# A Trainable Spectral-Spatial Sparse Coding Model for Hyperspectral Image Restoration

**Théo Bodrito**\*, **Alexandre Zouaoui**\*, **Jocelyn Chanussot, and Julien Mairal**
Inria, Univ. Grenoble Alpes, CNRS, Grenoble INP, LJK, 38000 Grenoble, France
`firstname.lastname@inria.fr`

## Abstract

Hyperspectral imaging offers new perspectives for diverse applications, ranging from the monitoring of the environment using airborne or satellite remote sensing, precision farming, food safety, planetary exploration, or astrophysics. Unfortunately, the spectral diversity of information comes at the expense of various sources of degradation, and the lack of accurate ground-truth "clean" hyperspectral signals acquired on the spot makes restoration tasks challenging. In particular, training deep neural networks for restoration is difficult, in contrast to traditional RGB imaging problems where deep models tend to shine. In this paper, we advocate instead for a hybrid approach based on sparse coding principles that retains the interpretability of classical techniques encoding domain knowledge with handcrafted image priors, while allowing to train model parameters end-to-end without massive amounts of data. We show on various denoising benchmarks that our method is computationally efficient and significantly outperforms the state of the art. [1]

## 1 Introduction

Hyperspectral imaging (HSI) enables measurements of the electromagnetic spectrum of a scene on multiple bands (typically about a hundred or more), which offers many perspectives over traditional color RGB imaging. For instance, the high-dimensional information present in a single pixel is sometimes sufficient to identify the signature of a particular material, which is of course infeasible in the RGB domain. Not surprisingly, hyperspectral imaging is then of utmost importance and has a huge number of scientific and technological applications such as remote sensing [8, 22, 47], quality evaluation of food products [16, 19, 36], medical imaging [2, 18, 38], agriculture and forestry [1, 37, 41], microscopy imaging in biology [25, 57], or exoplanet detection in astronomy [24].

Information contained in hyperspectral signals is much richer than in RGB images, but the price to pay is the need to deal with complex degradations that may arise from multiple sources, including sparse noise with specific patterns (stripes), in addition to photon and thermal noise [30, 52]. As a consequence, HSI denoising is a crucial pre-processing step to enhance the image quality before using data in downstream tasks such as semantic segmentation or spectral unmixing [31]. A second issue is the lack of large-scale collection of ground-truth high-quality signals and the large diversity of sensor types, which makes it particularly challenging to train machine learning models for restoration such as convolutional neural networks. To deal with the scarcity of ground-truth data, most successful approaches typically encode strong prior knowledge about data within the model architecture, which may be low-rank representations of input patches [17, 23, 54, 61, 71], sparse coding [13, 21, 23], or image self-similarities [40, 50, 72], which have proven to be very powerful in the RGB domain [9].

---

In this paper, we propose a fully interpretable machine learning model for hyperspectral images that may be seen as a hybrid approach between deep learning techniques, where parameters can be learned end to end with supervised data, and classical methods that essentially rely on image priors. Since designing an appropriate image prior by hand is very hard, our goal is to benefit from deep learning principles (here, differentiable programming [6]) while encoding domain knowledge and physical rules about hyperspectral data directly into the model architecture, which we believe is a key to develop robust approaches that do not require massive amounts of training data.

More precisely, we introduce a novel trainable spectral-spatial sparse coding model with two layers, which performs the following operations: (i) The first layer decomposes the spectrum measured at each pixel as a sparse linear combination of a few elements from a learned dictionary, thus performing a form of linear spectral unmixing per pixel, where dictionary elements can be seen as basis elements for spectral responses of materials present in the scene. (ii) The second layer builds upon the output of the first one, which is represented as a two-dimensional feature map, and sparsely encodes patches on a dictionary in order to take into account spatial relationships between pixels within small receptive fields. To further reduce the number of parameters to learn and leverage classical prior knowledge about spectral signals [61], we also assume that the dictionary elements admit a low-rank structure—that is, dictionary elements are near separable in the space and spectrum domains, as detailed later. Even though dictionary learning has been originally introduced for unsupervised learning [43, 48], we adopt an unrolled optimization procedure inspired by the LISTA algorithm [26], which has been very successful in imaging problems for training sparse coding models from supervised data [34, 35, 56, 64].

Our motivation for adopting a two-layer model is to provide a shared architecture for different HSI sensors, which often involve a different number of bands with different spectral responses. Our solution consists of learning sensor-specific dictionaries for the first layer, while the dictionary of second layer is shared across modalities. This allows training simultaneously on several HSI signals, the first layer mapping input data to a common space, before processing data by the second layer.

We experimentally evaluate our HSI model on standard denoising benchmarks, showing a significant improvement over the state of the art (including deep learning models and more traditional baselines), while being computationally very efficient at test time. Perhaps more important than pure quantitative results, we believe that our work also draws interesting conclusions for machine learning. First, by encoding prior knowledge within the model architecture directly, we obtain models achieving excellent results with a relatively small number of parameters to learn, a conclusion also shared by [34, 35] for RGB imaging; nevertheless, the effect is stronger in our work due to the scarcity of training data for HSI denoising and the difficulty to train deep learning models for this task. Second, we also show that interpretable architectures are useful: our model architecture can adapt to different noise levels per band and modify the encoding function at test time in a principled manner, making it well suited for solving blind denoising problems that are crucial for processing hyperspectral signals.

## 2    Related Work on Hyperspectral Image Denoising

**Learning-free and low-rank approaches.**    Classical image denoising methods such as BM3D [12] may be applied independently to each spectral band of HSI signals, but such an approach fails to capture relations between channels; Not surprisingly, multi-band techniques such as BM4D [40] have been shown to perform better for HSI, and other variants were subsequently proposed such as GLF [72]. Tensor-based methods such as LLRT [11] are able to exploit the underlying low-rank structure of HSI signals [17, 54, 68] and have shown particularly effective when combined with a non-local image prior as in NGMeet [28]. Finally, other approaches adapt traditional image processing priors such as total variation [66, 60], or wavelet sparsity [49, 53] but they tend to perform worse than GLF, LLRT, or NGMeet, see [32] for a survey on denoising techniques for HSI.

**Sparse coding models.**    Dictionary learning [48] is an unsupervised learning technique consisting of representing a signal as a linear combination of a few elements from a learned dictionary, which has shown to be very effective for various image restoration tasks [15, 44]. Several approaches have then combined dictionary learning and low-rank regularization. For instance, 3D patches are represented as tensors in [50] and are encoded by using spatial-spectral dictionaries [58]. In [71], 2D patches are extracted from the band-vectorized representation of the 3D HSI data and sparsely encoded on a dictionary, while encouraging low-rank representations with a trace norm penalty on

the reconstructed image. The low-rank constraint can also be enforced by designing the dictionary as the result of the matrix multiplication between spatial and spectral dictionaries learned by principal component analysis as in [21]. However, these methods typically compute sparse representations with an iterative optimization procedure, which may be computationally demanding at test time.

**Deep learning.** Like BM3D above, convolutional neural networks for grayscale image denoising (*e.g.*, DnCNN [69]) may also be applied to each spectral band, which is of course suboptimal. Because deep neural networks have been highly successful for RGB images with often low computational inference cost, there have been many attempts to design deep neural networks dedicated to HSI denoising. For instance, to account for the large number of hyperspectral bands, several approaches based on convolutional neural networks are operating on sliding windows in the spectral domain, [39, 55, 67], which allows training models on signals with different number of spectral bands, but the sliding window significantly increases the inference time. More precisely, attention layers are used in [55], while more traditional CNNs are used in [39], possibly with residual connections [67]. Recently, an approach based on recurrent architecture was proposed in [63] to process signals with an arbitrary number of bands, achieving impressive results for various denoising tasks.

**Hybrid approaches.** SMDS-Net [64] adopts a hybrid approach between sparse coding and deep learning models by adapting the RGB image restoration method of [35] to HSI images. The resulting pipeline however lacks interpretability: SMDS-Net first denoises the input image with non-local means [9], then performs subspace projection in the spectral domain using HySime[7], before sparsely encoding 3D patches (cubes) with a trainable version of Tensor-based ISTA [51]. Although this method reduces considerably the number of parameters in comparison to vanilla deep learning models, the spectral sliding windows approach lacks interpretability since the same denoising procedure is applied across different bands, which may not suffer from the same level of noise. In contrast, we propose a much simpler sparse coding model, which is physically consistent with the nature of hyperspectral signals, by introducing a novel differentiable low-rank sparse coding layer.

## 3 A Trainable Spectral-Spatial Sparse Coding Model (T3SC)

In this section, we introduce our trainable spectral-spatial sparse coding model dedicated to hyperspectral imaging, and start by presenting some preliminaries on sparse coding.

### 3.1 Background on Sparse Coding

**Image denoising with dictionary learning.** A classical approach introduced by Elad and Aharon [15] for image denoising consists in considering the set of small overlapping image patches (*e.g.*, $8 \times 8$ pixels) from a noisy image, and compute a sparse approximation of these patches onto a learned dictionary. The clean estimates for each patch are then recombined to produce the full image.

Formally, let us consider a noisy image $\mathbf{y}$ in $\mathbb{R}^{c \times h \times w}$ with $c$ channels and two spatial dimensions. We denote by $\mathbf{y}_1, \mathbf{y}_2, \cdots \mathbf{y}_n$ the $n$ overlapping patches from $\mathbf{y}$ of size $c \times s \times s$, which we represent as vectors in $\mathbb{R}^m$ with $m = cs^2$. Assuming that a dictionary $\mathbf{D} = [\mathbf{d}_1, \cdots, \mathbf{d}_p]$ in $\mathbb{R}^{m \times p}$ is given—we will discuss later how to obtain a "good" dictionary— each patch $\mathbf{y}_i$ is processed by computing a sparse approximation:

$$\min_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i\|^2 + \lambda \|\boldsymbol{\alpha}_i\|_1, \tag{1}$$

where $\| \cdot \|_1$ is the $l_1$-norm, which is known to induce sparsity in the problem solution [43], and $\boldsymbol{\alpha}_i$ is the sparse code representing the patch $\mathbf{y}_i$, while $\lambda$ controls the amount of regularization. Note that the $\ell_0$-penalty, which counts the number of non-zero elements, could also be used, leading to a combinatorial problem whose solution is typically approximated by a greedy algorithm. After solving the $n$ problems (1), each patch $\mathbf{y}_i$ admits a "clean" estimate $\mathbf{D}\boldsymbol{\alpha}_i$. Because each pixel belongs to several patches, the full denoised image $\hat{\mathbf{x}}$ is obtained by averaging these estimates.

Finding a good dictionary can be achieved in various manners. In classical dictionary learning algorithms, $\mathbf{D}$ is optimized such that the sum of the loss functions (1) is as small as possible, see [43] for a review. Adapting the dictionary with supervision is also possible [42], as discussed next.

