

# Spatial-Spectral Residuals Informed Diffusion Neural Operator for Pan-sharpening

## Supplementary Material

We organize the supplementary material as follows:

**Structure of the Auxiliary Predictor:** We introduce the structure of the auxiliary predictor in the main paper.

**Plug and Play Validation:** We proceed to validate the proposed spatial-spectral residual integration mechanism by applying it to a representative diffusion model.

**Additional Experimental Results:** We present additional experimental results on three datasets, including quantitative evaluations, qualitative comparisons, and feature maps.

### S1. Structure of the Auxiliary Predictor

As illustrated in Figure A1, the auxiliary predictor (corresponding to Figure 3 in our main paper) consists of 4 convolutional residual blocks. It is designed to extract the cross-modality feature  $\mathbf{V}$  and generate an initial HRMS estimate, with the latter serving as a replacement for the high-resolution reference  $\mathbf{H}$  in Stage I. Specifically, the auxiliary predictor takes as input the cross-modality PAN and LRMS images together with their spatial residual. To obtain this residual, we first transform the LRMS image  $\mathbf{L}$  into a single channel pseudo PAN image  $\hat{\mathbf{P}}$ , and then subtract it from the real PAN image  $\mathbf{P}$ , supplying the spatial details absent in  $\mathbf{L}$ .

### S2. Plug and Play Validation

We apply our proposed spatial-spectral residual integration strategy to another representative diffusion model [1] to validate its efficacy. As shown in Figure A2, integrating both spatial and spectral residuals leads to a notable improvement in PSNR over the original model. More importantly, this integration introduces only negligible overhead in terms of FLOPs, Memory, and Inference Time. These results strongly affirm the effectiveness and general applicability of the proposed residual integration approach.

### S3. Additional Experimental Results

#### S3.1. Evaluation on Reduced-resolution Scene

We begin by evaluating our model on the reduced resolution GF2 and QB datasets. As summarized in the left panels of Table 1 and Table 2, our method consistently outperforms state-of-the-art approaches across all metrics for both the GF2 and QB datasets. In particular, our method achieves the optimal results in PSNR and SAM on the GF2 and QB datasets, indicating superior spatial and spectral fusion quality. Overall, these quantitative results confirm our model’s superiority in preserving spectral fidelity and reconstructing spatial textures.

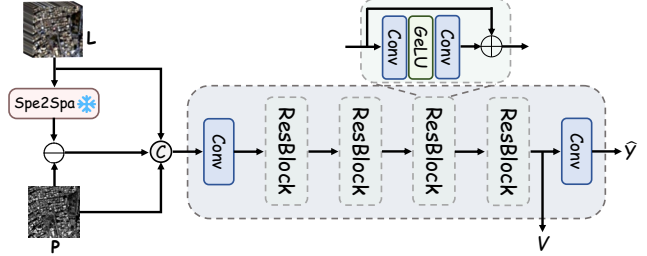


Figure A1. Structure of the auxiliary predictor in the main paper.

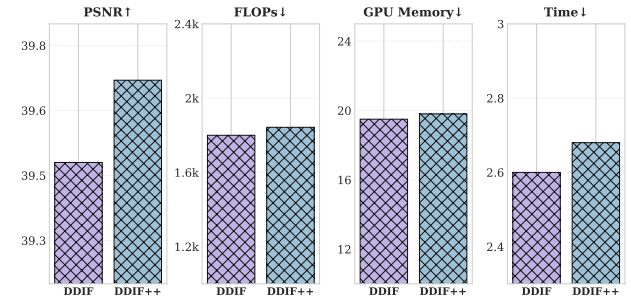


Figure A2. Ablation experiment evaluating the efficacy of our spatial-spectral residual integration strategy applied to DDIF [1]. “++” indicates the model integrating the spatial-spectral residuals.

