

SpectraLift: Physics-Guided Spectral-Inversion Network for Self-Supervised Hyperspectral Image Super-Resolution

Ritik Shah
University of Massachusetts
Amherst, MA 01003
rgshah@umass.edu

Marco F. Duarte
University of Massachusetts
Amherst, MA 01003
mduarte@umass.edu

July 18, 2025

Abstract

High-spatial-resolution hyperspectral images (HSI) are essential for applications such as remote sensing and medical imaging, yet HSI sensors inherently trade spatial detail for spectral richness. Fusing high-spatial-resolution multispectral images (HR-MSI) with low-spatial-resolution hyperspectral images (LR-HSI) is a promising route to recover fine spatial structures without sacrificing spectral fidelity. Most state-of-the-art methods for HSI-MSI fusion demand point spread function (PSF) calibration or ground truth high resolution HSI (HR-HSI), both of which are impractical to obtain in real world settings. We present SpectraLift, a fully self-supervised framework that fuses LR-HSI and HR-MSI inputs using only the MSI’s Spectral Response Function (SRF). SpectraLift trains a lightweight per-pixel multi-layer perceptron (MLP) network using (i) a synthetic low-spatial-resolution multispectral image (LR-MSI) obtained by applying the SRF to the LR-HSI as input, (ii) the LR-HSI as the output, and (iii) an ℓ_1 spectral reconstruction loss between the estimated and true LR-HSI as the optimization objective. At inference, SpectraLift uses the trained network to map the HR-MSI pixel-wise into a HR-HSI estimate. SpectraLift converges in minutes, is agnostic to spatial blur and resolution, and outperforms state-of-the-art methods on PSNR, SAM, SSIM, and RMSE benchmarks.

1 Introduction

Hyperspectral imaging (HSI) captures hundreds of narrow spectral bands, enabling precise material discrimination. However, dispersive optics and photon splitting impose a trade-off: richer spectral resolution entails coarser spatial detail, which blurs small features, exacerbates mixed-pixel effects, and degrades classification accuracy. Fusion-based hyperspectral super-resolution (HSI-SR) addresses this trade-off by combining a LR-HSI with a HR-MSI of the same scene. Early approaches (pan-sharpening, Bayesian inference, matrix factorization, tensor decomposition) leverage handcrafted priors but often struggle with complex spectral structures. Recent deep learning-based methods deliver high reconstruction fidelity but depend on large labeled datasets, intricate architectures, and opaque “black-box” mappings that hinder transparency.

Unsupervised and hybrid schemes incorporate physics priors (e.g., PSF models, sparsity, low-rank constraints) to reduce supervision and improve interpretability. Yet they typically demand precise PSF knowledge – impractical given variable sensor blur – or rely on scarce ground-truth HR-HSI. State-of-the-art unsupervised methods try to estimate the PSF instead of strictly requiring precise knowledge of it. However, blurring in hyperspectral sensors is caused due to unknown and variable optical characteristics such as lens diffraction and sensor-specific aberrations, resulting in

an inherent ambiguity. It is thus extremely difficult to estimate the PSF that caused the blur, which limits real-world deployment and raises concerns about model robustness and trustworthiness.

We propose SpectraLift, a self-supervised, lightweight framework that fuses LR-HSI and HR-MSI using only the multispectral sensor’s SRF. Note that the SRF is manufacturer-defined, well-documented, and routinely used for radiometric and atmospheric correction; common approximations such as Gaussian estimations also work well in practice. A synthetic LR-MSI is generated for training purposes by applying the SRF to the LR-HSI. This LR-MSI and the original LR-HSI serve (pixelwise) as the input-output pair for a compact MLP that we call the Spectral Inversion Network (SIN). The SIN learns an implicit MSI to HSI spectral mapping for the training images by minimizing an ℓ_1 loss between the estimated and true LR-HSI. At inference, the trained SIN applies its learned spectral mapping to output a HR-HSI from a HR-MSI input.

By formulating HSI-SR as per-pixel spectral inversion, SpectraLift can avoid PSF estimation and any HR-HSI supervision while being agnostic to spatial blur and resolution. SpectraLift training converges quickly and yields fully interpretable spectral mappings. Through extensive experiments on multiple benchmark datasets, we show that SpectraLift consistently outperforms state-of-the-art supervised and unsupervised methods in PSNR, SAM, SSIM, and RMSE, while remaining lightweight.

During the final preparation of this manuscript, we became aware of SSSR [1], which also employs an MLP-based spectral mapping optimized with an MSE loss. While there is a similarity between our approach and SSSR in the use of an MLP for mapping MSI to HSI, crucially, unlike SSSR, SpectraLift does not rely on scene-specific PSF calibration. Instead, it leverages only the known SRF of the MSI sensor, enabling practical deployment in real-world scenarios where PSFs are unknown or difficult to estimate. We have included SSSR as one of the baselines in the experiments of Section 3. We believe that SpectraLift is the first self-supervised HSI-SR method that is fully agnostic to PSF information while achieving a more robust formulation and superior performance.

2 Proposed Method

The MSI sensor’s SRF acts as a band compression operator, mapping high-dimensional HSI spectra to lower-dimensional MSI measurements. SpectraLift, illustrated in Figure 1, formulates HSI-SR as the self-supervised, per-pixel inversion of this spectral compression. Given only a pair of LR-HSI and HR-MSI and the SRF (or its approximation), a compact MLP learns to reverse the SRF-induced degradation and reconstruct the HSI from an MSI input. SpectraLift uses only the LR-HSI image for training and the HR-MSI image for inference.

Let $\mathbf{Y} \in \mathbb{R}^{h \times w \times C}$ be the observed LR-HSI with C bands, and $\mathbf{M} \in \mathbb{R}^{H \times W \times c}$ be the HR-MSI with $c < C$ bands. Under classical sensor models, these relate to the unknown HR-HSI $\mathbf{X} \in \mathbb{R}^{H \times W \times C}$ via spatial and spectral degradations: $\mathbf{Y} \approx \mathcal{D}(\mathcal{H}(\mathbf{X}, K), r)$ and $\mathbf{M} \approx \mathbf{X} \times \mathbf{R}$, where $\mathcal{H}(\mathbf{X}, K)$ denotes spatial convolution of the HSI image \mathbf{X} with the PSF K , $\mathcal{D}(\mathbf{X}, r)$ denotes spatial downsampling of the HSI \mathbf{X} by a factor r , and $\mathbf{X} \times \mathbf{R}$ denotes the product of the HSI \mathbf{X} with the MSI SRF $\mathbf{R} \in \mathbb{R}^{C \times c}$ along the spectral (third) mode, i.e., for $i = 1, \dots, H$, $j = 1, \dots, W$, $m = 1, \dots, c$, we have

$$\mathbf{M}_{ijm} = \sum_{n=1}^C \mathbf{X}_{ijn} \mathbf{R}_{nm}.$$

2.1 Spectral Inversion Network (SIN)

Denote a given input MSI pixel from the image \mathbf{M} by $\mathbf{m} = \mathbf{M}_{ij}$, and the corresponding output HSI pixel estimate by $\hat{\mathbf{x}}$. SIN implements a per-pixel mapping $\hat{\mathbf{x}} = f_{\Theta}(\mathbf{m})$, where Θ denotes the set of