3

**Differentiable programming for sparse coding.**  The proximal gradient descent method called ISTA [20] is a classical algorithm for solving the Lasso problem in Eq. (1), which consists of the following iterations

$$\boldsymbol{\alpha}_i^{(t+1)} = S_\lambda \left[ \boldsymbol{\alpha}_i^{(t)} + \eta \mathbf{D}^\top \left( \mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i^{(t)} \right) \right], \tag{2}$$

where $\eta > 0$ is a step-size and $S_\lambda[u] = \text{sign}(u) \max(|u| - \lambda, 0)$ is the soft-thresholding operator, which is applied pointwise to each entry of an input vector.

By noting that the above iteration can be seen as a sequence of affine transformations interleaved with pointwise non-linearities $S_\lambda$, it is then tempting to interpret $T$ iterations (2) as a multilayer feed-forward neural network with shared weights between the $T$ layers. Following such an insight, Gregor and LeCun have proposed the LISTA algorithm [26], where the parameters are learned such that the sequence approximates well the solution of the sparse coding problem (1).

Interestingly, the LISTA algorithm can also be used to train dictionaries for supervised learning tasks. This is the approach chosen in [35, 56] for image restoration, which considers the following iterations:

$$\boldsymbol{\alpha}_i^{(t+1)} = S_\lambda \left[ \boldsymbol{\alpha}_i^{(t)} + \mathbf{C}^\top \left( \mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i^{(t)} \right) \right], \tag{3}$$

which differs from (2) with the presence of a matrix $\mathbf{C}$ of the same size as $\mathbf{D}$. Even if the choice $\mathbf{C} = \eta\mathbf{D}$ (which recovers ISTA) is perfectly reasonable, using a different dictionary $\mathbf{C}$ has empirically shown to provide improvements in results quality [35], probably due to faster convergence of the LISTA iterations. Then, given some fixed parameters $\mathbf{C}, \mathbf{D}$, a clean estimate $\mathbf{W}\boldsymbol{\alpha}_i^{(T)}$ for each patch $\mathbf{y}_i$ is obtained by using a dictionary $\mathbf{W}$, where $T$ is the number of LISTA steps. The reason for allowing a different dictionary $\mathbf{W}$ than $\mathbf{D}$ is to correct the potential bias due to $\ell_1$-minimization.

Finally, the denoised image $\hat{\mathbf{x}}$ is reconstructed by averaging the patch estimates:

$$\hat{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^{n} \mathbf{R}_i \mathbf{W}\boldsymbol{\alpha}_i^{(T)}, \tag{4}$$

where $\mathbf{R}_i$ is the linear operator that places the patch $\hat{\mathbf{x}}_i$ at position $i$ in the image, and we assume—by neglecting border effects for simplicity—that each pixel admits the same number $m$ of estimates.

In contrast to classical restoration techniques based on dictionary learning, the LISTA point of view enables us to learn the model parameters $\mathbf{C}, \mathbf{D}, \mathbf{W}$ in a supervised fashion. Given a training set of pairs of noisy/clean images, we remark that the estimate $\hat{\mathbf{x}}$ is obtained from a noisy image $\mathbf{y}$ by a sequence of operations that are differentiable almost everywhere, as typical neural networks with rectified linear unit activation functions. A typical loss, which we optimize by stochastic gradient descent, is then

$$\min_{\mathbf{C}, \mathbf{D}, \mathbf{W}, \lambda} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \|\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}\|^2 \right],$$

where $(\mathbf{x}, \mathbf{y})$ is a pair of clean/noisy images drawn from some training distribution from which we can sample, and $\hat{\mathbf{x}}(\mathbf{y})$ is the clean estimate obtained from (4), given the noisy image $\mathbf{y}$.

## 3.2  A Trainable Low-Rank Sparse Coding Layer

We are now in shape to introduce a trainable layer encoding both sparsity and low-rank principles.

**Spatial-Spectral Representation.**  As shown in [10, 21], HSI patches can be well reconstructed by using only a few basis elements obtained by principal component analysis. The authors further decompose these into a Cartesian product of separate spectral and spatial dictionaries. In this paper, we adopt a slightly different approach, where we consider a single dictionary $\mathbf{D} = [\mathbf{d}_1, \ldots, \mathbf{d}_p]$ in $\mathbb{R}^{m \times p}$ as in the previous section with $m = cs^2$, but each element may be seen as a matrix of size $c \times s^2$ with low-rank structure. More precisely, we enforce the following representation

$$\forall j \in 1, \ldots, p, \ \mathbf{d}_j = \text{vec}\left(\mathbf{U}_j \times \mathbf{V}_j\right), \tag{5}$$

where $\mathbf{U}_j$ is in $\mathbb{R}^{s^2 \times r}$, $\mathbf{V}_j$ is in $\mathbb{R}^{r \times c}$, $r$ is the desired rank of the dictionary elements, and vec(.) is the operator than flattens a matrix to a vector. The hyperparameter $r$ is typically small with $r = 1, 2$ or 3. When $r = 1$, the dictionary elements are said to be separable in the spectral and spatial domains, which we found to be a too stringent condition to achieve good reconstruction in practice.
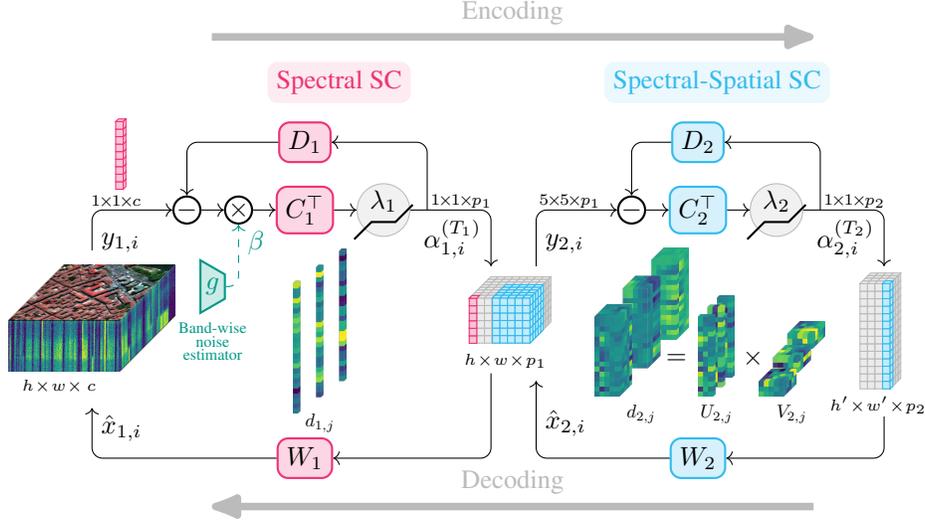
4

Figure 1: Architecture of T3SC : we propose a two-layer sparse coding model which is end-to-end trainable. The first layer performs a sensor-specific spectral decomposition, while the second layer encodes both spectral and spatial information.

The low-rank assumption allows us to build model with a reduced number of parameters, while encoding natural assumption about the data directly in the model architecture. Indeed, whereas a classical full-rank dictionary $\mathbf{D}$ admits $cs^2p$ parameters, the decomposition (5) yields dictionaries with $(s^2 + c)rp$ parameters only. Matrices $\mathbf{C}$ and $\mathbf{W}$ are parametrized in a similar manner.

**Convolutional variant and implementation tricks.** Whereas traditional sparse coding reconstructs local signals (patches) independently according to the iterations (3), another variant called convolutional sparse coding (CSC) represents the whole image by a sparse linear combination of dictionary elements placed at every possible location in the image [56]. From a mathematical point of view, the reconstruction loss for computing the codes $\boldsymbol{\alpha}_i$ given an input image $\mathbf{y}$ becomes

$$\min_{\{\boldsymbol{\alpha}_i \in \mathbb{R}^p\}_{i=1,\dots,n}} \frac{1}{2} \left\| \mathbf{y} - \frac{1}{m} \sum_{i=1}^{n} \mathbf{R}_i \mathbf{D} \boldsymbol{\alpha}_i \right\|^2 + \lambda \sum_{i=1}^{n} \|\boldsymbol{\alpha}_i\|_1. \tag{6}$$

An iterative approach for computing these codes can be obtained by a simple modification of (3) consisting of replacing the quantity $\mathbf{D}\boldsymbol{\alpha}_i^{(t)}$ by the $i$-th patch of the reconstructed image $\frac{1}{m} \sum_{i=1}^{n} \mathbf{R}_i \mathbf{D} \boldsymbol{\alpha}_i^{(t)}$. All of these operations can be efficiently implemented in standard deep learning frameworks, since the corresponding operations corresponds to a transposed convolution with $\mathbf{D}$, followed by convolution with $\mathbf{C}$, see [56] for more details. In this paper, we experimented with the CSC variant (6) and SC one (1), both with low-rank dictionaries, which were previously described. We observed that CSC was providing slightly better results and was thus adopted in our experiments. Following [35], another implementation trick we use is to consider a different $\lambda$ parameter per dictionary element, which slightly increases the number of parameters, while allowing to learn with a weighted $\ell_1$-norm in (6).

### 3.3 The Two-Layer Sparse Coding Model with Sensor-Specific Layer

One of the main challenge in hyperspectral imaging is to train a model that can generalize to several types of sensors, which typically admit different number of spectral bands. Whereas learning a model that is tuned to a specific sensor is perfectly acceptable in many contexts, it is often useful to learn a model that is able to generalize across different types of HSI signals. To alleviate this issue, several strategies have been adopted such as (i) projecting signals onto a linear subspace of fixed dimension, with no guarantee that representations within this subspace can be comparable between different signals, or (ii) processing input data using a sliding window across the spectral domain.

In this paper, we address this issue by learning a two-layer model, presented in Figure 1, where the first layer is tuned to a specific sensor, whereas the second layer could be generic. Note that the

second layer carries most of the model parameters (about $20\times$ more than in the first layer in our experiments). Formally, let us denote by $\boldsymbol{\alpha}$ in $\mathbb{R}^{p\times h\times w}$ the sparse encoding of an input tensor $\mathbf{y}$ in $\mathbb{R}^{c\times h\times w}$ as previously described. A sparse coding layer $\Phi$ naturally yields an encoder and a decoder such that:

$$\Phi^{enc} : \mathbf{y} \mapsto \boldsymbol{\alpha}, \quad \text{and} \quad \Phi^{dec} : \boldsymbol{\alpha} \mapsto \frac{1}{n} \sum_{i=1}^{n} \mathbf{R}_i \mathbf{W} \boldsymbol{\alpha}_i. \tag{7}$$

Given a noisy image $\mathbf{y}$, the denoising procedure described in the previous section with one layer can be written as

$$\hat{\mathbf{x}}(\mathbf{y}) = \Phi^{dec} \circ \Phi^{enc}(\mathbf{y}).$$

Then, a straightforward multilayer extension of the procedure may consist of stacking several sparse coding layers $\Phi_1, \ldots, \Phi_L$ together to form a multilayer sparse coding denoising model:

$$\hat{\mathbf{x}}(\mathbf{y}) = \Phi_1^{dec} \circ \cdots \circ \Phi_L^{dec} \circ \Phi_L^{enc} \circ \cdots \circ \Phi_1^{enc}(\mathbf{y}).$$

The model we propose is composed of two layers, as shown in Figure 1. The first layer encodes spectrally the input HSI image, meaning that it operates on $1 \times 1$ patches, whereas the second layer encodes both spectrally and spatially the output of the first layer.

