

A multi-scale attention (MSA) network for remote sensing scene images classification

Guokai Zhang, Weizhe Xu, Wei Zhao, Chenxi Huang, Eddie Ng Yk, Yongyong Chen*, Jian Su*

Abstract—The remote sensing scene images classification has been of great value to civil and military fields. Deep learning models, especially the convolutional neural network (CNN), have achieved great success in this task, however, they may suffer from two challenges: firstly, the sizes of the category objects are usually different, but the conventional CNN extracts the features with fixed convolution extractor which could cause the failure in learning the multi-scale features; secondly, some image regions may not be useful during the feature learning process, therefore, how to guide the network to select and focus on the most relevant regions is crucially vital for remote sensing scene image classification. To address these two challenges, we propose a multi-scale attention network (MSA-Network), which integrates a multi-scale (MS) module and a channel and position attention (CPA) module to boost the performance of the remote sensing scene classification. The proposed MS module learns multi-scale features by adopting various sizes of sliding windows from different depths' layers and receptive fields. The CPA module is composed of two parts: the channel attention (CA) module and the position attention (PA) one. The CA module learns the global attention features from channel-level, and the PA module extracts the local attention features from pixel-level. Thus, fusing both of those two attention features, the network is apt to focus on the more critical and salient regions automatically. Extensive experiments on UC Merced, AID, NWPU-RESISC45 datasets demonstrate that the proposed MSA-Network outperforms several state-of-the-art methods.

Index Terms—Remote sensing scene; multi-scale; attention; feature fusion

I. INTRODUCTION

The explosion of high-resolution remote sensing imaging technology has unleashed a veritable data deluge in investigating the land-use and land-cover scenes [9], [16], [19], [18]. Especially, the recognition and classification of the remote sensing scene images have been of great value to civil and military fields, due to its plentiful spatial and semantic information. Specifically, in this classification task, it uses pixel-based [3] or different levels of features to identify and label the images based on the image contents. However, the

remote sensing scene image usually contains complicated ground objects with different spatial distributions, such as roads, buildings, and rivers, and that makes it difficult to classify the specific category of the whole scene image.

To tackle this challenge, during the past years, many works have been proposed for scene classification. The early attempts for this task mainly focused on the hand-crafted features such as texture, shape, color, and spatial representations which are combined and fed into a classifier for prediction. For example, Yang et al. [65] designed a category of image descriptors based on saliency for remote sensing scene images classification. Luo et al. [40] proposed to extract the radiometric, Gaussian wavelet, Gabor, and Gray level co-occurrence matrix features with different spatial resolutions for indexing of remote sensing scene images. To extract more spatial features, the work in [77] developed an effective approach by fusing the local and global spatial features with multi-scale feature learning mode. Meanwhile, Huang et al. [24] utilized the patch-based multi-scale local binary pattern features and a Fisher vector for remote sensing scene images classification. Some other previous works also used the color features [49], [31], [57], [51], histogram of oriented gradients features [6] [7] [41], and bag-of-visual-words [74] [52] [12] for classification. However, the above methods based on the hand-crafted features may yield unsatisfactory performance since they require subjective and empirical feature definition and selection. Furthermore, these features are usually low-level and mid-level, which may become limited and inadequate for the complex remote sensing scene image feature learning.

Since the remarkable performance of the convolutional neural network (CNN) in many computer vision fields [61], [76], [25], [5], [73], [22], [17], [45], many researchers also used CNN to extract high-level features from the remote sensing scene images [15], [20]. Compared with the hand-crafted based feature learning, CNN could extract more semantic features by deep network layers in an end-to-end learning manner. Especially the feature representations from the deeper layers could provide more abstract and semantic information, which are more applicable to the remote sensing scene image classification task. There are many baseline CNN models such as AlexNet [26], VggNet [54], GoogleNet [56], ResNet [13], which have achieved impressive performance for the natural image classification on ImageNet dataset [27]. Inspired by these works, many works tried to explore the huge power of CNN for remote sensing scene image classification. For example, Han et al. [11] used the AlexNet as the backbone architecture and incorporated the spatial pyramid pooling layer to learn the multi-scale information of the remote sensing

Guokai Zhang, School of Optical-Electrical and Computer Engineering, University of Shanghai for Science and Technology, Shanghai China.

Weizhe Xu, School of Computer Science, The University of Manchester, Manchester UK.

Wei Zhao, School of Software Engineering, Tongji University, Shanghai China.

Chenxi Huang, School of Informatics, Xiamen University, Xiamen China.

Eddie Ng Yk, School of Mechanical and Aerospace Engineering, Nanyang Technological University, Singapore.

Yongyong Chen*, The School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen China.

Jian Su*, The School of Computer and Software, Nanjing University of Information Science and Technology, Nanjing China.

Corresponding authors: Yongyong Chen, Jian Su.

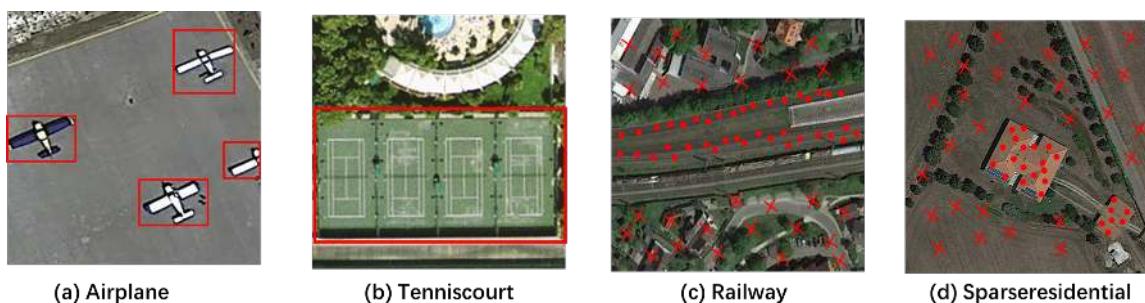


Fig. 1: The size of the target category airplane in scene image (a) is quite different from that of the target category tennis court in scene image (b). In the image (c) and (d), the red '.' denotes the target category key regions, and the red 'x' represents the useless regions of other objects, it is obvious that not all the image regions are useful during the feature learning process.

scene image. Castelluccio et al. [3] explored the CaffeNet [28] and GoogleNet [56] with widely different network settings for this classification task, in [50], [43], the deeper network architecture such as the ResNet was used to learn more semantic features. Meanwhile, some other hand-crafted features such as the texture, color, and Local Binary Patterns (LBP) [1] were combined with the deep features to further improve the classification performance. The previous works in [4], [69], [58], [39], [23] proposed different feature fusion strategies with various network structures to explore the effectiveness of different fusion modes. To alleviate the influence of limited remote sensing scene images, transfer learning or pre-trained methods [3], [4], [32], [21], [10], [34], [58], [44], [23] which pre-train the network on large amounts of datasets such as the ImageNet had been used for the improvement of the classification performance.

Although these works have achieved promising performance for remote sensing scene classification, they may encounter two main challenges. Firstly, the category objects of remote sensing scene images usually have different sizes as shown in Figure 1 (a) and (b), in which the size of the target category airplane in scene image (a) is quite different from that of the tennis court in (b). However, the conventional CNN extracts the features with a fixed convolution extractor which could be a handicap to learn the multi-scale features; secondly, as shown in Figure 1 (c) and (d), the red '.' denotes the target category key regions, and the red 'x' represents the useless regions of other objects, it is obvious that not all the image regions are useful during the feature learning process. Therefore, how to guide the network to select and focus on the most relevant regions is crucially vital for remote sensing scene image classification. To address these two challenges, in this paper, we propose a multi-scale attention network (MSA-Network) to achieve remote sensing scene images classification tasks. The proposed MSA-Network uses ResNet [13] as the backbone network and integrates a multi-scale (MS) module and a channel and position attention (CPA) module to further boost the classification performance. Inspired by the previous work [56], the designed MS module extracts multi-scale features from different receptive fields with various sizes of sliding windows. Besides, since different depths of layers may contain pyramidal scale features, we add the MS module behind each stage's last residual block to extract the multi-scale features

hierarchically. The CPA modules consist of two parts, i.e., the channel attention (CA) module, and position attention (PA) module, respectively. During the feature learning process, the CA module extracts the attention features from channel-level globally, while the PA module learns the attention features from pixel-level locally. By integrating those two attention features, it could guide the network to focus on more informative and critical regions globally and locally. The main contributions of this paper can be summarized as follows:

- A novel multi-scale module has been proposed to improve the network ability to capturing multi-scale features during the feature learning process.
- We design a channel and position attention module, which guides the network to focus on the informative critical regions globally and locally.
- Extensive experiments on three (UC Merced, AID, NWPU-RESISC45) datasets demonstrate our proposed MSA-Network has achieved competitive results over other state-of-the-art methods.

The rest of the paper is organized as follows. Section 2 introduces the proposed MSA-Network. The extensive experiments and results are presented in Section 3. In Section 4, we conduct a qualitative analysis of our designed model. In Section 5, we summarize this paper and give a brief discussion of the future work.

II. THE PROPOSED MSA-NETWORK

In this section, we give a detailed description of the designed model. As illustrated in Figure 2, the main backbone of our model is based on ResNet, which has achieved great success in many computer visions tasks. Meanwhile, for better extracting more multi-scale and discriminative features, we propose an MS module and a CPA module. Specifically, we add the MS module after each residual block to extract the multi-scale features hierarchically. Since the deeper layer contains more high-level and semantic features, we integrate the CPA module after the last MS module to focus on more informative and critical regions. For different depths of ResNet backbones, we list the parameters and network settings in Table I. The more detailed introductions of the proposed modules will be discussed in the following sections.

