

# Model Inspired Autoencoder for Unsupervised Hyperspectral Image Super-Resolution

Jianjun Liu, *Member, IEEE*, Zebin Wu, *Senior Member, IEEE*, Liang Xiao, *Member, IEEE*, and Xiao-Jun Wu, *Member, IEEE*

**Abstract**—This paper focuses on hyperspectral image (HSI) super-resolution that aims to fuse a low-spatial-resolution HSI and a high-spatial-resolution multispectral image to form a high-spatial-resolution HSI (HR-HSI). Existing deep learning-based approaches are mostly supervised that rely on a large number of labeled training samples, which is unrealistic. The commonly used model-based approaches are unsupervised and flexible but rely on hand-craft priors. Inspired by the specific properties of model, we make the first attempt to design a model inspired deep network for HSI super-resolution in an unsupervised manner. This approach consists of an implicit autoencoder network built on the target HR-HSI that treats each pixel as an individual sample. The nonnegative matrix factorization (NMF) of the target HR-HSI is integrated into the autoencoder network, where the two NMF parts, spectral and spatial matrices, are treated as decoder parameters and hidden outputs respectively. In the encoding stage, we present a pixel-wise fusion model to estimate hidden outputs directly, and then reformulate and unfold the model's algorithm to form the encoder network. With the specific architecture, the proposed network is similar to a manifold prior-based model, and can be trained patch by patch rather than the entire image. Moreover, we propose an additional unsupervised network to estimate the point spread function and spectral response function. Experimental results conducted on both synthetic and real datasets demonstrate the effectiveness of the proposed approach.

**Index Terms**—Super-resolution, hyperspectral image, autoencoder, unfolding, nonnegative matrix factorization.

## I. INTRODUCTION

**H**yperspectral image (HSI) is a kind of three dimensional image taken at different spectral bands, with its spectral range covering hundreds of contiguous and narrow bands that span the visible to infrared spectrum. The high spectral resolution of HSIs promotes various applications, such as material identification. Due to the limited incident energy, there is always a tradeoff between spectral resolution, spatial resolution and signal-to-noise ratio of images when designing the imaging sensors [1]–[6]. Thus, the spatial resolution of

HSIs is usually sacrificed, which impedes the subsequent tasks. Conversely, conventional multispectral images (MSIs) at much lower spectral resolution can be acquired with higher spatial resolution. An economical HSI super-resolution solution is to instead record a low-spatial-resolution HSI (LR-HSI) and a high-spatial-resolution MSI (HR-MSI), and to fuse them into a target high-spatial-resolution HSI (HR-HSI) [2], [3], [5].

HSI super-resolution that fuses a LR-HSI with a HR-MSI has attracted great attention [2], [3], [5]. This fusion problem arises from the Pansharpening problem that fuses a low-spatial-resolution MSI or HSI with a high-spatial-resolution panchromatic image [1], [4], [6]. Generally, the conventional approaches proposed for the Pansharpening problem can be extended to solve HSI super-resolution, but the fusion process of HSI super-resolution is relatively more complicated than that of Pansharpening due to the rich spectral information. Related fusion approaches can be roughly divided into four categories: component substitution [7], multiresolution analysis [8], model-based approaches and deep learning-based approaches. Among these categories, the research of model-based approaches is the most classic one, and deep learning-based approaches have been the most active one recently.

Model-based approaches consider building optimization models to obtain the target image. Given two observed images, they design fidelity terms and exploit spectral/spatial priors to enforce the desired result. Some approaches treat the target image as a variable and recover the target image entirely, such as group spectral embedding [9], clustering manifold structure [10], nonlocal patch tensor sparse representation [11], and structured sparse low-rank representation [12]. Most approaches consider separating the target image into parts and regenerating it via the recovered parts. There are many decomposition strategies by making assumptions about the target image. Examples are, that it lives in a low-dimensional subspace and the subspace-based models are solved by exploiting prior knowledge, such as piecewise smooth [13], dictionary learning [14], tensor multi-rank [15], and truncated matrix decomposition [16]; or that it can be represented linearly by pure spectral signatures and the endmember and abundance matrices are recovered simultaneously [17]–[20]; or that it can be sparsely represented by an over-complete spectral dictionary and different priors are used to obtain the spectral dictionary and coefficients [21]–[25]; or by approaches that separate the target image by tensor decomposition and update each component iteratively [26]–[29]. Moreover, there are some approaches that build models to estimate the point spread function (PSF) and spectral response function (SRF)

This work was supported in part by the National Natural Science Foundation of China under Grant No. 62071204, 61871226 and 61772274, by the Natural Science Foundation of Jiangsu Province under Grant No. BK20201338 and BK20180018, by the Jiangsu Provincial Social Developing Project under Grant No. BE2018727, by the China Postdoctoral Science Foundation under Grant No. 2021M691275, and by the Jiangsu Postdoctoral Research Funding Program under Grant No. 2021K148B.

Jianjun Liu and Xiao-Jun Wu are with the Jiangsu Provincial Engineering Laboratory for Pattern Recognition and Computational Intelligence, Jiangnan University, Wuxi, China (Email: liuofficial@163.com, wu\_xiaojun@jiangnan.edu.cn).

Zebin Wu and Liang Xiao are with the School of Computer Science, Nanjing University of Science and Technology, Nanjing, China (Email: zebin.wu@gmail.com, xiaoliang@mail.njust.edu.cn).

[13], [30]. The entire process of model-based approaches is unsupervised. Although these models are flexible and their theory is relatively complete, they rely on hand-craft priors and there are many empirical parameters to tune.

Deep learning-based approaches are data-driven. They build deep neural networks to solve the related fusion problems, and produce the target image by feeding observed images into the network. Some approaches enhance the ability to fuse images in the network structures, such as 3D convolutional neural networks (CNNs) [31], residual networks [32], multiscale structures [33], pyramid networks [34], attention networks [35], [36], cross-mode information [37], dense networks [38], [39], and adversarial network [40], [41]. Some approaches use detail information from high-spatial-resolution conventional images to improve performance [42]–[45]. Inspired by the specific properties of model, some form a hybrid of model- and deep learning-based approaches [46]–[48], and some use the deep unfolding technique to ease the construction of networks [49]–[51]. These approaches have shown good performance in exploiting the relationship between the observed and target images. However, they are mostly supervised that require plenty of labeled samples to train the networks, which limits their applications in many scenarios.

There are some deep learning-based approaches developed for HSI super-resolution that are performed in an unsupervised manner. For instance, Dian *et al.* [52] introduce a CNN denoiser to regularize the fusion model; Zhang *et al.* [53] integrate the deep image prior into the fusion model, and thereby present a unified unsupervised network for HSI super-resolution; Qu *et al.* [54] exploit an unsupervised approach composed of two autoencoder networks, which are coupled through a shared decoder; Wang *et al.* [55] propose a variational probabilistic autoencoder framework implemented by CNNs for HSI super-resolution; Yao *et al.* [56] propose a two-stream convolutional autoencoder framework inspired by coupled spectral unmixing, and introduce a cross-attention module to improve performance; Uezato *et al.* [57] design a network composed of an encoder-decoder network and a deep decoder network; Zheng *et al.* [58] propose a network consisting of three coupled autoencoder networks, inspired by coupled spectral unmixing, where the three autoencoder networks are coupled through two convolutional layers. Most approaches are built on the autoencoder architecture. Similar to model-based approaches, the construction of networks relies too much on human experience.

Inspired by the specific properties of model, we consider constructing an unsupervised network by referencing some models, and propose a model inspired autoencoder (MIAE) for unsupervised HSI super-resolution. Specifically, we perform nonnegative matrix factorization (NMF) on the target HR-HSI to maintain its intrinsic structure, and thereby propose an implicit autoencoder network for HR-HSI by integrating its NMF model. In the autoencoder network, each hyperspectral pixel is treated as an individual sample, and the two NMF parts of the target HR-HSI, i.e., spectral and spatial matrices, are treated as decoder parameters and hidden outputs respectively. Since the inputs of the autoencoder network are unknown, we take the two observed images as inputs and present a pixel-wise

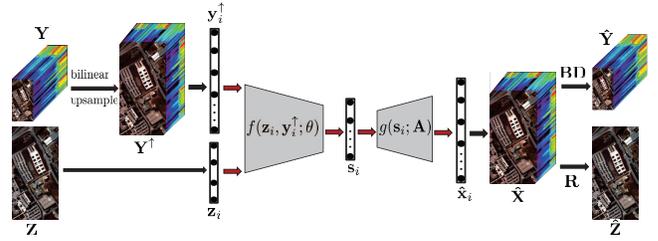


Fig. 1. The overall architecture of MIAE.

fusion model to estimate each hidden output vector directly. The pixel-wise fusion model is solved by the gradient descent algorithm, and the algorithm is reformulated and unfolded to form the encoder network. The loss function is just built on the mechanism of spectral and spatial degradations, and an additional blind estimation network is proposed to estimate the PSF and SRF. Compared with the existing HSI super-resolution approaches, some of the innovative characteristics of MIAE are highlighted as follows.

- 1) MIAE is an unsupervised deep learning-based approach that involves only one implicit autoencoder. The autoencoder network treats each pixel as an individual sample, and thus the proposed network can be treated as a kind of manifold prior-based model and can be trained patch by patch to accelerate the training process.
- 2) MIAE is constructed by referencing models, and thus the construction of the network is relatively concise. The NMF of the target HR-HSI is integrated into the autoencoder, and the encoder network is inspired by the pixel-wise fusion model.
- 3) An additional unsupervised network is proposed to estimate the PSF and SRF from the two observed images directly.

