

# Self-Attention-Based Deep Feature Fusion for Remote Sensing Scene Classification

Ran Cao, Leyuan Fang<sup>ID</sup>, *Senior Member, IEEE*, Ting Lu<sup>ID</sup>, *Member, IEEE*,  
and Nanjun He<sup>ID</sup>, *Student Member, IEEE*

**Abstract**—Remote sensing scene classification aims to assign automatically each aerial image a specific semantic label. In this letter, we propose a new method, called self-attention-based deep feature fusion (SAFF), to aggregate deep layer features and emphasize the weights of the complex objects of remote sensing scene images for remote sensing scene classification. First, the pretrained convolutional neural network (CNN) model is applied to extract the abstract multilayer feature maps from the original aerial imagery. Then, a nonparametric self-attention layer is proposed for spatial-wise and channel-wise weightings, which enhances the effects of the spatial responses of the representative objects and uses the infrequently occurring features more sufficiently. Thus, it can extract more discriminative features. Finally, the aggregated features are fed into a support vector machine (SVM) for classification. The proposed method is experimented on several data sets, and the results prove the effectiveness and efficiency of the scheme for remote sensing scene classification.

**Index Terms**—Deep feature fusion, deep learning, remote sensing scene classification, self-attention.

## I. INTRODUCTION

IN recent years, with the rapid development of the satellite imaging technology, a large number of high-resolution remote sensing images available, which contain a lot of scene-semantic information, enable us to measure the Earth surface with detailed structures [1]. For aerial-image interpretation, a lot of studies have been done [2]–[5], and remote sensing scene classification is the fundamental problem that aims to label a great amount of images without human intervention. However, the complex contents in the high-resolution images are also a big challenge for remote sensing scene classification.

In general, remote sensing scene-classification methods can be divided into three main classes: low-level methods, mid-level methods, and high-level methods [1]. Low-level methods [e.g., the scale-invariant feature transform (SIFT) [6]]

are supposed to use handcrafted features such as spectral, texture, and structure, which are the low-level visual features of the images, to distinguish the different classes of scenes. Mid-level methods rely on mid-level visual representation, which is encoded by the local features extracted from manually designed models [e.g., codebook in the bag of visual words (BoVW) [7]]. High-level methods are usually based on constructing deep learning networks [e.g., convolutional neural network (CNN) [8]] to obtain efficient visual information for scene classification. Recently, deep learning methods have achieved impressive performance on many computer-vision tasks, such as image classification [9] and image retrieval [10]. It has shown that deep learning can obtain deep and abstract feature representations, which has achieved some state-of-the-art results for remote sensing scene classification. While the low-level and mid-level methods just use the shallow local features that cannot be effectively adapted to various scenes, the high-level method can learn more discriminative and abstract semantic features through deep learning networks. Thus, deep-learning-based algorithms are often used in remote sensing classification. According to the characteristics of remote sensing scene data sets, training from scratch is possible to cause overfitting due to the small number of training images, and so, most works tend to use a pretrained model to extract features, such as VGG-VD16 [11], and AlexNet [8], which can easily get deep features without high computational cost. However, how to use efficiently the features obtained from the pretrained models is an urgent problem.

Meanwhile, attention mechanism, which essentially comes from the human visual attention mechanism, has shown good performance on various tasks, such as speech recognition [12] and object detection [13]. Attention mechanism is first proposed in neural machine translation [14], which helps to provide different weights for different areas of input and extract more meaningful information. Because of the complexity of the remote sensing scene data sets, it is hard to set weights for a specific region, which can be applied to various data sets. Building a network for weight computing can be costly and time-consuming, which may cause overfitting on the small-scale remote sensing scene data sets. Therefore, finding a way to use the attention mechanism is difficult but attractive.

