



Research Article

Hybrid Depth-Separable Residual Networks for Hyperspectral Image Classification

Cuijie Zhao,^{1,2} Hongdong Zhao¹, Guozhen Wang,³ and Hong Chen¹

¹College of Electronic Information Engineering, Hebei University of Technology, Tianjin 300401, China

²Tianjin University of Finance and Economics, Pearl River College, Tianjin 301811, China

³Department of Computer Science and Technology, Tianjin University Renai College, Tianjin 300000, China

Correspondence should be addressed to Hongdong Zhao; zhaohd@hebut.edu.cn

Received 13 July 2020; Revised 10 August 2020; Accepted 12 August 2020; Published 26 August 2020

Academic Editor: Zhihan Lv

Copyright © 2020 Cuijie Zhao et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

At present, the classification of the hyperspectral image (HSI) based on the deep convolutional network has made great progress. Due to the high dimensionality of spectral features, limited samples of ground truth, and high nonlinearity of hyperspectral data, effective classification of HSI based on deep convolutional neural networks is still difficult. This paper proposes a novel deep convolutional network structure, namely, a hybrid depth-separable residual network, for HSI classification, called HDSRN. The HDSRN model organically combines 3D CNN, 2D CNN, multiresidual network ROR, and depth-separable convolutions to extract deeper abstract features. On the one hand, due to the addition of multiresidual structures and skip connections, this model can alleviate the problem of over fitting, help the backpropagation of gradients, and extract features more fully. On the other hand, the depth-separable convolutions are used to learn the spatial feature, which reduces the computational cost and alleviates the decline in accuracy. Extensive experiments on the popular HSI benchmark datasets show that the performance of the proposed network is better than that of the existing prevalent methods.

1. Introduction

HSI has been widely used in environmental monitoring [1], mineral exploration [2], agricultural remote sensing [3], vegetation ecology [4], ocean remote sensing [5], and other earth observation tasks. In these applications, because HSI exhibits mixed land cover categories, resulting in high intraclass variability and interclass similarity, it is a huge challenge for any classification model. In order to improve the performance of HSI classification, traditional machine learning methods integrate spectral features and spatial features to achieve effective feature extraction such as random forest [6], SVM and its variants [7, 8], sparse self-representation [9], and artificial neural network [10]. However, these methods only extract surface features such as edges and textures of HSI, which will reduce the feature representation ability of hyperspectral images.

Deep learning [11, 12] method has been widely used in image processing, especially image classification [13] and

target recognition [14]. It can actively learn to extract features, with little manual intervention, and automatically find effective features. The deep model can also extract high-level abstract features by adding hierarchical abstractions, which are usually more robust to nonlinear processing. The basic network frameworks for deep learning include unsupervised neural networks [15], convolutional neural networks [16], cyclic neural networks [17], and recursive neural networks [18]. Among them, convolutional neural networks (CNNs) are classic models. This model is based on big data and deep network structures. It extracts rich deep features from the original hyperspectral data, ensuring the integrity of spatial and spectral information and avoiding the initiative and randomness of human feature extraction, and it can achieve better classification results than other deep learning models. In addition, the deep convolutional networks AlexNet [19], VGG [20], GoogleNet [21], ResNet [14], etc., perform well, which fully demonstrates the fact that convolutional neural networks are a good strategy for image classification.

In recent years, the classification of HSI based on the deep convolutional network has also made great progress. HSI is 3D data, which include 2D spatial information and 1D spectral information [22]. The classification of HSI based on CNN mainly uses 2D CNN or 3D CNN for hierarchical feature extraction. Chen et al. proposed to use 3D convolution to get spatial-spectral depth features [23]. A 3D CNN was designed by Li for spectral space classification [24]. Subsequently, Yang et al. proposed a 3D recurrent CNN [25]. Song et al. designed a deep feature fusion network to solve the hyperspectral classification problem [26], Fang et al. proposed a deep hash neural network [27], and Gong et al. introduced statistical metric methods for HSI spectral-spatial classification [28]. Zhong et al. designed the SSRN structure and introduced identity mapping residual blocks for spectral-spatial feature learning [22]. Liu et al. used 3D CNN and residual connection to construct a 12-layer deep network (Res-3D CNN) [29]. Lee and Kwon adopted the residual connection to make the network deeper and used 11 convolution kernels to learn the hierarchical features [30]. Recently, Roy et al. designed a concise model (HybridSN) [31] that combined 3D CNN and 2D CNN to extract spectral features, and by comparing with the most advanced models in the past, the classification effect is more excellent. Cao and Guo further introduced hybrid dilated convolutions (HDC) and the residual block based on SSRN and proposed a new end-to-end hybrid expansion residual deep convolutional network [32]. Wu et al. designed the 3D ResNeXt structure using feature fusion and label smoothing strategies [33].

It can be clearly seen from the above literature that only using 2D CNN for HSI classification cannot extract a good distinguishing feature map from the spectral dimension, and spectral information will be lost. Using only 3D CNN to extract features can enhance the accuracy, but the complexity will grow, and the performance is worse, when dealing with classes with similar textures in many spectral bands. The main reason for the above situation is that HSI is a 3D data image, with spatial dimensions and spectral dimensions. Using only 2D CNN cannot extract feature images with good discriminating ability from the spectral dimension. Similarly, deep 3D CNN complicates the model and greatly increases the amount of calculation. Moreover, many classes with similar textures, using 3D CNN alone, seem to be even worse.

In order to solve the above shortcomings, our model combines 3D CNN and 2D CNN that can fully extract spectral and spatial feature maps and overcome the shortcomings of single 2D CNN and 3D CNN. In addition, deep networks are difficult to train and are prone to problems such as overfitting, gradient disappearance, and gradient explosion. Therefore, it is natural to introduce the residual network in the model [34, 35] because the residual network can be improved by adding layers. With the depth of the network increasing, the operation costs in the model will also increase. Replacing the traditional 2D convolution with deep separable convolutions [36] can solve the problem of parameter and operation costs and further avoid overfitting. In this study, a new network is proposed by constructing a hybrid deep separable residual network. Firstly, the

framework designs a 3D residual module ROR to extract spatial-spectral mixing features. Subsequently, the feature information is converted from 3D data to a 2D feature map. Finally, the 2D depth-separable convolutions extract spatial features. Depth-separable convolution can enhance the feature learning ability of HSI and reduce the computational complexity. The multiresidual network ROR [37] in the network can enhance the learning ability. In addition, skip connections can extract spatial-spectral mixing features more effectively. The major contributions of this paper are listed as follows:

- (1) A hybrid depth-separable convolution residual network is proposed to enhance the feature learning ability of HSI. In this network, spatial-spectrum 3D CNN and spatial 2D CNN are combined into the model, which can better study deep-level spatial-spectral features.
- (2) We embed multiresidual network ROR in the 3D convolutional layer and the 2D convolutional layer that can greatly decrease the number of parameters, thereby simplifying the network structure and promoting the extraction of deep features.
- (3) In the 2D processing part, the use of deep separable convolutional layers reduces the number of parameters and avoids overfitting.

The rest of the article is organized as follows: in Section 2, we introduce the proposed framework. In Section 3, the HSI set and network configuration are explained, with experimental results and analyses. The conclusions and future directions are given in Section 4.

2. Methodology

2.1. Proposed Model. The HSI data can be regarded as a 3D cube, where the width of the cube is defined as W , the height is H , the spectral band is D , and the original input of hyperspectral data can be expressed as $R^{W \times H \times D}$. Each pixel contains D spectral measurements and forms a one-hot label vector $C = (c_1, c_2, \dots, c_n) \in R^{1 \times 1 \times n}$, where n represents the land cover categories. We design a HDSRN structure for HSI classification, including six parts, namely, PCA dimensionality reduction, 3D spectral and spatial mixed feature learning process, 3D to 2D deformed part, 2D learning process, average pooling layer, and FC layer. In Figure 1, the HDSRN classification framework is described in detail using the Indian Pines dataset as an example, where $W = H = 145$ and $D = 200$.

In the traditional 2D convolution operation, convolution is applied to the spatial dimension, and the 2D feature map is obtained. However, the HSI is 3D data, and it is necessary to capture the spectral feature; 2D CNN cannot process spectral information. The 3D CNN kernel can simultaneously extract spectral and spatial features, but it increases the computational complexity. In addition, when classifying a large number of features with similar textures on the spectral band, the performance is poor. In order to overcome the shortcomings of 2D CNN and 3D CNN and make full

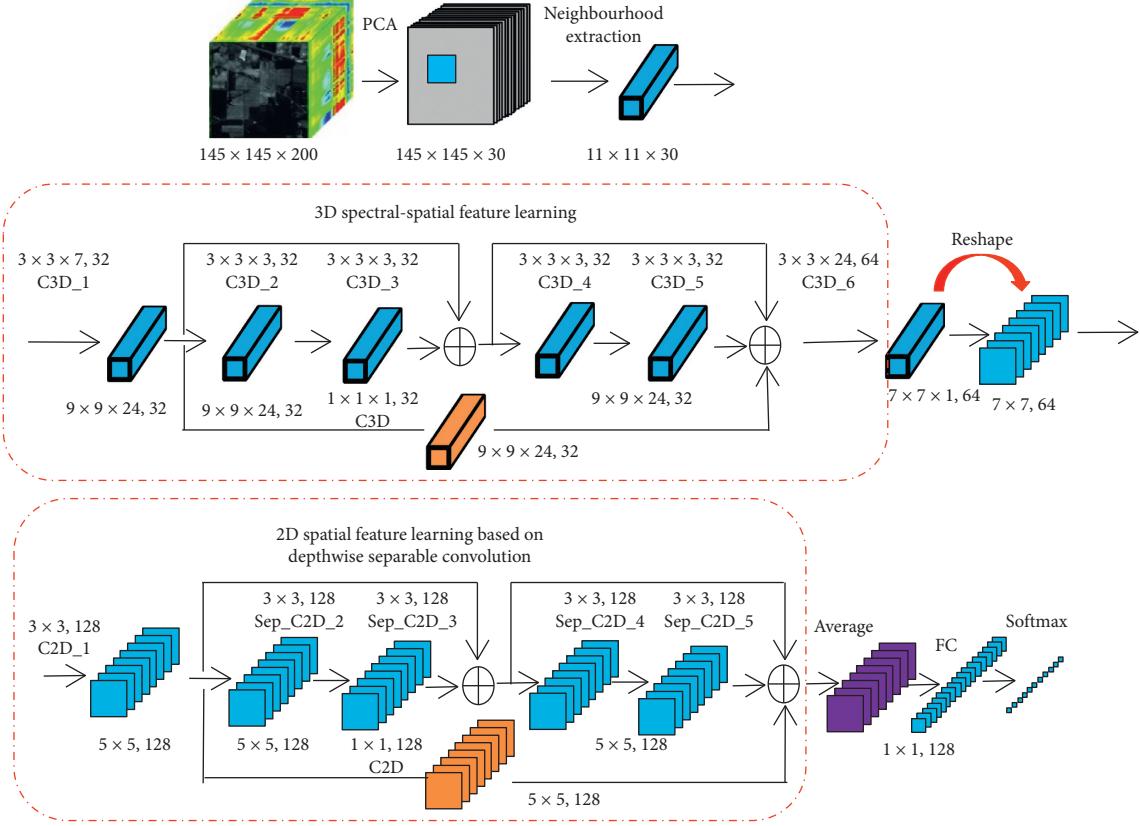


FIGURE 1: The HDSRN classification framework.

use of the automatic feature learning capabilities of 2D and 3D, the model uses a hybrid convolutional neural network framework to combine 3D CNN and 2D CNN into the model. The 3D spatial-spectral feature extraction part consists of two 3D convolutional layers and a set of multiresidual network ROR. The 2D spatial feature is composed of a 2D convolutional layer and a set of multiresidual network ROR. The reshape and average pooling modules are used to adjust the data size to meet the requirements of the next layer. The FC layer and the softmax layer are used to classify HSI.