### 3.4 Noise Adaptive Sparse Coding

An advantage of using a model based on a sparse coding objective (1) is to give the ability to encode domain knowledge within the model architecture. For instance, the Lasso problem (1) seen from a maximization a posteriori estimator implicitly assumes that the noise is i.i.d. If the noise variance is different on each spectral band, a natural modification to the model is to introduce weights and use a weighted-$\ell_2$ data fitting term (which could be applied as well to the CSC model of (6)):

$$\min_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} \frac{1}{2} \sum_{j=1}^{c} \beta_j \|\mathbf{M}_j(\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i)\|^2 + \lambda\|\boldsymbol{\alpha}_i\|_1, \tag{8}$$

where $\mathbf{M}_j$ is a linear operator that extracts band $j$ from a given HSI signal. From a probabilistic point of view, if $\sigma_j^2$ denotes the variance of the noise for band $j$, we may choose the corresponding weight $\beta_j$ to be proportional to $1/\sigma_j^2$. Yet, estimating accurately $\sigma_j^2$ is not always easy, and we have found it more effective to simply learn a parametric function $\beta_j = g(\mathbf{M}_j\mathbf{y})$—here, a very simple CNN with three layers, see supplementary material for details—which is applied independently to each band. It is then easy to modify the LISTA iterations accordingly to take into account these weights, and learn the model parameters jointly with those of the parametric function $g$.

### 3.5 Self-Supervised Learning: Blind-Band Denoising with No Ground Truth Data

Even though acquiring limited ground truth data for a specific sensor is often feasible, it is also interesting to be able to train models with no ground truth at all, *e.g.*, for processing images without physical access to the sensor. In such an unsupervised setting, deep neural networks are typically trained for RGB images by using blind-spot denoising techniques [33], consisting of predicting pixel values given their context. Here, we propose a much simpler approach exploiting the spectral redundancy between channels. More precisely, each time we draw an image for training, we randomly mask one band (ore more), and train the model to reconstruct the missing band from the available ones. Formally, the training objective becomes

$$\min_{\mathbf{C},\mathbf{D},\mathbf{W},\lambda} \mathbb{E}_{\mathbf{x},\mathbf{y},S} \left[ \sum_{j \notin S} \|\mathbf{M}_j(\hat{\mathbf{x}}_S(\mathbf{y}) - \mathbf{y})\|^2 \right], \tag{9}$$

where $S$ is the set of bands that are visible for computing the sparse codes $\boldsymbol{\alpha}_i$, leading to a reconstructed image that we denote by $\hat{\mathbf{x}}_S(\mathbf{y})$. Formally, it would mean considering the objective (8), but replacing the sum $\sum_{j=1}^{c}$ by $\sum_{j \in S}$. This is in spirit similar to blind-spot denoising, except that bands are masked instead of pixels, making the resulting implementation much simpler.

Table 1: Simplified comparison between learning-free and learning-based approaches.

| | *Data req.* | *training* | *inference* | *adapt. to new data* | *complex noise* |
|---|---|---|---|---|---|
| **learning-free** | no req. | no training | slow | easy | poor |
| **learning-based** | clean data | slow | fast | complicated | good perf. |

# 4 Experiments

We now present various experiments to demonstrate the effectiveness of our approach for HSI denoising, but first, we discuss the difficulty of defining the state of the art in this field. We believe indeed that it is not always easy to compare learning-free from approaches based on supervised learning. These two classes of approaches have very different requirements/characteristics, making one class more relevant than the other one in some scenarios, and less in others. Table 1 summarizes their characteristics, displaying advantages and drawbacks of both approaches.

**Benchmarked models.** Keeping in mind the previous dichotomy, we choose to compare our method to traditional methods such as bandwise BM3D [12] (implementation based on [45, 46]), BM4D [40], GLF [72], LLRT [11], NGMeet [28]. We also included deep learning models in our benchmark such as HSID-CNN [67], HSI-SDeCNN [39] 3D-ADNet [55], SMDS-Net [64] and QRNN3D [63]. Results of HSID-CNN, HSI-SDeCNN and 3D-ADNet on Washington DC Mall (available in the Appendix) are taken directly from the corresponding papers, as the train/test split is the same. Otherwise, the results were obtained by running the code obtained directly from the authors, except for SMDS-Net, where our implementation turned out to be slightly more effective. Note that the same architecture for our model was used in all our experiments (see Appendix).

**Datasets.** We evaluate our approach on two datasets with significantly different properties.

- *ICVL* [4] consists of 204 images of size $1392 \times 1300$ with 31 bands. We used 100 images for training and 50 for testing as in [63] but with a different train/test split ensuring that similar images—*e.g.*, picture from the same scene—are not used twice.
- *Washington DC Mall* is perhaps the most widely used dataset[2] for HSI denoising and consists of a high-quality image of size $1280 \times 307$ with 191 bands. Following [55], we split the image into two sub-images of size $600 \times 307$ and $480 \times 307$ for training and one sub-image of size $200 \times 200$ for testing. Even though the test image does not overlap with train images, they nevertheless share common characteristics. Interestingly, the amount of training data is very limited here.

Specific experiments were also conducted with the datasets APEX [29], Pavia[3], Urban[59] and CAVE [65], which appear in the supplementary material.

**Normalization.** Before denoising, HSI images are normalized to $[0, 1]$. For remote sensing datasets, we pre-compute the 2nd and 98th percentiles for each band, on the whole the training set. Then, normalization is performed on train and test images by clipping each band between those percentiles before applying bandwise min-max normalization, similar to [5, 39]. For the close-range dataset ICVL, we simply apply global min-max normalization as in [64, 63].

**Noise patterns.** We evaluate our model against different types of synthetic noise:

- *i.i.d Gaussian noise with known variance* $\sigma^2$, which is the same on all bands.
- *Gaussian noise with unknown band-dependent variance*: We consider Gaussian noise with different standard deviation $\sigma_j$ for each band, which is uniformly drawn in a fixed interval. These standard deviations change from an image to the other and are unknown at test time.
- *Noise with spectrally correlated variance*: We consider Gaussian noise with standard deviation $\sigma_j$ varying continuously across bands, following a Gaussian curve, see details in the appendix.
- *Stripes noise* : similar to [63], we applied additive stripes noise to 33% of bands. In those bands, 10-15% of columns are affected, meaning a value uniformly sampled in the interval $[-0.25, 0.25]$ is added to them. Moreover, all bands are disturbed by Gaussian noise with noise intensity $\sigma = 25$.

---

[2]`https://engineering.purdue.edu/~biehl/MultiSpec/hyperspectral.html`
[3]`http://www.ehu.eus/ccwintco/index.php?title=Hyperspectral_Remote_Sensing_Scenes`

**Metrics.** In order to assess the performances the previous methods, we used five different indexes widely used for HSI restoration, namely (i) Mean Peak Signal-to-Noise Ratio (MPSNR), which is the classical PSNR metric averaged across bands; (ii) Mean Structural Similarity Index Measurement (MSSIM), which is based on the SSIM metric [62]; (iii) Mean Feature Similarity Index Measurement (MFSIM) introduced in [70]; (iv) Mean ERGAS [14], and (v) Mean Spectral Angle Map (MSAM) [3]. We use MPSNR and MSSIM in the main paper and report the other metrics in the appendix.

**Implementation details.** We trained our network by minimizing the MSE between the groundtruth and restored images. For ICVL, we follow the training procedure described in [63]: we first center crop training images to size $1024 \times 1024$, then we extract patches of size $64 \times 64$ at scales 1:1, 1:2, and 1:4, with stride 64, 32 and 32 respectively. The number of extracted patches for ICVL amounts to 52962. For Washington DC Mall, we do not crop training images and the patches are extracted with stride 16, 8 and 8, for a total of 1650 patches. One epoch in Washington DC Mall corresponds to 10 iterations on the training dataset. Basic data augmentation schemes such as $90°$ rotations and vertical/horizontal flipping are performed. Code and additional details about optimization, implementation, computational resources, are provided in the supplementary material. As reported in Table 4, augmenting the number unrolled iterations improves the denoising performances at the expense of inference time. Since the Spectral-Spatial SC layer is the most time-consuming, the number of unrolled iterations chosen for the first and second layers are 12 and 5 respectively.

**Quantitative results on synthetic noise.** We present in Table 2 the results obtained on the ICVL dataset (results on DCMall are presented in the appendix). Our method uses the vanilla model of Section 3.3 for the experiments with constant $\sigma$ or correlated noise. For the blind denoising experiment with band-dependent $\sigma$ or for the stripe noise experiment, we use the variant of Section 3.4, which is designed to deal with unknown noise level per channel. The method "T3SC-SSL" implements the self-supervised learning approach of Section 3.5, which does not rely on ground-truth data.

- Our supervised approach achieves state-of-the-art results (or is close to the best performing baseline) on all settings. GLF performs remarkably well given that this baseline is learning-free.
- Our self-supervised method achieves a relatively good performance under i.i.d. Gaussian noise, but does not perform as well under more complex noise. This is a limitation of the approach which is perhaps expected and overcoming this limitation would require designing a different self-supervised learning scheme; this is an interesting problem, which is beyond the scope of this paper.

A visual result on ICVL is shown in Figure 2 for stripes noise. Inference times are provided in Table 3, showing that our approach is computationally efficient.

**Results on real noise.** We also conducted a denoising experiment on the Urban dataset, reporting a visual result in Figure 3. Deep models were pre-trained on the APEX dataset, which has the same number of channels as Urban (even though the sensors are different), with band-dependent noise with $\sigma \in [0-55]$. Please note that for this experiment we did not use Noise Adaptive Sparse Coding 3.4 for T3SC, as it is highly dependent on the type of sensor used for training. We show that learning-based models trained on synthetic noise are able to transfer to real data.

**Comments on the additional results presented in the appendix.** The appendix also contains (i) results on the DCMall dataset including additional baselines mentioned above; (ii) error bars for parts of our experimental results in order to assess their statistical significance; (iii) an experiment when learning simultaneously on several datasets with different types of sensors showing that the second layer can be generic and effective at the same time; (iv) additional visual results; (v) various ablation studies to illustrate the importance of different components of our method.

