# SpectralMamba: Efficient Mamba for Hyperspectral Image Classification

Jing Yao, *Member, IEEE,* Danfeng Hong, *Senior Member, IEEE,* Chenyu Li, and Jocelyn Chanussot, *Fellow, IEEE*

*Abstract*—Recurrent neural networks and Transformers have recently dominated most applications in hyperspectral (HS) imaging, owing to their capability to capture long-range dependencies from spectrum sequences. However, despite the success of these sequential architectures, the non-ignorable inefficiency caused by either difficulty in parallelization or computationally prohibitive attention still hinders their practicality, especially for large-scale observation in remote sensing scenarios. To address this issue, we herein propose SpectralMamba – a novel state space model incorporated efficient deep learning framework for HS image classification. SpectralMamba features the simplified but adequate modeling of HS data dynamics at two levels. First, in spatial-spectral space, a dynamical mask is learned by efficient convolutions to simultaneously encode spatial regularity and spectral peculiarity, thus attenuating the spectral variability and confusion in discriminative representation learning. Second, the merged spectrum can then be efficiently operated in the hidden state space with all parameters learned input-dependent, yielding selectively focused responses without reliance on redundant attention or imparallelizable recurrence. To explore the room for further computational downsizing, a piece-wise scanning mechanism is employed in-between, transferring approximately continuous spectrum into sequences with squeezed length while maintaining short- and long-term contextual profiles among hundreds of bands. Through extensive experiments on four benchmark HS datasets acquired by satellite-, aircraft-, and UAV-borne imagers, SpectralMamba surprisingly creates promising win-wins from both performance and efficiency perspectives. The code will be available at **https://github.com/danfenghong/SpectralMamba** for the sake of reproducibility.

*Index Terms*—Artificial intelligence, efficient, Mamba, hyperspectral image classification, state space model, spatial-spectral, transformer, remote sensing.

## I. INTRODUCTION

**T**HE emergent development of hyperspectral (HS) imaging remarkably empowers humans in observing the real world in greater detail and depth [1]. Unlike traditional photography which acquires images in a limited number of
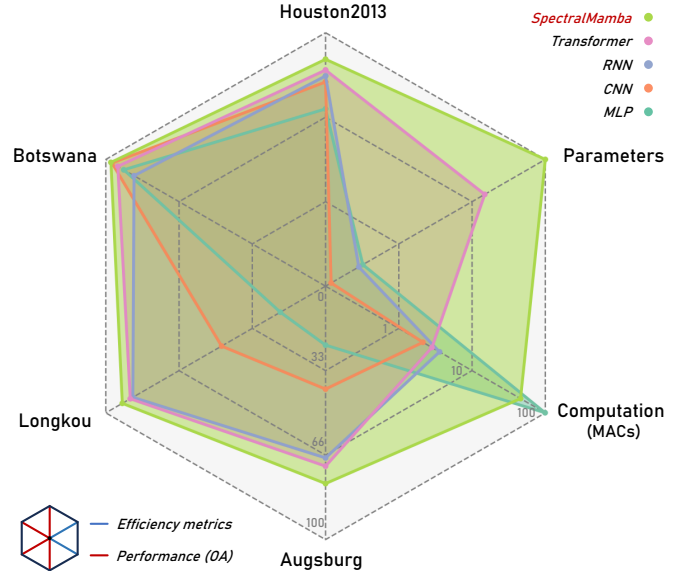
Fig. 1. The radar chart of experimental results of SpectralMamba and classic network architectures in terms of both hyperspectral image classification performance metric (OA) and mean efficiency metrics (number of parameters and MACs) on four benchmark datasets. To better visualize their differences, we set the lowest values of parameter number and MACs as the base score of 100, and customize a base-10 logarithmic scale on the MACs-axis. According to the chart, our SpectralMamba significantly outperforms its competitors along most metrics, showcasing its great potential as a novel efficient, and effective deep learning framework for hyperspectral data analysis.

broad spectral bands, the HS imaging system unprecedentedly achieves the simultaneous capture of both spatial and spectral information by measuring a spectrum of the energy at each pixel. The generated 3-dimensional (3-D) HS data cube contains the near-continuous spectral profile for each spatial resolution element, thus allowing for more accurate quantification, identification, and recognition of the imaged contents. Owing to the recent progress in aerospace and instrumental technology [2], HS imaging has gradually become an indispensable tool for remote sensing (RS). Of all the widespread applications it has found, HS image classification has drawn considerable attention across diverse fields, from environmental monitoring, urban planning, to military science, demonstrating its pervasive potential and cross-cutting importance [3], [4].

The primary objective of HS image classification in RS is to accurately identify the various land cover or land use types of interest within the image by leveraging the detailed spectral signatures associated with each pixel [5]. Despite the capability of HS imaging to capture hundreds of narrow wave-

length bands (typically ranging from the visible spectrum to the near-infrared region), providing in-depth characterization of the spectral properties of different materials, there remain two long-standing challenges in its practical applications.

1) The curse of dimensionality, also known as the Hughes phenomenon, is often encountered in HS image classification [6] when the classification accuracy undergoes an initial rise with more spectral bands observed but then decreases dramatically once a certain number of bands is reached [7]. The fundamental cause of this issue stems from the exponential growth of the feature space volume as the number of dimensions expands, rendering the computational processing and effective analysis of the HS data increasingly burdensome and challenging.

2) Spectral variability and spectral confusion are the other two prevalent phenomena that frequently manifest in HS data. The former issue refers to the situation where the same material displays varying spectral characteristics under different conditions, such as changes in illumination, atmospheric effects, or intrinsic variations, while the latter occurs when distinct materials exhibit similar spectral profiles [8].

In addition, several other issues always arise in conjunction with these challenges in HS data analysis, such as the limited availability of labeled training samples, and inevitable sensor noise with complex distributions, which makes it more challenging to precisely distinguish ground objects based solely on their spectral reflectances.

To tackle these challenges, researchers have dedicated significant efforts over the past decades to develop ever-advancing dimensionality reduction and feature extraction techniques for precise pixelwise recognition of HS images. In the early stages, researchers investigated the applicability of statistical approaches like principal component analysis, independent component analysis, kernel methods [9], and linear discriminant analysis [10], as well as machine learning and heuristics techniques encompassing subspace learning [11], manifold learning [12], ensemble methods [13], and active learning strategies [14], to effectively process and analyze HS data. During this period, shallow machine learning models, such as nearest neighbors, decision trees, and support vector machines, emerged as prevalent choices for effective backend classifiers complementing these feature extraction methods.

As deep learning (DL) has proliferated across numerous research fields since the last decade, the RS community has also embraced this powerful learning paradigm for HS data analysis, harnessing its capability to learn representations directly from the data, thereby mitigating the intrinsic cognitive bias imposed by inadequate mathematical modeling in conventional approaches [15], [16]. Among the various DL architectures for HS image classification, convolution neural network (CNN) has taken pride of place over a long period. Benefiting from the local receptive ability through shift-invariant convolutions, typical CNNs successfully realize hierarchical feature extraction and semantic abstraction in end-to-end training from input-target pairs. However, although this family of models is good at exploiting local contextual information, their inherent local connectivity and weight sharing inevitably restrict the modeling of long-range correlations and dynamics within and across data sequences, respectively.

At this point, sequence models, such as recurrent neural networks (RNNs) and Transformers, came into the public eye for their effectiveness in processing sequential data. By orderly unfolding the HS spectrum into a long sequence, RNNs and Transformers essentially enable the capture of long- and short-term spectral fingerprints through recurrent state modeling and attention mechanism, respectively. Besides, these sequence models have been widely proven to be more competent in handling non-linear data dynamics under complex RS scenes than CNNs. Despite these merits, they inherently suffer from parallel training difficulty and burdensome pairwise multiplicative computations, respectively. Although enormous efforts have been made to cope with these issues, as of yet, most of these variants can hardly avoid becoming increasingly sophisticated in either network structures or working flows in pursuing accuracy breakthroughs with improved representation, appearing to have reached a plateau that cannot avoid trading-off between the performance and computational efficiency.

Fortunately, recent advances in the state space model (SSM) make it broadly applicable and provide a novel avenue for sequentiality modeling. Building upon the theoretical foundation of Classical SSM from control theory and powerful modern DL advantages, the emerging deep SSMs unprecedentedly allow for efficient computation in learning very long-range dependencies over tens of thousands of timesteps and are dominating more and more benchmarks across various domains [17]. Nevertheless, existing SSMs are typically designed for causal learning on low-dimensional sequences such as audio and language, leaving their practicality on high-dimensional visual data such as HS imagery underexplored. Therefore, in this work, we excavate the potential of tailoring SSM to HS data by taking a deep dive into their traits. To be more specific, we propose SpectralMamba – an efficient and effective SSM-integrated DL framework for both pixelwise and patchwise input-based HS image classification. SpectralMamba leverages a simplified but adequate modeling of HS data dynamics in both spatial-spectral feature space and hidden state space, thereby alleviating the effect caused by spectral variability and spectral confusion. The underlying computational overhead caused by parameter size and computations is further reduced by a customized scanning strategy that can additionally enhance sequential representation while maintaining local spectral fingerprints of HS data. The main contributions of this article can be highlighted as follows.

