

H2TF for Hyperspectral Image Denoising: Where Hierarchical Nonlinear Transform Meets Hierarchical Matrix Factorization

Jiayi Li, Jinyu Xie, Yisi Luo, Xile Zhao, Jianli Wang

Abstract—Recently, tensor singular value decomposition (t-SVD) has emerged as a promising tool for hyperspectral image (HSI) processing. In the t-SVD, there are two key building blocks: (i) the low-rank enhanced transform and (ii) the accompanying low-rank characterization of transformed frontal slices. Previous t-SVD methods mainly focus on the developments of (i), while neglecting the other important aspect, i.e., the exact characterization of transformed frontal slices. In this letter, we exploit the potentiality in both building blocks by leveraging the Hierarchical nonlinear transform and the Hierarchical matrix factorization to establish a new Tensor Factorization (termed as H2TF). Compared to shallow counter partners, e.g., low-rank matrix factorization or its convex surrogates, H2TF can better capture complex structures of transformed frontal slices due to its hierarchical modeling abilities. We then suggest the H2TF-based HSI denoising model and develop an alternating direction method of multipliers-based algorithm to address the resultant model. Extensive experiments validate the superiority of our method over state-of-the-art HSI denoising methods.

Index Terms—Hyperspectral denoising, t-SVD, ADMM.

I. INTRODUCTION

HYPERSPECTRAL images (HSIs) inevitably contain mixed noise due to sensor failures or complex imaging conditions [1], [2], which seriously affects subsequent applications. Traditional hand-crafted HSI denoising methods, e.g., low-rankness [3], total variation (TV) [4], sparse representations [5], and non-local self-similarity [6], utilize interpretable domain knowledge to design generalizable regularizations for HSI denoising. Their representation abilities may be inferior to data-driven methods using deep neural networks (DNNs) [7]–[9], which can learn representative denoising mappings via supervised learning with abundant training pairs. However, supervised deep learning methods mostly neglect the prior information of HSIs, which sometimes results in generalization issues over different HSIs and various types of noise.

More recently, tensor singular value decomposition (t-SVD) attracts much attention in HSI denoising [10], [11]. The t-SVD views HSI as an implicit low-rank tensor and exploits the low-rankness in the transformed domain, which can more vividly characterize the structures of HSIs since it is flexible to select appropriate transforms and the accompanying low-rank characterization of the transformed frontal slices. Under such a framework, there are naturally two key building blocks: (i) The selection of the low-rank enhanced transform. A suitable transform can obtain a lower-rank transformed tensor and enhance

the recovery quality [12], [13]. (ii) The characterization of low-rankness of transformed frontal slices. The implicit low-rankness of HSIs is exploited by the low-rank modeling of frontal slices in the transformed domain.

Classical t-SVD-based methods mainly focused on the first building blocks, i.e., the design of different transforms. For example, the discrete Fourier transform (DFT) [14] was first used in the t-SVD, and then the discrete cosine transform (DCT) [15] was employed. Later methods exploited more representative and flexible transforms such as non-invertible transforms [16] and data-dependent transforms [17] to enhance the low-rankness of transformed frontal slices. These methods have achieved increasingly satisfactory results for HSI denoising [10], [11]. Nevertheless, these t-SVD methods pay less attention to the second building block, i.e., the exact characterization of transformed frontal slices. Specifically, they all employ shallow representations such as low-rank matrix factorization (MF) [13], QR factorization [18], and nuclear norm [12], [16] to characterize the transformed frontal slices.

In this work, we exploit a more representative formulation to capture complex structures of transformed frontal slices. Specifically, we leverage the hierarchical matrix factorization (HMF), which tailors a hierarchical formulation of learnable matrices along with nonlinear layers to capture each frontal slice in the transformed domain. The hierarchical modeling ability of HMF makes it more representative to capture the complex structures of HSIs. Meanwhile, we leverage the hierarchical nonlinear transform (HNT) to enhance the low-rankness of transformed frontal slices. With the Hierarchical nonlinear transform and Hierarchical matrix factorization, we develop a new Tensor Factorization method (termed as H2TF) under the t-SVD framework. Correspondingly, we develop the H2TF-based HSI denoising model. Attributed to the stronger representation abilities of HMF than shallow MF or its surrogates, our H2TF-based model can better capture fine details of the underlying clean HSI than conventional t-SVD-based methods. Thus, our model is expected to deliver better HSI denoising results. Meanwhile, the parameters of H2TF can be inferred from the observed noisy HSI in an unsupervised manner. In summary, the contributions of this letter are:

(i) We propose a new tensor factorization, i.e., the H2TF, which leverages the expressive power of two key building blocks—the HNT and the HMF, to respectively enhance the low-rankness of transformed data and characterize complex structures of transformed frontal slices. By virtue of their hierarchical modeling abilities, H2TF can faithfully capture

The authors are with the School of Mathematical Sciences, University of Electronic Science and Technology of China, Chengdu, China.

fine details of the clean HSI, and thus is beneficial for effectively removing heavy noise in the HSI.

(ii) We suggest an unsupervised H2TF-based HSI denoising model and develop an alternating direction method of multipliers (ADMM)-based algorithm. Extensive experiments on simulated and real-world data validate the superiority of our method over state-of-the-art (SOTA) HSI denoising methods, especially for details preserving and heavy noise removal.