**Broader Impact**

Our paper addresses the problem of denoising the signal, which is a key pre-processing step before using hyperspectral signals in concrete applications. As such, it is necessarily subject to dual use. For instance, HSI may be used for environmental monitoring, forestry, yield estimation in agriculture, natural disaster management planning, astronomy, archaeology, and medicine. Yet, HSI is also used by the petroleum industry for finding new oil fields, and has obvious military applications for surveillance. We believe the potential benefits of HSI for society are large enough to outweigh the

Table 2: Denoising performance on ICVL with various types of noise patterns. The first four rows correspond to i.i.d. Gaussian noise with fixed $\sigma$ per band. The next three rows corresponds to a noise level that depends on the band, taken uniformly on small interval. This is a blind-noise experiment since at test time, the noise level is unknown. The last two rows correspond to the scenarios with correlated $\sigma$ across bands, and with stripe noise, respectively. See main text for details.

| $\sigma$ | Metrics | Noisy | BM3D | BM4D | GLF | LLRT | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | MPSNR | 34.47 | 46.17 | 48.85 | 51.25 | 51.86 | **52.74** | 50.91 | 48.80 | <u>52.62</u> | 51.42 |
| | MSSIM | 0.7618 | 0.9843 | 0.9916 | 0.9949 | 0.9951 | **0.9960** | 0.9944 | 0.9918 | <u>0.9959</u> | 0.9952 |
| 25 | MPSNR | 21.44 | 37.86 | 39.89 | 43.16 | 43.43 | <u>44.74</u> | 42.83 | 44.20 | **45.38** | 44.73 |
| | MSSIM | 0.1548 | 0.9269 | 0.9510 | 0.9695 | 0.9746 | <u>0.9796</u> | 0.9700 | 0.9782 | **0.9825** | 0.9805 |
| 50 | MPSNR | 16.03 | 34.22 | 34.22 | 39.26 | 39.69 | 41.08 | 39.25 | <u>41.67</u> | **42.16** | 41.62 |
| | MSSIM | 0.0502 | 0.8654 | 0.8654 | 0.9197 | 0.9504 | 0.9602 | 0.9382 | <u>0.9655</u> | **0.9677** | 0.9646 |
| 100 | MPSNR | 10.85 | 30.43 | 32.47 | 34.79 | 36.39 | 37.55 | 35.64 | 37.19 | **38.99** | <u>38.50</u> |
| | MSSIM | 0.0144 | 0.7557 | 0.8155 | 0.7982 | 0.9182 | 0.9311 | 0.8815 | 0.9140 | **0.9439** | <u>0.9394</u> |
| [0-15] | MPSNR | 33.89 | 45.81 | 45.35 | 50.57 | 48.50 | 41.67 | 48.23 | <u>52.07</u> | **53.31** | 51.26 |
| | MSSIM | 0.6386 | 0.9767 | 0.9735 | 0.9948 | 0.9899 | 0.9078 | 0.9900 | <u>0.9957</u> | **0.9967** | 0.9955 |
| [0-55] | MPSNR | 23.36 | 39.06 | 38.43 | 44.22 | 41.13 | 32.94 | 41.76 | <u>47.13</u> | **48.64** | 46.82 |
| | MSSIM | 0.2601 | 0.9231 | 0.9074 | 0.9818 | 0.9580 | 0.7565 | 0.9620 | <u>0.9884</u> | **0.9911** | 0.9882 |
| [0-95] | MPSNR | 19.06 | 36.17 | 35.55 | 41.43 | 38.44 | 29.40 | 38.94 | 43.98 | **46.30** | <u>44.75</u> |
| | MSSIM | 0.1614 | 0.8760 | 0.8540 | 0.9674 | 0.9354 | 0.6609 | 0.9357 | 0.9753 | **0.9859** | <u>0.9822</u> |
| Corr. | MPSNR | 28.85 | 42.73 | 42.13 | 47.05 | 45.76 | 38.06 | 45.98 | <u>48.90</u> | **49.89** | 48.78 |
| | MSSIM | 0.4740 | 0.9599 | 0.9070 | 0.9881 | 0.9824 | 0.8536 | 0.9835 | <u>0.9911</u> | **0.9923** | 0.9911 |
| Strip. | MPSNR | 21.20 | 34.88 | 37.70 | 42.06 | 39.38 | 39.78 | 41.98 | <u>44.60</u> | **44.74** | 43.80 |
| | MSSIM | 0.1508 | 0.8641 | 0.9198 | 0.9628 | 0.9258 | 0.9333 | 0.9655 | **0.9806** | <u>0.9805</u> | 0.9773 |

Table 3: Inference time per image on ICVL with $\sigma = 50$; SMDS, QRNN3D and T3SC are using a V100 GPU; BM4D, GLF, LLRT and NGMeet are using an Intel(R) Xeon(R) CPU E5-1630 v4 @ 3.70GHz. Note that unlike GLF, NGMeet, and LRRT, learning-based approaches such as QRNN3D and our approach require a training procedure, which may be conducted offline. The cost of such a training step was about 13.5 hours for our method and 19 hours for QRNN3D on a V100 GPU.

| | BM3D | BM4D | GLF | LLRT | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|---|
| Inference time (s) | 1677 | 2382 | 5565 | 24384 | 2686 | 74.3 | **3.6** | <u>5.8</u> | 54.2 |



| (a) Groundtruth | (b) Noisy | (c) LLRT | (d) NGMeet |
|---|---|---|---|

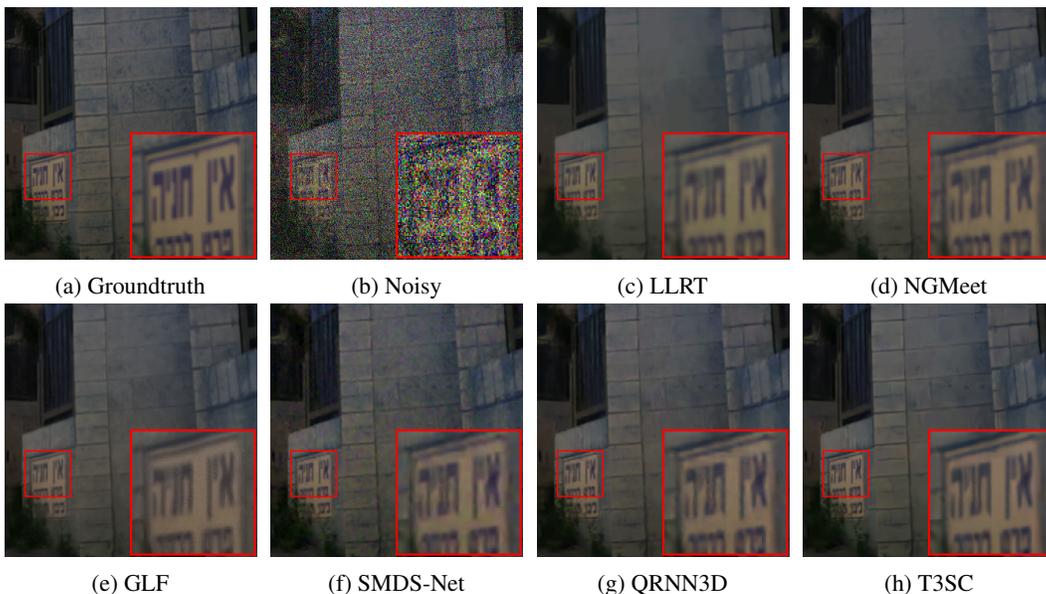| (e) GLF | (f) SMDS-Net | (g) QRNN3D | (h) T3SC |
|---|---|---|---|

Figure 2: Denoising results with Gaussian noise $\sigma = 25$ on ICVL with bands 9, 15, 28.

Table 4: Impact of the number of unrolled iterations per layer on denoising performances and inference time. This experiment was carried out on ICVL with $\sigma = 50$.

| Unrolled iterations per layer | 1 | 2 | 5 | 12 |
|---|---|---|---|---|
| MPSNR | 40.16 | 41.48 | 42.15 | 42.45 |
| Inference time (s) | 0.38 | 1.44 | 5.27 | 14.91 |



| (a) Input | (b) BM4D | (c) LLRT | (d) NGMeet |
|---|---|---|---|
| (e) GLF | (f) SMDS-Net | (g) QRNN3D | (h) T3SC |

Figure 3: Visual result on a real HSI denoising experiment on Urban dataset with bands 1, 108, 208.

potential harm. Nevertheless, we are planning to implement appropriate dissemination strategies to mitigate the risk of misuse for this work (notably with restrictive software licenses), while targeting a gold standard regarding the scientific reproducibility of our results.

# References

[1] T. Adão, J. Hruška, L. Pádua, J. Bessa, E. Peres, R. Morais, and J. J. Sousa. Hyperspectral imaging: A review on UAV-Based sensors, data processing and applications for agriculture and forestry. *Remote Sensing*, 9(11), 2017.

[2] H. Akbari, L. Halig, D. M. Schuster, B. Fei, A. Osunkoya, V. Master, P. Nieh, and G. Chen. Hyperspectral imaging and quantitative analysis for prostate cancer detection. *Journal of Biomedical Optics*, 17(7):1–11, 2012.

[3] L. Alparone, L. Wald, J. Chanussot, C. Thomas, P. Gamba, and L. M. Bruce. Comparison of pansharpening algorithms: Outcome of the 2006 grs-s data-fusion contest. *IEEE Transactions on Geoscience and Remote Sensing*, 45(10):3012–3021, 2007.

[4] B. Arad and O. Ben-Shahar. Sparse recovery of hyperspectral signal from natural RGB images. In *Proc. European Conference on Computer Vision (ECCV)*, 2016.

[5] N. Audebert, B. Saux, and S. Lefèvre. Deep learning for classification of hyperspectral data: A comparative review. *IEEE Geoscience and Remote Sensing Magazine*, 7(2):159–173, 2019.

[6] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research (JMLR)*, 18:1–43, 2018.

[7] J. M. Bioucas-Dias and J. M. P. Nascimento. Hyperspectral subspace identification. *IEEE Transactions on Geoscience and Remote Sensing*, 46(8):2435–2445, 2008.

[8] 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*, 1(2):6–36, 2013.

[9] A. Buades, B. Coll, and J.-M. Morel. A non-local algorithm for image denoising. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.

[10] A. Chakrabarti and T. Zickler. Statistics of real-world hyperspectral images. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011.

[11] Y. Chang, L. Yan, and S. Zhong. Hyper-laplacian regularized unidirectional low-rank tensor recovery for multispectral image denoising. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[12] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. Image Denoising by Sparse 3-D Transform-Domain Collaborative Filtering. *IEEE Transactions on Image Processing*, 16(8):2080–2095, 2007.

[13] C. F. Dantas, J. E. Cohen, and R. Gribonval. Hyperspectral image denoising using dictionary learning. In *2019 10th Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS)*, 2019.

[14] Q. Du, N. H. Younan, R. King, and V. P. Shah. On the performance evaluation of pan-sharpening techniques. *IEEE Geoscience and Remote Sensing Letters*, 4(4):518–522, 2007.

[15] M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image Processing*, 15(12):3736–3745, 2006.

[16] G. Elmasry, M. Kamruzzaman, D.-W. Sun, and Paul Allen. Principles and applications of hyperspectral imaging in quality evaluation of agro-food products: A review. *Critical Reviews in Food Science and Nutrition*, 52(11):999–1023, 2012.

[17] H. Fan, C. Li, Y. Guo, G. Kuang, and J. Ma. Spatial–Spectral total variation regularized low-rank tensor decomposition for hyperspectral image denoising. *IEEE Transactions on Geoscience and Remote Sensing*, 56(10):6196–6213, 2018.

