

Scale-Free Convolutional Neural Network for Remote Sensing Scene Classification

Jie Xie, *Student Member, IEEE*, Nanjun He^{id}, *Student Member, IEEE*, Leyuan Fang^{id}, *Senior Member, IEEE*, and Antonio Plaza^{id}, *Fellow, IEEE*

Abstract—Fine-tuning of pretrained convolutional neural networks (CNNs) has been proven to be an effective strategy for remote sensing image scene classification, particularly when a limited number of labeled data sets are available for training purposes. However, such a fine-tuning process often needs that the input images are resized into a fixed size to generate input vectors of the size required by fully connected layers (FCLs) in the pretrained CNN model. Such a resizing process often discards key information in the scenes and thus deteriorates the classification performance. To address this issue, in this paper, we introduce a scale-free CNN (SF-CNN) for remote sensing scene classification. Specifically, the FCLs in the CNN model are first converted into convolutional layers, which not only allow the input images to be of arbitrary sizes but also retain the ability to extract discriminative features using a traditional sliding-window-based strategy. Then, a global average pooling (GAP) layer is added after the final convolutional layer so that input images of arbitrary size can be mapped to feature maps of uniform size. Finally, we utilize the resulting feature maps to create a new FCL that is fed to a softmax layer for final classification. Our experimental results conducted using several real data sets demonstrate the superiority of the proposed SF-CNN method over several well-known classification methods, including pretrained CNN-based ones.

Index Terms—Free-scale convolutional neural networks (CNNs), fully connected layers (FCLs), remote sensing scene classification.

Manuscript received November 19, 2018; revised January 19, 2019 and March 9, 2019; accepted March 31, 2019. Date of publication April 25, 2019; date of current version August 27, 2019. This work was supported in part by the National Natural Science Fund of China for International Cooperation and Exchanges under Grant 61520106001, in part by the National Natural Science Foundation under Grant 61771192, in part by the National Natural Science Foundation of Hunan Province under Grant 2018JJ3077, in part by the Science and Technology Plan Project Fund of Hunan Province under Grant CX2018B171, Grant 2017RS3024, and Grant 2018TTP1013, in part by the Science and Technology Talents Program of Hunan Association for Science and Technology under Grant 2017TJ-Q09, in part by the China Postdoctoral Science Foundation under Project 2017T100597, and in part by the MINECO Project under Grant TIN2015-63646-C5-5-R. (Jie Xie and Nanjun He contributed equally to this work.) (Corresponding author: Leyuan Fang.)

J. Xie, N. He, and L. Fang 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: xj_xj@hnu.edu.cn; henanjun@hnu.edu.cn; fangleyuan@gmail.com).

A. Plaza is with the Hyperspectral Computing Laboratory, Department of Technology of Computers and Communications, Escuela Politécnica, University of Extremadura, E-10003 Cáceres, Spain (e-mail: aplaza@unex.es).

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

Digital Object Identifier 10.1109/TGRS.2019.2909695

I. INTRODUCTION

WITH the fast development of satellite sensor technology, remote sensing scene classification has drawn significant attention due to the wide range of applications that can now be addressed with such instruments, including urban planning [1], traffic flow prediction [2], and military reconnaissance [3]. The goal of remote sensing scene classification is to assign a specific label (e.g., *beach* or *bridge*) to a query remote sensing image.

In order to recognize a particular scene from a set of remotely sensed images, numerous feature extraction (FE) and classification methods have been proposed in the past decades, and an extensive review of them can be found in [4] and [5]. FE is a crucial part of the scene recognition process, which can be divided into three levels: 1) low-level FE; 2) midlevel FE; and 3) high-level FE. Early works on scene classification are based on low-level features, such as color [6], texture [7], and scale-invariant [8] features. However, low-level features may be too simple to describe the complex spatial layout of remote sensing scenes. To address this problem, many midlevel FE methods have been proposed. The bag-of-visual-words model [9] and the dictionary learning method [10] are two classic approaches for midlevel FE. Furthermore, to bridge the semantic gap, high-level FE techniques based on convolutional neural networks (CNNs) have been introduced in the field of remote sensing scene classification. The development process of these three different levels of FE methods reveals some important insights, as described in the following.

The key to low-level FE methods is the design of feature descriptors, such as color descriptors, textural descriptors, and scale-invariant descriptors. Specifically, in [11], an improved color code is introduced to accelerate scene classification performance by combining the advantages of digital image processing, geographical information systems, and data mining. In [12], morphological texture descriptors are applied to extract useful content from remotely sensed images. In [13], a method to extract scale and rotation-invariant features is proposed that greatly promotes feature generalization. Moreover, in [14], multiple kinds of feature descriptors are used to extract rich information from remotely sensed images, thus enhancing image representation and scene classification. However, the aforementioned methods only focus on relatively simple low-level features that cannot fully capture the rich information contained in remotely sensed images.

To solve this problem, many researchers have been focusing on how to efficiently represent remote sensing images using midlevel features. The main process of midlevel FE methods is to use a set of handcrafted feature descriptors (e.g., color or texture descriptors) to extract local image attributes from the original images and then build high-order statistical patterns by encoding these features [15]–[18]. The k -means clustering is a basic strategy to combine multiple kinds of features, and the bag-of-visual word-based methods [19], [20] are also quite popular in this context due to their effectiveness and simplicity. In [21], spatial pyramid matching is adopted to enhance the bag-of-visual-words model by combining local and global features. Furthermore, sparse coding [22], [23] has also been adopted for scene classification purposes by adding a sparsity constraint to the feature distributions that effectively reduces the complexity of the model and simplifies the associated learning tasks. In [24], a weighted deconvolutional sparse coding model is proposed for unsupervised extraction of edges and texture details from remotely sensed images. Nevertheless, these methods still heavily rely on low-level feature descriptors and cannot fully capture the high-level semantic information contained in the scenes.

In recent years, based on the excellent performance of CNNs in many image classification challenges (e.g., the ImageNet [25], Openimage [26], and Places365 [27]), numerous researchers have focused on adapting CNN-based methods to remote sensing problems [28]–[31]. CNN models can automatically achieve effective feature representation by means of hierarchical layers, where the shallow layers extract local low-level features and the deep layers extract global high-level semantic features [32], [33]. These high-level semantic features can be directly utilized to bridge the semantic gap between different scenes within the same class and thus achieve a better classification performance. Specifically, [34] proposes a gradient boosting random framework that utilizes CNN models, pretrained on a data set made up of natural images to accelerate the scene classification performance by integrating models with different structures. Considering the data shift problem between natural images and remote sensing images, a domain adaptation network is introduced in [35]. Resulting from the fact that different layers provide information with different degrees of effectiveness, [36] and [37] fuse features from different layers of a CNN model, pretrained on ImageNet, to increase the classification accuracy. Moreover, to solve the interclass similarity and intraclass diversity problems in remote sensing image data sets, off-the-shelf models are equipped with metric learning in [38] and [39] to enhance the scene classification results by changing the final feature space distribution. In general, fine-tuning CNN models that have been pretrained on a data set made up of natural images can provide an effective and also efficient strategy for remote sensing scene classification [40]–[42]. This is because only thousands of labeled images exist in remote sensing data sets (e.g., the UC Merced Land-Use data set [43], the Aerial Image data set, and the NWPU-RESISC45 data set) compared with millions of labeled images available in natural image data sets (e.g., ImageNet, Openimage, and Places365), and thus, it is

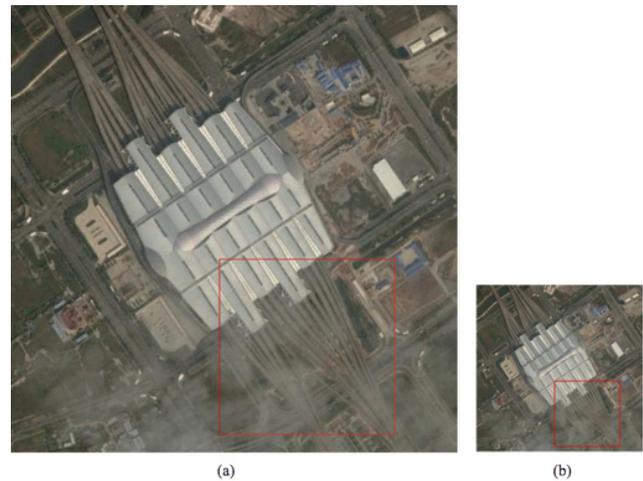


Fig. 1. (a) Original image with a size of 600×600 pixels. (b) Resized image with 224×224 pixels. The red rectangle marks an area that loses important details in the resizing process.

hard to train a CNN model from scratch for remote sensing scene classification.

Although fine-tuning the pretrained CNN models can lead to an acceptable performance in remote sensing scene classification, this strategy also has some limitations. Specifically, pretrained CNN models are learned on data sets of a relatively small and fixed size (e.g., 224×224 pixels in AlexNet [44]) due to the matrix operations performed in the fully connected layers (FCLs), and thus, it is required that the input images to be processed to have the same size as the images used for pretraining when the pretrained CNN models are fine-tuned. On the other hand, remote sensing images usually have higher sizes than the maximum ones allowed by pretrained CNN models. For example, in the widely used Aerial Image remote sensing scene data set [4], each image has a size of 600×600 pixels, which is much higher than the maximum image size allowed by AlexNet, i.e., 224×224 pixels. A common strategy to address this limitation is to resize the original image (e.g., from 600×600 to 224×224 pixels, as shown in Fig. 1) [45], [46]. However, some key information in the original image is inevitably lost during the preprocessing [notice the area marked with a red rectangle in Fig. 1(a) and (b)].

