# Hyperspectral Image Super-resolution via Deep Spatio-spectral Convolutional Neural Networks

Jin-Fan Hu, Ting-Zhu Huang, *Member, IEEE*, Liang-Jian Deng, *Member, IEEE*, Tai-Xiang Jiang, *Member, IEEE*, Gemine Vivone, *Senior Member, IEEE*, and Jocelyn Chanussot, *Fellow, IEEE*

*Abstract*—Hyperspectral images are of crucial importance in order to better understand features of different materials. To reach this goal, they leverage on a high number of spectral bands. However, this interesting characteristic is often paid by a reduced spatial resolution compared with traditional multispectral image systems. In order to alleviate this issue, in this work, we propose a simple and efficient architecture for deep convolutional neural networks to fuse a low-resolution hyperspectral image (LR-HSI) and a high-resolution multispectral image (HR-MSI), yielding a high-resolution hyperspectral image (HR-HSI). The network is designed to preserve both spatial and spectral information thanks to an architecture from two folds: one is to utilize the HR-HSI at a different scale to get an output with a satisfied spectral preservation; another one is to apply concepts of multi-resolution analysis to extract high-frequency information, aiming to output high quality spatial details. Finally, a plain mean squared error loss function is used to measure the performance during the training. Extensive experiments demonstrate that the proposed network architecture achieves best performance (both qualitatively and quantitatively) compared with recent state-of-the-art hyperspectral image super-resolution approaches. Moreover, other significant advantages can be pointed out by the use of the proposed approach, such as, a better network generalization ability, a limited computational burden, and a robustness with respect to the number of training samples.

*Index Terms*—Hyperspectral Image Super-resolution, Deep Convolutional Neural Network, Multiscale Structure, Image Fusion.

## I. INTRODUCTION

Traditional multispectral images (MSIs, *e.g.* RGB images) usually contain a reduced number of spectral bands providing a limited spectral information. It is well-known that, the more spectral bands we have, the better we would understand

J. -F. Hu, T. -Z. Huang and L. -J. Deng is with the School of Mathematical Sciences, University of Electronic Science and Technology of China, Chengdu, Sichuan, 611731, China (e-mail: hujf0206@163.com; tingzhuhuang@126.com; liangjian.deng@uestc.edu.cn).

T.-X. Jiang is with the FinTech Innovation Center, Financial Intelligence and Financial Engineering Research Key Laboratory of Sichuan province, School of Economic Information Engineering, Southwestern University of Finance and Economics, Chengdu, Sichuan, 610074, China (e-mail: taixiangjiang@gmail.com).

G. Vivone is with the Department of Information Engineering, Electrical Engineering and Applied Mathematics, University of Salerno, 84084 Fisciano, Italy and with the Institute of Methodologies for Environmental Analysis, CNR-IMAA, 85050 Tito Scalo, Italy (e-mails: gvivone@unisa.it; gemine.vivone@imaa.cnr.it).

J. Chanussot is with Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, Grenoble, 38000, France (e-mail: jocelyn.chanussot@gipsa-lab.grenoble-inp.fr).
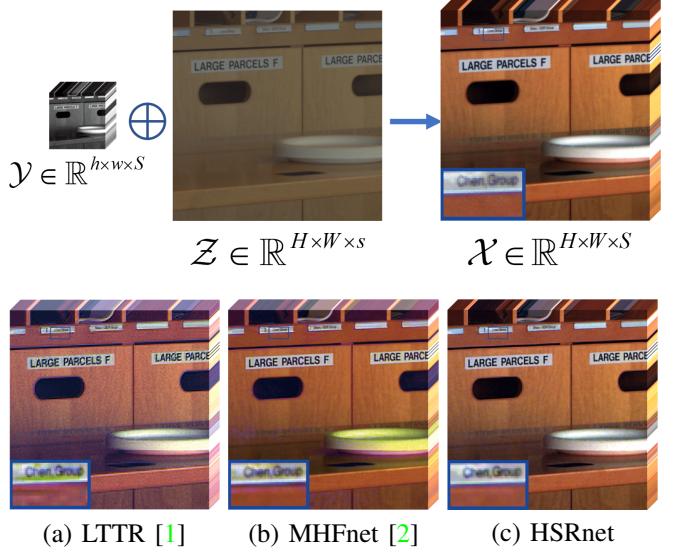


Fig. 1. First row: the schematic diagram of hyperspectral image resolution on a test image from the Harvard dataset ($h$ and $w$ represent the height and width of LR-HSI, $H$ and $W$ denote the height and width of HR-MSI, $s$ and $S$ denote the spectral band number of HR-MSI and LR-HSI, respectively). The right image is the ground-truth HR-HSI, $\mathcal{X}$. Second row: the results obtained by (a) LTTR (PSNR = 41.20dB), (b) MHFnet (PSNR = 38.70dB), and (c) the proposed HSRnet (PSNR = 43.93dB), where PSNR stands for the peak signal-to-noise ratio. Note that all the images are displayed with pseudo-color red, green, and blue (RGB) format using R = 28-th band, G = 12-th band, and B = 1-st band. Besides, MHFnet and HSRnet are both trained on the same CAVE dataset.

the latent spectral structure. Since hyperspectral imaging can obtain more spectral bands, it has become a non-negligible technology that is able to capture the intrinsic properties of different materials. However, due to the physical limitation of imaging sensors, there is a trade-off between the spatial resolution and the spectral resolution in a hyperspectral image (HSI), therefore it is burdensome to obtain an HSI with a high spatial resolution. In this condition, hyperspectral image super-resolution by fusing a low-resolution hyperspectral image (LR-HSI) with a high-resolution multispectral image (HR-MSI) is a promising way to address the problem.

Many researchers have focused on hyperspectral image super-resolution to increase the spatial resolution of LR-HSI proposing several algorithms. These latter are mainly based on the following models:

$$\mathbf{Y} = \mathbf{XBS}, \quad \mathbf{Z} = \mathbf{RX}, \tag{1}$$

where $\mathbf{Y} \in \mathbb{R}^{S \times hw}$, $\mathbf{Z} \in \mathbb{R}^{s \times HW}$ and $\mathbf{X} \in \mathbb{R}^{S \times HW}$ represent the mode-3 unfolding matrices of LR-HSI ($\mathcal{Y} \in \mathcal{R}^{h \times w \times S}$),

HR-MSI ($\mathcal{Z} \in \mathcal{R}^{H \times W \times s}$) and the latent HR-HSI ($\mathcal{X} \in \mathcal{R}^{H \times W \times S}$), respectively, $h$ and $w$ represent the height and width of LR-HSI, $H$ and $W$ denote the height and width of HR-MSI, $s$ and $S$ denote the spectral band number of HR-MSI and LR-HSI, respectively. Additionally, $\mathbf{B} \in \mathbb{R}^{HW \times HW}$ is the blur matrix, $\mathbf{S} \in \mathbb{R}^{HW \times hw}$ denotes the downsampling matrix, and $\mathbf{R} \in \mathbb{R}^{s \times S}$ represents the spectral response matrix. It is worth to be remarked that coherently with the notation adopted above, in this paper, we denote scalar, matrix, and tensor in non-bold case, bold upper case, and calligraphic upper case letters, respectively.

Based on the models in (1), many related approaches have been proposed. Different prior knowledge or regularization terms are integrated in those methods. However, the spectral response matrix $\mathbf{R}$ is usually unknown, thus the traditional methods need to select or estimate the matrix $\mathbf{R}$ and other involved parameters. Additionally, the related regularization parameters used in these kinds of approaches are often image-dependent.

Recently, with the tremendous development of neural networks, deep learning has become a promising way to deal with the hyperspectral image super-resolution problem. In [3], Dian *et al.* mainly focus on the spatial detail recovery learning image priors via a convolutional neural network (CNN). These learned priors have been included into a traditional regularization model to improve the final outcomes getting better image features than traditional regularization model-based methods. In [2], Xie *et al.* propose a model-enlightened deep learning method for hyperspectral image super-resolution. This method has exhibited an ability to preserve the spectral information and spatial details, thus obtaining state-of-the-art hyperspectral image super-resolution results.

However, deep learning-based approaches for hyperspectral image super-resolution also encounter some challenges. First of all, these methods sometimes have *complicated architectures* with millions of parameters to estimate. Second, due to the complicated architecture and large-scale training data, *expensive computation and storage* are usually involved. Third, deep learning-based methods are data-dependent, which usually holds a *weak network generalization*. Thus, the model trained on a specific dataset could poorly perform on a different kind of dataset. Instead, *the proposed network architecture can easily handle the above-mentioned drawbacks*.

In this paper, the proposed network architecture (called HSRnet from hereon) can be decomposed into two parts. One part is to preserve the spectral information of HR-HSI by upsampling the LR-HSI. The other part is mainly to get the spatial details of HR-HSI by training a convolutional neural network with the high-frequency information of HR-MSI and LR-HSI as inputs. By imposing the similarity between the network output and the reference (ground-truth) image, we can efficiently estimate the parameters involved in the network. In summary, this paper mainly consists of the following contributions:

1) The proposed network architecture is *simple* and *efficient*. As far as we know, it obtains better qualitative and quantitative performance than recent state-of-the-art hyperspectral image super-resolution methods. For example, our method shows significant improvements with respect to two state-of-the-art methods, one is deep learning-based [2] and the other one is regularization-based [1], see also Fig. 1. Besides, the proposed architecture involves fewer network parameters than other deep learning-based approaches thanks to our simple network design, more details are presented in Sec. IV-D.

2) The network architecture has a *promising generalization* ability to yield competitive results for different datasets even though the network is trained only on a specific dataset. This is due to the use of high-pass filters to feed the network with high-frequency spatial information. Extensive experiments corroborate this conclusion, see Fig. 10 and Tab. V.

