# Hybrid Spectral Denoising Transformer with Guided Attention

Zeqiang Lai[1],    Chenggang Yan[2],    Ying Fu[1†]

[1]Beijing Institute of Technology    [2] Hangzhou Dianzi University

{laizeqiang, fuying}@bit.edu.cn    cgyan@hdu.edu.cn

## Abstract

*In this paper, we present a Hybrid Spectral Denoising Transformer (HSDT) for hyperspectral image denoising. Challenges in adapting transformer for HSI arise from the capabilities to tackle existing limitations of CNN-based methods in capturing the global and local spatial-spectral correlations while maintaining efficiency and flexibility. To address these issues, we introduce a hybrid approach that combines the advantages of both models with a Spatial-Spectral Separable Convolution (S3Conv), Guided Spectral Self-Attention (GSSA), and Self-Modulated Feed-Forward Network (SM-FFN). Our S3Conv works as a lightweight alternative to 3D convolution, which extracts more spatial-spectral correlated features while keeping the flexibility to tackle HSIs with an arbitrary number of bands. These features are then adaptively processed by GSSA which performs 3D self-attention across the spectral bands, guided by a set of learnable queries that encode the spectral signatures. This not only enriches our model with powerful capabilities for identifying global spectral correlations but also maintains linear complexity. Moreover, our SM-FFN proposes the self-modulation that intensifies the activations of more informative regions, which further strengthens the aggregated features. Extensive experiments are conducted on various datasets under both simulated and real-world noise, and it shows that our HSDT significantly outperforms the existing state-of-the-art methods while maintaining low computational overhead. Code is at* https://github.com/Zeqiang-Lai/HSDT.

## 1. Introduction

Hyperspectral image (HSI) provides substantially more abundant spectral information than the ordinary color image, which makes it especially utilitarian in the field of remote sensing [3, 4], biometric authentication [60], detection [46], and geological science [19, 57]. Nevertheless, lim-
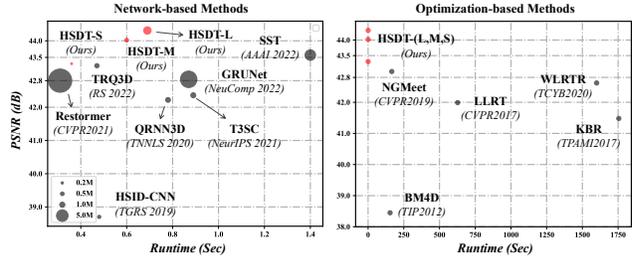
---

† Corresponding Author.



Figure 1: Our method achieves state-of-the-art performance while maintaining low computational overhead.

ited by imaging techniques, most existing HSI cameras still suffer from various types of noise that might degrade the performance of their applications, which urges the development of robust HSI denoising algorithms.

Motivated by the intrinsic properties of HSI, traditional HSI denoising approaches [70, 22] often exploit the optimization schemes with priors, *e.g.*, low rankness [75, 65], non-local similarities [45, 47], spatial-spectral correlation [52], and global correlation along the spectrum [66]. Whilst offering appreciable performance, the efficacy of these methods is largely dependent on the degree of similarity between the handcrafted priors and the real-world noise model, and these methods are often challenging to accelerate with modern hardwares due to the complex processing pipelines. Recent HSI denoising methods based on Convolutional Neural Network (CNN) [71, 35, 5] get rid of handcrafted regularizations with learning-based prior and often run faster with graphic accelerators and machine-learning frameworks [50]. However, these methods are still insufficient for exploring the characteristics of HSI, *e.g.*, global and local spectral-spatial correlations. For example, HSID-CNN [71] only considers the correlations between several adjacent spectral bands. QRNN3D [66] and GRUNet [35] model the global spectral correlations with quasi-recurrent units [6] but suffer from the problem of vanished correlations for long-range separate bands due to the recurrent multiplications of merging weights. Besides, recent methods [66, 74] tend to use 3D convolution to explore the local spectral-spatial correlations while maintaining the flexibil-

ity to handle different HSIs. This strategy, however, introduces substantially unwanted computation and parameters.

Starting from natural language processing [61], transformer architectures [17, 2] have recently been applied to various vision tasks including color image restoration [40, 72] and HSI processing [37, 7, 49]. With the multi-head self-attention of transformer, these methods enjoy stronger capabilities of capturing non-local similarity and long-range dependency over aforementioned CNN-based methods. Despite of that, they are still suboptimal and inflexible for diverse HSIs. On the one hand, existing attentions for HSIs apply along either spatial [37, 49] or 2D feature channel [37, 7] dimensions, which could introduce quadratic complexities or break down the structured spectral dependency. On the other hand, their 2D architectural designs also make their models specifically bound to one type of HSI, *e.g.*, HSI with 31 bands, and separate models have to be trained for other types, *e.g.*, HSI with 210 bands. This can be problematic since the amount of available datasets is unevenly distributed for different HSIs. Finally, HSIs often exhibit beneficial fixed structures, *e.g.*, relative intensity correlations of different bands for objects. Direct transfer of existing transformer blocks without considering this fact might lead to suboptimal performance for HSI denoising.

In this paper, we propose a novel Hybrid Spectral Denoising Transformer (HSDT) that effectively integrates the local spectral-spatial inductive bias of the convolution and the long-range spectral dependency modeling ability of the transformer. Unlike previous HSI transformers [37, 7, 49], HSDT is designed to be effective and flexible for handling diverse HSIs in a single model, which results in a variety of benefits, *e.g.*, the ability to jointly utilize HSIs with different numbers of bands for more sufficient training in data-insufficiency scenarios. To achieve it, (a) we first introduce a parallel Spectral-Spatial Separable Convolution (S3Conv) unit that efficiently extracts more spatial-spectral meaningful features than previous 3D [66] and separable convolutions [16]. (b) With the stronger local spectral-spatial inductive bias, the extracted features are then processed by a newly proposed Guided Spectral Self-Attention (GSSA) that performs the global self-attention along the 3D spectral rather than spatial [49] or 2D spectral/channel dimensions [37, 7] to selectively aggregate the information across different bands. This not only enriches our model with powerful capabilities for identifying long-range spectral correlations but also makes our model free of inflexibility of 2D spectral SA for dealing with different HSIs, quadratic complexity of spatial SA [17], the issue of vanished long-range dependency of QRNN [66]. (c) Besides, inspired the relatively stationary global patterns and statistics of features of different HSI bands, as shown in Fig. 2, we propose to enhance our GSSA with a set of learnable queries that encode the global spectral statistics for each band. We alterna-
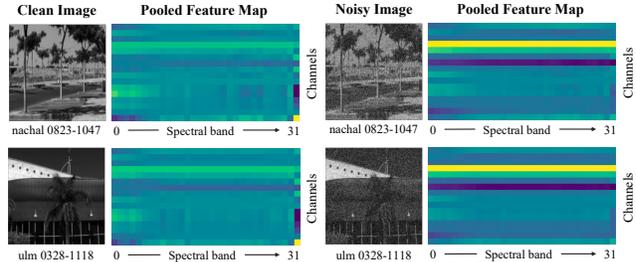


Figure 2: Visualization of feature maps of different bands by performing global average pooling on spatial locations. It can be observed that despite the difference of images, the pooled feature maps of clean and noisy images share different common patterns, which might be helpful for identifying, *e.g.*, cleaner bands, for denoising noisy bands.

tively switch between self-attention among spectral bands and cross-attention between spectral bands and learnable queries during the training, so that we can guide the GSSA to pay attention to features that are more discriminative and beneficial for denoising while keeping the flexibility. (d) Moreover, we propose a Self-Modulated Feed-Forward Network (SM-FFN) with a novel SM-branch to further strengthen the aggregated features of more informative regions. Extensive experiments on various datasets under different noise show that our HSDT consistently outperforms the existing state-of-the-art (SOTA) methods while maintaining low computational overhead, as shown in Fig. 1.

In summary, our contributions are that,

- We present HSDT, a 3D hybrid spectral denoising transformer that effectively captures the local spatial-spectral features and long-range global spectral correlations.
- We introduce GSSA guided by a set of learnable queries that encode the global statistics of HSIs, which models long-range spectral correlations along 3D spectrum instead of previous 2D spectral/channel dimensions.
- We propose SM-FFN with a novel self-modulated branch for driving the model to pay attention to more informative regions, along with a S3Conv for extracting spatial-spectral meaningful features.

## 2. Related Works

### 2.1. Hyperspectral Image Denoising

Traditional approaches for HSI denoising usually formulate the task as an optimization problem, which is solved by imposing different types of handcrafted regularizations [70, 48, 58, 45, 36]. Among these optimization-based methods, non-local similarity [45] has been widely utilized for its ability to integrate the image patches across the spectral and spatial locations. To reduce the computational burden, global spectral low-rank correlation [65, 58, 75] has also been heavily studied. Besides, different enhanced total

variation priors [51, 63, 70] are also adopted by considering the smoothness of local image patches. Though these methods could achieve favorable performance, most of them are computationally inefficient and can only address the noise satisfying the required assumptions, *e.g.*, Gaussian noise.

Meanwhile, recent works [71, 66, 35, 34] tend to exploit deep learning to learn denoising mapping purely in a data-driven manner. For most of these methods, the encoder-decoder U-Net [56, 16, 35, 66] architecture is the prominent choice due to its effectiveness for retaining both high- and low-level multi-scale representations. Residual learning [21] is also widely adopted to reduce learning difficulties from different perspectives, *e.g.*, residual image [9] and residual features [71, 35]. To consider the properties of HSIs, *e.g.*, spatial-spectral correlations, QRNN3D [66] proposes to use 3D convolution and quasi-recurrent unit [6]. Our work adopts techniques, residual learning, 3D convolution, and U-shape architecture, but our blocks, *e.g.*, S3Conv is more efficient than 3D convolution, and our GSSA could prevent vanished correlations for long-range spectral bands of QRU [65]. Separable convolution [28] is first introduced to replace 2D convolution. For HSI denoising, it is also adopted in [29, 16, 23] to reduce computational burden with similar motivations as ours. As a new alternative, our S3Conv is separable 3D instead of 2D convolution [29] and more effective than previous 3D variant [16].