To address the aforementioned problems, in this paper, we develop a new scale-free CNN (SF-CNN) architecture that can process remotely sensed images of arbitrary size, retaining the strong FE ability of pretrained CNN models. As it is the case with most pretrained CNN models, our proposed approach consists of two main parts: the convolutional layers and the FCLs. However, as opposed to traditional methods, in which the FCLs have a restriction that the sizes of the input image should be the same as those in the pretrained model, our SF-CNN addresses this issue by performing a convolution strategy on the FCLs while still efficiently extracting discriminative and highly representative features from the input images in a sliding-window manner. Furthermore, our newly developed SF-CNN adds a global average pooling (GAP) layer

after the final convolutional layer, mapping the input images of arbitrary size to feature maps of fixed size. The final output from the pooling layer is fed to the FCL and a softmax layer to obtain the final probability of each category. Our experimental results demonstrate that the proposed method can indeed fully utilize the remotely sensed images available (regardless of their size) and outperform the baseline methods and some state-of-the-art approaches on several publicly available benchmark data sets.

The remainder of this paper is organized as follows. Section II reviews some related works. Section III introduces the proposed SF-CNN. Section IV shows our experimental results conducted on several publicly available benchmark data sets. Section V concludes this paper with some remarks and hints at plausible future research lines.

II. RELATED WORKS

A. Structure of the CNN Model

With the rapid development of CNNs in the classification of natural images, many pretrained CNN models are now publicly available (e.g., AlexNet [44], GoogleNet [47], and VGGNet [48]). Generally, the most representative and discriminative features are captured by convolutional layers and FCLs. In this section, the mechanisms of these two types of layers are described in detail.

1) *Convolutional Layer*: The CNN model contains a group of cascaded convolutional layers, each comprising a set of convolutional kernels (also called filters), which are used to convolve the input data and then produce different kinds of output data. Let $\mathbf{X}_i = \{x_{l,l,l,i}, \dots, x_{w,h,c,i}, \dots, x_{W,H,C,i}\}$ —where W represents the width, H represents the height, and C represents the channel—be the input pattern in the i th convolutional layer. Let us also assume that there are a total of J kernels in the i th convolutional layer and that the size of each kernel is $K \times K \times C$, where K represents the width and height and C is the channel of each kernel. Let $\mathbf{W}_{j,i} = \{w_{l,l,l,j,i}, \dots, w_{k,k,c,j,i}, \dots, w_{K,K,C,j,i}\}$ —with $l \leq j \leq J$, $l \leq k \leq W$ and $l \leq c \leq H$ —be the j th kernel in the i th layer. The output of this convolutional layer $\mathbf{Y}_{j,i} = \{y_{l,l,j,i}, \dots, y_{w,h,j,i}, \dots, y_{W-K+1,H-K+1,j,i}\}$ can be obtained by

$$\mathbf{Y}_{j,i} = \mathbf{X}_i \otimes \mathbf{W}_{j,i} \quad (1)$$

where $\mathbf{Y}_i = \{\mathbf{Y}_{1,i}, \dots, \mathbf{Y}_{j,i}, \dots, \mathbf{Y}_{J,i}\}$ is the output of this layer and \otimes is the convolutional operation. Without padding, this operation is denoted as

$$y_{w,h,j,i} = \sum_{c=1}^C \sum_{n=1}^K \sum_{m=1}^K w_{m,n,c,j,i} x_{w+m-1,h+n-1,c,i}. \quad (2)$$

The mapping in (1) can also be defined as $\mathbf{Y}_i = \mathbf{f}(\mathbf{X}_i)$. When a CNN model is fine-tuned, and despite the fact that the size of $\mathbf{W}_{j,i} = \{\mathbf{W}_{1,i}, \dots, \mathbf{W}_{j,i}, \dots, \mathbf{W}_{J,i}\}$ is fixed, the size of \mathbf{X}_i is arbitrary. In conclusion, the mapping in (1) is not limited by the size of the input data.

2) *Fully Connected Layer*: Several FCLs follow the design of the final convolutional layer in the CNN model. The t th FCL has a mapping matrix $\mathbf{S}_t = \{s_{l,l,t}, \dots, s_{m,n,t}, \dots, s_{M,N,t}\}$ of size $M \times N$. It fully connects all the output data in the

previous layer and maps the data to a new vector $\mathbf{Z}_t = \{z_{l,t}, \dots, z_{n,t}, \dots, z_{N,t}\}$ with size $1 \times N$. Specifically, the output of the previous layer, which is also the input of the t th layer, needs to be reshaped to $\mathbf{A}_t = \{a_{l,t}, \dots, a_{m,t}, \dots, a_{M,t}\}$, of size $1 \times M$. In addition, the output from a convolutional layer can be represented as $\mathbf{X}_i = \{x_{l,l,l,i}, \dots, x_{w,h,c,i}, \dots, x_{W,H,C,i}\}$, where $M = W \times H \times C$. The relationship between \mathbf{A}_t and \mathbf{X}_i can be represented as $\mathbf{A}_t = \varphi(\mathbf{X}_i)$. The above-mentioned relationship can be denoted as

$$\mathbf{Z}_t = \mathbf{A}_t \mathbf{S}_t \quad (3)$$

where each $z_{n,t}$ of \mathbf{Z}_t can be obtained by

$$z_{n,t} = \sum_{m=1}^M a_{m,t} s_{m,n,t}. \quad (4)$$

In a transfer learning task, the FCLs in the pretrained model are essential to achieve high performance [49]. During the fine-tuning process of the CNN model, the size of \mathbf{S}_t ($M \times N$) is fixed, so that the size of \mathbf{A}_t should match this size. This imposes a limitation that the input images should have a fixed size. In addition, the CNN model used for classification purposes must use an FCL to generate the final label.

B. CNN-Based Scene Classification

CNNs exhibit powerful generalization ability and very good performance on natural image classification problems [50]. The great success of the CNN model is partly due to the huge amount of labeled training data sets available (e.g., the ImageNet, Openimage, and Places365 data sets have millions of labeled images). Recently, the CNN model has also been extended to remote sensing scene classification [51]–[54]. Since the number of labeled remote sensing images is still limited (e.g., the widely used Aerial Image data set only contains 10 000 labeled samples [4]), CNN models cannot be trained from scratch using these data sets. A popular strategy to alleviate this limitation is to adopt a transfer learning method, which utilizes the available remotely sensed images to fine-tune some CNN models (e.g., AlexNet, GoogleNet, or VGGNet) that have been already pretrained on some large-scale data set. Generally, fine-tuning of a pretrained CNN model takes advantage of the pretrained convolutional layers and FCLs to adapt the architecture to new classification tasks. This strategy has been shown to be effective for remote sensing scene classification purposes [34]–[42], [45].

III. SCALE-FREE CONVOLUTIONAL NEURAL NETWORK

Although fine-tuning CNNs can achieve the state-of-the-art scene classification performance in remote sensing problems, all available pretrained CNN models need to resize the input remotely sensed scene into a (lower) fixed size and thus inevitably discard some key information in the scene, which eventually deteriorates the scene recognition task. As described in Section II, this limitation results from the use of FCL matrices of fixed size, including the FCL matrix used to obtain the final label. In other words, the size of the input images must always match the size of the FCLs due to these

