

# Transforming Image Super-Resolution: A ConvFormer-based Efficient Approach

Gang Wu, Junjun Jiang, *Senior Member, IEEE*, Junpeng Jiang, and Xianming Liu, *Member, IEEE*

**Abstract**—Recent progress in single-image super-resolution (SISR) has achieved remarkable performance, yet the computational costs of these methods remain a challenge for deployment on resource-constrained devices. Especially for transformer-based methods, the self-attention mechanism in such models brings great breakthroughs while incurring substantial computational costs. To tackle this issue, we introduce the Convolutional Transformer layer (ConvFormer) and the ConvFormer-based Super-Resolution network (CFSR), which offer an effective and efficient solution for lightweight image super-resolution tasks. In detail, CFSR leverages the large kernel convolution as the feature mixer to replace the self-attention module, efficiently modeling long-range dependencies and extensive receptive fields with a slight computational cost. Furthermore, we propose an edge-preserving feed-forward network, simplified as EFN, to obtain local feature aggregation and simultaneously preserve more high-frequency information. Extensive experiments demonstrate that CFSR can achieve an advanced trade-off between computational cost and performance when compared to existing lightweight SR methods. Compared to state-of-the-art methods, e.g. ShuffleMixer, the proposed CFSR achieves  $0.39$  dB gains on Urban100 dataset for  $\times 2$  SR task while containing  $26\%$  and  $31\%$  fewer parameters and FLOPs, respectively. Code and pre-trained models are available at <https://github.com/Aitical/CFSR>.

**Index Terms**—Lightweight Image Super-Resolution, Large Kernel Convolution, Transformer, Self-attention.

## I. INTRODUCTION

**S**INGLE Image Super-Resolution (SISR) is fundamental task in computer vision that aims to enhance the resolution and quality of a single low-resolution image to a higher-resolution image. The goal is to generate a high-resolution image with fine details and improved visual perception from a single input image [1], [2]. The need for SISR arises from various real-world scenarios where high-resolution images are desired but limited by hardware capabilities or constraints. In many applications such as surveillance systems, medical imaging, satellite imagery, and digital photography, the acquisition of high-resolution images may be costly, time-consuming, or restricted. Therefore, SISR techniques provide a valuable solution by leveraging advanced algorithms and computational approaches to upsample low-resolution images.

In recent years, there has been considerable progress in the field of single-image super-resolution (SISR), largely attributed to the advent of deep learning techniques [2], [3]. The primary objective of SISR is to reconstruct a high-resolution image from its low-resolution counterpart, which

G. Wu, J. Jiang, J. Jiang, and X. Liu are with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China. E-mail: {gwu@hit.edu.cn, jiangjunjun@hit.edu.cn, 1190200226@stu.hit.edu.cn, csxm@hit.edu.cn}

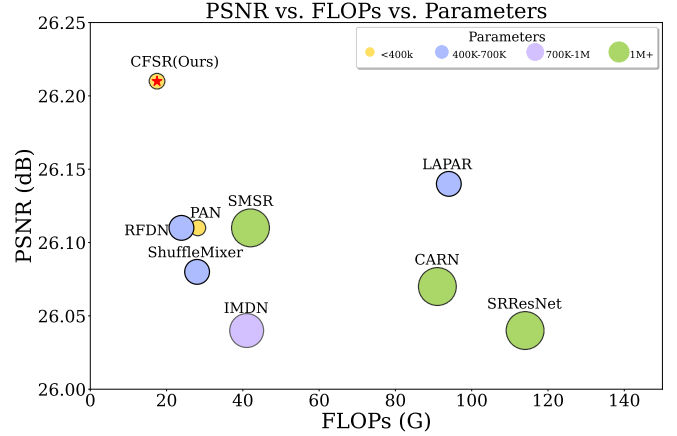


Fig. 1. Illustration of PSNR, FLOPs, and parameter counts of different SISR models on the Urban100 dataset for  $4\times$  SR task. The proposed CFSR approach achieves superior performance with less computational cost.

is a fundamental image processing task in the last decades. A groundbreaking study, SRCNN [4], introduced the concept of learning the mapping between low and high-resolution inputs using a convolutional neural network (CNN). This surpassed the performance of previous methods and led to further exploration of CNN-based methods. Subsequently, several studies have made great development of innovative SR models, such as deep and effective backbones [5]–[8] and attention mechanisms [9]–[12]. These CNN-based super-resolution methods have significantly advanced the state-of-the-art in SISR and demonstrate the power of deep learning and attention mechanisms in learning complex image representations and generating visually appealing high-resolution images from low-resolution inputs. However, despite the performance gains, these methods typically involve more complex models and higher computational complexity, which hinders their deployment on mobile and edge devices. To address this issue, the design of efficient and lightweight super-resolution models has become crucial. Many works have been proposed to reduce the amounts of parameters or floating point operations (FLOPs) to achieve lightweight models [13]–[21]. Zhao *et al.* [18] proposed the lightweight pixel-attention network (PAN) which replaces the standard residual or dense blocks with an efficient pixel-attention block. Sun *et al.* [21] proposed ShuffleMixer, which introduces the large kernel convolution into the lightweight SR network and significantly reduces the model complexity by the channel split-shuffle operation.

