A Spectral–Spatial Fusion Transformer Network for Hyperspectral Image Classification

Diling Liao¹⁰, Cuiping Shi¹⁰, Member, IEEE, and Liguo Wang¹⁰, Member, IEEE

Abstract-In the past, deep learning (DL) technologies have been widely used in hyperspectral image (HSI) classification tasks. Among them, convolutional neural networks (CNNs) use fixed-size receptive field (RF) to obtain spectral and spatial features of HSIs, showing great feature extraction capabilities, which are one of the most popular DL frameworks. However, the convolution using local extraction and global parameter sharing mechanism pays more attention to spatial content information, which changes the spectral sequence information in the learned features. In addition, CNN is difficult to describe the long-distance correlation between HSI pixels and bands. To solve these problems, a spectral-spatial fusion Transformer network (S2FTNet) is proposed for the classification of HSIs. Specifically, S2FTNet adopts the Transformer framework to build a spatial Transformer module (SpaFormer) and a spectral Transformer module (SpeFormer) to capture image spatial and spectral longdistance dependencies. In addition, an adaptive spectral-spatial fusion mechanism (AS²FM) is proposed to effectively fuse the obtained advanced high-level semantic features. Finally, a large number of experiments were carried out on four datasets, Indian Pines, Pavia, Salinas, and WHU-Hi-LongKou, which verified that the proposed S2FTNet can provide better classification performance than other the state-of-the-art networks.

Index Terms— Deep learning (DL), fusion, hyperspectral image (HSI), long-distance dependence.

I. INTRODUCTION

HYPERSPECTRAL images (HSIs) are captured by airborne imaging spectrometer and carry a lot of spectral and spatial information. In recent years, HSIs have played an important role in many fields, including health care [1], military [2], Earth exploration [3], and environmental protection [4]. Among them, HSI classification is an important stage of HSI processing and is one of the hot spots of image research. Specifically, HSI classification is to classify images pixel by pixel by learning prior knowledge [5], [6], [7].

In the early stage of research, classification methods paid more attention to the spectral feature extraction of

Diling Liao and Cuiping Shi are with the Department of Communication Engineering, Qiqihar University, Qiqihar 161000, China (e-mail: 2020910228@qqhru.edu.cn; shicuiping@qqhru.edu.cn).

Liguo Wang is with the College of Information and Communication Engineering, Dalian Nationalities University, Dalian 116000, China (e-mail: wangliguo@hrbeu.edu.cn).

Digital Object Identifier 10.1109/TGRS.2023.3286950

images, and many classical methods appeared, such as support vector machines (SVMs) [8], principal component analysis (PCA) [9], and composite kernels [10]. Although the above traditional methods can obtain the basic features of the image, the classification performance is not satisfactory. In addition, these methods have many disadvantages, for example, too much dependence on knowledge in professional fields, low generalization ability, and weak representation ability of acquired features. Therefore, deep learning (DL) technology is becoming more and more popular in computer vision tasks (such as classification [11], [12], [13], detection [14], [15], and segmentation [16]) because it can not only get rid of the constraints of manual but also adaptively learn high-level semantic information.

In recent years, many excellent frameworks have emerged for DL technology, including convolutional neural networks (CNNs) [17], generative adversarial networks (GANs) [18], [19], recurrent neural networks (RNNs) [20], [21], graph convolutional networks (GCNs) [22], [23], capsule network (CapsNet) [24], and vision Transformer (ViT) [25].

Among them, CNNs are one of the most popular DL methods, which improve the discriminative ability of features through local connection and global parameter sharing mechanism. Unlike other ordinary images, HSI contains rich spectral and spatial features, and the construction of CNN network can easily extract these two features of HSI. Hu et al. [26] used 1-D CNN to classify HSI pixel by pixel and verified that 1-D CNN is suitable for HSI classification tasks. In addition, the image has rich spatial information. In order to integrate the spatial information of the image, Zhao and Du [27] proposed 2-D CNN, which uses the adjacent pixels around the central classification pixel as training samples to perform classification tasks, improving the classification performance. However, only using 2-D CNN is not enough to extract spectral-spatial (SS) joint features of images. Therefore, Hamida et al. [28] cut the HSI into multiple 3-D cubes and constructed the 3-D CNN to extract the SS joint features of the image, verifying that the method can effectively improve the classification performance. Similarly, Roy et al. [29] designed an SS hybrid network based on 3-D CNN and 2-D CNN and proved its effectiveness. Shang et al. [30] proposed a classification method based on multiscale cross-branch response and second-order channel attention (MCRSCA), which considers the inherent spatial structure information of ground objects and avoids the loss of spatial details. With the gradual increase in 3-D CNN network depth, gradient disappearance and gradient explosion will occur [31], and the classifi-

1558-0644 © 2023 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See https://www.ieee.org/publications/rights/index.html for more information.

Manuscript received 26 April 2023; revised 2 June 2023; accepted 13 June 2023. Date of publication 26 June 2023; date of current version 30 June 2023. This work was supported in part by the National Natural Science Foundation of China under Grant 42271409 and Grant 62071084, in part by the Heilongjiang Science Foundation Project of China under Grant LH2021D022, and in part by the Fundamental Research Funds in Heilongjiang Provincial Universities of China under Grant 145209149. (*Corresponding author: Cuiping Shi.*)

cation accuracy will gradually decrease. In order to solve this problem, Zhong et al. [32] introduced the ResNet [11] structure into the designed spatial 3-D CNN module and spectral 3-D CNN module and extracted rich spatial and spectral features. In addition, Paoletti et al. [33] proposed a depth pyramid residual network for SS HSI classification by making better use of the potential of available information on each unit. In order to further improve the classification performance and alleviate the problem of overfitting, attention mechanism has been widely concerned and successfully applied in the HSI classification [34], [35], [36], [37]. For example, He et al. [38] proposed a dual global-local attention network (DGLANet). In order to reduce the spatial and spectral redundancy information of pixels, Mei et al. [39] proposed a network based on bidirectional long short-term memory (Bi-LSTM), which designed an SS attention mechanism and emphasized effective information. In addition, lightweight classification methods based on CNN are also popular. For example, a lightweight network [40] is constructed by using 3-D depthwise convolution, which reduces model parameters and computational overhead. Meng et al. [41] proposed a lightweight SS convolution module (LS²CM) as an alternative to the convolution layer. Kang et al. [42] proposed an SS classification framework based on edge preserving filtering (EPF). Zhong et al. [43] designed an iterative EPF (IEPF) method based on EPF and further improved the classification performance. In addition, they embedded an iterative strategy into SS classifiers and designed a new HSI classification method that combines multiple SS classifiers [44].

In the past, Transformer has received extensive attention in the field of natural language processing (NLP). It is worth noting that Transformer has recently been introduced into computer vision and successfully applied to image classification tasks [25]. Since the spectrum of HSIs is sequence data and usually contains hundreds of wavebands, He et al. [45] proposed a spatial-spectral Transformer (SST) network by combining transfer learning with the Transformer framework and proved that the Transformer can construct the correlation of spectral sequences. Similarly, Hong et al. [46] reconsidered Transformer from the perspective of spectral sequence attributes, proposed a spectral Transformer (SpeFormer) network, and confirmed that it has more significant advantages than classical ViT and advanced backbone networks. In general, CNN-based network access to high-level semantic features is relatively limited. Therefore, Sun et al. [47] proposed an SS feature tokenization Transformer (SSFTT) network to capture SS features and advanced semantic features. Similarly, Zhong et al. [48] proposed a new SS Transformer network (SSTN) to overcome the weak ability of CNNs to learn long-distance dependencies. Huang et al. [49] proposed a new 3-D swin Transformer-based hierarchical contrastive learning (3DSwinT-HCL) method based on 3-D swin Transformer. This method uses Transformer to effectively make up the shortcomings of CNNs lack of receptive field (RF) and inability to capture the order attribute of data. In order to solve the problem that the network is easily interfered by irrelevant information around the target pixel in the training phase, which

leads to inaccurate feature extraction, Bai et al. [50] proposed an HSI classification method based on the multibranch attention Transformer network. Zou et al. [51] proposed the local-enhanced SS Transformer (LESSForm) method, which alleviates the problem that Transformer-based classification methods usually generate inaccurate tag embedding from a single spectral or spatial dimension of the original HSI. Inspired by the bottleneck Transformer of computer vision, Song et al. [52] proposed a bottleneck spatial–spectral Transformer (BS2T) network, which uses Transformer to make the extracted features more spatial location aware and spectral aware. Mei et al. [53] proposed a group-aware hierarchical Transformer (GAHT) to solve the problem of overdispersion of features extracted by multihead self-attention (MHSA) in the Transformer.

Although the above DL methods have been widely used in the HSI classification, there are still some challenges. On the one hand, CNNs using the mechanism of local extraction and global parameter sharing pay more attention to spatial content information, thus distorting the spectrum sequence information in the learning features [51]. On the other hand, CNNs are difficult to describe the long-distance correlation between HSI pixels and bands. On the contrary, Transformer can not only effectively extract long-distance dependence but well maintain spectral sequence information. Therefore, this article proposed an HSI classification method based on spectral-spatial fusion Transformer network (S²FTNet). In particular, S²FTNet uses the Transformer framework to build the spatial Transformer (SpaFormer) module and SpeFormer module to capture the long-distance dependencies in image spatial and spectral. In addition, an adaptive spectral-spatial fusion mechanism (AS²FM) is proposed to effectively combine the obtained SS high-level semantic features.

The main contributions of this article are given as follows.