Visual comparisons further support these quantitative findings. As shown in Figure A4, the images generated by our model exhibit fewer visual artifacts and closer resemblance to the ground truth. This is reinforced by the mean absolute error (MAE) maps between the fused results and the ground truth, where enlarged residue regions contain notably fewer bright spots. Together, these visual observations align with and strengthen the quantitative conclusions.

We further visualize the feature maps across different denoiser layers using samples from three benchmark datasets. It is clear that the feature maps demonstrate a progressive enhancement in clarity and detail as the number of denoiser layers increases, as illustrated in Figure A3.

#### S3.2. Evaluation on Full-resolution Scene

We further evaluate our method under full resolution scenarios to verify its practical applicability and generalization ability. The results in the right panels of Table 1 and Table 2 show that our model maintains competitive performance on both the GF2 and QB datasets, mirroring the positive trends observed at reduced resolution. Notably, it achieves the highest HQNR scores on both datasets, outperforming

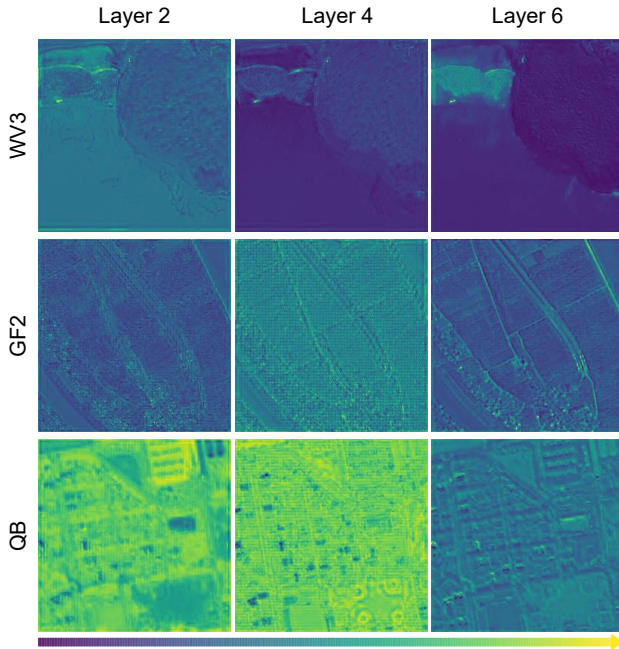


Figure A3. Feature maps across different denoiser layers.

traditional and deep-learning baselines, which indicates the superior spatial and spectral fidelity. These full resolution results highlight the robust generalization of our model in real-world applications.

Figure A5 presents visual comparisons on full resolution WV3, GF2 and QB samples. In the odd rows, our method produces images with finer spatial textures and visually consistent spectral characteristics. The even rows demonstrate the corresponding HQNR maps, where our method exhibits significantly fewer bright spots, consistent with its higher HQNR scores and indicating better spatial-spectral fidelity. These visual observations corroborate the quantitative results, jointly demonstrating the effectiveness of our approach in real-world high-resolution scenarios.