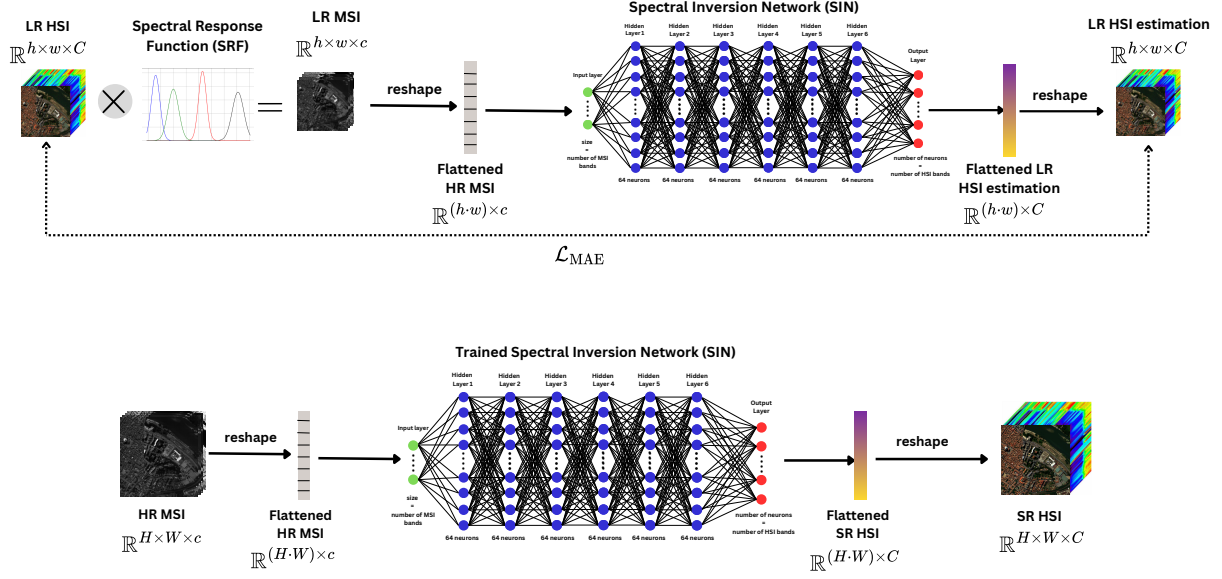


Figure 1: The SpectraLift pipelines. *Top*: self-supervised training of the Spectral Inversion Network (SIN) via SRF-based spectral inversion. *Bottom*: pixel-wise inference on HR-MSI with the SIN to produce the super resolved hyperspectral image (SR HSI)

SIN parameters, from an MSI spectrum $\mathbf{m} \in \mathbb{R}^c$ to an HSI spectrum $\hat{\mathbf{x}} \in \mathbb{R}^C$. The layers in SIN can be mathematically described as:

$$\begin{aligned}
\mathbf{x}^{(0)} &= \mathbf{m}, \\
\mathbf{x}^{(1)} &= \phi_1(\mathbf{x}^{(0)}), \\
\mathbf{x}^{(2)} &= \phi_2(\mathbf{x}^{(1)}) + \mathbf{x}^{(1)}, \\
\mathbf{x}^{(3)} &= \phi_3(\mathbf{x}^{(2)}), \\
\mathbf{x}^{(4)} &= \phi_4(\mathbf{x}^{(3)}) + \mathbf{x}^{(2)}, \\
\mathbf{x}^{(5)} &= \phi_5(\mathbf{x}^{(4)}), \\
\mathbf{x}^{(6)} &= \phi_6(\mathbf{x}^{(5)}) + \mathbf{x}^{(4)}, \\
\hat{\mathbf{x}} &= g_\theta(\mathbf{x}^{(6)}),
\end{aligned}$$

where each ϕ_i is a fully connected layer with 64 neurons and leaky ReLU activation, and g_θ is a fully connected output layer with linear activation. Residual/skip connections are incorporated in every other layer to enhance model expressiveness without increasing architectural complexity. Their inclusion leads to consistently improved spectral reconstruction and smoother convergence.

Training and Implementation: We synthesize the LR-MSI $\mathbf{Z} = \mathbf{Y} \times \mathbf{R}$ and optimize

$$\min_{\Theta} \frac{1}{hw} \sum_{i=1}^h \sum_{j=1}^w \|f_{\Theta}(\mathbf{Z}_{ij:}) - \mathbf{Y}_{ij:}\|_1,$$

training the SIN to invert the spectral degradation under an ℓ_1 loss. SpectraLift is implemented using the TensorFlow framework and trained with the Adam optimizer. For learning rate scheduling, we use the One-Cycle Learning Rate policy for the Washington DC Mall and Kennedy Space

Center datasets, as it accelerates convergence and promotes stability. For the Pavia University, Pavia Center, Botswana, and University of Houston datasets, we employ cosine annealing with restarts, which yielded superior convergence in these cases.

To achieve state-of-the-art performance, we tuned the scheduler parameters – specifically, the initial, maximum, and final learning rates for the One-Cycle policy, and the maximum and minimum learning rates for cosine annealing with restarts. However, SpectraLift achieves strong performance even with the default learning rate values provided in our pre-executed Jupyter notebooks. These notebooks, available in our GitHub repository <https://github.com/ritikgshah/SpectraLift>, contain the exact configurations used for each experiment.

Inference: At test time, the trained SIN is applied to each HR-MSI pixel $\mathbf{M}_{ij,:}$, yielding the pixels of the estimated HR-HSI $\hat{\mathbf{X}}$.

2.2 Rationale for Per-Pixel Formulation

SpectraLift’s design deliberately treats each pixel independently during spectral inversion. This choice is motivated by several practical considerations.

HSI-MSI temporal misalignment: Spatial priors typically assume co-registered HSI and MSI data, which is often unavailable in real-world scenarios due to temporal misalignments or platform differences, e.g., UAV vs. satellite acquisitions. Thus, moving objects may shift between acquisitions. This breaks the pixel-by-pixel correspondence assumed by many state of the art fusion methods and causes visible artifacts in the super-resolved outputs. SpectraLift entirely avoids this issue: it trains on synthetic LR-MSI derived solely from the LR-HSI and, at inference, processes each HR-MSI pixel independently, without ever requiring co-registration. Thus, no artifacts appear due to temporal offsets.

Dependence on PSF: Spatially coupled models are highly sensitive to unknown PSF variations, which can lead to artifacts when applied across diverse sensors or acquisition geometries. This design simplifies the training process and dramatically reduces computational overhead without sacrificing spectral fidelity, as evidenced by our results in Section 3.

Spatial Resolution and Blur Agnosticism: Because training uses only individual pixels, and both the input and target during training have very similar blurring, f_{Θ} has no mechanism to learn or depend on spatial structure. Consequently, the same learned mapping applies unchanged to any spatial sampling grid and optical acquisition setup, endowing SpectraLift with genuine spatial-resolution and quality agnosticism.