- In order to enhance the long-distance dependency of features and improve the representation ability of features, a Transformer block based on multihead double selfattention (MHD-SA) is proposed. Then, three improved Transformer blocks are constructed in parallel as a SpaFormer module to extract the long-distance dependence of images with different spatial dimensions.
- 2) In order to increase the RF of spectral extraction and learn more spectral sequence information, a SpeFormer module is designed. It uses convolution to replace the traditional Transformer's multilayer perceptron (MLP) and combines it with the proposed MHD-SA.
- Considering the different importance of high-level semantic features extracted by spatial branches and spectral branches, in order to combine them more effectively, an AS²FM is proposed.
- 4) Based on Transformer and CNN, we proposed an S²FTNet, which uses a dual-branch structure to extract spectral and spatial features and combines the features obtained from the two branches with an adaptive fusion mechanism. Extensive experiments have proved that our method has a better performance and potential compared with some state-of-the-art CNN-based and Transformer networks.

The rest of this article is arranged as follows. In Section II, the network structure of S^2FTNet proposed in this article is introduced in detail. In Section III, the parameter analysis of the model, quantitative analysis of comparative experiments, and visual evaluation are provided. Section IV gives the conclusion and prospect of this article.

II. METHODOLOGY

The method S²FTNet proposed in this article includes three main modules: SpaFormer, SpeFormer, and AS²FM. The overall network framework is shown in Fig. 1. Suppose that input HSI data are $X \in \mathbb{R}^{H \times W \times L}$, where W and H represent the width and height of the image, respectively, L represents the number of bands of the image, and the corresponding label set $Y_i \in \{1, 2, \dots, Class\}$. In order to facilitate feature extraction, X is first processed by an edgefilling strategy. Then, the new data obtained after filling are extracted in two ways. One is to extract the adjacent edge blocks of the pixel with the pixel to be classified as the center and reduce the spectral dimension by PCA to obtain data X_patch $\in \mathbb{R}^{s \times s \times b}$. The other is pixel-by-pixel extraction (PPE) to obtain data X_pixel $\in \mathbb{R}^{1 \times 1 \times L}$, where $s \times s$ represents the image space size after segmentation and b represents the number of spectral bands after PCA dimensionality reduction. Next, the two processed data are used as the input data of SpaFormer and SpeFormer modules, and the advanced semantic features extracted by the two modules are fused through an adaptive fuse mechanism. Finally, the fused feature vectors are transferred to the classifier for classification.

Then, the three main modules of S²FTNet proposed in this article are introduced in detail.

A. SpaFormer Module

In recent years, CNNs are one of the most classical DL frameworks and are widely used in HSI classification tasks. Convolution (Conv) of CNN uses a mechanism of local connection and global parameter sharing so that more attention is paid to the local features of the image during the feature extraction process. In contrast to Conv, Transformer can build long-distance dependencies, making up the shortcomings of Conv in feature extraction. Therefore, the SpaFormer uses the above two frameworks for modeling, and the structure is shown in Fig. 1. Next, this section will introduce the proposed SpaFormer module in detail.

First, the input image data $X_{\rm patch}$ passes through two Conv blocks, namely, 3-D convolution (Conv3D) and 2-D convolution (Conv2D), and each Conv block contains the convolution layer, the batch normalization (BN) layer, and the nonlinear activation layer. Specifically, $X_{\rm patch}$ extracts the SS joint information of the image through Conv3D, and the calculation process is

$$F_{3-D} = f\left(\delta_1\left(X_{\text{patch}}\Theta w^{3-D} + b^{3-D}\right)\right). \tag{1}$$

In Formula (1), w^{3-D} represents the weight offset of 3-D Conv, b^{3-D} represents the offset term, and F_{3-D} represents the output of Conv3D. Θ is a 3-D Conv operator, δ_1 is a 3-D BN operation, and $f(\cdot)$ is a nonlinear activation function ReLU. In order to further extract image spatial information, the module introduces Conv2D after Conv3D. The calculation principle of Conv2D is similar to that of Conv3D, and the formula is

$$F_{2-D} = f(\delta_2(F_{3-D} \odot w^{2-D} + b^{2-D})).$$
(2)

In Formula (2), $w^{2\text{-D}}$ represents the weight offset of 2-D Conv, $F_{2\text{-D}}$ represents the offset term, and $b^{2\text{-D}}$ represents the output of Conv2D. \odot is a 2-D Conv operator, and δ_2 is a 2-D BN operation. The module first extracts the SS joint and spatial features of the image by designing Conv3D and Conv2D, which provided complete shallow information for extracting high-level semantic features.

Then, three improved Transformer blocks are used for parallel connection to build the SpaFormer module, which is used to explore the long-distance dependency of images. As can be seen from Fig. 1, each Transformer block contains multiple components, including position embedding (PE), two layers of normalization (Norm), MHD-SA, and MLP.

To strengthen the correlation between positions, the Transformer block first introduced PE. To put it simply, all tokens $T = [T_1, T_2, ..., T_w]$ are connected to the learnable classification token T_0 , and the location information PE_{pos} is attached to all tokens, i.e.,

$$T_{\rm PE} = [T_0, T_1, T_2, \dots, T_w] + {\rm PE}_{\rm pos}.$$
 (3)

The proposed MHD-SA is the most important component of the entire Transformer, and its structure is shown in Fig. 2(a). At the same time, for the convenience of illustration, the single-head structure of MHD-SA is shown in Fig. 2(b). MHD-SA usually contains three feature inputs, namely, query (Q), key (K), and value (V), and Q, K, and V are obtained by linear mapping of three predefined weight matrices W_Q , W_K , and W_V . The self-attention score of single-headed double selfattention (DSA) is calculated by Q and K, and then, the score is weighted into V, i.e.,

$$SA = \operatorname{soft} \max\left(\frac{QK^{T}}{\sqrt{d_{K}}}\right)V \tag{4}$$

$$DSA = \text{soft} \max\left(\frac{L_{Q}(SA)L_{K}(SA)}{\sqrt{d_{L_{K}}}}\right)L_{V}(SA).$$
(5)

In (4) and (5), SA represents the self-attention value; $L_Q(\cdot)$, $L_K(\cdot)$, and $L_V(\cdot)$ represent the features obtained by SA through linear mapping; and d_K and d_{L_K} represent the feature dimensions of K and L_K , respectively. Generally, Transformer contains multiple-head self-attention, so the MHD-SA can be represented as

$$MHD-SA = Concat(DSA_1, DSA_2, \dots, DSA_h)W$$
(6)

where $Concat(\cdot)$ represents the cascade function, *h* represents the number of headers, and *W* represents the weight parameter.

Finally, MLP is introduced after MHD-SA to alleviate the problem of gradient explosion and gradient disappearance. The MLP structure contains two full connection layers, and a Gaussian error linear unit (GELU) is embedded between the two full connection layers.



Fig. 1. S²FTNet overall network framework.



Fig. 2. Overall structure of MHD-SA. (a) Multihead structure. (b) Single-head structure.

It is worth noting that SpaFormer contains three improved Transformer blocks. Although the three Transformer blocks have the same structure, the input data are different. It can be seen from Fig. 1 that the space size $s \times s$ of the input data of the three blocks performs pooling = false, pooling = 2, and pooling = 4 operations, and the output space size is $[s/pooling] \times [s/pooling]$, while [·] represents the upper rounding symbol. With different space sizes, Transformer blocks can be used to explore long-distance dependen-

cies of different spaces, which can enrich the diversity of features.

To sum up, the spatial branch contains two Conv blocks and SpaFormer modules. First, the SS joint and spatial features of the shallow layer are extracted through two Conv blocks to provide complete shallow information. Then, three improved Transformer blocks are paralleled, and different input space sizes are used to explore the long-distance dependency of features, which enriches the diversity of features.



Fig. 3. Overall structure of SpeFormer.

B. SpeFormer Module

HSI not only has rich spatial information but also contains hundreds of spectral bands. Extracting rich spectral features of images and taking full account of spectral sequence can improve the discrimination ability of features and classification performance. Therefore, inspired by [47], this article proposed a SpeFormer module. The overall structure is shown in Fig. 3.

It can be seen that the input data size is $\mathbb{R}^{1\times 1\times L}$ and *L* is the number of spectral bands of HSI. First, the input is dimensionally reduced by linear mapping and cascaded with learnable token T'_0 . Then, the results are embedded in position, and the obtained feature tensor T'_{PE} contains position and spectral order information. The calculation process is similar to the SpaFormer, i.e.,

$$T'_{\rm PE} = \left[T'_0, T'_1, T'_2, \dots, T'_w\right] + {\rm PE}'_{\rm pos}.$$
 (7)

Then, a Transformer block based on Conv is introduced, which fully considers the correlation between spectral sequences and can obtain the long-distance dependence between spectral sequences. The traditional MLP of Transformer includes two fully connected (FC) layers. Although the two layers of FC can extract spectral nonlinear features to a certain extent, it still lacks consideration of local spectral correlation. According to [54], the linear transformation at different positions in two FCs of the Transformer block is the same, but they use different parameters from one layer to another, which can be replaced by two 1×1 Conv. Therefore, in order to further explore the local spectral correlation and increase the convolution RF, SpeFormer uses two 3×3 Conv blocks (including a Conv layer and a BN layer) to replace FC in the traditional MLP block. This improved method can effectively increase the RF of spectral information extraction while avoiding the destruction of spectral order. Therefore, the improved Transformer block includes two-layer normalization, an MHD-SA, two Conv blocks, and a GELU. This process can be expressed as

SpeFormer =
$$\delta_2(f_2(g(\delta_1(f_1(\text{MHD-SA})))))).$$
 (8)

In Formula (8), $f(\cdot)$ represents the Conv function, $\delta(\cdot)$ represents the BN function, and SpeFormer represents the output result of improved Transformer block.

