
Augmented Deep Unrolling Networks for Snapshot Compressive Hyperspectral Imaging

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Snapshot compressive hyperspectral imaging requires reconstructing a hyperspec-
2 tral image from its snapshot measurement. This paper proposes an augmented deep
3 unrolling neural network for solving such a challenging reconstruction problem.
4 The proposed network is based on the unrolling of a proximal gradient descent
5 algorithm with two innovative modules for gradient update and proximal mapping.
6 The gradient update is modeled by a memory-assistant descent module motivated
7 by the momentum-based acceleration heuristics. The proximal mapping is mod-
8 eled by a sub-network with a cross-stage self-attention which effectively exploits
9 inherent self-similarities of a hyperspectral image along the spectral axis, as well
10 as enhancing the feature flow through the network. Moreover, a spectral geometry
11 consistency loss is proposed to encourage the model to concentrate more on the
12 geometric layer of spectral curves for better reconstruction. Extensive experiments
13 on several datasets showed the performance advantage of our approach over the
14 latest methods.

15 1 Introduction

16 Hyperspectral imaging captures a hyperspectral image (HSI) which is a 3D cube of intensities
17 that represents the integrals the radiance of a real scene across a wide range of spectral bands.
18 As an HSI provides rich spectral characteristics of objects of a scene, hyperspectral imaging has
19 found wide applications in many areas, *e.g.*, agriculture, industry, and science. Snapshot compressive
20 spectral imaging [1], often known as coded aperture snapshot spectral imaging (CASSI), is
21 a compressed-sensing-based technique for rapid and efficient acquisition of HSIs. In contrast to
22 traditional hyperspectral imaging techniques which use a sensor array for measuring the object at
23 many spectral bands, the CASSI only collects a single coded 2D snapshot, which measures the
24 object modulated by a physical mask and a disperser at the mixture of different wavelengths. A
25 reconstruction algorithm is then called to reconstruct the 3D HSI from its 2D compressive snapshot.

26 Let $\mathbf{X} \in \mathbb{R}^{M \times N \times \Lambda}$ denote an HSI with spatial indices m, n and spectral index λ . In general, the
27 snapshot from a CASSI device can be expressed as the following [1]:

$$\mathbf{Y}(m, n) = \sum_{\lambda=1}^{\Lambda} \rho(\lambda) \varphi(m - J(\lambda), n) \mathbf{X}(m - J(\lambda), n, \lambda) + \mathbf{N}(m, n), \quad (1)$$

28 where $\rho(\cdot)$ is the spectral response of the camera, $\varphi(\cdot, \cdot)$ the coded aperture pattern, $J(\cdot)$ the dispersive
29 function, and \mathbf{N} the measurement noise. For convenience, we re-express it in a matrix-vector form:

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{n}, \quad (2)$$

30 where Φ denotes the measurement matrix determined by ρ, ψ , and $\mathbf{x}, \mathbf{y}, \mathbf{n}$ are the vectorized form of
31 $\mathbf{X}, \mathbf{Y}, \mathbf{N}$, respectively. As Eq. (2) is an under-determined linear system with measurement noise,
32 HSI reconstruction needs to solve an ill-posed inverse problem,

33 In recent years, deep learning has become a prominent approach for developing powerful solutions to
 34 HSI reconstruction; see *e.g.* [2, 36, 3–6, 3, 7–10]. Most of them models the inverse mapping from
 35 the 2D snapshot to its corresponding HSI by a neural network (NN) trained over a dataset. Among
 36 existing designs of NN architecture, deep unrolling is the most popular one for HSI reconstruction, as
 37 it allows the inclusion of the physics of imaging. A typical deep unrolling network (DUN) unfolds
 38 an iterative scheme for solving some regularized variational model of (2), where the regularization-
 39 related parts are replaced by learnable NN modules. It can also be interpreted as a concatenation of
 40 the steps that alternates between an updating step and a refinement step: $\mathbf{x}^{(0)} \xrightarrow{\text{Update}} \mathbf{z}^{(0)} \xrightarrow{\text{Refine}}$
 41 $\mathbf{x}^{(1)} \xrightarrow{\text{Update}} \mathbf{z}^{(1)} \xrightarrow{\text{Refine}} \mathbf{x}^{(2)} \rightarrow \dots$. Despite extensive studies on HSI reconstruction, the
 42 practical need remains for the methods with better reconstruction accuracy.

43 The paper aims at developing a DUN for HSI reconstruction that brings noticeable performance
 44 improvement over existing deep NNs. The proposed DUN is based on the proximal gradient descent
 45 (PGD) algorithm [11, 12], one often seen iterative scheme for solving inverse problems in imaging.
 46 The PGD algorithm alternatively iterates between the following two steps:

- 47 1. A gradient descent step for updating the estimate of the image
- 48 2. A proximal mapping for refining the estimate via fitting some regularization term.

49 In comparison to existing DUNs for HSI reconstructions, there are three innovations in the design
 50 and training of the proposed one:

- 51 1. Updating step: Modeling the gradient descent step using an NN block with a momentum-
 52 motivated memory-assistant module which is implemented by long short-term memory.
- 53 2. Refinement step: Modeling the proximal mapping by a sub-NN with a across-stage self- attention
 54 module, for exploiting specific characteristics of HSIs and efficient feature flow.
- 55 3. Training loss: A spectral geometry consistency loss is proposed for regularizing the reconstruc-
 56 tion with better accuracy.

57 **Learnable memory-assistant module** In most existing DUNs for HSI reconstruction, the updating
 58 step usually is some pre-defined non-learnable process, *e.g.* gradient-based update. Gradient-based
 59 updates are in a zig-zag direction which slows down the movement to a minima. Also, the updates
 60 crawl near the minima or saddle points slowly as the gradient magnitude vanishes rapidly over there.
 61 A popular technique used for acceleration is the so-called *momentum* (*e.g.* RMSProp and Adam).
 62 Instead of using only the current gradient, momentum accumulates the gradients of the past steps to
 63 determine the direction to go, which helps move more quickly towards the minima as it dampens the
 64 zig-zag oscillations and builds the speed to quicken the convergence.

65 Motivated by the benefit brought by momentum in gradient-based update, we propose to learn an
 66 NN-based model for gradient-based update with the concept of momentum. As the effectiveness of
 67 momentum comes from its memory of the gradients of past steps, we propose an NN block with a
 68 memory-assistant mechanism such that it will leverage the gradient descents from previous stages,
 69 which is implemented using convolutional long short-term memory (ConvLSTM) units.

70 **Cross-stage self-attention module** An HSI has its specific physical characteristics. One is the
 71 self-similarity and strong correlation along the spectral axis, as the entries along the spectral axis
 72 measure the same object region but at different wavelengths. To exploit such specific physical
 73 property of HSIs, we propose a self-attention module along the spectral axis. While self-attention is
 74 not completely new in image reconstruction, our implementation is different from existing ones by
 75 defining in a cross-stage manner.

