

GEWDiff: Geometric Enhanced Wavelet-based Diffusion Model for Hyperspectral Image Super-resolution

Sirui Wang¹, Jiang He¹, Natàlia Blasco Andreo², Xiao Xiang Zhu^{1,3*}

¹Technical University of Munich, Arcisstraße 21, 80333 Munich, Germany

²Universitat Autònoma de Barcelona, 08193 Cerdanyola del Vallès, Barcelona, Spain

³Munich Center for Machine Learning, 80333 Munich, Germany

sirui.wang@tum.de, Natalia.Blasco@uab.cat, jiang.he@tum.de, xiaoxiang.zhu@tum.de

Abstract

Improving the quality of hyperspectral images (HSIs), such as through super-resolution, is a crucial research area. However, generative modeling for HSIs presents several challenges. Due to their high spectral dimensionality, HSIs are too memory-intensive for direct input into conventional diffusion models. Furthermore, general generative models lack an understanding of the topological and geometric structures of ground objects in remote sensing imagery. In addition, most diffusion models optimize loss functions at the noise level, leading to a non-intuitive convergence behavior and suboptimal generation quality for complex data. To address these challenges, we propose a Geometric Enhanced Wavelet-based Diffusion Model (GEWDiff), a novel framework for reconstructing hyperspectral images at 4-times super-resolution. A wavelet-based encoder-decoder is introduced that efficiently compresses HSIs into a latent space while preserving spectral-spatial information. To avoid distortion during generation, we incorporate a geometry-enhanced diffusion process that preserves the geometric features. Furthermore, a multi-level loss function was designed to guide the diffusion process, promoting stable convergence and improved reconstruction fidelity. Our model demonstrated state-of-the-art results across multiple dimensions, including fidelity, spectral accuracy, visual realism, and clarity.

Code — <https://github.com/zhu-xlab/GEWDiff>

Introduction

Hyperspectral images (HSIs) offer a unique perspective by capturing continuous spectral features of ground objects. Despite advancements in research, the high costs and low coverage of super-resolution (SR) hyperspectral data limit their applications. Currently, open-access hyperspectral airborne data are predominantly regional, typically focused on several cities, increasing their exclusivity and cost. Hyperspectral satellites have better coverage but suffer from insufficient spatial resolution. Improving the spatial resolution of hyperspectral satellite images is, therefore, crucial to allow for fully harnessing the potential of hyperspectral data in Earth observation (EO). Many fusion models that combine hyperspectral and multispectral images (MSIs) have been developed. However, the fusion model cannot generate

HSIs in any region of interest without any prior knowledge offered by VHR RGB data. Most fusion models can obtain 10 m resolution hyperspectral data with MSIs, such as Sentinel 2, but to get a spatial resolution from 10 to 2.5 m, we focused on the single-image-generation methods for HSIs at 4-times super-resolution.

Traditional hyperspectral image (HSI) super-resolution approaches typically rely on interpolation techniques, such as nearest neighbor or bilinear interpolation. While straightforward and computationally efficient, these methods fail to capture the complex, nonlinear relationships present in high-dimensional spectral data. Recently, deep learning-based methods, such as convolutional neural networks (CNNs) and generative adversarial networks (GANs) (Sidorov and Yngve Hardeberg 2019; Li, Wang, and Li 2020; Hou et al. 2022; Shi et al. 2022; Li et al. 2020), have emerged as powerful alternatives. These models are capable of learning expressive spectral-spatial representations directly from data, leading to significant improvements in reconstruction accuracy over classical methods. Recent studies have shown the strong potential of transformers (Zhang et al. 2023; Chen, Zhang, and Zhang 2023; Yu et al. 2023; Su et al. 2025) in vision tasks, due to their ability to model long-range spatial and spectral dependencies. However, a common limitation shared by the aforementioned methods is their difficulty in generating rich textures and complex spatial structures, despite their strong performance in preserving spectral fidelity.

Diffusion models have recently demonstrated remarkable success in generating high-quality natural images, as seen in models such as the Stable Diffusion (Rombach et al. 2021) and DiT (Peebles and Xie 2023) frameworks. Motivated by their strong generative capabilities and robustness in modeling complex distributions, researchers have begun exploring their applicability to hyperspectral image super-resolution. For instance, SpectralDiff (Chen et al. 2023) and HSR-Diff (Wu et al. 2023) extend the diffusion paradigm directly to the hyperspectral domain, aiming to better model spectral-spatial correlations. However, despite recent progress, adapting diffusion models to hyperspectral image generation remains a significant challenge. Unlike natural or multispectral images, HSI data typically exhibit a much lower signal-to-noise ratio and higher spectral dimensionality, making it difficult to design diffusion architectures that balance spatial fidelity, spectral accuracy, and visual re-

*Corresponding author.

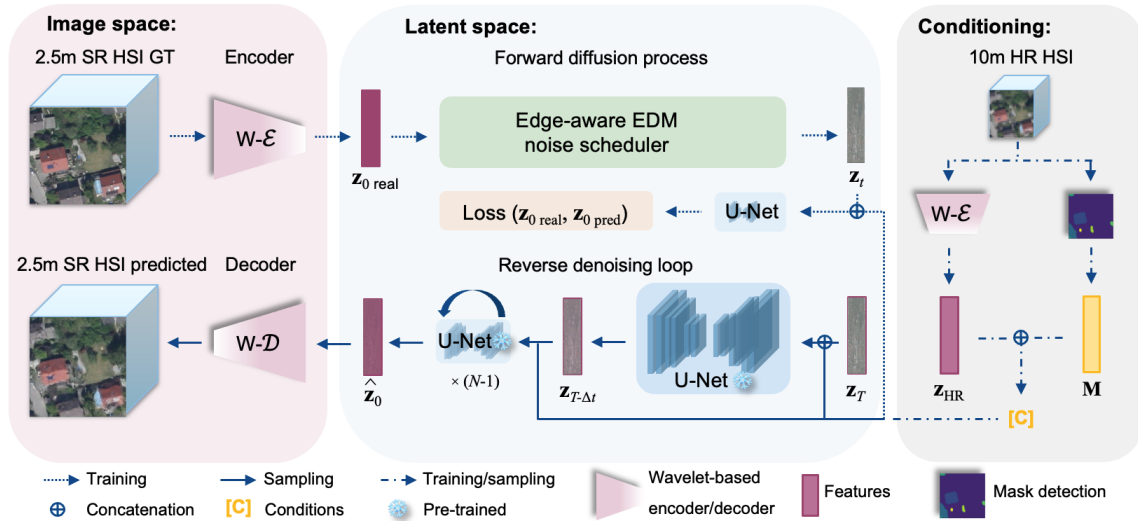


Figure 1: Illustration of the Geometric Enhanced Wavelet-based Diffusion Model pipeline.

alism. Existing architectures often suffer from a slow model convergence speed, excessive sampling steps, and a need for high computational GPU memory, which limits their practicality in real-world scenarios. We propose a Geometric Enhanced Wavelet-based Diffusion Model (GEWDiff) to address these challenges, as shown in Figure 1. The main contributions of this study are summarized as follows:

- **An efficient wavelet-based encoder-decoder** will almost losslessly transform hyperspectral data into a latent space by decomposing the input into multiple frequency levels. This method preserves spectral-spatial information while reducing channel dimensionality without long-term training.
- **Structural invariance control via the diffusion process:** An edge-aware noise scheduler was designed to improve generation efficiency and accuracy during training. Mask conditioning ensures preservation of the geometric integrity and prevents distortion.
- **A multi-level loss function comprising pixel-wise loss, perceptual loss, and gradient loss:** Each part of the loss function contributes to a balanced convergence speed. The loss function also enables the evaluation and alignment of predicted and ground-truth images on specific details and semantic understanding.

Related Work

Properties of wavelet-based diffusion models

Wavelet-based diffusion models have recently gained increasing attention due to their suitability for image generation. For example, WaveDiff (Phung, Dao, and Tran 2023) demonstrated that applying diffusion in the wavelet domain allows images to be compressed into a structured latent space, enabling almost-lossless reconstruction while significantly reducing computational overheads. This approach not only preserves high-fidelity details but also offers substantial memory savings, making it especially advantageous

for large-scale diffusion models. Building on this, Zhao et al. (Zhao et al. 2024) proposed a parallel diffusion strategy that separates high-frequency and low-frequency components for underwater image restoration. Shi et al. (Shi et al. 2024) proposed WaveDiffUR, a wavelet-domain diffusion model for remote sensing ultra-resolution (UR), which involved reformulating high-magnification SR as a conditional stochastic differential equation (SDE) solved via iterative wavelet decomposition, integrating pre-trained SR modules for scalability and a cross-scale pyramid (CSP) constraint to preserve spectral-spatial fidelity. Si et al. (Si et al. 2025) proposed CASSIDiff, the first diffusion model for CASSI hyperspectral reconstruction, which integrated a DWT-based feature fusion mechanism to reduce noise and a spectral-spatial attention module to capture spectral correlations. Despite the success of various wavelet-based diffusion models in natural and remote sensing images, their application to hyperspectral image generation remains unexplored.

Diffusion model for hyperspectral image super-resolution

Recent work has proposed adapting the diffusion process to better fit the characteristics of HSI data. As such, existing approaches can be broadly categorized into three paradigms: two-stage models, grouped autoencoder models, and end-to-end frameworks (Wang et al. 2023). Two-stage models decompose the super-resolution task into two separate subtasks handled by distinct networks. For example, HSI-Gene (Pang et al. 2024b) first generates high-resolution RGB bands from the input HSI, and then fuses them with the low-resolution HSI to reconstruct the final output. This modular design helps reduce computational complexity and allows for flexible training. Grouped models partition the spectral bands into groups and process them in parallel, often using autoencoder-style architectures. DMGASR (Wang et al. 2024), for instance, employs spectral grouping and trains separate VAE-based diffusion modules for different

groups, enabling scalable training across high-dimensional spectra while preserving inter-band correlations. End-to-end frameworks perform full-spectrum reconstruction in a single model, often incorporating various strategies to manage complexity and improve the expressiveness. For example, HIR-Diff (Pang et al. 2024a) integrates a codebook with singular value decomposition to compress and guide the generation process. MTLSC-Diff (Qu et al. 2024) uses classification maps as spatial priors to improve the generation accuracy. LSDiff (Cheng et al. 2024) applies the diffusion process in a compressed latent space, which allows reducing memory usage while maintaining the generation quality. Although some methods have been explored, most rely on two-stage training or fail to simultaneously ensure spectral fidelity and visual quality, while our model addresses both challenges effectively.