Recently, Transformer-based architectures have been proposed and attracted great attention in recent years due to their

impressive performance [22], [23]. Self-attention mechanism provides the promising long-range modeling and has achieved great breakthroughs in computer vision [22]. While the complexity of it is quadratic to the image size, which requires heavy computational resources. Liu *et al.* [23] proposed the Swin-Transformer, which performs self-attention within the large local window. Subsequently, many transformer-based SR methods have been advanced [24]–[27]. By leveraging self-attention mechanisms and hierarchical architectures, these methods can effectively obtain long-range relations and capture fine details. Transformer-based SR models showcase the promise of performance in addressing image super-resolution challenges that often surpass CNN-based models and achieve new state-of-the-art results. However, the application for lightweight models is limited due to the high computational cost substantial CPU or GPU memory requirements of self-attention mechanism.

In response to the challenge of computational efficiency in image super-resolution, we introduce a novel, self-attention-free approach that offers an excellent balance between computational cost and performance. This positions it as a viable solution for practical applications in lightweight image super-resolution. Specifically, we propose the Convolutional Transformer layer (ConvFormer) as a core component for effective and efficient feature extraction. Building on this foundation, we introduce the ConvFormer-based Super-Resolution network (CFSR) tailored for lightweight image super-resolution tasks. The transformer layer underpinning this approach, as outlined in [22], [23], [28], includes a feature mixer module and a feed-forward network. Drawing inspiration from the recent successes of CNN-based methods [29]–[32], our proposed feature mixer module employs large kernel convolutions as gate layers. This innovative design eliminates the need for self-attention in the feature mixer module, efficiently facilitating long-range relations and extensive receptive fields at a minimal additional computational cost. In addition, we introduce an edge-preserving feed-forward network (EFN) that refines the standard feed-forward network by incorporating enhanced edge extraction capabilities. Unlike the conventional feed-forward network (FFN) [22], which incorporates  $3 \times 3$  depth-wise convolutions for improved local feature aggregation in vision tasks [33], [34], our EFN integrates image gradient priors. This integration not only preserves high-frequency information but also introduces significant improvements for lightweight models without increasing complexity or parameter counts during inference, achieved through re-parameterization [19], [35], [36]. The architecture of CFSR, though straightforward and predominantly convolutional, is significantly more effective than previous methods. When benchmarked against existing methods for the  $\times 4$  SR task on the Urban100 dataset, as detailed in Figure 1, CFSR demonstrates superior performance. It notably excels in balancing reconstruction quality, model size, and computational efficiency, outperforming state-of-the-art methods with fewer parameters and reduced FLOPs.

We summarize the main contributions of our work as follows:

- 1) In traditional SISR models, the use of self-attention mechanisms for feature extraction has shown promising

results in enhancing super-resolution. However, it comes at a significant computational cost. To address this issue, this paper introduces ConvFormer, a feature mixer based on large kernel convolutions that replaces the self-attention module. ConvFormer efficiently captures long-range dependencies and extensive receptive fields while maintaining a lower computational complexity. This approach demonstrates superior performance and efficiency in lightweight image super-resolution tasks. It will shed new light on the design of CNN-based or hierarchical architecture for lightweight image super-resolution.

- 2) Conventional SISR algorithms often suffer from the loss of high-frequency information, resulting in a degradation of the quality of the generated super-resolved images. To tackle this challenge, this paper proposes EFN, an edge-preserving feed-forward network. EFN incorporates local feature aggregation through convolution layers with edge-preserving filters and preserves high-frequency information using skip connections and deconvolution layers. Compared to traditional methods, EFN achieves better preservation of image details and textures while maintaining high super-resolution performance.
- 3) It presents extensive experiments to verify the effectiveness of CFSR. Compared to the existing advanced methods, CFSR achieves superior performance with less computational cost. Detailed ablation studies are provided to analyze the impact of different components.

In the following section, we will first give some related work of lightweight image super-resolution methods and the progress of modern architectures in Section II. In Section III, we introduce and explain our proposed CFSR method in detail. Then, Section IV describes our training settings and experimental results including ablation analysis, where we compare the performance of our approach to other state-of-the-art methods. Finally, some conclusions are drawn in Section V.

## II. RELATED WORK

In this section, we will briefly introduce related literature including deep learning-based single-image super-resolution methods, transformer-based architectures, modern convolutional architectures, and the development of re-parameterizing methods.

### A. Single Image Super-Resolution

Recently, deep learning methods have achieved dramatic improvements in SISR tasks [2], [37], [38]. Especially for CNN-based models, various well-designed CNN architectures explore to further improve the SISR performance [5], [6], [39]. VDSR [5] introduces a very deep backbone to predict the residual construction between the LR input and the corresponding HR image. EDSR [6] incorporates residual blocks with skip connections, which allow direct propagation of information from earlier layers to later layers. Besides, attention mechanisms like the channel attention [40] has been introduced to SISR task as well [9]–[12]. Zhang *et al.* [9]

