

Review

# Survey of deep learning approaches for remote sensing observation enhancement

Grigorios Tsagkatakis <sup>1,2</sup> , Anastasia Aidini<sup>1,2</sup>, Konstantina Fotiadou<sup>1,2</sup>, Michalis Giannopoulos<sup>1,2</sup>, Anastasia Pentari<sup>1,2</sup>, Panagiotis Tsakalides <sup>1,2</sup>

<sup>1</sup> Signal Processing Lab (SPL), Institute of Computer Science, Foundation for Research and Technology - Hellas (FORTH), Crete, Greece

<sup>2</sup> Computer Science Department, University of Crete, Crete, Greece

\* Correspondence: greg@ics.forth.gr; Tel.: +30-2811 - 392725

Version September 4, 2019 submitted to Journal Not Specified

**Abstract:** Deep Learning, and Deep Neural Networks in particular, have established themselves as the new norm in signal and data processing, achieving state-of-the-art performance in image, audio, and natural language understanding. In remote sensing, a large body of research has been devoted to the application of deep learning for typical supervised learning tasks like classification. Less, yet equally important effort has also been allocated to addressing the challenges associated with the enhancement of low quality observations from remote sensing platforms. Addressing such channels is of paramount importance, both in itself, since high altitude imaging, environmental conditions and imaging systems trade-offs lead to low quality observation, as well as to facilitate subsequent analysis, like classification and detection. In this paper, we provide a comprehensive review of deep learning methods for the enhancement of remote sensing observations, focusing on critical tasks including single and multi-band super-resolution, denoising, restoration, pan-sharpening and fusion among others. In addition to the detailed analysis and comparison of recently presented approaches, different research avenues which could be explored in the future are also discussed.

**Keywords:** Deep Learning; Convolutional Neural Networks; Generative Adversarial Networks; Super-resolution; Denoising; Pan-sharpening; Fusion; Earth Observations; Satellite imaging

## 0. Introduction

Remote sensing of the environment, and especially Earth observation, is witnessing an explosion in terms of volume of available observations, which offer unprecedented capabilities towards the global-scale monitoring of natural and artificial processes [1,2]. However, the increase in volume, variety and complexity of measurements has led to a situation where data analysis is causing a bottleneck in the observation-to-knowledge pipeline. To address this challenge, machine learning approaches have been introduced for automating the analysis and facilitating the enhancement of remote sensing observations [3]. Unlike conventional machine learning approaches which first perform the extraction of appropriate hand-crafted features and then apply shallow classification techniques, deep machine learning and Deep Neural Networks (DNNs) in particular, have demonstrated astounding capabilities, which is primarily attributed to the automated extraction of meaningfully features, removing the need for identifying case-specific features [4]. The driving force behind the success of DNNs in image analysis can be traced to the following three key factors.

- More data available for training DNNs, especially for cases of supervised learning like classification, where annotations are typically provided by users.

- 31 • More processing power, and especially the explosion in availability of Graphical Processing  
32 Units (GPUs) which are optimized for high through processing of parallelizable problems like  
33 training DNNs.
- 34 • More advanced algorithms which have allowed DNNs to grow considerably both in terms of  
35 depth and output dimensions, leading to superior performance compared to more traditional  
36 shallow architectures.

37 The advent of DNNs has led to a paradigm shift in the remote sensing data analysis [5], where  
38 significant effort has been given in applying DNN method for supervised learning problems in  
39 imaging, like multi-class [6] and multi-label [7] image classification, autonomous vehicle operations  
40 [8], and accelerating magnetic resonance imaging [9] among others. While the problem of supervised  
41 classification of remote sensing observations has been under intense investigation in the past four  
42 years, in this work we consider the less explored situations involving remote sensing observation  
43 enhancement, which can be broadly considered as instances of inverse imaging problems.

44 In inverse imaging problems, the objective is to recover high quality signals from degraded, lower  
45 quality observations by inverting the effects of the degradation operation [10]. Indicative cases of  
46 inverse problems in optical remote sensing include estimation of higher spatial, spectral and temporal  
47 resolution observations, restoration through removal of noise, and enhancement of observation quality  
48 and characteristics through fusion of observations from different modalities. The challenge in inverse  
49 imaging problems lays in the fact that these problem are by nature ill-posed, i.e., a large number  
50 of different high resolution images (potentially infinite) can be mapped to the same low quality  
51 observations. As such, inverting the process amounts to properly exploiting any prior knowledge  
52 which might be available such as sparsity of representation or statistical priors [11].

53 In the past years, a large number of data-driven approaches have been presented for handling  
54 inverse problems where prior knowledge is automatically extracted through training machine learning  
55 systems with input-output pairs, while in the past few years, the machine learning approach of  
56 choice for this class of problems has been the deep learning framework, and prolific architectures like  
57 Convolutional Neural Networks (CNNs) [12]. Unlike typical natural images, in the context of remote  
58 sensing, a number of specific challenges are present and must be addressed including:

- 59 • High dimensionality of observations where Multispectral (HS) and Hyperspectral (HS) are often  
60 the input and thus exploiting spatial and spectral correlations is of paramount importance.
- 61 • Massive amounts of information encoded in each observation due to the large distance between  
62 sensor and scene, e.g. 400-600 km for low Earth orbit satellites, which implies that significantly  
63 more content is encoded in each image compared to a typical natural image.
- 64 • Sensor specific characteristics, including the radiometric resolution which unlike typical 8-bit  
65 imagery, in many cases involves observations of 12-bits per pixel.
- 66 • Challenging imaging conditions which are adverse affected by environmental conditions  
67 including the impact of the atmospheric effects such as clouds on the acquired imagery.

68 In this paper, we focus on passive imaging remote sensing platform, the majority of which  
69 involves multispectral and hyperspectral imaging instrument among spaceborne satellites and explore  
70 how DNNs can address associated inverse imaging problems targeting the quality enhancement of  
71 the acquired imagery. The rest of this paper is organized as follows: In Section 1, the prototypical  
72 DNN architectures considered in the related literature are presented including Convolutional Neural  
73 Networks (CNN), Autoencoders (AE), and Generative Adversarial Networks (GAN). In Section 2 we  
74 focus on the problem of super-resolution observation from either color, multispectral or hyperspectral  
75 observations while in Section 3 we explore methods related to the problem of pan-sharpening. Section 4  
76 present methodologies for denoising observations and estimating missing measurements while Section  
77 5 outlines the state-of-the-art in fusion of diverse sources, including the case of fusion of observation  
78 from different modalities. In Section 6, a set of challenges related to the problems at hand are presented,  
79 while future research endeavors and research endeavors which can not be directly mapped to existing

80 classes are presented. Note that in this work, we focus exclusively on the enhancement of remote  
 81 sensing observations, acquired either by an airborne, e.g. unmanned aerial vehicles, or spaceborne  
 82 satellites, and the observations are single images or image sequences (video) with one, three or multiple  
 83 spectral bands.

## 84 1. Deep Neural Network paradigms

85 In this section we provide an overview of the three major DNN approaches related to remote  
 86 sensing image enhancement, namely Autoencoders (AE), Convolutional Neural Networks (CNN), and  
 87 Generative Adversarial Networks (GANs).

### 88 1.1. Convolutional Neural Networks (CNNs)

89 CNNs are specific DNN architectures which involve the use of convolutions, instead of the more  
 90 traditional matrix multiplications, in a number of layers, and are an ideal tool for processing regularly  
 91 sampled data like 2D and 3D imagery. Prototypical CNN architectures for image analysis tasks are  
 92 composed of four key layers types, namely convolutional layers, non-linear activations, pooling and  
 93 fully connected layers [13].

#### 94 1.1.1. Key Components of CNNs

*Convolutional layers.* For the case of 2D imagery, the convolution operation applied on input  
 features  $I(i, j, k)$ , where  $(i, j)$  indicate the spatial location and  $k$  the input channel, is expressed as

$$Y^{(l)} = I^{(k)} * K^{(l,k)} = \sum_m \sum_n \sum_k I(m, n; k) K(i - m, j - n; l, k)$$

95 where  $K^{(l,k)} \in \mathbb{R}^{m \times n}$  is the convolution kernel of size  $m \times n$ , associated with input channel  $k$  and  
 96 output channel  $l$ . In typical setups, multiple kernels are learned for each convolutional layer. The  
 97 key motivation of utilizing convolutional layers is their ability to provide translation invariance with  
 98 respect to the detected features. Furthermore, the much smaller size of the kernels compared to the  
 99 input also promotes parameter sharing and sparsity in terms of connectivity, which also has a positive  
 100 impact in terms of storage requirements. Two additional parameters in the design of convolutional  
 101 layers, besides the size of the kernel, are the step-size of its application, a parameter known as stride  
 102 and the method of handling the boundaries.

103 *Activation function.* Each feature generated by convolving the inputs from the previous layer  
 104 with a kernel is then passed through a non-linear function called activation function. Although more  
 105 traditional CNN realizations employed activations like the sigmoid and hyperbolic tangent (tanh),  
 106 the majority of modern approaches utilize the Rectified linear unit, ReLU, which essentially imposes  
 107 a non-negativity constraint, i.e.  $\hat{x} = \max(0, x)$  and variations like the parametric PReLU [14] where  
 108  $\hat{x}_i = \max(0, x_i) + a_i \min(0, x_i)$  for the  $i$ -th channel. In addition to superior performance, ReLU variants  
 109 can also address the issues of vanishing gradient in deep architectures.

110 *Transposed convolution (or deconvolution) layers.* A requirement of CNN when applied in problems  
 111 like super-resolution is the increase in spatial resolution of the original input image, which has  
 112 been tackled through two approaches, (i) by performing an initial interpolation so that the DNN is  
 113 responsible for estimating the missing high frequencies, and (ii) by directly operating on the low spatial  
 114 resolution image and introducing deconvolution layers before producing the enhanced image [15].  
 115 While in convolutional layers, the stride, i.e. the step size for the application of the kernel, is a integer  
 116 number greater or equal to one, in deconvolutional layers a fractional stride value is utilized, which in  
 117 practice can still be realized through typical convolution operations using appropriate padding or by  
 118 reshuffling the outputs of the previous layer [16].

119 A key characteristic associated with CNN is the receptive field, a termed borrowed from the  
 120 biology which refers to the part of the sensed space that can elicit some neural response. Receptive fields

121 apply primarily to the case of CNN and refer to the spatial extent associated with each convolution  
122 kernel, while in the majority of CNN architectures for image enhancement, both generic and remote  
123 sensing ones, the convolution kernels are of size  $3 \times 3$ , and in some cases  $5 \times 5$  or even  $7 \times 7$ . However,  
124 due to the hierarchical structure of CNNs, each subsequent layer has a larger effective receptive field  
125 compared to the previous layer, allowing the CNN to progressively consider larger and larger parts of  
126 the input image. This process increases the receptive field by a factor of two for each additional layer.  
127 Dilated convolutions can also be utilized for increasing the receptive field by skipping pixels during  
128 their applications, e.g., a  $3 \times 3$  kernel is effectively applied on a  $5 \times 5$  region by padding, allowing for  
129 faster rates of increase. Note that pooling operators also increase the receptive field, however, they are  
130 not typically used in inverse problems due to the unwanted reduction in the output size. In terms of  
131 image enhancement, it has been shown that increasing the size of the receptive field, e.g., through the  
132 use of dilated convolutions, can lead to the extraction of more important features for generic image  
133 denoising [17] and super-resolution [18], as well as in despeckling SAR observations [19].

134 For completeness we note that in addition to the aforementioned layers, two other types of  
135 layers are common in the CNN literature, namely pooling and fully connected layers [13]. In order to  
136 reduce the size of the feature space and provide additional invariances, pooling layers, such as max  
137 pooling and average pooling, are introduced for aggregating the responses from small spatial regions  
138 and propagating the maximum or the average response respectively to subsequent layers. In fully  
139 connected layers, each node is connected with every other node from the previous layer, following the  
140 prescription of typical Multi-layer Perception architectures [20]. Fully connected and pooling layers are  
141 more frequently found in scenarios involving classification tasks where sequences of such layers are  
142 cascaded reaching the final output layer, e.g. [21–23], however, since reduction of dimensionality is not  
143 required in inverse imaging problems, they are not typically employed for observation enhancement  
144 tasks.

#### 145 1.1.2. Training and optimization of CNNs

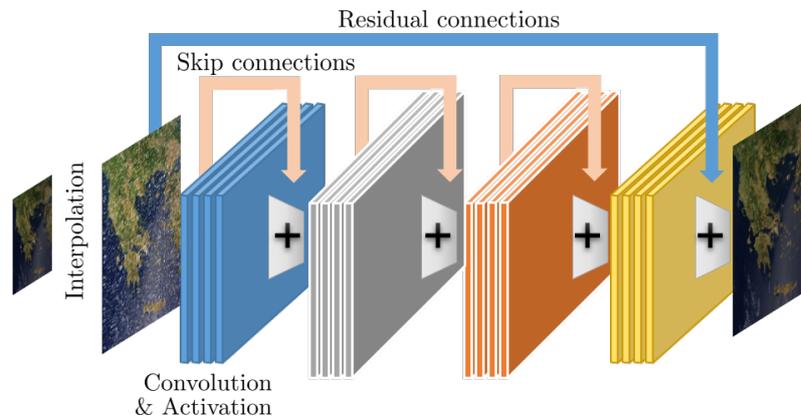
146 The process of training a CNN, i.e., identifying the optimal values for the kernel weights in the  
147 convolution and potentially fully connected or other layers, involves two major design choices, the  
148 selection of the appropriate loss function and the selection of the optimization process. Formally,  
149 given the CNN architecture parameters  $h_\theta$  and the input and target variables  $\mathbf{x}$ ,  $\mathbf{y}$ , the loss function  
150  $\mathcal{L} = \mathbb{E}(h_\theta(\mathbf{x}) - \mathbf{y})$  is typically the  $\ell_2$  norm between the predicted and the target high quality  
151 observation. The fundamental algorithm for training DNN is the Stochastic Gradient Descent (SGD)  
152 [24], however, state-of-the-art variations of SGD like the Adam [25] and the AdaDelta [26] offer very  
153 similar performance at a much lower computational cost.

154 While increasing the number of training examples typically leads to better performance, a  
155 challenging caveat is that the network can be too adapted to the training example and fail when  
156 new, unseen examples, are presented. This phenomenon is called overfitting and special measures  
157 must be introduced in order to combat it. To address this phenomenon, two important mechanisms  
158 have been proposed, namely dropout and regularization. In many cases, obtaining a sufficient number  
159 of training examples is not possible or it is very time and resource consuming. To address this situation,  
160 one approach involves training a DNN using a similar yet much larger dataset and then fine tuning the  
161 network with a much more limited number of case-specific examples. This process is known as transfer  
162 learning and its utilization in deep learning models [27] and has been instrumental for classification  
163 problems where a limited number of training examples is available.

164 Dropout [28] is a very successful method for addressing the problem of overfitting in terms of  
165 performance, yet extremely light in terms of computational requirements. This technique amounts to  
166 deactivating a number of nodes during each iteration of the training process. This forces the network  
167 to learn alternative paths for predicting the correct output, thus leading to greater flexibility in terms  
168 of generalization capacity. Although dropout can be introduced at any stage of training, the standard

169 practice is to introduce it in the fully connected layers since these layers account for the bulk of the  
 170 trainable parameters.

171 Batch Normalization is a local (per layer) normalizer, that operates on the node activations in a  
 172 way similar to the initial normalizing technique applied to the input data in the pre-processing step,  
 173 and studies [29] have shown that it can work very effectively as a method for addressing overfitting.  
 174 The primary goal of Batch Normalization is to enforce a zero mean and standard deviation of one for all  
 175 activations of the given layer and for each mini-batch. The main intuition behind Batch Normalization  
 176 lies in the fact that, as the neural network deepens, it becomes more probable that the neuronal  
 177 activations of intermediate layers might diverge significantly from desirable values and tend towards  
 178 saturation. This is known as Internal Covariate Shift [29] and Batch Normalization can play a crucial  
 179 role on mitigating its effects. Consequently, it can actuate the gradient descent operation to a faster  
 180 convergence, but it can also lead to overall higher accuracy and render the network more robust  
 181 against overfitting.



**Figure 1.** A typical CNN architecture for remote sensing image enhancement featuring convolutional and non-linear activation layers with residual and skip connections.

### 182 1.1.3. Prototypical CNN architectures

183 *Residual architectures.* An important extension of typical CNN architectures is the introduction of  
 184 the residual-learning concept [30], first proposed for single-image super-resolution [31], which has  
 185 allowed CNNs to grow much deeper without suffering the problem of vanishing/exploding gradients.  
 186 In essence, instead of learning a direct map from low quality inputs to high quality outputs, the CNN  
 187 is tasked with learning the residual, i.e., the difference between the low and high quality signals, which  
 188 typically represents missing high frequency information, at least for the case of super-resolution. In  
 189 order to allow networks to capture and extract features from multiple scale, *skip connections* between  
 190 different layers have also been considered and are now part of state-of-the-art approaches.

191 *Inception architectures.* An important design choice for a CNN is the appropriate kernel size, where  
 192 typical values are  $3 \times 3$  and  $5 \times 5$  for 2D image analysis. However, a single choice of kernel size  
 193 might lead to suboptimal performance. To address this challenge, the inception architectures [32,33]  
 194 introduce different size kernels at each convolutional layer and concatenate the filter outputs, while  
 195 the importance of each kernel is automatically adjusted during the training process. This way, feature  
 196 of different spatial extend are automatically selected for each layer, however, this benefit come at a  
 197 price of computational complexity due to the include in the number of trainable parameters.

198 *Hourglass architectures.* While in problems involving the classification of the input, the output space  
 199 is extremely smaller compared to the input space, for image enhancement problems, the dimensions  
 200 of the input and output spaces are typically the same, i.e., the size of the image itself. As such, CNN  
 201 architectures which progressively reduce the inputs to subsequent layers through pooling layers are  
 202 not appropriate. At the same time, critical features may reside at much higher dimensions compared

203 to the ambient dimensions of the input/output. A representative example of this architecture is the  
 204 SRCNN architecture by Dong et al. [15,34] which progressively shrinks and then expands the feature  
 205 extraction process through the introduction of deconvolution layers.

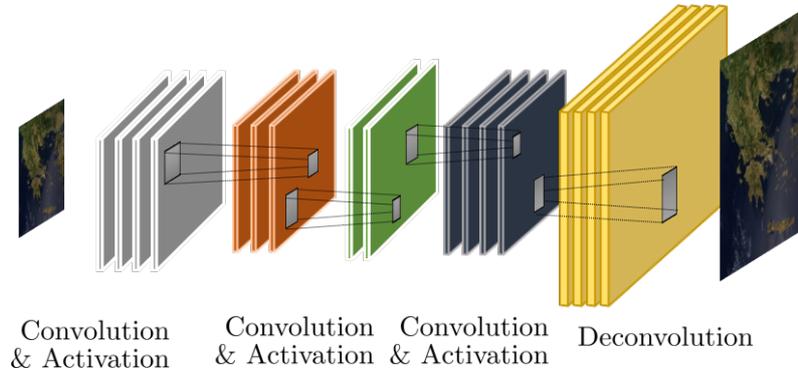


Figure 2. Hourglass shaped CNN architecture

206 *U-Net architectures.* Another very popular CNN architecture is the so-called U-Net architecture,  
 207 initially proposed for the problem of semantic segmentation in medical images [35]. In the U-Net  
 208 architecture, a symmetric design is followed where progressively convolutional and max pooling  
 209 layers is introduced leading to a very compact representation of the input image (contracting path),  
 210 which in sequence is expanded back to the dimensions of the input by convolution and upsampling  
 211 operators (expansive path).