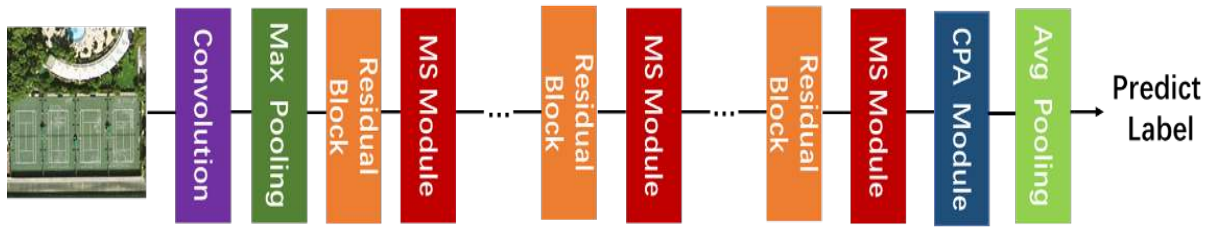


Fig. 2: An overview of the proposed architecture for remote sensing scene images classification. The main backbone of our network is based on ResNet, and we integrate the MS module and the CPA module with the network. The MS module is added behind each stage’s last residual block to extract the multi-scale features hierarchically. The CPA module is added behind the last MS module to guide the feature learning process to focus on more informative and critical regions. For different depth of ResNet backbones, we list detailed parameters and network settings in Table I.

TABLE I: Parameters setting of the network

layer name	18-layer	34-layer	50-layer	101-layer
Conv1	$7 \times 7, 64, \text{stride } 2$			
Pooling	$3 \times 3, \text{max pool, stride } 2$			
Conv2	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
MS module layer				
Conv2	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$
MS module layer				
Conv2	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$
MS module layer				
Conv2	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
MS module layer				
CPA module layer				
Output-layer	average pool, fc, softmax			

A. Multi-Scale Feature Learning

The remote sensing scene images usually compose of complex and diverse objects in the real world, and the main category of remote sensing scene object is often with various sizes. However, the conventional ResNet usually uses the fixed size convolutional layers to extract local features which could be a handicap for the network to learn multi-scale features. Thus, inspired by the previous work in [56], we use an MS feature learning module as a basic unit of the network to enhance the multi-scale feature learning. The detailed structure of the proposed MS module is illustrated in Figure 3. Denote the F_{i-1} as the input feature from the previous layer, and F_i as the output feature from the designed module. Instead of directly passing F_{i-1} into the next layer, we first apply multi-

scale feature learning with different convolution kernel sizes. Here, we use $\{1 \times 1, 3 \times 3, 6 \times 6, 9 \times 9\}$ as the basic units to generate four scale-level features. Next, an interlaced feature learning strategy is adopted for aggregating more contextual multi-scale information from the input features. Specifically, for the sizes of 1×1 and 3×3 kernels, they aim to learn more precise and subtle information. While for the sizes of 6×6 and 9×9 kernels, they are more likely to extract global and large-scale representations. The different sizes of convolution operations f_s can be formulated as:

$$f_s = \sum_{s \in S} (k_s * F_{i-1} + b_s), \quad (1)$$

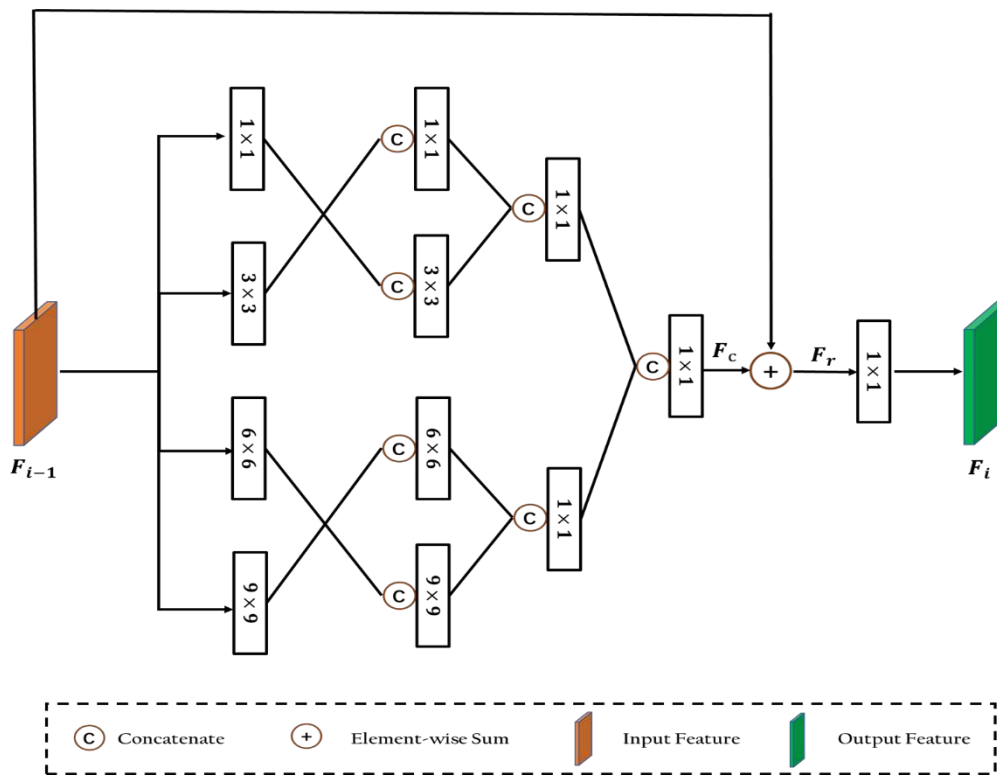


Fig. 3: The detailed structure of the MS module, F_{i-1} and F_i are the input and output feature maps of the i th layer, respectively.

where k_s is the kernel for the scale of $s \in \{1, 3, 6, 9\}$, and b_s is the bias for f_s . After the interlaced feature learning, two 1×1 convolution layers are utilized to aggregate the global and local features and then input to another 1×1 convolution layer to gain the squeezed feature F_c . Subsequently, we perform an element-wise sum F_c with F_{i-1} by residual learning to further improve the convergence ability of the network, and the output feature maps F_r can be given as:

$$F_r = F_c + F_{i-1}. \quad (2)$$

After that, the final F_i is gained by applying a 1×1 convolution operation on F_r . Notably, inspired by the fact that different depths of layers could contain pyramidal scale features [33], for our MAS-Network, we insert our MS module behind the last residual block of each stage to learn the scale-relevant feature pyramidally, and the detailed inserted position of this module is illustrated in Table I.

B. Channel and Position Attention Modules

Since the remote sensing scene images are captured from an overhead view, they usually contain complex and diverse objects. Thus, many objects are not useful for the image classification task. To handle this problem, we use the attention mechanism to guide the feature learning process to focus on more informative and critical regions. Conventionally, the attention mechanism has two categories: one is the hard attention which restricts the regions to 0 or 1, the other one is the soft attention which calculates the weight of the specific region. In this paper, we use the soft attention mechanism with two attention modules to select and learn more global and

local representations. Denote the output feature map from the last residual block as $F \in \mathbb{R}^{C \times H \times W}$, where C , H , W are the channel, height, and width of the output feature, separately. F_C is the channel attention feature map, and F_P is the position attention feature map. In the following subsections, we will give a detailed introduction to these two attention modules.

1) *Channel Attention Module*: Our channel attention module is inspired by the previous work squeeze and channel excitation block [29] and the spatial pyramid pooling [14]. The feature map of each channel contains different global and semantic responses, which is essential for remote sensing scene images understanding. Thus, in order to improve the ability to learn discriminative features from the channel level, we design a CA module to encode a wider of contextual and global representations from the channel level. The detailed structure of the proposed CA module is illustrated in Figure 4. For the input feature map $F \in \mathbb{R}^{H \times W \times C}$, we first perform three max-pooling operations with sizes of $\{2 \times 2, 4 \times 4, 8 \times 8\}$ to learn the contextual and global information, and it could be formulated as $\{F_{sp}, p \in \{2, 4, 8\}\}$. Then, a SE Block is utilized on F_{sp} to enable the network to focus on the most salient representations globally from the channel-level. Here, we denote each single F_{sp}^j as the j th channel feature map of F_{sp} where $j \in [1, C]$. In SE Block, a global-average pooling over F_{sp}^j is used to generate the global weight a_j of the j th channel. And the operation of the global-average pooling for the j th channel can be formulated as :

$$a_j = \frac{1}{H \times W} \sum_x \sum_y F_{sp}^j(x, y). \quad (3)$$

For better improving the generalization of the module, a fully connected (FC) layer is applied to the weight vector a ($a = \{a_1, a_2, \dots, a_j, \dots, a_C\}$) with a ReLU operator δ . Then, another FC layer with sigmoid activation σ is used to normalize outputs of the previous layer. These two FCs operations are defined as:

$$A_W = \sigma(W_2 \delta(W_1 a)), \quad (4)$$

where $W_1 \in \mathbb{R}^{C \times \frac{C}{r}}$ and $W_2 \in \mathbb{R}^{\frac{C}{r} \times C}$ are the respective weights of the two FCs, and the gained value of A_W represents the importance weight of channel feature map from the global-average level. Additionally, the variable of r is the bottleneck of the channel excitation, and we set it 4 empirically. Next, the channel attention output of \hat{F}_{sp} is gained by applying an multiplication between A_W and F_{sp} , which could be formulated as:

$$\hat{F}_{sp} = A_W \cdot F_{sp}. \quad (5)$$

After that, we concatenate those three channel attention outputs and apply a convolution layer with a size of 1×1 to squeeze the channel value to C . Therefore, the final output of channel attention module \hat{F} is calculated as:

$$F_C = Conv(Concat(\hat{F}_{s2}, \hat{F}_{s4}, \hat{F}_{s8})) \quad (6)$$

With the designed channel attention module, it could guide the network to focus on the globally crucial representations, which further improve the classification performance.

2) *Position Attention Module*: Different from the CA module, the goal of our PA module is to extract more subtle features from the pixel level. The detailed structure of the PA module is illustrated in Figure 5. In order to learn the importance of each feature pixel position, we first apply two position operations $\tau_M(\cdot)$ and $\tau_A(\cdot)$, which are utilized to calculate the max-value and average-value of each feature pixel position of the whole feature channels. After that, two sigmoid activations are employed to those two obtained features and gained the position weighted feature map F_M and F_A , respectively. Then, we multiply the input feature F with F_M and F_A to gain the position enhanced feature F'_M and F'_A , separately. And it could be formulated as follow :

$$F'_M = F \cdot F_M, F'_A = F \cdot F_A \quad (7)$$

Finally, the output feature map of the PA module is obtained by concatenating the F'_M and F'_A , and then applying a (1×1) convolution layer to aggregate those two features, and the aggregation process can be defined as:

$$F_P = Conv(Concat(F'_M, F'_A)) \quad (8)$$

3) *Channel and Local Feature Fusion*: In our designed model, the channel attention module could generate more global representations over the whole feature map, while the position attention module extracts the location attention feature from the pixel level. To incorporate more discriminative features from those two modules, we explore three combination methods to boost the remote sensing scene classification performance.