Method

Wavelet-based encoder and decoder

Our encoder-decoder is based on Regression wavelet analysis (RWA), first proposed for lossless hyperspectral image compression by (Amrani et al. 2016). RWA applies a predefined number of Haar wavelet (Haar 1910) decompositions intercalated with a linear regression of the spectral dimension to exploit the redundancy that still remains in the discrete wavelet transform (DWT) domain to further compress the data. RWA can compress HSIs with a more efficient, lossless, or near-lossless transform by storing the prediction error. The structure is shown in Figure 10.

Encoder. For the super-resolution task, RWA allows us to reduce the number of bands given to the diffusion model by using the Haar wavelet. Let \mathbf{I}_{LR} be the input 10 m high-resolution hyperspectral image, the J -th level RWA transform can be represented as:

$$(\mathbf{V}_{LR}^J, (\mathbf{w}_{LR}^j)^{1 \leq j \leq J}) = \text{RWA}(\mathbf{I}_{LR}, J), \quad (1)$$

$$\hat{\mathbf{w}}_i^j = \beta_{i,0}^j + \beta_{i,1}^j \mathbf{V}_1^j + \dots + \beta_{i,k}^j \mathbf{V}_k^j, \quad (2)$$

$$\min \|\mathbf{w}_i^j - \hat{\mathbf{w}}_i^j\|_2, \quad (3)$$

where \mathbf{w}_i^j represents the i -th details (high-coefficient) of the j -th level wavelet transform and $\hat{\mathbf{w}}_i^j$ its prediction; \mathbf{V}_k^j represents the k -th low-coefficient (main-coefficient) of the j -th level wavelet transform and $\beta_{i,k}^j$ the linear regression coefficients that will be learned to adjust the linear regression. Contrary to traditional RWA, where the residuals

$$\mathbf{W}_{LR}^j = \mathbf{w}_{LR}^j - \hat{\mathbf{w}}_{LR}^j, \quad (4)$$

are computed in order to fully recover the original signal, the proposed encoder will only store \mathbf{V}_{LR}^J and the weights of all the adjusted linear models $\mathbf{B}_{LR} = [\beta^{1 \leq j \leq J}]$, where J is the level of wavelet transforms applied. The main coefficients \mathbf{V}_{LR}^J , which contain the most critical information, are used as input for the principal component analysis (PCA). The following PCA transformation enables a more efficient compression of hyperspectral imagery (HSI) by achieving

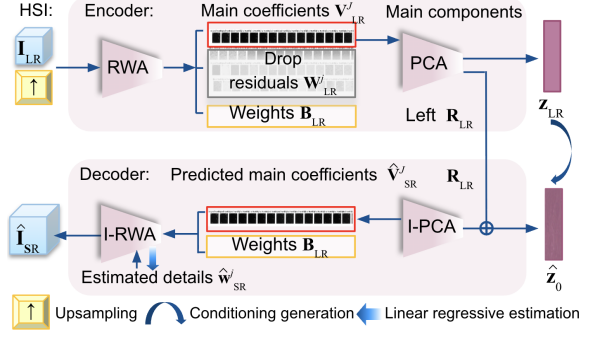


Figure 2: Illustration of the wavelet-based encoder-decoder.

a higher compaction factor while preserving more information. Furthermore, it can convert the sparse wavelet-based coefficients into a dense and orthogonal matrix, facilitating more coherent spectral analysis:

$$(\mathbf{z}_{LR}, \mathbf{R}_{LR}) = \text{PCA}(\mathbf{V}_{LR}^J), \quad (5)$$

where \mathbf{z}_{LR} is the feature that will be input to the diffusion latent space, and \mathbf{R}_{LR} represents the remaining components that will be kept and reused in the decoder equation (6).

Decoder. For the reconstruction, the inverse PCA recovers predicted features $\hat{\mathbf{z}}_0$ from the diffusion process to the super-resolution HSI main coefficients $\hat{\mathbf{V}}_{SR}^J$.

$$(\hat{\mathbf{V}}_{SR}^J) = \text{I-PCA}(\hat{\mathbf{z}}_0, \mathbf{R}_{LR}). \quad (6)$$

The final super-resolved image is obtained by an inverse RWA, having set the residuals \mathbf{W}_{LR}^j to zero, since these are not available for the super-resolution HSI image:

$$(\hat{\mathbf{I}}_{SR}) = \text{I-RWA}(\hat{\mathbf{V}}_{SR}^J, \mathbf{B}_{LR}, \mathbf{W}_{LR}^j, J). \quad (7)$$

The details $\hat{\mathbf{w}}_{SR}^j$ will be predicted by the adjusted linear regression model \mathbf{B}_{LR} to recover the information lost in the wavelet transform. Once the diffusion outputs the super-resolution components $\hat{\mathbf{V}}_{SR}^J$, inverse-RWA reconstructs the predicted main coefficients to the spectral dimension.

Geometric enhanced diffusion process

Hyperspectral image generation usually requires larger reverse sampling time steps. To solve this problem, we used EDM (Elucidating Diffusion Models) (Karras et al. 2022) as our baseline model. EDM adds noise in one step in the training process. A probability flow ordinary differential equation (ODE) then continuously increases the noise level of the image when moving forward in time (Karras et al. 2022). Instead of using a discrete time step to add noise, we used a concrete number σ to represent the strength of the noise:

$$\sigma \sim \exp(\mathcal{N}(P_{\text{mean}} = -1.2, P_{\text{std}} = 1.2)), \quad (8)$$

where P_{mean} represents the mean value, P_{std} is the standard deviation, and \mathcal{N} is a Gaussian distribution.

The “time variable” t used in our model is a continuous variable, which has the advantage of mapping the noise scale

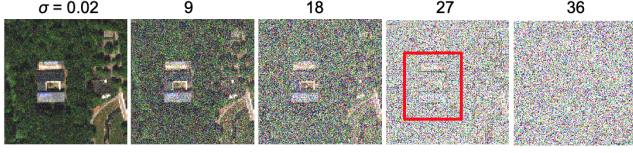


Figure 3: Edge perturbed noisy image over time.

to an approximately linear interval. In this way, the noise strength that will be added at the t moment can correspond to the size of t , and the relationship can be represented as:

$$t = -\log(\sigma_t). \quad (9)$$

Edge-aware noise scheduler. General diffusion models have an equal generation ability for each pixel. In the remote sensing scenario, we wanted to clarify the contour of buildings and other ground objects. Inspired by (Vandersanden et al. 2024), we designed an edge-aware noise scheduler in our training stage to increase the generation ability of the diffusion model for the pixels around the edges. The edge is preserved during the forward diffusion process. The noise around the edge is smaller than the general noise:

$$\mathbf{z}_t = \mathbf{z}_0 + \sigma_t \epsilon \odot (1 - \mathbf{E}(1 - \sigma_{\text{norm}}^2)\eta), \quad (10)$$

where \mathbf{z}_t represents the noisy features at moment t , $t \in (0, T)$, \mathbf{z}_0 is the features from the ground truth, σ_t represents the noise strength at moment t , $\epsilon \sim \mathcal{N}(0, 1)$ is random noise, \mathbf{E} is the binary edge map obtained from the input image, σ_{norm} is the normalized sigma at moment t , \odot highlights that the operation is a matrix multiplication, and $\eta = 0.5$ adjusts the perturbation strength influenced by the edge.

Mask controllable training and sampling. In our preliminary research, generating geometric objects without distortion was a big challenge. To address this, a mask is introduced as a condition that improves the ability to generate buildings. Segmentation is calculated from low-resolution RGB channels from hyperspectral images with a segment-anything model (Kirillov et al. 2023). We used one minus the average of NDVI (Kriegler et al. 1969) as the value of the mask to highlight the attention of buildings. Let S_s be the pixel number of the s th segmentation region, then the value M_s of the mask is defined as:

$$M_s = 1 - \frac{1}{|S_s|} \sum_{(x,y) \in S_s} \text{NDVI}_{\text{norm}}(x, y), \text{NDVI}_{\text{norm}} \in [0, 1]. \quad (11)$$

In the training stage, the predicted result $\hat{\mathbf{z}}_0$ is calculated in one step:

$$\hat{\mathbf{z}}_0 = f_\theta(\mathbf{z}_t, \mathbf{C}, \sigma_t), \mathbf{C} = [\mathbf{z}_{\text{LR}}, \mathbf{M}], \quad (12)$$

where $\hat{\mathbf{z}}_0$ represents the predicted features when $t = 0$; f_θ represents the objective function 3D U-Net with spectral fidelity enhancer (SFE) (Dong et al. 2021), as shown in figure 1; \mathbf{z}_{LR} represents the low-resolution condition; $\mathbf{M} \in (0, 1)^{H \times W}$ is the mask condition; and \mathbf{C} is the concatenation of all conditions.

