

Yuanhao Cai<sup>1,\*</sup>, Jing Lin<sup>1,\*</sup>, Zudi Lin<sup>2</sup>, Haoqian Wang<sup>1,†</sup>, Yulun Zhang<sup>3</sup>, Hanspeter Pfister<sup>2</sup>, Radu Timofte<sup>3,4</sup>, Luc Van Gool<sup>3</sup> 1 Shenzhen International Graduate School, Tsinghua University 2 Harvard University 3 ETH Zürich 4 JMU Würzburg

## Introduction

Existing leading methods for spectral reconstruction (SR) focus on designing deeper or wider convolutional neural networks (CNNs) to learn the end-to-end mapping from the RGB image to its hyperspectral image (HSI). These CNNbased methods achieve impressive restoration performance while showing limitations in capturing the long-range dependencies and self-similarity prior. To cope with this problem, we propose a novel Transformer-based method, Multistage Spectral-wise Transformer (MST++), for efficient spectral reconstruction. In particular, we employ Spectral-wise Multi-head Self-attention (S-MSA) that is based on the HSI spatially sparse while spectrally self-similar nature to compose the basic unit, Spectral-wise Attention Block (SAB). Then SABs build up Singlestage Spectral-wise Transformer (SST) that exploits a U-shaped structure to extract multi-resolution contextual information. Finally, our MST++, cascaded by several SSTs, progressively improves the reconstruction quality from coarse to fine. Comprehensive experiments show that our MST++ significantly outperforms other state-of-the-art methods. Our MST++ is based on our CVPR 2022 work MST and has won the first place in NTIRE 2022 Spectral Recovery Challenge.

Code for MST : https://github.com/caiyuanhao1998/MST Code for MST++ : https://github.com/caiyuanhao1998/MST-plus-plus

# Method



In figure (a), MST++ is cascaded by  $N_s$  Single-stage Spectral-wise Transformers (SSTs). In figure (b), SST adopts a three-level U-shaped architecture. Details of Spectral-wise Attention Block (SAB), Feed-Forward Network (FFN), and Spectralwise Multi-head Self-Attention (S-MSA) are shown in figure (c), (d), and (e).

#### > Overall Architecture

# MST++: Multi-stage Spectral-wise Transformer for Effcient Spectral Reconstruction

#### Spectral-wise Multi-head Self-Attention

Suppose the input tokens of S-MSA as X that is projected into Q, K, V as

 $\mathbf{Q} = \mathbf{X}\mathbf{W}^{\mathbf{Q}}, \mathbf{K} = \mathbf{X}\mathbf{W}^{\mathbf{K}}, \mathbf{V} = \mathbf{X}\mathbf{W}^{\mathbf{V}}$ 

Subsequently, **Q**, **K**, **V** are split into *N* heads along the spectral dimension and the self-attention is calculated inside each  $head_i$  as

$$\mathbf{A}_{j} = \operatorname{softmax}(\sigma_{j}\mathbf{K}^{\mathrm{T}}_{j}\mathbf{Q}_{j}), head_{j} = \mathbf{V}_{j}\mathbf{A}_{j}$$

Then the outputs of *N* heads are aggregated by a linear projection and is added with a position embedding that is produced by function  $f_p(\cdot)$  as

 $\mathbf{S} - \mathbf{MSA}(\mathbf{X}) = (\text{Concate}_{j=1}^{N}(head_{j}))\mathbf{W} + f_{p}(\mathbf{V})$ 

 $f_p(\cdot)$  consists of two depth-wise conv 3 × 3, a GELU activation, and reshape operation. Finally, we reshape the above results to get the output feature maps

### Discussion with Previous MSA Modules



We compare the computational complexity of our S-MSA, global spatial-wise MSA (G-MSA), and local window-based MSA (W-MSA) as

$$O(G - MSA) = 2(HW)^2C, \qquad O(W - MSA) = \frac{HW}{M^2}(2(M^2)^2C) = 2M^2HWC,$$
$$O(S - MSA) = N\left(\left(\frac{C}{N}\right)^2 HW + \left(\frac{C}{N}\right)^2 HW\right)^2 = \frac{2HWC^2}{N}$$

$$A) = 2(HW)^2 C, \qquad O(W - MSA) = \frac{HW}{M^2} (2(M^2)^2 C) = 2M^2 HWC,$$
$$O(S - MSA) = N \left( \left(\frac{C}{N}\right)^2 HW + \left(\frac{C}{N}\right)^2 HW \right)^2 = \frac{2HWC^2}{N}$$

### Ensemble Strategies

We adopt self ensemble, multi-scale ensemble, and top-k multi-model ensemble when testing on the test-challenge dataset. The top-k multi-model ensemble averages the outpuits of MIRNet, MPRet, Restormer, HINet, MST, and MST++.

# Experiment

#### Quantitative Results



### > Qualitative Results



# Conclusion

our contributions:

- The first Transformer MST++
- A novel S-MSA
- SOTA results
- A Strong Baseline



	NTIRE 2022 HSI Dataset - Valid						NTIRE 2022 HSI Dataset - Test		
	Params (M)	FLOPS (G)	MRAE	RMSE	PSNR	Username	MRAE	RMSE	
67]     ] 4]	4.65 31.70 2.42 4.04 2.66 5.21 3.75	304.45 163.81 158.32 270.61 173.81 31.04 42.95 93.77	0.3814 0.3476 0.3277 0.2500 0.2048 0.2032 0.1890 0.1833	0.0588 0.0550 0.0437 0.0367 0.0317 0.0303 0.0274 0.0274	26.36 26.89 28.29 31.22 32.13 32.51 33.29 33.40	pipixia uslab orange_dog askldklasfj HSHAJii ptdoge_hot test_pseudo gkdgkd	0.2434 0.2377 0.2377 0.2345 0.2308 0.2107 0.2036 0.1935	0.0411 0.0391 0.0376 0.0361 0.0364 0.0365 0.0324 0.0322	
5] ]	3.62 2.45 <b>1.62</b>	101.59 32.07 <b>23.05</b>	0.1833 0.1817 0.1772 0.1645	0.0274 0.0270 0.0256 0.0248	33.40 33.50 33.90 <b>34.32</b>	deeppf mialgo_ls	0.1935 0.1767 0.1247 0.1131	0.0322 0.0322 0.0257 0.0231	

Code and models are publicly available at