2.3 SIN Insights and Limitations

Exact SRF inversion: Under ideal conditions ($c \geq C$), \mathbf{R} admits a true inverse \mathbf{R}^{\dagger} (Moore-Penrose pseudo-inverse). In practice $c < C$, making inversion in \mathbb{R}^C ill-posed; however, real-world spectra lie on a low-dimensional non-linear manifold of intrinsic dimension $r \ll c, C$. The MLP learns a stable approximation $f_{\Theta} \approx \mathbf{R}^{\dagger}$ restricted to this manifold.

Physics-Guided Interpretability: Our sole physics assumption is the SRF \mathbf{R} , which is routinely provided by MSI manufacturers. This contrasts sharply with PSF-dependent unsupervised methods (and in particular [1]), since accurate PSF estimation is difficult or impossible in practice. With only a lightweight network and a clear ℓ_1 loss, SpectraLift converges rapidly and yields directly interpretable spectral inversion filters and avoids the opaque “black-box” nature of most current state-of-the-art models.

Limitations: A limitation of SpectraLift arises when the HR-MSI input contains only a single spectral band (e.g., a monochrome image). In this extreme case, predicting a full hyperspectral

signature from a scalar quantity is extremely ill-posed. In contrast, other state-of-the-art methods can exploit spatial context or impose assumptions about the similarity of spectra in neighboring pixels of the high-resolution output. While such assumptions can improve performance for highly degraded inputs, they may also reduce generalization when the input retains more spectral information. Furthermore, we note that single band HR images are rarely acquired in real world image fusion scenarios.

3 Experiments and Analysis

We evaluate SpectraLift on synthetic and real-world benchmarks. We compare against eight state-of-the-art baselines: four unsupervised (SSSR [1], MIAE [2], C2FF [3], SDP [4]) and four supervised (GuidedNet [5], FeINFN [6], FusFormer [7], MIMO-SST [8]). All datasets, precomputed results, and detailed instructions for replication are available at <https://github.com/ritikgshah/SpectraLift>. We have made every effort to ensure reproducibility by providing pre-executed Jupyter notebooks, Python scripts for end-to-end runs, configuration files, and environment replication scripts for seamless setup. This enables reviewers and practitioners to verify results and explore extensions with minimal friction.

Quality Metrics: For synthetic data, we evaluate SpectraLift using six widely adopted metrics that capture complementary aspects of reconstruction quality in hyperspectral image super-resolution:

- **Root Mean Squared Error (RMSE):** Measures the average pixel-wise difference between the reconstructed hyperspectral image $\hat{\mathbf{X}}$ and the ground truth \mathbf{X} . Lower RMSE indicates higher reconstruction fidelity.

$$\text{RMSE}(\mathbf{X}, \hat{\mathbf{X}}) = \sqrt{\frac{1}{HWC} \sum_{i=1}^H \sum_{j=1}^W \sum_{k=1}^C (\hat{\mathbf{X}}_{ijk} - \mathbf{X}_{ijk})^2}.$$

- **Peak Signal-to-Noise Ratio (PSNR):** Quantifies the ratio between the maximum possible pixel value and the power of the reconstruction error, expressed in decibels (dB). Higher PSNR indicates better perceptual quality.

$$\text{PSNR}(\mathbf{X}, \hat{\mathbf{X}}) = 10 \log_{10} \left(\frac{\text{MAX}^2}{\text{RMSE}(\mathbf{X}, \hat{\mathbf{X}})^2} \right),$$

where MAX is the maximum possible pixel value (e.g., $2^{16} - 1$ for 16-bit images).

- **Structural Similarity Index Measure (SSIM):** Assesses perceptual similarity by comparing luminance, contrast, and structure between $\hat{\mathbf{X}}$ and \mathbf{X} . Values close to 1 indicate high structural similarity.

$$\text{SSIM}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{(2\mu_X \mu_{\hat{\mathbf{X}}} + c_1)(2\sigma_{X\hat{\mathbf{X}}} + c_2)}{(\mu_X^2 + \mu_{\hat{\mathbf{X}}}^2 + c_1)(\sigma_X^2 + \sigma_{\hat{\mathbf{X}}}^2 + c_2)},$$

where μ , σ^2 , and $\sigma_{X\hat{\mathbf{X}}}$ are means, variances, and covariances, and c_1, c_2 are small constants to stabilize the denominator.

- **Universal Image Quality Index (UIQI):** Measures similarity in terms of luminance, contrast, and structure. Values range from -1 to 1 , with higher values indicating better quality.

$$\text{UIQI}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{4\sigma_{\mathbf{X}\hat{\mathbf{X}}}\mu_{\mathbf{X}}\mu_{\hat{\mathbf{X}}}}{(\sigma_{\mathbf{X}}^2 + \sigma_{\hat{\mathbf{X}}}^2)(\mu_{\mathbf{X}}^2 + \mu_{\hat{\mathbf{X}}}^2)}.$$

- **Erreur Relative Globale Adimensionnelle de Synthèse (ERGAS):** Provides a global indication of the relative error, normalized by the mean reflectance, and is commonly used in remote sensing. Lower ERGAS indicates higher reconstruction accuracy.

$$\text{ERGAS}(\mathbf{X}, \hat{\mathbf{X}}) = 100 \frac{r}{s} \sqrt{\frac{1}{C} \sum_{k=1}^C \frac{\text{RMSE}_k(\mathbf{X}, \hat{\mathbf{X}})^2}{\mu_k^2}},$$

where r/s is the ratio of spatial resolutions between $\hat{\mathbf{X}}$ and \mathbf{X} , and $\text{RMSE}_k(\mathbf{X}, \hat{\mathbf{X}})$ and μ_k are the RMSE and mean of band k , respectively.

- **Spectral Angle Mapper (SAM):** Computes the mean spectral angle (in degrees) between estimated and ground-truth spectral vectors at each pixel. It measures spectral similarity, with smaller angles indicating better fidelity. To ensure numerical stability, the arccos argument is clipped to avoid numerical issues.

$$\text{SAM}(\mathbf{X}, \hat{\mathbf{X}}) = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \left(\frac{180}{\pi} \arccos \left(\min \left(\frac{\langle \mathbf{X}_{ij\cdot}, \hat{\mathbf{X}}_{ij\cdot} \rangle}{\|\mathbf{X}_{ij\cdot}\|_2 \|\hat{\mathbf{X}}_{ij\cdot}\|_2 + \epsilon}, 1 - \delta \right) \right) \right),$$

where $\mathbf{X}_{ij\cdot}$ and $\hat{\mathbf{X}}_{ij\cdot}$ are the spectral vectors at pixel (i, j) , and ϵ and δ are small constants to avoid division by zero and numerical overflow, respectively (e.g., $\epsilon = 10^{-8}$, $\delta = 10^{-9}$).

We also profile execution time, number of parameters, inference peak GPU memory used and inference FLOPs to assess model complexity. For the inference GPU memory used and inference FLOPs we consider these for a single forward pass of the inputs through the model.