C. Adaptive Spectral–Spatial Fusion Mechanism

In this article, the proposed S^2 FTNet selects cross entropy as the loss function and optimizes the network through backpropagation where the expression of cross-entropy loss function is

Loss =
$$\frac{1}{C} \sum_{a=1}^{r} \left[-y'_a \log(y_a) - (1 - y'_a) \log(1 - y_a) \right].$$
 (9)

In Formula (9), y'_a and y_a represent real object labels and model prediction labels, respectively, *C* represents the total number of categories in the dataset, and Loss represents the average loss value of each mini-batch.

 S^2 FTNet includes two branches, SpaFormer branch and SpeFormer branch. Then, the high-level semantic features obtained from these two branches will be combined and sent to the classifier. In this section, we will introduce in detail how to effectively combine the features extracted from these two branches. Usually, two features are cascaded as follows:

$$F = \text{Concat}(F_{\text{Spa}}, F_{\text{Spe}}). \tag{10}$$

However, considering that the two important degrees of the features extracted from the two branches are different, we introduce the balance factor λ for score weighting, i.e.,

$$F = \text{Concat}(\lambda F_{\text{Spa}}, (1 - \lambda) F_{\text{Spe}}).$$
(11)

In the backpropagation process, the balance factor update can be expressed as

$$\lambda = \lambda_0 - \eta \frac{\partial}{\partial \lambda} \text{Loss}$$
(12)

where λ_0 is the random initial value of the balance factor and η is the learning rate. By adaptively determining the proportion of these two parts, the model has a stronger data representation ability than feature-weighted addition.

D. Algorithm Implementation Process

In this section, we give the implementation process of the proposed network S²FTNet, as shown in Table I. Take the Pavia dataset as an example, that is, input data $X_1 \in \mathbb{R}^{13 \times 13 \times 30}$. X performs edge filling and cuts and extracts cube by pixel, respectively, to obtain processed data $X_1 \in \mathbb{R}^{13 \times 13 \times 30}$ and $X_2 \in \mathbb{R}^{1 \times 1 \times 103}$. In the SpaFormer branch, first, select X_1 as the input data, and execute Conv3D and Conv2D. Among them, Conv3D and Conv2D, respectively, select eight convolution kernels with a size of $7 \times 7 \times 7$ and 64 convolution kernels with a size of 7×7 . Then, pooling = false, pooling = 2, and pooling = 4 operations are performed on the input image data space size $s \times s$. The space of the three images is 13×13 , 7×7 , and 4×4 . Then, in order to adapt to the improved Transformer blocks, they are reshaped and used as the input of three blocks. In the SpeFormer branch, first, select X_2 as the input data to reduce the complexity, and select $\dim = 64$ to linearly map the spectral dimensions of the data. Then, the linear mapping result is executed into a position embedded and improved Transformer block. It is worth noting that the advanced semantics extracted from the two branches are adaptively weighted by introducing the balance factor λ . Finally, the Softmax function is used for classification.

5515216

TABLE I IMPLEMENTATION PROCESS OF S²FTNET

Algorithm 1 S²FTNet implementation process

- **Input:** HSI image data $X \in \mathbf{R}^{H \times W \times L}$, label is $Y \in \mathbf{R}^{H \times W}$, PCA parameter b = 30, space size s = 13.
- **Output**: The classification accuracy and visual classification map of the marked samples.
- 1:The HSI data is edge filled and neighborhood and pixel by pixel cubes are extracted respectively. The data obtained are then processed by PCA to obtain $X_1 \in \mathbf{R}^{s \times s \times b}$ and $X_2 \in \mathbf{R}^{1 \times 1 \times L}$.
- 2:Set the GSD optimizer and learning rate r = 0.005, and select the batch size b = 64 and training iterations as e = 200.
- 3: **for** e = 1 to 200 do
- 4: Select X_1 as the input data of the SpeFormer, and execute Conv3D and Conv2D.
- 5: The space size $s \times s$ is operated by pooling = false, pooling = 2 and pooling=4 respectively to get three data.
- 6: A spatial Tranformer block that performs paralleling.
- 7: Select X_2 as the input data of the SpeFormer, and execute the improved SpeFormer block.
- 8: The outputs of the two branches are weighted by adaptive scores using balance factors λ .
- 9: Use the Softmax function to identify the label.
- 10:**end**

11:Save the parameters of the optimal model, and obtain the classification accuracy of the marker samples and the visual classification map of the ground object categories.

III. EXPERIMENTAL RESULTS AND ANALYSIS

In order to verify the effectiveness of the proposed method, a series of experiments is conducted. The experiments include network ablation experiments, parameter optimization, quantitative comparison, and visualization of classification results.

A. Dataset Description

In this article, three classical datasets and a newer dataset are selected for all experiments, Indian Pines, Pavia, Salinas, and WHU-Hi-LongKou datasets. Next, in this section, we will detail the category information of each dataset and the number of training samples for the proposed method; the specific information is shown in Table II.

1) Indian Pines Dataset: It was captured by airborne imaging spectrometer airborne visible infrared imaging spectrometer (AVIRIS) from an Indian Pine tree in Indiana in 1992. Among them, there are 16 land cover categories, mainly including corn, grass, soybean, and woods. The space size of the image is 145×145 , the spatial resolution is about 20 m, and the imaging wavelength range is $0.4-2.5 \ \mu$ m. It also contains 220 continuous spectral bands. In addition to the 104–108, 150–163, and 220 absorption bands, the remaining 200 bands were used for experiments.

2) Pavia Dataset: It was captured by the airborne imaging spectrometer ROSIS-03 over the University Pavia, Italy, in 2003. The space size of the image is 610×340 , with a spatial resolution of 1.3 m and 115 continuous spectral bands. Similarly, because individual bands cannot be reflected by water, there are only 103 bands left. Compared with the Indian Pines dataset, Pavia contains fewer land cover categories, including trees, asphalt roads, bricks, and meadows.

3) Salinas Dataset: It was captured by the imaging spectrometer AVIRIS over Salinas Valley, CA, USA. The space size is a total of 111104 pixels. In addition to background pixels, pixels remain for classification tasks. These pixels contain a total of 16 categories, including fallow and celery.

4) WHU-Hi-LongKou Dataset: It is collected from Longkou Town, Hubei Province, China, by the 8-mm focus (HNH) imaging sensor carried on the DJI Matrix 600 Pro (DJI M600 Pro) unmanned aerial vehicle (UAV) platform. The space size is 550×400 , the spatial resolution is about 0.463 m, the wavelength range is $0.4-1 \mu$ m, and 270 spectral bands are included. The number of land cover categories included in the WHU-Hi-LongKou dataset is the same as that in the Pavia dataset, which is a simple crop scenario. The main categories include water, broad leaf soybean, corn, rice, and cotton.

B. Experimental Setup

All experiments in this section are implemented on the platform of Intel¹ Core² i9-9900K CPU, NVIDIA GeForce RTX 2080Ti GPU, and 128-GB random access memory, and the language framework is Python. In addition, in order to better evaluate the classification performance of the model, we choose three common evaluation indicators: overall accuracy (OA), average accuracy (AA), and Kappa coefficient. Among them, OA represents the ratio of the number of accurately classified samples to the total number of samples, AA represents the average of the classification accuracy of each category, and Kappa is a measure of robustness.

The network constructed by combining CNN and Transformer is more inclined to spatial information of global context. In order to analyze the impact of different input space sizes *s* on the final classification performance, we selected 7– 15 input space sizes for experiments on four datasets. The adjacent space size interval is 2. The experimental results are shown in Fig. 4. It can be seen from Fig. 4 that the Indian Pines dataset is highly sensitive to different input space sizes. The classification accuracy OA of Pavia and WHU-Hi-LongKou datasets shows a trend of increasing first and then decreasing. For the Salinas dataset, with the increase in input space size *s*, OA increases first and then tends to be stable. It is worth noting that when s = 13, the four datasets have achieved the highest OA. Therefore, s = 13 is selected as the input space size of the proposed network.

In addition, different learning rates and batch sizes have a greater impact on the performance of the model. In order to explore the optimal learning rate and batch size of the proposed network, some relevant experiments were carried out, and the experimental results are shown in Fig. 5. Fig. 5(a)-(d)shows the results of experiments on Indian Pines, Pavia,

¹Registered trademark.

```
<sup>7</sup>, <sup>2</sup>Trademarked.
```

Authorized licensed use limited to: Harbin Engineering Univ Library. Downloaded on July 01,2023 at 08:01:03 UTC from IEEE Xplore. Restrictions apply.

