# Modeling Mask Uncertainty in Hyperspectral Image Reconstruction

Jiamian Wang*, Yulun Zhang†, Xin Yuan‡, Ziyi Meng§ and Zhiqiang Tao¶

## Abstract

Recently, hyperspectral imaging (HSI) has attracted increasing research attention, especially for the ones based on a coded aperture snapshot spectral imaging (CASSI) system. Existing deep HSI reconstruction models are generally trained on paired data to retrieve original signals upon 2D compressed measurements given by a particular optical hardware mask in CASSI, during which the mask largely impacts the reconstruction performance and could work as a "model hyperparameter" governing on data augmentations. This mask-specific training style will lead to a hardware miscalibration issue, which sets up barriers to deploying deep HSI models among different hardware and noisy environments. To address this challenge, we introduce mask uncertainty for HSI with a complete variational Bayesian learning treatment and explicitly model it through a mask decomposition inspired by real hardware. Specifically, we propose a novel Graph-based Self-Tuning (GST) network to reason uncertainties adapting to varying spatial structures of masks among different hardware. Moreover, we develop a bilevel optimization framework to balance HSI reconstruction and uncertainty estimation, accounting for the hyperparameter property of masks. Extensive experimental results and model discussions validate the effectiveness (over 33/30 dB) of the proposed GST method under two miscalibration scenarios and demonstrate a highly competitive performance compared with the state-of-the-art well-calibrated methods. Our code and pretrained model are available at https://github.com/Jiamian-Wang/mask_uncertainty_spectral_SCI

## 1 Introduction

Hyperspectral imaging (HSI) provides richer signals than the traditional RGB vision and has broad applications across agriculture [14, 32, 30], medical imaging [31, 18], remote sensing [5, 58, 62], etc. Various HSI systems have been built and studied in recent years, among which, the coded aperture snapshot spectral imaging (CASSI) system [12, 55] stands out due to its passive modulation property and has attracted increasing research attentions [50, 28, 57, 37, 17, 39, 59] in the computer vision community. The CASSI system adopts a hardware encoding & software decoding schema. It first utilizes an optical hardware mask to compress hyperspectral signals into a 2D measurement and then develops software algorithms to retrieve original signals upon the coded measurement conditioning on one particular mask used in the system. Therefore, the hardware mask generally plays a key role in reconstructing hyperspectral images and may exhibit a strongly-coupled (*one-to-one*) relationship with its reconstruction model.

While deep HSI networks [37, 17, 47, 49, 60] have recently shown a promising performance on high-fidelity reconstruction and real-time inference, they mainly treat the hardware mask as a fixed "model hyperparameter" (governing data augmentations on the compressed measurements) and train the reconstruction network on the paired hyperspectral images and measurements given the same mask. Empirically, this will cause a **hardware miscalibration** issue – the mask used in the pre-trained model is unpaired with the real captured measurement – when deploying a single deep reconstruction network among multiple hardware systems (usually with different masks, or different responses even using the same mask due to fabrication errors). As shown in Fig. 1, the performance of deep reconstruction networks pre-trained with one specific mask will badly degrade when applying in different unseen masks (*one-to-many*). Rather than re-training models on

---

*Santa Clara University. jwang16@scu.edu

†ETH Zürich. yulun100@gmail.com

‡Westlake University. xyuan@westlake.edu.cn

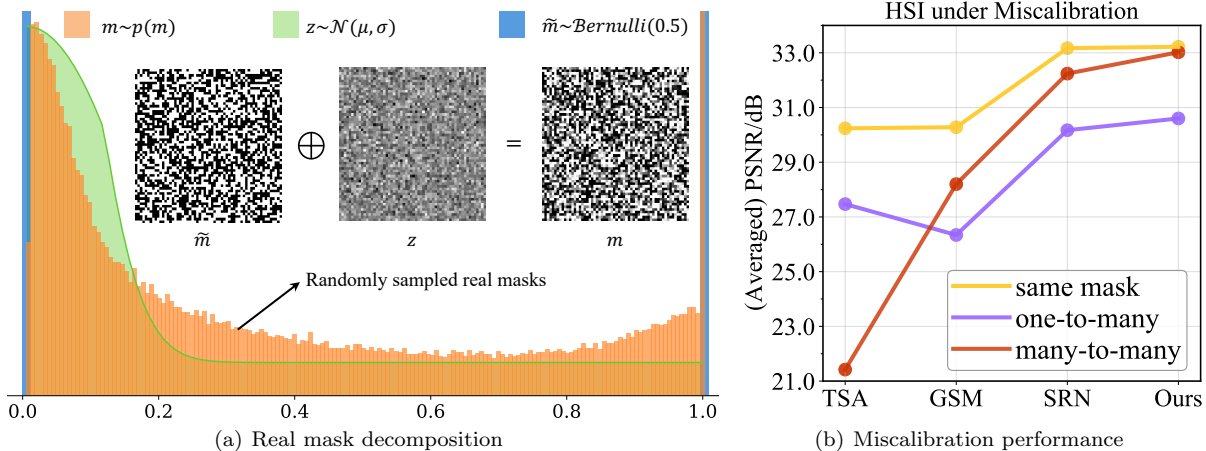§Kuaishou Technology. mengziyi@163.com

¶Santa Clara University. ztao@scu.edu

Figure 1: (a) A real mask $m \sim p(m)$ can be decomposed into an unknown clean mask $\tilde{m}$ plus random noise $z$. The mask distribution is plotted by realistic hardware mask values. Note that the distributions are demonstrated in a `symlog` scale. (b) Performance comparison under three different settings, including 1) the same mask for training/testing, 2) training on one mask and testing on multiple masks (one-to-many), and 3) training with random masks and testing on a held-out mask set (many-to-many).

each new mask, which is inflexible for practical usage, we are more interested in training a single model that could adapt to different hardware by exploring and exploiting uncertainties among masks[1].

One possible solution is to train a deep network over multiple CASSI systems, *i.e.*, using multiple masks and their corresponding encoded measurements following the deep ensemble [24] strategy. However, due to the distinct spatial patterns of each mask and its hyperparameter property, directly training the network with randomly sampled masks still cannot achieve a well-calibrated performance and sometimes performs even worse, *e.g.*, *many-to-many* of TSA-Net [37] in Fig. 1. Hence, we delve into one possible mask decomposition observed from the real hardware, which treats a mask as the unknown clean one plus random noise like Gaussian (see Fig. 1). We consider the noise stemming from two practical sources: 1) the hardware fabrication in real CASSI systems and 2) the functional mask values caused by different lighting environments. Notably, rather than modeling the entire mask distribution, which is challenging due to the high-dimensionality of a 2D map, we explicitly model the mask uncertainty as Gaussian noise centering around a given mask through its decomposition and resort to learn *self-tuning variances* adapting to different mask spatial patterns.

In this study, we propose a novel Graph-based Self-Tuning (GST) network to model mask uncertainty upon variational Bayesian learning and hyperparameter optimization techniques. On the one hand, we approximate the mask posterior distribution with variational inference under the given prior from real mask values, leading to a smoother mask distribution with smaller variance supported by empirical evidence. On the other hand, we leverage graph convolution neural networks to instantiate a stochastic encoder to reason uncertainties varying to different spatial structures of masks. Moreover, we develop a bilevel optimization framework (Fig. 2) to balance the HSI reconstruction performance and the mask uncertainty estimation, accounting for the high sensitive network responses to the mask changes. We summarize the contributions of this work as follows.

- We introduce *mask uncertainty* for CASSI to calibrate a single deep reconstruction network applying in multiple hardware, which brings a promising research direction to improve the robustness and flexibility of deploying CASSI systems to retrieve hyperspectral signals in real-world applications. To our best knowledge, this is the first work to explicitly explore and model mask uncertainty in the HSI reconstruction problem.

- A complete variational Bayesian learning framework has been provided to approximate the mask posterior distribution based on a mask decomposition inspired by real hardware mask observations. More-

---

[1]Note that in CASSI, the masks corresponding to different wavelengths are obtained by calibration with a single fixed wavelength to approximate a small range of band-limited signal. In [2], a high-order model was proposed to address this issue, which is different from the uncertainty problem in the mask considered in the present work.
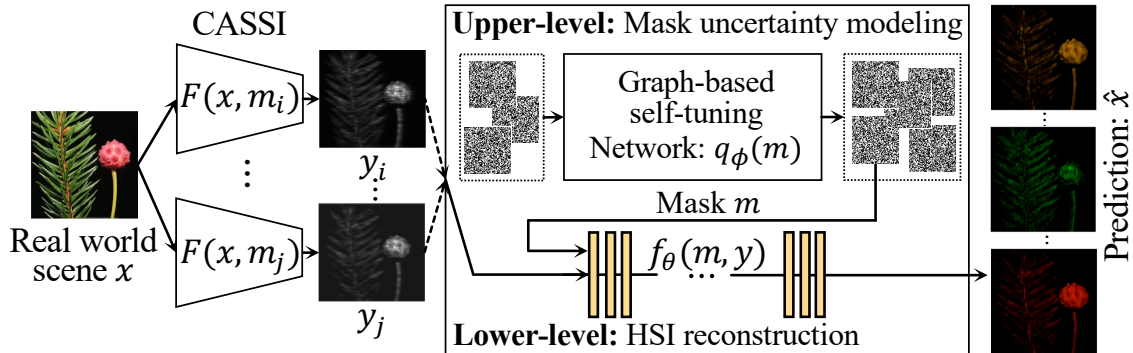
Figure 2: Illustration of the proposed method with a bilevel optimization framework. The upper-level employs the proposed Graph-based Self-Tuning (GST) network to model mask uncertainty through approximating the mask posterior distribution, while the lower-level adopts a reconstruction network $f_\theta(m, y)$ to retrieve the original hyperspectral images $\hat{x}$ by taking masks as hyperparameters. Our model could be applied in multiple CASSI systems using different unseen masks.

over, we design and develop a bilevel optimization framework (see Fig. 2) to jointly achieve high-fidelity HSI reconstruction and mask uncertainty estimation.

- We propose a novel Graph-based Self-Tuning (GST) network to automatically capture uncertainties varying to different spatial structures of 2D masks, leading to a smoother mask distribution over real samples and working as an effective data augmentation method.

- Extensive experimental results on both simulation data and real data demonstrate the effectiveness (over 33/30 dB) of our approach under two miscalibration cases. Besides, the proposed method also shows a highly competitive performance compared with state-of-the-art methods under the traditional well-calibrated setting.

## 2 Related Work

Recently, plenty of advanced algorithms have been designed from diverse perspectives to reconstruct the HSI data from measurements encoded by CASSI system. Among them, optimization-based methods solve the problem by introducing different priors, eg, GPSR [8], TwIST [3] and GAP-TV [54] etc., in which DeSCI [28] leads to the best results. Another novel stream is to empower optimization-based method by deep learning. For example, deep unfolding methods [15, 33, 48] and Plug-and-Play (PnP) structures [41, 56, 42, 39] have been raised. Although their well interpretability and robustness to masks to a certain degree, they suffer from low efficiency and unstable convergence.

Besides, a number of deep reconstruction networks [38, 40, 37, 47, 17] have been proposed, yielding SOTA performance with high reconstructive efficiency. Specifically, $\lambda$-net [40] is introduced under a generative adversarial framework. TSA-Net [37] reconstructs the HSI by considering spatial/spectral self-attentions, outperforming counterparts of the day. SRN [47] provides a lightweight reconstruction backbone based on residual learning. Recently, a Gaussian Scale Mixture (GSM) based method [17] is proposed, showing robustness on masks by enabling an approximation on different sensing matrices. However, all the above pre-trained networks perform unsatisfactorily on distinct unseen masks, raising a question that how to deploy a single reconstruction network among different hardware systems.