### 212 1.2. Autoencoders (AE)

213 A classical autoencoder (AE) is a deterministic DNN comprised of an input and an output  
 214 layer of the same size with a hidden layer in between, which is trained with back propagation in a  
 215 fully unsupervised manner, aiming to learn a faithful approximation of the input [36]. Specifically,  
 216 the formulation considers both as an input and as the output examples  $\mathbf{s} \in \mathbb{R}^N$ , and encodes the  
 217 information through a non-linear function  $\sigma : \mathbb{R}^N \rightarrow \mathbb{R}^M$ , such that each input vector is mapped to a  
 218 new feature space via  $M$  hidden units.

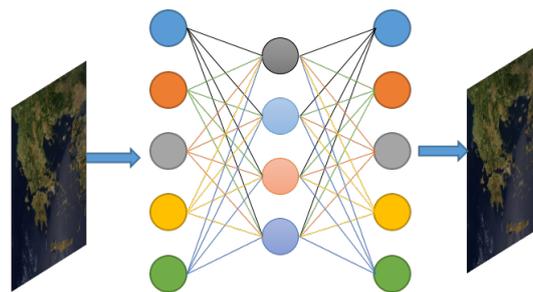


Figure 3. Typical structure of a single-layer autoencoder: This one-hidden layer structure learns the best possible compressed representation, so that the output is as close as possible to the input.

219 Formally, a single layer AE network consists of the input layer units  $\mathbf{s} \in \mathbb{R}^N$ , the hidden layer  
 220 unit,  $\mathbf{h} \in \mathbb{R}^M$ , and the output units  $\hat{\mathbf{s}} \in \mathbb{R}^N$ . and the objective is to learn a set of weights  $\mathbf{W} \in \mathbb{R}^{M \times N}$ ,  
 221 along with the associated encoding bias  $\mathbf{b}_1 \in \mathbb{R}^M$ , in order to generate compact and descriptive  
 222 features  $\mathbf{h} = \sigma(\mathbf{W}\mathbf{s} + \mathbf{b}_1)$  that can accurately reconstruct the input example. A typical example of  
 223 the function  $\sigma$  is the logistic sigmoid function. Afterwards, the decoding of  $\mathbf{h}$  is performed using  
 224 the weight matrix  $\mathbf{V} \in \mathbb{R}^{N \times M}$ , that connects the hidden layer with the output units  $\hat{\mathbf{s}} = \sigma(\mathbf{V}\mathbf{h} + \mathbf{b}_2)$ ,  
 225 where  $\mathbf{b}_2 \in \mathbb{R}^N$  stands for the decoding bias. In the majority of cases, tied weights are considered  
 226 such that  $\mathbf{W} = \mathbf{V}^T$ . Stacked Autoencoder (SSE) is considered the deep learning extension of AE where

multiple shallow (single layer) AE are stacked together and trained using greedy methods for each additional layer [37,38]. By applying a pooling operation after each layer, features of progressively larger input regions are essentially compressed into smaller ones, and thus are able to facilitate several classification or clustering tasks [39,40].

### 1.2.1. Sparse Autoencoders

An AE is closely related to the Principal Component Analysis (PCA), by performing an internal dimensionality reduction, an over-complete nonlinear mapping of the input vector  $\mathbf{s}$  can also be targeted by allowing more elements in the hidden layer, i.e.,  $M > N$ . In order to avoid trivial solutions like learning the identity function, additional constraints need to be imposed like sparsity of the internal representation. Consequently, to learn representative features, the error of the loss function:

$$\mathcal{L}(\mathbf{s}, \hat{\mathbf{s}}) = \frac{1}{2} \sum_{j=1}^N \|\hat{\mathbf{s}}_j - \mathbf{s}_j\|_2^2, \quad (1)$$

should be minimized, adhering to a sparsity constraint. In the aforementioned formulation,  $\mathbf{s}$  and  $\hat{\mathbf{s}}$  correspond to the input and the output data, respectively. In order to restrict the average activation to a small desired value, the sparsity constraint is imposed by introducing a Kullback-Leibler divergence regularization term such that:

$$\mathcal{L}(\mathbf{s}, \hat{\mathbf{s}}) = \frac{1}{2} \sum_{j=1}^N \|\hat{\mathbf{s}}_j - \mathbf{s}_j\|_2^2 + \beta \sum_{j=1}^M \text{KL}(\rho || \rho_j), \quad (2)$$

where  $\beta$  is a sparsity regularization parameter,  $M$  is the number of features,  $\rho$  is the average activation of  $\mathbf{h}$ ,  $\rho_j$  is the average activation of the  $\mathbf{h}_j$ -th vector over the input  $N$ -data, and  $KL$  denotes the Kullback-Leibler divergence regularization term defined as:

$$\text{KL}(\rho || \rho_j) = \rho \log \frac{\rho}{\rho_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \rho_j} \quad (3)$$

Consequently, the network learns weights such that only a few hidden nodes are activated by the given input.

### 1.2.2. Denoising Autoencoders

In a denoising autoencoder, the network is trained in order to reconstruct each data point from a corrupted input version [41]. Similar with the sparse autoencoder case, during the training process, the main objective is to learn both the encoder and decoder parameters. For this purpose, a noise distribution  $p(\tilde{\mathbf{x}}|\mathbf{x}, n)$  is defined, where the parameter  $n$  is associated to the amount of introduced noise. The autoencoder's weights are trained in order to reconstruct a random input from the training distribution. Formally, this process is summarized as follows:

$$(\theta^*, \theta'^*) = \arg \min_{\theta, \theta'} E_{(\mathbf{x}, \tilde{\mathbf{x}})} \left[ \mathcal{L}(\mathbf{X}, g_{\theta'}(f_{\theta}(\tilde{\mathbf{X}}))) \right], \quad (4)$$

where  $\mathcal{L}$  stands for the loss function.

One of the most crucial parameters in the denoising autoencoder is the noise distribution  $p$ . Gaussian noise is a common choice for continuous  $\mathbf{x}$ , defined as

$$p(\tilde{\mathbf{x}}|\mathbf{x}, n) = N(\tilde{\mathbf{x}}; \mathbf{x}, n), \quad (5)$$

256 while a masking noise can be used, when  $\mathbf{x}$  is a binary variable, as:

$$p(\tilde{\mathbf{x}}_i | \mathbf{x}_i, n) = \begin{cases} 0, & \text{with probability } n \\ \mathbf{x}_i, & \text{otherwise} \end{cases} \quad (6)$$

257 In any case, the amount of noise  $n$  affects the degree of corruption of the input. If  $n$  is high, the inputs  
 258 are more heavily corrupted during training. Additionally, the noise level provides a significant effect  
 259 on the representations that are learnt via the network. For instance, in scenarios where input data are  
 260 images, masking only a small amount of pixels will bias the process of learning the representation to  
 261 confront only local corruptions. On the other hand, masking larger areas will enforce the network to  
 262 utilize information from more distant regions.

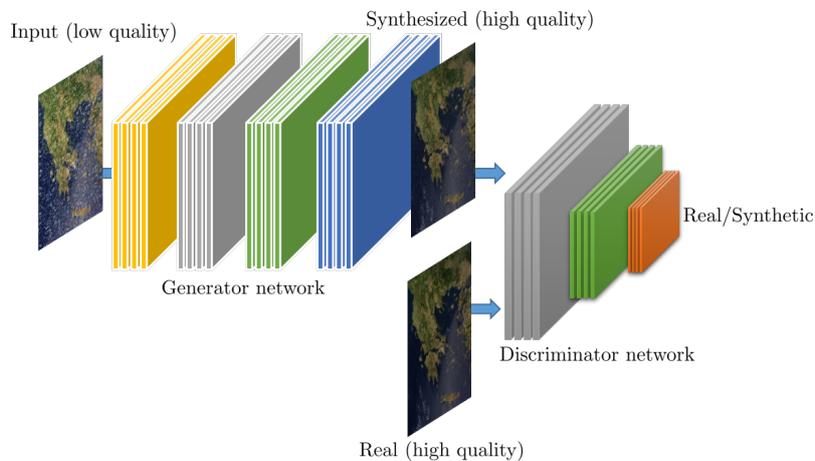
### 263 1.2.3. Variational Autoencoders

264 Nowadays, variational autoencoders (VAE) depict high performance in several generative  
 265 modeling tasks, such as handwritten digit recognition, face recognition, physical modeling of natural  
 266 scenes, and image segmentation among others [42]. Generative modeling focuses on learning models  
 267 of distributions  $P(X)$ , defined over data-points  $X$  that are spanned on a high-dimensional space  $\mathcal{X}$ .  
 268 However, the main limitation regarding the training phase of a generative model arises from the  
 269 complicated dependencies among the model's dimensions. In order to justify that the learnt model  
 270 represents efficiently the input dataset, we need to ensure that for every data-point  $X$  in the dataset,  
 271 there are exist latent variables able to synthesize successfully the input data-points.

272 Formally, let  $\mathbf{z}$  be the vector of latent variables lying on a high-dimensional space  $\mathcal{Z}$ , which can be  
 273 sampled according to some probability density function (PDF),  $P(\mathbf{z})$ , defined over  $\mathcal{Z}$ . Additionally, let  
 274  $f: \mathcal{Z} \times \Theta \rightarrow \mathcal{X}$ ,  $f(\mathbf{z}; \theta)$  be the family of deterministic functions parameterized by a vector  $\theta$  in some  
 275 space  $\Theta$ . The goal is to optimize  $\theta$  such that  $\mathbf{z}$  can be sampled from  $P(\mathbf{z})$ , and  $f(\mathbf{z}; \theta)$  can approach the  
 276 input data-points of  $\mathbf{X}$  with high probability. The aforementioned process is summarized as follows:

$$P(\mathbf{X}) = \int P(\mathbf{X} | \mathbf{z}; \theta) P(\mathbf{z}) d\mathbf{z} \quad (7)$$

277 Therefore, the main task of variational autoencoders (VAE) is the maximization of the probability  
 278 density function,  $P(\mathbf{X})$ . Unlike traditional sparse or denoising autoencoders, variational autoencoders  
 279 require no tuning parameters or sparsity penalties, while the output distribution is often chosen to be  
 280 Gaussian, such as:  $P(\mathbf{X} | \mathbf{z}; \theta) = \mathcal{N}(\mathbf{X} | f(\mathbf{z}; \theta), \sigma^2)$ .



**Figure 4.** A (conditional) Generative Adversarial Network architecture composed of a Generator and a discriminator network for image quality enhancement.

### 281 1.3. Generative Adversarial Networks

282 Generative Adversarial Networks (GANs) represent a radical new approach in DNN which  
 283 have been recently presented in [43] and have led to significant increase in terms of performance (e.g.  
 284 accuracy, quality), as well as have enabled the realization of new types of learning protocols, such as  
 285 the synthesis of extremely realistic images [44]. In the context of the GAN framework, instead of a  
 286 single DNN, training involves two DNNs, a “generator” and a “discriminator” network, where the  
 287 former synthesizes realistic signals given an input, and the later classifies inputs as real or synthetic.

288 In the original formulation of the GAN framework [43], the generator is seeded with randomized  
 289 noise input which leads to different output realizations, depending on its statistical characteristics.  
 290 For image enhancement problems, a variant of GANs, called conditional GANs (cGANs), is more  
 291 appropriate since the input to the generator is the image itself, although it could be very different from  
 292 the output, such as an edge map [45].

The prototypical GAN architecture for inverse imaging problems, involves an iterative training protocol that alternates between the synthesizing of a high quality image  $I^S$  given a low quality input image  $I^{IN}$ , performed by the generator  $G$ , and the classification of the high quality image as real  $I^R$  or synthetic  $I^S$ , executed by the discriminator  $D$ . Therefore, training a GAN corresponds to solving a min-max problem where the objective is to estimate the parameters (weights and biases) of the generator  $\theta_G$  and the discriminator  $\theta_D$ , given by

$$\min_{\theta_G} \max_{\theta_D} \mathbf{E}_{I^R \sim p_{train}(I^R)} [\log D_{\theta_D}(I^R)] + \mathbf{E}_{I^{IN} \sim p_G(I^{IN})} [\log(1 - D_{\theta_D}(G_{\theta_G}(I^{IN})))] \quad (8)$$

293 An important work that demonstrates the capabilities of GAN in typical inverse imaging problems,  
 294 is the Super Resolution GAN (SRGAN) architecture [46] which combines multiple loss functions,  
 295 including adversarial and perceptual, for single image super-resolution.

### 296 1.4. Performance evaluation

297 Obtaining examples for training a DNN model is in general is very challenging process, which for  
 298 the case of classification, requires significant labour and time. For the case of enhancement however,  
 299 the process can be substantially simplified by employing the Wald’s protocol [47]

#### 300 1.4.1. Evaluation metrics

- 301 • Peak Signal to Noise Ratio (PSNR) is a well established image quality metric, typically employed
- 302 for 8-bit images, an expresses image quality in decibels (dB), so that higher is better.
- 303 • Mean Squared Error (MSE) is a generic signal quality metric where lower values are sought.
- 304 • Structural Similarity Index (SSIM) is an image quality metric which considers the characteristics
- 305 of the human visual signal, such that values 1 indicate better performance.
- 306 • Signal-to-Reconstruction Error (SRE) is given in decibels (dB) and the higher values indicate
- 307 higher quality.
- 308 • Spectral Angle Mapper (SAM) is an image metric typically employed for MS and HS imagery
- 309 and values close to 0 indicate higher quality reconstruction.

#### 310 1.4.2. Remote sensing EO platforms

- 311 • Sentinel 2 satellite carrying the 12 band MSI instrument with spatial resolution between 10 and
- 312 20m for most bands in the visible, near and short-wave infrared.
- 313 • Landsat 7 and more recently Landsat 8 also acquire HS imagery over 11 bands in the visible to
- 314 infrared range at 30 and 100m resolution respectively.
- 315 • SPOT 6 and 7 provide PAN imagery at 2m spatial resolution and MS at 8m resolution offering a
- 316 swath of 60km
- 317 • Pleiades 1A and 1B acquire PAN imagery at 0.5m and MS at 2m, offering high agility for
- 318 responsive tasking.

- 319 • QuickBird and Ikonos are two decommissioned satellites offering 0.6m and 1m panchromatic  
320 (PAN) and 2.4 and 4m 4-band MS imaging respectively.
- 321 • Worldview-3 acquiring 0.31cm PAN and 1.24m eight-band MS imagery, as well as shortwave  
322 infrared imagery at 3.7m resolution.
- 323 • GaoFen-1 is equipped with two instrument sets, 2 cameras offering 2m PAN, 8m MS and four  
324 wide-field cameras acquiring MS imagery at 16m resolution.
- 325 • Jilin-1 satellite family by Chang Guang Satellite Technology Co., Ltd [48]
- 326 • The Airborne visible/infrared imaging spectrometer (AVIRIS) is a 224 band HS airborne  
327 instrument for which a number of annotated acquisitions are available.
- 328 • The Hyperion instrument aboard the EO-1 satellite is another well known HS instrument  
329 capturing imagery over 200 bands, which was decommissioned in 2017.
- 330 • ROSIS is a 102 band HS airborne instrument which is also associated with annotated acquisitions.

### 331 1.4.3. Typical datasets and data sources

- 332 • UC Merced contains 100 images of size  $256 \times 256$  from 21 different classes at 0.3m resolution.
- 333 • NWPU-RESIS45 [49] contains 700 images from 45 scene classes (31500 images in total).
- 334 • RSCNN7 contains 2800 images from seven different classes.
- 335 • NWPU-RESISC45 is a data set of aerial imagery of scenes from 45 categories, and each category  
336 contains 700 images of size of  $256 \times 256$  pixels.
- 337 • Indian Pines, Salinas, Cuprite and Kennedy Space Center from AVIRIS.
- 338 • Pavia Center and University from ROSIS offers 102 spectral bands over an urban area as 1.3m  
339 resolution.
- 340 • Botswana from EO-1 Hyperion offers 145 calibrated bands over different land cover types  
341 including seasonal swamps, occasional swamps, and drier woodlands.
- 342 • Kaggle Open Source Dataset and more specifically the Draper Satellite Image Chronology  
343 contains over 1000 high-resolution aerial photographs taken in southern California.
- 344 • Caltech, which is a publicly available aerial database that comprises an aircraft dataset and a  
345 vehicle dataset.

## 346 2. Super-resolution

### 347 2.1. Super-resolution of remote sensing observations

348 The problem of super-resolution is among the most studied problems in the context of DNN for  
349 remote sensing observation enhancement, where the objective is to increase the spatial and/or spectral  
350 resolution of observations, either as an end goal, or as an intermediate, in an effort to achieve higher  
351 classification accuracy. In order to training DNN architectures, the typical approach involves the  
352 generation of synthetic data through the application of a forward downsampling process by a given  
353 factor, and then the evaluation of the performance on the inverse upsampling process following the  
354 Wald's protocol [47]. Three of the most prolific CNN architectures, the Super-Resolution CNN (SRCNN)  
355 [34], the Very Deep Super Resolution (VDSR) [31] and the Enhanced Deep Super Resolution (EDSR)  
356 [50], have shown significant capacity in natural image super-resolution. In terms of methodological  
357 approach, the SRCNN is first CNN-based approach applied to this problem, the VDSR extends the  
358 SRCNN by introducing the notion of residual architectures and the EDSR introduces additional  
359 modifications making it among the highest performing methods for generic image super-resolution  
360 available today. We note that unlike typical CNN architectures, implementations that focus on  
361 super-resolution do not involve a pooling layer, since this operation reduces the spatial resolution,  
362 which is not desired in this context.

#### 363 2.1.1. Single image super-resolution

364 The use of DNN and CNN architectures for single remote sensing image super-resolution was  
365 first proposed by Liebel et al. [51]. In this work, the authors consider the paradigm of SRCNN [34]

366 for the spatial resolution enhancement of Sentinel-2 observations, but focus only on observations  
367 from the third band (559nm). The results demonstrate that using a pre-training SRNN network  
368 leads to significantly worse performance compared to "naive" methods like bicubic interpolation,  
369 while fine-tuning the network to the particular type of observation, leads to a marginal increase in  
370 performance (0.3dB). Similarly, in [52], the case of CNN for single-image super-resolution by factors  
371 of 2, 3 and 4 is investigated using observations from the SPOT 6 and 7 and the Pleiades 1A and 1B  
372 satellites. The authors compare the VDSR and the SRCNN architectures and report significant benefits  
373 of the VDSR compared to the SRCNN, however the gains offered by the VDSR are marginal compared  
374 to bicubic interpolation ( $\leq 1$ dB in PSNR). Huang et. al [53] also experimentally demonstrate that the  
375 VDSR applied to Sentinel-2, does not produce impressive results. Therefore, they introduce the Remote  
376 Sensing Deep Residual-Learning (RS-DRL) network, in which the features extracted by the SRCNN  
377 are introduced instead of the residual, achieving better performance than VDSR on Sentinel-2 images.