2.2. PCA Dimension Reduction and Data Preprocessing. In order to remove spectral redundancy, based on the lightweight design, we use PCA to reduce the dimensionality of the original hyperspectral data, retaining relatively few principal components. In Figure 1, PCA decreases the spectral band from D to B , and $B=30$, while the spatial dimension remains unchanged. In this case, the hyperspectral data are represented as $R^{W \times H \times B}$. The HSI image has large size and many bands, and direct processing requires high hardware and memory requirements. Therefore, before image classification and processing, we first divide the HSI image into small pieces. The hyperspectral image is divided into small overlapping 3D patches, expressed as $R^{S \times S \times B}$. The height and the width are both S , and the spectral band is B . The true value labels are decided by the label of the center pixel. The size of the patches cannot be too small because

being too small will result in a small receiving field and cannot fully extract image features. But, if the patches are too large, the amount of calculation will be large, and the training time and the test time will become longer. Through the comparative experiment in Section 3.2, we find that the classification effect is better when the S value is 11.

2.3. 3D Spectral Space Feature Learning. HSI data have the characteristics of a spectral-spatial 3D structure. Based on this feature, we construct a 3D convolution network suitable for HSI to extract spatial-spectral features. The 3D convolution operation [38] is achieved by convolving a 3D convolution kernel with 3D data. The input layer is a 3D image, which is composed of a spatial dimension and a spectral dimension. The 3D convolution kernel performs convolution operations on the two dimensions of the input 3D image and obtains a 3D feature map. The 3D convolution formula is as follows:

$$v_{i,j}^{x,y,z} = f \left(\sum_m \sum_{b=0}^{B_i-1} \sum_{h=0}^{H_i-1} \sum_{w=0}^{W_i-1} k_{i,j,m}^{b,h,w} v_{(i-1),m}^{(x+h),(y+w),(z+b)} + b_{i,j} \right), \quad (1)$$

where $v_{(i-1),m}^{(x+h),(y+w),(z+b)}$ is the value at the position $(x+h, y+w, z+b)$ of the m^{th} feature map output from the $i-1$ layer, $k_{i,j,m}^{b,h,w}$ is the value of the j^{th} convolution kernel of the i^{th} layer at the position (b, h, w) , B_i , H_i , and W_i are the kernel sizes along the spectral and spatial dimensions, respectively,

(b, h, w) is the index of the convolution kernel, (x, y, z) is the index of the feature map, (x, y) is the spatial dimension value, z is the spectral dimension value, $b_{i,j}$ is the deviation of the j^{th} feature map on the i^{th} neuron, and f is the activation function. In the hyperspectral classification task based on 3D convolution, we set the input data to $W \times H \times B$, C_1 , the size of the convolution kernel, is $k_1 \times k_2 \times k_3$, and the number of convolution kernels is p . If there is no padding and the step size is 1 in the convolution operation, the feature map size generated by 3D convolution is $(W - k_1 + 1) \times (H - k_2 + 1) \times (B - k_3 + 1)$, p . The number of weight parameters of the 3D convolutional layer is $p \times k_1 \times k_2 \times k_3 \times C_1$.

Figure 2 is the 3D spatial-spectral feature learning framework. This part is composed of two 3D convolutional layers and a set of multiresidual network ROR. The input data of the network are $11 \times 11 \times 30$, the size of the convolution kernel of the first layer is $3 \times 3 \times 7$, and the output is 32 feature maps of $9 \times 9 \times 24$ size. The second layer is a multiresidual network ROR, which uses 3D identity residual blocks to connect to deepen the network, avoiding weak signal loss and excessive fitting, which is conducive to improving efficiency and extracting better deep abstract features without introducing additional parameters. This part uses padding to ensure that the size of the output feature map is the same as the input size. In addition, the outermost skip connection realizes the fusion of feature data by summing the corresponding pixels. Skip connection alleviates the problem of gradient disappearance, contributes to gradient backpropagation, and can more fully extract features. On each convolutional layer connected by the residual block, we use 32 convolution kernels. The size of each kernel is $3 \times 3 \times 3$, from which we can get rich spectral and spatial features.

2.4. 3D to 2D Deformation. The 2D convolution operation is started, and spatial feature extraction is performed, after the 3D convolution operation. The network outputs 64 feature map data with a size of $7 \times 7 \times 1$. In order to learn the output features in the later 2D space, we reshape the 3D features into 64 2D feature maps of size 7×7 , as shown in Figure 3. After the reshaping operation, only the 2D spatial features need to be studied, which reduces the parameters and the operation cost compared with 3D convolution.

2.5. 2D Spatial Feature Learning Based on Depth-Separable Convolution. Convolutional layer is the core part of CNN to extract deep-level features. The main functional unit of the convolutional layer is a 2D convolution kernel, which acts on the input data of the previous layer to extract features and enhances the model's nonlinear feature extraction capability by adding an activation function, which is beneficial to the extraction of complex deep features. The 2D convolution is shown in the following equation:

$$v_{i,j}^{x,y} v_{(i-1),m}^{(x+h),(y+w)} = f \left(\sum_m \sum_{h=0}^{H_i-1} \sum_{w=0}^{W_i-1} k_{i,j,m}^{h,w} v_{(i-1),m}^{(x+h),(y+w)} + b_{i,j} \right), \quad (2)$$

where $v_{(i-1),m}^{(x+h),(y+w)}$ is the value at the position $(x + h, y + w)$ of the m^{th} feature map output from the $i - 1$ layer, $k_{i,j,m}^{h,w}$ is the value of the j^{th} convolution kernel of the i^{th} layer at the position (h, w) , B_i and H_i are the kernel sizes along the spatial dimensions, (h, w) is the index of the convolution kernel, (x, y) is the index of the feature map, $b_{i,j}$ is the deviation of the j^{th} feature map on the i^{th} neuron, and $f()$ is the activation function. In the hyperspectral classification task based on 2D convolution, we set the input data to $W \times H$, C_1 , the size of the convolution kernel is $k_1 \times k_2$, and the number of convolution kernels is p . Through the 2D convolution operation, the current layer can learn local features from the previous layer, and the size of the convolution kernel determines the size of the local space. Continuous convolutional layers can extract deeper and deeper features that are increasingly abstract. Continuous convolutional layer can extract deeper and more abstract features.

Unlike accustomed 2D convolution, depth-separable convolution performs a spatial convolution while keeping the channels independent, and then deep convolution is performed. The depth-separable convolution not only lessens the number of parameters and calculations in the network, but also speeds up the network training speed and reduces the odds of overfitting in HSI classification, as shown in Figure 4.

In general, the input image using convolution operation has the same length and width. Assuming that the input image size is $H \times H$ and the channel is D_1 , that is, the depth of the image, D_2 convolution kernel size $h \times h$ convolution operation is used. The specific steps are as follows.

The first step is the deep convolution operation, which uses D_1 convolution kernels of size $h \times h$ to perform the convolution operation. Each convolution kernel only convolves one channel of the input layer. If no padding is used and the step size is one, the mapping size obtained each time is shown in the following equation:

$$S_M = (H - h + 1)^2, \quad (3)$$

where H is the length and width of the input image and h is the convolution kernel size.

These maps are stacked together to create an image, and the size of the image is shown in the following equation:

$$S_h = (H - h + 1)^2 D_1, \quad (4)$$

where D_1 represents the number of channels. Finally, an output image of size S_h is obtained, and the depth of the image remains the same as the original.

The second step is to expand the depth convolution operation. We use D_1 convolution kernels of size 1×1 to perform the operation. Each convolution kernel is convolved with the input image of size S_h to obtain a size of S_M mapping. After D_2 times 1×1 convolution, the output image can be obtained. The size of the output image is shown in the following equation:

$$I_o = (H - h + 1)^2 D_2. \quad (5)$$

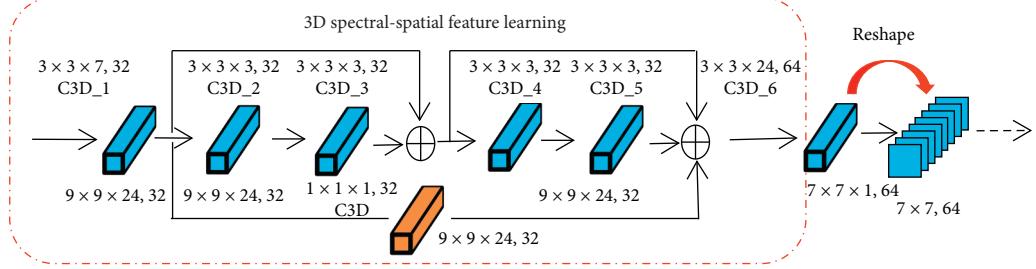


FIGURE 2: 3D spectral-spatial feature learning.

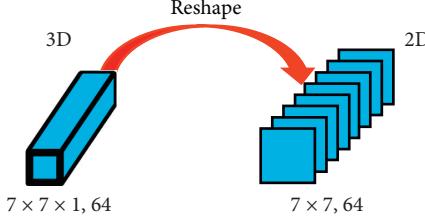


FIGURE 3: Framework of 3D to 2D deformable.