In this work, we enable a single well-trained deep reconstruction network to adapt to distinct real masks by proposing a variational approach for estimating the mask uncertainty. Popular uncertainty estimation methods include 1) Bayesian neural networks (BNN) [34, 4, 11] and 2) deep ensemble [24, 9, 27]. The former usually approximates the weight posterior distribution by using variational inference [4] or MC-dropout [11], while the latter generally trains a group of networks from random weight initialization. However, it is challenging to directly quantify mask uncertainty via BNNs or deep ensemble, since treating masks as model weights contraries to their hyperparameter properties.
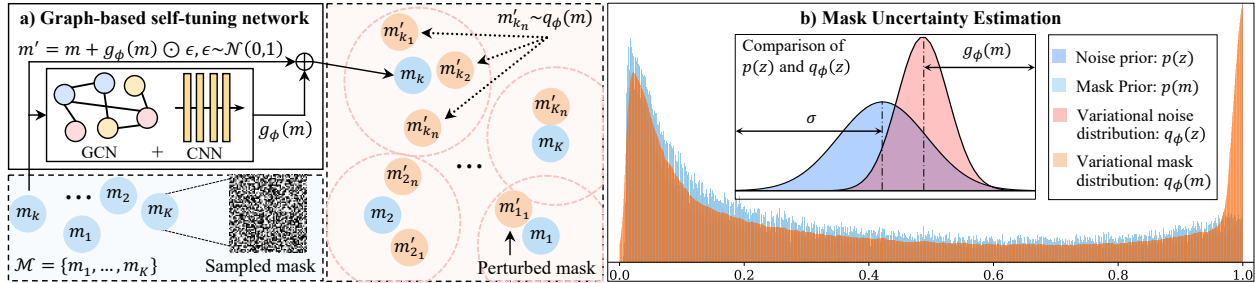
Figure 3: Illustration of modeling mask uncertainty with the proposed Graph-based Self-Tuning (GST) network. a) GST takes as input a real mask $m_k$ randomly sampled from different hardware masks $\mathcal{M}$ and obtains perturbed masks $m'_{k_n}$ by learning self-tuning variance centering on $m_k$. b) GST estimates mask uncertainty by approximating the mask posterior with a variational distribution $q_\phi(m)$, leading to a smoother mask distribution over the mask prior $p(m)$. More discussions are given in Section 4.2.

# 3 Methodology

## 3.1 Preliminaries

**HSI reconstruction.** The reconstruction based on the CASSI system [55, 37] generally includes a hardware-encoding forward process and a software-decoding inverse process. Let $x$ be a 3D hyperspectral image with the size of $H \times W \times \Lambda$, where $H$, $W$, and $\Lambda$ represent the height, width, and the number of spectral channels. The optical hardware encoder will compress the datacube $x$ into a 2D measurement $y$ upon a fixed physical mask $m$. The forward model of CASSI is

$$y = F(x; m) = \sum_\lambda^\Lambda \texttt{shift}(x)_\lambda \odot \texttt{shift}(m)_\lambda + \zeta, \tag{1}$$

where $\lambda$ refers to a spectral channel, $\odot$ represents the element-wise product, and $\zeta$ denotes the measurement noise. The shift operation is implemented by a single disperser as $\texttt{shift}(x)(u, v, i) = x(h, w + d(i - \lambda), i)$. In essence, the measurement $y$ is captured by spectral modulation[2] conditioning on the hardware mask $m$.

In the inverse process, we adopt a deep reconstruction network as the decoder: $\hat{x} = f_\theta(m, y)$ where $\hat{x}$ is the retrieved hyperspectral image, and $\theta$ represents all the learnable parameters. Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ be the dataset. The reconstruction network $f_\theta$ is generally trained to minimize an $\ell_1$ or $\ell_2$ loss as the following:

$$\min_\theta \sum_{x,y \in \mathcal{D}} \ell(f_\theta(m, y) - x) \quad \text{where} \quad y = F(x; m). \tag{2}$$

We instantiate $f_\theta$ as a recent HSI backbone model provided in [47], which benefits from nested residual learning and spatial/spectral-invariant learning. We employ this backbone for its lightweight structure to simplify the training[3].

**Hardware miscalibration.** As shown in Eq. (2), there is a *paired relationship* between the parameter $\theta$ and mask $m$ in deep HSI models. Thus, for different CASSI systems (i.e., distinct masks), multiple pairs $\{m_1; \theta_1\}, ..., \{m_K; \theta_K\}$ are expected for previous works. Empirically, the miscalibration between $m$ and $\theta$ will lead to obvious performance degradation. This miscalibration issue inevitably impairs the flexibility and robustness of deploying deep HSI models across real systems, considering the expensive training time and various noises existing in hardware. To alleviate such a problem, one straight-forward solution is to train the model $f_\theta$ with multiple masks, i.e., $\mathcal{M} = \{m_1, ..., m_K\}$, falling in a similar strategy to deep ensemble [24]. However, directly training a single network with random masks cannot provide satisfactory performance to unseen masks (see Section 4), since the lack of explicitly exploring the relationship between uncertainties and different mask structures.

---

[2]We used a two-pixel shift for neighbored spectral channels following [37, 47]. More details about spectral modulation could be found in [55].

[3]We provide more details toward SRN in Append B. More discussion on alternating reconstruction backbones could be found in Append C.

4

## 3.2 Mask Uncertainty

Modeling mask uncertainty is challenging due to 1) the high dimensionality of a 2D mask and the limited mask set size (i.e., $K$ for $\mathcal{M}$) and 2) the varying spatial structures (patterns) among masks. Hence, we propose to first estimate uncertainties around each mask through one possible mask decomposition in this section and then reason how the uncertainty will adapt to the change of mask structures with a self-tuning network in Section 3.3.

Inspired by the distribution of real mask values (Fig. 1 and Fig. 3), which renders two peaks at 0 and 1 and appears a Gaussian shape spreading over the middle, we decompose a mask as two components

$$m = \tilde{m} + z, \tag{3}$$

where we assume each pixel in $z$ follows a Gaussian distribution. For simplicity, we slightly abuse the notations by denoting the noise prior as $p(z) = \mathcal{N}(\mu, \sigma)$. The $\tilde{m}$ denotes the underlying clean binary mask with a specific spatial structure.

We estimate the mask uncertainty by approximating the mask posterior $p(m|X, Y)$ following [4, 10, 52], where $X = \{x_1, \ldots, x_N\}$ and $Y = \{y_1, \ldots, y_N\}$ indicate hyperspectral images and their corresponding measurements. To this end, we aim to learn a variational distribution $q_\phi(m)$ parameterized by $\phi$ to minimize the KL-divergence between $q_\phi(m)$ and $p(m|X, Y)$, i.e., $\min_\phi KL[q_\phi(m)||p(m|X, Y)]$, equivalent to maximizing the evidence lower bound (ELBO) [22, 16] as

$$\max_\phi \underbrace{\mathbb{E}_{q_\phi(m)}[\log p(X|Y, m)]}_{\text{reconstruction}} - \underbrace{KL[q_\phi(m)||p(m)]}_{\text{regularization}}, \tag{4}$$

where the first term measures the reconstruction (i.e., reconstructing the observations $X$ based on the measurements $Y$ and mask $m$ via $f_\theta(m, y)$), and the second term regularizes $q_\phi(m)$ given the mask prior $p(m)$. Inspired by Eq. (3), we treat the clean mask $\tilde{m}$ as a 2D constant with a particular structure and focus on mask uncertainties arising from the noise $z$. Thus, the variational distribution $q_\phi(m)$ is defined as a Gaussian distribution centering on a given $m \in \mathcal{M}$ by

$$q_\phi(m) = \mathcal{N}(m, g_\phi(m)), \tag{5}$$

where $g_\phi(m)$ learns *self-tuning variance* to model the uncertainty adapting to real masks sampled from $\mathcal{M}$. Correspondingly, the underlying variational noise distribution $q_\phi(z)$ follows Gaussian distribution with variance $g_\phi(m)$. We adopt the reparameterization trick [22] for computing stochastic gradients for the expectation w.r.t $q_\phi(m)$. Specifically, let $m' \sim q_\phi(m)$ be a random variable sampled from the variational distribution, we have

$$m' = t(\phi, \epsilon) = m + g_\phi(m) \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1). \tag{6}$$

Notably, we clamp all the pixel values of $m'$ in range $[0, 1]$.

The first term in Eq. (4) reconstructs $x$ with $p(x|y, m) \propto p(x|\hat{x} = f_\theta(m, y))$, yielding a squared error when $x|\hat{x}$ follows a Gaussian distribution [46]. Similar to AutoEncoders, we implement the negative log-likelihood $\mathbb{E}_{q_\phi(m)}[-\log p(X|Y, m)]$ as a $\ell_2$ loss and compute its Monte Carlo estimates with Eq. (6) as

$$\ell(\phi, \theta; \mathcal{D}) = \frac{N}{B} \sum_{i=1}^{B} \|f_\theta(y_i, t(\phi, \epsilon_i)) - x_i\|^2, \tag{7}$$

where $(x_i, y_i) \in \mathcal{D}$, $B$ denotes the mini-batch size, and $t(\phi, \epsilon_i)$ represents the $i$-th sample from $q_\phi(m)$. We leverage $t(\phi, \epsilon_i)$ to sample $B$ perturbed masks from $q_\phi(m)$ centering on one randomly sampled mask $m \in \mathcal{M}$ per batch.

Since $p(m)$ is unknown due to various spatial structures of masks, we resort to approximating the KL term in Eq. (4) with the entropy of $q_\phi(m)$. Eventually, we implement the ELBO($q(m)$) with the following loss:

$$\mathcal{L}(\phi, \theta; \mathcal{D}) = \ell(\phi, \theta; \mathcal{D}) + \beta \mathbb{H}[\log q_\phi(m)], \tag{8}$$

where $\mathbb{H}[\log q_\phi(m)]$ is computed by $\ln(g_\phi(m)\sqrt{2\pi e})$ and $\beta > 0$ interprets the objective function between variational inference and variational optimization [35, 20].
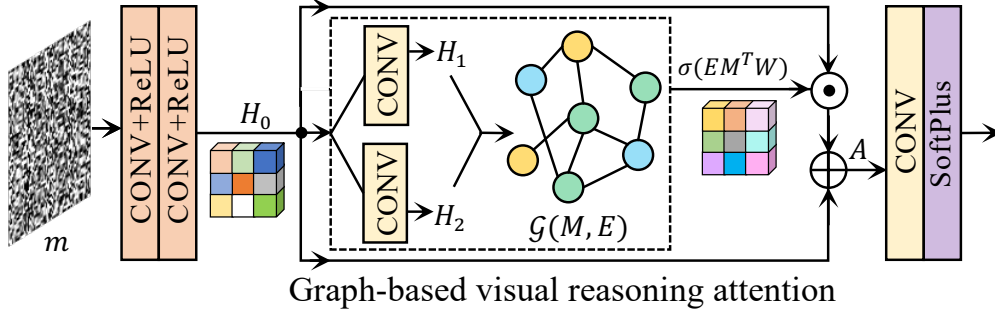
Graph-based visual reasoning attention

Figure 4: Structure of Graph-based Self-Tuning (GST) network. The model takes mask $m$ as input and outputs a 2D variance map, globally handling mask in a graph domain.

## 3.3 Graph-based Self-Tuning Network

We propose a graph-based self-tuning (GST) network to instantiate the variance model $g_\phi(m)$ in Eq. (5), which captures uncertainties around each mask and leads to a smoother mask distribution over real masks (see Fig. 3). The key of handling unseen masks (new hardware) is to learn how the distribution will change along with the varying spatial structures of masks. To this end, we implement the GST as a visual reasoning attention network [6, 25, 61]. It firstly computes pixel-wise correlations (visual reasoning) based on neural embeddings and then generates attention scores based on graph convolutional networks (GCN) [44, 23]. Unlike previous works [6, 25, 61], the proposed GST model is tailored for building a stochastic probabilistic encoder to capture the mask distribution.