In this letter, we propose a new method, called self-attention-based deep feature fusion (SAFF), to aggregate the multilayer features extracted from off-the-shelf models for remote sensing scene classification. The proposed framework includes three parts. The first part is to extract the convolutional features from the pretrained model. Second, SAFF is designed to make full use of the spatial and channel responses and emphasize the importance of the features that do not occur

Manuscript received September 9, 2019; revised December 28, 2019; accepted January 16, 2020. This work was supported in part by the National Natural Science Fund of China under Grant 61922029 and Grant 61520106001, in part by the Science and Technology Plan Project Fund of Hunan Province under Grant CX2018B171 and Grant 2018TP1013, and in part by the Natural Science Foundation of Hunan Province under Grant 2019JJ50079. (Ran Cao and Leyuan Fang are co-first authors.) (Corresponding author: Ting Lu.)

The authors are with the College of Electrical and Information Engineering, Hunan University, Changsha 410082, China, and also with the Key Laboratory of Visual Perception and Artificial Intelligence of Hunan Province, Hunan University, Changsha 410082, China (e-mail: cr@hnu.edu.cn; fangleyuan@gmail.com; tingluhnu@gmail.com; henanjun@hnu.edu.cn).

Color versions of one or more of the figures in this letter are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/LGRS.2020.2968550

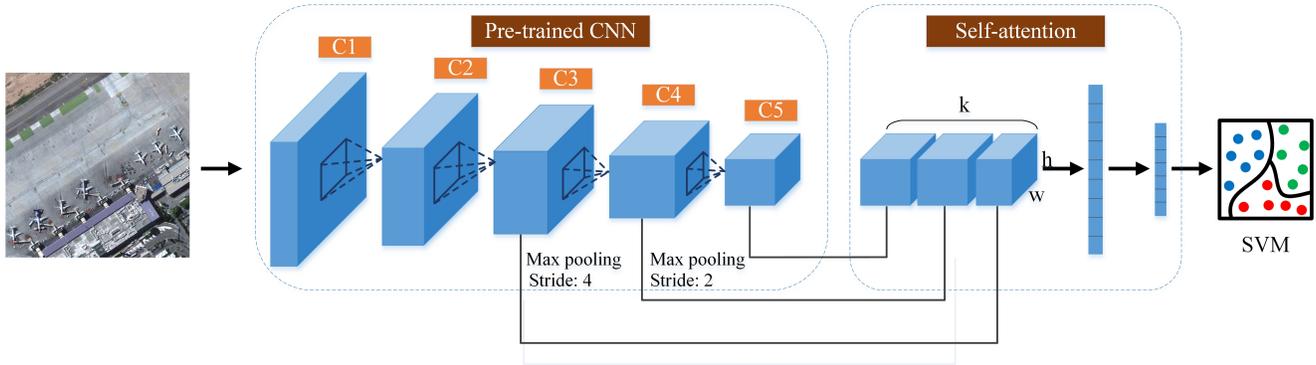


Fig. 1. Framework of the proposed SAFF for scene classification. It is composed of three parts: 1) multilayer convolutional feature extraction by a pretrained CNN, 2) SAFF, and 3) scene classification using SVM classifiers.

frequently by spatial-wise and channel-wise weightings, which converts the feature maps into aggregated vectors. Finally, we feed the vectors into a support vector machine (SVM) classifier with a linear kernel for scene classification. The proposed method associates the complex content with the semantic labels of the scene images without costly computation, which achieves great performance.

The rest of this letter is organized as follows. Section II briefly introduces the structure of the CNN. In Section III, the proposed method will be described. Section IV shows the results of the experiments on different data sets. Finally, we draw a conclusion and suggest some future works in Section V.

## II. CNNs

A CNN is an efficient deep learning network for extracting deep and abstract features, which transforms the input image into the final class scores [9]. The classic CNN model consists of convolutional layers, pooling layers, and fully connected (FC) layers.

### A. Convolutional Layers

The output of the convolutional layers is feature maps, where each element is obtained by multiple kernels. One kernel connects only to a local receptive field of the input feature maps, and its size can be manually adjusted. The parameters can be initialized by different ways, but different regions share the same weights, which is called the weight-sharing mechanism. Because of the receptive field and the weight-sharing mechanism, the parameters of the CNNs needed to be trained are much less than the FC networks. Calculating the dot product of the respective field and a set of weights will be sent to a nonlinear activation function to get the output feature maps.