The remainder of this paper is organized as follows. Section II proposes the proposed MAIE and its relationship to the model-based approaches, as well as the blind estimation network. In Section III, the effectiveness of MIAE is demonstrated through experiments on three synthetic datasets and one real dataset. Section IV provides concluding remarks.

## II. PROPOSED APPROACH

Fig. 1 illustrates the overall architecture of MIAE. The details of the proposed network are described as follows.

### A. NMF Inspired Autoencoder

NMF is a useful dimension reduction method [59]. It can capture the intrinsic structure of the data and represent the data in a sparse manner. The properties of NMF indicate that it can facilitate the inference process of super-resolution if we perform NMF on the target HR-HSI. Let us represent the target HR-HSI as a matrix  $\hat{\mathbf{X}} \in \mathbb{R}^{N_B \times N_H N_W}$ , where  $N_B$  denotes the spectral band, and  $N_H$  and  $N_W$  denote the spatial height and width respectively. NMF aims to factor  $\hat{\mathbf{X}}$  into two rank- $J$  ( $J < \min\{N_B, N_H N_W\}$ ) nonnegative matrices, i.e.,

$$\hat{\mathbf{X}} \approx \mathbf{AS}, \quad (1)$$

where the spectral matrix  $\mathbf{A} \in \mathbb{R}^{N_B \times J} \succeq 0$  and the spatial matrix  $\mathbf{S} \in \mathbb{R}^{J \times N_H N_W} \succeq 0$  with  $\succeq$  being a component-wise inequality.

NMF can be integrated into an autoencoder network [60]–[62]. Equation (1) can be rewritten as

$$\hat{\mathbf{x}}_i \approx \mathbf{A} \mathbf{s}_i, \forall i \quad (2)$$

where  $i = 1, 2, \dots, N_H N_W$ , and  $\hat{\mathbf{x}}_i \in \mathbb{R}^{N_B}$  and  $\mathbf{s}_i \in \mathbb{R}^J$  are the column vectors of  $\hat{\mathbf{X}}$  and  $\mathbf{S}$ , respectively. Let  $\hat{\mathbf{x}}_i$  represent the reconstructed vector and  $\mathbf{s}_i$  represent the hidden output vector, we can construct the following autoencoder network

$$\mathbf{x}_i \xrightarrow{f(\cdot)} \mathbf{s}_i \xrightarrow{g(\cdot)} \hat{\mathbf{x}}_i, \forall i \quad (3)$$

where the input data  $\mathbf{x}_i \in \mathbb{R}^{N_B}$  is the column vector of the observed HR-HSI  $\mathbf{X} \in \mathbb{R}^{N_B \times N_H N_W}$ . The network (3) consists of two networks  $f(\cdot)$  and  $g(\cdot)$ .  $f(\cdot)$  is the encoder network with  $\mathbf{s}_i = f(\mathbf{x}_i; \theta)$ , where  $\theta$  denotes all trainable parameters involved in the network.  $g(\cdot)$  is the decoder network with  $\hat{\mathbf{x}}_i = g(\mathbf{s}_i; \mathbf{A}) = \mathcal{C}_0^1(\mathbf{A} \mathbf{s}_i)$ , where  $\mathbf{A}$  is treated as the trainable weight matrix and  $\mathcal{C}_0^1(\cdot)$  is a clamp function that forces all elements of the input vector/matrix into the range  $[0, 1]$ . In hyperspectral unmixing [60]–[62], sum-to-one constraint is added to enforce the hidden vector, i.e.,  $\mathbf{1}_J^T \mathbf{s}_i = 1$  with  $\mathbf{1}_J \in \mathbb{R}^J$  being a vector of all 1s. We do not intend to finish the two tasks of fusion and unmixing at once and only use the nonnegative constraints. Specifically,  $\mathbf{s}_i$  and  $\mathbf{A}$  are enforced using  $\mathcal{C}_0^1(\mathbf{s}_i)$  and  $\mathcal{C}_0^1(\mathbf{A})$ , when designing the network. Then, if the input data  $\mathbf{X}$  is given, one can train the network (3) by feeding its  $N_H N_W$  hyperspectral pixels.

### B. Pixel-Wise Fusion Model Inspired Encoder Network

In HSI super-resolution, the input data  $\mathbf{X} \in \mathbb{R}^{N_B \times N_H N_W}$  is not given. One can not train the network (3) directly. Instead, we have two observed (i.e., degenerated) images of  $\mathbf{X}$ , a LR-HSI  $\mathbf{Y} \in \mathbb{R}^{N_B \times N_h N_w}$  and a HR-MSI  $\mathbf{Z} \in \mathbb{R}^{N_b \times N_H N_W}$ , where  $N_b < N_B$  is the multispectral band, and  $N_h$  and  $N_w$  are the spatial sizes. We assume that  $N_H = r N_h$  and  $N_W = r N_w$  with  $r > 1$  being the resolution ratio. The observations  $\mathbf{Y}$  and  $\mathbf{Z}$  can be modeled as spatially degraded and spectrally degraded versions of  $\mathbf{X}$ . Specifically, these two degeneration processes can be written as:

$$\mathbf{Y} \approx \mathbf{X} \mathbf{B} \mathbf{D} \quad (4)$$

$$\mathbf{Z} \approx \mathbf{R} \mathbf{X} \quad (5)$$

where the PSF  $\mathbf{B} \in \mathbb{R}^{N_H N_W \times N_H N_W}$  is the spatial blur,  $\mathbf{D} \in \mathbb{R}^{N_H N_W \times N_h N_w}$  is the spatial downsampling, and  $\mathbf{R} \in \mathbb{R}^{N_b \times N_B}$  is the SRF of multispectral sensor.

In (3), the network needs to be trained by feeding the input data pixel by pixel. It is the key to the success of autoencoder. For the degeneration processes, (5) can be rewritten as a pixel-wise formulation, i.e.,  $\mathbf{z}_i \approx \mathbf{R} \mathbf{x}_i$  with  $\mathbf{z}_i \in \mathbb{R}^{N_b}$  being the column vector of  $\mathbf{Z}$ , whereas (4) can't because of the coupling matrices  $\mathbf{B}$  and  $\mathbf{D}$ . We consider resizing  $\mathbf{Y}$  to the same size as  $\mathbf{X}$  using bilinear interpolation, in order to approximate  $\mathbf{X}$  pixel by pixel at the spectral level.  $\mathbf{x}_i$  can be obtained by solving

$$\min_{\mathbf{x}_i} \frac{1}{2} \|\mathbf{z}_i - \mathbf{R} \mathbf{x}_i\|_2^2 + \frac{\lambda}{2} \|\mathbf{y}_i^\uparrow - \mathbf{x}_i\|_2^2, \forall i \quad (6)$$

where  $\lambda > 0$  is the regularization parameter, and  $\mathbf{y}_i^\uparrow \in \mathbb{R}^{N_B}$  is the column vector of  $\mathbf{Y}^\uparrow \in \mathbb{R}^{N_B \times N_H N_W}$  with  $\mathbf{Y}^\uparrow$  being a bilinear interpolated version of  $\mathbf{Y}$ .

In order to construct the encoder network, it is unnecessary to solve  $\mathbf{x}_i$  and then design  $f(\cdot)$ , which will lead to error accumulation. We can treat  $\mathbf{x}_i$  as an implicit variable and design  $f(\cdot)$  by using  $\mathbf{z}_i$  and  $\mathbf{y}_i^\uparrow$  directly, i.e.,

$$(\mathbf{z}_i, \mathbf{y}_i^\uparrow) \xrightarrow{f(\cdot)} \mathbf{s}_i \xrightarrow{g(\cdot)} \hat{\mathbf{x}}_i, \forall i \quad (7)$$

Specifically, we want to obtain the hidden layer output  $\mathbf{s}_i$  by solving the following pixel-wise fusion model

$$\min_{\mathbf{s}_i} \frac{1}{2} \|\mathbf{z}_i - \mathbf{R} \mathbf{A} \mathbf{s}_i\|_2^2 + \frac{\lambda}{2} \|\mathbf{y}_i^\uparrow - \mathbf{A} \mathbf{s}_i\|_2^2, \forall i \quad (8)$$

and design  $f(\cdot)$  by unfolding all steps of its algorithm as network layers. Notably, in (8) both  $\mathbf{R}$  and  $\mathbf{A}$  are treated as new trainable parameters to facilitate the design of the encoder network. In the model-based HSI super-resolution, our pervious works [16], [48] have shown the effectiveness of (8).