We show the network structure of GST in Fig. 4. Given a real mask $m$, GST produces neural embedding $H_0$ by using two concatenated `CONV-ReLU` blocks. Then, we employ two `CONV` layers to convert $H_0$ into two different embeddings $H_1$ and $H_2$, and generate a graph representation by matrix multiplication $H_1^T H_2$, resulting in $\mathcal{G}(M, E)$, where the node matrix $M$ represents mask pixels and the edge matrix $E$ denotes the pixel-wise correlations. Let $W$ be the weight matrix of GCN. We obtain an enhanced attention cube by pixel-wise multiplication

$$A = H_0 \odot (\sigma(EM^TW) + \mathbf{1}), \tag{9}$$

where $\sigma$ is the sigmoid function. Finally, the self-tuning variance is obtained by

$$g_\phi(m) = \delta(\texttt{CONV}(A)), \tag{10}$$

where $\delta$ denotes the softplus function and $\phi$ denotes all the learnable parameters. Consequently, GST enables adaptive variance modeling to arbitrary real masks.

## 3.4 Bilevel Optimization

While it is possible to jointly train the HSI reconstruction network $f_\theta$ and the self-tuning network $g_\phi$ using the loss in Eq. (8), it is more proper to formulate the training of these two networks as a bilevel optimization framework accounting for two hyperparameter properties of masks. First, the deep reconstruction network is generally high-sensitive to the change/perturbation of masks. Thus, the model weight $\theta$ is largely subject to a mask $m$. Second, deep HSI methods [37, 47] usually employ a single mask and a group of shifting operations to lift the 2D measurement as a multi-channel input, where the mask works as a hyperparameter of data augmentation for training deep networks.

---

**Algorithm 1:** GST Training Algorithm

    **Input:** $\mathcal{D}^{trn}, \mathcal{D}^{val}, \mathcal{M}$; initialized $\theta, \phi$;
    **Output:** $\theta^*, \phi^*$
1  Pre-train $f_\theta(\cdot)$ on $\mathcal{D}^{trn}$ with $\alpha_0$ for $T^{init}$ epochs;
2  **while** *not converge* **do**
3      **for** $t = 1, ..., T^{trn}$ **do**
4         $\{(x_i, y_i)\}_{i=1}^B \sim \mathcal{D}^{trn}$;
5         $\theta \leftarrow \theta - \alpha_1 \frac{\partial}{\partial\theta}\ell(\phi, \theta; \mathcal{D}^{trn})$;
6      **end**
7      **for** $t = 1, ..., T^{val}$ **do**
8         $\{(x_i, y_i)\}_{i=1}^B \sim \mathcal{D}^{val}, m \sim \mathcal{M}, \epsilon \sim \mathcal{N}(0, 1)$;
9         $\phi \leftarrow \phi - \alpha_2 \frac{\partial}{\partial\phi}\mathcal{L}(\phi, \theta; \mathcal{D}^{val})$;
10     **end**
11 **end**

---

the mask works as a hyperparameter of data augmentation for training deep networks.

To be specific, we define the lower-level problem as HSI reconstruction and the upper-level problem as mask uncertainty estimation, and propose the final objective function of our GST model as the following:

$$\min_{\phi} \mathcal{L}(\phi, \theta^*; \mathcal{D}^{val}) \ \ \text{s.t.} \ \ \theta^* = \arg\min_{\theta} \ell(\phi, \theta; \mathcal{D}^{trn}), \tag{11}$$

where $\ell(\phi, \theta; \mathcal{D}^{trn})$ is provided in Eq. (7) with a training set and $\mathcal{L}(\phi, \theta^*; \mathcal{D}^{val})$ is given by Eq. (8) in a validation set. Upon Eq. (11), $f_\theta$ and $g_\phi$ are alternatively updated by computing gradients $\frac{\partial l}{\partial \theta}$ and $\frac{\partial \mathcal{L}}{\partial \phi}$. To better initializing the parameter $\theta$, we pre-train the reconstruction network $f_\theta(m, y)$ for several epochs. The entire training procedure of the proposed method is summarized in Algorithm 1. Notably, introducing Eq. (11) brings two benefits. 1) It could balance the solutions of HSI reconstruction and mask uncertainty estimation. 2) It enables the proposed GST as a hyperparameter optimization method, which could provide high-fidelity reconstruction even working on a single mask (see Table 3).

# 4 Experiments

**Simulation data.** We adapt the training set provided in [37]. Simulated measurements are computed by mimicking the compressing process of SD-CASSI optical system [37]. For both metric and perceptual comparison, we employ a benchmark test set that contains ten 256×256×28 HSIs following [17, 47, 39]. We build a validation set by splitting 40 HSIs from the training set.

**Real data.** We adopt five real 660×714 measurements provided in [37] for qualitative evaluation. We train the model on the expanded simulation training set by augmenting 37 HSIs originating from the KAIST dataset [7]. Also, the Gaussian noise ($\mathcal{N}(0, \varphi), \varphi \sim U[0, 0.05]$) is added on the simulated measurements during training, for the sake of mimicking practical measurement noise $\zeta$. All the other settings are kept the same as compared deep reconstruction methods.

**Mask set.** We adopt two 660×660 hardware masks in our experiment. Both are produced by the same fabrication process. For training, the mask set $\mathcal{M}$ is created by randomly cropping (256×256) from the mask provided in [37]. For testing, both masks are applied. In simulation, testing masks are differentiated from the training ones. For real HSI reconstruction, the second mask [38] is applied, indicating hardware miscalibration scenario.

**Implementation details.** The training procedures (Algorithm 1) for simulation and real case follow the same schedule: We apply the `xavier uniform` [13] initializer with `gain=1`. Before alternating, the reconstruction network is trained for $T^{init}$=20 epochs (learning rate $\alpha_0$=4× $10^{-4}$). Then, the reconstruction network $f_\theta(\cdot)$ is updated on training phase for $T^{trn}$=5 epochs ($\alpha_1$=4× $10^{-4}$) and the GST network is updated on validation phase for $T^{val}$=3 epochs ($\alpha_2$=1× $10^{-5}$). The learning rates are halved per 50 epochs and we adopt Adam optimizer [21] with default setting. In this work, we adopt SRN (`v1`) [47] as the reconstructive backbone, i.e., the full network without rescaling pairs. All the experiments were conducted on 4 NVIDIA GeForce RTX 3090 GPUs.

**Compared methods.** For hardware miscalibration, masks for data pair setup (i.e., CASSI compressing procedure) and network training should be different from those for testing. We specifically consider two scenarios: 1) many-to-many, i.e., training the model on mask set $\mathcal{M}$ and testing it by unseen masks; 2) One-to-many, i.e., training the model on single mask and testing it by diverse unseen masks, which brings more challenges. For quantitative performance comparison, in this work all the testing results are computed upon 100 testing trials (100 random unseen masks). We compare with four SOTA methods: TSA-Net [37], GSM-based method [17], SSI-ResU-Net (SRN) [47] and PnP-DIP [39], among which the first three are deep reconstruction networks and the last one is an iterative optimization-based method. Note that 1) PnP-DIP is a self-supervised method. We test it by feeding the data encoded by different masks in the testing mask set and compute the performance over all obtained results. 2) For real-world HSI reconstruction, all models are trained on the same mask while tested on the other. Specifically, the network inputs are initialized by testing mask for TSA-Net and SRN. For GSM, as demonstrated by the authors, we directly compute the sensing matrix of testing mask and replace the corresponding approximation in the network. We use PSNR and SSIM [51] as metrics for quantitative comparison.
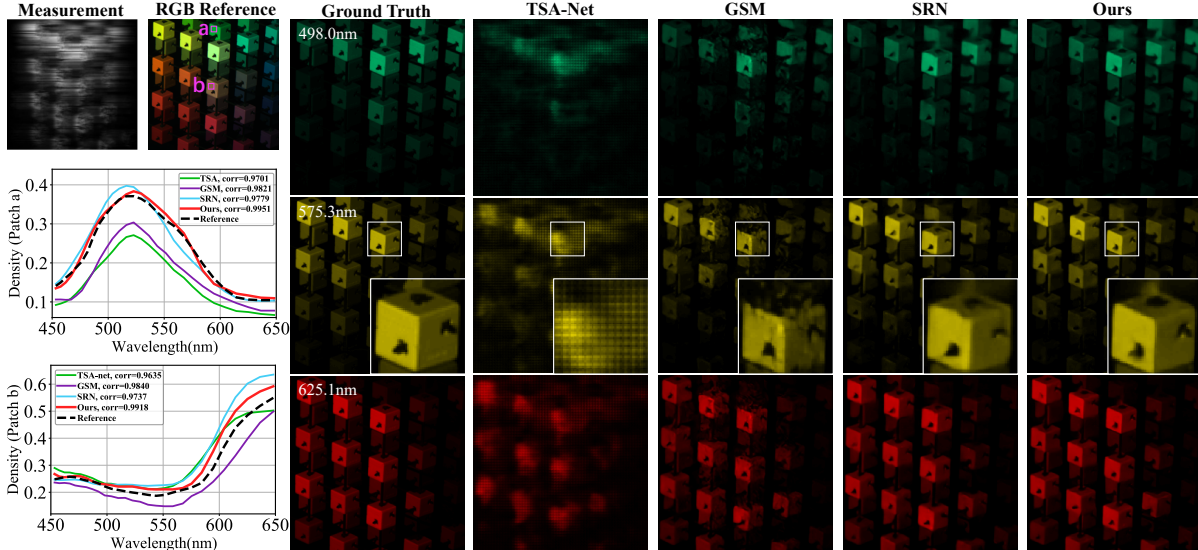
Figure 5: Reconstruction results on one simulation scene under hardware miscalibration (many-to-many). All methods are trained on the mask set $\mathcal{M}$ and tested by one unseen mask. Density curves computed on chosen patches are compared to analysis the spectra.

Table 1: PSNR(dB)/SSIM by different methods on 10 simulation scenes under the **many-to-many** hardware miscalibration. All the methods are trained with a mask set $\mathcal{M}$ and tested by random unseen masks. TSA-Net [37], GSM [17], and SRN [47] are obtained with a mask ensemble strategy. We report mean$_{\pm\texttt{std}}$ among 100 testing trials.