## 2.2. Vision Transformers

Transformer [61] has been first introduced as a parallel and purely attention-based alternative for recurrent neural networks [15, 26] in the literature of natural language processing. Though it is originally designed for modeling text, recent works such as ViT [17] and DeiT [59], have successfully transferred the transformer for high-level vision tasks. Recognizing the powerful representation abilities, this architecture is also expeditiously adapted for low-level tasks [40, 12, 39], such as natural image denoising. Among these methods, one of the key problems they attempt to overcome is the quadratic complexity of the Self-Attention (SA) mechanism in the transformer. To address it, SwinIR [40] is proposed as an adaption of Swin transformer [43] that replaces global attention with a more efficient shift-window-based attention. Similarly, Uformer [64] performs attention over non-overlapped patches and adopts U-Net architecture [56] to further increase efficiency. From a different perspective, Restormer [72] explores self-attention along the feature channels to realize the linear complexity. Despite their superior performance for various natural image restoration tasks, direct transfer of them for HSI can result in performance degradations since none of them consider the properties of HSI. Instead, our HSDT introduce S3Conv and GSSA that can extract more spectral correlated features, which is more suitable for HSI.

For HSI processing, the applications of transformer previously concentrate more on the classification [27, 25, 54] and spectral reconstruction [7, 41], but also recently extend to the HSI denoising [49, 37, 38, 68, 13]. For example, to exploit the spatial attention, TRQ3DNet [49] combines the QRU [66] and Uformer [64] block, while SST [37] and DSTrans [68] employ the Swin transformer block. To consider spectral correlations, most existing HSI transformers, including SST [37], DSTrans [68], and MST [7], utilize a 2D spectral/channel attention similar to Restormer [72]. Basically, these methods have three major drawbacks including (i) the spatial attention, *e.g.*, Uformer [64] and Swin [43], neither can model spectral relationships nor be computationally friendly for HSI, (ii) spectral attention as Restormer [72] and MST [7] essentially perform attention on feature channel dimension $C$ of 2D data $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ and $C$ could change for different layers, which makes them no longer spectral attention. (iii) Built based on 2D architectures for color images, these methods are usually not flexible for handle different HSIs in a single model. On the contrary, our GSSA is both effective for capturing global spectral correlations and flexible for diverse HSIs as we consider exact spectral attention along spectral $D$ dimension of 3D data $\mathbf{X} \in \mathbb{R}^{H \times W \times D \times C}$. Concurrent to our work, Hider [13] also consider 3D spectral attention. Different from it, our GSSA employs a simple pooling strategy instead of conv-reshape strategy as Restormer [72] to compute the query and key as well as learnable query, which makes it much more effective and efficient.

## 3. Hybrid Spectral Denoising Transformer

In this section, we present Hybrid Spectral Denoising Transformer (HSDT), a unified model for hyperspectral image denoising with an arbitrary number of bands. To achieve it effectively, our HSDT introduces several key designs, including (i) a powerful and lightweight spectral-spatial separable convolution as an alternative to 3D convolution, (ii) a guided spectral self-attention piloted by a set of learnable queries, and (iii) a self-modulated feed-forward network with an adaptive self-modulated branch.

The overall architecture of HSDT follows a U-shaped encoder-decoder with skip-connections [56], which is depicted in Fig. 3(a). Such hierarchical multi-scale design not only reduces the computational burden but also increases the receptive fields, which is different from conventional plain transformers [40, 17]. In general, HSDT is built by stacking a series of transformer blocks as,

$$\hat{\mathbf{X}} = \text{BN}(\text{S3Conv}(\mathbf{X})) \tag{1}$$

$$\mathbf{Y} = \text{SM-FFN}(\text{GSSA}(\hat{\mathbf{X}}) + \hat{\mathbf{X}}) \tag{2}$$

where $\mathbf{X} \in \mathbb{R}^{H \times W \times D \times C}, \mathbf{Y} \in \mathbb{R}^{\hat{H} \times \hat{W} \times D \times \hat{C}}$ are the input and output feature maps, $H, W$ denote spatial size, $D$ de-
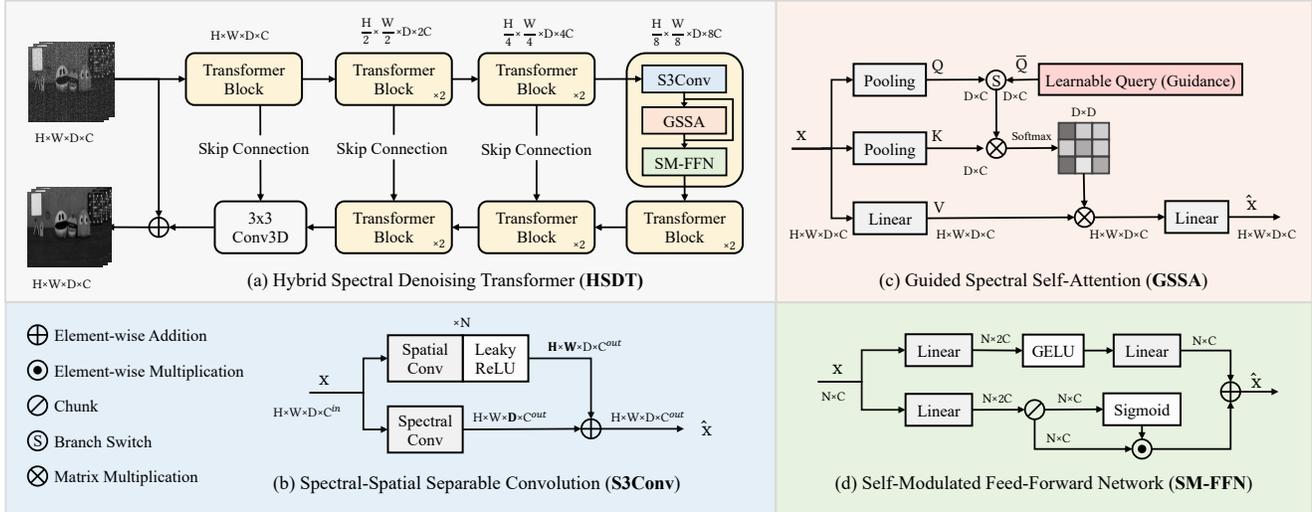
Figure 3: **Our HSDT architecture.** We adopt a hierarchical multiscale encoder-decoder (a) with each building transformer block stacks (d) a Spectral-Spatial Separable Convolution, (b) a Guided Spectral Self-Attention, and (c) a Self-Modulated Feed-Forward Network sequentially. We append batch normalization after each convolution and predict the residual image.

notes the number of spectral bands, $C$ denotes the number of feature maps, and BN denotes batch normalization [30]. More specifically, given the input noisy HSI, it is first projected into low-level features through a head transformer block and then passed through several transformer blocks to fuse the features along both spatial and spectral dimensions. The residual connection is added to the final output and the input noisy image. We use trilinear interpolations for upsampling and adopt additive skip connections in all levels of transformer blocks. Next, we illustrate the details of each network component.

## 3.1. Spectral-Spatial Separable Convolution

Modern HSI cameras are capable of capturing images with an exceedingly larger number of spectral bands, *e.g.*, 31 for Specim PS Kappa DX4 [20], and 224 for AVIRIS sensor [32]. However, it is still difficult to collect a large amount of training data for each device due to the complex and time-consuming imaging processes. Hence, it is one of the major concerns for an HSI denoising method if it can handle images with a different number of bands using a single model so that the dataset collected by different devices can be jointly used for more sufficient training.
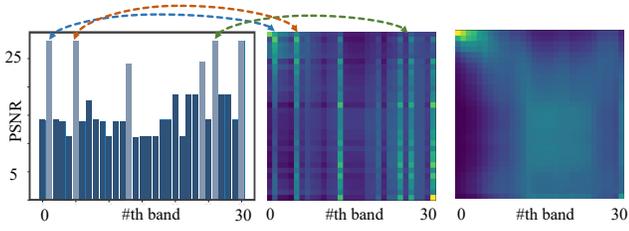
To address the issue, the sliding-window denoising strategy [71] and Conv3D [66] have been adopted to build band-flexible networks. However, their performance and efficiency are still limited by the local receptive field and heavy computation. In this work, we seek to preserve band flexibility and augment our denoising transformer with inductive bias while avoiding splitting HSI into windows. For this, we propose Spectral-Spatial Separable Convolution (S3Conv), a more lightweight and powerful variant of Conv3D that

parallel applies spectral and spatial convolution.

In detail, standard Conv3D filters spectral-spatial correlated features but introduces a heavy burden on the number of parameters and computation. To alleviate the computational burden, our S3Conv decouples the Conv3D into two parallel branches, which separately process the inputs along the spatial and spectral dimensions, as illustrated in Fig. 3(b). The spatial convolution extracts features with 2D filters for each band and the spectral convolution applies $1 \times 1$ projection to correlate the spectral information of all bands. To obtain the final features, we combine the output from two branches through element-wise addition. Instead of a direct extension of 2D separable Conv for 3D data, our S3Conv is spatial-spectral instead of spatial-channel separable, which makes it more suitable for extracting spectral correlation. Besides, our S3Conv is parallel instead of sequentially separated, which makes it easily accelerated.

## 3.2. Guided Spectral Self-Attention

Despite the spatial self-attention [40, 56] improves the model performance by considering spatial interaction and non-local similarities, it is computationally demanding and might be difficult to deal with HSIs with a different number of bands. In this work, we propose an efficient Guided Spectral Self-Attention (GSSA) that applies 3D SA along the spectral than spatial nor channel dimensions. Our GSSA is intuitively supported by the spectral correlations of HSI and has linear complexity and long-range relation modeling abilities. This makes our model extremely more powerful at locating the informative regions to assist the denoising than existing spectral integration techniques [66, 35].

(a) SNR of each band     (b) Attn w/ LQ     (c) Attn w/o LQ

Figure 4: The Learnable Query (LQ) guides the model to pay attention (Attn) to more informative bands with higher signal-to-noise ratio (SNR).



(a) Noisy HSI     (b) Feature map     (c) Modulation Weight

Figure 5: Our Self-Modulated FFN amplifies the features in high-information-density regions, *e.g.*, edges, with an element-wise modulation weight. Yellow regions denote higher weight.