1) We propose a novel SSM-based backbone network termed SpectralMamba that makes a further step towards performant and computation-friendly HS image classification from the perspective of sequence modeling. To the best of our knowledge, this is the first work that well tailors the deep SSM for HS data and its analysis.

2) Targeting the high dimensionality of HS data, and the issues of spectral variability and confusion, we propose the strategies of piece-wise sequential scanning (PSS) and gated spatial-spectral merging (GSSM) to fully encode the underlying spatial regularity and spectral
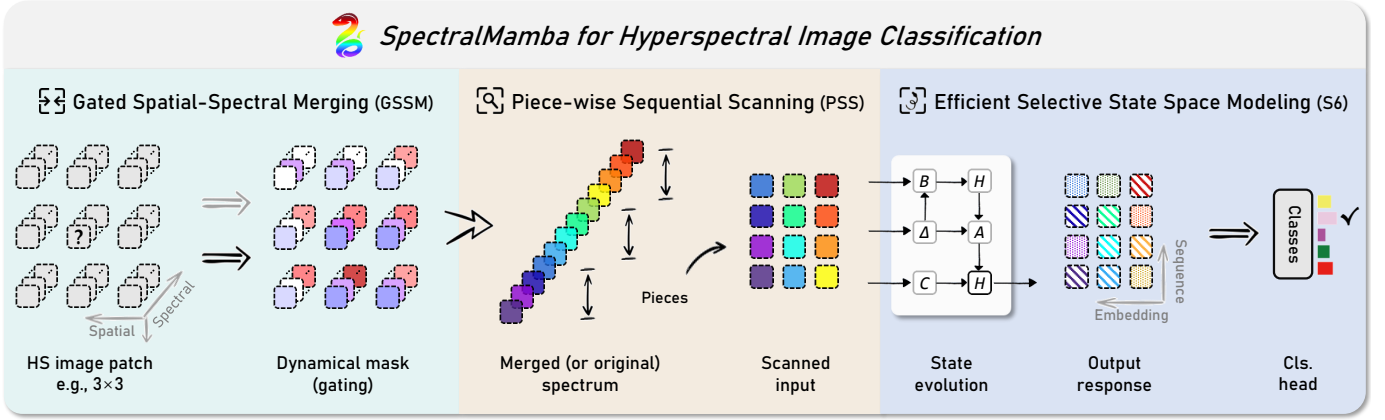
Fig. 2. SpectralMamba mainly consists of three components, i.e., gated spatial-spectral merging (GSSM) module, piece-wise sequential scanning (PSS) strategy, and efficient selective state space (S6) modeling. Its pixelwise counterpart functions by directly operating the original spectrum from the middle stage. $A$, $B$, $C$, and $\Delta$ denote all learnable parameters in the hidden state space, where $H$ records the selectively embedded sequence for the final output.

peculiarity, yielding a more robust discriminative representation learned via a fully lightweight architecture.

3) Through extensive experimental comparison on four benchmark HS datasets acquired from satellite-, aircraft-, and UAV-borne platforms, our SpectralMamba significantly outperforms the representative competitors with classic backbones at a generally minimal computational resource cost (shown in Fig. 1). The ablation studies further verify the effectivity of our key components, such as PSS maximally brings approximately 4% improvement in OA while reducing 60% parameters and 40% computations than our baseline.

The remainder of this article is organized as follows. Section II introduces the preliminary elements for state space models, elaboration of our SpectralMamba, and method analysis. Section III details the experiments, including descriptions of the datasets and implementation, evaluation of both performance and computational cost, comparison results and analysis, and ablation studies. Finally, Section IV concludes the work and points out plausible future directions.

## II. METHODOLOGY

### A. Preliminaries

*1) State Space Model:* Inspired from classical SSMs [18] and modern DL advances, most notably CNNs, RNNs, and Transformers, structured state space sequence models (S4) [17], [19] have recently emerged and garnered considerable attention for modeling sequential data. This class of models typically originates from a continuous-time system that maps an input function or sequence $x(t) \in \mathbb{R}^M$ to an output response signal $y(t) \in \mathbb{R}^O$ through an implicit latent state $h(t) \in \mathbb{R}^N$, which can be mathematically formulated using the following ordinary differential equations,

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t), \qquad (1)$$
$$y(t) = \mathbf{C}h(t) + \mathbf{D}x(t), \qquad (2)$$

where $\mathbf{A} \in \mathbb{R}^{N \times N}$ and $\mathbf{C} \in \mathbb{R}^{O \times N}$ control how the current state evolves over time and translates to the output,

$\mathbf{B} \in \mathbb{R}^{N \times M}$ and $\mathbf{D} \in \mathbb{R}^{O \times M}$ depict how the input influences the state and the output, respectively. Herein we consider the case of a single-input single-output system with $O = M = 1$ and omit $\mathbf{D}x(t)$ term by treating it as a skip connection as the S4 model does [17].

*2) Discretization:* The first step in applying SSMs to discrete signals such as language, audio, and images, is to transform the system parameters into their "discretized" counterpart. A commonly adopted discretization method is the zero-order hold rule, by which the reparameterization holds as follows,

$$\bar{\mathbf{A}} = \exp(\Delta\mathbf{A}), \qquad (3)$$
$$\bar{\mathbf{B}} = (\Delta\mathbf{A})^{-1}(\bar{\mathbf{A}} - \mathbf{I})(\Delta\mathbf{B})$$
$$\approx (\Delta\mathbf{A})^{-1}(\Delta\mathbf{A})(\Delta\mathbf{B}) \qquad (4)$$
$$= \Delta\mathbf{B},$$

where the first-order Taylor series approximation is used by following [20]. The timescale parameter $\Delta$ denotes the sampling step, i.e., $x_k = x(k\Delta)$, which also trades off the state and current input during the evolution. Then the discrete SSM can be formulated into the following recurrent representation,

$$h_k = \bar{\mathbf{A}}h_{k-1} + \bar{\mathbf{B}}x_k, \qquad (5)$$
$$y_k = \mathbf{C}h_k, \qquad (6)$$

which can be computed similarly to RNNs. To better accommodate GPU acceleration for efficient training, S4 also unrolls the above linear recurrence, yielding its global convolutional representation as

$$\mathbf{y} = \mathbf{x} * \bar{\mathbf{K}}, \qquad (7)$$

where $\bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \ldots, \mathbf{C}\bar{\mathbf{A}}^{L-1}\bar{\mathbf{B}})$ represents the SSM convolution kernel, and $L$ is the sequence length of input.

*3) Mamba:* Besides the linearity, another simplifying assumption of the above systems is the time invariance, that is, all the system parameters are defined as time-independent. Recently, a novel class of selective SSMs (S6) breaks this constraint by parameterizing $(\Delta, \mathbf{B}, \mathbf{C})$ as functions of input $\mathbf{x}$, thus endowing SSMs with additional selection ability to focus on the important or ignore the unimportant. As a
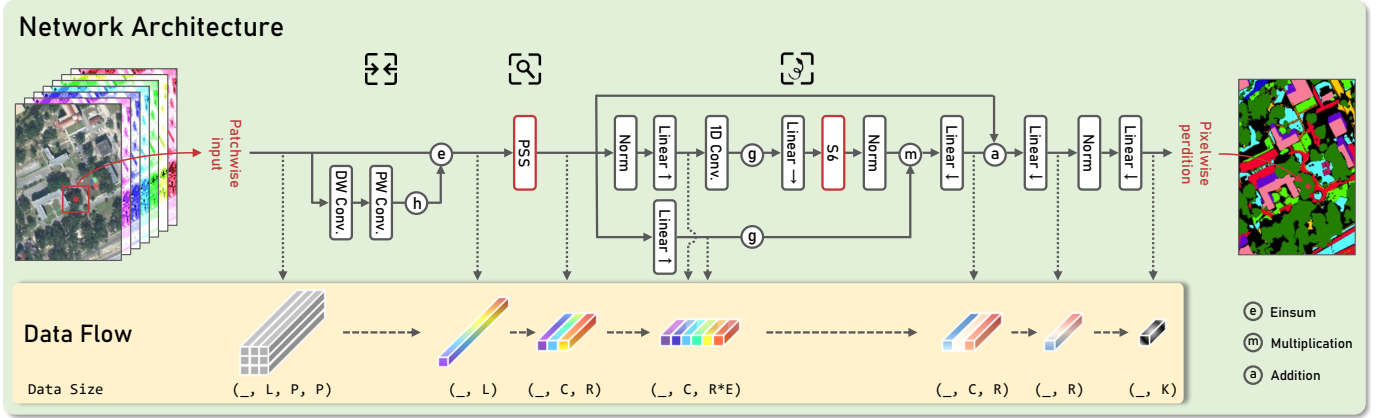
Fig. 3. Detailed architectural design and data processing pipeline of our proposed SpectralMamba, exemplified using a patchwise input under batch training for hyperspectral image classification, where PSS and S6 refer to the piece-wise sequential scanning and selective state space model, respectively.

common fashion, a simplified neural network architecture consists of linear layers, a convolution layer, a residual connection, nonlinear transformations, and most importantly the S6 kernel together to form the Mamba block. A hardware-aware optimization is also proposed to guarantee its efficient implementation.