| Scene | TSA-Net [37] | | GSM [17] | | PnP-DIP† [39] | | SRN [47] | | Ours | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| 1 | $23.45_{\pm0.29}$ | $0.6569_{\pm0.0051}$ | $31.38_{\pm0.20}$ | $0.8826_{\pm0.0032}$ | $29.24_{\pm0.98}$ | $0.7964_{\pm0.0532}$ | $33.26_{\pm0.16}$ | $0.9104_{\pm0.0018}$ | $\mathbf{33.99}_{\pm0.14}$ | $\mathbf{0.9258}_{\pm0.0013}$ |
| 2 | $18.52_{\pm0.12}$ | $0.5511_{\pm0.0049}$ | $25.94_{\pm0.22}$ | $0.8570_{\pm0.0041}$ | $25.73_{\pm0.54}$ | $0.7558_{\pm0.0117}$ | $29.86_{\pm0.23}$ | $0.8809_{\pm0.0029}$ | $\mathbf{30.49}_{\pm0.17}$ | $\mathbf{0.9002}_{\pm0.0022}$ |
| 3 | $18.42_{\pm0.30}$ | $0.5929_{\pm0.0127}$ | $26.11_{\pm0.20}$ | $0.8874_{\pm0.0034}$ | $29.61_{\pm0.45}$ | $0.8541_{\pm0.0125}$ | $31.69_{\pm0.20}$ | $0.9093_{\pm0.0020}$ | $\mathbf{32.63}_{\pm0.16}$ | $\mathbf{0.9212}_{\pm0.0013}$ |
| 4 | $30.44_{\pm0.15}$ | $0.8940_{\pm0.0043}$ | $34.72_{\pm0.35}$ | $0.9473_{\pm0.0023}$ | $38.21_{\pm0.66}$ | $0.9280_{\pm0.0078}$ | $39.90_{\pm0.22}$ | $0.9469_{\pm0.0012}$ | $\mathbf{41.04}_{\pm0.23}$ | $\mathbf{0.9667}_{\pm0.0014}$ |
| 5 | $20.89_{\pm0.23}$ | $0.5648_{\pm0.0077}$ | $26.15_{\pm0.24}$ | $0.8256_{\pm0.0061}$ | $28.59_{\pm0.79}$ | $0.8481_{\pm0.0183}$ | $30.86_{\pm0.16}$ | $0.9232_{\pm0.0019}$ | $\mathbf{31.49}_{\pm0.17}$ | $\mathbf{0.9379}_{\pm0.0017}$ |
| 6 | $23.04_{\pm0.19}$ | $0.6099_{\pm0.0060}$ | $30.97_{\pm0.29}$ | $0.9224_{\pm0.0025}$ | $29.70_{\pm0.51}$ | $0.8484_{\pm0.0186}$ | $34.20_{\pm0.23}$ | $0.9405_{\pm0.0014}$ | $\mathbf{34.89}_{\pm0.29}$ | $\mathbf{0.9545}_{\pm0.0009}$ |
| 7 | $15.97_{\pm0.14}$ | $0.6260_{\pm0.0042}$ | $22.58_{\pm0.24}$ | $0.8459_{\pm0.0054}$ | $27.13_{\pm0.31}$ | $0.8666_{\pm0.0079}$ | $27.27_{\pm0.16}$ | $0.8515_{\pm0.0026}$ | $\mathbf{27.63}_{\pm0.16}$ | $\mathbf{0.8658}_{\pm0.0024}$ |
| 8 | $22.64_{\pm0.18}$ | $0.6366_{\pm0.0066}$ | $29.76_{\pm0.22}$ | $0.9059_{\pm0.0021}$ | $28.38_{\pm0.35}$ | $0.8325_{\pm0.0203}$ | $32.35_{\pm0.22}$ | $0.9320_{\pm0.0015}$ | $\mathbf{33.02}_{\pm0.26}$ | $\mathbf{0.9471}_{\pm0.0013}$ |
| 9 | $18.91_{\pm0.11}$ | $0.5946_{\pm0.0083}$ | $27.23_{\pm0.11}$ | $0.8899_{\pm0.0021}$ | $33.63_{\pm0.26}$ | $0.8779_{\pm0.0073}$ | $32.83_{\pm0.13}$ | $0.9205_{\pm0.0016}$ | $\mathbf{33.45}_{\pm0.13}$ | $\mathbf{0.9317}_{\pm0.0013}$ |
| 10 | $21.90_{\pm0.18}$ | $0.5249_{\pm0.0110}$ | $28.05_{\pm0.21}$ | $0.8877_{\pm0.0055}$ | $27.24_{\pm0.43}$ | $0.7957_{\pm0.0226}$ | $30.25_{\pm0.14}$ | $0.9053_{\pm0.0019}$ | $\mathbf{31.49}_{\pm0.15}$ | $\mathbf{0.9345}_{\pm0.0015}$ |
| *Avg.* | $21.42_{\pm0.07}$ | $0.6162_{\pm0.0030}$ | $28.20_{\pm0.01}$ | $0.8852_{\pm0.0001}$ | $29.66_{\pm0.38}$ | $0.8375_{\pm0.0093}$ | $32.24_{\pm0.10}$ | $0.9121_{\pm0.0010}$ | $\mathbf{33.02}_{\pm0.01}$ | $\mathbf{0.9285}_{\pm0.0001}$ |

†PnP-DIP is a mask-free method which reconstructs from measurements encoded by random masks.

## 4.1 HSI Reconstruction Performance

We evaluate our method under different settings on both simulation and real data. More visualizations are provided in the supplementary material.

**Miscalibration (many-to-many).** Training the deep reconstruction networks with a mask ensemble strategy could improve the generalization ability and might be a potential solution to deploy a single network across different hardware, such as training TSA-Net [37], GSM [17], and SRN [47] on a mask set. However, as shown in Table 1 and Table 3, these methods generally suffer from a clear performance degradation under miscalibration compared with their well-calibrated performance. For example, TSA-Net [37] drops around 8db on PSNR. Benefiting from modeling mask uncertainty, our approach achieves high-fidelity results (over 33dB) on both cases, with only a 0.2db drop.

Fig. 5 compares the reconstructive results of different methods perceptually. Our method retrieves more details at different spectral channels. We randomly choose two regions, corresponding to two colors separately (patch a and b in RGB reference), to analyze results regarding spectra. Density curves of spectral correlations to ground truth are compared in the bottom left of Fig. 5.

**Miscalibration (one-to-many).** In Table 2, all the methods are trained on a single mask and tested on

Table 2: PSNR(dB)/SSIM by different methods on 10 simulation scenes under the **one-to-many** hardware miscalibration. All the methods are trained by a single mask and tested by random unseen masks. We report `mean`$_{\pm\texttt{std}}$ among 100 testing trials.

| Scene | TSA-Net [37] | | GSM [17] | | PnP-DIP$^\dagger$ [39] | | SRN [47] | | Ours | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| 1 | $28.49_{\pm0.58}$ | $0.8520_{\pm0.0081}$ | $28.20_{\pm0.95}$ | $0.8553_{\pm0.0185}$ | $29.24_{\pm0.98}$ | $0.7964_{\pm0.0532}$ | $31.24_{\pm0.77}$ | $0.8878_{\pm0.0117}$ | $\mathbf{31.72}_{\pm0.76}$ | $\mathbf{0.8939}_{\pm0.0119}$ |
| 2 | $24.96_{\pm0.51}$ | $0.8332_{\pm0.0064}$ | $24.46_{\pm0.96}$ | $0.8330_{\pm0.0189}$ | $25.73_{\pm0.54}$ | $0.7558_{\pm0.0117}$ | $27.87_{\pm0.82}$ | $0.8535_{\pm0.0131}$ | $\mathbf{28.22}_{\pm0.85}$ | $\mathbf{0.8552}_{\pm0.0144}$ |
| 3 | $26.14_{\pm0.76}$ | $0.8829_{\pm0.0108}$ | $23.71_{\pm1.18}$ | $0.8077_{\pm0.0221}$ | $29.61_{\pm0.45}$ | $0.8541_{\pm0.0125}$ | $28.31_{\pm0.88}$ | $\mathbf{0.8415}_{\pm0.0213}$ | $28.77_{\pm1.13}$ | $0.8405_{\pm0.0257}$ |
| 4 | $35.67_{\pm0.47}$ | $0.9427_{\pm0.0028}$ | $31.55_{\pm0.75}$ | $0.9385_{\pm0.0074}$ | $38.21_{\pm0.66}$ | $0.9280_{\pm0.0078}$ | $\mathbf{37.93}_{\pm0.72}$ | $\mathbf{0.9476}_{\pm0.0057}$ | $37.60_{\pm0.81}$ | $0.9447_{\pm0.0071}$ |
| 5 | $25.40_{\pm0.59}$ | $0.8280_{\pm0.0108}$ | $24.44_{\pm0.96}$ | $0.7744_{\pm0.0291}$ | $28.59_{\pm0.79}$ | $0.8481_{\pm0.0183}$ | $27.99_{\pm0.79}$ | $0.8680_{\pm0.0194}$ | $\mathbf{28.58}_{\pm0.79}$ | $\mathbf{0.8746}_{\pm0.0208}$ |
| 6 | $29.32_{\pm0.60}$ | $0.8796_{\pm0.0047}$ | $28.28_{\pm0.92}$ | $0.9026_{\pm0.0094}$ | $29.70_{\pm0.51}$ | $0.8484_{\pm0.0186}$ | $32.13_{\pm0.87}$ | $\mathbf{0.9344}_{\pm0.0061}$ | $\mathbf{32.72}_{\pm0.79}$ | $0.9339_{\pm0.0061}$ |
| 7 | $22.80_{\pm0.65}$ | $0.8461_{\pm0.0101}$ | $21.45_{\pm0.79}$ | $0.8147_{\pm0.0162}$ | $27.13_{\pm0.31}$ | $0.8666_{\pm0.0079}$ | $24.84_{\pm0.73}$ | $0.7973_{\pm0.0150}$ | $\mathbf{25.15}_{\pm0.76}$ | $\mathbf{0.7935}_{\pm0.0173}$ |
| 8 | $28.09_{\pm0.43}$ | $0.8738_{\pm0.0043}$ | $28.08_{\pm0.76}$ | $0.9024_{\pm0.0089}$ | $28.38_{\pm0.35}$ | $0.8325_{\pm0.0203}$ | $31.32_{\pm0.59}$ | $0.9324_{\pm0.0043}$ | $\mathbf{31.84}_{\pm0.56}$ | $\mathbf{0.9323}_{\pm0.0042}$ |
| 9 | $27.75_{\pm0.55}$ | $0.8865_{\pm0.0054}$ | $26.80_{\pm0.78}$ | $0.8773_{\pm0.0144}$ | $33.63_{\pm0.26}$ | $0.8779_{\pm0.0073}$ | $31.06_{\pm0.66}$ | $0.8997_{\pm0.0091}$ | $\mathbf{31.11}_{\pm0.72}$ | $\mathbf{0.8988}_{\pm0.0104}$ |
| 10 | $26.05_{\pm0.48}$ | $0.8114_{\pm0.0072}$ | $26.40_{\pm0.77}$ | $0.8771_{\pm0.0124}$ | $27.24_{\pm0.43}$ | $0.7957_{\pm0.0226}$ | $29.01_{\pm0.61}$ | $0.9028_{\pm0.0092}$ | $\mathbf{29.50}_{\pm0.68}$ | $\mathbf{0.9030}_{\pm0.0098}$ |
| Avg. | $27.47_{\pm0.46}$ | $0.8636_{\pm0.0060}$ | $26.34_{\pm0.06}$ | $0.8582_{\pm0.0012}$ | $29.66_{\pm0.38}$ | $0.8375_{\pm0.0093}$ | $30.17_{\pm0.63}$ | $0.8865_{\pm0.0108}$ | $\mathbf{30.60}_{\pm0.08}$ | $\mathbf{0.8881}_{\pm0.0013}$ |

$^\dagger$PnP-DIP is a mask-free method which reconstructs from measurements encoded by random masks.

Table 3: PSNR (dB) and SSIM values by different algorithms on the simulation dataset under the well-calibrated setting (training/test on the **same mask**). We adopt the same 256×256 real mask provided in previuous works [37, 17] for a fair comparison.

| Scene | $\lambda$-net [40] | | HSSP [48] | | TSA-Net [37] | | GSM [17] | | PnP-DIP [39] | | SRN [47] | | Ours | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| 1 | 30.82 | 0.8492 | 31.07 | 0.8577 | 31.26 | 0.8920 | 32.38 | 0.9152 | 31.99 | 0.8633 | 34.13 | 0.9260 | **34.19** | **0.9292** |
| 2 | 26.30 | 0.8054 | 26.30 | 0.8422 | 26.88 | 0.8583 | 27.56 | 0.8977 | 26.56 | 0.7603 | 30.60 | 0.8985 | **31.04** | **0.9014** |
| 3 | 29.42 | 0.8696 | 29.00 | 0.8231 | 30.03 | 0.9145 | 29.02 | 0.9251 | 30.06 | 0.8596 | 32.87 | 0.9221 | **32.93** | **0.9224** |
| 4 | 37.37 | 0.9338 | 38.24 | 0.9018 | 39.90 | 0.9528 | 36.37 | 0.9636 | 38.99 | 0.9303 | **41.27** | **0.9687** | 40.71 | 0.9672 |
| 5 | 27.84 | 0.8166 | 27.98 | 0.8084 | 28.89 | 0.8835 | 28.56 | 0.8820 | 29.09 | 0.8490 | 31.66 | 0.9376 | **31.83** | **0.9415** |
| 6 | 30.69 | 0.8527 | 29.16 | 0.8766 | 31.30 | 0.9076 | 32.49 | 0.9372 | 29.68 | 0.8481 | 35.14 | **0.9561** | 35.14 | 0.9543 |
| 7 | 24.20 | 0.8062 | 24.11 | 0.8236 | 25.16 | 0.8782 | 25.19 | 0.8860 | 27.68 | 0.8639 | 27.93 | **0.8638** | 28.08 | 0.8628 |
| 8 | 28.86 | 0.8307 | 27.94 | 0.8811 | 29.69 | 0.8884 | 31.06 | 0.9234 | 29.01 | 0.8412 | 33.14 | **0.9488** | 33.18 | 0.9486 |
| 9 | 29.32 | 0.8258 | 29.14 | 0.8676 | 30.03 | 0.8901 | 29.40 | 0.9110 | 33.35 | 0.8802 | 33.49 | 0.9326 | **33.50** | **0.9332** |
| 10 | 27.66 | 0.8163 | 26.44 | 0.8416 | 28.32 | 0.8740 | 30.74 | 0.9247 | 27.98 | 0.8327 | 31.43 | **0.9338** | 31.59 | 0.9311 |
| Avg. | 29.25 | 0.8406 | 28.93 | 0.8524 | 30.24 | 0.8939 | 30.28 | 0.9166 | 30.44 | 0.8529 | 33.17 | 0.9288 | **33.22** | **0.9292** |

multiple unseen masks. We pose this setting to further demonstrate the hardware miscalibration challenge. As can be seen, except for the mask-free method PnP-DIP, the methods usually experience large performance difference compared with using mask ensemble in Table 1 and the well-calibrated case in Table 3. This observation supports the motivation of modeling mask uncertainty – 1) simply using mask ensemble may aggravate the miscalibration (TSA-Net using ensemble performs even worse) and 2) the model trained with a single mask cannot be effectively deployed in different hardware.