(1) **Concatenation** The concatenation operation is to concatenate two features at the same location x, y along the specific channel direction d into a vector:

$$y_{concat}^{x,y,d} = concat(F_C^{x,y,d}, F_P^{x,y,d}), \quad (9)$$

where $d \in \{0, 1, 2\}$, and $1 \leq x \leq H$, $1 \leq y \leq W$. $y_{concat}^{x,y,d}$ is the output feature by concatenation operation. Although the concatenation operation increases the complexity with stacking more channels, there is no information loss during the fusion process.

(2) **Addition** The addition operation is to compute the sum of two input features at the same location. The addition operation is formulated as:

$$y_{add}^{x,y,c} = F_C^{x,y,c} + F_P^{x,y,c}, \quad (10)$$

where $1 \leq x \leq H$, $1 \leq y \leq W$, $1 \leq c \leq C$, and $y_{add}^{x,y,c}$ is the output feature by addition operation.

(3) **Nonlinear Fusion** The strategy of nonlinear fusion is similar to the concatenation fusion, except that we apply an FC layer with the nonlinear activation function to each input feature before the concatenation operation:

$$y_{nonlinear}^{x,y,d} = concat(W_C(F_C^{x,y,d}), W_P(F_P^{x,y,d})). \quad (11)$$

Here, $1 \leq x \leq H$, $1 \leq y \leq W$ and $d \in \{0, 1, 2\}$ is the specific channel direction, W_C and W_P is the corresponding weight of the FC layer, $y_{nonlinear}^{x,y,c}$ is the output feature by nonlinear fusion operation.

III. EXPERIMENTS AND RESULTS

In this section, we evaluate the proposed MSA-Network model on three different aerial scene public datasets. We first give a brief description of all datasets, then the detailed experimental setting of the model is introduced. Finally, we conduct extensive experiments to further validate the performance of the proposed MSA-Network model.

A. Introduction of All Datasets

1) *UC Merced Dataset*: The UC Merced dataset, which was extracted from the USGS National Map Urban Area Imagery collection, consists of 21 land-use scene classes, as shown in Figure 6. There are 100 images for each class, and each image has a 1-foot spatial resolution, measuring 256×256 pixels.

2) *AID Dataset*: The AID dataset is a new large-scale image dataset that contains sample images collected from Google Earth imagery. There are 10000 images within 30 classes, as shown in Figure 6. Each image of the AID dataset has a spatial resolution ranging from 8m to half a meter, and there are about 220 to 400 samples measuring 600×600 pixels in each class. For each class of the AID dataset, all the images are selected from different countries and regions around the world. As a consequence, the dataset has high intra-classes diversities.

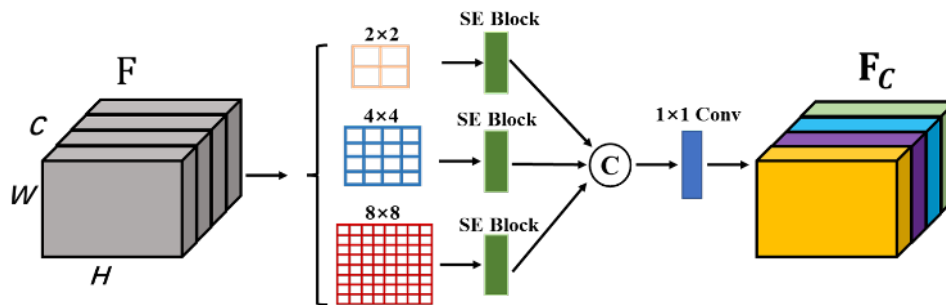


Fig. 4: The structure of our CA module, F denotes the input feature map with dimension of $C \times H \times W$, and F_C is the final output from the CA module.

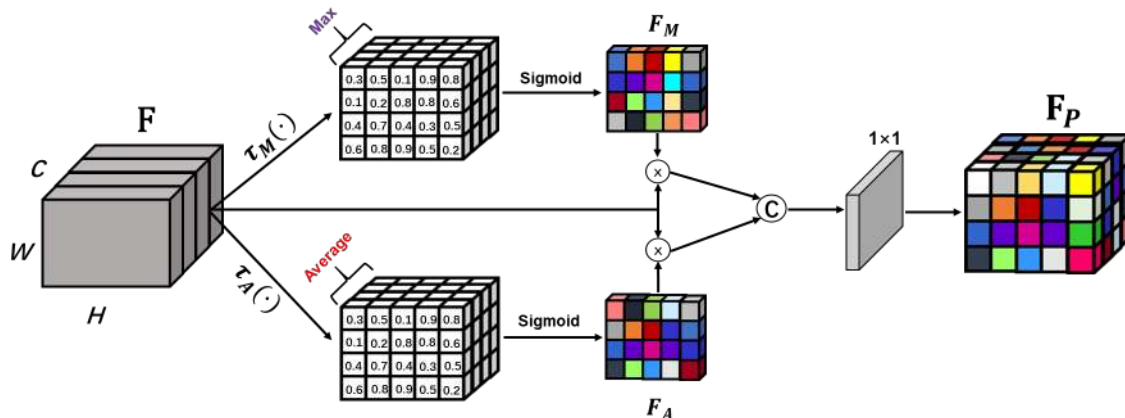


Fig. 5: The structure of PA module, F denotes the input feature maps with the dimension of $C \times H \times W$, F_P is the final output from PA module.

3) *NWPU-RESISC45 Dataset*: The NWPU-RESISC45 dataset, created by Northwestern Polytechnical University (NWPUP), contains 31500 images for remote sensing image scene classification. The dataset consists of 45 scene classes, and there are 700 images measuring 256×256 pixels in each class. The 45 classes are shown in Figure 6.

From Figures 6, we can see that the sizes of the categories are various with different distributions, which could be a challenge for the conventional neural network to learn the semantic features.

B. Implementation Details

In this paper, we use Tensorflow as the basic framework to implement the proposed MSA-Network model. The main backbone of our network is based on the ResNet [13]. Following the experimental setting of [71] [59] [47], for the UC Merced dataset, we use 80% and 50% as the training sample ratios, respectively. For the AID dataset, we use 50% and 20% as the training sample ratios, respectively. For the NWPU-RESISC45 dataset, we use 20% and 10% training sample ratios, respectively. During the training process, we use Adam optimization [30] to make the network convergent. The initial learning rate is set as 1.0×10^{-4} , and then we reduce it by 0.1 factor after the val-loss not improved 10 epochs. The batch-size of the proposed MSA-Network is 32, and we use real-time data augmentation on the training dataset such as the

random rotation, flip, and cropping. We pre-train our model on the ImageNet dataset, and dropout is adopted to avoid the network being over-fitting. We conduct our experiments on Ubuntu 14.04 operating system with 64GB memory, and an NVIDIA GTX 1080 graphics processing unit (GPU) has been used to accelerate the training process.

C. Evaluation Protocol

The overall accuracy (OA) and confusion matrix are selected as the criterion to evaluate the performance of the MSA-Network. The value of OA is calculated by the ratio between the correct numbers of classified images and the total number of images. It is one of the basic evaluation metrics for the classification task, and the higher value of OA denotes the more accurate classification performance. The confusion matrix is a table layout that describes the errors and confusion of each class. Each row of the table denotes the predicted category instance, and the column represents the actual category instance. In addition, the diagonal of the table are the numbers of all the classes correctly classified. The final confusion matrix is calculated by the best classification result of all the training ratios for each dataset.

D. The Performance of Different Network Depths

Different depths of networks could extract various representations. For a deeper network, it is liable to learn more

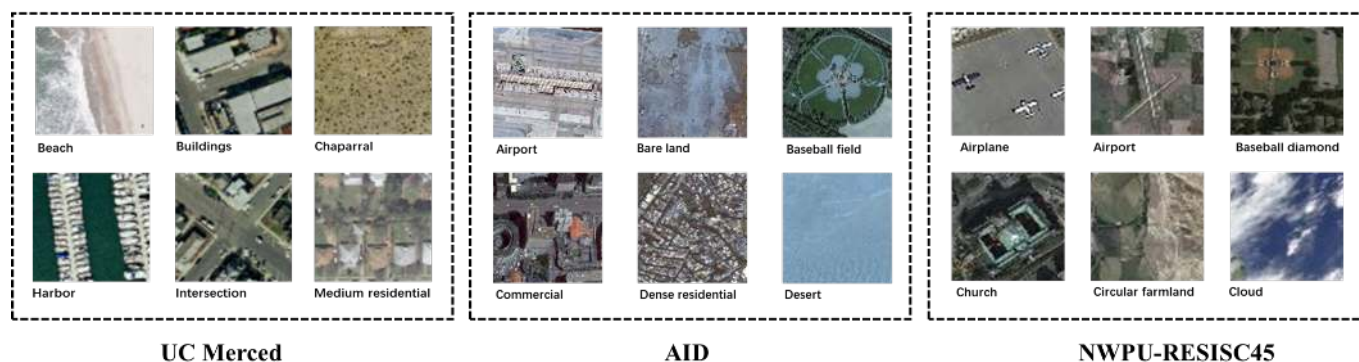


Fig. 6: Example images of UC Merced/AID/NWPU-RESISC45 dataset.

semantic and high-level features, for a shallower network, it would extract tinier and more detailed information. Thus, in this section, we first explore the performance of four network depths (18-layer, 34-layer, 50-layer, 101-layer) on the three datasets with our designed architecture. The depth of the network is deepened by adding the convolution and residual blocks as previous work [13]. The results are reported in Figure 7. We use the 80%, 50%, and 20% of training data for the UC Merced dataset, AID dataset, and NWPU-RESISC45 dataset, respectively. On the UC Merced dataset, the proposed MSA-Network achieves ($97.73 \pm 0.45\%$), ($97.82 \pm 0.38\%$), ($98.59 \pm 0.42\%$), and ($98.96 \pm 0.21\%$) classification accuracy for 18-layers, 34-layers, 50-layers, and 101-layers, separately. On the AID dataset, the results are ($94.86 \pm 0.51\%$), ($95.32 \pm 0.21\%$), ($95.89 \pm 0.33\%$), and ($96.01 \pm 0.43\%$) for 18-layers, 34-layers, 50-layers, and 101-layers, respectively. The performance of 18-layers, 34-layers, 50-layers, and 101-layers on NWPU-RESISC45 dataset is ($90.31 \pm 0.42\%$), ($90.78 \pm 0.51\%$), ($91.67 \pm 0.43\%$), and ($92.52 \pm 0.23\%$). The experimental results demonstrate that the deeper network layers could achieve better performance compared with shallower network layers. The best result on three datasets is by using the 101-layers.

