

# Cross-Spatial Pixel Integration and Cross-Stage Feature Fusion Based Transformer Network for Remote Sensing Image Super-Resolution

Yuting Lu, Lingtong Min, Binglu Wang<sup>†</sup>, *Member, IEEE*, Le Zheng, *Senior Member, IEEE*, Xiaoxu Wang, *Member, IEEE*, Yongqiang Zhao, *Member, IEEE*, and Teng Long, *Fellow, IEEE*

**Abstract**—Remote sensing image super-resolution (RSISR) plays a vital role in enhancing spatial details and improving the quality of satellite imagery. Recently, Transformer-based models have shown competitive performance in RSISR. To mitigate the quadratic computational complexity resulting from global self-attention, various methods constrain attention to a local window, enhancing its efficiency. Consequently, the receptive fields in a single attention layer are inadequate, leading to insufficient context modeling. Furthermore, while most transform-based approaches reuse shallow features through skip connections, relying solely on these connections treats shallow and deep features equally, impeding the model’s ability to characterize them. To address these issues, we propose a novel transformer architecture called Cross-Spatial Pixel Integration and Cross-Stage Feature Fusion Based Transformer Network (SPIFFNet) for RSISR. Our proposed model effectively enhances global cognition and understanding of the entire image, facilitating efficient integration of features cross-stages. The model incorporates cross-spatial pixel integration attention (CSPIA) to introduce contextual information into a local window, while cross-stage feature fusion attention (CSFFA) adaptively fuses features from the previous stage to improve feature expression in line with the requirements of the current stage. We conducted comprehensive experiments on multiple benchmark datasets, demonstrating the superior performance of our proposed SPIFFNet in terms of both quantitative metrics and visual quality when compared to state-of-the-art methods.

**Index Terms**—remote sensing image super-resolution, transformer network, cross-spatial pixel integration, cross-stage feature fusion

## I. INTRODUCTION

REMOTE sensing imaging technology is of paramount importance in numerous fields, including environmental monitoring [1]–[3], disaster management [4], [5], urban planning [6], [7], and object detection [8]–[10]. Therefore,

Yuting Lu, Xiaoxu Wang and Yongqiang Zhao are with School of Automation, Northwestern Polytechnical University, Xi’an 710072, China (e-mail:lyt1996@mail.nwpu.edu.cn, woayofly1982@nwpu.edu.cn, zhaoyq@nwpu.edu.cn). Lingtong Min is with School of Electronics and Information, Northwestern Polytechnical University, Xi’an 710072, China (e-mail:minlingtong@nwpu.edu.cn).

Binglu Wang, Le Zheng and Teng Long are with the Radar Research Laboratory, School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China (e-mail: wbl921129@gmail.com, le.zheng.cn@gmail.com, longteng@bit.edu.cn).

<sup>†</sup>Corresponding author: Binglu Wang.

This work is supported by the Postdoctoral Science Foundation of China under Grant 2022M710393, the Fourth Special Grant of China Postdoctoral Science Foundation (in front of the station) 2022TQ0035 and the Shaanxi Science Fund for Distinguished Young Scholars 2022JC-49.

the acquisition of high-resolution remote sensing images is imperative for the effective implementation and analysis of remote sensing image applications. However, challenges arise due to factors such as sensor noise, optical distortion and environmental interference, which can significantly degrade the image quality. Image super-resolution (SR) is a typical computer vision task that involves reconstructing high-resolution (HR) images from low-resolution (LR) images. The primary objective of SR is to mitigate the detrimental impact of acquisition equipment and environmental factors on remote sensing imaging outcomes, thereby enhancing the resolution of remote sensing images. As an alternative to developing physical imaging technologies, SR has gained significant attention in recent years for its ability to effectively generate high-resolution remote sensing images [11].

Traditional RSISR methods often rely on interpolation-based techniques, such as bicubic interpolation [12] or Lanczos interpolation. While these methods are simple, they may yield limited performance due to their inability to capture high-frequency details and structural information in the generated images. Recent advancements in deep learning [13] have led to the emergence of convolutional neural networks (CNNs) as powerful tools for various image processing tasks, including RSISR. CNN-based methods [14]–[20] have demonstrated promising results in learning complex representations from large datasets. However, despite the success of CNNs, they still possess certain limitations when employed for RSISR. CNNs typically operate locally with fixed receptive fields, which may hinder their ability to effectively capture long-range dependencies. As a result, they may have limited modeling capacity for remote sensing scenes with large spatial extent [21].

Transformer-based architectures, initially developed for natural language processing tasks, have emerged as a suitable solution for addressing the limitations of CNNs and have demonstrated impressive performance in diverse computer vision tasks, including image classification [22]–[24] and object detection [25]–[27]. The pioneering vision transformer model [22] employs a redundant attention mechanism, resulting in quadratic computation complexity relative to the image size. This high computational complexity poses challenges for its application in high-resolution predictions for RSISR tasks. To mitigate this issue, recent proposals have explored the use of self-attention within small spatial regions [28]–[30]. However, since remote sensing images usually cover a large range of

arXiv:2307.02974v1 [cs.CV] 6 Jul 2023



Fig. 1. Examples of similar windows in remote sensing images. These windows are similar in texture, shape and semantic features but far apart in space.

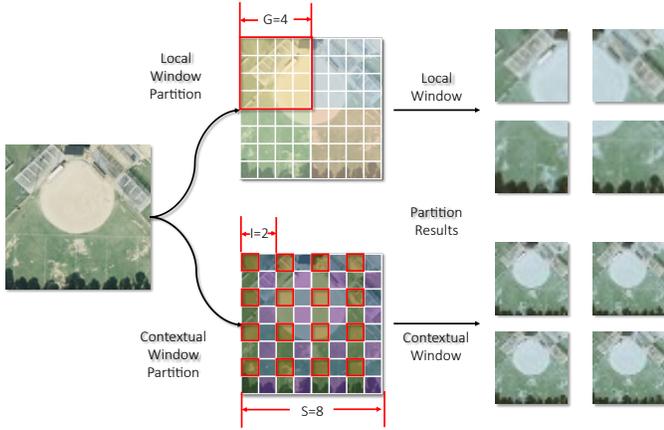


Fig. 2. Diagrams of difference between local window partition and contextual window partition. The pixels in the red box belong to the same window.

areas, ground objects and landforms with similar features are far apart in space. As shown in Fig. 1, the texture, shape and semantic features of the areas in the boxes are similar to each other but far apart in space. As a result, these methods that partition an image into fixed-size windows and employ self-attention within those windows to model pixel dependencies fail to capture interactions between distant but similar pixels, which is crucial for attaining optimal performance [31]. Furthermore, the majority of current transformer-based RSISR methods primarily rely on skip connections to transmit shallow features to deep features. However, treating these reused shallow features equally impedes the representational capability of transformers, despite the proven effectiveness of skip connections in RSISR [32]–[34].

To address these limitations, we propose two components: cross-space pixel Fusion attention (CSPIA) and cross-stage feature fusion attention (CSFFA). CSPIA allows the local window to perceive the contextual window (refer to Fig. 2) by maximizing the similarity between image pairs. This, in turn, effectively enlarges the receptive field, as depicted in Fig. 3, enabling the utilization of valuable context information from the image. In parallel, CSFFA enhances feature expression by adaptively integrating features cross-stages. CSPIA consists of three main steps: Space Division (SD), Local-Context Match-

ing (LCM) and Cross Attention (CA). The SD is responsible for obtaining local windows and contextual windows through different spatial partitioning strategies. Then, LCM is used to obtain the most matched contextual window for each pair of local window. Finally, the most similar contextual window is selected to conduct CA with corresponding local window. In this way, context information can be integrated into current local window efficiently, as shown in Fig 5. Building upon CSPIA, we construct a cross-spatial pixel integration block (CSPIB). Following the fusion of pixels from contextual windows, we apply the standard multihead self-attention (MSA) to capture local-range dependencies within the refined area of the local window and a local  $3 \times 3$  convolution further handles local details. Subsequently, CSFFA calculates cross-covariance of feature channels across stages to generate cross-stage attention map based on both shallow features and deep features (after projection of key and query). CSFFA enables the model to adaptively adjust the channel-wise feature maps at cross-stage of the network to enhance the informative multiscale feature representation ability. Furthermore, to enhance the flow of complementary features and allow subsequent network layers to focus on finer image details, we integrate a feed forward network (FFN) [35] into our model. Utilizing CSFFA and FFN, we construct a cross-stage feature fusion block (CSFFB). By combining CSPIBs and CSFFBs, we develop a transformer network, named SPIFFNet, which incorporates cross-spatial pixel integration and cross-stage feature fusion, specifically designed for RSISR. This architecture is depicted in Fig. 4. Furthermore, our experimental findings unequivocally demonstrate the superiority of the proposed SPIFFNet model over state-of-the-art methods.