76 One additional function for such a cross-stage self-attention module is to exploit the similarity of the
 77 features learned over different stages by forming a path between two different stages. Such similarities
 78 among the featured learned at different stages come from the fact that the role of refinement step is
 79 supposed to the same across different stages. The benefit of utilizing such similarity is two-fold. One
 80 is for more efficient feature delivery across the full stages, and the other is for enabling interactions
 81 among the features at different stages during the training.

82 **Loss on spectral geometry consistency** In addition to the standard ℓ_1 loss, a spectral geometry
 83 consistency loss is proposed for training the DUN for HSI reconstruction. Such a loss encourages the

84 model to concentrate more on the profile of spectral changes during reconstruction, which helps to
 85 improve the reconstruction accuracy as empirically observed.

86 2 Related Work

87 By imposing certain priors on HSIs, regularization is a widely-used approach to solving the problem
 88 of HSI reconstruction. The priors for natural images have been extended to HSIs, *e.g.*, sparsity prior
 89 in image gradients used in total variation [13, 14], sparsity prior under a learned dictionary [2, 15],
 90 and non-local self-similarity prior in the form of low-rankness for spatial-spectral patches [16–19].
 91 These pre-defined priors are often insufficient for the HSIs with complex and diverse structures.

92 There is an increasing trend to use the implicit image prior encoded in a pre-trained or untrained
 93 NN for regularization. Plug-and-play methods [14, 20, 21] employ the NNs pre-trained on the
 94 denoising tasks of HSIs or natural images to regularize the reconstruction process. However, pre-
 95 trained denoising NNs are usually not very effective to handle the noise and artifacts generated in the
 96 iterative reconstruction process. Self-supervised learning methods [22, 23] use an untrained NN to
 97 re-parameterize the latent HSI and train it to match the observed snapshot. Such an online learning
 98 scheme is computationally expensive and cannot leverage the knowledge from external data.

99 It has been a prominent approach that to end-to-end train a DNN that maps a snapshot to the latent
 100 HSI; see *e.g.* [24, 5, 25, 26, 9, 8]. Many existing studies employ the DUN architecture *e.g.* [3, 6, 7, 4].
 101 Recall that a DUN often consists of pairs of steps: one step for updating the estimate of the latent HSI
 102 and the other step for refining the estimate with a learnable prior. Most existing works focus on the
 103 latter, which can be viewed as a denoising NN that exploits different image priors, *e.g.*, spatial-spectral
 104 prior [3], non-local self-similarity prior [6], and patch-level Gaussian scale mixture prior [7].

105 **Learning updating steps in DUNs** Zhang *et al.* [4] replaced the operators Φ , Φ^\top appearing in the
 106 gradient descent step of PGD by convolutions and residual blocks, with a channel attention block to
 107 estimate the step size in PGD from the estimate output by the previous stage. Different from that, we
 108 do not learn those operators but utilize them to have a better update step. Working on natural image
 109 recovery rather than on HSI reconstruction, Mou *et al.* [27] used a residual block to estimate the
 110 gradient descent step. In comparison, we use an LSTM to leverage the dependency between different
 111 stages for estimating the updating step.

112 **Self-attention for HSI reconstruction** Self-attention (SA) has been exploited in existing works for
 113 HSI reconstruction. Miao *et al.* [5] used a generative adversarial network with SA for the initial stage
 114 in the NN. Meng *et al.* [28] used three spatial-spectral SA modules to exploit the spatial-spectral
 115 correlation of an HSI. Hu *et al.* [9] develops a spatial-spectral attention module with efficient feature
 116 fusion. In comparison to these methods, ours treats spectral maps as tokens for SA and calculates the
 117 SA along the spectral dimension. This shares a similar idea with a parallel work [8] which also treats
 118 spectral maps as tokens in a transformer-based model. Different from it, we use SA in a cross-stage
 119 manner which enhances the feature flow at the same time.

120 **Training loss for HSI reconstruction** Most existing NNs for HSI reconstruction are trained by
 121 the standard mean-squared-error loss or ℓ_1 loss. Hu *et al.* [9] introduced a frequency-domain loss
 122 to narrow the frequency-domain discrepancy between network predictions and ground truths. In
 123 comparison, the loss we proposed narrows the discrepancy in terms of spectral geometric changes.

124 3 Proposed Approach

125 The proposed DUN for HSI reconstruction is based on the PGD algorithm [11, 12] for the following
 126 optimization model regularized by the functional \mathcal{R} :

$$\min_{\mathbf{x}} \|\mathbf{y} - \Phi\mathbf{x}\|_2^2 + \lambda\mathcal{R}(\mathbf{x}), \quad \lambda \in \mathbb{R}^+, \quad (3)$$

127 The PGD algorithm for solving Eq. (3) alternately iterates between two steps: gradient-descent (GD)
 128 step for updating the estimate, and proximal mapping (PM) step for refining the estimate by fitting
 129 the functional \mathcal{R} with encoded image prior: For $k = 1, \dots, K$,

$$[\text{GD}]: \quad \mathbf{u}^{(k)} = \mathbf{x}^{(k-1)} + \gamma^{(k)} \Phi^\top (\mathbf{y} - \Phi\mathbf{x}^{(k-1)}), \quad (4)$$

$$[\text{PM}]: \quad \mathbf{x}^{(k)} = \text{Prox}_{\mathcal{R}}(\mathbf{u}^{(k)}) \triangleq \underset{\mathbf{x}'}{\text{argmin}} \|\mathbf{x} - \mathbf{u}^{(k)}\|_2^2 + 2\gamma^{(k)}\mathcal{R}(\mathbf{u}^{(k)}). \quad (5)$$

130 where $\gamma^{(k)}$ denotes step size. Most existing DUNs focus on modeling the PM step (5) by an NN for a
 131 data-driven prior. The GD step (4) usually is kept unchanged with the learnable parameter $\gamma^{(k)}$.

132 We propose a Memory-Assistant Descent (MAD) block to model the GD step (4) and a Cross-stage
 133 Attentive Proximal (CAP) sub-network to model the PM step (5). The former functions as gradient
 134 descent across different stages for momentum-motivated acceleration, which leads to a more efficient
 135 update than that only using the gradient at current stage. The latter is to utilize the self-similarities
 136 existing in an HSI with a cross-stage manner, which enable us to exploit special characteristics of
 137 HSIs and fasten feature flow through the DNN. In short, the proposed NN, called MadcapNet, is the
 138 concatenation of K stages, each of which contains a pair of a MAD block and a CAP sub-network;
 139 see Figure 1 for the diagram of MadcapNet.