3) Multi-scale information is integrated into our network architecture, which significantly improves the performance of the proposed method. The effectiveness of a multi-scale module has been proven in many computer vision works [4]–[9] and further discussed in Sec. IV-C.

4) The network shows a good *robustness to the number of training samples*, which indicates that our method could get very high performance with a different amount of training data. Furthermore, *shorter training and testing times* compared with a state-of-the-art deep learning-based approach (see Tab. IX) have been remarked.

The rest of the paper is outlined as follows. Section II presents the related works about the hyperspectral super-resolution problem. Section III introduces the proposed network architecture. In Section IV, extensive experiments are conducted to assess the effectiveness of the proposed architecture. Furthermore, some discussions about the image spectral response, the network generalization, the computational burden, and the use of the multi-scale module are provided to the readers.

## II. RELATED WORKS

Hyperspectral image super-resolution is a popular topic, which is receiving more and more attention. In particular, the combination of hyperspectral data with higher spatial resolution multispectral images is representing a fruitful scheme leading to satisfying results. Recent fusion or super-resolution approaches can be roughly categorized into two families: model-based approaches and deep learning-based methods.

Model-based approaches are classic solutions. Indeed, many works have been already published [1], [10]–[32] for the super-resolution problem. For instance, Dian *et al.* [27] exploit the spectral correlations and the non-local similarities by clustering the HR-MSI in order to create clusters with similar structures. Low tensor-train rank prior is used in [33], the so-called [1] method. The tensor train (TT) rank consists of ranks of matrices formed by a well-balanced matricization scheme. The effectiveness of low TT rank (LTTR) prior has been utilized in [34], which shows ability in image and video reconstruction. Compared to normal matrix ranks, the tensor rank keeps more abundant information about the data cube. Then, they regard the super-resolution as an optimization problem that, with the help of low tensor-train rank constraint,
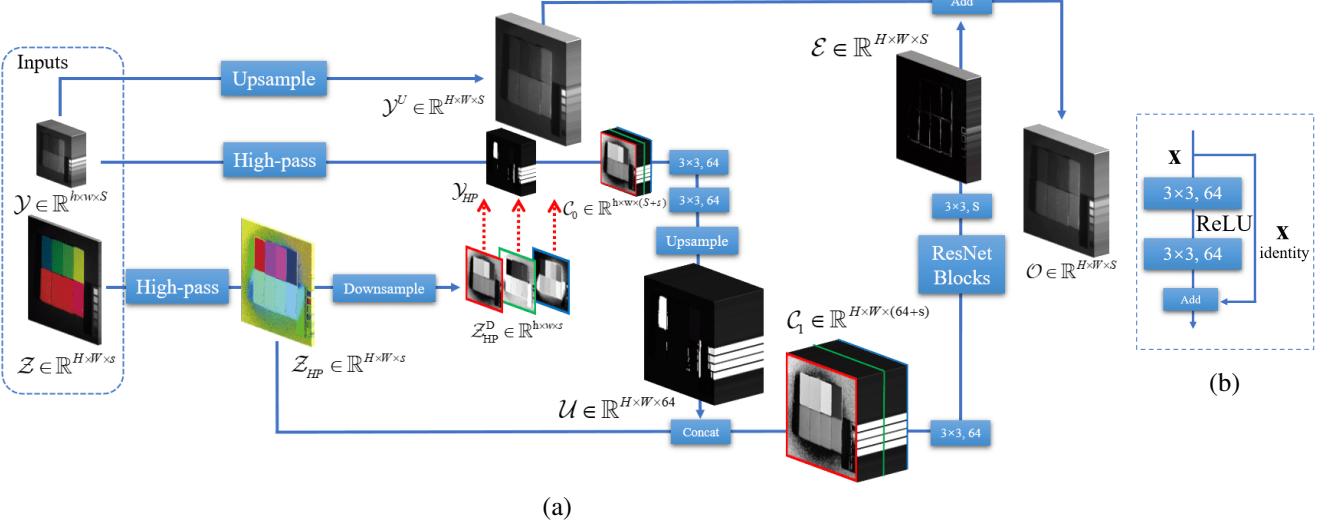
Fig. 2. The flowchart of the proposed network architecture (HSRnet). (a) Architecture of our HSRnet. LR-HSI $\mathcal{Y}$ and HR-MSI $\mathcal{Z}$ are the two inputs, and the $\mathcal{O}$ is the final output. (b) Illustration of one ResNet block with two layers and 64 kernels (size $3 \times 3$) for each layer.

has a satisfying solution under the well-known alternating direction multipliers minimization (ADMM) [35] framework.

Deep learning-based methods have recently showed exceptional performance in the field of image super-resolution, see *e.g.* [2], [3], [36]–[51]. A powerful example is provided by the so-called PanNet developed in [40]. Here, Yang *et al.* designed a new architecture training the deep-learning network with high-pass filtered details rather than original images. This is done in order to simultaneously preserve the spatial and spectral structures. Thanks to the use of high-pass filters, a greater generalization capability is observed. Another instance of deep learning-based methods for solving the hyperspectral image super-resolution issue is provided in [2], where a model-based deep learning method is proposed. The method exhibits a great ability to preserve structures and details, as well as it obtains state-of-the-art results. Unlike other deep learning-based methods that mainly regard the image super-resolution issue as a simple regression problem, this approach is based on the generation mechanism of the HSI and the MSI to build a novel fusion model. It adopts the low rankness knowledge along with the spectral mode of the HR-HSI under analysis. Instead of solving the model by traditional alternating iterative algorithms, the authors design a deep network learning the proximal operators and model parameters by exploiting CNNs.

## III. THE PROPOSED HSRNET

In this section, we introduce first the regularization-based model for the hyperspectral image super-resolution problem. Motivated by the above-mentioned model, we propose our network architecture that will be detailed in Sec. III-B.

### A. Problem Formulation

Estimating the HR-HSI from LR-HSI and HR-MSI is an ill-posed inverse problem. Thus, prior knowledge is introduced

exploiting regularization terms under the maximum a posteriori (MAP) framework. Those methods can be formulated as:

$$\min_{\mathbf{X}} \mathcal{L} = \lambda_1 f_1(\mathbf{X}, \mathbf{Y}) + \lambda_2 f_2(\mathbf{X}, \mathbf{Z}) + R(\mathbf{X}), \quad (2)$$

where $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ are the mode-3 unfolding matrices of tensor HR-HSI, LR-HSI, and HR-MSI, respectively, which have been introduced in Sec. I. $\lambda_1$ and $\lambda_2$ represent two regularization parameters, $f_1$ and $f_2$ force the spatial and spectral consistency, respectively, and $R$ stands for the regularization term depending on the prior knowledge. In general, $f_1$ and $f_2$ are defined based on the relations in (1), *i.e.* ,

$$\begin{aligned} f_1(\mathbf{X}, \mathbf{Y}) &= \|\mathbf{Y} - \mathbf{XBS}\|_F^2, \\ f_2(\mathbf{X}, \mathbf{Z}) &= \|\mathbf{Z} - \mathbf{RX}\|_F^2, \end{aligned} \quad (3)$$

where $\|\mathbf{X}\|_F = \sqrt{\sum \sum x_{ij}^2}$ is the Frobenius norm. In particular, the regularization term $R$ is crucial for regularization-based methods.

Deep learning can be viewed as an estimation problem of a function mapping input data with ground-truth (labeled) data. In our case, starting from the input images (*i.e.*, LR-HSI and HR-MSI), we can estimate the mapping function $f$ by minimizing the following expression:

$$\min_{\Theta} \ \mathcal{L} = \|f_{\Theta}(\mathbf{Y}, \mathbf{Z}) - \mathbf{X}\|_F^2 , \quad (4)$$

where $\mathbf{Y}$ and $\mathbf{Z}$ are the LR-HSI and the HR-MSI, respectively, and $\mathbf{X}$ is the reference (ground-truth) HR-HSI. The mapping function $f$ can be viewed as a deep convolutional neural network, thus $\Theta$ represents the parameters of the network. Besides, the prior knowledge can be viewed as being implicitly expressed by the learned parameters. In the next subsection, we will present the network architecture recasting the problem as in (4), where the function $f$ is estimated thanks to several examples provided to the network during the training phase.
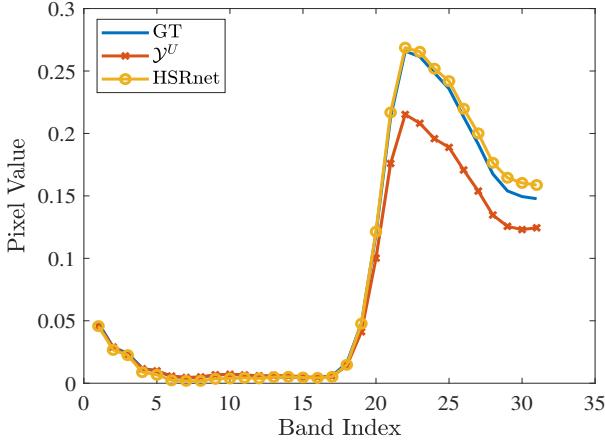
Fig. 3. Sampled spectral signatures for the object at pixel (175, 400) as obtained by the (ground-truth) HR-HSI, the upsampled LR-HSI $\mathcal{Y}^U$, and the estimated version of the high resolution HSI exploiting the proposed HSRnet.



Fig. 4. The residual maps: (a) $\mathcal{E} = \mathcal{O} - \mathcal{Y}^U$ and (b) $\mathcal{E}_{gt} = \mathcal{X} - \mathcal{Y}^U$.

### B. Network Architecture

Fig. 2 shows the proposed HSRnet for the hyperspectral image super-resolution problem. From the figure, it is easy to see that we decompose the network into two parts, such that the two parts can preserve the most crucial characteristics of a hyperspectral image, *i.e.,* the spectral information and the spatial details.