### B. Pooling Layers

The pooling layers are usually placed behind the convolutional layers and used to reduce the spatial dimension of the feature maps. The downsampling is usually operated by average-pooling or max-pooling on a local region to reduce the computational cost.

### C. FC Layers

The FC layers are the last few layers to transform the pooled feature maps to a vector, and the last fully connected layer is a softmax layer that computes the scores for each class.

## III. PROPOSED METHOD

Fig. 1 illustrates the overall framework of the proposed SAFF for scene classification. As can be seen, the framework is made up of three parts: 1) multilayer convolutional feature extraction by a pretrained CNN; 2) SAFF; and 3) scene classification using linear SVM classifiers. The details of this architecture are described as follows.

### A. Multilayer Convolutional Feature Extraction

The CNN model can be simply explained by the underlying formula

$$f(X) = f_l(\dots f_2(f_1(X; w_1); w_2) \dots, w_l) \quad (1)$$

where  $X$  is the input image,  $w_1, w_2 \dots w_l$ , represent the weights of the respective layers, and  $f_l$  is the activation function of layer  $l$ . The output of each layer will become the input of the next layer, and, ultimately, the input image is transformed into a feature vector that can be classified by a specific classifier. The weights have been pretrained on other data set that is similar to the remote sensing scene images and adopting pretrained weights enables us to save a lot of calculation cost and time. The adopted CNN models are VGG-VD16 and AlexNet. In the proposed method, we only use the feature maps of three convolutional layers (conv3-3, conv4-3, and conv5-3) of VGG-VD16 and three convolutional layers (conv3, conv4, and conv5) of AlexNet. Since we do not use the FC layers, the size of the input images is arbitrary, and more information can be kept. According to different models, we extract the feature maps of three respective layers and implement the max-pooling operation on them to obtain feature maps with the same size as the last layers. Thus, the final number of feature maps is the sum of the three-layer feature maps.

### B. SAFF

With the pretrained CNN model, we have got a stacked  $K$ -dimensional feature maps of three convolutional layers. Next, SAFF, as shown in Fig. 2, is constructed for deep feature aggregation, which is made up of two steps: spatial-wise weighting and channel-wise weighting.

1) *Spatial-Wise Weighting*: Different from traditional sum-pooling [10], weighted sum-pooling is a transformation that gives spatial-wise weights to stress the characteristics of complex objects of the scene images. To get the weight

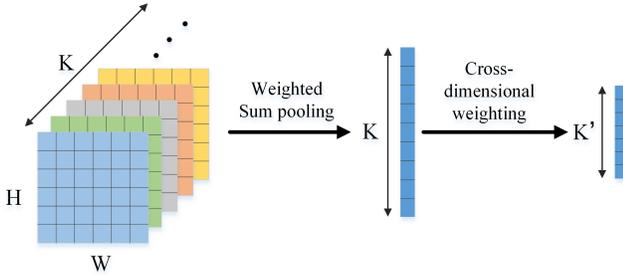


Fig. 2. SAFF.

map  $S'$ , we just use the ratio of each element in the concatenated feature maps to represent the importance of each pixel. Because it is complicated to find out the common and concrete region that is favorable for scene classification, we focus on the extracted feature maps. For the concatenated feature maps  $X \in \mathbb{R}^{H \times W \times K}$ , we calculate the sum of  $X$  on the third dimension to get the  $S \in \mathbb{R}^{H \times W}$ , which is shown as follows:

$$S = \sum_k X^k \quad (2)$$

where  $S$  is the sum of the concatenated feature map on the third dimension and  $X^k$  represents the feature map of channel  $k$ .

Then, we compute the ratio  $S_{xy}'$  of each element  $S_{xy}$  to the sum of  $S$  and replace  $S_{xy}$  with  $S_{xy}'$  to get a weight map  $S'$ , which will be separately applied to the location of each channel, as follows:

$$S_{xy}' = \left( \frac{S_{xy}}{(\sum_{m,n} S_{mn}^a)^{1/a}} \right)^{1/b} \quad (3)$$

where  $S_{xy}'$  is the weight of the spatial coordinate  $(x, y)$  and  $S_{xy}$  is the element of  $S$ .  $(m, n)$  represents the location on the feature maps.  $a, b$  are the parameters of normalization, which are separately set to be 0.5, 2 in the experiments.

