

NODiff: Neural Operator Diffusion for Multispectral Image Fusion

Junming Hou^{1*}, Ran Ran^{2*}, Sixing Chen¹, Zihao Chen¹, Xiaofeng Cong¹, Junling Li^{1†}, Liang-Jian Deng^{2, 3†}

¹Southeast University, Nanjing 211189, China

²University of Electronic Science and Technology of China, Chengdu 611731, China

³Multi-Hazard Early Warning Key Laboratory of Sichuan Province

junming_hou@seu.edu.cn, ranran@std.uestc.edu.cn, sixingchen@seu.edu.cn, 213233966@seu.edu.cn, cxf_svip@163.com, junlingli@seu.edu.cn, liangjian.deng@uestc.edu.cn

Abstract

Pansharpening is a powerful technique for generating high-resolution multispectral (HRMS) images by fusing currently available image pairs of low-resolution multispectral (LRMS) and texture-rich panchromatic (PAN) data, effectively addressing the physical constraints of satellite sensors. While recent generative diffusion models have demonstrated impressive performance gains in this domain, their prohibitive computational demands and training costs hinder practicality in resource-constrained remote sensing satellite systems. In this work, we propose NODiff, a novel diffusion framework that replaces the conventional attention-based denoising backbone with a neural operator, seamlessly integrating operator learning and generative modeling into an efficient yet effective solution for pansharpening. In practice, we implement our approach through a two-stage learning paradigm: First, we pretrain the proposed Neural Operator-based diffusion model to learn the high-resolution texture priors essential for pansharpening. Afterward, we freeze the pretrained parameters, and design a lightweight conditional detail guidance adapter to enable efficient fine-tuning for generating desired HRMS images. Meanwhile, a time-aware low-rank adaptation is introduced to dynamically refine high-frequency details potentially affected by spectral mode truncation. Extensive experiments on multiple benchmark datasets demonstrate that NODiff achieves competitive pansharpening performance while significantly reducing training and inference costs. Beyond pansharpening, our method provides new insights into building resource-efficient generative models.

Code — <https://github.com/coder-JMHou/NODiff>

Introduction

High-resolution remote sensing multispectral imagery has found widespread application across a broad spectrum of civilian and military domains, such as precision agriculture, environmental monitoring, mapping services, and target recognition (Meng et al. 2019; Vivone et al. 2024). Due to physical imaging constraints, current satellite multispectral sensors cannot directly acquire high-resolution multispectral (HRMS) data. As a common solution, they often ob-

*Contributed equally.

†Corresponding authors.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

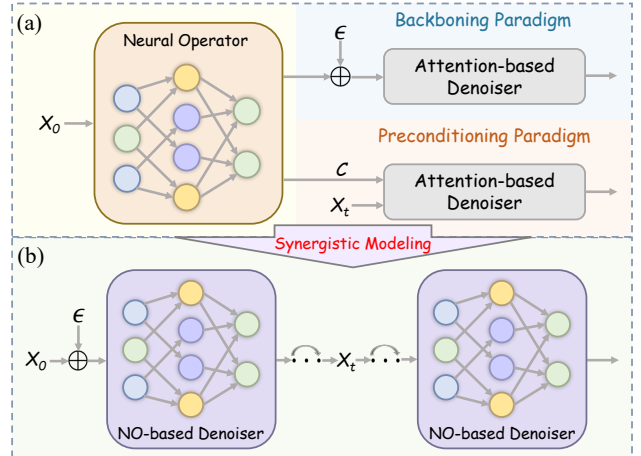


Figure 1: Comparison of integration schemes between Neural Operators and Diffusion Models: (a) Existing approaches (Liu and Tang 2025; Xu et al. 2025), and (b) Our design. Rather than the simple integration, our design delves into synergizing the strengths of two emerging neural architectures, resulting in an efficient and effective framework.

serve a low-resolution (LR) MS image and its paired texture-rich panchromatic (PAN) image of the same scene, providing complementary spectral and spatial information, respectively. Pansharpening technology is developed to fuse the complementary information from PAN and MS modalities to generate HRMS images, thereby benefiting the downstream remote sensing applications.

To date, the pansharpening field has undergone substantial methodological development, evolving from traditional model-driven algorithms to deep learning paradigms. Traditional approaches primarily include Component Substitution (CS), Multi-Resolution Analysis (MRA), and Variational Optimization (VO), according to their underlying fusion principles and algorithmic paradigms. However, these approaches often heavily rely on manually designed priors and unrealistic assumptions, leading to suboptimal fusion quality. With the growing availability of large-scale datasets, recent years have witnessed substantial advances in deep learning-based pansharpening methods, which significantly outperform traditional algorithms. These approaches exploit

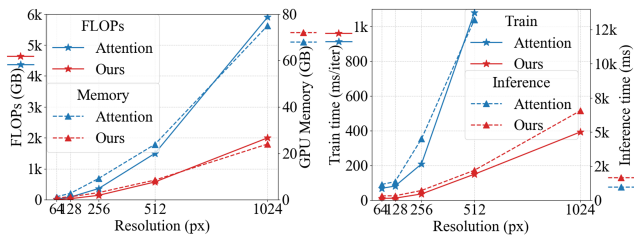


Figure 2: Comparison of computational cost and efficiency between our design and conventional attention-based denoising backbone, in terms of FLOPs, GPU memory usage, training time per iteration, and inference time.

rich training data to directly capture the complex non-linear mapping from input PAN–LRMS pairs to their corresponding HRMS outputs by training deep neural networks. In pursuit of improved performance, a variety of advanced network architectures have been proposed, spanning from convolutional neural networks (CNNs) to global transformers.

Very recently, diffusion models have been applied to pansharpening, yielding superior performance over conventional deep learning models. However, such performance gains come at the expense of substantial computational demands and training costs, posing significant challenges for practical applications, particularly in resource-constrained remote sensing satellites. Neural operators (NOs) are a recently developed neural architecture designed to solve partial differential equations (PDEs). They directly learn mappings between infinite-dimensional function spaces, providing resolution-invariant paradigms that generalize across varying input resolutions. Recent efforts have been made to integrate NOs and diffusion models for image super-resolution, as summarized in Figure 1: (1) Preconditioning: NOs act as preprocessors (*e.g.*, up-sample operator) to produce conditional detail guidance for diffusion model (Liu and Tang 2025); (2) Backboning: NOs serve as the super-resolution backbone, refined subsequently by diffusion processes (Xu et al. 2025). While promising, these approaches inherit the computational burden and inefficiency of diffusion models, limiting their practical applicability.