### B. SpectralMamba: Overview

To break through the performance and efficiency bottlenecks of existing methods based on CNN, RNN, or Transformer backbones, we propose the SpectralMamba, a Mamba – S6 model incorporated DL solution to tackle HS image classification. The key of SpectralMamba lies in its simultaneous modeling of HS data dynamics by gated spectrum merging in the spatial-spectral space and selective sequence learning in hidden state space through a minimally parameterized network architecture. Furthermore, a novel sequential scanning strategy tailored for HS data is proposed to make the framework more computation-friendly by unraveling the hundreds-of-bands spectrum into pieces. The patchwise SpectralMamba can also be flexibly transformed into its pixelwise counterpart by skipping the spatial-spectral encoding stem. Fig. 2 gives an illustration of how our proposed SpectralMamba works, and the detailed network architecture as well as its data flow are shown in Fig. 3.

### C. SpectralMamba: Key Components

Let the 1-D vector $\mathbf{x}^{pixel} = [\mathbf{x}_1^{pixel}, \ldots, \mathbf{x}_L^{pixel}] \in \mathbb{R}^{1 \times L}$ denote a given pixel of an HS image, where $L$ is the band number of the spectrum. We consider the state space modeling within the spectral domain, that is, our aim becomes to find its output response $\mathbf{y}^{pixel} \in \mathbb{R}^L$ through a well-defined S6 model. However, simply treating the reflectance values at each band as the representation may limit the mining of sequential patterns. Therefore we expand the model dimension by a factor, $E = 8$ in our case, to enlarge the state space size. As the core architectural part in Fig. 3 depicts, the Mamba block designed for HS data consists of three streams. Its mainstream comprises two distinct LayerNorm layers in input

and expanded state spaces, three linear layers for expanding, keeping, and compressing feature dimensions, a SiLU nonlinear activation, and the S6 block. The other two streams are the skip connection and an excitation-like multiplication to adaptively transform original information across layers [21]. Note that the use of skip connection and nonlinearity is crucial for stable training with fast convergence, while the practical performance appears to be not that sensitive to the choice of normalization and activation function.

*1) Piece-wise Sequential Scanning:* The aforementioned state space modeling is empowered to pay attention or forget features at particular wavelengths through an input-dependent parameterizing way [20]. When applied to HS data with hundreds of near-continuous bands, its high spectral redundancy drives us to rethink the input manner. Unlike the very recent efforts in modifying Mamba for natural image processing by considering spatial multi-directional scanning [22], we propose a novel piece-wise sequential scanning (PSS) along the spectral dimension to fully leverage the reflectance characteristics of different types of ground objects.

In specific, we can formulate the PSS module as

$$PSS(\mathbf{x}^{pixel}) = [S_1 \mathbf{x}^{pixel}, \ldots, S_R \mathbf{x}^{pixel}], \quad (8)$$

where $S_r \mathbf{x}^{pixel} = [\mathbf{x}_{(r-1)C+1}^{pixel}, \ldots, \mathbf{x}_{rC}^{pixel}]^\top \in \mathbb{R}^{C \times 1}$ scans continuous pieces from the original spectrum $\mathbf{x}^{pixel}$ for $r = 1, \ldots, R$, and $R$ is the number of pieces with length $C$. This also acts similarly as a resampling with features at each position in the sequence enriched. By applying PSS before our Mamba block, the response accordingly turns from a 1-D $L$-length vector to a 2-D output of shape $C \times R$. We then add one more pre-layer before the common softmax-based classification head to finally obtain $K$-length categorical logits.

*2) Gated Spatial-Spectral Merging:* It is crucially important to consider spatial information for discriminative representation learning. However, conventional spatial-spectral feature extraction methods commonly treat each patch equally with fixed convolution kernels. Inspired by the S6 model that learns interactions along the sequence in an input-dependent way, we propose to further increase the content-awareness by introducing a dynamical gate function for adaptive spatial-

TABLE I
SUMMARY OF FOUR INVESTIGATED HS DATASETS, INCLUDING THE TAGS OF DATA ACQUISITION INFORMATION, CATEGORICAL INFORMATION OF GROUND OBJECTS, AND THE CORRESPONDING NUMBERS OF TRAIN AND TEST SAMPLES.

| Dataset | Houston2013 | | | Augsburg | | | Longkou | | | Botswana | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Sensor | ITRES CASI-1500 | | | HySpex | | | Headwall Nano-Hyperspec | | | Hyperion | | |
| Platform | Aircraft-borne | | | Aircraft-borne | | | UAV-borne | | | Satellite-borne | | |
| Loc. & Time | America, 2012 | | | Europe, 2018 | | | Asia, 2018 | | | Africa, 2001 | | |
| GSD | 10 m | | | 30 m | | | 0.463 m | | | 30 m | | |
| Wavelength | 380 nm - 1050 nm | | | 400 nm - 2500 nm | | | 400 nm - 1000 nm | | | 400 nm - 2500 nm | | |
| Data Size | $349 \times 1905 \times 144$ | | | $332 \times 485 \times 180$ | | | $550 \times 400 \times 270$ | | | $1476 \times 256 \times 145$ | | |
| Class No. | Class Name | Train | Test | Class Name | Train | Test | Class Name | Train | Test | Class Name | Train | Test |
| 1 | Healthy Grass | 198 | 1053 | Forest | 80 | 13427 | Corn | 54 | 34457 | Water | 24 | 246 |
| 2 | Stressed Grass | 190 | 1064 | Residential Area | 80 | 30249 | Cotton | 53 | 8321 | Hippo Grass | 20 | 81 |
| 3 | Synthetic Grass | 192 | 505 | Industrial Area | 80 | 3771 | Sesame | 54 | 2977 | Floodplain Grasses 1 | 21 | 230 |
| 4 | Tree | 188 | 1056 | Low Plants | 80 | 26777 | Broad-Leaf Soybean | 53 | 63159 | Floodplain Grasses 2 | 24 | 191 |
| 5 | Soil | 186 | 1056 | Allotment | 80 | 495 | Narrow-Leaf Soybean | 54 | 4097 | Reeds | 20 | 249 |
| 6 | Water | 182 | 143 | Commercial Area | 80 | 1565 | Rice | 50 | 11804 | Riparian | 20 | 249 |
| 7 | Residential | 196 | 1072 | Water | 81 | 1449 | Water | 53 | 67003 | Fires Car | 20 | 239 |
| 8 | Commercial | 191 | 1053 | - | - | - | Roads and Houses | 53 | 7071 | Island Interior | 20 | 183 |
| 9 | Road | 193 | 1059 | - | - | - | Mixed Weed | 51 | 5178 | Acacia Woodlands | 21 | 293 |
| 10 | Highway | 191 | 1036 | - | - | - | - | - | - | Acacia Shrub Lands | 22 | 226 |
| 11 | Railway | 181 | 1054 | - | - | - | - | - | - | Acacia Grasslands | 20 | 285 |
| 12 | Parking Lot 1 | 192 | 1041 | - | - | - | - | - | - | Short Mopane | 22 | 159 |
| 13 | Parking Lot 2 | 184 | 285 | - | - | - | - | - | - | Mixed Mopane | 20 | 248 |
| 14 | Tennis Court | 181 | 247 | - | - | - | - | - | - | Exposes Soils | 20 | 75 |
| 15 | Running Track | 187 | 473 | - | - | - | - | - | - | - | - | - |
| Total | - | 2832 | 12197 | - | 561 | 77733 | - | 475 | 204067 | - | 294 | 2954 |

spectral embedding. With the proposed gated spatial-spectral merging (GSSM), we can substitute the $\mathbf{x}^{pixel}$ in Eq. (8) with merged spectrum computed as follows,

$$GSSM(\mathbf{x}^{patch}) = h(f_{PW}(f_{DW}(\mathbf{x}^{patch}))) \otimes \mathbf{x}^{patch}, \quad (9)$$

where $h$ is the sigmoid activation, $f$ represents the composite function composed of depthwise (DW) convolution and pointwise (PW) convolution, $\otimes$ denotes the Einstein summation along the spatial dimension that combines two tensors of shape $L \times P \times P$ into a 1-D vector with length $L$. Through GSSM, we hope to adaptively encode the semantic relationships between the center pixel and its neighborhoods in learning a more discriminative "spectrum", thereby attenuating spectral variability and confusion effect.