*1) Spectral preservation:* The LR-HSI $\mathcal{Y} \in \mathbb{R}^{h \times w \times S}$ [1] has the same spectral band number as the ground-truth HR-HSI $\mathcal{X} \in \mathbb{R}^{H \times W \times S}$. Indeed, most of the spectral information of the HR-HSI is contained in the LR-HSI (the remaining part is due to the spectral information of the high resolution spatial details). In order to corroborate it, we plot the sampled spectral signatures obtained by the ground-truth HR-HSI $\mathcal{X}$ and by the corresponding upsampled LR-HSI $\mathcal{Y}^U \in \mathbb{R}^{H \times W \times S}$ in Fig. 3. It is easy to be noted that the plots are very close to each other indicating that $\mathcal{Y}^U$ holds most of the spectral content of $\mathcal{X}$. Therefore, in order to guarantee a spectral preservation, we simply upsample $\mathcal{Y}$ getting $\mathcal{Y}^U$ (as shown in the top part of Fig. 2(a)).

Admittedly, $\mathcal{Y}^U$ is able to preserve the spectral information, but many spatial details are lost (which can retain part of the spectral information). Instead, the proposed HSRnet can learn the spectral information of the HR-HSI, even preserving the spatial counterpart. As a result, the final outcome of the proposed HSRnet clearly shows an almost perfect spectral preservation, see Fig. 3.

*2) Spatial preservation:* Since the HR-MSI $\mathcal{Z} \in \mathbb{R}^{H \times W \times s}$ contains high spatial resolution information, we aim to use $\mathcal{Z}$ to extract spatial details injecting them into the final hyperspectral super-resolution image. Moreover, $\mathcal{Y}$ still contains some spatial details, thus we also consider employing $\mathcal{Y}$ to extract them. However, we do not simply concatenate $\mathcal{Z}$ and $\mathcal{Y}$ together taking them into the network because that will not lead to a satisfying detail preservation. Indeed, we calculate first the spatial details at the LR-HSI scale, called $\mathcal{Y}_{HP}$ in Fig.

---

[1] We use three coordinates format to better represent the 3D hyperspectral image, *i.e.,* $h \times w \times S$.
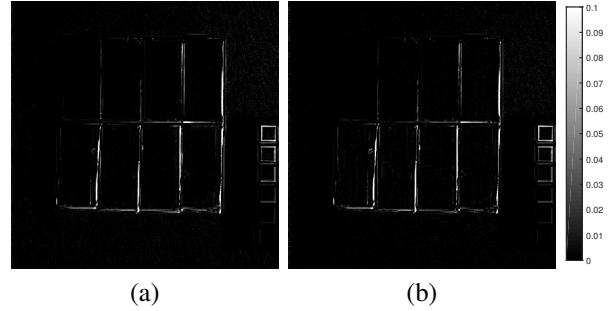
2. In particular, we simply take them from high-pass filtering the LR-HSI. Moreover, we add other details at the same scale by extracting them from the HR-MSI $\mathcal{Z}$. This is done by filtering and then downsampling the HR-MSI $\mathcal{Z}$ getting $\mathcal{Z}_{HP}^D$, see Fig. 2 again. This information has the advantage to occupy less memory and to require less computational burden to be processed compared to the original information in $\mathcal{Z}$. Finally, we concatenate this information, *i.e.* $\mathcal{Y}_{HP}$ and $\mathcal{Z}_{HP}^D$, to get $\mathcal{C}_0 \in \mathbb{R}^{h \times w \times (S+s)}$.

In order to complete the multiresolution analysis, thus introducing a multi-scale module in our network, the details at the HR-MSI scale are also extracted. This is performed by simply filtering them using a properly designed high-pass filter. These details, denoted as $\mathcal{Z}_{HP}$, can be concatenated with $\mathcal{C}_0$ (*i.e.*, the details at the lower scale) after this latter is properly convoluted and upsampled to the HR-MSI scale. Thus, $\mathcal{C}_1 \in \mathbb{R}^{H \times W \times (64+s)}$ indicates the concatenation of the details at two different scales (the LR-HSI one and the HR-MRI one). This represents the input of the ResNet implementing the well-known concept of multi-resolution analysis often considered in previously developed researches (*e.g.* [5]–[9]) either by designing diverse kernel sizes for convolution [5], [6] or extracting different spatial resolutions by filtering input data [7]–[9].

It is worth to be remarked that the high-pass filtering step is realized by the subtraction of the original image and its low-pass version, which is obtained by an average filter with a kernel size equal to $6 \times 6$. The upsampled operation is implemented by deconvolving with a kernel of size $6 \times 6$. Moreover, the concatenation operator is about adding the multispectral bands with high spatial resolution (3 bands, RGB image) into the hyperspectral bands (as shown in Fig. 2). In this work, the red, the green, and the blue slices of $\mathcal{Z}_{HP}^U$ and $\mathcal{Z}_{HP}$ are inserted as the head, the middle, and the tail frontal slices to complement the spectral information of the hyperspectral image.

Fig. 4 shows a comparison between $\mathcal{E}$ and $\mathcal{E}_{gt}$. From the figure, it is clear that $\mathcal{E}$ (*i.e.*, the details extracted by the proposed approach) and $\mathcal{E}_{gt}$ (*i.e.*, the details extracted by using the reference image) are very close to each other validating the effectiveness of the proposed network design. This result is only obtained thanks to the use of a multi-scale module combining details at two different scales guaranteeing a better spatial detail content in input of the ResNet.
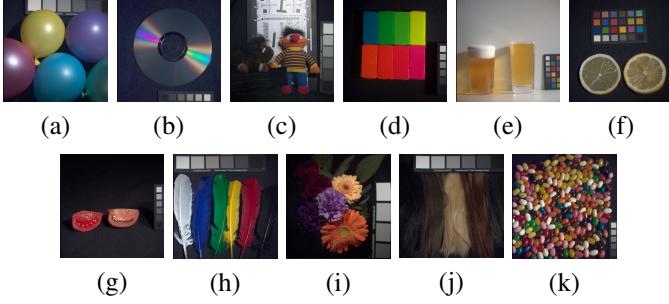
Fig. 5. The 11 testing images from the CAVE dataset. (a) *balloons*, (b) *cd*, (c)*chart and stuffed toy*, (d) *clay*, (e) *fake and real beers*, (f) *fake and real lemon slices*, (g) *fake and real tomatoes*,(h)*feathers*,(i)*flowers*,(j)*hairs*,(k)*jelly beans*.

*3) Loss function:* After obtaining the spectral preserved $\mathcal{Y}^U$ image and the spatial preserved $\mathcal{E}$ image from the ResNet fed by the image cube $\mathcal{C}_1$, we subsequently add the two outputs together to get the outcome. Thus, the loss function exploited during the training phase to drive the estimation of the function mapping in (4) can be defined as

$$\min_{\Theta} \mathcal{L} = \left\| f_{\Theta}(\mathcal{Y}_{HP}, \mathcal{Z}_{HP}^D, \mathcal{Z}_{HP}) + \mathcal{Y}^U - \mathcal{X} \right\|_F^2, \quad (5)$$

where $f_{\Theta}(\cdot)$ is the mapping function that has as input the details at the two different scales used to estimate the spatial preserved image $\mathcal{E}$ and the upsampled LR-HSI $\mathcal{Y}^U$. The loss function imposes the similarity between the network output $f_{\Theta}(\mathcal{Y}_{HP}^U, \mathcal{Z}_{HP}^D, \mathcal{Z}_{HP}) + \mathcal{Y}^U$ and the reference (ground-truth) $\mathcal{X}$ image.

## C. Network Training

*1) Training data:* In the work, we mainly use the CAVE dataset [52] for training the network. It contains 32 hyperspectral images with size $512 \times 512$ and 31 spectral bands. Additionally, each hyperspectral image also has a corresponding RGB image with size $512 \times 512$ and 3 spectral bands (*i.e.,* the HR-MSI image). We selected 20 images [2]for training the network, and the other 11 images to be considered for testing[3], as done for the MHFnet in [2]. The CAVE test images are shown in Fig. 5.

*2) Data simulation:* We extracted 3920 overlapped patches with a size of $64 \times 64 \times 31$ from the 20 images of the CAVE dataset used as ground-truth, thus forming the HR-HSI patches. Accordingly, the LR-HSI patches are generated starting from the HR-HSI by applying a Gaussian blur with kernel size equal to $3 \times 3$ and standard deviation equal to 0.5 and then downsampling the blurred patches to the size of $16 \times 16$, *i.e.,* with a downsampling factor of 4. Moreover, the HR-MSI patches (*i.e.,*, the RGB patches) are generated similarly as for the HR-HSI patches, but using the corresponding (already available) RGB data. Thus, other 3920 patches of size of $64 \times 64 \times 3$ are available to represent the HR-MSI. Following these indications, the patches for the training phase

[2]We selected the same 20 images as for the training of the MHFnet.
[3]One image, *i.e.*, "Watercolors", is discarded as it is unavailable for use.
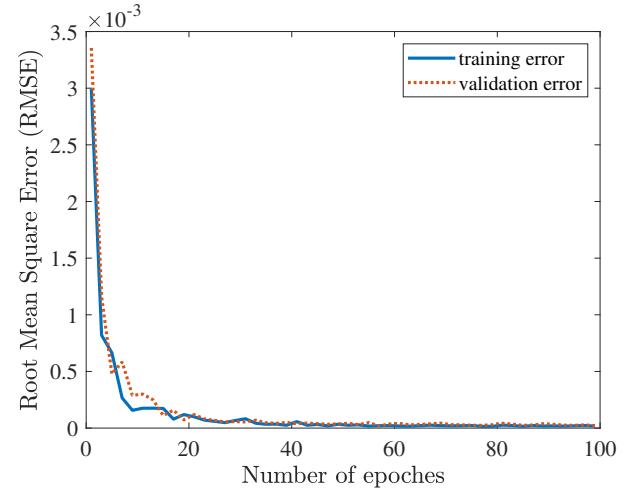


Fig. 6. Training and validation errors for the proposed HSRnet.

are the $80\%$ of the whole dataset and the rest (*i.e.*, the $20\%$) is used for the testing phase.