**Same mask (one-to-one).** Table 3 reports the well-calibrated performance for all the methods, *i.e.*, training/testing models on the same real mask. While our approach is specially designed for training with multiple masks, it still consistently outperforms all the competitors by leveraging bilevel optimization.

**Results on real data.** Fig. 6 visualizes reconstruction results on the real dataset, where the *left* provides results using the same mask and the *right* is under the one-to-many miscalibration setting. For the same mask, the proposed method is supposed to perform comparably. For the miscalibration, we train all the models on a single real mask provided in [37] and test them on the other one [38]. As shown, the proposed method produces plausible results and improves over other methods visually. Besides, spectral fidelity is also demonstrated.
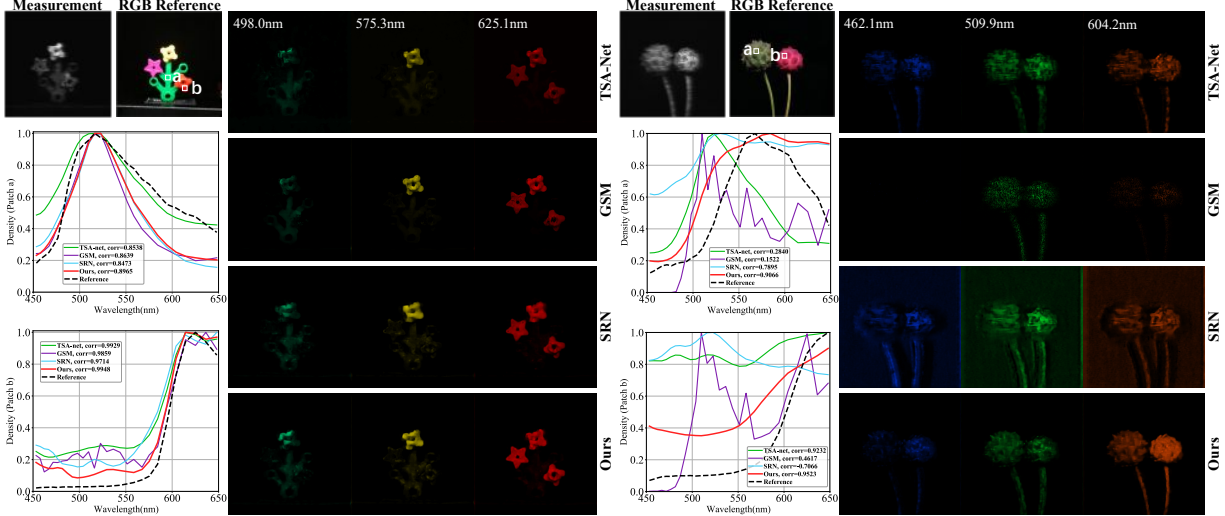
Figure 6: Real HSI reconstruction. **Left**: same mask (one-to-one) reconstruction, i.e., all methods are trained and tested on the same 660×660 real mask. **Right**: miscalibration (one-to-many) setting, i.e., all methods are trained on a single mask and tested by unseen masks (Here we adopt another 660×660 real mask). Density curves computed on chosen patches are compared to analysis the spectra.

Table 4: Ablation study and complexity analysis. All the methods are tested on simulation test set under the many-to-many setting with one NVIDIA RTX 3090 GPU. We report the PSNR (dB)/SSIM among 100 testing trials, the total training time, and the test time per sample. PnP-DIP is self-supervised, thus no training is required.

| Settings | PSNR | SSIM | #params (M) | FLOPs (G) | Training (day) | Test (sec.) |
|---|---|---|---|---|---|---|
| TSA-Net [37] | $21.42_{\pm 0.07}$ | $0.6162_{\pm 0.0030}$ | 44.25 | 110.06 | 1.23 | 0.068 |
| GSM [17] | $28.20_{\pm 0.01}$ | $0.8852_{\pm 0.0001}$ | 3.76 | 646.35 | 6.05 | 0.084 |
| PnP-DIP [39] | $29.66_{\pm 0.38}$ | $0.8375_{\pm 0.0093}$ | 33.85 | 64.26 | − | 482.78 |
| w/o GST | $32.24_{\pm 0.10}$ | $0.9121_{\pm 0.0010}$ | 1.25 | 81.84 | 1.14 | 0.061 |
| w/o Bi-Opt | $32.43_{\pm 0.02}$ | $0.9206_{\pm 0.0001}$ | 1.27 | 82.87 | 1.83 | 0.061 |
| w/o GCN | $32.82_{\pm 0.01}$ | $0.9262_{\pm 0.0001}$ | 1.27 | 82.78 | 1.63 | 0.062 |
| Ours (full model) | $33.02_{\pm 0.01}$ | $0.9285_{\pm 0.0001}$ | 1.27 | 82.87 | 2.56 | 0.062 |

## 4.2 Model Discussion

**Ablation study**. Table 4 compares the performance and complexity of the proposed `full model` with three ablated models as follows. 1) The model `w/o GST` is equivalent to training the reconstruction backbone SRN [47] with a mask ensemble strategy. 2) The model `w/o Bi-Opt` is implemented by training the proposed method without using the bilevel optimization framework. 3) In the model `w/o GCN`, we replace the GCN module in GST with convolutional layers carrying a similar size of parameters. The bilevel optimization achieves 0.59dB improvement without overburdening the complexity. The GCN contributes 0.2dB with 0.09G FLOPs increase. Overall, the proposed GST yields 0.8dB improvement with negligible cost (i.e., +0.02M #params, +1.03G FLOPs, and +1.14 days training), and could be used in multiple unseen masks without re-training.

**Complexity comparison**. In Table 4, we further compare the complexity of the proposed method with several recent HSI methods. The proposed method possess one of the smallest model size. Besides, our method shows a comparable FLOPs and training time with others. Notably, given $M$ distinct masks, TSA-Net, GSM, and SRN require $M\times$ training time as reported to achieve well-calibrated performance. Instead, the proposed method only needs to be trained one time to provide calibrated reconstructions over multiple unseen masks.

**Self-tuning variance under different priors.** We first validate the effectiveness of the self-tuning variance by comparing with the fix-valued variance, i.e., scalars from 0 to 1. As shown by the green curve in Fig. 7 (a), fixed variance only achieves no more than 32dB performance. And the best performance by 0.005 indicates a strong approximation nature to the mask noise. By comparison, we explore the behaviour of the
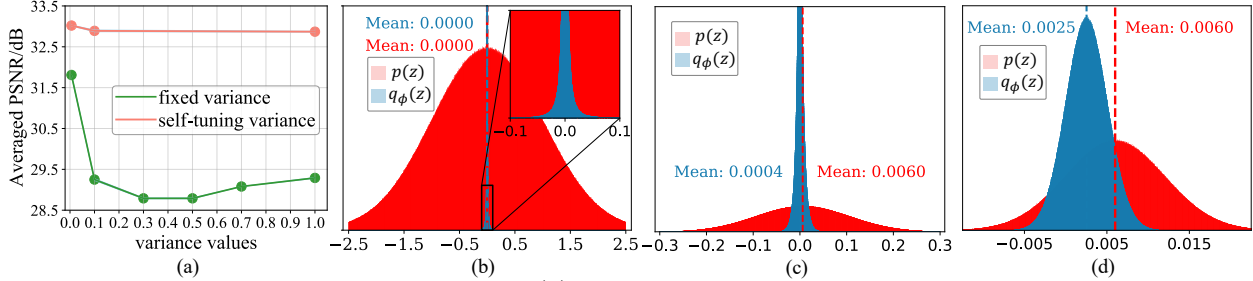
Figure 7: Discussion on self-tuning variance. (a) Performance comparison between self-tuning variance and fixed ones. (b) The standard normal prior $\mathcal{N}(0,1)$. (c) Set the prior as $\mathcal{N}(0.006, 0.1)$ by observing real masks. (d) Set the prior as $\mathcal{N}(0.006, 0.005)$ by observing real masks and the performance curve in (a).
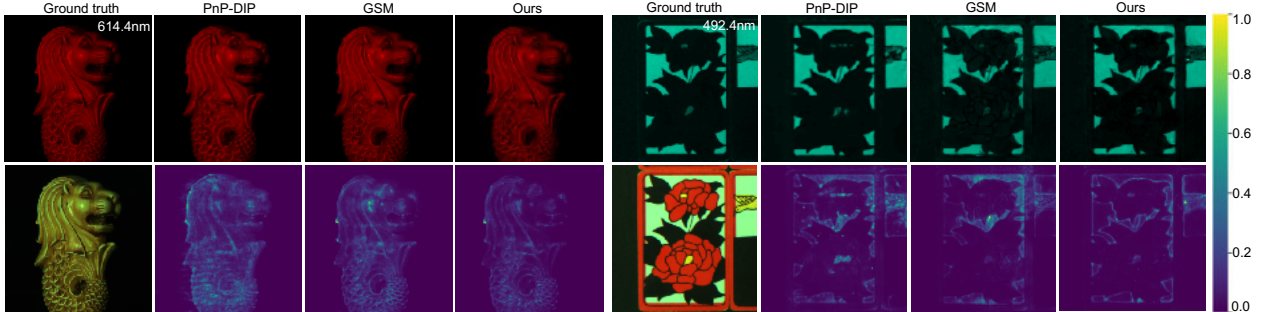


Figure 8: Illustration of epistemic uncertainty induced by multiple masks. For each block, the first row shows the averaged reconstruction results of selected channels given by different methods and the second demonstrates the corresponding epistemic uncertainty.

self-tuning variance upon different noise priors and achieve no less than 32.5dB performance (red curve in Fig. 7 (a)). Specifically, we implement the noise prior $p(z)$ by exchanging the standard normal distribution of auxiliary variable $\epsilon$ in Eq. (6). We start from $\mathcal{N}(0,1)$, which is so broad that the GST network tries to centralize variational noise and restricting the randomness as Fig. 7 (b) shown. Then, we constraint the variance and approximate the mean value by the minimum of the real mask histogram to emphasize the near-zero noise, proposing $\mathcal{N}(0.006, 0.1)$. The corresponding variational noise distribution deviates from the prior as shown in Fig. 7 (c), indicating the underlying impact of GST network. Finally, we further combine the previous fixed-variance observation and propose $\mathcal{N}(0.006, 0.005)$. As the red curve in Fig. 7 (a) indicated, the best reconstruction performance is also obtained. In summary, the proposed method effectively restricts the posited noise prior, leading to the variational noise distribution with a reduced range.