The article presents three key contributions that can be summarized as follows:

- 1) We propose cross-spatial pixel integration attention (CSPIA) to introduce contextual information into the local window. The contextual information of the image enhances the global cognition and understanding of the entire image. By incorporating the context information within the local window, the model gains a better understanding of the relationship between ground features and the surrounding environment, thereby enhancing the consistency and accuracy of the RSISR results.

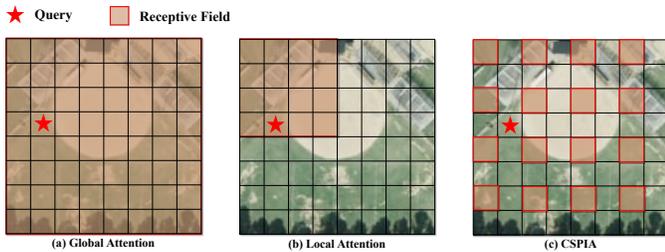


Fig. 3. Comparison of CSPIA with other Transformer models. The red star denote the query, and masks with solid line boundaries denote the regions to which the queries attend. (a) Global Attention [36] adopts full attention for all queries. (b) Local Attention [29] uses partitioned window attention. (c) CSPIA learns the context information region that is most similar to current query.

2) We propose cross-stage feature fusion attention (CSFFA), which facilitates effective feature representation by modeling the interdependencies among different channels across stages. By dynamically assigning weights to different channels, this mechanism enhances the model’s capacity to capture essential image features while suppressing irrelevant ones, thereby producing higher-quality super-resolved images.

3) Based on CSPIA and CSFFA, we propose SPIFFNet, a cross-spatial pixel integration and cross-stage feature fusion based transformer network for remote sensing image super-resolution. SPIFFNet effectively captures contextual information to enhance the global perception ability of the local window and adaptively fuses information from previous stages to enhance feature representation. Experimental results on benchmark datasets validate that SPIFFNet achieves state-of-the-art performance in terms of objective metrics as well as visual quality.

The rest of this article is organized as follows. Section II presents the related works on SR. Section III introduces the proposed SPIFFNet model, including the CSPIA and the CSFFA. Section IV presents the experimental results and analysis. Finally, Section V concludes this paper.

## II. RELATED WORK

### A. Deep Learning-based Methods for SR

Deep learning-based super-resolution (SR) methods predominantly rely on standard convolutional neural networks (CNNs) owing to their robust nonlinear representation capabilities. Typically, these methods approach super-resolution as an image-to-image regression task, with the objective of learning the direct mapping from LR to HR images. SRCNN [37] first uses three convolution layers to map the low-resolution images to high-resolution images. Building upon SRCNN, Kim *et al.* extended the network depth in their work called DRCN [38], resulting in considerable performance improvements over SRCNN [37]. FSRCNN [39] achieves high computational efficiency without compromising restoration quality through a redesigned architecture of SRCNN. VDSR [40] addressed the challenge of handling multi-scale images within a unified framework by incorporating residual learning, gradient cropping, and an increased number of network layers. EDSR [41] achieves superior performance by streamlining the

model architecture, eliminating redundant modules from the conventional ResNet framework. For remote sensing images, LGCNet [42] stands as the pioneering CNN-based model for super-resolution, introducing the concept of local and global contrast features to enhance the preservation of details and clarity in the reconstructed images. Haut *et al.* [43] coordinates several different improvements in network design to achieve the most advanced performance on the RSISR task. A novel single-path feature reuse approach and a second-order learning mechanism are proposed by Dong *et al.* [44], which aim to effectively utilize both small and large difference features. Although these methods have achieved impressive results, the limited receptive field of CNNs cannot capture the long-range dependencies between pixels, thereby limiting their performance.

### B. Transformer-based Methods

The Transformer network, initially proposed in 2017 for machine translation tasks [45], has gained popularity in computer vision due to its remarkable performance in image processing [22], [46]. Since its inception, numerous visual models based on the Transformer have been proposed [28], [30], [35], [36], [47]. As an example, Chen *et al.* introduced the Image Processing Transformer (IPT) as a novel pre-trained model for low-level computer vision tasks. To fully leverage the potential of the transformer, a substantial amount of corrupted image pairs is generated using the ImageNet dataset. The IPT model adapts to diverse image processing tasks through multi-head and multi-tail training, along with the incorporation of contrastive learning techniques. Another well-known image restoration method called Uformer [30] utilizes the Locally Enhanced Window Transformer block, which reduces computational requirements by utilizing a non-overlapping window-based self-attention mechanism. Moreover, three skip connection schemes are explored to facilitate efficient information transfer from the encoder to the decoder. Furthermore, the Restoration Transformer (Restormer) network [35] proposes an effective Transformer model that captures remote pixel interactions and is suitable for large image, which is achieved through key design choices in the building blocks, including multi-head attention and feedforward networking. Efficient Super-Resolution Transformer (ESRT) [47] presents a hybrid model that combines a lightweight CNN backbone (LCB) with a lightweight transformer backbone (LTB) can dynamically adjust the feature map size, achieving competitive results with low computational cost. SwinIR [29] introduced a robust baseline model for image restoration that utilizes the Swin Transformer architecture. In the context of RSISR, TransENet [36] proposes a multilevel enhancement architecture based on the Transformer framework, which can be integrated with the conventional super-resolution (SR) framework to effectively merge multi-scale high- and low-dimensional features.

## III. METHODOLOGY

In this section, we introduce the proposed SPIFFNet for RSISR. The overall framework of SPIFFNet is presented in Section III-A and the SPIFFNet group that integrates CSPIA

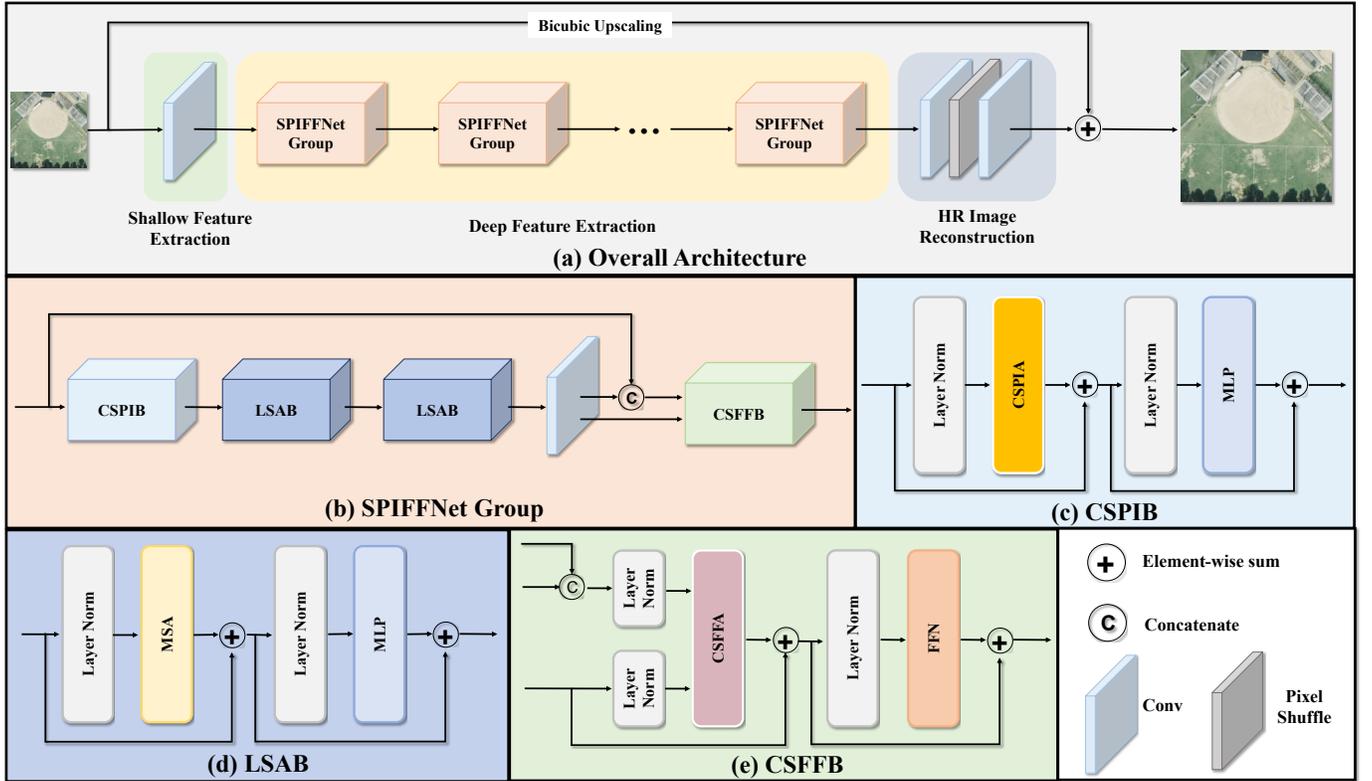


Fig. 4. (a) Architecture of SPIFFNet for high-resolution image restoration. (b) The main components of the model: Cross-Spatial Pixel Integration Block (CSPIB), Local Spatial Attention Block (LSAB), Local  $3 \times 3$  Convolution and Cross-Stage Feature Fusion Block (CSFFB). (c) Cross-Spatial Pixel Integration Block (CSPIB) that implements the injection of context information into a local window. (d) Local Spatial Attention Block (LSAB) model local-range dependencies in the fine area of a local window. (e) Cross-Stage Feature Fusion Block (CSFFB) performs channel attention across stages for feature fusion.