378 In order to encode information at different scales, such as those present in remote sensing imagery,  
379 Lei et al. [54] propose a modification of the VDSR architecture by directly propagating features  
380 from different convolutional layers to a "local-global" combination layer, the so-called "multifork"  
381 structure. The capabilities of the proposed LGCNet architecture are evaluated on RGB imagery from  
382 the UC-merced dataset and the results demonstrate that: (i) generic single-image super-resolution  
383 methods like the LGCNET and other DL methods like SRCNN can achieve higher quality estimation  
384 compared to bicubic for scale factor 3 on aerial imagery (1-2dB) and (ii) the proposed LGCNet method  
385 does not achieves any noteworthy improvement compared to generic DL approaches like [34] and [15].  
386 The extraction of features from multiple scales is also addressed by Xu et al. [55], who propose the  
387 DMCN architecture, a symmetric *hourglass*-structured CNN with multiple skip connections, termed  
388 memory connections, for single-image super-resolution. The reported results demonstrate an increase  
389 of less than 0.5dB and 0.02 in terms of PSNR and SSIM respectively on the NWPU-RESISC45, UC  
390 Merced and GaoFen1 datasets compared to the LGCNet [54] architecture. An interesting remark of this  
391 work is that the introduction of downsampling and corresponding upsampling layers lead to around  
392 53% and 67% reduction in memory footprint and inference time respectively.

393 Multi-scale feature extraction is also discussed in [56], where a joint wavelet and CNN based  
394 method is proposed for aerial image super-resolution. The method involves training multiple CNNs  
395 in multiple frequency bands generated by performing wavelet decomposition on aerial images, in an  
396 effort to approximate the wavelet multiscale representation and restore frequency features at different  
397 directions. For inference, each CNN is responsible for estimating the corresponding representation at  
398 the specific scale and the high resolution image is obtained through wavelet synthesis. The evaluation  
399 results indicate that the proposed method achieves marginal improvements of around 0.1dB compared  
400 to the VDSR architecture on noise-free images from the RSSCN7 dataset, while for the case of noisy  
401 imagery, a similar performance to the VDSR is reported.

402 Another super-resolution approach incorporating the wavelet transform and CNN equipped  
403 with local and global residual connections is proposed in Ma et al. [57] which addresses the issue by  
404 focusing on the frequency domain. Formally, the super resolution CNN accepts as inputs four subband  
405 representations of the input image obtained by a 2D discrete wavelet transform and the original low  
406 resolution input. The the fours input images are introduced to the proposed recursive residual CNN  
407 architecture and the outputs of the network are combined by an inverse discrete wavelet transform to  
408 generate the final high spatial resolution image. Evaluation of the method on the airplane images from  
409 the NWPU-RESISC45 dataset demonstrate marginal improvements ( $\leq 0.3$ dB) compared to DRNN [58],  
410 a recently presented DNN scheme for natural image super-resolution. Extraction of features across  
411 multiple scales for single satellite image super-resolution is also explored in the work by Lu et al. [59]  
412 where a multi-scale residual DNN is considered.

413 While moderate increase in the spatial resolution, e.g.  $\times 2, \times 3$ , of single remote sensing RGB  
414 images has been to a large extend addressed, the case of higher scale factors, e.g.  $\times 4, \times 8$  is still  
415 extremely challenging. This challenge is picked up by the authors of the Fast Residual Dense

416 BackProjection Network (FRDBPN) [60], which is heavily inspired by the method in [61], originally  
417 developed for natural color image super-resolution. The proposed methods employs upprojection and  
418 downscaling units, which consist of specific sequences of convolution and deconvolutions (dilated  
419 convolutions), as well as both global and local residual connections. Experimental results on the  
420 UC Merced dataset demonstrate a moderate ( $\leq 0.5\text{dB}$ ) increase in quality estimation compared to other  
421 DNN based natural single-image super-resolution methods

422 The Generative Adversarial Network (GAN) framework, which involves the interaction between  
423 a generator network and a discriminator network in order to produce high quality and realistic  
424 reconstructions, has been recently considered for the super-resolution of remote sensing images. Haut  
425 et al. [62] focus on a generative model, proposing an *hourglass* architecture network which is trained in  
426 an unsupervised way, i.e., without using an external dataset of examples. Specifically, the generator  
427 network starts from random noise as input and iteratively refines the high-resolution image estimation  
428 by performing a downsampling step that generates a low spatial resolution data which is fed to the  
429 second upsampling step. Compared to a large class of unsupervised state-of-the-art methods, the  
430 proposed model achieves certain gains from the UC-Merced, the SRCNN7 and the NWPU-RESIS45  
431 datasets, especially for  $\times 4$  resolution enhancement, which is more than 0.5dB for specific classes.  
432 Another GAN based approach is also explored in [63], which involves a series of modifications to the  
433 prototypical SRGAN [46] architecture. A major novelty of this work is that training of the architectures  
434 is conducted under a transfer learning protocol using a dataset of natural images before fine-tuning  
435 the network on aerial images from the UC Merced dataset. The proposed approaches demonstrated an  
436 average increase of 0.5dB compared to the generic SRGAN framework, which itself achieves a gain of  
437 1.3dB compared to the SRCNN framework.

438 Last, in [64], a GAN architecture is proposed which specifically targets the estimation of high  
439 frequency edge information from remote sensing imagery for video acquiring platforms. The  
440 generator network is composed of two sub-networks, namely the ultradense subnetwork which  
441 aims at reconstructing a baseline high resolution image and the edge-enhancement subnetwork which  
442 addresses the issue of noise contaminating the estimated edges. Comparisons with state-of-the-art  
443 super-resolution techniques for natural images demonstrate significant improvements, where for the  
444 Kaggle Open Source dataset, the proposed method surpasses SRCNN [34] by more than 2dB, the  
445 VDSR [65] by 1.5dB and the GAN based SRGAN [46] method by 1dB.

#### 446 2.1.2. Multispectral and Hyperspectral image super-resolution

447 In addition to single color image super-resolution, approaches for the enhancement of both spatial  
448 and spectral resolution of MS and HS observations have also been proposed. Efforts such as the ones  
449 by Hu et al. [66] as well as by C. Wang et al. [67] were among the first to consider DNNs for HS  
450 super-resolution, however, their investigation focused on natural scene, not remote sensing ones. One  
451 of the first DNN based approaches for remote sensing HS super-resolution was proposed by Yuan et al.  
452 [68] where the authors utilize the single-image super-resolution SRCNN architecture, pre-trained on  
453 natural images, for enhancing each band independently, and then employ a collaborative non-negative  
454 matrix factorization in order to enforce spectral consistency through the encoding of end-member  
455 information. In [69], another CNN architecture similar to the SRCNN [34] is considered for the  
456 automated spatial resolution enhancement of low-resolution MS observation, exploiting overlapping  
457 observation between two 4-band sensors aboard different satellites, for generating training examples.  
458 The authors specifically consider the fusion of low-spatial resolution images from the 740km swath  
459 Advanced Wide Field Sensor sensors and high-spatial resolution images from the 140km swath Linear  
460 Imaging Self Scanner sensors aboard the Indian Space Research Organisation's Resourcesat 1 and 2.

461 A two-step approach is considered in [70], where spatial resolution enhancement is achieved  
462 through an appropriate combination of deep Laplacian Pyramid Networks (LPNs) [71] and spectral  
463 enhancement is addressed through a dictionary learning with non-negativity constraint approach.  
464 Specifically, the LPN is trained using pairs of low and high resolution natural imagery and then applied

465 for the super-resolution of each band individually. Then a dictionary encoding the spectral content,  
 466 trained on low spatial resolution HS observations, is employed for increasing the spectral resolution  
 467 of the estimated bands. Experimental results with observations from the Indian Pines and the Pavia  
 468 Center dataset demonstrate substantial gains, 3.4dB and 2.2dB respectively, compared to the SRCNN  
 469 [34] method which achieved the second best performance.

470 An issue related to the previous CNN-based super-resolution methods is that quality, and therefore  
 471 the loss function, is quantified through the  $\ell_2$  norm, i.e.,  $\sum_i (\mathbf{x}_i - \hat{\mathbf{x}}_i)^2$ , where  $i$  indexes the examples.  
 472 Although this loss function may be appropriate for single band imagery, for the MS/HS case, it can lead  
 473 to spectral inconsistencies. To address this challenge, Zheng et al. [72] proposed a multi-losses function  
 474 Network for simultaneously enforcing spatial and spectral consistency. The composite loss function  
 475 in this case measures the normalized  $\ell_1$  norm between bands for capturing spatial quality and the  
 476 cosine plus the  $\ell_1$  between spectral pixels for spectral quality. Evaluation on the Pavia Center dataset  
 477 demonstrated a substantial increase of 1.8dB compared to the VDSR [31] and of 0.5dB compared to  
 478 EDSR [50], single image/band approaches.

479 A specific situation of MS super-resolution, which is however of great importance, involves  
 480 increasing the resolution of observation from the ESA Sentinel 2 satellite. A CNN based  
 481 super-resolution approach, inspired by the state-of-the-art single natural image EDSR method  
 482 [50], is considered for Sentinel 2 imagery [73], targeting the estimation of 10m ground sampling  
 483 distance images from 20m and 60m images. Training examples are generated through a simulated  
 484 downsampling processing of the available observations, based on the assumption that this degradation  
 485 process is closely related to the actual reduction in the observed spatial resolution. An important  
 486 observation, made by the authors, is that the proposed CNN is able to be generalize to any location  
 487 over the Earth, even though it is trained using examples from a small number of representative regions.  
 488 The results indicated that a very deep CNN architecture is able to achieve significantly higher quality  
 489 estimation, in the order of 6dB for the SRE metric compared to competing methods.

490 The problem of Sentinel 2 super-resolution is also investigated in [74], where a residual-structured  
 491 CNN is considered. A detailed investigation of the impact of different hyper-parameters on the  
 492 performance of observation enhancement from two locations, an urban and a coastal scene, revealed  
 493 that: (i) increasing the number of training patches can lead to worse performance due to overfitting  
 494 for certain scenes, (ii) increasing the number of training epoch offers minimal improvement, (iii) the  
 495 number of residual blocks is not a critical design parameter, and (iv) optimal values exist related to  
 496 patch size, and deviations from this size lead to worse performance. Gargiulo et al. [75] also consider  
 497 Sentinel 2 observations but focus on the super-resolution of the 11<sup>th</sup> SWIR band by posing the problem  
 498 as a pan-sharpening one and employing a CNN architecture similar to [76].

499 A limitation of typical DNN based on SR methods applied to HS observations is that when 2D  
 500 convolutions are applied, the convolutional layers will tend to fuse information across spectral bands.  
 501 For many cases, different bands are characterized by different SNR, therefore the fusion operation  
 502 can introduce noise to the extracted features. To address this issue, a Separable-Spectral convolution  
 503 and Inception Network (SSIN) is proposed in [77] in order to both extract band-specific features and  
 504 fuse features across bands. Validation of the Pavia center dataset demonstrated in increase of 0.8dB  
 505 compared to the VDSR [31] method for upscaling factor  $\times 2$  but almost no gains for  $\times 4$ .

506 The *simultaneous* encoding of both spatial and spectral information for the purpose of satellite  
 507 HS super-resolution was first proposed in [78]. To achieve this objective, a 3D CNN, i.e. a CNN  
 508 with three-dimensional convolutional kernels, is employed. The proposed architectures assumes an  
 509 interpolated image as input and the application of sequences of 3D convolution for estimating the high  
 510 quality output, measured in terms of MSE. Experimental validation using observations from the ROSIS  
 511 and HYDICE sensors demonstrate that using 3D convolutions lead to higher quality compared to 2D  
 512 methods like [51] while another observation is that the optimal performance is achieved by a moderate  
 513 number of convolutional layers and that more layers lead to marginal decrease in performance.

514 The potential of simultaneous spatio-spectral super-resolution in addition to sub-pixel  
515 classification of MS observations from aerial vehicles (drones) is examined in [79]. To achieve both  
516 tasks, a network inversion based architecture is proposed where a CNN architecture accepts both  
517 the low resolution image, as well as image-derived features like histograms-of-oriented-gradients.  
518 Experiments on airborne datasets such as Indian Pines and Salinas, indicate that this super-resolution  
519 approach better preserves the spectral fidelity when compared to sparse coding approaches which  
520 generally in corrupted dictionaries and sparse codes.

521 Ran et al. [80] consider a hierarchical CNN architecture which learns an end-to-end mapping  
522 from low resolution to high-resolution images, by combining feature extraction and edge enhancement  
523 in the hierarchical layers. Extensions of this approach based on residual learning and multi-scale  
524 version are also investigated for further improvements. The performance of the proposed approach  
525 is compared with bicubic interpolation and dictionary based methods on the areal color, MS and HS  
526 images suggest that better recovery quality and reduced computational requirements compared to  
527 dictionary based methods.

528 Last, we note two interesting approaches for spectral super-resolution using CNN. The first  
529 one explores the use of various CNN architecture for estimating extended spectral resolution MS  
530 imagery from 3 channel RGB images, without resorting to additional hardware components or imaging  
531 protocols, by exploiting the correlations between the RGB values and the corresponding HS radiance  
532 values. The HSCNN+ method [81] proposes deep residual and densely connected CNN architectures,  
533 which do not require an explicit upsampling pre-processing operator. The proposed method achieved  
534 the first place in the NTIRE 2018 Spectral Reconstruction Challenge for both the “Clean” and “Real  
535 World” tracks, which however consider natural scenes, not remote sensing ones. Similar approaches  
536 are also considered in [82,83].

537 The second approach involves the use of DNN in the framework of Compressed Sensing [84].  
538 Unlike traditional MS/HS architecture, approaches based on Compressed Sensing employ randomized  
539 sampling processes in order to acquire some form of compressed observations and subsequently  
540 employ optimization algorithms for estimating the image. By doing so, such approaches are able to  
541 overcome different trade-offs related to acquisition, including producing imagery with higher spatial  
542 resolution compared to the number of detectors used or provide a full spectral resolution MS from  
543 a single observation e.g. [85,86]. In the context of Compressed Sensing, Yuann et al. [87] consider  
544 the super-resolution of observations acquired by a Compressed Sensing architectures, specifically  
545 the CASSI [88], using CNN architectures, which allow demonstrates promising performance under  
546 simulated conditions. A similar approach involving the recovery of HS observations from Compressed  
547 Sensing measurements is also considered in [89], while in [90], in addition to the recovery process,  
548 optimization of the acquisition process is also explored.

### 549 2.1.3. Video super-resolution

550 While the majority of remote sensing platforms acquire still imagery, grayscale (panchromatic),  
551 RGB, or hyperspectral, a new breed of platforms focuses on the acquisition of video imagery, including  
552 the Jilin-1 satellites. To achieve however high temporal resolution, the imaging architecture typically  
553 needs to sacrifice spatial resolution, mandating the application of super-resolution for providing high  
554 quality spatio-temporal observations. In [91], a CNN architecture based on the VDSR [31] is explored  
555 for super-resolution of video satellite frames. Training examples from a high-spatial resolution static  
556 imaging satellite (Gaofen 2) are utilized, while the learning process is applied for super-resolving  
557 single frames from video sequences, acquired by another platform (Jilin 1). To address issues related  
558 to inability of existing methods to enhance image boundaries, the authors proposed an appropriate  
559 padding of the boundaries. The proposed method achieves a gain of 0.5dB when trained on high  
560 resolution satellite imagery compared to VDSR trained on natural imagery.

561 Another approach for satellite video super-resolution is proposed in [92] in which extracted video  
562 frames are used for both training and testing the DNN architecture. The proposed network is a variant

563 of the SRCNN method [34] where a deconvolution layer is introduced for producing high resolution  
 564 outputs, thus removing the need for pre-processing. Validation for an upscale factor of 3 on a diverse  
 565 set of scenes acquired by the Jilin-1 platform demonstrates an increase between 0.5 and 1.7dB compared  
 566 to SRCNN. Jiang et al. [93] consider a two stage process for video frame super-resolution, where a  
 567 pre-trained CNN with an architecture similar to [16] produced multiple features at the first stage,  
 568 which are then combined through a densely connected network at the second stage. The proposed  
 569 PECNN architecture achieve a minor increase in terms of no-reference quality metrics, like the average  
 570 gradient of the naturalness image quality when evaluated on aerial imagery from the Kaggle dataset,  
 571 however, the method demonstrates a substantial reduction in terms of processing time per image.

572 Satellite video super-resolution is also explored in [94] where a CNN variant termed deep  
 573 distillation recursive network (DDRNN) is introduced, featuring (i) groups of ultra-dense residual blocks  
 574 encoding multiple short and long skip connections without substantially increasing the computation  
 575 and memory requirements, (ii) a multi-scale purification unit for the recovery of high-frequency  
 576 components and (iii) a reconstruction module which corresponds to a shuffling operator for producing  
 577 the final high resolution image. Evaluation on observations from Jilin-1 satellite, as well as from the  
 578 aerial imagery available through the Kaggle Open Source Dataset, demonstrates a gain of more than  
 579 1dB compared to the VDSR [31] method for an upscale factor of 4 when trained on the same dataset.

## 580 2.2. Discussion

581 The number of applications of DNN in remote sensing super-resolution has increased dramatically  
 582 in the past 2-3 years, demonstrating superior performance compared to other state-of-the-art method.  
 583 In Table 1 we provide a high-level grouping of the different methods discussed above.

Methods	Observation type	Approach
[51–53]	Single Image/band	single-scale CNN
[54–57,59]	Single Image/band	multi-scale CNN
[62–64]	Single Image	GAN
[68]	HS/MS	multi-scale CNN (LPN)
[72]	HS/MS	multiple loss CNN
[77]	HS/MS	band-specific CNN
[78]	HS/MS	3D-CNN
[91–93]	Video frames	CNN variants

Table 1. Listing of representative super-resolution approaches

584 In addition to the qualitative grouping of the different super-resolution approaches presented  
 585 in the Table 1, Tables 2 and 3 provide some quantitative performance results for single band and  
 586 multi-band image super-resolution respectively. In both cases, we report the gains offers by each  
 587 method compared to the bicubic interpolation measured in terms of PSNR, for different requested  
 588 super-resolution factors. Regarding the case of single band images and the case of  $\times 4$  super-resolution,  
 589 the reported results indicate that the use of GAN based architectures, pre-trained on generic imagery,  
 590 lead to the higher increase in terms of visual quality. For the case of multispectral imagery, the increase  
 591 in performance is significantly lower, which the results suggest that exploiting both spatial and spectral  
 592 information lead to the best performance.

593 Based on our analysis of the current state-of-the-art, a number of observations can be made  
 594 regarding the problem of remote sensing super-resolution.

- 595 • In terms of single image, single band super-resolution, approaches based on the GAN framework  
 596 appear to achieve the highest quality image estimation. This observation for the case of remote  
 597 sensing is in-line with results for methods applied in generic natural image.
- 598 • For the case of spatial-spectral super-resolution, the best performing approaches pay special  
 599 attention of simultaneously increasing the spatial resolution without introducing unwanted

Method	Dataset	Performance gain		
		$\times 2$	$\times 3$	$\times 4$
LGCNet [54]	UC Merced	+9%	+7%	+5%
WMCNN [56]	RSSCN7	+8%		+3%
[62]	UCMERCED/RSSCN7/NWPU-RESIS45	+9%		+7%
TGAN [63]	UC Merced (airplanes) pre-trained on the DIV2K			+16%
DMCN [55]	UC Merced	+10%	+8%	+7%
MRNN [59]	NWPU-RESISC45			+13%
FRDBPN [60]	UC Merced			+5%
EEGAN [64]	Kaggle	+14%	4%	+13%

**Table 2.** Relative performance gains (with respect to PSNR) for different DNN based *single image* super-resolution methods compared to bicubic interpolation for scale factors , and

Method	Dataset	Performance gain		
		$\times 2$	$\times 3$	$\times 4$
[68]	UC Merced			+2%
MLFN [72]	UC Merced	+11%		
SSIN [77]	UC Merced	+8%	+6%	+4%
3D-FCNN [78]	UC Merced	+9%	+5%	+3%

**Table 3.** Relative performance gains (with respect to mean PSNR) for different DNN based *multispectral* super-resolution methods compared to bicubic interpolation for scale factors , and

artifacts in the spectral domain, unlike early approaches which consider each spectral channel independently.