## References

- [1] Zihan Cao, Shiqi Cao, Liang-Jian Deng, Xiao Wu, Junming Hou, and Gemine Vivone. Diffusion model with disentangled modulations for sharpening multispectral and hyper-spectral images. *Information Fusion*, 104:102158, 2024. 1
- [2] Liang-Jian Deng, Gemine Vivone, Cheng Jin, and Jocelyn Chanussot. Detail injection-based deep convolutional neural networks for pansharpening. *IEEE Transactions on Geoscience and Remote Sensing*, 59(8):6995–7010, 2020. 3, 4
- [3] Junming Hou, Zihan Cao, Naishan Zheng, Xuan Li, Xiaoyu Chen, Xinyang Liu, Xiaofeng Cong, Danfeng Hong, and Man Zhou. Linearly-evolved transformer for pansharpening. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 1486–1494, 2024. 3, 4
- [4] Jie Huang, Haorui Chen, Jiakuan Ren, Siran Peng, and Liangjian Deng. A general adaptive dual-level weighting mechanism for remote sensing pansharpening. *arXiv preprint arXiv:2503.13214*, 2025. 3, 4
- [5] Zhi-Rui Jin, Tian-Jing Zhang, Tai-Xiang Jiang, Gemine Vivone, and Liang-Jian Deng. Lagconv: Local-context adaptive convolution kernels with global harmonic bias for pansharpening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1113–1121. AAAI Press, 2022. 3, 4
- [6] Simone Lolli, Luciano Alparone, Andrea Garzelli, and Gemine Vivone. Haze correction for contrast-based multispectral pansharpening. *IEEE Geoscience and Remote Sensing Letters*, 14(12):2255–2259, 2017. 3, 4
- [7] Qingyan Meng, Wenxu Shi, Sijia Li, and Linlin Zhang. Pandiff: A novel pansharpening method based on denoising diffusion probabilistic model. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–17, 2023. 3, 4
- [8] Jiangtong Tan, Jie Huang, Naishan Zheng, Man Zhou, Keyu Yan, Danfeng Hong, and Feng Zhao. 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*, pages 25922–25931, 2024. 3, 4
- [9] Gemine Vivone. Robust band-dependent spatial-detail approaches for panchromatic sharpening. *IEEE transactions on Geoscience and Remote Sensing*, 57(9):6421–6433, 2019. 3, 4
- [10] Zhong-Cheng Wu, Ting-Zhu Huang, Liang-Jian Deng, Jie Huang, Jocelyn Chanussot, and Gemine Vivone. Lrtcfpan: Low-rank tensor completion based framework for pansharpening. *IEEE Transactions on Image Processing*, 32:1640–1655, 2023. 3, 4
- [11] Man Zhou, Jie Huang, Chun-Le Guo, and Chongyi Li. Fourmer: An efficient global modeling paradigm for image restoration. In *International conference on machine learning*, pages 42589–42601. PMLR, 2023. 3, 4
- [12] Man Zhou, Naishan Zheng, Xuanhua He, Danfeng Hong, and Jocelyn Chanussot. Probing synergistic high-order interaction for multi-modal image fusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 3, 4