and CSFFB is carefully discussed in Section III-B. Furthermore, in Section III-C, we will provide a concise overview of the implementation details.

#### A. Overview of SPIFFNet

In this section, we introduce the framework of our method, as shown in Fig. 4. Given an input image  $I_{LR} \in \mathbb{R}^{H \times W \times 3}$ , where  $H$  and  $W$  are the image height and width. Then, the input  $I_{LR}$  undergoes a transformation into the feature space through a  $3 \times 3$  convolutional layer

$$F_0 = \text{Conv}(I_{LR}) \quad (1)$$

where the Conv denotes  $3 \times 3$  convolution and the  $F_0 \in \mathbb{R}^{H \times W \times C}$  represents the shallow features.

Then, several SPIFFNet groups, each of which involves CSPIB, LSAB, local  $3 \times 3$  convolution and CSFFB are set up after the convolutional layer for deep feature extraction. We extract deep feature  $F_{DF} \in \mathbb{R}^{H \times W \times C}$  from  $F_0$  as

$$F_{DF} = H_{DF}(F_0) \quad (2)$$

where  $H_{DF}$  represents the deep feature extraction module which contains  $K$  SPIFFNet groups. Specifically, the intermediate features  $F_1, F_2, \dots, F_K$  and the final deep feature  $F_{DF}$  are extracted sequentially

$$F_i = H_i(F_{i-1}), \quad i = 1, 2, \dots, K \quad (3)$$

where  $H_i$  denotes the  $i$ -th SPIFFNet group.

Finally, the deepest features  $F_{DF}$  are reconstructed using a  $3 \times 3$  convolutional layer and pixel-shuffle upsampling operations [48] to generate SR image  $I_r$ . In addition, a bilinear interpolation of the LR image  $I_b$  is incorporated in the summation process to aid in the recovery process for the super-resolution output  $I_r$

$$I_{SR} = I_r + I_b \quad (4)$$

We train the proposed model using the L1 loss function. The loss function is obtained by comparing the LR images  $I_{LR}$  with their corresponding HR reference images  $I_{HR}$

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \left\| I_{HR}^{(i)} - I_{SR}^{(i)} \right\|_1 \quad (5)$$

where  $\theta$  represents the parameters of the SPIFFNet, and  $N$  denotes the number of training samples.

#### B. SPIFFNet Group

In this section, we introduce the SPIFFNet group, a crucial component of our SPIFFNet model. Each SPIFFNet group consists of four components: cross-spatial pixel integration attention block (CSPIB), Local Spatial Attention Block (LSAB), local  $3 \times 3$  convolution and cross-stage feature fusion block (CSFFB). The Cross-Spatial Pixel Integration Block (CSPIB) expands the model's receptive field, allowing it to capture

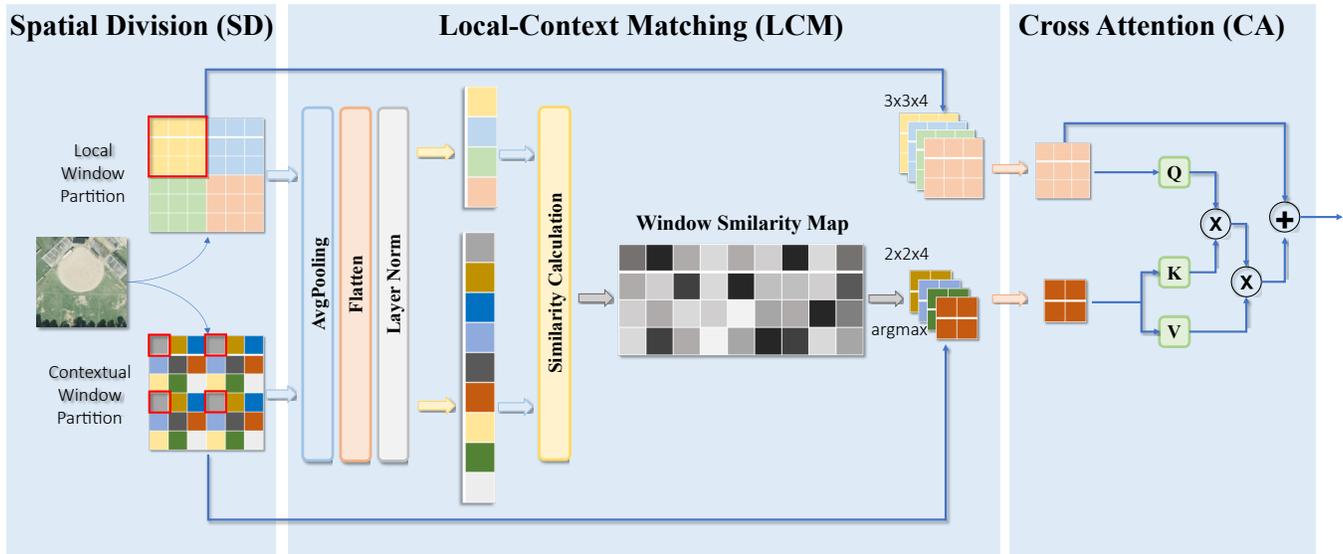


Fig. 5. Illustration of cross-spatial pixel integration attention (CSPIA).  $\otimes$  and  $\oplus$  are matrix multiplication and element-wise addition operations, respectively.

long-range dependencies and contextual information from the input feature maps. This expansion is facilitated by the utilization of the Cross-Spatial Pixel Integration Attention (CSPIA). The LSAB is responsible for capturing the correlations of local spatial information. Local  $3 \times 3$  convolution deals with local details in a fine-grained manner. CSFFB is designed to adaptive integration of information from the previous stage according to the needs of the characteristics of the current stage rather than treating them equally. This combination of local and global information, along with the cross stage adaptive information fusion, enables the model to capture complex spatial dependencies and contextual information effectively.

1) *Cross-Spatial Pixel Integration Block (CSPIB)*: Previous methods often overlooked the interaction between local and contextual features. The CSPIB is designed to expand the local spatial window to capture more context information, as shown in Fig. 4(c). CSPIB contains two sequential modules, the cross-spatial pixel integration attention (CSPIA) is designed to capture contextual information for local windows and the MLP module for feature projection.