- Exploiting temporal information, which for many remote sensing cases, can be directly translated to access to multi-view imagery given the regular imaging protocols in satellites imaging.
- Super-resolution of video sequences require the exploitation of both spatial and temporal information, such as the approach proposed in [95] which encodes spatio-temporal information using a recurrent residual network and is applied to the restoration of generic video sequences. This methodology however has not been explored in the context of remote sensing, where methods primarily consider single frames.
- Although numerous approach have been proposed for addressing the combination of two resolution dimensions, methods which can simultaneously enhance remote sensing images along their spatial, spectral and temporal dimensions have not been considered.

The above observation can be utilized in order to outline an exemplary architecture of MS/HS super-resolution shown in Figure 5. This exemplary architecture captures the essence of the most successful paradigms in the area, utilizing 3D convolutions, skip connections over multiple layers, dilated convolution for upsampling and the use of multi-objective loss functions.

### 3. Pan-sharpening

While both high spatial and spectral resolution observations are desired in remote sensing, imaging systems have to sacrifice the resolution in one domain in order to achieve the requested performance in the other. As a result, the majority of spaceborne remote sensing imaging systems acquire two types of observations, high spatial resolution single band panchromatic (PAN) images, or moderate to high spectral resolution multispectral (MS) or hyperspectral (HS) observations at much lower spatial resolution. For instance, the SPOT 6 and 7 satellites acquire PAN imagery at 1.5m spatial resolution and 6m spatial resolution for the 4 band MS, while the WorldView-4<sup>1</sup> acquires PAN imagery

<sup>1</sup> at least before a failed gyroscope rendered it inoperable in January 2019

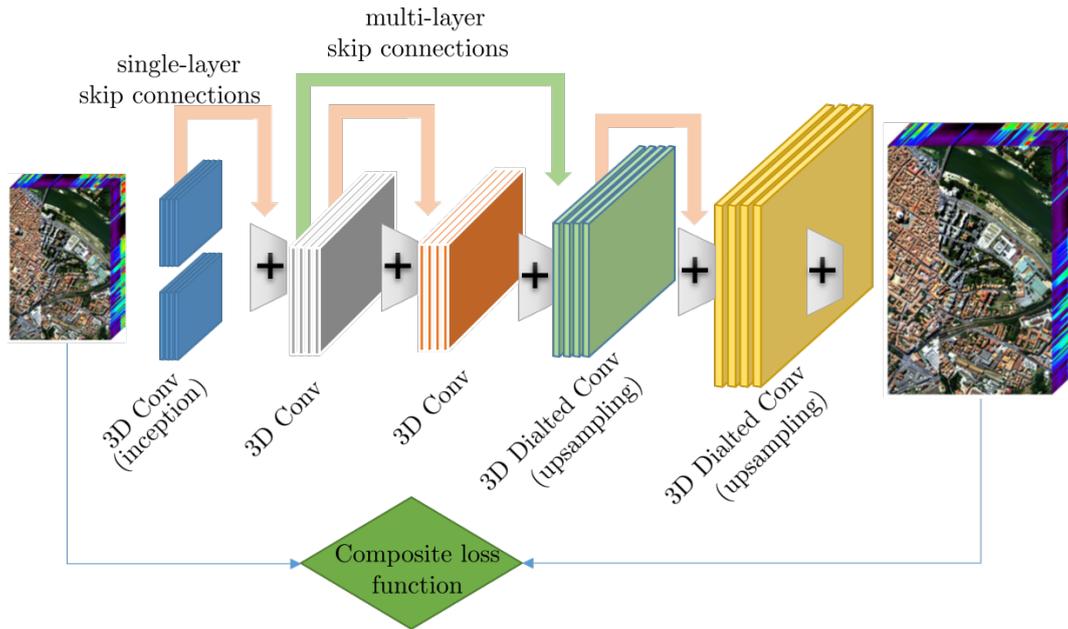


Figure 5. An exemplary CNN architecture for MS/HS super-resolution.

624 at 0.31m and 4 bands MS at 1.24m. Pan-sharpening methods explore how the two types of observations  
 625 can be combined in order to produce high spatio-spectral observations [96]. For this problem, a number  
 626 of DNN approaches have been recently proposed.

### 627 3.1. AE based approaches

628 Huang et al. [97] are among the first to consider DNNs for the problem of pan-sharpening remote  
 629 sensing observation. Specifically, they propose a modified Sparse Denoising Autoencoder (MSDA)  
 630 architecture, aiming at learning the patch-level relationship between low and high spatial resolution  
 631 MS and PAN images. The authors consider a series of pre-trained MSDAs which form a stacked MSDA,  
 632 and subsequently fine-tune the entire network to infer high spatial resolution image patches for each  
 633 MS band, which are then appropriately averaged to produce the final image/cube. The method is  
 634 evaluated on two different satellite data-sets from the IKONOS and QuickBird satellites, obtaining  
 635 higher quality imagery compared to conventional algorithms.

636 In order to extract the most representative features from each source, a multi-level neural network  
 637 approach for the Pan-sharpening problem is proposed in [98]. The authors propose a network  
 638 composed of four individual Stacked Autoencoders (SAE) networks, and specifically a variation  
 639 called Denoising SAE, which are coupled together in order to derive the desired high-resolution  
 640 MS image. Concerning the role that each of the separate network play, two of them represent the  
 641 high-resolution PAN image and the up-sampled MS image in a new feature space via SAE, while  
 642 the subsequent third one builds the connection (i.e. mapping) between these two representations.  
 643 Finally, a fine-tuning fourth network is introduced at the highest level in order to ameliorate the  
 644 learned parameters of the networks via joint optimization through back-propagation. The proposed  
 645 method is evaluated on the QuickBird dataset, and managed to outperform not only conventional  
 646 pan-sharpening techniques but deep learning based ones (i.e. [99]) as well.

647 A multi-resolution analysis framework for confronting the Pan-sharpening problem is proposed  
 648 in [100] where a two-stage DNN scheme based on the SAE framework is investigated. In the first  
 649 stage a low-resolution PAN image is constructed by its high-resolution counterpart via the help of  
 650 two autoencoders (the second one is trained with input the features derived by the first one). In the  
 651 second stage the desired high-resolution MS image is reconstructed via its low-resolution version  
 652 according to the respective relationship between PAN images. Once both AE are trained alone, the

653 whole network (i.e. both of them together) is fine-tuned in order to optimize its hyper-parameters. The  
654 proposed approach is tested in two separate data-sets (namely, QuickBird and WorldView-3), where it  
655 outperformed its competitive algorithms. A possible shortcoming of the method could be considered  
656 the absence of comparison with other deep learning-based methods, which seem to be state-of-the-art  
657 in the field the last years.

658 A recently approach utilizing AE for pan-sharpening is proposed by Azarang et al. [101] where  
659 the author consider a convolutional AE architecture, composed of an encoding and the associated  
660 decoding stage from improving the spatial information of the low resolution MS bands. The objective  
661 is achieved by learning the nonlinear relationship between a PAN image and its spatially degraded  
662 version at a patch level, and utilizes the trained model to increase the spatial resolution of each MS  
663 band independently. Evaluation on data from QuickBird, Pleiades and GeoEye-1 demonstrate some  
664 gains compared to other non-DL based methods.

### 665 3.2. CNN based approaches

666 In addition to AE, CNNs have also been considered for pan-sharpening. One of the first attempts  
667 to employ CNNs in pan-sharpening is reported in [102] where the authors consider a CNN architecture  
668 based on the SRCNN [34] architecture for increasing the spatial resolution of MS observations in a  
669 band-by-band fashion and then use the spatially enhanced band corresponding to the PAN image for  
670 adjusting the original PAN image before subsequently applying a Gram–Schmidt transform to fuse  
671 the MS with the PAN observations. Training of the model is performed on natural images from the  
672 2013 version of ImageNet dataset while for the validation, observations from QuickBird satellite are  
673 employed, demonstrating higher quality estimation compared to competing non-DL methods.

674 Another early attempt in using CNN for pan-sharpening is explored in [76] which also borrows  
675 ideas from the SRCNN super-resolution architecture and extends them to tackle the pan-sharpening  
676 problem. The key novelty of this work involves the upsampling of the low spatial resolution 4 band  
677 MS observations to the resolution of the PAN image and stacking both observations into a 5 component  
678 input. The authors further propose a non-linear radiometric parameter (typical in remote sensing  
679 applications) injection as a surplus input to the network, which also contributes to an increase in the  
680 output image quality. The performance of the proposed approach, termed PNN, is validated using 3  
681 different datasets from the IKONOS, GeoEye-1, and WorldView 2 satellites, where it outperforms the  
682 competitive non-DL algorithms with respect to various full-reference and no-reference metrics.

683 Another early attempt in the application of CNNs for par-sharpening is proposed by Li et al. [103]  
684 who, inspired by the VDSR architecture [31], propose the introduction of convolution/deconvolution  
685 layers and residual connections. The proposed method, called DCNN, is capable of estimating the full  
686 resolution output, not by averaging different patches but n a band/channel-wise fashion, achieving  
687 promising results with respect to the competing methods in two different data-sets from the QuickBird  
688 and GeoEye-1 satellites.

689 Wei et al. [99], [104] also propose a residual CNN architecture, in order to evaluate the performance  
690 of deeper networks as opposed to shallow architectures. To overcome the problem of vanishing  
691 gradient characterizing very deep architectures, the proposed approach, called DRPNN, considers  
692 a residual architecture with skip connections, inspired by the VDSR method for super-resolution,  
693 allowing the network to accept low resolution images in the input, i.e., without the need for a  
694 upsampling pre-processing step. The method is tested in two separate data-sets, namely QuickBird  
695 and WorldView-2, and demonstrates significant gains compared to traditional non-DL approaches, as  
696 well as well as a CNN one [76]. A similar residual CNN approach is also considered in [105] which  
697 is shown to outperform conventional algorithms, as well as the CNN based method [76] in terms of  
698 performance on observations from the LANDSAT 7 ETM+.

699 In a similar vein, i.e., using residual connections for training deeper CNNs, the authors in [106]  
700 consider an architecture based on the ResNet [30] architecture. The proposed scheme incorporates  
701 domain specific knowledge in order to preserve both spatial and spectral information, by adding

702 up-sampled MS images to the output and training the network in the high-frequency domain.  
703 Experimental results demonstrate the superiority of the proposed method compared to conventional  
704 pan-sharpening methods as well as CNN based [76]) in the WorldView-3 dataset, as well as greater  
705 generalization capabilities in new datasets, namely, WorldView-2 and WorldView-3, where the  
706 proposed approach is capable of producing high quality results without the need to be retrained.

707 Another CNN approach for pan-sharpening is proposed by [107,108] where the authors introduce  
708 a multi-scale feature extraction by considering three different size convolutional filter at each layer to  
709 build a respective multi-scale network, inspired by the inception architecture [33]. By concatenating  
710 the derived features across the spectral dimension, they retained as much spatial information as  
711 possible, while minimizing spectral distortion. To build a deeper network architecture, they adopt  
712 residual connection approach leading to sparse and informative features. The proposed architecture is  
713 tested in two different datasets from the QuickBird and WorldView-2 satellites, where it outperformed  
714 competitive approaches, including the CNN baseline [76], by 0.7 and 0.3dB respectively.

715 Building upon the method proposed in [76], which became a baseline for most of the  
716 aforementioned deep learning-based approaches for Pan-sharpening, the authors propose an updated  
717 approach in [109]. More specifically, they modified their initial approach by taking into account  
718 research that had taken place after their initial publication, by adopting the following changes: (i)  
719 use of  $\ell_1$  loss function-instead of  $\ell_2$ ; (ii) processing and learning in the image residual domain-rather  
720 than in the raw image one; (iii) use of deeper network architectures (in the [76], the network consisted  
721 of only three convolutional layers). The aforementioned modifications are tested in the GeoEye-1  
722 dataset and compared to the initial setup proposed in [76], and proved to be quite beneficial in terms  
723 of performance as well as training time reduction.

724 Capitalizing on these observations, the authors extended the approach in [110] by making it  
725 adaptive to every new image that needed to be pan-sharpened via the trained network. For doing so,  
726 every new image fed a fine-tuning process on the pre-trained deep neural network (until convergence),  
727 and right afterwards it entered the trained model in order to be pan-sharpened at its output. To  
728 prove the validity of their -updated- approach, the authors tested it in four different datasets, namely  
729 IKONOS, GeoEye-1, WorldView-2 and WorldView-3, and compared it with their initial approach ([76])  
730 as well as with conventional pan-sharpening techniques. The reported results in all data-sets verified  
731 the merits of the new approach both in terms of performance metrics as well as of visual inspection.

732 Another approach involving CNNs for pan-sharpening is the RSIFNN architecture [111] which  
733 poses the problem as a fusion of PAN and MS inputs. Specifically, the RSIFNN introduces a two-stream  
734 fusion paradigm such that features from the PAN and MS inputs are first independently extracted and  
735 subsequently fused in order to produce high spatio-spectral observations, while a residual connection  
736 between the low resolution MS and the outputs is also introduced. The authors considered PAN  
737 and MS (4 bands) observations the from QuickBird and Gaofen-1 satellites and generate synthetic  
738 examples through appropriate downsampling (Wald's protocol) for comparing the performance of  
739 the proposed and various state-of-the-art methods, including the CNN-based method by Massi et al.  
740 [76] as well as the SRCNN method. The authors report a gain of around 0.3dB compared to [76] and  
741 significant gains using additional metrics likes SAM, while also demonstrating a dramatic decrease  
742 ( $\times 6$ ) in running time for inference compared to other CNN based approaches. In a similar spirit, Zhang  
743 et al. [112] consider two stream bidirectional Pyramid networks which is able to extract and infuse  
744 spatial information to MS from PAN imagery along multiple scales. Validation using GF2, IKONOS,  
745 QuickBird, and WorldView3 observations demonstrate some minor improvements in performance.

746 Wei Yao et al. [113] propose a pixel-level regression for encoding the relationship between the  
747 pixel values in the low and high resolution observation. For that case, they employed a CNN inspired  
748 by the U-Net architecture [35], but refined with respect to the problem at hand and trained with image  
749 patches. The proposed method obtained quite improved performance with respect to localization  
750 accuracy as well as feature semantic levels, compared to other DNN based methods. At the same time,  
751 the proposed approach outperforms other sophisticated Pan-sharpening methods in two different

752 remote sensing datasets in various indices of quality, leading to sharpened visual results both in terms  
 753 of spatial details as well as spectral characteristics.

754 Increasing the spatial resolution of MS observation using PAN imagery using CNNs is also  
 755 explored in [114] where a novel loss function penalizing the discrepancy in the spectral domain  
 756 is proposed. The proposed scheme employs two networks, the first one called 'Fusion Network'  
 757 which introduces high spatial resolution information to upsampled MS images and the second one,  
 758 the 'Spectral compensation' network which tries to minimize spectral distortions. Validation on  
 759 observations from Pleiades, Worldview-2 and GeoEye-1 demonstrate the the proposed scheme achieves  
 760 superior performance compared to [76] and [99] in almost all cases.

761 In the recently proposed method by Guo et al. [115], the authors utilize the latest CNN design  
 762 innovations including single and multi-scale dilated convolutions, weight regularization (weight  
 763 decay) and residual connections for pan-sharpening Worldview-2 and IKONOS imagery. Compared  
 764 to state-of-the-art methods, including [76] and [106], the proposed method achieved both superior  
 765 performance in terms of quality (more than 0.5 in SAM compared to competing CNN approaches) and  
 766 computational requirements for inference.

767 Building on the disruptive GAN framework, the authors in [116] proposed a two-stream CNN  
 768 architecture (one stream for the PAN images and the other one for the MS ones) as the generator  
 769 of high-quality pan-sharpened images, accompanied by a fully-convolutional discriminator. In  
 770 this way, the fusion of images is performed in the feature level (of the generator), rather than in  
 771 pixel-level, reducing in this way the spectral distortion. The aforementioned approach is tested on two  
 772 different datasets (namely, QuickBird and GaoFen-1), where it outperformed conventional and deep  
 773 learning-based ([76]) pan-sharpening methods, even if only the generator component is used solely. A  
 774 similar idea is also explored in [117] where additional contains related to spectral features are imposed.

### 775 3.3. Discussion

method	Approach
[97,98,100,101]	AE variants
[76,102]	CNN based on SRCNN architecture
[99,103–106]	CNN with residual connections
[107,108]	Inception-like CNN for multi-scale feature extraction
[110]	Target-adapted CNN
[111,112]	Two-stream CNN architecture
[113]	CNN-based pixel-level regression over multiple scales
[114]	Two CNNs associated with spectral and spectral dimension
[116]	Two-stream GAN-CNN architecture

**Table 4.** Listing of representative pan-sharpening approaches

776 In order to offer some insights into the performance of each method, Table XX report the gains  
 777 of each method, measured in terms of Spectral Angle Mapper (SAM), compared a state-of-the-art  
 778 non-DNN method of MTF-GLP [118] and one the first DNN based approach, the PNN [76]. The  
 779 results suggest that currently, the highest achieving DNN method is [115], based on its performance  
 780 on pan-sharpening WorldView-2 observations. This observation suggest that utilizing all available  
 781 machinery during the realization of a CNN architecture offers significant performance gains.

782 A quick overview of the major approaches related to pan-sharpening remote sensing observation  
 783 using DNN methods is presented in Table 4. Since in all cases the input is a PAN and an MS image, the  
 784 table reports the different approaches groups along the specific DNN methodology employed. Based  
 785 on the analysis of existing literature, a number of observations can be made

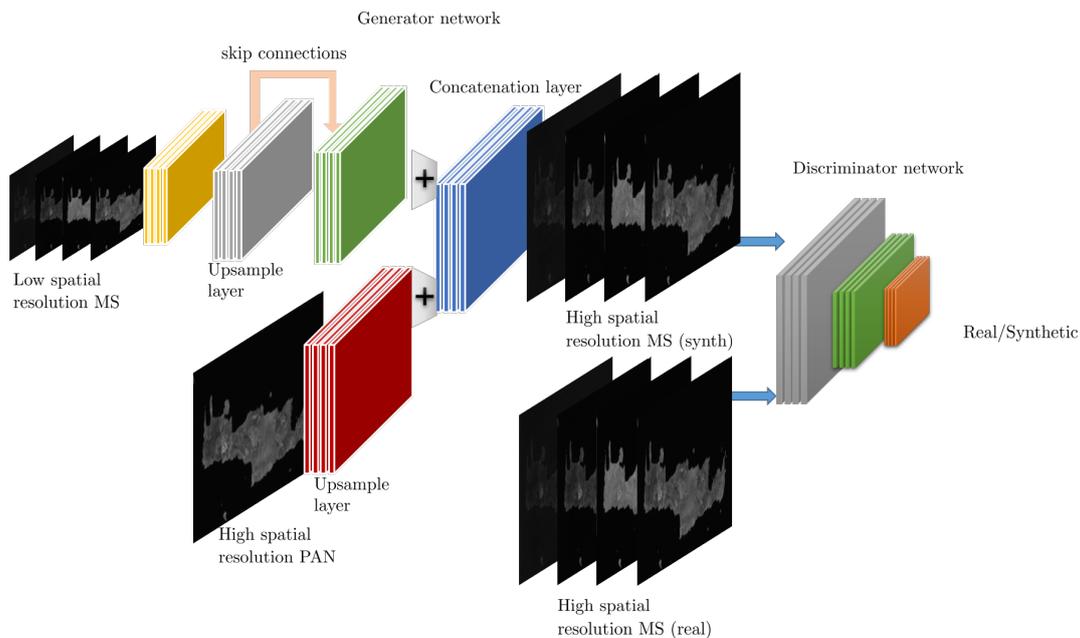
- 786 • The majority of state-of-the-art performing methods rely on some variant of CNNs, typically  
 787 employing fully convolutional architectures, while AE approaches are significantly less utilized  
 788 for this problem.