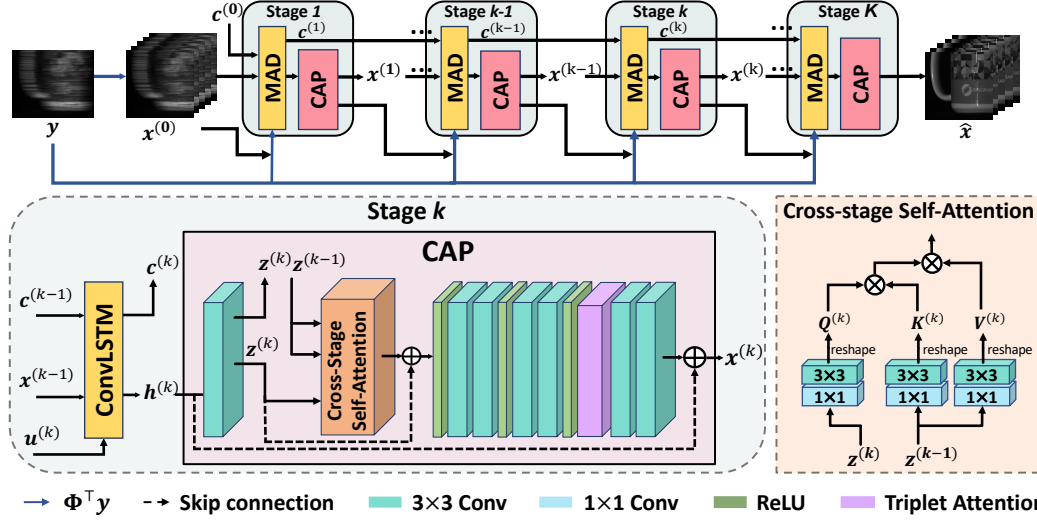


Figure 1: Diagram of the proposed augmented deep unrolling neural network for HSI reconstruction.

140 3.1 Memory-Assistant Descent Blocks

141 The MAD blocks are a set of ConvLSTM units [29] placed at each stage of the NN, which utilizes
 142 the long-range dependencies among all cascading stages for momentum-assistant gradient update. In
 143 each MAD block, the gradient map is defined by

$$\mathbf{u}^{(k)} = \Phi^\top (\mathbf{y} - \Phi \mathbf{x}^{(k-1)}) \quad (6)$$

144 is taken as the input for the k -th ConvLSTM unit, which introduces information on gradient descent.
 145 Let $\mathbf{h}^{(k)}$, $\mathbf{c}^{(k)}$ denote the hidden state and cell state in the ConvLSTM at the k -th stage respectively,
 146 where $\mathbf{h}^{(k)}$ is of the same size as $\mathbf{x}^{(k)}$. The MAD block is defined as

$$[\mathbf{h}^{(k)}, \mathbf{c}^{(k)}] = \text{ConvLSTM}(\mathbf{u}^{(k)}, \mathbf{x}^{(k-1)}, \mathbf{c}^{(k-1)}), \quad (7)$$

147 for $k = 1, \dots, K$. Different from original ConvLSTM units which use the previous hidden state
 148 $\mathbf{h}^{(k-1)}$ as input, we replace $\mathbf{h}^{(k-1)}$ by $\mathbf{x}^{(k-1)}$, the output from the CAP sub-network of the previous
 149 stage. The motivation behind is to utilize the current gradient descent defined over $\mathbf{x}^{(k-1)}$. Then, $\mathbf{h}^{(k)}$
 150 is used as the input of the CAP sub-network and $\mathbf{c}^{(k)}$ is fed to the MAD block at the next stage as an
 151 accumulator of state information.

152 In the k -th stage, the ConvLSTM unit calculates \mathbf{h}_k , \mathbf{c}_k by the following rules

$$\mathbf{c}^{(k)} = \mathbf{f}_k \odot \mathbf{c}^{(k-1)} + \mathbf{i}^{(k)} \odot \tanh(\mathbf{g}^{(k)}), \quad (8)$$

$$\mathbf{h}^{(k)} = \mathbf{o}^{(k)} \odot \tanh(\mathbf{c}^{(k)}), \quad (9)$$

153 where \odot denotes Hadamard product, and \mathbf{i}_k , \mathbf{f}_k , \mathbf{o}_k , \mathbf{g}_k denote the input gate, forget gate, output
 154 gate, and the intermediate result, respectively, which are calculated as follows:

$$\mathbf{i}^{(k)} = \text{sigmoid}(\mathbf{W}_{\text{mi}}\mathbf{u}^{(k)} + \mathbf{W}_{\text{xi}}\mathbf{x}^{(k-1)} + \mathbf{b}_i), \quad (10)$$

$$\mathbf{f}^{(k)} = \text{sigmoid}(\mathbf{W}_{\text{mf}}\mathbf{u}^{(k)} + \mathbf{W}_{\text{xf}}\mathbf{x}^{(k-1)} + \mathbf{b}_f), \quad (11)$$

$$\mathbf{g}^{(k)} = \mathbf{W}_{\text{mg}}\mathbf{u}^{(k)} + \mathbf{W}_{\text{xg}}\mathbf{x}^{(k-1)} + \mathbf{b}_g, \quad (12)$$

$$\mathbf{o}^{(k)} = \text{sigmoid}(\mathbf{W}_{\text{mo}}\mathbf{u}^{(k)} + \mathbf{W}_{\text{xo}}\mathbf{x}^{(k-1)} + \mathbf{b}_o), \quad (13)$$

155 where \mathbf{W}_{**} are implemented by 3×3 convolutional layers with bias terms \mathbf{b}_* .

156 3.2 Cross-stage Attentive Proximal Sub-networks

157 The CAP blocks function as a learnable PM step (5) which refines the estimate from the MAD
 158 block. It can be understood as a denoising NN by interpreting the estimation residual as noise. Given
 159 $\mathbf{h}^{(k)}$ (of the same size as \mathbf{x}) from the MAD block as input, we map it to a feature tensor $\mathbf{z}^{(k)}$ via a
 160 convolutional layer, which is then processed by a cross-stage SA module. Afterward, the results are
 161 fed to a sequence of convolutional layers with rectified linear units (ReLUs) and a triplet attention [30].
 162 The output with the same size as \mathbf{x} is combined with the input $\mathbf{h}^{(k)}$ via a skip connection, yielding
 163 the reconstructed HSI $\mathbf{x}^{(k)}$ at the current stage. See Figure 1 for the details.