[18] B. Fei. Chapter 3.6 - Hyperspectral imaging in medical applications. In J. M. Amigo, editor, *Hyperspectral Imaging*, volume 32 of *Data Handling in Science and Technology*, pages 523–565. Elsevier, 2020.

[19] Y.-Z. Feng and Da-Wen Sun. Application of hyperspectral imaging in food safety inspection and control: A review. *Critical Reviews in Food Science and Nutrition*, 52(11):1039–1058, 2012.

[20] M. Figueiredo and R. Nowak. An EM algorithm for wavelet-based image restoration. *IEEE Transactions on Image Processing*, 12(8):906–916, 2003.

[21] Y. Fu, A. Lam, I. Sato, and Y. Sato. Adaptive spatial-spectral dictionary learning for hyperspectral image denoising. In *Proc. International Conference on Computer Vision (ICCV)*, 2015.

[22] A. F. Goetz. Three decades of hyperspectral remote sensing of the Earth: A personal view. *Remote Sensing of Environment*, 113:S5–S16, 2009.

[23] X. Gong, W. Chen, and J. Chen. A low-rank tensor dictionary learning method for hyperspectral image denoising. *IEEE Transactions on Signal Processing*, 68:1168–1180, 2020.

[24] C. A. G. Gonzalez, O. Absil, and M. van Droogenbroeck. Supervised detection of exoplanets in high-contrast imaging sequences. *Astronomy & Astrophysics*, 613:A71, May 2018.

[25] A. A. Gowen, Y. Feng, E. Gaston, and V. Valdramidis. Recent applications of hyperspectral imaging in microbiology. *Talanta*, 137:43–54, 2015.

[26] K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In *Proc. International Conference on Machine Learning (ICML)*, 2010.

[27] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.

[28] W. He, Q. Yao, C. Li, N. Yokoya, and Q. Zhao. Non-local meets global: An integrated paradigm for hyperspectral denoising. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[29] K. I. Itten, F. Dell'Endice, A. Hueni, M. Kneubühler, D. Schläpfer, D. Odermatt, F. Seidel, S. Huber, J. Schopfer, T. Kellenberger, et al. Apex-the hyperspectral esa airborne prism experiment. *Sensors*, 8(10):6235–6259, 2008.

[30] J. Kerekes and J. Baum. Hyperspectral imaging system modeling. *Lincoln Laboratory Journal*, 14(1):117–130, 2003.

[31] N. Keshava and J. F. Mustard. Spectral unmixing. *IEEE signal processing magazine*, 19(1):44–57, 2002.

[32] Z. Kong, X. Yang, and L. He. A comprehensive comparison of multi-dimensional image denoising methods, 2020.

[33] S. Laine, T. Karras, J. Lehtinen, and T. Aila. High-quality self-supervised deep image denoising. In *Adv. Neural Information Processing Systems (NeurIPS)*, 2019.

[34] B. Lecouat, J. Ponce, and J. Mairal. Designing and learning trainable priors with non-cooperative games. In *Adv. Neural Information Processing Systems (NeurIPS)*, 2020.

[35] B. Lecouat, J. Ponce, and J. Mairal. Fully trainable and interpretable non-local sparse models for image restoration. In *Proc. European Conference on Computer Vision (ECCV)*, 2020.

[36] Y. Liu, H. Pu, and D.-W. Sun. Hyperspectral imaging technique for evaluating food quality and safety during various processes: A review of recent applications. *Trends in Food Science & Technology*, 69:25–35, 2017.

[37] B. Lu, P. D. Dao, J. Liu, Y. He, and J. Shang. Recent advances of hyperspectral imaging technology and applications in agriculture. *Remote Sensing*, 12(16), 2020.

[38] G. Lu and B. Fei. Medical hyperspectral imaging: A review. *Journal of Biomedical Optics*, 19(1):1–24, 2014.

[39] A. Maffei, J. M. Haut, M. E. Paoletti, J. Plaza, L. Bruzzone, and A. Plaza. A single model CNN for hyperspectral image denoising. *IEEE Transactions on Geoscience and Remote Sensing*, 58(4):2516–2529, 2020.

[40] M. Maggioni, V. Katkovnik, K. Egiazarian, and A. Foi. Nonlocal Transform-Domain Filter for Volumetric Data Denoising and Reconstruction. *IEEE Transactions on Image Processing*, 22(1):119–133, 2013.

[41] S. Mahesh, D. Jayas, J. Paliwal, and N. White. Hyperspectral imaging to classify and monitor quality of agricultural materials. *Journal of Stored Products Research*, 61:17–26, 2015.

[42] J. Mairal, F. Bach, and J. Ponce. Task-driven dictionary learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 34(4):791–804, 2012.

[43] J. Mairal, F. Bach, J. Ponce, et al. Sparse modeling for image and vision processing. *Foundations and Trends in Computer Graphics and Vision*, 8(2-3):85–283, 2014.

[44] J. Mairal, M. Elad, and G. Sapiro. Sparse representation for color image restoration. *IEEE Transactions on image processing*, 17(1):53–69, 2007.

[45] Y. Mäkinen, L. Azzari, and A. Foi. Exact transform-domain noise variance for collaborative filtering of stationary correlated noise. In *IEEE International Conference on Image Processing (ICIP)*, 2019.

[46] Y. Mäkinen, L. Azzari, and A. Foi. Collaborative filtering of correlated noise: Exact transform-domain variance for improved shrinkage and patch matching. *IEEE Transactions on Image Processing*, 29:8339–8354, 2020.

[47] D. G. Manolakis, R. B. Lockwood, and T. W. Cooley. *Hyperspectral Imaging Remote Sensing: Physics, Sensors, and Algorithms*. Cambridge University Press, Cambridge, 2016.

[48] B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.

[49] H. Othman and S.-E. Qian. Noise reduction of hyperspectral imagery using hybrid spatial-spectral derivative-domain wavelet shrinkage. *IEEE Transactions on Geoscience and Remote Sensing*, 44(2):397–408, 2006.

[50] Y. Peng, D. Meng, Z. Xu, C. Gao, Y. Yang, and B. Zhang. Decomposable nonlocal tensor dictionary learning for multispectral image denoising. In *Proc. Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.

[51] N. Qi, Y. Shi, X. Sun, and B. Yin. TenSR: Multi-dimensional tensor sparse representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[52] B. Rasti, P. Scheunders, P. Ghamisi, G. Licciardi, and J. Chanussot. Noise Reduction in Hyperspectral Imagery: Overview and Application. *Remote Sensing*, 10(3):482, Mar. 2018.

[53] B. Rasti, J. R. Sveinsson, M. O. Ulfarsson, and J. A. Benediktsson. Hyperspectral image denoising using first order spectral roughness penalty in wavelet domain. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 7(6):2458–2467, 2014.

[54] B. Rasti, M. O. Ulfarsson, and P. Ghamisi. Automatic hyperspectral image restoration using sparse and low-rank modeling. *IEEE Geoscience and Remote Sensing Letters*, 14(12):2335–2339, 2017.

[55] Q. Shi, X. Tang, T. Yang, R. Liu, and L. Zhang. Hyperspectral image denoising using a 3-D attention denoising network. *IEEE Transactions on Geoscience and Remote Sensing*, 2021.

[56] D. Simon and M. Elad. Rethinking the CSC model for natural images. In *Adv. Neural Information Processing Systems (NeurIPS)*, 2019.

[57] V. Studer, J. Bobin, M. Chahid, H. S. Mousavi, E. Candes, and M. Dahan. Compressive fluorescence microscopy for biological and hyperspectral imaging. *Proceedings of the National Academy of Sciences*, 109(26):E1679–E1687, 2012.

[58] L. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.

[59] G. R. L. (U.S.). Hypercube sample data set: Hydice sensor imagery urban.

[60] M. Wang, Q. Wang, J. Chanussot, and D. Hong. $L_0$-$l_1$ hybrid total variation regularization and its applications on hyperspectral image mixed noise removal and compressed sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2021.

[61] M. Wang, Q. Wang, J. Chanussot, and D. Li. Hyperspectral image mixed noise removal based on multidirectional low-rank modeling and spatial-spectral total variation. *IEEE Transactions on Geoscience and Remote Sensing*, 59(1):488–507, 2021.

[62] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.

[63] K. Wei, Y. Fu, and H. Huang. 3-D quasi-recurrent neural network for hyperspectral image denoising. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.

[64] F. Xiong, J. Zhou, J. Lu, and Y. Qian. SMDS-Net: Model Guided Spectral-Spatial Network for Hyperspectral Image Denoising. *arXiv:2012.01829*, 2020.

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

[66] Q. Yuan, L. Zhang, and H. Shen. Hyperspectral image denoising with a Spatial–Spectral view fusion strategy. *IEEE Transactions on Geoscience and Remote Sensing*, 52(5):2314–2325, 2014.

[67] Q. Yuan, Q. Zhang, J. Li, H. Shen, and L. Zhang. Hyperspectral Image Denoising Employing a Spatial–Spectral Deep Residual Convolutional Neural Network. *IEEE Transactions on Geoscience and Remote Sensing*, 57(2):1205–1218, Feb. 2019.

[68] H. Zhang, W. He, L. Zhang, H. Shen, and Q. Yuan. Hyperspectral image restoration using low-rank matrix recovery. *IEEE Transactions on Geoscience and Remote Sensing*, 52(8):4729–4743, 2014.

[69] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang. Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017.

[70] L. Zhang, L. Zhang, X. Mou, and D. Zhang. Fsim: A feature similarity index for image quality assessment. *IEEE Transactions on Image Processing*, 20(8):2378–2386, 2011.

[71] Y.-Q. Zhao and J. Yang. Hyperspectral image denoising via sparse representation and low-rank constraint. *IEEE Transactions on Geoscience and Remote Sensing*, 53(1):296–308, 2015.

[72] L. Zhuang and J. M. Bioucas-Dias. Hyperspectral image denoising based on global and non-local low-rank factorizations. In *IEEE International Conference on Image Processing (ICIP)*, 2017.

# Supplementary Material

## A  Implementation details

In this section, we provide additional implementation details, which are useful to reproduce our experiments (note that the code is also provided).

**Noise with spectrally correlated variance.**  For each band $i \in [\![0, c-1]\!]$, the standard deviation of the Gaussian noise is defined as :

$$\sigma_i = \beta \, exp \left[ -\frac{1}{4\eta^2} \left( \frac{i}{c} - \frac{1}{2} \right)^2 \right]$$

with $\beta = 23.08$ and $\eta = 0.157$.

**Preprocessing.**  A basic centering step is used for each input patch of our model. More precisely, for the first layer, each band of the input hyperspectral image is centered independently prior to patches extraction, and means are added back after decoding. For the second layer, patches are centered independently for each band (and similarly, the means are added back after decoding).