Method	Dataset	Baseline approaches	
		MTF-GLP	PNN
PUNET [113][110]	Ikonos / WorldView-2	-	+15% / +16%
DRPNN [99]	QuickBird / WorldView-2	+35% / +25%	+8% / +5%
MSDCNN [107]	QuickBird / WorldView-2	+35% / +26%	+7% / +5%
RSIFNN [111]	Gaofen-1 / QuickBird	-	+11% / +8%
L1-RL-FT [110]	WorldView-2 / WorldView-3	+38% / +29%	+14% / +37%
BDPN [112]	QuickBird / WorldView-3	-	+25% / + 11%
[115]	WorldView-2	-	+33%
[114]	Pleiades / WorldView-2	+50% / +26%	+46% / +20%
PNN [76]	WorldView-2 / Ikonos	+25% / +21%	-
PanNet [106]	WorldView-3	-	+15%
PSGAN [116]	QuickBird / GF-1	-	+18% / +35%

**Table 5.** Relative performance gains (with respect to SAM) for different DNN pan-sharpening methods compared to non-DNN MTF-GLP [118] method and baseline DNN based PNN [76] method

- 789
- The introduction of successful CNN components like inception modules and residual/skip connections lead to higher quality outputs.
  - Among the best performing methods are the ones that employ two-stream architectures in which cases the PAN and MS are first analyzed independently by two separate paths and the extracted features are then jointly processed in order to produce the output.
  - GANs has also been introduced in the problem of pan-sharpening with some promising initial results, however, more detailed analysis is required in order to justify any potential benefits.
- 790
- 791
- 792
- 793
- 794
- 795

796 Given the analysis of existing methods, Figure 6 presents a graphical illustration of an exemplary  
 797 GAN architecture for pan-sharpening employing CNNs for both the generator and the discriminator  
 798 network. This exemplary architecture captures the latest and most successfully architecture design  
 799 considerations including the notion of skip connections, the use of upsampling layers using dilated  
 800 convolutions, all the context of a GAN paradigm which has been shown to lead to the highest quality  
 801 estimations.



**Figure 6.** An exemplary two-stream GAN-CNN architecture for pan-sharpening.

## 802 4. Restoration

### 803 4.1. Denoising and deblurring

804 Although MS and HS observations collected from remote sensing platforms provide valuable  
805 information, they are often inevitably contaminated with noise injected during acquisition and  
806 processing like compression, which can jeopardize the performance of subsequent analysis.  
807 Observation denoising is thus a crucial process which has been extensively studied by the remote  
808 sensing community [119], while in recent years, several DNN architectures have been proposed for  
809 observation denoising.

810 In [120], the authors introduce a CNN with trainable non-linear activation functions for HS  
811 denoising, where both the convolutional kernels' weights as well as the non-linear functions are  
812 learned from overlapping clean-noisy HS patch-pairs serving as training examples. In this work,  
813 principal component analysis is initially applied in order to extract the major components from the  
814 highly correlated spectra, which are subsequently used for denoising the remaining components. Then,  
815 the major components and the remaining denoised components are inversely transformed in order to  
816 obtain a denoised HS. Compared with the state-of-the-art HS denoising methods, the experimental  
817 results on both synthetic and real remote sensing observations from the AVIRIS sensor (Salinas valley)  
818 demonstrate that the proposed method can achieve faster and more effective performance, at the cost,  
819 however, of not preserving the spectral information.

820 A HS denoising method which can properly remove spatial noise while accurately preserving  
821 the spectral information is proposed in [121], where a residual CNN with batch normalization is  
822 employed for learning the appropriate mapping between the noisy and the clean HS observations. In  
823 the denoising stage of the method, the learned mapping is adopted to reconstruct the clean spectral  
824 difference. Meanwhile, a key band is selected based on a principal component transformation matrix  
825 and denoised, acting as a starting point for the correction of other spectral bands. Comparative  
826 analyses on observations acquired by the ROSIS and the the Moon Mineralogy Mapper (M3) acquired  
827 over the Aristarchus Plateau demonstrated that the method achieved higher performance compared to  
828 non-DL method, due to the fact that textures are retained in both large and small scales. The authors  
829 in [122] further extend the idea of denoising key bands similar to [121] while also utilizing trainable  
830 non-linearity functions introduced in [120] for denoising the remaining non-key bands. Experiments  
831 on HYDICE imagery demonstrate that this model enjoys superior performance in spatial recovery and  
832 spectral preservation compared to other non-DNN approaches.

833 Unlike traditional CNNs that are limited by the requirements of the training sample size, Chen  
834 et al. in [123] introduce an efficient deep CNN model for the denoising of aerial images, which is  
835 appropriate for small training datasets. The learning algorithm of this model employs a multi-scale  
836 residual-learning approach with batch normalization and dropout, which not only speeds up the  
837 training process but also improves the denoising performance of the model. Experimental results  
838 on images of the Caltech database showed that the proposed method achieves better performance  
839 compared to state-of-the-art denoising methods using small training datasets, while the comparison  
840 also included deep learning based methods like PNN [76] and SRCNN [34], which however, are not  
841 designed for image denoising.

842 Another DNN method for HS denoising is proposed in [124], by learning a non-linear end-to-end  
843 mapping between the noisy and clean signals with a deep spatial-spectral CNN. More specifically, a  
844 2D-CNN is used to enhance the feature extraction ability of single bands, and a 3D-CNN is employed  
845 to simultaneously utilize the correlation between adjacent bands, offering greater flexibility with  
846 respect to the available bands. To capture both the multi-scale spatial and spectral features which  
847 characterize remote sensing imagery, different convolutional kernel sizes are employed. Experimental  
848 results from the Indian Pines and Pavia University datasets demonstrate that this method outperforms  
849 many of the mainstream non-DNN methods for HS denoising. A 3D CNN based architecture is also  
850 explored in [125] for HS image denoising where the notion of atrous convolution is considered, a term

851 which has being replaced by dilated convolution, which is able to expand the receptive field associated  
852 with each convolution kernel without increasing the number of parameters. The atrous convolutions,  
853 along with the introduction of skip connections and multi-scale inception like filters, demonstrated an  
854 increase in quality, in the order to 0.5dB for some cases, compared to other non-DNN approaches.

855 While the previously reported approaches consider only one type of noise, typically Gaussian, in  
856 many real life situation, especially in the case of HS imaging, one has to deal with multiple types of  
857 noise, i.e., hybrid noise. The use of DNN for addressing hybrid noise in HS is considered by Zhang et al.  
858 [126] by simultaneously exploiting both spatial and spectral information through learning of spatial and  
859 spectral gradients at multiple scales. The proposed SSGN architecture further utilizes a spatial-spectral  
860 loss function for reducing distortions. Experiments performed on HS imagery from HYDICE, AVIRIS  
861 and Hyperion EO-1, confirmed that this method performs better at mixed noise cases (Gaussian noise +  
862 stripe noise + deal lines) than state-of-the-art non-DNN denoising algorithms. Addressing the removal  
863 of mixed noise types in HS is also considered by Chang et al. [127] where the proposed HSI-DeNet is  
864 a GAN based approach which directly operates on the 3D data and incorporates high performance  
865 machinery like residual learning and dilated convolutions. Experiments on HYDICE images indicate  
866 that this method outperforms the state-of-the-art in terms of both speed and performance. However,  
867 larger training datasets are needed to improve the generalization of this model to accommodate all  
868 kinds of complex noise categories and HS data.

869 A GAN based methodology is also explored in [128] in order to deblur degraded remote sensing  
870 images in the context of image restoration. A common approach to solve this problem is the  
871 incorporation of various priors into a restoration procedure as constrained conditions that often lead  
872 to inaccurate results. In contrary, in [128] an end-to-end kernel-free blind deblurring learning method  
873 is presented that does not need any prior assumptions for the blurs. During training, a discriminator  
874 network is defined, with a gradient penalty which is shown to be robust to the choice of generator  
875 architecture, and a perceptual loss which is a simple  $\ell_2$  loss based on the difference of the generated  
876 and target image CNN feature maps that focuses on restoring general content instead of texture details.  
877 This method can handle blur caused by camera shake and object movement, using fewer parameters  
878 compared to multi-scale CNN. The experimental results on the Google map dataset showed that  
879 this method can obtain competitive results compared to other state-of-the-art methods. However,  
880 checkerboard artifacts still exist in some deblurred images, produced by the deconvolution-transposed  
881 convolution operation.

882 Finally, in [129], a novel deep memory-connected DNN with a large receptive field is introduced  
883 for remote sensing image restoration. Inspired by neuroscience concepts, the authors therein propose  
884 local and global memory connections to combine image detail information in lower layers with  
885 global information in higher layers of the network. In addition, to alleviate the computational  
886 burden and accelerate the training process, down-sampling units are proposed to shrink the spatial  
887 size of the feature map. Experiments on three remote sensing datasets, namely from UC Merced,  
888 NWPU-RESISC45 and GaoFen-1, indicate that this method yields promising improvements and better  
889 visual performance over the current state-of-the-art, achieving high-quality image recovery for both  
890 known and unknown noise level at the same time.

#### 891 4.2. Missing data recovery

892 Remote sensing observations are often hindered by missing measurements, which is attributed to  
893 poor atmospheric conditions, such as thick clouds, as well as internal satellite sensors malfunction  
894 such as dead pixels and scan lines. Thus, an issue of major importance is the estimation of the missing  
895 observations, also known as gap filling when considering time-series of observations, by properly  
896 exploiting the spatial, temporal and spectral information [130]. Most of the earlier techniques (such as  
897 interpolation) deal with each domain separately, losing in that way crucial information, as most of them  
898 take into account linear correlations between the data. On the other hand, the recent breakthrough of

899 the deep CNNs methods gives researchers the opportunity not only to reconstruct better the missing  
900 information but also incorporates knowledge from all aspects in order to address this challenge

901 In [131] a CNN based approach which exploits information under a unified  
902 spatial–temporal–spectral model is presented for missing measurements recovery, focusing  
903 on three representative cases, namely removing dead lines from the 6<sup>th</sup> band of Aqua MODIS  
904 instrument, scan lines in the Landsat Enhanced Thematic Mapper Plus and thick clouds. The proposed  
905 architecture consider es two sources of data, that is spatial data with missing values and auxiliary  
906 spectral or temporal data without missing values, where initial features are independently extracted  
907 and subsequently fused in deeper layers of the CNN. The proposed CNN is augmented with a number  
908 of enhancements including multiscale feature extraction using dilated convolutions and residual  
909 learning using skip connections. Experimental results under various conditions using both simulated  
910 and real, as well as additional sources of noise like registration error, demonstrate the very high  
911 performance and robustness of the proposed method.

912 A critical aspect of missing observation recovery, especially when temporal analysis is considered,  
913 is the preservation of causality. In their work, Das et al. [132] consider a deep learning based approach  
914 for reconstructing missing time series data, which relies on making use of observations from both  
915 preceding and subsequent time-instances while at the same time maintaining the causality constraint.  
916 The proposed network architecture comprises of four different modules, namely feature preparation,  
917 feature learning, prediction and data tuning, and specifically, which are organized in an ensemble of  
918 forward prediction DNNs based on the Deep-STEP architecture [133], each of which aims at predicting  
919 the immediate next missing observation in a sequence given previous time instances. Experiments  
920 performed using products (NDVI) derived from the Landsat-7 TM-5 satellite over the Barddhaman  
921 district, West Beghal-India, show clear superiority of the proposed approach over six competing  
922 algorithms under a number of quantitative metrics (i.e. MAE, RMSE, PSNR and SSIM).

923 A GAN-based method for estimating arbitrary shaped occluded regions of sea surface temperature  
924 images by exploiting historical observations is proposed in [134]. The authors consider the  
925 deep convolutional generative adversarial network [135] in order to address the unavailability of  
926 observations due to cloud occlusion, where the generator network is tasked with producing realistic  
927 estimation while the discriminator network must classify the inputs as real or synthesized. Estimating  
928 of missing regions is subsequently carried out by trying to estimating the closest vector representation  
929 of the uncorrupted image through the minimization of a loss function encapsulating the  $\ell_2$  norm on  
930 the available observation, the adversarial error and the deviation of a average values. The results  
931 demonstrate a significant gain offered by the proposed approach compared to competing methods  
932 including the deep learning method in [136] which is designed for generic imagery.

933 Conditional GANs are also considered in [137] for estimation of missing multispectral imagery  
934 by combining optical and radar multitemporal observations. The proposed methodology involves  
935 introducing pairs of registered optical and SAR patches from multiple dates for training the generation  
936 which is then tasked to produce an optical patch, given the same data SAR and optical/SAR from a  
937 different day, while the discriminator is trained to judge is such synthesized data is realistic or not.  
938 The performance of the scheme is validated using optical data from Landat 8 OLI instrument and  
939 Sentinel 2 MSI instrument and SAR data for the Sentinel 1 SAR over regions of diverse land cover  
940 types which as used as input to a Random Forest classifier. The experimental results demonstrate that  
941 using synthesized imagery in addition to using the available SAR data can lead to almost the same  
942 accuracy as the case where both types of observation is available Similarly to [137], in [138] conditional  
943 GANs are also considered for missing observations recovery using multitemporal optical and radar  
944 (SAR) observations from Sentinels 2 and 1 respectively.

945 **4.3. Discussion**

946 An overview of exemplary applications of DNN in the problem of remote sensing observation  
 947 restoration, which includes the removal of different types of noise and the estimation of missing  
 948 observations, is presented in Table 6.

Method	Observation type	Approach
[120]	HSI + Gaussian noise	CNN with trainable activation function
[121,122]	HSI + Gaussian noise	Residual CNN with key band selection
[123]	Aerial RGB + Gaussian noise	Multi-scale CNN for learning residual noise component
[124]	HSI + Gaussian noise	Multi-scale 3D CNN using neighboring bands
[125]	HSI + Gaussian/Poisson noise	Multi-scale 3D CNN w/ atrous layers
[127]	HSI + mixed noise	2D CNN with residual connection
[127,128]	Noise & Blur	GAN architecture
[131]	Missing spatial measurements	CNN based fusion using auxiliary observations
[132,133]	Missing temporal measurements	CNN for recovery via temporal prediction
[134]	Cloud occlusions	GAN
[137]	Missing MS observations	cGAN using multitemporal MS and SAR observation

**Table 6.** Major approaches in DNN based remote sensing observation restoration.

949 A quantified comparison of different denoising methods is presented in Table 7, where gains  
 950 compared to the state-of-the-art non-DNN BM4D [139] method are presented in terms of PSNR,  
 951 averaged over the bands. The results indicate that the GAN based method reported in [121] offers the  
 952 highest performance gains, while being able to handle different types of noise.

Method	Dataset	comparison to BM4D [139]
HDnTSDL [122]	Washington DC Mall (HYDICE)	+12% / +10% / +13% (Gaussian w/ $\sigma=15/25/50$ )
SSDRN [121]	PAVIA (ROSIS)	+2% / +3% / +4% (Gaussian w/ $\sigma=15/25/50$ )
SSGN [126]	Pavia (ROSIS), Washington DC Mall (HYDICE), Indian Pines (AVIRIS)	+20% / +17% (Gaussian, Gaussian+stripe)
[124]	Washington DC Mall (HYDICE)	+6% / +8% / +11% (Gaussian w/ $\sigma=25/50/100$ )
[125]	Pavia University (ROSIS), Indian Pines (AVIRIS)	+1% / +4% (Gaussian w/ $\sigma=10$ )

**Table 7.** Relative performance gains (with respect to mean PSNR) for different DNN denoising methods compared to non-DNN BM4D [139] method

- 953 • Highest performing methods for restoration involve CNN architectures with multi-scale feature  
 954 extraction capabilities, while GANs architectures also appear promising.  
 955 • While the case of denoising with known noise characteristics has been explored, further research  
 956 involving unknown or mixed distribution noise is required.  
 957 • The utilization of multiple sources of observations with different quality characteristics targeting  
 958 restoration is also another topics of great importance.  
 959 • Estimating missing observation, especially the cases involving cloud occlusion of optical  
 960 imagery, using radar observations which are impervious to this problem also warrants further  
 961 investigation.

962 **5. Fusion**

963 Fusion refers to the simultaneous encoding of observations from multiple modalities and/or  
 964 sensors with diverse characteristics, targeting the generation of higher quality or new types of

965 observations [140]. The framework of fusion is closely related to the pan-sharpening one, however,  
966 it considers more generic cases such as the joint analysis of high-spectral-low-spatial HS resolution  
967 observations with low-spectral-high-spatial MS ones, or observation from optical and radar sensors.

### 968 5.1. Fusion of MS and HS observations

969 Enhancing the spatial resolution of HS observations through fusion with high spatial resolution  
970 MS measurements is a major line of research, since HS and MS instruments are characterized by  
971 different imaging capabilities due to the inherent trade-offs in the system design [141]. Palsson et al.  
972 [142] were among the first to consider a 3D CNN for increasing the spatial resolution HS observation,  
973 by exploiting the high spatial resolution of registered MS observations. The method operates on image  
974 patches and involves the singular value decomposition to the HS observations so that only the spectral  
975 loadings are super-resolved, while preserving the spectral singular values. To train and validate the  
976 performance of the CNN, HS observations from the ROSIS Pavia center data set are appropriately  
977 downsampled. Experimental comparison with a wavelet-based approach demonstrate that higher quality  
978 reconstruction is possible, even when the input observations are affected by noise.

979 Yang et al. [143] recently proposed a two-branch CNN architecture where features from the HS  
980 and the MS are initially extracted independently and then are fused together, effectively enhancing the  
981 spectral dimension of MS by transferring the spectral features from the HS. The overall architecture  
982 involves sequences of convolutional layers, 1D for the HS and 2D for the MS, and the subsequent  
983 concatenation of their outputs which are used as input to a sequence of fully connected layers producing  
984 the super-resolved spectrum for each spatial location. The performance is quantified using synthetic  
985 observations from the AVIRIS dataset and EnMAP, as well as real observations captured by Sentinel 2  
986 for the MS and the EO-1 Hyperion sensor for HS.