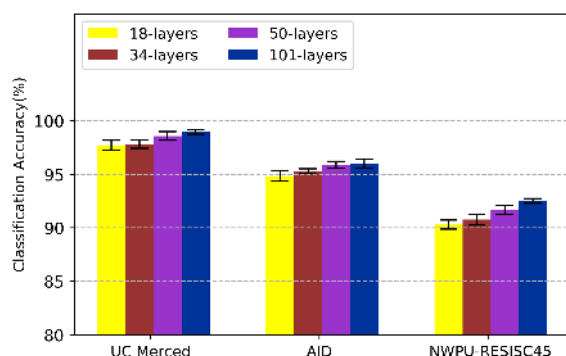


Fig. 7: The comparison results with different network depths.

E. The Effectiveness of Different Modules

In this paper, we propose an MS module and two attention modules (CA module and PA module) to extract the multi-scale, global, and local representations more effectively. In

order to evaluate the effectiveness of different modules, we conduct relevant experiments on those three datasets. The detailed comparison result is shown in Table II. From the experimental results, we can observe that adding one of the MS, CA, and PA module could efficiently improve the classification performance of the network on different datasets, that further proves the effectiveness of our proposed modules. Overall, the best performance is achieved by using the MS module, CA module, and PA module simultaneously, and it further proves that the multi-scale, channel, and position attention features are crucially important in the remote sensing images classification task.

F. Compared with Different Fusion Methods

The designed CA module generates the features from the channel level, and the PA module generates the features from pixel level. Thus, fusing both of the two features could enrich the discrimination of the feature learning model. In this section, we explore the effectiveness of three fusion methods (Addition, Concatenation, Nonlinear Fusion) on all datasets. The detailed results are shown in Figure 8. On the UC Merced dataset, the OA is ($98.37 \pm 0.39\%$), ($98.96 \pm 0.21\%$), and ($97.25 \pm 0.64\%$) for addition, concatenation, and nonlinear fusion methods, respectively. Meanwhile, the OA on the AID dataset is ($95.77 \pm 0.51\%$), ($96.01 \pm 0.43\%$), and ($95.32 \pm 0.37\%$). On NWPU-RESISC45 dataset, it achieves ($92.01 \pm 0.46\%$), ($92.52 \pm 0.23\%$), and ($91.79 \pm 0.56\%$) with the three fusion methods, respectively. The performance on these three datasets demonstrates that the concatenation fusion method performs best compared with the other two fusion ones. Furthermore, the addition method is slightly better than the nonlinear fusion method, and it could be the reason that some spatial correspondences between the two generated attention features are decreased at the FC layers. Since the concatenation fusion method performs best, we use it as the final fusion method of our model.

G. Classification of the UC Merced Dataset

In order to evaluate the effectiveness of the proposed MSA-Network, we compare it with other state-of-the-art methods on the UC Merced dataset. The detailed results are presented in Table III. In [3], [4], [32], [21], [10], [34], [58], [44],

TABLE II: The effectiveness of different modules on the UC Merced, AID and NWPU-RESISC45 dataset

Dataset	MS module	CA module	PA module	Accuracy(%)
UC Merced	×	×	×	96.88 ± 0.32
	✓	×	×	97.09 ± 0.22
	×	✓	×	97.21 ± 0.63
	×	×	✓	97.13 ± 0.32
	✓	✓	×	98.16 ± 0.13
	✓	×	✓	97.83 ± 0.25
	×	✓	✓	98.35 ± 0.42
	✓	✓	✓	98.96 ± 0.21
AID	×	×	×	93.02 ± 0.27
	✓	×	×	93.13 ± 0.26
	×	✓	×	94.33 ± 0.37
	×	×	✓	93.87 ± 0.44
	✓	✓	×	95.31 ± 0.51
	✓	×	✓	94.73 ± 0.25
	×	✓	✓	95.75 ± 0.32
	✓	✓	✓	96.01 ± 0.43
NWPU-RESISC45	×	×	×	89.83 ± 0.31
	✓	×	×	90.01 ± 0.34
	×	✓	×	90.32 ± 0.45
	×	×	✓	90.89 ± 0.22
	✓	✓	×	91.25 ± 0.37
	✓	×	✓	91.74 ± 0.65
	×	✓	✓	92.03 ± 0.33
	✓	✓	✓	92.52 ± 0.23

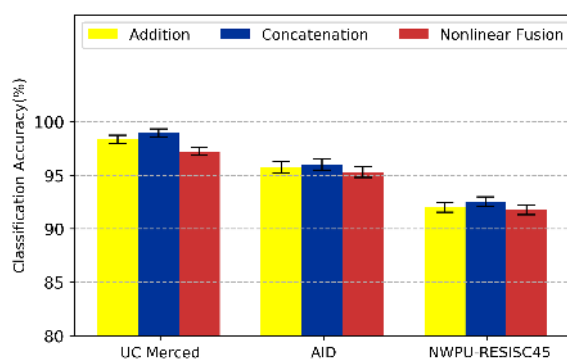


Fig. 8: The comparison results with different fusion methods.

[23], researchers used pre-trained CNNs to boost the performance of the models. In [74], [47], [64], [23], bag-of-visual-words (BoVW) was applied to land-use classification. In [68], [58], the extreme learning machine classifier was applied for final classification with fused features. In [34], [66], they constructed a pyramid image to characterize both the photometric and geometric aspects of an image. In [1],

[60], [23], images with the LBP based texture, local, and global information was employed to convey the comprehensive message to the models. From this table, we can see that the proposed MSA-Network has gained (98.96 ± 0.21%) and (97.80 ± 0.33%) classification accuracy for 80% and 50% training ratios, separately. It is noteworthy that the overall accuracy of 80% training ratio is better than the 50% training data ratio, it can be explained that with more training data samples, the network could learn more high-level and abstract features which further enhance the classification performance. Overall, the experimental results demonstrate that our MSA-Network could gain state-of-the-art classification performance on the UC Merced dataset. The detailed confusion matrix with the training ratio of 80% is illustrated in Figure 9. The confusion matrix result shows that most of the categories have achieved 100% classification accuracy.

H. Classification of the AID Dataset

Table IV shows the results with other state-of-the-art methods on the AID dataset. The FACNN in [39] combined the

TABLE III: Comparisons of overall accuracy (%) on the UC Merced dataset

Method	80% Training Ratio	50% Training Ratio
SPCK [66]	73.14	/
SPCK+ [66]	76.05	/
SPCK++ [66]	77.38	/
BoVW [3]	76.81	/
BoVW + SCK [64]	77.71	/
Saliency-Guided [72]	82.72 ± 1.18	/
VGG-VD-16 [59]	95.21 ± 1.20	94.14 ± 0.69
CaffeNet [59]	95.02 ± 0.81	93.98 ± 0.67
SRSCNN [37]	95.57	/
CNN-ELM [58]	95.62	/
salM ³ LBP-CLM [60]	95.75 ± 0.80	94.21 ± 0.75
TEX-Net-LF [1]	97.72 ± 0.54	96.91 ± 0.36
Appearance-based [48]	96.05 ± 0.62	/
LGFBOW [74]	96.88 ± 1.32	/
GoogLeNet [59]	94.31 ± 0.89	92.70 ± 0.60
Fusion by addition [4]	/	97.42 ± 1.79
DCA fusion [4]	/	96.90 ± 0.09
CCP-net [47]	97.52 ± 0.97	/
Two-Stream Fusion [68]	98.02 ± 1.03	96.97 ± 0.75
DSFATN [10]	98.25	/
Deep CNN Transfer [21]	98.49	/
ResNet [50]	98.50 ± 1.40	/
Aggregate strategy [34]	97.40	96.25
SHHTFM [75]	98.33 ± 0.98	/
VGG-VD16+AlexNet [32]	98.81 ± 0.38	/
FACNN [39]	98.81 ± 0.24	/
CTFCNN [23]	98.44 ± 0.58	/
SAFF [2]	97.02 ± 0.78	/
SSRL [63]	94.05 ± 0.12	/
GBNet [55]	98.57 ± 0.48	97.05 ± 0.19
InceptionV3+Xception[46]	97.86 ± 0.59	96.60 ± 0.65
MSA-Network	98.96 ± 0.21	97.80 ± 0.33

feature learning, feature aggregation, and classifier for joint training. In [67], they utilized a multi-level fusion method, which can make a judgment by incorporating different levels' information. In [4], combined with the SIFT feature, the deep learning feature can get a discriminative image presentation which overcoming the scale and rotation variability. The result in Table IV shows that our model has achieved (93.53 ± 0.21%) classification accuracy with 20% training data, and (96.01 ± 0.43%) classification accuracy with 50% training data, which is the best performance compared with all competing methods. Overall, the result demonstrates that our MSA-Network could achieve state-of-the-art classification performance on the AID dataset. The detailed confusion matrix

with a fixing training ratio of 50% is illustrated in Figure 10. The result shows that most of the categories have gained high classification accuracy, which is above 90%.