Synthetic Datasets: Following Wald’s protocol [9], for each ground-truth HSI (GT): Washington DC Mall (DC), Kennedy Space Center (KSC), Botswana, Pavia University (Pavia U), and Pavia Center (Pavia), we generate LR-HSI by first applying spatial convolution of the GT with one of ten PSFs shown in Figure 2 (Gaussian, Kolmogorov, Airy, Moffat, Sinc, Lorentzian Squared, Hermite, Parabolic, Gabor, Delta), all with kernel size of (15,15). We then downsample the convolution output by $r \in \{4, 8, 16, 32\}$, and finally add Gaussian white noise with SNR matched to r : (4, 35 dB), (8, 30 dB), (16, 25 dB), (32, 20 dB), yielding $10 \times 4 = 40$ LR-HSIs for each GT. HR-MSI are synthesized by applying SRFs for $b \in \{1, 3, 4\}$ (IKONOS), $b = 8$ (WorldView-2), and $b = 16$ (WorldView-3), then adding 40 dB Gaussian white noise, producing 5 HR-MSIs for each GT.

We consider 80 LR-HSI/HR-MSI pairs per GT: 10 PSFs \times 8 representative (r, b) configurations $\{(4,4), (8,4), (16,4), (32,4), (8,1), (8,3), (8,8), (8,16)\}$. We apply all the degradations to the GT normalized between $[0,1]$. Supervised methods are trained on 75% crops and tested on 25% crops for each LR-HSI-HR-MSI input pair. Spatial dimensions for the training and testing crops are constrained to be multiples of 32, ensuring compatibility with models that require integer downsampling without residual pixels. Unsupervised methods use full images but report metrics on the same 25% test regions. To support a more fair comparison across supervised and unsupervised

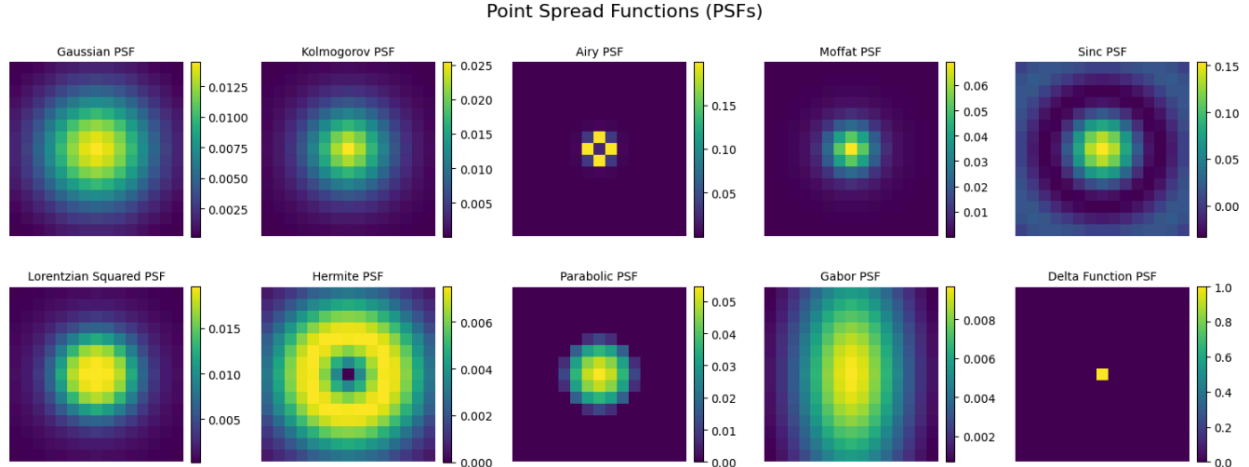


Figure 2: Point Spread Functions used for Synthetic LR HSI generation

approaches, the latter are also given access to the PSF and SRF used for LR-HSI and HR-MSI generation.

Tables 1-5 report mean quality and complexity metrics over these 80 LR-HSI/HR-MSI pairs for the Washington DC Mall, Kennedy Space Center, Pavia University, Pavia Center, and Botswana benchmarks, respectively, while Tables 6-10 show the same quality metrics when the single-band HR image is omitted. We do not include the model complexity metrics in tables 6-10 since they have almost no difference compared to the ones reported in tables 1-5. Across all five synthetic datasets, SpectraLift outperforms unsupervised baselines on almost all quality measures. Supervised methods enjoyed a modest advantage – they were trained on a 75% crop of the same scene (distinct from the 25% test region), allowing them to learn scene-specific spectral behaviors. Despite lacking this advantage, on naturally occurring scenes (Washington DC Mall, Kennedy Space Center, Botswana) SpectraLift even surpasses most supervised approaches, attaining the highest PSNR (35.96dB) and lowest SAM (3.22) on DC, state-of-the-art RMSE (0.04472), PSNR (27.10dB), SSIM (0.925), ERGAS (8.41) and SAM (8.70) on KSC, and the highest SSIM (0.974) and lowest SAM (1.46) on Botswana. Only MIMO-SST edges ahead, but with a more complex method relying on "black box" models with knowledge of the scene specific spectra. On urban scenes with more man-made materials (Pavia University, Pavia Center), scene-specific spectral information gives supervised models like MIMO-SST and FeINFN an edge, but SpectraLift still outperforms unsupervised baselines, demonstrating its advantage when comparing methods with equal training knowledge. This advantage would explain their subpar metrics on DC, KSC, and Botswana combined with their good performance against unsupervised methods on the Pavia University and Pavia Center datasets. When we omit the single-band HR image – reflecting more realistic operational scenarios – SpectraLift rises to the top on the natural scenes: it secures best RMSE, PSNR, SSIM and SAM on DC (Table 6), best results for all metrics (except UIQI) on KSC (Table 7), and best SSIM, UIQI, and SAM on Botswana (Table 10), while remaining highly competitive on the urban benchmarks. These results underscore its effectiveness when the HR input has more than a single band.

Table 1: Quality measures for the Washington DC Mall dataset: Mean value of 80 LR-HSI/HR-MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓	Time (s) ↓	Params (M) ↓	FLOPs (G) ↓	GPU Mem (MB) ↓
MIAE	0.03611	30.28	0.934	0.968	5.33	5.21	212.69	0.0218	8.46	707.54
C2FF	0.02681	34.74	0.960	<i>0.975</i>	<i>4.73</i>	3.53	<i>78.54</i>	0.0979	<i>7.29</i>	948.88
SDP	0.03029	31.99	0.939	0.907	16.59	4.29	799.75	6.4546	511.52	5370.37
SSSR	0.05024	27.98	0.908	0.938	7.49	5.83	67.24	<i>0.0334</i>	0.000067	286.31
GuidedNet	0.03494	29.98	0.919	0.872	30.30	4.81	398.07	6.7111	178.41	81.52
FeINFN	0.02558	32.66	0.966	0.956	7.83	3.79	3018.80	3.7021	251.70	982.47
FusFormer	0.02555	32.55	0.876	0.688	17.34	3.89	11110.07	0.1883	946.28	6328.36
MIMO-SST	0.02193	<i>35.21</i>	0.969	0.982	2.92	<i>3.29</i>	376.34	2.1879	56.37	393.74
SpectraLift	<i>0.02344</i>	35.96	<i>0.967</i>	0.969	5.87	3.22	91.20	0.0336	26.28	<i>286.31</i>