164 Recall that SA [31] relates input feature tokens to compute a refined feature representation. It first
 165 generates a key/query/value vector of length d from each token, and all the key/query/value vectors
 166 are stored as \mathbf{K} , \mathbf{Q} , \mathbf{V} respectively. Then, SA is calculated as follows:

$$\text{SA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{1}{\sqrt{d}}\mathbf{Q}\mathbf{K}^\top\right)\mathbf{V}. \quad (14)$$

167 We treat each feature channel as a token so as to exploit the self-similarities among feature channels.
 168 Such tokens are aligned due to natural alignment of spectral slices of an HSI. In the k th stage, rather
 169 than use the feature $\mathbf{z}^{(k)}$ at current stage to calculate $\mathbf{K}^{(k)}$, $\mathbf{Q}^{(k)}$, $\mathbf{V}^{(k)}$, we only use $\mathbf{z}^{(k)}$ for $\mathbf{Q}^{(k)}$
 170 while using the feature $\mathbf{z}^{(k-1)}$ of previous stage for $\mathbf{K}^{(k)}$, $\mathbf{V}^{(k)}$. Concretely, we calculate

$$\mathbf{Q}^{(k)} = \mathbf{W}_{\text{Qd}}^{(k)}\mathbf{W}_{\text{Qp}}^{(k)}\mathbf{z}^{(k)}, \mathbf{K}^{(k)} = \mathbf{W}_{\text{Kd}}^{(k)}\mathbf{W}_{\text{Kp}}^{(k)}\mathbf{z}^{(k-1)}, \mathbf{V}^{(k)} = \mathbf{W}_{\text{Vd}}^{(k)}\mathbf{W}_{\text{Vp}}^{(k)}\mathbf{z}^{(k-1)}, \quad (15)$$

171 where $\mathbf{W}_{(*p)}^{(k)}$, $\mathbf{W}_{(*d)}^{(k)}$ are 1×1 convolutions and 3×3 depth-wise convolutions respectively for better
 172 encoding spatial-channel context.

173 The motivation of the cross-stage strategy is as follows. The DUN architecture alternates between
 174 the update and the refinement. Since the CAP sub-networks at different stages play the same role
 175 of refinement, their extracted features should be highly correlated and the features extracted from
 176 the previous stage provide good initials for the corresponding ones at the next stage. However, the
 177 aforementioned pipeline does not utilize such correlations for more efficient training, which may
 178 result in a bottleneck for features flowing through the whole DUN. The proposed cross-stage SA
 179 scheme forms a path between two stages, which allows efficient feature transmission during inference
 180 and enhances feature interactions during training.

181 The mutli-head strategy [31] is adopted for the cross-stage SA. First, we split the key/query/value ma-
 182 trices into H heads along channel dimension: $\mathbf{Q}^{(k)} = [\mathbf{Q}_1^{(k)}, \dots, \mathbf{Q}_H^{(k)}]$, $\mathbf{K}^{(k)} = [\mathbf{K}_1^{(k)}, \dots, \mathbf{K}_H^{(k)}]$,
 183 and $\mathbf{V}^{(k)} = [\mathbf{V}_1^{(k)}, \dots, \mathbf{V}_H^{(k)}]$. Then, the output is calculated as

$$\mathbf{O}^{(k)} = \cup_{j=1}^H \text{SA}(\mathbf{Q}_j^{(k)}, \mathbf{K}_j^{(k)}\mathbf{V}_j^{(k)}), \quad (16)$$

184 which is reshaped for subsequent processing.

185 3.3 Loss function for Training

186 To better train a NN for HSI reconstruction, we propose an additional loss called spectral geometry
 187 consistency (SGC) loss. For an HSI $\mathbf{X} \in \mathbb{R}^{M \times N \times \Lambda}$, we define the geometry map $\mathcal{D}(\mathbf{x})$ as follows.

$$\mathcal{D}(\mathbf{X}) = \nabla_c(\text{sign}(\nabla_c\mathbf{X})) \in \{-1, 0, 1\}^{M \times N \times \Lambda}, \quad (17)$$

188 where ∇_c calculates the gradient along the spectral axis, and $\text{sign}(\cdot)$ denotes element-wise sign
 189 function. For a spatial location (m_0, n_0) , $\mathcal{D}(\mathbf{X})[m_0, n_0, \cdot]$ indicates the wavelengths where the
 190 monotony of spectral values changes, which is one geometrical property of the spectral curve. Based
 191 on \mathcal{D} , the SGC loss emphasizes the geometrical layout consistency between the reconstructed HSI
 192 and ground truth.

193 Considering HSIs exhibit high spatial sparsity, the irrelevant dark regions are omitted for robustness.
 194 This is achieved by constructing a mask $\mathbf{M}_\mathbf{X}$ that thresholds the max density along spectral dimension:
 195 $\mathbf{M}_\mathbf{X}(m, n, \lambda) = 1$ if $\max_\lambda \mathbf{X}(m, n, \lambda) \geq \alpha$; and 0 otherwise. Let $\mathbf{X}, \widehat{\mathbf{X}}$ denote the reconstructed
 196 HSI and its ground truth respectively. The SGC loss is defined as

$$\mathcal{L}_{\text{sgc}} \triangleq \|\mathbf{M}_\mathbf{X} \odot \mathcal{D}(\mathbf{X}) - \mathbf{M}_{\widehat{\mathbf{X}}} \odot \mathcal{D}(\widehat{\mathbf{X}})\|_1. \quad (18)$$

197 By minimizing \mathcal{L}_{sgc} , the HSI predicted by the NN is biased to the one with the same wavelength-
 198 density trends of ground truths, which helps to alleviate possible over-fitting. Then, the overall loss is
 199

$$\mathcal{L} \triangleq \mathcal{L}_1 + \gamma \mathcal{L}_{\text{sgc}} = \|\mathbf{X} - \widehat{\mathbf{X}}\|_1 + \gamma \|\mathbf{M}_\mathbf{X} \odot \mathcal{D}(\mathbf{X}) - \mathbf{M}_{\widehat{\mathbf{X}}} \odot \mathcal{D}(\widehat{\mathbf{X}})\|_1, \quad \gamma \in \mathbb{R}^+. \quad (19)$$

200 4 Experiments

201 We implement MadcapNet with PyTorch. The stage number K is set to 6. On all convolutional
 202 layers, the kernel sizes are all set to 3×3 , and both the stride and padding number are set to 1. The
 203 head number H for the self-attention in CAP blocks is set to 8. Regarding the training loss, we
 204 set $\alpha = \frac{5}{255}$ for $\mathbf{M}_\mathbf{X}$ and $\gamma = 0.5$ for Eq. (19). The training is done via the Adam optimizer with
 205 a fixed learning rate of 10^{-4} and a maximal epoch number of 200. The same data augmentation
 206 scheme as [7] is adopted, including rotation and flipping. All the models are trained and tested
 207 on an NVIDIA GeForce RTX 1080Ti GPU. Our code will be released on GitHub. upon paper’s
 208 acceptance. Following [7], Peak-Signal-to-Noise-Ratio (PSNR) and Structured SIMilarity (SSIM)
 209 index are adopted as the metrics to evaluate the reconstruction results quantitatively.