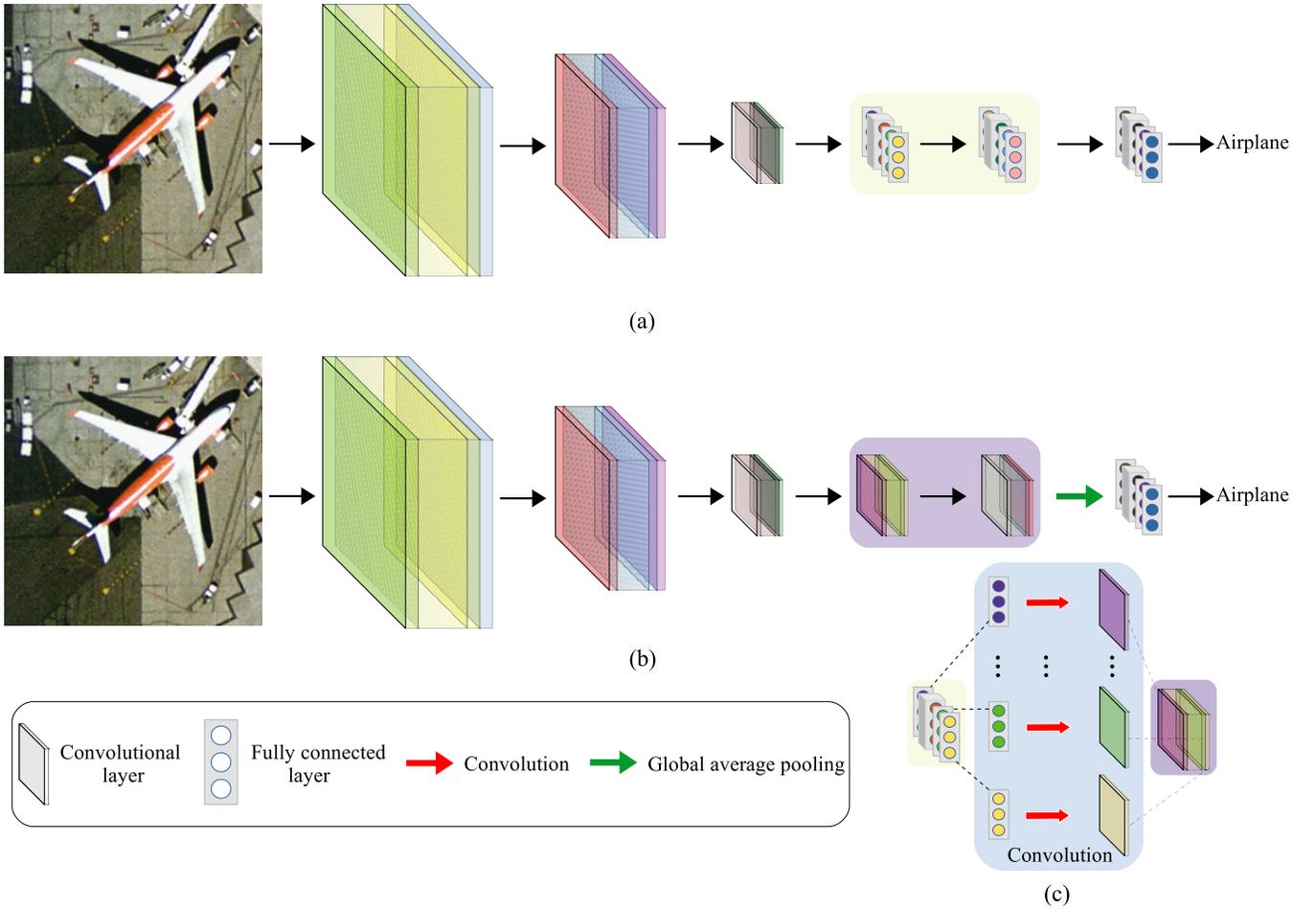


Fig. 2. Graphical representation of the architecture of (a) original model, (b) our SF-CNN model, and (c) FCL convolution.

matrix operations. In order to address this limitation and allow the input scenes to be of arbitrary size, we introduce a new SF-CNN in this section. Note that our method does not reduce the powerful FE ability of pretrained CNN models, since we have used an equivalent structure in our newly proposed architecture, as described in the following.

A. Architecture of the SF-CNN

Fig. 2(b) shows the architecture of the proposed SF-CNN, in which the parameters of the convolutional layer are directly transferred from a pretrained CNN model on ImageNet. Specifically, the proposed SF-CNN contains two main components: 1) FCLs' convolution and 2) extra GAP layer. With these two components, the proposed SF-CNN allows the input remote sensing scenes to be of arbitrary size. Note that the first component is crucial to retain the ability of the pretrained CNN model to extract effective features for scene classification. The second component matches the input data of the FCL to the size of the FCL needed to obtain the final label. In the following, these two key components (and the optimization process of our newly developed SF-CNN) are described in detail.

1) *FCLs' Convolution*: As described in Section II-A2, although FCLs are very important for transfer learning, they

require that the input images have a fixed size. Generally, the data streams flowing into the FCLs can be divided into two main categories. The first one is the output of a convolutional layer $\mathbf{X}_i = \{x_{1,1,i}, \dots, x_{w,h,c,i}, \dots, x_{W,H,C,i}\}$, which is vectorized as $\mathbf{A}_t = \{a_{1,t}, \dots, a_{m,t}, \dots, a_{M,t}\}$ before feeding it into the t th FCL by regularization, denoted by $\mathbf{A}_t = \varphi(\mathbf{X}_i)$. The second one is the output of the other FCL $\mathbf{Z}_t = \{z_{1,t}, \dots, z_{n,t}, \dots, z_{N,t}\}$. In the t th FCL, the input data $\mathbf{A}_t = \{a_{1,t}, \dots, a_{m,t}, \dots, a_{M,t}\}$ are linearly mapped to a vector $\mathbf{Z}_t = \{z_{1,t}, \dots, z_{n,t}, \dots, z_{N,t}\}$ by a mapping matrix $\mathbf{S}_t = \{s_{1,1,t}, \dots, s_{m,n,t}, \dots, s_{M,N,t}\}$. The size of the input images must be fixed during the fine-tuning of a pretrained CNN model due to the mapping matrix \mathbf{S}_t , which is of fixed size. Therefore, the remotely sensed scenes used to fine-tune CNN models normally need to be resized, which may discard key information in the scenes. However, as mentioned in Section II-A1, input images with arbitrary size can now be directly fed into the convolutional layers in a transfer learning task. In other words, the convolutional layers have no limitation regarding the size of the input images. Hence, it is effective to modify the FCLs to match the convolutional layers [55] so that the remote sensing scenes with arbitrary size can be directly fed into the pretrained CNN model. A feasible strategy is FCLs' convolution, achieved by converting the mapping matrix to

TABLE I
SETTING OF THE HYPERPARAMETERS USED FOR OPTIMIZATION

Models	Iterations	Batch sizes	Learning rates (convolutional layer)	Learning rates (FCL)	Weight decay	Momentum	Optimization algorithm
AlexNet	15000	128	0.001	0.01	0.0005	0.9	Nesterov
GoogleNet	15000	128	0.001	0.01	0.0005	0.9	Nesterov
VGGNet	15000	50	0.001	0.01	0.0005	0.9	Nesterov

the convolutional kernels. The FCLs' convolution not only efficiently extracts features from the input data but also eliminates the limitation that the input data need to have a fixed size. Here, to this end, the n th column of the mapping matrix $\mathbf{S}_{n,t} = \{s_{1,n,t}, \dots, s_{m,n,t}, \dots, s_{M,n,t}\}$ is converted into the n th convolutional kernel $\mathbf{W}_{n,i} = \{w_{1,1,1,n,i}, \dots, w_{w,h,c,n,i}, \dots, w_{W,H,C,n,i}\}$ by a mapping regulation $\mathbf{S}_{n,t} = \varphi(\mathbf{W}_{n,i})$. On the one hand, the original input data from the convolutional layer $\mathbf{X}_i = \{x_{1,1,1,i}, \dots, x_{w,h,c,i}, \dots, x_{W,H,C,i}\}$ are directly fed into this FCL, and C of $\mathbf{W}_{n,i}$ is regarded as the C of \mathbf{X}_i . On the other hand, when the original input data come from the FCL $\mathbf{A}_t = \{a_{1,t}, \dots, a_{m,t}, \dots, a_{M,t}\}$, the width and height of the input data can be regarded as $\mathbf{X}_i = \{x_{1,1,1,i}, \dots, x_{1,1,c,i}, \dots, x_{1,1,C,i}\}$, and W and H of $\mathbf{W}_{n,i}$ are all set to 1. $\mathbf{W}_i = \{w_{1,1,1,n,i}, \dots, w_{1,1,c,n,i}, \dots, w_{1,1,C,n,i}\}$. In this context, C of \mathbf{X}_i is equal to M of \mathbf{A}_t . As a matter of fact, the FCL can be regarded as a special kind of convolutional layer, where the size of the convolutional kernels equals the size of the input data. The equivalence of this transformation process is demonstrated in Section III-A2.

2) *Global Average Pooling Layer*: The GAP layer has been used by some available CNN architectures to reduce the model size and address overfitting issues [56], [57]. Compared with the global max pooling (GMP), the GAP is more suitable for classification tasks [58], [59], especially for scene classification tasks, where some categories (e.g., center and school) require global information to classify. By contrast, the GMP is suitable for object localization tasks [60], [61] due to its robustness to the local spatial variation [62]–[64]. Our method incorporates a GAP layer right after the final convolutional layer. Specifically, the size of the input data $\mathbf{X}_i = \{x_{1,1,1,i}, \dots, x_{w,h,c,i}, \dots, x_{W,H,C,i}\}$ in the i th layer is $W \times H \times C$, and the size of the output data $\mathbf{G}_i = \{g_{1,i}, \dots, g_{c,i}, \dots, g_{C,i}\}$ in the i th layer is $1 \times C$. The operation conducted by the GAP layer is given by

$$g_{c,i} = \frac{\sum_{w=1}^W \sum_{h=1}^H x_{w,h,c,i}}{W \times H}. \quad (5)$$

This operation is used to obtain the average value of each channel, so that the arbitrary size of the output data is only related to the number of channels, which depends on the value of j in the last convolutional layer.