Although (8) has an analytic solution, it is not suitable as the encoder network and is difficult to implement by a network. (8) can be solved by the gradient descent algorithm as

$$\mathbf{s}_i^k = \mathbf{s}_i^{k-1} - \eta (\bar{\mathbf{A}}^T \bar{\mathbf{A}} \mathbf{s}_i^{k-1} - \bar{\mathbf{A}}^T \mathbf{z}_i + \lambda \bar{\mathbf{A}}^T \mathbf{A} \mathbf{s}_i^{k-1} - \lambda \bar{\mathbf{A}}^T \mathbf{y}_i^\uparrow), \quad (9)$$

where  $\bar{\mathbf{A}} = \mathbf{R} \mathbf{A} \in \mathbb{R}^{N_b \times J}$ ,  $\eta > 0$  is the step, and  $k = 1, 2, \dots, K$  represents the  $k$ -th iteration. To better design the network, input data and intermediate variables are distinguished by rewriting (9) as

$$\mathbf{s}_i^k = (\mathbf{I} - \eta \bar{\mathbf{A}}^T \bar{\mathbf{A}} - \eta \lambda \bar{\mathbf{A}}^T \mathbf{A}) \mathbf{s}_i^{k-1} + \eta \bar{\mathbf{A}}^T \mathbf{z}_i + \eta \lambda \bar{\mathbf{A}}^T \mathbf{y}_i^\uparrow, \quad (10)$$

where  $\mathbf{I}$  represents the identity matrix. According to the  $K$  iterations of (10), the proposed encoder network is mainly a structure of  $K$  stages. Fig. 2 illustrates the details of  $f(\cdot)$  when  $K = 3$ . In (10), all variables  $\mathbf{s}_i^{k-1}$ ,  $\mathbf{z}_i$  and  $\mathbf{y}_i^\uparrow$  are left multiplied by a matrix. This process is implemented using a fully connected layer followed by a Leaky ReLU, and (10) can be rewritten as

$$\mathbf{s}_i^k = f_s^k(\mathbf{s}_i^{k-1}; \theta_s^k) + \eta f_z(\mathbf{z}_i; \theta_z) + \eta \lambda f_y(\mathbf{y}_i^\uparrow; \theta_y), \quad (11)$$

where  $\{f_s^k\}_{k=2}^K$ ,  $f_z$  and  $f_y$  represent the modules designed for the multiplication of matrix and vector, and  $\{\theta_s^k\}_{k=2}^K$ ,  $\theta_z$  and  $\theta_y$  represent the trainable parameters involved in the corresponding networks. The red dotted boxes in Fig. 2 show the layers of  $\{f_s^k\}_{k=2}^K$ ,  $f_z$  and  $f_y$ , where two fully connected layers are used in  $f_y$  for the purpose of feature extraction. In (11), the three modules are combined linearly. To break the fixed format of optimization model and provide more flexibility, the linear combination is implemented by concatenating these modules and performing a fully connection and a leaky ReLU. (11) can be rewritten as

$$\mathbf{s}_i^k = f^k(f_s^k(\mathbf{s}_i^{k-1}; \theta_s^k), f_z(\mathbf{z}_i; \theta_z), f_y(\mathbf{y}_i^\uparrow; \theta_y); \theta^k), \quad (12)$$

where  $\{f^k\}_{k=1}^K$  and  $\{\theta^k\}_{k=1}^K$  represent the modules and trainable parameters for the linear combination. Finally, by performing  $f^k$  from 1 to  $K$ , we can obtain the hidden output by  $\mathbf{s}_i = \mathcal{C}_0^1(\mathbf{s}_i^K)$  and have the trainable parameters  $\theta = \{\theta_z, \theta_y, \{f_s^k\}_{k=2}^K, \{\theta^k\}_{k=1}^K\}$ .

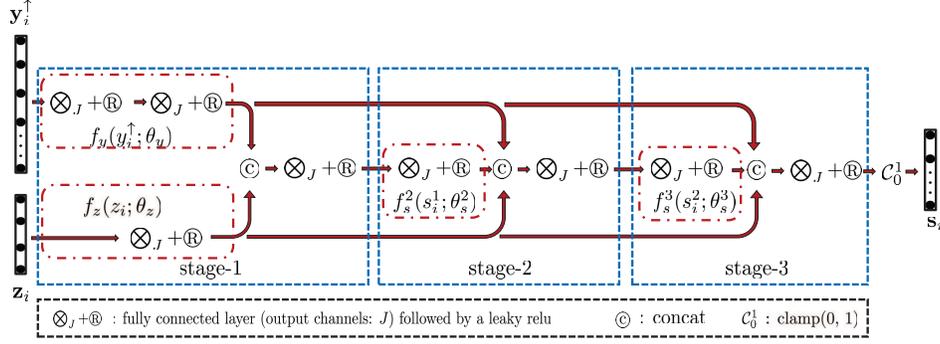


Fig. 2. Details of the proposed encoder network when  $K = 3$ .

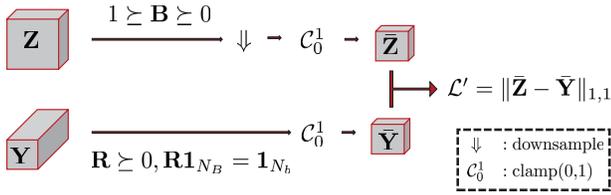


Fig. 3. Details of the blind estimation network.

### C. Loss Function and Training Strategy

To train the autoencoder network (7), we can't build loss function by using the target  $\hat{\mathbf{x}}_i$  directly, since the input data  $\mathbf{x}_i$  is just an implicit variable. Instead, we have two observation images  $\mathbf{Y}$  and  $\mathbf{Z}$  to work with. Similar to (4) and (5), the outputs LR-HSI  $\hat{\mathbf{Y}}$  and HR-MSI  $\hat{\mathbf{Z}}$  can be modeled as  $\hat{\mathbf{Y}} \approx \hat{\mathbf{X}}\mathbf{B}\mathbf{D}$  and  $\hat{\mathbf{Z}} \approx \mathbf{R}\hat{\mathbf{X}}$ . To measure the difference between the outputs and observations, the  $l_1$ -norm is used because it is more robust to outliers than the  $l_2$ -norm. Then, the overall loss function can be written as

$$\mathcal{L} = \|\mathbf{Z} - \hat{\mathbf{Z}}\|_{1,1} + \|\mathbf{Y} - \hat{\mathbf{Y}}\|_{1,1}, \quad (13)$$

where  $\|\cdot\|_{1,1}$  represents the absolute sum of all the matrix elements.

When training the network (7) using the loss function (13), one has to combine all pixels  $\hat{\mathbf{x}}_i$  into an image  $\hat{\mathbf{X}}$ , since  $\hat{\mathbf{Y}}$  and  $\hat{\mathbf{Z}}$  are coupled together by  $\mathbf{B}\mathbf{D}$ . In other words, one has to train (7) by feeding  $\mathbf{Z}$  and  $\mathbf{Y}^\uparrow$  entirely. In spite of this, we can still train (7) by using small patches to accelerate the training process. Specifically, we divide  $\mathbf{Z}$  and  $\mathbf{Y}^\uparrow$  into overlapped patches so that the patches cover all pixels, and then discard the pixels affected by spatial blur  $\mathbf{B}$  at the boundaries of these patches, when computing the loss function (13).

### D. Blind Estimation Network

The PSF  $\mathbf{B}$  and SRF  $\mathbf{R}$  are required to train the proposed autoencoder network. By combing (4) and (5), we have

$$\mathbf{Z}\mathbf{B}\mathbf{D} \approx \mathbf{R}\mathbf{Y}. \quad (14)$$

By imposing some physical constraints, one can obtain  $\mathbf{B}$  and  $\mathbf{R}$  by solving

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{R}} \quad & \|\mathbf{Z}\mathbf{B}\mathbf{D} - \mathbf{R}\mathbf{Y}\|_{1,1} \\ \text{s.t.} \quad & \mathbf{1} \succeq \mathbf{B} \succeq \mathbf{0}, \mathbf{R} \succeq \mathbf{0}, \mathbf{R}\mathbf{1}_{N_b} = \mathbf{1}_{N_b} \end{aligned} \quad (15)$$

Problem (15) can be solved by some optimization algorithms. Instead, we solve (15) by training a network. Specifically, we treat  $(\mathbf{Z}, \mathbf{Y})$  as inputs and  $(\mathbf{B}, \mathbf{R})$  as trainable parameters. Then, the blind estimation network can be constructed by using the following loss function

$$\mathcal{L}' = \|\bar{\mathbf{Z}} - \bar{\mathbf{Y}}\|_{1,1}, \quad (16)$$

where  $\bar{\mathbf{Z}} = \mathcal{C}_0^1(\mathbf{Z}\mathbf{B}\mathbf{D})$  and  $\bar{\mathbf{Y}} = \mathcal{C}_0^1(\mathbf{R}\mathbf{Y})$  represent the output data. The details of the blind estimation network are illustrated in Fig. 3.

### E. Relationship to Model-Based Approaches

The proposed MIAE can be regarded as a kind of specific fusion model. By combing (1), (7) and (13), MIAE can be rewritten as

$$\begin{aligned} \min_{\mathbf{A}, \theta} \quad & \|\mathbf{Z} - \mathbf{R}\mathbf{A}\mathbf{S}\|_{1,1} + \|\mathbf{Y} - \mathbf{A}\mathbf{S}\mathbf{B}\mathbf{D}\|_{1,1} \\ \text{s.t.} \quad & \mathbf{s}_i = f(\mathbf{z}_i, \mathbf{y}_i^\uparrow; \theta), \forall i \\ & \mathbf{1} \succeq \mathbf{A} \succeq \mathbf{0}, \mathbf{1} \succeq \mathbf{S} \succeq \mathbf{0}, \mathbf{1} \succeq \mathbf{A}\mathbf{S} \succeq \mathbf{0} \end{aligned} \quad (17)$$

In (17),  $f(\cdot)$  can be thought of as a nonlinear mapping function, and each  $\mathbf{s}_i$  is only associated with the input  $\mathbf{z}_i$  and  $\mathbf{y}_i^\uparrow$  that correspond to its spatial position. The constraint  $\mathbf{s}_i = f(\mathbf{z}_i, \mathbf{y}_i^\uparrow; \theta)$  acts as a manifold regularization that embeds the combination of  $\mathbf{z}_i$  and  $\mathbf{y}_i^\uparrow$  into a low-dimensional space  $\mathbb{R}^J$ .  $\mathbf{s}_i = f(\mathbf{z}_i, \mathbf{y}_i^\uparrow; \theta)$  also acts as a self-supervised deep prior regularization that only uses itself as training data.