Table 2: Quality measures for the Kennedy Space Center dataset: Mean value of 80 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓	Time (s) ↓	Params (M) ↓	FLOPs (G) ↓	GPU Mem (MB) ↓
MIAE	0.04855	26.41	0.907	0.962	8.54	8.88	189.24	0.0996	31.12	711.06
C2FF	<i>0.04540</i>	<i>26.97</i>	<i>0.922</i>	0.956	<i>8.50</i>	8.90	<i>77.24</i>	0.0913	<i>5.42</i>	707.95
SDP	0.04659	26.73	0.906	0.949	8.63	9.65	643.65	6.1322	388.55	4209.39
SSSR	0.05767	25.05	0.851	0.930	9.54	10.08	70.70	0.0284	0.000057	211.06
GuidedNet	0.04951	26.11	0.913	0.956	9.17	9.25	374.44	6.0301	176.43	71.32
FeINFN	0.05013	26.02	0.902	0.932	9.07	10.77	2206.34	3.6520	261.73	984.15
FusFormer	0.05124	25.82	0.888	0.921	9.26	11.57	6976.85	0.1811	787.98	6322.09
MIMO-SST	0.04660	26.68	0.917	0.941	8.70	9.47	305.21	2.1361	56.79	363.68
SpectraLift	0.04472	27.10	0.925	<i>0.957</i>	8.41	8.70	92.24	<i>0.0327</i>	20.42	<i>211.06</i>

Table 3: Quality measures for the Pavia University dataset: Mean value of 80 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓	Time (s) ↓	Params (M) ↓	FLOPs (G) ↓	GPU Mem (MB) ↓
MIAE	0.03775	29.62	0.906	0.984	3.64	5.01	99.99	0.0880	18.11	445.91
C2FF	0.03618	31.70	0.914	0.976	3.36	4.63	<i>73.24</i>	0.0590	<i>2.25</i>	293.47
SDP	0.03088	31.88	0.917	0.982	3.14	4.29	385.58	4.6211	192.66	2491.03
SSSR	0.05851	26.31	0.829	0.950	5.46	5.68	66.48	0.0099	0.000020	81.49
GuidedNet	0.03060	30.80	0.920	<i>0.986</i>	3.36	4.19	243.72	3.4673	68.62	22.63
FeINFN	<i>0.02789</i>	31.82	<i>0.933</i>	0.982	3.04	<i>4.18</i>	1187.35	3.4082	127.63	381.94
FusFormer	0.03017	31.00	0.929	0.983	3.32	4.22	5556.52	0.1460	471.06	6231.11
MIMO-SST	0.02236	34.34	0.945	0.989	2.49	3.33	221.63	1.8836	22.93	153.75
SpectraLift	0.03073	<i>32.79</i>	0.928	0.980	<i>2.99</i>	<i>4.18</i>	91.03	<i>0.0279</i>	11.52	<i>81.49</i>

Table 4: Quality measures for the Pavia Center dataset: Mean value of 80 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓	Time (s) ↓	Params (M) ↓	FLOPs (G) ↓	GPU Mem (MB) ↓
MIAE	0.04344	28.67	0.901	0.984	3.58	6.99	314.86	0.0878	68.30	1677.58
C2FF	0.02861	33.38	0.956	0.988	2.46	4.90	<i>85.77</i>	0.0586	<i>8.43</i>	1094.10
SDP	0.02979	32.11	0.952	0.988	2.61	5.58	1085.04	4.6011	724.77	9375.14
SSSR	0.06399	25.67	0.839	0.958	5.26	6.28	67.55	0.0098	0.000020	<i>304.91</i>
GuidedNet	0.02929	31.10	0.954	<i>0.992</i>	2.68	4.98	723.75	3.4408	263.18	92.44
FeINFN	0.02532	32.48	<i>0.964</i>	<i>0.992</i>	2.33	4.83	5131.20	3.4048	490.85	1459.95
FusFormer	<i>0.02262</i>	33.58	0.968	0.993	<i>2.09</i>	<i>4.49</i>	19432.89	0.1455	1884.16	6321.71
MIMO-SST	0.02168	35.02	0.963	0.993	2.01	4.19	706.33	1.8802	87.91	575.90
SpectraLift	0.02716	<i>33.88</i>	0.958	0.988	2.36	4.80	90.84	<i>0.0278</i>	43.41	<i>304.91</i>

Table 5: Quality measures for the Botswana dataset: Mean value of 80 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓	Time (s) ↓	Params (M) ↓	FLOPs (G) ↓	GPU Mem (MB) ↓
MIAE	0.01838	35.06	0.954	<i>0.997</i>	1.75	1.77	162.78	0.0190	7.09	552.00
C2FF	0.01572	36.51	0.970	0.998	<i>1.48</i>	1.59	74.82	0.0776	<i>5.48</i>	716.48
SDP	0.01863	34.82	0.951	0.995	2.29	1.97	698.01	5.4788	416.77	4835.41
SSSR	0.02050	34.20	0.953	0.996	1.93	1.69	<i>86.12</i>	<i>0.0194</i>	0.000039	208.97
GuidedNet	0.01970	34.25	0.927	0.986	4.11	1.93	439.14	4.7894	162.62	68.65
FeINFN	0.02250	33.44	0.946	0.995	1.95	2.53	2718.46	3.5485	266.91	919.15
FusFormer	0.02937	31.10	0.780	0.815	11.23	3.35	8333.27	0.1662	944.11	6299.68
MIMO-SST	0.01346	37.76	<i>0.973</i>	0.998	1.35	<i>1.54</i>	365.49	2.0289	53.87	367.98
SpectraLift	<i>0.01452</i>	<i>37.04</i>	0.974	<i>0.997</i>	1.60	1.46	93.97	0.0306	23.03	<i>209.00</i>

Table 6: Quality measures for the Washington DC dataset without 1 band HR image: Mean value of 70 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓
MIAE	0.03239	31.14	0.948	0.971	5.03	4.92
C2FF	0.01569	36.89	<i>0.988</i>	0.989	<i>3.74</i>	<i>2.17</i>
SDP	0.02214	33.49	0.961	0.914	16.29	3.49
SSSR	0.04269	29.16	0.926	0.947	7.01	4.75
Guided Net	0.03173	30.72	0.928	0.870	31.35	4.43
FeINFN	0.02142	33.72	0.978	0.959	7.58	3.42
Fus Former	0.02193	33.51	0.886	0.689	17.30	3.56
MIMO-SST	<i>0.01457</i>	<i>37.00</i>	0.991	0.989	2.16	2.32
Spectra Lift	0.01266	38.22	0.991	<i>0.980</i>	5.13	1.90