**Cross-Spatial Pixel Integration Attention (CSPIA)**: In this section, our objective is to expand the receptive field of local windows, enabling them to capture context information from the input feature maps. Previous studies have demonstrated that SR networks with a wider effective receptive field achieve superior performance [31]. The challenge lies in enabling the network to model global connectivity while preserving computational efficiency. Due to the fixed partitioning of windows at the layer level, there are no direct connections between windows. One straightforward approach is to exhaustively combine the information from every window pair. However, this approach is unnecessary and inefficient since many windows are irrelevant and uninformative. Additionally, redundant interactions may introduce noise that impairs the model’s performance. Based on these observations, we introduce an innovative technique called cross-spatial pixel

integration attention (CSPIA), where each local window adaptively integrates pixels with the most correlated global window. Specifically, as shown in Fig. 5, the CSPIA consists of three steps: Spatial Division (SD), Local-Context Matching (LCM) and Cross Attention (CA). Through SD, we split the feature map into two parts: local windows and global windows. In the case of local windows, adjacent embeddings of size  $G \times G$  are grouped together. An example is illustrated in Fig. 2 with  $G = 4$ . In the case of global windows, where the input size is  $S \times S$ , the feature map is sampled at a fixed interval  $I$ . Fig. 2 demonstrates an example with  $I = 2$ , where embeddings with a red border belong to a window. The height or width of the group for global windows is calculated as  $G = S/I$ . All windows can be processed in parallel, after which the outputs are pasted to their original location in window aggregation module. It is worthwhile to note that such cropping strategy is adaptive to arbitrary input size, which means no padding pixels are needed. Then, local windows and contextual windows are spatially pooled into one-dimensional tokens. These tokens encode the characteristic of the windows, which are later used for similarity calculation and window matching. This process can be expressed as

$$\bar{X}_i = \arg \max_{X_j} L(X_i)^T L(X_j), j \neq i \quad (6)$$

where  $\bar{X}_i$  is the best-matching global window with current local window  $X_i$ , and  $L(\cdot)$  is the average pooling function along spatial dimension followed by flatten operation and layer normalization. Since the *argmax* operation is non-differentiable, we replace it with Gumbel-Softmax operation [49] during training so as to make it possible to train end-to-end. After that, pixel information of  $\bar{X}_i$  are fused into  $X_i$  via Cross-Attention (CA)

$$X_i = CA(X_i, \bar{X}_i) \quad (7)$$

As illustrated in Fig. 5, CA works in a similar way to the standard self-attention [22], but the key and value are calcu-

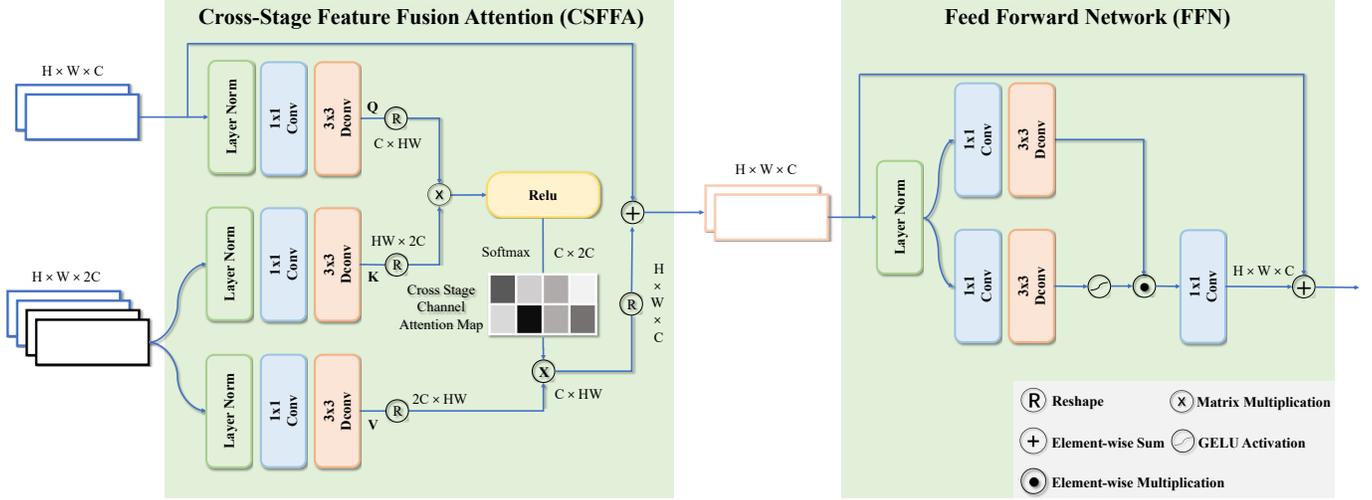


Fig. 6. Illustration of cross-stage feature fusion block (CSFFB). It consists of two parts: cross-stage feature fusion attention (CSFFA) and FFN.

lated using  $\bar{X}_i$ . As a result, CSPIA can enable contextual pixel integration while introducing little computational overhead.

2) *Local Spatial Attention Block (LSAB)*: As shown in Fig. 4(d), LSAB adopts the standard multihead self-attention (MSA) paradigm [22], with two modifications. Firstly, LSAB operates at the window level instead of the image level. Secondly, positional embedding is omitted due to the introduction of the convolutional layer, which implicitly learns positional relationships and enhances the network’s efficiency and conciseness. LSAB is designed to model local-range dependencies within a window, facilitating the comprehensive utilization of contextual information.

Specifically, for feature  $X \in \mathbb{R}^{P^2 \times C}$ , the corresponding *query*, *key* and *value* matrices  $Q \in \mathbb{R}^{P^2 \times d}$ ,  $K \in \mathbb{R}^{P^2 \times d}$ ,  $V \in \mathbb{R}^{P^2 \times C}$  are computed as

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V \quad (8)$$

where the weight matrices  $W_Q$ ,  $W_K$  and  $W_V$  are shared across windows,  $P$  is the window size. By comparing the similarity between  $Q$  and  $K$ , we obtain a attention map of size  $\mathbb{R}^{P^2 \times P^2}$  and multiply it with  $V$ . Overall, the calculation of Multi-head Self-Attention (MSA) can be expressed as

$$MSA(X) = softmax(QK^T / \sqrt{d})V \quad (9)$$

Here  $\sqrt{d}$  is used to control the magnitude of  $QK^T$  before applying the softmax function.

Similar to the conventional transformer layer [22], the MLP is employed after MSA module to further transform features. MLP contains two fully-connected layers, and one GELU nonlinearity is applied after the first linear layer.

3) *Local 3 × 3 Convolution*: By adding a local 3 × 3 convolutional layer after feature extraction, the Transformer-based network is infused with the inherent inductive bias of convolution operations. This enhances the foundation for aggregating shallow and deep features in subsequent stages.

4) *Cross-Stage Feature Fusion Block (CSFFB)*: Skip Connections are commonly used to propagate shallow features to deeper layers. However, the long-term information from the

shallow stages tends to be attenuated. Although the shallow features can be reused through skip connections, they are treated indiscriminately with the deep features across different stages, thereby impeding the representational capacity of CNNs. To address this concern, we introduce CSFFB, illustrated in Fig. 4(e). CSFFB comprises two consecutive components: the Cross-Stage Feature Fusion Attention (CSFFA) for adaptive feature fusion across stages, and the FFN for feature transformation.

**Cross-Stage Feature Fusion Attention (CSFFA)**: Figure 6 demonstrates the operation of the CSFFA, which calculates attention scores at the channel level using feature maps cross-stages. These scores are then applied to the feature maps, enabling the weighting of each channel’s contribution and the integration of features from previous stages with the current input. This process promotes the fusion of channel information across stages, facilitating the model to capture both low-level details and high-level contextual information. As a result, the approach enables more precise and effective image super-resolution.

Specifically, the features of the current stage, denoted as  $X_{cur} \in \mathbb{R}^{H \times W \times C}$ , and the previous stage, denoted as  $X_{pre} \in \mathbb{R}^{H \times W \times C}$ , are concatenated along the feature dimension to obtain  $Y \in \mathbb{R}^{H \times W \times 2C}$ . Our CSFFA then gets *query*, *key*, and *value* projections:  $Q = W_d^Q W_p^Q X_{cur}$ ,  $K = W_d^K W_p^K Y$ , and  $V = W_d^V W_p^V Y$ . Where  $W_p^{(\cdot)}$  denotes the 1 × 1 point-wise convolution, and  $W_d^{(\cdot)}$  represents the 3 × 3 depth-wise convolution. Subsequently, we reshape the  $Q$  and  $K$  to facilitate their dot-product interaction, resulting in a cross-stage channel attention map  $A$  with dimensions  $\mathbb{R}^{C \times 2C}$ . The CSFFA process can be defined as

$$X_{out} = W_p Attn(Q, K, V) + X_{cur} \quad (10)$$

$$Attn(Q, K, V) = V \cdot Softmax(ReLu(Q \cdot K / \alpha))$$

where  $X_{cur}$  and  $X_{out}$  are the input and output feature maps, respectively;  $Q \in \mathbb{R}^{C \times HW}$ ;  $K \in \mathbb{R}^{HW \times 2C}$ ; and  $V \in \mathbb{R}^{2C \times HW}$  matrices are obtained by the  $X_{cur} \in \mathbb{R}^{H \times W \times C}$  and  $Y \in \mathbb{R}^{H \times W \times 2C}$ , respectively. To enhance GELU control