210 4.1 Evaluation on Synthetic Data

211 **CAVE and KAIST datasets** Following [28, 7], we use the CAVE dataset [32] containing 32 HSIs
 212 with 31 spectral bands for training, and 10 scenes with 31 spectral bands from the KAIST dataset [14]
 213 for test. All these HSIs are cropped into patches with a spatial size of 256×256 and reduced to 28
 214 wavelengths ranging from 450nm to 650nm via spectral interpolation. The snapshot measurements
 215 are generated by the 256×256 mask of CASSI used in [28].

216 Ten existing methods are chosen for comparison, including (a) two conventional methods: GAP-
 217 TV [13] and DeSCI [17]; (b) one self-supervised learning-based method: PnP-DIP [22]; and (c)
 218 seven supervised learning-based methods: λ -Net [5] HSSP [3], DNU [6], TSA-Net [28], DGSM [7],
 219 HDNet [9], and MST-L [8]. The HSSP, DNU and DGSM are based on DUNs. The HDNet and
 220 MST-L are from two latest works accepted in an upcoming conference.

221 The quantitative results are listed in Table 1, which are quoted from [8, 9] whenever possible and
 222 otherwise obtained with released codes. It can be seen that our approach significantly outperforms the
 223 compared ones. Specifically, MadcapNet shows remarkable superior performance over other DUNs.
 224 It also surpasses MST-L and HDNet (*i.e.* two latest methods) with an average PSNR gain of more
 225 than 1dB and 2dB respectively. Table 1 also compares the model complexity of different methods in
 226 terms of number of parameters and number of Giga Floating-point Operations Per Second (GFLOPS).
 227 Although our model contains ConvLSTM and self-attention blocks, it is still kept compact to maintain
 228 a relatively-low model complexity. Among all compared methods, our MadcapNet has the smallest
 229 number of GLOPS, and it is smaller than all other models except DNU. These results show the
 230 practicability of MadcapNet for real applications. To conclude, our approach can achieve the best
 231 trade-off between performance and model complexity.

232 **ICVL and Harvard datasets** We also conduct experiments on the ICVL dataset [33] and the
 233 Harvard dataset [34], respectively. The ICVL dataset consists of 201 HSIs of real-world objects, each
 234 with 31 spectral bands collected from 400nm to 700 nm at a 10nm step. The Harvard dataset consists
 235 of 50 outdoor scenes, each with 31 spectral bands collected from 420nm to 720nm at a 10nm step.

Table 1: Quantitative results in PSNR(dB) (even rows) and SSIM (odd rows) on KAIST dataset.

Method	#Param.	#GFLOPS	Scene#1	#2	#3	#4	#5	#6	#7	#8	#9	#10	Mean
GAP-TV	-	-	26.82	22.89	26.31	30.65	23.64	21.85	23.76	21.98	22.63	23.10	24.36
			0.754	0.61	0.802	0.852	0.703	0.663	0.688	0.655	0.682	0.584	0.669
DeSCI	-	-	27.13	23.04	26.62	34.96	23.94	22.38	24.45	22.03	24.56	23.59	25.27
			0.748	0.62	0.818	0.897	0.706	0.683	0.743	0.673	0.732	0.587	0.721
λ -net	62.64M	117.98	30.10	28.49	27.73	37.01	26.19	28.64	26.47	26.09	27.50	27.13	28.53
			0.849	0.805	0.870	0.934	0.817	0.853	0.806	0.831	0.826	0.816	0.841
HSSP	-	-	31.48	31.09	28.96	34.56	28.53	30.83	28.71	30.09	30.43	28.78	30.35
			0.858	0.842	0.823	0.902	0.808	0.877	0.824	0.881	0.868	0.842	0.852
DNU	1.19M	163.48	31.72	31.13	29.99	35.34	29.03	30.87	28.99	30.13	31.03	29.14	30.74
			0.863	0.846	0.845	0.908	0.833	0.887	0.839	0.885	0.876	0.849	0.863
PnP-DIP	33.85M	64.42	32.68	27.26	31.30	40.54	29.79	30.39	28.18	29.44	34.51	28.51	31.26
			0.890	0.833	0.914	0.962	0.900	0.877	0.913	0.874	0.927	0.851	0.894
TSA-Net	44.25M	110.06	32.03	31.00	32.25	39.19	29.39	31.44	30.32	29.35	30.01	29.59	31.46
			0.892	0.858	0.915	0.953	0.884	0.908	0.878	0.888	0.890	0.874	0.894
DGSMF	3.76M	646.65	33.26	32.09	33.06	40.54	28.86	33.08	30.74	31.55	31.66	31.44	32.63
			0.915	0.898	0.925	0.964	0.882	0.937	0.886	0.923	0.911	0.925	0.917
HDNet	2.35M	154.00	34.95	32.52	34.52	43.00	32.49	35.96	29.18	34.00	34.56	32.22	34.34
			0.948	0.953	0.957	0.981	0.957	0.965	0.937	0.961	0.958	0.950	0.957
MST-L	2.03M	28.15	35.40	35.87	36.51	42.27	32.77	34.80	33.66	32.67	35.39	32.50	35.18
			0.941	0.944	0.953	0.973	0.947	0.955	0.925	0.948	0.949	0.941	0.948
<u>MadcapNet</u>	1.51M	24.24	36.24	37.49	37.07	42.85	34.09	35.61	35.37	33.96	36.67	33.12	36.32
			0.951	0.961	0.963	0.981	0.962	0.966	0.949	0.962	0.960	0.948	0.961

236 Following the protocol of [3, 35], 50 HSIs in the ICVL dataset and 9 HSIs in the Harvard dataset are
 237 used for test respectively, and the rest samples for training. All HSIs for training and test are cropped
 238 into patches with a spatial size of 48×48 , while keeping the band number unchanged. The snapshot
 239 measurements are generated by the 48×48 mask of CASSI used in [3].

240 Six existing methods are selected for comparison, including (a) a conventional method: SSNR [16];
 241 and (b) six supervised learning-based methods: HSCNN [36], λ -Net [5], DNU [6], DTLP [37], and
 242 HDNet [9]. The DNU and DTLP use DUNs, and the HDNet is a latest method.

243 See Table 2 for the quantitative comparison. The results of the compared methods are cited from [37].
 244 The proposed one outperformed all other methods, with more than 0.85db PSNR improvement on
 245 both datasets. Such noticeable performance gains of MadcapNet over other DUNs again demonstrated
 246 the effectiveness of our network architecture.

Table 2: Quantitative results in PSNR(dB) and SSIM on ICVL and Harvard datasets.