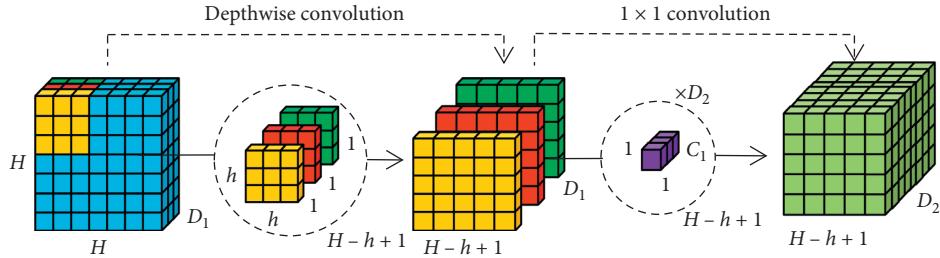


FIGURE 4: Depth-separable convolution.

The depth-separable convolution converts an input layer of $H \times H \times D_1$ into an output layer of I_o .

We compare the number of parameters used in the two methods. The number of traditional 2D convolution operation parameters is shown in the following equation:

$$P_t = D_1 h^2 D_2. \quad (6)$$

The number of depth-separable convolution parameters is shown in the following equation:

$$P_d = D_1 h^2 + D_1 D_2, \quad (7)$$

where $D_1 h^2$ represents the number of operation parameters in the first step and $D_1 D_2$ is the number of parameters in the second step.

The ratio of the depth-separable convolution and the traditional 2D convolution parameter is shown in equation (8). It can be found that the use of deep separable convolution can greatly reduce the number of parameters and improve operational efficiency:

$$R_p = \frac{1}{D_2} + \frac{1}{h^2}. \quad (8)$$

Figure 5 is a spatial feature learning framework based on depth-separable convolution. This part consists of a

convolutional layer and a set of multiresidual network ROR. In Sep_C2D_2, Sep_C2D_3, Sep_C2D_4, and Sep_C2D_5, the depth-separable convolution replaces the traditional 2D convolution operation, which greatly reduces the number of parameters and the operation cost. The input data of the network are 64 feature maps with a size of 7×7 . Firstly, this part of the network uses 128 convolution kernels of size 3×3 to realize convolution operations on the input data. Secondly, it uses identity residual blocks to connect and deepen the network to extract better deep abstract features. We use padding to keep the size of the output feature map unchanged from the input size. Similarly, the outermost skip connection realizes the fusion of feature data. On each convolutional layer connected by the residual block, we use 128 convolution kernels, and each kernel size is 3×3 . Finally, through the average pooling layer, the 128 output feature maps are converted into 128 special maps with a size of 1×1 .

2.6. Residual Network. In deep learning, the shallow network cannot significantly enhance the classification effect of the network, and the deep network can better learn abstract features [11, 14]. But, the deeper the network, the more obvious the phenomenon of gradient disappearing, and the

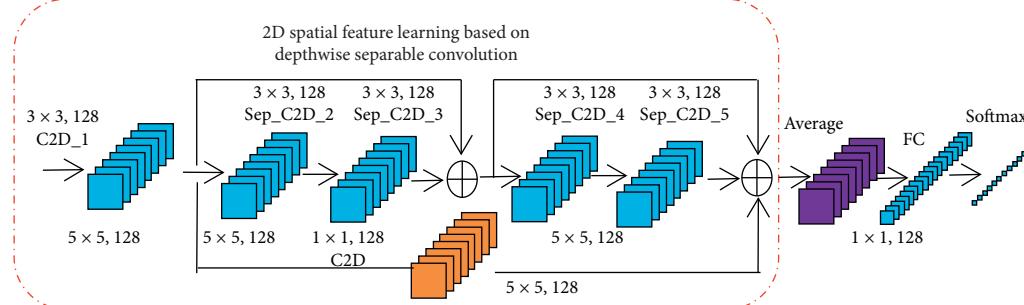


FIGURE 5: 2D spatial feature learning based on depth-separable convolution.

training effect of the network will not be very good. The proposed classification model uses residual connection to deepen the network to solve the gradient dispersion. In addition, the use of residual connections does not introduce additional parameters. According to whether the input and output sizes are the same, ResNet [14] divides the residual connection into identical residual connection and non-identity residual connection, as shown in Figure 6. In the identity residual connection, the input data are X , and the identity connection is used to inject X directly into the downstream of the network. The dimension of the output data has not changed, as shown in Figure 6(a). In the nonidentity residual connection, the size of the input data and the output data does not match. A convolutional layer is added to the shortcut path to adjust the input data x to a suitable size to match the data size of the main path, as shown in Figure 6(b). He et al. [14] show that the identity residual connection can effectively solve network degradation.

Due to the high spectral resolution and high spatial correlation of HSI, we design two consecutive identical residual blocks in the proposed HDSRN model, as shown in Figure 7. In this structure, the gradient in the upper layer can quickly propagate back to the lower layer, thereby facilitating and standardizing the model training process. In addition, we added a skip connection to fuse feature data at the outermost layer of two consecutive identical residual blocks. Skip connection can alleviate the problem of gradient disappearance, help the backpropagation of the gradient, and more fully extract features.

3. Experimental Results and Discussion

In this section, we first select three popular HSI benchmark datasets for our experiment and describe evaluation indices and experimental settings. Then, we discuss the impact of input spatial dimensions and proportions of training samples on classification performance. Finally, we compare the HDSRN model with the existing state-of-the-art methods, such as SVM-RBF [39], 2D CNN [40], 3D CNN [24], SSRN [22], and HybridSN [31].

3.1. Datasets. We use Indian Pines (IN), University of Pavia (UP), and Salinas Scene (SA) datasets [39] to verify the model.

IN dataset [41, 42] was collected by AVIRIS sensors in northwestern Indiana. The spatial resolution is about 20m, and the image has a 145×145 spatial dimension and 224 spectral bands. Among them, 24 spectral bands covering the water absorption area are discarded, so it becomes 3D data of size $145 \times 145 \times 200$. The data have 21,025 pixels in total, and the characteristic pixels are 10,249 pixels. The planting area is all agricultural crops with a total of 16 categories.

UP dataset [41, 42] was collected by ROSIS sensors in Pavia University. The spatial resolution is 1.3 m, the image has a 610×340 spatial dimension and 103 spectral bands, and its wavelength range is between 430 and 860 nm. The UP dataset contains a total of 2,207,400 pixels, and the feature pixels are only 42,776 and contains 9 categories.

SA dataset [41, 42] was collected by the AVIRIS sensors in Salinas Valley, California, USA. The spatial resolution is 3.7 m, and the image has a 512×217 spatial dimension and 224 spectral bands. 20 spectral bands covering the water absorption area are discarded, so it becomes 3D data of size $512 \times 512 \times 204$. The SA dataset included a total of 111104 pixels, 56975 pixels are background pixels, and the feature pixels are 54129 pixels with a total of 16 categories.

3.2. Evaluation Indices. In this paper, we use the confusion matrix to represent the classification accuracy of HSI. Confusion matrix is an indicator to judge the results of the classification model and is used to judge the quality of the classifier. The judgment indexes customarily used are overall accuracy (OA), average accuracy (AA), and kappa coefficient (kappa).

OA refers to the ratio of the number of correctly classified samples to the overall number of samples, as shown in the following equation:

$$OA = \frac{\sum_{i=1}^n M_{ii}}{N}, \quad (9)$$

where N is the overall number of samples, M_{ii} is the diagonal element of the corresponding confusion matrix, and n is the number of categories.

AA is the average value of the classification accuracy of each category, and its calculation process is shown in the following equation:

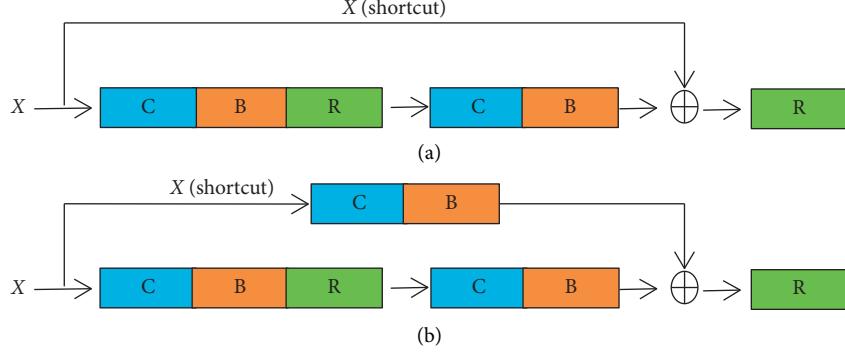


FIGURE 6: The residual connections: (a) identity connection; (b) nonidentity connection.

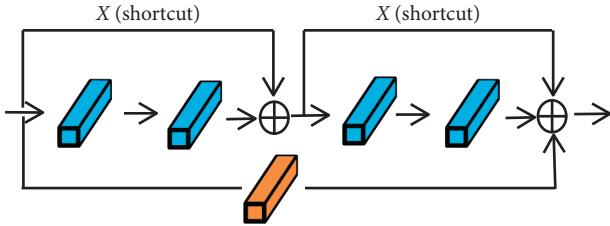


FIGURE 7: Spectral-spatial residual network.

$$AA = \frac{\sum_{i=1}^n \left(M_{ii} / \sum_{j=1}^n M_{ij} \right)}{n}. \quad (10)$$

Kappa coefficient is an important indicator to measure classification performance, and its calculation process is shown in the following equation:

$$\kappa = \frac{N \sum_{i=1}^n M_{ii} - \sum_{i=1}^n (M_{i+} + M_{+i})}{N^2 - \sum_{i=1}^n (M_{i+} + M_{+i})}, \quad (11)$$

where M_{i+} is the sum of the i^{th} row in the confusion matrix and M_{+i} represents the sum of the i^{th} column in the confusion matrix.

3.3. Experimental Settings. We studied the impact of the sample ratio and the input space size on classification performance. For specific experimental details, see Section 3.3. In the end, the most suitable space size we chose is 11×11 . For the fairness of comparison, we collected the same spatial dimensions 11×11 in the input 3D patches of different datasets. Taking the Indian Pines dataset as an example, the detailed network parameter settings of the model proposed in this paper are shown in Table 1. We applied the Adam algorithm and chose the optimal learning rate of 0.001. The experimental hardware platform is i5-7500 CPU and GTX960 GPU. For these three datasets, in order to get more accurate statistical results, each experiment was repeated 5 times, and the average value of the classification indicators was used as the final result.