Finally, we compute the output of weighted sum-pooling by the following equation:

$$\Phi = \sum_{y=1}^H \sum_{x=1}^W S_{xy}' f(x, y) \quad (4)$$

where  $\Phi$  is the output of weighted sum-pooling,  $H$  and  $W$  are the height and width of the aggregated feature maps,  $S_{xy}'$  is the spatial-wise weights, and  $f(x, y)$  is the value of the corresponding location  $(x, y)$  in the concatenated feature maps  $X$ . The spatial-wise weighting enhances the importance of the salient objects and decreases the effects of unsalient objects [15].

2) *Channel-Wise Weighting*: The second step is cross-dimension weighting [15]. As indicated in [15], the images of the common classes have highly correlated channel sparsity, the sparsity pattern of the channels contains discriminative information, and the sparsity of the feature maps is used to show the importance of infrequently occurring features. This step is also operated on the stacked convolutional feature map  $X \in \mathbb{R}^{H \times W \times K}$ . We first obtain the sparsity by calculating the property of nonzero elements of each channel, and  $\Xi_k$  is the sparsity of every channel, shown

as follows:

$$\Omega = \frac{1}{WH} \sum_{ij} \mathbb{1}[\lambda^{ij} > 0] \quad (5)$$

$$\Xi_k = 1 - \Omega_k \quad (6)$$

where  $\mathbb{1}(\cdot)$  is the indicator function—if  $\cdot$  is True, return 1; else, return 0. Through weighted sum-pooling, the channels with frequently occurring features are already activated. Therefore, we are supposed to improve the ratio of the channels with infrequently occurring features. The channel-wise weights are defined as the following equation:

$$\omega_k = \log \left( \frac{K\delta + \sum_h \Omega_h}{\delta + \Omega_k} \right) \quad (7)$$

where  $\omega_k$  is the weight of channel  $k$ ,  $K$  is the number of channels of the concatenated feature maps  $X$ , and  $\delta$  is a small constant if the formula is senseless, which is infinitely close to 0.

After the calculation of the weights, we apply the spatial-wise weighting maps to the location of each channel and compute the sum of each channel. Therefore, we get a feature vector that has the same number with channels. Then, we apply the cross-dimension weights to the feature vector. Finally, principal components analysis (PCA) whitening [16] is used to eliminate the redundant information and reduce the dimension of the feature vector, which is also a part of channel-wise weighting. Channel-wise weighting stresses the differences among the images by increasing the weights of infrequently occurring features.

### C. SVM Classification

Through the self-attention layer, the feature vectors with a dimension of  $K'$  are derived for final classification. Before scene image classification, normalization is used to avoid overfitting. For the whole data set, the common operation is applied. Then, we randomly select feature vectors as training samples to train a linear SVM, while the rest of the feature vectors are used to test the effectiveness of our method.

## IV. EXPERIMENT

### A. Data Sets

In order to evaluate the feasibility and effectiveness of our method, we experiment this method on several popular remote sensing scene data sets.

1) *UC Merced Land Use Data Set*: There are 2100 images in the UC Merced land use (UC) data set [7] labeled into 21 scene classes. Each class consists of 100 images with the size of  $256 \times 256$  pixels in the RGB space. The pixel spatial resolution is 1 ft.

2) *Aerial Image Dataset*: This data set has a number of 10000 images within 30 classes. Each class consists of images ranging from 220 up to 420. The size of each aerial image is fixed to be  $600 \times 600$  pixels in the RGB space. Aerial image dataset (AID) [1] has multiresolutions, and the pixel resolution changes from about 8 m to about half a meter.