I. Classification of the NWPU-RESISC45 Dataset

For the NWPU-RESISC45 dataset, the result over other state-of-the-art methods is shown in Table V. In [36], the lecture focused on the four new loss functions to achieve better performance. Based on a residual network and dense convolutional networks in [42], it achieved a competitive result. Using a metric learning regularization term, the Siamese

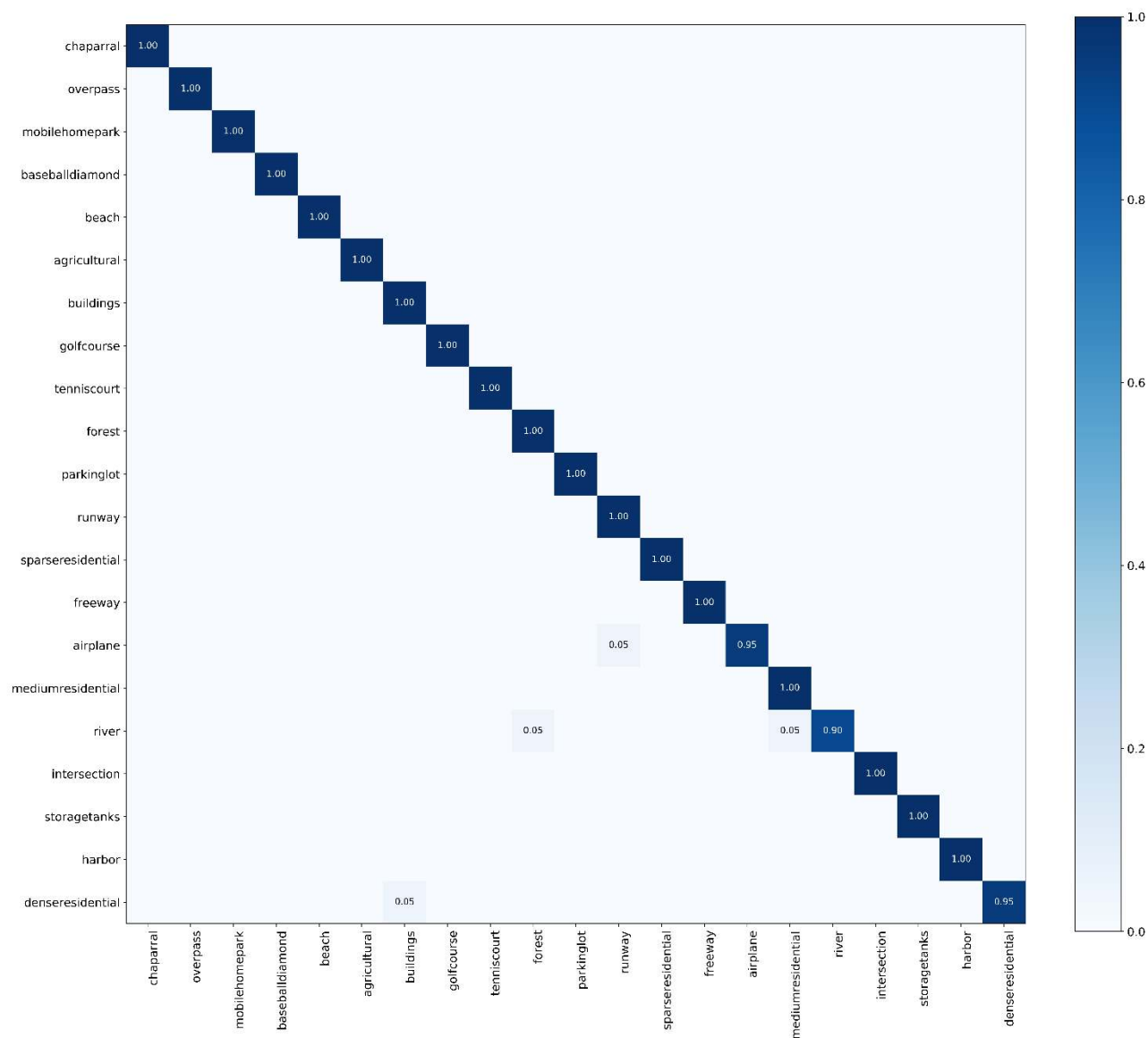


Fig. 9: The confusion matrix on UC Merced dataset with training data of 80%.

CNN proposed in [35] was also very robust and effective. it notes that the proposed architecture in [35] replaces the final fully connected layer with a convolutional layer to predict the corresponding label of each class. The result demonstrates that our method achieves the highest classification accuracy compared with other state-of-the-art methods. For the 20% training data, our model achieves $(93.52 \pm 0.21\%)$ classification accuracy, and for the 10% training data, it has gained the classification accuracy of $(90.38 \pm 0.17\%)$. Especially, compared with the other state-of-the-art network architectures (e.g., AlexNet, GoogLeNet, VGG-16), the proposed MSA-Network could gain better performance. It could be explained from two folds: first, our designed network is based on ResNet, which has a deeper network architecture to learn more high-

level and semantic features; second, with our designed MS module and CPA module, the model is apt to extract more multi-scale and salient features, which further improve the classification performance of the model. The confusion matrix of 20% training data is illustrated in Figure 11, indicating that our model achieves competitive classification results in most of the categories.

IV. QUALITATIVE ANALYSIS

A. Multi-scale Feature Map Visualization

In this section, we visualize some examples of the MS module feature maps on different datasets in Figure 12. Since

TABLE IV: Comparisons of overall accuracy (%) on the AID dataset

Method	50% Training Ratio (%)	20% Training Ratio
BoVW [60]	/	78.66 ± 0.52
MS-CLBP+FV [60]	/	86.48 ± 0.27
CaffeNet [59]	89.53 ± 0.31	86.86 ± 0.47
GoogLeNet [59]	86.39 ± 0.55	83.44 ± 0.40
VGG-VD-16 [59]	89.64 ± 0.36	86.59 ± 0.29
DCA fusion [4]	89.71 ± 0.33	/
<i>salM³LBP - CLM</i> [60]	89.76 ± 0.45	/
TEX-Net-LF [1]	95.73 ± 0.16	/
Fusion by addition [4]	91.87 ± 0.36	/
Bidirectional adaptive feature fusion [38]	93.56	/
Multilevel fusion [67]	94.17 ± 0.32	/
Two-Stream Fusion [69]	94.58 ± 0.25	92.32 ± 0.21
FACNN [39]	95.45 ± 0.11	/
CTFCNN [23]	94.91 ± 0.24	/
LCNN-BFF [53]	94.62 ± 0.16	91.66 ± 0.48
SAFF [2]	93.83 ± 0.28	90.25 ± 0.29
GBNet [55]	95.52	/
MSA-Network	96.01 ± 0.43	93.53 ± 0.21

TABLE V: Comparisons of overall accuracy (%) on the NPWU-RESISC45 dataset

Method	20% Training Ratio (%)	10% Training Ratio
LASC-CNN[70]	84.30	81.37
Two-Stream Fusion [69]	83.16 ± 0.18	80.22 ± 0.22
BoCF [8]	84.32 ± 0.17	82.65 ± 0.31
AlexNet [8]	85.16 ± 0.18	81.22 ± 0.19
GoogLeNet [8]	86.02 ± 0.18	82.57 ± 0.12
VGG-16 [8]	90.36 ± 0.18	87.15 ± 0.45
Triple networks [36]	92.33 ± 0.20	/
ResNet [42]	91.96 ± 0.71	89.24 ± 0.75
Siamese_ResNet [35]	92.28	/
LCNN-BFF [53]	91.73 ± 0.17	86.53 ± 0.15
SAFF [2]	87.86 ± 0.14	84.38 ± 0.19
SSRL [63]	83.12 ± 0.26	/
MF ² Net [62]	92.73 ± 0.21	90.17 ± 0.25
MSA-Network	93.52 ± 0.21	90.38 ± 0.17

different depths of layers could contain various semantic features, we select the feature maps from the four depth layers of our designed architecture. Here, we denote the 'MS_module layer_1', 'MS_module layer_2', 'MS_module layer_3', and 'MS_module layer_4' as the layers from the first, second, third, and fourth of the MS module layer, respectively. From the visualization results, we can see that the MS module could extract diverse features through different sizes of filters,

and the shallower depth layer tends to extract more edge and subtle representations, while the deeper depth layer is liable to learn more high-level and abstract features. Overall, with the designed MS module, the designed network could encode more discriminative features from different levels, which further improve the classification performance of the model.

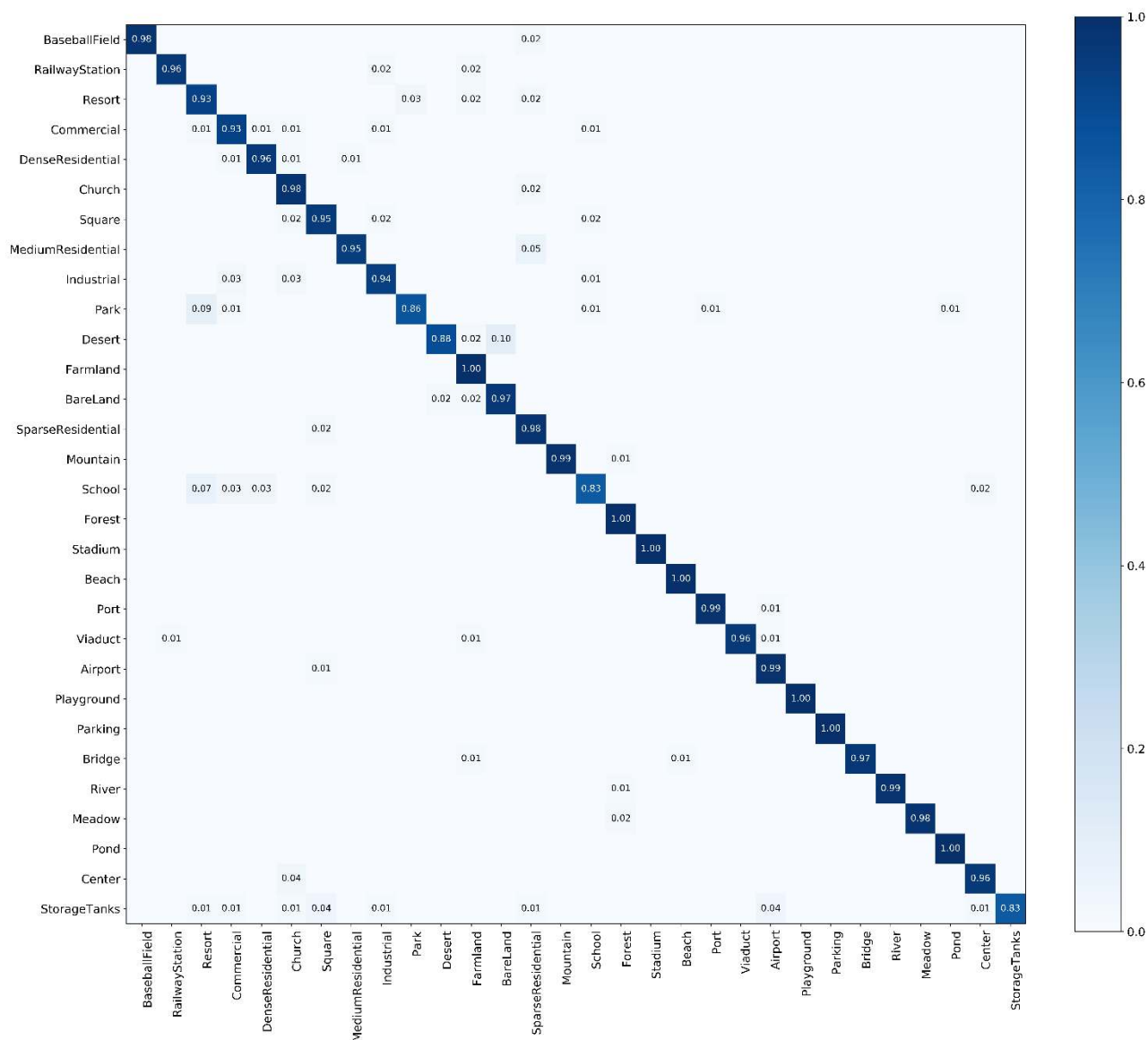


Fig. 10: The confusion matrix on AID dataset with training data of 50%.