3.4. Experimental Parameter Discussion. The ratio of training samples and the input space size are two important factors that affect the performance of HSI classification. We use a series of comparative experiments to determine the training sample ratio and the size of the input space. When the training sample ratios of the datasets IN, UP, and SA are 20%, 10%, and 10%, respectively, the classification performance is optimal. Through the experiment in Section 3.4.2, the classification performance is the best when the input space size is 11×11 .

3.4.1. The Impact of the Training Dataset Proportion. In order to select the appropriate training sample ratio for the dataset, we conducted a comparative experiment with different training sample ratios. The proportions of the training samples we selected are 2%, 5%, 10%, 15%, and 20%. The same spatial dimension is extracted for different datasets. For example, the spatial dimension of IN is $11 \times 11 \times 30$, the spatial dimension of UP is $11 \times 11 \times 15$, and the spatial dimension of SA is $11 \times 11 \times 15$. Each experiment was repeated 5 times, and the average value of the classification indicators was used as observation objects. The changes in OA, AA, and kappa values at different ratios are shown in Table 2. It can be found that as the ratio of input training samples increases, the accuracy gradually improves. When the proportion of training samples in the three datasets reaches 20%, the classification accuracy exceeds 99%, especially the accuracy of the SA dataset is close to 100%. Figure 8 is a classification accuracy chart of three datasets with different training sample proportions. We can clearly see that when the training set is between 2% and 15%, the accuracy is significantly improved. When our training set reaches 15%, the accuracy began to grow slowly. The training sample ratio reaches 20%, and the classification accuracy improves to 99.6%. For the IN dataset, a 20% training sample ratio is enough to train the network. It can be clearly seen from Figures 8(b) and 8(c) that when the training set is between 2% and 10%, the classification accuracy increases significantly. As the training set reaches 10%, the increasing rate of accuracy starts to slow down. The UP classification accuracy reaches 99.73%, and the SA reaches 99.98%. For UP and SA datasets, 10% of the training samples are selected to train the network.

TABLE 1: The architecture of the proposed network.

Layer	Output	Convolution kernel size	Convolution kernel number
Input	11, 11, 30, 1		
C3D_1	9, 9, 24, 32	(3, 3, 7)	32
C3D_2	9, 9, 24, 32	(3, 3, 3)	32
C3D_3	9, 9, 24, 32	(3, 3, 3)	32
C3D_4	9, 9, 24, 32	(3, 3, 3)	32
C3D_5	9, 9, 24, 32	(3, 3, 7)	32
C3D	9, 9, 24, 32	(1, 1, 1)	32
C3D_6	7, 7, 1, 64	(3, 3, 24)	64
Reshape	7, 7, 64		
C2D_1	5, 5, 128	(3, 3)	128
Sep_C2D_2	5, 5, 128	(3, 3)	128
Sep_C2D_3	5, 5, 128	(3, 3)	128
Sep_C2D_4	5, 5, 128	(3, 3)	128
Sep_C2D_5	5, 5, 128	(3, 3)	128
C2D	5, 5, 128	(1, 1)	128
Average pooling	1, 1, 128		
FC	16		
Output	Categories of land cover		

TABLE 2: Classification accuracy of different training sample ratios.

Ratios (%)	IN datasets			UP datasets			SA datasets		
	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)
2	86.69	82.73	85.87	98.39	96.63	97.67	99.19	98.83	99.11
5	94.99	95.02	93.73	99.40	98.71	98.88	99.81	99.41	99.63
10	98.49	98.05	98.16	99.86	99.83	99.81	99.98	99.98	99.98
15	99.03	99.00	98.97	99.87	99.85	99.86	99.99	99.98	99.98
20	99.72	99.60	99.70	99.89	99.86	99.87	100	100	100

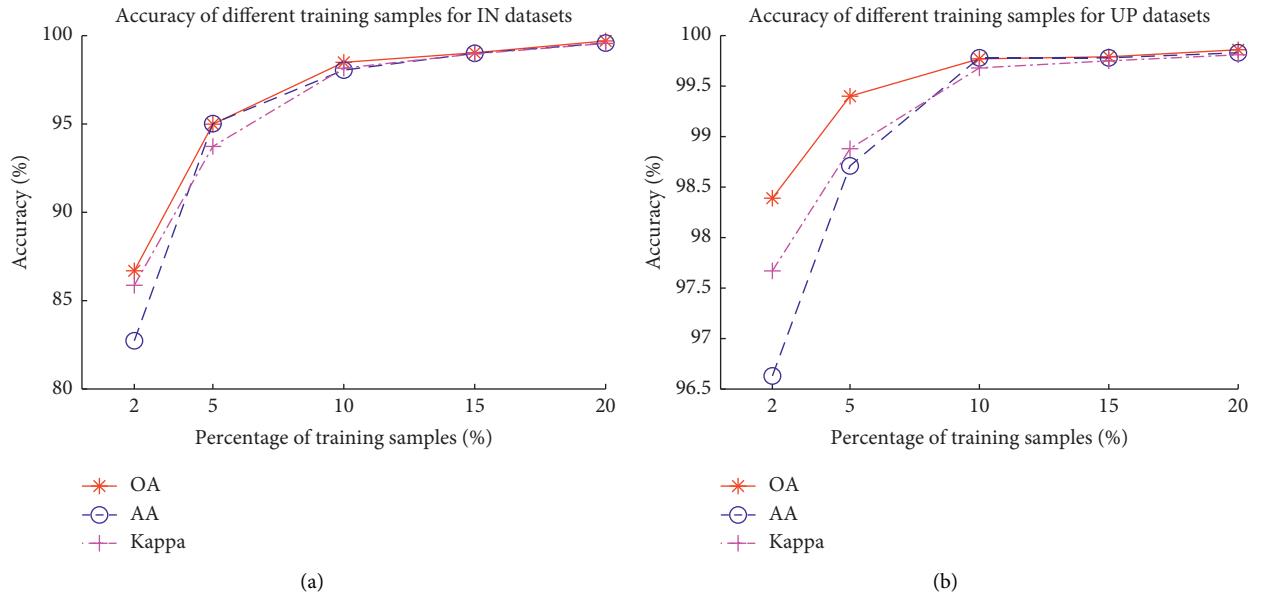


FIGURE 8: Continued.

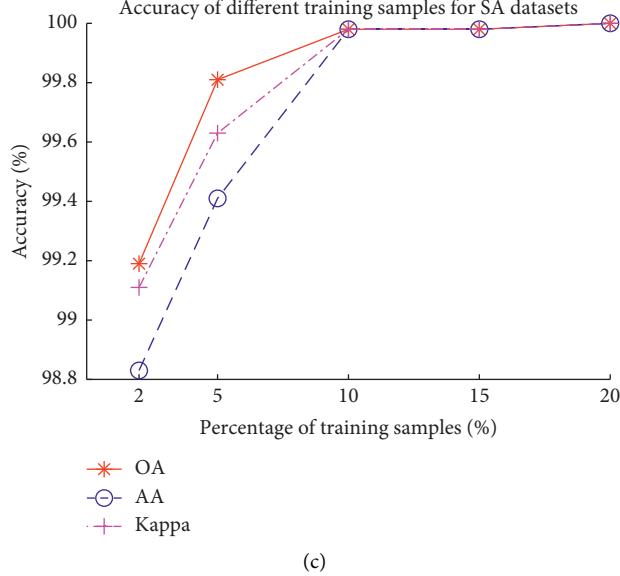


FIGURE 8: Accuracy of different training samples for three datasets: (a) IN; (b) UP; (c) SA.

TABLE 3: Accuracy of different spatial sizes.

Spatial size	IN datasets			UP datasets			SA datasets		
	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)
5 × 5	95.86	96.08	95.09	98.36	98.20	98.25	98.89	98.88	98.85
7 × 7	98.96	98.25	98.64	99.29	99.35	99.27	99.58	99.45	99.57
9 × 9	99.30	99.31	99.24	99.68	99.66	99.58	99.87	99.86	99.88
11 × 11	99.72	99.60	99.70	99.86	99.83	99.81	99.98	99.98	99.98
13 × 13	99.65	99.61	99.60	99.87	99.83	99.85	100	100	100

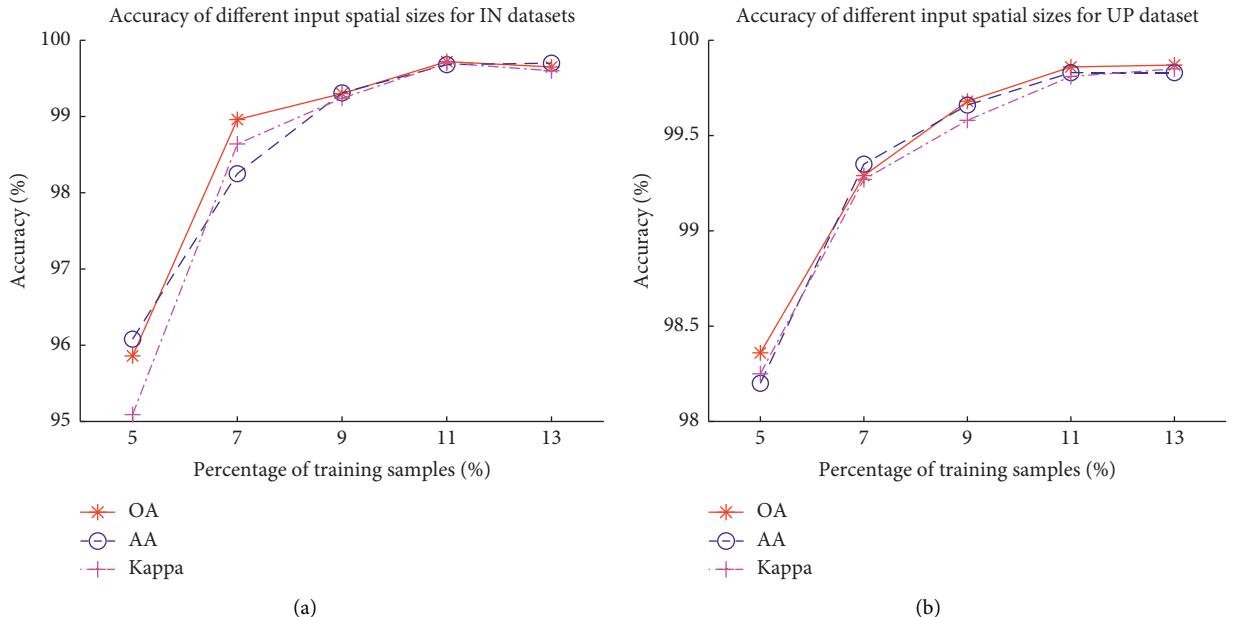


FIGURE 9: Continued.