	Indian Pines				Salinas		
Class	Class name	Train	Test	Class	Class name	Train	Test
1	Alfalfa	4	42	1	Brocoil-green-weeds_1	100	1909
2	Corn-notill	142	1286	2	Brocoil-green-weeds_2	186	3540
3	Corn-mintill	82	748	3	Fallow	98	1878
4	Corn	23	214	4	Fallow-rough-plow	69	1325
5	Grass-pasture	48	435	5	Fallow-smooth	133	2545
6	Grass-trees	72	658	6	Stubble	197	3762
7	Grass-pasture-mowed	3	25	7	Celery	178	3401
8	Hay-windrowed	47	431	8	Grapes-untrained	563	10708
9	Oats	3	17	9	Soil-vinyard-develop	310	5893
10	Soybean-notill	97	875	10	Corn-senesced-green-weeds	163	3115
11	Soybean-mintill	245	2210	11	Lettuce-romaine-4wk	53	1015
12	Soybean-clean	59	534	12	Lettuce-romaine-5wk	96	1831
13	Wheat	20	185	13	Lettuce-romaine-6wk	45	871
14	Woods	126	1139	14	Lettuce-romaine-7wk	53	1017
15	Bldg-Grass-Tree-Drivers	38	348	15	Vinyard-untrained	363	6905
16	Stone-Steel-Towers	9	84	16	Vinyard-vertical-trellis	90	1717
/	Total	1018	9231	/	Total	2697	51432
	Pavia				WHU-Hi-LongKou		
Class	Class name	Train	Test	Class	Class name	Train	Test
1	Asphalt	331	6300	1	Corn	172	34339
2	Meadows	932	17717	2	Cotton	41	8333
3	Gravel	104	1995	3	Sesame	15	3016
4	Trees	153	2911	4	Broad-leaf soybean	316	62896
5	Painted metal sheets	67	1278	5	Narrow-leaf soybean	20	4131
6	Bare Soil	251	4778	6	Rice	59	11795
7	Bitumen	66	1264	7	Water	335	66721
8	Self-blocking bricks	184	3498	8	Roads and houses	35	7089
9	Shadows	47	900	9	Mixed weed	26	5203
/	Total	2135	40641	/	Total	1019	203523

TABLE II DETAILED CATEGORY INFORMATION OF FOUR DATASETS



Fig. 4. Impact of different input space sizes on OA.

Salinas, and WHU-Hi-LongKou datasets, respectively. Among them, different contour colors represent different ranges of OA values, and red and blue represent a gradual decrease in OA values. It can be found that the OA value of the same dataset is more sensitive to different learning rates and batch sizes of the model. Especially for the Indian Pines and WHU-Hi-LongKou datasets, due to the small number of training data samples used in the training process, the learning rate has a significant impact on them.

Specifically, for the Indian Pines dataset, as shown in Fig. 5(a), the optimal learning rate and batch size are $5e^{-4}$ and 64, respectively. For the Pavia dataset, as shown in Fig. 5(b),



Fig. 5. Effect of different learning rates and batch sizes on performance accuracy OA. (a) Experimental results on Indian Pines dataset. (b) Experimental results on Pavia dataset. (c) Experimental results on Salinas dataset. (d) Experimental results on WHU-Hi-LongKou dataset.

when the batch size is 64 or 128, the learning rate has little impact on the performance of OA. Similarly, for the Salinas dataset, as shown in Fig. 5(c), when the learning rate is large and the batch size is large, better OA values can often be obtained. For the WHU-Hi-LongKou dataset, as shown in

Fig. 5(d), when the batch size is 64, the selected learning rate can achieve better classification results. Therefore, through the parameter experiment of the model, it can be found that the best learning rate and batch size of the classification network proposed in this article are $5e^{-3}$ and 64, respectively.

C. Ablation Experiments

In the proposed method, the network mainly includes four parts, Conv2D&3D, SpaFormer, SpeFormer, and AS²FM. In order to better verify the impact of each part on the classification performance OA value. We conducted ablation experiments on them in four datasets, and the experimental results are shown in Table III. Among them, " $\sqrt{}$ " indicates that the module is available, and "-" indicates that the module is not used. There are five cases in total. It can be seen from Table III that Case 1 only includes Conv2D and Conv3D, and the OA value obtained is low. In Cases 2 and 3, SpaFormer and Spe-Former are added based on Conv2D&Conv3D, respectively. It can be found that the accuracy of OA is worth improving greatly. Generally, the features extracted from the two branches will be combined in a cascade (Cat) manner, as in Case 4. In order to better combine these two features, we introduce a balance factor to fuse the features obtained from the two branches. The experiment shows that the OA value of Case 5 is higher than that of Case 4 on the four datasets, which fully proved the effectiveness of this adaptive combination method.

In addition, we also conducted experiments on the impact of different b of PCA (b = 30, b = 60, and b = 90) on classification performance. The experimental results are shown in Fig. 6, and Fig. 6(a)-(c) shows the impact of different b values on the OA (%), running time (s), and parameter (k). From Fig. 6(a), it can be seen that different b values have little impact on the Salinas dataset, and they slowly decrease as b increases in the other three datasets. We infer that this is due to the dimensionality disaster caused by the high-dimensional characteristics of HSIs and the inclusion of redundant features, which leads to a small reduction in classification accuracy. As shown in Fig. 6(b) and (c), it can be seen that with the increase in b, the running time and parameter increase exponentially. Therefore, in our proposed method, we choose b = 30 as the optimal dimensionality reduction parameter for PCA.

D. Analysis of Experimental Results

In order to verify the superiority of the proposed classification network, we have selected a classifier [iterative support vector machine (ISVM)] [44] and a variety of state-of-theart networks based on CNN and Transformer, including 2-D CNN [27], 3-D CNN [28], Hybrid-SN [29], PyResNet [33], LiteDepthwiseNet [40], MCRSCA [30], ViT [25], SF [46], SSFTT [47], SSTN [48], and GAHT [53].

1) Quantitative Analysis: The classification accuracy of OA, AA, Kappa, and each category of all methods on the four datasets is shown in Tables IV–VII. The best classification results are in bold. As can be seen from the Tables IV–VII, CNN-based methods have achieved relatively good classification results due to their strong ability to extract context features. However, due to the limited advanced global features

obtained by CNN, it is easy to fall into the performance bottleneck. In addition, although Transformer-based methods show great potential by building long-distance dependencies, the classification performance of networks built only using Transformer frameworks is not satisfactory, such as ViT and SF. However, the classification network constructed by combining CNN and Transformer framework has achieved good classification results, such as SSFTT, SSTN, GAHT, and the proposed method. It is worth noting that ISVM based on classifier design has also obtained competitive classification results.

In general, the classification accuracy of the proposed classification method is better than that of other comparison methods on the four datasets. This result not only benefits from the proposed method S²FTNet, which combines the advantages of CNN and Transformer, but also benefits from the effective fusion of the extracted SS high-level semantic features. More specifically, compared with the best CNN method among the comparison methods (MCRSCA), the OA value of S²FTNet is 0.38%, 0.39%, 2.78%, and 1.04% higher on the Indian Pines, Pavia, Salinas, and WHU-Hi-LongKou datasets, respectively. Compared with the best Transformer method in the comparison method (SSFTT), the OA value of S²FTNet is 1.07%, 0.23%, 0.39%, and 0.40% higher on the Indian Pines, Pavia, Salinas, and WHU-Hi-LongKou datasets, respectively. Compared with ISVM classifiers, the OA values of S²FTNet on four datasets are 0.90%, 3.36%, 0.03%, and 0.42% higher. It is worth noting that our method achieves 100% accuracy for individual categories in some datasets, for example, category 1 (Alfalfa), category 3 (Corn-mintill), category 7 (Grass pace moved), category 8 (Hay windowed), and category 9 (Oats) on the Indian Pines dataset; category 5 (Painted metal sheets), category 6 (Bare Soil), and category 7 (Bitumen) on Pavia dataset; and category 1 (Brocool-green-weeds_1), category 7 (Celery), category 10 (Corn-senced-green-weeds), category 11 (Lettuce-remaine-4wk), category 12 (Lettuce-remaine-5wk), and category 13 (Lettuce-remaine-6wk) on the Salinas dataset.

2) Visual Evaluation: Figs. 7-10 show the classification results of all methods on four datasets. It can be clearly seen that the visual effect of the proposed method is closer to the real ground object map. On the Indian Pines dataset, the CNN-based classification method has a poor classification effect on edge categories, while the classification method combining CNN and Transformer has better classification results than CNN, which also benefits from more abundant features extracted, including global and local features. The Pavia dataset contains fewer bands, and the distribution of buildings is more complex. The proposed S²FTNet method has less noise in the classification result map, while most of the comparison methods have more classification errors in the category "Meadows." For the Salinas dataset, two categories that are easy to observe, Vinyard untrained and Grapes untrained, our method has the best visual effect, followed by SSFTT. Among them, 2-D CNN, 3-D CNN, ViT, and SF in the comparison method have serious misclassification. For the WHU-Hi-LongKou dataset, the images mainly include crops with similar spectra. Our method combines CNN and Transformer to build a spatial and spectral extraction module,

 TABLE III

 Impact of Different Modules on Network OA Value (%)

Case	Conv2D&3D	SpaFormer	SpeFormer	Cat	AS ² FM	Indian Pines	Pavia	Salinas	WHU-Hi-LongKou
1	\checkmark	-	-	-	-	79.38	96.96	98.64	94.12
2	\checkmark	\checkmark	-	-	-	96.80	98.43	99.27	98.89
3		-	\checkmark	-	-	97.30	98.03	99.43	98.71
4	\checkmark	\checkmark	\checkmark		-	97.85	98.97	99.76	98.92
5	\checkmark	\checkmark	\checkmark	-		98.50	99.38	99.80	99.39



Fig. 6. Comparison of different *b* values of PCA. (a) Impact of different *b* values on OA. (b) Impact of different *b* values on running time. (c) Impact of different *b* values on parameter.