987 The use of CNN for MS and HS image fusion is further explored in [144] where two sub-networks,  
988 an encoder and a pyramid fully convolutional decoder network, are considered for progressively  
989 reconstructing the high spatial resolution HS image. The proposed architecture seeks to extract features  
990 from multiple scales by introducing different spatial resolution HR imagery during the estimation  
991 process, while to further improve the reconstruction quality, the typical  $\ell_2$  loss is augmented with a  
992 term focusing on the gradient difference for generating enhanced details. Experimental validation of  
993 the method using simulated observations produced by real Hyperion, HYDICE and ROSIS sensor data  
994 demonstrate significant improvements as high as 1dB for Hyperion, 2.5dB for HYDICE and 1.5dB for  
995 ROSIS compared to [142] which achieve the second best performance. Furthermore, the results also  
996 demonstrate that the introduction of a more appropriate loss function always have a positive effect in  
997 reconstruction quality.

998 In addition to the fusion of observations with different spatio-spectral characteristics,  
999 fusion is also considered for generating high spatio-temporal resolution observations [145].  
1000 The proposed architecture is composed of three parts, the spatial resolution enhancement of  
1001 high-temporal-low-spatial resolution imagery, the extraction of high-frequency components from  
1002 low-temporal-high-spatial resolution observations and the final fusion of the extracted features.  
1003 Effectively, the architecture is trained by considering pairs of different resolution observations from  
1004 the same date and using the learned model to enhance the spatial resolution of a low-spatial resolution  
1005 observation from a future date. The method is validated using daily observations at 500m spatial  
1006 resolution observation from MODIS and 30m resolution from Landsat 8 OLI sensor which has a 16 day  
1007 revisit cycle, operating on single bands.

1008 Building on the fusion context, the authors of [146] propose a blind image quality  
1009 assessment-flavoured model for deriving high-resolution short-wave infrared (SWIR) images, via  
1010 the help of pan-sharpening and hyper-sharpening techniques. More precisely, they build four  
1011 different architectures which combine PAN and SWIR images in order to derive very high-resolution  
1012 images obtained from the WorldView-3 satellite dataset. These approaches (combining widely used  
1013 pan-sharpening and hyper-sharpening algorithms) differ from each other to the "nature" of the

1014 employed fusion process (i.e. sequential, parallel, sequential-parallel, parallel-sequential), and are  
1015 evaluated with respect to a new image quality assessment measure which weights the level of spectral  
1016 and spatial distortions (capitalizing on the famous SSIM and Natural Image Quality Evaluator (NIQE)  
1017 image processing performance metrics). The four above approaches (each one comprising of 10  
1018 different algorithmic schemes for the pan-sharpening) are compared in terms of the aforementioned  
1019 quality measure, as well as of the required computational time, ending up with quite promising results  
1020 in the context of blind assessment of pan-sharpened images in the case of no-overlapping between the  
1021 SWIR bands and the PAN one.

## 1022 5.2. Fusion of spectral and radar observations

1023 Extending the majority of fusion approaches which consider MS and HS observations, a more  
1024 generic fusion framework which explores the fusion of optical sequences, synthetic aperture radar  
1025 (SAR) sequences and digital elevation model is proposed by Scarpa et al. [147] in an effort to estimate  
1026 missing optical features, typically due to cloud coverage. The proposed CNN architecture accepts very  
1027 small patches of cloud-free co-registered optical and SAR images from the Sentinel 2 and 1 platforms  
1028 respectively and is tasked at estimating biophysical parameters like then normalized difference  
1029 vegetation index (NDVI). To validate the performance, the quality of the NDVI estimation is quantified  
1030 over an agricultural region of Burkina Faso, West Africa, and over an extended period of time. Results  
1031 demonstrate both that high-quality estimation of spectral features is possible from radar data and that  
1032 CNN-based architectures can handle large temporal gaps between observations.

1033 A key challenge when trying to fusion observations from multiple sources is the registration or  
1034 equivalently the identification of corresponding patches, between different modalities. To that end,  
1035 the fusion of optical and radar (SAR) observations is explored in [148] where a so-called "siamese"  
1036 CNN network architecture is employed for predicting if patches from the two sources are a match  
1037 or not. To validate the performance of the system, the authors employ an automated framework  
1038 called "SARptical" [149] which is able to extract 3D point clouds from optical and SAR observation and  
1039 perform the matching in this highly accurate 3D point cloud space. Despite the significant challenges  
1040 due to the difference in viewing angles and the imaging protocol employed by SAR systems, the results  
1041 that the quality of matching is comparable to state-of-the-art hand crafted features.

1042 In [150], Quan et al. consider the use of GANs for generating examples for registration of optical  
1043 and SAR imaging data. Since both optical and SAR observations of a scene may not even be available,  
1044 traditional data augmentation techniques for deep learning-based models -employing geometric  
1045 transformations for generating abundant training data- are out of the question for the problem at hand.  
1046 To overcome such an obstacle, the authors propose the use of a GAN which is charged with the task  
1047 of producing coupled training data with respect to its input (i.e. derive pseudo SAR data when fed  
1048 with optical data, and vice-versa). Subsequently, a CNN is employed in order to infer the labels of  
1049 the generated multi-modal image data (in a patch-wise level). At the test phase, the trained model  
1050 in conjunction with three different kind of constraints that need to be met, is used for predicting the  
1051 labels of the test images and the transform matrix alongside with the registered images as well. The  
1052 proposed method is tested with several registration methods in different multi-modal image data  
1053 (i.e. optical-and-SAR, optical-and-LIDAR), outperforming them in several qualitative and machine  
1054 learning based measures.

1055 The generation of SAR observations which although highly realistic do not match corresponding  
1056 optical examples is explored in [151], in an effort to generate challenging negative examples which can  
1057 lead to higher matching accuracy and lower false positive rates between SAR and optical observations  
1058 presented in [148]. The Conditional GAN framework is also explored in [152] for the improvement of  
1059 the geolocalization accuracy of optical imagery through the automatic extraction of ground control  
1060 points from SAR derived imagery. The approach involved using cGANs for generating SAR like images  
1061 from optical imagery, which can then be compared to actual SAR observation and used to extract  
1062 ground control points which are used by matching. Once the matching parameters are found, the can

1063 be directly applied to the original optical imagery to achieve high quality registration. The proposed  
 1064 method is validated through the demonstration of increased matching accuracy on observations from  
 1065 TerraSAR-X and PRISM measurements.

### 1066 5.3. Discussion

1067 Table 8 provides a quick outlook at the different DNN-based approaches in remote sensing  
 1068 observation fusion. We note that we focus on the case where observations from multiple sources are  
 1069 fused in order to provide high quality and/or more informative observations. Given the existing  
 1070 state-of-the-art, the following observation can be made.

Method	Inputs	Objective	Approach
[142]	MS and HS	spatial/spectral resolution	3D-CNN using concatenated observations
[143,144]	MS and HS	spatial/spectral resolution	Fusion using two-stream CNNs
[145]	MS and HS	spatial/temporal resolution	CNN sub-networks fusion
[147]	MS and SAR	NDVI estimation	CNN using concatenated observations
[148,150]	RGB and SAR	registration	GAN-based framework

Table 8. Listing on key approaches in DNN based observation fusion.

1071 In Table 9, the performance of three MS/HS fusion methods, two specifically developed for  
 1072 fusion [142,144], and a generic pan-sharpening one [76] with respect to PSNR are presented. These  
 1073 results indicate the problem of HS/MS fusion is different compared to pan-sharpening and that more  
 1074 case-specific methods, for example [144] which considers the fusion of features from both modalities,  
 1075 offer better performance.

Method	Dataset		
	Botswana (NASA EO-1)	Washington DC Mall (HYDICE)	Pavia University (AVIRIS)
PFCN+GDL [144]	35.36	38.38	41.80
3D-CNN [142]	34.03	36.48	39.93
PNN [76]	30.30	28.71	36.51

Table 9. Relative performance gains (with respect to PSNR) for three DNN based MS/HS fusion approaches

- 1076 • Almost exclusively, CNN architectures has been used for HS and MS fusion, the majority of  
 1077 which follow the same principles as the case of pan-sharpening.
- 1078 • Although different approaches have been proposed for addressing pairs of resolution dimensions,  
 1079 i.e., spatial-spectral and spatial-temporal, no approach has been put forward for increasing the  
 1080 resolution along spatial, spectral and temporal resolution.
- 1081 • There is limited investigation in fusion of observation from different modalities, i.e., optical and  
 1082 radar. We believe this domain to be extremely promising and hence more research needs to be  
 1083 conducted.

## 1084 6. Challenges and perspectives

1085 The impact of machine learning in remote sensing observation analysis has been of paramount  
 1086 importance, since the first application of such methods in the later 90's, primarily focusing on the  
 1087 automated extraction of information like land cover estimation [153] and detection of oil spills [154].  
 1088 Since the early 2000, machine learning has been gaining attention for problems related to observation  
 1089 enhancement like super-resolution [155], while in the past five years, this problem has been radically  
 1090 addressed through DNN approaches, demonstrating that machine learning and DNN is particular  
 1091 have still significant unexplored potential in the remote sensing domain.

1092 Karpatne et al. [156] present a comprehensive investigation of the potential and challenges of  
1093 machine learning in geosciences, a key application domain of remote sensing, which involve: (i)  
1094 inherent challenges associated with geoscience processes like spatio-temporal characteristics and  
1095 multi-variate nature of phenomena, (ii) challenges associated with geoscience data collection including  
1096 different sampling resolutions, missing measurements and highly varying quality and (iii) scarcity of  
1097 training examples and associated ground truth.

1098 One of the major benefit of DNN is their ability in *simultaneously addressing multiple problems* [157].  
1099 For example in [158], the authors propose a multi-scale fully convolutional network for multi-task  
1100 problems and more specifically for the simultaneous super-resolution and colorization of remote  
1101 sensing imagery. Another potential area of innovation involves the exploration of enhancing satellite  
1102 derived products like land surface temperature and soil moisture. As an illustrative example, while  
1103 typical remote sensing super-resolution approaches are applied to imagery, may that be PAN, MS  
1104 or HS, CNNs [34] have also been recently considered for super-resolving climate variables like daily  
1105 precipitation measurements [159] surpassing in performance more established approaches.

1106 From an architectural point-of-view, undoubtedly, one of the most interesting and innovative  
1107 approaches is the concept of GANs. In the previous sections, the application of GANs in different  
1108 problems are presented, including super-resolution [64], pan-sharpening [116] and restoration [127].  
1109 In addition to the potential of GANs in well-known problems, GANs have also been proposed  
1110 for addressing emerging problems. A particular instance involves the generating realistic training  
1111 examples. as in the case of [160], where Lin et al. explore the potential of generating realistic remote  
1112 sensing imagery by training a GAN to map images for segmentation purposes, while in [161], GANs  
1113 are employed for generating ground-level views from aerial and satellite imagery.

1114 While utilizing machine learning for the automated extraction of actionable information from a  
1115 single source has been extensively studied, substantially less effort has been allocated on approaches  
1116 for jointly analyzing observations from *multiple modalities* and *different sources*. Compared to the  
1117 single-modality case, handling observations from multiple instruments can lead to a much higher  
1118 quality estimation of geophysical parameters, exploiting each instrument's unique capabilities.  
1119 Achieving this objective requires the simultaneous encoding of diverse types of observations where  
1120 each spatial location on the Earth for example, must be associated with a spectral profile encoding the  
1121 response over multiple wavelengths acquired with MS, and reflected signals of different polarizations  
1122 acquired by SAR. Enabling this capability requires the ability to simultaneously encode high  
1123 dimensional observation, more than 3 which is currently the state-of-the-art.

1124 Another, equally important challenge characterizing existing approaches, is the inability of  
1125 integrating observations associated with *diverse sampling scales*. In the case of remote sensing, this  
1126 translates to integrating low spatial resolution global-scale satellite data with high-accuracy localized  
1127 in-situ sensor network measurements. A consequence of paramount importance related to this issue is  
1128 that the systematically-acquired satellite observations require human experts for providing annotations,  
1129 which, although of high quality, cannot be effectively scaled-up. To address this challenge, coarse  
1130 satellite observations will need to be automatically annotated using "ground-truth" measurements  
1131 from in-situ networks and ground-based surveys, seamlessly handling the variety in spatial and  
1132 temporal scales, and the irregularity of sensing patterns. Some initial attempts have been explored for  
1133 remote sensing observation classification e.g. [7], however, more research is needed in order to enable  
1134 the automated accurate and timely estimation of key environmental parameters at global scales.

1135 Analyzing observations from a single time instance cannot provide the necessary insight into the  
1136 temporal dynamics of phenomena of interest. *Time-series analysis* can provide concrete solutions to  
1137 problems like estimating measurements from regions where no up-to-date observations are available,  
1138 effectively increasing the temporal resolution. Furthermore, data-driven time-series processing can  
1139 enable the prompt identification of anomalies, a crucial issue since subtle changes within normal  
1140 variation, which cannot be identified using pre-defined thresholds, can indicate the onset of major  
1141 environmental events. Multi-temporal remote sensing observation classification is gaining traction,

1142 e.g. [162], however, addressing the challenges associated with observation enhancement is still in its  
 1143 infancy and can only be achieved by automatically exploiting correlations across modalities and time.

1144 A last yet quite important point, is the incorporation of prior knowledge in DNN models, a  
 1145 very active research topics which seeks ways of introducing physics guided modeling into the design  
 1146 of neural networks. Indeed, the use of highly accurate and physical plausibility constraints are  
 1147 key ingredients when trying to utilize DNN for scientific discovery, especially for the case of Earth  
 1148 monitoring [163]. Currently, the few research papers published on the topics consider the case of  
 1149 dynamical systems where in addition to a typical loss function, an additional regularization penalizing  
 1150 physical inconsistency is introduced. For example, in the work by Karpatne et al. [164,165], the output  
 1151 of a physics-driven model of lake temperature is considered as extra features, in addition to actual  
 1152 drivers like meteorological conditions, in the DNN algorithm, while a physics-based loss function is  
 1153 also considered during the optimization process. Although this framework has not been considered  
 1154 for the case of remote sensing observation analysis and enhancement, we expect that significant  
 1155 benefits can results from the integration of data-driven and physics-driven models, where for example,  
 1156 one could consider the physical processes governing soil dynamics during a dynamic observation  
 1157 super-resolution process.

1158 **Author Contributions:** conceptualization, G.T. and P.T.; methodology, G.T., K.F., A.A., A.P., M.G., P.T.; formal  
 1159 analysis, G.T., K.F., A.A., A.P., M.G., P.T.; writing—original draft preparation, G.T., A.A. K.F., M.G., A.A.;  
 1160 writing—review and editing, P.T.; visualization, G.T.; supervision, P.T.; project administration, P.T.; funding  
 1161 acquisition, P.T.

1162 **Funding:** The research work was funded by Greece and the European Union (European Social Fund) in the  
 1163 context of the Youth Employment Initiative through the Operational Program for Human Resources Development,  
 1164 Education and Lifelong Learning, under grant no. MIS 5004457.

1165 **Conflicts of Interest:** The authors declare no conflict of interest.

## 1166 Abbreviations

1167 The following abbreviations are used in this manuscript:

1168	(S) AE	(Stacked) Autoencoder
	CNN	Convolutional Neural Network
	DNN	Deep Neural Network
	EO	Earth Observation
	GAN	Generative Adversarial Network
	HS	Hyperspectral
	NDVI	Normalized Difference Vegetation Index
	MS	Multispectral
1169	MSE	Mean-Squared-Error
	RGB	Red-Green-Blue color channels
	PAN	Panchromatic
	PSNR	Peak Signal-to-Noise Ratio
	SAE	Stacked Autoencoders
	SAM	Spectral Angle Mapper
	SAR	Synthetic Aperture Radar
	SWIR	Short-Wave Infra-Red
	SSIM	Structural SIMilarity (Index)

## 1170 References

- 1171 1. Nativi, S.; Mazzetti, P.; Santoro, M.; Papeschi, F.; Craglia, M.; Ochiai, O. Big data challenges in building the  
 1172 global earth observation system of systems. *Environmental Modelling & Software* **2015**, *68*, 1–26.
- 1173 2. Ma, Y.; Wu, H.; Wang, L.; Huang, B.; Ranjan, R.; Zomaya, A.; Jie, W. Remote sensing big data computing:  
 1174 Challenges and opportunities. *Future Generation Computer Systems* **2015**, *51*, 47–60.

- 1175 3. Lary, D.J.; Alavi, A.H.; Gandomi, A.H.; Walker, A.L. Machine learning in geosciences and remote sensing.  
1176 *Geoscience Frontiers* **2016**, *7*, 3–10.
- 1177 4. Chen, X.W.; Lin, X. Big data deep learning: challenges and perspectives. *IEEE access* **2014**, *2*, 514–525.
- 1178 5. Zhang, L.; Zhang, L.; Du, B. Deep learning for remote sensing data: A technical tutorial on the state of the  
1179 art. *IEEE Geoscience and Remote Sensing Magazine* **2016**, *4*, 22–40.
- 1180 6. Romero, A.; Gatta, C.; Camps-Valls, G. Unsupervised deep feature extraction for remote sensing image  
1181 classification. *IEEE Transactions on Geoscience and Remote Sensing* **2016**, *54*, 1349–1362.
- 1182 7. Karalas, K.; Tsagkatakis, G.; Zervakis, M.; Tsakalides, P. Deep learning for multi-label land cover  
1183 classification. *Image and Signal Processing for Remote Sensing XXI. International Society for Optics  
1184 and Photonics*, 2015, Vol. 9643, p. 96430Q.
- 1185 8. Maqueda, A.I.; Loquercio, A.; Gallego, G.; García, N.; Scaramuzza, D. Event-based vision meets deep  
1186 learning on steering prediction for self-driving cars. *Proceedings of the IEEE Conference on Computer  
1187 Vision and Pattern Recognition*, 2018, pp. 5419–5427.
- 1188 9. Wang, S.; Su, Z.; Ying, L.; Peng, X.; Zhu, S.; Liang, F.; Feng, D.; Liang, D. Accelerating magnetic resonance  
1189 imaging via deep learning. *2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI). IEEE*,  
1190 2016, pp. 514–517.
- 1191 10. Ribes, A.; Schmitt, F. Linear inverse problems in imaging. *IEEE Signal Processing Magazine* **2008**, *25*, 84–99.
- 1192 11. Yu, G.; Sapiro, G.; Mallat, S. Solving inverse problems with piecewise linear estimators: From Gaussian  
1193 mixture models to structured sparsity. *IEEE Transactions on Image Processing* **2012**, *21*, 2481–2499.
- 1194 12. Jin, K.H.; McCann, M.T.; Froustey, E.; Unser, M. Deep convolutional neural network for inverse problems  
1195 in imaging. *IEEE Transactions on Image Processing* **2017**, *26*, 4509–4522.
- 1196 13. Goodfellow, I.; Bengio, Y.; Courville, A.; Bengio, Y. *Deep learning*; Vol. 1, MIT press Cambridge, 2016.
- 1197 14. He, K.; Zhang, X.; Ren, S.; Sun, J. Delving deep into rectifiers: Surpassing human-level performance on  
1198 imagenet classification. *Proceedings of the IEEE international conference on computer vision*, 2015, pp.  
1199 1026–1034.
- 1200 15. Dong, C.; Loy, C.C.; Tang, X. Accelerating the super-resolution convolutional neural network. *European  
1201 Conference on Computer Vision. Springer*, 2016, pp. 391–407.
- 1202 16. Shi, W.; Caballero, J.; Huszár, F.; Totz, J.; Aitken, A.P.; Bishop, R.; Rueckert, D.; Wang, Z. Real-time single  
1203 image and video super-resolution using an efficient sub-pixel convolutional neural network. *Proceedings  
1204 of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1874–1883.
- 1205 17. Wang, T.; Sun, M.; Hu, K. Dilated deep residual network for image denoising. *2017 IEEE 29th International  
1206 Conference on Tools with Artificial Intelligence (ICTAI). IEEE*, 2017, pp. 1272–1279.
- 1207 18. Lin, G.; Wu, Q.; Qiu, L.; Huang, X. Image super-resolution using a dilated convolutional neural network.  
1208 *Neurocomputing* **2018**, *275*, 1219–1230.
- 1209 19. Zhang, Q.; Yuan, Q.; Li, J.; Yang, Z.; Ma, X. Learning a dilated residual network for SAR image despeckling.  
1210 *Remote Sensing* **2018**, *10*, 196.
- 1211 20. Haykin, S.; Network, N. A comprehensive foundation. *Neural networks* **2004**, *2*, 41.
- 1212 21. Stivaktakis, R.; Tsagkatakis, G.; Tsakalides, P. Deep Learning for Multilabel Land Cover Scene  
1213 Categorization Using Data Augmentation. *IEEE Geoscience and Remote Sensing Letters* **2019**.
- 1214 22. Hamida, A.B.; Benoit, A.; Lambert, P.; Amar, C.B. 3-D Deep learning approach for remote sensing image  
1215 classification. *IEEE Transactions on Geoscience and Remote Sensing* **2018**, *56*, 4420–4434.
- 1216 23. Fotiadou, K.; Tsagkatakis, G.; Tsakalides, P. Deep convolutional neural networks for the classification of  
1217 snapshot mosaic hyperspectral imagery. *Electronic Imaging* **2017**, 2017, 185–190.
- 1218 24. Bottou, L. Online learning and stochastic approximations. *On-line learning in neural networks* **1998**, *17*, 142.
- 1219 25. Kingma, D.P.; Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* **2014**.
- 1220 26. Zeiler, M.D. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701* **2012**.
- 1221 27. Bengio, Y. Deep learning of representations for unsupervised and transfer learning. *Proceedings of ICML  
1222 Workshop on Unsupervised and Transfer Learning*, 2012, pp. 17–36.
- 1223 28. Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; Salakhutdinov, R. Dropout: a simple way to  
1224 prevent neural networks from overfitting. *The Journal of Machine Learning Research* **2014**, *15*, 1929–1958.
- 1225 29. Ioffe, S.; Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal  
1226 covariate shift. *arXiv preprint arXiv:1502.03167* **2015**.