II. THE PROPOSED H2TF

A. The t-SVD framework

We first introduce the general formulation of t-SVD. Suppose that the noisy HSI $\mathcal{Y} \in \mathbb{R}^{h \times w \times b}$ admits $\mathcal{Y} = \mathcal{X} + \mathcal{N}$, where \mathcal{X} denotes the clean HSI and \mathcal{N} denotes noise. To infer the underlying clean HSI \mathcal{X} from the observed \mathcal{Y} , t-SVD method generally formulates the following model:

$$\min_{\mathcal{Z}, \theta} L(\mathcal{Y}, \mathcal{X}) + \sum_k \psi(\mathcal{Z}^{(k)}), \text{ where } \mathcal{X} = \phi_\theta(\mathcal{Z}). \quad (1)$$

Here, $L(\mathcal{Y}, \mathcal{X})$ denotes the fidelity term and $\psi(\mathcal{Z}^{(k)})$ represents the low-rank characterization of $\mathcal{Z}^{(k)}$ (which denotes the k -th frontal (spatial) slice of $\mathcal{Z} \in \mathbb{R}^{h \times w \times b}$ [16]). $\phi_\theta(\cdot) : \mathbb{R}^{h \times w \times b} \rightarrow \mathbb{R}^{h \times w \times b}$ denotes a transform with learnable parameters θ , which transforms the low-rank representation \mathcal{Z} into the original domain. Sometimes the transform $\phi_\theta(\cdot)$ may not be learnable (e.g., the fixed DFT [14]), and in those situations the optimization variable only includes \mathcal{Z} .

The philosophy of the t-SVD model (1) is to minimize the rank in the transformed domain, which can model the implicit low-rankness of HSI. There are naturally two key building blocks for exactly modeling the implicit low-rankness, i.e., the selection of the transform $\phi_\theta(\cdot)$ and the exact low-rank characterization $\psi(\cdot)$ of the transformed frontal slice $\mathcal{Z}^{(k)}$. Most t-SVD-based methods focus on the design of different transforms $\phi_\theta(\cdot)$ (see examples in [13], [16], [17]), but all of them pay less attention to the exact characterization of the transformed frontal slice. They mostly adopt shallow representations to characterize $\mathcal{Z}^{(k)}$, e.g., MF [13], [19], QR factorization [18], and nuclear norm [15], [16]. However, these shallow representations may not be expressive enough to capture fine details of the clean HSI. Therefore, more representative methods are desired to enhance the representation abilities of the model in the transformed domain.

B. HMF for Characterizing $\mathcal{Z}^{(k)}$

To cope with this challenge, we leverage the HMF (hierarchical matrix factorization) to characterize $\mathcal{Z}^{(k)}$. The hierarchical modeling ability of HMF helps it more faithfully capture complex structures of the transformed frontal slice $\mathcal{Z}^{(k)}$ than shallow counter partners, e.g., SVD, MF, and QR factorization.

The standard MF used in previous t-SVD methods [13], [19] decomposes a low-rank matrix $\mathbf{Z} \in \mathbb{R}^{h \times w}$ into two factors as $\mathbf{Z} = \mathbf{W}_2 \mathbf{W}_1$, where $\mathbf{W}_2 \in \mathbb{R}^{h \times r}$, $\mathbf{W}_1 \in \mathbb{R}^{r \times w}$, and r is the rank. To model the hierarchical structures of \mathbf{Z} , we extend the MF to the product of multiple matrix factors $\{\mathbf{W}_d\}_{d=1}^l$:

$$\mathbf{Z} = \mathbf{W}_l \mathbf{W}_{l-1} \cdots \mathbf{W}_1, \quad (2)$$

where $\mathbf{W}_d \in \mathbb{R}^{r_d \times r_{d-1}}$, $r_l = h$, and $r_0 = w$. It was shown in [20] that such a linear HMF can induce an implicit low-rank

regularization on \mathbf{Z} when using gradient-based optimization. Generally, the larger l is (i.e., adding depth to the HMF), the tendency towards low-rank solutions goes stronger and oftentimes leads to better matrix recovery performances. Thus, the HMF is suitable to play the role of low-rank regularization in the t-SVD model (1).

Nevertheless, the linear HMF (2) may not be sufficient to capture nonlinear interactions inside HSIs. It motivates us to utilize the nonlinear HMF [21], [22] to model the low-rank matrix \mathbf{Z} via $\mathbf{Z} = \mathbf{W}_l \sigma(\mathbf{W}_{l-1} \cdots \mathbf{W}_3 \sigma(\mathbf{W}_2 \mathbf{W}_1))$, where $\sigma(\cdot)$ is a nonlinear scalar function. Classical HMF-based methods [20], [21] only utilize HMF to tackle the two-dimensional matrix. However, matrixing the HSI inevitably destroys its high-dimensional data structures. Therefore, we suggest tailoring b nonlinear HMFs to model the transformed tensor \mathcal{Z} by using each HMF to represent one of the frontal slices of \mathcal{Z} . Formally, we represent each frontal slice of \mathcal{Z} by

$$\mathcal{Z}^{(k)} = \mathcal{W}_l^{(k)} \sigma(\mathcal{W}_{l-1}^{(k)} \cdots \mathcal{W}_3^{(k)} \sigma(\mathcal{W}_2^{(k)} \mathcal{W}_1^{(k)})), k = 1, 2, \dots, b.$$

The above HMFs can be equivalently formulated as the tensor formulation $\mathcal{Z} = \mathcal{W}_l \Delta \sigma(\mathcal{W}_{l-1} \Delta \cdots \mathcal{W}_3 \Delta \sigma(\mathcal{W}_2 \Delta \mathcal{W}_1))$, where Δ is the tensor face-wise product [23] and $\{\mathcal{W}_d \in \mathbb{R}^{r_d \times r_{d-1} \times b}\}_{d=1}^l$ are some factor tensors.