In this work, we build upon a novel perspective to synergize the strengths of these two emerging neural paradigms. Instead of simple integration, our approach assimilates the neural operator as the denoising backbone of the diffusion model by replacing conventional attention-based architectures, resulting in an efficient yet effective NO-based diffusion framework. As shown in Figure 2, our NO-based design demonstrates competitive computational efficiency compared to conventional attention-based counterparts. In practice, we implement it through a two-stage training paradigm: First, we pretrain the proposed NO-based diffusion model to steer the diffusion generative process toward learning high-resolution texture priors, which are pivotal for pansharpening. Second, we introduce a straightforward conditional detail guidance adapter that takes LRMS and PAN image pairs as inputs to fine-tune the pretrained NO-based diffusion model for generating corresponding

HRMS images. Meanwhile, a lightweight time-aware low-rank adaptation (LoRA) technique is employed to further refine high-frequency details that may be affected due to spectral mode truncation. Our method, built upon these two core processes, achieves competitive pansharpening performance compared to existing advanced approaches. Notably, compared to conventional attention-based diffusion models, our design achieves a favorable trade-off between efficiency and effectiveness by integrating the complementary strengths of two emerging and promising neural paradigms. We summarize the main contributions as follows:

- We propose NODiff, a novel diffusion framework that replaces costly attention-based denoiser architectures with resolution-invariant neural operator, thereby synergizing operator learning with generative modeling into an efficient yet effective solution for pansharpening.
- We introduce a lightweight conditional detail guidance adapter for parameter-efficient fine-tuning of the pre-trained knowledge-rich diffusion model, complemented by a time-aware low-rank adaptation mechanism that dynamically refines high-frequency textures affected by spectral mode truncation, enabling the generation of high-quality HRMS images.
- Extensive experiments on multiple benchmark datasets demonstrate that our method outperforms current advanced pansharpening approaches, and also delivers a promising performance-efficiency balance compared to conventional attention-based diffusion counterparts.

Related Works

Deep Learning-based Pansharpening. In recent years, deep learning has emerged as the prevailing paradigm in the pansharpening domain owing to its powerful nonlinear learning capability, greatly surpassing traditional algorithms (Dian et al. 2021; Vivone et al. 2024). The implementation of PNN indicates the superiority of deep learning methods, sparking the development of diverse advanced network architectures, including convolutional neural networks (CNNs) (Yang et al. 2017; Deng et al. 2020; Hu et al. 2021; Hou et al. 2025), global transformers (Zhou et al. 2022; Hou et al. 2024; Liu, Dian, and Li 2025; Wu et al. 2025), and recently emerging generative diffusion models (Meng et al. 2023; Cao et al. 2024; Kim et al. 2025).

Image Super-resolution Neural Operator. As newly developed data-driven solutions for PDEs, neural operators directly learn mappings between infinite-dimensional function spaces (Kovachki et al. 2023; Aizzadenesheli et al. 2024). By optimizing network training in function spaces, NOs achieve superior nonlinear approximation and exhibit invariance to discretization (Lu et al. 2021; Li et al. 2021; Tran et al. 2023). More importantly, this resolution-invariant property enables NOs to generalize well across discretizations, thus making them particularly well-suited for resolution-varying tasks. For example, SRNO first presents a NO-based continuous image super-resolution approach (Wei and Zhang 2023). HiNOTE introduces a hierarchical NO framework for arbitrary-scale super-resolution of scientific data (Luo, Qian, and Yoon 2024). He et al. apply this novel

neural architecture to continuous-resolution hyperspectral pansharpening (He et al. 2025a). Very recently, several studies have integrated neural operators with diffusion models to harness their complementary strengths, yielding promising results (Liu and Tang 2025; Xu et al. 2025). Nonetheless, these approaches still inherit the substantial computational overhead characteristic of diffusion models. By contrast, we employ NOs as highly efficient and effective noise predictors by replacing the conventional attention-based counterparts, achieving a desirable performance-efficiency balance. Most importantly, this offers a novel perspective for combining these two emerging neural architectures.

Diffusion Models for Pansharpening. As a prevailing generative model, denoising diffusion probabilistic models (DDPMs) (Ho, Jain, and Abbeel 2020) have found wide applications in various domains, such as text-to-image generation (Saharia et al. 2022; Gu et al. 2022) and image editing (Kawar et al. 2023; Couairon et al. 2022; Shi et al. 2024). In addition, DDPMs have also shown notable efficacy in image processing tasks. Building on this, Song et al. (Song, Meng, and Ermon 2021) propose the denoising diffusion implicit model (DDIM) which employs a non-Markovian sampling process to accelerate diffusion model inference. In recent years, diffusion models have attracted growing interest in the pansharpening field, with PAN and LRMS serving as conditioning information (Kim et al. 2025; Cao et al. 2024). Although these methods exhibit superior fusion capabilities over conventional neural networks, they often require greater computational resources and training costs, which hinder their broader applications. This motivates us to explore more efficient variants that are suitable for resource-constrained remote sensing satellites.

Methodology

Overview of the Proposed Method

Accurately reconstructing fine-grained texture details is a key objective of pansharpening learning, as required by downstream remote sensing applications. Very recently, diffusion models have been introduced to learn the high-resolution residuals between LRMS and HRMS images, demonstrating impressive pansharpening performance improvements (Cao et al. 2024; Kim et al. 2025). However, these methods require substantial computational resources and training costs, limiting their practicality on resource-constrained remote sensing satellites. Our approach seeks to remedy this by incorporating high-resolution texture-specific priors into a computationally efficient diffusion model. To achieve this, we implement our method through a two-stage learning paradigm illustrated in Figure 3: (1) High-resolution Texture Prior Learning: First, we pretrain the proposed Neural Operator-based conditional diffusion model to learn high-resolution texture priors. (2) Conditional Detail Guidance Fine-Tuning: Following that, we use a straightforward conditional detail guidance adapter to fine-tune the pretrained diffusion model for generating corresponding HRMS images, aided by a time-aware LoRA. Our method synergizes the unique strengths of two emerging and powerful neural architectures, offering new insights for

building efficient yet effective generative pansharpening approaches suitable for resource-constrained satellites.

Learning High-Resolution Diffusion Texture Priors

In this stage, our primary objective is to integrate high-resolution texture-specific priors into a diffusion model, enhancing its ability to reconstruct fine-grained textures while maintaining computational efficiency. Motivated by recent advances in integrating diffusion models with neural operators for image super-resolution (Liu and Tang 2025; Xu et al. 2025), we seek to explore this promising synergy to construct an efficient generative pansharpening model. To this end, we revisit the unique strengths of these two emerging neural architectures and identify the key aspects particularly suited for super-resolution scenarios: 1) the Neural Operators’ function-space mapping capacity and resolution-invariance; and 2) the diffusion model’s generative ability to produce fine-grained details. Building upon these insights, we develop an efficient yet effective Neural Operator-based diffusion model, where the denoising network is constructed by cascaded Fourier Neural Operator (FNO) layers (Li et al. 2021) shown in Figure 3, replacing conventional attention-heavy architectures.