During the sampling stage, DPM-Solver++ (Lu et al. 2022) accelerates the generation by employing a second-order approximation to solve the underlying ODE, while utilizing adaptive time stepping to significantly reduce the number of function evaluations. In our sampler, $t \in [0, T]$ will be separated into N steps. The step size is $\Delta t = t_{n+1} - t_n, n = 0, 1, \dots, N - 1$. The initial noisy image and noise strength at step n can then be calculated with:

$$\mathbf{z}_T = \sigma_T \cdot \epsilon, \quad (13)$$

$$\sigma_n = \left(\sigma_{\text{max}}^{1/\rho} + \frac{n}{N-1} (\sigma_{\text{min}}^{1/\rho} - \sigma_{\text{max}}^{1/\rho}) \right)^\rho, \quad (14)$$

where ρ is the scheduling curvature parameter, and $\sigma_{\text{max/min}}$ is the maximum/minimum noise strength. The $n+1$ step denoised features \mathbf{z}_{n+1} can be calculated with:

$$\gamma = -\frac{1}{2} \cdot \frac{t_{n+1} - t_n}{t_n - t_{n-1}}, \quad (15)$$

$$\tilde{f}_\theta = (1 - \gamma)f_\theta(\hat{\mathbf{z}}_n, \mathbf{C}, \sigma_n) + \gamma f_\theta(\hat{\mathbf{z}}_{n-1}, \mathbf{C}, \sigma_{n-1}), \quad (16)$$

$$\mathbf{z}_{n+1} = \frac{\sigma_{n+1}}{\sigma_n} \hat{\mathbf{z}}_n - \sigma_{n+1}(e^{-\Delta t} - 1) \cdot \tilde{f}_\theta. \quad (17)$$

Multi-level loss function

Multi-level loss equations, such as equation 18, can ensure that the generated image is accurate in all aspects.

$$\mathcal{L} = \lambda(t) \cdot (\lambda_1 \mathcal{L}_{\text{pixel}} + \lambda_2 \mathcal{L}_{\text{perc}} + \lambda_3 \mathcal{L}_{\text{grad}}), \quad (18)$$

To balance the convergence speed, we set $\lambda_1 = 0.8, \lambda_2 = 0.1, \lambda_3 = 0.1$. $\lambda(t)$ indicates the loss weighting based on t . Considering pixel loss can ensure that the absolute value of the spectral information of each pixel is accurate. Here we use the combination of L2 norm loss and Spectral Angle Mapper (SAM) (Yuhas, Goetz, and Boardman 1992) loss:

$$\mathcal{L}_{\text{pixel}} = (\|\mathbf{z}_0 - \hat{\mathbf{z}}_0\|^2 + \text{SAM}(\mathbf{z}_0, \hat{\mathbf{z}}_0))/2. \quad (19)$$

Perceptron loss (Johnson, Alahi, and Fei-Fei 2016) can ensure that the generated image is similar in high-level feature space:

$$\mathcal{L}_{\text{perc}} = \|\phi \text{VGG}(\hat{\mathbf{z}}_0) - \phi \text{VGG}(\mathbf{z}_0)\|_2^2. \quad (20)$$

Gradient loss (Lu and Chen 2022) can ensure that the image gradient information is consistent with the real image. The images generated by DPM Solver++ have high-contrast characteristics. Gradient loss is defined as follows:

$$\mathcal{L}_{\text{grad}} = \frac{1}{2} (\|\nabla_x \hat{\mathbf{z}}_0 - \nabla_x \mathbf{z}_0\|^1 + \|\nabla_y \hat{\mathbf{z}}_0 - \nabla_y \mathbf{z}_0\|^1), \quad (21)$$

where $\nabla_{x/y}$ represents the gradient of the image in the x/y direction.

	Metric	MCNet	MSDFormer	ESSAFormer	DMGASR	HIR Diff	SNLSR	Ours
(a)	PSNR↑	28.300±0.0480	28.284±0.0000	27.483±0.1392	26.986±0.2052	24.833±0.1079	28.531±0.0001	28.863±0.2940
	SSIM↑	0.6658±0.0025	0.6592±0.0000	0.5915±0.0191	0.5831±0.0118	0.6401±0.0024	0.6718±0.0000	0.7104±0.0212
	SAM↓	8.3332±0.1243	8.7442±0.0000	9.2114±0.2482	11.340±0.0743	8.9538±0.0053	7.8911±0.0000	8.4283±0.3073
	CC↑	0.7440±0.0000	0.7645±0.0000	0.7374±0.0004	0.6767±0.0128	0.7543±0.0011	0.7527±0.0000	0.7945±0.0165
	RMSE↓	0.0557±0.0002	0.0544±0.0000	0.0560±0.0002	0.0627±0.0006	0.0810±0.0015	0.0552±0.0000	0.0548±0.0030
	FID↓	116.14±0.0856	103.74±0.0000	97.438±14.779	49.026±2.3984	50.596±1.0436	125.75±0.0691	44.464±17.627
	LV↑	0.0004±0.0000	0.0004±0.0000	0.0004±0.0000	0.0037±0.0004	0.0021±0.0002	0.0003±0.0000	0.0041±0.0022
(b)	PSNR↑	24.216±0.0157	24.359±0.0000	24.103±0.0132	23.021±0.0146	21.567±0.5351	24.305±0.0000	24.933±0.0079
	SSIM↑	0.5355±0.0008	0.5536±0.0000	0.5210±0.0010	0.4925±0.0025	0.4987±0.0133	0.5404±0.0000	0.6337±0.0106
	SAM↓	11.663±0.0482	11.912±0.0000	12.166±0.0156	16.158±0.0117	12.348±0.1154	11.418±0.0001	11.323±0.0456
	CC↑	0.7050±0.0004	0.7238±0.0000	0.7077±0.0004	0.6575±0.0007	0.7347±0.0003	0.7102±0.0000	0.7771±0.0003
	RMSE↓	0.0685±0.0001	0.0669±0.0000	0.0682±0.0002	0.0779±0.0011	0.0940±0.0061	0.0680±0.0000	0.0668±0.0020
	FID↓	257.45±0.1949	272.78±0.0000	288.85±7.1990	120.06±11.769	375.96±23.877	267.88±0.0032	64.333±4.2810
	LV↑	0.0007±0.0000	0.0006±0.0000	0.0007±0.0000	0.0034±0.0005	0.0005±0.0000	0.0005±0.0000	0.0087±0.0003
(c)	PSNR↑	33.389±0.2641	28.709±0.0000	25.504±0.1340	32.864±0.2049	34.473±0.0069	35.734±0.0000	35.837±0.1176
	SSIM↑	0.7441±0.0029	0.4766±0.0000	0.4120±0.0312	0.6802±0.0159	0.7362±0.0017	0.7525±0.0000	0.7747±0.0045
	SAM↓	8.5500±0.0896	12.213±0.0000	18.724±0.5899	11.476±0.3843	8.3601±0.0446	7.6613±0.0000	7.4735±0.0532
	CC↑	0.6495±0.0145	0.6300±0.0000	0.6326±0.0090	0.5001±0.0161	0.7102±0.0001	0.7733±0.0000	0.7906±0.0055
	RMSE↓	0.0476±0.0006	0.0525±0.0000	0.0690±0.0006	0.0542±0.0013	0.0420±0.0001	0.0471±0.0000	0.0468±0.0006
	FID↓	464.13±5.9021	738.62±0.0000	701.35±16.290	245.38±63.176	363.23±6.9705	470.34±0.0000	238.12±16.970
	LV↑	0.0003±0.0000	0.0003±0.0000	0.0010±0.0000	0.0031±0.0015	0.0002±0.0000	0.0002±0.0000	0.0011±0.0000
(d)	Tr time (s)	1.33×10^4	3.99×10^4	7.65×10^4	3.16×10^5	—	2.80×10^4	3.10×10^5
	Te time (s)	18.13	10.40	7.98	334.00	212.90	4.10	28.70
	NFE	256 ²	256 ²	256 ²	20 × 8	20	256 ²	50
	Model size	6.50 MB	57.7 MB	3.70 MB	1.18 GB	1.56 GB	7.70 MB	4.55 GB

Table 1: Quantitative comparison with SOTA SR models of PSNR, SSIM, SAM, CC, RMSE, FID, and LV on (a) MDAS sample 1, (b) MDAS sample 2, and (c) WDC dataset. (d) Model efficiency was evaluated with the training/testing time, number of function evaluations (NFE), and model size. (Best performance value is highlighted in bold. Noise-affected values are underlined.)



Figure 4: 4-times visual comparisons with SOTA SR models on (a) MDAS sample 1, (b) MDAS sample 2, and (c) WDC dataset.

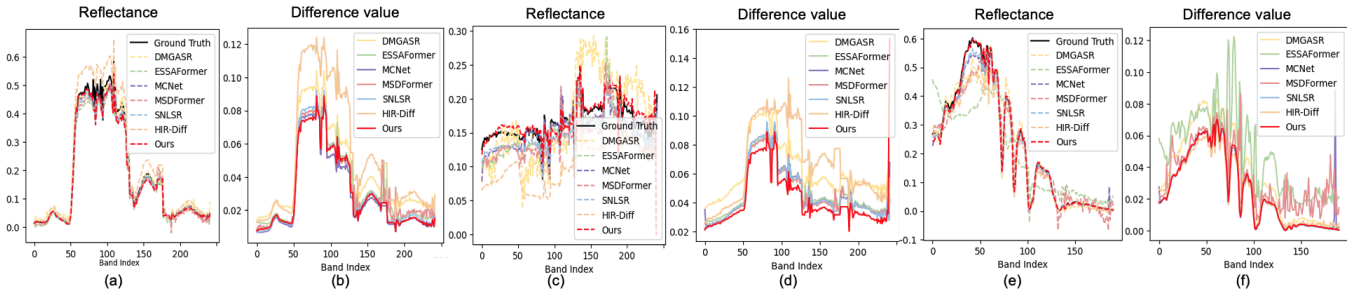


Figure 5: Spectral profile of a random pixel and mean difference value in each band of (a–b) MDAS sample 1, (c–d) MDAS sample 2, and (e–f) WDC dataset.

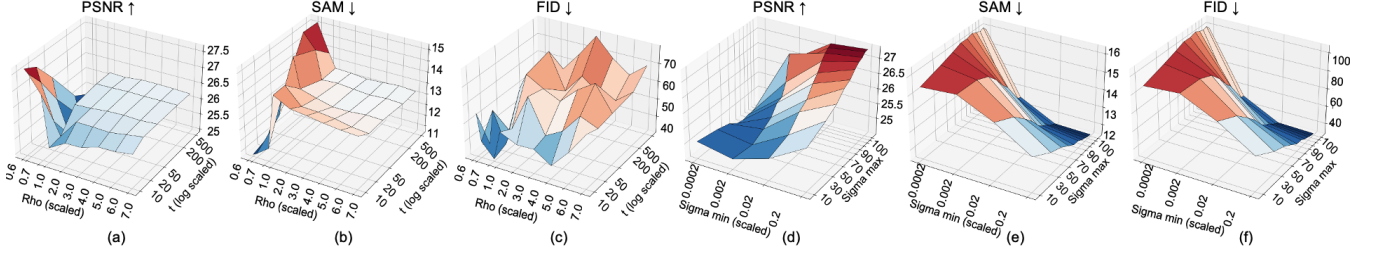


Figure 6: Qualitative comparison with different numbers of ρ and time steps for (a) PSNR, (b) SAM, and (c) FID. Qualitative comparison with different numbers of σ_{\max} and σ_{\min} for (d) PSNR, (e) SAM, and (f) FID on validation dataset.