B. Visualization of the Attention Map

In this section, we show some examples of attention maps on different datasets. For better comparison, we compare our model with the ResNet architecture (depth with 101). The detailed visualization results are illustrated in Figure 13. Compared with the ResNet visualization results, it is obvious that our designed model can attend to the more crucial regions, especially on the category relevant regions. Meanwhile, the designed model could provide more diverse and detailed region features which could further enhance the final classification performance. We suggest that adding the CA module, could guide the model to focus on more global regions while adding the PA module, tends to guide the model

to focus on more subtle regions, thus, fusing both the CA module and PA module could efficiently improve the model ability to learn more crucial and salient features.

V. CONCLUSION

In this paper, we propose an MSA-Network to handle the remote sensing scene image classification task, in which the MSA-Network uses the ResNet and incorporates an MS module and CPA module to further improve the performance of the designed architecture. The proposed MS module extracts multi-scale features from different receptive fields with various sizes of sliding windows. Moreover, we add the MS module behind each stage's last residual block to extract the multi-scale features hierarchically. The CPA module consists of two

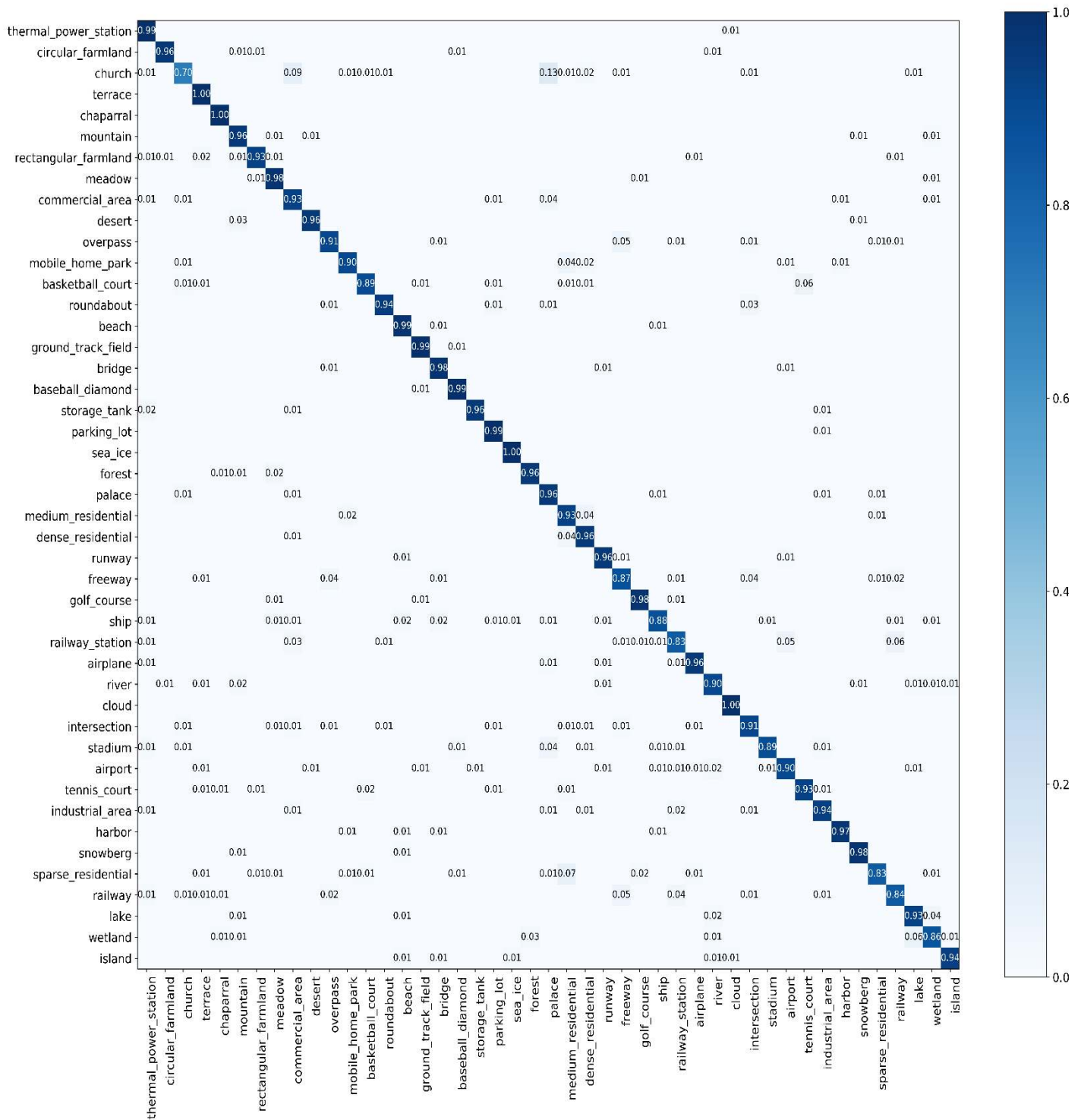


Fig. 11: The confusion matrix on NPWU-RESISC45 dataset with training data of 20%.

parts: the CA module and the PA module. The CA module aims to extract the attention features from channel-level globally, while the PA module learns the attention features

from pixel-level locally. With those two attention modules, the proposed MSA-Network could focus on more informative and critical regions globally and locally. Experimental

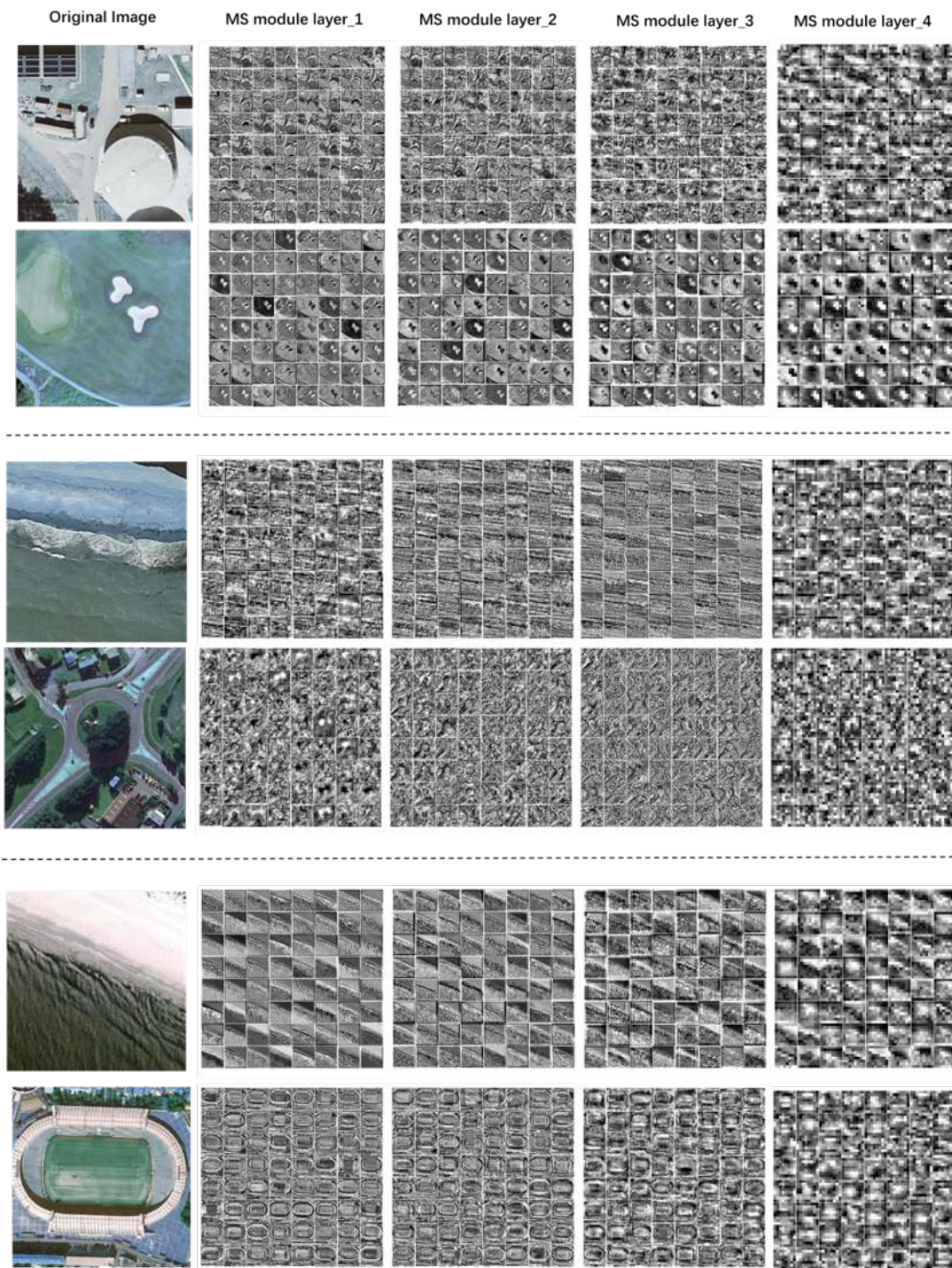


Fig. 12: Some examples of the MS module feature maps on different datasets, the images on the top two rows are samples from the UC Merced dataset, the images on the middle two rows are samples from the NPWU-RESISC45 dataset, and the images on the bottom two rows are the samples from AID dataset. The 'MS_module layer_1', 'MS_module layer_2', 'MS_module layer_3', and 'MS_module layer_4' as the layers from the first, second, third, and fourth of the MS module layer, respectively.

results on UC Merced, NPWU-RESISC45, and AID datasets demonstrate that the proposed MSA-Network could achieve better performance over several state-of-the-arts on the overall classification accuracy. Furthermore, we also conduct relevant experiments on our designed modules to further analyze the effectiveness of each module. In future work, we will try

to explore the designed MS module and CPA module with deeper network architecture such as DenseNet to validate the effectiveness of our modules. Meanwhile, some other techniques such as ensemble more models with early fusion or late fusion may further improve the classification accuracy.