Let $\mathbf{L} \in \mathbb{R}^{H \times W \times s}$ denote the up-sampled LRMS image, and $\mathbf{H} \in \mathbb{R}^{H \times W \times s}$ the corresponding reference HRMS image. We define the residual image $\mathbf{X}_0 \in \mathbb{R}^{H \times W \times s}$ as $\mathbf{H} - \mathbf{L}$, which contains rich texture details absent in \mathbf{L} . Similar to (Cao et al. 2024; Kim et al. 2025), the forward diffusion process adds Gaussian noise to \mathbf{X}_0 over T iterations in the Markov steps:

$$q(\mathbf{X}_t | \mathbf{X}_{t-1}) = \mathcal{N}(\mathbf{X}_t; \sqrt{1 - \beta_t} \mathbf{X}_{t-1}, \beta_t I), \quad (1)$$

where \mathbf{X}_t denotes the noisy residual image at step t , and β_t controls the noise variance. Through reparameterization, \mathbf{X}_t can be obtained directly by adding noise in one step to \mathbf{X}_0 :

$$\mathbf{X}_t = \sqrt{\bar{\alpha}_t} \mathbf{X}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad (2)$$

where ϵ denotes a standard Gaussian distribution and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$, $\alpha_t = 1 - \beta_t$. The reverse process samples \mathbf{X}_0 from pure noise $\mathbf{X}_T \sim \mathcal{N}(0, I)$ step-by-step. The posterior distribution of the denoising process is written as:

$$q(\mathbf{X}_{t-1} | \mathbf{X}_t, \mathbf{X}_0) = \mathcal{N}\left(\mathbf{X}_{t-1}; \mu_t(\mathbf{X}_t, \mathbf{X}_0), \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t I\right),$$

$$\mu_t(\mathbf{X}_t, \mathbf{X}_0) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{X}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right), \quad (3)$$

where ϵ_t is the t -th step added noise in the forward process.

To address the efficiency-performance trade-off in conventional attention-based denoising architectures, we design an efficient FNO-based denoiser ϵ_θ that predicts the noise at each step. We use the \mathbf{H} as high-resolution conditional guidance by concatenating it with the noisy image \mathbf{X}_t :

$$\mathbf{F}_t = \text{Conv}([\mathbf{X}_t, \mathbf{H}]). \quad (4)$$

Subsequently, the resulting \mathbf{F}_t is processed by the denoising network composed of cascaded FNO layers, as follows:

$$\mathbf{F}_t^l = \sigma \left(\underbrace{\mathcal{F}^{-1}(\mathcal{P}_{\mathcal{K}}(\mathcal{F}(\mathbf{F}_t^{l-1})))}_{\text{Global Spectral Operator}} \odot \mathbf{R} + \underbrace{\mathbf{W} \mathbf{F}_t^{l-1}}_{\text{Linear Transform}} \right), \quad (5)$$

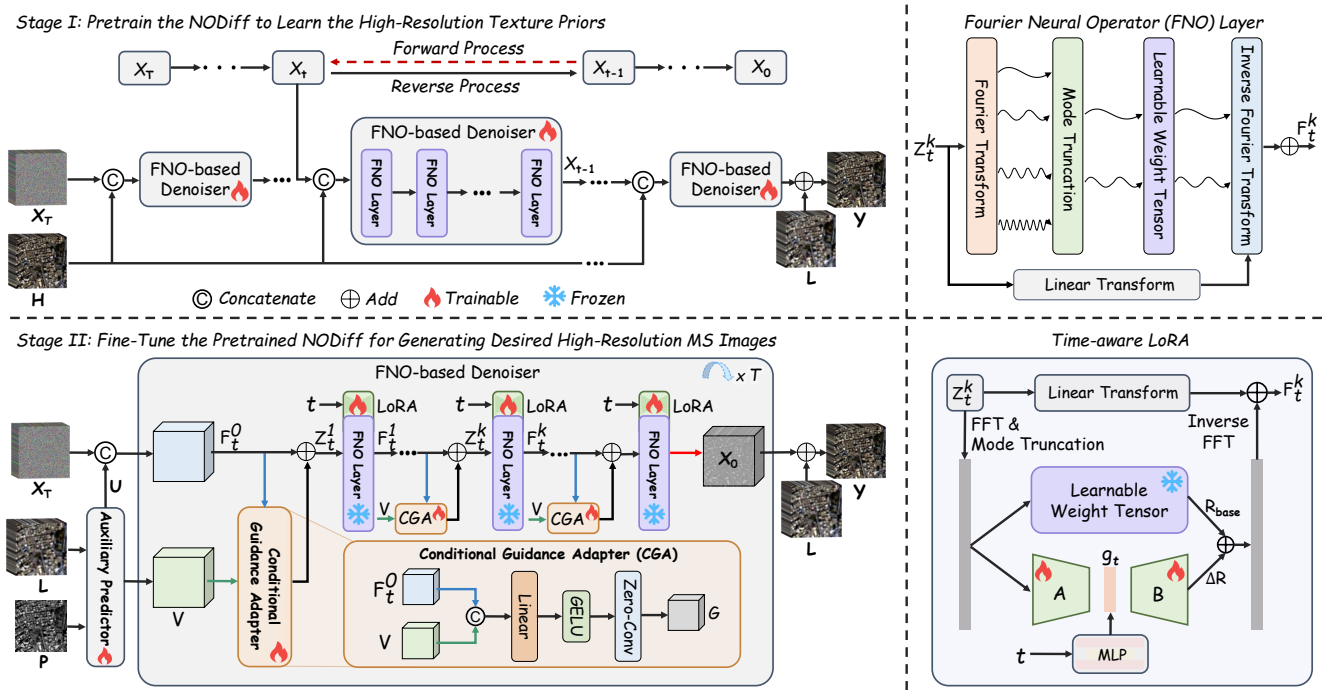


Figure 3: The framework of the proposed NODiff which is implemented through two main procedures. Firstly, we pretrain the NODiff to learn high-resolution diffusion texture priors conditioned on the HRMS images. Subsequently, we introduce a lightweight conditional detail guidance adapter using LRMS and PAN image pairs as inputs to fine-tune the pretrained NODiff for generating corresponding HRMS images, assisted by a time-aware LoRA design.

where $\mathbf{F}_t^l (l = 1, 2, \dots, m)$ represents the features of the l -th FNO layer. $\mathcal{F}(\cdot)$ and $\mathcal{F}^{-1}(\cdot)$ denote the Fourier and inverse Fourier transforms, respectively. $\mathcal{P}_{\mathcal{K}}(\cdot)$ signifies the spectral mode truncation operation. R is a learnable complex weight tensor. W and $\sigma(\cdot)$ represent the linear transformation and nonlinear activation.