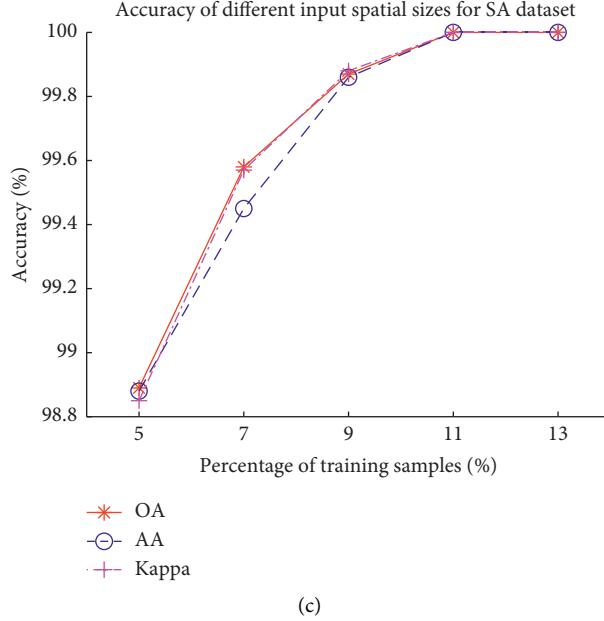


FIGURE 9: Accuracy of different spatial sizes: (a) IN; (b) UP; (c) SA.

TABLE 4: Training time and test time of different spatial sizes.

Spatial size	IN datasets		UP datasets		SA datasets	
	Training time (s)	Test time (s)	Training time (s)	Test time (s)	Training time (s)	Test time (s)
5 × 5	94.9	1.6	99.9	2.3	100.0	3.1
7 × 7	173.5	1.9	225.7	3.5	351.1	5.2
9 × 9	323.7	2.6	500.7	5.3	650.3	6.2
11 × 11	577.1	3.0	787.3	6.8	926.3	8.2
13 × 13	995.1	6.2	1296.6	11.6	1537.0	14.3

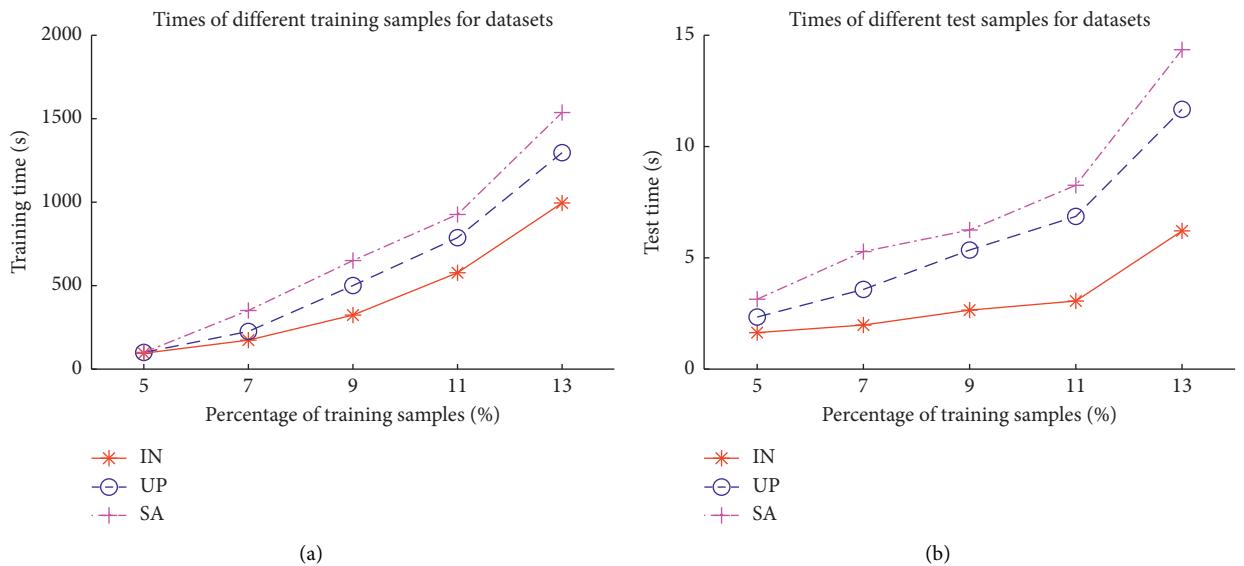


FIGURE 10: (a) Training time and (b) test time of different spatial window sizes for three datasets.

TABLE 5: Accuracy of different residual block combinations.

Combination	IN datasets			UP datasets			SA datasets		
	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)
1 + 1	99.48	99.58	99.40	99.73	99.70	99.72	99.85	99.78	99.82
1 + 2	99.46	99.60	99.41	99.42	99.50	99.37	99.57	99.45	99.53
2 + 1	99.50	99.31	99.44	99.44	99.30	99.42	99.52	99.56	99.47
2 + 2	99.72	99.60	99.70	99.86	99.83	99.81	99.98	99.98	99.98
2 + 3	99.41	99.02	99.32	99.42	99.35	99.34	99.56	99.45	99.54
3 + 2	98.35	97.73	98.07	99.74	99.70	99.73	99.86	99.83	99.85
3 + 3	99.72	99.51	99.65	99.86	99.84	99.83	99.98	99.95	99.96

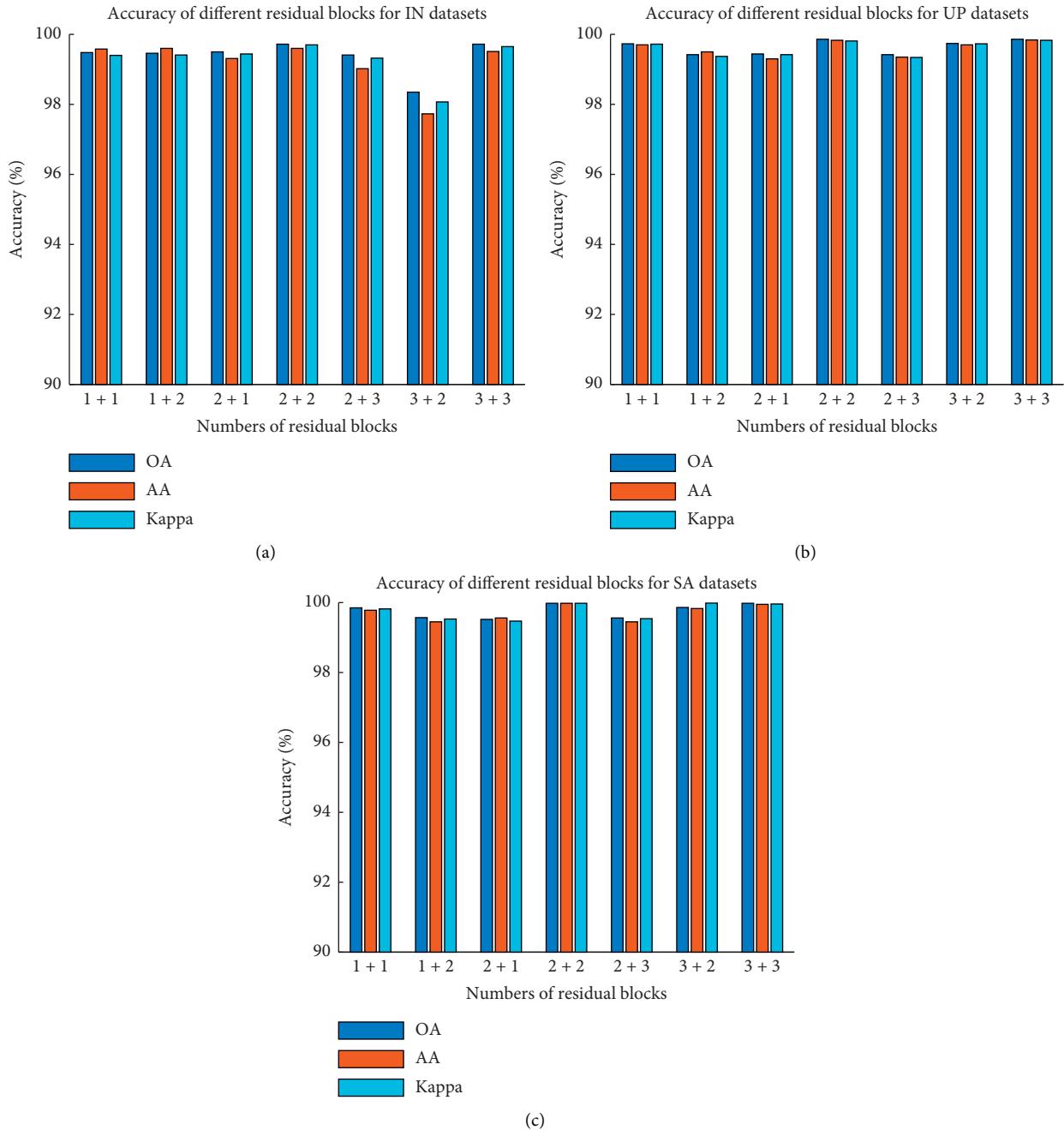


FIGURE 11: Accuracy comparison chart of different combinations: (a) IN; (b) UP; (c) SA.

TABLE 6: Comparative experiment for depth-separable convolution.

Compare items	IN datasets		UP datasets		SA datasets	
	HDSRN	Model A	HDSRN	Model A	HDSRN	Model A
Params.	714368	1233920	439904	959456	439904	959456
Training time (s)	577.1	841.8	787.3	945.7	926.3	1362.5
Testing time (s)	3.0	4.5	6.8	7.5	8.2	9.5
OA (%)	99.72	96.24	99.86	97.36	99.98	98.55

TABLE 7: Classification accuracy of different models.

Model	IN datasets			UP datasets			SA datasets		
	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)	OA (%)	AA (%)	Kappa (%)
HDSRN	99.72	99.60	99.70	99.86	99.83	99.81	99.98	99.98	99.98
Model B	99.49	99.36	93.48	99.74	99.71	99.68	99.82	99.11	99.83

TABLE 8: The comparison experiments in the IN dataset.