*3) Training platform and parameters setting:* The proposed network is trained on Python 3.7.4 with Tensorflow 1.14.0 and Linux operating system with NVIDA GPU GeForce GTX 2080Ti. We use Adam optimizer with a learning rate equal to 0.0001 in order to minimize the loss function (5) by 100,000 iterations and 32 batches. The ResNet block in our network architecture is crucial. Indeed, we use 6 ResNet blocks (each one with two layers and 64 kernels of size $3 \times 3$ for each layer. See Fig. 2). Fig. 6 shows the training and validation errors of the proposed HSRnet confirming the convergence of the proposed convolutional neural network using the above-mentioned parameters setting.

## IV. EXPERIMENTAL RESULTS

In this section, we compare the proposed HSRnet with several state-of-the-art methods for the hyperspectral super-resolution problem. In particular, the benchmark consists of the CNMF method[4] [21], the FUSE approach[5] [53], the GLP-HS method[6] [19], the LTTR technique[7] [1], the LTMR approach[8] [22], the MHFnet[9] [2], and the proposed HSRnet approach. For a fair comparison, the MHFnet is trained on the same training data as the proposed approach. Furthermore, the batch size and the training iterations of the MHFnet are set to 32 and 100,000, respectively, as for the proposed approach.

Two widely used benchmark datasets, *i.e.,* the CAVE database[10] [52] and the Harvard database[11] [54], are selected.

For quantitative evaluation, we adopt four quality indexes (QIs), *i.e.,* the peak signal-to-noise ratio (PSNR), the spectral

[4]http://naotoyokoya.com/Download.html
[5]http://wei.perso.enseeiht.fr/publications.html
[6]http://openremotesensing.net/knowledgebase/hyperspectral-and-multispectral-data-fusion/
[7]https://github.com/renweidian
[8]https://github.com/renweidian
[9]https://github.com/XieQi2015/MHF-net
[10]http://www.cs.columbia.edu/CAVE/databases/multispectral/
[11]http://vision.seas.harvard.edu/hyperspec/download.html

TABLE I
AVERAGE QIS AND RELATED STANDARD DEVIATIONS OF THE RESULTS ON
100 PATCHES EXTRACTED FROM THE TESTING IMAGES ON THE CAVE
DATASET. THE BEST VALUES ARE HIGHLIGHTED IN BOLDFACE.

| Method | PSNR | SAM | ERGAS | SSIM |
|---|---|---|---|---|
| CNMF | 31.4±3.3 | 5.95±4.0 | 8.19±19.2 | 0.96±0.04 |
| FUSE | 28.9±3.3 | 10.36±6.4 | 7.23±6.0 | 0.91±0.07 |
| GLP-HS | 30.8±4.0 | 6.28±3.9 | 5.9±4.8 | 0.94±0.05 |
| LTTR | 31.1±3.6 | 7.36±3.5 | 6.55±5.3 | 0.94±0.04 |
| LTMR | 30.6±3.5 | 7.61±3.6 | 7.25±6.8 | 0.93±0.04 |
| MHFnet | 35.1±5.9 | 7.29±7.2 | 30.7±146.4 | 0.96±0.03 |
| HSRnet | **38.2**±5.3 | **2.94**±1.8 | **2.99**±3.6 | **0.99**±0.01 |
| Best value | +∞ | 0 | 0 | 1 |

TABLE II
AVERAGE QIS AND RELATED STANDARD DEVIATIONS OF THE RESULTS ON
11 TESTING IMAGES ON THE CAVE DATASETS. THE BEST VALUES ARE
HIGHLIGHTED IN BOLDFACE.

| Method | PSNR | SAM | ERGAS | SSIM |
|---|---|---|---|---|
| CNMF | 32.2±4.5 | 14.96±5.2 | 8.79±4.8 | 0.911±0.04 |
| FUSE | 31.5±2.5 | 17.71±7.8 | 9.07±6.2 | 0.870±0.05 |
| GLP-HS | 35.4±2.7 | 7.91±3.0 | 5.61±3.6 | 0.946±0.02 |
| LTTR | 36.8±2.8 | 6.65±2.5 | 5.66±2.8 | 0.957±0.03 |
| LTMR | 36.2±2.7 | 7.66±2.9 | 5.70±2.7 | 0.949±0.03 |
| MHFnet | 43.3±2.8 | 4.34±1.5 | 2.33±1.4 | 0.989±0.01 |
| HSRnet | **44.0**±2.9 | **3.09**±1.0 | **1.93**±1.0 | **0.992**±0.00 |
| Best value | +∞ | 0 | 0 | 1 |

angle mapper (SAM) [55], the erreur relative globale adi-mensionnelle de synthèse (ERGAS) [56], and the structure similarity (SSIM) [57]. The SAM measures the average angle between the spectral vectors of the target and of the reference image. Instead, the ERGAS represents the fidelity of the image based on the weighted sum of mean squared errors. The ideal value in both the cases is zero. The lower the index, the better the quality. Finally, PSNR and SSIM are widely used to evaluate the similarity between the target and the reference image. The higher the index, the better the quality. The ideal value for SSIM is one.

### A. Results on CAVE Dataset

In order to point out the effectiveness of all the methods on different kinds of scenarios, we divide first the remaining 11 testing images on the CAVE dataset into small patches of size $128 \times 128$. Then, 100 patches are randomly selected. We exhibit the average QIs and corresponding standard deviations of the results for the different methods on these patches in Table I. From Table I, we can find that the proposed HSRnet significantly outperforms the compared methods. In particular, the SAM value of our method is much lower than that of the compared approaches (about the half with respect to the best compared method). This is in agreement with our previously developed analysis, namely that the proposed HSRnet is able to preserve the spectral features of the acquired scene.

Afterwards, we conduct the experiments on the whole 11 testing images. Table II presents the average QIs on the 11 testing images. To ease the readers' burden, we only show the visual results on *balloons*, *clay*, and *fake and real bears*. Table III lists the specific QIs of the results on these two images for the different methods. The proposed method outperforms the compared approaches. Furthermore, the running time of

TABLE III
QIS OF THE RESULTS BY DIFFERENT METHODS AND THE RUNNING TIMES
ON (A) *balloons*, (D) *clay*, AND (E) *fake and real beers* ON THE CAVE
DATASET. G INDICATES THAT THE METHOD IS RUNNING ON THE GPU
DEVICE, WHILE C DENOTES THE USE OF THE CPU. THE BEST VALUES
ARE HIGHLIGHTED IN BOLDFACE.

| Method | CNMF | FUSE | GLPHS | LTTR | LTMR | MHFnet | HSRnet |
|---|---|---|---|---|---|---|---|
| (a) 512 × 512 | | | | | | | |
| PSNR | 31.26 | 32.02 | 39.73 | 39.13 | 39.21 | 45.24 | **49.51** |
| SAM | 9.89 | 10.56 | 3.29 | 3.29 | 4.15 | 2.91 | **1.64** |
| ERGAS | 4.57 | 4.30 | 1.81 | 2.11 | 2.11 | 1.06 | **0.59** |
| SSIM | 0.926 | 0.928 | 0.975 | 0.980 | 0.980 | 0.992 | **0.996** |
| (d) 512 × 512 | | | | | | | |
| PSNR | 31.35 | 32.18 | 37.59 | 37.09 | 37.06 | 43.09 | **45.06** |
| SAM | 17.56 | 17.68 | 10.68 | 7.00 | 7.64 | 7.71 | **4.60** |
| ERGAS | 7.19 | 9.25 | 4.78 | 5.20 | 5.23 | 2.94 | **2.06** |
| SSIM | 0.926 | 0.900 | 0.963 | 0.976 | 0.973 | 0.986 | **0.993** |
| (e) 512 × 512 | | | | | | | |
| PSNR | 30.41 | 35.98 | 37.57 | 38.99 | 38.66 | 41.97 | **45.97** |
| SAM | 4.81 | 3.97 | 1.25 | 1.97 | 2.18 | 1.62 | **0.94** |
| ERGAS | 2.19 | 1.70 | 1.23 | 1.25 | 1.26 | 0.76 | **0.42** |
| SSIM | 0.965 | 0.962 | 0.969 | 0.975 | 0.972 | 0.986 | **0.992** |
| Average time(s) | 27.1(C) | 1.9(C) | 4.6(C) | 767.8(C) | 271.3(C) | 4.4(G) | **1.7**(G) |

the HSRnet is also the lowest one. In Fig. 7, we display the pseudo-color images of the fusion results and the corresponding error maps on three images. From the error maps in Fig. 7, it can be observed that the proposed HSRnet approach has a better reconstruction of the high resolution details with respect to the compared methods, thus clearly reducing the errors in the corresponding error maps.

The spectral fidelity is of crucial importance when the fusion of hyperspectral images is considered. In order to illustrate the spectral reconstruction provided by the different methods, we plot the spectral vectors for two exemplary cases, see Fig. 8. It is worth to be remarked that the spectral vectors estimated by our method and the ground-truth ones are very close to each other.

### B. Results on Harvard Dataset

The Harvard dataset is a public dataset that has 77 HSIs of indoor and outdoor scenes including different kinds of objects and buildings. Every HSI has a spatial size of 1392×1040 with 31 spectral bands, and the spectral bands are acquired at an interval of 10nm in the range of 420-720nm. 10 images are randomly selected for testing. The test images are shown in Fig. 9.