Compared to shallow counter partners, e.g., MF, QR factorization, and nuclear norm, the above nonlinear HMF can better capture complex hierarchical structures of HSIs due to its nonlinear hierarchical modeling abilities, which helps to better recover fine details of HSI and remove heavy noise.

C. The Proposed H2TF

Next, we introduce our H2TF. Recall that two key building blocks in the t-SVD are the selection of the transform $\phi_\theta(\cdot)$ and the characterization of the transformed frontal slice $\mathcal{Z}^{(k)}$. We have leveraged the HMF to characterize $\mathcal{Z}^{(k)}$, and we further leverage the HNT (hierarchical nonlinear transform) as the first building block $\phi_\theta(\cdot)$:

$$\phi_\theta(\mathcal{Z}) := \sigma(\cdots \sigma(\mathcal{Z} \times_3 \mathbf{H}_1) \times_3 \cdots \times_3 \mathbf{H}_{m-1}) \times_3 \mathbf{H}_m,$$

where $\sigma(\cdot)$ is a nonlinear scalar function, $\theta := \{\mathbf{H}_p \in \mathbb{R}^{b \times b}\}_{p=1}^m$ are learnable parameters of HNT, and \times_3 is the mode-3 tensor-matrix product [24]. It was throughout demonstrated [13] that the HNT can effectively enhance the low-rankness of transformed tensor and thus obtain a better low-rank representation than shallow transforms (e.g., DFT [14] and DCT [15]), which benefits the implicit low-rank modeling.

Definition 1 (H2TF). Finally, we can define the following factorization modality of a certain low-rank tensor \mathcal{X} parameterized by $\{\mathcal{W}_d\}_{d=1}^l$ and $\{\mathbf{H}_p\}_{p=1}^m$:

$$\mathcal{X} = \phi_\theta \left(\underbrace{\mathcal{W}_l \Delta \sigma(\mathcal{W}_{l-1} \Delta \cdots \mathcal{W}_3 \Delta \sigma(\mathcal{W}_2 \Delta \mathcal{W}_1))}_{\text{Hierarchical matrix factorization}} \right), \quad (3)$$

$$\phi_\theta(\mathcal{Z}) := \underbrace{\sigma(\cdots \sigma(\mathcal{Z} \times_3 \mathbf{H}_1) \times_3 \cdots \times_3 \mathbf{H}_{m-1}) \times_3 \mathbf{H}_m}_{\text{Hierarchical nonlinear transform}},$$

which we call the H2TF representation of \mathcal{X} .

A general illustration of the proposed H2TF is shown in Fig. 1. H2TF benefits from the HMF to exploit complex hierarchical information of transformed frontal slices and the HNT to

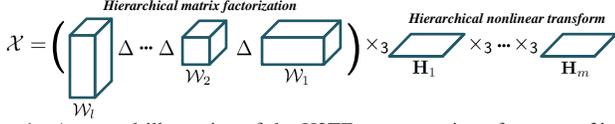


Fig. 1. A general illustration of the H2TF representation of a tensor \mathcal{X} . The nonlinear layer $\sigma(\cdot)$ is omitted for space consideration.

enhance the low-rankness in the transformed domain. With the hierarchical modeling abilities of H2TF, the characterization of HSIs would be more accurate. Therefore, H2TF can more faithfully capture fine details and rich textures of HSIs and remove heavy mixed noise. Now, we discuss the connections between H2TF and some popular matrix/tensor factorizations.

Remark 1. By changing the layer number of hierarchical matrix factorization (i.e., l) and the layer number of hierarchical nonlinear transform (i.e., m), H2TF includes many matrix/tensor factorizations as special cases:

(i) When $l = 2$, i.e., the HMF degenerates into the MF, our H2TF degenerates into the hierarchical low-rank tensor factorization [13]. (ii) When $m = 1$ and \mathbf{H}_m is an identity matrix (i.e., the transform $\phi_\theta(\cdot)$ is an identical mapping), our H2TF degenerates into the plain HMFs [21], [22] applied on each frontal slice of the tensor separately. In the following, we interpret this case as “ $m = 0$ ” since the transform is neglected. (iii) When $l = 2$ and $m = 1$ with \mathbf{H}_m being the fixed inverse DFT matrix, our H2TF degenerates into the classical low-tubal-rank tensor factorization [19], [25].

Moreover, H2TF can explicitly preserve the low-rankness of the tensor when omitting some nonlinearity, as stated below.

Lemma 1. Suppose that $\mathcal{X} = \phi(\mathcal{W}_l \Delta (\mathcal{W}_{l-1} \Delta \dots \Delta \mathcal{W}_1)) \in \mathbb{R}^{h \times w \times b}$, where $\{\mathcal{W}_d \in \mathbb{R}^{r_d \times r_{d-1} \times b}\}_{d=1}^l$ ($r_l = h$ and $r_0 = w$) are factor tensors, $\phi(\mathcal{Z}) := \mathcal{Z} \times_3 \mathbf{F}^{-1}$ is the inverse DFT, and \mathbf{F}^{-1} is the inverse DFT matrix (which is a special case of H2TF). Then we have $\text{rank}_t(\mathcal{X}) \leq \min\{r_0, r_1, \dots, r_l\}$, where $\text{rank}_t(\cdot)$ denotes the tensor tubal-rank [12]–[14].