Dataset	Metric	SSNR	HSCNN	λ -Net	DNU	DTLP	HDNet	<u>MadcapNet</u>
ICVL	PSNR	30.40	28.45	29.01	32.61	34.53	36.38	37.60
	SSIM	0.943	0.934	0.946	0.966	0.977	0.981	0.985
Harvard	PSNR	31.14	27.60	29.37	31.11	32.43	34.02	34.88
	SSIM	0.942	0.895	0.909	0.929	0.941	0.950	0.956

247 **Visual inspection** See Figure 2 for the visualization of HSI reconstruction results on two samples
 248 from the KAIST and Harvard datasets respectively. The spectral curves (density versus wavelength)
 249 correspond to the points marked by green boxes in the RGB references. In the legends of both
 250 figures, we provide the curve correlation value between the result of a compared method and the
 251 ground truth. Those values show that the HSIs reconstructed by the proposed MadcapNet have
 252 the highest correlation to the ground truths. We also visualize three spectral channels of an entire
 253 reconstructed HSI and zoom in the selected regions marked by yellow boxes. It can be seen that the
 254 results of MadcapNet are more visually pleasing than that of other compared methods, with a better
 255 reconstruction of structures.

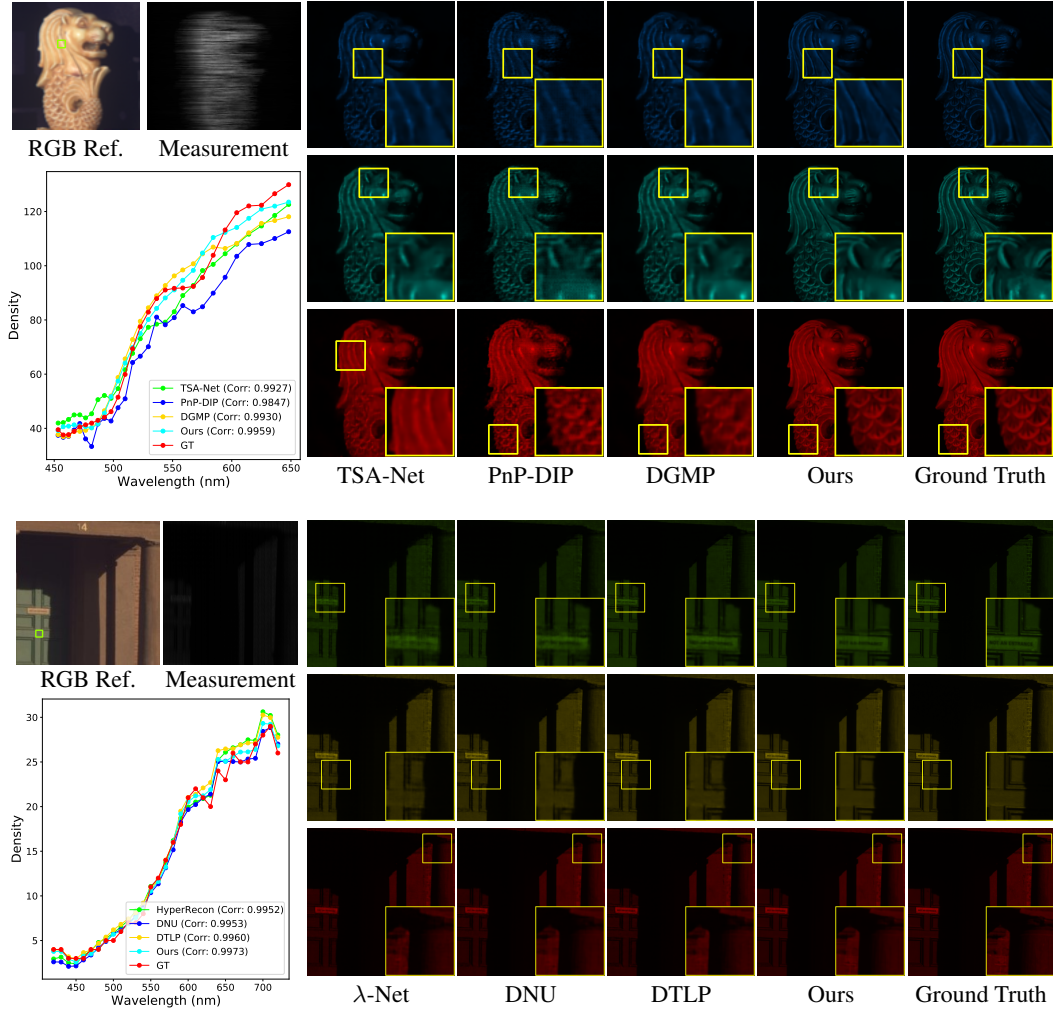


Figure 2: Visual comparison of HSI reconstruction on two samples from KAIST and Harvard datasets respectively. Left: spectra curves of the selected regions marked by green boxes. Right: reconstruction on the spectral channels.

256 4.2 Evaluation on Real Data

257 We also conduct a test on the real snapshots of spatial size 660×714 from [7, 28], which are captured
 258 by a real system with 28 wavelengths ranging from 450nm to 650nm and with 54-pixel dispersion
 259 in the column dimension. Following [7, 28], we use the mask associated with that real system to
 260 generate snapshots on both the CAVE and KAIST datasets, and then we inject 11-bit shot noise to the
 261 snapshots for simulating real situations. The resulting data is used to retrain our model. Due to the
 262 lack of ground truths in test data, we only compare the qualitative results of different methods. See
 263 Figure 3 for the reconstruction results on a real scene, and see more in the supplementary materials.
 264 The performance of MadcapNet is also good on the real data. This indeed has demonstrated the good
 265 generalization performance of our model.

266 4.3 Ablation Studies

267 Ablation studies are conducted on the KAIST dataset. We form some baselines by removing one or
 268 more main components of our approach. Concretely, we consider (a) replace the MAD blocks by
 269 the GD steps (4); (b) replace the cross-stage SA in the CAP network with the inner-stage SA which
 270 uses the features at current stage to calculate $\mathbf{K}^{(k)}$, $\mathbf{Q}^{(k)}$, $\mathbf{V}^{(k)}$ in (15); (c) replace the cross-stage SA

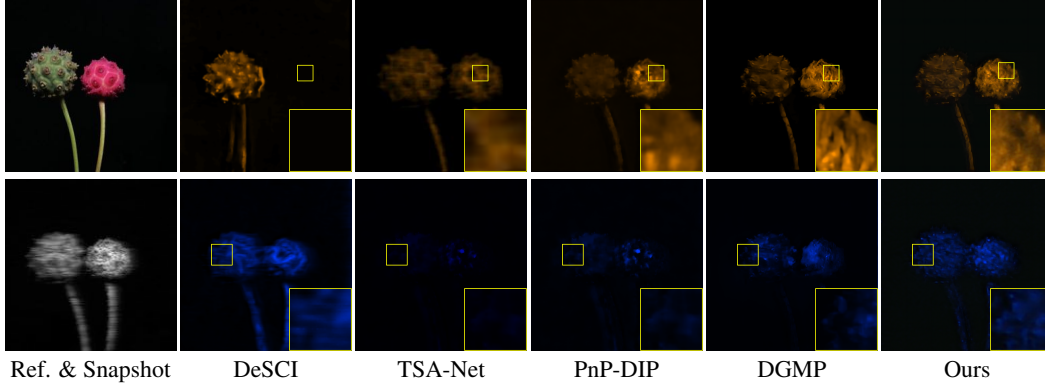


Figure 3: Visual comparison of HSI reconstruction on real data, in terms of two spectral channels.