**Spectral Self Attention.** GSSA takes 3D feature maps from previous S3Conv, *i.e.*, $\mathbf{X} \in \mathbb{R}^{H \times W \times D \times C}$, and performs 3D attention on $D$ dimension, instead of $C$ dimension of $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ of previous 2D spectral attention [7, 37]. To perform attention, we first convert it into query, key, and value. Unlike conventional attention block [61], only value $\hat{\mathbf{V}} \in \mathbb{R}^{H \times W \times D \times C}$ is linearly projected from $\mathbf{X}$ in GSSA. Linear projections for query and key are not necessary according to our experiments, so we omit it for simplicity. Instead, we directly perform the global average pooling on input $\mathbf{X}$ along the spatial dimensions to obtain the global features of each band, *i.e.*, $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{D \times C}$. This pooling strategy is not only parameter-free but it also differs from previous reshape strategy [72, 7] that causes larger computation in the following dot-product attention.
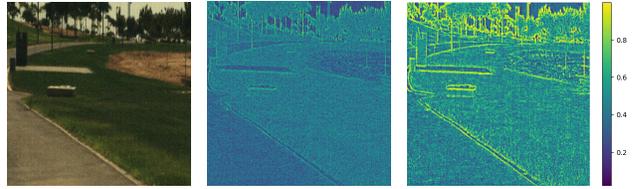
Then, the transposed attention map $\mathbf{A}$ in the shape of $\mathbb{R}^{D \times D}$ is obtained via dot-product between key $\mathbf{K}$ and query $\mathbf{Q}$ with softmax normalization. Finally, we multiply the attention map with the value $\hat{\mathbf{V}}$ to dynamically select the essential features across the spectrum for each band, *i.e.*,

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \hat{\mathbf{V}}) = \hat{\mathbf{V}} \cdot \text{Softmax}(\mathbf{K} \cdot \mathbf{Q}) \quad (3)$$

$$\hat{\mathbf{X}} = \mathbf{W}\,\text{Attention}(\mathbf{Q}, \mathbf{K}, \hat{\mathbf{V}}) + \mathbf{X}. \quad (4)$$

To better transform the features, we perform another linear projection $\mathbf{W} \in R^{C \times C}$ on the fused features. Beside, we learn the residual features to stabilize the training.

**Learnable Query.** Different spectral bands of HSIs often exhibit some fixed relative relationships due to the physical spectral constraints. Inspired by this, we investigate the global spectral feature maps of different HSIs, as shown in Fig. 2. Perhaps surprisingly, we found that clean and noisy HSIs exhibit clearly different shared patterns, which indicates the possibilities to utilize them for identifying more useful bands in GSSA. To achieve it, we propose to model these global statistics of each band with a set of learnable queries, as shown in Fig. 3(c). We perform cross-attention (CA) between the pooled keys and learnable queries to filter out most discriminative features for each band. The learnable queries are jointly trained with the other part of the

model and not restricted to the input feature maps so that they can better fit the statistics of clean HSIs. As a results, the attention with learnable queries produce a remarkably better attention map with clear interpretability that identifies more informative bands with higher SNR to assist the denoising of more noisy ones, as shown in Fig. 4.

**Alternative Training Strategy.** The learnable queries are very useful to generate more discriminative attention, but the number of queries has to be predefined to the number of bands of the training dataset. It would disable the model from handling HSIs with a different number of bands simultaneously. To address the issue, we propose an alternative training strategy, in which we randomly switch between SA and CA with learnable queries during the training. Since the pooling operation in SA also captures global information to some extent, the proposed learnable query with this strategy can then be viewed as guidance leading the training of self-attention on more descriptive bands, which consequently leads to better performance.

### 3.3. Self-Modulated Feed-Forward Network

Feed-Forward Network (FFN) is one of the most essential parts of transformer architectures, and it has been reported that it might be the key to construct the meta structure of transformer than SA [69]. Traditional FFN [61] processes the output features from the SA layer with two linear projections and a non-linear activation between them.

In this work, we propose Self-Modulated FFN (SM-FFN) with a fundamental augmentation of the vanilla FFN using self-modulation. As shown in Fig. 3(d), we first expand the input features channels $\mathbf{X} \in \mathbb{R}^{H \times W \times D \times C}$ with a scale factor of 2 through a linear projection, $\mathbf{Y} = \mathbf{W_3}\mathbf{X}, \mathbf{W_3} \in \mathbb{R}^{C \times 2C}$, then we split (also known as *chunk* operation) the expanded features $\mathbf{Y}$ into two parts $\mathbf{F}, \mathbf{W} \in \mathbb{R}^{H \times W \times D \times C}$. We treat one part $\mathbf{F}$ as the candidate features and another part $\mathbf{W}$ after the sigmoid normalization as an element-wise modulation weight. With the vanilla FFN branch, our SM-FFN can be described as,

$$\text{SM-Branch}(\mathbf{X}) = \mathbf{F} \odot \text{Sigmoid}(\mathbf{W}), \quad (5)$$
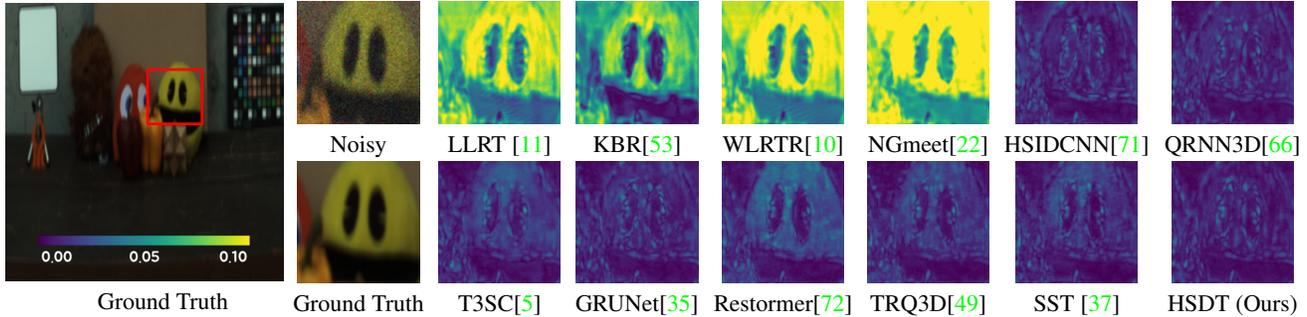
Figure 6: Visual comparison with error maps for simulated Gaussian noise removal on ICVL dataset. ($\sigma = 70$)

| Methods | Params (M) | Runtime (s) | $\sigma = 30$ | | | $\sigma = 50$ | | | $\sigma = 70$ | | | $\sigma = $ blind | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | PSNR | SSIM | SAM | PSNR | SSIM | SAM | PSNR | SSIM | SAM | PSNR | SSIM | SAM |
| Noisy | - | - | 18.59 | 0.110 | 0.807 | 14.15 | 0.046 | 0.991 | 11.23 | 0.025 | 1.105 | 17.34 | 0.114 | 0.859 |
| LLRT [11] | - | 627 | 41.99 | 0.967 | 0.056 | 38.99 | 0.945 | 0.075 | 37.36 | 0.930 | 0.087 | 40.97 | 0.956 | 0.064 |
| KBR [53] | - | 1755 | 41.48 | 0.984 | 0.088 | 39.16 | 0.974 | 0.100 | 36.71 | 0.961 | 0.113 | 40.68 | 0.979 | 0.080 |
| WLRTR [10] | - | 1600 | 42.62 | 0.988 | 0.056 | 39.72 | 0.978 | 0.073 | 37.52 | 0.967 | 0.095 | 41.66 | 0.983 | 0.064 |
| NGmeet [22] | - | 166 | 42.99 | 0.989 | 0.050 | 40.26 | 0.980 | 0.059 | 38.66 | 0.974 | 0.067 | 42.23 | 0.985 | 0.053 |
| HSID-CNN [71] | 0.40 | 0.48 | 41.72 | 0.987 | 0.067 | 39.39 | 0.980 | 0.083 | 37.77 | 0.972 | 0.096 | 40.95 | 0.984 | 0.072 |
| QRNN3D [66] | 0.86 | 0.44 | 42.22 | 0.988 | 0.062 | 40.15 | 0.982 | 0.074 | 38.30 | 0.974 | 0.094 | 41.37 | 0.985 | 0.068 |
| T3SC [5] | 0.83 | 0.95 | 42.36 | 0.986 | 0.079 | 40.47 | 0.980 | 0.087 | 39.05 | 0.974 | 0.096 | 41.52 | 0.983 | 0.085 |
| GRUNet [35] | 14.2 | 0.87 | 42.84 | 0.989 | 0.052 | 40.75 | 0.983 | 0.062 | 39.02 | 0.977 | 0.080 | 42.03 | 0.987 | 0.057 |
| Restormer [72] | 26.2 | 0.31 | 42.80 | 0.990 | 0.062 | 41.03 | 0.985 | 0.062 | 39.62 | 0.980 | 0.069 | 41.99 | 0.987 | 0.064 |
| TRQ3D [49] | 0.68 | 0.47 | 43.25 | 0.990 | 0.046 | 41.30 | 0.985 | 0.053 | 39.86 | 0.980 | 0.061 | 42.47 | 0.988 | 0.054 |
| SST [37] | 4.14 | 1.40 | 43.57 | 0.991 | 0.045 | 41.41 | 0.986 | 0.052 | 39.89 | 0.980 | 0.058 | 42.81 | 0.988 | 0.047 |
| SERT [38] | 1.91 | 0.44 | 43.99 | 0.991 | 0.043 | 41.82 | 0.986 | 0.051 | 40.28 | 0.981 | 0.059 | 43.20 | 0.988 | 0.047 |
| HSDT-S(Ours) | 0.13 | 0.36 | 43.31 | 0.990 | 0.047 | 41.16 | 0.985 | 0.055 | 39.66 | 0.980 | 0.064 | 42.57 | 0.988 | 0.051 |
| HSDT-M(Ours) | 0.52 | 0.60 | 44.02 | 0.991 | 0.041 | 41.82 | 0.986 | 0.049 | 40.33 | 0.981 | 0.055 | 43.32 | 0.989 | 0.045 |
| HSDT-L(Ours) | 2.09 | 0.69 | **44.31** | **0.992** | **0.041** | **42.09** | **0.987** | **0.048** | **40.59** | **0.982** | **0.054** | **43.59** | **0.989** | **0.044** |

Table 1: Gaussian denoising results on ICVL. *Blind* denotes Gaussian noise with random sigma (ranged from 30 to 70). HSDT-M doubles the width/channels of HSDT-S. HSDT-L increases the network depth of HSDT-M with an additional transformer block. The runtime is averaged over 10 runs.

$$\text{SM-FFN}(\mathbf{X}) = \mathbf{W_1}(\text{GELU}(\mathbf{W_2 X})) + \text{SM-Branch}(\mathbf{X}), \quad (6)$$

where $\mathbf{W_2} \in \mathbb{R}^{C \times 2C}, \mathbf{W_1} \in \mathbb{R}^{2C \times C}$. Intuitively, our SM branch is designed and works as a soft max-pooling that amplifies the activation of regions with higher information density. As reflected in Fig. 5, this makes our network more robust by emphasizing the features in regions that are more important for improving the reconstruction quality, *e.g.*, edges and corners.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets.** We use three public HSI datasets with natural scenes, *i.e.*, ICVL [1], CAVE [67], RealHSI [74], and one remotely sensed dataset Urban [32]. ICVL and CAVE are unpaired datasets with solely clean images. RealHSI and Urban are paired and unpaired datasets with real-world noise. ICVL and CAVE share the same spectral resolution of 31, while RealHSI and Urban have different spectral res-

olutions of 34 and 210. We train and test the models on ICVL under simulated noise. For other datasets, we test the zero-shot performance with the models trained on ICVL. We use 100 images from ICVL for training, and we randomly select the main $512 \times 512$ regions of 50 images from ICVL, and the entire regions of 12, 15 and 1 images from CAVE, RealHSI, and Urban for testing.

