# SpectralDiff: A Generative Framework for Hyperspectral Image Classification with Diffusion Models

Ning Chen, Jun Yue, Leyuan Fang, *Senior Member, IEEE,* and Shaobo Xia

*Abstract*—Hyperspectral Image (HSI) classification is an important issue in remote sensing field with extensive applications in earth science. In recent years, a large number of deep learning-based HSI classification methods have been proposed. However, existing methods have limited ability to handle high-dimensional, highly redundant, and complex data, making it challenging to capture the spectral-spatial distributions of data and relationships between samples. To address this issue, we propose a generative framework for HSI classification with diffusion models (SpectralDiff) that effectively mines the distribution information of high-dimensional and highly redundant data by iteratively denoising and explicitly constructing the data generation process, thus better reflecting the relationships between samples. The framework consists of a spectral-spatial diffusion module, and an attention-based classification module. The spectral-spatial diffusion module adopts forward and reverse spectral-spatial diffusion processes to achieve adaptive construction of sample relationships without requiring prior knowledge of graphical structure or neighborhood information. It captures spectral-spatial distribution and contextual information of objects in HSI and mines unsupervised spectral-spatial diffusion features within the reverse diffusion process. Finally, these features are fed into the attention-based classification module for per-pixel classification. The diffusion features can facilitate cross-sample perception via reconstruction distribution, leading to improved classification performance. Experiments on three public HSI datasets demonstrate that the proposed method can achieve better performance than state-of-the-art methods. For the sake of reproducibility, the source code of SpectralDiff will be publicly available at https://github.com/chenning0115/SpectralDiff.

*Index Terms*—Deep neural network, hyperspectral image classification, spectral-spatial diffusion, feature extraction, diffusion models, deep generative model.

## I. INTRODUCTION

Ning Chen is with the Institute of Remote Sensing and Geographic Information System, Peking University, Beijing 100871, China (e-mail: chenning0115@pku.edu.cn).

Jun Yue is with the School of Automation, Central South University, Changsha 410083, China (e-mail: jyue@pku.edu.cn).

Leyuan Fang is with the College of Electrical and Information Engineering, Hunan University, Changsha 410082, China, and also with the Peng Cheng Laboratory, Shenzhen 518000, China (e-mail: fangleyuan@gmail.com).

Shaobo Xia is with the Department of Geomatics Engineering, Changsha University of Science and Technology, Changsha 410114, China (e-mail: shaobo.xia@csust.edu.cn).

HYPERSPECTRAL imaging is a cutting-edge technology that enables the acquisition of high-resolution spectral information of objects. By integrating spatial and spectral reflectance information, each pixel in a hyperspectral image (HSI) corresponds to a unique spectral curve, providing rich information for identification and differentiation of diverse materials and surfaces. Beyond the limited perception of human eyes, the spectral detection range of hyperspectral imaging empowers a comprehensive understanding of nature [1], [2]. This technology has demonstrated significant potential in numerous fields, such as environmental management, agriculture, land management, ecology, geology, urban planning, and oceanography. HSI classification, which involves assigning pixels to specific land cover classes, such as soil and grass, is one of the most significant applications of hyperspectral imaging. As a fundamental component of hyperspectral data processing, HSI classification plays an indispensable role in most hyperspectral imaging applications [3], [4].

The high dimensionality of HSI poses a significant challenge for the accurate classification of pixels. With hundreds of spectral bands and massive amounts of data, it can be difficult to identify relevant features. To overcome this challenge, researchers have developed a range of methods to map spectral vectors from high-dimensional space to low-dimensional feature space in order to extract effective spectral features. These methods include classic statistical transformation techniques such as principal component analysis (PCA) [5], minimum noise fraction (MNF) [6], local preserving projection (LPP) [7], linear discriminant analysis (LDA) [8], independent component analysis (ICA) [9], and sparse preserving projection (SPP) [10]. However, the spatial heterogeneity and homogeneity of HSIs make it difficult to fully utilize them by extracting spectral features alone. To address this limitation, researchers have proposed a series of methods to jointly extract spatial and spectral features, such as extended morphological profile (EMP) [11] and extended attribute profile (EAP) [12].

With the successful introduction and rapid development of deep neural networks (DNN) [13], [14], it has achieved outstanding performance in image classification [15], [16], image segmentation [17], [18], [19], instance segmentation [20], image vectorization [21], and object detection [22], [23]. DNN offers a solution by offering an adaptable and potent framework for automatically learning complicated features and relationships in data [24]. The HSI classification accuracy has been steadily improved by a number of HSI spectral feature extraction techniques based on DNN. Due to the fact that HSI

has both spectral and spatial features, some spectral-spatial feature extraction techniques combining DNN have been proposed in order to fully explore the three-dimensional data characteristics [25], including stacked auto-encoders [26], deep fully convolutional network [25], deep prototypical network [27], spatial pyramid pooling [28], and spectralformer [29].

Despite achieving favorable outcomes in HSI classification, DNN-based methods remain limited in their ability to model spectral-spatial relationships across samples. Current approaches primarily rely on graph neural networks (GNNs) for modeling sample relationships [30], [31], [32], [33]. However, using GNNs to measure the relationships between all samples requires designing graph structures or neighborhood information, which increases the complexity of the design and implementation processes and introduces subjectivity. From a data distribution standpoint, GNNs-based methods are ineffective in capturing the spectral-spatial distribution of data and do not fully represent the data generation process. As a result, these methods have limited perceptiveness towards contextual features when constructing relationships between samples.

To address the aforementioned challenge, we propose a generative framework based on diffusion models, named SpectralDiff. The proposed framework constructs the data generation process through iterative denoising, thereby obtaining spectral-spatial distribution information of the data and better capturing the relationships between samples. The framework consists of two modules, namely the spectral-spatial diffusion module and the attention-based classification module. In the spectral-spatial diffusion module, we use hyperspectral cube data with spatial and spectral dimensions and construct a hyperspectral channel distribution in the spectral-spatial latent feature space through the spectral-spatial diffusion process, which is a Markov process composed of forward and reverse processes. In the forward process, Gaussian noise is added to the hyperspectral channel, while in the reverse process, the noise is removed through multiple time steps with a spectral-spatial denoising network. By constructing the distribution of samples and the explicit sample generation process, the relationships between samples are constructed through the hidden variables of the spectral-spatial denoising network. We extract the spectral-spatial diffusion features that aggregate these latent variables and feed the generated features into the attention-based classification module. This module directly generates per-pixel classification results for hyperspectral data, thus achieving cross-sample perception and improving classification performance.

The primary innovative aspect of our study is adopting a generative perspective, wherein we demonstrate the process of modeling sample generation to acquire spectral-spatial features infused with contextual information of the spectral-spatial distribution. This is achieved without the need for pre-defined graph structures or neighborhood information. It is worth noting that our proposed generative framework is loosely coupled between each module, which can evolve independently in the future, thus boosting the development of HSI classification. The main contributions of this paper can be summarized as follows.

- We formulate the construction of relationships between samples from a generative perspective, designated as a spectral-spatial diffusion process. To the best of our knowledge, this is the first study to apply diffusion models to HSI classification.
- We present an HSI classification framework based on forward and reverse diffusion processes, which harnesses the iterative generative process of constructing samples for obtaining spectral-spatial features endowed with contextual information, all without the need for pre-defined graph structures or neighborhood information.
- Numerous experiments demonstrate the superiority of the proposed method over existing state-of-the-art techniques. Moreover, through ablation experiments, we validate the effectiveness of spectral-spatial diffusion features.

The rest of this article is structured as follows. In Section II, we provide a brief overview of the existing work on HSI classification and diffusion models. Section III outlines our proposed SpectralDiff in detail, highlighting the key components of our approach. In Section IV, we present the experimental setup and discuss our results from both qualitative and quantitative perspectives. Additionally, we validate the effectiveness of our approach through ablation experiments. Finally, in Section V, we summarize the conclusions of our study and provide insights into the future directions of research in this field.

## II. RELATED WORK

In this section, we provide an introduction to the background and related work associated with our proposed method. Specifically, we discuss existing techniques for HSI classification based on traditional spectral-spatial feature extraction, as well as those based on deep learning. Furthermore, we present an overview of the diffusion model and its development context, highlighting the key advancements and challenges in this field. Through this exposition, we aim to establish a solid foundation for our proposed method and contextualize our contribution within the broader scope of HSI classification research.

### A. Hyperspectral Image Classification

Hyperspectral sensor captures both spatial and spectral information. In an HSI image, every pixel vector corresponds to a unique spectral curve. The primary objective of HSI image classification is to assign each pixel to a specific land cover class, such as river, forest, lake, farmland, building, grassland, mineral, road and rock. Classification is a pivotal step in the application of HSI and has significant implications for environment, geology, mining, ecology, forestry, agriculture, and other areas [34], [4]. HSI classification enables quick and precise acquisition of ground feature, empowering informed decision-making and management. For instance, in agriculture, HSI classification can be used to track disease progression and crop development, improving crop quality and productivity. In the mining sector, HSI classification can make it easier to find and identify minerals, increasing the effectiveness of mining natural resources [3], [35].

Spectral feature extraction based on continuous spectral signals is widely used in HSI feature extraction, which considers each pixel of the HSI as an independent spectral vector

and ignores the spatial relationships between pixels when generating the ground objects' features [36]. Linear feature extraction methods such as principal component analysis (PCA) [5], linear discriminant analysis (LDA) [8], and minimum noise fraction (MNF) [6], [37] are among the commonly used spectral feature extraction techniques. Nonetheless, the nonlinear nature of HSI has led to the development of numerous nonlinear feature extraction approaches in recent years, complementing the statistical transformation methods based on prior knowledge. In recent years, many nonlinear HSI feature extraction methods for have been proposed [38], [1], including local Fisher discriminant analysis (LFDA) [39], manifold learning [40], sparsity preserving projections (SPP) [41], improved manifold coordinate representations [42] and locality preserving projections (LPP) [7].

The HSI analysis technique involves extracting separable features of patterns from the spectral signal to reduce the dimensionality of the data. Despite its efficacy, this method alone cannot fully utilize the characteristics of HSI due to its spatial homogeneity and heterogeneity. Consequently, incorporating image texture and structural features is necessary to overcome this limitation. Researchers have proposed several methods to extract spatial structure and texture information from HSIs to obtain spatial features such as three-dimensional gray-level cooccurrence [43] and discriminative Gabor feature selection [44]. In the domain of HSI feature extraction and classification, the classic spectral-spatial joint feature extraction methods include extended morphological profile [11], spatial and spectral regularized local discriminant embedding [45], extended attribute profile [12], spatial–spectral manifold alignment [46], directional morphological profiles [11], and spectral-spatial locality preserving projection [47].

Deep learning has emerged as a significant breakthrough in the field of machine learning, providing automatic feature learning from data [13]. HSI classification, in particular, benefits from deep learning models due to their ability to handle complex and nonlinear relationships between input data and output classes [16], leading to improved classification accuracy compared to traditional machine learning methods [48], [37]. To fully utilize the spectral-spatial features of HSI, joint spectral-spatial HSI feature extraction methods based on deep learning, such as recurrent neural network [49], deep residual network [50], [51], capsule networks [52], and transformer [29], have been proposed. While these methods have yielded satisfactory results, challenges remain, such as the inability to capture global relationships between samples. To address this issue, researchers have proposed using graph neural network to model the relationship between samples [30]. However, the use of graph neural network for global relationship modeling incurs high computing and memory costs, and slow gradient decline [32]. In this paper, we propose using diffusion models to model the global relationship of samples and verify the effectiveness of this method through extensive comparative experiments.

## B. Diffusion Models

Diffusion models, also known as diffusion probabilistic models [53], belong to the class of latent variable models (LVM) in machine learning [54]. They draw inspiration from non-equilibrium thermodynamics and construct a Markov chain that gradually introduces random noise into the input data [55]. Subsequently, the Markov chain is trained utilizing variational inference [56], delivering remarkable performance across various domains such as natural language processing [57], [58], time series forecasting [59], [60], and molecular graph modeling [61], [62]. Currently, research on diffusion models is generally based on three primary paradigms [63], namely score-based generative models [64], [65], [66], stochastic differential equations (SDE) [67], and denoising diffusion probabilistic models (DDPM) [68], [53].