271 with a same number of convolutional layers; (d) remove the SGC loss \mathcal{L}_{sgc} . For a fair comparison,
 272 each baseline is configured to have (nearly) the same number of parameters as the original model, by
 273 uniformly increasing the channel numbers of convolutional layers. The results are listed in Table 3.

274 It can be seen that each main component in our approach plays an important role. Using the MAD
 275 blocks as an alternate to GD steps can improve PSNR by almost 1db. It also brings improvement across
 276 all baseline settings. Benefiting from the power of SA, the cross-stage SA brings noticeable PSNR
 277 gain. In addition, the SA utilized in the cross-stage manner leads to around 0.36dB improvement in
 278 PSNR over that utilized in the inner-stage manner. The SGC loss also has a solid contribution to the
 279 performance. See Figure 4 for an illustration of the effect of the SGC loss, where training with \mathcal{L}_{sgc}
 280 makes the tendency of the predicted spectral curves closer to ground truths. See also supplementary
 281 materials for more results.

Table 3: Results in ablation studies on KAIST dataset.

Metric	w/o MAD	w/o CAP	Cross→Inner	w/o \mathcal{L}_{sgc}	Original
PSNR(dB)	35.35	35.80	35.96	35.53	36.32
SSIM	0.947	0.956	0.958	0.951	0.961

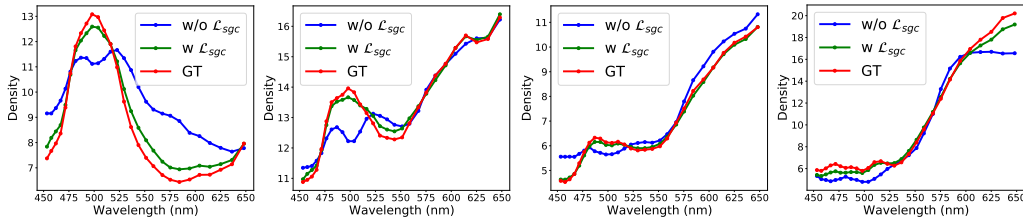


Figure 4: Spectra of selected regions on Scene#1 (first two) and Scene#5 (last two) of KAIST dataset.

282 5 Conclusion

283 In this paper, we proposed an augmented DUN for CASSI-based hyperspectral imaging. The proposed
 284 DUN is based on the unfolding of PGD, with three-fold augmentations: momentum-motivated
 285 ConvLSTM-assistant module for improving the gradient descent steps, a sub-network with cross-stage
 286 self-attention for exploiting self-similarities of an HSI and enhancing feature flow simultaneously,
 287 and a loss to induce predictions biased to spectral geometry consistency. The combination of these
 288 augmentations leads to noticeable performance improvement in HSI reconstruction, which were
 289 demonstrated by extensive experiments. The proposed DUN also sees its potential application to
 290 other compressive imaging problems. We will study it in the future.

References

- 291
292 [1] Michael E Gehm, Renu John, David J Brady, Rebecca M Willett, and Timothy J Schulz. Single-shot
293 compressive spectral imaging with a dual-disperser architecture. *Opt. Lett.*, 15(21):14013–14027, 2007.
- 294 [2] Xing Lin, Yebin Liu, Jiamin Wu, and Qionghai Dai. Spatial-spectral encoded compressive hyperspectral
295 imaging. *ACM Trans. Graph.*, 33(6):1–11, 2014.
- 296 [3] Lizhi Wang, Chen Sun, Ying Fu, Min H Kim, and Hua Huang. Hyperspectral image reconstruction using a
297 deep spatial-spectral prior. In *Proc. CVPR*, pages 8032–8041, 2019.
- 298 [4] Xuanyu Zhang, Yongbing Zhang, Ruiqin Xiong, Qilin Sun, and Jian Zhang. Herosnet: Hyperspectral
299 explicable reconstruction and optimal sampling deep network for snapshot compressive imaging. *Proc.*
300 *CVPR*, 2022.
- 301 [5] Xin Miao, Xin Yuan, Yunchen Pu, and Vassilis Athitsos. λ -net: Reconstruct hyperspectral images from a
302 snapshot measurement. In *Proc. CVPR*, pages 4059–4069, 2019.
- 303 [6] Lizhi Wang, Chen Sun, Maoqing Zhang, Ying Fu, and Hua Huang. Dnu: deep non-local unrolling for
304 computational spectral imaging. In *Proc. CVPR*, pages 1661–1671, 2020.
- 305 [7] Tao Huang, Weisheng Dong, Xin Yuan, Jinjian Wu, and Guangming Shi. Deep gaussian scale mixture
306 prior for spectral compressive imaging. In *Proc. CVPR*, pages 16216–16225, 2021.
- 307 [8] Yuanhao Cai, Jing Lin, Xiaowan Hu, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc
308 Van Gool. Mask-guided spectral-wise transformer for efficient hyperspectral image reconstruction. *Proc.*
309 *CVPR*, 2022.
- 310 [9] Xiaowan Hu, Yuanhao Cai, Jing Lin, Haoqian Wang, Xin Yuan, Yulun Zhang, Radu Timofte, and Luc
311 Van Gool. Hdnet: High-resolution dual-domain learning for spectral compressive imaging. *Proc. CVPR*,
312 2022.
- 313 [10] Ying Fu, Tao Zhang, Lizhi Wang, and Hua Huang. Coded hyperspectral image reconstruction using deep
314 external and internal learning. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.
- 315 [11] Patrick L Combettes and Jean-Christophe Pesquet. Proximal splitting methods in signal processing. In
316 *Fixed-point algorithms for inverse problems in science and engineering*, pages 185–212. Springer, 2011.
- 317 [12] Atsushi Nitanda. Stochastic proximal gradient descent with acceleration techniques. *Proc. NeurIPS*, 27,
318 2014.
- 319 [13] Xin Yuan. Generalized alternating projection based total variation minimization for compressive sensing.
320 In *Proc. ICIP*, pages 2539–2543. IEEE, 2016.
- 321 [14] Inchang Choi, Daniel S. Jeon, Giljoo Nam, Diego Gutierrez, and Min H. Kim. High-quality hyperspectral
322 reconstruction using a spectral prior. *ACM Trans. Graph.*, 36(6):218:1–13, 2017.
- 323 [15] Xin Yuan, Tsung-Han Tsai, Ruoyu Zhu, Patrick Lull, David Brady, and Lawrence Carin. Compressive
324 hyperspectral imaging with side information. *IEEE J. Sel. Top Signal Process*, 9(6):964–976, 2015.
- 325 [16] Ying Fu, Yinqiang Zheng, Imari Sato, and Yoichi Sato. Exploiting spectral-spatial correlation for coded
326 hyperspectral image restoration. In *Proc. CVPR*, pages 3727–3736, 2016.
- 327 [17] Yang Liu, Xin Yuan, Jinli Suo, David J Brady, and Qionghai Dai. Rank minimization for snapshot
328 compressive imaging. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(12):2990–3006, 2018.
- 329 [18] Shipeng Zhang, Lizhi Wang, Ying Fu, Xiaoming Zhong, and Hua Huang. Computational hyperspectral
330 imaging based on dimension-discriminative low-rank tensor recovery. In *Proc. CVPR*, pages 10183–10192,
331 2019.
- 332 [19] Yong Chen, Ting-Zhu Huang, Wei He, Naoto Yokoya, and Xi-Le Zhao. Hyperspectral image compressive
333 sensing reconstruction using subspace-based nonlocal tensor ring decomposition. *IEEE Transactions on*
334 *Image Processing*, 29:6813–6828, 2020.
- 335 [20] Xin Yuan, Yang Liu, Jinli Suo, and Qionghai Dai. Plug-and-play algorithms for large-scale snapshot
336 compressive imaging. In *Proc. CVPR*, pages 1447–1457, 2020.
- 337 [21] Haiquan Qiu, Yao Wang, and Deyu Meng. Effective snapshot compressive-spectral imaging via deep
338 denoising and total variation priors. In *Proc. CVPR*, pages 9127–9136, 2021.