### D. SpectralMamba: Method Analysis

It is non-trivial to directly extend SSM to applications of HS data. Building upon the insights on SSM and structural prior knowledge of HS data, our SpectralMamba offers a viable SSM-based baseline to address the dense prediction application of HS images. The proposed PSS strategy not only enables the model to uncover local characteristics of the spectral profile but also further improves efficiency by narrowing the width of core operation networks. Moreover, the GSSM module is designed based on the observation that the semantic relationships between the central pixel and its neighboring pixels usually vary spatially and spectrally across the scene. The widely-existed mixed pixel phenomenon can also differ for pixels within a local patch, especially for those falling on boundaries. We hope to efficiently capture such highly spatial-spectral-variant HS data dynamics through a lightweight mask learner, thus yielding a merged spectrum with increased discriminative ability for the subsequent sequential learning in the state space.

The connections between our proposed SpectralMamba and related works are also worth noting. On one hand, although CasRNN – an RNN-based representative for HS image classification – has considered the similar spectral redundancy issue by hierarchically learning from adjacent spectral bands to nonadjacent ones, their imparallelizable recurrence trait still accumulated both the computations and parameters in pursuing a stable training [23]. Also, the unbearable quadratic complexity of self-attention in conventional Transformers has significantly magnified the computational burden as the number of neighboring spectral bands considered in a so-called groupwise embedding increases [24]. In contrast, the proposed PSS in our SpectralMamba perfectly matches the efficient feature selection via state evolution in S6, simultaneously preserving local spectral patterns and enlarging the feature dimension of sequence, while further improving computational efficiency. What's more, our practice verifies that a non-overlapped scanning, i.e., $R = L/C$, is enough to produce a promising performance at a lower computational overhead. On the other hand, different from conventional gated convolutions for natural images [25], our GSSM provides a lightweight gating mechanism in capturing the highly spatial-spectral-variant HS data dynamics. It not only fits the setting to maintain the spectral sequentiality but also complements the subsequent content-aware learning in the state space with efficient spatial-spectral rectification. In the following experiments section, we will demonstrate how SpectralMamba outperforms these predecessors by enhancing the efficacy in interpreting HS data on the abovementioned merits, meanwhile, maintaining high efficiency with low computational resource requirements.

## III. EXPERIMENTS

### A. Datasets

We select four benchmark HS datasets to conduct experiments, covering all types of acquisition platforms, that are, aircraft-borne, unmanned aerial vehicle (UAV)-borne, and satellite-borne, hoping to give a comprehensive and faithful
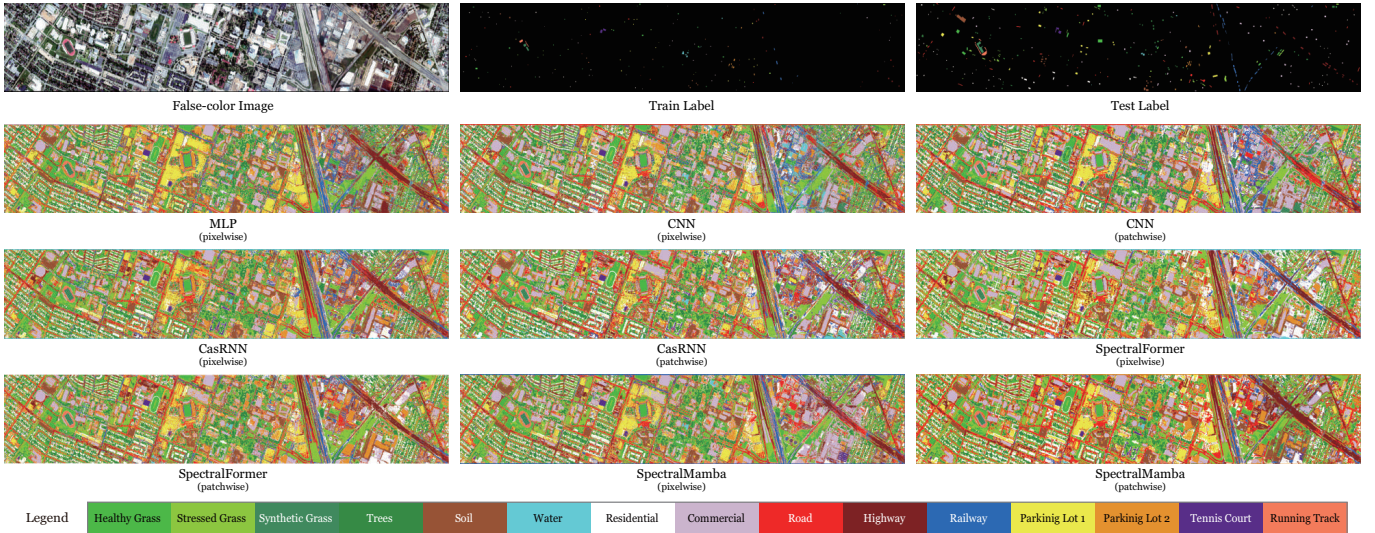
Fig. 4. An illustration of the false-color image, train and test labels, and classification maps obtained by compared methods on the Houston2013 HS dataset.

TABLE II
QUANTITATIVE COMPARISON WITH RELATED REPRESENTATIVES ON HOUSTON2013 HS DATASET. THE BEST AND SECOND-BEST OVERALL RESULTS ARE SHOWN IN BOLD AND UNDERLINED, RESPECTIVELY.

| Method | MLP [27] | CNN [28] | | CasRNN [23] | | SpectralFormer [24] | | SpectralMamba | |
|---|---|---|---|---|---|---|---|---|---|
| Implementation | Pixelwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise |
| Healthy Grass | 83.29 | 84.24 | 82.24 | 83.67 | 82.62 | 82.81 | 82.53 | 93.83 | **95.92** |
| Stressed Grass | 59.59 | 68.05 | 86.84 | 97.74 | 98.31 | 98.68 | 96.71 | 95.30 | **98.87** |
| Synthetic Grass | 98.81 | 91.09 | 96.24 | **100.00** | 96.44 | 99.60 | 98.22 | 99.80 | 98.61 |
| Tree | 83.71 | 80.59 | **99.72** | 98.20 | 98.96 | 98.48 | 98.20 | 96.50 | 99.05 |
| Soil | 46.12 | 49.62 | 85.51 | 97.73 | 97.92 | 96.02 | 98.77 | 98.86 | **99.91** |
| Water | 93.71 | 92.31 | **95.80** | 95.10 | 95.10 | 95.10 | 93.01 | 95.10 | 94.41 |
| Residential | 59.51 | 78.64 | 78.54 | 86.75 | 85.26 | 85.91 | **91.32** | 82.46 | 87.03 |
| Commercial | 69.42 | 64.01 | **82.24** | 60.21 | 66.29 | 50.71 | 67.43 | 77.30 | 64.86 |
| Road | 76.30 | 78.56 | 70.73 | 74.22 | 71.48 | 73.56 | 71.95 | **82.06** | 80.64 |
| Highway | 53.38 | 39.77 | 50.10 | 81.37 | 69.40 | 88.51 | 86.29 | 60.33 | **98.17** |
| Railway | 81.69 | 56.45 | 85.77 | **88.80** | 83.02 | 75.81 | 81.02 | 81.40 | 80.36 |
| Parking Lot 1 | 52.83 | 42.94 | 62.73 | 68.01 | 57.44 | 65.03 | 65.90 | 74.26 | **81.56** |
| Parking Lot 2 | 67.02 | 57.89 | 81.40 | 75.79 | **84.91** | 72.63 | 69.82 | 72.63 | 79.65 |
| Tennis Court | 92.31 | 94.33 | 99.60 | 99.60 | 98.79 | 99.60 | 97.17 | 99.19 | **100.00** |
| Running Track | 96.19 | 98.10 | 94.50 | 97.89 | 95.35 | 98.73 | **99.15** | 98.31 | 98.52 |
| OA (%) ↑ | 69.94 | 67.58 | 80.57 | 85.21 | 82.94 | 83.32 | 85.27 | <u>85.64</u> | **89.52** |
| AA (%) ↑ | 74.26 | 71.77 | 83.46 | 87.01 | 85.42 | 85.41 | 86.50 | <u>87.16</u> | **90.50** |
| κ ↑ | 0.6749 | 0.6501 | 0.7894 | 0.8397 | 0.8149 | 0.8193 | 0.8402 | <u>0.8442</u> | **0.8864** |
| MACs (M) ↓ | **19.64** | 26.83 | 640.09 | 548.45 | 549.95 | 667.04 | 681.19 | <u>23.52</u> | 36.21 |
| Params (K) ↓ | 305.94 | 418.77 | 1192.96 | 350.61 | 353.49 | 72.56 | 74.10 | **14.23** | <u>36.55</u> |

verification of our proposed SpectralMamba. The details are summarized in Table I.

*1) Houston2013 HS Dataset:* The first HS dataset Houston2013 consists of HS imagery acquired by ITRES CASI-1500 HS imager with 144 spectral bands ranging from 380 nm to 1050 nm. The investigated scene has 349 × 1905 pixels at a ground sampling distance (GSD) of 10 m and 15 LULC categories, covering the University of Houston campus and the surrounding urban area. It was provided by the IEEE Geoscience and Remote Sensing Society data fusion contest in 2013 and has been widely used in research and competitions related to HS image analysis and pixel-based classification [26].