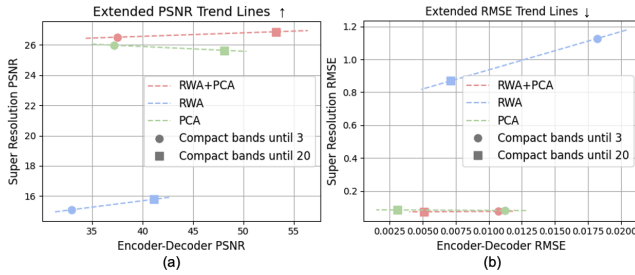


Figure 7: Encoder-decoder reconstruction quality compared with the performance of the whole SR process for the validation dataset.

Experiments

Datasets and implementation details

A group of EeteS simulated EnMap hyperspectral data, called the EnMap Campaign (realistic EnMAP-like high-resolution data), was used for training. EnMap is a hyperspectral satellite managed by the DLR Earth Observation Center, offering 30 m spatial resolution since June 2022. Additionally, the EnMAP Campaign Portal (access via supplementary material) captures aerial hyperspectral imagery and simulates EnMap-like data, boasting a spatial resolution of 2.5–4 m and containing data from 2009 to 2016. We gathered 8000 pairs of $256 \times 256 \times 242$ patches of super-resolution HSI images and aligned them with 4-times downsampled images from EnMap Campaign and MDAS (Hu et al. 2023). This dataset covers 15 cities in Europe and the Americas, representing diverse ground objects. Among 8,000 pairs, one

pair was split into the validation dataset and used for parameter selection, and a group of ablation study experiments. Two pairs are selected for EnMap simulation testing. The WDC (Biehl et al. 2015) dataset is used as one of the test datasets to see the transmission on other datasets. We deployed the model on four NVIDIA A100 GPUs with a learning rate of 1×10^{-4} and trained it for 200 epochs.

Quantitative metrics

Fidelity was evaluated using two standard metrics: the Peak Signal-to-Noise Ratio (PSNR) (Huynh-Thu and Ghanbari 2012), which measures the pixel-wise quality, and the Structural Similarity Index (SSIM) (Wang et al. 2004), which captures structural consistency. Both were averaged over spectral bands. **Realism and clarity** were addressed using the Fréchet Inception Distance (FID) (Heusel et al. 2017), noting that it is computed in an RGB-trained feature space and thus was only used for relative comparisons. We also included Local Variation (LV) (Pertuz, Puig, and Garcia 2013) to evaluate the local texture sharpness. **Spectral accuracy** was measured via Spectral Angle Mapper (SAM) (Yuhas, Goetz, and Boardman 1992), Cross-Correlation (CC), and Root Mean Square Error (RMSE). These metrics assess the angular, correlational, and pixel-wise spectral alignment.

State-of-the-art image generation

We evaluated the performance and efficiency of our model by comparing it with the six SOTA models: MCNet (Li, Wang, and Li 2020), MSDFormer (Chen, Zhang, and Zhang 2023), ESSAFormer (Zhang et al. 2023), DMGASR (Wang et al. 2024), HIR Diff (Pang et al. 2024a), and SNLSR (Hu et al. 2024). We trained these models on our dataset with the

Method	Baseline	A	B	C	D	E	F	G	H	I	J	Ours
w/RWA		✓		✓	✓	✓	✓	✓	✓	✓	✓	✓
w/PCA			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
w/Mask		✓	✓	✓		✓	✓	✓	✓	✓		✓
w/Edge		✓	✓		✓		✓		✓	✓	Inverse	✓
w/L pix	✓	✓	✓	✓		✓			✓	✓	✓	✓
w/L perc		✓	✓	✓	✓		✓		✓	✓	✓	✓
w/L geo		✓	✓	✓	✓			✓	✓	✓	✓	✓
w/Unet3D		✓	✓	✓	✓	✓	✓	✓	✓		✓	✓
w/SFE		✓	✓	✓	✓	✓	✓	✓		✓	✓	✓
PSNR↑	2.0476	15.788	25.640	26.579	26.681	26.101	21.664	26.467	26.388	26.181	26.8713	27.013
SSIM↑	-0.0150	-0.0173	0.5267	0.6503	0.6530	0.6140	0.1580	0.6443	0.6240	0.6131	0.6667	0.6573
SAM↓	124.15	85.239	15.125	11.766	12.160	12.804	25.528	12.445	12.788	14.085	11.7688	11.501
CC↑	0.2643	0.4763	0.5530	0.6803	0.6882	0.6441	0.0145	0.6741	0.6681	0.6500	0.6999	0.7008
RMSE↓	1.1879	0.8687	0.0850	0.0758	0.0749	0.0804	0.1342	0.0769	0.0783	0.0803	0.0739	0.0726
FID↓	5019.1	484.16	83.627	43.445	36.269	40.303	701.94	40.195	47.475	65.985	34.9402	30.110
LV↑	7.5344	0.6231	0.0160	0.0079	0.0077	0.0096	0.1684	0.0093	0.0074	0.0089	0.0080	0.0083

Table 2: Ablation study. Quantitative comparison for the effect of each module on the validation dataset. Results are averaged over 4 runs. The baseline used an EDM backbone and DPM-Solver++ sampler on 242 bands. A: no PCA in the encoder; B: no RWA in the encoder; C: no edge perturbation; D: no mask conditioning; E/F/G: retained pixel/perceptual/geometric loss; H: no spectral fidelity module in the encoder; I: used 2D U-Net instead of 3D; J: more noise on edge (Zhang et al. 2025).

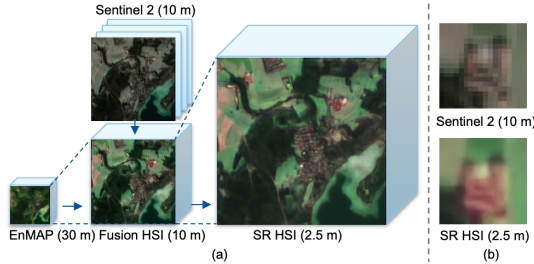


Figure 8: (a) Visualization of EnMAP, Sentinel 2, Fusion image, Super resolution image. (b) Effect of super-resolution.

implementation details provided in each paper, except for the unsupervised method HIR Diff. The HIR Diff model provides a pre-trained checkpoint based on their fully prepared HSI dataset. Our model demonstrated strong performance in generating medium to large-scale ground objects. In Figure 4 and Table 1, we can see that both the visualization result and quantitative result outperformed the other models. The WDC dataset result shows our model can adapt to another dataset with a different spectral profile. However, our model may struggle with accuracy when the input conditional image lacks sufficient semantic information; for example, the rooftop of the MDAS sample 1 reconstructed image.

Effect of the encoder-decoder alone vs. with super-resolution. We compared our encoder-decoder (RWA + PCA) against RWA-only and PCA-only baselines (Figure 7), under 3- and 20-band compression. Without super-resolution, our encoder-decoder achieved almost lossless reconstruction with a PSNR up to 56. With super-resolution, using RWA + PCA compression to 20 bands outperformed all the baselines, demonstrating strong spectral compaction and generation. As shown, retaining more bands improves the quality but increases the cost. We chose 20 bands to ensure a balance between performance and efficiency.

Effect of ρ , number of timesteps N , σ_{\max} and σ_{\min} on sampling. The EDM noise schedule is shaped by ρ , which controls how noise levels decay from σ_{\max} to σ_{\min} via N steps. Higher ρ sharpens early denoising but increases the randomness in later steps, potentially degrading spectral consistency. We found that lower ρ values [0.6–0.7] yielded smoother noise schedules via 50 steps, which also better preserve spectral fidelity (Figure 6). We also studied the effects of σ_{\max} and σ_{\min} . A larger σ_{\max} allows more diversity but risks over-noising; a smaller one limits the detail. We set $\sigma_{\max} = 80$ for a balance. For σ_{\min} , smaller values extend the denoising phase, improving the detail but increasing artifacts. We found the best results when $\sigma_{\min} \in [0.02, 0.2]$.

Real-world application. We tried to combine EnMAP and Sentinel-2 imagery to 10 m resolution hyperspectral data via the unsupervised method HySure (Simões et al. 2015). Leveraging our 4-times super-resolution generation model, GEWDiff, we could finally produce EnMAP hyperspectral images with a 2.5-meter resolution. The no-reference image quality assessment MetaIQA (Zhu et al. 2020) was improved from 0.1997 to 0.2029 (Figure 8 (b)).

Ablation study

The results of our ablation studies, presented in Table 2, offer significant insights into the contributions of various components of the GEWDiff model. The use of a suitable encoder plays a crucial role in our model design. The multi-level loss function achieves better performance than an L2 loss. The design of a 3D objective function, U-Net, and its spectral fidelity enhancer makes progress toward stabilizing the results. The geometric enhancement strategies, such as edge perturbation and mask conditioning, did not show a significant improvement in the global metrics. However, some effects could be observed from Figure 4, whereby the edges are clear, and there was no obvious building distortion.

Conclusion

We proposed GEWDiff, which improves the spatial resolution of hyperspectral images by a factor of 4. Our method integrates wavelet-domain transforms and geometric priors to effectively preserve both spectral fidelity and spatial textures while accelerating convergence. The experimental results showed that GEWDiff outperformed the current SOTA baselines. One limitation of GEWDiff is that the result relies too much on the input conditions. This limitation could be addressed in future work through the integration of classifier-free guidance, which would enable the model to better generalize under weak or ambiguous conditioning. Moreover, we would also contribute to model distillation for further lightweight alternatives.

Acknowledgments

This manuscript has been accepted for publication in AAAI 2026. Please refer to the original version via [link: <https://aaai.org/conference/aaai/aaai-26/>].

This work was supported in part by [the place to write support project]; and in part by Munich Center for Machine Learning. The authors sincerely thank GFZ Helmholtz-Zentrum for providing the EnMAP Campaign datasets (Buddenbaum and Hill 2020; Beamish et al. 2020; Brell et al. 2020; Milewski et al. 2020; Cooper et al. 2020; Hank et al. 2015; Jarmer and Siegmann 2017; Neumann, Weiss, and Itzerott 2015; Okujeni, van der Linden, and Hostert 2016; Foerster et al. 2015; Boesche et al. 2016) used in this paper, and Dr. Jingliang Hu for providing the MDAS dataset.