- 339 [22] Ziyi Meng, Zhenming Yu, Kun Xu, and Xin Yuan. Self-supervised neural networks for spectral snapshot
340 compressive imaging. *Proc. CVPR*, 2021.
- 341 [23] Yuhui Quan, Xinran Qin, Mingqin Chen, and Yan Huang. High-quality self-supervised snapshot hyper-
342 spectral imaging. In *proc. ICASSP*, pages 1526–1530. IEEE, 2022.
- 343 [24] Lizhi Wang, Tao Zhang, Ying Fu, and Hua Huang. Hyperreconnet: Joint coded aperture optimization and
344 image reconstruction for compressive hyperspectral imaging. *IEEE Trans. Image Process.*, 28(5):2257–
345 2270, 2018.
- 346 [25] Tao Zhang, Ying Fu, Lizhi Wang, and Hua Huang. Hyperspectral image reconstruction using deep external
347 and internal learning. In *Proc. CVPR*, pages 8559–8568, 2019.
- 348 [26] Kouhei Yorimoto and Xian-Hua Han. Hypermixnet: Hyperspectral image reconstruction with deep mixed
349 network fryom a snapshot measurement. In *Proc. ICCV Workshop*, pages 1184–1193, 2021.
- 350 [27] Chong Mou, Qian Wang, and Jian Zhang. Deep generalized unfolding networks for image restoration.
351 *Proc. CVPR*, 2022.
- 352 [28] Ziyi Meng, Jiawei Ma, and Xin Yuan. End-to-end low cost compressive spectral imaging with spatial-
353 spectral self-attention. In *Proc. ECCV*, pages 187–204. Springer, 2020.
- 354 [29] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Con-
355 volutional lstm network: A machine learning approach for precipitation nowcasting. *Proc. NeurIPS*, 28,
356 2015.
- 357 [30] Diganta Misra, Trikey Nalamada, Ajay Uppili Arasanipalai, and Qibin Hou. Rotate to attend: Convolutional
358 triplet attention module. In *Proc. WACV*, pages 3139–3148, 2021.
- 359 [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
360 Kaiser, and Illia Polosukhin. Attention is all you need. *Proc. NeurIPS*, 30, 2017.
- 361 [32] Fumihito Yasuma, Tomoo Mitsunaga, Daisuke Iso, and Shree K Nayar. Generalized assorted pixel
362 camera: postcapture control of resolution, dynamic range, and spectrum. *IEEE Trans. Image Process.*,
363 19(9):2241–2253, 2010.
- 364 [33] Boaz Arad and Ohad Ben-Shahar. Sparse recovery of hyperspectral signal from natural rgb images. In
365 *Proc. ECCV*, pages 19–34. Springer, 2016.
- 366 [34] Ayan Chakrabarti and Todd Zickler. Statistics of real-world hyperspectral images. In *Proc. CVPR*, pages
367 193–200. IEEE, 2011.
- 368 [35] Shipeng Zhang, Lizhi Wang, Lei Zhang, and Hua Huang. Learning tensor low-rank prior for hyperspectral
369 image reconstruction. In *Proc. CVPR*, pages 12006–12015, 2021.
- 370 [36] Zhiwei Xiong, Zhan Shi, Huiqun Li, Lizhi Wang, Dong Liu, and Feng Wu. Hscnn: Cnn-based hyperspectral
371 image recovery from spectrally undersampled projections. In *Proc. ICCV Workshop*, pages 518–525, 2017.
- 372 [37] Siming Zheng, Yang Liu, Ziyi Meng, Mu Qiao, Zhishen Tong, Xiaoyu Yang, Shensheng Han, and Xin Yuan.
373 Deep plug-and-play priors for spectral snapshot compressive imaging. *Photonics Res.*, 9(2):B18–B29,
374 2021.

375 Checklist

- 376 1. For all authors...
- 377 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
378 contributions and scope? [Yes]
- 379 (b) Did you describe the limitations of your work? [Yes] See supplemental material.
- 380 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
381 supplemental material.
- 382 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
383 them?[Yes]
- 384 2. If you are including theoretical results...
- 385 (a) Did you state the full set of assumptions of all theoretical results? [N/A]

- 386 (b) Did you include complete proofs of all theoretical results? [N/A]
- 387 3. If you ran experiments...
- 388 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 389 mental results (either in the supplemental material or as a URL)? [No] But we promise
- 390 to release all our codes upon paper's acceptance, as stated in Section 4.
- 391 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 392 were chosen)? [Yes]
- 393 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 394 ments multiple times)? [N/A]
- 395 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 396 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.
- 397 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 398 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 399 (b) Did you mention the license of the assets? [N/A]
- 400 (c) Did you include any new asset either in the supplemental material or as a URL? [N/A]
- 401 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 402 using/curating? [N/A]
- 403 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 404 information or offensive content? [N/A]
- 405 5. If you used crowdsourcing or conducted research with human subjects...
- 406 (a) Did you include the full text of instructions given to participants and screenshots, if
- 407 applicable? [N/A]
- 408 (b) Did you describe any potential participant risks, with links to Institutional Review
- 409 Board (IRB) approvals, if applicable? [N/A]
- 410 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 411 spent on participant compensation? [N/A]