*2) Longkou HS Dataset:* The third Longkou dataset, also known as WHU-Hi-LongKou, is a specific UAV-borne HS dataset acquired in Longkou Town, Hubei, China. This dataset was captured using the Headwall Nano-Hyperspec sensor mounted on a UAV at a flight altitude of 500 meters. We

use the opened data that has 550 × 400 pixels at a GSD of 0.463 m, and 270 spectral bands ranging from 400 nm to 1000 nm. More than 204 thousand labelings with 6 crop types and 3 LULC categories were used for evaluation. This dataset is part of the WHU-Hi dataset collection, which also includes other two datasets HanChuan and HongHu [29].

*3) Augsburg HS Dataset:* The second HS dataset Augsburg consists of HS imagery acquired by the HySpex – an airborne imaging spectrometer system operated by the the Remote Sensing Technology Institute of the German Aerospace Center [11]. The pre-processed HS imagery has 180 bands ranging from 400 nm to 2500 nm with high quality. Our selected sub-region includes 332 × 485 pixels at a GSD of 30m, covering the city of Augsburg, Germany. The ground reference map comprising 7 categories was elaborately produced by manual labeling based on OpenStreetMap product.

*4) Botswana HS Dataset:* The last Botswana dataset is a collection of HS imagery acquired by the NASA EO-1 satellite
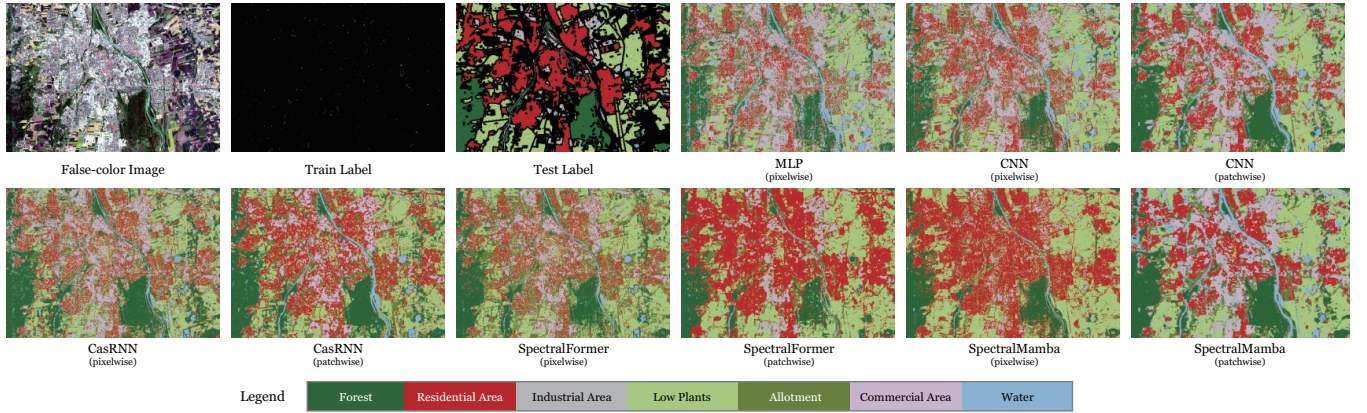
Fig. 5. An illustration of the false-color image, train and test labels, and classification maps obtained by compared methods on the Augsburg HS dataset.

TABLE III
QUANTITATIVE COMPARISON WITH RELATED REPRESENTATIVES ON AUGSBURG HS DATASET. THE BEST AND SECOND-BEST OVERALL RESULTS ARE SHOWN IN BOLD AND UNDERLINED, RESPECTIVELY.

| Method | MLP [27] | CNN [28] | | CasRNN [23] | | SpectralFormer [24] | | SpectralMamba | |
|---|---|---|---|---|---|---|---|---|---|
| Implementation | Pixelwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise |
| Forest | 13.27 | 12.15 | 53.77 | 94.30 | 96.31 | 90.23 | 92.74 | 87.96 | **98.31** |
| Residential Area | 23.61 | 32.71 | 35.84 | 47.53 | 63.76 | 45.50 | 53.58 | 57.33 | **73.79** |
| Industrial Area | 34.37 | 25.40 | 32.67 | 30.76 | 9.71 | 24.85 | **54.65** | 35.38 | 48.71 |
| Low Plants | 22.90 | 29.86 | 38.45 | 57.23 | 66.52 | 71.43 | **84.46** | 77.39 | 77.60 |
| Allotment | 62.63 | 48.48 | 80.00 | 71.11 | 90.30 | 82.22 | 75.35 | 77.37 | **94.14** |
| Commercial Area | 42.56 | 30.22 | 45.30 | 45.30 | **58.34** | 53.04 | 36.17 | 48.12 | 58.21 |
| Water | 51.83 | 50.38 | 61.08 | 53.42 | 63.35 | 49.55 | 64.67 | 61.90 | **71.98** |
| OA (%) ↑ | 23.26 | 28.20 | 40.63 | 58.35 | 67.77 | 61.62 | <u>71.03</u> | 68.49 | **77.90** |
| AA (%) ↑ | 35.88 | 32.75 | 49.59 | 57.09 | 64.04 | 59.54 | <u>65.95</u> | 63.64 | **74.68** |
| κ ↑ | 0.0942 | 0.1094 | 0.2692 | 0.4645 | 0.5673 | 0.5018 | <u>0.6157</u> | 0.5864 | **0.7048** |
| MACs (M) ↓ | **20.10** | 41.74 | 1038.34 | 685.38 | 687.24 | 832.60 | 850.29 | <u>30.96</u> | 50.55 |
| Params (K) ↓ | 313.10 | 651.61 | 1861.63 | 348.55 | 352.15 | 72.04 | 73.58 | **13.56** | <u>47.94</u> |

over the Okavango Delta area of Botswana between 2001 and 2004. This dataset was captured using the well-known Hyperion sensor, which collects original data at a 30-meter spatial resolution in 242 bands covering the 400-2500 nm portion of the spectrum. After pre-processing by removing uncalibrated and noisy bands, 145 bands are commonly used for classification experiments. The ground reference comprises 14 identified classes representing different land cover types in the investigated region.

### B. Experimental Setup

To ensure the reproduction of our experimental results and a faithful verification of the effectiveness of our method, we herein present the necessary details in implementing all compared methods on investigated datasets.

*1) Train and Test Set Split:* Besides the data quality, the split of the train and test set has a significant effect on assessing the model performance, particularly the deep models. We strive to reason the popularity of the Houston2013 dataset and conclude four criteria for sample selection with easy practicality and wide acceptability.

1) First, the samples selected for training are better evenly distributed in the scene to ameliorate the spectral variability phenomenon that may deteriorate most models.
2) Second, it is important to maintain a moderate train set size to prevent the evaluation that is either too

simplistic or too challenging, therefore striking a balance in accurately reflecting the model's performance and generalization ability.
3) Third, rather than pixelwise sampling, image segment-based sampling or labeling tends to be more efficient in progressively constructing the train set to a preset size.
4) Last but not least, class-balanced sampling is always preferable, especially for unseen scenes with no prior knowledge of the class distribution.

Based on the above criteria, we propose the following steps to determine the train and test set for those datasets without benchmark split. The first step is to apply the SLIC superpixel method to over-segment the image into a large number of segments [30], most of which can then be made of the same ground object. The next step is straightforward to randomly collect those homogeneous segments progressively until the budget of training samples for each class is reached. The classwise budgets for the last three HS datasets we adopt are empirically set to 80, 50, and 20, respectively. This method can also be used to organize a separate validation set [31].

*2) Evaluation Metrics:* Besides the conventional performance metrics such as classwise accuracy (CA), overall accuracy (OA), average accuracy (AA), and Kappa coefficient (κ), we introduce two more metrics to evaluate the efficacy of different methods, which are multiply-accumulate operations (MACs) and parameters (Params) of each network. As the name suggests, MACs refer to the number of multiply-
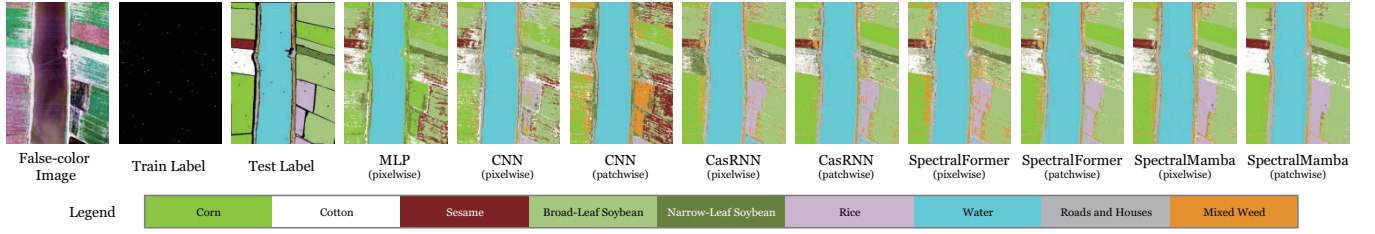
Fig. 6. An illustration of the false-color image, train and test labels, and classification maps obtained by compared methods on the Longkou HS dataset.