Lemma 1 indicates that H2TF can preserve the low-rankness in the linear special case, where the degree of low-rankness (the upper bound of tubal-rank) is conditioned on the sizes of factor tensors. Therefore, we can readily control the degree of low-rankness by tuning the sizes of factor tensors in H2TF.

D. H2TF for HSI Denoising

H2TF is a potential tool for multi-dimensional data analysis and processing. We consider HSI denoising as a representative real-world application. By applying the H2TF representation (3) into (1), we can obtain the following HSI denoising model:

$$\min_{\{\mathcal{W}_d\}_{d=1}^l, \{\mathbf{H}_p\}_{p=1}^m} L(\mathcal{Y}, \mathcal{X}),$$

where $\mathcal{X} = \phi_\theta(\mathcal{W}_l \Delta \sigma(\mathcal{W}_{l-1} \Delta \dots \Delta \mathcal{W}_3 \Delta \sigma(\mathcal{W}_2 \Delta \mathcal{W}_1)))$.

In the HSI denoising problem, we consider the fidelity term as $L(\mathcal{Y}, \mathcal{X}) = \|\mathcal{Y} - \mathcal{X} - \mathcal{S}\|_F^2 + \alpha_1 \|\mathcal{S}\|_{\ell_1}$, where $\|\cdot\|_F^2$ denotes the Frobenius norm and we introduce $\mathcal{S} \in \mathbb{R}^{h \times w \times b}$ to represent sparse noise (often contains impulse noise and stripes). The ℓ_1 -norm enforces the sparsity on \mathcal{S} so that the sparse noise can be eliminated. Here, α_1 is a trade-off parameter.

Meanwhile, our H2TF can be readily combined with other proven techniques to enhance the denoising abilities. Here,

we consider the hybrid spatial-spectral TV (HSSTV) regularization [26] to further capture spatial and spatial-spectral local smoothness of HSIs. The HSSTV is formulated as $\|\mathcal{X}\|_{\text{HSSTV}} := \alpha_2 \|\mathcal{X}\|_{\text{TV}} + \alpha_3 \|\mathcal{X}\|_{\text{SSSTV}}$, where $\|\mathcal{X}\|_{\text{TV}} := \|\nabla_x \mathcal{X}\|_{\ell_1} + \|\nabla_y \mathcal{X}\|_{\ell_1}$, $\|\mathcal{X}\|_{\text{SSSTV}} := \|\nabla_x (\nabla_z \mathcal{X})\|_{\ell_1} + \|\nabla_y (\nabla_z \mathcal{X})\|_{\ell_1}$, and α_i ($i = 2, 3$) are trade-off parameters. Here, the derivative operators are defined as $(\nabla_x \mathcal{X})_{(i,j,k)} := \mathcal{X}_{(i+1,j,k)} - \mathcal{X}_{(i,j,k)}$, $(\nabla_y \mathcal{X})_{(i,j,k)} := \mathcal{X}_{(i,j+1,k)} - \mathcal{X}_{(i,j,k)}$, and $(\nabla_z \mathcal{X})_{(i,j,k)} := \mathcal{X}_{(i,j,k+1)} - \mathcal{X}_{(i,j,k)}$, where $\mathcal{X}_{(i,j,k)}$ denotes the (i, j, k) -th element of \mathcal{X} .

Based on the formulations of fidelity term and HSSTV, the proposed H2TF-based HSI denoising model is formulated as

$$\min_{\{\mathcal{W}_d\}_{d=1}^l, \{\mathbf{H}_p\}_{p=1}^m, \mathcal{S}} \|\mathcal{Y} - \mathcal{X} - \mathcal{S}\|_F^2 + \alpha_1 \|\mathcal{S}\|_{\ell_1} + \|\mathcal{X}\|_{\text{HSSTV}},$$

where $\mathcal{X} = \phi_\theta(\mathcal{W}_l \Delta \sigma(\mathcal{W}_{l-1} \Delta \dots \Delta \mathcal{W}_3 \Delta \sigma(\mathcal{W}_2 \Delta \mathcal{W}_1)))$. (4)

Compared to previous t-SVD-based HSI denoising methods [10], [11], H2TF has powerful representation abilities brought from the hierarchical structures and thus could better capture fine details of HSIs. Besides, the parameters of H2TF are unsupervisedly inferred from the noisy HSI by optimizing (4) without the requirement of training process.

E. ADMM-Based Algorithm

To tackle the problem (4), we develop an ADMM-based algorithm. By introducing auxiliary variables \mathcal{V}_i ($i = 1, 2, 3, 4$), (4) can be equivalently formulated as

$$\min_{\{\mathcal{W}_d\}_{d=1}^l, \{\mathbf{H}_p\}_{p=1}^m, \mathcal{S}, \{\mathcal{V}_i\}_{i=1}^4} \|\mathcal{Y} - \mathcal{X} - \mathcal{S}\|_F^2 + \alpha_1 \|\mathcal{S}\|_{\ell_1} + \alpha_2 \|\mathcal{V}_1\|_{\ell_1} + \alpha_2 \|\mathcal{V}_2\|_{\ell_1} + \alpha_3 \|\mathcal{V}_3\|_{\ell_1} + \alpha_3 \|\mathcal{V}_4\|_{\ell_1},$$

s.t. $\mathcal{V}_1 = \nabla_x \mathcal{X}$, $\mathcal{V}_2 = \nabla_y \mathcal{X}$, $\mathcal{V}_3 = \nabla_x (\nabla_z \mathcal{X})$, $\mathcal{V}_4 = \nabla_y (\nabla_z \mathcal{X})$, where $\mathcal{X} = \phi_\theta(\mathcal{W}_l \Delta \sigma(\mathcal{W}_{l-1} \Delta \dots \Delta \mathcal{W}_3 \Delta \sigma(\mathcal{W}_2 \Delta \mathcal{W}_1)))$. The corresponding augmented Lagrangian function is