## III. EXPERIMENTAL RESULTS AND ANALYSIS

In this section, experiments on both synthetic and real datasets are conducted to demonstrate the performance of the proposed MIAE. Before the following experiments, all datasets are scaled to the range  $[0, 1]$ . The quality of the fused images in the synthetic datasets are assessed with root mean squared error (RMSE), peak signal-noise-ratio (PSNR), spectral angle mapper (SAM), relative dimensionless global error in synthesis (ERGAS), and universal image quality index (UIQI) [2], [3], [5].

### A. Synthetic Datasets and Implementation Details

Three real-life HSI datasets, University of Paiva (PaviaU), Kennedy Space Center (KSC), and Washington DC Mall (DC) are manipulated to use as synthetic reference images for the simulation experiments.

- 1) The PaviaU dataset is acquired by the Reflective Optics System Imaging Spectrometer (ROSIS), with a spectral range of 0.43 to 0.86  $\mu\text{m}$ . The ROSIS sensor is characterized by 115 spectral bands and 103 remained after removal of noisy bands. This image, with size of  $610 \times 340$  pixels, has spatial resolution of 1.3  $m$  per pixel. We select the up-left  $512 \times 256$ -pixel part as the reference image.
- 2) The KSC dataset is acquired by the Airborne Visible/Infrared Imaging Spectrometer (AVIRIS), with a spectral range of 0.4 to 2.5  $\mu\text{m}$ . The AVIRIS sensor is characterized by 224 spectral bands and the number of spectral bands is reduced to 176 by removing water absorption bands. The size of this image is  $512 \times 614$  with a spatial resolution of 18  $m$ . We select the up-left  $512 \times 512$ -pixel part as the reference image.
- 3) The DC dataset is acquired by the Hyperspectral digital imagery collection experiment (HYDICE) image, with a spectral range of 0.4 to 2.4  $\mu\text{m}$ . The HYDICE sensor is characterized by 210 spectral bands, and bands in the region where the atmosphere is opaque have been removed, leaving 191 bands. This image, with size of  $1208 \times 307$  pixels, has a spatial resolution of about 2.8  $m$ . We select a  $512 \times 256$ -pixel part as the reference image.

For each reference image, we generate the two observation images, LR-HSI and HR-MSI, according to the Wald’s protocol [63]. To generate the LR-HSI, we spatially blur the reference image and then downsample it by a factor of 8 ( $r = 8$ ) in each direction. A Gaussian blur of  $15 \times 15$  pixels, with a mean of 0 and a standard deviation of 3.40, is applied to each band of the reference image. To generate the HR-MSI,  $\mathbf{R}$  is derived from the spectral response of the IKONOS satellite. We generate a 4-band image by averaging the bands of the reference image according to the spectral response profiles of the RGB and NIR bands. To account for ubiquitous noise or error, moderate Gaussian noise is added to the LR-HSI (SNR = 30 dB) and the HR-MSI (SNR = 40 dB).

We implement and train the proposed network and blind estimation network using the PyTorch framework. As discussed in Section II-C, we divide the observation images into patches to accelerate the training process. Take the HR-MSI as a reference,  $40 \times 40$ -pixel overlapping patches with a stride of 24 are extracted for training. The batch sizes are 25 for the PaviaU and DC datasets, and 50 for the KSC dataset. An Adam optimizer is used to train the network for 10000 iterations. The learning rate is initialized as  $5 \times 10^{-3}$  and gradually decayed by multiplying  $1 - \frac{1}{9000} \max(0, \textit{iteration} - 1000)$ , where ‘*iteration*’ represents the current number of iterations. As for the blind estimation network, it is trained by feeding the observed images entirely, the total number of iterations is 5000, and the learning rate is set as  $5 \times 10^{-5}$ .

TABLE I  
QUALITY MEASURES OF NONBLIND AND BLIND MIAE

	PaviaU			KSC		DC	
	best	nonblind	blind	nonblind	blind	nonblind	blind
RMSE	0	0.0169	0.0172	0.0426	0.0427	0.0127	0.0145
PSNR	$+\infty$	37.57	37.33	34.21	34.05	37.48	35.90
SAM	0	2.41	2.43	6.98	7.00	1.56	1.84
ERGAS	0	0.647	0.656	3.122	3.129	14.184	14.224
UIQI	1	0.988	0.988	0.882	0.882	0.983	0.975

### B. Influence of Parameters

Two parameters, rank  $J$  and stage  $K$ , need to be given when constructing the proposed network. In this set of experiments, we investigate them and show how they impact quality measures of MIAE. Fig. 4 illustrates the PSNR results of MIAE as a function of  $J$  when  $K = 1, 2, \dots, 5$ . It can be seen that, for all datasets, the PSNR performance improves as  $J$  increases, but a large  $J$  will cause overfitting or performance degradation. Compared with small values of  $K$ , a moderate  $K$  is better and a too large  $K$  is prone to overfitting. Thus,  $K$  is eventually set as 3 for all datasets, and  $J$  is eventually set as 80 for the PaviaU and KSC datasets and 30 for the DC dataset.

### C. Experiment Results on Synthetic Datasets

1) *Blind and Nonblind*: Section II-D presents a blind estimation network for estimating the PSF and SRF. This experiment is used to evaluate the estimated  $\mathbf{B}$  and  $\mathbf{R}$ . Table I shows the quality measures of the proposed MIAE using the exact and estimated  $\mathbf{B}$  and  $\mathbf{R}$ , that is, nonblind and blind cases. It can be seen that, for the PaviaU and KSC datasets, the performance degradation caused by blind estimation is very small when compared with the nonblind estimation; and for the DC dataset, the performance degradation is also not significant.

2) *Influence of LR-HSI Interpolation*: For the proposed MIAE, bilinear interpolation is used to upsample the LR-HSI to the same size of the target HR-HSI. This experiment shows how the interpolation method affects the performance of MIAE. Four interpolation methods are considered, i.e., bilinear interpolation, nearest interpolation, bicubic interpolation and cubic spline interpolation. The quality measures to assess the different interpolation methods are given in Table II. In most cases, there is no obvious difference between these interpolation methods. The nearest interpolation performs slightly worse on the KSC dataset, and the bicubic interpolation on the DC dataset.

3) *Comparison With the State of the Art*: Nine unsupervised methods, which can be divided into model- and deep learning-based approaches, are compared to evaluate the performance of MIAE. The model-based approaches consist of six methods. The first method is the baseline one (denoted by SLYV) that solves a Sylvester equation [64], and the next five methods are coupled NMF (CNMF) [17], coupled spectral unmixing (CSU) [18], NSSR [21], HySure [13], and NPTSR [11]. The deep learning-based approaches are CNNFUS [52] and three autoencoder-based methods, i.e., uSDN [54], HyCoNet

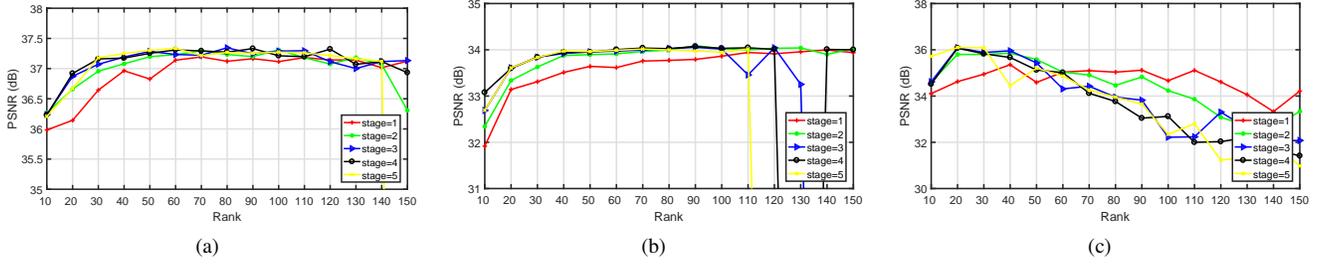


Fig. 4. PSNR as a function of rank  $J$  when using different stages  $K$ . (a) PaviaU dataset. (b) KSC dataset. (c) DC dataset.