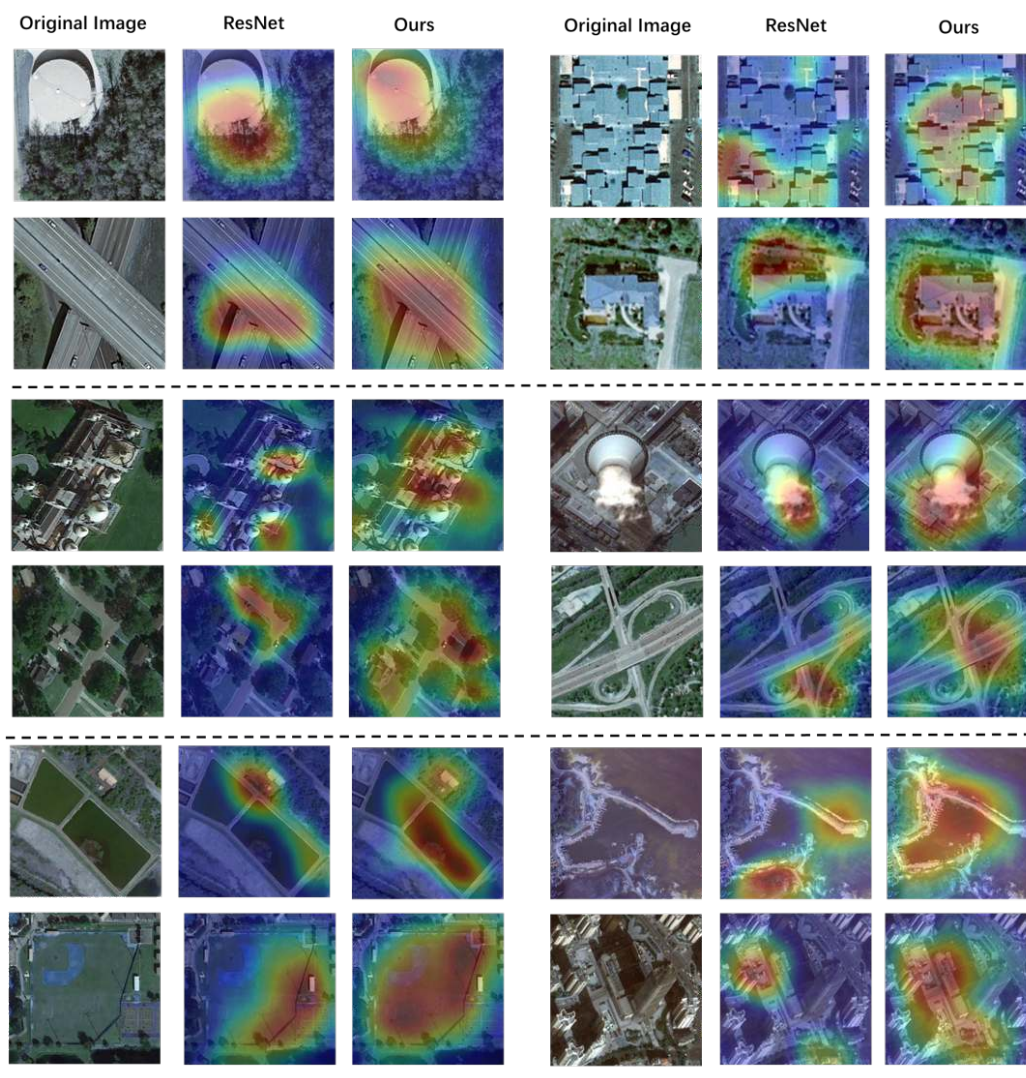


Fig. 13: Some examples of the attention maps on different datasets, the images on the top two rows are samples from the UC Merced dataset, the images on the middle two rows are samples from the NPWU-RESISC45 dataset, and the images on the bottom two rows are the samples from the AID dataset.

ACKNOWLEDGMENT

The authors would like to thank the paper reviewers. We also thank the associate editor for supplying valuable comments and suggestions. This work was funded by the Shenzhen College Stability Support Plan under Grant GXWD20201230155427003-20200824113231001, and the General Program of National Natural Science Foundation of China (NSFC) under Grant 62106063 and 62102259.

REFERENCES

- [1] Rao Muhammad Anwer, Fahad Shahbaz Khan, Joost van de Weijer, Matthieu Molinier, and Jorma Laaksanen. Binary patterns encoded convolutional neural networks for texture recognition and remote sensing scene classification. *ISPRS journal of photogrammetry and remote sensing*, 138:74–85, 2018.
- [2] Ran Cao, Leyuan Fang, Ting Lu, and Nanjun He. Self-attention-based deep feature fusion for remote sensing scene classification. *IEEE Geoscience and Remote Sensing Letters*, 18(1):43–47, 2020.
- [3] Marco Castelluccio, Giovanni Poggi, Carlo Sansone, and Luisa Verdoliva. Land use classification in remote sensing images by convolutional neural networks. *arXiv preprint arXiv:1508.00092*, 2015.
- [4] Souleyman Chaib, Huan Liu, Yanfeng Gu, and Hongxun Yao. Deep feature fusion for vhr remote sensing scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(8):4775–4784, 2017.
- [5] Liang-Chieh Chen, Yi Yang, Jiang Wang, Wei Xu, and Alan L Yuille. Attention to scale: Scale-aware semantic image segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3640–3649, 2016.
- [6] G. Cheng, J. Han, L. Guo, Z. Liu, S. Bu, and J. Ren. Effective and efficient midlevel visual elements-oriented land-use classification using vhr remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 53(8):4238–4249, 2015.
- [7] G. Cheng, P. Zhou, J. Han, L. Guo, and J. Han. Auto-encoder-based shared mid-level visual dictionary learning for scene classification using very high resolution remote sensing images. *Computer Vision Iet*, 9(5):639–647, 2015.
- [8] Gong Cheng, Zhenpeng Li, Xiwen Yao, Lei Guo, and Zhongliang Wei. Remote sensing image scene classification using bag of convolutional features. *IEEE Geoscience and Remote Sensing Letters*, 14(10):1735–1739, 2017.
- [9] Luis Gómez-Chova, Devis Tuia, Gabriele Moser, and Gustau Camps-Valls. Multimodal classification of remote sensing images: A review and future directions. *Proceedings of the IEEE*, 103(9):1560–1584, 2015.
- [10] Xi Gong, Zhong Xie, Yuanyuan Liu, Xuguo Shi, and Zhuo Zheng. Deep