**From mask uncertainty to epistemic uncertainty**. The hardware mask plays a similar role to model hyperparameter and largely impacts the weights of reconstruction networks. Thus, marginalizing over the mask posterior distribution will induce the epistemic uncertainty (also known as model uncertainty [10, 19]) and reflect as pixel-wise variances (the second row in Fig. 8) of the reconstruction results over multiple unseen masks. As can be seen, the mask-free method PnP-DIP [39] still produces high uncertainties given measurements of the same scene coded by different hardware masks. While employing a deep ensemble strategy could alleviate this issue, such as training GSM [17] with mask ensemble, it lacks an explicit way to quantify mask uncertainty and may lead to unsatisfactory performance (see Table 1). Differently, the proposed GST method models mask uncertainty by approximating the mask posterior through a variational Bayesian treatment, exhibiting high-fidelity reconstruction result with low epistemic uncertainties across different masks as shown in Fig. 8.

# 5 Conclusions

In this work, we have explored a practical hardware miscalibration issue when deploying deep HSI models in real CASSI systems. Our solution is to calibrate a single reconstruction network via modeling mask uncertainty. We proposed a complete variational Bayesian learning treatment upon one possible mask

decomposition inspired by observations on real masks. Bearing the objectives of variational mask distribution modeling and HSI retrieval, we introduced and implemented a novel Graph-based Self-Tuning (GST) network that proceeds HSI reconstruction and uncertainty reasoning under a bilevel optimization framework. The proposed method enabled a smoothed distribution and achieved promising performance under two different miscalibration scenarios. We hope the proposed insight will benefit future work in this novel research direction.

# Appendix

## A  Overview

In this appendix, we present additional results and analyses about the proposed method as follows.

- **Reconstruction Backbone SRN**: A detailed introduction to the reconstruction backbone network used in the manuscript (Section B).
- **Alternating Reconstruction Backbones**: More analyses on alternating reconstruction backbones employed in the proposed method (Section C).
- **Spectral Fidelity Analysis**: Evaluation on the spectral fidelity of the reconstruction results by the proposed method (Section D).
- **Epistemic Uncertainty Analysis**: More visualization and discussion on epistemic uncertainty of the proposed method (Section E).
- **Complementary Ablation Studies**: Ablation studies under one-to-many miscalibration and the same mask setting (one-to-one) (Section F).
- **Self-tuning variance Analysis**: More discussions about the proposed self-tuning variance. Specifically, more results for fixed variance are provided. Also, we demonstrate the convergence of $g_\phi(m)$, variational noise distributions given distinct noise priors (Section G).
- **Datasets**: More illustrations on the dataset, includes training data, validation data, testing data, and mask set (Section H).

## B  Reconstruction Backbone: SRN

In the manuscript, we adopt a recent deep reconstruction network, SRN [47] as the backbone $f_\theta(\cdot)$. Specifically, the network input $x_{in} \in \mathbb{R}^{H \times W \times \Lambda}$ is initialized by the measurement $y \in \mathbb{R}^{H \times (W+\Lambda-1) \times \Lambda}$ and the mask $m \in \mathbb{R}^{H \times W}$

$$x_{in}[:,:,\lambda] = \texttt{shift}(y)_\lambda \odot m, \tag{12}$$

where $\odot$ is a Hadamard product and the $\texttt{shift}$ is the reverse operation applied in the forward process (see Eq. (1) in the manuscript for more details).

The model is composed of a 1) *main body*, which is simultaneously bridged by a global skip connection, 2) a *head* operation and 3) a *tail* operation, both of which are conducted by a $\texttt{CONV-ReLU}$ structure. Let $x_{head}$ and $x_{body}$ be the output of the *head* operation and the *main body*, respectively. We have

$$x_{body} = f_{res}^J(f_{res}^{J-1}(...(f_{res}^1(x_{head}))...)), \tag{13}$$

where $J = 16$ concatenated residual blocks share the same structure, i.e., $f_{res}(x) = x + (\texttt{CONV}(\texttt{ReLU}(\texttt{CONV}(x))))$.

## C  Alternating Reconstruction Backbones

The performance and the robustness toward masks of the deep reconstruction networks largely depend on their constructions. Thus, we validate the effectiveness of the proposed method upon backbones with different architectures.

Table 5: Averaged PSNR(dB)/SSIM of the different models. We consider the miscalibration many-to-many scenario for a fair comparison. For three types of backbones, this is implemented by training upon a mask ensemble and testing on random masks. The *mean* and *std* are obtained upon 100 testing trials.

| Models | PSNR (dB) | SSIM |
|---|---|---|
| SRN [47] | $32.24_{\pm 0.10}$ | $0.9121_{\pm 0.0010}$ |
| Spectral ViT [1] | $31.62_{\pm 0.09}$ | $0.9282_{\pm 0.0010}$ |
| SwinIR [26] | $33.49_{\pm 0.10}$ | $0.9501_{\pm 0.0010}$ |
| SRN+GST (Ours) | $\mathbf{33.02}_{\pm 0.01}$ | $\mathbf{0.9285}_{\pm 0.0001}$ |
| Spectral ViT+GST (Ours) | $\mathbf{32.15}_{\pm 0.01}$ | $\mathbf{0.9330}_{\pm 0.0001}$ |
| SwinIR+GST (Ours) | $\mathbf{34.15}_{\pm 0.01}$ | $\mathbf{0.9548}_{\pm 0.0001}$ |

**SwinIR Backbone.** In this supplementary material, we consider transformer architectures as the backbone. Specifically, transformer acquires modeling ability from attention mechanism [45], which has been proved to behave quite differently from the traditional ConvNets [43].

Given the initialized input $x_{in}$ by Eq. (12), we firstly implement the backbone by Swin transformer structure [29], which computes the spatial self-attention. It is composed of three modules: (1) shallow feature extraction by a CONV3×3 layer, i.e., $x_{SF} = \text{CONV}(x_{in})$, (2) deep feature extraction module consisting of $K$ concatenated residual Swin transformer blocks, i.e, $x_{DF} = f_{DF}(x_{SF})$ where $f_{DF}(\cdot) = f_{RSTB}^{K}(f_{RSTB}^{K-1}(...(f_{RSTB}^{1}(\cdot))...))$, and (3) a reconstruction module by a CONV3×3 layer, i.e., $\widehat{x} = \text{CONV}(x_{DF})$.

For each residual Swin transformer block $f_{RSTB}(\cdot)$, we have $L$ Swin transformer layers, which conducts window-based MSA and MLP

$$x = f_{\text{W-MSA}}(f_{\text{LN}}(x)) + x, \quad x = f_{\text{MLP}}(f_{\text{LN}}(x)) + x, \tag{14}$$

where the details of the $f_{\text{W-MSA}}(\cdot)$, $f_{\text{MLP}}(\cdot)$, and $f_{\text{LN}}(\cdot)$ could be found in [29]. In the experiment, we set the $K = 4$, $L = 6$. For all the blocks, we let the embedding dimension to be 60 and number of heads to be 6.[4]

**Spectral ViT.** We also provide another type of vision transformer, which exchanges the previous spatial self-attention with the spectral self-attention. Specifically, it treats the feature map of each embedding channel as a token. Given the query $\mathbf{Q}$, key $\mathbf{K}$ and value $\mathbf{V}$, we have output $\mathbf{X}$

$$\mathbf{X} = \mathbf{V}\text{Attn}(\mathbf{K}, \mathbf{Q}), \quad where \quad \text{Attn}(\mathbf{K}, \mathbf{Q}) = \text{softmax}(\mathbf{K}^T\mathbf{Q}/\delta), \tag{15}$$

where $\delta$ denotes a learnable scalar. For more details please refer to [1].

**Comparison.** We summarize the performance of different backbones under miscalibration scenario many-to-many in Table 5. For detailed illustration of the miscalibration setting, please refer to manuscript Section 4.1. Notably, the integration of the backbone into our method is implemented by training the full model upon a mask dataset $\mathcal{M}$ in a bilevel optimization framework.

By comparison, one can draw the following conclusions. (1) For the metric comparison, our method brings performance gain $\Delta\text{PSNR} = 0.78\text{dB}$, $0.53\text{dB}$, and $0.66\text{dB}$, respectively, for different backbones. (2) The proposed method enables high-fidelity reconstruction with the highest confidence. Specifically, in many-to-many case, $\text{PSNR}_{std} > 0.1\text{dB}$. Our method achieves $10\times$ the randomness control, indicating a better epistemic modeling capacity (Please see Section E).

# D   Spectral Fidelity.

In this work, we adopt two methods to demonstrate the spectral fidelity of the reconstruction results.

Firstly, given the prediction $\widehat{x} \in \text{R}^{H \times W \times N_\lambda}$, we treat each spectral channel as a R.V. (random variable) of $HW$ dimensions and calculate the channel-wise correlations. For each hyperspectral image with $N_\lambda = 28$, a correlation matrix of $28 \times 28$ could be visualized. We compare these matrices by the references and the predictions in Fig. 9. The more consistent they are, the higher spectral fidelity we achieve. By observation,

---

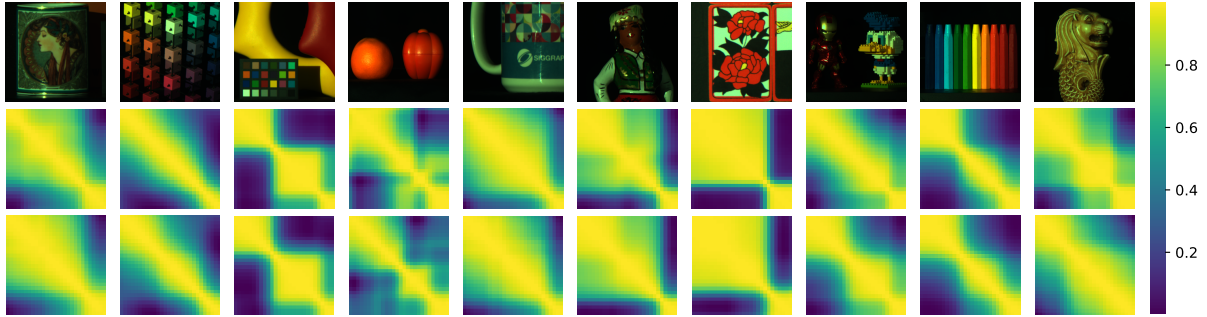[4]For implementation please refer to https://github.com/JingyunLiang/SwinIR.

Figure 9: RGB references of the benchmark simulation test set (top line) and spectral correlation coefficient visualizations by the reference (middle line) as well as the proposed method (bottom line). Each correlation coefficient map is of the size 28×28.
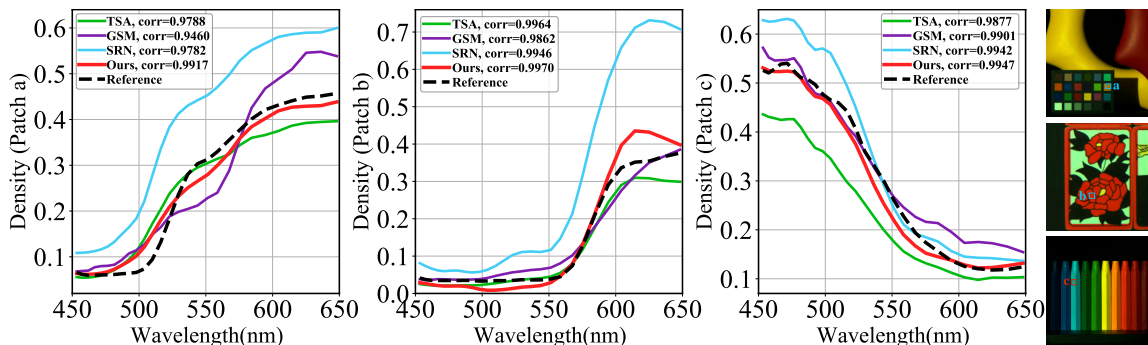


Figure 10: Spectral correlation to the ground truth on exampled locations. Spatial patches (patch `a,b,c` as plotted in the right most column RGB references, please zoom in for better visualization) are chosen to ensure the monochromaticity. Density curves are computed upon the predictions by different methods within the chosen patch.

the correlation matrices by the predictions show highly similar visual patterns as the reference, indicating that the proposed method effectively captures long-range spectra dependencies. We also notice that there might be minor differences at the centers of the matrices between the visualizations. Rectifying this part is pretty challenging as the model need to precisely distinguish the difference between each adjacent spectral pairs.