TABLE II  
QUALITY MEASURES OF MIAE USING DIFFERENT INTERPOLATION METHODS

	PaviaU					KSC				DC			
	best	bilinear	nearest	bicubic	spline	bilinear	nearest	bicubic	spline	bilinear	nearest	bicubic	spline
RMSE	0	0.0172	0.0172	0.0171	0.0172	0.0427	0.0431	0.0428	0.0426	0.0145	0.0151	0.0148	0.0142
PSNR	$+\infty$	37.33	37.28	37.30	37.27	34.05	33.88	34.02	34.04	35.90	35.61	34.77	35.78
SAM	0	2.43	2.43	2.41	2.43	7.00	7.11	7.01	7.01	1.84	1.93	1.83	1.82
ERGAS	0	0.656	0.656	0.652	0.653	3.129	3.164	3.135	3.130	14.224	14.204	14.349	14.188
UIQI	1	0.988	0.988	0.988	0.988	0.882	0.879	0.882	0.881	0.975	0.974	0.975	0.975

TABLE III  
QUALITY MEASURES FOR THE PAVIAU DATASET USING DIFFERENT METHODS (THE BEST VALUES ARE HIGHLIGHTED)

Method	SLYV	CNMF	CSU	NSSR	HySure	NPTSR	CNNFUS	uSDN	HyCoNet	MIAE
RMSE	0.1072	0.0196	0.0231	0.0236	0.0194	0.0186	0.0237	0.0258	0.0188	<b>0.0172</b>
PSNR	23.79	35.75	33.87	33.93	36.09	36.71	35.24	32.84	36.67	<b>37.33</b>
SAM	12.62	2.62	2.89	3.21	2.70	2.64	3.16	3.49	2.66	<b>2.43</b>
ERGAS	3.646	0.741	0.849	0.871	0.728	0.699	0.825	0.905	0.720	<b>0.656</b>
UIQI	0.853	0.987	0.983	0.983	0.986	0.987	0.984	0.982	0.987	<b>0.988</b>

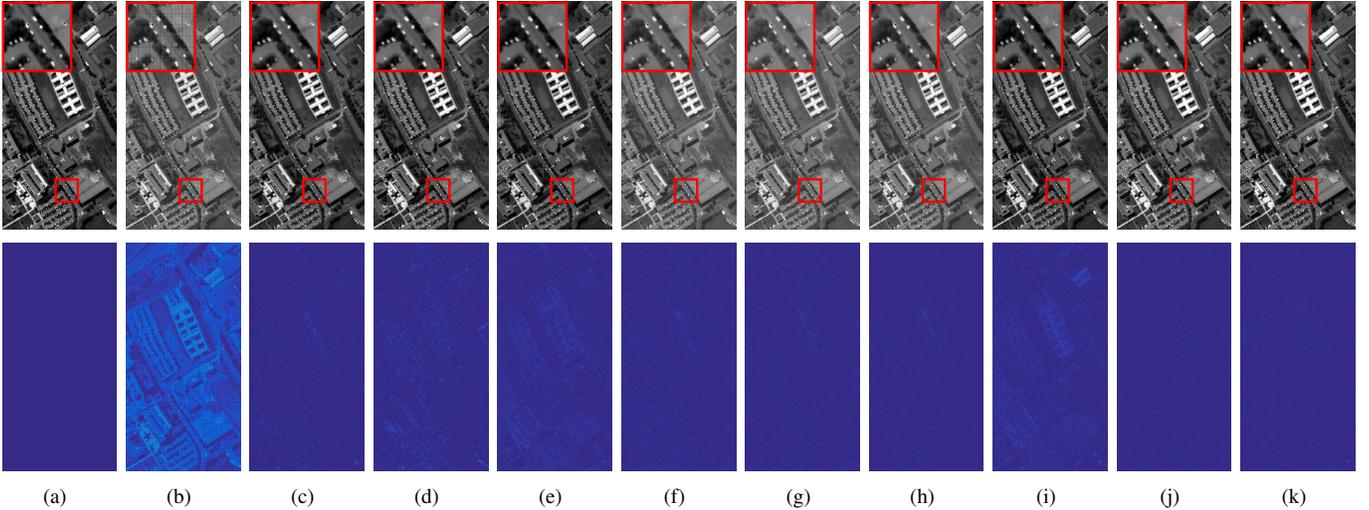


Fig. 5. Images (with a meaningful region marked and zoomed in 3 times for easy observation) and error maps at band 30 of HSI super-resolution results when applied to the PaviaU dataset. (a) Reference image. (b) SLYV. (c) CNMF. (d) CSU. (e) NSSR. (f) HySure. (g) NPTSR. (h) CNNFUS. (i) uSDN. (j) HyCoNet. (k) MIAE.

TABLE IV  
QUALITY MEASURES FOR THE KSC DATASET USING DIFFERENT METHODS (THE BEST VALUES ARE HIGHLIGHTED)

Method	SLYV	CNMF	CSU	NSSR	HySure	NPTSR	CNNFUS	uSDN	HyCoNet	MIAE
RMSE	0.1574	0.0454	0.0465	0.0513	0.0453	0.0450	0.0534	0.0504	0.0441	<b>0.0427</b>
PSNR	19.33	32.70	31.17	30.77	32.95	33.30	30.95	30.06	33.49	<b>34.05</b>
SAM	23.23	7.78	8.02	8.65	7.64	7.29	9.01	9.14	7.22	<b>7.00</b>
ERGAS	8.753	3.497	3.405	3.738	3.336	3.328	3.954	3.717	3.262	<b>3.129</b>
UIQI	0.506	0.870	0.855	0.836	0.881	<b>0.887</b>	0.843	0.857	0.878	0.882

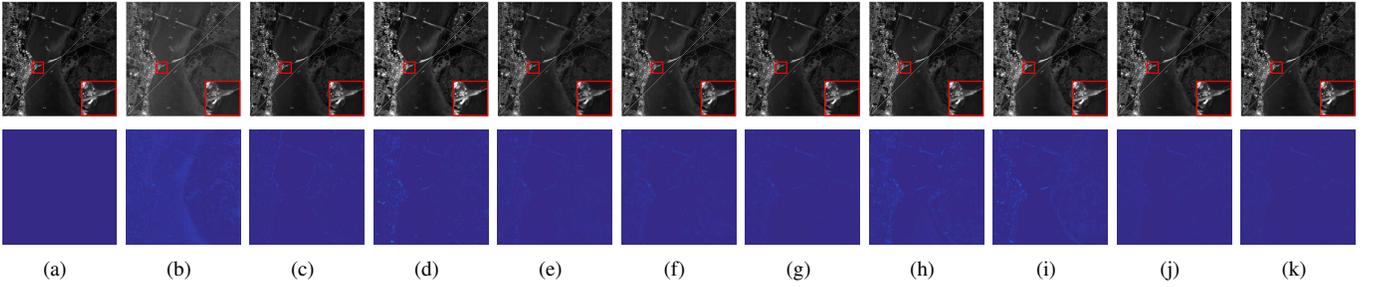


Fig. 6. Images (with a meaningful region marked and zoomed in 3 times for easy observation) and error maps at band 30 of HSI super-resolution results when applied to the KSC dataset. (a) Reference image. (b) SLYV. (c) CNMF. (d) CSU. (e) NSSR. (f) HySure. (g) NPTSR. (h) CNNFUS. (i) uSDN. (j) HyCoNet. (k) MIAE.

TABLE V  
QUALITY MEASURES FOR THE DC DATASET USING DIFFERENT METHODS (THE BEST VALUES ARE HIGHLIGHTED)

Method	SLYV	CNMF	CSU	NSSR	HySure	NPTSR	CNNFUS	uSDN	HyCoNet	MIAE
RMSE	0.1685	0.0283	0.0220	0.0366	0.0275	0.0261	0.0332	0.0278	0.0197	<b>0.0145</b>
PSNR	18.53	32.61	32.48	29.90	32.15	32.77	32.39	29.24	31.63	<b>35.90</b>
SAM	26.56	3.82	2.82	5.33	3.48	3.54	3.70	3.53	<b>1.83</b>	1.84
ERGAS	37.085	15.182	14.406	14.447	13.797	13.886	14.300	17.345	<b>10.943</b>	14.224
UIQI	0.469	0.930	0.929	0.905	0.930	0.944	0.954	0.856	0.914	<b>0.975</b>

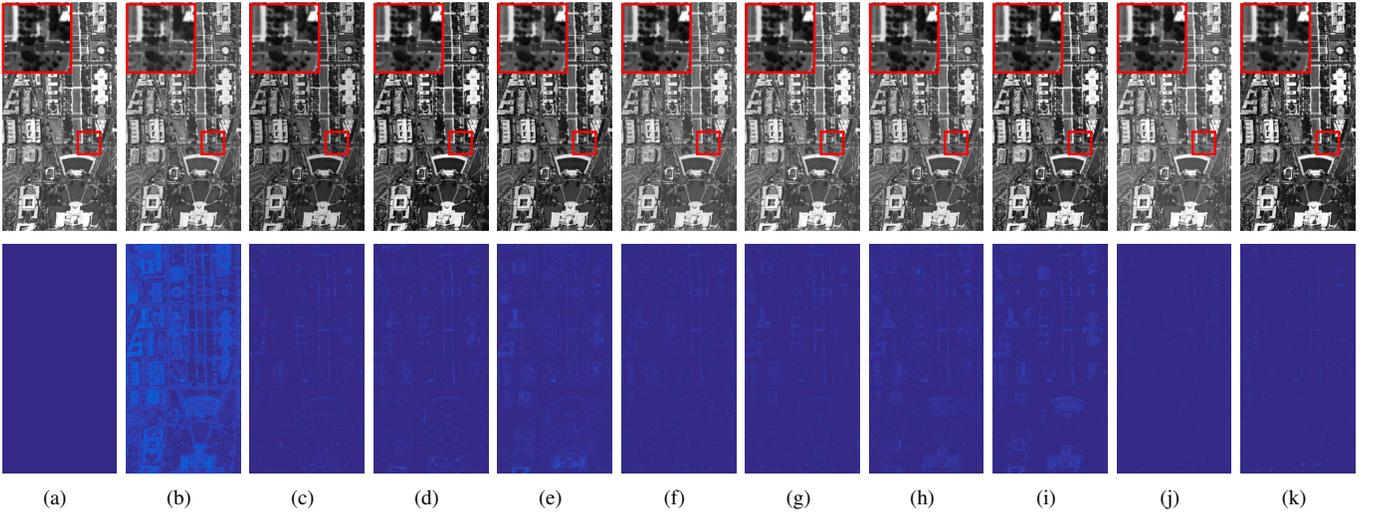


Fig. 7. Images (with a meaningful region marked and zoomed in 3 times for easy observation) and error maps at band 30 of HSI super-resolution results when applied to the DC dataset. (a) Reference image. (b) SLYV. (c) CNMF. (d) CSU. (e) NSSR. (f) HySure. (g) NPTSR. (h) CNNFUS. (i) uSDN. (j) HyCoNet. (k) MIAE.