The diffusion model aims to infer the underlying structure of a dataset by modeling the diffusion of data points in the latent space [69], which consists of two core processes: forward and reverse. During the forward process, the input data is gradually perturbed by introducing Gaussian noise in multiple time steps. Conversely, the reverse process aims to restore the original input data by reducing the discrepancy between the predicted noise and the increased noise in multiple reverse time steps. In the field of computer vision, this entails training a neural network to perform denoising of images that are blurred by Gaussian noise, by learning the reverse diffusion process [70]. The diffusion model has gained popularity in recent times, owing to its remarkable flexibility and robustness, and has been successfully employed to address a diverse range of intricate visual challenges [56], such as image inpainting [71], image generation [55], [68], [67], [69], image-to-image translation [72], [73], image fusion [70], and image super-resolution [74], [75], [76].

The feature representation learned from the diffusion models has been demonstrated to be highly effective in various discriminating tasks, such as image classification [77], object detection [78] and image segmentation [79], [80]. By constructing the distribution of input data, the diffusion models can effectively capture the underlying patterns and relationships between samples. In this paper, we propose a generative framework for hyperspectral classification based on the diffusion models, which explicitly constructs the data generation process through hierarchical denoising. The proposed framework is able to build the distribution of multi-channel input data and effectively exploit the distributional information of high-dimensional and highly redundant data, thereby improving the performance of HSI classification.

## III. METHODOLOGY

In this section, we provide a detailed description of the proposed SpectralDiff, as shown in Fig. 1. By training the spectral-spatial denoising network to estimate the noise added in the forward spectral-spatial diffusion process, we establish the relationship among all samples.

### A. Spectral-Spatial Diffusion Module

Given an HSI $\mathcal{H} \in \mathbb{R}^{H \times W \times B}$, where $H$ and $W$ denote the height and width of $\mathcal{H}$, respectively. $B$ represents the number of spectral channels. To learn the joint latent structure of images with hundreds of channels, we add noise to the
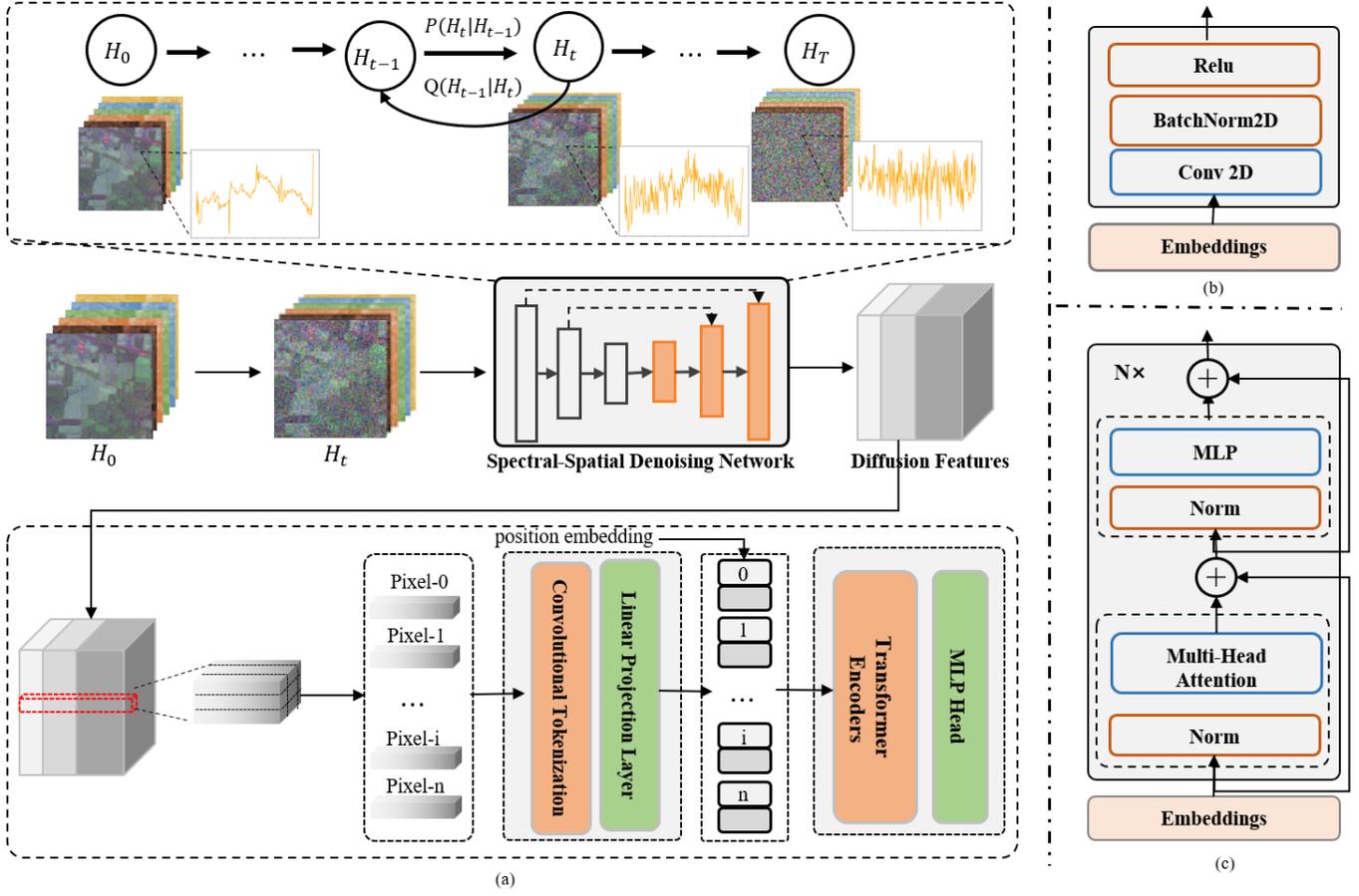
Fig. 1: Overview of the proposed SpectralDiff. (a) The architecture of the generative framework with diffusion process. $\mathcal{H}_0$ and $\mathcal{H}_t$ represent hyperspectral images of timestep 0 and timestep $T$, respectively. $P(\cdot|\cdot)$ and $Q(\cdot|\cdot)$ represent the forward and reverse spectral-spatial diffusion processes, respectively. (b) Structure of Convolutional Tokenization. (c) Structure of Transformer Encoders.

spectral-spatial instance in the forward process, and eliminate the noise added by the forward process by training a spectral-spatial denoising network in the reverse process [68]. The purpose of training the diffusion model with forward and reverse processes is to learn the joint latent structure between hyperspectral channels by simulating the diffusion of hyperspectral channels in the latent space [66], [81].

*1) Forward Spectral-Spatial Diffusion Process:* The forward spectral-spatial diffusion process, drawing inspiration from the principles of non-equilibrium thermodynamics, can be considered as Markov chain [55], [68]. Its progressive stages entail the gradual incorporation of Gaussian noise into the data. In this study, we extract spectral-spatial instance $\mathcal{HI} \in \mathbb{R}^{K \times K \times B}$ using a $K \times K$ neighborhood from the whole HSI $\mathcal{H}$. The forward spectral-spatial diffusion process is distinguished by the "no memory" property, whereby the probability distribution of the HSI at a given time $t + 1$ is exclusively determined by its state at time $t$. At the $t$th step, the spectral-spatial instance imbued with noise is expressed as follows:

$$P(\mathcal{HI}_t|\mathcal{HI}_{t-1}) = \mathcal{N}(\mathcal{HI}_t; \sqrt{\alpha_t}\mathcal{HI}_{t-1}, (1 - \alpha_t)I) \quad (1)$$

where $\mathcal{HI}_{t-1}$ and $\mathcal{HI}_t$ denote the noisy hyperspectral instances at timestep $t - 1$ and $t$, respectively. $I$ stands for the standard normal distribution. $\sqrt{\alpha_t}\mathcal{HI}_{t-1}$ and $(1 - \alpha_t)I$ denote the mean and variance of $P(\mathcal{HI}_t|\mathcal{HI}_{t-1})$, respectively. The variance of Gaussian noise added at timestep $t$ is controlled by a variance schedule referred to as $\alpha_t$. The magnitude of the added Gaussian noise decreases as the value of $\alpha_t$ increases. By Eq. (1), we can derive the expression of $\mathcal{HI}_1 \in \mathbb{R}^{K \times K \times B}$ during the first diffusion as follows:

$$\mathcal{HI}_1 = \sqrt{\alpha_1}\mathcal{HI}_0 + \sqrt{1 - \alpha_1}\epsilon \quad (2)$$

where $\mathcal{HI}_0$ stands for the hyperspectral instance before diffusion. $\epsilon \in \mathbb{R}^{K \times K \times B}$ is the added Gaussian noise. By using Eq. (1) and Eq. (2), the expression of $\mathcal{HI}_t$ can be derived as follows:

$$\begin{cases} \mathcal{HI}_t = \sqrt{\bar{\alpha}_t}\mathcal{HI}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \\ \bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i \end{cases} \quad (3)$$

where $\bar{\alpha}_t$ represents the product of $\alpha_1$ to $\alpha_t$. The computation of $\mathcal{HI}_t$ in the context of forward spectral-spatial diffusion process hinges upon the timestep $t$, variance schedule $\alpha_1, ..., \alpha_t$, and the noise sampled from the standard normal
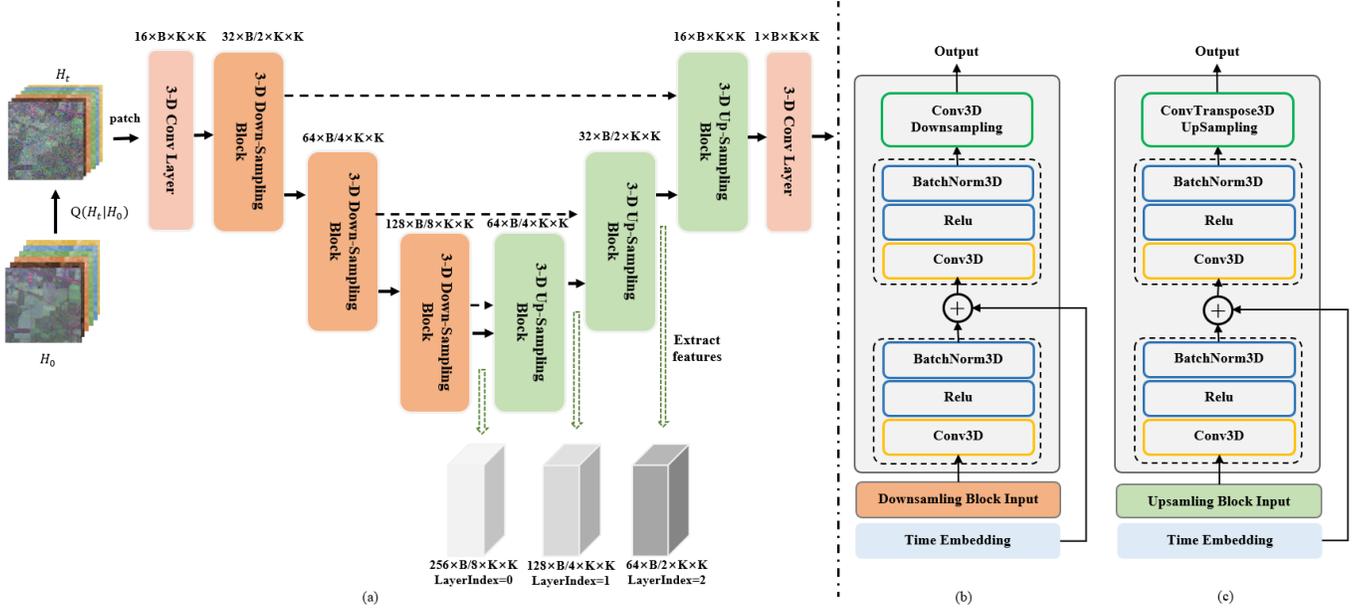
Fig. 2: Structure of the Spectral-Spatial Denoising Network. (a) The original HSI data is processed through a forward diffusion process to obtain image data with added Gaussian noise. Subsequently, the data is passed through an encoder-decoder structure and finally output as a tensor with the same shape as the original HSI data. This output represents the network's prediction for the initial noise $\epsilon$ added. (b) The details of the 3-D down-sampling block. (c) The details of the 3-D up-sampling block.

distribution. Given these inputs, the hyperspectral instance at timestep $t$ can be directly generated by Eq. (3).