As in the previous settings, the original data is regarded as the ground-truth HR-HSI. The LR-HSI data is simulated as in Sec. III-C. Instead, the HR-MSI (not already available for this dataset) is obtained by applying the method provided by [58], where the spectral response functions are obtained from CIE[12].

We would like to remark that both our method and the MHFnet are trained on the CAVE dataset, and we directly test them on the Harvard dataset without any retraining or
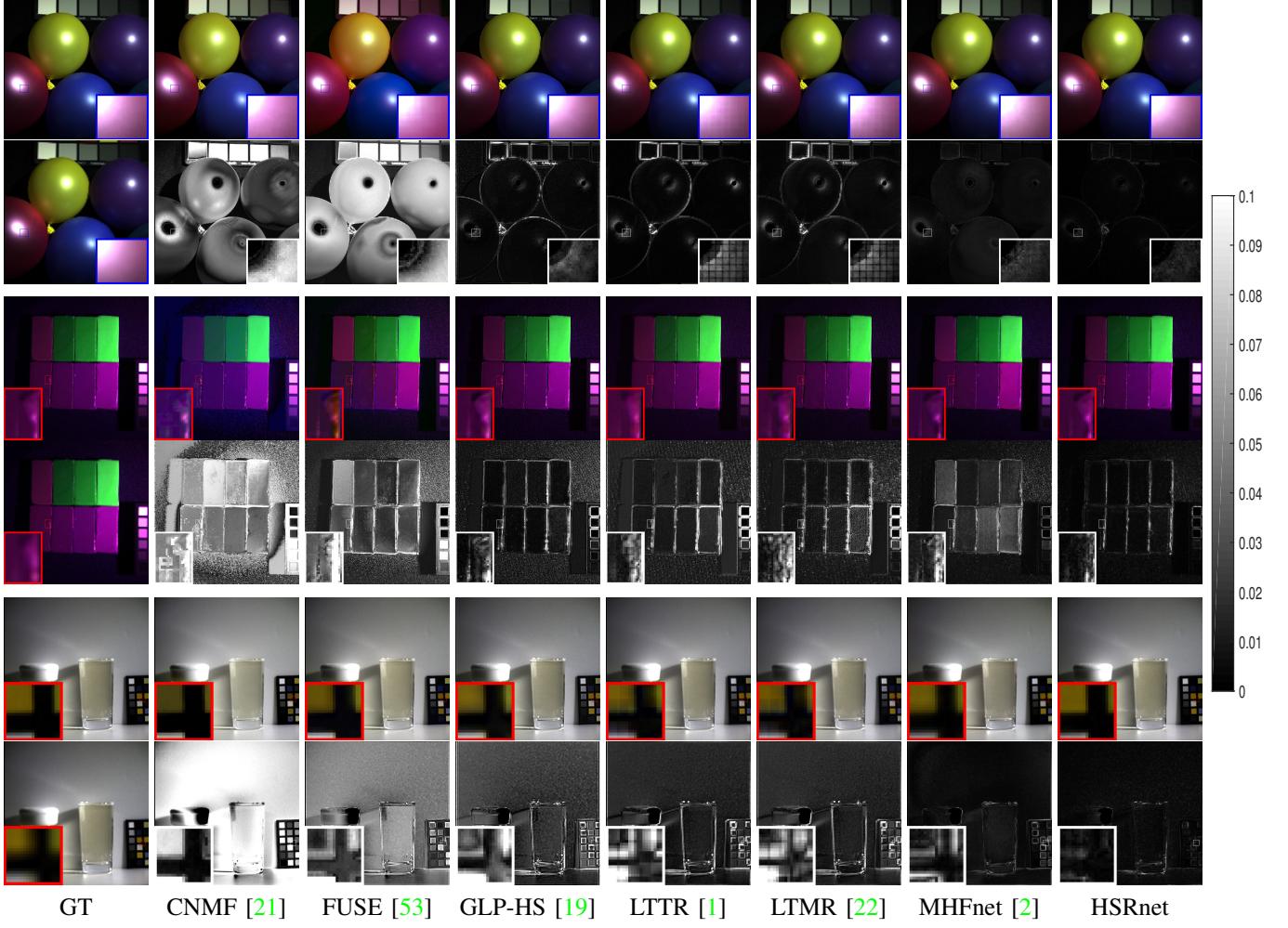
---
[12]http://www.cvrl.org

Fig. 7. The first column: the true pseudo-color images from the original CAVE dataset and the corresponding LR-HSI images of *balloons* (R-23, G-18, B-7) (1st-2nd rows), *clay* (R-3, G-16, B-2) (3rd-4th rows), and *fake and real beers* (R-24, G-23, B-18) (5th-6th rows). 2nd-8th columns: the true pseudo-color fused products and the corresponding residuals for the different methods in the benchmark pointing out some close-ups to facilitate the visual analysis.

TABLE IV
AVERAGE QIs AND RELATED STANDARD DEVIATIONS OF THE RESULTS ON
100 PATCHES EXTRACTED FROM THE IMAGES ON THE HARVARD DATASET.
THE BEST VALUES ARE HIGHLIGHTED IN BOLDFACE.

| Method | PSNR | SAM | ERGAS | SSIM |
|---|---|---|---|---|
| CNMF | 27.6±3.7 | 3.62±2.2 | 3.86±4.1 | 0.95±0.05 |
| FUSE | 26.7±3.7 | 5.40±4.1 | 4.07±4.0 | 0.94±0.06 |
| GLP-HS | 26.0±3.4 | 4.74±3.3 | 4.26±3.3 | 0.93±0.06 |
| LTTR | 27.4±3.5 | 4.65±2.5 | 4.87±3.1 | 0.94±0.06 |
| LTMR | 26.9±3.7 | 6.06±3.0 | 4.29±3.3 | 0.92±0.07 |
| MHFnet | 26.6±5.2 | 8.09±4.6 | 62.18±178.2 | 0.88±0.11 |
| HSRnet | **29.3**±4.4 | **3.44**±2.0 | **3.5**±2.2 | **0.97**±0.03 |
| Best value | +∞ | 0 | 0 | 1 |

TABLE V
AVERAGE QIs AND RELATED STANDARD DEVIATIONS OF THE RESULTS
FOR 10 TESTING IMAGES ON THE HARVARD DATASET. THE BEST VALUES
ARE HIGHLIGHTED IN BOLDFACE.

| Method | PSNR | SAM | ERGAS | SSIM |
|---|---|---|---|---|
| CNMF | 34.3±3.8 | 4.72±2.3 | 4.37±2.4 | 0.94±0.02 |
| FUSE | 32.9±3.8 | 7.48±3.5 | 4.79±2.0 | 0.93±0.03 |
| GLP-HS | 35.0±4.8 | 4.87±2.2 | 4.26±1.6 | 0.93±0.04 |
| LTTR | 36.1±5.4 | 6.06±2.3 | 6.19±2.2 | 0.90±0.07 |
| LTMR | 37.2±4.5 | 6.13±2.3 | 4.82±3.1 | 0.93±0.05 |
| MHFnet | 36.4±5.5 | 7.03±4.0 | 16.57±14.6 | 0.91±0.08 |
| HSRnet | **39.5**±4.7 | **3.38**±1.1 | **3.27**±1.5 | **0.97**±0.02 |
| Best value | +∞ | 0 | 0 | 1 |

fine-tuning. Thus, the performance on the Harvard dataset of these two methods could reflect their generalization abilities.

Moreover, we firstly divide these 10 testing images into patches of size $128 \times 128$ randomly selecting 100 patches. Table IV shows the QIs of the results for the different methods on these 100 patches. We can observe that our method is still the best method for all the different QIs, while the margins between our method and the MHFnet become larger than those in Table I. Particularly, the ERGAS value of the MHFnet

ranks last place. Thus, this test corroborates that the proposed approach has a better generalization ability than the compared deep learning-based method (*i.e.*, the MHFnet).

Table V records the average QIs and the corresponding standard deviations for the different methods using the 10 testing images. Table VI gives the QIs and the running times for three specific datasets of the Harvard dataset. The proposed method ranks first with the lowest running time. Finally, considering the details in the pseudo-color images in Fig. 10,

(a) *balloons* (276, 277)
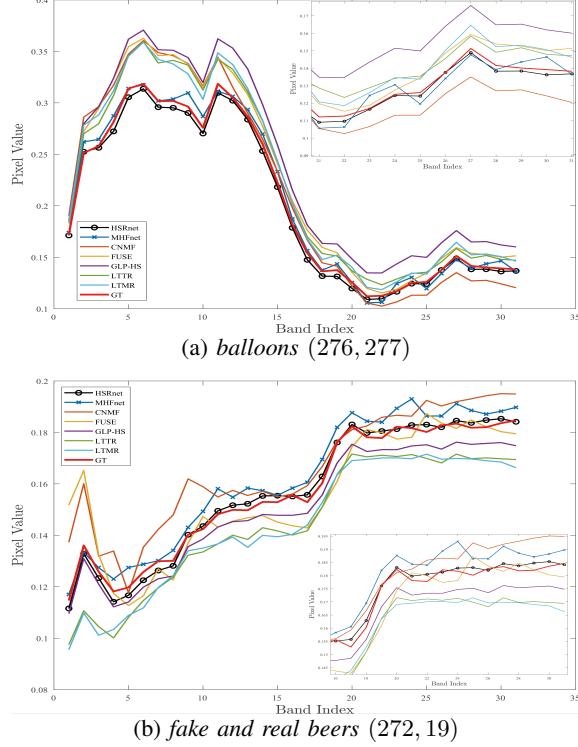


(b) *fake and real beers* (272, 19)

Fig. 8. Selected spectral vectors for the outcomes coming from the different fusion methods and the ground-truth (GT). The indications of the specific dataset and the location of the pixel under analysis are also provided.
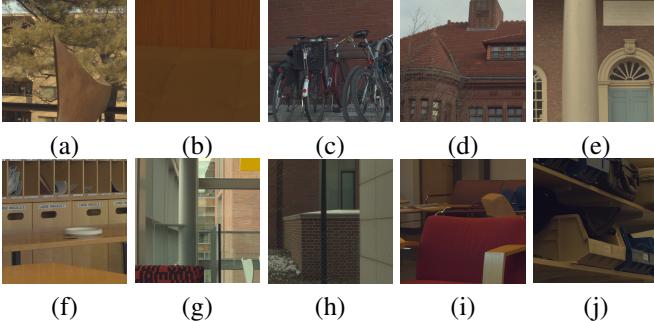


Fig. 9. The 10 testing images from the Harvard dataset. (a) *tree*, (b) *cushion*, (c)*bikes*, (d) *roof*, (e) *door*, (f) *cabinet*, (g) *window*, (h) *wall*, (i) *chairs*, (j) *baskets*.

we can see that the results of our method get the highest qualitative performance, thus obtaining error maps that are very dark (*i.e.*, with errors that tend to zero everywhere).