**Noise Settings.** We consider simulated Gaussian and complex noise as well as real-world noise provided by RealHSI [74] and Urban [32]. The Gaussian noise is sampled from zero-mean i.i.d Gaussian distribution with different variances and we evaluate different methods on 30, 50, 70, and blind (range from 10 to 70) noise strengths. Following [66], the complex noise is composed of non-i.i.d Gaussian noise and one or several types of complex noise, including, stripe, deadline, and impulse noise.

**Implementation Details.** We implement the proposed method with PyTorch [50]. The network is trained by optimizing the mean-square-root error between predicted im-
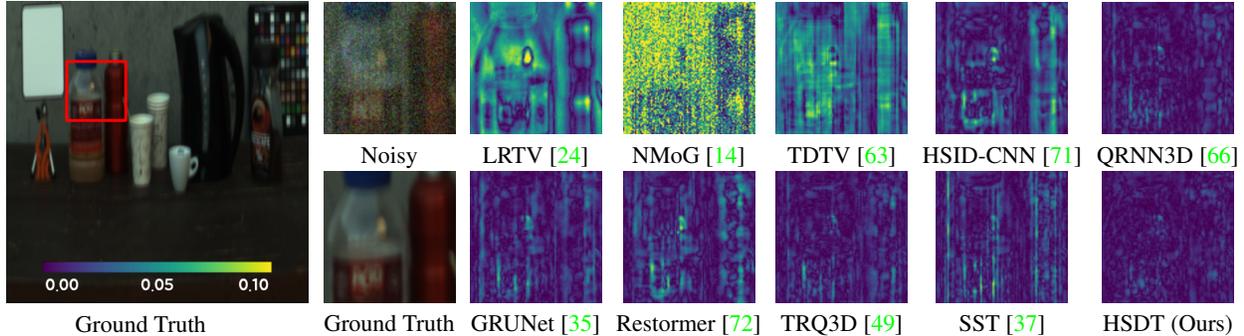
Figure 7: Visual comparison with error maps for simulated complex noise removal on ICVL dataset. (*mixture* case)

Top row labels: Noisy, LRTV [24], NMoG [14], TDTV [63], HSID-CNN [71], QRNN3D [66]
Bottom row labels: Ground Truth, GRUNet [35], Restormer [72], TRQ3D [49], SST [37], HSDT (Ours)
Left image label: Ground Truth

| Methods | non-iid | | | stripe | | | deadline | | | impulse | | | mixture | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | SAM | PSNR | SSIM | SAM | PSNR | SSIM | SAM | PSNR | SSIM | SAM | PSNR | SSIM | SAM |
| Noisy | 18.25 | 0.168 | 0.898 | 17.80 | 0.159 | 0.910 | 17.61 | 0.155 | 0.917 | 14.80 | 0.114 | 0.926 | 14.08 | 0.099 | 0.944 |
| LRTV [24] | 33.62 | 0.905 | 0.077 | 33.49 | 0.905 | 0.078 | 32.37 | 0.895 | 0.115 | 31.56 | 0.871 | 0.242 | 30.47 | 0.858 | 0.287 |
| NMoG [14] | 34.51 | 0.812 | 0.187 | 33.87 | 0.799 | 0.265 | 32.87 | 0.797 | 0.276 | 28.60 | 0.652 | 0.486 | 27.31 | 0.632 | 0.513 |
| TDTV [63] | 38.14 | 0.944 | 0.075 | 37.67 | 0.940 | 0.081 | 36.15 | 0.930 | 0.099 | 36.67 | 0.935 | 0.094 | 34.77 | 0.919 | 0.113 |
| HSID-CNN [71] | 40.14 | 0.984 | 0.067 | 39.53 | 0.983 | 0.720 | 39.49 | 0.983 | 0.071 | 36.69 | 0.959 | 0.156 | 35.36 | 0.954 | 0.169 |
| QRNN3D [66] | 42.79 | 0.978 | 0.052 | 42.35 | 0.976 | 0.055 | 42.23 | 0.976 | 0.056 | 39.23 | 0.945 | 0.109 | 38.25 | 0.938 | 0.108 |
| T3SC [5] | 41.28 | 0.987 | 0.065 | 40.85 | 0.986 | 0.072 | 39.54 | 0.983 | 0.096 | 36.06 | 0.952 | 0.203 | 34.48 | 0.946 | 0.228 |
| GRUNet [35] | 42.89 | 0.992 | 0.047 | 42.39 | 0.991 | 0.050 | 42.11 | 0.991 | 0.050 | 40.70 | 0.985 | 0.067 | 38.51 | 0.980 | 0.081 |
| Restormer [72] | 40.81 | 0.987 | 0.050 | 40.49 | 0.986 | 0.052 | 39.89 | 0.986 | 0.055 | 37.60 | 0.972 | 0.746 | 36.21 | 0.968 | 0.752 |
| TRQ3D [49] | 43.34 | 0.992 | 0.042 | 43.05 | 0.992 | 0.043 | 42.70 | 0.992 | 0.045 | 41.22 | 0.983 | 0.075 | 40.27 | 0.983 | 0.075 |
| SST [37] | 43.43 | 0.993 | 0.042 | 43.02 | 0.992 | 0.044 | 42.95 | 0.992 | 0.044 | 41.27 | 0.985 | 0.064 | 39.19 | 0.983 | 0.067 |
| SERT [38] | 44.10 | 0.993 | 0.038 | 43.80 | 0.993 | 0.040 | 43.55 | 0.993 | 0.041 | 40.56 | 0.977 | 0.080 | 39.25 | 0.977 | 0.079 |
| HSDT-S(Ours) | 43.46 | 0.992 | 0.044 | 43.13 | 0.992 | 0.047 | 42.97 | 0.991 | 0.047 | 41.11 | 0.976 | 0.110 | 40.22 | 0.974 | 0.116 |
| HSDT-M(Ours) | 44.56 | 0.994 | 0.039 | 44.29 | 0.993 | 0.040 | 44.18 | 0.993 | 0.040 | 41.28 | 0.977 | 0.101 | 40.46 | 0.976 | 0.106 |
| HSDT-L(Ours) | 44.94 | 0.994 | 0.036 | 44.69 | 0.994 | 0.038 | 44.55 | 0.994 | 0.037 | 42.02 | 0.980 | 0.101 | 41.07 | 0.980 | 0.101 |

Table 2: Complex denoising results on ICVL. *non-iid* denotes Non i.i.d Gaussian noise. *stripe*, *deadline*, *impulse* denote the combination of *non-iid* and corresponding complex noise. *mixture* denotes the combination of all the mentioned noise.

ages and the ground truth. We adopt Adam [33] optimizer with a multi-step learning scheduler whose initial learning rate is set to $1 \times 10^{-3}$ and decayed by a factor when it reaches the predefined milestones. Following [66], we use a multi-stage training strategy to train models for Gaussian and complex noise. The learning rate warmup is used between training stages. The batch size and training patch size are set to 16 and $64 \times 64$. All our models for Gaussian noise are trained for 80 epochs, and the models for complex noise are obtained by another 30 epochs of fine-tuning. We use pretrained complex denoising models for real noise.

### 4.2. Main Results

**Gaussian Denoising.** We evaluate our method against SOTA optimization-based methods (*i.e.* LLRT [11], KBR [53], WLRTR [10], and NGmeet [22]), and deep-learning-based HSI denoising methods (*i.e.*, HSID-CNN [71], QRNN3D [66], GRUNet [35], T3SC [5], TRQ3D [49]), SST [37], and SERT [38]. The SOTA RGB denoising transformer, *i.e.*, Restormer [72], is also included by adjusting the input and output channels and training with the same

HSI dataset as ours. As the quantitative results provided in Tab. 1. we can observe that our method achieves superior performance outperforming the other ones by over 1 dB improvement on PSNR. Besides, we achieve the lowest SAM metric, which indicates that our method is better at maintaining spectral consistency. The synthetic color image is given in Fig. 6 for the visual comparison.

**Complex Denoising.** While Gaussian denoising might be useful for some scenarios, *e.g.*, as a plug-and-play denoiser [35], it is not common for real-world images. We, therefore, also compare our method with several recently developed methods on simulated complex noise. As most optimization-based methods only perform well on the noise settings they can solve, we compare our method with a different set of optimization-based methods including, LRTV [24], NMoG [14] and TDTV [63], in addition with the same set of deep-learning-based methods as Gaussian denoising. As shown in Tab. 2, our method again achieves the best performance with up to 1.7 dB improvement on PSNR. It demonstrates the stronger modeling ability of our method.
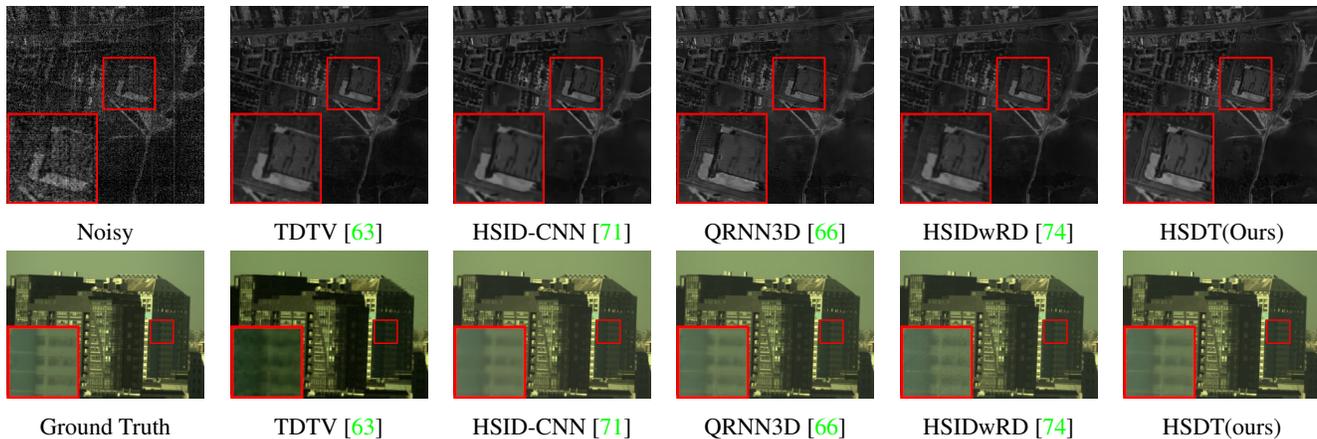
Figure 8: Visual results for real-world noise removal on Urban Dataset (Top) and Real [74] Dataset (Bottom).