*2) Reverse Spectral-Spatial Diffusion Process:* In the process of reverse spectral-spatial diffusion, a spectral-spatial denoising network is trained to gradually denoise the noisy hyperspectral instance to obtain the original hyperspectral instance $\mathcal{HI_0}$ [79]. In the $t$-th step of the reverse diffusion process, denoising operation is performed on the noisy hyperspectral instance $\mathcal{HI_t}$ to obtain the hyperspectral instance of the previous step, which is $\mathcal{HI_{t-1}}$. Given the hyperspectral instance at step $t$, the conditional probability of the hyperspectral instance at step $t-1$ can be expressed as follows [68], [82]:

$$Q(\mathcal{HI_{t-1}}|\mathcal{HI_t}) = \mathcal{N}(\mathcal{HI_{t-1}}; \mu_\theta(\mathcal{HI_t}, t), \sigma_t^2 \boldsymbol{I}) \quad (4)$$

where $\mu_\theta(\mathcal{HI_t}, t)$ and $\sigma_t^2$ denote the mean and variance of the condition distribution $Q(\mathcal{HI_{t-1}}|\mathcal{HI_t})$, respectively. The variance is determined by the variance schedule, which can be expressed as follows:

$$\sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}(1 - \alpha_t) \quad (5)$$

To obtain the mean of the conditional distribution $Q(\mathcal{HI_{t-1}}|\mathcal{HI_t})$, we need to train the network to predict the added noise. The mean can be represented as follows:

$$\mu_\theta(\mathcal{HI_t}, t) = \frac{1}{\sqrt{\alpha_t}}(\mathcal{HI_t} - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathcal{HI_t}, t)) \quad (6)$$

where $\epsilon_\theta(\cdot, \cdot)$ denote the spectral-spatial denoising network whose input is the timestep $t$ and the noisy hyperspectral instance $\mathcal{HI_t}$ at timestep $t$.

*B. Loss Function of Spectral-Spatial Diffusion Process*

For the spectral-spatial denoising network training, the hyperspectral instance $\mathcal{HI_0^i} \in \mathbb{R}^{K \times K \times B}$ are first extracted, followed by sampling the noise matrix $\epsilon^i \in \mathbb{R}^{K \times K \times B}$ of equivalent size from the standard normal distribution. Subsequently, the timestep $t$ is sampled from a uniform distribution $U(\{1, ..., T\})$. With the above-mentioned sampling completed, the noisy hyperspectral instance at timestep $t$ can be calculated by Eq. (3). The noisy hyperspectral instance and timestep are fed into the spectral-spatial denoising network to generate the predicted noise. The loss function of the spectral-spatial diffusion process can be expressed as follows:

$$\begin{aligned}
\mathcal{L}_{ssdp} &= \sum_{i=1}^{N} \left\| \boldsymbol{\epsilon^i} - \epsilon_\theta(\mathcal{HI_t^i}, t) \right\|_1 \\
&= \sum_{i=1}^{N} \left\| \boldsymbol{\epsilon^i} - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathcal{HI_0^i} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|_1
\end{aligned} \quad (7)$$

where $N$ and $i$ represent the total number of hyperspectral instances and the index of a hyperspectral instance, $\epsilon_\theta(\cdot, \cdot)$ represents the spectral-spatial denoising network discussed in the previous section. This network operates by taking two input components: the time step $t$ and the corresponding hyperspectral image $\mathcal{HI_t^i}$ with Gaussian noise added. The symbol $\|\cdot\|_1$ denotes the L1 norm. This loss function quantifies the discrepancy between the initial noise $\boldsymbol{\epsilon^i}$ and the estimated noise. Consequently, the spectral-spatial denoising network has the capability to predict the magnitude of noise introduced to the original image.

## C. Structure of the Spectral-Spatial Denoising Network

The structure of the spectral-spatial denoising network adopts a similar architecture to U-Net [74], as shown in Fig. 2. The network takes the time step $t$ and the corresponding hyperspectral image $\mathcal{H}_t$ with added Gaussian noise as inputs. Considering the potentially large size of remote sensing images, we used the patch form of the image $\mathcal{HI}_t$ as the input instead. The network outputs predictions for the Gaussian noise $\epsilon$. The loss function is shown in Eq. (7), as described in previous section.

Specifically, the input undergoes one 3D convolution layer and three 3D down-sampling blocks to achieve spectral and spatial feature encoding. Subsequently, it is decoded via three 3D up-sampling blocks and another 3D convolution layer to generate the output. Unlike diffusion model processing RGB images, the Spectral-Spatial Denoising Network employs 3D convolution structures for the efficient extraction of spectral and spatial information simultaneously during down-sampling and up-sampling processes. To extract the diffusion feature of every pixel in both encoder and decoder processes, we preserved the spatial dimensions while compressing and restoring the spectral dimensions. The internal structures of 3D down-sampling and 3D up-sampling blocks are depicted in Fig. 2(b) and Fig. 2(c), respectively, consisting of two-layer blocks that incorporate a 3D-convolution layer, a BatchNorm3D layer, and a ReLU activation layer. The down-sampling structure applies a 3D convolution operation with a stride of 2 to gradually reduce the spectral dimension, while the up-sampling block utilizes a 3D deconvolutional layer with a stride of 2 to gradually restore the spectral dimension.

## D. Attention-based Classification Module

Once the training of the spectral spatial denoising network is completed, we will employ the network to generate spectral-spatial diffusion features. Specifically, we utilize the high-hyperspectral image $\mathcal{H}_t$ corrupted by Gaussian noise at time step $t$ as the network's input. During the network's inference process, we extract the activation tensors from the down-sampling and up-sampling blocks as diffusion features. As depicted in Fig. 2(a), in practical implementation, the U-net's presence of shortcut connections enables us to directly obtain the input of each up-sampling layer for the corresponding feature extraction. Importantly, the obtained diffusion features may differ in information across different timestamps and up-sampling layers, potentially affecting the final classification performance.

Fig. 1(a) illustrates the fundamental architecture of the attention-based classification module. Here, we denote the extracted diffusion features as $F \in R^{H \times W \times L}$, where $H$, $W$, and $L$ represent height, width, and the feature channel, respectively. It is worth noting that the spatial dimension of the diffusion features is consistent with that of the raw features. For each pixel $i$, the corresponding channel dimension not only contains the original spectral information, but also includes the deep features fused through the spectral-spatial diffusion module. we use the patch $Z_i \in R^{k \times k \times L}$, which is centered

on the pixel, as the input to the classifier. The spectral-spatial features of $Z_i$ are mapped via Conv2D structure and linear mapping layer. This process produces a token sequence that is appropriate for processing with Transformer. During this process, the region surrounding pixel $i$ is separated into distinct spatial tokens, which are then consolidated through the attention mechanism. The kernel of the Conv2D is set to $3 \times 3$. To provide positional information, the model adds position embedding to the input information before feeding it to the transformer encoder. The transformer encoder structure is depicted in Fig. 1(b). It comprises multiple identical substructures, each of which contains a multi-head attention layer and an MLP layer. Finally, the model maps the output to the classification through the MLP layer. The attention-based classification module combines the CNN and transformer structures to form an effective classifier. This approach utilizes a CNN structure to perform spectral-spatial feature mapping, and a Transformer structure for deep feature extraction, which leads to exceptional classification performance.

## IV. EXPERIMENTS

### A. Experimental Settings

*1) Datasets:* In order to verify the effectiveness of our algorithm, we applied the algorithm to three public datasets: the Indian Pines dataset, the Pavia University dataset, and the Salinas dataset.

The Indian Pines (IP) dataset was collected in 1992 using the Airborne Visible/Infrared Imaging Spectrometer (AVIRIS) Sensor, covering the northwestern region of Indiana in the United States. Comprising of 224 spectral bands, the uncorrected dataset spans a range of 0.4 to $2.5\mu$m and is composed of 145×145 pixels with a spatial resolution of 20 meters. It contains 16 different land-cover classes. For the purpose of experimentation, 24 water-absorption bands and noise bands were removed, and a subset of 200 bands were selected for our analysis.



Fig. 3: The false-color composite image, the corresponding ground-truth map, and the legend of the Indian Pines dataset.

The Pavia University (PU) dataset was collected in 2001 by the Reflective Optics System Imaging Spectrometer (ROSIS) Sensor, covering the Pavia University in Northern Italy. It contains of 115 spectral bands while 12 noise bands were removed so that 103 bands were used. It is composed of 610 × 340 pixels with a spatial resolution of 1.3 m. The dataset convers Nine categories.

TABLE I
TRAINING AND TEST SAMPLE NUMBERS IN THE INDIAN PINES DATASET, THE PAVIA UNIVERSITY DATASET, AND THE SALINAS DATASET

| No. | Indian Pines | | | Pavia University | | | Salinas | | |
|---|---|---|---|---|---|---|---|---|---|
| | Class | Train. | Test. | Class | Train. | Test. | Class | Train. | Test. |
| 1 | Alfalfa | 30 | 16 | Asphalt | 30 | 6601 | BrocoliGreenWeeds1 | 30 | 1979 |
| 2 | CornMotill | 30 | 1398 | Meadows | 30 | 18619 | BrocoliGreenWeeds2 | 30 | 3696 |
| 3 | CornMintill | 30 | 800 | Gravel | 30 | 2069 | Fallow | 30 | 1946 |
| 4 | Corn | 30 | 207 | Trees | 30 | 3034 | FallowRoughPlow | 30 | 1364 |
| 5 | GrassPasture | 30 | 453 | PaintedMetalSheets | 30 | 1315 | FallowSmooth | 30 | 2648 |
| 6 | GrassTrees | 30 | 700 | BareSoil | 30 | 4999 | Stubble | 30 | 3929 |
| 7 | GrassPastureMowed | 14 | 14 | Bitumen | 30 | 1300 | Celery | 30 | 3549 |
| 8 | HayWindrowed | 30 | 448 | SelfBlockingBricks | 30 | 3652 | GrapesUntrained | 30 | 11241 |
| 9 | Oats | 10 | 10 | Shadows | 30 | 917 | SoilVinyardDevelop | 30 | 6173 |
| 10 | SoybeanNotill | 30 | 942 | | | | ComSenescedGreenWeeds | 30 | 3248 |
| 11 | SoybeanMintill | 30 | 2425 | | | | LettuceRomaine4wk | 30 | 1038 |
| 12 | SoybeanClean | 30 | 563 | | | | LettuceRomaine5wk | 30 | 1897 |
| 13 | Wheat | 30 | 175 | | | | LettuceRomaine6wk | 30 | 886 |
| 14 | Woods | 30 | 1235 | | | | LettuceRomaine7wk | 30 | 1040 |
| 15 | BuildingsGrassTreesDrives | 30 | 356 | | | | VinyardUntrained | 30 | 7238 |
| 16 | StoneSteelTowers | 30 | 63 | | | | VinyardVerticalTrellis | 30 | 1777 |
| | Total | 444 | 9805 | Total | 270 | 42506 | Total | 480 | 53649 |



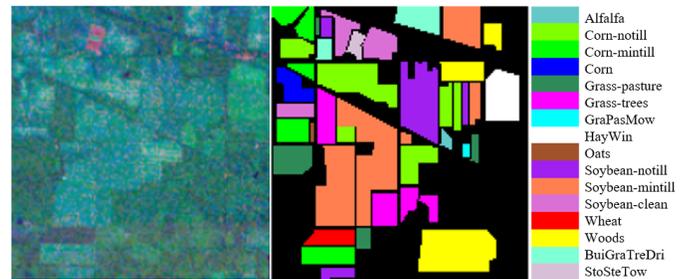Fig. 4: The false-color composite image, the corresponding ground-truth map, and the legend of the Pavia University dataset.
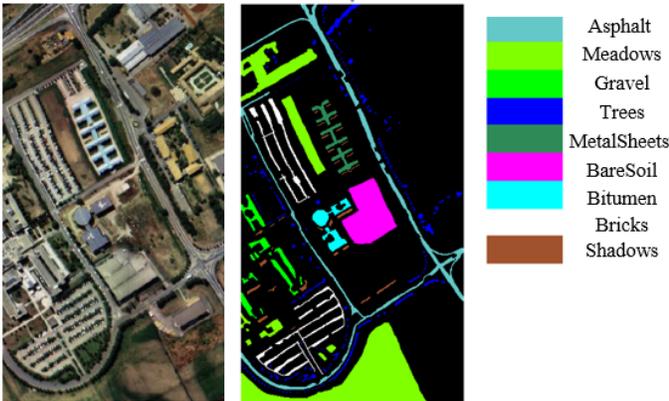


Fig. 5: The false-color composite image, the corresponding ground-truth map, and the legend of the Salinas dataset.