TABLE IV Classification Accuracy of OA, AA, Kappa, and Various Categories of All Methods on the Indian Pines Dataset. The Best Classification Results Are in Bold

	Classifier				CNNs					Tran	sformers		
Methods	ISVM	2DCNN	3DCNN	Hybrid-SN	PyResNet	LiteDepthwiseNet	MCRSCA	ViT	SF	SSFTT	SSTN	GAHT	Proposed
	[44]	[27]	[28]	[29]	[33]	[40]	[30]	[25]	[46]	[47]	[48]	[53]	Toposed
OA (%)	97.60	82.04	81.15	94.31	92.86	96.28	98.12	79.73	88.54	97.43	95.43	83.00	98.50
AA (%)	98.46	89.09	87.15	94.32	92.15	95.03	96.31	83.36	91.81	93.85	84.54	86.22	97.65
$\kappa \times 100$	97.26	79.26	78.26	93.51	91.87	95.75	97.86	76.75	86.88	97.07	94.78	80.49	98.30
1	97.83	100	100	96.91	98.25	92.51	90.54	99.00	100	85.71	39.41	85.56	100
2	91.04	76.01	72.54	90.87	93.47	96.60	98.17	72.83	82.89	98.13	96.91	80.11	98.52
3	99.40	77.24	69.29	92.05	87.13	96.24	98.10	69.39	84.39	96.66	95.55	76.80	100
4	99.58	96.07	95.83	91.32	99.77	95.42	96.79	82.71	91.44	96.73	96.46	88.60	99.01
5	99.59	93.85	95.23	98.73	95.33	95.43	97.33	85.30	93.50	98.39	94.99	91.29	98.78
6	99.86	79.85	81.47	97.87	96.91	98.16	98.24	85.22	91.74	99.09	96.98	88.04	99.68
7	96.43	80.00	60.00	91.32	69.90	81.41	91.36	92.86	100	100	31.12	91.67	100
8	100	97.04	95.89	98.34	97.05	97.08	99.84	92.06	95.34	98.38	96.71	90.06	100
9	100	100	100	95.54	85.19	92.46	86.88	79.50	100	58.82	40.00	87.03	100
10	97.84	83.60	87.14	91.61	89.02	94.35	97.53	76.73	87.01	99.31	88.37	79.48	96.86
11	96.95	75.16	75.66	95.23	93.95	96.42	98.56	77.39	88.19	98.10	96.29	78.98	99.04
12	99.49	82.73	77.70	92.74	85.23	94.88	95.81	68.13	79.56	89.70	92.87	74.75	99.01
13	100	100	99.86	99.52	98.36	99.63	98.54	93.56	95.37	99.46	98.79	93.59	92.00
14	99.72	92.99	93.23	96.79	96.83	97.55	99.76	92.06	93.92	98.68	98.79	92.67	98.14
15	97.70	90.86	91.10	94.74	91.83	94.59	96.31	82.41	91.14	93.97	94.17	81.98	98.78
16	100	100	99	85.57	96.16	97.75	97.16	84.62	94.33	90.47	95.24	98.94	82.50

which fused spectral information and spatial information well. The obtained classification result has a better edge effect and less intraclass noise.

In order to more clearly illustrate the effectiveness of the proposed S²FTNet method, we compared T-SNE visualization of features obtained by various methods (including 3-D CNN, Hybrid-SN, and SSTN) on four datasets. The experimental results are shown in Figs. 11–14. Different colors represent labels of different categories. From left to right, they are the category distribution results of methods 3-D CNN, Hybrid-SN, SSTN, and proposed. More specifically, on the Indian Pines dataset, both 3-D CNN and SSTN methods have serious label

mixing. Although Hybrid-SN has obtained a better intraclass distance than 3-D CNN and SSTN, the interclass distance is still not satisfactory. However, our method has a more obvious cluster, showing better intraclass and interclass distance. For the Pavia dataset, 3-D CNN and SSTN methods performed poorly, and category 2 (yellow), category 4 (gray), and category 9 (yellow) were still seriously mixed. Compared with the Indian Pines dataset, Hybrid-SN performs better. However, our approach is still significantly better. For the Salinas dataset, the category distribution of 3-D CNN, SSTN, and Hybrid-SN is mostly in a strip shape, with a large gap in the distance within the category. However, most of the categories of our

					KE	SULTS ARE IN BOL	.D						
	Classifier				CNNs			Tran	sformers				
Methods	ISVM	2DCNN	3DCNN	Hybrid-SN	PyResNet	LiteDepthwiseNet	MCRSCA	ViT	SF	SSFTT	SSTN	GAHT	Proposed
	[44]	[27]	[28]	[29]	[33]	[40]	[30]	[25]	[46]	[47]	[48]	[53]	rioposeu
OA (%)	96.02	94.55	93.69	97.99	97.72	98.86	98.99	94.35	95.89	99.15	97.20	94.68	99.38
AA (%)	97.70	93.55	93.38	97.49	97.00	98.73	98.38	92.15	93.64	98.62	96.75	94.40	98.89
K	94.80	92.74	91.56	97.33	96.98	98.49	98.66	92.48	94.55	98.87	96.27	92.93	99.18
1	97.39	91.07	88.64	97.19	96.63	99.50	98.37	90.74	93.23	99.67	95.89	91.86	99.56
2	93.53	97.18	96.27	99.22	99.38	99.89	99.87	97.57	98.96	99.99	98.83	97.03	99.96
3	94.76	83.47	83.55	97.66	93.46	99.22	95.52	82.37	82.55	98.59	92.36	88.98	94.14
4	99.09	99.31	98.99	99.08	97.81	99.41	98.42	99.14	99.83	93.71	97.21	98.65	98.04
5	99.85	99.84	99.44	98.18	99.11	99.79	99.97	96.87	99.73	99.84	96.13	100	100
6	97.97	95.81	95.83	98.68	99.17	99.86	99.36	94.43	97.79	99.54	99.95	91.94	100
7	99.70	89.88	91.04	97.36	99.39	99.95	95.96	79.37	80.51	99.53	99.08	92.16	100
8	98.29	86.76	87.87	92.16	91.07	91.49	98.04	89.59	90.31	98.00	92.89	88.92	98.94
9	98.73	98.66	98.76	97.87	96.94	99.51	99.94	99.30	99.88	98.67	98.41	100	99.33

TABLE V Classification Accuracy of OA, AA, Kappa, and Various Categories of All Methods on the Pavia Dataset. The Best Classification Results Are in Bold

TABLE VI

CLASSIFICATION ACCURACY OF OA, AA, KAPPA, AND VARIOUS CATEGORIES OF ALL METHODS ON THE SALINAS DATASET. THE BEST CLASSIFICATION RESULTS ARE IN BOLD

	Classifier				CNNs					Tran	sformers		
Methods	ISVM	2DCNN	3DCNN	Hybrid-SN	PyResNet	LiteDepthwiseNet	MCRSCA	ViT	SF	SSFTT	SSTN	GAHT	Proposed
	[44]	[27]	[28]	[29]	[33]	[40]	[30]	[25]	[46]	[47]	[48]	[53]	Toposed
OA (%)	99.67	96.01	96.62	98.99	96.57	96.88	97.09	97.87	97.72	99.41	94.03	95.56	99.80
AA (%)	99.41	98.02	98.22	99.29	97.27	97.97	98.30	99.43	98.85	99.37	98.08	97.29	99.74
K imes 100	99.63	95.55	96.24	98.88	96.18	96.52	96.75	97.55	97.46	99.34	93.40	95.05	99.78
1	100	99.81	99.72	99.86	87.44	100	99.20	97.87	99.40	99.95	99.37	99.83	100
2	99.60	99.62	99.31	99.98	99.88	99.99	99.08	99.43	99.97	99.92	99.86	99.85	99.75
3	100	99.39	98.51	99.82	80.83	96.37	99.41	97.55	98.75	99.89	98.52	97.05	99.31
4	97.85	99.51	99.42	99.63	99.82	98.34	99.30	98.98	99.79	99.85	99.03	97.07	98.49
5	99.74	99.10	99.02	99.45	99.86	99.42	98.77	98.68	99.13	98.66	98.85	98.80	99.76
6	99.90	99.95	99.96	99.87	96.92	99.98	99.93	99.92	99.93	99.79	99.91	99.89	99.92
7	99.64	99.62	99.76	99.81	99.86	99.82	98.87	99.84	99.96	99.97	99.92	99.01	100
8	99.87	90.84	92.18	98.15	98.42	94.98	94.74	92.45	94.86	98.55	94.63	91.91	99.76
9	99.98	99.73	99.93	99.88	100	99.73	99.84	99.39	99.93	99.88	99.45	99.93	99.98
10	99.42	95.84	96.64	99.29	99.27	96.15	97.75	98.16	98.25	99.13	99.64	95.71	100
11	99.44	97.98	98.69	99.05	98.97	98.09	97.54	94.14	99.11	100	98.80	97.39	100
12	100	99.58	99.59	99.87	99.43	99.38	99.83	99.08	99.75	100	99.91	98.74	100
13	98.91	98.95	99.06	99.07	99.84	97.44	99.99	98.55	100	99.20	99.42	99.03	100
14	96.45	99.17	99.18	98.02	98.13	98.65	98.90	97.35	99.69	95.67	99.60	97.82	99.90
15	99.71	89.75	90.77	97.15	97.79	89.80	90.40	89.02	93.05	99.93	83.00	85.42	99.86
16	100	99.54	99.85	99.73	99.88	99.32	99.19	99.50	99.96	99.48	100	99.10	99.18

TABLE VII

CLASSIFICATION ACCURACY OF OA, AA, KAPPA, AND VARIOUS CATEGORIES OF ALL METHODS ON THE WHU-HI-LONGKOU DATASET. THE BEST CLASSIFICATION RESULTS ARE IN BOLD