No. of classes	Train/test	SVM-RBF	2D CNN	3D CNN	SSRN	HybridSN	HDSRN
1	9/37	61.5	85.88	94.63	97.82	97.98	99.78
2	286/1142	78.68	91.31	93.9	99.17	98.37	99.59
3	166/664	73.41	91.07	94.85	99.53	99.48	99.82
4	47/190	71.58	80.38	93.48	97.79	97.38	99.82
5	97/386	80.38	91.89	93.56	99.24	99.23	99.38
6	146/584	92.27	99.01	94.2	99.51	99.14	99.89
7	6/22	79.52	82.59	89.73	98.7	99	99.32
8	96/382	87.38	100	96.01	99.85	100	100
9	4/16	85.87	66.55	95	98.5	99.01	100
10	194/778	77.58	86.38	94.55	98.74	98.76	99.61
11	491/1964	83.75	90.47	93.87	99.3	99.64	99.89
12	119/474	83.21	82.89	91.52	98.43	99.12	98.55
13	41/164	84.64	99.06	93.89	100	100	99.86
14	253/1012	98.01	97.86	91.77	99.31	100	100
15	77/309	94.3	90.52	95.03	99.2	99.35	99.89
16	19/74	61.43	98.94	93.57	97.82	97.26	98.57
OA (%)		82.83	90.89	94.07	99.19	99.26	99.72
AA (%)		80.84	89.68	93.72	98.93	98.98	99.62
Kappa × 100		82.23	88.56	93.87	99.07	99.09	99.70

TABLE 9: The comparison of experiments in the UP dataset.

No. of classes	Train/test	SVM-RBF	2D CNN	3D CNN	SSRN	HybridSN	HDSRN
1	663/5968	93.68	97.37	97.4	99.75	99.76	99.78
2	1865/16784	97.02	99.26	94.73	99.79	99.78	99.8
3	210/1889	82.41	80.73	95.05	98.29	99.01	99.82
4	306/2758	96.51	95.54	98.04	99.52	99.53	99.56
5	135/1211	98.38	99.75	99.01	99.82	99.87	99.9
6	503/4526	90.01	93.14	98.62	99.77	99.77	99.81
7	133/1197	85.92	91.65	97.02	99.65	99.69	100
8	368/3314	88.08	92.39	98.23	99.05	99.21	99.8
9	95/852	99.85	99.09	99.29	99.78	99.80	100
OA (%)		82.67	96.89	99.07	99.62	99.72	99.86
AA (%)		80.84	95.79	98.75	99.49	99.60	99.83
Kappa × 100		81.21	96.56	98.87	99.50	99.64	99.81

3.4.2. The Effect of Space Size. In the deep convolutional neural network, the larger the size of the input image, the larger the number of operation parameters and the higher the computational complexity. In addition, if the size of the

input image is too small, the receiving field will be too small to obtain a good classification result. In our experiment, we tried 5 different spatial input sizes, namely, 5×5 , 7×7 , 9×9 , 11×11 , and 13×13 , to evaluate the influence of the network

TABLE 10: The comparison of experiments in the SA dataset.

No. of classes	Train/test	SVM-RBF	2D CNN	3D CNN	SSRN	HybridSN	HDSRN
1	201/1808	92.56	95.51	95.19	99.78	100.00	100.00
2	373/3353	94.78	95.87	96.12	99.78	100.00	100.00
3	198/1778	94.22	95.64	95.89	99.78	100.00	100.00
4	139/1255	97.36	99.99	97.18	99.99	100.00	100.00
5	268/2410	93.38	94.51	94.70	99.78	100.00	100.00
6	396/3563	95.36	96.71	95.63	99.78	100.00	100.00
7	358/3221	95.79	95.51	95.36	99.78	100.00	100.00
8	1127/10144	80.21	86.93	86.33	99.78	100.00	100.00
9	620/5583	97.67	97.42	97.19	99.78	100.00	100.00
10	328/2950	88.99	90.93	90.94	92.69	100.00	100.00
11	107/961	90.20	94.24	90.20	90.68	98.24	99.79
12	193/1734	96.42	99.91	96.24	99.95	99.61	100.00
13	92/824	95.53	95.53	95.31	99.85	99.59	100.00
14	107/963	91.26	92.97	94.18	99.88	100.00	100.00
15	727/6541	72.67	90.93	81.39	99.12	99.32	99.84
16	181/1626	89.37	93.43	94.09	99.78	100.00	100.00
OA (%)		92.67	95.34	94.02	99.64	99.80	99.98
AA (%)		91.61	94.75	93.49	98.76	99.80	99.98
Kappa × 100		92.21	94.93	93.57	99.60	99.80	99.98

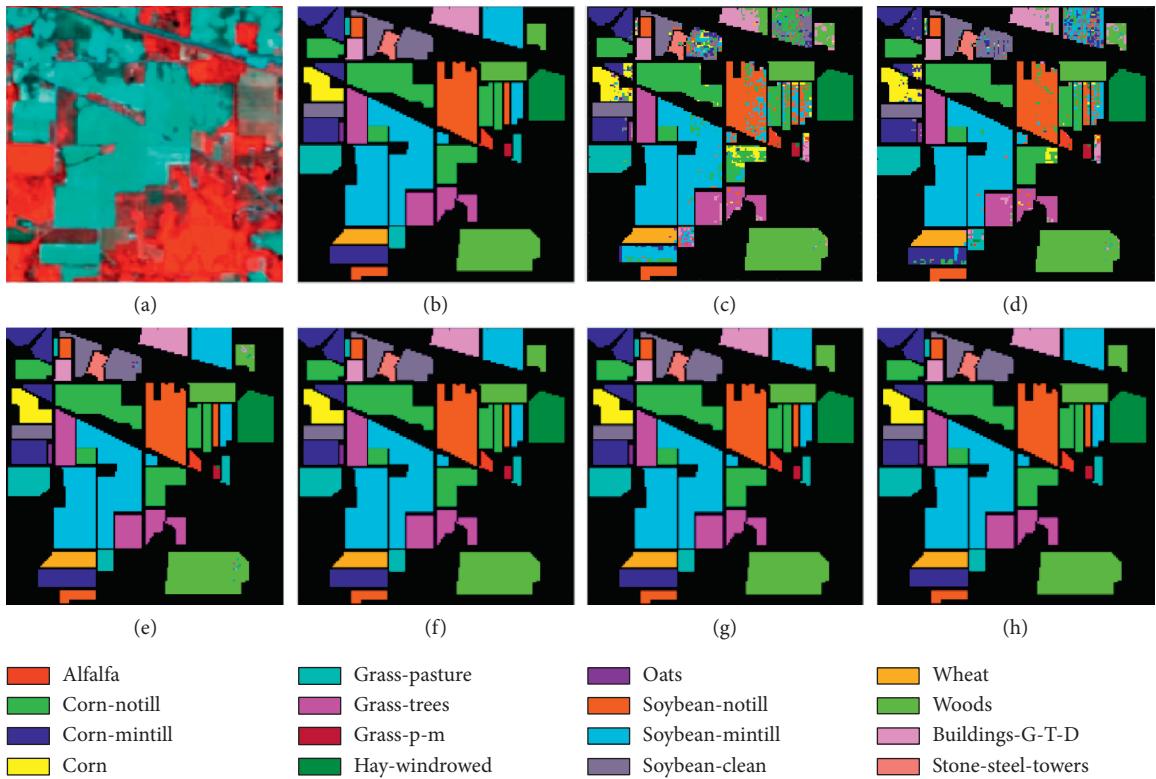


FIGURE 12: Classification results of the models in comparison with the UP dataset. (a) False color image, (b) ground truth, and (c)-(h) predicted classification maps for SVM-RBF, 2D CNN, 3D CNN, SSRN, HybridSN, and proposed HDSRN.

input size. For the IN, UP, and SA datasets, the training dataset ratios are 20%, 10%, and 10%, each experiment was repeated 5 times, and the average value of the classification indicators was used as the observation object. Table 3 and Figure 9 show the changes in OA, AA, and kappa in the IN, UP, and SA datasets under different spatial sizes. Figure 9(a) shows that as the size of the input space increases, the classification accuracy of the IN dataset begins to change

slowly after 11×11 . When the space input size is 13×13 , there is a slight downward trend. Figure 9(b) and Figure 9(c) show that as the size of the input space gradually increases, the classification accuracy of UP and SA datasets is significantly improved, and when the space input size reaches 11×11 , it starts to change slowly.

Table 4 and Figure 10 show the changes in training time and test time under different spatial sizes. We can see from