Table 6: Averaged spectral correlations (↑) to the reference. For each scene, we compute the averaged correlation values upon density curves corresponding to five selected patches. Please refer to examples in Fig. 10 for a detailed computational procedure.

| Methods | Scene1 | Scene2 | Scene3 | Scene4 | Scene5 | Scene6 | Scene7 | Scene8 | Scene9 | Scene10 | *Avg.* |
|---------|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|--------|
| GSM [17] | 0.9903 | 0.9674 | 0.9751 | 0.9435 | 0.9891 | 0.9904 | 0.9951 | 0.9700 | 0.9870 | 0.8926 | 0.9701 |
| SRN [47] | 0.9910 | 0.9701 | 0.9753 | 0.9429 | 0.9898 | 0.9910 | 0.9959 | 0.9865 | 0.9870 | 0.9000 | 0.9730 |
| Ours | **0.9956** | **0.9910** | **0.9921** | **0.9458** | **0.9925** | **0.9982** | **0.9969** | 0.9867 | **0.9903** | **0.9004** | **0.9790** |

Secondly, we quantitatively compare the spectral fidelity of different methods upon density curves. As three examples demonstrated in Fig. 10, we first crop a small spatial patch from the prediction (exampled by the RGB reference on the right most column), then we draw the density curve by pixel intensities in that small patch. Finally, the correlations between the reference curve and that from the predictions are computed. Higher correlation values indicate a higher spectral fidelity for the cropped patch. The small patch is chosen to ensure the monochromaticity (wavelength). For example, if we choose the patch whose color lies in blue∼cyan range (bottom-right RGB reference in Fig. 10), the energy of the density would concentrate in the 450nm∼500nm.

To globally compare the spectral fidelity, we randomly choose five monochromatic patch of each scene and compute an averaged correlation value upon five density curves. We report the correlation values of ten
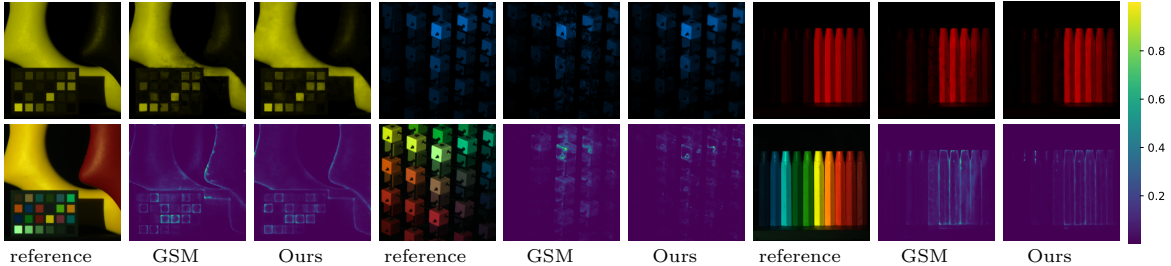
Figure 11: Three exampled epistemic uncertainty visualizations by GSM [17] and the proposed method. For each example, we demonstrate averaged reconstruction results on selected wavelengths (i.e., 567.5nm, 471.6nm, and 614.4nm) in top line, and the epistemic uncertainty in the bottom line.

Table 7: Ablation study of the proposed method under the traditional setting (one-to-one).

| Settings | PSNR(dB) | SSIM |
|---|---|---|
| w/o GST | 33.17 | 0.9288 |
| w/o Bi-Opt | 32.73 | 0.9193 |
| w/o GCN | 33.23 | 0.9286 |
| Ours (full model) | 33.22 | 0.9292 |

scenes in Table 6 and demonstrate the superiority of the proposed method, as compared to GSM [17] and SRN [47].

# E    Epistemic Uncertainty

As illustrated in Section 4.2 of the manuscript, the proposed method demonstrates low epistemic uncertainty by approximating the mask distribution. In this section, we provide more visualizations and analyses on epistemic uncertainty in Fig. 11. Specifically, we test the well-trained models upon random real masks and repeat 100 trials. Both GSM and the proposed method are trained upon the same mask set $\mathcal{M}$ for a fair comparison. For each exampled hyperspectral images, we compare the averaged reconstruction and epistemic uncertainty on a selected spectral channel. Notably, in low-frequency regions, both methods show high confidence, while in high-frequency regions (i.e., edges), the proposed method presents a much-lower epistemic uncertainty, which would potentially benefit the down-stream applications like object detection or segmentation upon hyperspectral images.

# F    Ablation Study

In this paper, three scenarios are introduced: 1) one-to-one setting, which is the traditional setting considered by previous reconstruction methods, 2) one-to-many miscalibration, 3) many-to-many miscalibration. Notably, the third scenario enables a complete mask distribution modeling, for which reason we put more emphasize on it and provide the ablation study accordingly in the manuscript. Following that, Table 7 conducts the same ablation experiments under the traditional setting (one-to-one). For miscalibration (one-to-many), we also do verification and report the performance in Table 8. The ablated models include

- w/o GST: we remove the graph-based self-tuning (GST) network from the proposed method. Actually it degrades into the reconstruction backbone SRN [47] applied under corresponding scenarios.

- w/o Bi-Opt: we simultaneously optimize all of the parameters by the lower-level loss function, i.e., Eq. (7) in the manuscript, based on the original training set.

- w/o GCN: for the self-tuning network, we exchange the GCN with a convolutional layer carrying more parameters for a fair comparison.

Table 8: Ablation study of the proposed method under the setting of miscalibration (one-to-many) among 100 testing trials.

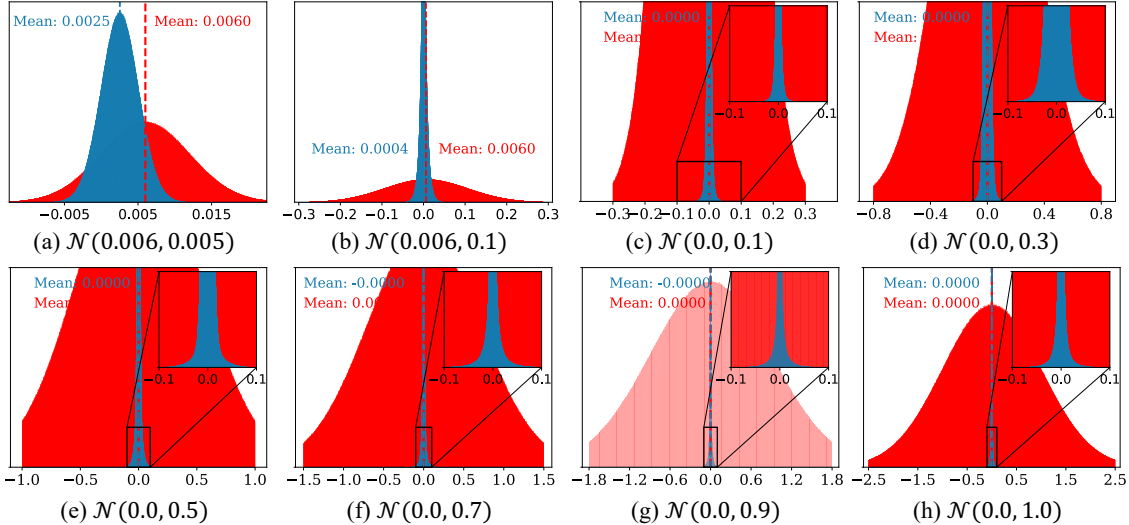| Settings | PSNR(dB) | SSIM |
|---|---|---|
| w/o GST | 30.17±0.63 | 0.8865±0.0108 |
| w/o Bi-Opt | 30.30±0.06 | 0.8843±0.0011 |
| w/o GCN | 30.13±0.07 | 0.8849±0.0011 |
| Ours (full model) | 30.60±0.08 | 0.8881±0.0013 |



Figure 12: Learned variational noise distribution under different priors. Eight different priors (red) are adopted in the experiment. By comparison, variational noise distributions (blue) are characterized by smaller variance. Please refer to the red curves in Fig. 13 for corresponding reconstruction performance comparison.

As shown in the Table 7 and Table 8, both the GST module and bilevel optimization strategy contribute significantly for the final performance boost. While the ablated model `w/o GCN` in self-tuning network works comparably with the `Ours (full model)` under the traditional setting by PSNR, it falls behind regarding SSIM, indicating a sub-optimal reconstruction ability.

# G    Model Discussion

**Fixed variance.** In the manuscript, we showcase that the fixed variance with distinct values only achieves sub-optimal performance compared with the self-tuning variance. In Fig. 13 (b), we also plot the SSIM curve (green) with different values of the fixed variance. Besides, the original PSNR curve (green) is shown in Fig. 13 (a).

**Self-tuning variance under different priors.** Due to the limitation of time and computational resource, we only discussed the self-tuning variance under three most representative noise priors in the manuscript, i.e., $\mathcal{N}(0.006, 0.005)$, $\mathcal{N}(0.006, 0.1)$ and $\mathcal{N}(0.0, 1.0)$. In this section, we demonstrate additional results corresponding to more noise priors. Similar to the performance curves plotted in Fig. 7 (a) in the manuscript, we report the corresponding performances by red curves in Fig. 13, verifying the superiority of the self-tuning variance. In Fig. 12, we visualize both the noise prior and the subsequent variational noise distributions. The similarity between all subplots – the learned variational noise distribution gets a smaller variance than the given prior – validates the effectiveness of mask uncertainty modeling.

In Fig. 14, we explore the convergence of the $g_\phi(m)$ during the training phase of our best model. For pre-training phase of reconstruction network (first 20 epochs as mentioned in the manuscript), the range of $g_\phi(m)$ remains invariant. An interesting observation is that the fluctuation of the $g_\phi(m)$ value is accompanied by
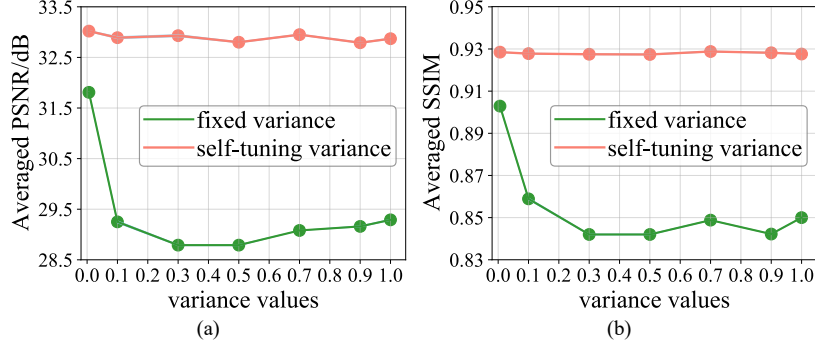
Figure 13: Performance comparison between fixed variance (green) and self-tuning variance (red). The PSNR is compared in (a) and SSIM is compared in (b). Reconstruction using self-tuning variance outperforms that using fixed variance with different values.
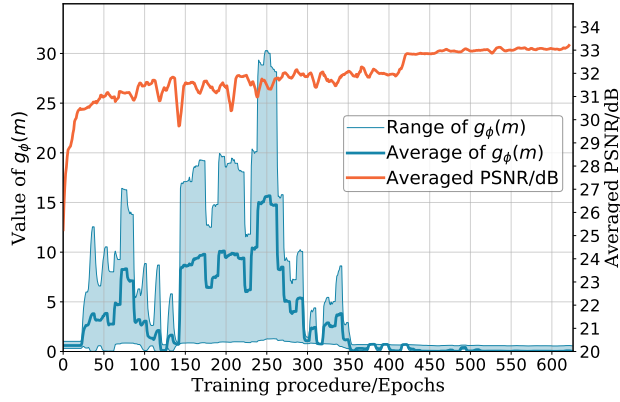


Figure 14: Observation on $g_\phi(m)$ during training. The $g_\phi(m)$ gradually converge to a smaller range with more epochs of training. Meanwhile, a better reconstruction performance can be observed. Both training epochs and validation epochs are jointly counted.

fluctuation of the reconstruction performance, indicating the underlying impact of the self-tuning variance. During the last 200 epochs (including training and validation epochs), a converged $g_\phi(m)$ contributes to a steady performance improvement.