TABLE I  
MEAN PSNR (DB) AND SSIM OVER THE UCMERGED TEST DATASET AND THE AID TEST DATASET

Dataset	Scale	Metric	SRCNN	LGCNet	VDSR	DCM	HSENet	TransENet	SPIFFNet Ours
UCMerced	×2	PSNR	32.84	33.48	33.87	33.65	34.22	34.03	<b>34.68</b>
		SSIM	0.9152	0.9235	0.9280	0.9274	0.9327	0.9301	<b>0.9354</b>
	×3	PSNR	28.66	29.28	29.76	29.52	30.00	29.92	<b>30.43</b>
		SSIM	0.8038	0.8238	0.8354	0.8394	0.8420	0.8408	<b>0.8510</b>
	×4	PSNR	26.78	27.02	27.54	27.22	27.73	27.77	<b>28.09</b>
		SSIM	0.7219	0.7333	0.7522	0.7528	0.7623	0.7630	<b>0.7733</b>
AID	×2	PSNR	34.49	34.80	35.05	35.21	35.30	35.28	<b>35.66</b>
		SSIM	0.9286	0.9320	0.9346	0.9366	0.9377	0.9374	<b>0.9397</b>
	×3	PSNR	30.55	30.73	31.15	31.31	31.39	31.45	<b>31.70</b>
		SSIM	0.8372	0.8417	0.8522	0.8561	0.8581	0.8595	<b>0.8631</b>
	×4	PSNR	28.40	28.61	28.99	29.17	29.34	29.38	<b>29.54</b>
		SSIM	0.7561	0.7626	0.7753	0.7824	0.7881	0.7909	<b>0.7938</b>

TABLE II  
ABLATION EXPERIMENTS ON UCMERGED-X4

CSPIA	CSFFA	Param	PSNR	SSIM
✗	✗	1.459M	28.009	0.7711
✗	✓	1.510M	28.039	0.7718
✓	✗	1.537M	28.077	0.7723
✓	✓	1.659M	28.086	0.7732

and promote the development of sophisticated image attributes, we introduce a ReLU non-linearity function before the softmax normalization. The ReLU non-linearity function applies sparse constraints to the cross-stage attention map promotes the model’s focus on the most informative regions and mitigates the influence of noisy or irrelevant features.

To transform features, we use two parallel  $1 \times 1$  convolutions and  $3 \times 3$  convolutions to the feature map. Subsequently, a SimpleGate activation function [50] multiplies one of the branches to regulate the flow of complementary features and facilitate feature transformation. To be specific, when provided with an input tensor  $X \in \mathbb{R}^{H \times W \times C}$ , the FFN can be expressed as

$$\hat{X} = W_p^0 \text{Gating}(X) + X$$

$$\text{Gating}(X) = \phi(W_d^1 W_p^1(\text{LN}(X))) \odot W_d^2 W_p^2(\text{LN}(X)) \quad (11)$$

where the symbol  $\odot$  denotes element-wise multiplication,  $\phi$  refers to the GELU non-linearity function, and LN denotes layer normalization [51].

### C. Implementation Details

This article focuses on RSISR at three magnification factors:  $\times 2$ ,  $\times 3$ , and  $\times 4$ . During the training, we randomly sample LR remote sensing images and their corresponding HR reference windows as  $48 \times 48$  windows. To augment the training samples, we apply random rotations ( $90^\circ$ ,  $180^\circ$ , and  $270^\circ$ ) and horizontal flipping. The proposed SPIFFNet consists of 10 blocks, each of which is designed with a local window size and a global window size of 16 and a feature dimension of 64. Additionally, the SPIFFNet employs 4 attention heads. Further details and experimental analyses are presented in Section IV.

We employ the Adam optimizer [52] for model optimization, setting  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ , and  $\varepsilon = 10^{-8}$ . We set the

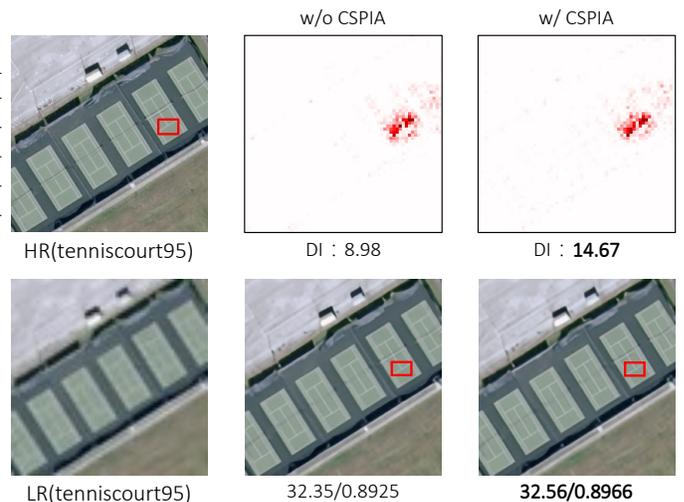


Fig. 7. LAM [31] comparison between the full SPIFFNet (w/o CSPIA) and the variant without GPA (w/ CSPIA) for  $\times 4$  SR. The first column shows the low remote sensing resolution image and the corresponding high remote sensing resolution image. The second and third columns show SPIFFNet and its variants and the corresponding effective receptive fields.

initial learning rate to  $4 \times 10^{-4}$  and the mini-batch size to 8. We train the model for a total of 2000 epochs, gradually reducing the learning rate to  $5 \times 10^{-7}$  at epoch 2000 using the cosine annealing schedule [53].

## IV. EXPERIMENTAL RESULTS AND ANALYSES

### A. Experimental Datasets and Metrics

1) *Datasets*: This study employs two publicly available remote sensing datasets, namely UCMerced [54] and AID [55]. The UCMerced dataset comprises 21 classes, each containing 100 remote sensing scenes images that have dimensions of  $256 \times 256$  pixels. We partitioned the dataset into two subsets: one for training purposes and the other for testing. Each subset comprises 1050 images. The AID dataset consists of 10,000 images representing 30 classes of remote sensing scenes which have dimensions of  $600 \times 600$  pixels. In the case of the AID dataset, 80% of the total dataset is randomly assigned as the training set, while the remaining images are allocated for testing.

TABLE III  
MEAN PSNR(DB) OF EACH CLASS FOR UPSCALING FACTOR 3 ON UCMERGED TEST DATASET

Class NO.	Class name	LGCNet	VDSR	DCM	HSENet	TransENet	SPIFFNet (Ours)
1	agricultural	27.66	27.75	<b>29.06</b>	27.64	28.02	28.29
2	airplane	29.12	29.76	<b>30.77</b>	30.09	29.94	30.45
3	baseballdiamond	34.72	34.98	33.76	35.05	35.04	<b>35.61</b>
4	beach	37.37	37.57	36.38	37.69	37.53	<b>37.90</b>
5	buildings	27.81	28.53	28.51	28.95	28.81	<b>29.57</b>
6	chaparral	26.39	26.61	26.81	26.70	26.69	<b>26.87</b>
7	denseresidential	28.25	28.88	28.79	29.24	29.11	<b>29.66</b>
8	forest	28.44	28.52	28.16	28.59	28.59	<b>28.67</b>
9	freeway	29.52	30.21	30.45	30.63	30.38	<b>31.13</b>
10	golfcourse	36.51	36.57	34.43	36.62	36.68	<b>37.05</b>
11	harbor	23.63	24.36	<b>26.55</b>	24.88	24.72	25.47
12	intersection	28.29	28.83	29.28	29.21	29.03	<b>29.67</b>
13	mediumresidential	27.76	28.26	27.21	28.55	28.47	<b>28.94</b>
14	mobilehomepark	24.59	25.24	26.05	25.70	25.64	<b>26.25</b>
15	overpass	26.58	27.70	27.77	28.22	27.83	<b>28.55</b>
16	parkinglot	23.69	24.12	24.95	24.66	24.45	<b>25.37</b>
17	river	29.12	29.23	28.89	29.22	29.25	<b>29.54</b>
18	runway	31.15	31.41	<b>32.53</b>	31.15	31.25	31.69
19	sparseresidential	30.53	31.47	29.81	31.64	31.57	<b>31.89</b>
20	storagetanks	32.17	32.72	29.02	32.95	32.71	<b>33.23</b>
21	tenniscourt	31.58	32.29	30.76	32.71	32.51	<b>33.19</b>
	avg	29.28	29.76	29.52	30.00	29.92	<b>30.43</b>

2) *Metrics*: We select PSNR and SSIM [56] as the evaluation metrics for RSISR, and assess all super-resolution results on the RGB channels.

### B. Ablation Studies

We conducted experiments in this section to validate the components of our method. All experiments were performed using the same experimental setup, with the UCMerged dataset and a uniform magnification factor of 4 was applied. We start with a naive baseline by removing both components. Then we add CSPIA and CSFFA to the baseline, respectively. At last, both components are employed to compose our final version of method. The results are reported in Table II.