TABLE IV
QUANTITATIVE COMPARISON WITH RELATED REPRESENTATIVES ON LONGKOU HS DATASET. THE BEST AND SECOND-BEST OVERALL RESULTS ARE SHOWN IN BOLD AND UNDERLINED, RESPECTIVELY.

| Method | MLP [27] | CNN [28] | | CasRNN [23] | | SpectralFormer [24] | | SpectralMamba | |
|---|---|---|---|---|---|---|---|---|---|
| Implementation | Pixelwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise |
| Corn | 37.36 | 12.68 | 4.64 | 89.38 | 95.98 | 85.47 | 85.47 | 91.60 | **97.48** |
| Cotton | 28.16 | 25.12 | 33.29 | 30.50 | 70.46 | 69.53 | 69.53 | 74.49 | **95.01** |
| Sesame | 29.63 | 30.50 | 39.13 | 51.63 | 65.91 | 60.36 | 60.36 | 78.17 | **92.17** |
| Broad-Leaf Soybean | 30.95 | 32.32 | 47.81 | 73.85 | 77.15 | 73.44 | 73.44 | 79.27 | **82.87** |
| Narrow-Leaf Soybean | 32.63 | 25.24 | 66.49 | 76.42 | 86.45 | 83.26 | 83.26 | **90.85** | 90.68 |
| Rice | 0.00 | 22.29 | 13.20 | 80.77 | 83.62 | 79.67 | 79.67 | 95.41 | **98.29** |
| Water | 7.47 | 6.17 | 80.98 | 99.83 | **99.99** | **99.99** | **99.99** | 98.66 | 99.60 |
| Roads and Houses | 0.00 | 13.53 | 7.27 | 79.86 | 86.17 | 77.90 | 77.90 | 61.76 | **88.81** |
| Mixed Weed | 0.00 | 23.99 | 33.45 | 48.05 | 54.63 | 53.13 | 53.13 | 68.75 | **73.74** |
| OA (%) ↑ | 20.57 | 18.51 | 47.30 | 82.92 | 87.69 | 84.03 | <u>88.86</u> | 87.80 | **92.48** |
| AA (%) ↑ | 18.47 | 21.32 | 36.25 | 70.03 | 80.04 | 75.86 | <u>82.29</u> | 82.11 | **90.96** |
| κ ↑ | 0.0575 | 0.0693 | 0.3558 | 0.7815 | 0.8417 | 0.7964 | <u>0.8566</u> | 0.8432 | **0.9028** |
| MACs (M) ↓ | **21.61** | 93.74 | 2333.70 | 954.68 | 957.48 | 1276.52 | 1303.65 | <u>52.79</u> | 93.18 |
| Params (K) ↓ | 336.65 | 1499.14 | 4186.11 | 349.07 | 354.47 | <u>72.17</u> | 73.71 | **18.68** | 94.55 |

accumulate operations during actual network training. In our experiments, we set the batch size for all methods as 64 to fairly compute and compare their MACs for one batch. The less value the MACs and Params take, the less computational resource the corresponding model costs.

*3) Compared Methods:* The major aim of our experiments is to assess whether the proposed SpectralMamba can be deemed as a revolutionized DL tool for the HS image classification task. Therefore, we select four prevalent types of DL-based solutions for comparison, which are MLP, CNN, RNN-based, and Transformer-based models. The details of our competitors are as follows.

1) The MLP contains an input block, two hidden blocks, and a classification head. Each of the first three blocks consists of a fully-connected (FC) layer, a 1-D batch normalization (BN) layer, and a ReLU activation [27].
2) The CNN contains two convolution blocks and a classification head. The convolution block consists of 1-D or 2-D convolution layers depending on pixelwise or patchwise input [28], corresponding 1-D or 2-D BN layers, and ReLU.
3) The CasRNN is selected to represent the RNN-based models [23]. It has pixelwise and patchwise versions. The pixelwise CasRNN mainly contains cascaded gated recurrent units, while the patchwise one additionally introduces two separable convolution blocks and a max-pooling layer before the pixelwise CasRNN[1].
4) The pixelwise SpectralFormer and patchwise Spectral-Former are selected to represent the Transformer-based models [24]. The former one is equipped with the

groupwise spectral embedding to enhance local spectral details while the latter adds an FC layer to encode the spatial information from flattened image patches[2].

*4) Implementation Details:* All of our experiments are conducted mainly based on the PyTorch framework using a workstation with an Nvidia GeForce GTX 3080 GPU card. We tuned the learning rate and weight decay hyperparameters via a rough grid search on intervals of $\{10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}\}$. The epoch is set to 500. A StepLR scheduler is adopted to shrink the learning rate by multiplying a factor of 0.9 for each 20 epochs.

Another commonly overlooked issue is information leakage, that is, more testing samples will be seen as the patch size of training samples increases [32]. To mitigate the effect caused by this issue, we set the patch size for all patchwise models as 3. We claim that this setup is enough for assessment from the perspective of tendency, which also consumes less energy and is aligned with sustainable development goals.

*C. Results and Analysis*

We summarize the quantitative results of all compared methods on six metrics, i.e., performance metrics as CA, OA, AA, and κ, and efficacy metrics as MACs and Params, in Tables II to V for Houston2013, Augsburg, Longkou, and Botswana datasets, respectively. The classification maps of full scenes are correspondingly visualized in Figs. 4 to 7.

From the tables, we can draw several conclusions that are consistent on four datasets. Patchwise implementations generally tend to produce better performance than pixelwise ones by

---

[1]https://github.com/RenlongHang/CasRNN

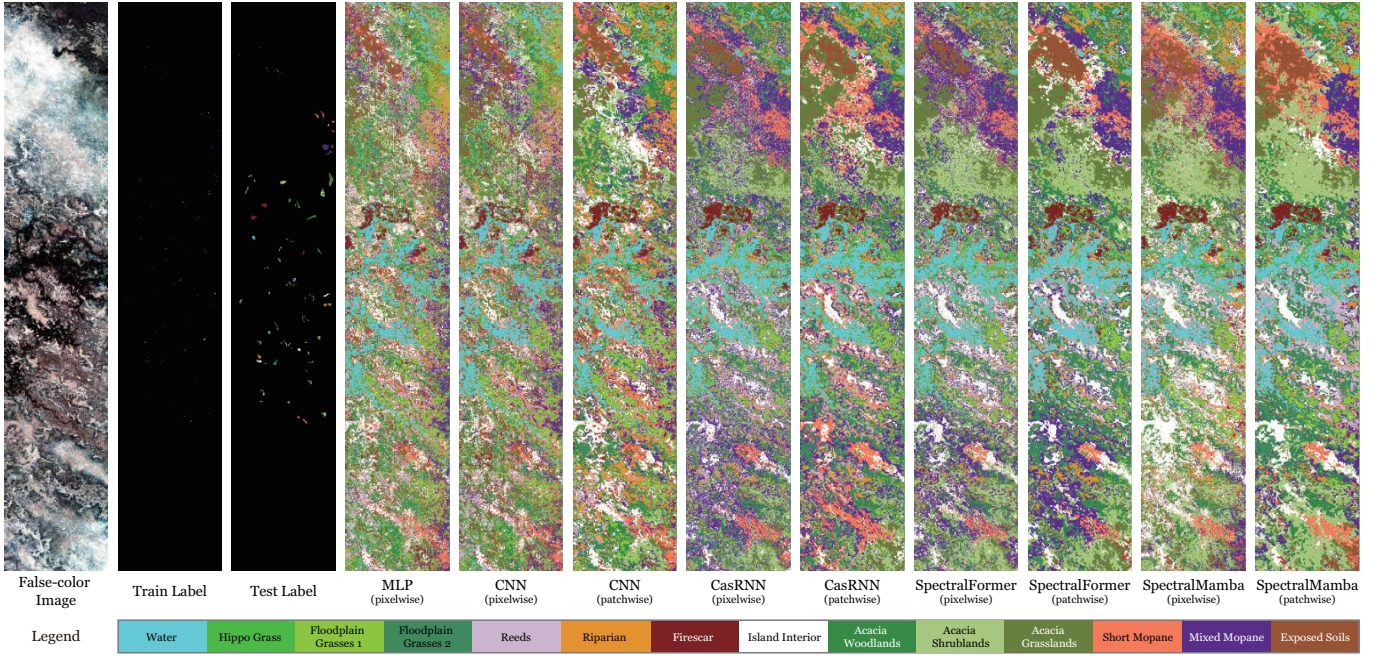[2]https://github.com/danfenghong/IEEE_TGRS_SpectralFormer

Fig. 7. An illustration of the false-color image, train and test labels, and classification maps obtained by compared methods on the Botswana HS dataset.