3) *NWPU-RESISC45 Data Set*: The NWPU-RESISC45 (NWPU) data set [17] is made up of 31500 images divided into 45 scene classes. Each class consists of 700 images with the size of  $256 \times 256$  pixels in the RGB space. The spatial resolution of the most images changes from about 30 to 0.2 m per pixel.

## B. Experimental Setup

In our experiment, two pretrained CNN models, VGG-VD16 and AlexNet, are applied to extract the multilayer deep features. For different models, various feature maps are used for aggregation (conv3, conv4, and conv5 of AlexNet; conv3-3, conv4-3, and conv5-3 of VGG-VD16). To preserve the spatial information, we only use the outputs of the convolutional layers, and thus, it is unnecessary to adapt the size of our scene images to the fixed input size of the pretrained models. In order to evaluate the performance of our method, the experiments are repeated ten times with the training set consisting of randomly selected samples updated every time. Specifically, these data sets are divided to training samples and testing samples with different proportions. The training ratio of the UC data set is 80%. For the AID, the training ratios are set to 50% and 20%, while the training ratios are 20% and 10% in the NWPU data set. The final overall accuracy and standard deviation are determined by the average of ten-time results.

To validate the effectiveness of the proposed method, we implement some ablation experiments. First, we set the spatial-wise weights and channel-wise weights as 1, which is traditional sum-pooling, to aggregate the last three convolutional features maps. Here, these methods are called by the AlexNet + sum-pooling and VGG\_VD16 + sum-pooling, respectively. Second, our method is only applied on the last convolutional layer, i.e., AlexNet + single-layer SAFF and VGG\_VD16 + single-layer SAFF. Moreover, the proposed method will be compared with several classical ways for classification. The feature vectors obtained from the second FC layer of the pretrained models is fed into an SVM classifier for classification, which is regarded as the baseline. The discriminant correlation analysis (DCA) [18], the multiscale CNN (MCNN) [19], the spatial pyramid pooling (SPP)-Net [20], and the bag of convolutional feature (BoCF) [17] are applied to different data sets to get the best performance. All these experiments are conducted on a desktop with an Intel Core i7-8700K CPU at 3.70 GHz.

## C. Parameter Setting

In our method, only one parameter needs to be adjusted for the best performance: the dimension of the feature vectors,  $K'$ . Through the operation of SAFF, the feature vectors that have the common dimension with the channel of stacked feature maps are resized to  $K'$ -dimension vectors, while  $K'$  is confined between 0 and the number of the channels of the stacked feature maps. First, in order to get reasonable performance, the redundant information is ignored, and the feature vectors are simply set to be the maximum without reduction in the dimension. Then, we reduce the dimension step by step to compare the results of different dimensions of the feature vectors. The parameter  $K'$  is then analyzed on the different data sets. Fig. 3 shows the overall accuracy changing with the dimension of the feature vectors on the AID with the training ratios of 20% and 50%. The horizontal axis represents the dimension of the feature vector, and the vertical axis shows the overall accuracy of classification. In general, the overall accuracy is improved, while the length of the feature vector is increased. When the feature vectors are

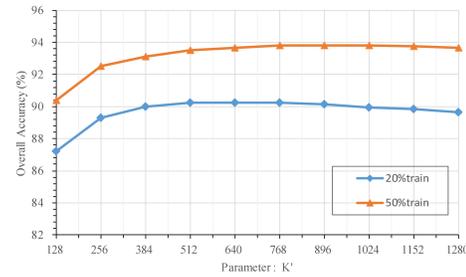


Fig. 3. Effect of different dimensions of feature vector on the AID with the training ratios of 20% and 50%.