After iterative refinement, we project the last FNO layer's output to obtain the predicted noise at step t :

$$\hat{\epsilon}_t = \text{Conv}(\mathbf{F}_t^m). \quad (6)$$

With the predicted noise $\hat{\epsilon}_t$, the reverse process is parameterized as:

$$\begin{aligned} \mathbf{X}_{t-1} &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{X}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(\mathbf{X}_t, \mathbf{H}, t) \right) \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{X}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon}_t \right). \end{aligned} \quad (7)$$

After T -step denoising process, we can reconstruct the high-resolution multispectral image $\mathbf{Y} \in \mathbb{R}^{H \times W \times s}$ by adding the predicted residual image \mathbf{X}_0 to \mathbf{L} . Accordingly, the optimization objective is defined as:

$$\mathcal{L}_{\text{Stage I}} = \mathbb{E} \left\| \epsilon - \epsilon_\theta \left(\sqrt{\alpha_t}(\mathbf{H} - \mathbf{L}) + \sqrt{1 - \alpha_t} \epsilon, \mathbf{H}, t \right) \right\|_1. \quad (8)$$

Conditional Detail Guidance Fine-Tuning

The pretraining stage enables NODiff to encapsulate rich prior knowledge, which we aim to fully retain and leverage during downstream pansharpening. To achieve this, we

freeze the parameters of the pretrained FNO-based denoiser and introduce a flexible and parameter-efficient adaptation strategy implemented through the cross-modality detail guidance adapter and time-aware low-rank adaptation illustrated in Figure 3. Specifically, we first introduce a straightforward auxiliary predictor consisting of several convolution blocks, which takes the Sobel-filtered PAN and LRMS as inputs to produce: 1) The pseudo super-resolved (SR) image to replace the \mathbf{H} in the Stage I; and 2) The cross-modality detail features as the input of the conditional detail guidance adapter along with the latent noisy feature \mathbf{X}_t :

$$\begin{aligned} \mathbf{V} &= \Phi(\mathcal{S}(\mathbf{P}, \mathbf{L})), \\ \mathbf{U} &= \text{Conv}(\mathbf{V}) + \mathbf{L}, \end{aligned} \quad (9)$$

where $\mathcal{S}(\cdot)$ denotes the mapping of Sobel operator, and $\Phi(\cdot)$ represents the backbone of the auxiliary predictor, where $\mathbf{V} \in \mathbb{R}^{H \times W \times C}$ is the extracted cross-modality detail features. \mathbf{U} is the initially predicted pseudo SR used to replace the high-resolution reference image \mathbf{H} .

At the diffusion step t , we first concatenate \mathbf{U} with the noisy feature \mathbf{X}_t along the channel dimension, and then fuse them through a convolution layer, resulting in the input feature $\mathbf{F}_t^0 \in \mathbb{R}^{H \times W \times C}$ of the FNO-based denoiser:

$$\mathbf{F}_t^0 = \text{Conv}([\mathbf{X}_t, \mathbf{U}]). \quad (10)$$

Then, we inject the cross-modality detail feature \mathbf{V} into each intermediate layer of the denoiser through a conditional detail guidance adapter, which is realized through a bidirectional information flow interaction. Specifically, for the l -th

Method	Reduced Resolution				Full Resolution		
	PSNR(\pm std)	SAM(\pm std)	ERGAS(\pm std)	Q2 ⁿ (\pm std)	D_λ (\pm std)	D_s (\pm std)	HQNR(\pm std)
BDS-PC (Vivone 2019)	32.969 \pm 2.784	5.429 \pm 1.823	4.698 \pm 1.617	0.829 \pm 0.097	0.063 \pm 0.024	0.073 \pm 0.036	0.870 \pm 0.053
BT-H (Lolli et al. 2017)	33.080 \pm 2.880	4.920 \pm 1.425	4.579 \pm 1.496	0.832 \pm 0.094	0.057 \pm 0.023	0.081 \pm 0.037	0.867 \pm 0.054
LRTCFFan (Wu et al. 2023)	33.613 \pm 2.839	4.737 \pm 1.412	4.315 \pm 1.442	0.846 \pm 0.091	0.018 \pm 0.007	0.053 \pm 0.026	0.931 \pm 0.031
FusionNet (Deng et al. 2020)	38.042 \pm 2.592	3.325 \pm 0.698	2.467 \pm 0.645	0.904 \pm 0.090	0.024 \pm 0.009	0.036 \pm 0.014	0.941 \pm 0.020
LAGConv (Jin et al. 2022)	38.568 \pm 2.789	3.104 \pm 0.559	2.300 \pm 0.613	0.910 \pm 0.091	0.037 \pm 0.015	0.042 \pm 0.015	0.923 \pm 0.025
Fourmer (Zhou et al. 2023)	38.268 \pm 2.727	3.236 \pm 0.681	2.419 \pm 0.665	0.911 \pm 0.090	0.022 \pm 0.010	0.035 \pm 0.004	0.944 \pm 0.013
HFIN (Tan et al. 2024)	38.534 \pm 2.786	3.088 \pm 0.635	2.306 \pm 0.557	0.912 \pm 0.089	0.025 \pm 0.008	0.043 \pm 0.017	0.934 \pm 0.024
HOIF (Zhou et al. 2024)	38.352 \pm 2.855	3.186 \pm 0.643	2.385 \pm 0.668	0.913 \pm 0.086	0.039 \pm 0.021	0.039 \pm 0.010	0.924 \pm 0.026
PanDiff (Meng et al. 2023)	37.860 \pm 2.773	3.316 \pm 0.674	2.492 \pm 0.656	0.906 \pm 0.088	0.027 \pm 0.012	0.054 \pm 0.026	0.920 \pm 0.036
PanMamba (He et al. 2025b)	39.012 \pm 2.818	2.914 \pm 0.592	2.184 \pm 0.521	<u>0.920\pm0.085</u>	0.018 \pm 0.007	0.031 \pm 0.011	0.952 \pm 0.015
ADWM (Huang et al. 2025)	39.170 \pm 2.878	<u>2.914\pm0.589</u>	<u>2.145\pm0.531</u>	0.919 \pm 0.086	0.024 \pm 0.010	0.029\pm0.015	0.948 \pm 0.021
NODiff (ours)	39.314\pm2.898	2.859\pm0.585	2.109\pm0.526	0.923\pm0.083	0.014\pm0.005	<u>0.030\pm0.010</u>	0.956\pm0.012
Ideal value	∞	0	0	1	0	0	1

Table 1: Quantitative results for reduced and full resolution WV3 samples, comparing several representative state-of-the-art methods. Bold: Best; Underline: Second best.

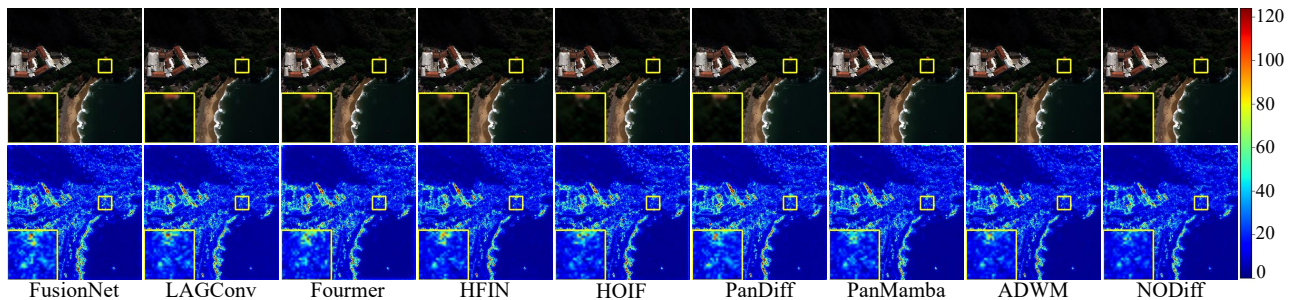


Figure 4: The visual results (top) and mean absolute error maps (bottom) of all compared DL-based methods on a reduced resolution WV3 testing sample.