# H  Dataset

**HSI data set.** We adopt the training set provided in [37] and follow the same data augmentation operations. Specifically, the training set contains 205 1024×1024×28 training samples, all of which sources from the CAVE dataset [53]. Our model is trained on 256×256×28 patches randomly cropped from these 205 samples. For a fair comparison with the other deep reconstruction networks, we create a validation data set by randomly splitting 40 hyperspectral images from the above 205 samples. Therefore, *no new HSI data is introduced for our model training.* For the model testing, ten simulation hyperspectral images corresponding to 10 scenes shown in Fig. 9 are used for quantitative and perceptual comparison, following previous works [37, 17, 47, 39, 36].

**Mask set.** Two 660×660 real masks following the same fabrication process are employed in this work. Fig. 15 demonstrates the histograms of both masks. As mentioned in the manuscript, the training mask set $\mathcal{M}$ is built by randomly cropping 256×256 patches from the first real mask [37]. For simulation data, testing masks are collected from both real masks. Notably, **there is no overlap between training and testing mask sets.** For real HSI reconstruction, no testing mask set is available. The second 660×660 real mask [38] is directly applied for testing purpose, indicating the miscalibration scenario.
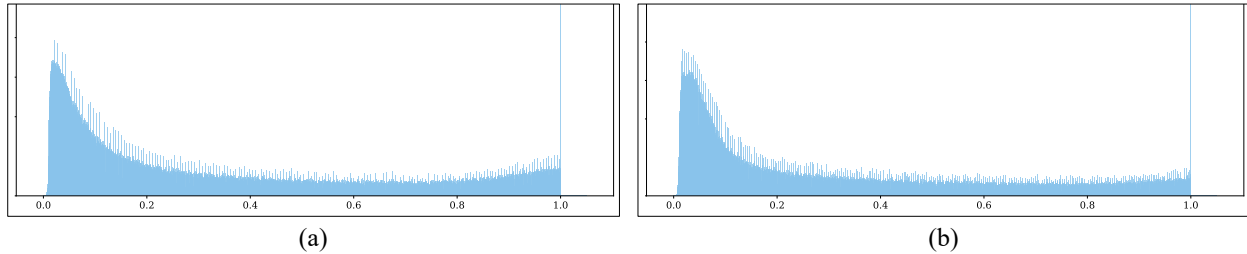
17

Figure 15: Histograms of two real masks applied in this work. (a) sources from [37] and (b) sources from [38]. Both masks are produced by the same fabrication process. Bin number is set to 2000 for both histograms.

# References

[1] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. In *NeurIPS*, 2021. 13

[2] Henry Arguello, Hoover Rueda, Yuehao Wu, Dennis W Prather, and Gonzalo R Arce. Higher-order computational model for coded aperture spectral imaging. *Applied optics*, 52(10):D12–D21, 2013. 2

[3] José M Bioucas-Dias and Mário AT Figueiredo. A new twist: Two-step iterative shrinkage/thresholding algorithms for image restoration. *TIP*, 2007. 3

[4] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *ICML*, 2015. 3, 5

[5] Marcus Borengasser, William S Hungate, and Russell Watkins. *Hyperspectral remote sensing: principles and applications*. CRC press, 2007. 1

[6] Xinlei Chen, Li-Jia Li, Li Fei-Fei, and Abhinav Gupta. Iterative visual reasoning beyond convolutions. In *CVPR*, 2018. 6

[7] Inchang Choi, MH Kim, D Gutierrez, DS Jeon, and G Nam. High-quality hyperspectral reconstruction using a spectral prior. Technical report, 2017. 7

[8] Mário AT Figueiredo, Robert D Nowak, and Stephen J Wright. Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems. *IEEE Journal of selected topics in signal processing*, 2007. 3

[9] Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757*, 2019. 3

[10] Yarin Gal. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016. 5, 11

[11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 1050–1059, 2016. 3

[12] Michael E Gehm, Renu John, David J Brady, Rebecca M Willett, and Timothy J Schulz. Single-shot compressive spectral imaging with a dual-disperser architecture. *Optics express*, 2007. 1

[13] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, 2010. 7

[14] Aoife A Gowen, Colm P O'Donnell, Patrick J Cullen, Gérard Downey, and Jesus M Frias. Hyperspectral imaging–an emerging process analytical tool for food quality and safety control. *Trends in food science & technology*, 18(12):590–598, 2007. 1

[15] John R Hershey, Jonathan Le Roux, and Felix Weninger. Deep unfolding: Model-based inspiration of novel deep architectures. *arXiv preprint arXiv:1409.2574*, 2014. 3

[16] Matthew D Hoffman and Matthew J Johnson. Elbo surgery: yet another way to carve up the variational evidence lower bound. In *NeurIPS Workshop*, 2016. 5

[17] Tao Huang, Weisheng Dong, Xin Yuan, Jinjian Wu, and Guangming Shi. Deep gaussian scale mixture prior for spectral compressive imaging. In *CVPR*, 2021. 1, 3, 7, 8, 9, 10, 11, 14, 15, 17

[18] William R Johnson, Daniel W Wilson, Wolfgang Fink, Mark S Humayun, and Gregory H Bearman. Snapshot hyperspectral imaging in ophthalmology. *Journal of biomedical optics*, 12(1):014036, 2007. 1

[19] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems*, volume 30, 2017. 11

[20] Mohammad Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and scalable bayesian deep learning by weight-perturbation in adam. In *ICML*, 2018. 5

[21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 7

[22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 5

[23] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017. 6

[24] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017. 2, 3, 4

[25] Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. Visual semantic reasoning for image-text matching. In *ICCV*, 2019. 6

[26] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *CVPR*, 2021. 13

[27] Jeremiah Zhe Liu, John Paisley, Marianthi-Anna Kioumourtzoglou, and Brent Coull. Accurate uncertainty estimation and decomposition in ensemble learning. *arXiv preprint arXiv:1911.04061*, 2019. 3

[28] Yang Liu, Xin Yuan, Jinli Suo, David J Brady, and Qionghai Dai. Rank minimization for snapshot compressive imaging. *TPAMI*, 2018. 1, 3

[29] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *CVPR*, 2021. 13

[30] Delia Lorente, Nuria Aleixos, JUAN Gómez-Sanchis, Sergio Cubero, OSCAR LEONARDO García-Navarrete, and José Blasco. Recent advances and applications of hyperspectral imaging for fruit and vegetable quality assessment. *Food and Bioprocess Technology*, 5(4):1121–1142, 2012. 1

[31] Guolan Lu and Baowei Fei. Medical hyperspectral imaging: a review. *Journal of biomedical optics*, 19(1):010901, 2014. 1

[32] Renfu Lu and Yud-Ren Chen. Hyperspectral imaging for safety inspection of food and agricultural products. In *Pathogen Detection and Remediation for Safe Eating*, volume 3544, pages 121–133. International Society for Optics and Photonics, 1999. 1

[33] Jiawei Ma, Xiao-Yang Liu, Zheng Shou, and Xin Yuan. Deep tensor admm-net for snapshot compressive imaging. In *ICCV*, 2019. 3

[34] David J. C. MacKay. *Bayesian Methods for Adaptive Models*. PhD thesis, California Institute of Technology, 1992. 3

[35] Matthew MacKay, Paul Vicol, Jon Lorraine, David Duvenaud, and Roger Grosse. Self-tuning networks: Bilevel optimization of hyperparameters using structured best-response functions. *arXiv preprint arXiv:1903.03088*, 2019. 5

[36] Ziyi Meng, Shirin Jalali, and Xin Yuan. Gap-net for snapshot compressive imaging. *arXiv preprint arXiv:2012.08364*, 2020. 17

[37] Ziyi Meng, Jiawei Ma, and Xin Yuan. End-to-end low cost compressive spectral imaging with spatial-spectral self-attention. In *ECCV*, 2020. 1, 2, 3, 4, 6, 7, 8, 9, 10, 17, 18

[38] Ziyi Meng, Mu Qiao, Jiawei Ma, Zhenming Yu, Kun Xu, and Xin Yuan. Snapshot multispectral endomicroscopy. *Optics Letters*, 2020. 3, 7, 9, 17, 18

[39] Ziyi Meng, Zhenming Yu, Kun Xu, and Xin Yuan. Self-supervised neural networks for spectral snapshot compressive imaging. In *ICCV*, 2021. 1, 3, 7, 8, 9, 10, 11, 17

[40] Xin Miao, Xin Yuan, Yunchen Pu, and Vassilis Athitsos. l-net: Reconstruct hyperspectral images from a snapshot measurement. In *ICCV*, 2019. 3, 9

[41] Mu Qiao, Xuan Liu, and Xin Yuan. Snapshot spatial–temporal compressive imaging. *Optics letters*, 2020. 3

[42] Mu Qiao, Ziyi Meng, Jiawei Ma, and Xin Yuan. Deep learning for video compressive sensing. *APL Photonics*, 2020. 3

[43] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? In *NeurIPS*, 2021. 13

[44] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 2008. 6

[45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 13

[46] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *JMLR*, 2010. 5

[47] Jiamian Wang, Yulun Zhang, Xin Yuan, Yun Fu, and Zhiqiang Tao. A new backbone for hyperspectral image reconstruction. *arXiv preprint arXiv:2108.07739*, 2021. 1, 3, 4, 6, 7, 8, 9, 10, 12, 13, 14, 15, 17

[48] Lizhi Wang, Chen Sun, Ying Fu, Min H Kim, and Hua Huang. Hyperspectral image reconstruction using a deep spatial-spectral prior. In *CVPR*, 2019. 3, 9

[49] Lizhi Wang, Chen Sun, Maoqing Zhang, Ying Fu, and Hua Huang. Dnu: Deep non-local unrolling for computational spectral imaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1661–1671, 2020. 1

[50] Lizhi Wang, Zhiwei Xiong, Guangming Shi, Feng Wu, and Wenjun Zeng. Adaptive nonlocal sparse representation for dual-camera compressive hyperspectral imaging. *TPAMI*, 2017. 1

[51] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *TIP*, 2004. 7

[52] Andrew Gordon Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. *arXiv preprint arXiv:2002.08791*, 2020. 5

[53] F. Yasuma, T. Mitsunaga, D. Iso, and S.K. Nayar. Generalized Assorted Pixel Camera: Post-Capture Control of Resolution, Dynamic Range and Spectrum. Technical report, 2008. 17

[54] Xin Yuan. Generalized alternating projection based total variation minimization for compressive sensing. In *ICIP*, 2016. 3

[55] Xin Yuan, David J. Brady, and Aggelos K. Katsaggelos. Snapshot compressive imaging: Theory, algorithms, and applications. *IEEE Signal Processing Magazine*, 2021. 1, 4

[56] Xin Yuan, Yang Liu, Jinli Suo, and Qionghai Dai. Plug-and-play algorithms for large-scale snapshot compressive imaging. In *CVPR*, 2020. 3

[57] Xin Yuan, Yang Liu, Jinli Suo, Fredo Durand, and Qionghai Dai. Plug-and-play algorithms for video snapshot compressive imaging. *TPAMI*, 2021. 1

[58] Yuan Yuan, Xiangtao Zheng, and Xiaoqiang Lu. Hyperspectral image superresolution by transfer learning. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 10(5):1963–1974, 2017. 1

[59] Shipeng Zhang, Lizhi Wang, Ying Fu, Xiaoming Zhong, and Hua Huang. Computational hyperspectral imaging based on dimension-discriminative low-rank tensor recovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10183–10192, 2019. 1

[60] Tao Zhang, Ying Fu, Lizhi Wang, and Hua Huang. Hyperspectral image reconstruction using deep external and internal learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8559–8568, 2019. 1

[61] Yulun Zhang, Kai Li, Kunpeng Li, and Yun Fu. Mr image super-resolution with squeeze and excitation reasoning attention network. In *CVPR*, 2021. 6

[62] Yunhao Zou, Ying Fu, Yinqiang Zheng, and Wei Li. Csr-net: Camera spectral response network for dimensionality reduction and classification in hyperspectral imagery. *Remote Sensing*, 12(20):3294–3314, 2020. 1