References

- Amrani, N.; Serra-Sagristà, J.; Laparra, V.; Marcellin, M. W.; and Malo, J. 2016. Regression Wavelet Analysis for Lossless Coding of Remote-Sensing Data. *IEEE Transactions on Geoscience and Remote Sensing*, 54(9): 5616–5627.
- Beamish, A.; Chabrillat, S.; Brell, M.; Heim, B.; and Sachs, T. 2020. Toolik Lake Research Natural Area AISA-Eagle hyperspectral Mosaic.
- Biehl, L.; Maud, A. R.; Hsu, W.-K.; and Yeh, T. T. 2015. MultiSpec.
- Boesche, N. K.; Mielke, C.; Segl, K.; Chabrillat, S.; Rogass, C.; Thomson, D.; Lundeen, S.; Brell, M.; and Guanter, L. 2016. EnGeoMAP Test Data: Simulated EnMAP Satellite Data for Mountain Pass, USA and Rodalquilar, Spain.
- Brell, M.; Spengler, D.; Ruhtz, T.; Ward, K.; Chabrillat, S.; Segl, K.; Foerster, S.; and Itzerott, S. 2020. Demmin, Germany (October 2015) - an EnMAP Preparatory Flight Campaign.
- Buddenbaum, H.; Dotzler, S.; and Hill, J. 2015a. Donnersberg, 2014-07-03 - An EnMAP Preparatory Flight Campaign (Datasets).
- Buddenbaum, H.; Dotzler, S.; and Hill, J. 2015b. Nationalpark Hunsrück-Hochwald, 2014-05-05 - An EnMAP Preparatory Flight Campaign (Datasets).
- Buddenbaum, H.; and Hill, J. 2020. Gerolstein, 2016-09-08 - An EnMAP Preparatory Flight Campaign.
- Chen, N.; Yue, J.; Fang, L.; and Xia, S. 2023. SpectralDiff: A generative framework for hyperspectral image classification with diffusion models. *IEEE Transactions on Geoscience and Remote Sensing*, 61: 1–16.
- Chen, S.; Zhang, L.; and Zhang, L. 2023. MSDformer: Multiscale deformable transformer for hyperspectral image super-resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 61: 1–14.
- Cheng, Y.; Ma, Y.; Fan, F.; Ma, J.; Yao, Y.; and Mei, X. 2024. Latent spectral-spatial diffusion model for single hyperspectral super-resolution. *Geo-spatial Information Science*, 1–16.
- Cooper, S.; Okujeni, A.; Jänicke, C.; Segl, K.; van der Linden, S.; and Hostert, P. 2020. 2013 Simulated EnMAP Mosaics for the San Francisco Bay Area, USA.
- Dong, W.; Hou, S.; Xiao, S.; Qu, J.; Du, Q.; and Li, Y. 2021. Generative dual-adversarial network with spectral fidelity and spatial enhancement for hyperspectral pansharpening. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12): 7303–7317.
- Foerster, S.; Brosinsky, A.; Wilczok, C.; and Bauer, M. 2015. Isábena 2011 - An EnMAP Preparatory Flight Campaign (Datasets).
- Haar, A. 1910. Zur Theorie der orthogonalen Funktionensysteme. *Mathematische Annalen*, 69: 331–371.
- Hank, T. B.; Locherer, M.; Richter, K.; and Mauser, W. 2015. Neusling (Landau a.d. Isar) 2012 - A Multitemporal and Multisensoral Agricultural EnMAP Preparatory Flight Campaign (Datasets).
- Hank, T. B.; Richter, K.; and Mauser, W. 2015. Neusling (Landau a.d. Isar) 2009 - An Agricultural EnMAP Preparatory Flight Campaign Using the HyMap Instrument (Datasets).
- Heusel, M.; Ramsauer, H.; Unterthiner, T.; Nessler, B.; and Hochreiter, S. 2017. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising Diffusion Probabilistic Models. arXiv:2006.11239.
- Hou, J.; Zhu, Z.; Hou, J.; Zeng, H.; Wu, J.; and Zhou, J. 2022. Deep posterior distribution-based embedding for hyperspectral image super-resolution. *IEEE Transactions on Image Processing*, 31: 5720–5732.
- Hu, J.; Liu, R.; Hong, D.; Camero, A.; Yao, J.; Schneider, M.; Kurz, F.; Segl, K.; and Zhu, X. X. 2023. MDAS: a new multimodal benchmark dataset for remote sensing. *Earth System Science Data*, 15(1): 113–131.
- Hu, Q.; Wang, X.; Jiang, J.; Zhang, X.-P.; and Ma, J. 2024. Exploring the Spectral Prior for Hyperspectral Image Super-Resolution. *IEEE Transactions on Image Processing*, 33: 5260–5272.
- Huynh-Thu, Q.; and Ghanbari, M. 2012. The accuracy of PSNR in predicting video quality for different video scenes and frame rates. *Telecommunication Systems*, 49(1): 13–27.

- Jarmer, T.; and Siegmann, B. 2017. Köthen 2011/2012 - An EnMAP Preparatory Flight Campaign.
- Johnson, J.; Alahi, A.; and Fei-Fei, L. 2016. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. *arXiv:1603.08155*.
- Karras, T.; Aittala, M.; Aila, T.; and Laine, S. 2022. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35: 26565–26577.
- Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; Dollár, P.; and Girshick, R. 2023. Segment Anything. *arXiv:2304.02643*.
- Kriegler, F. J.; Malila, W. A.; Nalepka, R. F.; and Richardson, W. 1969. Preprocessing transformations and their effects on multispectral recognition. In *Proceedings of the Sixth International Symposium on Remote Sensing of Environment*, 97–131.
- Li, J.; Cui, R.; Li, B.; Song, R.; Li, Y.; Dai, Y.; and Du, Q. 2020. Hyperspectral Image Super-Resolution by Band Attention Through Adversarial Learning. *IEEE Transactions on Geoscience and Remote Sensing*, 58(6): 4304–4318.
- Li, Q.; Wang, Q.; and Li, X. 2020. Mixed 2D/3D convolutional network for hyperspectral image super-resolution. *Remote sensing*, 12(10): 1660.
- Lu, C.; Zhou, Y.; Bao, F.; Chen, J.; Li, C.; and Zhu, J. 2022. DPM-Solver++: Fast solver for guided sampling of diffusion probabilistic models. *ArXiv preprint, arXiv:2211.01095*.
- Lu, Z.; and Chen, Y. 2022. Single image super-resolution based on a modified U-net with mixed gradient loss. *signal, image and video processing*, 16(5): 1143–1151.
- Milewski, R.; Chabrilat, S.; Brell, M.; Behling, R.; and Eichstaedt, H. 2020. Omongwa Pan, Namibia (June 2015) - an EnMAP Preparatory Flight Campaign.
- Neumann, C.; Weiss, G.; and Itzerott, S. 2015. Döberitzer Heide 2008/2009 - An EnMAP Preparatory Flight Campaign (Datasets).
- Okujeni, A.; van der Linden, S.; and Hostert, P. 2016. Berlin-Urban-Gradient dataset 2009 - An EnMAP Preparatory Flight Campaign (Datasets).
- Pang, L.; Rui, X.; Cui, L.; Wang, H.; Meng, D.; and Cao, X. 2024a. HIR-Diff: Unsupervised hyperspectral image restoration via improved diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3005–3014.
- Pang, L.; Tang, D.; Xu, S.; Meng, D.; and Cao, X. 2024b. HSiGene: A Foundation Model For Hyperspectral Image Generation. *arXiv:2409.12470*.
- Peebles, W.; and Xie, S. 2023. Scalable Diffusion Models with Transformers. *arXiv:2212.09748*.
- Pertuz, S.; Puig, D.; and Garcia, M. A. 2013. Analysis of focus measure operators for shape-from-focus. *Pattern Recognition*, 46(5): 1415–1432.
- Phung, H.; Dao, Q.; and Tran, A. 2023. Wavelet diffusion models are fast and scalable image generators. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10199–10208.
- Qu, J.; Xiao, L.; Dong, W.; and Li, Y. 2024. MTLSC-Diff: Multitask learning with diffusion models for hyperspectral image super-resolution and classification. *Knowledge-Based Systems*, 303: 112415.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. *arXiv:2112.10752*.
- Segl, K.; Küster, T.; Rogaß, C.; Kaufmann, H.; Sang, B.; Mogulsky, V.; and Hofer, S. 2012. EeteS: An end-to-end image simulation tool applied to THE EnMAP hyperspectral mission. In *2012 IEEE International Geoscience and Remote Sensing Symposium*, 5025–5028.
- Shi, Y.; Han, L.; Dancy, D.; and Han, L. 2024. WaveDiffUR: A diffusion SDE-based solver for ultra magnification super-resolution in remote sensing images. *arXiv:2412.18996*.
- Shi, Y.; Han, L.; Han, L.; Chang, S.; Hu, T.; and Dancey, D. 2022. A Latent Encoder Coupled Generative Adversarial Network (LE-GAN) for Efficient Hyperspectral Image Super-Resolution. *IEEE Transactions on Geoscience and Remote Sensing*, 60: 1–19.
- Si, Y.; Lin, Z.; Wang, X.; and He, S. 2025. A New Hyperspectral Reconstruction Method With Conditional Diffusion Model for Snapshot Spectral Compressive Imaging. *IEEE Transactions on Instrumentation and Measurement*, 74: 1–14.
- Sidorov, O.; and Yngve Hardeberg, J. 2019. Deep hyperspectral prior: Single-image denoising, inpainting, super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 0–0.
- Simões, M.; Bioucas-Dias, J.; Almeida, L. B.; and Chanussot, J. 2015. A Convex Formulation for Hyperspectral Image Superresolution via Subspace-Based Regularization. *IEEE Transactions on Geoscience and Remote Sensing*, 53(6): 3373–3388.
- Su, X.; Shen, X.; Wan, M.; Nie, J.; Chen, L.; Liu, H.; and Zhou, X. 2025. EigenSR: Eigenimage-Bridged Pre-Trained RGB Learners for Single Hyperspectral Image Super-Resolution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 7033–7041.
- Vandersanden, J.; Holl, S.; Huang, X.; and Singh, G. 2024. Edge-preserving noise for diffusion models. *ArXiv preprint, arXiv:2410.01540*.
- Wang, X.; Hu, Q.; Cheng, Y.; and Ma, J. 2023. Hyperspectral image super-resolution meets deep learning: A survey and perspective. *IEEE/CAA Journal of Automatica Sinica*, 10(8): 1668–1691.
- Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4): 600–612.
- Wang, Z.; Li, D.; Zhang, M.; Luo, H.; and Gong, M. 2024. Enhancing hyperspectral images via diffusion model and group-autoencoder super-resolution network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 5794–5804.