1) *Effects of CSPIA*: Table II summarizes the results of this ablation study. We can see that the model with CSPIA significantly outperforms the baseline model, indicating the effectiveness of the expandable window mechanism in capturing both global and local features.

To better understand the main reason of the improvement brought by CSPIA, we utilize LAM [31] to visualize the effective receptive field of a input window. As shown in Fig. 7, the window benefits from a global range of useful pixels by using CSPIA. The results indicate the effectiveness of the proposed CSPIA in improving PSNR and SSIM performances.

2) *Effects of CSFFA*: Table II summarizes the results of this ablation study. The model incorporating CSFFA demonstrates a significant performance improvement over the baseline model, providing strong evidence of the effectiveness of CSFFA.

### C. Comparisons with Other Methods

This section presents a comparative analysis of the proposed method with several deep learning-based SR methods, namely SRCNN [37], VDSR [40], LGCNet [42], DCM [43], HSENet [57], and TransENet [36].

1) *Quantitative Results on UCMerged Dataset*: The performance of these approaches on the UCMerged dataset is reported in Table I. The best result is indicated in bold font. Notably, some results have been reported in multiple published articles [43], [58]. To ensure consistency, we retrained these comparison methods using the open-source code, subjecting all methods to the same testing conditions. The results demonstrate that our SPIFFNet achieves the highest values in terms of PSNR and SSIM. Table III provides a summary of the PSNR for each class in the UCMerged dataset at an upscale factor of 3. We observed that SPIFFNet significantly outperforms HSENet [57] and TransENet [36] on the buildings, harbors and parking lot classes, which require precise local context information for object discrimination and image detail reconstruction. These notable results serve as further evidence of the effectiveness of our proposed method.

2) *Quantitative Results on AID Dataset*: We conducted additional experiments on the AID dataset to further validate the effectiveness of SPIFFNet. Table I presents the PSNR and SSIM results of SPIFFNet compared to other methods on this dataset. Compared to other methods, SPIFFNet achieves the best results. Furthermore, Table [55] presents the test results for each category at a magnification factor of 4 on the AID dataset. The consistent superior performance of SPIFF in various scenarios further demonstrates the effectiveness of our approach.

3) *Qualitative Results*: Fig. 8 displays several super-resolved examples from the UCMerged dataset, such as scenes depicting "agricultural", "buildings", and "overpass". Similarly, Fig. 9 showcases examples from the AID dataset, including scenes of "playground" and "parking". Our method demonstrates superior performance compared to other methods in challenging areas such as texture and edge, as evident from the visual results. This observation provides further evidence of the effectiveness of our approach.

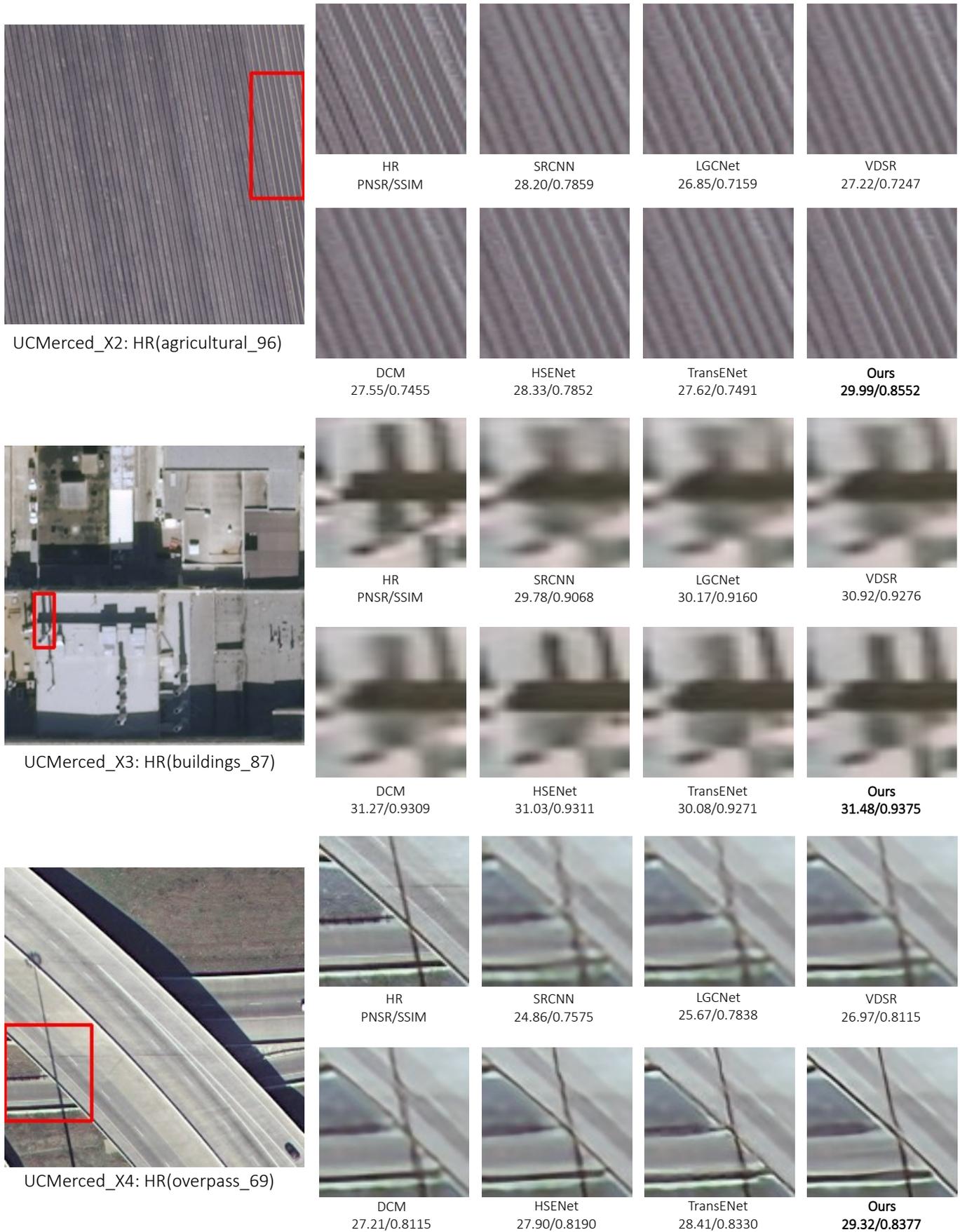


Fig. 8. Result comparisons on UCMerced dataset with different methods.