$$\begin{aligned} & \mathcal{L}_\mu(\{\mathcal{W}_d\}_{d=1}^l, \{\mathbf{H}_p\}_{p=1}^m, \mathcal{S}, \{\mathcal{V}_i\}_{i=1}^4, \{A_i\}_{i=1}^4) \\ &= \|\mathcal{Y} - \mathcal{X} - \mathcal{S}\|_F^2 + \alpha_1 \|\mathcal{S}\|_{\ell_1} + \alpha_2 \|\mathcal{V}_1\|_{\ell_1} + \alpha_2 \|\mathcal{V}_2\|_{\ell_1} + \\ & \alpha_3 \|\mathcal{V}_3\|_{\ell_1} + \alpha_3 \|\mathcal{V}_4\|_{\ell_1} + \frac{\mu}{2} \|\nabla_x \mathcal{X} + \frac{A_1}{\mu} - \mathcal{V}_1\|_F^2 + \\ & \frac{\mu}{2} \|\nabla_y \mathcal{X} + \frac{A_2}{\mu} - \mathcal{V}_2\|_F^2 + \frac{\mu}{2} \|\nabla_x (\nabla_z \mathcal{X}) + \frac{A_3}{\mu} - \mathcal{V}_3\|_F^2 + \\ & \frac{\mu}{2} \|\nabla_y (\nabla_z \mathcal{X}) + \frac{A_4}{\mu} - \mathcal{V}_4\|_F^2, \end{aligned}$$

where μ is the penalty parameter, A_i ($i = 1, 2, 3, 4$) are multipliers, and \mathcal{X} is defined as in (3). The joint minimization problem can be decomposed into easier subproblems, followed by the update of Lagrangian multipliers.

The \mathcal{V}_i ($i = 1, 2, 3, 4$) subproblems are

$$\begin{cases} \min_{\mathcal{V}_1} \frac{\mu}{2} \|\nabla_x \mathcal{X}^t + \frac{A_1}{\mu} - \mathcal{V}_1\|_F^2 + \alpha_2 \|\mathcal{V}_1\|_{\ell_1} \\ \min_{\mathcal{V}_2} \frac{\mu}{2} \|\nabla_y \mathcal{X}^t + \frac{A_2}{\mu} - \mathcal{V}_2\|_F^2 + \alpha_2 \|\mathcal{V}_2\|_{\ell_1} \\ \min_{\mathcal{V}_3} \frac{\mu}{2} \|\nabla_x (\nabla_z \mathcal{X}^t) + \frac{A_3}{\mu} - \mathcal{V}_3\|_F^2 + \alpha_3 \|\mathcal{V}_3\|_{\ell_1} \\ \min_{\mathcal{V}_4} \frac{\mu}{2} \|\nabla_y (\nabla_z \mathcal{X}^t) + \frac{A_4}{\mu} - \mathcal{V}_4\|_F^2 + \alpha_3 \|\mathcal{V}_4\|_{\ell_1}, \end{cases}$$

which can be exactly solved by $\mathcal{V}_1^{t+1} = \text{Soft}_{\frac{\alpha_2}{\mu}}(\nabla_x \mathcal{X}^t + \frac{A_1}{\mu})$, $\mathcal{V}_2^{t+1} = \text{Soft}_{\frac{\alpha_2}{\mu}}(\nabla_y \mathcal{X}^t + \frac{A_2}{\mu})$, $\mathcal{V}_3^{t+1} = \text{Soft}_{\frac{\alpha_3}{\mu}}(\nabla_x (\nabla_z \mathcal{X}^t) + \frac{A_3}{\mu})$, and $\mathcal{V}_4^{t+1} = \text{Soft}_{\frac{\alpha_3}{\mu}}(\nabla_y (\nabla_z \mathcal{X}^t) + \frac{A_4}{\mu})$, where $(\text{Soft}_v(\mathcal{X}))_{(i,j,k)} := \text{sign}(\mathcal{X}_{(i,j,k)}) \max\{|\mathcal{X}_{(i,j,k)}| - v, 0\}$.

The \mathcal{S} subproblem is $\min_{\mathcal{S}} \|\mathcal{Y} - \mathcal{X}^t - \mathcal{S}\|_F^2 + \alpha_1 \|\mathcal{S}\|_{\ell_1}$, which can be exactly solved by $\mathcal{S}^{t+1} = \text{Soft}_{\frac{\alpha_1}{2}}(\mathcal{Y} - \mathcal{X}^t)$.

The \mathcal{X} subproblem is

$$\min_{\{\mathcal{W}_d\}_{d=1}^l, \{\mathbf{H}_p\}_{p=1}^m} \|\mathcal{Y} - \mathcal{X} - \mathcal{S}^t\|_F^2 + \frac{\mu}{2} (\|\nabla_x \mathcal{X} - \mathcal{D}_1^t\|_F^2 + \|\nabla_y \mathcal{X} - \mathcal{D}_2^t\|_F^2 + \|\nabla_x(\nabla_z \mathcal{X}) - \mathcal{D}_3^t\|_F^2 + \|\nabla_y(\nabla_z \mathcal{X}) - \mathcal{D}_4^t\|_F^2),$$

where $\mathcal{D}_i^t := \mathcal{V}_i^t - \frac{\Lambda_i^t}{\mu}$ ($i = 1, 2, 3, 4$). The clean HSI \mathcal{X} is parameterized by $\{\mathcal{W}_d\}_{d=1}^l$ and $\{\mathbf{H}_p\}_{p=1}^m$, as presented in Eq. (3). To tackle the nonlinear and nonconvex \mathcal{X} subproblem, we apply the adaptive moment estimation (Adam) algorithm [27]. In each iteration of the ADMM-based algorithm, we employ one step of the Adam to update $\{\mathcal{W}_d\}_{d=1}^l$ and $\{\mathbf{H}_p\}_{p=1}^m$.