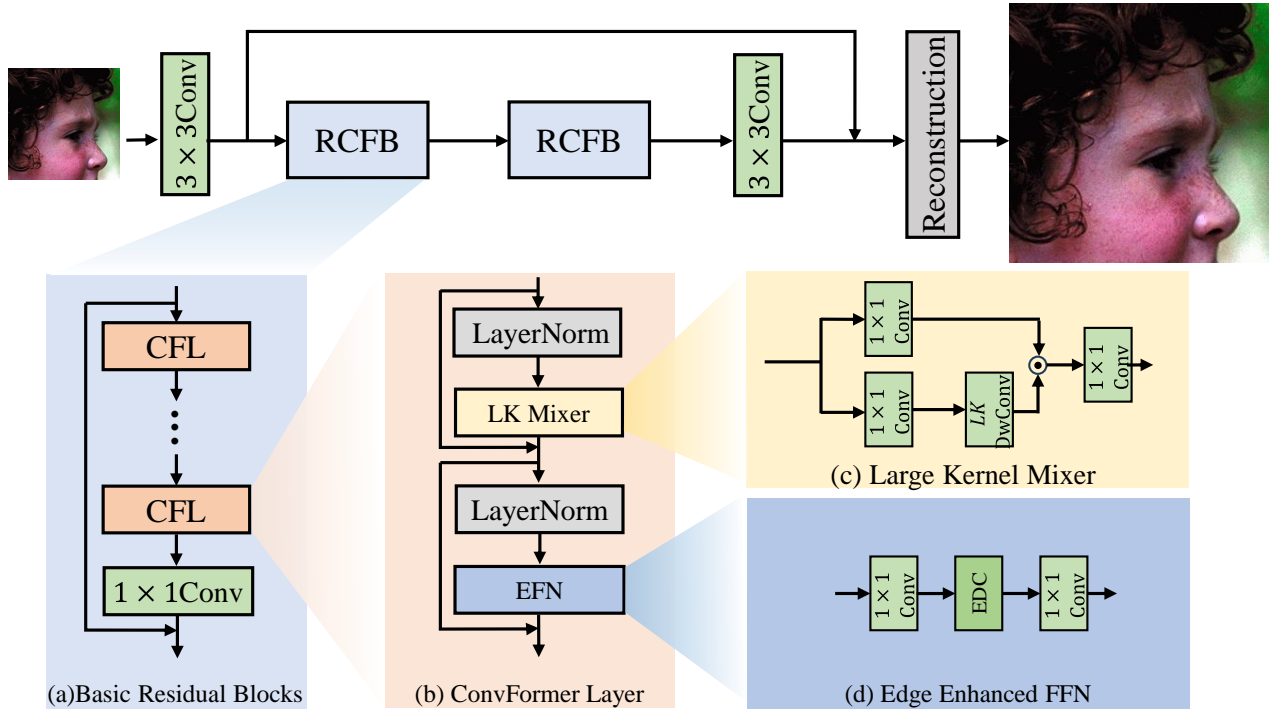


Fig. 2. Detailed implementation and different components in the proposed CFSR. The architecture of CFSR is mainly stacked by the basic residual block, which contains several ConvFormer layers. The ConvFormer Block plays a pivotal role, containing the proposed large kernel feature mixer (LK Mixer) and edge-preserving feed-forward network (EFN).

introduced the channel attention mechanism and proposed the RCAN model, which extend the backbone into over 400 layers.

In contrast to achieving advancing performance with a rapidly increased number of parameters and computational cost, many lightweight SISR models have been exploited by reducing parameters, especially for resource-limited devices [14]–[17], [19], [20], [41]. Hui *et al.* proposed a deep information distillation network (IDN) [41] and extended it into the information multi-distillation network (IMDN) [15], which won the AIM2020 challenge. Zhang *et al.* [19] proposed an edge-oriented Convolution Block (ECB) for real-time inference, which extracts 1st-order and 2nd-order spatial derivatives from intermediate features. Most recently, Sun *et al.* [21] introduced the large kernel convolution into the lightweight SR and proposed the ShuffleMixer. By channel split-shuffle operation, it efficiently reduced latent projection features.

In this paper, we proposed a pretty simple ConvNet for lightweight SR, named CFSR. Compared to the previous work [21], CFSR achieves  $0.39$  dB gains on Urban100 dataset for  $\times 2$  SR task while containing 26% and 31% fewer parameter counts and FLOPs, respectively.

### B. Transformer-based Architectures

Most recently, vision transformers have attracted great attention [22], [23] and many work have been proposed to explore transformer-based architectures for image restoration [24], [25], [27], [42], [43]. Various pre-trained models with full or local window-based attention are exploited and applied to target restoration tasks [24], [42]–[44]. Liang *et al.* [25] firstly introduced the Swin-Transformer into the image restoration

tasks and proposed a hierarchical architecture SwinIR. Cai *et al.* [44] proposed the hierarchical patch-based Transformer architecture, which significantly enhances single image super-resolution by progressively recovering high-resolution images through a hierarchy of patch partitions. This multi-stage model begins with small patch sizes, merging them into full resolution in later stages, and includes a unique attention-based position encoding and a multi-receptive field attention module. On the other hand, several works have been studied for lightweight transformer-based models [25], [45]–[47]. Wu *et al.* [45] proposed the lightweight model TCSR which introduces a sliding-window-based self-attention mechanism. One advantage of using transformers for SISR is their ability to capture global context information, which can be beneficial for generating high-quality HR images. However, compared to CNN-based models, transformer-based methods usually requires much more computational resources, even with a small model capacity, such as the SwinIR-light [25].

### C. Modern CNN-based Architectures

Several works investigated some modern CNN-based architectures [29], [31], [32], [48], [49]. On the one hand, large kernel convolution has been revisited [31], [32], [48]. Liu *et al.* [32] explored a modern CNN-based architecture and introduced larger kernels that utilize  $7 \times 7$  kernel size. Building upon this work, Ding *et al.* [31] further brought the kernel size up to 31. Subsequently, Liu *et al.* [48] further extended the kernel size up to 51 by sparse training. These advancements in kernel size have shown improvements in capturing complex image details and enhancing the super-resolution performance.

On the other hand, many work paid more attention to the hierarchical architecture combined convolution and transformer [29], [49]. These architectures leverage the strengths of both convolutional and transformer networks to capture local and global information effectively. Inspired by these findings, in this paper, we exploit a simple transformer-like ConvNet for lightweight SR tasks, where we replace the self-attention module with a large kernel-based mixer and improve the feed-forward network to preserve more high-frequency information. By leveraging the advantages of large kernel convolutions and transformer-like architectures, the proposed method holds potential for achieving better results in SISR while maintaining computational efficiency.