Wu, C.; Wang, D.; Bai, Y.; Mao, H.; Li, Y.; and Shen, Q. 2023. HSR-Diff: Hyperspectral image super-resolution via conditional diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7083–7093.

Yu, D.; Li, Q.; Wang, X.; Zhang, Z.; Qian, Y.; and Xu, C. 2023. Dstrans: Dual-stream transformer for hyperspectral image restoration. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 3739–3749.

Yuhas, R. H.; Goetz, A. F. H.; and Boardman, J. W. 1992. Discrimination among semi-arid landscape endmembers using the Spectral Angle Mapper (SAM) algorithm. In *Summaries of the Third Annual JPL Airborne Geoscience Workshop*, volume 1, 147–149. Pasadena, CA: JPL.

Zhang, L.; You, W.; Shi, K.; and Gu, S. 2025. Uncertainty-guided Perturbation for Image Super-Resolution Diffusion Model. arXiv:2503.18512.

Zhang, M.; Zhang, C.; Zhang, Q.; Guo, J.; Gao, X.; and Zhang, J. 2023. Essaformer: Efficient transformer for hyperspectral image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 23073–23084.

Zhao, C.; Cai, W.; Dong, C.; and Hu, C. 2024. Wavelet-based Fourier Information Interaction with Frequency Diffusion Adjustment for Underwater Image Restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 8281–8291.

Zhu, H.; Li, L.; Wu, J.; Dong, W.; and Shi, G. 2020. MetaIQA: Deep Meta-learning for No-Reference Image Quality Assessment. arXiv:2004.05508.

Appendix

There is a growing demand for high-quality, global-scale hyperspectral data. As shown in Figure 9, we compare the global coverage of EnMap data, Sentinel-2 data, and high-resolution airborne hyperspectral imagery. It is evident that higher data quality typically comes at the expense of reduced spatial coverage. In this work, we aim to develop a cost-efficient AI-based model that enhances the quality of hyperspectral satellite imagery (e.g., EnMap), thereby making high-quality data more accessible for downstream applications such as environmental monitoring, land cover classification, and precision agriculture.

Research review for HSI super-resolution.

A Mixed Convolutional Network (MCNet) (Li, Wang, and Li 2020) is proposed for hyperspectral image super-resolution (SR), which introduces a Mixed Convolutional Module (MCM) with 3D CNNs. A novel deep learning method named PDE-Net is proposed (Hou et al. 2022) to formulate hyperspectral (HS) image embedding as an approximation of the posterior distribution of spatial-spectral features. Generative adversarial networks (GANs) have gained attention as a promising class of generative models for hyperspectral image super-resolution. A Latent-Encoder GAN (LE-GAN) (Shi et al. 2022) introduces a latent encoder to

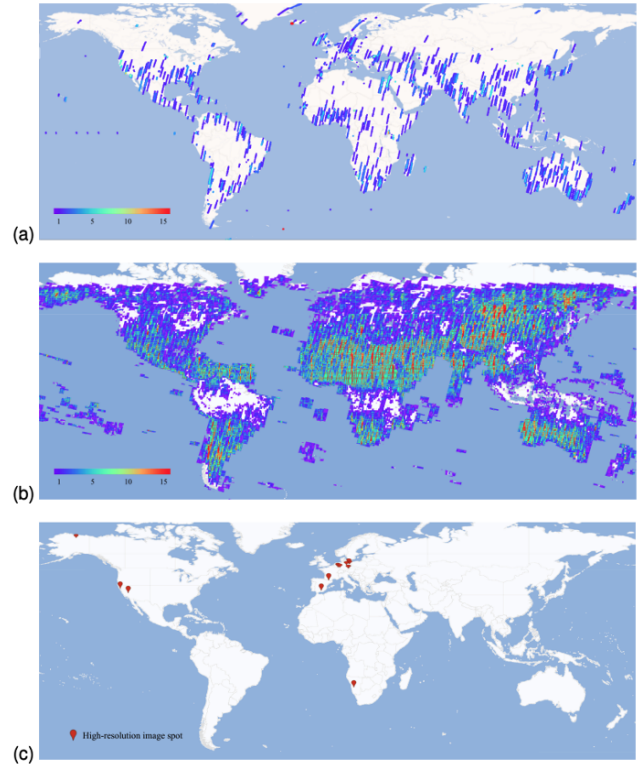


Figure 9: Satellite trajectory and number of repeat times for (a) EnMap ($GSD = 30m$) and (b) Sentinel-2 ($GSD = 10m$) satellite from 01.2024 to 03.2024. (Sentinel-2 images are filtered by cloud coverage $< 10\%$) (c) Location of high-resolution EnMap Campaign ($GSD \approx 2.5m$) hyperspectral airborne images.

map spectral-spatial features into a latent space. A GAN-based method with band attention (Li et al. 2020) leverages full-band input and spectral-spatial constraints to reduce texture blur and spectral distortion. However, they often suffer from limitations such as mode collapse and low sample diversity, which can limit their ability to fully capture the complex distribution of hyperspectral data. Notable examples include ESSAformer (Zhang et al. 2023), which introduces enhanced spectral-spatial attention mechanisms; MSDFormer (Chen, Zhang, and Zhang 2023) leverages multi-scale dynamic fusion; DStrans (Yu et al. 2023), which incorporates dynamic spectral transformers; and EigenSR (Su et al. 2025), which explores low-rank representations within a transformer framework. Transformers show architectures’ growing versatility and effectiveness in modeling complex hyperspectral data. While these models effectively maintain spectral consistency across bands, they often fail to generate fine-grained spatial details and high-frequency textures that are critical for accurate hyperspectral reconstruction.

Challenges in hyperspectral image generation.

Despite recent progress, hyperspectral image (HSI) generation still faces several practical and technical challenges, we explained the details for challenges and the solution we give

Name	Application	Sensor	Product-Level	Date
Gerobstein (Germany)	Forest	HySpec	L2	2016-09-08
Tonlik Lake (Alaska, USA)	Vegetation	AISA (Eagle)	L2	2016-08-27
Drenmin (Germany)	Soil	HySpex VNIR-1600	L2	2015-10-01
Namäla	Soil	HySpex SWIR320m-e	L2	2015-06-06
		HySpex VNIR-1600		
		HySpex SWIR320m-e		
Donnersberg (DE)	Forest	HySpex VNIR-1600	L2	2014-07-02
		HySpex SWIR320m-e		
		HySpex VNIR-1600		
Mountain Pass, USA	Geology	AVIRIS-NG	L2	2014-06-21
Redalquillar, Spain	Geology	HyMap	L2	2014-06-15
Hamrdick-Hochwald (DE)	Forest	HySpex VNIR-1600	L2	2014-05-04
		HySpex SWIR320m-e		
San Francisco Bay Area (USA)	Vegetation	AVIRIS-classic	L2	2013-11-22
San Francisco Bay Area (USA)	Vegetation	AVIRIS-classic	L2	2013-06-07
San Francisco Bay Area (USA)	Vegetation	AVIRIS-classic	L2	2013-04-10
Neurling (DE)	Agriculture	AVIS-3	L2	2012-09-07
Neurling (DE)	Agriculture	HySpex VNIR-1600	L2	2012-08-11
		HySpex SWIR320m-e		
MDAS (DE)	Urban/Urban-Rural	HySpec	L2	2018-5-7
Neurling (DE)	Agriculture	AVIS-3	L2	2012-06-15
Neurling (DE)	Agriculture	AVIS-3	L2	2012-05-24
Köthen (DE)	Agriculture	asiaDual	L2	2012-05-23
Neurling (DE)	Agriculture	HySpex VNIR-1600	L2	2012-05-07
		HySpex SWIR320m-e		
Neurling (DE)	Agriculture	AVIS-3	L2	2012-04-27
Neurling (DE)	Agriculture	APEX	L2	2011-09-09
Isabean (ES)	Environmental gradients	asiaEagle	L2	2011-08-08
Köthen (DE)	Agriculture	asiaDual	L2	2011-06-26
Köthen (DE)	Agriculture	asiaDual	L2	2011-05-09
Isabean (ES)	Environmental gradients	asiaEagle	L2	2011-04-01
Döbertizer Heide (DE)	Vegetation	HyMap	L2	2009-08-19
		HyMap		
Berlin (DE)	Urban/Urban-Rural	HyMap	L2	2009-08-19
Berlin (DE)	Urban/Urban-Rural	HyMap	L2	2009-08-19
Neurling (DE)	Agriculture	HyMap	L2	2009-07-26
Steinbeifen (DE)	Agriculture	HyMap	L2	2009-07-26
Döbertizer Heide (DE)	Vegetation	HyMap	L2	2008-08-06

Table 3: Data Collection Summary

in our paper:

- **Computational cost.** Diffusion models such as DDPM (Ho, Jain, and Abbeel 2020) are memory-intensive, requiring up to 1.3 GB of GPU memory per spectral channel, making them hard to scale to high-bandwidth data. In contrast, our proposed method with 3D U-Net reduces this requirement to a maximum of 1.0 GB per channel. Additionally, we can provide the code for a lightweight version of the model with a 2D U-Net, which operates in a 3-channel latent space. This configuration supports a batch size of 16 without significant performance degradation and is fully trainable on a single NVIDIA RTX 3090 GPU.
- **Sampling speed.** Generating HSIs typically requires thousands of timesteps (e.g., 5000 for DDPM), which significantly slows down inference. Our proposed sampling scheduler reduces to 50 timesteps, which can be tested on CPU devices.
- **Low signal-to-noise ratio.** Maintaining strong signal integrity across many channels is difficult, especially under noise-aware training objectives. Hyperspectral image generation is particularly challenging due to the inherently low signal-to-noise ratio. To address this, we employ an encoder-decoder architecture that transforms the

complex high-dimensional image space into a more compact and structured latent space.

- **Slow convergence.** Diffusion models tend to converge very slowly when processing a large number of channels simultaneously. To alleviate this issue, our encoder-decoder architecture effectively reduces the number of channels to 20 in the latent space.
- **Geometric distortion.** Preserving geometric correctness across spectral bands is critical but challenging, especially when using spatially-aware denoising processes. Conventional natural image generation models often struggle to capture the unique geometric and structural patterns required for remote sensing image synthesis. To address this, we propose a mask-controllable, edge-aware architecture that enhances geometric fidelity and reduces spatial distortion in the generated hyperspectral images.
- **Real-world applicability.** For satellite-based HSI (e.g., EnMap), a two-stage resolution enhancement is necessary. First, spatial resolution is increased to Sentinel-2 level (10 m) using a fusion model trained on EeteS [2] simulated data. Then, a diffusion model further upsamples it to 2.5 m, enabling high-resolution, global-scale hyperspectral generation.