### C. Ablation Study

*1) High-pass filters:* In order to investigate the effects of the use of high-pass filters, we compare our HSRnet with its variant that is similar to the original HSRnet but without any high-pass filtering. After removing the high-pass filters, the data cube $\mathcal{C}_0$ in Fig. 2 is obtained by concatenating the LR-HSI $\mathcal{Y}$ and the downsampled version of the HR-MSI, *i.e.*, $\mathcal{Z}^D$. The network is trained on the same training data of the HSRnet with the same training settings. Table VII presents the average QIs of these two networks on the 11 testing images for the CAVE dataset and the 10 testing images for the Harvard dataset. As we can see from Table VII, the mean

TABLE VI
QIS OF THE RESULTS FOR THE DIFFERENT METHODS AND THE RUNNING TIMES ON (A) *trees* ,(C) *bikes*, AND (H) *wall* FOR THE HARVARD DATASET. G INDICATES THAT THE METHOD IS RUNNING ON THE GPU DEVICE, WHILE C DENOTES THE USE OF THE CPU. THE BEST VALUES ARE HIGHLIGHTED IN BOLDFACE.

| (a) $1000 \times 1000$ | | | | | | |
|---|---|---|---|---|---|---|
| Method | CNMF | FUSE | GLPHS | LTTR | LTMR | MHFnet | HSRnet |
| PSNR | 32.54 | 31.27 | 34.03 | 31.63 | 32.99 | 35.25 | **37.54** |
| SAM | 5.21 | 7.95 | 5.41 | 7.73 | 7.30 | 5.00 | **3.01** |
| ERGAS | 3.91 | 5.39 | 4.10 | 8.68 | 4.79 | 29.05 | **3.1** |
| SSIM | 0.912 | 0.882 | 0.911 | 0.829 | 0.867 | 0.917 | **0.961** |

| (c) $1000 \times 1000$ | | | | | | |
|---|---|---|---|---|---|---|
| Method | CNMF | FUSE | GLPHS | LTTR | LTMR | MHFnet | HSRnet |
| PSNR | 33.74 | 31.67 | 33.52 | 34.26 | 36.77 | 38.24 | **39.25** |
| SAM | 4.15 | 7.75 | 5.12 | 6.19 | 6.06 | 5.00 | **3.56** |
| ERGAS | 3.28 | 3.80 | 3.87 | 4.56 | 2.90 | 8.05 | **2.38** |
| SSIM | 0.938 | 0.924 | 0.879 | 0.908 | 0.938 | 0.957 | **0.974** |

| (h) $1000 \times 1000$ | | | | | | |
|---|---|---|---|---|---|---|
| Method | CNMF | FUSE | GLPHS | LTTR | LTMR | MHFnet | HSRnet |
| PSNR | 39.69 | 3.007 | 39.33 | 42.55 | 41.90 | 43.97 | **44.76** |
| SAM | 5.04 | 9.11 | 5.83 | 5.94 | 6.65 | 5.27 | **3.91** |
| ERGAS | 7.56 | 7.26 | 7.01 | 7.44 | 6.61 | 14.33 | **3.77** |
| SSIM | 0.921 | 0.942 | 0.959 | 0.974 | 0.972 | 0.977 | **0.989** |
| Average time(s) | 102.1(C) | 7.3(C) | 16.1(C) | 2049.5(C) | 864.3(C) | 6.8(G) | **2.4**(G) |

TABLE VII
AVERAGE QIS AND RELATED STANDARD DEVIATIONS OF THE RESULTS ON THE CAVE AND THE HARVARD DATASETS USING THE PROPOSED METHOD WITH AND WITHOUT THE HIGH-PASS (HP) FILTERS. THE BEST VALUES ARE HIGHLIGHTED IN BOLDFACE.

| CAVE | | | | |
|---|---|---|---|---|
| Method | PSNR | SAM | ERGAS | SSIM |
| w/o HP | 39.4±3.3 | 3.88±1.3 | 3.60±2.4 | 0.98±0.01 |
| with HP | **44.0±2.9** | **3.09±1.0** | **1.93±1.0** | **0.99±0.00** |
| Harvard | | | | |
| Method | PSNR | SAM | ERGAS | SSIM |
| w/o HP | 32.8±5.6 | 4.96±2.6 | 8.12±6.6 | 0.90±0.07 |
| with HP | **39.5±4.7** | **3.38±1.3** | **3.27±1.5** | **0.97±0.02** |

values and standard deviations of the proposed network are much better than that of the one without the high-pass filters. This demonstrates that the use of high-pass filters lead to better and more stable performance. In particular, the QIs of the Harvard testing images prove that the filters significantly provide better generalization ability. Thus, the high-pass filters are of crucial importance for competitive performance of the proposed HSRnet.

*2) Multi-scale module:* Concatenating multi-scale images is a key part of our network architecture. This leads to the extraction of several details at two different scales, which represent useful information for the super-resolution processing. To prove the strength of this module, we compare our original HSRnet and the simpler architecture that only uses the main scale, $\mathcal{C}_1$ in proposed Network is replaced by the one in Fig. 11. The results of the two compared approaches are reported in Table VIII. The QI values show the necessity of the multi-scale module in our HSRnet representing a part of the proposed architecture that is less important than the high-pass filtering, but relevant in order to improve the performance measured by some QIs, see *e.g.* the SAM and the ERGAS.
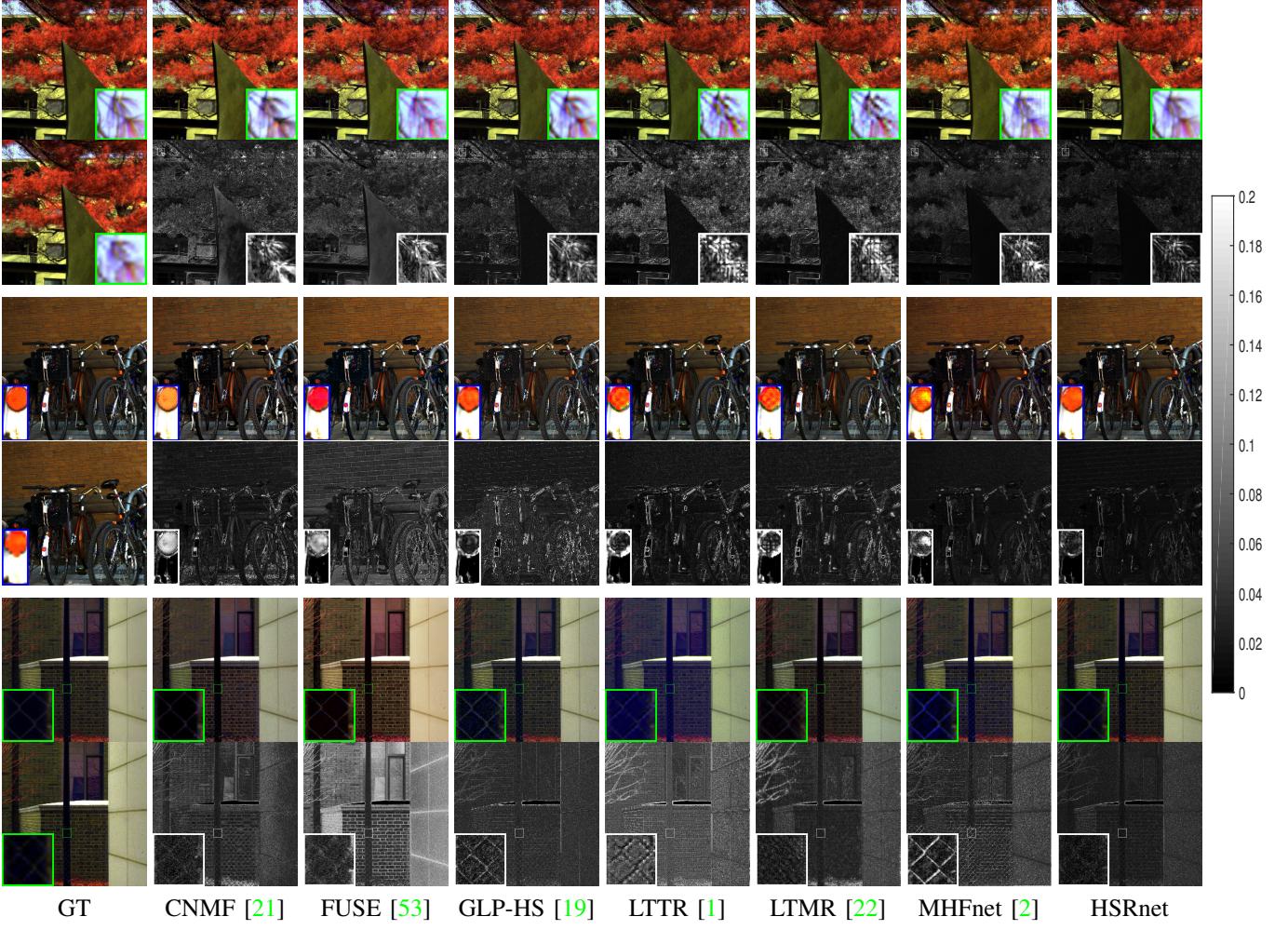
Fig. 10. The first column: the true pseudo-color images from the original Harvard dataset and the corresponding LR-HSI images of *tree* (R-30, G-27, B-7) (1st-2nd rows), *bikes* (R-31, G-18, B-9) (3rd-4th rows), and (h) *window* (R-31, G-28, B-1) (5th-6th rows). 2nd-8th columns: the true pseudo-color fused products and the corresponding residuals for the different methods in the benchmark pointing out some close-ups to facilitate the visual analysis.
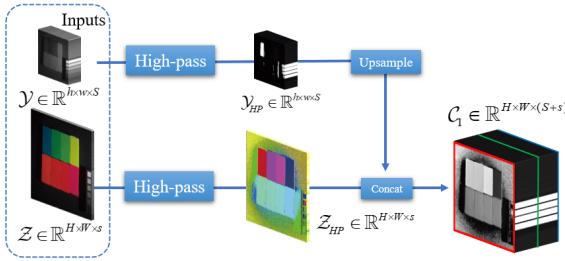