TABLE I  
COMPARISON OF OVERALL CLASSIFICATION ACCURACY (%)  
WITH THE UC DATA SET

Method	OA
GoogleNet [1]	94.31±0.89
VGG_VD16 [1]	95.21±1.20
DCA [18]	96.90±0.77
MCNN [19]	96.66±0.90
SPP-Net [20]	96.67±0.94
AlexNet+ single-layer SAFF	89.67±1.39
VGG_VD16+ single-layer SAFF	88.79±0.87
AlexNet+ sum pooling [10]	94.10±0.93
VGG_VD16+ sum pooling [10]	91.67±1.40
AlexNet+ SAFF (ours)	96.13±0.97
VGG_VD16+ SAFF (ours)	<b>97.02±0.78</b>

TABLE II  
COMPARISON OF OVERALL CLASSIFICATION ACCURACY (%)  
WITH THE AID

Method	Training ratios	
	20%	50%
GoogleNet [1]	83.44±0.40	86.39±0.55
VGG_VD16 [1]	86.59±0.29	89.64±0.36
DCA [18]	-	89.71±0.33
Fusion by concatenation [18]	-	91.87±0.36
MCNN [19]	-	91.80±0.22
SPP-Net [20]	87.44±0.45	91.45±0.38
AlexNet+ single-layer SAFF	79.17±0.36	85.06±0.36
VGG_VD16+ single-layer SAFF	79.85±0.51	85.77±0.43
AlexNet+ sum pooling [10]	80.51±0.59	86.09±0.34
VGG_VD16+ sum pooling [10]	75.34±0.67	82.81±0.69
AlexNet+ SAFF (ours)	87.51±0.36	91.83±0.27
VGG_VD16+ SAFF (ours)	<b>90.25±0.29</b>	<b>93.83±0.28</b>

very short, due to insufficient information, the performance is not good. However, by increasing the length of the feature vector, the feature is more sufficient, and the precision greatly improves until it reaches a steady value varying in a small range. After reaching best results, the accuracy decreases a little due to some redundant information. As can be seen, when the dimension of the feature vector  $K'$  is 896, the performance is optimal.

## D. Experimental Results

Through a series of experiments, we make sure the suitable parameters for each data set while using different pretrained models, and show the overall accuracies of different methods in Tables I–III. In the tables, the performance of sum-pooling and SAFF on single convolutional features is not as good as the proposed method, which verifies that the SAFF and combining features of multiple layers contribute to final performance. When comparing with the baseline, our method gets better

TABLE III  
COMPARISON OF OVERALL CLASSIFICATION ACCURACY (%)  
WITH THE NWPU DATA SET

Method	Training ratios	
	10%	20%
AlexNet [17]	76.69±0.21	79.85±0.13
GoogleNet [17]	76.19±0.38	78.48±0.26
VGG_VD16 [17]	76.47±0.18	79.79±0.15
AlexNet+ BoCF [17]	55.22±0.39	59.22±0.18
GoogleNet+ BoCF [17]	78.92±0.17	80.97±0.17
VGG_VD16+ BoCF [17]	82.65±0.31	84.32±0.17
SPP-Net [20]	82.13±0.30	84.64±0.23
AlexNet+ single-layer SAFF	71.00±0.26	76.17±0.21
VGG_VD16+ single-layer SAFF	69.67±0.32	74.42±0.29
AlexNet+ sum pooling [10]	74.39±0.25	79.83±0.13
VGG_VD16+ sum pooling [10]	70.14±0.47	76.23±0.32
AlexNet+ SAFF (ours)	80.05±0.29	84.00±0.17
VGG_VD16+ SAFF (ours)	<b>84.38±0.19</b>	<b>87.86±0.14</b>

results that can be reflected more obviously on the NWPU-RESISC45 data set, and the overall accuracy is 8% higher than the CNN method. In [17], BoCF is also proposed to use the convolutional feature maps for classification. The overall classification accuracy achieved by our method is 3% higher than the BoCF on the NWPU data set. The performance of DCA [18], which fuses the last two FC layers of the CNN model, is about 89.71%, while our method is 93.83% on the AID with 50% training sample. Compared with SPP-Net [20], which used SPP to aggregate the last convolutional features, our method achieves better performance on all data sets. In [19], the MCNN constructs two branch networks that make the network available for the images with arbitrary scales, whose overall accuracy is about 2% lower than the proposed method on the AID. On the UC data set, our method shows almost the same performance with the DCA method. Although the UC data set is a classic data set for remote sensing scene classification, the construction of the imagery is simpler than others. Our method aims to make use of the complex content of the image to obtain better results; thus, we cannot greatly improve the overall accuracy on this data set. From the tables, we can observe that the VGGNet-16 achieves better performance in our method compared with that uses the feature maps of AlexNet. The difference is caused by the depth of different models. The deeper the network is, the more distinguishable the obtained features will be. The AlexNet has just five convolutional layers, while VGGNet-16 has 12 convolutional layers. Therefore, the VGGNet-16 is more competitive for classification. Through a series of experiments with different pretrained models and different training ratios, the results prove the effectiveness of our method.