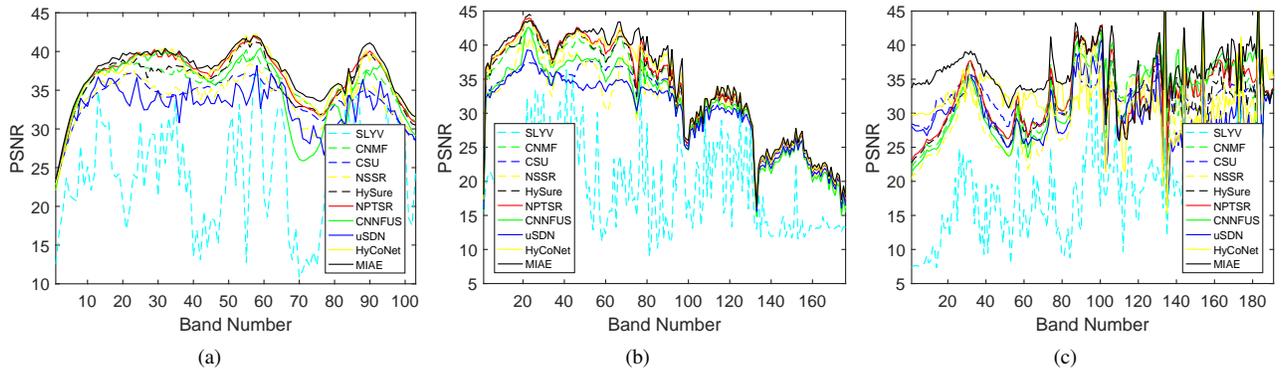


Fig. 8. PSNR as a function of spectral band. (a) PaviaU dataset. (b) KSC dataset. (c) DC dataset.

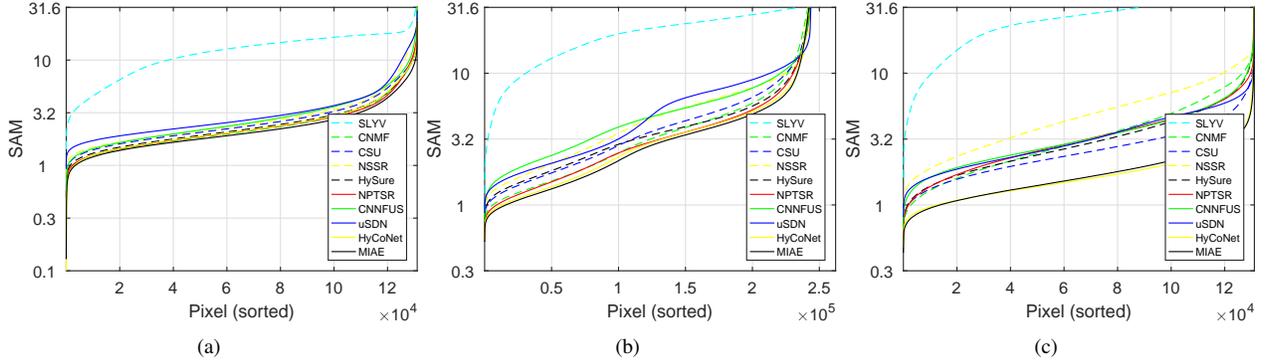


Fig. 9. SAM (plotted in a  $\log_{10}(\cdot)$  scale) as a function of sorted pixel. (a) PaviaU dataset. (b) KSC dataset. (c) DC dataset.

[58] and also MIAE. The free parameters of the compared methods are tuned to be optimal with the test datasets, and the default training strategies are used for the deep learning-based approaches. All of the compared methods are blind, where the estimated  $\mathbf{B}$  and  $\mathbf{R}$  are used. For those methods that do not involve blind estimation procedures,  $\mathbf{B}$  and  $\mathbf{R}$  are estimated by the proposed blind estimation network.

The five quantitative results of the compared methods for the PaviaU dataset are shown in Table III with the best values marked in bold. It can be seen that, all methods outperform the baseline method SLYV, and the proposed MIAE method gives the best quantitative results followed by NPTSR and HyCoNet. Both the model- and deep learning-based approaches can yield good results. Fig. 5 illustrates the reference image and the fusion results of the compared methods in form of the 30th band gray and error images. Visually, it can be observed that the baseline method SLYV has severe spatial distortion and all other methods outperform it. MIAE and HyCoNet perform better than the other methods in terms of both zoomed region and error map. Fig. 8 (a) shows the PSNR as a function of spectral band for the compared methods. It can be seen that the proposed MIAE method performs best in almost all bands followed by HyCoNet and NPTSR. Fig. 9 (a) shows the SAM between the reference image and the fusion results for each pixel using the compared methods, with the pixels sorted by ascending error. As illustrated in this figure, MIAE consistently outperforms the others at the pixel level.

Table IV reports the five quality measures of the compared methods for the KSC dataset. From this table, we can see that the baseline method SLYV performs the worst, NPTSR gives the best UIQI result, and the proposed MIAE method performs best for the remaining four quality measures. NPTSR and HyCoNet are only inferior to MIAE. In Fig. 6, we show the reference image and the fusion results of the compared methods in form of the 30th band gray and error images. Visually, it can be observed that the reconstructed results of HyCoNet and MIAE are better than the others, and the baseline method SLYV gives the worst images. Fig. 8 (b) gives PSNR as a function of the spectral band for the compared methods. MIAE, HyCoNet and NPTSR achieve high results in most bands. Fig. 9 (b) gives the SAMs for each pixel between the reference image and the fusion results, with the pixels sorted in order of ascending error. It can be observed that

MIAE is the best followed by HyCoNet and NPTSR.

Table V summarizes the five quality measures of the compared methods for the DC dataset. From this table, we can see that MIAE gives three best quantitative results and one second best, and HyCoNet gives two best. The 30th band gray and error images of the reference image and the fusion results of the compared methods are given in Fig. 7. Through visual inspection, we can see that MIAE and HyCoNet exhibit good reconstructed results. PSNR and SAM, as functions of the spectral band and by pixel sorted on error, are shown in Figs. 8 (c) and 9 (c), respectively. It can be seen that MIAE outperforms the others in terms of band-level PSNR, and MIAE and HyCoNet achieve higher results than the others in terms of pixel-level SAM.

4) *Computational Efficiency*: All experiments are carried out using a desktop computer with an Intel Core i9-7900X CPU, a GeForce GTX 2080Ti GPU, and 64-GB memory. The first seven methods SLYV, CNMF, CSU, NSSR, HySure, NPTSR and CNMFUS are performed using MATLAB, and the remaining three autoencoder-based methods uSDN, HyCoNet and MIAE are implemented by the PyTorch framework. Table VI summarizes the running times of the first seven methods and the training times of the autoencoder-based methods, and the number of trainable parameters for each autoencoder network is reported in Table VII. It can be seen that, MIAE takes less time to train the network than the other two autoencoder-based methods, and its trainable parameters are much less than HyCoNet. Ignoring the platform, MIAE is comparable to the model-based approaches.

#### D. Experiment Results on Real Data

The University of Houston (UH) dataset released by the 2018 IEEE GRSS Data Fusion Contest [65] is used to evaluate MIAE in practical application. The original data is acquired by the National Center for Airborne Laser Mapping (NCALM), covering the UH campus and its surrounding urban areas. This experiment selects a LR-HSI and a high-resolution RGB (HR-RGB) image from this multi-modal optical remote sensing datasets. The LR-HSI collected by ITRES CASI-1500 sensor contains  $4172 \times 1202$  pixels with a spatial resolution of 1 m and 48 spectral bands with a spectral range of 0.38 to 1.05  $\mu m$ . The HR-RGB image collected by DiMAC ULTRALIGHT+ sensor contains  $83440 \times 24040$  pixels. Take the LR-HSI as a

TABLE VI  
RUNNING/TRAINING TIMES (IN SECONDS) OF THE COMPARED METHODS

Method	SLYV	CNMF	CSU	NSSR	HySure	NPTSR	CNNFUS	uSDN	HyCoNet	MIAE
PaviaU	1.6	16.4	156.9	181.9	167.6	1339.0	9.4	504.5	562.6	186.5
KSC	5.6	36.5	302.7	499.0	336.1	4599.2	14.1	827.9	970.7	396.4
DC	3.0	16.5	172.2	290.6	171.1	2424.3	8.8	477.8	577.9	239.1

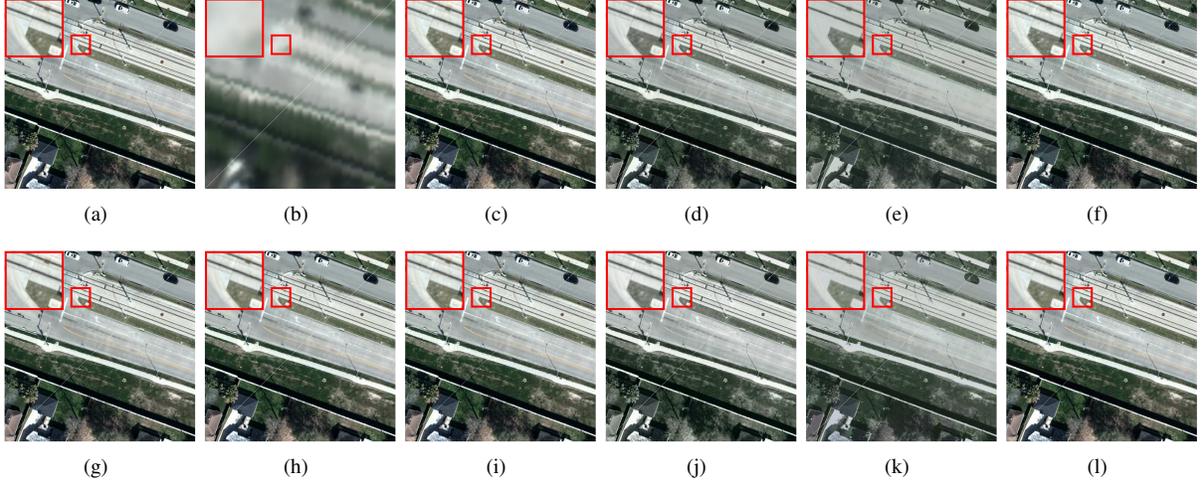


Fig. 10. RGB images (with a meaningful region marked and zoomed in 3 times for easy observation) of HSI super-resolution results when applied to real dataset. (a) HR-RGB image. (b) LR-HSI. (c) SLYV. (d) CNMF. (e) CSU. (f) NSSR. (g) HySure. (h) NPTSR. (i) CNNFUS. (j) uSDN. (k) HyCoNet. (l) MIAE.