- 1227 30. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. Proceedings of the IEEE  
1228 conference on computer vision and pattern recognition, 2016, pp. 770–778.
- 1229 31. Kim, J.; Kwon Lee, J.; Mu Lee, K. Accurate image super-resolution using very deep convolutional networks.  
1230 Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1646–1654.
- 1231 32. Szegedy, C.; Liu, W.; Jia, Y.; Sermanet, P.; Reed, S.; Anguelov, D.; Erhan, D.; Vanhoucke, V.; Rabinovich,  
1232 A. Going deeper with convolutions. Proceedings of the IEEE conference on computer vision and pattern  
1233 recognition, 2015, pp. 1–9.
- 1234 33. Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; Wojna, Z. Rethinking the inception architecture for  
1235 computer vision. Proceedings of the IEEE conference on computer vision and pattern recognition, 2016,  
1236 pp. 2818–2826.
- 1237 34. Dong, C.; Loy, C.C.; He, K.; Tang, X. Image super-resolution using deep convolutional networks. *IEEE*  
1238 *transactions on pattern analysis and machine intelligence* **2016**, *38*, 295–307.
- 1239 35. Ronneberger, O.; Fischer, P.; Brox, T. U-net: Convolutional networks for biomedical image segmentation.  
1240 International Conference on Medical image computing and computer-assisted intervention. Springer, 2015,  
1241 pp. 234–241.
- 1242 36. Hinton, G.E.; Zemel, R.S. Autoencoders, minimum description length and Helmholtz free energy.  
1243 Advances in neural information processing systems, 1994, pp. 3–10.
- 1244 37. Bengio, Y.; Lamblin, P.; Popovici, D.; Larochelle, H. Greedy layer-wise training of deep networks. Advances  
1245 in neural information processing systems, 2007, pp. 153–160.
- 1246 38. Vincent, P.; Larochelle, H.; Lajoie, I.; Bengio, Y.; Manzagol, P.A. Stacked denoising autoencoders: Learning  
1247 useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*  
1248 **2010**, *11*, 3371–3408.
- 1249 39. Shin, H.C.; Orton, M.R.; Collins, D.J.; Doran, S.J.; Leach, M.O. Stacked autoencoders for unsupervised  
1250 feature learning and multiple organ detection in a pilot study using 4D patient data. *IEEE transactions on*  
1251 *pattern analysis and machine intelligence* **2013**, *35*, 1930–1943.
- 1252 40. Glorot, X.; Bordes, A.; Bengio, Y. Domain adaptation for large-scale sentiment classification: A deep  
1253 learning approach. Proceedings of the 28th international conference on machine learning (ICML-11), 2011,  
1254 pp. 513–520.
- 1255 41. Vincent, P.; Larochelle, H.; Bengio, Y.; Manzagol, P.A. Extracting and composing robust features with  
1256 denoising autoencoders. Proceedings of the 25th international conference on Machine learning. ACM,  
1257 2008, pp. 1096–1103.
- 1258 42. Doersch, C. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908* **2016**.
- 1259 43. Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y.  
1260 Generative adversarial nets. Advances in neural information processing systems, 2014, pp. 2672–2680.
- 1261 44. Wang, T.C.; Liu, M.Y.; Zhu, J.Y.; Tao, A.; Kautz, J.; Catanzaro, B. High-resolution image synthesis and  
1262 semantic manipulation with conditional gans. Proceedings of the IEEE Conference on Computer Vision  
1263 and Pattern Recognition, 2018, pp. 8798–8807.
- 1264 45. Isola, P.; Zhu, J.Y.; Zhou, T.; Efros, A.A. Image-to-image translation with conditional adversarial networks.  
1265 Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1125–1134.
- 1266 46. Ledig, C.; Theis, L.; Huszár, F.; Caballero, J.; Cunningham, A.; Acosta, A.; Aitken, A.; Tejani, A.; Totz, J.;  
1267 Wang, Z.; others. Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network.  
1268 Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on. IEEE, 2017, pp. 105–114.
- 1269 47. Zeng, Y.; Huang, W.; Liu, M.; Zhang, H.; Zou, B. Fusion of satellite images in urban area: Assessing the  
1270 quality of resulting images. 2010 18th International Conference on Geoinformatics. IEEE, 2010, pp. 1–4.
- 1271 48. Fan, J. Chinese Earth Observation Program and Policy. In *Satellite Earth Observations and Their Impact on*  
1272 *Society and Policy*; Springer, 2017; pp. 105–110.
- 1273 49. Cheng, G.; Han, J.; Lu, X. Remote sensing image scene classification: Benchmark and state of the art.  
1274 *Proceedings of the IEEE* **2017**, *105*, 1865–1883.
- 1275 50. Lim, B.; Son, S.; Kim, H.; Nah, S.; Lee, K.M. Enhanced deep residual networks for single image  
1276 super-resolution. The IEEE conference on computer vision and pattern recognition (CVPR) workshops,  
1277 2017, Vol. 1, p. 4.

- 1278 51. Liebel, L.; Körner, M. SINGLE-IMAGE SUPER RESOLUTION FOR MULTISPECTRAL REMOTE SENSING  
1279 DATA USING CONVOLUTIONAL NEURAL NETWORKS. *International Archives of the Photogrammetry,*  
1280 *Remote Sensing & Spatial Information Sciences* **2016**, *41*.
- 1281 52. Tuna, C.; Unal, G.; Sertel, E. Single-frame super resolution of remote-sensing images by convolutional  
1282 neural networks. *International Journal of Remote Sensing* **2018**, *39*, 2463–2479.
- 1283 53. Huang, N.; Yang, Y.; Liu, J.; Gu, X.; Cai, H. Single-Image Super-Resolution for Remote Sensing Data Using  
1284 Deep Residual-Learning Neural Network. International Conference on Neural Information Processing.  
1285 Springer, 2017, pp. 622–630.
- 1286 54. Lei, S.; Shi, Z.; Zou, Z. Super-resolution for remote sensing images via local–global combined network.  
1287 *IEEE Geoscience and Remote Sensing Letters* **2017**, *14*, 1243–1247.
- 1288 55. Xu, W.; Guangluan, X.; Wang, Y.; Sun, X.; Lin, D.; Yirong, W. High Quality Remote Sensing Image  
1289 Super-Resolution Using Deep Memory Connected Network. IGARSS 2018-2018 IEEE International  
1290 Geoscience and Remote Sensing Symposium. IEEE, 2018, pp. 8889–8892.
- 1291 56. Wang, T.; Sun, W.; Qi, H.; Ren, P. Aerial Image Super Resolution via Wavelet Multiscale Convolutional  
1292 Neural Networks. *IEEE Geoscience and Remote Sensing Letters* **2018**, *15*, 769–773.
- 1293 57. Ma, W.; Pan, Z.; Guo, J.; Lei, B. Achieving Super-Resolution Remote Sensing Images via the Wavelet  
1294 Transform Combined With the Recursive Res-Net. *IEEE Transactions on Geoscience and Remote Sensing* **2019**.
- 1295 58. Tai, Y.; Yang, J.; Liu, X. Image super-resolution via deep recursive residual network. Proceedings of the  
1296 IEEE conference on computer vision and pattern recognition, 2017, pp. 3147–3155.
- 1297 59. Lu, T.; Wang, J.; Zhang, Y.; Wang, Z.; Jiang, J. Satellite Image Super-Resolution via Multi-Scale Residual  
1298 Deep Neural Network. *Remote Sensing* **2019**, *11*, 1588.
- 1299 60. Pan, Z.; Ma, W.; Guo, J.; Lei, B. Super-Resolution of Single Remote Sensing Image Based on Residual Dense  
1300 Backprojection Networks. *IEEE Transactions on Geoscience and Remote Sensing* **2019**.
- 1301 61. Haris, M.; Shakhnarovich, G.; Ukita, N. Deep back-projection networks for super-resolution. Proceedings  
1302 of the IEEE conference on computer vision and pattern recognition, 2018, pp. 1664–1673.
- 1303 62. Haut, J.M.; Fernandez-Beltran, R.; Paoletti, M.E.; Plaza, J.; Plaza, A.; Pla, F. A new deep generative network  
1304 for unsupervised remote sensing single-image super-resolution. *IEEE Transactions on Geoscience and Remote*  
1305 *Sensing* **2018**, pp. 1–19.
- 1306 63. Ma, W.; Pan, Z.; Guo, J.; Lei, B. Super-Resolution of Remote Sensing Images Based on Transferred  
1307 Generative Adversarial Network. IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing  
1308 Symposium. IEEE, 2018, pp. 1148–1151.
- 1309 64. Jiang, K.; Wang, Z.; Yi, P.; Wang, G.; Lu, T.; Jiang, J. Edge-Enhanced GAN for Remote Sensing Image  
1310 Superresolution. *IEEE Transactions on Geoscience and Remote Sensing* **2019**.
- 1311 65. Shuai, Y.; Wang, Y.; Peng, Y.; Xia, Y. Accurate Image Super-Resolution Using Cascaded Multi-Column  
1312 Convolutional Neural Networks. 2018 IEEE International Conference on Multimedia and Expo (ICME).  
1313 IEEE, 2018, pp. 1–6.
- 1314 66. Li, Y.; Hu, J.; Zhao, X.; Xie, W.; Li, J. Hyperspectral image super-resolution using deep convolutional neural  
1315 network. *Neurocomputing* **2017**, *266*, 29–41.
- 1316 67. Wang, C.; Liu, Y.; Bai, X.; Tang, W.; Lei, P.; Zhou, J. Deep Residual Convolutional Neural Network for  
1317 Hyperspectral Image Super-Resolution. International Conference on Image and Graphics. Springer, 2017,  
1318 pp. 370–380.
- 1319 68. Yuan, Y.; Zheng, X.; Lu, X. Hyperspectral image superresolution by transfer learning. *IEEE Journal of*  
1320 *Selected Topics in Applied Earth Observations and Remote Sensing* **2017**, *10*, 1963–1974.
- 1321 69. Collins, C.B.; Beck, J.M.; Bridges, S.M.; Rushing, J.A.; Graves, S.J. Deep learning for multisensor image  
1322 resolution enhancement. Proceedings of the 1st Workshop on Artificial Intelligence and Deep Learning for  
1323 Geographic Knowledge Discovery. ACM, 2017, pp. 37–44.
- 1324 70. He, Z.; Liu, L. Hyperspectral Image Super-Resolution Inspired by Deep Laplacian Pyramid Network.  
1325 *Remote Sensing* **2018**, *10*, 1939.
- 1326 71. Lai, W.S.; Huang, J.B.; Ahuja, N.; Yang, M.H. Deep laplacian pyramid networks for fast and accurate  
1327 superresolution. IEEE Conference on Computer Vision and Pattern Recognition, 2017, Vol. 2, p. 5.
- 1328 72. Zheng, K.; Gao, L.; Zhang, B.; Cui, X. Multi-Losses Function Based Convolution Neural Network for Single  
1329 Hyperspectral Image Super-Resolution. 2018 Fifth International Workshop on Earth Observation and  
1330 Remote Sensing Applications (EORSA). IEEE, 2018, pp. 1–4.

- 1331 73. Lanaras, C.; Bioucas-Dias, J.; Galliani, S.; Baltasvias, E.; Schindler, K. Super-resolution of Sentinel-2 images:  
1332 Learning a globally applicable deep neural network. *ISPRS Journal of Photogrammetry and Remote Sensing*  
1333 **2018**, *146*, 305 – 319. doi:https://doi.org/10.1016/j.isprsjprs.2018.09.018.
- 1334 74. Palsson, F.; Sveinsson, J.; Ulfarsson, M. Sentinel-2 Image Fusion Using a Deep Residual Network. *Remote*  
1335 *Sensing* **2018**, *10*, 1290.
- 1336 75. Gargiulo, M.; Mazza, A.; Gaetano, R.; Ruello, G.; Scarpa, G. A CNN-Based Fusion Method for  
1337 Super-Resolution of Sentinel-2 Data. IGARSS 2018-2018 IEEE International Geoscience and Remote  
1338 Sensing Symposium. IEEE, 2018, pp. 4713–4716.
- 1339 76. Masi, G.; Cozzolino, D.; Verdoliva, L.; Scarpa, G. Pansharpening by convolutional neural networks. *Remote*  
1340 *Sensing* **2016**, *8*, 594.
- 1341 77. Zheng, K.; Gao, L.; Ran, Q.; Cui, X.; Zhang, B.; Liao, W.; Jia, S. Separable-spectral convolution and inception  
1342 network for hyperspectral image super-resolution. *International Journal of Machine Learning and Cybernetics*  
1343 **2019**, pp. 1–15.
- 1344 78. Mei, S.; Yuan, X.; Ji, J.; Zhang, Y.; Wan, S.; Du, Q. Hyperspectral Image Spatial Super-Resolution via 3D  
1345 Full Convolutional Neural Network. *Remote Sensing* **2017**, *9*, 1139.
- 1346 79. Arun, P.V.; Herrmann, I.; Budhiraju, K.M.; Karnieli, A. Convolutional network architectures for  
1347 super-resolution/sub-pixel mapping of drone-derived images. *Pattern Recognition* **2019**, *88*, 431–446.
- 1348 80. Ran, Q.; Xu, X.; Zhao, S.; Li, W.; Du, Q. Remote sensing images super-resolution with deep convolution  
1349 networks. *Multimedia Tools and Applications* **2019**, pp. 1–17.
- 1350 81. Shi, Z.; Chen, C.; Xiong, Z.; Liu, D.; Wu, F. Hscnn+: Advanced cnn-based hyperspectral recovery from rgb  
1351 images. IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 2018,  
1352 Vol. 3, p. 5.
- 1353 82. Han, X.H.; Shi, B.; Zheng, Y. Residual HSRCNN: Residual Hyper-Spectral Reconstruction CNN from an  
1354 RGB Image. 2018 24th International Conference on Pattern Recognition (ICPR). IEEE, 2018, pp. 2664–2669.
- 1355 83. Shi, Z.; Chen, C.; Xiong, Z.; Liu, D.; Zha, Z.J.; Wu, F. Deep residual attention network for spectral image  
1356 super-resolution. Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 0–0.
- 1357 84. Eldar, Y.C.; Kutyniok, G. *Compressed sensing: theory and applications*; Cambridge university press, 2012.
- 1358 85. Golbabaee, M.; Vandergheynst, P. Hyperspectral image compressed sensing via low-rank and joint-sparse  
1359 matrix recovery. 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).  
1360 Ieee, 2012, pp. 2741–2744.
- 1361 86. Tsagakatakis, G.; Tsakalides, P. Compressed hyperspectral sensing. Image Sensors and Imaging Systems  
1362 2015. International Society for Optics and Photonics, 2015, Vol. 9403, p. 940307.
- 1363 87. Yuan, H.; Yan, F.; Chen, X.; Zhu, J. Compressive Hyperspectral Imaging and Super-resolution. 2018 IEEE  
1364 3rd International Conference on Image, Vision and Computing (ICIVC). IEEE, 2018, pp. 618–623.
- 1365 88. Arce, G.R.; Brady, D.J.; Carin, L.; Arguello, H.; Kittle, D.S. Compressive coded aperture spectral imaging:  
1366 An introduction. *IEEE Signal Processing Magazine* **2013**, *31*, 105–115.
- 1367 89. Xiong, Z.; Shi, Z.; Li, H.; Wang, L.; Liu, D.; Wu, F. Hscnn: Cnn-based hyperspectral image recovery from  
1368 spectrally undersampled projections. Proceedings of the IEEE International Conference on Computer  
1369 Vision Workshops, 2017, Vol. 2.
- 1370 90. Wang, L.; Zhang, T.; Fu, Y.; Huang, H. HyperReconNet: Joint Coded Aperture Optimization and Image  
1371 Reconstruction for Compressive Hyperspectral Imaging. *IEEE Transactions on Image Processing* **2018**,  
1372 *28*, 2257–2270.
- 1373 91. Luo, Y.; Zhou, L.; Wang, S.; Wang, Z. Video Satellite Imagery Super Resolution via Convolutional Neural  
1374 Networks. *IEEE Geoscience and Remote Sensing Letters* **2017**, *14*, 2398–2402.
- 1375 92. Xiao, A.; Wang, Z.; Wang, L.; Ren, Y. Super-Resolution for “Jilin-1” Satellite Video Imagery via a  
1376 Convolutional Network. *Sensors* **2018**, *18*, 1194.
- 1377 93. Jiang, K.; Wang, Z.; Yi, P.; Jiang, J. A progressively enhanced network for video satellite imagery  
1378 superresolution. *IEEE Signal Processing Letters* **2018**, *25*, 1630–1634.
- 1379 94. Jiang, K.; Wang, Z.; Yi, P.; Jiang, J.; Xiao, J.; Yao, Y. Deep Distillation Recursive Network for Remote Sensing  
1380 Imagery Super-Resolution. *Remote Sensing* **2018**, *10*, 1700.
- 1381 95. Yang, W.; Feng, J.; Xie, G.; Liu, J.; Guo, Z.; Yan, S. Video super-resolution based on spatial-temporal  
1382 recurrent residual networks. *Computer Vision and Image Understanding* **2018**, *168*, 79–92.