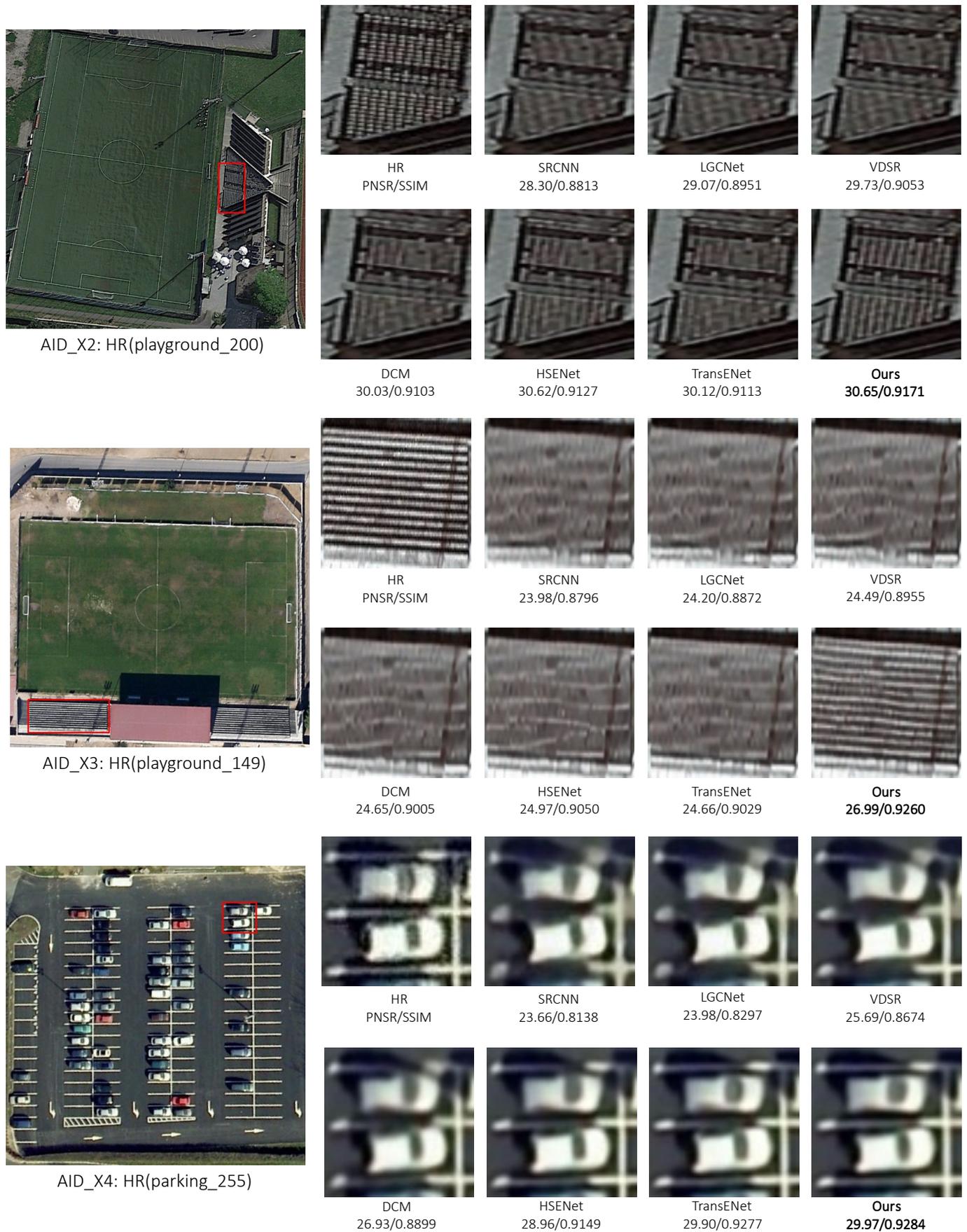


Fig. 9. Result comparisons on AID dataset with different methods.

TABLE IV  
MEAN PSNR(DB) OF EACH CLASS FOR UPSCALING FACTOR 4 ON AID TEST DATASET

Class NO.	Class Name	LGCNet	VDSR	DCM	HSENet	TransENet	SPIFFNet (Ours)
1	airport	28.39	28.82	28.99	29.15	29.23	<b>29.34</b>
2	bareland	35.78	35.98	36.17	36.25	36.20	<b>36.48</b>
3	baseballfield	30.75	31.18	31.36	31.52	<b>31.59</b>	31.48
4	beach	32.08	32.29	32.45	32.54	32.55	<b>32.78</b>
5	bridge	30.67	31.19	31.39	31.57	31.63	<b>31.87</b>
6	center	26.92	27.48	27.72	27.95	28.03	<b>28.21</b>
7	church	23.68	24.12	24.29	24.47	24.51	<b>24.63</b>
8	commercial	27.24	27.62	27.78	27.94	27.97	<b>28.07</b>
9	denseresidential	24.33	24.70	24.87	25.06	25.13	<b>25.20</b>
10	desert	39.06	39.13	39.27	39.37	39.31	<b>39.58</b>
11	farmland	33.77	34.20	34.42	34.56	34.58	<b>34.79</b>
12	forest	28.20	28.36	28.47	28.56	28.56	<b>28.62</b>
13	industrial	26.24	26.72	26.92	27.13	27.21	<b>27.32</b>
14	meadow	32.65	32.77	32.88	32.94	32.94	<b>33.14</b>
15	mediumresidential	27.63	28.06	28.25	28.44	28.45	<b>28.57</b>
16	mountain	28.97	29.11	29.18	29.22	29.28	<b>29.33</b>
17	park	27.37	27.69	27.82	27.95	28.01	<b>28.16</b>
18	parking	24.40	25.21	25.74	26.27	26.40	<b>26.56</b>
19	playground	29.04	29.62	29.92	30.20	30.30	<b>30.64</b>
20	pond	30.00	30.26	30.39	30.49	30.53	<b>30.86</b>
21	port	26.02	26.43	26.62	26.84	26.91	<b>26.99</b>
22	railwaystation	27.76	28.19	28.38	28.55	28.61	<b>28.74</b>
23	resort	27.32	27.71	27.88	28.04	28.08	<b>28.13</b>
24	river	30.60	30.82	30.91	30.98	31.00	<b>31.11</b>
25	school	26.34	26.78	26.94	27.15	27.22	<b>27.35</b>
26	sparseresidential	26.27	26.46	26.53	26.63	26.63	<b>26.68</b>
27	square	28.39	28.91	29.13	29.33	29.39	<b>29.59</b>
28	stadium	26.37	26.88	27.10	27.32	27.41	<b>27.54</b>
29	storagetanks	25.48	25.86	26.00	26.16	26.20	<b>26.26</b>
30	viaduct	27.26	27.74	27.93	28.13	28.21	<b>28.36</b>
	avg	28.61	28.99	29.17	29.34	29.38	<b>29.54</b>

## V. CONCLUSION

This paper introduces a novel transformer-based method called Cross-Spatial Pixel Integration and Cross-Stage Feature Fusion Based Transformer Network (SPIFFNet) for RSISR. The aim of SPIFFNet is to enhance the global perception ability of local windows by introducing context information and to improve the representation ability of features by integrating cross-stage features. SPIFFNet consists of two key components: CSPIA, which introduces context information into the reconstruction of local windows to enhance global awareness, and CSFFA, which enables adaptive aggregation of features across different stages of the network, resulting in more effective information fusion and superior super-resolution performance. We conducted extensive experiments on benchmark datasets to validate the effectiveness of SPIFFNet.

## REFERENCES

- [1] Z. Yin, F. Ling, X. Li, X. Cai, H. Chi, X. Li, L. Wang, Y. Zhang, and Y. Du, "A cascaded spectral-spatial cnn model for super-resolution river mapping with modis imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2021.
- [2] S. Yang, M. Wang, P. Li, L. Jin, B. Wu, and L. Jiao, "Compressive hyperspectral imaging via sparse tensor and nonlinear compressed sensing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 11, pp. 5943–5957, 2015.
- [3] J. Xue, Y. Zhao, S. Huang, W. Liao, J. C.-W. Chan, and S. G. Kong, "Multilayer sparsity-based tensor decomposition for low-rank tensor completion," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 11, pp. 6916–6930, 2021.
- [4] F. Min, L. Wang, S. Pan, and G. Song, "D 2 unet: Dual decoder u-net for seismic image super-resolution reconstruction," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–13, 2023.
- [5] F. Wang, J. Li, Q. Yuan, and L. Zhang, "Local-global feature-aware transformer based residual network for hyperspectral image denoising," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–19, 2022.
- [6] Y. Zhang, R. Zong, L. Shang, and D. Wang, "On coupling classification and super-resolution in remote urban sensing: An integrated deep learning approach," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–17, 2022.
- [7] J. Yang, L. Xiao, Y.-Q. Zhao, and J. C.-W. Chan, "Variational regularization network with attentive deep prior for hyperspectral-multispectral image fusion," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–17, 2021.
- [8] Y. Liu, Z. Xiong, Y. Yuan, and Q. Wang, "Distilling knowledge from super resolution for efficient remote sensing salient object detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [9] J. Xue, Y. Zhao, W. Liao, and J. C.-W. Chan, "Nonlocal low-rank regularized tensor decomposition for hyperspectral image denoising," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 7, pp. 5174–5189, 2019.
- [10] J. Yang, L. Xiao, Y.-Q. Zhao, and J. C.-W. Chan, "Unsupervised deep tensor network for hyperspectral-multispectral image fusion," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [11] X. Wang, J. Yi, J. Guo, Y. Song, J. Lyu, J. Xu, W. Yan, J. Zhao, Q. Cai, and H. Min, "A review of image super-resolution approaches based on deep learning and applications in remote sensing," *Remote Sensing*, vol. 14, no. 21, p. 5423, 2022.
- [12] R. Keys, "Cubic convolution interpolation for digital image processing," *IEEE transactions on acoustics, speech, and signal processing*, vol. 29, no. 6, pp. 1153–1160, 1981.
- [13] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [14] S. Jia, Z. Wang, Q. Li, X. Jia, and M. Xu, "Multiattention generative adversarial network for remote sensing image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022.
- [15] S. Lei, Z. Shi, and Z. Zou, "Coupled adversarial training for remote sensing image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 5, pp. 3633–3643, 2019.