	Classifier				CNNs					Tran	sformers		
Methods	ISVM	2DCNN	3DCNN	Hybrid-SN	PyResNet	LiteDepthwiseNet	MCRSCA	ViT	SF	SSFTT	SSTN	GAHT	Proposed
	[44]	[27]	[28]	[29]	[33]	[40]	[30]	[25]	[46]	[47]	[48]	[53]	TToposed
OA (%)	98.97	89.95	95.12	98.60	97.73	98.87	98.35	86.48	92.43	98.99	97.88	96.89	99.39
AA (%)	96.03	80.61	89.72	96.73	96.56	98.52	94.46	73.71	82.60	97.81	96.13	92.55	99.20
$\kappa_{ imes 100}$	98.65	86.56	93.57	98.16	97.01	98.51	97.83	82.07	90.05	98.68	97.21	95.91	99.46
1	99.14	96.03	99.75	99.20	98.83	98.88	99.67	83.68	96.11	98.24	99.86	99.58	99.97
2	99.94	59.83	66.72	96.53	95.35	98.37	94.67	39.75	72.39	99.86	90.14	87.33	99.25
3	81.16	96.67	98.92	94.00	99.68	99.93	90.91	59.38	83.48	99.07	99.41	92.14	99.57
4	99.30	85.20	95.10	98.63	99.53	98.25	98.85	87.93	93.71	99.60	97.18	96.12	99.51
5	89.64	56.31	76.30	95.13	96.01	99.11	82.43	35.03	67.89	97.57	97.81	79.07	99.64
6	99.64	91.99	97.81	98.74	99.03	99.73	99.20	95.80	94.80	99.45	99.79	99.51	99.66
7	99.95	98.86	98.80	99.69	96.64	99.90	99.96	99.25	97.6	99.82	99.76	99.89	99.98
8	98.13	67.55	82.63	95.29	91.76	93.47	92.20	92.09	80.60	90.93	86.26	86.13	90.45
9	97.32	73.03	91.40	93.36	92.19	99.00	92.23	70.50	57.07	95.79	94.94	93.24	98.13

methods are clustered and have large intraclass distances. Due to the large number of sample categories in the WHU-Hi-LongKou dataset, its category distribution visualization effect is relatively full, but it is not difficult to see that there are some mixed categories in 3-D CNN, SSTN, and Hybrid-SN, and the category distribution is relatively scattered. On the contrary, LIAO et al.: SPECTRAL-SPATIAL FUSION TRANSFORMER NETWORK FOR HSI CLASSIFICATION



Fig. 7. Classification visualization maps of all methods on the Indian Pines dataset. (a) Real ground feature map. (b)–(n) Classification map of ISVM, 2-D CNN, 3-D CNN, Hybrid-SN, PyResNet, LiteDepthwiseNet, MCRSCA, ViT, SF, SSFTT, SSTN, GAHT, and proposed, respectively.



Fig. 8. Classification visualization maps of all methods on the Pavia dataset. (a) Real ground feature map. (b)–(n) Classification map of ISVM, 2-D CNN, 3-D CNN, Hybrid-SN, PyResNet, LiteDepthwiseNet, MCRSCA, ViT, SF, SSFTT, SSTN, GAHT, and proposed, respectively.

our method obtains that the features of the same category are more clustered, and the distribution of different categories is more dispersed. In general, the proposed method S²FTNet has better interclass distance and minimized intraclass distance and plays an important role in capturing the relationship between HSI classification samples.

3) Model Hyperparametric Analysis: In the designed network, considering the different importance of features extracted from spatial and spectral branches and their different contributions to the final classification results, we introduced a balance factor λ into the network and weighted the two

branches by fractions. It will be updated gradually with the change in loss value during the training. In order to observe the changes of balance factor λ and loss value, we selected two datasets for the experiment, Indian Pines and WHU-Hi-LongKou datasets. The experimental results are shown in Fig. 15(a) and (b). The red dot represents the balance factor λ value, and the blue dot represents the loss value. The abscissa represents the training epoch. The left and right ordinates have different magnitudes. The left ordinate is the loss value, and the right ordinate is the balance factor value. It can be found that, on the one hand, the training epoch of the two datasets is



Fig. 9. Classification visualization maps of all methods on the Salinas dataset. (a) Real ground feature map. (b)–(n) Classification map of ISVM, 2-D CNN, 3-D CNN, Hybrid-SN, PyResNet, LiteDepthwiseNet, MCRSCA, ViT, SF, SSFTT, SSTN, GAHT, and proposed, respectively.



Fig. 10. Classification visualization maps of all methods on the WHU-Hi-LongKou dataset. (a) Real ground feature map. (b)–(n) Classification map of ISVM, 2-D CNN, 3-D CNN, Hybrid-SN, PyResNet, LiteDepthwiseNet, MCRSCA, ViT, SF, SSFTT, SSTN, GAHT, and proposed, respectively.

about 40, and the loss value is close to 0, which shows that the combination of these two branch features can achieve faster convergence. On the other hand, the balance factor λ updated slowly and tends to be stable with the increase in epoch, and the stable value is about 0.590. The above results show that the features extracted by the SpaFormer branch and the SpaFormer branch are different in importance, and the SpaFormer branch

accounts for a larger proportion than the SpeFormer branch, and the SS features obtained are more abundant. Finally, the classification performance can be effectively improved by the adaptive fusion of these two features. For the Indian Pines and WHU-Hi-LongKou datasets, the former has many categories, while the latter has a large spatial resolution. The long-distance spectral and spatial features extracted by the SpaFormer branch



Fig. 11. T-SNE visualization of different methods on the Indian Pines dataset. (a) 3-D CNN. (b) Hybrid-SN. (c) SSTN. (d) Proposed.



Fig. 12. T-SNE visualization of different methods on the Pavia dataset. (a) 3-D CNN. (b) Hybrid-SN. (c) SSTN. (d) Proposed.



Fig. 13. T-SNE visualization of different methods on the Salinas dataset. (a) 3-D CNN. (b) Hybrid-SN. (c) SSTN. (d) Proposed.



Fig. 14. T-SNE visualization of different methods on the WHU-Hi-LongKou dataset. (a) 3-D CNN. (b) Hybrid-SN. (c) SSTN. (d) Proposed.

TABLE VIII Comparison of Classification OA (%) Results of S²FTNet Combined With Different Classifiers

datasets	Indian Pines	Pavia	Salinas	WHU-Hi-LongKou
SVM_EPF	88.95	76.21	94.38	97.33
SVM_IEPF	97.72	99.56	99.63	98.97
S ² FTNet	98.50	99.38	99.80	99.39
S ² FTNet_EPF	97.34	78.39	99.07	98.86
S ² FTNet_IEPF	98.96	99.90	99.91	99.59

contribute greatly to the classification results of the two datasets.

4) Combining Different Classifiers: In this section, we select two spatial and spectral classifiers, EPF and IEPF, and combine our method with these two classifiers for exper-

iments. The experimental results are shown in Table VIII. From the results, we can see that the method with the lowest classification accuracy OA value is EPF. IEPF has improved EPF and greatly improved the classification accuracy. In addition, our method combines EPF and IEPF classifiers, and it can be found that compared to EPF, the classification accuracy of S²FTNet_EPF has been improved on all four datasets. Similarly, compared to IEPF, S²FTNet_IEPF can also effectively improve the classification performance. This shows that our proposed method can effectively extract spatial and spectral features.

5) Model Efficiency Analysis: In order to evaluate the running efficiency of the proposed methods, this article conducts running efficiency test experiments for all methods, and Table IX shows the results of the experiments. As can be seen from Table IX, compared with the method SSFTT, which



Fig. 15. Changes of balance factor λ and loss value on different datasets. (a) Indian Pines dataset. (b) WHU-Hi-LongKou dataset.

 TABLE IX

 Comparison of Running Time of All Methods on Four Datasets

datasets	time	2DCNN	3DCNN	Hybrid-SN	PyResNet	LiteDepthwiseNet	MCRSCA	ViT	SF	SSFTT	SSTN	GAHT	Proposed
Indian Pines	Train.(s)	0.36	1.34	0.38	2.39	1.72	1.61	2.98	3.76	0.17	0.41	0.25	0.42
	Test.(s)	6.73	32	2.84	6.98	3.53	5.00	10.95	71.46	1.06	1.19	0.81	1.30
Pavia	Train.(s)	0.44	1.79	0.37	4.94	1.75	3.93	2.31	3.59	0.31	0.76	0.42	0.79
Pavia	Test.(s)	15.4	78	12.76	32.87	7.95	43.13	17.53	40.92	3.23	5.10	2.85	5.15
Salinas	Train.(s)	1.36	4.38	0.27	6.15	4.60	3.24	6.43	6.25	0.19	1.06	0.63	1.12
Sainas	Test.(s)	38.6	41.8	16.21	41.39	21.37	14.10	24.56	33.30	2.76	7.16	4.27	7.29
WHU-Hi-LongKou	Train.(s)	0.42	2.25	1.85	1.63	2.33	1.14	1.38	1.59	0.08	0.44	0.38	0.39
WHU-HI-LongKou –	Test.(s)	6.8	7.0	79.10	171.8	130.42	30.00	38.13	44.25	8.68	32.25	21.70	35.98



Fig. 16. Comparison of different training sample percentages. (a) Pavia dataset. (b) Salinas.

requires the shortest training time and test time, the training time and test time required for the method S^2FTNet proposed in this article are slightly longer. This is because the proposed method is a two-branch Transformer structure. Compared with other Transformer-based methods, S^2FTNet generally requires less running time. In addition, compared with the CNN-based method, the Transformer-based method requires much less training time and testing time. In general, the efficiency of Transformer-based method is significantly higher than that of CNN-based method. Compared with other methods, the running time of the proposed S^2FTNet is relatively close to

that of the optimal method. The experiment fully shows that S^2FTNet not only has good classification accuracy but also has satisfactory operation efficiency.