FNO layer of the denoiser, the adapter first integrates the $(l-1)$ -th output noisy feature $\mathbf{F}_t^{l-1} \in \mathbb{R}^{H \times W \times C}$ ($l = 1, 2, \dots, m$) with \mathbf{V} to produce a residual guidance feature $\mathbf{G} \in \mathbb{R}^{H \times W \times C}$, which is then added to \mathbf{F}_t^{l-1} to form the input feature $\mathbf{Z}_t^l \in \mathbb{R}^{H \times W \times C}$ for the l -th FNO layer:

$$\begin{aligned} \mathbf{G} &= \phi([\mathbf{V}, \mathbf{F}_t^{l-1}]), \\ \mathbf{Z}_t^l &= \mathbf{F}_t^{l-1} + \mathbf{G}, \end{aligned} \quad (11)$$

where $\phi(\cdot)$ denotes the mapping of the detail guidance adapter as illustrated in Figure 3.

In parallel, we introduce a time-aware low-rank adaptation (Ta-LoRA) technique into the spectral convolution of the FNO layer to further refine high-frequency details that may be affected owing to mode truncation. Let \mathbf{R}_{base} be the complex weight tensor corresponding to the retained mode set (frozen). A timestep embedding is first mapped into a rank-wise vector $\mathbf{g}_t \in \mathbb{R}^r$ using an MLP layer:

$$\mathbf{g}_t = 1 + \text{Tanh}(\text{MLP}(\text{Embed}(t))). \quad (12)$$

Building on the diagonalized form $\mathbf{D}_t = \text{diag}(\mathbf{g}_t) \in \mathbb{R}^{r \times r}$, the low-rank factors $\mathbf{A} \in \mathbb{R}^{r \times C_{\text{in}}}$ and $\mathbf{B} \in \mathbb{C}^{C_{\text{out}} \times r \times m_x \times m_y}$, and a scaling factor α , the time-modulated low-rank increment is computed as:

$$\Delta \mathbf{R} = \alpha \cdot \mathbf{A} \mathbf{D}_t \mathbf{B}. \quad (13)$$

Thus, the final effective spectral weights are updated as:

$$\mathbf{R} = \mathbf{R}_{\text{base}} + \Delta \mathbf{R}. \quad (14)$$

Notably, the increment $\Delta \mathbf{R}$ is applied exclusively to the retained spectral modes including positive/negative bands, leaving the pretrained FNO dataflow unchanged. Thus, this Ta-LoRA incurs negligible additional parameters and computational overhead. Moreover, the time-aware vector \mathbf{g}_t dynamically adapts the LoRA increment across different noise levels, making the later iterations more sensitive and effective for high-frequency details. The optimization objective is formally identical to that of Stage I and is expressed as:

$$\mathcal{L}_{\text{Stage II}} = \mathbb{E} \left\| \epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_t} (\mathbf{H} - \mathbf{L}) + \sqrt{1 - \bar{\alpha}_t} \epsilon, \mathbf{P}, \mathbf{L}, t \right) \right\|_1. \quad (15)$$

Experiments

Datasets, Metrics and Implementation Details. We experiment on three popular benchmark datasets, including the WorldView-3 (WV3), GaoFen-2 (GF2), and QuickBird (QB), to evaluate our model’s efficacy. The training datasets are simulated using observations from the original satellite imagery according to Wald’s protocol. All the training and test data are publicly available at the PanCollection (Deng et al. 2022). The reduced resolution results are evaluated

Method	Reduced Resolution-GF2				Reduced Resolution-QB			
	PSNR(\pm std)	SAM(\pm std)	ERGAS(\pm std)	Q2 ⁿ (\pm std)	PSNR(\pm std)	SAM(\pm std)	ERGAS(\pm std)	Q2 ⁿ (\pm std)
BDS-PC (Vivone 2019)	35.180 \pm 2.317	1.681 \pm 0.360	1.667 \pm 0.445	0.892 \pm 0.035	32.547 \pm 3.202	8.089 \pm 1.980	7.515 \pm 0.800	0.831 \pm 0.090
BT-H (Lolli et al. 2017)	36.054 \pm 2.236	1.649 \pm 0.360	1.528 \pm 0.409	0.918 \pm 0.025	32.648 \pm 3.273	7.194 \pm 1.552	7.401 \pm 0.838	0.833 \pm 0.088
LRTCFFan (Wu et al. 2023)	37.599 \pm 2.331	1.298 \pm 0.312	1.272 \pm 0.343	0.935 \pm 0.030	33.260 \pm 3.272	7.187 \pm 1.711	6.928 \pm 0.812	0.855 \pm 0.087
FusionNet (Deng et al. 2020)	39.639 \pm 2.270	0.974 \pm 0.212	0.988 \pm 0.222	0.964 \pm 0.009	37.532 \pm 2.518	4.923 \pm 0.908	4.159 \pm 0.321	0.925 \pm 0.090
LAGConv (Jin et al. 2022)	42.735 \pm 1.447	0.786 \pm 0.148	0.687 \pm 0.113	0.980 \pm 0.009	38.181 \pm 2.456	4.547 \pm 0.830	3.826 \pm 0.420	0.934 \pm 0.088
Fourmer (Zhou et al. 2023)	40.670 \pm 1.903	0.976 \pm 0.209	0.885 \pm 0.185	0.970 \pm 0.011	36.797 \pm 2.597	5.079 \pm 1.009	4.494 \pm 0.410	0.924 \pm 0.080
HFIN (Tan et al. 2024)	42.189 \pm 1.752	0.843 \pm 0.148	0.735 \pm 0.126	0.977 \pm 0.011	38.247 \pm 2.403	4.542 \pm 0.805	3.813 \pm 0.322	0.934 \pm 0.085
HOIF (Zhou et al. 2024)	40.982 \pm 1.802	0.943 \pm 0.205	0.841 \pm 0.162	0.974 \pm 0.009	38.242 \pm 2.119	4.521 \pm 0.811	3.825 \pm 0.510	0.933 \pm 0.094
PanDiff (Meng et al. 2023)	42.326 \pm 1.635	0.875 \pm 0.133	0.727 \pm 0.115	0.981 \pm 0.008	<u>38.538\pm2.401</u>	4.524 \pm 0.792	<u>3.688\pm0.343</u>	<u>0.937\pm0.084</u>
PanMamba (He et al. 2025b)	42.907 \pm 1.811	0.743 \pm 0.156	0.684 \pm 0.129	0.982 \pm 0.008	37.356 \pm 2.570	4.625 \pm 0.904	4.277 \pm 0.779	0.929 \pm 0.085
ADWM (Huang et al. 2025)	43.884 \pm 1.714	0.672 \pm 0.130	0.597 \pm 0.107	0.985 \pm 0.006	38.466 \pm 2.420	4.450 \pm 0.809	3.705 \pm 0.346	0.937 \pm 0.085
NODiff (ours)	44.262\pm1.833	0.645\pm0.129	0.573\pm0.110	0.986\pm0.006	38.700\pm2.446	4.372\pm0.787	3.607\pm0.328	0.938\pm0.085
Ideal value	∞	0	0	1	∞	0	0	1

Table 2: Quantitative results for reduced resolution GF2 and QB samples, comparing several representative state-of-the-art methods. Bold: Best; Underline: Second best.