| Method         | Reduced Resolution                 |                                   |                                   |                                   | Full Resolution                   |                                   |                                   |
|----------------|------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
|                | PSNR( $\pm$ std)                   | SAM( $\pm$ std)                   | ERGAS( $\pm$ std)                 | $Q2^n$ ( $\pm$ std)               | $D_\lambda$ ( $\pm$ std)          | $D_s$ ( $\pm$ std)                | HQNR( $\pm$ std)                  |
| BDS-PC [9]     | 35.180 $\pm$ 2.317                 | 1.681 $\pm$ 0.360                 | 1.667 $\pm$ 0.445                 | 0.892 $\pm$ 0.035                 | 0.076 $\pm$ 0.030                 | 0.155 $\pm$ 0.028                 | 0.781 $\pm$ 0.041                 |
| BT-H [6]       | 36.054 $\pm$ 2.236                 | 1.649 $\pm$ 0.360                 | 1.528 $\pm$ 0.409                 | 0.918 $\pm$ 0.025                 | 0.060 $\pm$ 0.025                 | 0.131 $\pm$ 0.019                 | 0.817 $\pm$ 0.031                 |
| LRTCFFPan [10] | 37.599 $\pm$ 2.331                 | 1.298 $\pm$ 0.312                 | 1.272 $\pm$ 0.343                 | 0.935 $\pm$ 0.030                 | 0.033 $\pm$ 0.027                 | 0.090 $\pm$ 0.014                 | 0.881 $\pm$ 0.023                 |
| FusionNet [2]  | 39.639 $\pm$ 2.270                 | 0.974 $\pm$ 0.212                 | 0.988 $\pm$ 0.222                 | 0.964 $\pm$ 0.009                 | 0.035 $\pm$ 0.012                 | 0.101 $\pm$ 0.013                 | 0.867 $\pm$ 0.018                 |
| LAGConv [5]    | 42.735 $\pm$ 1.447                 | 0.786 $\pm$ 0.148                 | 0.687 $\pm$ 0.113                 | 0.980 $\pm$ 0.009                 | 0.028 $\pm$ 0.013                 | 0.079 $\pm$ 0.014                 | 0.895 $\pm$ 0.020                 |
| Fourmer [11]   | 40.670 $\pm$ 1.903                 | 0.976 $\pm$ 0.209                 | 0.885 $\pm$ 0.185                 | 0.970 $\pm$ 0.011                 | 0.047 $\pm$ 0.039                 | <b>0.038<math>\pm</math>0.010</b> | 0.917 $\pm$ 0.035                 |
| HFIN [8]       | 42.189 $\pm$ 1.752                 | 0.843 $\pm$ 0.148                 | 0.735 $\pm$ 0.126                 | 0.977 $\pm$ 0.011                 | 0.027 $\pm$ 0.020                 | 0.062 $\pm$ 0.009                 | 0.912 $\pm$ 0.018                 |
| HOIF [12]      | 40.982 $\pm$ 1.802                 | 0.943 $\pm$ 0.205                 | 0.841 $\pm$ 0.162                 | 0.974 $\pm$ 0.009                 | 0.029 $\pm$ 0.015                 | 0.051 $\pm$ 0.011                 | 0.922 $\pm$ 0.019                 |
| PanDiff [7]    | 42.326 $\pm$ 1.635                 | 0.875 $\pm$ 0.133                 | 0.727 $\pm$ 0.115                 | 0.981 $\pm$ 0.008                 | 0.028 $\pm$ 0.020                 | 0.073 $\pm$ 0.010                 | 0.902 $\pm$ 0.021                 |
| LFormer [3]    | 44.196 $\pm$ 1.800                 | <u>0.648<math>\pm</math>0.130</u> | <u>0.578<math>\pm</math>0.112</u> | 0.985 $\pm$ 0.007                 | <u>0.021<math>\pm</math>0.010</u> | 0.050 $\pm$ 0.001                 | <u>0.930<math>\pm</math>0.013</u> |
| ADWM [4]       | 43.884 $\pm$ 1.714                 | 0.672 $\pm$ 0.130                 | 0.597 $\pm$ 0.107                 | <u>0.985<math>\pm</math>0.006</u> | 0.022 $\pm$ 0.012                 | 0.052 $\pm$ 0.011                 | 0.928 $\pm$ 0.014                 |
| SRINO (ours)   | <b>44.228<math>\pm</math>1.709</b> | <b>0.646<math>\pm</math>0.129</b> | <b>0.573<math>\pm</math>0.101</b> | <b>0.985<math>\pm</math>0.006</b> | <b>0.019<math>\pm</math>0.010</b> | <u>0.044<math>\pm</math>0.008</u> | <b>0.938<math>\pm</math>0.011</b> |
| Ideal value    | $\infty$                           | <b>0</b>                          | <b>0</b>                          | <b>1</b>                          | <b>0</b>                          | <b>0</b>                          | <b>1</b>                          |

Table 1. Quantitative results for reduced and full resolution GF2 samples, comparing several representative state-of-the-art methods. Bold: Best; Underline: Second best.