TABLE V
QUANTITATIVE COMPARISON WITH RELATED REPRESENTATIVES ON BOTSWANA HS DATASET. THE BEST AND SECOND-BEST OVERALL RESULTS ARE
SHOWN IN BOLD AND UNDERLINED, RESPECTIVELY.

| Method | MLP [27] | CNN [28] | | CasRNN [23] | | SpectralFormer [24] | | SpectralMamba | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Implementation | Pixelwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise | Pixelwise | Patchwise |
| Water | **100.00** | **100.00** | **100.00** | **100.00** | **100.00** | **100.00** | **100.00** | **100.00** | **100.00** |
| Hippo Grass | **100.00** | **100.00** | **100.00** | 97.53 | 97.53 | **100.00** | 98.77 | 95.06 | **100.00** |
| Floodplain Grasses 1 | 94.78 | 98.26 | **100.00** | 81.30 | 96.52 | 90.87 | 95.65 | 95.65 | **100.00** |
| Floodplain Grasses 2 | 96.86 | 96.86 | 94.76 | 88.48 | 94.76 | 93.72 | 92.15 | 92.15 | **97.91** |
| Reeds | 90.36 | 89.16 | 90.76 | 75.90 | 73.49 | 79.52 | 87.55 | 86.75 | **93.17** |
| Riparian | 71.89 | 82.73 | **97.59** | 71.08 | 68.27 | 65.86 | 87.15 | 77.91 | 94.78 |
| Firescar | 97.49 | 96.23 | **100.00** | 96.65 | 89.96 | 87.87 | 99.58 | 97.91 | **100.00** |
| Island Interior | 99.45 | 98.91 | **100.00** | 91.26 | 97.81 | 97.27 | **100.00** | 95.63 | **100.00** |
| Acacia Woodlands | 84.98 | 74.74 | 90.10 | 72.35 | 72.35 | 74.74 | 87.71 | 86.35 | **96.25** |
| Acacia Shrublands | 95.13 | 93.36 | **97.79** | 73.45 | 74.34 | 92.48 | 95.13 | 96.02 | **97.79** |
| Acacia Grasslands | 87.37 | 92.28 | **98.25** | 85.96 | 92.28 | 84.91 | 97.19 | 85.61 | 94.04 |
| Short Mopane | 96.23 | 95.60 | **100.00** | 99.37 | **100.00** | 96.86 | **100.00** | 91.19 | 99.37 |
| Mixed Mopane | 89.52 | 91.94 | 98.39 | 76.61 | 89.92 | 88.71 | 93.55 | 97.58 | **99.60** |
| Exposed Soils | **100.00** | 98.67 | 98.67 | 94.67 | 98.67 | 98.67 | **100.00** | **100.00** | **100.00** |
| OA (%) ↑ | 91.81 | 92.21 | <u>97.19</u> | 84.19 | 87.14 | 87.44 | 94.55 | 91.88 | **97.66** |
| AA (%) ↑ | 93.15 | 93.48 | <u>97.59</u> | 86.05 | 88.99 | 89.39 | 95.32 | 92.70 | **98.06** |
| κ ↑ | 0.9112 | 0.9156 | <u>0.9695</u> | 0.8287 | 0.8606 | 0.8639 | 0.9409 | 0.9119 | **0.9747** |
| MACs (M) ↓ | **19.64** | 27.19 | 658.12 | 570.80 | 572.31 | 671.63 | 685.89 | <u>23.39</u> | 36.25 |
| Params (K) ↓ | 305.93 | 424.43 | 1209.34 | 350.35 | 353.25 | 72.49 | 74.03 | **12.33** | <u>34.96</u> |

exploitation of spatial information (except for CasRNN on the Houston2013 dataset, which is probably caused by the unstable training issue). In specific, pixelwise MLP and pixelwise CNN involve fewer computations than CasRNN and SpectralFormer. Still, patchwise CNN encounters an evident explosion in terms of MACs due to its use of inseparable 2-D convolution. If we examine it further through the lens of computational efficiency, although SpectralFormer owns much fewer parameters than those of CasRNN, its self-attention operations result in higher MACs, particularly for dealing with longer sequences in the middle two datasets. Most importantly, our SpectralMamba unsurprisingly achieves the best CAs in most classes, the best OA, AA, and κ, and meanwhile the best Params, and the MACs that are only second to those of the simplistic MLP with a negligible margin.

Some other interesting observations are also worth noting. If we set the quantitative tendency on the Houston2013 dataset as the reference, both the conventional pixelwise MLP and pixelwise CNN behave contrarily on Botswana and the other two datasets. This could possibly be explained by that the train and test sample distributions are much closer to each other on the Botswana dataset than on the other two. Patchwise CNN raises the performance of pixelwise CNN by 12.43% and 28.79% in terms of OA on the middle two datasets, but still at a low level. CasRNN, SpectralFormer, and our SpectralMamba, no matter whether pixelwise or patchwise implementation, consistently show a stable classification performance at a desirable level, which verifies the importance and superiority of their sequential modeling capability. Furthermore, our SpectralMamba successfully addresses the computational weaknesses that lie

TABLE VI
QUANTITATIVE RESULTS OF ABLATION BY SWITCHING ON AND OFF THREE KEY MODULES IN THE PROPOSED SPECTRALMAMBA ON HOUSTON2013 HS DATASET. THE RELATIVE IMPROVEMENT AND DEGRADATION ARE SHOWN IN CYAN AND RED, RESPECTIVELY.

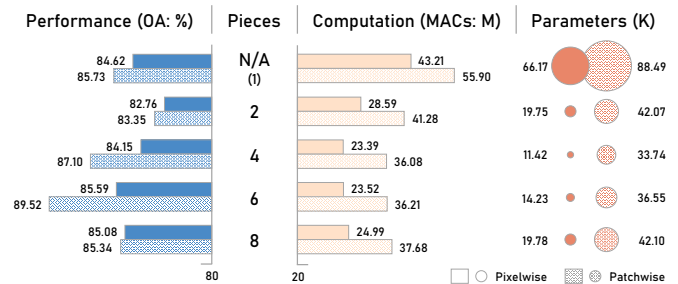| Module | Implementation | | | | | |
|---|---|---|---|---|---|---|
| | Pixelwise | | Patchwise | | | |
| GSSM | - | - | ✗ | ✓ | ✓ | ✓ |
| PSS | ✗ | ✓ | ✓ | ✗ | ✓ | ✓ |
| Mamba | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ |
| OA (%) ↑ | 84.62 | 85.59 (+0.97) | 86.74 (+2.12) | 85.73 (+1.11) | 84.40 (-0.22) | 89.52 (+4.90) |
| AA (%) ↑ | 85.56 | 87.23 (+1.67) | 88.40 (+2.84) | 87.55 (+1.99) | 86.57 (+1.01) | 90.50 (+4.94) |
| κ ↑ | 0.8330 | 0.8438 (+0.0108) | 0.8561 (+0.0231) | 0.8451 (+0.0121) | 0.8308 (-0.0008) | 0.8864 (+0.0534) |
| MACs (M) ↓ | 43.21 | 23.25 (-19.96) | 23.60 (-19.61) | 55.90 (+12.69) | 12.73 (-30.48) | 36.21 (-7.00) |
| Params (K) ↓ | 66.17 | 14.23 (-51.94) | 14.24 (-51.93) | 88.49 (+22.32) | 22.75 (-43.42) | 36.55 (-29.62) |



Fig. 8. Ablation results of our SpectralMamba on Houston2013 HS Dataset by varying the number of pieces for the proposed piece-wise sequential scanning.

in the former two prevalent sequence models by leveraging appropriate spatial-spectral merging, sequential scanning, and state space modeling.

We can draw clues from the visualization results that exhibit trends close to those of the above-mentioned quantitative results. For example, our SpectralMamba offers more reliable predictions for categories such as *road* and *highway* in the Houston2013 scene. As for the Augsburg scene, although CasRNN and SpectralFormer also show *forest* as good as ours, they typically confuse *residential area* with other classes. Moreover, the trajectory of the river and the shape of paddy fields, i.e., *water* category, can be clearly depicted by our SpectralMamba. The superiority of our methods is more obvious in the Longkou scene. Take *rice* area as an instance, in contrast to other methods that are easy to mix it with *water* or *mixed weed*, ours discover a region that is most regular and homogeneous used to plant *rice*. Although the labels in the Botswana scene are more sparse, we can still tell that our patchwise SpectralMamba generates classification maps with less salt and pepper noises-like patterns that are more likely in accord with the real distribution of ground objects.

### D. Ablation Studies

*1) Module Effectivity:* We first conduct ablation experiments on three key components in our SpectralMamba, that are, GSSM, PSS, and Mamba. We make a quantitative comparison by removing one or two constituents are summarized in Table VI. In most cases, an evident relative improvement can be observed on the basis of a single use of Mamba under pixelwise implementation. For example, our proposed scanning strategy brings nearly 1% and 3% accuracy improvement under pixelwise and patchwise settings, respectively. What's more, it can lead to 20M and 50K decreases in MACs and Params respectively, compared to the scenarios where only this module is switched off. We also observe a significant improvement in terms of 2.78% OA and 2.10% AA by activating the GSSM module in patchwise case. Note that the Mamba can still be deemed the most valuable module as more than 5% OA decrease occurs when removing it.