| Methods | Params(M) | RealHSI [74] | | CAVE [67] | |
|---|---|---|---|---|---|
| | | PSNR | SAM | PSNR | SAM |
| Noisy | - | 23.31 | 0.257 | 18.99 | 0.901 |
| NMoG [14] | - | 30.90 | 1.762 | 30.84 | 37.86 |
| TDTV [63] | - | 31.14 | 1.853 | 33.14 | 22.34 |
| HSID-CNN [71] | 0.40 | 31.05 | 0.096 | 36.09 | 0.318 |
| QRNN3D [66] | 0.86 | 31.13 | 0.094 | 37.80 | 0.247 |
| GRUNet [35] | 14.2 | 31.03 | 0.091 | 37.33 | 0.288 |
| HSIDwRD [74] | 23.6 | 31.23 | 0.092 | 39.37 | 0.188 |
| HSDT-L(Ours) | 0.52 | **31.42** | **0.091** | **39.80** | **0.174** |

Table 3: Additional results on RealHSI and CAVE.

The visual comparison is given in Fig. 7. Similar to the results of Gaussian noise, our reconstruct better images.

**Real World Denoising.** We also conduct experiments on the recently developed RealHSI dataset [74] using the models trained on ICVL. We compare our methods with the model in [74] and several leading competing methods of complex denoising. The methods that cannot generalize to HSIs with 34 bands, *i.e.*, T3SC [5] and Restormer [72], SST [37], SERT [38], TRQ3D [49], are not included. The quantitative results are given in Tab. 3. It can be seen that our method achieves better performance with fewer parameters. The visual results are provided in the bottom row of Fig. 8. It can be observed that our method produces cleaner and sharper results while the others are blurry or could not completely remove the noise.

### 4.3. Generalization to Other Datasets

We have demonstrated the stronger generalization abilities of our method on real-world denoising dataset [74] with models trained purely on ICVL and simulated noise. With the same models, we provide results on more datasets with different scenes and the number of bands.

| SSA | Guidance | SM-FFN | Conv | Params | PSNR | SAM |
|---|---|---|---|---|---|---|
| ✗ | ✗ | ✗ | Conv3D | 0.43M | 38.28 | 0.101 |
| ✓ | ✗ | ✗ | Conv3D | 0.55M | 41.34 | 0.057 |
| ✓ | ✓ | ✗ | Conv3D | 0.55M | 41.46 | 0.054 |
| ✓ | ✓ | ✓ | Conv3D | 0.58M | 41.62 | 0.052 |
| ✓ | ✓ | ✓ | S3Conv | 0.52M | **41.82** | **0.049** |

Table 4: Break-down ablation on the proposed components.

**Natural HSI.** The CAVE [67] is another widely used HSI denoising dataset that contains more indoor scenes, *e.g.*, different materials and objects. We evaluate the performance under mixture complex noise with the models trained on ICVL. As shown in Tab. 3, our model still achieves the best performance against the others with a large margin, which proves the stronger generalization capabilities of our method to handle out-of-distribution images.

**Remotely Sensed HSI.** We demonstrate that our model trained on ICVL with 31 bands can be directly transferred to datasets with a totally different number of bands, *e.g.*, Urban with 210 bands, without any fine-tuning. This is largely supported by the flexibility of the proposed GSSA and S3Conv units, while such flexibilities are not presented in models, *e.g.*, Restormer [72] and TRQ3D [49]. We present the visual comparison in Fig. 8. It can be observed that QRNN3D [66] could not completely remove the row noise and lose the details of the roof. HSID-CNN [71] and TDTV [24] eliminate the most noise but produce more blurry results. Our method instead removes most noise while maintaining the sharper details.

### 4.4. Ablation Studies

We adopt Gaussian denoising ($\sigma = 50$) on ICVL to conduct the ablation studies. The baseline model is derived by removing our S3Conv, GSSA, and SM-FNN from HSDT-M. We also conduct the per-component ablation studies, which are provided in the supplementary material.

**Break-down Ablation.** We provide the results of breakdown ablations in Tab. 4, in which we gradually add the proposed blocks back to the baseline model. It can be seen that GSSA provides the most performance gain, which can be attributed to the importance of spectral correlations for HSI denoising. The LQ guidance and SM-FFN improve the models with negligible parameter growth, while S3Conv improves the performance with even fewer parameters.

# 5. Conclusion

We present HSDT, an effective and flexible transformer for hyperspectral image denoising. Built upon the hybrid hierarchical architecture, our HSDT is equipped with a novel S3Conv, GSSA, and SM-FFN module to effectively integrate the local spectral-spatial inductive bias and the long-range spectral dependency modeling. Our S3Conv extracts highly correlated local spatial-spectral features without harming efficiency and flexibility, while the GSSA provides stronger capabilities for capturing global spectral correlations, guided by a set of learnable queries that encode the band-wise spectral signatures. With the SM-FFN to further strengthen the aggregated features of more informative regions, our model outperforms the existing SOTA methods on various datasets under simulated and real-world noise.

# Acknowledgement

# References

[1] Boaz Arad and Ohad Ben-Shahar. Sparse recovery of hyperspectral signal from natural rgb images. In *Eur. Conf. Comput. Vis.*, pages 19–34. Springer, 2016. 6

[2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 2

[3] José M Bioucas-Dias, Antonio Plaza, Gustavo Camps-Valls, Paul Scheunders, Nasser Nasrabadi, and Jocelyn Chanussot. Hyperspectral remote sensing data analysis and future challenges. *IEEE Trans. Geosci. Remote Sens.*, 1(2):6–36, 2013. 1

[4] George Alan Blackburn. Hyperspectral remote sensing of plant pigments. *Journal of Experimental Botany*, 58(4):855–867, 2007. 1

[5] Théo Bodrito, Alexandre Zouaoui, Jocelyn Chanussot, and Julien Mairal. A trainable spectral-spatial sparse coding model for hyperspectral image restoration. *Adv. Neural Inform. Process. Syst.*, 34:5430–5442, 2021. 1, 6, 7, 8

[6] James Bradbury, Stephen Merity, Caiming Xiong, and Richard Socher. Quasi-recurrent neural networks. *arXiv preprint arXiv:1611.01576*, 2016. 1, 3

[7] Yuanhao Cai, Jing Lin, Xiaowan Hu, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc Van Gool. Mask-guided spectral-wise transformer for efficient hyperspectral image reconstruction. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 17502–17511, 2022. 2, 3, 5, 13, 14, 15

[8] Stanley H Chan, Xiran Wang, and Omar A Elgendy. Plug-and-play admm for image restoration: Fixed-point convergence and applications. *IEEE Trans. Comput. Imaging.*, 3(1):84–98, 2016. 15

[9] Yi Chang, Luxin Yan, Houzhang Fang, Sheng Zhong, and Wenshan Liao. Hsi-denet: Hyperspectral image restoration via convolutional neural network. *IEEE Trans. Geosci. Remote Sens.*, 57(2):667–682, 2018. 3

[10] Yi Chang, Luxin Yan, Xi-Le Zhao, Houzhang Fang, Zhijun Zhang, and Sheng Zhong. Weighted low-rank tensor recovery for hyperspectral image restoration. *IEEE Trans. Cybern.*, 50(11):4558–4572, 2020. 6, 7

[11] Yi Chang, Luxin Yan, and Sheng Zhong. Hyper-laplacian regularized unidirectional low-rank tensor recovery for multispectral image denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4260–4268, 2017. 6, 7

[12] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 12299–12310, 2021. 3

[13] Hongyu Chen, Guangyi Yang, and Hongyan Zhang. Hider: A hyperspectral image denoising transformer with spatial–spectral constraints for hybrid noise removal. *IEEE Transactions on Neural Networks and Learning Systems*, 2022. 3, 13, 14, 15

[14] Yang Chen, Xiangyong Cao, Qian Zhao, Deyu Meng, and Zongben Xu. Denoising hyperspectral image with non-iid noise structure. *IEEE Trans. Cybern.*, 48(3):1054–1066, 2017. 7, 8

[15] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014. 3

[16] Weisheng Dong, Huan Wang, Fangfang Wu, Guangming Shi, and Xin Li. Deep spatial–spectral representation learning for hyperspectral image denoising. *IEEE Trans. Comput. Imaging.*, 5(4):635–648, 2019. 2, 3, 14

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2, 3

[18] Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018. 15

[19] James Ellis. Searching for oil seeps and oil-impacted soil with hyperspectral imagery. *Earth Observation Magazine*, 1:25–28, 2001. 1

[20] Ying Fu, Zhiyuan Liang, and Shaodi You. Bidirectional 3d quasi-recurrent neural networkfor hyperspectral image super-resolution. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, 14:2674–2688, 2021. 4, 16

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016. 3

[22] Wei He, Quanming Yao, Chao Li, Naoto Yokoya, and Qibin Zhao. Non-local meets global: An integrated paradigm for hyperspectral denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6868–6877, 2019. 1, 6, 7

[23] Wei He, Quanming Yao, Naoto Yokoya, Tatsumi Uezato, Hongyan Zhang, and Liangpei Zhang. Spectrum-aware and transferable architecture search for hyperspectral image restoration. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIX*, pages 19–37. Springer, 2022. 3

[24] Wei He, Hongyan Zhang, Liangpei Zhang, and Huanfeng Shen. Total-variation-regularized low-rank matrix factorization for hyperspectral image restoration. *IEEE Trans. Geosci. Remote Sens.*, 54(1):178–188, 2015. 7, 8

[25] Xin He, Yushi Chen, and Zhouhan Lin. Spatial-spectral transformer for hyperspectral image classification. *Remote Sens.*, 13(3):498, 2021. 3

[26] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 3

[27] Danfeng Hong, Zhu Han, Jing Yao, Lianru Gao, Bing Zhang, Antonio Plaza, and Jocelyn Chanussot. Spectralformer: Rethinking hyperspectral image classification with transformers. *IEEE Trans. Geosci. Remote Sens.*, 60:1–15, 2021. 3