**Code and patch sizes**  The hyperparameters of our model are presented in Table 5.

| Layer | Patches size | Code size | Unrolled iterations | Rank |
|:---:|:---:|:---:|:---:|:---:|
| Spectral SC | $1 \times 1$ | 64 | 12 | 1 |
| Spectral-Spatial SC | $5 \times 5$ | 1024 | 5 | 3 |

Table 5: Architecture of our model

Table 14 shows that the combination of both layers is more performant than each layer independently.

**Initialization**  All parameters are initialized with He initialization [27].

**Blocks inference**  In order to apply our model to large images, we split them into blocks of size $256 \times 256$ with an overlap of 6 pixels. Each block is denoised independently. The output image is obtained by aggregating the denoised blocks. Pixels comprised in several blocks are averaged.

**Weights estimator**  For complex noise such as Gaussian noise with band-dependent variance or stripes noise, our model uses a CNN to estimate the weights $\beta_j$ associated with each band. The CNN operates on centered patches of size $56 \times 56$ both during training (random crops) and inference (blocks inference), and its architecture is described in Table 6.

| Layer | Kernel size | Stride | #filters | Output size |
|:---:|:---:|:---:|:---:|:---:|
| Inputs | | | | $1 \times 56 \times 56$ |
| Conv2D + ReLU | 5 | 2 | 64 | $64 \times 26 \times 26$ |
| MaxPooling2D | | | | $64 \times 13 \times 13$ |
| Conv2D + ReLU | 3 | 2 | 128 | $128 \times 6 \times 6$ |
| MaxPooling2D | | | | $128 \times 3 \times 3$ |
| Conv2D + Sigmoid | 3 | 1 | 1 | $1 \times 1 \times 1$ |

Table 6: CNN architecture for estimating $\beta_j$

The ablation study presented in Table 15 shows that this extension improves performances substantially for complex noise.

**Optimization**  Our models are trained with batch size of 16 for 60 epochs. We use the Adam optimizer, the initial learning rate is $3 \times 10^{-4}$, and is divided by two at epoch 30 and 45.

**Self-Supervised Learning**   During the training, $n$ bands randomly selected are masked simultaneously, and reconstructed from the available ones. The MSE loss is applied on the masked bands only. For testing, the $n$ masked bands are evenly distributed along the spectral dimension. All bands are reconstructed after $\lceil c/n \rceil$ iterations, where $c$ denotes the total number of bands and $\lceil \cdot \rceil$ is the ceiling operator. We used $n = 4$ for ICVL and $n = 16$ for Washington DC Mall.

For the SSL setting to be realistic, the noise is added to the clean image before patches are extracted. Otherwise, the model would have indirect access to the groundtruth by seeing the same patch with different noise realizations. As a consequence, the denoising task is much harder on complex noise when data is limited, as shown in Table 7.

# B   Additional Quantitative Results.

**Washington DC Mall dataset.**   Results for this dataset are presented in Table 7. Additional baselines are presented in Table 8. The conclusions are similar to those already drawn in the main paper.

Table 7: Denoising performances on Washington DC Mall.

| $\sigma$ | Metrics | Noisy | BM3D | BM4D | GLF | LLRT | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | MPSNR | 34.31 | 35.10 | 41.13 | 39.57 | 41.83 | 37.57 | 42.83 | <u>43.42</u> | **43.85** | 42.56 |
| | MSSIM | 0.9821 | 0.9875 | 0.9962 | 0.9953 | 0.9968 | 0.9928 | 0.9971 | <u>0.9973</u> | **0.9978** | 0.9967 |
| 25 | MPSNR | 20.70 | 24.51 | 31.08 | 35.25 | 34.95 | 35.38 | 35.64 | 35.04 | **36.74** | <u>35.92</u> |
| | MSSIM | 0.7688 | 0.8859 | 0.9690 | 0.9883 | 0.9863 | 0.9886 | 0.9889 | 0.9864 | **0.9912** | <u>0.9894</u> |
| 50 | MPSNR | 15.25 | 20.80 | 26.69 | 31.77 | 30.94 | 31.88 | 31.76 | 31.72 | **33.12** | <u>31.96</u> |
| | MSSIM | 0.5314 | 0.7508 | 0.9220 | 0.9761 | 0.9704 | 0.9759 | <u>0.9765</u> | 0.9741 | **0.9819** | 0.9762 |
| 100 | MPSNR | 10.48 | 17.65 | 22.51 | 27.81 | 26.82 | 27.86 | 28.02 | 27.41 | **29.48** | <u>28.04</u> |
| | MSSIM | 0.2888 | 0.5427 | 0.8141 | 0.9475 | 0.9322 | 0.9460 | <u>0.9491</u> | 0.9375 | **0.9618** | 0.9460 |
| [0-15] | MPSNR | 33.32 | 34.62 | 37.22 | 39.89 | 40.04 | 37.40 | 40.77 | **43.72** | <u>41.83</u> | 38.16 |
| | MSSIM | 0.9551 | 0.9746 | 0.9903 | 0.9950 | 0.9951 | 0.9926 | 0.9958 | **0.9971** | <u>0.9968</u> | 0.9917 |
| [0-55] | MPSNR | 22.45 | 26.11 | 29.04 | 38.37 | 33.36 | 32.55 | 34.31 | <u>38.44</u> | **39.28** | 31.93 |
| | MSSIM | 0.7450 | 0.8683 | 0.9504 | <u>0.9934</u> | 0.9811 | 0.9780 | 0.9859 | 0.9925 | **0.9945** | 0.9781 |
| [0-95] | MPSNR | 18.18 | 23.06 | 25.77 | <u>36.98</u> | 30.07 | 29.21 | 30.80 | 35.84 | **37.20** | 27.79 |
| | MSSIM | 0.5889 | 0.7688 | 0.9033 | <u>0.9914</u> | 0.9643 | 0.9589 | 0.9718 | 0.9877 | **0.9920** | 0.9561 |
| Corr. | MPSNR | 28.48 | 30.50 | 33.69 | 37.96 | 37.77 | 36.56 | 38.54 | <u>39.84</u> | **40.79** | 39.61 |
| | MSSIM | 0.9085 | 0.9515 | 0.9637 | 0.9928 | 0.9921 | 0.9911 | 0.9934 | <u>0.9944</u> | **0.9960** | <u>0.9944</u> |
| Strip. | MPSNR | 20.47 | 24.08 | 29.07 | 35.27 | 34.13 | 34.94 | 35.24 | <u>35.25</u> | **36.34** | 34.50 |
| | MSSIM | 0.7621 | 0.8672 | 0.9433 | 0.9877 | 0.9833 | <u>0.9876</u> | <u>0.9876</u> | 0.9874 | **0.9906** | 0.9853 |

Table 8: Denoising performances on Washington DC Mall with additional baselines.

| $\sigma$ | Metrics | Noisy | BM3D | BM4D | GLF | LLRT | NGMeet | 3D-ADNet | HSID-CNN | HSI-SDeCNN | SMDS-Net | QRNN3D | T3SC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | MPSNR | 34.31 | 35.10 | 41.13 | 39.57 | 41.83 | 37.57 | 42.08 | 41.68 | 39.98 | 42.83 | 43.42 | **43.85** |
| | MSSIM | 0.9821 | 0.9875 | 0.9962 | 0.9953 | 0.9968 | 0.9928 | 0.9968 | 0.9966 | 0.9954 | 0.9971 | 0.9973 | **0.9978** |
| 25 | MPSNR | 20.70 | 24.51 | 31.08 | 35.25 | 34.95 | 35.38 | 33.78 | 33.05 | 33.44 | 35.64 | 35.04 | **36.74** |
| | MSSIM | 0.7688 | 0.8859 | 0.9690 | 0.9883 | 0.9863 | 0.9886 | 0.9825 | 0.9813 | 0.9822 | 0.9889 | 0.9864 | **0.9912** |
| 50 | MPSNR | 15.25 | 20.80 | 26.69 | 31.77 | 30.94 | 31.88 | 29.73 | 28.96 | 29.61 | 31.76 | 31.72 | **33.12** |
| | MSSIM | 0.5314 | 0.7508 | 0.9220 | 0.9761 | 0.9704 | 0.9759 | 0.9587 | 0.9536 | 0.9608 | 0.9765 | 0.9741 | **0.9819** |
| 100 | MPSNR | 10.48 | 17.65 | 22.51 | 27.81 | 26.82 | 27.86 | 24.74 | 25.29 | 25.75 | 28.02 | 27.41 | **29.48** |
| | MSSIM | 0.2888 | 0.5427 | 0.8141 | 0.9475 | 0.9322 | 0.9460 | 0.9064 | 0.9014 | 0.9121 | 0.9491 | 0.9375 | **0.9618** |

**Study of statistical significance for the ICVL dataset.**   In order to evaluate the statistical significance of our results, we present some results in Table 9 for some of our models and baselines, by running models with five different random seeds. Note that we did not conduct such a study for all results in this paper in order to keep the computational cost of the project reasonable. The conclusions of the paper remain unchanged.