Table 7: Quality measures for the Kennedy Space Center dataset without 1 band HR image: Mean value of 70 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓
MIAE	0.04736	26.62	0.915	0.963	8.40	8.75
C2FF	<i>0.04247</i>	<i>27.45</i>	<i>0.936</i>	0.959	<i>8.14</i>	<i>8.38</i>
SDP	0.04386	27.16	0.923	0.953	8.31	9.16
SSSR	0.05608	25.31	0.857	0.933	9.29	9.55
Guided Net	0.04851	26.28	0.920	0.957	9.06	9.09
FeINFN	0.04926	26.16	0.908	0.932	8.98	10.74
Fus Former	0.05059	25.93	0.892	0.921	9.19	11.54
MIMO-SST	0.04450	27.03	0.933	0.943	8.46	9.09
Spectra Lift	0.04178	27.58	0.939	<i>0.960</i>	8.05	8.15

Table 8: Quality measures for the Pavia University dataset without 1 band HR image: Mean value of 70 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓
MIAE	0.03348	30.50	0.924	0.986	3.29	4.81
C2FF	0.02485	33.55	0.947	0.987	2.63	3.55
SDP	0.02170	33.47	0.950	0.990	2.54	3.49
SSSR	0.05040	27.38	0.855	0.959	4.99	4.79
Guided Net	0.02662	31.67	0.933	0.987	3.17	3.90
FeINFN	0.02412	32.74	0.948	0.983	2.85	3.94
Fus Former	0.02616	31.89	0.943	0.985	3.14	3.97
MIMO-SST	0.01632	35.85	0.965	0.992	2.12	2.75
Spectra Lift	<i>0.01871</i>	<i>34.79</i>	<i>0.963</i>	<i>0.991</i>	<i>2.21</i>	<i>3.07</i>

Table 9: Quality measures for the Pavia Center dataset without 1 band HR image: Mean value of 70 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓
MIAE	0.03932	29.47	0.917	0.986	3.27	6.67
C2FF	0.01770	35.36	0.978	0.996	1.73	3.43
SDP	0.02081	33.72	0.971	0.995	2.02	4.31
SSSR	0.05790	26.56	0.849	0.963	4.89	5.02
Guided Net	0.02563	31.95	0.962	0.994	2.46	4.63
FeINFN	0.02208	33.36	0.970	0.993	2.14	4.59
Fus Former	0.01920	34.57	0.975	0.994	1.88	4.24
MIMO-SST	0.01474	36.73	0.981	<i>0.996</i>	1.56	3.19
Spectra Lift	<i>0.01607</i>	<i>35.93</i>	<i>0.980</i>	0.997	<i>1.61</i>	<i>3.30</i>

Table 10: Quality measures for the Botswana dataset without 1 band HR image: Mean value of 70 LR HSI-HR MSI configurations (**best in bold**), *second best in italics*.

Method	RMSE ↓	PSNR ↑	SSIM ↑	UIQI ↑	ERGAS ↓	SAM ↓
MIAE	0.01725	35.53	0.962	<i>0.997</i>	1.65	1.71
C2FF	0.01428	37.19	<i>0.974</i>	0.998	<i>1.34</i>	<i>1.46</i>
SDP	0.01742	35.32	0.955	0.995	2.21	1.88
SSSR	0.01953	34.61	0.956	0.996	1.83	1.58
Guided Net	0.01918	34.48	0.929	0.985	4.17	1.89
FeINFN	0.02250	33.52	0.948	0.995	1.94	2.56
Fus Former	0.02868	31.26	0.787	0.814	11.29	3.33
MIMO-SST	0.01250	38.31	0.979	0.998	1.26	<i>1.46</i>
Spectra Lift	<i>0.01290</i>	<i>37.80</i>	0.979	0.998	1.47	1.31

While SpectraLift achieves strong performance on most quality metrics (PSNR, SSIM, SAM, RMSE), it performs comparatively less well on ERGAS and UIQI in certain cases. This difference reflects the nature of these metrics: ERGAS and UIQI emphasize global luminance and contrast consistency, making them sensitive to small biases and variance shifts. In contrast, RMSE and SAM directly assess pixel-wise accuracy and spectral shape preservation, while SSIM captures spatial detail but focuses on local structural similarity. To maintain a lightweight, physics-grounded, and PSF-agnostic formulation, SpectraLift forgoes scene-specific PSF calibration and black-box spatial modeling—choices that can boost global metrics (ERGAS and UIQI) under synthetic conditions but often reduce robustness in real-world scenarios. Future work can extend this foundation to better address global consistency while preserving SpectraLift’s core strengths.

Beyond reconstruction accuracy, SpectraLift offers clear advantages in efficiency. Its training time is comparable to the fastest methods and substantially lower than that of more complex architectures, while its inference cost in FLOPs is close to the most lightweight approaches. This balance of high performance and low complexity underscores the practicality of SpectraLift for real-world hyperspectral super-resolution tasks.

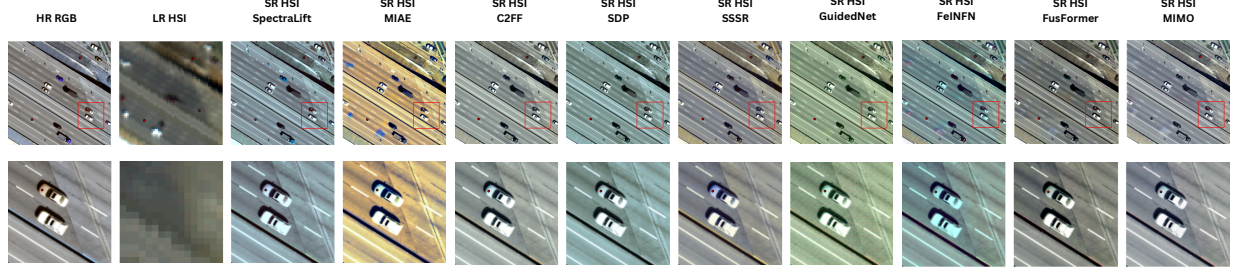
Real-World Data: We further evaluate SpectraLift on the University of Houston (UH) dataset from the 2018 IEEE GRSS Data Fusion Contest [10], which provides an LR-HSI ($4172 \times 1202 \times 50$) and an HR-RGB ($83440 \times 24040 \times 3$). We extract two distinct $64 \times 64 \times 50$ regions from the LR-HSI and their aligned $1280 \times 1280 \times 3$ HR-RGB patches, downsample the RGB to $512 \times 512 \times 3$ ($r = 8$), and apply all unsupervised methods (including SpectraLift) to recover a $512 \times 512 \times 50$ HR-HSI. Testing is restricted to a small part of the UH dataset in order to ensure that baseline methods could process the inputs in a single batch; if batch processing were used, the output of each batch would have to be stitched together and visible seam artifacts would appear (as can be seen for FusFormer in Fig. 3a and 3c, which could not process each region in a single batch). SpectraLift’s pixel-wise mapping is immune to such effects when inferred using batch processing.