#### D. Re-parameterizing Methods

Ding *et al.* [35] proposed the RepVGG, which provides a practical network architecture and the concept of re-parameterizing operation. By re-parameterizing, multiple linear convolution learned in training process can be merged as one single convolution in inference instead of introducing extra costs. For SISR task, Wang *et al.* [50] proposed RepSR, a plain architecture for SISR task by re-parameterization. The authors analyzed the impact of the BatchNorm (BN) operation in SISR task and successfully re-introduced BN into SR. Zhang *et al.* [19] proposed the ECBSR, which introduces more 1st-order and 2nd-order gradient information into the vanilla convolution by re-parameterizing. Wang *et al.* [51] proposed DDistill-SR, which combines re-parameterization with dynamic convolution [52] to extend the learnable landscape while introducing less complexity in inference.

In this study, we introduce the Convolutional Transformer-based Super-Resolution network (CFSR), a simple yet effective model for lightweight image super-resolution. Drawing inspiration from large kernel methods like ConvNet [32] and re-parameterization strategies exemplified by RepVGG [35], our approach is particularly informed by developments in lightweight super-resolution models such as PAN [18] and ECBSR [19]. CFSR marks a departure from the pixel-attention reliance of PAN, revisiting and streamlining the convolution and self-attention feature extraction mechanisms. We introduce an advanced large kernel feature mixer, engineered to deliver exceptional performance and a significantly expanded receptive field. Additionally, our model capitalizes on the successes of Transformer-based architectures, integrating a novel edge-enhanced re-parameterization operation into the feed-forward network. This innovation not only enhances local feature extraction but also preserves a greater amount of high-frequency information, a significant advancement over the convolution-centric approach of ECBSR [19]. By extending the re-parameterizing branch and focusing on high-frequency detail retention, CFSR stands at the forefront of the current wave of convolutional technique advancements, uniquely tailoring these approaches for the intricate demands of super-resolution tasks.

### III. PROPOSED METHOD

In this paper, we propose a novel lightweight SISR method called CFSR network to leverage large kernel convolutions

as gate layers and replace the self-attention module present in transformers. This design enables efficient handling of long-range dependencies and extensive receptive fields while maintaining a lightweight computational cost. Additionally, we introduce the edge-preserving feed-forward network (EFN). EFN incorporates significant image gradient prior, thereby providing more high-frequency information. Furthermore, by re-parameterizing, EFN is free to improve the performance without any extra costs. In this section, we will present the detailed implementation of the propose CFSR.

#### A. Network Architecture

In Figure 2, we illustrate the proposed CFSR framework, which encompasses three pivotal stages: shallow feature extraction, deep feature extraction, and the image reconstruction module. The shallow feature extraction phase is designed to distill low-level image features, such as edges, textures, and fine-grained details, from the input image and map them into a latent space. These features are important for preserving the local structure and details of the image. In the deep feature extraction stage, the model extracts higher-level, more compact feature representations. It captures texture and structure information, which is essential for recovering lost details and enhancing the clarity of the image. Utilizing both shallow and deep features, the image reconstruction module is capable of generating high-resolution images. In the following, we will present the details of these three components.

**Shallow feature extraction.** Given a low-resolution (LR) input image  $I_{LR} \in \mathbb{R}^{H \times W \times 3}$ , where  $H$  and  $W$  are the height and width of this image. The shallow feature extraction utilizes a  $3 \times 3$  convolution layer to map  $I_{LR}$  into the latent feature space, It can be formulated as:

$$F_{sf} = H_{sf}(I_{LR}), \quad (1)$$

where  $H_{sf}(\cdot)$  denotes the convolutional layer for shallow feature extraction,  $F_0 \in \mathbb{R}^{H \times W \times C}$  is the output shallow feature, and  $C$  is the number of channels.

**Deep feature extraction.** Then we use a stack of two basic residual blocks (BRB), which contains several ConvFormer layers (CFL), and a  $3 \times 3$  convolution layer is added at the end of the BRB to aggregate the local features. Specifically, the detailed implementation of BRB and CFL are presented in Figure 2(a) and Figure 2(b), respectively. This process of deep feature extraction is formulated as:

$$F_k = \text{BRB}_k(F_{k-1}), \quad (2)$$

where  $\text{BRB}_k(\cdot)$  denotes the  $k$ -th BRB.  $F_{k-1}$  and  $F_k$  are the input feature and the output feature of the  $k$ -th BRB, respectively. Finally, the total deep feature extraction is:

$$F_{df} = H_{df}(F_{sf}), \quad (3)$$

where  $H_{df}(\cdot)$  presents the general deep feature extraction of the proposed CFSR network, and  $F_{df}$  presents the output of the deep backbone.

**Image reconstruction.** Image reconstruction module aims to reconstruct a high-resolution image  $I_{SR} \in \mathbb{R}^{rH \times rW \times 3}$  by aggregating both shallow and deep features as:

$$I_{SR} = H_{rec}(F_{sf} + F_{df}), \quad (4)$$

where  $r$  is a scale factor.  $H_{rec}(\cdot)$  represents the reconstruction module, which comprises of a  $3 \times 3$  convolution layer and a pixel-shuffle operation.

**Loss function.** The optimization of CFSR parameters is achieved through the minimization of the  $L_1$  pixel loss, which can be formulated as:

$$L_1 = \|I_{SR} - I_{HR}\|_1, \quad (5)$$

where  $I_{HR}$  is the corresponding ground-truth high-resolution image.