[28] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 3

[29] Ryuji Imamura, Tatsuki Itasaka, and Masahiro Okuda. Zero-shot hyperspectral image denoising with separable image prior. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3

[30] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proc Int. Conf. Mach Learn.*, pages 448–456. PMLR, 2015. 4

[31] Junjun Jiang, He Sun, Xianming Liu, and Jiayi Ma. Learning spatial-spectral prior for super-resolution of hyperspectral imagery. *IEEE Trans. Comput. Imaging.*, 6:1082–1096, 2020. 16

[32] Linda S Kalman and Edward M Bassett III. Classification and material identification in an urban environment using hydice hyperspectral data. In *Imaging Spectrometry III*, volume 3118, pages 57–68. SPIE, 1997. 4, 6

[33] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7

[34] Zeqiang Lai and Ying Fu. Mixed attention network for hyperspectral image denoising. *arXiv preprint arXiv:2301.11525*, 2023. 3

[35] Zeqiang Lai, Kaixuan Wei, and Ying Fu. Deep plug-and-play prior for hyperspectral image restoration. *Neurocomputing*, 481:281–293, 2022. 1, 3, 4, 6, 7, 8, 15, 16

[36] Zeqiang Lai, Kaixuan Wei, Ying Fu, Philipp Härtel, and Felix Heide. ∇-prox: Differentiable proximal algorithm modeling for large-scale optimization. *ACM Transactions on Graphics (TOG)*, 42(4):1–19, 2023. 2

[37] Miaoyu Li, Ying Fu, and Yulun Zhang. Spatial-spectral transformer for hyperspectral image denoising. *arXiv preprint arXiv:2211.14090*, 2022. 2, 3, 5, 6, 7, 8, 13, 14, 15

[38] Miaoyu Li, Ji Liu, Ying Fu, Yulun Zhang, and Dejing Dou. Spectral enhanced rectangle transformer for hyperspectral image denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5805–5814, 2023. 3, 6, 7, 8

[39] Yawei Li, Yuchen Fan, Xiaoyu Xiang, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. Efficient and explicit modelling of image hierarchies for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18278–18289, 2023. 3

[40] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Int. Conf. Comput. Vis.*, pages 1833–1844, 2021. 2, 3, 4

[41] Jing Lin, Yuanhao Cai, Xiaowan Hu, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc Van Gool. Coarse-to-fine sparse transformer for hyperspectral image reconstruction. *arXiv preprint arXiv:2203.04845*, 2022. 3

[42] Yang Liu, Xin Yuan, Jinli Suo, David J. Brady, and Qionghai Dai. Rank minimization for snapshot compressive imaging. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(12):2990 – 3006, 2019. 16

[43] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Int. Conf. Comput. Vis.*, pages 10012–10022, 2021. 3

[44] Patrick Llull, Xuejun Liao, Xin Yuan, Jianbo Yang, David Kittle, Lawrence Carin, Guillermo Sapiro, and David J Brady. Coded aperture compressive temporal imaging. *Optics express*, 21(9):10526–10545, 2013. 15

[45] Matteo Maggioni, Vladimir Katkovnik, Karen Egiazarian, and Alessandro Foi. Nonlocal transform-domain filter for volumetric data denoising and reconstruction. *IEEE Trans. Image Process.*, 22(1):119–133, 2012. 1, 2

[46] Nishir Mehta, Sushant P Sahu, Shahensha Shaik, Ram Devireddy, and Manas Ranjan Gartia. Dark-field hyperspectral imaging for label free detection of nano-bio-materials. *Wiley Interdisciplinary Reviews: Nanomedicine and Nanobiotechnology*, 13(1):e1661, 2021. 1

[47] Yiqun Mei, Yuchen Fan, and Yuqian Zhou. Image super-resolution with non-local sparse attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3517–3526, 2021. 1

[48] Hisham Othman and Shen-En Qian. Noise reduction of hyperspectral imagery using hybrid spatial-spectral derivative-domain wavelet shrinkage. *IEEE Trans. Geosci. Remote Sens.*, 44(2):397–408, 2006. 2

[49] Li Pang, Weizhen Gu, and Xiangyong Cao. Trq3dnet: A 3d quasi-recurrent and transformer based network for hyperspectral image denoising. *Remote Sens.*, 14(18):4598, 2022. 2, 3, 6, 7, 8

[50] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Adv. Neural Inform. Process. Syst.*, 32, 2019. 1, 6

[51] J. Peng, Q. Xie, Q. Zhao, Y. Wang, L. Yee, and D. Meng. Enhanced 3dtv regularization and its applications on hsi denoising and compressed sensing. *IEEE Trans. Image Process.*, 29:7889–7903, 2020. 3

[52] Yi Peng, Deyu Meng, Zongben Xu, Chenqiang Gao, Yi Yang, and Biao Zhang. Decomposable nonlocal tensor dictionary learning for multispectral image denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2949–2956, 2014. 1

[53] Xie Qi, Zhao Qian, Deyu Meng, and Zongben Xu. Kronecker-basis-representation based tensor sparsity and its applications to tensor recovery. *IEEE Trans. Pattern Anal. Mach. Intell.*, PP(99):1–1, 2017. 6, 7

[54] Yuhao Qing, Wenyi Liu, Liuyan Feng, and Wanjia Gao. Improved transformer net for hyperspectral image classification. *Remote Sens.*, 13(11):2216, 2021. 3

[55] Haiquan Qiu, Yao Wang, and Deyu Meng. Effective snapshot compressive-spectral imaging via deep denoising and total variation priors. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9127–9136, 2021. 15, 16

[56] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015. 3, 4

[57] Randall B Smith. Introduction to hyperspectral imaging with tmips. *MicroImages Tutorial Web site*, 14, 2006. 1

[58] Le Sun, Byeungwoo Jeon, Yuhui Zheng, and Zebin Wu. Hyperspectral image restoration using low-rank representation on spectral difference image. *IEEE Geosci. Remote. Sens. Lett.*, 14(7):1151–1155, 2017. 2

[59] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *Proc Int. Conf. Mach Learn.*, pages 10347–10357. PMLR, 2021. 3

[60] Muhammad Uzair, Arif Mahmood, and Ajmal Mian. Hyperspectral face recognition with spatiospectral information fusion and pls regression. *IEEE Trans. Image Process.*, 24(3):1127–1137, 2015. 1

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Adv. Neural Inform. Process. Syst.*, 30, 2017. 2, 3, 5

[62] Ashwin A Wagadarikar, Nikos P Pitsianis, Xiaobai Sun, and David J Brady. Video rate spectral imaging using a coded aperture snapshot spectral imager. *Optics Express*, 17(8):6368–6388, 2009. 15

[63] Yao Wang, Jiangjun Peng, Qian Zhao, Yee Leung, Xi-Le Zhao, and Deyu Meng. Hyperspectral image restoration via total variation regularized low-rank tensor decomposition. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, 11(4):1227–1243, 2017. 3, 7, 8

[64] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 17683–17693, 2022. 3

[65] Kaixuan Wei and Ying Fu. Low-rank bayesian tensor factorization for hyperspectral image denoising. *Neurocomputing*, 331:412–423, 2019. 1, 2, 3

[66] Kaixuan Wei, Ying Fu, and Hua Huang. 3-d quasi-recurrent neural network for hyperspectral image denoising. *IEEE Trans Neural Netw Learn Syst.*, 32(1):363–375, 2021. 1, 2, 3, 4, 6, 7, 8, 14, 15, 16

[67] F. Yasuma, T. Mitsunaga, D. Iso, and S.K. Nayar. Generalized Assorted Pixel Camera: Post-Capture Control of Resolution, Dynamic Range and Spectrum. Technical report, Nov 2008. 6, 8

[68] Dabing Yu, Qingwu Li, Xiaolin Wang, Zhiliang Zhang, Yixi Qian, and Chang Xu. Dstrans: Dual-stream transformer for hyperspectral image restoration. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3739–3749, 2023. 3

[69] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10819–10829, 2022. 5

[70] Qiangqiang Yuan, Liangpei Zhang, and Huanfeng Shen. Hyperspectral image denoising employing a spectral–spatial adaptive total variation model. *IEEE Trans. Geosci. Remote Sens.*, 50(10):3660–3677, 2012. 1, 2, 3

[71] Qiangqiang Yuan, Qiang Zhang, Jie Li, Huanfeng Shen, and Liangpei Zhang. Hyperspectral image denoising employing a spatial–spectral deep residual convolutional neural network. *IEEE Trans. Geosci. Remote Sens.*, 57(2):1205–1218, 2018. 1, 3, 4, 6, 7, 8

[72] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5728–5739, 2022. 2, 3, 5, 6, 7, 8, 13, 14, 15

[73] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44:6360–6376, 2021. 15

[74] Tao Zhang, Ying Fu, and Cheng Li. Hyperspectral image denoising with realistic data. In *Int. Conf. Comput. Vis.*, pages 2248–2257, 2021. 1, 6, 8

[75] XiLe Zhao, Hao Zhang, TaiXiang Jiang, Michael K Ng, and Xiong-Jun Zhang. Fast algorithm with theoretical guarantees for constrained low-tubal-rank tensor recovery in hyperspec-

tral images denoising. *Neurocomputing*, 413:397–409, 2020.
1, 2

# A. More Details about GSSA

**Computational Complexity of GSSA.** Given an input $\mathbf{X} \in \mathbb{R}^{H \times W \times D \times C}$ where $H, W$ denote height and width, $D$ denotes the number of spectral bands, $C$ denotes the features channels, the computational complexity of each step of GSSA is summarized in Tab. 5. Since the feature channels are typically larger than the number of spectral bands, the asymptotic computational complexity of GSSA is dominated by two linear transformations, *i.e.*, *Linear for V* and *Post linear*.

| Step | Complexity |
|---|---|
| Linear for $V$ | $(H \times W \times D) \times C^2$ |
| Pooling for $Q,K$ | $2 \times (H \times W) \times (D \times C)$ |
| Compute attention matrix | $D \times D \times C$ |
| Feature aggregation | $H \times W \times C \times D \times D$ |
| Post linear | $(H \times W \times D) \times C^2$ |
| Total | $O((H \times W \times D) \times C^2)$ |

Table 5: The computational complexity of GSSA. The overall complexity of GSSA is linear with respect to image size.