Finally, the Lagrange multipliers are updated by $\Lambda_1^{t+1} = \Lambda_1^t + \mu(\nabla_x \mathcal{X}^t - \mathcal{V}_1^t)$, $\Lambda_2^{t+1} = \Lambda_2^t + \mu(\nabla_y \mathcal{X}^t - \mathcal{V}_2^t)$, $\Lambda_3^{t+1} = \Lambda_3^t + \mu(\nabla_x(\nabla_z \mathcal{X}^t) - \mathcal{V}_3^t)$, and $\Lambda_4^{t+1} = \Lambda_4^t + \mu(\nabla_y(\nabla_z \mathcal{X}^t) - \mathcal{V}_4^t)$.

III. EXPERIMENTS

A. Experimental Settings

We compare H2TF with SOTA model-based methods LRTDTV [28], SSTV-LRTF [11], RCTV [4], and HLRTF [13] and deep learning methods HSID-CNN [9] and SDeCNN [8]. We use the pre-trained models of HSID-CNN and SDeCNN provided by authors. All hyperparameters of these methods are carefully adjusted based on authors' suggestions to achieve the best results. We report the peak-signal-to-noise-ratio (PSNR) and structural similarity (SSIM). For more implementation details, please refer to supplementary materials.

We include four HSIs and three multi-spectral images (MSIs) as simulated datasets. The HSIs are *WDC* ($256 \times 256 \times 32$), *PaviaC* ($256 \times 256 \times 32$), *PaviaU* ($256 \times 256 \times 32$), and *Indian* ($145 \times 145 \times 32$). The MSIs are *Beads* ($256 \times 256 \times 31$), *Cloth* ($256 \times 256 \times 31$), and *Cups* ($256 \times 256 \times 31$) in the CAVE dataset [29]. The noise settings of simulated data are explained as below. **Case 1:** All bands are added with Gaussian noise of standard deviation 0.2. **Case 2:** The Gaussian noise for Case 1 is kept. Besides, all bands are added with impulse noise with sampling rate 0.1. **Case 3:** The same as Case 2 plus 50% of bands corrupted by deadlines. The number of deadlines for each chosen band is generated randomly from 6 to 10, and their spatial width is chosen randomly from 1 to 3. **Case 4:** The same as Case 2 plus 40% of bands corrupted by stripes. The number of stripes in each corrupted band is chosen randomly from 6 to 15. **Case 5:** The same as Case 2 plus both the deadlines in Case 3 and the stripes in Case 4. To test our method in real scenarios, we choose two real-world noisy HSIs *Shanghai* ($300 \times 300 \times 32$) and *Urban* ($307 \times 307 \times 32$) as real-world experimental datasets.

B. Experimental Results

1) *Results:* The quantitative results on simulated data are reported in Table I. Our H2TF obtains better quantitative results than other competitors. H2TF outperforms other TV and tensor factorization-based methods (LRTDTV, SSTV-LRTF, RCTV, and HLRTF), which shows the stronger representation abilities of H2TF than existing shallow tensor factorizations thanks to the hierarchical structures of H2TF. Some visual results on simulated and real data are shown in Figs. 2-3.

TABLE I
AVERAGE QUANTITATIVE DENOISING RESULTS BY DIFFERENT METHODS.

Dataset	Method	Case 1		Case 2		Case 3		Case 4		Case 5	
		PSNR	SSIM								
HSIs	LRTDTV	30.88	0.888	29.55	0.849	28.29	0.825	29.35	0.843	28.31	0.824
	SSTV-LRTF	30.77	0.887	30.35	0.879	28.32	0.839	29.51	0.859	27.42	0.812
<i>WDC</i>	HSID-CNN	29.61	0.863	22.89	0.691	21.98	0.661	22.16	0.669	21.22	0.635
	SDeCNN	30.26	0.873	23.97	0.735	23.33	0.725	23.41	0.723	22.63	0.714
<i>PaviaC</i>	RCTV	29.53	0.853	29.05	0.839	26.75	0.782	28.47	0.826	26.34	0.772
<i>PaviaU</i>	HLRTF	30.12	0.868	29.65	0.855	29.58	0.853	29.15	0.846	29.03	0.841
<i>Indian</i>	H2TF	32.51	0.919	31.41	0.900	31.34	0.899	30.83	0.893	30.84	0.892
MSIs	LRTDTV	28.85	0.889	27.05	0.838	26.31	0.828	26.83	0.830	26.13	0.819
	SSTV-LRTF	27.64	0.878	27.48	0.864	26.25	0.855	26.79	0.844	25.11	0.823
<i>Beads</i>	HSID-CNN	25.86	0.827	21.22	0.660	20.97	0.645	20.68	0.646	20.34	0.626
	SDeCNN	28.43	0.886	22.04	0.715	22.32	0.709	21.53	0.706	21.70	0.698
<i>Cloth</i>	RCTV	28.15	0.869	27.49	0.866	25.77	0.839	26.98	0.854	25.46	0.829
<i>Cups</i>	HLRTF	29.21	0.884	28.73	0.886	28.67	0.884	28.10	0.870	28.03	0.868
	H2TF	31.51	0.940	29.46	0.906	29.47	0.901	29.22	0.908	29.03	0.896

H2TF generally outperforms other competitors in two aspects. First, H2TF can more effectively remove heavy mixed noise. Second, H2TF preserves fine details of HSIs better than other methods. The superior performances of H2TF are mainly due to its hierarchical modeling abilities, which help to better characterize fine details of HSI and robustly capture the underlying structures of HSI under extremely heavy noise. More visual results can be found in supplementary.