Since UH lacks HR-HSI ground truth and its PSF/SRF are unknown, we use the IKONOS RGB SRF and let each unsupervised method estimate the PSF. Supervised baselines were trained on Pavia Center HSI with inputs being generated by following Wald’s protocol for $(r, b) = (8, 3)$ via three overlapping $512 \times 512 \times 3$ HR-RGB patches and the corresponding LR-HSI patches, and then evaluated on the UH test pairs. For FusFormer, 53 non-overlapping HR-RGB patches of size $128 \times 128 \times 3$ and the corresponding LR-HSI patches had to be used to meet memory limits.

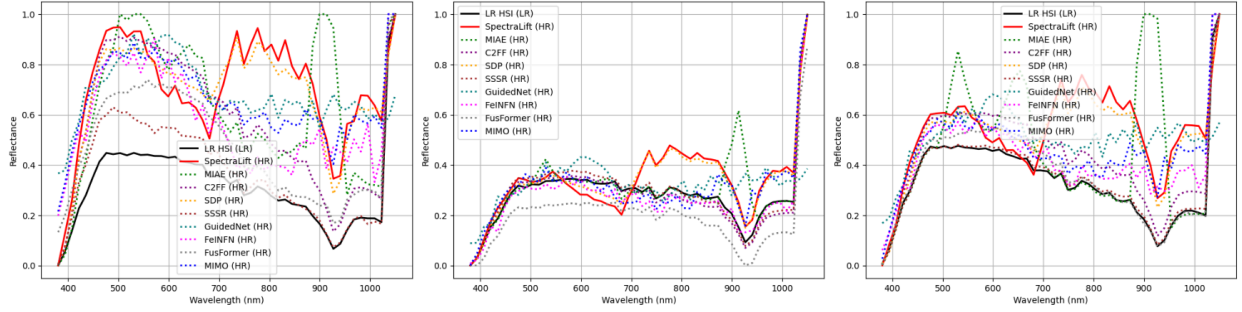
It is important to note that we do not report proxy metrics such as Quality with No Reference (QNR) for the UH dataset. These metrics, while useful in certain controlled scenarios, are highly sensitive to the assumed PSF and SRF characteristics. In our case, no information regarding the PSF or exact SRF of the MSI acquisition is available, making such metrics unreliable and potentially misleading. Instead, we assess real-world performance through visual comparisons (Fig. 3a, 3c) and spectral fidelity plots (Fig. 3b, 3d).

Figure 3b plots the spectral signatures at three manually selected pixels – car, bare earth, and highway – marked in red in Figure 3a; Figure 3d does the same for concrete (residential building), metal (car), and grass pixels in Figure 3c. Because the spatial downsampling factor is $r = 8$, each LR-HSI pixel $\mathbf{Y}_{i,j,:}$ corresponds to an 8×8 patch in the HR grid. To compare spectra, we extract the spectrum of the LR pixel at (i, j) and pair it with the spectrum of the HR-HSI at the top-right corner of its corresponding block, namely $(8i, 8j)$. For example, the LR pixel at $(8, 8)$ maps to the HR position $(64, 64)$, so in Figures 3b and 3d we plot LR-HSI pixel $(8, 8)$ against HR-HSI pixel $(64, 64)$ for a direct spectral fidelity comparison.

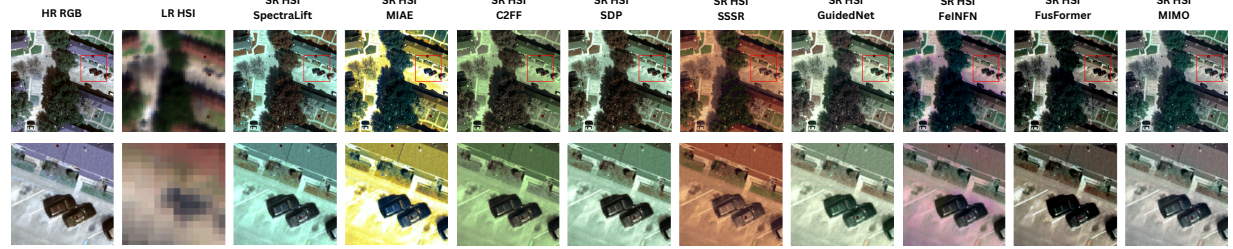
SpectraLift delivers crisp, high-resolution reconstructions in both UH test scenes, whereas competing methods exhibit various spatial and color artifacts. In Figure 3c, several baselines (MIAE, C2FF, SSSR, and FeINFN) introduce noticeable color shifts, producing HR-HSIs with unnatural



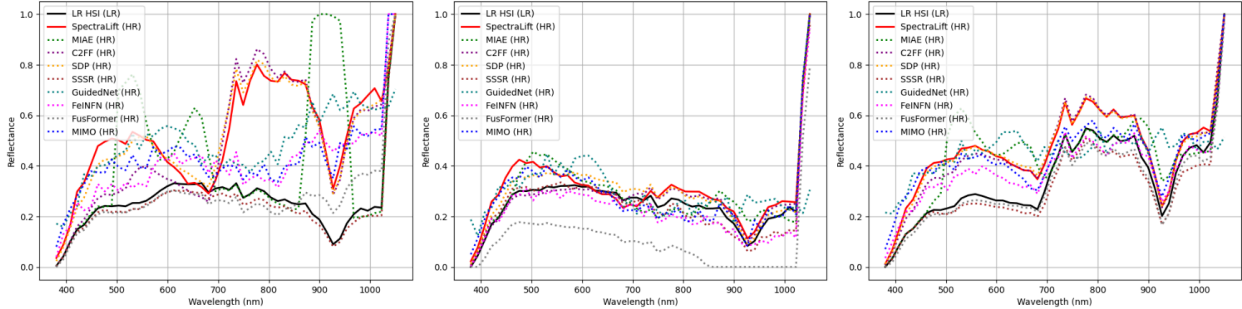
(a) UH super-resolved images for test scene 1. Second row shows a zoomed-in crop of the region marked in red.



(b) UH spectra for test scene 1. Left: Car, Center: Bare earth, Right: Highway.



(c) UH super-resolved images for test scene 2. Second row shows a zoomed-in crop of the region marked in red.



(d) UH spectra for test scene 2. Left: Building roof, Center: Car, Right: Grass.

Figure 3: University of Houston super-resolved results and corresponding spectra for two test scenes. (a, c) Super-resolved images with zoomed-in crops. (b, d) Spectral plots for selected regions.

tints – while SpectraLift faithfully preserves the true colors of the underlying HR-RGB. Some baselines produce HR-HSIs with spatial inconsistencies as well, such as Guided Net introducing small

Table 11: Ablation study of SpectraLift on the DC dataset using only Gaussian PSF (**best in bold**).