### B. ConvFormer Layer

This section introduces a streamlined, fully convolutional network backbone for lightweight computing, aimed at reducing computational complexity and memory usage. Firstly, we introduce a large kernel convolution-based feature mixer module, which requires less computational costs while provide effective large receptive fields. Secondly, we introduce the proposed EFN, which induces more high-frequency information for SR model.

TABLE I

COMPARISON OF COMPUTATIONAL COSTS BETWEEN GLOBAL SELF-ATTENTION (SA), LOCAL WINDOW SELF-ATTENTION (LWSA), AND THE PROPOSED LARGE KERNEL MIXER (LK), WHERE  $K$  IS THE WINDOW/KERNEL SIZE.

Module	Complexity	Parameters
SA	$\mathcal{O}(4HWC^2 + 2H^2W^2C)$	$4C^2$
LWSA	$\mathcal{O}(4HWC^2 + 2HWC^2K)$	$4C^2$
LK	$\mathcal{O}(3HWC^2 + HWC^2K)$	$3C^2 + CK^2$

**Large kernel mixer.** Self-attention is a powerful feature extractor, but its high computational cost makes it impractical for real-time application. Recent studies have demonstrated the effectiveness of employing large kernel convolutions [31]. In Table I, a detailed comparison of the computational complexity between multi-head self-attention (MHSA) [22], local window self-attention (LWSA) [23], and large kernel convolution (LK) [31] is exhibited.

Here the  $H, W, C$ , and  $K$  represent the height, width, channel dimension, and kernel (local window) size, respectively. The comparison reveals that, given identical window and kernel sizes, the computational complexity of the LK token mixer is significantly less than that of MHSA. The complexity of LK approximates to half that of LWSA, while simultaneously exhibiting a reduction in parameter count, where  $K \ll C$ . Consequently, large kernel convolutions offer a more resource-efficient choice for the design of lightweight models. In this paper, we propose a simple yet effective feature mixer module as presented in Figure 2(c). Here we take a  $1 \times 1$  convolution followed by a large kernel convolution as the feature mixing gate. Feature extraction of our LK mixer is formulated as follows:

$$\begin{aligned} V &= \text{Conv}_{1 \times 1}(F), \\ F_{gate} &= \text{Conv}_{1 \times 1}(\text{DwConv}_{k \times k}(F)), \\ F_{out} &= V \odot F_{gate}, \end{aligned} \quad (6)$$

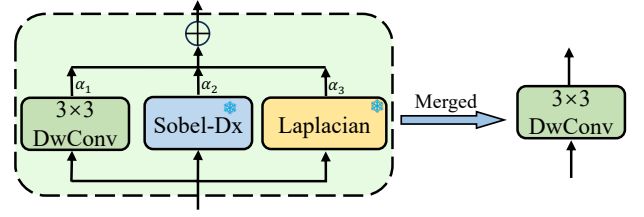


Fig. 3. Illustration of the edge-preserving depth-wise convolution (EDC). It contains a multi-branch structure with pre-defined gradient kernels and is equivalent to a single  $3 \times 3$  depth-wise convolution in inference by re-parameterizing.

where  $F$  is the input of the ConvFormer layer,  $F_{out}$  is output of it,  $\odot$  presents per-pixel production,  $k$  is the kernel size, and  $\text{DwConv}$  is depth-wise convolution.

**Edge-preserving feed-forward network.** In the feed-forward network (FFN) in a vanilla Transformer unit, previous studies have enhanced it by integrating depth-wise convolution, thereby improving the local feature ensemble [33], [34]. Considering that SISR, being an ill-posed problem, aims at learning an inversion from LR to HR, where the high-frequency information is crucial to it [53]. To obtain more high-frequency information in latent features, we propose an edge-preserving feed-forward network (EFN), by our Edge-preserving Depth-Wise Convolution (EDC), as illustrated in Figure 2(d). This allows a best of both words for a richer high-frequency information while maintaining local feature ensemble. The implementation of the proposed EFN is as follows:

$$F_1 = \text{Conv}_{1 \times 1}(F_{in}), \quad (7)$$

$$F_2 = \text{GELU}(F_1), \quad (8)$$

$$F_{EDC} = \text{EDC}(F_2), \quad (9)$$

$$F_3 = \text{ReLU}(F_{EDC}), \quad (10)$$

$$F_4 = \text{Conv}_{1 \times 1}(F_{EDC}), \quad (11)$$

where  $\text{GELU}(\cdot)$  is the activation function.

Detailed implementation of our EDC is presented in Figure 3. It takes a multi-branch structure, containing a standard depth-wise convolution (DwConv) and three DwConvs with pre-defined gradient kernels. Denote  $K_{3 \times 3} \in \mathbb{R}^{C \times 1 \times 3 \times 3}$  and  $B_{3 \times 3}$  the learnable kernel weights and bias for the vanilla DwConv, where  $C$  presents the output channels and 3 is the spatial size. The feature extraction is formulated as:

$$F_{3 \times 3} = K_{3 \times 3} * F_2 + B_{3 \times 3}, \quad (12)$$

where  $*$  presents the depth-wise convolution operation.