Table 9: Denoising performances on ICVL with multiple seeds

| $\sigma$ | Metrics | Noisy | GLF | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|
| 5 | MPSNR | $34.47 \pm 0.01$ | $51.25 \pm 0.01$ | $\mathbf{52.74 \pm 0.01}$ | $50.78 \pm 0.09$ | $49.54 \pm 1.28$ | $\underline{52.62 \pm 0.01}$ | $51.37 \pm 0.03$ |
|  | MSSIM | $0.7619 \pm 0.0001$ | $0.9951 \pm 0.0001$ | $\mathbf{0.9961 \pm 0.0001}$ | $0.9943 \pm 0.0001$ | $0.9924 \pm 0.0021$ | $\underline{0.9960 \pm 0.0001}$ | $0.9952 \pm 0.0001$ |
| 25 | MPSNR | $21.43 \pm 0.01$ | $43.16 \pm 0.01$ | $\underline{44.74 \pm 0.01}$ | $42.63 \pm 0.11$ | $44.20 \pm 0.16$ | $\mathbf{45.37 \pm 0.02}$ | $44.70 \pm 0.02$ |
|  | MSSIM | $0.1548 \pm 0.0002$ | $0.9696 \pm 0.0001$ | $0.9797 \pm 0.0001$ | $0.9687 \pm 0.0009$ | $0.9780 \pm 0.0009$ | $\mathbf{0.9825 \pm 0.0001}$ | $\underline{0.9805 \pm 0.0001}$ |
| 50 | MPSNR | $16.03 \pm 0.01$ | $39.26 \pm 0.01$ | $41.09 \pm 0.01$ | $39.09 \pm 0.08$ | $41.47 \pm 0.14$ | $\mathbf{42.16 \pm 0.01}$ | $\underline{41.62 \pm 0.01}$ |
|  | MSSIM | $0.0503 \pm 0.0001$ | $0.9198 \pm 0.0002$ | $0.9603 \pm 0.0001$ | $0.9359 \pm 0.0012$ | $0.9639 \pm 0.0012$ | $\mathbf{0.9677 \pm 0.0001}$ | $\underline{0.9648 \pm 0.0001}$ |
| 100 | MPSNR | $10.85 \pm 0.01$ | $34.78 \pm 0.01$ | $37.55 \pm 0.01$ | $35.59 \pm 0.04$ | $38.38 \pm 0.60$ | $\mathbf{38.99 \pm 0.01}$ | $\underline{38.51 \pm 0.01}$ |
|  | MSSIM | $0.0144 \pm 0.0001$ | $0.7981 \pm 0.0004$ | $0.9312 \pm 0.0001$ | $0.8781 \pm 0.0017$ | $0.9370 \pm 0.0114$ | $\mathbf{0.9439 \pm 0.0002}$ | $\underline{0.9397 \pm 0.0001}$ |
| [0-15] | MPSNR | $33.94 \pm 0.09$ | $50.68 \pm 0.11$ | $41.57 \pm 0.14$ | $48.00 \pm 0.13$ | $\underline{52.10 \pm 0.12}$ | $\mathbf{53.10 \pm 0.12}$ | $51.21 \pm 0.11$ |
|  | MSSIM | $0.6381 \pm 0.0013$ | $0.9950 \pm 0.0001$ | $0.9065 \pm 0.0022$ | $0.9899 \pm 0.0001$ | $\underline{0.9958 \pm 0.0001}$ | $\mathbf{0.9966 \pm 0.0001}$ | $0.9955 \pm 0.0002$ |
| [0-55] | MPSNR | $23.41 \pm 0.09$ | $44.41 \pm 0.12$ | $32.93 \pm 0.09$ | $41.42 \pm 0.18$ | $\underline{47.26 \pm 0.12}$ | $\mathbf{48.57 \pm 0.28}$ | $46.47 \pm 0.23$ |
|  | MSSIM | $0.2621 \pm 0.0025$ | $0.9820 \pm 0.0004$ | $0.7534 \pm 0.0031$ | $0.9593 \pm 0.0015$ | $\underline{0.9889 \pm 0.0004}$ | $\mathbf{0.9915 \pm 0.0005}$ | $0.9856 \pm 0.0024$ |
| [0-95] | MPSNR | $19.11 \pm 0.09$ | $41.62 \pm 0.11$ | $29.40 \pm 0.12$ | $38.86 \pm 0.06$ | $\underline{44.07 \pm 0.08}$ | $\mathbf{46.24 \pm 0.24}$ | $43.98 \pm 0.46$ |
|  | MSSIM | $0.1644 \pm 0.0031$ | $0.9667 \pm 0.0007$ | $0.6601 \pm 0.0051$ | $0.9352 \pm 0.0004$ | $\underline{0.9758 \pm 0.0003}$ | $\mathbf{0.9863 \pm 0.0005}$ | $0.9735 \pm 0.0049$ |

**CAVE dataset.** We report denoising performances of T3SC on the CAVE Dataset in Table 10 To evaluate T3SC, the dataset was divided in four splits : three were used for training and one for testing. The values reported for T3SC are averaged across all rotations of the test split.

Table 10: Denoising performances on CAVE dataset with Gaussian noise.

| $\sigma$ | Metrics | Noisy | NGMeet | T3SC |
|---|---|---|---|---|
| 5 | MPSNR | 35.05 | 47.96 | **49.16** |
| 25 | MPSNR | 21.99 | 42.44 | **42.77** |
| 50 | MPSNR | 16.37 | 38.89 | **39.7** |
| 100 | MPSNR | 10.96 | 34.99 | **36.48** |

**Joint training across heterogeneous datasets.** In Table 11, we study the problem of training a single model on three different datasets, APEX, DC Mall, and Pavia, involving a different number of channels. As mentioned in the paper, this model involves a common second layer and a spectral dictionary per dataset. These result show that most of the model parameters (which are present in the second layer) can in fact be shared across datasets without significant loss of accuracy when compared to the training of three different models (thus involving three times more parameters).

Table 11: Results for joint training experiment

| Training procedure | Model | Metrics | APEX | DC Mall | Pavia Center |
|---|---|---|---|---|---|
| Independant trainings | QRNN3D | MPSNR | 33.19 | 31.72 | 30.56 |
|  |  | MSSIM | 0.9619 | 0.9741 | 0.9569 |
|  | T3SC | MPSNR | **34.91** | **33.12** | **31.32** |
|  |  | MSSIM | **0.9730** | **0.9819** | **0.9617** |
| Joint training | QRNN3D | MPSNR | 31.95 | 30.97 | 29.12 |
|  |  | MSSIM | 0.9501 | 0.9690 | 0.9428 |
|  | T3SC | MPSNR | **34.74** | **33.08** | **31.30** |
|  |  | MSSIM | **0.9711** | **0.9819** | **0.9616** |

**Additional metrics.** Additional metrics are provided for the ICVL and DCMall datasets, respectively in Tables 12 and 13. The conclusions of the paper are unchanged.

| $\sigma$ | Metrics | Noisy | BM3D | BM4D | GLF | LLRT | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | MFSIM | 0.9953 | 0.9978 | 0.9986 | 0.9994 | <u>0.9995</u> | **0.9996** | 0.9993 | 0.9987 | **0.9996** | <u>0.9995</u> |
| 5 | MERGAS | 6.18 | 1.48 | 1.10 | 0.84 | 0.7740 | **0.69** | 0.87 | 1.14 | <u>0.70</u> | 0.83 |
| | MSAM | 0.2460 | 0.0518 | 0.0390 | 0.0267 | 0.0229 | **0.0211** | 0.0307 | 0.0412 | <u>0.0223</u> | 0.0286 |
| | MFSIM | 0.9218 | 0.9773 | 0.9829 | 0.9944 | 0.9942 | 0.9954 | 0.9921 | <u>0.9967</u> | **0.9970** | 0.9965 |
| 25 | MERGAS | 27.33 | 3.86 | 3.21 | 2.13 | 2.19 | <u>1.77</u> | 2.20 | 1.86 | **1.65** | 1.79 |
| | MSAM | 0.5989 | 0.1286 | 0.1005 | 0.0595 | 0.0459 | **0.0384** | 0.0717 | 0.0537 | <u>0.0406</u> | 0.0501 |
| | MFSIM | 0.8100 | 0.9488 | 0.9488 | 0.9851 | 0.9851 | 0.9863 | 0.9782 | **0.9928** | <u>0.9925</u> | 0.9914 |
| 50 | MERGAS | 51.48 | 5.88 | 5.88 | 3.33 | 3.92 | 2.71 | 3.33 | <u>2.50</u> | **2.40** | 2.56 |
| | MSAM | 0.7546 | 0.1964 | 0.1964 | 0.1029 | 0.0682 | **0.0505** | 0.1033 | 0.0571 | <u>0.0549</u> | 0.0663 |
| | MFSIM | 0.6471 | 0.8942 | 0.9008 | 0.9679 | 0.9637 | 0.9661 | 0.9456 | **0.9835** | <u>0.9824</u> | 0.9805 |
| 100 | MERGAS | 95.97 | 9.11 | 7.96 | 5.59 | 6.22 | 4.08 | 5.04 | <u>4.20</u> | **3.46** | 3.66 |
| | MSAM | 0.8619 | 0.2984 | 0.2228 | 0.1847 | 0.0919 | **0.0679** | 0.1441 | 0.1009 | <u>0.0761</u> | 0.0889 |
| | MFSIM | 0.9876 | 0.9954 | 0.9963 | 0.9991 | 0.9985 | 0.9965 | 0.9984 | <u>0.9995</u> | **0.9996** | 0.9993 |
| [0-15] | MERGAS | 10.11 | 1.91 | 2.07 | 0.98 | 1.17 | 4.53 | 1.20 | <u>0.79</u> | **0.69** | 0.91 |
| | MSAM | 0.3412 | 0.0680 | 0.0672 | 0.0328 | 0.0311 | 0.1772 | 0.0408 | <u>0.0265</u> | **0.0234** | 0.0322 |
| | MFSIM | 0.9087 | 0.9743 | 0.9768 | 0.9950 | 0.9900 | 0.9755 | 0.9890 | <u>0.9984</u> | **0.9985** | 0.9978 |
| [0-55] | MERGAS | 33.34 | 4.17 | 4.73 | 2.07 | 3.02 | 14.69 | 2.50 | <u>1.39</u> | **1.20** | 1.52 |
| | MSAM | 0.6478 | 0.1443 | 0.1412 | 0.0687 | 0.0636 | <u>0.4086</u> | 0.0784 | 0.0427 | **0.0370** | 0.0502 |
| | MFSIM | 0.8291 | 0.9524 | 0.9560 | 0.9911 | 0.9798 | 0.9536 | 0.9772 | <u>0.9969</u> | **0.9972** | 0.9962 |
| [0-95] | MERGAS | 54.92 | 5.83 | 6.73 | 2.86 | 4.64 | 24.82 | 3.46 | 2.17 | **1.58** | <u>1.92</u> |
| | MSAM | 0.7720 | 0.2001 | 0.1928 | 0.0992 | 0.0813 | 0.5574 | 0.1042 | 0.0622 | **0.0471** | <u>0.0615</u> |
| | MFSIM | 0.9704 | 0.9902 | 0.9923 | 0.9981 | 0.9968 | 0.9919 | 0.9969 | <u>0.9990</u> | **0.9991** | 0.9988 |
| Corr. | MERGAS | 14.20 | 2.61 | 3.74 | 1.46 | 1.63 | 6.37 | 1.55 | <u>1.12</u> | **1.02** | 1.15 |
| | MSAM | 0.4617 | 0.0934 | 0.1540 | 0.0468 | 0.0416 | 0.2550 | 0.0515 | <u>0.0316</u> | **0.0291** | 0.0367 |
| | MFSIM | 0.9068 | 0.9579 | 0.9736 | 0.9926 | 0.9871 | 0.9880 | 0.9900 | **0.9968** | <u>0.9965</u> | 0.9956 |
| Strip. | MERGAS | 28.14 | 7.65 | 4.65 | 2.52 | 4.34 | 4.32 | 2.44 | <u>1.78</u> | **1.77** | 2.00 |
| | MSAM | 0.6067 | 0.2197 | 0.1442 | 0.0764 | 0.1272 | 0.1298 | 0.0790 | **0.0439** | <u>0.0534</u> | 0.0631 |