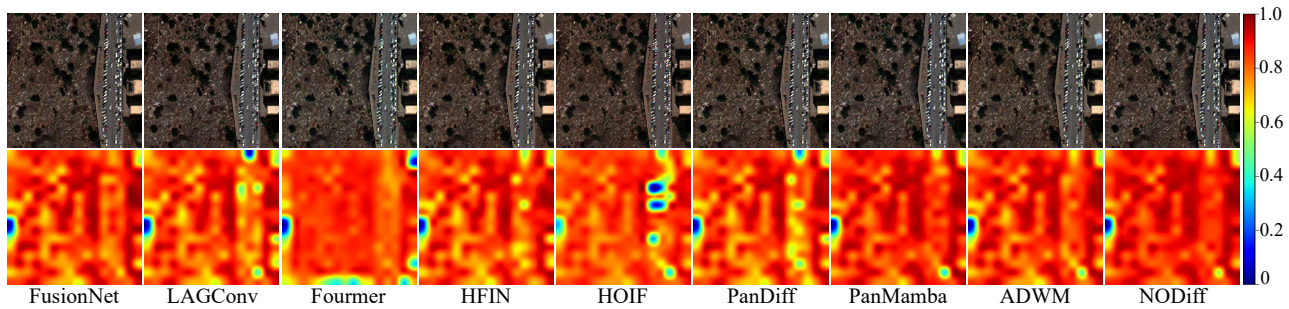


Figure 5: The visual results (top) and the corresponding HQNR maps (bottom) of all compared DL-based methods on a full resolution WV3 testing sample.

using four widely adopted reference-based metrics: PSNR, SAM (Yahas, Goetz, and Boardman 1992), ERGAS (Wald 2002), and Q2n (Zhou, Civco, and Silander 1998). The real-world full resolution performance is compared using three non-reference indicators: D_λ , D_s , and HQNR (Arienzo et al. 2022). All experiments are conducted using the PyTorch framework on an NVIDIA GeForce RTX 4090 GPU. Additional training details and experimental results are provided in the supplementary material due to page limitations.

Baselines. We compare NODiff with several state-of-the-art DL-based pansharpening methods, including FusionNet (Deng et al. 2020), LAGConv (Jin et al. 2022), Fourmer (Zhou et al. 2023), HFIN (Tan et al. 2024), HOIF (Zhou et al. 2024), PanDiff (Meng et al. 2023), PanMamba (He et al. 2025b), and ADWM (Huang et al. 2025). Additionally, three classical algorithms are selected for a comprehensive comparison: BT-H (Lolli et al. 2017), BDS-PC (Vivone 2019), and LRTCFFan (Wu et al. 2023).

Comparison With SOTA Methods

Evaluation on Reduced Resolution Scene. We begin by quantitatively evaluating the similarity between the fused images and ground truth in the reduced resolution setting. As presented in Table 1 (the left panel) and Table 2, our model consistently outperforms the compared cutting-edge pansharpening methods across all evaluation metrics on three

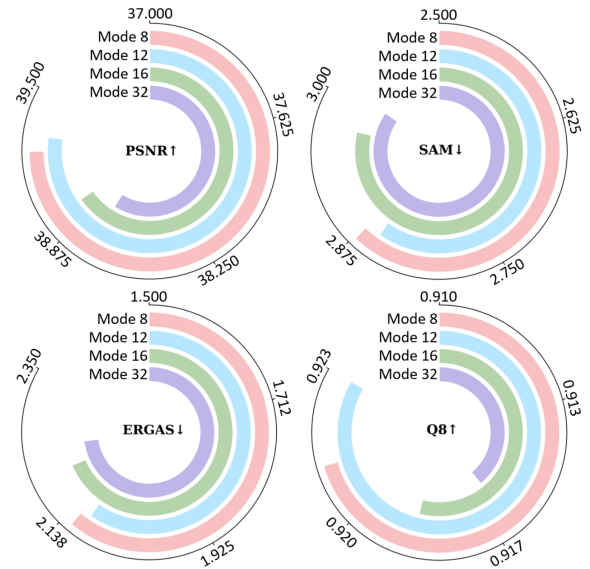


Figure 6: Ablation study on the spectral mode truncation.

datasets. Figure 4 provides the RGB visualizations and the corresponding Mean Absolute Error (MAE) residual maps compared to the ground truth for a WV3 example. Our

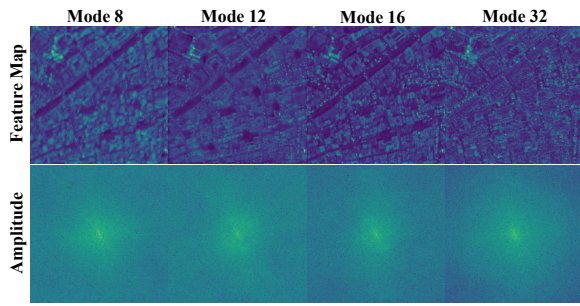


Figure 7: Features and Amplitude under mode truncation.

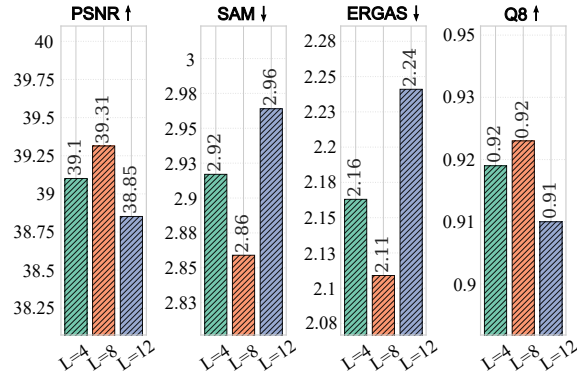


Figure 8: Ablation study on the FNO layer number.

model produces images with minimal aberrations, as evidenced by its sparse MAE residues. These results provide compelling evidence of the superiority of our method.