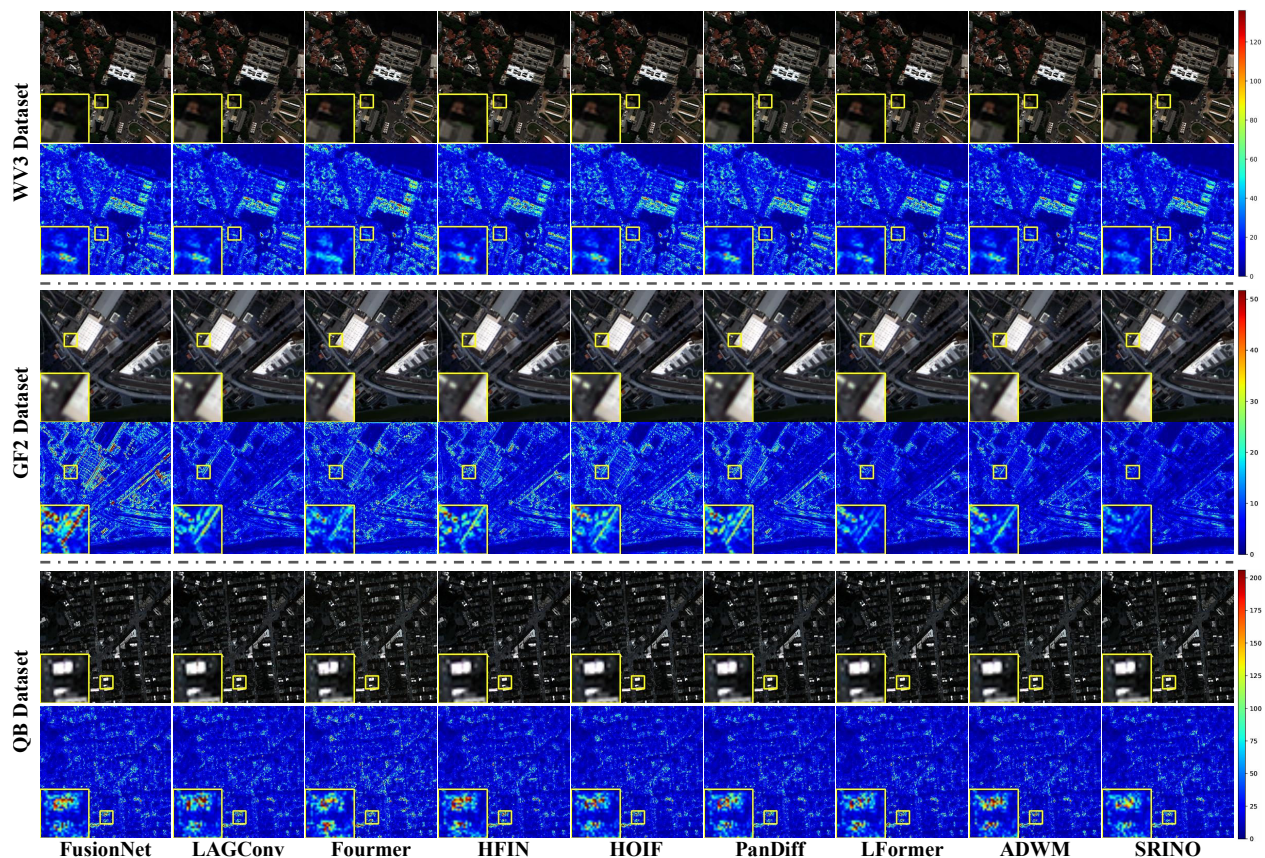


Figure A4. The visual results (odd rows) and the corresponding mean absolute error (MAE) maps (even rows) of all compared DL-based methods on reduced resolution samples from the WV3, GF2 and QB datasets, respectively.