6) Comparison of Different Training Sample Percentages: The percentage of training samples plays a decisive role in the HSI classification. However, the lack of labeled samples limits the training of the model. Therefore, it is necessary to verify the effectiveness of the method under small training samples. In this article, we selected 0.5%, 1%, 5%, and 10% of Pavia and Salinas datasets for experiments. The experimental results are shown in Fig. 16. The abscissa represents the percentage of training samples, and the ordinate represents the OA value. It can be seen that our method has obtained the best results under different training sample percentages. In addition, the suboptimal methods for Pavia and Salinas datasets are LiteDepthwiseNet and SSFTT, respectively. It is worth noting that our method has an OA value exceeding 95% on both datasets at a 0.5% sample percentage. Through small sample experiments, we verify that the proposed method can also achieve better classification accuracy under limited training samples.

IV. CONCLUSION

In this article, we proposed an S²FTNet method, which fully considers the spectral sequence and long-distance dependence of HSI data. Different from the traditional CNN-based methods, the proposed method combines CNN and Transformer frameworks, making up the disadvantage that CNN is difficult to describe HSI long-distance correlation. Specifically, the proposed S²FTNet includes two branches, SpaFormer branch and SpeFormer branch. Among them, the SpaFormer branch adopts CNN and the improved Transformer block to establish the long-distance dependence of spectral and spatial, which enriches the SS features. The SpeFormer branch adopts the method of preserving spectral sequence, combined with the improved MHD-SA and Conv, to explore the long-distance dependence between different spectral bands. Due to the different importance of the extracted features, in order to balance the high-level semantic features extracted from the two branches, this article also proposed an AS²FM. Finally, in order to verify the advantages of the proposed method, three classical datasets and a new dataset are chosen and a series of experiments is carried out, which verified the effectiveness of the proposed method.

In the future, we will further explore the HSI classification method based on Transformer and extract more representative semantic features through a small number of labeled samples to reduce the demand of the model on the number of training samples.

ACKNOWLEDGMENT

The authors would like to thank the Editor-in-Chief, the Associate Editor, and the reviewers for their insightful comments and suggestion. They would also like to thank Prof. Zhong's team from the State Key Laboratory of Surveying, Mapping and Remote Sensing Information Engineering, Wuhan University, Wuhan, China, for the WHU-Hi-LongKou dataset.

REFERENCES

- Q. Huang, W. Li, B. Zhang, Q. Li, R. Tao, and N. H. Lovell, "Blood cell classification based on hyperspectral imaging with modulated Gabor and CNN," *IEEE J. Biomed. Health Informat.*, vol. 24, no. 1, pp. 160–170, Jan. 2020.
- [2] G. A. Lampropoulos, T. Liu, S.-E. Qian, and C. Fei, "Hyperspectral classification fusion for classifying different military targets," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Boston, MA, USA, Jul. 2008, pp. 262–265.

- [3] D. Hong et al., "Interpretable hyperspectral artificial intelligence: When nonconvex modeling meets hyperspectral remote sensing," *IEEE Geosci. Remote Sens. Mag. Replaces Newsletter*, vol. 9, no. 2, pp. 52–87, Jun. 2021.
- [4] D. M. Tratt, K. N. Buckland, E. R. Keim, and P. D. Johnson, "Urban-industrial emissions monitoring with airborne longwave-infrared hyperspectral imaging," in *Proc. 8th Workshop Hyperspectral Image Signal Process., Evol. Remote Sens. (WHISPERS)*, Los Angeles, CA, USA, Aug. 2016, pp. 1–5.
- [5] L. Mou and X. X. Zhu, "Learning to pay attention on spectral domain: A spectral attention module-based convolutional network for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 1, pp. 110–122, Jan. 2020.
- [6] C. Yu, R. Han, M. Song, C. Liu, and C. Chang, "A simplified 2D-3D CNN architecture for hyperspectral image classification based on spatial–spectral fusion," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 13, pp. 2485–2501, 2020.
- [7] J. He, L. Zhao, H. Yang, M. Zhang, and W. Li, "HSI-BERT: Hyperspectral image classification using the bidirectional encoder representation from transformers," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 1, pp. 165–178, Jan. 2020.
- [8] F. Melgani and L. Bruzzone, "Classification of hyperspectral remote sensing images with support vector machines," *IEEE Trans. Geosci. Remote Sens.*, vol. 42, no. 8, pp. 1778–1790, Aug. 2004.
- [9] M. D. Farrell and R. M. Mersereau, "On the impact of PCA dimension reduction for hyperspectral detection of difficult targets," *IEEE Geosci. Remote Sens. Lett.*, vol. 2, no. 2, pp. 192–195, Apr. 2005.
- [10] G. Camps-Valls, L. Gomez-Chova, J. Munoz-Mari, J. Vila-Frances, and J. Calpe-Maravilla, "Composite kernels for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 3, no. 1, pp. 93–97, Jan. 2006.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* (CVPR), Jun. 2016, pp. 770–778.
- [12] Y. Chen, J. Li, H. Xiao, X. Jin, S. Yan, and J. Feng, "Dual path networks," in *Proc. NIPS*, 2017, pp. 1–9.
- [13] S. Sabour, N. Frosst, and G. E Hinton, "Dynamic routing between capsules," 2017, arXiv:1710.09829.
- [14] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2017.
- [15] C. Li, R. Cong, J. Hou, S. Zhang, Y. Qian, and S. Kwong, "Nested network with two-stream pyramid for salient object detection in optical remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 9156–9166, Nov. 2019.
- [16] 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.
- [17] 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 Trans. Geosci. Remote Sens.*, vol. 54, no. 10, pp. 6232–6251, Oct. 2016.
- [18] Z. Zhong, J. Li, D. A. Clausi, and A. Wong, "Generative adversarial networks and conditional random fields for hyperspectral image classification," *IEEE Trans. Cybern.*, vol. 50, no. 7, pp. 3318–3329, Jul. 2020.
- [19] J. Wang, F. Gao, J. Dong, and Q. Du, "Adaptive DropBlock-enhanced generative adversarial networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 6, pp. 5040–5053, Jun. 2021.
- [20] R. Hang, Q. Liu, D. Hong, and P. Ghamisi, "Cascaded recurrent neural networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 8, pp. 5384–5394, Aug. 2019.
- [21] S. Hao, W. Wang, and M. Salzmann, "Geometry-aware deep recurrent neural networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 3, pp. 2448–2460, Mar. 2021.
- [22] H. Zhang, J. Zou, and L. Zhang, "EMS-GCN: An end-to-end mixhop superpixel-based graph convolutional network for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2022.
- [23] Y. Ding, Y. Chong, S. Pan, Y. Wang, and C. Nie, "Spatial–spectral unified adaptive probability graph convolutional networks for hyperspectral image classification," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Sep. 23, 2021, doi: 10.1109/TNNLS.2021.3112268.

- [24] M. E. Paoletti et al., "Capsule networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 4, pp. 2145–2160, Apr. 2019.
- [25] A. Dosovitskiy et al., "An image is worth 16×16 words: Transformers for image recognition at scale," 2020, arXiv:2010.11929.
- [26] W. Hu, Y. Huang, L. Wei, F. Zhang, and H. Li, "Deep convolutional neural networks for hyperspectral image classification," *J. Sensors*, vol. 2015, pp. 1–12, Jun. 2015.
- [27] W. Zhao and S. Du, "Spectral-spatial feature extraction for hyperspectral image classification: A dimension reduction and deep learning approach," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 8, pp. 4544–4554, Aug. 2016.
- [28] A. Ben Hamida, A. Benoit, P. Lambert, and C. Ben Amar, "3-D deep learning approach for remote sensing image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 8, pp. 4420–4434, Aug. 2018.
- [29] 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 Geosci. Remote Sens. Lett.*, vol. 17, no. 2, pp. 277–281, Feb. 2020.
- [30] R. Shang, H. Chang, W. Zhang, J. Feng, Y. Li, and L. Jiao, "Hyperspectral image classification based on multiscale cross-branch response and second-order channel attention," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5532016.
- [31] D. Hong, L. Gao, J. Yao, B. Zhang, A. Plaza, and J. Chanussot, "Graph convolutional networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 7, pp. 5966–5978, Jul. 2021.
- [32] Z. Zhong, J. Li, Z. Luo, and M. Chapman, "Spectral-spatial residual network for hyperspectral image classification: A 3-D deep learning framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 2, pp. 847–858, Feb. 2018.
- [33] M. E. Paoletti, J. M. Haut, R. Fernandez-Beltran, J. Plaza, A. J. Plaza, and F. Pla, "Deep pyramidal residual networks for spectral-spatial hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 2, pp. 740–754, Feb. 2019.
- [34] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 2011–2023, Aug. 2020.
- [35] S. Woo, "CBAM: Convolutional block attention module," in Proc. Eur. Conf. Comput. Vis. (ECCV), 2018, pp. 1–17.
- [36] J. Fu et al., "Dual attention network for scene segmentation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 3141–3149.
- [37] L. Wang, J. Peng, and W. Sun, "Spatial-spectral squeeze-and-excitation residual network for hyperspectral image classification," *Remote Sens.*, vol. 11, no. 7, p. 884, Apr. 2019.
- [38] K. He et al., "A dual global-local attention network for hyperspectral band selection," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5527613.
- [39] S. Mei, X. Li, X. Liu, H. Cai, and Q. Du, "Hyperspectral image classification using attention-based bidirectional long short-term memory network," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5509612.
- [40] B. Cui, X. Dong, Q. Zhan, J. Peng, and W. Sun, "LiteDepthwiseNet: A lightweight network for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5502915.
- [41] Z. Meng, L. Jiao, M. Liang, and F. Zhao, "A lightweight spectralspatial convolution module for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, 2022, Art. no. 5505105.
- [42] X. Kang, S. Li, and J. A. Benediktsson, "Spectral-spatial hyperspectral image classification with edge-preserving filtering," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 5, pp. 2666–2677, May 2014.
- [43] S. Zhong, C. Chang, and Y. Zhang, "Iterative edge preserving filtering approach to hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 1, pp. 90–94, Jan. 2019.
- [44] S. Zhong, S. Chen, C. Chang, and Y. Zhang, "Fusion of spectral-spatial classifiers for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 6, pp. 5008–5027, Jun. 2021.
- [45] X. He, Y. Chen, and Z. Lin, "Spatial-spectral transformer for hyperspectral image classification," *Remote Sens.*, vol. 13, no. 3, p. 498, Jan. 2021.
- [46] D. Hong et al., "SpectralFormer: Rethinking hyperspectral image classification with transformers," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5518615.
- [47] L. Sun, G. Zhao, Y. Zheng, and Z. Wu, "Spectral-spatial feature tokenization transformer for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5522214.