*2) Scanning Pieces:* The other ablation experiment worth conducting is to investigate the impact of varying the number of scanning pieces, i.e., $R$. The quantitative results are

illustrated in Fig. 8 by setting $R$ to values from 2 to 8 with an interval of 2. From the figure, we can observe that the baseline without using a scanning strategy can still attain an acceptable performance, while its computation and parameter scales are much higher than those using scanning. Specifically, in both pixelwise and patchwise cases, the MACs and Params metrics simultaneously show an initial decrease followed by an increase as $R$ increases, in contrast to the trend for performance. The underlying reason behind this phenomenon lies in that too many pieces inevitably widen the networks for SSM and increase the difficulty in pattern recognition from resulting spectral sequences with short lengths. Fortunately, by trading off all metrics, we can find the optimal $R = 6$ for the Houston2013 dataset, whose computation and parameters are maximally 46% and 79% less than w/o scanning while bringing 1% to 4% improvements in OA approximately.

## IV. CONCLUSION

The issues of spectral redundancy and spectral variability have long been plaguing humans from precisely sensing and perceiving the world via HS imaging. Existing intelligent methods for solving the HS image classification task, based on either conventional CNNs or trending sequence models like RNNs and Transformers, can hardly address these issues simultaneously and efficiently. Building upon the foundations laid by recent advances in SSMs, we establish the first HS data-oriented deep SSM model, i.e., SpectralMamba, by elaborately proposing PSS and GSSM to ease the sequentiality learning in the state domain and rectify the spectrum in the spatial-spectral domain, respectively. SpectralMamba enjoys a surprisingly simplified and lightweight architecture and achieves credibly superior HS image classification performance in most aspects. In the future, we will endeavor to unveil its potential for more critical applications using HS data under resource-constrained scenarios.

## REFERENCES

[1] J. M. Bioucas-Dias, A. Plaza, G. Camps-Valls, P. Scheunders, N. Nasrabadi, and J. Chanussot, "Hyperspectral remote sensing data analysis and future challenges," *IEEE Geoscience and Remote Sensing Magazine*, vol. 1, no. 2, pp. 6–36, 2013.

[2] B. Xi, J. Li, Y. Li, R. Song, D. Hong, and J. Chanussot, "Few-shot learning with class-covariance metric for hyperspectral image classification," *IEEE Transactions on Image Processing*, vol. 31, pp. 5079–5092, 2022.

[3] D. Hong, B. Zhang, X. Li, Y. Li, C. Li, J. Yao, N. Yokoya, H. Li, P. Ghamisi, X. Jia, A. Plaza, P. Gamba, J. A. Benediktsson, and J. Chanussot, "Spectralgpt: Spectral remote sensing foundation model," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. DOI:10.1109/TPAMI.2024.3362475.

[4] S. Chang and P. Ghamisi, "Changes to captions: An attentive network for remote sensing change captioning," *IEEE Transactions on Image Processing*, 2023.

[5] M. Ahmad, S. Shabbir, S. K. Roy, D. Hong, X. Wu, J. Yao, A. M. Khan, M. Mazzara, S. Distefano, and J. Chanussot, "Hyperspectral image classification—traditional to deep models: A survey for future prospects," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 968–999, 2022.

[6] 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 Transactions on Cybernetics*, vol. 51, no. 7, pp. 3602–3615, 2021.

[7] W. Ma, C. Gong, Y. Hu, P. Meng, and F. Xu, "The hughes phenomenon in hyperspectral classification based on the ground spectrum of grasslands in the region around qinghai lake," in *International Symposium on Photoelectronic Detection and Imaging 2013: Imaging Spectrometer Technologies and Applications*, vol. 8910, pp. 363–373, SPIE, 2013.

[8] J. Theiler, A. Ziemann, S. Matteoli, and M. Diani, "Spectral variability of remotely sensed target materials: Causes, models, and strategies for mitigation and robust exploitation," *IEEE Geoscience and Remote Sensing Magazine*, vol. 7, no. 2, pp. 8–30, 2019.

[9] G. Camps-Valls and L. Bruzzone, "Kernel-based methods for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 43, no. 6, pp. 1351–1362, 2005.

[10] G. Camps-Valls, D. Tuia, L. Bruzzone, and J. A. Benediktsson, "Advances in hyperspectral image classification: Earth monitoring with statistical learning methods," *IEEE Signal Processing Magazine*, vol. 31, no. 1, pp. 45–54, 2013.

[11] D. Hong, J. Hu, J. Yao, J. Chanussot, and X. X. Zhu, "Multimodal remote sensing benchmark datasets for land cover classification with a shared and specific feature learning model," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 178, pp. 68–80, 2021.

[12] D. Hong, N. Yokoya, N. Ge, J. Chanussot, and X. X. Zhu, "Learnable manifold alignment (lema): A semi-supervised cross-modality learning framework for land cover and land use classification," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 147, pp. 193–205, 2019.

[13] A. Samat, P. Du, S. Liu, J. Li, and L. Cheng, "$E^2$lms: Ensemble extreme learning machines for hyperspectral image classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 4, pp. 1060–1069, 2014.

[14] J. M. Haut, M. E. Paoletti, J. Plaza, J. Li, and A. Plaza, "Active learning with convolutional neural networks for hyperspectral image classification using a new bayesian approach," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 11, pp. 6440–6461, 2018.

[15] C. Li, B. Zhang, D. Hong, J. Yao, and J. Chanussot, "Lrr-net: An interpretable deep unfolding network for hyperspectral anomaly detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023. DOI:10.1109/TGRS.2023.3279834.

[16] Y. Xu, W. Yu, P. Ghamisi, M. Kopp, and S. Hochreiter, "Txt2img-mhn: Remote sensing image generation from text using modern hopfield networks," *IEEE Transactions on Image Processing*, 2023.

[17] A. Gu, K. Goel, and C. Ré, "Efficiently modeling long sequences with structured state spaces," *International Conference on Learning Representations*, 2022.

[18] R. E. Kalman, "A new approach to linear filtering and prediction problems," 1960.

[19] E. Nguyen, K. Goel, A. Gu, G. Downs, P. Shah, T. Dao, S. Baccus, and C. Ré, "S4nd: Modeling images and videos as multidimensional signals with state spaces," *Advances in Neural Information Processing Systems*, vol. 35, pp. 2846–2861, 2022.

[20] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.

[21] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7132–7141, 2018.

[22] Y. Liu, Y. Tian, Y. Zhao, H. Yu, L. Xie, Y. Wang, Q. Ye, and Y. Liu, "Vmamba: Visual state space model," *arXiv preprint arXiv:2401.10166*, 2024.

[23] R. Hang, Q. Liu, D. Hong, and P. Ghamisi, "Cascaded recurrent neural networks for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 8, pp. 5384–5394, 2019.

[24] D. Hong, Z. Han, J. Yao, L. Gao, B. Zhang, A. Plaza, and J. Chanussot, "Spectralformer: Rethinking hyperspectral image classification with transformers," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022. Doi: 10.1109/TGRS.2021.3130716.

[25] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-form image inpainting with gated convolution," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4471–4480, 2019.

[26] J. Yao, B. Zhang, C. Li, D. Hong, and J. Chanussot, "Extended vision transformer (exvit) for land use and land cover classification: A multimodal deep learning framework," *IEEE Transactions on Geoscience and Remote Sensing*, 2023. 10.1109/TGRS.2023.3284671.

[27] D. Hong, B. Zhang, H. Li, Y. Li, J. Yao, C. Li, M. Werner, J. Chanussot, A. Zipf, and X. X. Zhu, "Cross-city matters: A multimodal remote sensing benchmark dataset for cross-city semantic segmentation using high-resolution domain adaptation networks," *Remote Sensing of Environment*, vol. 299, p. 113856, 2023.

[28] J. Yao, X. Cao, D. Hong, X. Wu, D. Meng, J. Chanussot, and Z. Xu, "Semi-active convolutional neural networks for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022. DOI:10.1109/TGRS.2022.3206208.

[29] Y. Zhong, X. Hu, C. Luo, X. Wang, J. Zhao, and L. Zhang, "Whu-hi: Uav-borne hyperspectral with high spatial resolution (h2) benchmark datasets and classifier for precise crop identification based on deep convolutional neural network with crf," *Remote Sensing of Environment*, vol. 250, p. 112012, 2020.

[30] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012.

[31] J. Yao, D. Hong, H. Wang, H. Liu, and J. Chanussot, "Ucsl: Towards unsupervised common subspace learning for cross-modal image classification," *IEEE Transactions on Geoscience and Remote Sensing*, 2023. Doi: 10.1109/TGRS.2023.3282951.

[32] N. Audebert, B. Le Saux, and S. Lefèvre, "Deep learning for classification of hyperspectral data: A comparative review," *IEEE Geoscience and Remote Sensing Magazine*, vol. 7, no. 2, pp. 159–173, 2019.