| Method        | Reduced Resolution                 |                                   |                                   |                                   | Full Resolution                   |                                   |                                   |
|---------------|------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
|               | PSNR( $\pm$ std)                   | SAM( $\pm$ std)                   | ERGAS( $\pm$ std)                 | $Q2^n$ ( $\pm$ std)               | $D_\lambda$ ( $\pm$ std)          | $D_s$ ( $\pm$ std)                | HQNR( $\pm$ std)                  |
| BDS-PC [9]    | 32.547 $\pm$ 3.202                 | 8.089 $\pm$ 1.980                 | 7.515 $\pm$ 0.800                 | 0.831 $\pm$ 0.090                 | 0.198 $\pm$ 0.033                 | 0.164 $\pm$ 0.048                 | 0.672 $\pm$ 0.058                 |
| BT-H [6]      | 32.648 $\pm$ 3.273                 | 7.194 $\pm$ 1.552                 | 7.401 $\pm$ 0.838                 | 0.833 $\pm$ 0.088                 | 0.230 $\pm$ 0.072                 | 0.165 $\pm$ 0.017                 | 0.643 $\pm$ 0.065                 |
| LRTCFFan [10] | 33.260 $\pm$ 3.272                 | 7.187 $\pm$ 1.711                 | 6.928 $\pm$ 0.812                 | 0.855 $\pm$ 0.087                 | <b>0.023<math>\pm</math>0.012</b> | 0.071 $\pm$ 0.035                 | 0.909 $\pm$ 0.044                 |
| FusionNet [2] | 37.532 $\pm$ 2.518                 | 4.923 $\pm$ 0.908                 | 4.159 $\pm$ 0.321                 | 0.925 $\pm$ 0.090                 | 0.057 $\pm$ 0.018                 | 0.052 $\pm$ 0.009                 | 0.894 $\pm$ 0.021                 |
| LAGConv [5]   | 38.181 $\pm$ 2.456                 | 4.547 $\pm$ 0.830                 | 3.826 $\pm$ 0.420                 | 0.934 $\pm$ 0.088                 | 0.086 $\pm$ 0.024                 | 0.068 $\pm$ 0.014                 | 0.852 $\pm$ 0.018                 |
| Fourmer [11]  | 36.797 $\pm$ 2.597                 | 5.079 $\pm$ 1.009                 | 4.494 $\pm$ 0.410                 | 0.924 $\pm$ 0.080                 | <u>0.033<math>\pm</math>0.013</u> | 0.050 $\pm$ 0.010                 | <u>0.920<math>\pm</math>0.018</u> |
| HFIN [8]      | 38.247 $\pm$ 2.403                 | 4.542 $\pm$ 0.805                 | 3.813 $\pm$ 0.322                 | 0.934 $\pm$ 0.085                 | 0.067 $\pm$ 0.025                 | 0.078 $\pm$ 0.019                 | 0.860 $\pm$ 0.018                 |
| HOIF [12]     | 38.242 $\pm$ 2.119                 | 4.521 $\pm$ 0.811                 | 3.825 $\pm$ 0.510                 | 0.933 $\pm$ 0.094                 | 0.077 $\pm$ 0.026                 | 0.059 $\pm$ 0.025                 | 0.869 $\pm$ 0.045                 |
| PanDiff [7]   | 38.538 $\pm$ 2.401                 | 4.524 $\pm$ 0.792                 | 3.688 $\pm$ 0.343                 | 0.937 $\pm$ 0.084                 | 0.058 $\pm$ 0.022                 | 0.064 $\pm$ 0.025                 | 0.882 $\pm$ 0.041                 |
| LFormer [3]   | 38.674 $\pm$ 2.370                 | <u>4.382<math>\pm</math>0.789</u> | <u>3.616<math>\pm</math>0.314</u> | <u>0.939<math>\pm</math>0.082</u> | 0.037 $\pm$ 0.016                 | 0.089 $\pm$ 0.026                 | 0.877 $\pm$ 0.021                 |
| ADWM [4]      | 38.466 $\pm$ 2.420                 | 4.450 $\pm$ 0.809                 | 3.705 $\pm$ 0.346                 | 0.937 $\pm$ 0.085                 | 0.064 $\pm$ 0.025                 | <b>0.024<math>\pm</math>0.016</b> | 0.914 $\pm$ 0.036                 |
| SRINO (ours)  | <b>38.864<math>\pm</math>2.346</b> | <b>4.351<math>\pm</math>0.773</b> | <b>3.542<math>\pm</math>0.313</b> | <b>0.940<math>\pm</math>0.085</b> | 0.041 $\pm$ 0.013                 | <u>0.032<math>\pm</math>0.013</u> | <b>0.927<math>\pm</math>0.024</b> |
| Ideal value   | $\infty$                           | <b>0</b>                          | <b>0</b>                          | <b>1</b>                          | <b>0</b>                          | <b>0</b>                          | <b>1</b>                          |