- [48] Z. Zhong, Y. Li, L. Ma, J. Li, and W. Zheng, "Spectral-spatial transformer network for hyperspectral image classification: A factorized architecture search framework," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5514715.
- [49] X. Huang, M. Dong, J. Li, and X. Guo, "A 3-D-swin transformerbased hierarchical contrastive learning method for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5411415.
- [50] J. Bai et al., "Hyperspectral image classification based on multibranch attention transformer networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5535317.
- [51] J. Zou, W. He, and H. Zhang, "LESSFormer: Local-enhanced spectralspatial transformer for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5535416.
- [52] R. Song, Y. Feng, W. Cheng, Z. Mu, and X. Wang, "BS2T: Bottleneck spatial-spectral transformer for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5532117.
- [53] S. Mei, C. Song, M. Ma, and F. Xu, "Hyperspectral image classification using group-aware hierarchical transformer," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5539014.
- [54] M. Lin, Q. Chen, and S. Yan, "Network in network," 2013, arXiv:1312.4400.



Diling Liao received the bachelor's degree from the Zhuhai College, Jilin University, Zhuhai, China, in 2019. He is currently pursuing the master's degree with Qiqihar University, Qiqihar, China.

His research interests include hyperspectral image processing and machine learning.



Cuiping Shi (Member, IEEE) received the M.S. degree from Yangzhou University, Yangzhou, China, in 2007, and the Ph.D. degree from the Harbin Institute of Technology (HIT), Harbin, China, in 2016. From 2017 to 2020, she held a post-doctoral

research position at the College of Information and Communications Engineering, Harbin Engineering University, Harbin. She is currently a Professor with the Department of Communication Engineering, Qiqihar University, Qiqihar, China. She has published two academic books about remote sensing

image processing and more than 60 articles in journals and conference proceedings. Her main research interests include remote sensing image processing, pattern recognition, and machine learning.

Dr. Shi's doctoral dissertation won the Nomination Award of Excellent Doctoral Dissertation of HIT in 2016.



Liguo Wang (Member, IEEE) received the M.S. and Ph.D. degrees in signal and information processing from the Harbin Institute of Technology, Harbin, China, in 2002 and 2005, respectively.

From 2006 to 2008, he held a post-doctoral research position at the College of Information and Communications Engineering, Harbin Engineering University, Harbin, where he is currently a Professor. Since 2020, he has worked with the College of Information and Communication Engineering, Dalian Nationalities University, Dalian, China. He has pub-

lished two books about hyperspectral image processing and more than 130 articles in journals and conference proceedings. His main research interests include remote sensing image processing and machine learning.



一、检索要求

- 1. 委托人: 石翠萍 Shi, CP (Shi, Cuiping)
- 2. 委托单位:齐齐哈尔大学
- 3. 检索目的:论文被 SCI-E 收录情况

二、检索范围

Science Citation Index Expanded (SCI-EXPANDED)	1990-present	网络版
JCR-(Journal Citation Reports)	2022	网络版
中国科学院文献情报中心期刊分区表(升级版)	2022	网络版

检索报告

三、检索结果

委托人提供的1篇论文被SCI-E收录,论文被收录、所在期刊的JCR影响因子、中科院期刊分区(升级版)情况见附件一。

特此证明!







1

目初后

附件一: SCI-E收录情况

标题: A Spectral-Spatial Fusion Transformer Network for Hyperspectral Image Classification 作者: Liao, DL (Liao, Diling); <u>Shi, CP (Shi, Cuiping)</u>; Wang, LG (Wang, Liguo) 来源出版物: IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING 卷: 61 文献号: 5515216 DOI: 10.1109/TGRS.2023.3286950 出版年: 2023 Web of Science 核心合集中的 "被引频次": 0 被引频次合计: 0 使用次数 (最近 180 天): 2 使用次数 (2013 年至今): 2

引用的参考文献数:54

摘要: In the past, deep learning (DL) technologies have been widely used in hyperspectral image (HSI) classification tasks. Among them, convolutional neural networks (CNNs) use fixed-size receptive field (RF) to obtain spectral and spatial features of HSIs, showing great feature extraction capabilities, which are one of the most popular DL frameworks. However, the convolution using local extraction and global parameter sharing mechanism pays more attention to spatial content information, which changes the spectral sequence information in the learned features. In addition, CNN is difficult to describe the long-distance correlation between HSI pixels and bands. To solve these problems, a spectral-spatial fusion Transformer network (S2FTNet) is proposed for the classification of HSIs. Specifically, S2FTNet adopts the Transformer framework to build a spatial Transformer module (SpaFormer) and a spectral Transformer module (SpeFormer) to capture image spatial and spectral long-distance dependencies. In addition, an adaptive spectral-spatial fusion mechanism (AS(2)FM) is proposed to effectively fuse the obtained advanced high-level semantic features. Finally, a large number of experiments were carried out on four datasets, Indian Pines, Pavia, Salinas, and WHU-Hi-LongKou, which verified that the proposed S2FTNet can provide better classification performance than other the state-of-the-art networks.

入藏号: WOS:001022708100025

语言: English

文献类型: Article

作者关键词: Deep learning (DL); fusion; hyperspectral image (HSI); long-distance dependence KeyWords Plus: RESIDUAL NETWORK

地址: [Liao, Diling; Shi, Cuiping] Qiqihar Univ, Dept Commun Engn, Qiqihar 161000, Peoples R China.

[Wang, Liguo] Dalian Nationalities Univ, Coll Informat & Commun Engn, Dalian 116000, Peoples R China.

通讯作者地址: Shi, CP (通讯作者), Qiqihar Univ, Dept Commun Engn, Qiqihar 161000, Peoples R China.

电子邮件地址: 2020910228@qqhru.edu.cn; shicuiping@qqhru.edu.cn; wangliguo@hrbeu.edu.cn Affiliations: Qiqihar University; Dalian Minzu University

出版商: IEEE-INST ELECTRICAL ELECTRONICS ENGINEERS INC

出版商地址: 445 HOES LANE, PISCATAWAY, NJ 08855-4141 USA

Web of Science Index: Science Citation Index Expanded (SCI-EXPANDED)

Web of Science 类别: Geochemistry & Geophysics; Engineering, Electrical & Electronic; Remote Sensing; Imaging Science & Photographic Technology

研究方向: Geochemistry & Geophysics; Engineering; Remote Sensing; Imaging Science & Photographic Technology

IDS 号: L4AS8

ISSN: 0196-2892

eISSN: 1558-0644

29 字符的来源出版物名称缩写: IEEE T GEOSCI REMOTE

ISO 来源出版物缩写: IEEE Trans. Geosci. Remote Sensing

来源出版物页码计数:16

基金资助致谢:

基金资助机构 授权号

National Natural Science Foundation of China

42271409

62071084

Heilongjiang Science Foundation Project of China

LH2021D022

Fundamental Research Funds in Heilongjiang Provincial Universities of China 145209149

This work was supported in part by the National Natural Science Foundation of China under Grant 42271409 and Grant 62071084, in part by the Heilongjiang Science Foundation Project of China under Grant LH2021D022, and in part by the Fundamental Research Funds in Heilongjiang Provincial Universities of China under Grant 145209149.

输出日期: 2023-08-30

期刊影响因子 ™ 2022: 8.2

中国科学院文献情报中心期刊分区(升级版, 2022)截图如下:

IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING

刊名	IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING		
年份	2022		
ISSN	0196-2892		
Review	否		
Open Access	否		
Web of Science	SCIE		
	学科	分区	Top期刊
大类	工程技术	1	是
	GEOCHEMISTRY & GEOPHYSICS 地球化学与地球物理	1	H
1. **	ENGINEERING, ELECTRICAL & ELECTRONIC 工程:电子与电气	2	2
小交	IMAGING SCIENCE & PHOTOGRAPHIC TECHNOLOGY 成像科学与照相技术	2	5
	REMOTE SENSING 遥感	2	5

The End