The Salinas (SA) hyperspectral dataset was captured using an airborne sensor over Salinas Valley, California, USA. It consists of 512x217 pixels and 224 spectral bands ranging from 400 to 2500 nm, with a spatial resolution of about 3.7 meters. The dataset includes 16 crop types and has been widely utilized in classification, clustering, and feature extraction. Its availability has advanced hyperspectral imaging and agricultural monitoring research.

Table I shows the size of the training and testing datasets used in this experiment, including the sample distribution of each land cover category.

*2) Evaluation Metrics:* We will compare the effectiveness of our algorithm with other algorithms from four aspects, mainly including Overall Accuracy (OA), Average Accuracy (AA), kappa coefficient ($\kappa$), and the classification accuracy of each land cover category itself.

*3) Training Details:* The PyTorch training framework was utilized to implement and train the model, with the basic hardware environment comprising an AMD EPYC 7543 production-grade CPU, 128GB of memory, and two NVIDIA GeForce RTX 3090 GPUs, each with 24GB of memory.

The diffusion model was optimized using the Adam optimizer, with a learning rate of 1e-4, a batch size of 256, and a patch size of 64. The selected patch size for the experiment is relatively large compared to the overall size of the remote sensing image. In this study, given the available computational resources, it is advisable to maximize the patch size in order to model a larger range of pixel relationships. In order to enhance the fitting performance of diffusion, it is recommended to use a large number of training epochs. For the IP, PU and SA datasets, this study utilized training results with more than 30000 epochs. The classification model was trained using an Adam optimizer, with a learning rate of 1e-3 and a smaller

batch size of 64. Empirically, it was observed that this model converged rapidly, and fewer than 50 epochs were sufficient to achieve the desired classification accuracy for all three datasets.

## B. Performance Analysis

To assess the effectiveness of the proposed approach, various representative algorithms, namely CNN1D, CNN2D, SF, miniGCN, SSRN, SSFTT, DMVL(+SVM), and SSGRN, were selected for the control experiments.

1) CNN1D[38], consists of 5 layers: a 1-D convolutional layer, a batch normalization (BN) layer, a rectified linear unit (ReLU) layer, an average pooling layer, and an output layer. CNN1D is one of the classical CNN algorithms, and our proposed algorithm also includes a CNN structure.

2) CNN2D[83], consists of three 2-D convolutional blocks, each of which contains a 2-D convolutional layer, a BN layer, and a ReLU activation function. Each 2-D convolutional block has 8, 16, and 32 $3 \times 3$ 2-D filters, respectively. CNN2D is one of the classical CNN algorithms, and our proposed algorithm also includes a CNN structure.

3) SF[29], SpectralFormer is one of the classical transformer-based algorithms, and our structure also incorporates a transformer structure.

4) miniGCN[32], a variant of graph convolutional networks (GCN) designed to train large-scale GCNs in a mini-batch fashion, making it more efficient and flexible than traditional GCNs. miniGCN is one of the classical algorithms for constructing sample relationships.

5) SSRN[50], Spectral-Spatial Residual Network, enhances the effectiveness of CNN-based models by fusing and extracting spectral-spatial information through the Spectral-Spatial Residual Network. SSRN represents one of the state-of-the-art CNN-based methods.

6) SSFTT[84], Spectral-Spatial Feature Tokenization Transformer, effectively integrates CNN and Transformers, representing high-level semantic features, and achieving state-of-the-art performance.

7) DMVL(+SVM)[85], Deep Multiview Learning(+SVM) performs unsupervised feature extraction followed by classification using an SVM classifier. Similar to our proposed algorithm, DMVL belongs to the two-stage algorithms, and its classification performance reaches one of the state-of-the-art levels.

8) SSGRN[33], Spectral-Spatial Graph Reasoning Network, employs intermediate features to generate super-pixels, facilitating the creation of homogeneous regions and graph structures. SSGRN is one of the state-of-the-art graph-based algorithms.

9) Ours: Our proposed model follows the basic architecture described earlier, comprising two stages: diffusion model training and classification. During diffusion model training, the input data is divided into $64 \times 64$ patches, which are first passed through a 3D convolution layer with a kernel size of (3,3,3) and padding size

of (1,1,1). Subsequently, three 3D down-sampling and three 3D up-sampling layers are employed for further spatial and spectral feature extraction. Finally, another 3D convolution layer with a kernel size of (3,3,3) and padding size of (1,1,1) will be passed. To extract diffusion features, we performed spatial and spectral reconstruction with noise images and timestamps as inputs. We used the output of U-net down-sampling and up-sampling layers while reducing its dimensionality to serve as diffusion features. In this study, PCA was used. In the classification stage, the diffusion features go through a series of layers, including the Convolutional Tokenization layer, Linear Projection layer, transformer encoders, and MLP layer, to ultimately obtain the classification results. At this stage, the patch size is set to 13.

*1) Quantitative Results:* Tables II-IV respectively presents the classification performance of the IP, PU and SA datasets, including OA, AA, $\kappa$, and the classification accuracy of each category. The best performance is highlighted in bold. It is worth noting that our model exhibits the best overall performance on these three datasets.

Through detailed analysis, it can be observed from the table that traditional algorithms such as CNN1D and CNN2D exhibit relatively poor performance in classification. This is primarily because these models have weak feature extraction capabilities. Especially for HSIs, simple models struggle to effectively capture both spatial and spectral features, particularly when dealing with abundant spectral data. On the other hand, the SpectralFormer and miniGCN models are more complex, but their performance is relatively inferior in limited sample situations.

The SSRN, SSFTT, DMVL and SSGRN algorithms perform relatively well. The SSRN and SSFTT models deeply extract spectral-spatial information, leading to a significant improvement in classification performance. The DMVL algorithm effectively utilizes unsupervised information. Additionally, the SSGRN algorithm enhances model performance through additional modeling of global relationships. Our proposed algorithm enhances the modeling of the spectral-spatial relationship through the diffusion process, surpassing the aforementioned algorithms. Furthermore, it integrates spectral-spatial features more effectively using the transformer architecture, resulting in superior performance. In comparison to the second-best algorithm, our approach achieved an improvement of $0.81\%$, $0.93\%$, and $1.97\%$ in overall accuracy (OA) on the IP, PU and SA datasets, respectively.

From an analysis of the classification performance in different subcategories, it can be concluded that our algorithm demonstrates relatively balanced results across various subcategories in the IP and SA datasets, with higher values of AA and $\kappa$. However, on the PU dataset, our algorithm exhibits relatively weaker performance in achieving balanced classification results among different subcategories compared to models such as SSRN, Nevertheless, the overall AA remains high.

Fig. 9 presents a comprehensive comparison of the overall accuracy among corresponding models using the IP, PU, and

TABLE II
CLASSIFICATION RESULTS OF EXPERIMENTS ON THE INDIAN PINES DATASET

| Class | CNN1D [38] | CNN2D [83] | SF [29] | miniGCN [32] | SSRN [50] | SSFTT [84] | DMVL[85] | SSGRN[33] | Ours |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 93.750 | **100.000** | **100.000** | 12.500 | **100.000** | **100.000** | 68.657 | **100.000** | **100.000** |
| 2 | 58.512 | 63.877 | 49.070 | 65.880 | **95.494** | 80.472 | 92.023 | 86.409 | 85.479 |
| 3 | 57.500 | 58.000 | 71.500 | 42.625 | 84.875 | 89.250 | 89.074 | 92.375 | **94.125** |
| 4 | 78.744 | 79.227 | 96.618 | 60.870 | 98.068 | 99.034 | 80.068 | **100.000** | **100.000** |
| 5 | 86.313 | 80.795 | 77.042 | 84.989 | 79.470 | **99.338** | 96.742 | 95.144 | 96.026 |
| 6 | 90.571 | 89.857 | 85.857 | 96.714 | 93.429 | 98.714 | 85.697 | 98.143 | **99.857** |
| 7 | 92.857 | **100.000** | 85.714 | 92.857 | 92.857 | **100.000** | 31.111 | **100.000** | **100.000** |
| 8 | 87.946 | 94.643 | 99.107 | 98.214 | **100.000** | **100.000** | 97.951 | **100.000** | **100.000** |
| 9 | 80.000 | 70.000 | **100.000** | 60.000 | **100.000** | **100.000** | 74.074 | **100.000** | **100.000** |
| 10 | 54.246 | 69.958 | 80.149 | 63.376 | **96.178** | 91.720 | 88.447 | 86.837 | 84.501 |
| 11 | 51.093 | 46.887 | 52.371 | 37.237 | 81.402 | 88.454 | **97.728** | 90.680 | 91.959 |
| 12 | 69.805 | 58.792 | 52.753 | 58.792 | 91.652 | 85.968 | 84.301 | 95.204 | **95.560** |
| 13 | 97.143 | **100.000** | 98.857 | 97.714 | **100.000** | **100.000** | 93.488 | **100.000** | **100.000** |
| 14 | 76.923 | 81.377 | 79.109 | 89.555 | **99.757** | 99.028 | 97.709 | 92.793 | 97.085 |
| 15 | 71.910 | 63.483 | 62.079 | 53.652 | 85.674 | 98.596 | 84.171 | **100.000** | **100.000** |
| 16 | 90.476 | 88.889 | **100.000** | 90.476 | **100.000** | **100.000** | 39.732 | **100.000** | 98.413 |
| OA(%) | 66.007 | 66.966 | 67.782 | 63.916 | 90.658 | 91.566 | 90.487 | 92.336 | **93.146** |
| AA(%) | 77.362 | 77.861 | 80.639 | 69.091 | 93.678 | 95.661 | 81.311 | 96.099 | **96.438** |
| $\kappa \times 100$ | 61.658 | 62.807 | 63.820 | 59.439 | 89.386 | 90.392 | 89.216 | 91.254 | **92.175** |

TABLE III
CLASSIFICATION RESULTS OF EXPERIMENTS ON THE PAVIA UNIVERSITY DATASET

| Class | CNN1D [38] | CNN2D [83] | SF [29] | miniGCN [32] | SSRN [50] | SSFTT [84] | DMVL[85] | SSGRN[33] | Ours |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 71.156 | 66.444 | 61.521 | 60.627 | **91.668** | 86.972 | 76.361 | 82.397 | 86.154 |
| 2 | 85.273 | 84.758 | 87.975 | 92.207 | 91.664 | 97.728 | 96.995 | 96.026 | **98.394** |
| 3 | 76.027 | 73.659 | 87.240 | 79.894 | 94.490 | 92.943 | 90.013 | 93.717 | **96.568** |
| 4 | 77.192 | 80.488 | 94.100 | 82.367 | 91.595 | 91.760 | 66.530 | **98.484** | 76.401 |
| 5 | 99.468 | 99.163 | **100.000** | 99.163 | 98.327 | 99.924 | 86.293 | 99.468 | 99.544 |
| 6 | 88.138 | 86.437 | 65.573 | 65.613 | 99.480 | 98.720 | 97.113 | **100.000** | 98.300 |
| 7 | 90.308 | 86.462 | 90.000 | 96.154 | 99.846 | 99.000 | 94.895 | **99.923** | 99.077 |
| 8 | 79.107 | 76.999 | 78.614 | 80.312 | 96.166 | 74.535 | 97.243 | 87.431 | **98.987** |
| 9 | **100.000** | 99.891 | 99.455 | **100.000** | 95.202 | 99.455 | 69.598 | 92.475 | 91.167 |
| OA(%) | 82.772 | 81.424 | 81.511 | 82.355 | 93.636 | 93.667 | 89.768 | 93.850 | **94.775** |
| AA(%) | 85.185 | 83.811 | 84.942 | 84.037 | **95.382** | 93.448 | 86.116 | 94.436 | 93.843 |
| $\kappa \times 100$ | 77.606 | 75.880 | 75.795 | 76.620 | 91.706 | 91.617 | 86.580 | 91.919 | **93.061** |

SA datasets with varying sample sizes. The figure clearly demonstrates that increasing the sample size leads to a continuous improvement in the overall classification accuracy for each model. Remarkably, our proposed model consistently outperforms other models across a wide range of sample sizes. When the number of samples per-class drops to 10, our model still demonstrates the best performance on the PU dataset, while it does not exhibit significant advantages on the IP and SA datasets. We consider this may be attributed to the larger number of feature channels extracted through the diffusion process, resulting in a certain degree of overfitting in the classification model.