2) *Discussions:* The HMF is an important building block in H2TF. We test the influence of the layer number of HMF (i.e., l); see Fig. 4 (a). A suitable layer number of HMF (e.g., $l = 5$) can obtain both good performances and a lightweight model. The HNT is another important building block. We change the layer number of HNT to test its influence; see Fig. 4 (b). Also, a proper layer number of HNT (e.g., $m = 2$) can bring good performances. According to Lemma 1, the sizes of factor tensors in HMF, i.e., $\{r_d\}_{d=1}^4$, determine the degree of low-rankness. Hence, we test such connections by changing the sizes of factor tensors; see Fig. 4 (c) (Here, r_0 and r_5 are fixed as the sizes of observed data and $\{r_d\}_{d=1}^4$ are selected in $\{(1, 2, 4, 8), (2, 4, 8, 16), (3, 6, 12, 24), \dots, (20, 40, 80, 160)\}$). When the sizes (rank) are too small, the model lacks representation abilities and when the sizes (rank) are too large, the model overfits. Nevertheless, our method is quite robust w.r.t. $\{r_d\}_{d=1}^4$.

IV. CONCLUSIONS

We propose the H2TF for HSI denoising. Our H2TF simultaneously leverages the hierarchical matrix factorization and the hierarchical nonlinear transform to compactly represent HSIs with powerful representation abilities, which can more faithfully capture fine details of HSIs than classical tensor factorization methods. Comprehensive experiments validate the superiority of H2TF over SOTA methods, especially for HSI details preserving and heavy noise removal.

REFERENCES

- [1] N. Liu, W. Li, R. Tao, and J. E. Fowler, "Wavelet-domain low-rank/group-sparse destriping for hyperspectral imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 12, pp. 10310–10321, 2019.
- [2] W. He, H. Zhang, L. Zhang, W. Philips, and W. Liao, "Weighted sparse graph based dimensionality reduction for hyperspectral images," *IEEE Geoscience and Remote Sensing Letters*, vol. 13, no. 5, pp. 686–690, 2016.
- [3] F. Xu, Y. Chen, C. Peng, Y. Wang, X. Liu, and G. He, "Denoising of hyperspectral image using low-rank matrix factorization," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 7, pp. 1141–1145, 2017.

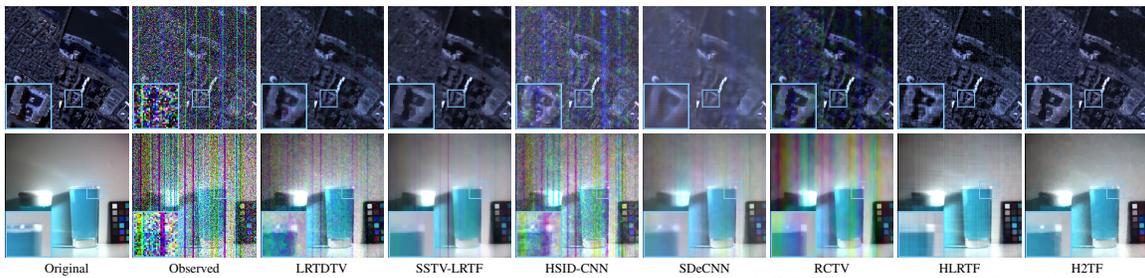


Fig. 2. Pseudo-color images of HSI denoising results by different methods on simulated data *PaviaC* Case 4 (first row) and *Cups* Case 5 (second row).

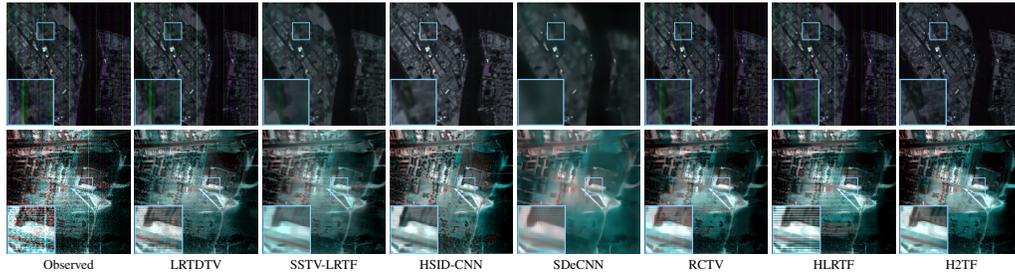


Fig. 3. Pseudo-color images of HSI denoising results by different methods on real-world data *Shanghai* (first row) and *Urban* (second row).

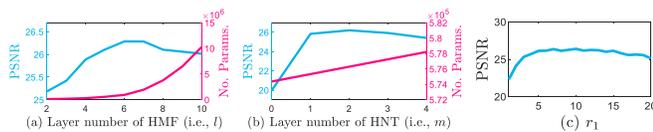


Fig. 4. Results on *Beads* Case 5 with (a) different layer number of HMF, (b) different layer number of HNT, and (c) different sizes of factor tensors.