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

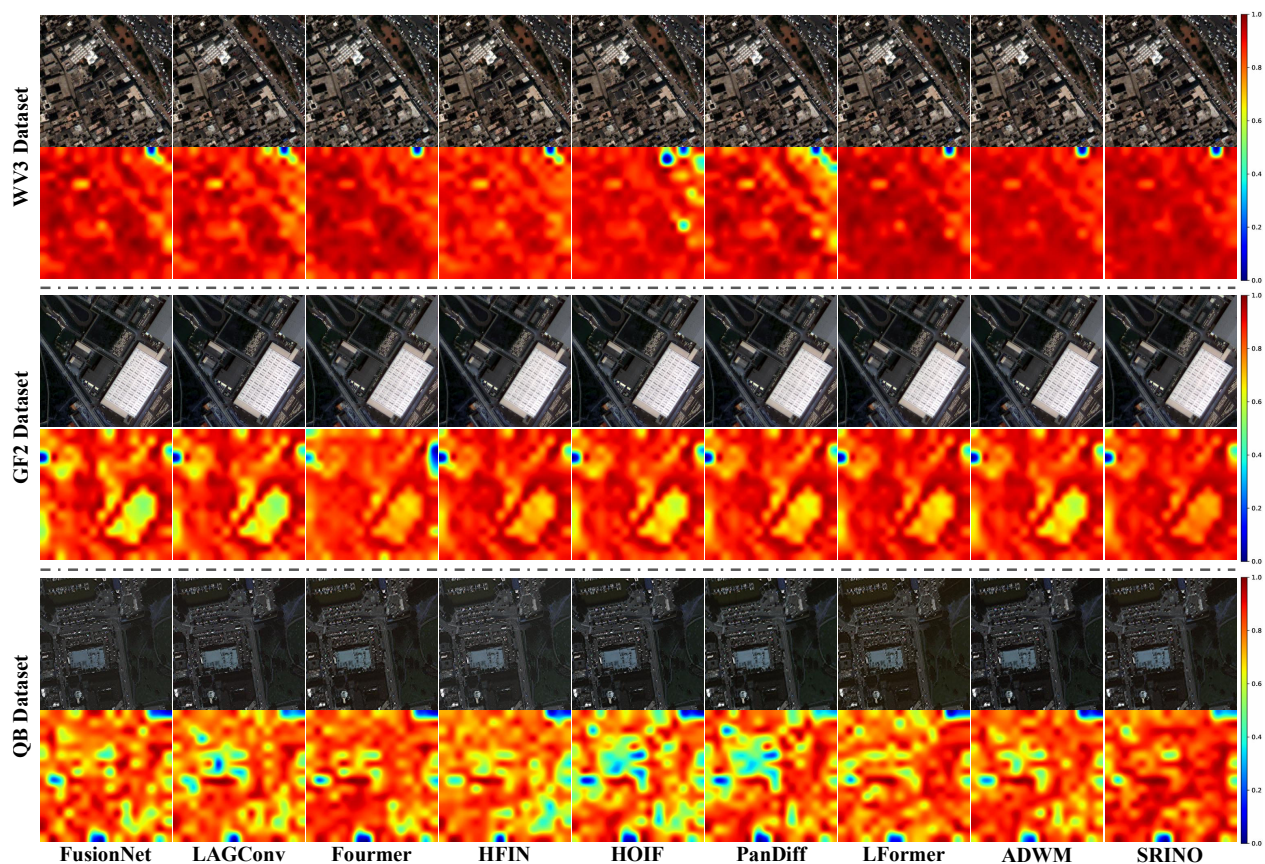


Figure A5. The visual results (odd rows) and the corresponding HQNR maps (even rows) of all compared DL-based methods on full-resolution samples from the WV3, GF2 and QB datasets, respectively.