3) *Optimization*: The output of the final FCL $\mathbf{Z}_t = \{z_{1,t}, \dots, z_{n,t}, \dots, z_{N,t}\}$ is fed to the softmax layer in order to obtain the probability that of each image belonging to each class. This probability is denoted as follows [65]:

$$P_k = \frac{e^{z_{k,t}}}{\sum_{n=1}^N e^{z_{n,t}}}, \quad k = 1, \dots, n, \dots, N \quad (6)$$

where N is the number of classes and $\mathbf{P} = [P_1, \dots, P_k, \dots, P_N]^T$. Then, the class with the maximal probability is used as the estimated label δ_n for each image. Based on the estimated labels, the loss function L_f can be obtained via a combination of logistic loss and an additional weight decay term for regularization

$$\min L_f = \min \left(- \sum_{batch} \theta \cdot \log(\mathbf{P}) + \lambda \left(\|\mathbf{W}^{(\cdot)}\|_F^2 + \|\mathbf{S}^{(\cdot)}\|_F^2 \right) \right) \quad (7)$$

where θ represents a vector that uses 1 as a true label and 0 otherwise, the $\mathbf{W}^{(\cdot)}$ represents the set of all parameters in the convolutional layers, the $\mathbf{S}^{(\cdot)}$ represents the set of all parameters in the FCLs, and λ is the weight decay coefficient of the SF-CNN. To minimize the loss function L_f , the backward propagation algorithm [66] is adopted to update the aforementioned parameters $\mathbf{W}^{(\cdot)}$ and $\mathbf{S}^{(\cdot)}$. Specifically, it propagates the predicted error from the last layer to the first one and modifies the parameters according to the gradient of the propagated error at each layer. In general, the stochastic gradient descent (SGD) algorithm is applied to achieve this goal. Table I summarizes the hyperparameters used for optimization. Note that the batch size of VGGNet is set to 50 due to the limitations in the memory of the graphical processing unit (GPU) used for implementing our approach.

B. Equivalence Proof of the FCLs' Convolution

In this section, our proposed FCLs' convolution is proved to be an equivalent transformation, with no effects on the training and testing process, compared with the original FCLs. These two processes consist of forward propagation and backpropagation. Specifically, in order to make the proof more concise, we define $\mathbf{S}_{n,t} = \{s_{1,n,t}, \dots, s_{m,n,t}, \dots, s_{M,n,t}\}$, and a new mapping regulation ψ is defined as

$$\mathbf{X}_i = \psi(\varphi(\mathbf{X}_i)). \quad (8)$$

This mapping regulation ψ is considered an inverse mapping φ . In the following, we detail the equivalence proof on the forward propagation and backpropagation phases.

1) *Equivalence Proof on the Forward Propagation Phase*: As described in Section II, in the i th FCL, the relationship between the input data $\mathbf{A}_t = \{a_{1,t}, \dots, a_{m,t}, \dots, a_{M,t}\}$ and the output data $\mathbf{Z}_t = \{z_{1,t}, \dots, z_{n,t}, \dots, z_{N,t}\}$ can be represented as $\mathbf{Z}_t = \mathbf{A}_t \mathbf{S}_t$. According to Sections II and III-A1, the

convolution transformation can be denoted as

$$Y_i = \psi(A_t) \otimes [\psi(S_{1,t}), \dots, \psi(S_{n,t}), \dots, \psi(S_{N,t})] \quad (9)$$

$$Y_i = [X_i \otimes W_{1,i}, \dots, X_i \otimes W_{n,i}, \dots, X_i \otimes W_{N,i}] \quad (10)$$

$$Y_i = [Y_{1,i}, \dots, Y_{n,i}, \dots, Y_{N,i}]. \quad (11)$$

On the one hand, when the sizes of X_i and $W_{n,i}$ are the same, the number of elements in $Y_{n,i}$ is 1, which means that W and H of $Y_{n,i}$ are 1. From (2), (11) can be obtained by

$$Y_{n,i} = \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W w_{w,h,c,n,i} x_{w,h,c,i}. \quad (12)$$

On the other hand, when the matrices (X_i and $W_{n,i}$) are obtained from vectors (A_t and $s_{m,t}$) by the same regulation (ψ), $w_{w,h,c,n,i} x_{w,h,c,i}$ and $a_{m,t} s_{m,n,t}$ have a one-to-one relationship, and an equality relation can be obtained as follows:

$$Y_{n,i} = z_{n,t}, \quad (13)$$

$$Y_i = Z_t. \quad (14)$$

From (14), no changes are found in the output after the FCLs' convolution, which indicates that this operation retains its ability to extract discriminative features.

2) *Equivalence Proof on the Backpropagation Phase:* Because of the existing equivalence between Y_i and Z_t [from (14)], to make the process more concise, the mappings between P and the output Y_i or Z_t are defined as

$$P = F(Y_i) \quad (15)$$

or

$$P = F(Z_t). \quad (16)$$

During the backpropagation process, S_t is updated by

$$S_t = S_t - \alpha \frac{\partial L_f}{\partial S_t}. \quad (17)$$

From (3), (6), and (7), $\frac{\partial L_f}{\partial S_t}$ of (17) can be solved as

$$\frac{\partial L_f}{\partial S_t} = \frac{-\partial(\sum_z \theta \cdot \log(P))}{\partial S_t} + \frac{\lambda \partial(\|W^{(\cdot)}\|_F^2 + \|S^{(\cdot)}\|_F^2)}{\partial S_t} \quad (18)$$

$$\frac{\partial L_f}{\partial S_t} = -\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} \frac{\partial P}{\partial Z_t} \frac{\partial Z_t}{\partial S_t} + \lambda \frac{\partial(\|S^{(\cdot)}\|_F^2)}{\partial S_t} \quad (19)$$

$$\frac{\partial L_f}{\partial S_t} = -\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Z_t) \cdot A_t + \lambda \frac{\partial(\sum_{n=1}^N \sum_{m=1}^M s_{m,n,t}^2)}{\partial S_t} \quad (20)$$

$$S_t = S_t - \alpha \left(-\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Z_t) \cdot A_t + \lambda \frac{\partial(\sum_{n=1}^N \sum_{m=1}^M s_{m,n,t}^2)}{\partial S_t} \right) \quad (21)$$

$$s_{m,n,t} = s_{m,n,t} - \alpha \left(-\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Z_t) \cdot a_{m,t} + 2\lambda s_{m,n,t} \right). \quad (22)$$

After the convolution transformation, the optimization is now denoted as follows:

$$W_i = W_i - \alpha \frac{\partial L_f}{\partial W_i}. \quad (23)$$

From (1), (2), (6), (7), (11), and (12), $(\partial L_f / \partial W_i)$ of (23) can be solved as

$$\frac{\partial L_f}{\partial W_i} = \frac{-\partial(\sum_z \theta \cdot \log(P))}{\partial W_i} + \frac{\lambda \partial(\|W^{(\cdot)}\|_F^2 + \|S^{(\cdot)}\|_F^2)}{\partial W_i} \quad (24)$$

$$\frac{\partial L_f}{\partial W_i} = -\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} \frac{\partial P}{\partial Y_i} \frac{\partial Y_i}{\partial W_i} + \lambda \frac{\partial(\|W^{(\cdot)}\|_F^2)}{\partial W_i} \quad (25)$$

$$\frac{\partial L_f}{\partial W_i} = -\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Y_i) \cdot X_i + \lambda \frac{\partial(\sum_{n=1}^N \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W w_{w,h,c,n,i}^2)}{\partial W_i} \quad (26)$$

$$W_i = W_i - \alpha \left(-\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Y_i) \cdot X_i + \lambda \frac{\partial(\sum_{n=1}^N \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W w_{w,h,c,n,i}^2)}{\partial W_i} \right) \quad (27)$$

$$w_{w,h,c,n,i} = w_{w,h,c,n,i} - \alpha \left(-\frac{\partial(\sum_z \theta \cdot \log(P))}{\partial P} F'(Y_i) \cdot x_{w,h,c,i} + 2\lambda w_{w,h,c,n,i} \right). \quad (28)$$

Since $W_i = \psi(S_i)$, $X_i = \psi(A_i)$ and ψ is a linear mapping with a coefficient of 1 and a bias of 0, relationship between the $w_{w,h,c,n,i}$, $x_{w,h,c,i}$, $s_{m,n,i}$, and $a_{m,i}$ is denoted as

$$\forall s_{m,n,i}, \exists! w_{w,h,c,n,i} = s_{m,n,t} \text{ and } a_{m,t} = x_{w,h,c,i}. \quad (29)$$

From (14)–(16), this expression can be solved as

$$F'(Z_t) = F'(Y_i). \quad (30)$$

Therefore, we conclude that (22) and (28) are equivalent. Based on the aforementioned description, the convolution transformation is indeed equivalent during the backpropagation phase.

IV. EXPERIMENTAL RESULTS

A. Data Sets' Description

To validate the effectiveness of our newly developed SF-CNN model, we perform a set of comprehensive experiments on three publicly available benchmark remote sensing scene data sets that are the UC Merced Land-Use data set [43], the Aerial Image data set [4], and the NWPU-RESISC45 data set [5].

- 1) The UC Merced Land-Use data set consists of 2100 images divided into 21 land-use classes, including agricultural, airplane, baseball diamond, beach, buildings, chaparral, dense residential, forest, freeway, golf course, harbor, intersection, medium residential,



Fig. 3. Some examples of scenes that are easily misclassified in the UC Merced Land-Use data set.



Fig. 4. Some examples of scenes with high interclass similarity in the Aerial Image data set.

mobile home park, overpass, parking lot, river, runway, sparse residential, storage tanks, and tennis courts. Each class contains 100 aerial images with 256×256 pixels, and each pixel has a spatial resolution of 0.3 m in the red–green–blue (RGB) color space. Fig. 3 shows some examples of scenes in the UC Merced Land-Use data set which are easily misclassified.