Next, we take the 1st-order and 2nd-order gradient kernels, such as Sobel filters and Laplacian filters. Denote the  $K_{D_x}$ ,  $K_{D_y}$  the horizontal and vertical Sobel filters:

$$K_{D_x} = \begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix}, \quad K_{D_y} = \begin{bmatrix} +1 & +2 & +1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}. \quad (13)$$

To align with the shape of kernel  $K_{3 \times 3}$  in depth-wise convolution, we simply expand and repeat the Sobel filters

into the  $C \times 1 \times 3 \times 3$  size. The 1st-order gradient of latent feature map is extracted as:

$$F_{Sobel} = K_{D_x} * F_2 + B_{D_x} + K_{D_y} * F_2 + B_{D_y}, \quad (14)$$

where  $B_{D_x}$  and  $B_{D_y}$  are biases.

Moreover, Laplacian filter  $K_{Lap}$  is utilized to extract 2nd-order gradient, where we take the 4-neighborhood and 8-neighborhood Laplacian operator as follows:

$$K_{Lap4} = \begin{bmatrix} 0 & +1 & 0 \\ +1 & -4 & +1 \\ 0 & +1 & 0 \end{bmatrix}, K_{Lap8} = \begin{bmatrix} +1 & +1 & +1 \\ +1 & -8 & +1 \\ +1 & +1 & +1 \end{bmatrix}, \quad (15)$$

and the same reshaping operation is adopted as aforementioned to extract the 2nd-order intermediate feature:

$$F_{Lap} = K_{Lap4} * F_2 + B_{Lap4} + K_{Lap8} * F_2 + B_{Lap8}. \quad (16)$$

The full feature extraction in EDC layer is:

$$F_{EDC} = \alpha_1 F_{3 \times 3} + \alpha_2 F_{Sobel} + \alpha_3 F_{Lap}, \quad (17)$$

where the parameters  $\alpha_1, \alpha_2, \alpha_3$  function as learnable competition coefficients for each branch. These coefficients are regulated by a straightforward softmax function, which aids in maintaining a higher retention of high-frequency feature information within the EFN framework.

**Merged EDC by re-parameterization in inference.** Following [19], [35], we can merge the multi-branch EDC layer into one single DwConv in inference without introducing additional complexity. Denote  $K$  and  $B$  the merged kernel weight and bias of the vanilla DwConv in inference. They can be achieved as follows:

$$K = \alpha_1 K_{3 \times 3} + \alpha_2 (K_{D_x} + K_{D_y}) + \alpha_3 (K_{Lap4} + K_{Lap8}), \quad (18)$$

$$B = \alpha_1 B_{3 \times 3} + \alpha_2 (B_{D_x} + B_{D_y}) + \alpha_3 (B_{Lap4} + B_{Lap8}). \quad (19)$$

Finally, we merge the five branches into one single DwConv operation by re-parameterizing, and the feature extraction of EDC layer in inference is:

$$F_{EDC} = K * F_2 + B. \quad (20)$$

## IV. EXPERIMENTS AND ANALYSIS

In this section we will describe the detailed evaluation experiments. Firstly, we introduce the experiment settings and comparison methods. Then quantitative and qualitative results are reported on some public datasets of SOTA lightweight methods and our proposed method. Lastly, to verify the technical contribution of the proposed method, we present the performance of different variants of the proposed method through some ablation studies.

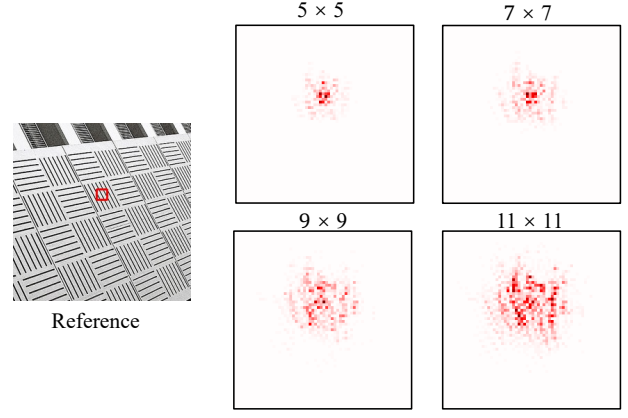


Fig. 4. LAM [54] attributions of different kernel sizes. From left to right, there is the reference input in the first column and LAM attributions of 5, 7, 9, 11 kernel sizes in the second and third columns.

### A. Experimental Setup

**Datasets and evaluation metrics.** Following comparison methods [16], [18], [21], [59], we train our model on the DIV2K [63] and Flickr2K [6] datasets, which contain 3450 high-quality images. We test the performance of CFSR on five benchmark test datasets, including Set5 [64], Set14 [65], BSD100 [66], Urban100 [67] and Manga109 [68]. We evaluate the Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM) on Y channel of transformed YCbCr space.

**Training details.** The channel number, RCFB number and CFL number in each RCFB are set to 48, 2 and 6, respectively. The sizes  $k_1$  and  $k_2$  of two depth-wise convolution in CFL are set to 9 and 3. During training, we randomly crop the image patches with the fixed size of  $64 \times 64$  and set the batch size to 16 for training. We employ randomly rotating  $90^\circ$ ,  $180^\circ$ ,  $270^\circ$  and horizontal flip for data augmentation. We use ADAM [69] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.99$  to optimize  $L_1$  loss. The initial learning rate is  $2e-4$ . The CFSR is implemented by PyTorch [70] and trained with Nvidia RTX A4000 GPU.

**Comparison Methods** We compare the proposed CFSR with state-of-the-art efficient SR approaches, including SR-CNN [4], VDSR [5], LapSRN [55], IDN [41], CARN [14], SRResNet [56], sLWSR [57], IMDN [15], LatticeNet [62], LAPAR [16], SMSR [17], ECBSR [19], DR SAN [58], PAN [18], DDistill [59], RFDN [60], ShuffleMixer [21], and some Transformer-based methods, including SwinIR-light [25] and ELAN [61].