Dataset

Hyperspectral aerial images that have been used in the framework are called EnMAP Campaign (Buddenbaum and Hill 2020; Beamish et al. 2020; Brell et al. 2020; Milewski et al. 2020; Buddenbaum, Dotzler, and Hill 2015a,b; Cooper et al. 2020; Hank et al. 2015; Jarmer and Siegmann 2017; Neumann, Weiss, and Itzerott 2015; Hank, Richter, and Mauser 2015; Okujeni, van der Linden, and Hostert 2016; Foerster et al. 2015; Boesche et al. 2016). It is a preparatory science program to support method and application development in the prelaunch phase of the EnMAP satellite mission. This project aims to use aerial HSI data to simulate EnMap images. We also used another EnMap-like aerial dataset called the MDAS dataset (Hu et al. 2023). All datasets are shown in Table. 1.

We use Level-2 reflectance products with digital number (DN) values ranging from 0 to 10,000, and a scale factor of 10,000. Each sample consists of a low-resolution input patch of size $64 \times 64 \times 242$ and a corresponding high-resolution ground truth patch of $256 \times 256 \times 242$, resulting in 8000 paired samples for training and evaluation. We use the MDAS dataset as the reference spectral profile, as it has been preprocessed using the EeteS (Segl et al. 2012) tool to simulate EnMap spectral characteristics. To ensure consistency across datasets from different sensors, we apply nearest-neighbor interpolation to align their spectral profiles with that of MDAS.

Methodology

In this section, we include two paragraphs that were not shown in our paper for a better understanding.

Information loss of wavelet-based encoder-decoder.

Our encoder-decoder is designed to compact hyperspectral images into a latent space by operating solely on the spectral dimension. However, minimizing information loss during this spectral compression is critical for accurate HSI reconstruction. In this section, we evaluate the effectiveness of our approach in achieving an almost lossless transformation, demonstrating its ability to preserve essential spectral information.

RWA significantly improves the redundancy processing capability of DWT by introducing a regression model. During the inverse transformation, the detail coefficients are predicted from the main coefficients through the regression model, and the original image is reconstructed using the saved weights. The results of the wavelet transform are usually sparse, even though they compact the HSI efficiently. In contrast, the results of PCA are usually dense, and the principal components are orthogonal. PCA is more suitable for processing global features in the diffusion model. As a result, we combine two transformers.

In Figure 10, we present reconstructed images using RWA+PCA, RWA-only, and PCA-only methods in line (a). The compression to 3–4 bands is intentional to better visualize reconstruction errors. The RWA-only method struggles to recover full color information, while the PCA-only approach produces overly smooth, blurry results. In contrast,

RWA+PCA yields clearer latent features (line (b)) and more faithful reconstructions. In Figure 11, all methods demonstrate strong recovery of spectral profiles, even when compressed to just 3–4 bands. In practice, we compact to 20 bands to ensure an almost lossless spectral transformation.

Table. 2 presents a quantitative comparison of various RWA levels combined with different numbers of retained PCA components. The first row shows the final configuration used in our experiments, chosen to balance GPU memory usage and encoder-decoder efficiency.

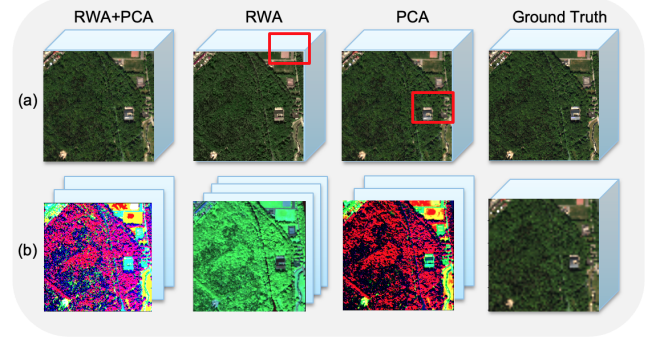


Figure 10: (a) Recovered encoder → decoder (RGB visualization), (b) Latent space (False color visualization for first three bands).

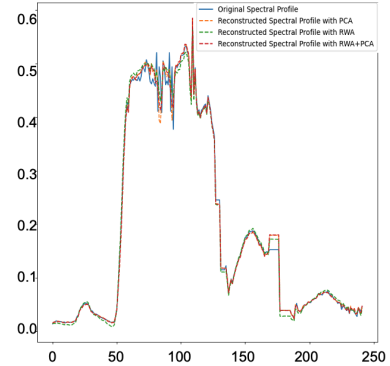


Figure 11: Spectral profile of pixel (125,125) recovered encoder → decoder original vs. reconstructed by RWA, PCA and RWA+PCA.

Structure of U-Net in GEWDiff framework.

In previous experiments, we found that the simple model did not have enough ability to preserve spectral information. 3D U-Net performs convolutions in both spatial and spectral dimensions. Thus, it can capture joint spatial-spectral features from PCA components in latent space, which is critical for accurately modeling hyperspectral images.

Additionally, we introduce a Spectral Fidelity Enhancer (SFE) module into the diffusion model to further preserve the spectral structure of hyperspectral images. The SFE can be intuitively understood as a spectral attention mechanism

RWA compact until	PCA compact until	MPSNR \uparrow	MSSIM \uparrow	ERGAS \downarrow	SAM \downarrow	CrossCorrelation \uparrow	RMSE \downarrow
121 bands \checkmark	20 bands \checkmark	53.26	0.9962	0.9505	1.7577	0.9982	0.0051
31 bands	20 bands	52.14	0.995	1.1208	2.3755	0.9974	0.0068
61 bands	20 bands	52.77	0.9956	0.9889	2.1529	0.9979	0.0062
61 bands	10 bands	49.81	0.9939	1.3837	2.3702	0.9967	0.0068
16 bands	10 bands	48.76	0.9933	1.4816	2.5506	0.9961	0.0073
121 bands	10 bands	49.42	0.9942	1.5159	2.0848	0.9965	0.006
31 bands	10 bands	50.01	0.9938	1.3397	2.4749	0.9967	0.0071
242 bands	10 bands	49.14	0.994	1.6179	1.8881	0.9962	0.0055
16 bands	6 bands	46.47	0.9898	1.8661	2.8588	0.9936	0.0082
31 bands	6 bands	47.27	0.9912	1.8545	2.799	0.9945	0.0081
61 bands	6 bands	47.31	0.9912	1.8261	2.8198	0.9946	0.0081
121 bands	6 bands	47.24	0.991	1.8519	2.8949	0.9944	0.0084
61 bands	4 bands	45.56	0.986	2.0299	3.1919	0.9914	0.0092
121 bands	4 bands	45.62	0.9859	2.0476	3.2506	0.9913	0.0094
31 bands	4 bands	45.47	0.9859	2.0517	3.1815	0.9913	0.0092
16 bands	4 bands	45.21	0.9858	2.0456	3.1974	0.9913	0.0092
121 bands	3 bands	42.54	0.9768	4.2612	3.6212	0.975	0.0109
61 bands	3 bands	42.51	0.9771	4.2225	3.5594	0.9755	0.0107
31 bands	3 bands	42.47	0.977	4.2165	3.5467	0.9755	0.0107
16 bands	3 bands	42.37	0.9768	4.2571	3.5744	0.9751	0.0108

Table 4: Reconstruction encoder \rightarrow decoder results on different numbers of RWA compacted bands and PCA compacted bands.

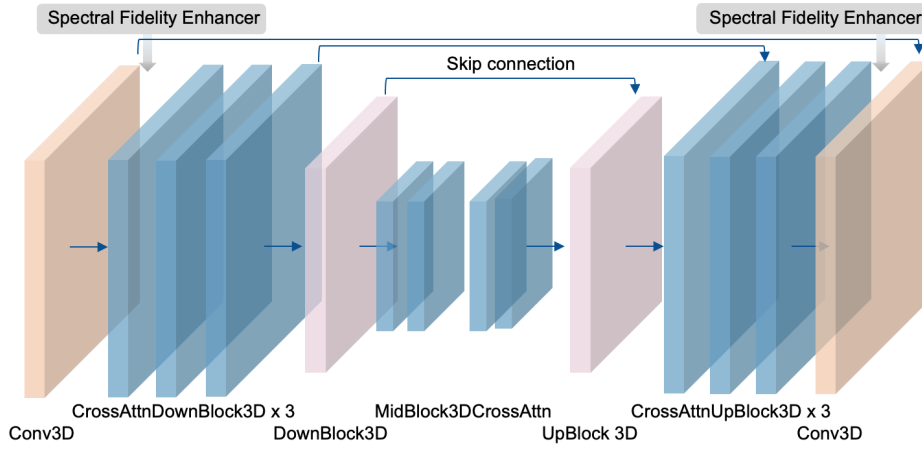


Figure 12: Illustration of three-dimensional U-Net structure.

that emphasizes key spectral features during the generation process. By guiding the model to focus on critical wavelength information, the SFE helps reduce spectral distortion and enhances the fidelity of the reconstructed spectral profiles.

Hyperparameter analysis in sampling stage

In this section, we give a more detailed explanation of the hyperparameters that were used in our method.

Effect of ρ on sampling and spectral fidelity.

In EDM, the parameter ρ governs the curvature of the noise schedule that interpolates between the maximum noise level σ_{\max} and the minimum noise level σ_{\min} across sampling steps. The ρ of EDM does not decrease linearly from $\sigma_{\max} \rightarrow \sigma_{\min}$, but is calculated by the following formula:

$$\sigma_n = \left(\sigma_{\max}^{1/\rho} + \frac{n}{N-1} \left(\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho} \right) \right)^\rho, \quad (22)$$

Specifically, ρ determines whether these σ decrease linearly ($\rho = 1$) or nonlinearly ($\rho > 1$), as shown in Figure. 13

ρ controls how the noise levels decay, with larger values resulting in a schedule that rapidly decreases at early steps and flattens at later steps. This structure allows more refinement in the low-noise regime, which typically enhances perceptual detail and image fidelity. However, in the context of conditional hyperspectral image synthesis, maintaining spectral consistency across all bands is crucial. High values of ρ tend to introduce increased stochasticity in the denoising process, especially in high-frequency components, potentially leading to unwanted spectral artifacts or incoherent band-wise variations. We designed a group of experiments on the validation dataset. To mitigate this, we adopt a relatively low value $\rho \in [0.6, 0.7]$, which produces a more linear and uniform noise schedule. This setting constrains the denoising trajectory, suppressing random high-frequency noise while promoting better alignment with the spectral distribution of the conditioning data. As a result, it helps preserve

spectral fidelity and structural coherence in the reconstructed hyperspectral images.