- 2) The Aerial Image data set consists of 10 000 images divided into 30 scene classes, including airport, bare land, baseball field, beach, bridge, center, church, commercial, dense residential, desert, farmland, forest, industrial, meadow, medium residential, mountain, park, parking, playground, pond, port, railway station, resort, river, school, sparse residential, square, stadium, storage tanks, and viaduct. Each class contains hundreds of aerial images (from 220 to 420) with 600×600 pixels in the RGB color space. The spatial resolution of these image ranges from 8 to 0.5 m/pixel. Fig. 4 shows some examples of the Aerial Image data set. As it can be seen in Fig. 4, some classes exhibit a very high interclass similarity (e.g., *bare land* and *desert*), which is the main difficulty for the classification of scenes in this data set.
- 3) The NWPU-RESISC45 data set consists of 31 500 images divided into 45 classes, including airplane, airport, baseball diamond, basketball court, beach, bridge, chaparral, church, circular farmland, cloud, commercial

area, dense residential, desert, forest, freeway, golf course, ground track field, harbor, industrial area, intersection, island, lake, meadow, medium residential, mobile home park, mountain, overpass, palace, parking lot, railway, railway station, rectangular farmland, river, roundabout, runway, sea ice, ship, snowberg, sparse residential, stadium, storage tank, tennis court, terrace, thermal power station, and wetland. Each class contains 700 images with spatial resolution ranging from about 30 to 0.2 m/pixel and a size of 256×256 pixels in the RGB color space. This is one of the largest remote sensing scene data sets in terms of the number of images and the categories, which leads to larger intraclass differences and higher interclass similarities than the ones observed in the two aforementioned data sets. Some examples are given in Fig. 5.

B. Experimental Setup

For the UC Merced Land-Use data set, a training proportion of 80% ($Pr = 80\%$) randomly selected samples is considered for training, and the remaining 20% of the labeled samples are used for testing. For the Aerial Image data set, the considered training proportions are $Pr = 20\%$ and $Pr = 50\%$. For the NWPU-RESISC45 data set, the considered training proportions are $Pr = 10\%$ and $Pr = 20\%$. These proportions



Fig. 5. Some classes with large intraclass difference and high interclass similarity in the NWPU-RESISC45 data set.

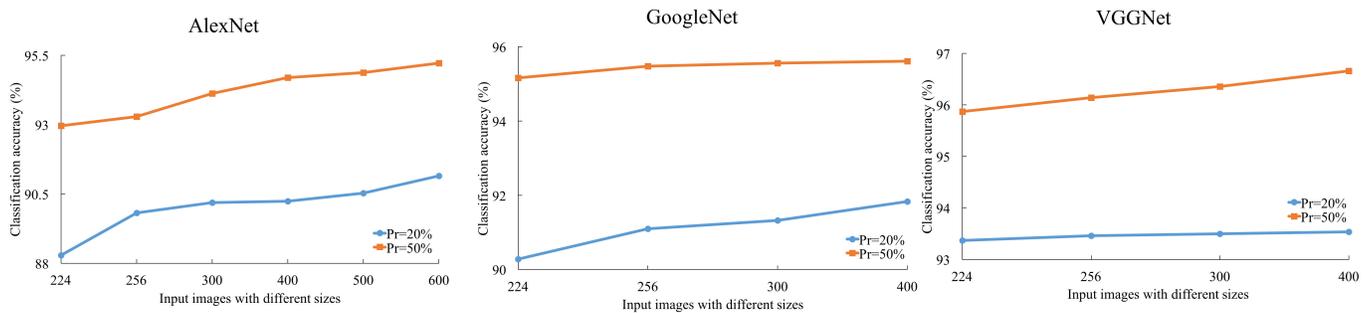


Fig. 6. Graphical illustration of the benefits obtained using remote sensing images without a fixed size for scene classification using our newly proposed SF-CNN model. Three different pretraining strategies are considered. (a) AlexNet. (b) GoogleNet. (c) VGGNet.

have been selected in accordance with the previous studies in the literature, in order to facilitate our comparisons with state-of-the-art approaches.

To evaluate the results of the proposed SF-CNN for scene classification, the average accuracy (AA), the Kappa coefficient (Kappa), the overall accuracy (OA), and the confusion matrix are adopted as evaluation metrics in the following experiments. Three classic pretrained CNN models (i.e., AlexNet, GoogleNet, or VGGNet) are utilized to analyze the generalization ability of our proposed framework on the three publicly available remote sensing scene data sets. The pretrained FCLs of these CNN models are convolutionalized, and a large-margin Gaussian mixture loss is added to obtain more representative and discriminative features [67]. In addition, all our experimental results are obtained as the average of ten repeated experiments using different, randomly selected training samples. Our experiments are conducted on a PC with CPU i7-7700K, 16 GB of RAM, and a GPU (GTX 1080 Ti).

C. Benefits of Using Input Images With Different Sizes

To illustrate that the exploitation of remotely sensed images without any limitations regarding their size can significantly enhance the classification accuracies obtained by the proposed SF-CNN model, in this experiment, we consider input images with different sizes. We take the images in the Aerial Image data set as a baseline due to their original size of 600×600 pixels. These images are resized to 224×224 , 256×256 , 300×300 , 400×400 , and 500×500 pixels. Note that, due to the memory limitations of the GPU (GTX 1080 Ti) used in our experiments, the SF-CNN models pretrained with

GoogleNet and VGGNet are only resized to 224×224 , 256×256 , 300×300 , and 400×400 pixels, with the batch sizes of 128 and 50, respectively. As it can be easily observed from Fig. 6, the proposed SF-CNN models with pretrained AlexNet, GoogleNet, and VGGNet exhibit better classification accuracies when remote sensing images with larger size are fed to the networks. In addition, Table II shows that the size of the input images has a more significant impact on the classification accuracy achieved by the SF-CNN on AlexNet, which suggests that the AlexNet is more sensitive to the size of the input images. Specifically, from Table II, the accuracies obtained on the VGGNet in the categories, *school* (73.54–84.94 increase with $Pr = 20\%$) and *center* (92.31–97.75 increase $Pr = 50\%$), show significant improvements. Moreover, since the SF-CNN model based on VGGNet extracts more representative features than the other two considered models, it also exhibits the highest classification accuracy among the three considered models.

D. Comparisons With Other Methods

The proposed SF-CNN is expected to take better advantage of input images with larger spatial resolution. In this section, we compare the proposed approach with the baseline model and also with some pretrained CNN models already available in the literature [37], [38], which consider input images of fixed size. In [38], two metric functions are adopted in a new discriminative CNN (D-CNN) model to handle the problems of interclass similarity and intraclass diversity in remote sensing scene classification, but the images considered in the study must be artificially selected as a group of input data.

TABLE II
CLASSIFICATION PERFORMANCE OF THE PROPOSED METHOD OBTAINED BY USING INPUT IMAGES WITH DIFFERENT SIZES

Models	Pr=20%						Pr=50%					
	AlexNet		GoogleNet		VGGNet		AlexNet		GoogleNet		VGGNet	
Sizes	200	600	200	400	200	400	200	600	200	400	200	400
Airport	85.76	89.20	81.94	90.12	92.92	93.34	91.11	96.34	97.22	97.13	98.89	99.50
Bare land	92.74	95.12	96.77	97.65	98.19	97.12	97.42	97.09	96.13	97.97	98.71	98.77
Baseball field	93.75	97.12	97.73	92.93	99.07	96.53	98.18	98.76	100.00	99.00	99.09	99.21
Beach	95.94	96.21	95.63	97.44	98.96	98.69	100.00	99.17	100.00	99.91	100.00	100.00
Bridge	90.97	95.10	94.10	94.28	97.09	96.12	96.67	94.67	98.89	97.69	96.11	98.95
Center	79.81	80.25	80.77	81.17	81.46	81.19	93.08	91.98	91.54	90.68	92.31	97.75
Church	76.56	86.94	83.85	90.46	90.31	87.96	89.17	93.00	92.50	93.24	95.00	95.95
Commercial	75.71	85.67	80.36	80.88	92.35	91.37	92.57	99.10	95.43	97.62	98.29	98.92
Dense residential	92.99	91.42	92.99	93.87	97.47	93.54	98.05	95.77	97.07	97.47	97.07	97.19
Desert	95.42	89.54	94.17	93.38	93.54	94.52	97.33	95.00	96.00	97.24	98.67	99.39
Farmland	93.92	93.54	94.26	94.12	98.52	97.58	96.22	93.15	97.84	97.21	97.84	98.44
Forest	98.00	98.96	99.00	99.88	99.21	98.94	95.20	95.67	98.40	99.91	96.80	98.46
Industrial	78.53	83.29	76.60	80.69	85.79	89.68	89.23	89.41	86.15	84.53	92.31	93.39
Meadow	96.88	98.62	98.21	98.20	97.53	96.82	98.57	98.96	99.29	99.20	99.29	98.63
Medium residential	82.33	84.87	82.76	86.22	88.57	87.87	95.86	97.60	95.86	96.46	97.93	98.68
Mountain	97.43	97.02	98.16	99.41	99.47	99.94	98.82	99.67	98.82	99.91	99.41	100.00
Park	82.86	84.60	87.50	83.38	86.64	90.30	78.86	86.53	89.14	87.91	90.29	92.06
Parking	99.04	98.68	99.36	99.60	99.57	99.62	99.49	99.67	100.00	99.91	100.00	100.00
Playground	92.57	97.26	97.97	98.85	97.51	98.59	95.68	96.97	97.84	99.91	99.46	98.98
Pond	94.05	95.50	95.83	96.42	96.94	96.37	96.67	96.81	98.10	98.96	98.57	99.11
Port	90.79	94.37	93.75	96.60	96.59	97.31	94.74	98.09	96.32	97.80	95.79	96.38
Railway station	81.25	87.46	87.02	88.38	92.52	89.36	90.77	96.59	95.38	95.29	93.85	93.14
Resort	69.40	75.82	67.67	70.28	76.93	69.77	74.48	85.19	77.93	78.53	79.31	81.44
River	90.85	93.56	91.77	95.70	95.94	96.59	92.20	96.26	96.10	96.50	96.59	97.13
School	59.17	70.79	68.33	72.96	73.54	84.94	69.33	84.34	84.67	79.91	82.00	86.06
Sparse residential	95.83	97.88	97.08	99.63	93.96	96.19	98.67	97.00	98.67	99.24	99.33	97.45
Square	70.08	75.34	78.03	79.67	77.86	84.03	78.18	81.38	83.64	85.97	84.24	84.30
Stadium	94.40	96.94	93.10	94.85	92.88	93.04	93.10	92.77	96.55	97.15	95.86	96.61
Storage tanks	91.32	95.45	95.49	97.06	97.09	97.86	93.89	98.00	97.22	98.24	97.22	98.95
Viaduct	97.62	99.07	98.81	100.00	99.91	99.94	98.57	98.72	100.00	99.91	99.52	100.00
AA	87.86	90.85	89.97	91.48	92.94	93.17	92.74	94.79	95.09	95.35	95.66	96.50
OA	88.30	91.15	90.28	91.83	93.31	93.60	92.98	94.93	95.26	95.53	95.86	96.66
Kappa	87.88	90.83	89.93	91.51	93.06	93.38	92.73	94.76	95.09	95.37	95.71	96.54