TABLE VII  
NUMBER OF TRAINABLE PARAMETERS

	uSDN	HyCoNet	MIAE
PaviaU	37.9K	377.7K	87.8K
KSC	48.9K	389.5K	99.5K
DC	51.1K	391.9K	21.7K

reference, we select an area of  $64 \times 64 \times 48$  as our observation data, and downsample the corresponding area of the HR-RGB to be a  $512 \times 512 \times 3$ -size image. That is, the resolution ratio is  $r = 8$ . RGB images of the real dataset and the fusion results of the compared methods mentioned in Section III-C3 are given in Fig. 10. Visually, it can be seen that MIAE, NPTSR and HySure give the good color and brightness results, and the result of the proposed MIAE is much closer to the HR-RGB image.

#### IV. CONCLUSION

This paper has proposed an unsupervised MIAE network for HSI super-resolution. The proposed MIAE involves an implicit autoencoder network and the structures are concise. Firstly, inspired by that performing NMF on the target HR-HSI can facilitate the inference process of super-resolution, the implicit autoencoder network is built on the target HR-HSI by integrating its NMF model, where the two NMF parts, spectral and spatial matrices, are treated as decoder parameters and hidden outputs respectively. The autoencoder network treats each hyperspectral pixel of the target HR-HSI as an individual sample, that is, the network is trained pixel by pixel. Secondly, the ‘implicit’ indicates the input pixel of the

autoencoder network is unknown, and thus a pixel-wise fusion model taken the two observed images as inputs is presented to estimate the hidden layer vector directly. The pixel-wise fusion model is simple and effective. Specifically, the LR-HSI is resized to the same size of the target HR-HSI using bilinear interpolation, in order to feed the network pixel by pixel. To break the fixed format of model and provide more flexibility, the gradient descent algorithm is used to solve the pixel-wise fusion model, and the algorithm is reformulated and unfolded to form the encoder network. Finally, the loss function is built on the relationship between the target HR-HSI and the two observed images. With the specific pixel-wise architecture, MIAE can be treated as a kind of manifold prior-based model and can be trained patch by patch to accelerate the training process. Moreover, a blind estimation network is proposed to estimate the PSF and SRF in an unsupervised manner. MIAE has been experimentally tested using three synthetic datasets and one real dataset, and the experimental results demonstrate its effectiveness. Although the results obtained by MIAE are very encouraging, further improvements such as the application of convolutional autoencoder should be pursued in future.

#### ACKNOWLEDGMENT

The authors would like to thank the authors of [11], [13], [17], [18], [21], [52], [54], [58], [64] for providing their codes. They would like to thank NCALM and the Hyperspectral Image Analysis Laboratory at UH for providing the UH datasets, and the Image Analysis and Data Fusion Technical Committee of the IEEE GRSS for supporting the annual Data Fusion Contest.

## REFERENCES

- [1] L. Alparone, L. Wald, J. Chanussot, C. Thomas, P. Gamba, and L. M. Bruce, "Comparison of pansharpening algorithms: Outcome of the 2006 grss data-fusion contest," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 45, no. 10, pp. 3012–3021, 2007.
- [2] L. Loncan, L. B. De Almeida, J. M. Bioucas-Dias, X. Briottet, J. Chanussot, N. Dobigeon, S. Fabre, W. Liao, G. A. Licciardi, M. Simoes *et al.*, "Hyperspectral pansharpening: A review," *IEEE Geoscience and remote sensing magazine*, vol. 3, no. 3, pp. 27–46, 2015.
- [3] N. Yokoya, C. Grohnfeldt, and J. Chanussot, "Hyperspectral and multispectral data fusion: A comparative review of the recent literature," *IEEE Geoscience and Remote Sensing Magazine*, vol. 5, no. 2, pp. 29–56, 2017.
- [4] X. Meng, H. Shen, H. Li, L. Zhang, and R. Fu, "Review of the pansharpening methods for remote sensing images based on the idea of meta-analysis: Practical discussion and challenges," *Information Fusion*, vol. 46, pp. 102–113, 2019.
- [5] R. Dian, S. Li, B. Sun, and A. Guo, "Recent advances and new guidelines on hyperspectral and multispectral image fusion," *Information Fusion*, vol. 69, pp. 40–51, 2021.
- [6] G. Vivone, M. D. Mura, A. Garzelli, R. Restaino, G. Scarpa, M. Ulfarsson, L. Alparone, and J. Chanussot, "A new benchmark based on recent advances in multispectral pansharpening: Revisiting pansharpening with classical and emerging pansharpening methods," *IEEE Geoscience and Remote Sensing Magazine*, vol. 9, pp. 53–81, 2021.
- [7] T.-M. Tu, S.-C. Su, H.-C. Shyu, and P. S. Huang, "A new look at ihs-like image fusion methods," *Information fusion*, vol. 2, no. 3, pp. 177–186, 2001.
- [8] F. Nencini, A. Garzelli, S. Baronti, and L. Alparone, "Remote sensing image fusion using the curvelet transform," *Information fusion*, vol. 8, no. 2, pp. 143–156, 2007.
- [9] K. Zhang, M. Wang, and S. Yang, "Multispectral and hyperspectral image fusion based on group spectral embedding and low-rank factorization," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 3, pp. 1363–1371, 2017.
- [10] L. Zhang, W. Wei, C. Bai, Y. Gao, and Y. Zhang, "Exploiting clustering manifold structure for hyperspectral imagery super-resolution," *IEEE Transactions on Image Processing*, vol. 27, no. 12, pp. 5969–5982, 2018.
- [11] Y. Xu, Z. Wu, J. Chanussot, and Z. Wei, "Nonlocal patch tensor sparse representation for hyperspectral image super-resolution," *IEEE Transactions on Image Processing*, vol. 28, no. 6, pp. 3034–3047, 2019.
- [12] J. Xue, Y. Zhao, Y. Bu, W. Liao, J. Chan, and W. Philips, "Spatial-spectral structured sparse low-rank representation for hyperspectral image super-resolution," *IEEE Transactions on Image Processing*, vol. 30, pp. 3084–3097, 2021.
- [13] M. Simões, J. Bioucas-Dias, L. B. Almeida, and J. Chanussot, "A convex formulation for hyperspectral image superresolution via subspace-based regularization," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 6, pp. 3373–3388, 2015.
- [14] Q. Wei, J. Bioucas-Dias, N. Dobigeon, and J.-Y. Tourneret, "Hyperspectral and multispectral image fusion based on a sparse representation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 7, pp. 3658–3668, 2015.
- [15] R. Dian and S. Li, "Hyperspectral image super-resolution via subspace-based low tensor multi-rank regularization," *IEEE Transactions on Image Processing*, vol. 28, no. 10, pp. 5135–5146, 2019.
- [16] J. Liu, Z. Wu, L. Xiao, J. Sun, and H. Yan, "A truncated matrix decomposition for hyperspectral image super-resolution," *IEEE Transactions on Image Processing*, vol. 29, pp. 8028–8042, 2020.
- [17] N. Yokoya, T. Yairi, and A. Iwasaki, "Coupled nonnegative matrix factorization unmixing for hyperspectral and multispectral data fusion," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 2, pp. 528–537, 2012.
- [18] C. Lanaras, E. Baltasavias, and K. Schindler, "Hyperspectral super-resolution by coupled spectral unmixing," in *IEEE International Conference on Computer Vision*, 2015, pp. 3586–3594.
- [19] C.-H. Lin, F. Ma, C.-Y. Chi, and C.-H. Hsieh, "A convex optimization-based coupled nonnegative matrix factorization algorithm for hyperspectral and multispectral data fusion," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 3, pp. 1652–1667, 2018.
- [20] R. Wu, W.-K. Ma, X. Fu, and Q. Li, "Hyperspectral super-resolution via global-local low-rank matrix estimation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 10, pp. 7125–7140, 2020.
- [21] W. Dong, F. Fu, G. Shi, X. Cao, J. Wu, G. Li, and X. Li, "Hyperspectral image super-resolution via non-negative structured sparse representation," *IEEE Transactions on Image Processing*, vol. 25, no. 5, pp. 2337–2352, 2016.
- [22] M. A. Veganzones, M. Simoes, G. Licciardi, N. Yokoya, J. M. Bioucas-Dias, and J. Chanussot, "Hyperspectral super-resolution of locally low rank images from complementary multisource data," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 274–288, 2016.
- [23] X.-H. Han, B. Shi, and Y. Zheng, "Self-similarity constrained sparse representation for hyperspectral image super-resolution," *IEEE Transactions on Image Processing*, vol. 27, no. 11, pp. 5625–5637, 2018.
- [24] C. Yi, Y.-Q. Zhao, and J. C.-W. Chan, "Hyperspectral image super-resolution based on spatial and spectral correlation fusion," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 7, pp. 4165–4177, 2018.
- [25] X. Han, J. Yu, J.-H. Xue, and W. Sun, "Hyperspectral and multispectral image fusion using optimized twin dictionaries," *IEEE Transactions on Image Processing*, vol. 29, pp. 4709–4720, 2020.
- [26] S. Li, R. Dian, L. Fang, and J. M. Bioucas-Dias, "Fusing hyperspectral and multispectral images via coupled sparse tensor factorization," *IEEE Transactions on Image Processing*, vol. 27, no. 8, pp. 4118–4130, 2018.
- [27] C. I. Kanatsoulis, X. Fu, N. D. Sidiropoulos, and W.-K. Ma, "Hyperspectral super-resolution: A coupled tensor factorization approach," *IEEE Transactions on Signal Processing*, vol. 66, no. 24, pp. 6503–6517, 2018.
- [28] Y. Xu, Z. Wu, J. Chanussot, and Z. Wei, "Hyperspectral images super-resolution via learning high-order coupled tensor ring representation," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 11, pp. 4747–4760, 2020.
- [29] Y. Chen, J. Zeng, W. He, X.-L. Zhao, and T.-Z. Huang, "Hyperspectral and multispectral image fusion using factor smoothed tensor ring decomposition," *IEEE Transactions on Geoscience and Remote Sensing*, pp. 1–17, 2021.
- [30] L. Bungert, D. A. Coomes, M. J. Ehrhardt, J. Rasch, R. Reisenhofer, and C.-B. Schönlieb, "Blind image fusion for hyperspectral imaging with the directional total variation," *Inverse Problems*, vol. 34, no. 4, p. 044003, 2018.
- [31] F. Palsson, J. R. Sveinsson, and M. O. Ulfarsson, "Multispectral and hyperspectral image fusion using a 3-d-convolutional neural network," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 5, pp. 639–643, 2017.
- [32] G. Scarpa, S. Vitale, and D. Cozzolino, "Target-adaptive cnn-based pansharpening," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 9, pp. 5443–5457, 2018.
- [33] Q. Yuan, Y. Wei, X. Meng, H. Shen, and L. Zhang, "A multiscale and multidepth convolutional neural network for remote sensing imagery pan-sharpening," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 3, pp. 978–989, 2018.
- [34] Y. Zhang, C. Liu, M. Sun, and Y. Ou, "Pan-sharpening using an efficient bidirectional pyramid network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 8, pp. 5549–5563, 2019.
- [35] Y. Zheng, J. Li, Y. Li, J. Guo, X. Wu, and J. Chanussot, "Hyperspectral pansharpening using deep prior and dual attention residual network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 11, pp. 8059–8076, 2020.
- [36] J. Jiang, H. Sun, X. Liu, and J. Ma, "Learning spatial-spectral prior for super-resolution of hyperspectral imagery," *IEEE Transactions on Computational Imaging*, vol. 6, pp. 1082–1096, 2020.
- [37] X. Zhang, W. Huang, Q. Wang, and X. Li, "Ssr-net: Spatial-spectral reconstruction network for hyperspectral and multispectral image fusion," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, pp. 5953–5965, 2021.
- [38] X. Dong, X. Sun, X. Jia, Z. Xi, L. Gao, and B. Zhang, "Remote sensing image super-resolution using novel dense-sampling networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 2, pp. 1618–1633, 2021.
- [39] D. Liu, J. Li, and Q. Yuan, "A spectral grouping and attention-driven residual dense network for hyperspectral image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 9, pp. 7711–7725, 2021.
- [40] J. Li, R. Cui, B. Li, R. Song, Y. Li, Y. Dai, and Q. Du, "Hyperspectral image super-resolution by band attention through adversarial learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 6, pp. 4304–4318, 2020.
- [41] W. Dong, S. Hou, S. Xiao, J. Qu, Q. Du, and Y. Li, "Generative dual-adversarial network with spectral fidelity and spatial enhancement for hyperspectral pansharpening," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2021.