Fig. 11. Concatenation strategy with single scale. If we use this simple and single scale structure to replace the multiscale concatenation $\mathcal{C}_1$ of our HSRnet in Fig. 2, it will get worse outcome than our HSRnet, which validates the importance of our multi-scale concatenation.

TABLE VIII
AVERAGE QIS AND RELATED STANDARD DEVIATIONS OF THE RESULTS ON THE CAVE AND THE HARVARD DATASETS USING THE PROPOSED METHOD WITH A DIFFERENT NUMBER OF SCALES. THE BEST VALUES ARE HIGHLIGHTED IN BOLDFACE.

| CAVE | | | | |
|---|---|---|---|---|
| Method | PSNR | SAM | ERGAS | SSIM |
| one scale | 42.9±3.3 | 3.20±1.1 | 2.18±1.2 | **0.99**±0.00 |
| HSRnet | **44.0**±2.9 | **3.09**±1.0 | **1.93**±1.0 | **0.99**±0.00 |
| Harvard | | | | |
| Method | PSNR | SAM | ERGAS | SSIM |
| one scale | 38.8±4.4 | 3.66±1.9 | 3.64±1.8 | **0.97**±0.02 |
| HSRnet | **39.5**±4.7 | **3.38**±1.3 | **3.27**±1.5 | **0.97**±0.02 |

*D. Comparison with MHFnet*

To our knowledge, the MHFnet developed by Xie *et al*. [2] outperforms the state-of-the-art of the model-based and the deep learning-based methods, actually representing the best way to address the HSI super-resolution problem. Due to the fact that the MHFnet and our HSRnet are both deep learning-based methods, in this subsection, we keep on discussing about

the HSRnet comparing it with the MHFnet.

*1) Sensitivity to the number of training samples:* We train the MHFnet and our HSRnet with different numbers of training samples to illustrate their sensitivity with respect to this parameter. We randomly select 500, 1000, 2000, and 3136 samples from the training data. Testing data consists of 7 testing images on the CAVE dataset and 10 testing images on the Harvard dataset. Table IX reports the average QIs of the results obtained by the MHFnet and by our HSRnet varying

TABLE IX
RESULTS OF THE TWO DEEP LEARNING-BASED METHODS VARYING THE
NUMBER OF THE TRAINING SAMPLES. THE BEST VALUES ARE
HIGHLIGHTED IN BOLDFACE.

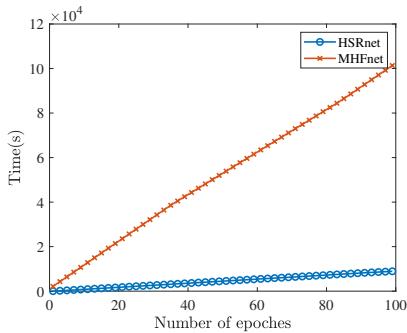| Datasets | # training data | Methods | PSNR | SAM | ERGAS | SSIM |
|---|---|---|---|---|---|---|
| CAVE | 3136 | MHFnet | 43.27 | 4.34 | 2.33 | 0.989 |
| | | HSRnet | **44.00** | **3.09** | **1.93** | **0.992** |
| | 2000 | MHFnet | 43.37 | 4.50 | 2.39 | 0.988 |
| | | HSRnet | **43.91** | **3.03** | **1.96** | **0.992** |
| | 1000 | MHFnet | **43.42** | 4.47 | 2.34 | 0.988 |
| | | HSRnet | 43.40 | **3.16** | **2.08** | **0.991** |
| | 500 | MHFnet | **42.74** | 4.77 | **2.50** | 0.987 |
| | | HSRnet | 40.99 | **3.65** | 2.89 | 0.987 |
| Harvard | 3136 | MHFnet | 36.41 | 7.03 | 16.57 | 0.915 |
| | | HSRnet | **39.53** | **3.38** | **3.27** | **0.970** |
| | 2000 | MHFnet | 36.54 | 6.93 | 18.42 | 0.912 |
| | | HSRnet | **39.87** | **3.40** | **3.33** | **0.970** |
| | 1000 | MHFnet | 36.16 | 6.99 | 26.49 | 0.916 |
| | | HSRnet | **39.44** | **3.47** | **3.54** | **0.968** |
| | 500 | MHFnet | 36.18 | 7.41 | 25.95 | 0.903 |
| | | HSRnet | **38.69** | **3.55** | **3.81** | **0.966** |



Fig. 12. The comparison of the training times for the MHFnet and the proposed HSRnet.

the number of the training samples. From the results on the CAVE dataset in Table IX, we can note that the MHFnet performs well when the training samples are less. This can be attributed to its elaborately designed network structure. Our method steadily outperforms the MHFnet in the cases of 2000 and 3196 training samples. Instead, from the results on the Harvard dataset, we can remark that the generalization ability of our method is robust with respect to changes of the numbers of the training samples (due to the use of the high-pass filters in the architecture). Whereas the MHFnet shows poor performance due to its manual predefined parameters that are sensitive to scene changes.

*2) Network generalization:* In the above content, MHFnet and our HSRnet are both trained with CAVE data. We can find that our HSRnet outperforms the MHFnet in all the experiments on the testing data provided by the Harvard dataset. This shows the remarkable generalization ability of our network. To further corroborate it, we retrain these two networks on training samples provided by the Harvard dataset. Namely, we extract from the Harvard dataset 3763 training samples, in which the HR-MSI is of size $64 \times 64$ and the LR-HSI is of size $16 \times 16$. As previously done, we select the same 11 images from the CAVE dataset and the same 10 images from the Harvard dataset to build the testing set. We

show the QIs of the results for these two networks trained on the Harvard dataset in Table X. It can be seen that the generalization ability of the MHFnet is still limited. Instead, the proposed approach still shows an excellent generalization ability when it is used on CAVE data but trained on the Harvard samples.

TABLE X
AVERAGE QIS AND RELATED STANDARD DEVIATIONS OF THE RESULTS
FOR THE NETWORKS TRAINED ON THE HARVARD DATASET. THE BEST
VALUES ARE HIGHLIGHTED IN BOLDFACE.

| CAVE | | | | |
|---|---|---|---|---|
| Method | PSNR | SAM | ERGAS | SSIM |
| MHFnet | 34.9±2.5 | 13.15±4.2 | 5.73±2.4 | 0.93±0.02 |
| HSRnet | **40.5**±2.8 | **4.21**±1.6 | **3.20**±1.6 | **0.98**±0.01 |
| Harvard | | | | |
| Method | PSNR | SAM | ERGAS | SSIM |
| MHFnet | **41.0**±5.3 | 3.36±1.6 | 3.33±1.9 | 0.97±0.02 |
| HSRnet | 40.1±5.5 | **3.06**±1.1 | **2.49**±1.0 | **0.98**±0.02 |

*3) Parameters and training time:* MHFnet contains 3.6 million parameters, instead, 2.1 million parameters have to be learned by our HSRnet. In Fig. 12, we plot the training time with respect to the epochs. We can find that our network needs much less training time than MHFnet. Actually, from Tables III and VI, the testing time of our HSRnet is also less than that of the MHFnet. Indeed, fewer parameters result in less training and testing times, making our method more practical.

## V. CONCLUSIONS

In this paper, a simple and efficient deep network architecture has been proposed for addressing the hyperspectral image super-resolution issue. The network architecture consists of two parts: $i$) a spectral preservation module and $ii$) a spatial preservation module that has the goal to reconstruct image spatial details starting from multi-resolution versions of input data. The combination of these two parts is performed to get the final network output. This latter is compared with the reference (ground-truth) image under the Frobenius norm based loss function. This is done with the aim of estimating the network parameters during the training phase.

Extensive experiments demonstrated the superiority of our HSRnet with respect to recent state-of-the-art hyperspectral image super-resolution approaches. Additionally, advantages of our HSRnet have been reported also from other points of view, such as, the network generalization, the limited computational burden, and the robustness with respect to the number of training samples.