*2) Qualitative results:* Figs. 6-8 display the classification results of IP, PU and SA datasets. Visually, our model exhibits lower noise levels and closely aligns with the Ground Truth. Traditional classification algorithms such as CNN1D, CNN2D, SF and miniGCN produce noisy classification maps with discontinuous land cover blocks and rough classification

results. Meanwhile, SSRN, SSFTT, DMVL and SSGRN algorithms show improved performance, with a reduction in noisy points. Furthermore, from a detailed perspective, our proposed algorithm exhibits better overall segmentation of land cover.

*C. Model Analysis*

*1) Ablation study:* We further analyzed the benefits of the proposed SpectraiDiff, which involves feeding the original spectral features and the features extracted from the pre-training process with the diffusion model into our proposed attention-based classification model. As shown in Table V, the classification results on IP, PU and SA datasets are presented. The results demonstrate that using diffusion features as input significantly outperforms the use of raw features on all three datasets, resulting in improved OA, AA, and $\kappa$ metrics. It is worth noting that using diffusion features as input, compared to using raw features as input, results in an overall classification accuracy improvement (OA) of 0.78%,

TABLE IV

CLASSIFICATION RESULTS OF EXPERIMENTS ON THE SALINAS DATASET

| Class | CNN1D [38] | CNN2D [83] | SF [29] | miniGCN [32] | SSRN [50] | SSFTT [84] | DMVL[85] | SSGRN[33] | Ours |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 99.596 | 99.747 | 99.192 | 99.141 | **100.000** | **100.000** | 96.958 | 92.976 | **100.000** |
| 2 | 99.729 | 99.729 | 99.513 | 99.946 | 95.996 | **100.000** | 99.946 | 95.860 | **100.000** |
| 3 | 97.431 | 98.972 | 84.275 | 97.996 | 81.449 | **100.000** | 99.697 | 98.561 | **100.000** |
| 4 | 99.267 | 98.974 | 98.754 | 99.560 | **100.000** | 99.927 | 66.098 | 93.109 | **100.000** |
| 5 | 97.432 | 99.131 | 92.674 | 98.678 | 91.843 | 98.263 | **99.875** | 95.582 | 98.036 |
| 6 | 99.440 | 99.211 | 99.491 | 99.822 | 99.466 | 99.949 | 98.555 | 99.746 | **100.000** |
| 7 | 99.606 | 99.436 | 99.127 | 99.577 | 99.944 | 99.972 | **100.000** | 98.929 | 99.239 |
| 8 | 84.610 | 65.599 | 68.268 | 77.493 | 86.932 | 82.929 | 96.086 | **97.545** | 96.975 |
| 9 | 98.753 | 99.141 | 97.732 | 99.887 | **100.000** | **100.000** | **100.000** | 99.417 | **100.000** |
| 10 | 89.070 | 85.345 | 88.208 | 92.857 | 96.121 | 98.060 | **100.000** | 88.762 | 98.615 |
| 11 | 98.555 | 90.944 | 96.435 | 95.857 | **100.000** | **100.000** | 96.703 | 97.399 | 99.807 |
| 12 | 99.895 | 99.736 | 99.895 | 99.789 | **100.000** | **100.000** | 99.353 | 99.051 | 99.947 |
| 13 | 98.984 | 98.307 | 99.887 | 98.984 | **100.000** | 97.178 | 94.731 | 92.099 | **100.000** |
| 14 | 91.250 | 96.058 | 97.596 | 96.250 | 99.615 | 99.519 | 86.829 | 99.231 | **100.000** |
| 15 | 64.424 | 72.271 | 81.044 | 53.868 | 90.564 | 97.513 | **99.090** | 93.396 | 98.826 |
| 16 | 97.805 | 93.697 | 90.096 | 99.043 | **100.000** | 99.043 | 99.724 | **100.000** | **100.000** |
| OA(%) | 90.540 | 87.338 | 88.248 | 88.181 | 94.352 | 95.789 | 97.005 | 96.539 | **98.971** |
| AA(%) | 94.740 | 93.519 | 93.262 | 94.297 | 96.371 | 98.272 | 95.853 | 96.354 | **99.465** |
| $\kappa \times 100$ | 89.451 | 85.931 | 86.973 | 86.823 | 93.722 | 95.322 | 96.668 | 96.144 | **98.854** |

5.87%, and 3.17% on the IP, PU, and SA datasets, respectively. Fig. 10 demonstrates the improvement in performance under different sample sizes. Within the framework of the proposed Attention-based Classification Module, the variation in the number of training samples per class from 20 to 80 reveals a consistent trend: OA achieved with diffusion features as input consistently surpasses that obtained with raw features. However, when the sample size decreases to 10 per class, the inclusion of diffusion features does not yield any performance improvement on both the IP and PU datasets. This occurrence can be attributed to the potential overfitting of the classification model, which arises due to the high dimensionality of the generated diffusion features.

TABLE V

CLASSIFICATION PERFORMANCE ANALYSIS BETWEEN DIFFERENT FEATURE INPUT ON THE INDIAN PINES DATASET, PAVIA UNIVERSITY DATASET AND SALINAS DATASET

| Dataset | Feature Input | Metrics | | |
|---|---|---|---|---|
| | | OA(%) | AA(%) | $\kappa$*100 |
| IP | raw features | 92.37 | 95.61 | 91.27 |
| | diffusion features | 93.15 | 96.44 | 92.17 |
| PU | raw features | 88.91 | 89.77 | 85.51 |
| | diffusion features | 94.78 | 93.84 | 93.06 |
| SA | raw features | 95.80 | 98.26 | 95.34 |
| | diffusion features | 98.97 | 99.47 | 98.85 |

*2) Diffusion Model Analysis:* We further analyzed the effects of the diffusion model on the final classification model. Firstly, to verify the diffusion model effectiveness, we used the well-trained diffusion model to recover and reconstruct the spectral curves of the hyperspectral data. we utilized HIS images corrupted with Gaussian noise as the input, and performed the Reverse Spectral-Spatial Diffusion Process step by step for each timestamp. As an example, Fig. 11 presents the restoration effects corresponding to different timestamps on the Indian Pines dataset using a false-color image. From a visual perspective, the diffusion model basically reconstructs the original remote sensing image content. Additionally, Fig. 12 shows the reconstruction process for spectral curve of one land-cover type. As the timestamps change, the spectral curve gradually reconstructs to the original shape of the land-cover types from a state similar to white noise. This indicates that the diffusion model has embedded the spectral curve information into the model parameters, providing a data foundation for using the diffusion feature for land-cover classification.

When extracting features using the Diffusion model, there are two crucial influencing parameters to consider, namely the Timestamp and the Layerindex. Timestamp refers to the number of denoising steps the Diffusion model takes to restore noisy images additionally. Layerindex refers to the location of the U-Net output used as a feature layer in the Diffusion model. We have conducted classification experiments on various Timestamp and Layerindex values, and the results are presented in the Table VI. In the case of the IP dataset, there are some fluctuations in classification performance for different Timestamp and Layerindex values, but no significant changes. We postulate that this may be attributed to the relatively small size of the IP dataset. However, for the PU and SA datasets, there is a certain correlation between classification performance and Timestamp/Layerindex. When considering the Timestamp dimension, a decreasing trend in classification performance is observed when using features with larger Timestamp, and the optimal performance generally occurs in smaller Timestamp groups (Timestamp = 5). We believe that when the Timestamp is larger, the number of iterations for denoising using the Diffusion model is lower, leading to
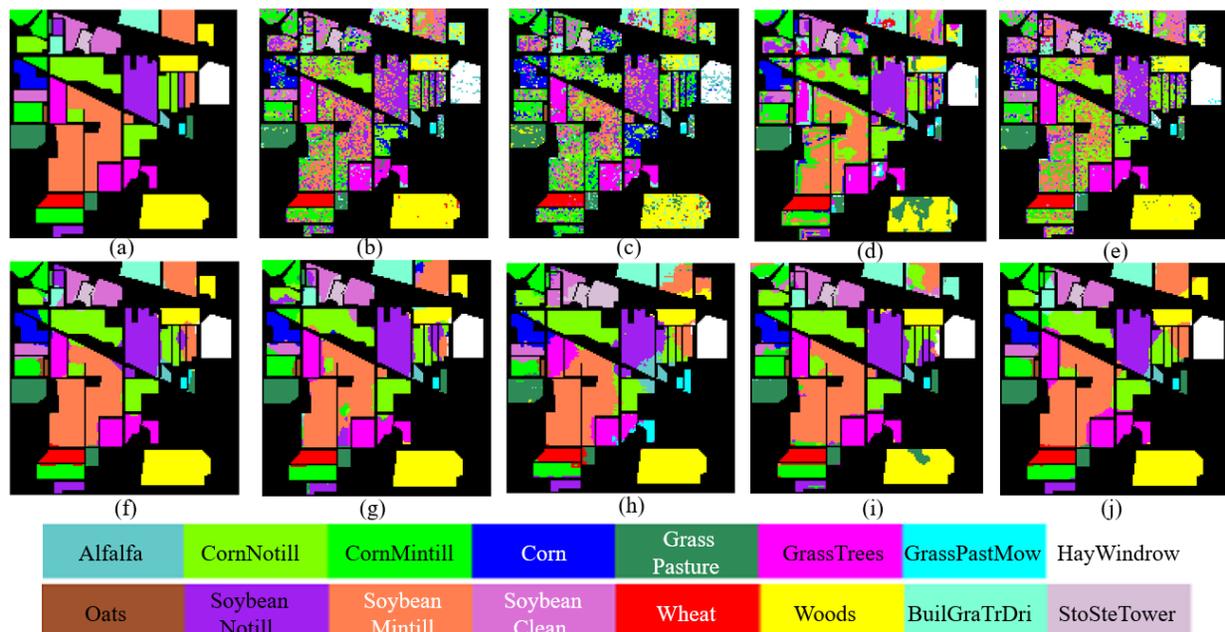
Fig. 6: Classification maps of the Indian Pines dataset. (a) Ground-truth map. (b) CNN1D (OA = 66.007%). (c) CNN2D (OA = 66.966%). (d) SF (OA = 67.782%). (e) miniGCN (OA = 63.916%). (f) SSRN (OA = 90.658%). (g) SSFTT (OA = 91.566%). (h) DMVL (OA = 90.487%). (i) SSGRN (OA = 92.336%). (j) Ours (OA = 93.146%).