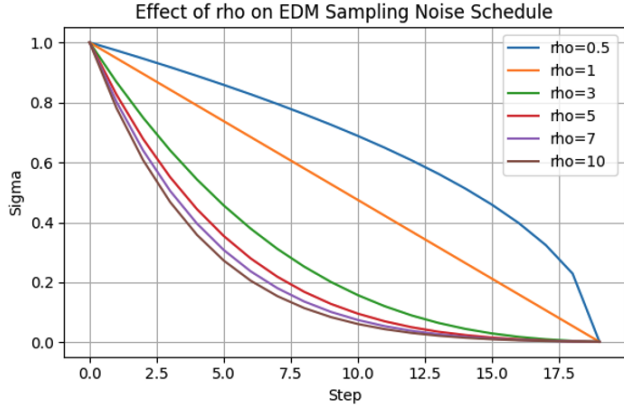


Figure 13: Effect of ρ on EDM sampling noise schedule.

Effect of σ_{\max} and σ_{\min} on sampling and spectral fidelity.

σ_{\max} and σ_{\min} are the two parameters that control the maximum noise strength and minimum noise strength in the sampling stage. At this point, the diffusion model has learned the ability to denoise within a specific noise scale range. We test the effect of these two parameters on the sampling stage. σ_{\max} represents the initial loss, which will affect the “divergence” of the image. A higher σ_{\max} causes higher initial noise, higher image freedom. The result image may contain more details, but with large amounts of noise. A lower σ_{\max} causes more conservative results, and may lose color information. We adopted $\sigma_{\max} = 80$ as the final setting. σ_{\min} represents the end noise, which will affect the “convergence” of the image. The σ_{\min} is large, the convergence is shorter, and the image is more blurry. Smaller σ_{\min} causes longer convergence process. Reconstruction images contain more details but also have more artifacts. The best performance σ_{\min} range is [0.002, 0.2].

Training and testing settings.

Our training and testing pipeline involves several key hyperparameters that influence both performance and computational efficiency. For spectral dimensionality reduction, we first compact the hyperspectral data using a wavelet transform, retaining 121 bands (compact bands), followed by a PCA transformation during training to further reduce the spectral dimension to 20 bands (PCA bands). Training is conducted with a batch size of 1 and typically runs for 200 epochs. The number of GPUs (num processes) can vary from 1 to 4, depending on available hardware. To improve learning quality, we apply a combination of loss functions: pixel loss (l1 lambda = 0.8), perceptual loss (l2 lambda = 0.1), and gradient loss (l3 lambda = 0.1). We also include optional modules such as mask conditioning (mask) and edge perturbation (edge) to enhance robustness. During inference, the sampling process is controlled by the number of diffu-

sion steps (timesteps), typically set to 50 for a balance between quality and speed. Noise scheduling is governed by σ_{\min} (0.002–0.2), σ_{\max} (80–90), and σ_{data} (0.5), while the parameter ρ (0.6–0.7) adjusts the curvature of the sampling trajectory. Additionally, the recall option allows resuming training from any specific epoch. These parameters are tuned to achieve a balance between computational cost and reconstruction fidelity.

Table 3 and Figure 14 present a summary of hyperparameter effects and visualize the outcomes of suboptimal parameter settings. Before the camera-ready version is finalized, we will release the model checkpoints to support reproducibility. We encourage users to adjust these hyperparameters for their specific datasets, eliminating the need for retraining.

Robustness under imperfect conditioning.

To verify the robustness of GEWDiff under imperfect conditions, we conducted controlled perturbation experiments on mask input and low-resolution (LR) images to simulate real segmentation and sensor noise. Specifically, we applied two methods: (1) random mask erosion/dilation (1–3 pixels) combined with random spatial translation (± 1 –2 pixels) to assess its sensitivity to geometric boundary errors; and (2) introducing 1% additive Gaussian noise into the low-resolution input to test its spectral stability. Quantitative evaluation results in Table 6 show that segmentation perturbation leads to a 0.03 dB decrease in PSNR and an increase in FID by 11 (MDAS1), while noisy low-resolution input leads to a 0.2 dB decrease in PSNR and an increase in FID by 1.5. Notably, the structural fidelity remains essentially unchanged under these perturbations, indicating that GEWDiff’s cross-scale geometric prior and diffusion denoising process effectively mitigates moderate mask errors and input noise. Although extreme segmentation failures may still affect the reconstruction results, these results demonstrate that GEWDiff can generalize beyond ideal mask conditions and maintain stable performance under real noise and boundary deviations. Moreover, the mask and NDVI used in our work will not generate extra data demand, as they were derived from low-resolution images. In future work, we will explore uncertainty-aware mask conditionalization and adaptive attention weighting to further improve their robustness to highly ambiguous or noisy inputs.

Practical utility on downstream land-cover classification task

To further evaluate the practicality of GEWDiff, we assessed its impact on downstream hyperspectral land cover classification tasks. We randomly sampled 30 points covering four semantic categories (trees, buildings, bare soil, impermeable surfaces) and trained a random forest classifier using super-resolution output (SR) as input. For comparison, we applied the same classifier configuration and sampling scheme to the original 10-meter low-resolution (LR) input. Experiments were conducted on two independent samples (MDAS 1 and MDAS 2), and the classification results were benchmarked against pseudo-ground values extracted from high-resolution data, as shown in Figure 15.

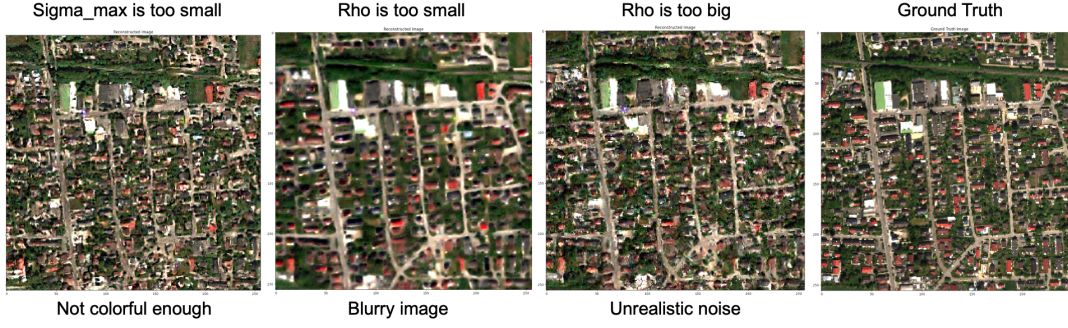


Figure 14: Three typical effects of unsuitable parameters on validation set (informative only).

Parameter	Higher	Lower	Influence
σ_{\max}	Larger initial noise, high image freedom, more details, but with larger noise.	The initial noise is small, more conservative, and may lose color information.	This will affect the "divergence" of the image.
σ_{\min}	The end noise is large, the convergence is poor, and the image is blurry.	Smaller ending noise and cleaner convergence.	Smaller is cleaner, but also more risky.
ρ	The number of steps is denser in the early stage and sparser in the late stage.	The number of steps is sparse in the early stage and denser in the late stage.	Control the curvature of the "sampling curve".
t	Larger and more precise, but slower to compute.	The sampling speed is fast, but it is easy to be blurry.	Affects the total number of steps for denoising.

Table 5: Summary of hyperparameter effects (informative only)

Dataset	Methods	Mask perturbation		Noise perturbation	
MDAS 1	PSNR \uparrow	28.8372	0.0544	28.6	0.12
	SSIM \uparrow	0.7286	0.0032	0.719	0.005
	SAM \downarrow	8.5045	0.0323	8.684	0.008
	CC \uparrow	0.8113	0.0025	0.803	0.005
	RMSE \downarrow	0.05246	0.00028	0.05313	0.00067
	FID \downarrow	56.8863	0.8352	46.05	5.3
	LV \uparrow	0.003698	0.000022	0.0037	0.00009

Table 6: Conditions perturbation experiments

In both samples, GEWDiff consistently improved the accuracy of semantic classification. On MDAS1, SR improved the average IoU from 0.4312 to 0.4486 and the average F1 score from 0.5553 to 0.5761, achieving more accurate class boundary delimitation while maintaining the same overall accuracy (OA) as the LR baseline (0.9018). On the more challenging MDAS2 samples, SR significantly improved OA (0.6996 \rightarrow 0.7222), mIoU (0.3967 \rightarrow 0.4181), and mF1 (0.5319 \rightarrow 0.5572). Visual comparison (Appendix Figure X) reveals a clearer distinction between built-up areas and vegetation, as well as a more refined detection of impermeable structures. These results demonstrate that GEWDiff does not merely upsample spectral data, but generates representations capable of distinguishing more class features and improving downstream discriminative capabilities. Therefore, even with limited supervision information and a small number of training samples, this model can bring tangible benefits to practical remote sensing applications.

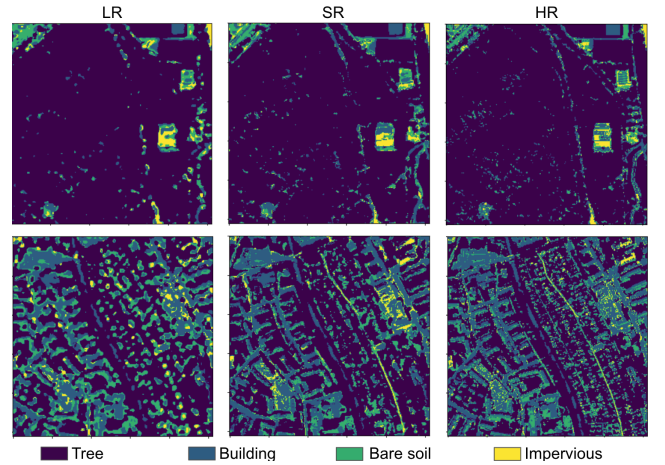


Figure 15: Classification map of low-resolution image, super-resolution image, and high-resolution image.