Evaluation on Full Resolution Scene. We evaluate the real-world applicability of NODiff using full resolution data. As shown in the right portion of Table 1, it consistently achieves leading performance across most metrics on the WV3 dataset, reaffirming the encouraging trends observed at reduced resolution assessment. In particular, it yields the highest HQNR score compared to both traditional and deep learning baselines, indicating superior fusion quality with respect to spectral fidelity and spatial detail. This is corroborated by the HQNR visualization shown in Figure 5, where our result exhibits a uniformly dark tone with sparse bright spots, indicating minimal spatial-spectral distortion. These full resolution results further substantiate the robust generalization capability of our method in real-world scenarios.

Ablation Study

Effect of the Spectral Mode Truncation. We first investigate the impact of mode truncation in spectral convolution. Specifically, we compare four configurations with mode values set to 8, 12, 16, and 32. As shown in Figure 6, our method with the mode=12 consistently achieves the best performance across all evaluation metrics. As visualized in Figure 7, the feature map of our method exhibits clear and finer textures, along with a more balanced frequency distribution. In comparison, higher spectral modes (*e.g.*, 16 and 32) introduce unwanted high-frequency noise, leading to noticeable

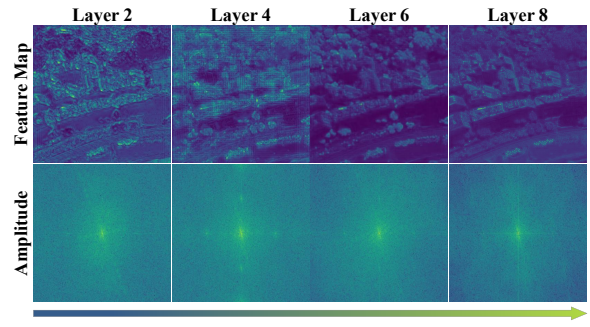


Figure 9: Features and Amplitude across FNO layers.

Method	PSNR	SAM	ERGAS	Q8
w/o Stage I	38.865	2.992	2.200	0.917
w/o CGA	38.681	3.028	2.356	0.911
w/o Ta-LoRA	<u>39.032</u>	<u>2.957</u>	<u>2.185</u>	<u>0.919</u>
NODiff (Ours)	39.314	2.859	2.109	0.923

Table 3: Effect of the training and fine-tuning designs.

artifacts in the feature maps; whereas a lower mode (8) results in overly smooth features and loss of important details.

Effect of the FNO Layer Number. We examine the effect of varying the number of FNO layers. As summarized in Figure 8, the model achieves optimal performance with 8 layers, and the feature maps become progressively clearer as depth increases (Figure 9), highlighting our model’s expressiveness. Further increasing the depth leads to a performance drop, possibly because the amplified high-frequency deviation caused by spectral truncation diminishes sensitivity to finer textures. Notably, the model with only 4 layers also performs competitively over other cutting-edge approaches.

Effect of the Fine-tuning Designs. We analyze the impact of different training paradigms and fine-tuning designs. As reported in Table 3, the two-stage pretraining-fine-tuning pipeline significantly outperforms the one-stage training scheme. Moreover, incorporating both CGA and Ta-LoRA yields the best performance across all metrics, demonstrating the effectiveness of each fine-tuning component.

Conclusion

We propose an efficient yet effective Neural Operator-based diffusion framework that seamlessly synergizes operator learning and generative modeling to learn high-resolution diffusion texture priors for pansharpening. To enable task-specific adaptation, we introduce a conditional detail guidance adapter for parameter-efficient fine-tuning of texture generation, complemented by a time-aware low-rank adaptation to remedy the potential negative effects of spectral mode truncation. Building on these ingredients, our method achieves superior performance, and offers a scalable solution for resource-efficient generative modeling.

Acknowledgements

The work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 62301151, in part by the Project of the Department of Science and Technology of Sichuan Province under Grant 2025YFNH0001.

References

- Arienzo, A.; Vivone, G.; Garzelli, A.; Alparone, L.; and Chanussot, J. 2022. Full-resolution quality assessment of pansharpening: Theoretical and hands-on approaches. *IEEE Geoscience and Remote Sensing Magazine*, 10(3): 168–201.
- Azizzadenesheli, K.; Kovachki, N.; Li, Z.; Liu-Schiaffini, M.; Kossaiji, J.; and Anandkumar, A. 2024. Neural operators for accelerating scientific simulations and design. *Nature Reviews Physics*, 6(5): 320–328.
- Cao, Z.; Cao, S.; Deng, L.-J.; Wu, X.; Hou, J.; and Vivone, G. 2024. Diffusion model with disentangled modulations for sharpening multispectral and hyperspectral images. *Information Fusion*, 104: 102158.
- Couairon, G.; Verbeek, J.; Schwenk, H.; and Cord, M. 2022. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*.
- Deng, L.-J.; Vivone, G.; Jin, C.; and Chanussot, J. 2020. Detail injection-based deep convolutional neural networks for pansharpening. *IEEE Transactions on Geoscience and Remote Sensing*, 59(8): 6995–7010.
- Deng, L.-J.; Vivone, G.; Paoletti, M. E.; Scarpa, G.; He, J.; Zhang, Y.; Chanussot, J.; and Plaza, A. 2022. Machine learning in pansharpening: A benchmark, from shallow to deep networks. *IEEE Geoscience and Remote Sensing Magazine*, 10(3): 279–315.
- Dian, R.; Li, S.; Sun, B.; and Guo, A. 2021. Recent advances and new guidelines on hyperspectral and multispectral image fusion. *Information Fusion*, 69: 40–51.
- Gu, S.; Chen, D.; Bao, J.; Wen, F.; Zhang, B.; Chen, D.; Yuan, L.; and Guo, B. 2022. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10696–10706.
- He, L.; Fang, Z.; Li, J.; Ye, H.; and Plaza, A. 2025a. Arbitrary-Resolution Hyperspectral Pansharpening Neural Operators. *IEEE Transactions on Geoscience and Remote Sensing*, 63: 1–19.
- He, X.; Cao, K.; Zhang, J.; Yan, K.; Wang, Y.; Li, R.; Xie, C.; Hong, D.; and Zhou, M. 2025b. Pan-mamba: Effective pan-sharpening with state space model. *Information Fusion*, 115: 102779.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Hou, J.; Cao, Z.; Zheng, N.; Li, X.; Chen, X.; Liu, X.; Cong, X.; Hong, D.; and Zhou, M. 2024. Linearly-evolved Transformer for Pan-sharpening. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 1486–1494.
- Hou, J.; Liu, X.; Wu, C.; Cong, X.; Huang, C.; Deng, L.-J.; and You, J. W. 2025. Bidomain uncertainty gated recursive network for pan-sharpening. *Information Fusion*, 118: 102938.
- Hu, J.-F.; Huang, T.-Z.; Deng, L.-J.; Jiang, T.-X.; Vivone, G.; and Chanussot, J. 2021. Hyperspectral image super-resolution via deep spatio-spectral attention convolutional neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12): 7251–7265.
- Huang, J.; Chen, H.; Ren, J.; Peng, S.; and Deng, L. 2025. A General Adaptive Dual-level Weighting Mechanism for Remote Sensing Pansharpening. *arXiv preprint arXiv:2503.13214*.
- Jin, Z.-R.; Zhang, T.-J.; Jiang, T.-X.; Vivone, G.; and Deng, L.-J. 2022. LAGConv: Local-Context Adaptive Convolution Kernels with Global Harmonic Bias for Pansharpening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 1113–1121. AAAI Press.
- Kawar, B.; Zada, S.; Lang, O.; Tov, O.; Chang, H.; Dekel, T.; Mosseri, I.; and Irani, M. 2023. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 6007–6017.
- Kim, S.; Do, J.; Lee, J.; and Kim, M. 2025. U-Know-DiffPAN: An Uncertainty-aware Knowledge Distillation Diffusion Framework with Details Enhancement for PAN-Sharpening. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 23069–23079.
- Kovachki, N.; Li, Z.; Liu, B.; Azizzadenesheli, K.; Bhat-tacharya, K.; Stuart, A.; and Anandkumar, A. 2023. Neural operator: Learning maps between function spaces with applications to pdes. *Journal of Machine Learning Research*, 24(89): 1–97.
- Li, Z.; Kovachki, N. B.; Azizzadenesheli, K.; Liu, B.; Bhat-tacharya, K.; Stuart, A.; and Anandkumar, A. 2021. Fourier Neural Operator for Parametric Partial Differential Equations. In *International Conference on Learning Representations*.
- Liu, X.; and Tang, H. 2025. DiffFNO: Diffusion Fourier Neural Operator. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 150–160.
- Liu, Y.; Dian, R.; and Li, S. 2025. Low-rank transformer for high-resolution hyperspectral computational imaging. *International Journal of Computer Vision*, 133(2): 809–824.
- Lolli, S.; Alparone, L.; Garzelli, A.; and Vivone, G. 2017. Haze correction for contrast-based multispectral pansharpening. *IEEE Geoscience and Remote Sensing Letters*, 14(12): 2255–2259.
- Lu, L.; Jin, P.; Pang, G.; Zhang, Z.; and Karniadakis, G. E. 2021. Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators. *Nature machine intelligence*, 3(3): 218–229.
- Luo, X.; Qian, X.; and Yoon, B.-J. 2024. Hierarchical Neural Operator Transformer with Learnable Frequency-aware Loss Prior for Arbitrary-scale Super-resolution. In *Proceedings of the 41st International Conference on Machine Learning*.