## REFERENCES

[1] R. Dian, S. Li, and L. Fang, "Learning a low tensor-train rank representation for hyperspectral image super-resolution," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 9, pp. 2672 –2683, 2019. 1, 2, 5, 7, 9

[2] Q. Xie, M. Zhou, Q. Zhao, D. Meng, W. Zuo, and Z. Xu, "Multispectral and hyperspectral image fusion by MS/HS fusion net," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1585–1594. 1, 2, 3, 5, 7, 9

[3] 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. 2, 3

[4] X. Fu, B. Liang, Y. Huang, X. Ding, and J. Paisley, "Lightweight pyramid networks for image deraining," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–14, 2019. 2

[5] J. Li, F. Fang, K. Mei, and G. Zhang, "Multi-scale residual network for image super-resolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 517–532. 2, 4

[6] W. Ren, S. Liu, H. Zhang, J. Pan, X. Cao, and M.-H. Yang, "Single image dehazing via multi-scale convolutional neural networks," in *European conference on computer vision*. Springer, 2016, pp. 154–169. 2, 4

[7] H. Zhang, V. Sindagi, and V. M. Patel, "Multi-scale single image dehazing using perceptual pyramid deep network," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 902–911. 2, 4

[8] C. Yeh, C. Huang, and L. Kang, "Multi-scale deep residual learning-based single image haze removal via image decomposition," *IEEE Transactions on Image Processing*, vol. 29, pp. 3153–3167, 2020. 2, 4

[9] X. Wang, H. Ma, X. Chen, and S. You, "Edge preserving and multi-scale contextual neural network for salient object detection," *IEEE Transactions on Image Processing*, vol. 27, no. 1, pp. 121–134, 2018. 2, 4

[10] F. Fang, F. Li, C. Shen, and G. Zhang, "A variational approach for pan-sharpening," *IEEE Transactions on Image Processing*, vol. 22, no. 7, pp. 2822–2834, 2013. 2

[11] T. Wang, F. Fang, F. Li, and G. Zhang, "High-quality bayesian pan-sharpening," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 227–239, 2019. 2

[12] X. Fu, Z. Lin, Y. Huang, and X. Ding, "A variational pan-sharpening with local gradient constraints," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 265–10 274. 2

[13] C. I. Kanatsoulis, X. Fu, N. D. Sidiropoulos, and W.-K. Ma, "Hyperspectral super-resolution: Combining low rank tensor and matrix structure," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 3318–3322. 2

[14] K. Zhang, M. Wang, S. Yang, and L. Jiao, "Spatial–spectral-graph-regularized low-rank tensor decomposition for multispectral and hyperspectral image fusion," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 4, pp. 1030–1040, 2018. 2

[15] R. A. Borsoi, T. Imbiriba, and J. C. M. Bermudez, "Super-resolution for hyperspectral and multispectral image fusion accounting for seasonal spectral variability," *IEEE Transactions on Image Processing*, vol. 29, pp. 116–127, 2020. 2

[16] Y. Xing, M. Wang, S. Yang, and K. Zhang, "Pansharpening with multiscale geometric support tensor machine," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 5, pp. 2503–2517, 2018. 2

[17] B. Aiazzi, S. Baronti, and M. Selva, "Improving component substitution pansharpening through multivariate regression of MS + Pan data," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 45, no. 10, pp. 3230–3239, 2007. 2

[18] X. Han, J. Luo, J. Yu, and W. Sun, "Hyperspectral image fusion based on non-factorization sparse representation and error matrix estimation," in *2017 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2017, pp. 1155–1159. 2

[19] M. Selva, B. Aiazzi, F. Butera, L. Chiarantini, and S. Baronti, "Hyper-sharpening: A first approach on SIM-GA data," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 8, no. 6, pp. 3008–3024, 2015. 2, 5, 7, 9

[20] T. Akgun, Y. Altunbasak, and R. M. Mersereau, "Super-resolution reconstruction of hyperspectral images," *IEEE Transactions on Image Processing*, vol. 14, no. 11, pp. 1860–1875, 2005. 2

[21] 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. 2, 5, 7, 9

[22] 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. 2, 5, 7, 9

[23] Z. Pan and H. Shen, "Multispectral image super-resolution via RGB image fusion and radiometric calibration," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1783–1797, 2019. 2

[24] K. Zhang, M. Wang, and S. Yang, "Multispectral and hyperspectral image fusion based on group spectral embedding and low-rank factor-ization," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 3, pp. 1363–1371, 2016. 2

[25] 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. 2

[26] 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. 2

[27] R. Dian, S. Li, L. Fang, and Q. Wei, "Multispectral and hyperspectral image fusion with spatial-spectral sparse representation," *Information Fusion*, vol. 49, pp. 262–270, 2019. 2

[28] K. X. Rong, L. Jiao, S. Wang, and F. Liu, "Pansharpening based on low-rank and sparse decomposition," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 12, pp. 4793–4805, 2014. 2

[29] K. Rong, S. Wang, X. Zhang, and B. Hou, "Low-rank and sparse matrix decomposition-based pan sharpening," in *2012 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 2012, pp. 2276–2279. 2

[30] H. A. Aly and G. Sharma, "A regularized model-based optimization framework for pan-sharpening," *IEEE Transactions on Image Processing*, vol. 23, no. 6, pp. 2596–2608, 2014. 2

[31] P. Liu, L. Xiao, and T. Li, "A variational pan-sharpening method based on spatial fractional-order geometry and spectral–spatial low-rank priors," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 3, pp. 1788–1802, 2017. 2

[32] M. A. Veganzones, M. Simes, 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. 2

[33] I. V. Oseledets, "Tensor-train decomposition," *SIAM Journal on Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, 2011. 2

[34] J. A. Bengua, H. N. Phien, H. D. Tuan, and M. N. Do, "Efficient tensor completion for color image and video recovery: Low-rank tensor train," *IEEE Transactions on Image Processing*, vol. 26, no. 5, pp. 2466–2479, 2017. 2

[35] J. Bioucasdias and M. Figueiredo, "Alternating direction algorithms for constrained sparse regression: Application to hyperspectral unmixing," *In Workshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing*, pp. 1–4, 2010. 3

[36] S. Mei, J. J. X. Yuan, S. Wan, J. Hou, and Q. Du, "Hyperspectral image super-resolution via convolutional neural network," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 4297–4301. 3

[37] J. Hu, Y. Li, X. Zhao, and W. Xie, "A spatial constraint and deep learning based hyperspectral image super-resolution method," in *IEEE International Geoscience and Remote Sensing Symposium*, 2017, pp. 5129–5132. 3

[38] Q. Huang, W. Li, T. Hu, and R. Tao, "Hyperspectral image super-resolution using generative adversarial network and residual learning," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2019, pp. 3012–3016. 3

[39] F. Palsson, J. Sveinsson, and M. 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. 3

[40] J. Yang, X. Fu, Y. Hu, Y. Huang, X. Ding, and J. Paisley, "PanNet: A deep network architecture for Pan-Sharpening," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1753–1761. 3

[41] J. Yang, Y.-Q. Zhao, and J. C.-W. Chan, "Hyperspectral and multispectral image fusion via deep two-branches convolutional neural network," *Remote Sensing*, vol. 10, no. 5, p. 800, 2018. 3

[42] Z. Shao and J. Cai, "Remote sensing image fusion with deep convolutional neural network," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 5, pp. 1656–1669, 2018. 3

[43] Y. Li, J. Hu, X. Zhao, W. Xie, and J. Li, "Hyperspectral image super-resolution using deep convolutional neural network," *Neurocomputing*, vol. 266, pp. 29–41, 2017. 3

[44] W. Yao, Z. Zeng, C. Lian, and H. Tang, "Pixel-wise regression using u-net and its application on pansharpening," *Neurocomputing*, vol. 312, pp. 364–371, 2018. 3

[45] S. Vitale, "A cnn-based pansharpening method with perceptual loss," in *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 2019, pp. 3105–3108. 3

[46] F. Palsson, J. R. Sveinsson, and M. O. Ulfarsson, "Sentinel-2 image fusion using a deep residual network," *Remote Sensing*, vol. 10, no. 8, p. 1290, 2018. 3

[47] X. Han, j. Yu, J. Luo, and W. Sun, "Hyperspectral and multispectral image fusion using cluster-based multi-branch bp neural networks," *Remote Sensing*, vol. 11, no. 10, p. 1173, 2019. 3

[48] Y. Liu, X. Chen, Z. Wang, Z. Wang, R. Ward, and X. Wang, "Deep learning for pixel-level image fusion: Recent advances and future prospects," *Information Fusion*, vol. 42, pp. 158–173, 2018. 3

[49] W. Huang, L. Xiao, Z. Wei, H. Liu, and S. Tang, "A new pan-sharpening method with deep neural networks," *IEEE Geoscience and Remote Sensing Letters*, vol. 12, no. 5, pp. 1037–1041, 2015. 3

[50] Y. Rao, L. He, and J. Zhu, "A residual convolutional neural network for pan-shaprening," in *2017 International Workshop on Remote Sensing with Intelligent Processing (RSIP)*. IEEE, 2017, pp. 1–4. 3

[51] X. Y. Liu, Y. Wang, and Q. Liu, "Psgan: a generative adversarial network for remote sensing image pan-sharpening," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 873–877. 3

[52] F. Yasuma, T. Mitsunaga, D. Iso, and S. 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. 5

[53] 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. 5, 7, 9

[54] A. Chakrabarti and T. Zickler, "Statistics of real-world hyperspectral images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2011, pp. 193–200. 5

[55] R. Yuhas, J. Boardman, and A. Goetz, "Determination of semi-arid landscape endmembers and seasonal trends using convex geometry spectral unmixing techniques," *The 4th Annual JPL Airborne Geoscience Workshop*, 1993. 6

[56] L. Wald, *Data fusion: definitions and architectures: fusion of images of different spatial resolutions*. Presses des MINES, 2002. 6

[57] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004. 6

[58] CIE, "Fundamental chromaticity diagram with physiological axes–part 1," *Technical Report*, 2006. 6