TABLE III
CLASSIFICATION PERFORMANCE OBTAINED BY DIFFERENT METHODS ON THE UC MERCED LAND-USE DATA SET

Methods	OA
Baseline	77.71
D-CNN with AlexNet	96.67±0.10
MSCP+MRA with AlexNet	97.32±0.52
SF-CNN with AlexNet	96.98±0.24
D-CNN with GoogleNet	97.07±0.12
SF-CNN with GoogleNet	98.10±0.25
D-CNN with VGGNet	98.93±0.10
MSCP+MRA with VGGNet	98.40±0.34
SF-CNN with VGGNet	99.05±0.27

In [37], a covariance-based multilayer fusion strategy (MSCP) is proposed to exploit the highly correlated and complementary information from different layers using a multiresolution analysis (MRA) to enhance the obtained results (we will refer to this technique hereinafter as MSCP+MRA). All the results obtained after our detailed comparison (including OAs and standard deviations) are presented in Tables III–V, where we can observe that the use of remote sensing images with higher resolution helps the proposed SF-CNN model to outperform the previously developed methods for scene classification.

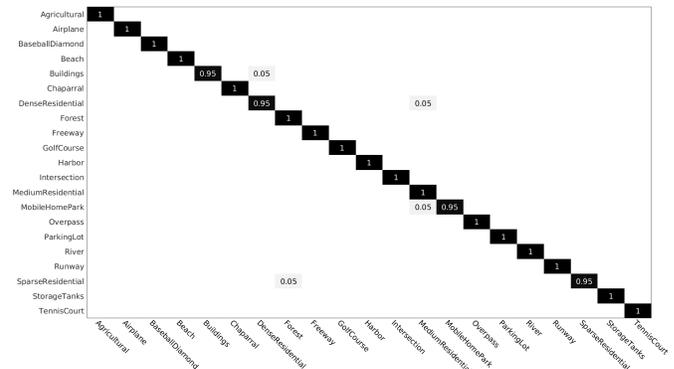


Fig. 7. Confusion matrix for the UC Merced Land-Use data set using the proposed SF-CNN pretrained with VGGNet (Pr = 80%).

1) *Experiment 1: UC Merced Land-Use Data Set:* First, we perform an experiment with the UC Merced Land-Use data set. As it can be observed from Table III, fine-tuning pretrained CNN models offers a practical strategy for the classification of small data sets. The proposed SF-CNN achieves the highest classification performance with OA superior to 99%. Furthermore, all the samples in 17 categories are classified correctly. As it can be observed in Fig. 7, there are four categories that

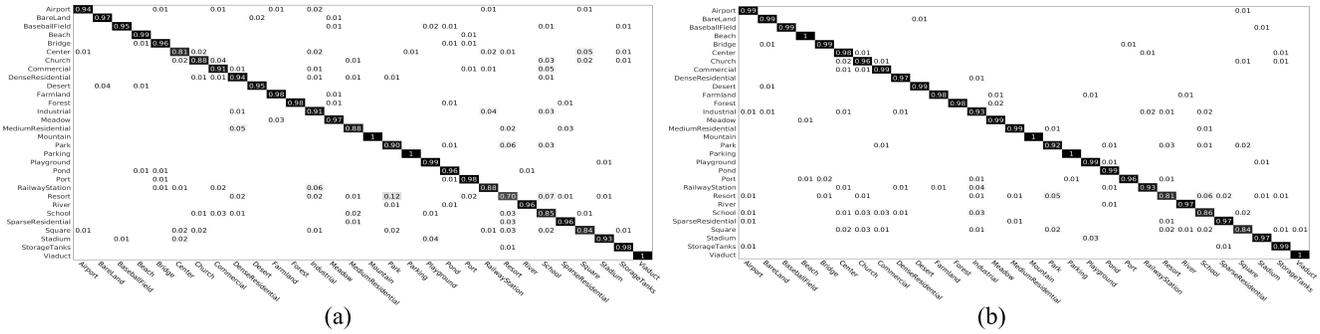


Fig. 8. Confusion matrix for the Aerial Image data set using the proposed SF-CNN with pretrained VGGNet. (a) Pr = 20%. (b) Pr = 50%.

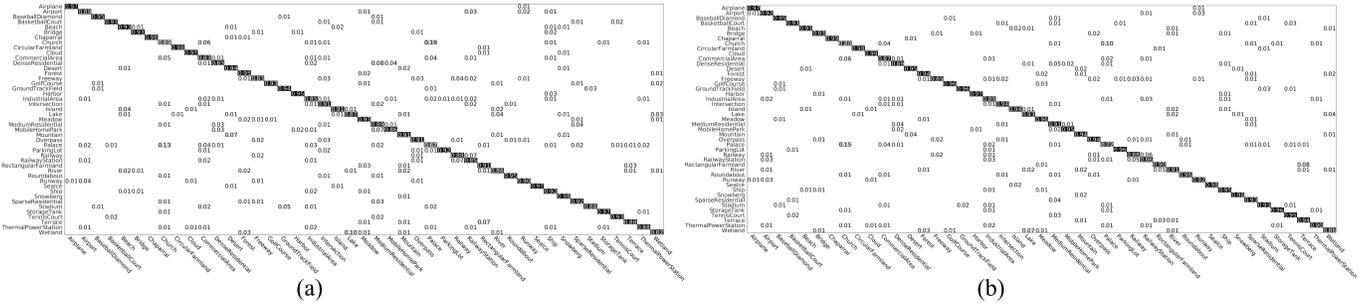


Fig. 9. Confusion matrix for the NWPU-RESISC45 data set using the proposed SF-CNN pretrained with VGGNet. (a) Pr = 10%. (b) Pr = 20%.

TABLE IV

CLASSIFICATION PERFORMANCE OBTAINED BY DIFFERENT METHODS ON THE AERIAL IMAGE DATA SET

Methods	Pr	
	20%	50%
Baseline on AlexNet	86.86±0.47	89.53±0.31
D-CNN with AlexNet	85.62±0.10	94.47±0.12
MSCP+MRA with AlexNet	90.65±0.19	94.11±0.15
SF-CNN with AlexNet	91.15±0.13	94.93±0.13
Baseline on GoogleNet	86.59±0.29	89.64±0.36
D-CNN with GoogleNet	88.79±0.10	96.22±0.10
SF-CNN with GoogleNet	91.83±0.11	95.53±0.09
Baseline on VGGNet	83.44±0.40	86.39±0.55
D-CNN with VGGNet	90.82±0.16	96.89±0.10
MSCP+MRA with VGGNet	92.21±0.17	96.56±0.18
SF-CNN with VGGNet	93.60±0.12	96.66±0.11

TABLE V

CLASSIFICATION PERFORMANCE OBTAINED BY DIFFERENT METHODS ON THE NWPU-RESISC45 DATA SET

Methods	Pr	
	10%	20%
Baseline on AlexNet	81.22±0.19	85.16±0.18
D-CNN with AlexNet	85.56±0.20	87.24±0.12
MSCP+MRA with AlexNet	83.31±0.23	87.05±0.23
SF-CNN with AlexNet	84.24±0.19	87.78±0.17
Baseline on GoogleNet	82.57±0.12	86.02±0.18
D-CNN with GoogleNet	88.89±0.10	90.49±0.15
SF-CNN with GoogleNet	87.43±0.13	90.51±0.13
Baseline on VGGNet	87.57±0.45	90.36±0.18
D-CNN with VGGNet	89.22±0.50	91.89±0.22
MSCP+MRA with VGGNet	88.07±0.18	90.81±0.13
SF-CNN with VGGNet	89.89±0.16	92.55±0.14

are misclassified (i.e., *buildings* is misclassified as *dense residential*, *dense residential* is misclassified as *medium residential*, *mobile home park* is misclassified as *medium residential*, and *sparse residential* is misclassified as *forest*), with the test errors that all equal to 0.05. Note that the number of test samples for each category is 20 in the UCM21 data set. Thus, actually, only one sample ($20 \times 0.05 = 1$) is misclassified for each of the four categories, and four samples in total are misclassified for the whole test data set.