- [16] X. Dong, X. Sun, X. Jia, Z. Xi, L. Gao, and B. Zhang, "Remote sensing image super-resolution using novel dense-sampling networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 2, pp. 1618–1633, 2020.
- [17] S. Wang, T. Zhou, Y. Lu, and H. Di, "Contextual transformation network for lightweight remote-sensing image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2021.
- [18] X. Jin, J. He, Y. Xiao, and Q. Yuan, "Learning a local-global alignment network for satellite video super-resolution," *IEEE Geoscience and Remote Sensing Letters*, vol. 20, pp. 1–5, 2023.
- [19] J. Liu, Z. Wu, L. Xiao, and X.-J. Wu, "Model inspired autoencoder for unsupervised hyperspectral image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–12, 2022.
- [20] Y. Li, Q. Xu, Z. He, and W. Li, "Progressive task-based universal network for raw infrared remote sensing imagery ship detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [21] D. He and Y. Zhong, "Deep hierarchical pyramid network with high-frequency-aware differential architecture for super-resolution mapping," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–15, 2023.
- [22] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [23] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," pp. 568–578, 2021.
- [24] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," pp. 22–31, 2021.
- [25] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," pp. 213–229, 2020.
- [26] J. Beal, E. Kim, E. Tzeng, D. H. Park, A. Zhai, and D. Kislyuk, "Toward transformer-based object detection," *arXiv preprint arXiv:2012.09958*, 2020.
- [27] Z. Sun, S. Cao, Y. Yang, and K. M. Kitani, "Rethinking transformer-based set prediction for object detection," pp. 3611–3620, 2021.
- [28] H. Chen, Y. Wang, T. Guo, C. Xu, Y. Deng, Z. Liu, S. Ma, C. Xu, C. Xu, and W. Gao, "Pre-trained image processing transformer," pp. 12 299–12 310, 2021.
- [29] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "Swinir: Image restoration using swin transformer," pp. 1833–1844, 2021.
- [30] Z. Wang, X. Cun, J. Bao, W. Zhou, J. Liu, and H. Li, "Uformer: A general u-shaped transformer for image restoration," pp. 17 683–17 693, 2022.
- [31] J. Gu and C. Dong, "Interpreting super-resolution networks with local attribution maps," pp. 9199–9208, 2021.
- [32] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," pp. 286–301, 2018.
- [33] R. Lan, L. Sun, Z. Liu, H. Lu, Z. Su, C. Pang, and X. Luo, "Cascading and enhanced residual networks for accurate single-image super-resolution," *IEEE transactions on cybernetics*, vol. 51, no. 1, pp. 115–125, 2020.
- [34] R. Lan, L. Sun, Z. Liu, H. Lu, C. Pang, and X. Luo, "Madnet: a fast and lightweight network for single-image super resolution," *IEEE transactions on cybernetics*, vol. 51, no. 3, pp. 1443–1453, 2020.
- [35] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," pp. 5728–5739, 2022.
- [36] S. Lei, Z. Shi, and W. Mo, "Transformer-based multistage enhancement for remote sensing image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2021.
- [37] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 295–307, 2015.
- [38] J. Kim, J. K. Lee, and K. M. Lee, "Deeply-recursive convolutional network for image super-resolution," pp. 1637–1645, 2016.
- [39] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," pp. 391–407, 2016.
- [40] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," pp. 1646–1654, 2016.
- [41] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," pp. 136–144, 2017.
- [42] S. Lei, Z. Shi, and Z. Zou, "Super-resolution for remote sensing images via local-global combined network," *IEEE Geoscience and Remote Sensing Letters*, vol. 14, no. 8, pp. 1243–1247, 2017.
- [43] J. M. Haut, M. E. Paoletti, R. Fernández-Beltrán, J. Plaza, A. Plaza, and J. Li, "Remote sensing single-image superresolution based on a deep compendium model," *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 9, pp. 1432–1436, 2019.
- [44] X. Dong, L. Wang, X. Sun, X. Jia, L. Gao, and B. Zhang, "Remote sensing image super-resolution using second-order multi-scale networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 4, pp. 3473–3485, 2020.
- [45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [46] N. Parmar, A. Vaswani, J. Uszkoreit, Ł. Kaiser, N. Shazeer, A. Ku, and D. Tran, "Image transformer," pp. 4055–4064, 2018.
- [47] Z. Lu, J. Li, H. Liu, C. Huang, L. Zhang, and T. Zeng, "Transformer for single image super-resolution," pp. 457–466, 2022.
- [48] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," pp. 1874–1883, 2016.
- [49] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," *arXiv preprint arXiv:1611.01144*, 2016.
- [50] L. Chen, X. Chu, X. Zhang, and J. Sun, "Simple baselines for image restoration," pp. 17–33, 2022.
- [51] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [52] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [53] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," *arXiv preprint arXiv:1608.03983*, 2016.
- [54] Y. Yang and S. Newsam, "Bag-of-visual-words and spatial extensions for land-use classification," pp. 270–279, 2010.
- [55] G.-S. Xia, J. Hu, F. Hu, B. Shi, X. Bai, Y. Zhong, L. Zhang, and X. Lu, "Aid: A benchmark data set for performance evaluation of aerial scene classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 3965–3981, 2017.
- [56] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [57] S. Lei and Z. Shi, "Hybrid-scale self-similarity exploitation for remote sensing image super-resolution," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–10, 2021.
- [58] M. Qin, S. Mavromatis, L. Hu, F. Zhang, R. Liu, J. Sequeira, and Z. Du, "Remote sensing single-image resolution improvement using a deep gradient-aware network with image-specific enhancement," *Remote Sensing*, vol. 12, no. 5, p. 758, 2020.



**Yuting Lu** is pursuing the Ph.D. degree in control science and engineering from the Northwestern Polytechnic University, Xi'an, China. His research interests include image super-resolution and computer vision.



**Lingtong Min** received the B.S. degree from Northeastern University, Shenyang, China, in 2012, and the Ph.D. degree from Zhejiang University, Hangzhou, in 2019. He is an Associate Professor with Northwestern Polytechnical University. His main research interests are computer vision, pattern recognition, and remote sensing image understanding.



**Binglu Wang** (M'21) received the Ph.D. degree in Control Science and Engineering with the School of Automation at Northwestern Polytechnic University, Xi'an, China, in 2021. He is currently a Post-doctoral with the Department of Electrical Engineering, Beijing Institute of Technology, Beijing, China. His research interests include Computer Vision, Digital Signal Processing and Deep Learning.



**Le Zheng** (Senior Member, IEEE) received the B.Eng. degree from Northwestern Polytechnical University (NWPU), Xi'an, China, in 2009 and Ph.D degree from Beijing Institute of Technology (BIT), Beijing, China in 2015, respectively. He has previously held academic positions in the Electrical Engineering Department of Columbia University, New York, U.S., first as a Visiting Researcher from 2013 to 2014 and then as a Postdoc Research Fellow from 2015 to 2017. From 2018 to 2022, he worked at Aptiv (formerly Delphi), Los Angeles, as a Principal

Radar Systems Engineer, leading projects on the next-generation automotive radar products. Since July 2022, he has been a Full Professor with the School of Information and Electronics, BIT. His research interests lie in the general areas of radar, statistical signal processing, wireless communication, and high-performance hardware, and in particular in the area of automotive radar and integrated sensing and communications (ISAC).



**Xiaoxu Wang** (M'10) received the M.S. and Ph.D. degrees from the School of Automation, Harbin Engineering University, Harbin, China, in 2008 and 2010, respectively. He was as a Postdoctoral Researcher from 2010 to 2012, and as an Associate Professor from 2013 to October 2018 in Automation School of Northwestern Polytechnical University. He is currently a professor with the Northwestern Polytechnical University. His main research interests include deep learning, inertial navigation and non-linear estimation.



**Yongqiang Zhao** (M'05) received the B.S., M.S., and Ph.D. degrees in control science and engineering from the Northwestern Polytechnic University, Xi'an, China. From 2007 to 2009, he was as a Post-Doctora Researcher with McMaster University, Hamilton, ON, Canada, and Temple University, Philadelphia, PA, USA. He is currently a Professor with the Northwestern Polytechnical University. His research interests include polarization vision, hyperspectral imaging, and pattern recognition.



**Teng Long** (Fellow IEEE) was born in Fujian, China, in 1968. He received the M.S. and Ph.D. degrees in electrical engineering from the Beijing Institute of Technology, Beijing, China, in 1991 and 1995, respectively. He was a Visiting Scholar with Stanford University, California, in 1999, and University College London, in 2002. He has been a Full Professor with the Department of Electrical Engineering, Beijing Institute of Technology, since 2000. He has authored or co-authored more than 300 articles. His research interests include synthetic aperture

radar systems and real-time digital signal processing, with applications to radar and communication systems. Dr. Long is a Fellow of the Institute of Electronic and Technology and the Chinese Institute of Electronics. He was the recipient of many awards for his contributions to research and invention in China. He has been a member of the Chinese Engineering Academy since 2021.