Table 12: Additional metrics on ICVL

| $\sigma$ | Metrics | Noisy | BM3D | BM4D | GLF | LLRT | NGMeet | SMDS | QRNN3D | T3SC | T3SC-SSL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | MFSIM | 0.9534 | 0.9578 | 0.9772 | 0.9824 | 0.9817 | 0.9785 | 0.9802 | **0.9824** | <u>0.9814</u> | 0.9804 |
| 5 | MERGAS | 3.12 | 2.84 | 1.50 | 1.96 | 1.46 | 2.50 | 1.38 | <u>1.26</u> | **1.19** | 1.54 |
| | MSAM | 0.0862 | 0.0775 | 0.0427 | 0.0495 | 0.0395 | 0.0569 | 0.0373 | <u>0.0349</u> | **0.0329** | 0.0425 |
| | MFSIM | 0.8213 | 0.8676 | 0.9394 | <u>0.9661</u> | 0.9629 | 0.9655 | 0.9639 | 0.9614 | **0.9673** | 0.9648 |
| 25 | MERGAS | 14.96 | 9.50 | 4.55 | 2.91 | 3.31 | 2.94 | 2.87 | 3.08 | **2.50** | <u>2.77</u> |
| | MSAM | 0.3087 | 0.1753 | 0.1044 | 0.0684 | 0.0726 | 0.0671 | 0.0676 | 0.0709 | **0.0599** | <u>0.0668</u> |
| | MFSIM | 0.7174 | 0.7861 | 0.8974 | <u>0.9495</u> | 0.9439 | 0.9484 | 0.9464 | 0.9487 | **0.9542** | 0.9465 |
| 50 | MERGAS | 28.00 | 14.51 | 7.44 | 4.24 | 4.89 | 4.28 | 4.45 | 4.38 | **3.68** | <u>4.23</u> |
| | MSAM | 0.4785 | 0.2175 | 0.1438 | 0.0890 | 0.0925 | <u>0.0864</u> | 0.0944 | 0.0880 | **0.0768** | 0.0934 |
| | MFSIM | 0.6000 | 0.6821 | 0.8240 | 0.9188 | 0.9065 | <u>0.9209</u> | 0.9170 | 0.9100 | **0.9329** | 0.9163 |
| 100 | MERGAS | 48.42 | 20.83 | 11.98 | 6.54 | 7.58 | 6.66 | <u>6.52</u> | 7.01 | **5.51** | 6.52 |
| | MSAM | 0.6566 | 0.2700 | 0.1939 | 0.1183 | 0.1193 | <u>0.1147</u> | 0.1205 | 0.1297 | **0.0977** | 0.1265 |
| | MFSIM | 0.9338 | 0.9455 | 0.9690 | **0.9831** | 0.9774 | 0.9761 | 0.9787 | <u>0.9828</u> | 0.9782 | 0.9679 |
| [0-15] | MERGAS | 5.42 | 4.29 | 2.29 | 2.10 | 1.89 | 2.53 | 1.68 | **1.36** | <u>1.48</u> | 2.34 |
| | MSAM | 0.1358 | 0.1052 | 0.0610 | 0.0509 | 0.0487 | 0.0582 | 0.0438 | **0.0368** | <u>0.0395</u> | 0.0624 |
| | MFSIM | 0.8196 | 0.8642 | 0.9261 | **0.9766** | 0.9554 | 0.9523 | 0.9603 | 0.9714 | <u>0.9748</u> | 0.9507 |
| [0-55] | MERGAS | 18.46 | 10.41 | 5.56 | <u>2.37</u> | 3.86 | 4.19 | 3.22 | <u>2.37</u> | **2.05** | 4.73 |
| | MSAM | 0.3563 | 0.1879 | 0.1171 | <u>0.0572</u> | 0.0798 | 0.0961 | 0.0731 | 0.0581 | **0.0518** | 0.1226 |
| | MFSIM | 0.7471 | 0.8057 | 0.8837 | **0.9725** | 0.9377 | 0.9339 | 0.9473 | 0.9613 | <u>0.9689</u> | 0.9370 |
| [0-95] | MERGAS | 29.42 | 14.25 | 8.15 | <u>2.68</u> | 5.36 | 6.14 | 4.60 | 3.07 | **2.50** | 6.95 |
| | MSAM | 0.4899 | 0.2274 | 0.1466 | <u>0.0632</u> | 0.0973 | 0.1262 | 0.0962 | 0.0719 | **0.0604** | 0.1520 |
| | MFSIM | 0.9028 | 0.9229 | 0.9519 | <u>0.9783</u> | 0.9713 | 0.9693 | 0.9721 | **0.9790** | 0.9768 | 0.9744 |
| Corr. | MERGAS | 8.25 | 5.91 | 4.07 | 2.29 | 2.44 | 2.67 | 2.10 | <u>1.92</u> | **1.65** | 1.97 |
| | MSAM | 0.2049 | 0.1368 | 0.1106 | 0.0559 | 0.0593 | 0.0661 | 0.0540 | <u>0.0481</u> | **0.0436** | 0.0527 |
| | MFSIM | 0.8177 | 0.8621 | 0.9365 | **0.9663** | 0.9604 | 0.9649 | 0.9639 | 0.9619 | <u>0.9651</u> | 0.9582 |
| Strip. | MERGAS | 15.38 | 10.20 | 4.84 | 3.00 | 3.55 | 3.09 | <u>2.99</u> | 3.02 | **2.62** | 3.26 |
| | MSAM | 0.3152 | 0.1886 | 0.1101 | <u>0.0698</u> | 0.0794 | 0.0705 | 0.0700 | 0.0702 | **0.0623** | 0.0795 |

Table 13: Additional metrics on DCMall

**Ablation studies.** In this paragraph, we present different ablation studies, demonstrating in Table 14 that our two-layer model outperforms single-layer models. We also demonstrate the usefulness of our variant with weights $\beta_j$ in Table 15 when the noise variance varies a lot between different bands.

| Metrics | Noisy | Spec | SpecSpat | Spec + SpecSpat |
|---------|-------|------|----------|-----------------|
| MPSNR | 16.03 | 30.96 | 40.13 | **42.17** |
| MSSIM | 0.0502 | 0.6884 | 0.9533 | **0.9677** |
| MFSIM | 0.8100 | 0.9708 | 0.9849 | **0.9925** |
| MERGAS | 51.48 | 8.84 | 3.00 | **2.39** |
| MSAM | 0.7546 | 0.1300 | 0.1021 | **0.0547** |

Table 14: Combination of sparse coding layers: we denote by *Spec* the Spectral Sparse Coding layer and by *SpecSpat* the Spectral-Spatial Sparse Coding layer. This experiment was run on ICVL with $\sigma = 50$.

Table 15: Our model without/with band-wise noise estimator (NE) on ICVL with band-dependent Gaussian noise and stripes noise

|  | Metrics | T3SC | T3SC + NE |
|--------|---------|------|-----------|
| [0-15] | MPSNR | 52.85 | **53.31** |
|  | MSSIM | 0.9963 | **0.9967** |
| [0-55] | MPSNR | 47.39 | **48.64** |
|  | MSSIM | 0.9890 | **0.9911** |
| [0-95] | MPSNR | 44.92 | **46.30** |
|  | MSSIM | 0.9821 | **0.9859** |
| Strip. | MPSNR | 44.68 | **44.74** |
|  | MSSIM | 0.9801 | **0.9805** |

# C  Visual Examples

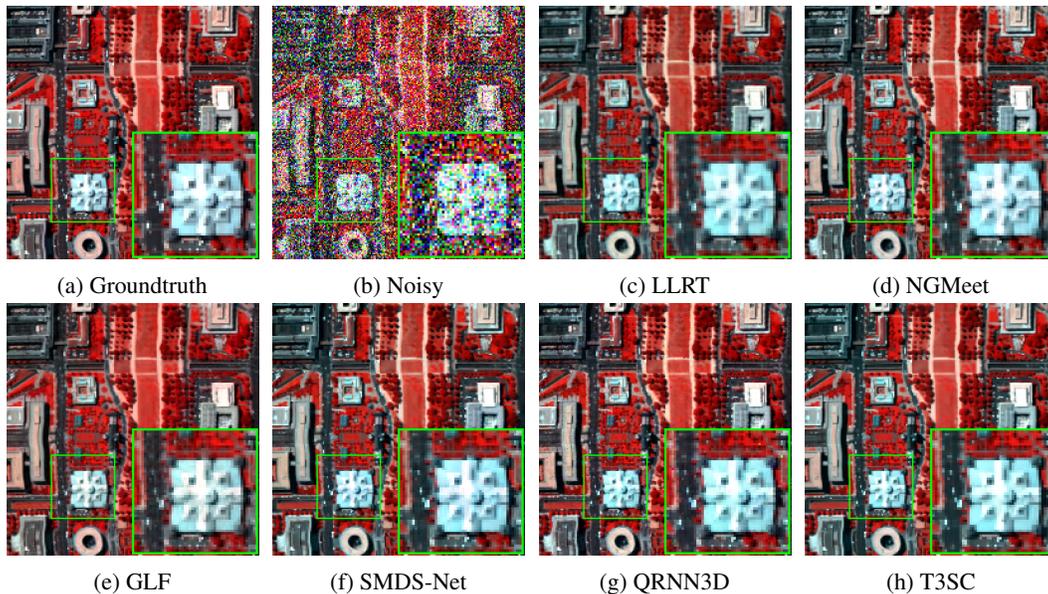Finally, we show additional visual examples in Figure 4 and 5.



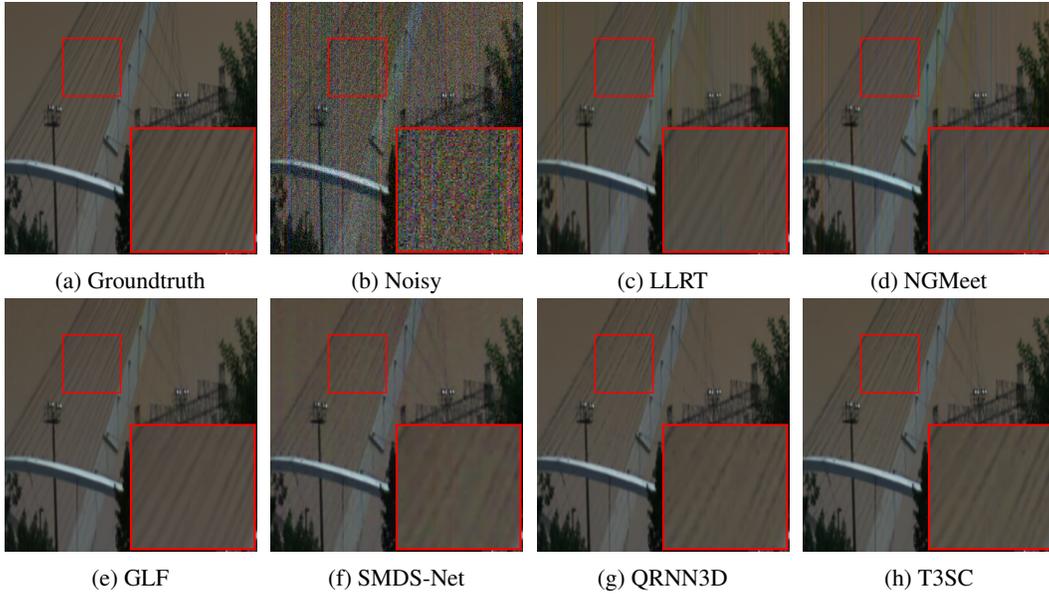Figure 4: Simulated Gaussian noise ($\sigma = 100$) on DCMall

Figure 5: Visual results for the denoising experiment with stripes noise on ICVL with bands 9, 15, 28.

# D GPU ressources

The total number of GPU hours involved in this project is around 19k hours on NVIDIA Tesla V100 16Go, including preliminary experiments, model design, final experiments and running baseline methods.