2) *Experiment 2: Aerial Image Data Set*: Our second experiment is performed on the Aerial Image data set. Table IV shows the OAs and standard deviations obtained by three different pretrained models. In this case, we can observe that the use of additional training samples is quite helpful for improving the classification accuracies in each pretrained

CNN model. In this case, when Pr = 20%, our SF-CNN obtains the state-of-the-art results, with OA above 91% on the AlexNet and GoogleNet and OA above 93% on the VGGNet. In addition, the use of images with higher spatial resolution greatly improves the classification accuracies, especially under limited training samples. Fig. 8 shows that the following classes, *resort*, *school*, and *square*, are easily misclassified, which is due to the high interclass similarities exhibited by those classes. It should be noted that these images are also difficult to label for humans. Although the classification performance of the D-CNN based on VGGNet is better than the one achieved by the proposed method, the D-CNN needs to select image pairs as the input manually, which can be very time-consuming. In our SF-CNN, the sequence of training samples is just randomly shuffled.

3) *Experiment 3: NWPU-RESISC45 Data Set*: Our third experiment is conducted on the NWPU-RESISC45 data set. As shown in Table V, the best classification performance is obtained by the proposed SF-CNN method, which demonstrates that higher spatial resolution in the input scenes helps the CNN models to extract more discriminative features for scene classification purposes. The proposed method is the only one to obtain an OA above 92% with $Pr = 20\%$. As it can be observed in Fig. 9, some pairs of categories (e.g., *church* and *palace*, *rectangular farmland* and *terrace*, and *lake* and *wetland*) are easy to be confused, as a result of their high interclass similarity. On the other hand, some categories, such as *freeway*, *palace*, and *thermal power station*, are hard to classify, owing to their high interclass diversity. Specifically, the classification accuracy of the *palace* category is below 75%, which hinders the OA obtained for the NWPU-RESISC45 data set.

V. CONCLUSION

In this paper, a new SF-CNN model has been developed for remotely sensed scene classification purposes. The main advantage of the proposed method is that it allows the input remote sensing images to be of arbitrary sizes and does not require any resizing of such images prior to the processing. This preserves key information in high spatial resolution images, which is greatly beneficial to ultimately achieve better classification performance. Specifically, the proposed method first transfers the FCLs in the pretrained CNN model to convolutional layers and then uses a GAP layer after the final convolutional layer. Our experiments using three classic pretrained CNN models on three publicly available data sets verify the effectiveness of the proposed method when compared with other state-of-the-art approaches.

As with any new approach, there are some unresolved issues that may present challenges over time. In our method, the input images in a minibatch must have the same size. This is expected to be solved by setting the minibatch size to 1 and fine-tuning the pretrained CNN models with batch normalization layers, which is expected to require larger training and testing times that can be dealt with by developments in the GPU technology. Moreover, the lack of a sufficient number of labeled images is one of the biggest obstacles in the domain of scene classification, which can easily lead to the problem of overfitting for some complicated CNN models. In this regard, we are working on a new design of CNN models that will allow the input data to be multistructural, which may be helpful for integrating already available off-the-shelf data sets (and also for the collection of new data sets).

ACKNOWLEDGMENT

The authors would like to thank the editors and the anonymous reviewers for their valuable comments and suggestions, which greatly helped them to enhance the technical quality and presentation of this paper.

REFERENCES

- [1] N. Longbotham, C. Chaapel, L. Bleiler, C. Padwick, W. J. Emery, and F. Pacifici, "Very high resolution multiangle urban classification analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 4, pp. 1155–1170, Apr. 2012.
- [2] Y. Lv, Y. Duan, W. Kang, Z. Li, and F.-Y. Wang, "Traffic flow prediction with big data: A deep learning approach," *IEEE Trans. Intell. Transp. Syst.*, vol. 16, no. 2, pp. 865–873, Apr. 2015.
- [3] J. Leitloff, S. Hinz, and U. Stilla, "Vehicle detection in very high resolution satellite images of city areas," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 7, pp. 2795–2806, Jul. 2010.
- [4] 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.
- [5] G. Cheng, J. Han, and X. Lu, "Remote sensing image scene classification: Benchmark and state of the art," *Proc. IEEE*, vol. 105, no. 10, pp. 1865–1883, Oct. 2017.
- [6] M. J. Swain and D. H. Ballard, "Color indexing," *Int. J. Comput. Vis.*, vol. 7, no. 1, pp. 11–32, 1991.
- [7] S. Bhagavathy and B. S. Manjunath, "Modeling and detection of geospatial objects using texture motifs," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 12, pp. 3706–3715, Dec. 2006.
- [8] Y. Yang and S. Newsam, "Comparing SIFT descriptors and Gabor texture features for classification of remote sensed imagery," in *Proc. 15th IEEE Int. Conf. Image Process.*, Oct. 2008, pp. 1852–1855.
- [9] J. Fan, T. Chen, and S. Lu, "Unsupervised feature learning for land-use scene recognition," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 4, pp. 2250–2261, Apr. 2017.
- [10] B. Tu, X. Zhang, X. Kang, G. Zhang, J. Wang, and J. Wu, "Hyperspectral image classification via fusing correlation coefficient and joint sparse representation," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 3, pp. 340–344, Mar. 2018.
- [11] H. Li, H. Gu, Y. Han, and J. Yang, "Object-oriented classification of high-resolution remote sensing imagery based on an improved colour structure code and a support vector machine," *Int. J. Remote Sens.*, vol. 31, no. 6, pp. 1453–1470, Mar. 2010.
- [12] E. Aptoula, "Remote sensing image retrieval with global morphological texture descriptors," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 5, pp. 3023–3034, May 2014.
- [13] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.*, vol. 60, no. 2, pp. 91–110, 2004.
- [14] B. Luo, S. Jiang, and L. Zhang, "Indexing of remote sensing images with different resolutions by multiple features," *IEEE J. Sel. Topics Appl. Earth Observat. Remote Sens.*, vol. 6, no. 4, pp. 1899–1912, Aug. 2013.
- [15] N. He, L. Fang, S. Li, and A. J. Plaza, "Covariance matrix based feature fusion for scene classification," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Jul. 2018, pp. 3587–3590.
- [16] L. Fang, N. He, S. Li, P. Ghamisi, and J. A. Benediktsson, "Extinction profiles fusion for hyperspectral images classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 3, pp. 1803–1815, Mar. 2018.
- [17] P. Zhong and R. Wang, "Learning conditional random fields for classification of hyperspectral images," *IEEE Trans. Image Process.*, vol. 19, no. 7, pp. 1890–1907, Jul. 2010.
- [18] L. Fang, N. He, S. Li, A. J. Plaza, and J. Plaza, "A new spatial-spectral feature extraction method for hyperspectral images using local covariance matrix representation," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 6, pp. 3534–3546, Jun. 2018.
- [19] Y. Yang and S. Newsam, "Geographic image retrieval using local invariant features," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 2, pp. 818–832, Feb. 2013.
- [20] 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.
- [21] S. Chen and Y. Tian, "Pyramid of spatial relations for scene-level land use classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 4, pp. 1947–1957, Apr. 2015.
- [22] A. M. Cheriyyadat, "Unsupervised feature learning for aerial scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 1, pp. 439–451, Jan. 2014.
- [23] J. Zou, W. Li, C. Chen, and Q. Du, "Scene classification using local and global features with collaborative representation fusion," *Inf. Sci.*, vol. 348, pp. 209–226, Jun. 2016.
- [24] X. Lu, X. Zheng, and Y. Yuan, "Remote sensing scene classification by unsupervised representation learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 9, pp. 5148–5157, Sep. 2017.
- [25] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 248–255.
- [26] Z. Xu, L. Zhu, and Y. Yang, "Few-shot object recognition from machine-labeled Web images," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jul. 2017, pp. 5358–5366.