### B. Main Results

The proposed CFSR achieves promising performance with less model complexity in both quantitative and qualitative results.

**Quantitative evaluation.** Table II presents quantitative comparisons for the upscaling factors of 2x, 3x, and 4x on five test datasets. Additionally, we also enumerate the parameter counts and FLOPs for each method. Remarkably, our proposed CFSR outperforms existing advanced CNN-based methods in terms of both PSNR and SSIM across all scales and datasets,



Fig. 5. Visual comparisons for SR( $\times 4$ ) methods on Set14, Manga109, and Urban100 datasets (**Zoom in for more details**).

TABLE II  
 QUANTITATIVE COMPARISON WITH SOME STATE-OF-THE-ART SR APPROACHES ON FIVE WIDELY USED BENCHMARK DATASETS. MULT-ADDS IS EVALUATED ON A  $1280 \times 720$  HR IMAGE. RESULTS OF OURS ARE IN **BOLD**. ONE CAN FIND THAT OUR CFSR ACHIEVES SOTA PERFORMANCE AMONG EXISTING ADVANCED CNN-BASED METHODS AND SIGNIFICANTLY BRIDGE THE GAP BETWEEN TRANSFORMER-BASED METHODS.

Method	Scale	Params	FLOPs	Set5	Set14	BSD100	Urban100	Manga109
				PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM
SRCNN [4]		8K	52.7G	36.66/0.9067	32.45/0.9067	31.36/0.8879	29.50/0.8946	35.60/0.9663
VDSR [5]		666K	612.6G	37.53/0.9587	33.03/0.9124	31.90/0.8960	30.76/0.9140	37.22/0.9750
LapSRN [55]		251K	29.9G	37.52/0.9591	32.99/0.9124	31.80/0.8952	30.41/0.9103	37.27/0.9740
IDN [41]		553K	124.6G	37.83/0.9600	33.30/0.9148	32.08/0.8985	31.27/0.9196	38.01/0.9749
CARN [14]		1592K	222.8G	37.76/0.9590	33.52/0.9166	32.09/0.8978	31.92/0.9256	38.36/0.9765
SRResNet [56]		1,370K	341.7G	38.05/0.9607	33.64/0.9178	32.22/0.9002	32.23/0.9295	38.05/0.9607
sLWSR [57]		322K	20.6G	31.87/0.8910	28.06/0.7740	27.46/0.7320	25.65/0.772	-/
IMDN [15]		694K	158.8G	38.00/0.9605	33.63/0.9177	32.19/0.8996	32.17/0.9283	38.88/0.9774
LAPAR-A [16]	×2	548K	171.0G	38.01/0.9605	33.62/0.9183	32.19/0.8999	32.10/0.9283	38.67/0.9772
SMSR [17]		985K	131.6G	38.00/0.9601	33.64/0.9179	32.17/0.8990	32.19/0.9284	38.76/0.9771
ECBSR [19]		596K	137.3G	37.90/0.9615	33.34/0.9178	32.10/0.9018	31.71/0.9250	-/
PAN [18]		261K	70.5G	38.00/0.9605	33.59/0.9181	32.18/0.8997	32.01/0.9273	38.70/0.9773
DRSAN [58]		370K	85.5G	37.99/0.9606	33.57/0.9177	32.16/0.8999	32.10/0.9279	-/
DDistill-SR [59]		414K	128.0G	38.03/0.9606	33.61/0.9182	32.19/0.9000	32.18/0.9286	38.94/0.9777
RFDN [60]		534K	95.0G	38.05/0.9606	33.68/0.9184	32.16/0.8994	32.12/0.9278	38.88/0.9773
ShuffleMixer [21]		394K	91.0G	38.01/0.9606	33.63/0.9180	32.17/0.8995	31.89/0.9257	38.83/0.9774
<b>CFSR (Ours)</b>		291K	62.6G	<b>38.07/0.9607</b>	<b>33.74/0.9192</b>	<b>32.24/0.9005</b>	<b>32.28/0.9300</b>	<b>39.00/0.9778</b>
SwinIR-light [25]		878K	195.6G	38.14/0.9611	33.86/0.9206	32.31/0.9012	32.76/0.9340	39.12/0.9783
ELAN-light [61]		582K	168.4G	38.17/0.9611	33.94/0.9207	32.30/0.9012	32.76/0.9340	39.11/0.9782
VDSR [5]		666K	612.6G	33.66/0.9213	29.77/0.8314	28.82/0.7976	27.14/0.8279	32.01/0.9340
LapSRN [55]		502K	90.0G	33.81/0.9220	29.79/0.8325	28.82/0.7980	27.07/0.8275	32.21/0.9350
CARN [14]		1,592K	118.8G	34.29/0.9255	30.29/0.8407	29.06/0.8034	28.06/0.8493	33.50/0.9440
SRResNet [56]		1,554K	190.2G	34.41/0.9274	30.36/0.8427	29.11/0.8055	28.20/0.8535	33.54/0.9448
IMDN [15]		703K	72.0G	34.36/0.9270	30.32/0.8417	29.09/0.8046	28.17/0.8519	33.61/0.9445
LatticeNet [62]		765K	76.3G	34.40/0.9272	30.32/0.8416	29.10/0.8049	28.19/0.8513	-/
LAPAR-A [16]	×3	594K	114.0G	34.36/0.9267	30.34/0.8421	29.11/0.8054	28.15/0.8523	33.51/0.9441