**Fast Implementation.** With the simplification of pixel-wise attention via global average pooling, our GSSA can be efficiently implemented with a depth-wise convolution by treating the shared attention map as a convolution filter and swapping the spectral and channel dimensions. The speed comparison is shown in Tab. 6, and it can be seen that the Conv-based implementation is approximately 20% faster than the naive `Matmul`-based one.

| Implementation | Runtime (s) | PSNR |
|---|---|---|
| Matmul-based | 0.60 | 41.82 |
| Conv-based | 0.47 | 41.82 |

Table 6: Speed of different implementations of GSSA. Our Conv-based implementation reduces the running time without harming the performance.

## A.1. Comparison against other Attention.

Here, we provide a more detailed explanation regarding the differences between our GSSA and existing channel or spectral attention mechanisms. *We highlight that our GSSA is significantly different from previous attention mechanisms in a variety aspects.* Since GSSA performs attention along spectral rather than spatial dimensions, we here compare it with four previous attention mechanisms that apply along spectral or channel dimensions including:

Fig. 9 illustrates the structures of the aforementioned attention mechanisms. It is worth noting that all previous methods are essentially variants of MDTA proposed in Restormer, whereas our GSSA is fundamentally distinct from them. In the following, we will provide a detailed

| Attention | Method | Task |
|---|---|---|
| MDTA | Restormer [72] | Color image restoration |
| MS-MSA | MST [7] | Spectral Reconstruction |
| GSA | SST [37] | HSI denoising |
| MGSA | Hider [13] | HSI denoising |

Table 7: The competing attention mechanisms.

explanation of the main differences between the previous methods and our GSSA.

**3D vs 2D Data Format.** The first notable difference, which can be easily confused with previous work, is that *our GSSA performs attention on the spectral dimension*, i.e., the $D$ dimension of a 5D data cube $x \in \mathbb{R}^{B \times C \times D \times H \times W}$. In contrast, previous works, such as MST, and SST, even though they refer to their attention mechanisms as spectral attention, essentially apply channel attention along the $C$ dimension of a 4D data cube $x \in \mathbb{R}^{B \times C \times H \times W}$, which is the same as MDTA. Our 3D approach provides the flexibility to handle HSIs with different bands within a single model. Additionally, it achieves superior performance by preserving the structures of different bands, *i.e.*, each band possesses its own feature set, and their relationship remains unchanged across layers of the entire model.

**QKV Projection.** The second key difference pertains to the projections used for the query, key, and value. Conventional attention mechanisms typically employ three linear projections to project the input into query, key, and value. This approach is utilized in all of the compared methods, with the exception of our GSSA. Instead, *our GSSA applies linear projection solely for the value, which greatly simplifies the design*. By contrast, MDTA needs a extra 3x3 depth-wise convolution after the linear projection. MGSA is identical to MDTA, except that it employs a 3D convolution. MS-MSA and GSA are the same and solely utilize linear projections, with the exception that MS-MSA employs an additional mask attention specifically designed for spectral reconstruction.

**Pooling vs Reshape.** The third difference is that our GSSA uses global average pooling to obtain feature maps for each band. This differs from previous methods that adopt a reshape approach. Our method is significantly more computationally efficient compared to previous approaches. Previous methods reshape the Q, K, and V tensors from a shape of $H \times W \times C$ into $HW \times C$, treating $HW$ as the features for each channel. This leads to a time complexity of dot-product attention that is linear with respect to the image size, *i.e.*, $D \times D \times HWC$. In contrast, our GSSA approach only has a constant time complexity $D \times D \times C$, where $D$ denotes the number of bands.

**Learnable Query.** The fourth notable difference is the introduction of the learnable query (LQ), which is motivated by the fixed patterns of pixel values across different

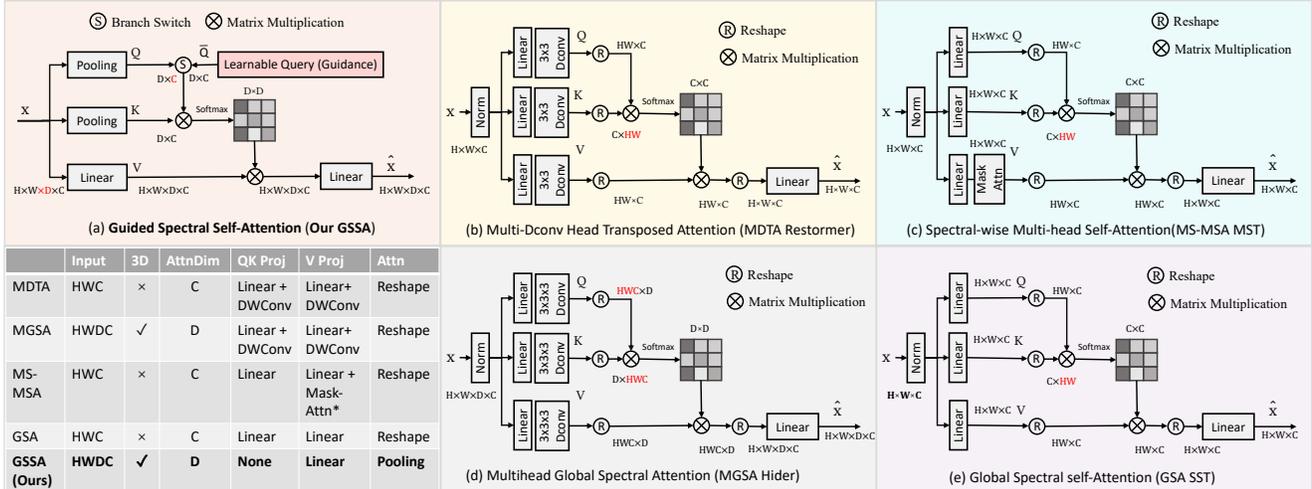| | Input | 3D | AttnDim | QK Proj | V Proj | Attn |
|---|---|---|---|---|---|---|
| MDTA | HWC | × | C | Linear + DWConv | Linear + DWConv | Reshape |
| MGSA | HWDC | ✓ | D | Linear + DWConv | Linear+ DWConv | Reshape |
| MS-MSA | HWC | × | C | Linear | Linear + Mask-Attn* | Reshape |
| GSA | HWC | × | C | Linear | Linear | Reshape |
| **GSSA (Ours)** | **HWDC** | **✓** | **D** | **None** | **Linear** | **Pooling** |

Figure 9: Comparison of different spectral/channel attention mechanisms. Our GSSA is significantly different from previous attention mechanisms. We could observe MDTA and MGSA are almost identical; MS-MSA and GSA are almost identical. Besides, MS-MSA and GSA are also basically simpler version of MDTA without depthwise convolution. Please refer to the text for detailed explanation.

bands. For example, the values of band 100 nm and 200 nm are correlated. Our LQ helps to identify these correlations and the alternative training strategy enables improvements without any extra cost on the number of parameters, inference time, and the flexibility to handle HSIs with different bands.

## B. More Ablation Studies

To evaluate the effectiveness of the proposed components, we conduct a series of experiments to explore the different design choices for each part of our HSDT architecture. Specifically, we compare the proposed blocks, which include GSSA, S3Conv, and SM-FNN, by separately replacing them with existing blocks that share the same functionality, *e.g.*, replacing S3Conv with Conv3D. We use HSDT-M as the base model and evaluate the performance of the different blocks by replacing them one at a time. For blocks that cannot be incorporated into our 3D architectural design of HSDT, such as 2D spectral attention [37], we report the results obtained using their respective models.

**Spatial-Spectral Separable Convolution.** We evaluate several variants of our S3Conv. The most straightforward variant, S3Conv-S, sets the number of spatial convolutions to 1, while the S3Conv variant that we adopt uses 2. Another variant, S3Conv-Seq, applies spatial and spectral convolutions sequentially instead of in parallel. As shown in Table 8, both variants achieve comparable performance with roughly 60% of the parameters used by Conv3D. Our adopted version achieves a 0.2 dB PSNR gain with only 80% of the parameters used by Conv3D. Notably, our S3Conv approach significantly outperforms previous HSI

| Model | #P(Conv) | #P(Total) | PSNR | SAM |
|---|---|---|---|---|
| Conv3D | 0.43M | 0.58M | 41.62 | 0.052 |
| Sep3D [16] | 0.37M | 0.53M | 41.44 | 0.054 |
| S3Conv-S | 0.26M | 0.42M | 41.47 | 0.052 |
| S3Conv-Seq | 0.26M | 0.42M | 41.58 | 0.052 |
| S3Conv | 0.36M | 0.52M | **41.82** | **0.049** |

Table 8: Comparison of different S3Conv variants against 3D convolution and previous separable convolution. Our S3Conv achieves significant better performance with fewer parameters. Our methods are highlighted as gray . #P denotes the model parameters.

separable convolution approaches [16], achieving over 0.4 dB PSNR improvement with even fewer parameters.

**Guided Spectral Self-Attention.** We compare the proposed GSSA approach with existing spectral fusion techniques, including QRU [66], GSA [37], MS-MSA [7], MDTA [72], and MGSA [13]. It is worth noting that although GSA and MS-MSA are named as spectral attention, they are essentially channel attentions derived from MDTA, as discussed earlier. Furthermore, GSA, MS-MSA, and MDTA are all 2D attention approaches that work with 4D data formats instead of the 5D data format used by HSDT. Therefore, we report the results of their models when compared with GSA, MS-MSA, and MDTA. For 3D spectral fusion techniques such as QRU and MGSA, we report the results of models that replace the GSSA of HSDT-M with them. Table 9 presents the results of different attention mechanisms. Our GSSA approach achieves the best results against the other approaches. Notably, our GSSA outper-

forms previous GSA and MGSA approaches (which are also designed for HSI denoising) by a large margin, demonstrating the effectiveness of our designs.

| Model | Params | PSNR | SAM |
|---|---|---|---|
| QRU [66] | 0.57M | 41.31 | 0.064 |
| GSA [37] & MS-MSA [7] | 4.14M | 41.41 | 0.052 |
| MDTA [72] | 26.2M | 41.03 | 0.062 |
| MGSA [13] | 0.50M | 39.74 | 0.102 |
| GSSA | 0.52M | **41.82** | **0.049** |

Table 9: Results of our GSSA in comparison with other attention blocks. Our GSSA achieves a prominent improvement against QRU by over 0.5 PSNR improvement, while previous HSI denoising transformer with GSA only outperforms QRU by only 0.1 PSNR.