[4] J. Peng, H. Wang, X. Cao, X. Liu, X. Rui, and D. Meng, "Fast noise removal in hyperspectral images via representative coefficient total variation," *IEEE Transactions on Geoscience and Remote Sensing*, 2022. doi=10.1109/TGRS.2022.3229012.

[5] B. Rasti, M. O. Ulfarsson, and P. Ghamisi, "Automatic hyperspectral image restoration using sparse and low-rank modeling," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 12, pp. 2335–2339, 2017.

[6] W. He, Q. Yao, C. Li, N. Yokoya, Q. Zhao, H. Zhang, and L. Zhang, "Non-local meets global: An iterative paradigm for hyperspectral image restoration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 4, pp. 2089–2107, 2022.

[7] Y. Chang, L. Yan, H. Fang, S. Zhong, and W. Liao, "HSI-DeNet: Hyperspectral image restoration via convolutional neural network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 2, pp. 667–682, 2019.

[8] 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*, vol. 58, no. 4, pp. 2516–2529, 2020.

[9] 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*, vol. 57, no. 2, pp. 1205–1218, 2019.

[10] Y.-B. Zheng, T.-Z. Huang, X.-L. Zhao, T.-X. Jiang, T.-H. Ma, and T.-Y. Ji, "Mixed noise removal in hyperspectral image via low-fibered-rank regularization," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 1, pp. 734–749, 2020.

[11] 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*, vol. 56, no. 10, pp. 6196–6213, 2018.

[12] J.-L. Wang, T.-Z. Huang, X.-L. Zhao, T.-X. Jiang, and M. K. Ng, "Multi-dimensional visual data completion via low-rank tensor representation under coupled transform," *IEEE Transactions on Image Processing*, vol. 30, pp. 3581–3596, 2021.

[13] Y. Luo, X. Zhao, D. Meng, and T. Jiang, "HLRTF: Hierarchical low-rank tensor factorization for inverse problems in multi-dimensional imaging," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19281–19290, 2022.

[14] C. Lu, J. Feng, Y. Chen, W. Liu, Z. Lin, and S. Yan, "Tensor robust principal component analysis with a new tensor nuclear norm," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 4, pp. 925–938, 2020.

[15] C. Lu, X. Peng, and Y. Wei, "Low-rank tensor completion with a new tensor nuclear norm induced by invertible linear transforms," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5989–5997, 2019.

[16] T.-X. Jiang, M. K. Ng, X.-L. Zhao, and T.-Z. Huang, "Framelet representation of tensor nuclear norm for third-order tensor completion," *IEEE Transactions on Image Processing*, vol. 29, pp. 7233–7244, 2020.

[17] H. Kong, C. Lu, and Z. Lin, "Tensor Q-rank: new data dependent definition of tensor rank," *Machine Learning*, vol. 110, pp. 1867–1900, 2021.

[18] Y. Zheng and A.-B. Xu, "Tensor completion via tensor QR decomposition and L2,1-norm minimization," *Signal Processing*, vol. 189, p. 108240, 2021.

[19] P. Zhou, C. Lu, Z. Lin, and C. Zhang, "Tensor factorization for low-rank tensor completion," *IEEE Transactions on Image Processing*, vol. 27, no. 3, pp. 1152–1163, 2018.

[20] S. Arora, N. Cohen, W. Hu, and Y. Luo, "Implicit regularization in deep matrix factorization," in *Advances in Neural Information Processing Systems*, vol. 32, pp. 7413–7424, 2019.

[21] Z. Li, T. Sun, H. Wang, and B. Wang, "Adaptive and implicit regularization for matrix completion," *SIAM Journal on Imaging Sciences*, vol. 15, no. 4, pp. 2000–2022, 2022.

[22] J. Fan and J. Cheng, "Matrix completion by deep matrix factorization," *Neural Networks*, vol. 98, pp. 34–41, 2018.

[23] E. Kerfeld, M. Kilmer, and S. Aeron, "Tensor-tensor products with invertible linear transforms," *Linear Algebra and its Applications*, vol. 485, pp. 545–570, 2015.

[24] M. E. Kilmer, K. Braman, N. Hao, and R. C. Hoover, "Third-order tensors as operators on matrices: A theoretical and computational framework with applications in imaging," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 1, pp. 148–172, 2013.

[25] X.-Y. Liu, S. Aeron, V. Aggarwal, and X. Wang, "Low-tubal-rank tensor completion using alternating minimization," *IEEE Transactions on Information Theory*, vol. 66, no. 3, pp. 1714–1737, 2019.

[26] S. Takeyama, S. Ono, and I. Kumazawa, "Mixed noise removal for hyperspectral images using hybrid spatio-spectral total variation," in *IEEE International Conference on Image Processing*, pp. 3128–3132, 2019.

[27] D. Kingma and J. Ba, "ADAM: A method for stochastic optimization," in *International Conference on Learning Representations*, 2014.

[28] Y. Wang, J. Peng, Q. Zhao, Y. Leung, X.-L. Zhao, and D. Meng, "Hyperspectral image restoration via total variation regularized low-rank tensor decomposition," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 4, pp. 1227–1243, 2017.

[29] F. Yasuma, T. Mitsunaga, D. Iso, and S. K. Nayar, "Generalized assorted pixel camera: Postcapture control of resolution, dynamic range, and spectrum," *IEEE Transactions on Image Processing*, vol. 19, no. 9, pp. 2241–2253, 2010.