- [42] J. Yang, X. Fu, Y. Hu, Y. Huang, X. Ding, and J. Paisley, "Pannet: A deep network architecture for pan-sharpening," in *IEEE International Conference on Computer Vision*, 2017, pp. 5449–5457.
- [43] W. Xie, J. Lei, Y. Cui, Y. Li, and Q. Du, "Hyperspectral pansharpening with deep priors," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 5, pp. 1529–1543, 2019.
- [44] X. Fu, W. Wang, Y. Huang, X. Ding, and J. Paisley, "Deep multiscale detail networks for multiband spectral image sharpening," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 5, pp. 2090–2104, 2021.
- [45] L.-J. Deng, G. Vivone, C. Jin, and J. Chanussot, "Detail injection-based deep convolutional neural networks for pansharpening," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 8, pp. 6995–7010, 2021.
- [46] R. Dian, S. Li, A. Guo, and L. Fang, "Deep hyperspectral image sharpening," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 11, pp. 5345–5355, 2018.
- [47] H. Shen, M. Jiang, J. Li, Q. Yuan, Y. Wei, and L. Zhang, "Spatial-spectral fusion by combining deep learning and variational model," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 8, pp. 6169–6181, 2019.
- [48] D. Shen, J. Liu, Z. Xiao, J. Yang, and L. Xiao, "A twice optimizing net with matrix decomposition for hyperspectral and multispectral image fusion," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 13, pp. 4095–4110, 2020.
- [49] Q. Xie, M. Zhou, Q. Zhao, D. Meng, W. Zuo, and Z. Xu, "Multispectral and hyperspectral image fusion by ms/hs fusion net," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1585–1594.
- [50] W. Wei, J. Nie, Y. Li, L. Zhang, and Y. Zhang, "Deep recursive network for hyperspectral image super-resolution," *IEEE Transactions on Computational Imaging*, vol. 6, pp. 1233–1244, 2020.
- [51] W. Dong, C. Zhou, F. Wu, J. Wu, G. Shi, and X. Li, "Model-guided deep hyperspectral image super-resolution," *IEEE Trans Image Process*, vol. 30, pp. 5754–5768, 2021.
- [52] R. Dian, S. Li, and X. Kang, "Regularizing hyperspectral and multi-spectral image fusion by cnn denoiser," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, pp. 1124–1135, 2021.
- [53] L. Zhang, J. Nie, W. Wei, Y. Li, and Y. Zhang, "Deep blind hyperspectral image super-resolution," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, pp. 2388–2400, 2021.
- [54] Y. Qu, H. Qi, and C. Kwan, "Unsupervised sparse dirichlet-net for hyperspectral image super-resolution," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2511–2520.
- [55] Z. Wang, B. Chen, R. Lu, H. Zhang, and P. K. Varshney, "Fusionnet: An unsupervised convolutional variational network for hyperspectral and multispectral image fusion," *IEEE Transactions on Image Processing*, vol. 29, pp. 7565–7577, 2020.
- [56] J. Yao, D. Hong, J. Chanussot, D. Meng, X. Zhu, and Z. Xu, "Cross-attention in coupled unmixing nets for unsupervised hyperspectral super-resolution," in *European Conference on Computer Vision*, 2020, pp. 208–224.
- [57] T. Uezato, D. Hong, N. Yokoya, and W. He, "Guided deep decoder: Unsupervised image pair fusion," in *European Conference on Computer Vision*, 2020, pp. 87–102.
- [58] K. Zheng, L. Gao, W. Liao, D. Hong, B. Zhang, X. Cui, and J. Chanussot, "Coupled convolutional neural network with adaptive response function learning for unsupervised hyperspectral super resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 3, pp. 2487–2502, 2021.
- [59] D. D. Lee and H. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, pp. 788–791, 1999.
- [60] Y. Qu and H. Qi, "udas: An untied denoising autoencoder with sparsity for spectral unmixing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, pp. 1698–1712, 2019.
- [61] Y. Qian, F. Xiong, Q. Qian, and J. Zhou, "Spectral mixture model inspired network architectures for hyperspectral unmixing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 10, pp. 7418–7434, 2020.
- [62] B. Palsson, M. O. Ulfarsson, and J. R. Sveinsson, "Convolutional autoencoder for spectral-spatial hyperspectral unmixing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 1, pp. 535–549, 2021.
- [63] T. Ranchin and L. Wald, "Fusion of high spatial and spectral resolution images: the arsis concept and its implementation," *Photogrammetric engineering and remote sensing*, vol. 66, no. 1, pp. 49–61, 2000.
- [64] Q. Wei, N. Dobigeon, and J.-Y. Tourneret, "Fast fusion of multi-band images based on solving a sylvester equation," *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 4109–4121, 2015.
- [65] Y. Xu, B. Du, L. Zhang, D. Cerra, M. Pato, E. Carmona, S. Prasad, N. Yokoya, R. Hansch, and B. Le Saux, "Advanced multi-sensor optical remote sensing for urban land use and land cover classification: Outcome of the 2018 ieeec grss data fusion contest," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 12, no. 6, pp. 1709–1724, 2019.