**Self-Modulated Feed-Forward Network.** The proposed SM-Branch can be used without additional conventional FFN. As shown in Tab. 10, the sole use of SM-Branch also outperforms the conventional FFN, and the combination of them both yields the best results with very few extra parameters. The GDFN [72] developed for RGB restoration performs poorly and might be unsuitable for our model.

| Model | Params | PSNR | SAM |
|---|---|---|---|
| FFN | 0.49M | 41.67 | 0.050 |
| GDFN [72] | 0.49M | 37.38 | 0.094 |
| SM-Branch | 0.45M | 41.74 | 0.051 |
| SM-FFN | 0.52M | **41.82** | **0.049** |

Table 10: Comparison of the existing FFN with our SM-FFN and SM-Branch.

## C. More Discussions

**Visualization of S3Conv.** To demonstrate the effectiveness of our S3Conv. We provide a comparison of the features map between S3Conv and conventional 3D convolution. As shown in Fig. 10, our S3Conv extracts more spatial meaningful features.

**Analysis of SM-FFN.** The proposed SM-FFN is designed for strengthening the features with higher activation via a self-modulation operation. The improvement provided by SM-FFN could be intuitively explained by the emphasis on more informative regions that typically have higher activation. In the following, we provide some possible relations between our SM-FFN and the SiLU [18] activation, which might further imply why our SM-FFN works better. Specifically, The SiLU activation is,

$$y = x \odot \mathrm{sigmoid}(x), \qquad (7)$$



(a) Input  (c) Feature maps produced by Conv3D
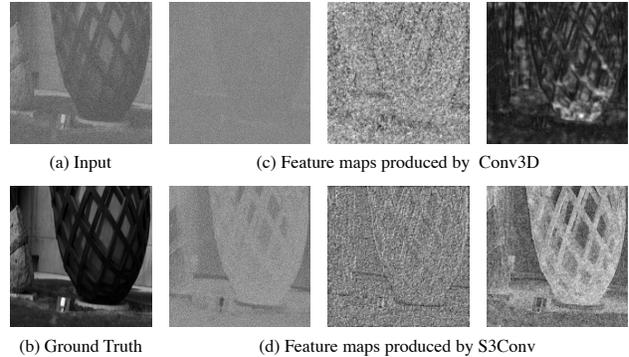(b) Ground Truth  (d) Feature maps produced by S3Conv

Figure 10: Comparison of the feature maps extracted by conventional 3D convolution and our S3Conv.

where $x$ and $y$ are the input and output feature maps. It can be observed that SiLU could be treated as a kind of self-modulation where the modulation weight is computed from the input itself. However, such homogeneous self-modulation might be limited in expressive abilities. Instead, our SM-FFN employs a heterogeneous self-modulation,

$$y = \mathrm{Linear}_1(x) \odot \mathrm{sigmoid}(\mathrm{Linear}_2(x)), \qquad (8)$$

where we adopt two extra linear projections to project input x into two different spaces. This removes the restriction of SiLU where the input $x$ should simultaneously play two roles of features and modulation weight. Thus, our SM-FFN can obtain the advantages of SiLU, *e.g.*, training stability and implicit regularization while maintaining more representation capability. Consequently, it leads to better performance than conventional FFN.

## D. Extension as Plug-and-Play Prior

Considering the superior performance of our method on the Gaussian denoising task, we demonstrate that HSDT can be used a plug-and-play (PnP) prior [8] to solve general HSI restoration tasks with proximal algorithms, *e.g.*, ADMM and HQS.

**Experimental Setup.** We adopt PnP-ADMM [35] to extend our method to the tasks of compressive sensing, and super-resolution. To meet the requirements of PnP algorithms, *i.e.*, Gaussian denoiser for continuous noise strengths, we retrain our model, *i.e.*, HSDT-M, with an additional noise level map [73] on simulated Gaussian noise ranged from 0 to 70. We run 40 iterations for compressive sensing and 24 iterations for super-resolution. The hyperparameters of the algorithms are manually tuned to achieve the best performance.

**Compressive Sensing.** We conduct the simulated experiments on CASSI [62] system. Following [55], the shifting random binary mask [44] is used in our simulation. We

| Method | PSNR | SSIM |
|---|---|---|
| 2DTV | 25.26 | 0.863 |
| 3DTV | 28.46 | 0.910 |
| DeSCI [42] | 26.62 | 0.912 |
| SCI-TV-FFDNet [55] | 29.35 | 0.925 |
| DPHSIR [35] | 30.56 | 0.945 |
| PnP-HSDT (ours) | **31.64** | **0.948** |

(a) Results on the task of compressive sensing.

| | 2x | | 4x | |
|---|---|---|---|---|
| Method | PSNR | SSIM | PSNR | SSIM |
| Bicubic | 35.13 | 0.9575 | 35.12 | 0.954 |
| SSPSR [31] | 47.55 | 0.995 | 39.19 | 0.979 |
| Bi-3DQRNN [20] | 42.53 | 0.989 | 39.56 | 0.979 |
| DPHSIR [35] | 48.75 | 0.996 | 40.95 | 0.980 |
| PnP-HSDT (ours) | **49.76** | **0.996** | **41.56** | **0.982** |

(b) Results on the task of super-resolution.

Table 11: Experimental results of our PnP extension on the task of compressive sensing and super-resolution.

provide the results on CAVE `Toy`, which is obtained from [42]. We compare several recent methods, including DPH-SIR [35], SCI-TV-FFDNet [55], DeSCI [42], and traditional methods, *i.e.*, 2DTV and 3DTV. The quantitative results are shown in Tab. 1a. It can be seen that our method obtains the best performance with over 1 dB improvement on PSNR. Specifically, the improvement is purely obtained through the superior denoising ability of our model, which means our model can also be integrated into other more advanced PnP methods for further improvement, *e.g.*, [55].

**Super-Resolution.** We also provide results on the task of HSI super-resolution. Following [35], we first blur the high-resolution HSI via an $8 \times 8$ Gaussian blur kernel with $\sigma = 3$, and then downsample the image to obtain the low-resolution HSI. We provide the results on ICVL with a scale factor of 2 and 4. The competing methods include several recently developed methods, *e.g.*, SSPSR [31], Bi3DQRNN [20], and DPHSIR [35] . As shown in Tab. 1b, our method achieves the best performance. In particular, our method only needs the pretrained Gaussian denoising model, which is the same as [35]. The improvement against [35] comes from the better PnP denoising prior, which further demonstrates the stronger denoising ability of our method.

# E. More Implementation Details

**Setup of the Learning Rate.** In this part, we provide more details about the multi-step learning rate scheduler that we used for training our simulated Gaussian and complex denoising models. Specifically, we use a multi-stage training strategy to train the models for Gaussian noise and complex noise. The learning rate is set up as shown in Tab. 3a. We use learning rate warmup to gradually increase the learning rate from 0 to $1 \times 10^{-3}$ for the first epoch of the second stage.

**Details of the Simulated Complex Noise.** We follow [66] for constructing simulated complex noise. In details, we consider the non-independent and non-identically distributed (non-i.i.d) Gaussian noise, stripe noise, deadline noise, impulse noise, and the combination of the aforementioned noise (denoted as mixture noise). The details about these five cases of noise are listed as follows,

- **Non-i.i.d noise**. The non-independent and non-identically distributed Gaussian is added to every pixel of each HSI. The noise strength is randomly selected from 10, 30, 50, and 70.
- **Stripe noise**. Stripe noise (5% to 15% percentages of columns) is added to randomly selected one-third of bands. Non-i.i.d. Gaussian noise is added to All bands.
- **Deadline noise**. Deadline noise is added to randomly selected one-third of bands. Non-i.i.d. Gaussian noise is added to All bands.
- **Impulse noise**. Impulse noise with intensity ranging from 10% to 70% is added to randomly selected one-third of bands. Non-i.i.d. Gaussian noise is added to All bands.
- **Mixture noise**. Each band is randomly corrupted by at least one kind of noise mentioned above.

**System Configuration.** In the main paper, we compare the running time of different methods. All the comparisons are performed with an Nvidia GeForce RTX 3090, and an Intel(R) Core(TM) i9-10850K CPU @ 3.60GHz on Ubuntu 20.04.1 LTS. All the CNN-based methods are implemented and tested with PyTorch 1.7.1. All the optimization-based methods are implemented and tested with Matlab. We test the running time on ICVL with an image size of $512 \times 512$ by repeating the test 10 times and averaging the results.

# F. Future work.

In this work, we propose a transformer architecture, *i.e.*, HSDT for hyperspectral image denoising. We introduce several effective and generalizable components to better explore the spatial-spectral and global spectral correlations of HSI. Specifically, it is worthwhile to explore the applications of the proposed S3Conv and HSDT for more network architectures and tasks. Furthermore, our learnable queries could also be extended to condition on some external information for more explicit guidance. For example, we might be able to inject the Gaussian noise strength into the network with learnable queries, through an embedding layer. This is helpful for a PnP Gaussian denoiser, where the noise strength is known.

| Stage 1 | Gaussian Noise $\sigma = 50$ | | | | |
|---|---|---|---|---|---|
| Epoch | 0 - 20 | 20 - 30 | | | |
| LR | $1 \times 10^{-3}$ | $1 \times 10^{-4}$ | | | |
| **Stage 2** | **Gaussian Noise $\sigma = 10, 30, 50, 70$** | | | | |
| Epoch | 30 - 45 | 45 - 55 | 55 - 60 | 60 - 65 | 65 - 75 | 75 - 80 |
| LR | $1 \times 10^{-3}$ | $1 \times 10^{-4}$ | $5 \times 10^{-5}$ | $1 \times 10^{-5}$ | $5 \times 10^{-6}$ | $1 \times 10^{-6}$ |
| **Stage 3** | **Complex Noise** | | | | |
| Epoch | 80 - 90 | 90 - 95 | 95 - 100 | 100 - 105 | 105 - 110 |
| LR | $1 \times 10^{-3}$ | $5 \times 10^{-4}$ | $1 \times 10^{-4}$ | $5 \times 10^{-5}$ | $1 \times 10^{-5}$ |

(a) Our multi-step learning rate scheduler.

| System | Ubuntu 20.04.1 LTS |
|---|---|
| GPU | Nvidia GeForce RTX 3090 |
| CPU | Intel(R) Core(TM) i9-10850K CPU |
| Framework | PyTorch 1.7.1 |
| Driver | Cuda 11.2 |
| Software | Matlab 2020 |
| Dataset | ICVL |
| Image Size | $512 \times 512$ |
| Repeat times | 10 |

(b) System configuration for the speed test.

Table 12: More implementation details. (a) We adopt a multi-stage training strategy with the learning warmup setup for the first epoch. (b) We provide the system configuration as the results of the speed test are strongly correlated with the configuration.

## G. Broader Impacts

Our work has no ethical issues or broader impacts.