TABLE VI

THE PERFORMANCE OF DIFFERENT LAYERINDEX AND TIMESTAMP IN THE INDIAN PINES DATASET, THE PAVIA UNIVERSITY DATASET, AND SALINAS DATASET

| LayerIndex | Timestamp | Indian Pines | | | Pavia University | | | Salinas | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | OA(%) | AA(%) | $\kappa*100$ | OA(%) | AA(%) | $\kappa*100$ | OA(%) | AA(%) | $\kappa*100$ |
| 0 | 5 | 92.73 | 95.96 | 91.69 | **94.78** | 93.84 | **93.06** | **98.97** | **99.47** | **98.85** |
| | 10 | 92.19 | 96.08 | 91.08 | 93.60 | 93.14 | 91.57 | 98.64 | 99.35 | 98.49 |
| | 50 | 92.67 | 96.10 | 91.62 | 93.48 | 91.79 | 91.37 | 97.73 | 99.05 | 97.47 |
| | 100 | 92.43 | 95.93 | 91.35 | 92.86 | 92.04 | 90.53 | 97.80 | 98.87 | 97.56 |
| | 200 | 91.39 | 95.56 | 90.18 | 91.62 | 90.84 | 88.96 | 97.70 | 98.72 | 97.44 |
| 1 | 5 | 92.97 | 96.43 | 91.98 | 94.31 | **94.19** | 92.53 | 98.69 | 99.22 | 98.54 |
| | 10 | 92.41 | 96.11 | 91.34 | 94.72 | 93.83 | 93.00 | 98.32 | 99.20 | 98.13 |
| | 50 | 91.90 | 95.92 | 90.75 | 92.73 | 93.03 | 90.46 | 97.99 | 98.94 | 97.76 |
| | 100 | 92.12 | 95.80 | 90.99 | 94.12 | 92.79 | 92.19 | 97.98 | 99.08 | 97.75 |
| | 200 | 92.02 | 95.89 | 90.89 | 91.84 | 90.55 | 89.27 | 96.73 | 98.42 | 96.36 |
| 2 | 5 | 93.15 | **96.44** | **92.17** | 92.02 | 91.65 | 89.52 | 98.02 | 99.03 | 97.80 |
| | 10 | 92.85 | 96.39 | 91.84 | 91.26 | 91.02 | 88.55 | 97.74 | 98.77 | 97.49 |
| | 50 | 92.93 | 96.21 | 91.93 | 91.10 | 90.66 | 88.23 | 95.51 | 97.62 | 95.01 |
| | 100 | **93.16** | 96.32 | **92.17** | 86.51 | 85.12 | 82.25 | 95.96 | 97.47 | 95.51 |
| | 200 | 92.60 | 96.11 | 91.55 | 81.54 | 81.72 | 75.88 | 93.53 | 96.37 | 92.80 |

1 LayerIndex=0, 1, 2 respectively represent the inputs of the three up-sampling layers in the U-Net model, with larger numbers closer to the output layer.
2 A smaller Timestamp indicates that the diffusion denoising process is closer to the original image position.
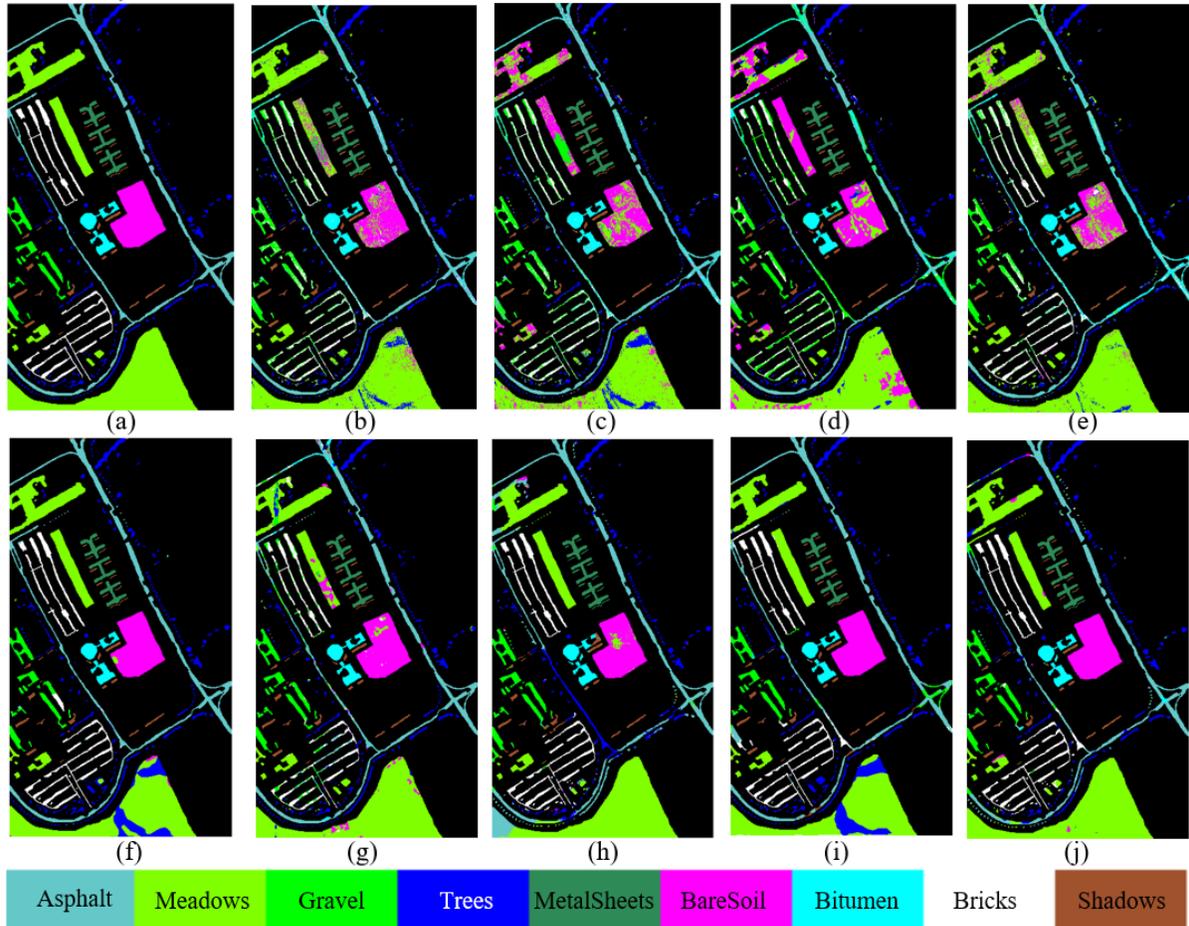
Fig. 7: Classification maps of the Pavia University dataset. (a) Ground-truth map. (b) CNN1D (OA = 82.772%). (c) CNN2D (OA = 81.424%). (d) SF (OA = 81.511%). (e) miniGCN (OA = 82.355%). (f) SSRN (OA = 93.636%). (g) SSFTT (OA = 93.667%). (h) DMVL (OA = 89.768%). (i) SSGRN (OA = 93.850%). (j) Ours (OA = 94.775%).
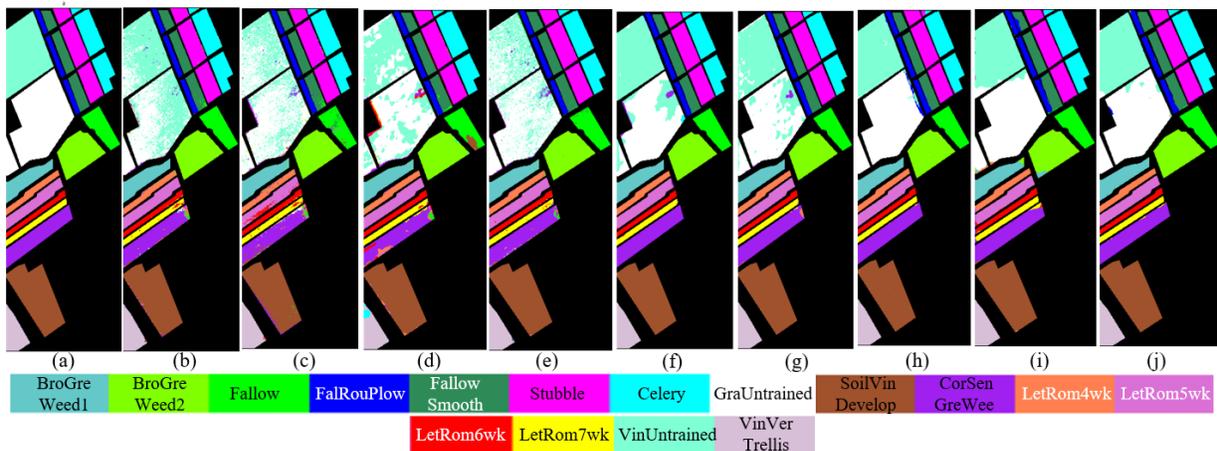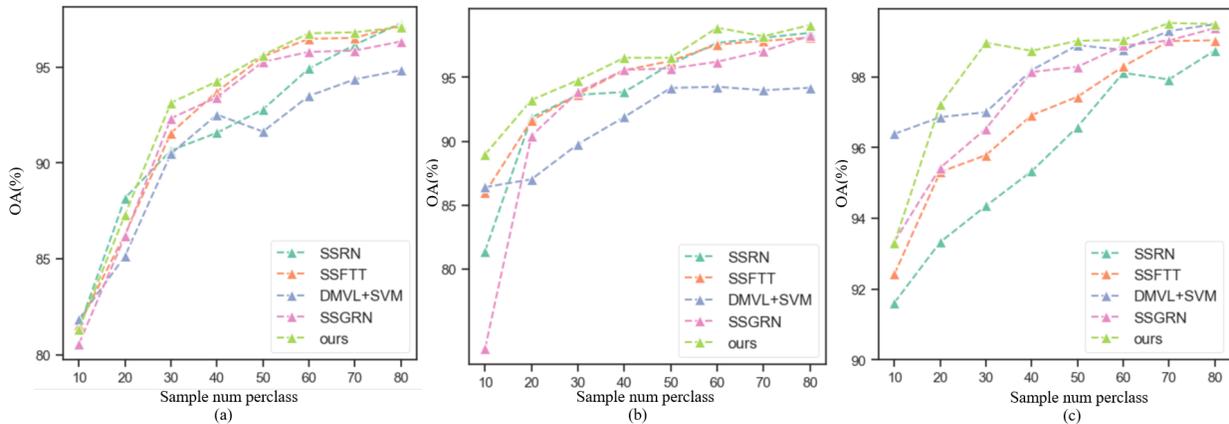


Fig. 8: Classification maps of the Salinas dataset. (a) Ground-truth map. (b) CNN1D (OA = 90.540%). (c) CNN2D (OA = 87.338%). (d) SF (OA = 88.248%). (e) miniGCN (OA = 88.181%). (f) SSRN (OA = 94.352%). (g) SSFTT (OA = 95.789%). (h) DMVL (OA = 97.005%). (i) SSGRN (OA = 96.539%). (j) Ours (OA = 98.971%).

Fig. 9: Evolution of OA as a function of number of training samples per class. (a) Indian Pines. (b) University of Pavia.(c) Salinas.
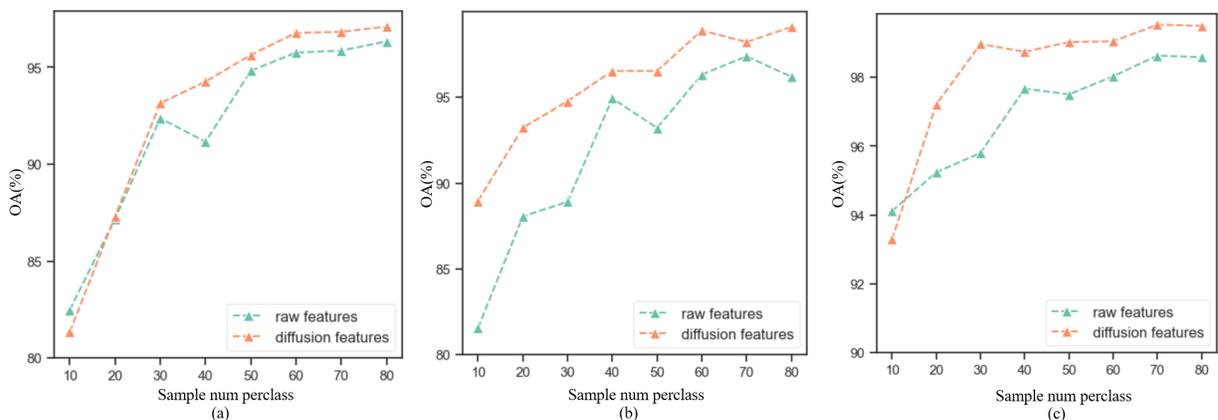


Fig. 10: Overall Accuracy with Varying Sample Sizes: Raw Features Inputs vs. Diffusion Features Inputs. (a) Indian Pines. (b) University of Pavia.(c) Salinas.

TABLE VII
GPU INFERENCE TIME (S) OF THE PROPOSED METHOD
AND THE COMPARISON METHODS

| Method | IP | PU | SA |
|---|---|---|---|
| CNN1D | 9.74 | 37.94 | 20.01 |
| CNN2D | 14.65 | 45.84 | 24.74 |
| SF | 29.81 | 163.34 | 141.30 |
| miniGCN | 0.95 | 2.40 | 2.13 |
| SSRN | 48.02 | 248.06 | 222.80 |
| SSFTT | 31.58 | 110.08 | 55.96 |
| SSGRN | 1.13 | 4.99 | 2.98 |
| Ours | 48.53 | 256.23 | 242.16 |

relatively more noise information in the input and resulting in a deviation in the classification performance. Considering the Layerindex dimension, both datasets (PU & SA) showed better performance at Layerindex 0 than at Layerindex 1 and 2.