Model	RMSE ↓	PSNR ↑	SSIM ↑	SAM ↓	Params (M) ↓	FLOPs (G) ↓
Baseline	0.0243	35.50	0.9652	3.3500	0.0336	26.28
No skip connections	0.0247	35.29	0.9643	3.4155	0.0336	26.21
No learning rate scheduler	0.0256	34.84	0.9622	3.4887	0.0336	26.28
MSE loss	0.0244	35.14	0.9603	3.4103	0.0336	26.28
Cosine similarity loss	0.1549	17.15	0.7537	3.7929	0.0336	26.28
ReLU	0.0248	35.15	0.9649	3.3980	0.0336	26.28
GeLU	0.0250	35.04	0.9662	3.4001	0.0336	26.88
8 hidden layers	0.0244	35.44	0.9650	3.3585	0.0420	32.79
4 hidden layers	0.0244	35.47	0.9652	3.3690	0.0253	19.77
2 hidden layers	0.0248	35.19	0.9661	3.3898	0.0170	13.25
1 hidden layer	0.0267	34.53	0.9647	3.5704	0.0128	9.98
32 hidden layer size	0.0246	35.16	0.9658	3.3561	0.0118	9.15
128 hidden layer size	0.0276	34.68	0.9562	3.7911	0.1080	84.68
Linear map	0.0336	31.60	0.9623	4.2414	0.0012	0.90

checkerboard artifacts that are noticeable when zoomed in. Temporal misalignment between the LR-HSI and HR-RGB is most pronounced in Figure 3a, where moving cars occupy different positions in the two inputs. Methods that fuse the modalities jointly (MIAE, FeINFN, FusFormer, and MIMO) produce ghosting and smearing around those vehicles. Conversely, SpectraLift’s HR-HSIs are sharp, align strongly with the HR-RGB, and preserve the scene’s spatial structure.

Across both UH scenes, SpectraLift’s per-pixel inversion recovers spectral shapes that mirror those of the LR-HSI, particularly for naturally occurring materials such as grass and bare earth (Figures 3d, 3b). While absolute reflectance values deviate – an expected consequence of the ill-posed spectral unmixing – the overall spectral structure is faithfully preserved, supporting downstream analysis tasks that rely on relative band signatures. For certain man-made structures like highways and building rooftops, other baselines produce HR-HSIs whose spectra better aligns with that of the LR-HSI.

4 Ablation

To quantify the impact of our architectural and training choices, we evaluate a suite of SpectraLift/SIN variants on the Washington DC Mall benchmark using only the Gaussian PSF; results are given in Table 11. We perform the experiments using the same procedure as described in Section 3. All variants of SIN share the same hyperparameters and differ only in the component under test. Removing our selective residual links degrades reconstruction quality, confirming that skip connections stabilize the ill-posed spectral inversion by biasing the network when appropriate at almost no additional complexity cost. Disabling the learning-rate scheduler also harms convergence, while swapping the MAE loss for MSE causes a small drop in fidelity. By contrast, replacing MAE with a cosine-similarity loss collapses performance, underlining that we must penalize direct magnitude errors in the spectral domain. Changing activations – Leaky ReLU to ReLU or GeLU – has only a minor effect, showing robustness to the choice of nonlinearity.

SIN depth and width both matter: our chosen six-layer MLP achieves the best balance of accuracy and efficiency. An eight-layer variant offers no additional gains at substantially higher compute, whereas a four-layer or two-layer version remains competitive while having fewer param-

eters and FLOPs. Just using a MLP with a single hidden layer is also capable of achieving decent quality metrics at substantially lower computational cost. Reducing hidden-unit size further cuts complexity with only slight quality loss, while doubling hidden-unit size does not offer any benefits at much higher model complexity. A purely linear map (single dense layer without nonlinear activation) performs much worse, empirically validating that modeling hyperspectral signals on a near-linear manifold – rather than as a strictly linear mapping – is essential.

Overall, the baseline configuration is good for balancing the fidelity-efficiency trade-off, and users with tighter resource budgets can select the shallower or narrower variants to save compute with only marginal performance degradation.

5 Conclusion

We have introduced SpectraLift, a fully self-supervised HSI-SR framework that leverages the well-defined, readily available SRF (or its Gaussian approximations in cases where the exact SRF is not available) to directly invert spectral compression with a compact, per-pixel MLP. By casting fusion as an ℓ_1 spectral-inversion problem grounded in a clear physics model, SpectraLift avoids black-box approaches, and delivers interpretable mappings. Extensive evaluations across synthetic and real-world benchmarks show SpectraLift consistently matches or exceeds leading supervised methods on key quality metrics, while outperforming state-of-the-art unsupervised baselines. Its lightweight design and robustness to spatial blur and temporal misalignment further highlight its practical utility. SpectraLift establishes a foundation for lightweight, physics-guided spectral inversion and paves the way for future advances in efficient, transparent HSI-MSI fusion models for HSI-SR.

References

- [1] A. Rajaei, E. Abiri, and M. Helfroush, “Self-supervised spectral super-resolution for a fast hyperspectral and multispectral image fusion,” *Sci. Rep.*, vol. 14, no. 1, 2024.
- [2] J. Liu, Z. Wu, L. Xiao, and X.-J. Wu, “Model inspired autoencoder for unsupervised hyperspectral image super-resolution,” *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–12, 2022.
- [3] J. Li, K. Zheng, W. Liu, Z. Li, H. Yu, and L. Ni, “Model-guided coarse-to-fine fusion network for unsupervised hyperspectral image super-resolution,” *IEEE Geosci. Remote Sens. Letters*, vol. 20, pp. 1–5, 2023.
- [4] J. Liu, Z. Wu, and L. Xiao, “A spectral diffusion prior for unsupervised hyperspectral image super-resolution,” *IEEE Trans. Geosci. Remote Sens.*, vol. 62, pp. 1–13, 2024.
- [5] R. Ran, L.-J. Deng, T.-X. Jiang, J.-F. Hu, J. Chanussot, and G. Vivone, “GuidedNet: A general CNN fusion framework via high-resolution guidance for hyperspectral image super-resolution,” *IEEE Trans. Cybern.*, vol. 53, no. 7, pp. 4148–4161, 2023.
- [6] Y.-J. Liang, Z. Cao, S. Deng, H.-X. Dou, and L.-J. Deng, “Fourier-enhanced implicit neural fusion network for multispectral and hyperspectral image fusion,” in *Adv. Neural Inf. Proc. Syst.*, vol. 37, 2024, pp. 63 441–63 465.
- [7] J.-F. Hu, T.-Z. Huang, L.-J. Deng, H.-X. Dou, D. Hong, and G. Vivone, “Fusformer: A transformer-based fusion network for hyperspectral image super-resolution,” *IEEE Geosci. Remote Sens. Letters*, vol. 19, pp. 1–5, 2022.

- [8] J. Fang, J. Yang, A. Khader, and L. Xiao, “MIMO-SST: Multi-input multi-output spatial-spectral transformer for hyperspectral and multispectral image fusion,” *IEEE Trans. Geosci. Remote Sens.*, vol. 62, pp. 1–20, 2024.
- [9] T. Ranchin and L. Wald, “Fusion of high spatial and spectral resolution images: The ARSIS concept and its implementation,” *Photogramm. Eng. Remote Sens.*, vol. 66, no. 1, pp. 49–61, Jan. 2000.
- [10] Y. Xu, B. Du, L. Zhang, *et al.*, “Advanced multi-sensor optical remote sensing for urban land use and land cover classification: Outcome of the 2018 IEEE GRSS Data Fusion Contest,” *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 12, no. 6, pp. 1709–1724, 2019.