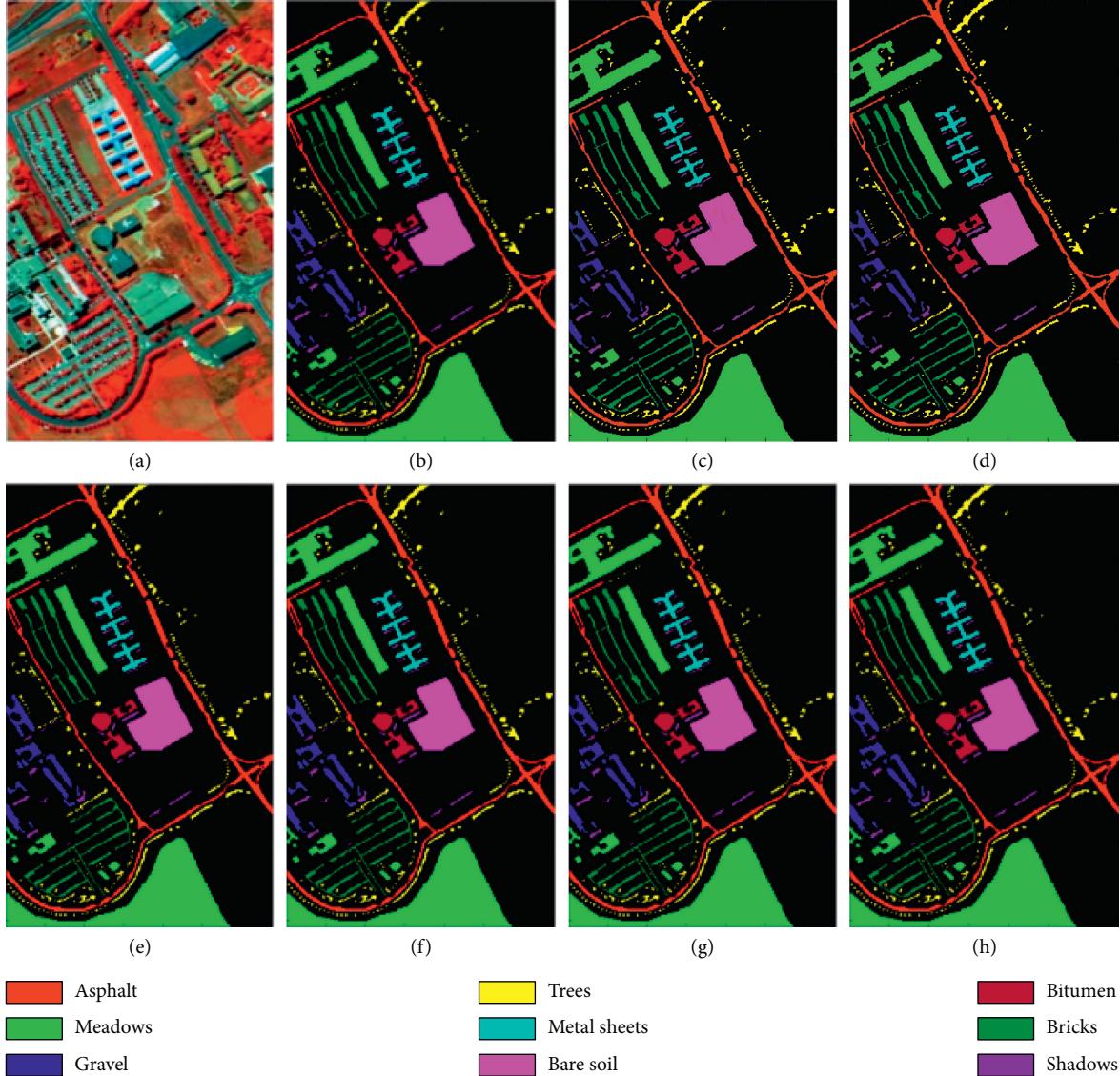


FIGURE 13: Classification maps for the IN dataset. (a) False color image, (b) ground truth, and (c)-(h) predicted classification maps for SVM-RBF, 2D CNN, 3D CNN, SSRN, HybridSN, and proposed HDSRN.

Figure 10 that as the size of the input space increases, the number of calculation parameters also gradually increases, and the training time and test time both increase sharply, resulting in a significant jump in the calculation cost. Through the analysis of the effect of the above space size on the proposed model, we find that using a space size of 11×11 is the most suitable for our model. Therefore, the input space size we chose in the comparison experiment is 11×11 . At this time, the accuracy of the three datasets is 99.72% (IN), 99.86% (UP), and 99.98% (SA).

3.4.3. The Effect of Residual Networks. We conducted a comparative experiment with different combinations of residual blocks. We tested the combination of 3D residual blocks and 2D residual blocks of $1+1$, $1+2$, $2+1$, $2+2$, $2+3$, $3+2$, and $3+3$ and gave each group's accuracy, as shown in Table 5. Figure 11 is a comparison bar chart of the

classification results. We can find that the classification effect of the three datasets is the best, when the combination is 2 + 2. In addition, we can also find from Figure 11 that the accuracy has not continued to improve, but has decreased, as the number of residual blocks increases. The first reason may be limited training samples, and the second reason may be that a deeper network increases the complexity of feature extraction.

3.4.4. The Effect of Depth-Separable Convolution. We conducted a comparative experiment to test the impact of depth-separable convolution. The traditional 2D convolution is replaced by the depth-separable convolution to form a comparison model A. The other settings were consistent with HDSRN. The training sample ratios of the IN, UP, and SA datasets are 20%, 10%, and 10%. The same spatial dimension is extracted for different datasets. For example, the

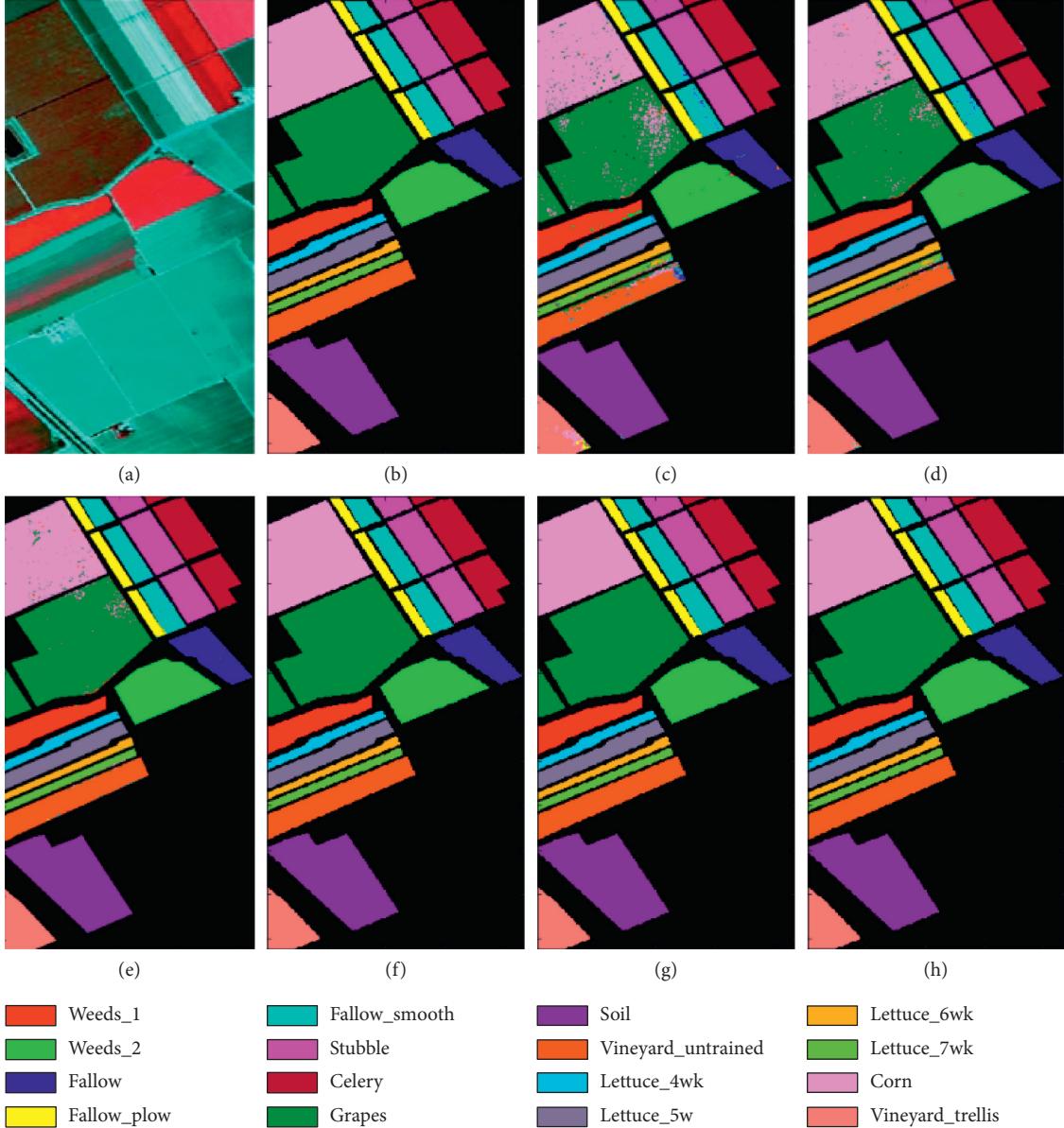


FIGURE 14: Classification maps for the SA dataset. (a) False color image, (b) ground truth, and (c)-(h) classification maps for SVM-RBF, 2D CNN, 3D CNN, SSRN, HybridSN, and proposed HDSRN.

spatial dimension of IN is $11 \times 11 \times 30$, the spatial dimension of UP is $11 \times 11 \times 15$, and the spatial dimension of SA is $11 \times 11 \times 15$. For three datasets, we conducted 5 experiments. Table 6 shows the comparison between traditional 2D convolution and depth-separable convolution in parameters, sample training time, testing time, and overall classification accuracy. We can find that the depth-separable convolution reduces the number of parameters and operation time, avoids overfitting, and improves the performance of HSI classification.

3.4.5. The Effect of Skip Connections. The HDSRN framework included a skip layer. In order to test the effect of skip connection, we designed a framework without skip connection as model B and conducted related comparative

experiments on three datasets. We can see from Table 7 that the outermost skip connection can improve the classification accuracy of HSI because skip connection alleviates the problem of gradient disappearance, helps gradients propagate backward, and can extract features more fully.

3.5. The Comparative Experiment with Popular Methods. So as to evaluate the HSI classification capability of HDSRN, we compared the model with the popular methods, such as SVM-RBF [37], 2D CNN [38], 3D CNN [24], SSRN [22], and HybridSN [31]. We used some public codes to train and test the data, which can be accessed online at <https://github.com/eecn/Hyperspectral-Classification> and <https://github.com/gokriznastic/HybridSN>. For the fairness of the experiment, through the comparative experiment in Section 3.4,

we set the input space to the same size. The input sizes of IN, UP, and SA datasets are $11 \times 11 \times 30$, $11 \times 11 \times 15$, and $11 \times 11 \times 15$, respectively. The training data ratios of IN, UP, and SA datasets are 20%, 10%, and 10%, respectively. We conducted 5 repeated trials with the mean classification metrics as final results.

Tables 8–10 show the comparison of experiments using different methods. We can see that the HDSRN method performs well, and the classification effect on the three datasets is better than other methods. One possible reason is that the proposed network model is based on the spatial spectrum 3D CNN and 2D CNN hierarchical framework, and they are complementary. This design method helps to capture more contextual information. Another possible reason is that the multiresidual network ROR is embedded in the proposed model, which can extract deeper abstract features. The last reason is that the depth-separable convolution reduces the computational cost and avoids overfitting.

From the analysis of Tables 8–10, we can see that in the IN dataset, the classification accuracy of HDSRN is about 0.5% higher than SSRN and about 0.4% higher than HybridSN. In the UP and SA datasets, the classification accuracy of HDSRN is also significantly improved. We can also find that the accuracy of network classification that only focuses on spectral or spatial features is usually less than 97%. Therefore, the classification method based on the combination of spectrum and spatial features is significantly better than the traditional method. In the IN dataset, HDSRN obtained significantly better classification results than SSRN and HybridSN, in the classification of categories 2, 5, 9, and 16. In addition, because there are fewer training samples for classes 1, 7, 9, and 16, the classification accuracy of these types of features is unstable and significantly lower than other categories, which has a greater impact on the overall classification accuracy. However, HDSRN can still classify these categories, and the classification accuracy is higher than 97%.