- 1383 96. Zhang, Y.; Mishra, R.K. A review and comparison of commercially available pan-sharpening techniques for  
1384 high resolution satellite image fusion. 2012 IEEE International Geoscience and Remote Sensing Symposium.  
1385 IEEE, 2012, pp. 182–185.
- 1386 97. Huang, W.; Xiao, L.; Wei, Z.; Liu, H.; Tang, S. A new pan-sharpening method with deep neural networks.  
1387 *IEEE Geoscience and Remote Sensing Letters* **2015**, *12*, 1037–1041.
- 1388 98. Cai, W.; Xu, Y.; Wu, Z.; Liu, H.; Qian, L.; Wei, Z. Pan-Sharpener Based on Multilevel Coupled Deep  
1389 Network. IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium. IEEE, 2018,  
1390 pp. 7046–7049.
- 1391 99. Wei, Y.; Yuan, Q.; Shen, H.; Zhang, L. Boosting the accuracy of multispectral image pansharpening by  
1392 learning a deep residual network. *IEEE Geoscience and Remote Sensing Letters* **2017**, *14*, 1795–1799.
- 1393 100. Azarang, A.; Ghassemian, H. A new pansharpening method using multi resolution analysis framework  
1394 and deep neural networks. Pattern Recognition and Image Analysis (IPRIA), 2017 3rd International  
1395 Conference on. IEEE, 2017, pp. 1–6.
- 1396 101. Azarang, A.; Manoochehri, H.E.; Kehtarnavaz, N. Convolutional Autoencoder-Based Multispectral Image  
1397 Fusion. *IEEE Access* **2019**, *7*, 35673–35683.
- 1398 102. Zhong, J.; Yang, B.; Huang, G.; Zhong, F.; Chen, Z. Remote sensing image fusion with convolutional neural  
1399 network. *Sensing and Imaging* **2016**, *17*, 10.
- 1400 103. Li, N.; Huang, N.; Xiao, L. Pan-sharpening via residual deep learning. Geoscience and Remote Sensing  
1401 Symposium (IGARSS), 2017 IEEE International. IEEE, 2017, pp. 5133–5136.
- 1402 104. Wei, Y.; Yuan, Q. Deep residual learning for remote sensed imagery pansharpening. 2017 International  
1403 Workshop on Remote Sensing with Intelligent Processing (RSIP). IEEE, 2017, pp. 1–4.
- 1404 105. Rao, Y.; He, L.; Zhu, J. A residual convolutional neural network for pan-sharpening. Remote Sensing with  
1405 Intelligent Processing (RSIP), 2017 International Workshop on. IEEE, 2017, pp. 1–4.
- 1406 106. Yang, J.; Fu, X.; Hu, Y.; Huang, Y.; Ding, X.; Paisley, J. PanNet: A deep network architecture for  
1407 pan-sharpening. Proc. IEEE Int. Conf. Comput. Vis.(ICCV), 2017, pp. 1753–1761.
- 1408 107. Yuan, Q.; Wei, Y.; Meng, X.; Shen, H.; Zhang, L. A multiscale and multidepth convolutional neural network  
1409 for remote sensing imagery pan-sharpening. *IEEE Journal of Selected Topics in Applied Earth Observations and*  
1410 *Remote Sensing* **2018**, *11*, 978–989.
- 1411 108. Wei, Y.; Yuan, Q.; Meng, X.; Shen, H.; Zhang, L.; Ng, M. Multi-scale-and-depth convolutional neural  
1412 network for remote sensed imagery pan-sharpening. 2017 IEEE International Geoscience and Remote  
1413 Sensing Symposium (IGARSS). IEEE, 2017, pp. 3413–3416.
- 1414 109. Masi, G.; Cozzolino, D.; Verdoliva, L.; Scarpa, G. CNN-based pansharpening of multi-resolution  
1415 remote-sensing images. 2017 Joint Urban Remote Sensing Event (JURSE). IEEE, 2017, pp. 1–4.
- 1416 110. Scarpa, G.; Vitale, S.; Cozzolino, D. Target-adaptive CNN-based pansharpening. *IEEE Transactions on*  
1417 *Geoscience and Remote Sensing* **2018**, pp. 1–15.
- 1418 111. Shao, Z.; Cai, J. Remote Sensing Image Fusion With Deep Convolutional Neural Network. *IEEE Journal of*  
1419 *Selected Topics in Applied Earth Observations and Remote Sensing* **2018**, *11*, 1656–1669.
- 1420 112. Zhang, Y.; Liu, C.; Sun, M.; Ou, Y. Pan-Sharpener Using an Efficient Bidirectional Pyramid Network.  
1421 *IEEE Transactions on Geoscience and Remote Sensing* **2019**.
- 1422 113. Yao, W.; Zeng, Z.; Lian, C.; Tang, H. Pixel-wise Regression using U-Net and its Application on  
1423 Pansharpening. *Neurocomputing* **2018**.
- 1424 114. Eghbalian, S.; Ghassemian, H. Multi spectral image fusion by deep convolutional neural network and new  
1425 spectral loss function. *International Journal of Remote Sensing* **2018**, *39*, 3983–4002.
- 1426 115. Guo, Y.; Ye, F.; Gong, H. Learning an efficient convolution neural network for pansharpening. *Algorithms*  
1427 **2019**, *12*, 16.
- 1428 116. Liu, X.; Wang, Y.; Liu, Q. PSGAN: A Generative Adversarial Network for Remote Sensing Image  
1429 Pan-Sharpener. *arXiv preprint arXiv:1805.03371* **2018**.
- 1430 117. Zhang, Y.; Li, X.; Zhou, J. SFTGAN: a generative adversarial network for pan-sharpening equipped with  
1431 spatial feature transform layers. *Journal of Applied Remote Sensing* **2019**, *13*, 026507.
- 1432 118. Aiuzzi, B.; Alparone, L.; Baronti, S.; Garzelli, A.; Selva, M. MTF-tailored multiscale fusion of high-resolution  
1433 MS and Pan imagery. *Photogrammetric Engineering & Remote Sensing* **2006**, *72*, 591–596.

- 1434 119. Ghamisi, P.; Yokoya, N.; Li, J.; Liao, W.; Liu, S.; Plaza, J.; Rasti, B.; Plaza, A. Advances in hyperspectral  
1435 image and signal processing: A comprehensive overview of the state of the art. *IEEE Geoscience and Remote*  
1436 *Sensing Magazine* **2017**, *5*, 37–78.
- 1437 120. Xie, W.; Li, Y. Hyperspectral imagery denoising by deep learning with trainable nonlinearity function.  
1438 *IEEE Geoscience and Remote Sensing Letters* **2017**, *14*, 1963–1967.
- 1439 121. Xie, W.; Li, Y.; Jia, X. Deep convolutional networks with residual learning for accurate spectral-spatial  
1440 denoising. *Neurocomputing* **2018**.
- 1441 122. Xie, W.; Li, Y.; Hu, J.; Chen, D.Y. Trainable spectral difference learning with spatial starting for hyperspectral  
1442 image denoising. *Neural Networks* **2018**, *108*, 272–286.
- 1443 123. Chen, C.; Xu, Z. Aerial-Image Denoising Based on Convolutional Neural Network with Multi-Scale  
1444 Residual Learning Approach. *Information* **2018**, *9*, 169.
- 1445 124. Yuan, Q.; Zhang, Q.; Li, J.; Shen, H.; Zhang, L. Hyperspectral image denoising employing a spatial-spectral  
1446 deep residual convolutional neural network. *IEEE Transactions on Geoscience and Remote Sensing* **2018**, pp.  
1447 1–14.
- 1448 125. Liu, W.; Lee, J. A 3-D Atrous Convolution Neural Network for Hyperspectral Image Denoising. *IEEE*  
1449 *Transactions on Geoscience and Remote Sensing* **2019**.
- 1450 126. Zhang, Q.; Yuan, Q.; Li, J.; Liu, X.; Shen, H.; Zhang, L. Hybrid Noise Removal in Hyperspectral Imagery  
1451 With a Spatial-Spectral Gradient Network. *arXiv preprint arXiv:1810.00495* **2018**.
- 1452 127. Chang, Y.; Yan, L.; Fang, H.; Zhong, S.; Liao, W. HSI-DeNet: Hyperspectral image restoration via  
1453 convolutional neural network. *IEEE Transactions on Geoscience and Remote Sensing* **2018**, pp. 1–16.
- 1454 128. Zhang, Y.; Xiang, Y.; Bai, L. Generative Adversarial Network for Deblurring of Remote Sensing Image.  
1455 2018 26th International Conference on Geoinformatics. IEEE, 2018, pp. 1–4.
- 1456 129. Xu, W.; Xu, G.; Wang, Y.; Sun, X.; Lin, D.; Wu, Y. Deep Memory Connected Neural Network for Optical  
1457 Remote Sensing Image Restoration. *Remote Sensing* **2018**, *10*, 1893.
- 1458 130. Shen, H.; Li, X.; Cheng, Q.; Zeng, C.; Yang, G.; Li, H.; Zhang, L. Missing information reconstruction of  
1459 remote sensing data: A technical review. *IEEE Geoscience and Remote Sensing Magazine* **2015**, *3*, 61–85.
- 1460 131. Zhang, Q.; Yuan, Q.; Zeng, C.; Li, X.; Wei, Y. Missing Data Reconstruction in Remote Sensing image with a  
1461 Unified Spatial-Temporal-Spectral Deep Convolutional Neural Network. *IEEE Transactions on Geoscience*  
1462 *and Remote Sensing* **2018**, pp. 1–15.
- 1463 132. Das, M.; Ghosh, S.K. A deep-learning-based forecasting ensemble to predict missing data for remote  
1464 sensing analysis. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* **2017**,  
1465 *10*, 5228–5236.
- 1466 133. Das, M.; Ghosh, S.K. Deep-STEP: A deep learning approach for spatiotemporal prediction of remote  
1467 sensing data. *IEEE Geoscience and Remote Sensing Letters* **2016**, *13*, 1984–1988.
- 1468 134. Dong, J.; Yin, R.; Sun, X.; Li, Q.; Yang, Y.; Qin, X. Inpainting of Remote Sensing SST Images With  
1469 Deep Convolutional Generative Adversarial Network. *IEEE Geoscience and Remote Sensing Letters* **2019**,  
1470 *16*, 173–177.
- 1471 135. Yu, Y.; Gong, Z.; Zhong, P.; Shan, J. Unsupervised Representation Learning with Deep Convolutional  
1472 Neural Network for Remote Sensing Images. International Conference on Image and Graphics. Springer,  
1473 2017, pp. 97–108.
- 1474 136. Pathak, D.; Krahenbuhl, P.; Donahue, J.; Darrell, T.; Efros, A.A. Context encoders: Feature learning by  
1475 inpainting. Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp.  
1476 2536–2544.
- 1477 137. Bermudez, J.D.; Happ, P.N.; Feitosa, R.Q.; Oliveira, D.A. Synthesis of Multispectral Optical Images From  
1478 SAR/Optical Multitemporal Data Using Conditional Generative Adversarial Networks. *IEEE Geoscience*  
1479 *and Remote Sensing Letters* **2019**.
- 1480 138. Grohnfeldt, C.; Schmitt, M.; Zhu, X. A Conditional Generative Adversarial Network to Fuse Sar And  
1481 Multispectral Optical Data For Cloud Removal From Sentinel-2 Images. IGARSS 2018-2018 IEEE  
1482 International Geoscience and Remote Sensing Symposium. IEEE, 2018, pp. 1726–1729.
- 1483 139. Maggioni, M.; Katkovnik, V.; Egiazarian, K.; Foi, A. Nonlocal transform-domain filter for volumetric data  
1484 denoising and reconstruction. *IEEE transactions on image processing* **2012**, *22*, 119–133.
- 1485 140. Simone, G.; Farina, A.; Morabito, F.C.; Serpico, S.B.; Bruzzone, L. Image fusion techniques for remote  
1486 sensing applications. *Information fusion* **2002**, *3*, 3–15.

- 1487 141. Liu, Y.; Chen, X.; Wang, Z.; Wang, Z.J.; Ward, R.K.; Wang, X. Deep learning for pixel-level image fusion:  
1488 recent advances and future prospects. *Information Fusion* **2018**, *42*, 158–173.
- 1489 142. Palsson, F.; Sveinsson, J.R.; Ulfarsson, M.O. Multispectral and hyperspectral image fusion using a  
1490 3-D-convolutional neural network. *IEEE Geoscience and Remote Sensing Letters* **2017**, *14*, 639–643.
- 1491 143. Yang, J.; Zhao, Y.Q.; Chan, J.C.W. Hyperspectral and Multispectral Image Fusion via Deep Two-Branches  
1492 Convolutional Neural Network. *Remote Sensing* **2018**, *10*, 800.
- 1493 144. Zhou, F.; Hang, R.; Liu, Q.; Yuan, X. Pyramid Fully Convolutional Network for Hyperspectral and  
1494 Multispectral Image Fusion. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*  
1495 **2019**.
- 1496 145. Tan, Z.; Yue, P.; Di, L.; Tang, J. Deriving high spatiotemporal remote sensing images using deep  
1497 convolutional network. *Remote Sensing* **2018**, *10*, 1066.
- 1498 146. Kwan, C.; Budavari, B.; Bovik, A.C.; Marchisio, G. Blind quality assessment of fused worldview-3 images  
1499 by using the combinations of pansharpening and hypersharpening paradigms. *IEEE Geoscience and Remote*  
1500 *Sensing Letters* **2017**, *14*, 1835–1839.
- 1501 147. Scarpa, G.; Gargiulo, M.; Mazza, A.; Gaetano, R. A CNN-Based Fusion Method for Feature Extraction from  
1502 Sentinel Data. *Remote Sensing* **2018**, *10*, 236.
- 1503 148. Hughes, L.H.; Schmitt, M.; Mou, L.; Wang, Y.; Zhu, X.X. Identifying corresponding patches in sar and  
1504 optical images with a pseudo-siamese cnn. *IEEE Geoscience and Remote Sensing Letters* **2018**, *15*, 784–788.
- 1505 149. Wang, Y.; Zhu, X.X.; Zeisl, B.; Pollefeys, M. Fusing meter-resolution 4-D InSAR point clouds and optical  
1506 images for semantic urban infrastructure monitoring. *IEEE Transactions on Geoscience and Remote Sensing*  
1507 **2017**, *55*, 14–26.
- 1508 150. Quan, D.; Wang, S.; Liang, X.; Wang, R.; Fang, S.; Hou, B.; Jiao, L. Deep Generative Matching Network  
1509 for Optical and SAR Image Registration. IGARSS 2018-2018 IEEE International Geoscience and Remote  
1510 Sensing Symposium. IEEE, 2018, pp. 6215–6218.
- 1511 151. Hughes, L.; Schmitt, M.; Zhu, X. Mining Hard Negative Samples for SAR-Optical Image Matching Using  
1512 Generative Adversarial Networks. *Remote Sensing* **2018**, *10*, 1552.
- 1513 152. Merkle, N.; Auer, S.; Müller, R.; Reinartz, P. Exploring the potential of conditional adversarial networks for  
1514 optical and SAR image matching. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote*  
1515 *Sensing* **2018**, *11*, 1811–1820.
- 1516 153. DeFries, R.; Chan, J.C.W. Multiple criteria for evaluating machine learning algorithms for land cover  
1517 classification from satellite data. *Remote Sensing of Environment* **2000**, *74*, 503–515.
- 1518 154. Kubat, M.; Holte, R.C.; Matwin, S. Machine learning for the detection of oil spills in satellite radar images.  
1519 *Machine learning* **1998**, *30*, 195–215.
- 1520 155. Lu, Y.; Inamura, M. Spatial resolution improvement of remote sensing images by fusion of subpixel-shifted  
1521 multi-observation images. *International Journal of Remote Sensing* **2003**, *24*, 4647–4660.
- 1522 156. Karpatne, A.; Ebert-Uphoff, I.; Ravela, S.; Babaie, H.A.; Kumar, V. Machine learning for the geosciences:  
1523 Challenges and opportunities. *IEEE Transactions on Knowledge and Data Engineering* **2018**.
- 1524 157. Rick Chang, J.; Li, C.L.; Póczos, B.; Vijaya Kumar, B.; Sankaranarayanan, A.C. One Network to Solve Them  
1525 All—Solving Linear Inverse Problems Using Deep Projection Models. Proceedings of the IEEE Conference  
1526 on Computer Vision and Pattern Recognition, 2017, pp. 5888–5897.
- 1527 158. Liu, H.; Fu, Z.; Han, J.; Shao, L.; Liu, H. Single satellite imagery simultaneous super-resolution and  
1528 colorization using multi-task deep neural networks. *Journal of Visual Communication and Image Representation*  
1529 **2018**, *53*, 20–30.
- 1530 159. Vandal, T.; Kodra, E.; Ganguly, S.; Michaelis, A.; Nemani, R.; Ganguly, A.R. DeepSD: Generating high  
1531 resolution climate change projections through single image super-resolution. Proceedings of the 23rd ACM  
1532 SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2017, pp. 1663–1672.
- 1533 160. Lin, D.Y.; Wang, Y.; Xu, G.L.; Fu, K. Synthesizing remote sensing images by conditional adversarial  
1534 networks. Geoscience and Remote Sensing Symposium (IGARSS), 2017 IEEE International. IEEE, 2017, pp.  
1535 48–50.
- 1536 161. Deng, X.; Zhu, Y.; Newsam, S. What Is It Like Down There? Generating Dense Ground-Level Views and  
1537 Image Features From Overhead Imagery Using Conditional Generative Adversarial Networks. *arXiv*  
1538 *preprint arXiv:1806.05129* **2018**.

- 1539 162. Ji, S.; Zhang, C.; Xu, A.; Shi, Y.; Duan, Y. 3D convolutional neural networks for crop classification with  
1540 multi-temporal remote sensing images. *Remote Sensing* **2018**, *10*, 75.
- 1541 163. Reichstein, M.; Camps-Valls, G.; Stevens, B.; Jung, M.; Denzler, J.; Carvalhais, N.; others. Deep learning  
1542 and process understanding for data-driven Earth system science. *Nature* **2019**, *566*, 195.
- 1543 164. Karpatne, A.; Watkins, W.; Read, J.; Kumar, V. Physics-guided neural networks (pgnn): An application in  
1544 lake temperature modeling. *arXiv preprint arXiv:1710.11431* **2017**.
- 1545 165. Jia, X.; Willard, J.; Karpatne, A.; Read, J.; Zwart, J.; Steinbach, M.; Kumar, V. Physics guided RNNs for  
1546 modeling dynamical systems: A case study in simulating lake temperature profiles. Proceedings of the  
1547 2019 SIAM International Conference on Data Mining. SIAM, 2019, pp. 558–566.

1548 © 2019 by the authors. Submitted to *Journal Not Specified* for possible open access  
1549 publication under the terms and conditions of the Creative Commons Attribution (CC BY) license  
1550 (<http://creativecommons.org/licenses/by/4.0/>).