We compared the inference time between different algorithms. It is noteworthy that our algorithm, being a two-stage algorithm, only accounts for the inference time at the classification stage in this analysis. Table. VII shows the testing time for each algorithm. Our algorithm takes longer compared to simple CNN algorithms such as CNN1D, CNN2D, and GCN algorithms. However, the increase in time is not significant when compared to complex CNN algorithms like SSRN and transformer algorithms, as they are still at the same level.

## V. CONCLUSION

In this study, a novel approach is proposed for constructing the spectral-spatial distribution of HSI data from a generative perspective and capturing the spectral-spatial features. The proposed method provides a unique viewpoint for the spectral-spatial diffusion process and plays a critical role in establishing relationships between samples. With the proposed SpectralDiff, sample relationships can be adaptively constructed without prior knowledge of graph structure or neighborhood information. This approach captures the data distribution and contextual information of objects in HSI, achieving cross-sample perception. Experimental results demonstrate that this method outperforms state-of-the-art techniques.

Looking towards the future, an exciting avenue of exploration lies in studying the potential of diffusion models for out-
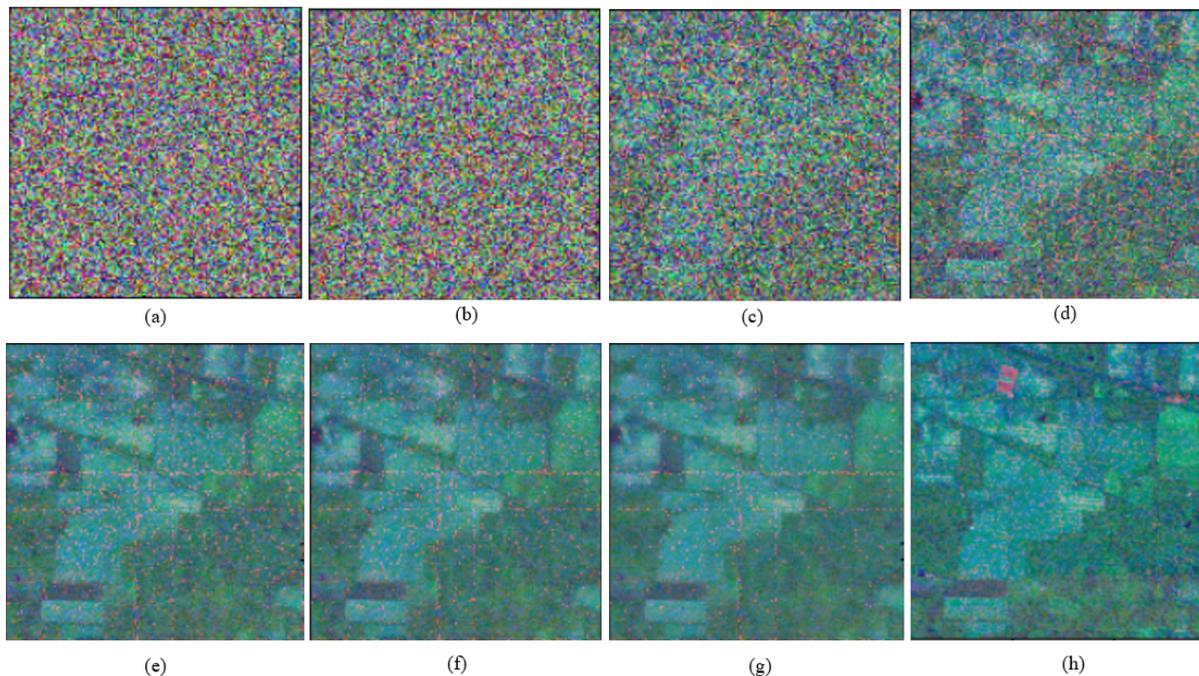
Fig. 11: False-color images of the reconstructed IndianPines dataset corresponding to different timestamps by Reverse Spectral-Spatial Diffusion Process. (a) $t = 400$. (b) $t = 200$. (c) $t = 100$. (d) $t = 50$. (e) $t = 10$. (f) $t = 5$. (g) $t = 0$. (h) ground truth.
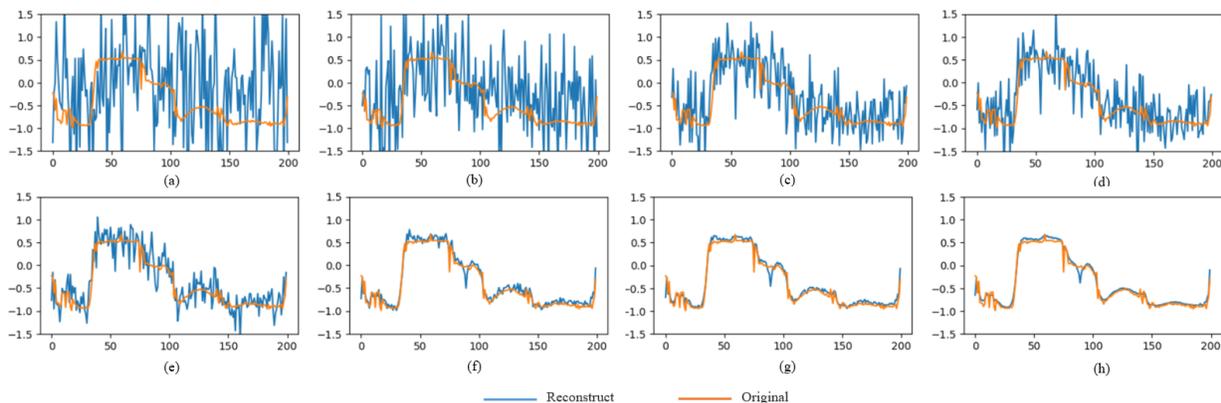


Fig. 12: The spectral curves of the reconstructed Indian Pines dataset corresponding to different timestamps by Reverse Spectral-Spatial Diffusion Process, class=woods, x-axis:spectral band number, y-axis: normalized spectral values. (a) $t = 400$. (b) $t = 200$. (c) $t = 100$. (d) $t = 80$. (e) $t = 50$. (f) $t = 10$. (g) $t = 5$. (h) $t = 0$.

of-distribution generalization and detection in the context of hyperspectral imaging, building upon the generative paradigm. It is expected that diffusion models will continue to advance these areas, capturing underlying data manifolds through diffusion processes, thereby learning to generalize well to unseen examples lying outside the training distribution. Additionally, diffusion models demonstrate strong detection performance in identifying out-of-distribution samples. Going forward, there is significant potential for diffusion models to further contribute to the fields of out-of-distribution generalization and detection. With further research, developments in leveraging the power of diffusion models to analyze complex and high-dimensional hyperspectral data are expected to continue, leading to exciting opportunities for future applications in diverse areas.

REFERENCES

[1] S. Li, W. Song, L. Fang, Y. Chen, P. Ghamisi, and J. A. Benediktsson, "Deep learning for hyperspectral image classification: An overview," IEEE Trans. Geosci. Remote Sens., vol. 57, no. 9, pp. 6690–6709, 2019.
[2] L. Zhang, L. Zhang, and B. Du, "Deep learning for remote sensing data: A technical tutorial on the state of the art," IEEE Geosci. Remote Sens. Mag., vol. 4, no. 2, pp. 22–40, 2016.
[3] X. X. Zhu, D. Tuia, L. Mou, G.-S. Xia, L. Zhang, F. Xu, and F. Fraundorfer, "Deep learning in remote sensing: A comprehensive review and list of resources," IEEE Geosci. Remote Sens. Mag., vol. 5, no. 4, pp. 8–36, 2017.
[4] P. Ghamisi, J. Plaza, Y. Chen, J. Li, and A. J. Plaza, "Advanced spectral classifiers for hyperspectral images: A review," IEEE Geosci. Remote Sens. Mag., vol. 5, no. 1, pp. 8–32, 2017.
[5] S. Prasad and L. M. Bruce, "Limitations of principal components analysis for hyperspectral target recognition," IEEE Geosci. Remote Sens. Lett., vol. 5, no. 4, pp. 625–629, 2008.

[6] M.-D. Yang, K.-S. Huang, Y. F. Yang, L.-Y. Lu, Z.-Y. Feng, and H. P. Tsai, "Hyperspectral image classification using fast and adaptive bidimensional empirical mode decomposition with minimum noise fraction," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 12, pp. 1950–1954, 2016.

[7] Y. Zhai, L. Zhang, N. Wang, Y. Guo, Y. Cen, T. Wu, and Q. Tong, "A modified locality-preserving projection approach for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 8, pp. 1059–1063, 2016.

[8] S. D. Fabiyi, P. Murray, J. Zabalza, and J. Ren, "Folded lda: Extending the linear discriminant analysis algorithm for feature extraction and data reduction in hyperspectral remote sensing," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 14, pp. 12 312–12 331, 2021.

[9] J. Wang and C.-I. Chang, "Independent component analysis-based dimensionality reduction with applications in hyperspectral image analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 6, pp. 1586–1600, 2006.

[10] L. Xie, M. Yin, X. Yin, Y. Liu, and G. Yin, "Low-rank sparse preserving projections for dimensionality reduction," *IEEE Trans. Image Process.*, vol. 27, no. 11, pp. 5261–5274, 2018.

[11] A. Plaza, P. Martinez, J. Plaza, and R. Perez, "Dimensionality reduction and classification of hyperspectral image data using sequences of extended morphological transformations," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 3, pp. 466–479, 2005.

[12] P. R. Marpu, M. Pedergnana, M. Dalla Mura, J. A. Benediktsson, and L. Bruzzone, "Automatic generation of standard deviation attribute profiles for spectral–spatial classification of remote sensing data," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 2, pp. 293–297, 2013.

[13] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *science*, vol. 313, no. 5786, pp. 504–507, 2006.

[14] J. Yue, L. Fang, P. Ghamisi, W. Xie, J. Li, J. Chanussot, and A. J. Plaza, "Optical remote sensing image understanding with weak supervision: Concepts, methods, and perspectives," *IEEE Geosci. Remote Sens. Mag.*, pp. 2–21, 2022.

[15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, June 2016.

[16] S. Wang, J. Yue, J. Liu, Q. Tian, and M. Wang, "Large-scale few-shot learning via multi-modal knowledge discovery," in *ECCV*. Springer, 2020, pp. 718–734.

[17] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, 2018.

[18] Z. Tian, H. Zhao, M. Shu, Z. Yang, R. Li, and J. Jia, "Prior guided feature enrichment network for few-shot segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 2, pp. 1050–1065, 2022.

[19] L. Wu, L. Fang, X. He, M. He, J. Ma, and Z. Zhong, "Querying labeled for unlabeled: Cross-image semantic consistency guided semi-supervised semantic segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 7, pp. 8827–8844, 2023.

[20] L. Fang, Y. Jiang, Y. Yan, J. Yue, and Y. Deng, "Hyperspectral image instance segmentation using spectral–spatial feature pyramid network," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–13, 2023.

[21] L. Fang, Y. Yan, J. Yue, and Y. Deng, "Toward the vectorization of hyperspectral imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–14, 2023.

[22] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 142–158, 2016.

[23] 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, 2017.

[24] L. Fang, D. Zhu, J. Yue, B. Zhang, and M. He, "Geometric-spectral reconstruction learning for multi-source open-set classification with hyperspectral and lidar data," *IEEE/CAA J. Autom. Sin.*, vol. 9, no. 10, pp. 1892–1895, 2022.

[25] J. Yue, D. Zhu, L. Fang, P. Ghamisi, and Y. Wang, "Adaptive spatial pyramid constraint for hyperspectral image classification with limited training samples," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–14, 2021.

[26] Y. Chen, Z. Lin, X. Zhao, G. Wang, and Y. Gu, "Deep learning-based classification of hyperspectral data," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 7, no. 6, pp. 2094–2107, 2014.

[27] C. Zhang, J. Yue, and Q. Qin, "Global prototypical network for few-shot hyperspectral image classification," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 13, pp. 4748–4759, 2020.

[28] J. Yue, S. Mao, and M. Li, "A deep learning framework for hyperspectral image classification using spatial pyramid pooling," *Remote Sens. Lett.*, vol. 7, no. 9, pp. 875–884, 2016.

[29] D. Hong, Z. Han, J. Yao, L. Gao, B. Zhang, A. Plaza, and J. Chanussot, "Spectralformer: Rethinking hyperspectral image classification with transformers," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022.