- salient feature based anti-noise transfer network for scene classification of remote sensing imagery. *Remote Sensing*, 10(3):410, 2018.
- [11] X. Han, Y. Zhong, L. Cao, and L. Zhang. Pre-trained alexnet architecture with pyramid pooling and supervision for high spatial resolution remote sensing image scene classification. *Remote Sensing*, 9(8):848, 2017.
- [12] S. Hao, S. Xian, H. Wang, L. Yu, and X. Li. Automatic target detection in high-resolution remote sensing images using spatial sparse coding bag-of-words model. *IEEE Geoscience Remote Sensing Letters*, 9(1):109–113, 2011.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *IEEE*, 2016.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1904–1916, 2015.
- [15] Danfeng Hong, Lianru Gao, Jing Yao, Bing Zhang, Antonio Plaza, and Jocelyn Chanussot. Graph convolutional networks for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [16] Danfeng Hong, Xin Wu, Pedram Ghamisi, Jocelyn Chanussot, Naoto Yokoya, and Xiao Xiang Zhu. Invariant attribute profiles: A spatial-frequency joint feature extractor for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 58(6):3791–3808, 2020.
- [17] Danfeng Hong, Naoto Yokoya, Jocelyn Chanussot, Jian Xu, and Xiao Xiang Zhu. Learning to propagate labels on graphs: An iterative multitask regression framework for semi-supervised hyperspectral dimensionality reduction. *ISPRS journal of photogrammetry and remote sensing*, 158:35–49, 2019.
- [18] Danfeng Hong, Naoto Yokoya, Jocelyn Chanussot, and Xiao Xiang Zhu. An augmented linear mixing model to address spectral variability for hyperspectral unmixing. *IEEE Transactions on Image Processing*, 28(4):1923–1938, 2018.
- [19] Danfeng Hong, Naoto Yokoya, Nan Ge, Jocelyn Chanussot, and Xiao Xiang Zhu. Learnable manifold alignment (lema): A semi-supervised cross-modality learning framework for land cover and land use classification. *ISPRS journal of photogrammetry and remote sensing*, 147:193–205, 2019.
- [20] Danfeng Hong, Naoto Yokoya, Gui-Song Xia, Jocelyn Chanussot, and Xiao Xiang Zhu. X-modalnet: A semi-supervised deep cross-modal network for classification of remote sensing data. *ISPRS Journal of Photogrammetry and Remote Sensing*, 167:12–23, 2020.
- [21] Fan Hu, Gui-Song Xia, Jingwen Hu, and Liangpei Zhang. Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery. *Remote Sensing*, 7(11):14680–14707, 2015.
- [22] Yan Hu. Design and implementation of abnormal behavior detection based on deep intelligent analysis algorithms in massive video surveillance. *Journal of Grid Computing*, 18(2):227–237, 2020.
- [23] Hong Huang and Kejie Xu. Combing triple-part features of convolutional neural networks for scene classification in remote sensing. *Remote Sensing*, 11(14):1687, 2019.
- [24] Longhui Huang, Chen Chen, Wei Li, and Qian Du. Remote sensing image scene classification using multi-scale completed local binary patterns and fisher vectors. *Remote Sensing*, 8(6):483, 2016.
- [25] Zhenhua Huang, Jinyi Tang, Guangxu Shan, Juan Ni, Yunwen Chen, and Cheng Wang. An efficient passenger-hunting recommendation framework with multitask deep learning. *IEEE Internet of Things Journal*, 6(5):7713–7721, 2019.
- [26] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5mb model size. 2016.
- [27] D. Jia, D. Wei, R. Socher, L. J. Li, L. Kai, and F. F. Li. Imagenet: A large-scale hierarchical image database. *Proc of IEEE Computer Vision Pattern Recognition*, pages 248–255, 2009.
- [28] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678, 2014.
- [29] H. Jie, S. Li, and S. Gang. Squeeze-and-excitation networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [30] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *Computer Science*, 2014.
- [31] L. Lei and S. Ning. Object-oriented classification of high-resolution remote sensing image using structural feature. *IEEE*, 2010.
- [32] Erzhu Li, Junshi Xia, Peijun Du, Cong Lin, and Alim Samat. Integrating multilayer features of convolutional neural networks for remote sensing scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(10):5653–5665, 2017.
- [33] T. Y. Lin, P. Dollar, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. *IEEE Computer Society*, 2017.
- [34] Na Liu, Lihong Wan, Yu Zhang, Tao Zhou, Hong Huo, and Tao Fang. Exploiting convolutional neural networks with deeply local description for remote sensing image classification. *IEEE access*, 6:11215–11228, 2018.
- [35] X. Liu, Y. Zhou, J. Zhao, R. Yao, Bing Liu, and Y. Zheng. Siamese convolutional neural networks for remote sensing scene classification. *IEEE Geoscience and Remote Sensing Letters*, pages 1–5, 2019.
- [36] Y. Liu and C. Huang. Scene classification via triplet networks. *IEEE Journal of Selected Topics in Applied Earth Observations Remote Sensing*, 11(1):220–237, 2018.
- [37] Yanfei Liu, Yanfei Zhong, Feng Fei, Qiqi Zhu, and Qianqing Qin. Scene classification based on a deep random-scale stretched convolutional neural network. *Remote Sensing*, 10(3):444, 2018.
- [38] Xiaoqiang Lu, Weijun Ji, Xuelong Li, and Xiangtao Zheng. Bidirectional adaptive feature fusion for remote sensing scene classification. *Neurocomputing*, 328:135–146, 2019.
- [39] Xiaoqiang Lu, Hao Sun, and Xiangtao Zheng. A feature aggregation convolutional neural network for remote sensing scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 57(10):7894–7906, 2019.
- [40] Bin Luo, Shujing Jiang, and Liangpei Zhang. Indexing of remote sensing images with different resolutions by multiple features. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 6(4):1899–1912, 2013.
- [41] M. L. Mekhali, F. Melgani, Y. Bazi, and N. Alajlan. Land-use classification with compressive sensing multifeature fusion. *IEEE Geoscience Remote Sensing Letters*, 12(10):2155–2159, 2015.
- [42] R. Minetto, M. P. Segundo, and S. Sarkar. Hydra: an ensemble of convolutional neural networks for geospatial land classification. *IEEE Transactions on Geoscience and Remote Sensing*, 2018.
- [43] Rodrigo Minetto, Maurício Pamplona Segundo, and Sudeep Sarkar. Hydra: An ensemble of convolutional neural networks for geospatial land classification. *IEEE Transactions on Geoscience and Remote Sensing*, 57(9):6530–6541, 2019.
- [44] L. Na, X. Lu, L. Wan, H. Hong, and F. Tao. Improving the separability of deep features with discriminative convolution filters for rsi classification. *ISPRS International Journal of Geo-Information*, 7(3):95, 2018.
- [45] Mohammad Nauman, Hafeez Ur Rehman, Gianfranco Politano, and Alfredo Benso. Beyond homology transfer: Deep learning for automated annotation of proteins. *Journal of Grid Computing*, 17(2):225–237, 2019.
- [46] Biserka Petrovska, Eftim Zdravevski, Petre Lameski, Roberto Corizzo, Ivan Štajduhar, and Jonatan Lerga. Deep learning for feature extraction in remote sensing: A case-study of aerial scene classification. *Sensors*, 20(14):3906, 2020.
- [47] K. Qi, Q. Guan, Y. Chao, F. Peng, S. Shen, and H. Wu. Concentric circle pooling in deep convolutional networks for remote sensing scene classification. *Remote Sensing*, 10(6):934, 2018.
- [48] Kunlun Qi, Chao Yang, Qingfeng Guan, Huayi Wu, and Jianya Gong. A multiscale deeply described correlations-based model for land-use scene classification. *Remote Sensing*, 9(9):917, 2017.
- [49] Jad Santos, Oab Penatti, and Rds Torres. Evaluating the potential of texture and color descriptors for remote sensing image retrieval and classification. In *VISAPP 2010 - Proceedings of the Fifth International Conference on Computer Vision Theory and Applications, Angers, France, May 17-21, 2010 - Volume 2*, 2010.
- [50] Grant J Scott, Matthew R England, William A Starms, Richard A Marcum, and Curt H Davis. Training deep convolutional neural networks for land-cover classification of high-resolution imagery. *IEEE Geoscience and Remote Sensing Letters*, 14(4):549–553, 2017.
- [51] Z. Shao, W. Zhou, L. Zhang, and J. Hou. Improved color texture descriptors for remote sensing image retrieval. *Journal of Applied Remote Sensing*, 8(1):083584, 2014.
- [52] X. Sheng, F. Tao, D. Li, and S. Wang. Object classification of aerial images with bag-of-visual words. *IEEE Geoscience Remote Sensing Letters*, 7(2):366–370, 2010.
- [53] Cuiqing Shi, Tao Wang, and Liguang Wang. Branch feature fusion convolution network for remote sensing scene classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 13:5194–5210, 2020.

- [54] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *Computer Science*, 2014.
- [55] Hao Sun, Siyuan Li, Xiangtao Zheng, and Xiaoqiang Lu. Remote sensing scene classification by gated bidirectional network. *IEEE Transactions on Geoscience and Remote Sensing*, 58(1):82–96, 2019.
- [56] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [57] P. Tokarczyk, J. D. Wegner, S. Walk, and K. Schindler. Features, color spaces, and boosting: New insights on semantic classification of remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 53(1):280–295, 2014.
- [58] Qian Weng, Zhengyuan Mao, Jiawen Lin, and Wenzhong Guo. Land-use classification via extreme learning classifier based on deep convolutional features. *IEEE Geoscience and Remote Sensing Letters*, 14(5):704–708, 2017.
- [59] G. S. Xia, J. Hu, F. Hu, B. Shi, and L. Zhang. Aid: A benchmark data set for performance evaluation of aerial scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(7):3965–3981, 2017.
- [60] Xiaoyong, Bian, Chen, Chen, Long, Tian, Qian, and Du. Fusing local and global features for high-resolution scene classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 10(6):2889–2901, 2017.
- [61] Gaowei Xu, Min Liu, Zhuofu Jiang, Dirk Söfker, and Weiming Shen. Bearing fault diagnosis method based on deep convolutional neural network and random forest ensemble learning. *Sensors*, 19(5):1088, 2019.
- [62] Kejie Xu, Hong Huang, Yuan Li, and Guangyao Shi. Multilayer feature fusion network for scene classification in remote sensing. *IEEE Geoscience and Remote Sensing Letters*, 17(11):1894–1898, 2020.
- [63] Peiyao Yan, Feng He, Yajie Yang, and Fei Hu. Semi-supervised representation learning for remote sensing image classification based on generative adversarial networks. *IEEE Access*, 8:54135–54144, 2020.
- [64] Y. Yang and S. Newsam. Bag-of-visual-words and spatial extensions for land-use classification. In *Sigspatial International Conference on Advances in Geographic Information Systems*, page 270, 2010.
- [65] Yi Yang and Shawn Newsam. Comparing sift descriptors and gabor texture features for classification of remote sensed imagery. In *2008 15th IEEE international conference on image processing*, pages 1852–1855. IEEE, 2008.
- [66] Y. Yi and S. Newsam. Spatial pyramid co-occurrence for image classification. In *IEEE International Conference on Computer Vision*, pages 1465–1472, 2011.
- [67] Y. Yu and F. Liu. Aerial scene classification via multilevel fusion based on deep convolutional neural networks. *IEEE Geoscience and Remote Sensing Letters*, PP(2):1–5, 2018.
- [68] Y. Yu and F. Liu. A two-stream deep fusion framework for high-resolution aerial scene classification. *Computational Intelligence and Neuroscience*, 2018(4-5):1–13, 2018.
- [69] Yunlong Yu and Fuxian Liu. A two-stream deep fusion framework for high-resolution aerial scene classification. *Computational intelligence and neuroscience*, 2018, 2018.
- [70] Baohua Yuan, Shijin Li, and Ning Li. Multiscale deep features learning for land-use scene recognition. *Journal of Applied Remote Sensing*, 12(1):015010, 2018.
- [71] F Zhang, B. Du, and L. Zhang. Scene classification via a gradient boosting random convolutional network framework. *IEEE Transactions on Geoscience Remote Sensing*, 54(3):1793–1802, 2016.
- [72] Fan Zhang, Bo Du, and Liangpei Zhang. Saliency-guided unsupervised feature learning for scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 53(4):2175–2184, 2014.
- [73] Shuai Zhao, Dong Ming Zhang, and Hong Wei Huang. Deep learning-based image instance segmentation for moisture marks of shield tunnel lining. *Tunnelling and Underground Space Technology*, 95:103156, 2020.
- [74] Q. Zhu, Y. Zhong, B. Zhao, and G. S. Xia. Bag-of-visual-words scene classifier with local and global features for high spatial resolution remote sensing imagery. *IEEE Geoscience and Remote Sensing Letters*, 2016.
- [75] Qiqi Zhu, Yanfei Zhong, Siqi Wu, Liangpei Zhang, and Deren Li. Scene classification based on the sparse homogeneous-heterogeneous topic feature model. *IEEE Transactions on Geoscience and Remote Sensing*, 56(5):2689–2703, 2018.
- [76] Tehseen Zia and Saad Razzaq. Residual recurrent highway networks for learning deep sequence prediction models. *Journal of Grid Computing*, 18(1):169–176, 2020.
- [77] Jinyi Zou, Wei Li, Chen Chen, and Qian Du. Scene classification using local and global features with collaborative representation fusion. *Information Sciences*, 348:209–226, 2016.