- [27] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6, pp. 1452–1464, Jun. 2018.
- [28] Y. Yu, Z. Gong, C. Wang, and P. Zhong, "An unsupervised convolutional feature fusion network for deep representation of remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 1, pp. 23–27, Jan. 2018.
- [29] P. Zhong, Z. Gong, S. Li, and C.-B. Schönlieb, "Learning to diversify deep belief networks for hyperspectral image classification," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 55, no. 6, pp. 3516–3530, Jun. 2017.
- [30] X. Yao, J. Han, G. Cheng, X. Qian, and L. Guo, "Semantic annotation of high-resolution Satellite images via weakly supervised learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 6, pp. 3660–3671, Jun. 2016.
- [31] R. Dian, S. Li, A. Guo, and L. Fang, "Deep hyperspectral image sharpening," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 11, pp. 5345–5355, Nov. 2018.
- [32] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proc. Eur. Conf. Comput. Vis.*, Sep. 2014, pp. 818–833.
- [33] L. Fang, G. Liu, S. Li, P. Ghamisi, and J. A. Benediktsson, "Hyperspectral image classification with squeeze multibias network," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 3, pp. 1291–1301, Mar. 2019.
- [34] F. Zhang, B. Du, and L. Zhang, "Scene classification via a gradient boosting random convolutional network framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 3, pp. 1793–1802, Mar. 2016.
- [35] E. Othman, Y. Bazi, F. Melgani, H. Alhichri, N. Alajlan, and M. Zuur, "Domain adaptation network for cross-scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 8, pp. 4441–4456, Aug. 2017.
- [36] E. Li, J. Xia, P. Du, C. Lin, and A. Samat, "Integrating multilayer features of convolutional neural networks for remote sensing scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 10, pp. 5653–5665, Oct. 2017.
- [37] N. He, L. Fang, S. Li, A. Plaza, and J. Plaza, "Remote sensing scene classification using multilayer stacked covariance pooling," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 12, pp. 6899–6910, Dec. 2018.
- [38] G. Cheng, C. Yang, X. Yao, L. Guo, and J. Han, "When deep learning meets metric learning: Remote sensing image scene classification via learning discriminative CNNs," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 5, pp. 2811–2821, May 2018.
- [39] Z. Gong, P. Zhong, Y. Yu, and W. Hu, "Diversity-promoting deep structural metric learning for remote sensing scene classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 1, pp. 371–390, Jan. 2018.
- [40] W. Zhao and S. Du, "Scene classification using multi-scale deeply described visual words," *Int. J. Remote Sens.*, vol. 37, no. 17, pp. 4119–4131, Sep. 2016.
- [41] F. P. S. Luus, B. P. Salmon, F. van den Bergh, and B. T. J. Maharaj, "Multiview deep learning for land-use classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 12, pp. 2448–2452, Dec. 2015.
- [42] 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.
- [43] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," in *Proc. 18th SIGSPATIAL Int. Conf. Adv. Geographic Inf. Syst.*, Nov. 2010, pp. 270–279.
- [44] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 1097–1105.
- [45] 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.
- [46] G. Wang, B. Fan, S. Xiang, and C. Pan, "Aggregating rich hierarchical features for scene classification in remote sensing imagery," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 10, no. 9, pp. 4104–4115, Sep. 2017.
- [47] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 1–9.
- [48] 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.
- [49] C.-L. Zhang, J.-H. Luo, X.-S. Wei, and J. Wu, "In defense of fully connected layers in visual representation transfer," in *Proc. Pacific Rim Conf. Multimedia*, Sep. 2017, pp. 807–817.
- [50] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jul. 2017, pp. 4700–4708.
- [51] D. Marmanis, M. Datcu, T. Esch, and U. Stilla, "Deep learning earth observation classification using ImageNet pretrained networks," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 1, pp. 105–109, Jan. 2016.
- [52] K. Nogueira, O. A. B. Penatti, and J. A. dos Santos, "Towards better exploiting convolutional neural networks for remote sensing scene classification," *Pattern Recognit.*, vol. 61, pp. 539–556, Jan. 2017. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0031320316301509>
- [53] N. He *et al.*, "Feature extraction with multiscale covariance maps for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 2, pp. 755–769, Feb. 2019.
- [54] J. Zhu, L. Fang, and P. Ghamisi, "Deformable convolutional neural networks for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 15, no. 8, pp. 1254–1258, Aug. 2018. doi: 10.1109/LGRS.2018.2830403.
- [55] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 640–651, Apr. 2017.
- [56] M. Lin, Q. Chen, and S. Yan. (Dec. 2013). "Network in network." [Online]. Available: <https://arxiv.org/abs/1312.4400>
- [57] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2016, pp. 770–778.
- [58] Y. Hou, Q. Kong, J. Wang, and S. Li. (Nov. 2018). "Polyphonic audio tagging with sequentially labelled data using CRNN with learnable gated linear units." [Online]. Available: <https://arxiv.org/abs/1811.07072>
- [59] K. Jarrett, K. Kavukcuoglu, M. Ranzato, and Y. LeCun, "What is the best multi-stage architecture for object recognition?" in *Proc. IEEE 12th Int. Conf. Comput. Vis.*, Sep./Oct. 2009, pp. 2146–2153.
- [60] W. Deng, L. Zheng, Q. Ye, G. Kang, Y. Yang, and J. Jiao, "Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 994–1003.
- [61] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Is object localization for free?—Weakly-supervised learning with convolutional neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2015, pp. 685–694.
- [62] Y.-L. Boureau, F. Bach, Y. LeCun, and J. Ponce, "Learning mid-level features for recognition," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 2559–2566.
- [63] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 1794–1801.
- [64] Y.-L. Boureau, J. Ponce, and Y. Lecun, "A theoretical analysis of feature pooling in visual recognition," in *Proc. 27th Int. Conf. Mach. Learn.*, 2010, pp. 111–118.
- [65] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [66] Y. LeCun *et al.*, "Handwritten digit recognition with a back-propagation network," in *Proc. Adv. Neural Inf. Process. Syst.*, 1990, pp. 396–404.
- [67] W. Wan, Y. Zhong, T. Li, and J. Chen, "Rethinking feature distribution for loss functions in image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 9117–9126.



Jie Xie (S'18) received the B.Sc. degree from the Hunan University of Science and Technology, Xiangtan, China, in 2015. He is currently pursuing the Ph.D. degree in control science and engineering with Hunan University, Changsha, China.

His research interests include hyperspectral image processing, remote sensing images processing, and deep learning.



Nanjun He (S'17) received the B.S. degree from the Central South University of Forestry and Technology, Changsha, China, in 2013. He is currently pursuing the Ph.D. degree with the Laboratory of Vision and Image Processing, Hunan University, Changsha.

From 2017 to 2018, he was a Visiting Ph.D. Student with the Hyperspectral Computing Laboratory, University of Extremadura, Cáceres, Spain, supported by the China Scholarship Council. His research interests include remote sensing image classification and remote sensing object detection.



Leyuan Fang (S'10–M'14–SM'17) received the B.S. and Ph.D. degrees from the College of Electrical and Information Engineering, Hunan University, Changsha, China, in 2008 and 2015, respectively.

From 2011 to 2012, he was a Visiting Ph.D. Student with the Department of Ophthalmology, Duke University, Durham, NC, USA, supported by the China Scholarship Council. Since 2017, he has been an Associate Professor with the College of Electrical and Information Engineering, Hunan University. His research interests include sparse representation and

multiresolution analysis in remote sensing and medical image processing.

Dr. Fang received the Scholarship Award for Excellent Doctoral Student granted by the Chinese Ministry of Education in 2011.



Antonio Plaza (M'05–SM'07–F'15) received the M.Sc. and Ph.D. degrees in computer engineering from the Hyperspectral Computing Laboratory, Department of Technology of Computers and Communications, University of Extremadura, Cáceres, Spain, in 1999 and 2002, respectively.

He is currently the Head of the Hyperspectral Computing Laboratory. He has authored more than 600 publications, including over 200 JCR journal papers (over 160 in the IEEE journals), 23 book chapters, and around 300 peer-reviewed conference

proceeding papers. His research interests include hyperspectral data processing and parallel computing of remote sensing data.

Dr. Plaza was a member of the Editorial Board of the *IEEE Geoscience and Remote Sensing Newsletter* from 2011 to 2012 and the *IEEE Geoscience and Remote Sensing Magazine* in 2013. He was also a member of the Steering Committee of the IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATIONS AND REMOTE SENSING (JSTARS). He was a recipient of the Recognition of Best Reviewers of IEEE GEOSCIENCE AND REMOTE SENSING LETTERS in 2009, the IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING in 2010, the Recognition as an Outstanding Associate Editor of the IEEE ACCESS in 2017, the Best Column Award of the *IEEE Signal Processing Magazine* in 2015, the 2013 Best Paper Award of the IEEE JSTARS, the Most Highly Cited Paper Award of the *Journal of Parallel and Distributed Computing* from 2005 to 2010, and the Best Paper Awards from the IEEE International Conference on Space Technology and the IEEE Symposium on Signal Processing and Information Technology. He has reviewed more than 500 manuscripts for over 50 different journals. He has guest edited 10 special issues on hyperspectral remote sensing for different journals. He served as an Associate Editor for the IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING from 2007 to 2012. He is currently an Associate Editor of the IEEE ACCESS. He served as the Director of Education Activities for the IEEE Geoscience and Remote Sensing Society (GRSS) from 2011 to 2012 and the President of the Spanish Chapter of the IEEE GRSS from 2012 to 2016. He served as the Editor-in-Chief of the IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING from 2013 to 2017.