[30] A. Qin, Z. Shang, J. Tian, Y. Wang, T. Zhang, and Y. Y. Tang, "Spectral–spatial graph convolutional networks for semisupervised hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 16, no. 2, pp. 241–245, 2019.

[31] S. Wan, C. Gong, P. Zhong, S. Pan, G. Li, and J. Yang, "Hyperspectral image classification with context-aware dynamic graph convolutional network," *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 1, pp. 597–612, 2021.

[32] 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, 2021.

[33] D. Wang, B. Du, and L. Zhang, "Spectral-spatial global graph reasoning for hyperspectral image classification," *IEEE Trans. Neural Networks Learn. Syst.*, pp. 1–14, 2023.

[34] M. Fauvel, Y. Tarabalka, J. A. Benediktsson, J. Chanussot, and J. C. Tilton, "Advances in spectral-spatial classification of hyperspectral images," *Proc. IEEE*, vol. 101, no. 3, pp. 652–675, 2013.

[35] S. Liu, Q. Shi, and L. Zhang, "Few-shot hyperspectral image classification with unknown classes using multitask deep learning," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–18, 2020.

[36] J. Yue, L. Fang, H. Rahmani, and P. Ghamisi, "Self-supervised learning with adaptive distillation for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–13, 2021.

[37] J. Yue, L. Fang, and M. He, "Spectral-spatial latent reconstruction for open-set hyperspectral image classification," *IEEE Trans. Image Process.*, vol. 31, pp. 5227–5241, 2022.

[38] B. Rasti, D. Hong, R. Hang, P. Ghamisi, X. Kang, J. Chanussot, and J. A. Benediktsson, "Feature extraction for hyperspectral imagery: The evolution from shallow to deep: Overview and toolbox," *IEEE Geosci. Remote Sens. Mag.*, vol. 8, no. 4, pp. 60–88, 2020.

[39] W. Li, S. Prasad, J. E. Fowler, and L. M. Bruce, "Locality-preserving discriminant analysis in kernel-induced feature spaces for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 8, no. 5, pp. 894–898, 2011.

[40] H. Huang, G. Shi, H. He, Y. Duan, and F. Luo, "Dimensionality reduction of hyperspectral imagery based on spatial–spectral manifold learning," *IEEE Trans. Cybern.*, vol. 50, no. 6, pp. 2604–2616, 2020.

[41] N. H. Ly, Q. Du, and J. E. Fowler, "Sparse graph-based discriminant analysis for hyperspectral imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 7, pp. 3872–3884, 2014.

[42] C. Bachmann, T. Ainsworth, and R. Fusina, "Improved manifold coordinate representations of large-scale hyperspectral scenes," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 10, pp. 2786–2803, 2006.

[43] F. Tsai and J.-S. Lai, "Feature extraction of hyperspectral image cubes using three-dimensional gray-level cooccurrence," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 6, pp. 3504–3513, 2013.

[44] L. Shen, Z. Zhu, S. Jia, J. Zhu, and Y. Sun, "Discriminative gabor feature selection for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 1, pp. 29–33, 2013.

[45] Y. Zhou, J. Peng, and C. L. P. Chen, "Dimension reduction using spatial and spectral regularized local discriminant embedding for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 2, pp. 1082–1095, 2015.

[46] D. Hong, N. Yokoya, J. Chanussot, J. Xu, and X. X. Zhu, "Joint and progressive subspace analysis (jpsa) with spatial–spectral manifold alignment for semisupervised hyperspectral dimensionality reduction," *IEEE Trans. Cybern.*, vol. 51, no. 7, pp. 3602–3615, 2021.

[47] A. Kianisarkaleh and H. Ghassemian, "Spatial-spectral locality preserving projection for hyperspectral image classification with limited training samples," *Int. J. Remote Sens.*, vol. 37, no. 21, pp. 5045–5059, 2016.

[48] J. Yue, W. Zhao, S. Mao, and H. Liu, "Spectral-spatial classification of hyperspectral images using deep convolutional neural networks," *Remote Sens. Lett.*, vol. 6, no. 6, pp. 468–477, 2015.

[49] L. Mou, P. Ghamisi, and X. X. Zhu, "Deep recurrent neural networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 7, pp. 3639–3655, 2017.

[50] 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, 2018.

[51] 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, 2019.

[52] M. E. Paoletti, J. M. Haut, R. Fernandez-Beltran, J. Plaza, A. Plaza, J. Li, and F. Pla, "Capsule networks for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 4, pp. 2145–2160, 2019.

[53] A. Q. Nichol and P. Dhariwal, "Improved denoising diffusion probabilistic models," in *ICML*, vol. 139, 2021, pp. 8162–8171.

[54] D. Kingma, T. Salimans, B. Poole, and J. Ho, "Variational diffusion models," in *NeurIPS*, vol. 34, 2021, pp. 21 696–21 707.

[55] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *ICML*. PMLR, 2015, pp. 2256–2265.

[56] F.-A. Croitoru, V. Hondru, R. T. Ionescu, and M. Shah, "Diffusion models in vision: A survey," *arXiv preprint arXiv:2209.04747*, 2022.

[57] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, J. Ho, D. J. Fleet, and M. Norouzi, "Photorealistic text-to-image diffusion models with deep language understanding," in *NeurIPS*, vol. 35, 2022, pp. 36 479–36 494.

[58] J. Richter, S. Welker, J.-M. Lemercier, B. Lay, and T. Gerkmann, "Speech enhancement and dereverberation with diffusion-based generative models," *IEEE/ACM Trans. Audio, Speech, Language Process.*, pp. 1–13, 2023.

[59] J. M. L. Alcaraz and N. Strodthoff, "Diffusion-based time series imputation and forecasting with structured state space models," *arXiv preprint arXiv:2208.09399*, 2022.

[60] K. Rasul, C. Seward, I. Schuster, and R. Vollgraf, "Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting," in *ICML*. PMLR, 2021, pp. 8857–8868.

[61] M. Zhang, M. Qamar, T. Kang, Y. Jung, C. Zhang, S.-H. Bae, and C. Zhang, "A survey on graph diffusion models: Generative ai in science for molecule, protein and material," *arXiv preprint arXiv:2304.01565*, 2023.

[62] B. Jing, G. Corso, J. Chang, R. Barzilay, and T. Jaakkola, "Torsional diffusion for molecular conformer generation," *arXiv preprint arXiv:2206.01729*, 2022.

[63] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, Y. Shao, W. Zhang, B. Cui, and M.-H. Yang, "Diffusion models: A comprehensive survey of methods and applications," *arXiv preprint arXiv:2209.00796*, 2022.

[64] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," *NeurIPS*, vol. 32, 2019.

[65] A. Vahdat, K. Kreis, and J. Kautz, "Score-based generative modeling in latent space," in *NeurIPS*, vol. 34. Curran Associates, Inc., 2021, pp. 11 287–11 302.

[66] Y. Song and S. Ermon, "Improved techniques for training score-based generative models," in *NeurIPS*, vol. 33. Curran Associates, Inc., 2020, pp. 12 438–12 448.

[67] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," in *ICLR*, 2021.

[68] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *NeurIPS*, vol. 33. Curran Associates, Inc., 2020, pp. 6840–6851.

[69] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," in *NeurIPS*, vol. 34, 2021, pp. 8780–8794.

[70] J. Yue, L. Fang, S. Xia, Y. Deng, and J. Ma, "Dif-fusion: Towards high color fidelity in infrared and visible image fusion with diffusion models," *arXiv preprint arXiv:2301.08072*, 2023.

[71] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool, "Repaint: Inpainting using denoising diffusion probabilistic models," in *CVPR*, 2022, pp. 11 461–11 471.

[72] C. Saharia, W. Chan, H. Chang, C. Lee, J. Ho, T. Salimans, D. Fleet, and M. Norouzi, "Palette: Image-to-image diffusion models," in *ACM SIGGRAPH*, 2022.

[73] M. Zhao, F. Bao, C. Li, and J. Zhu, "EGSDE: Unpaired image-to-image translation via energy-guided stochastic differential equations," in *NeurIPS*, 2022.

[74] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, "Image super-resolution via iterative refinement," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 4, pp. 4713–4726, 2023.

[75] M. Daniels, T. Maunu, and P. Hand, "Score-based generative neural networks for large-scale optimal transport," in *NeurIPS*, vol. 34, 2021, pp. 12 955–12 965.

[76] H. Chung, B. Sim, and J. C. Ye, "Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction," in *CVPR*, June 2022, pp. 12 413–12 422.

[77] R. S. Zimmermann, L. Schott, Y. Song, B. A. Dunn, and D. A. Klindt, "Score-based generative classifiers," *arXiv preprint arXiv:2110.00473*, 2021.

[78] S. Chen, P. Sun, Y. Song, and P. Luo, "Diffusiondet: Diffusion model for object detection," *arXiv preprint arXiv:2211.09788*, 2022.

[79] D. Baranchuk, A. Voynov, I. Rubachev, V. Khrulkov, and A. Babenko, "Label-efficient semantic segmentation with diffusion models," in *ICLR*, 2021, pp. 1–15.

[80] T. Amit, E. Nachmani, T. Shaharbany, and L. Wolf, "Segdiff: Image segmentation with diffusion probabilistic models," *arXiv preprint arXiv:2112.00390*, 2022.

[81] S. Gu, D. Chen, J. Bao, F. Wen, B. Zhang, D. Chen, L. Yuan, and B. Guo, "Vector quantized diffusion model for text-to-image synthesis," in *CVPR*, June 2022, pp. 10 696–10 706.

[82] W. G. C. Bandara, N. G. Nair, and V. M. Patel, "Ddpm-cd: Remote sensing change detection using denoising diffusion probabilistic models," *arXiv preprint arXiv:2206.11892*, 2022.

[83] 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, 2016.

[84] 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.

[85] B. Liu, A. Yu, X. Yu, R. Wang, and W. Guo, "Deep multiview learning for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. PP, no. 99, pp. 1–15, 2020.

**Ning Chen** received the B.S. degree from the School of Earth and Space Sciences, Peking University, Beijing, China, in 2016, and the M.S. degree in GIS from the School of Earth and Space Sciences, Peking University, Beijing, China in 2019.

He is currently an Assistant Engineer with the Institute of Remote Sensing and Geographic Information System, Peking University. His research interests include satellite image understanding, recommendation system, large-scale sparse learning and pattern recognition.

**Jun Yue** received the B.Eng. degree in geodesy from Wuhan University, Wuhan, China, in 2013 and the Ph.D. degree in GIS from Peking University, Beijing, China, in 2018.

He is currently an Assistant Professor with the School of Automation, Central South University. His research interests include satellite image understanding, pattern recognition, and few-shot learning. Dr. Yue serves as a reviewer for IEEE Transactions on Image Processing, IEEE Transactions on Neural Networks and Learning Systems, IEEE Transactions on Geoscience and Remote Sensing, ISPRS Journal of Photogrammetry and Remote Sensing, IEEE Geoscience and Remote Sensing Letters, IEEE Transactions on Biomedical Engineering, Information Fusion, Information Sciences, etc.

**Leyuan Fang** (Senior Member, IEEE) received the Ph.D. degree from the College of Electrical and Information Engineering, Hunan University, Changsha, China, in 2015.

From August 2016 to September 2017, he was a Postdoc Researcher with the Department of Biomedical Engineering, Duke University, Durham, NC, USA. He is currently a Professor with the College of Electrical and Information Engineering, Hunan University. His research interests include sparse representation and multi-resolution analysis in remote sensing and medical image processing. He is the associate editors of IEEE Transactions on Image Processing, IEEE Transactions on Geoscience and Remote Sensing, IEEE Transactions on Neural Networks and Learning Systems, and Neurocomputing. He was a recipient of one 2nd-Grade National Award at the Nature and Science Progress of China in 2019.

**Shaobo Xia** received the bachelor's degree in geodesy and geomatics from the School of Geodesy and Geomatics, Wuhan University, Wuhan, China, in 2013, the master's degree in cartography and geographic information systems from the Institute of Remote Sensing and Digital Earth, Chinese Academy of Sciences, Beijing, China, in 2016, and the Ph.D. degree in geomatics from the University of Calgary, Calgary, AB, Canada, in 2020.

He is an Assistant Professor with the Department of Geomatics Engineering, Changsha University of Science and Technology, Changsha, China. His research interests include point cloud processing and remote sensing.