- ing, volume 235 of *Proceedings of Machine Learning Research*, 33466–33485. PMLR.
- Meng, Q.; Shi, W.; Li, S.; and Zhang, L. 2023. PanDiff: A novel pansharpening method based on denoising diffusion probabilistic model. *IEEE Transactions on Geoscience and Remote Sensing*, 61: 1–17.
- Meng, X.; Shen, H.; Li, H.; Zhang, L.; and Fu, R. 2019. Review of the pansharpening methods for remote sensing images based on the idea of meta-analysis: Practical discussion and challenges. *Information Fusion*, 46: 102–113.
- Saharia, C.; Chan, W.; Saxena, S.; Li, L.; Whang, J.; Denton, E. L.; Ghasemipour, K.; Gontijo Lopes, R.; Karagol Ayan, B.; Salimans, T.; et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35: 36479–36494.
- Shi, Y.; Xue, C.; Liew, J. H.; Pan, J.; Yan, H.; Zhang, W.; Tan, V. Y.; and Bai, S. 2024. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8839–8849.
- Song, J.; Meng, C.; and Ermon, S. 2021. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*.
- Tan, J.; Huang, J.; Zheng, N.; Zhou, M.; Yan, K.; Hong, D.; and Zhao, F. 2024. Revisiting Spatial-Frequency Information Integration from a Hierarchical Perspective for Panchromatic and Multi-Spectral Image Fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 25922–25931.
- Tran, A.; Mathews, A.; Xie, L.; and Ong, C. S. 2023. Factorized Fourier Neural Operators. In *The Eleventh International Conference on Learning Representations*.
- Vivone, G. 2019. Robust band-dependent spatial-detail approaches for panchromatic sharpening. *IEEE transactions on Geoscience and Remote Sensing*, 57(9): 6421–6433.
- Vivone, G.; Deng, L.-J.; Deng, S.; Hong, D.; Jiang, M.; Li, C.; Li, W.; Shen, H.; Wu, X.; Xiao, J.-L.; et al. 2024. Deep Learning in Remote Sensing Image Fusion: Methods, protocols, data, and future perspectives. *IEEE Geoscience and Remote Sensing Magazine*.
- Wald, L. 2002. *Data fusion: definitions and architectures: fusion of images of different spatial resolutions*. Presses des MINES.
- Wei, M.; and Zhang, X. 2023. Super-resolution neural operator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18247–18256.
- Wu, X.; Cao, Z.-H.; Huang, T.-Z.; Deng, L.-J.; Chanussot, J.; and Vivone, G. 2025. Fully-Connected Transformer for Multi-Source Image Fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(3): 2071–2088.
- Wu, Z.-C.; Huang, T.-Z.; Deng, L.-J.; Huang, J.; Chanussot, J.; and Vivone, G. 2023. LRTCFFpan: Low-rank tensor completion based framework for pansharpening. *IEEE Transactions on Image Processing*, 32: 1640–1655.
- Xu, Z.; Tang, Y.; Xu, B.; and Li, Q. 2025. NeurOp-Diff: Continuous Remote Sensing Image Super-Resolution via Neural Operator Diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 12491–12501.
- Yang, J.; Fu, X.; Hu, Y.; Huang, Y.; Ding, X.; and Paisley, J. 2017. PanNet: A deep network architecture for pansharpening. In *Proceedings of the IEEE international conference on computer vision*, 5449–5457.
- Yuhas, R. H.; Goetz, A. F.; and Boardman, J. W. 1992. Discrimination among semi-arid landscape endmembers using the spectral angle mapper (SAM) algorithm. In *JPL, Summaries of the Third Annual JPL Airborne Geoscience Workshop. Volume 1: AVIRIS Workshop*.
- Zhou, J.; Civco, D. L.; and Silander, J. A. 1998. A wavelet transform method to merge Landsat TM and SPOT panchromatic data. *International journal of remote sensing*, 19(4): 743–757.
- Zhou, M.; Huang, J.; Fang, Y.; Fu, X.; and Liu, A. 2022. Pan-sharpening with customized transformer and invertible neural network. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, 3553–3561.
- Zhou, M.; Huang, J.; Guo, C.-L.; and Li, C. 2023. Fourmer: An efficient global modeling paradigm for image restoration. In *International conference on machine learning*, 42589–42601. PMLR.
- Zhou, M.; Zheng, N.; He, X.; Hong, D.; and Chanussot, J. 2024. Probing synergistic high-order interaction for multimodal image fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.