## V. CONCLUSION

In this letter, we presented an SAFF method for remote sensing scene classification by aggregating the convolutional features extracted from the pretrained CNN models. In order to reduce the computation cost, the proposed method is dedicated to reflect the importance of elaborating surface objects and focuses more on the infrequently occurring features. Compared with the benchmark and other competitive ways, the aggregated feature processed by the proposed method contains

more discriminative properties for classification. The results experimented on different data sets verify the effectiveness of our method. However, there are still some works to be completed. In the future, we plan to design an end-to-end network that can automatically calculate the weights without human intervention.

## REFERENCES

- [1] G.-S. Xia *et al.*, "AID: A benchmark data set for performance evaluation of aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 7, pp. 3965–3981, Jul. 2017.
- [2] X. Bai, C. Liu, P. Ren, J. Zhou, H. Zhao, and Y. Su, "Object classification via feature fusion based marginalized kernels," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 1, pp. 8–12, Jan. 2015.
- [3] X. Bai, H. Zhang, and J. Zhou, "VHR object detection based on structural feature extraction and query expansion," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 10, pp. 6508–6520, Oct. 2014.
- [4] J. Xie, N. He, L. Fang, and A. Plaza, "Scale-free convolutional neural network for remote sensing scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 9, pp. 6916–6928, Sep. 2019.
- [5] N. He, L. Fang, S. Li, J. Plaza, and A. Plaza, "Skip-connected covariance network for remote sensing scene classification," *IEEE Trans. Neural Netw. Learn. Syst.*, to be published, doi: 10.1109/tnnls.2019.2920374.
- [6] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, Nov. 2004.
- [7] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. 18th SIGSPATIAL Int. Conf. Adv. Geograph. Inf. Syst. (GIS)*, 2010, pp. 270–279.
- [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Conf. Adv. Neural Inf. Process. Syst. (NIPS)*, 2012, pp. 1097–1105.
- [9] F. Hu, G.-S. Xia, J. Hu, and L. Zhang, "Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery," *Remote Sens.*, vol. 7, no. 11, pp. 14680–14707, Nov. 2015.
- [10] A. B. Yandex and V. Lempitsky, "Aggregating local deep features for image retrieval," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1269–1277.
- [11] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015, pp. 1–13.
- [12] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Proc. Conf. Adv. Neural Inf. Process. Syst. (NIPS)*, 2015, pp. 577–585.
- [13] C. Wang, X. Bai, S. Wang, J. Zhou, and P. Ren, "Multiscale visual attention networks for object detection in VHR remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 2, pp. 310–314, Feb. 2019.
- [14] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015, pp. 1–15.
- [15] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2016, pp. 685–701.
- [16] H. Jégou and O. Chum, "Negative evidences and co-occurrences in image retrieval: The benefit of PCA and whitening," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2012, pp. 774–787.
- [17] G. Cheng, Z. Li, X. Yao, L. Guo, and Z. Wei, "Remote sensing image scene classification using bag of convolutional features," *IEEE Geosci. Remote Sens. Lett.*, vol. 14, no. 10, pp. 1735–1739, Oct. 2017.
- [18] S. Chaib, H. Liu, Y. Gu, and H. Yao, "Deep feature fusion for VHR remote sensing scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 8, pp. 4775–4784, Aug. 2017.
- [19] Y. Liu, Y. Zhong, and Q. Qin, "Scene classification based on multiscale convolutional neural network," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 12, pp. 7109–7121, Dec. 2018.
- [20] 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 Sens.*, vol. 9, no. 8, p. 848, Aug. 2017.