Figures 12–14 are the visualization of the classification by the HDSRN model and the comparison network. These maps include false color images, ground truth, and visualization images of different comparison methods. We can see that SVM-RBF has the worst visual effect among these models. The generated visual image is relatively rough, the classification accuracy is low, and the noise is obvious. This may be due to the fact that traditional methods cannot effectively extract spatial-spectral features, resulting in unsatisfactory classification results. The second is 2D CNN, and the classification effect is relatively poor. This may be because 2D CNN is unable to extract a good identification feature map from the spectral dimension, and a lot of spectral information is lost. The visual images of 3D CNN are relatively smooth. The visual images of SSRN, HybridSN, and HDSRN are smoother. However, SSRN has misclassified noise in categories 7 and 9 of the UP dataset. Compared with other methods, HDSRN has a higher classification accuracy. In addition, the classification map provided by HDSRN is the most accurate, and the edge contours of features are clearer than other methods.

4. Conclusions

In this paper, we propose a hybrid separable convolutional residual model for HSI classification. The model integrates spatial-spectral residual blocks, spatially separable convolutional residual blocks, and the outermost skip connections. The spatial-spectral 3D feature and the spatial 2D feature can be continuously extracted. Concretely, the HDSRN model can extract spectral feature information and spatial feature information, and these feature information is complementary. Additionally, we embed a multilevel residual network ROR in the model to learn spectral and spatial features, which improves the network optimization and learning capabilities. Finally, in spatial feature learning, we use separable convolution to extract valuable features, reducing the number of parameters and computing time and alleviating the decline in accuracy. The comparative experimental results on the IN, UP, and SA benchmark datasets verify the superiority of the HDSRN method. The focus of future work will be to use transfer learning methods to further solve the HSI classification problem.

Data Availability

The data used to support the findings of this study are included within the article.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This research was funded by the Foundation of Science and Technology on Electro-Optical Information Security Control Laboratory, grant no. 614210701041705.

References

- [1] N. Ma, Y. Peng, S. Wang, and P. Leong, “An unsupervised deep hyperspectral anomaly detector,” *Sensors*, vol. 18, no. 3, p. 693, 2018.
- [2] G. Notesco, Y. Ogen, and E. Ben-Dor, “Mineral classification of makhtesh ramon in Israel using hyperspectral longwave infrared (LWIR) remote-sensing data,” *Remote Sensing*, vol. 7, no. 9, pp. 12282–12296, 2015.
- [3] T.-H. Hsieh and J.-F. Kiang, “Comparison of CNN algorithms on hyperspectral image classification in agricultural lands,” *Sensors*, vol. 20, no. 6, p. 1734, 2020.
- [4] J. M. Bioucas-Dias, A. Plaza, G. Camps-Valls, P. Scheunders, N. Nasrabadi, and J. Chanussot, “Hyperspectral remote sensing data analysis and future challenges,” *IEEE Geoscience and Remote Sensing Magazine*, vol. 1, no. 2, pp. 6–36, 2013.
- [5] D. Prasad and K. Agarwal, “Classification of hyperspectral or trichromatic measurements of ocean color data into spectral classes,” *Sensors*, vol. 16, no. 3, p. 413, 2016.
- [6] X. Li, W. Chen, Q. Zhang, and L. Wu, “Building Auto-Encoder Intrusion Detection System based on random forest feature selection,” *Computers & Security*, vol. 95, p. 101851, 2020.
- [7] Y. Chen, Z. Lin, X. Zhao, G. Wang, and Y. Gu, “Deep learning-based classification of hyperspectral data,” *IEEE*

- Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 6, pp. 2094–2107, 2014.
- [8] C. Zhao, H. Zhao, G. Wang, and H. Chen, “Improvement SVM classification performance of hyperspectral image using chaotic sequences in artificial bee colony,” *IEEE Access*, vol. 8, pp. 73947–73956, 2020.
 - [9] P. Hu, X. Liu, Y. Cai, and Z. Cai, “Band selection of hyperspectral images using multiobjective optimization-based sparse self-representation,” *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 3, pp. 452–456, 2019.
 - [10] C. Hernández-Espinosa, M. Fernández-Redondo, and J. Torres-Sospedra, “Some experiments with ensembles of neural networks for classification of hyperspectral images,” *Advances in Neural Networks-ISNN 2004*, pp. 912–917, 2004.
 - [11] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
 - [12] D. Heaven, “Why deep-learning AIs are so easy to fool,” *Nature*, vol. 574, no. 7777, pp. 163–166, 2019.
 - [13] M. H. Hesamian, S. Mashohor, M. I. Saripan, and W. A. Wan Adnan, “Effect of image resolution on intensity based scene illumination classification using neural network,” *The Imaging Science Journal*, vol. 63, no. 8, pp. 433–439, 2015.
 - [14] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Las Vegas, NV, USA, June 2016.
 - [15] Y. Wei, X. Luo, L. Hu, Y. Peng, and J. Feng, “An improved unsupervised representation learning generative adversarial network for remote sensing image scene classification,” *Remote Sensing Letters*, vol. 11, no. 6, pp. 598–607, 2020.
 - [16] J. Gu, Z. Wang, J. Kuen et al., “Recent advances in convolutional neural networks,” *Pattern Recognition*, vol. 77, pp. 354–377, 2018.
 - [17] P. Ghamisi, N. Yokoya, J. Li et al., “Advances in hyperspectral image and signal processing: a comprehensive overview of the state of the art,” *IEEE Geoscience and Remote Sensing Magazine*, vol. 5, no. 4, pp. 37–78, 2017.
 - [18] L. Dong, F. Wei, K. Xu, S. Liu, and M. Zhou, “Adaptive multi-compositionality for recursive neural network models,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 3, pp. 422–431, 2016.
 - [19] S. A. Amirshahi, M. Pedersen, and S. X. Yu, “Image quality assessment by comparing CNN features between images,” *Journal of Imaging Science and Technology*, vol. 60, p. 60410, 2016.
 - [20] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” <http://arxiv.org/abs/1409.1556>.
 - [21] C. Szegedy, W. Liu, Y. Jia et al., “Going deeper with convolutions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Boston, MA, USA, June 2015.
 - [22] Z. Zhong, J. Li, Z. Luo, and M. Chapman, “Spectral-spatial residual network for hyperspectral image classification: a 3-D deep learning framework,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 2, pp. 847–858, 2018.
 - [23] Y. Chen, H. Jiang, C. Li, X. Jia, and P. Ghamisi, “Deep feature extraction and classification of hyperspectral images based on convolutional neural networks,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 10, pp. 6232–6251, 2016.
 - [24] Y. Li, H. Zhang, and Q. Shen, “Spectral-spatial classification of hyperspectral imagery with 3D convolutional neural network,” *Remote Sensing*, vol. 9, no. 1, p. 67, 2017.
 - [25] X. Yang, Y. Ye, X. Li, R. Y. K. Lau, X. Zhang, and X. Huang, “Hyperspectral image classification with deep learning models,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 9, pp. 5408–5423, 2018.
 - [26] W. Song, S. Li, L. Fang, and T. Lu, “Hyperspectral image classification with deep feature fusion network,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 6, pp. 3173–3184, 2018.
 - [27] L. Fang, Z. Liu, and W. Song, “Deep hashing neural networks for hyperspectral image feature extraction,” *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 9, pp. 1412–1416, 2019.
 - [28] Z. Gong, P. Zhong, W. Hu, Z. Xiao, and X. Yin, “A novel statistical metric learning for hyperspectral image classification,” <http://arxiv.org/abs/1905.05087>.
 - [29] B. Liu, X. Yu, P. Zhang, and X. Tan, “Deep 3D convolutional network combined with spatial-spectral features for hyperspectral image classification,” *Acta Geodaetica et Cartographica Sinica*, vol. 48, pp. 53–63, 2019.
 - [30] H. Lee and H. Kwon, “Going deeper with contextual CNN for hyperspectral image classification,” *IEEE Transactions on Image Processing*, vol. 26, no. 10, pp. 4843–4855, 2017.
 - [31] S. K. Roy, G. Krishna, S. R. Dubey, and B. B. Chaudhuri, “HybridSN: exploring 3-D-2-D CNN feature hierarchy for hyperspectral image classification,” *IEEE Geoscience and Remote Sensing Letters*, vol. 17, no. 2, pp. 277–281, 2020.
 - [32] F. Cao and W. Guo, “Deep hybrid dilated residual networks for hyperspectral image classification,” *Neurocomputing*, vol. 384, pp. 170–181, 2020.
 - [33] P. Wu, Z. Cui, Z. Gan, and F. Liu, “Three-dimensional ResNeXt network using feature fusion and label smoothing for hyperspectral image classification,” *Sensors*, vol. 20, no. 6, p. 1652, 2020.
 - [34] G. Li, H. Tang, Y. Sun et al., “Hand gesture recognition based on convolution neural network,” *Cluster Computing*, vol. 22, no. S2, pp. 2719–2729, 2019.
 - [35] W. Cheng, Y. Sun, G. Li, G. Jiang, and H. Liu, “Jointly network: a network based on CNN and RBM for gesture recognition,” *Neural Computing and Applications*, vol. 31, no. S1, pp. 309–323, 2019.
 - [36] F. Chollet, “Deep learning with depthwise separable convolutions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, HI, USA, July 2017.
 - [37] K. Zhang, M. Sun, T. X. Han, X. Yuan, L. Guo, and T. Liu, “Residual networks of residual networks: multilevel residual networks,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, pp. 1303–1314, 2017.
 - [38] S. Ji, W. Xu, M. Yang, and K. Yu, “3D convolutional neural networks for human action recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 221–231, 2013.
 - [39] B. Kuo, H. Ho, C. Li, C. Hung, and J. Taur, “A kernel-based feature selection method for SVM with RBF kernel for hyperspectral image classification,” *Applied Earth Observations and Remote Sensing*, vol. 7, pp. 317–326, 2014.
 - [40] Q. Wang, J. Gao, and Y. Yuan, “A joint convolutional neural networks and context transfer for street scenes labeling,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 5, pp. 1457–1470, 2018.
 - [41] F. Luo, B. Du, L. Zhang, L. Zhang, and D. Tao, “Feature learning using spatial-spectral hypergraph discriminant analysis for hyperspectral image,” *IEEE Transactions on Cybernetics*, vol. 49, pp. 2406–2419, 2018.
 - [42] Computational Intelligence Group of the Basque University (UPV/EHU), *Hyperspectral Remote Sensing Scenes*, Computational Intelligence Group of the Basque University (UPV/EHU), San Sebastian, Spain, 2020, http://www.ehu.eus/ccwintco/index.php/Hyperspectral_Remote_Sensing_Scenes.