SMSR [17]		993K	67.8G	34.40/0.9270	30.33/0.8412	29.10/0.8050	28.25/0.8536	33.68/0.9445
PAN [18]		261K	39.0G	34.40/0.9271	30.36/0.8423	29.11/0.8050	28.11/0.8511	33.61/0.9448
DRSAN [58]		410K	43.2G	34.41/0.9272	30.27/0.8413	29.08/0.8056	28.19/0.8529	-/
Distill-SR [59]		414K	57.4G	34.37/0.9275	30.34/0.8420	29.11/0.8053	28.19/0.8528	33.69/0.9451
ShuffleMixer [21]		415K	43.0G	34.40/0.9272	30.37/0.8423	29.12/0.8051	28.08/0.8498	33.69/0.9448
<b>CFSR (Ours)</b>		298K	28.5G	<b>34.50/0.9279</b>	<b>30.44/0.8437</b>	<b>29.16/0.8066</b>	<b>28.29/0.8553</b>	<b>33.86/0.9462</b>
SwinIR-light [25]		886K	87.2G	34.62/0.9289	30.54/0.8463	29.20/0.8082	28.66/0.8624	33.98/0.9478
ELAN-light [61]		590K	75.7G	34.61/0.9288	30.55/0.8463	29.21/0.8081	28.69/0.8624	34.00/0.9478
SRCNN [4]		8K	52.7G	30.48/0.8626	27.50/0.7513	26.90/0.7101	24.52/0.7221	27.58/0.8555
VDSR [5]		666K	612.6G	31.35/0.8838	28.01/0.7674	27.29/0.7251	25.18/0.7524	28.83/0.8870
LapSRN [55]		813K	149.4G	31.54/0.8852	28.09/0.7700	27.32/0.7275	25.21/0.7562	29.09/0.8900
IDN [41]		553K	32.3G	31.82/0.8903	28.25/0.7730	27.41/0.7297	25.41/0.7632	29.41/0.8942
CARN [14]		1592K	90.9G	32.13/0.8937	28.60/0.7806	27.58/0.7349	26.07/0.7837	30.47/0.9084
SRResNet [56]		1,518K	146G	32.17/0.8951	28.61/0.7823	27.59/0.7365	26.12/0.7871	30.48/0.9087
IMDN [15]		715K	40.9G	32.21/0.8948	28.58/0.7811	27.56/0.7353	26.04/0.7838	30.45/0.9075
LatticeNet [62]		777K	43.6G	32.18/0.8943	28.61/0.7812	27.57/0.7355	26.14/0.7844	-/
LAPAR-A [16]	×4	659K	94.0G	32.15/0.8944	28.61/0.7818	27.61/0.7366	26.14/0.7871	30.42/0.9074
SMSR [17]		1006K	41.6G	32.12/0.8932	28.55/0.7808	27.55/0.7351	26.11/0.7868	30.54/0.9085
ECBSR [19]		603K	34.7G	31.92/0.8946	28.34/0.7817	27.48/0.7393	25.81/0.7773	-/
PAN [18]		272K	28.2G	32.13/0.8948	28.61/0.7822	27.59/0.7363	26.11/0.7854	30.51/0.9095
DRSAN [58]		410K	30.5G	32.15/0.8935	28.54/0.7813	27.54/0.7364	26.06/0.7858	-/
DDistill-SR [59]		434K	33.0G	32.23/0.8960	28.62/0.7823	27.58/0.7365	26.20/0.7891	30.48/0.9090
RFDN [60]		550K	23.9G	32.24/0.8952	28.61/0.7819	27.57/0.7360	26.11/0.7858	30.58/0.9089
ShuffleMixer [21]		411K	28.0G	32.21/0.8953	28.66/0.7827	27.61/0.7366	26.08/0.7835	30.65/0.9093
<b>CFSR (Ours)</b>		307K	17.5G	<b>32.33/0.8964</b>	<b>28.73/0.7842</b>	<b>27.63/0.7381</b>	<b>26.21/0.7897</b>	<b>30.72/0.9111</b>
SwinIR-light [25]		897K	49.6G	32.44/0.8976	28.77/0.7858	27.69/0.7406	26.47/0.7980	30.92/0.9151
ELAN-light [61]		601K	43.2G	32.43/0.8975	28.78/0.7858	27.69/0.7406	26.54/0.7982	30.92/0.9150





Fig. 6. LAM [54] comparisons between SOTA methods and the proposed CFSSR. Results of two samples from Urban100 and B100 datasets are presented. One can find that the proposed CFSSR outperforms other advanced models with larger receptive fields and richer textures (**Zoom in for more details**).

TABLE III  
ABLATION ON THE SIZE OF LARGE KERNEL CONVOLUTION FOR  $\times 4$  SR.  
WE TEST THE RESULTS ON URBAN100 AND MANGA109 DATASETS.

Kernel Size	Params	FLOPs	Urban100 PSNR/SSIM	Manga109 PSNR/SSIM
$5 \times 5$	274K	15.6G	26.02/0.7834	30.47/0.9082
$7 \times 7$	288K	16.4G	26.08/0.7854	30.54/0.9089
$9 \times 9$	307K	17.5G	26.13/0.7875	30.60/0.9098
$11 \times 11$	330K	18.9G	26.16/0.7880	30.64/0.9102

and significantly bridges the gap between Transformer-based methods [25], [61]. In detail, when compared to DDistill-SR [59], CFSSR maintains superior performance across all scales and datasets, while exhibiting approximately half the computational complexity. Furthermore, let us focus on the 3x super-resolution tasks. In comparison to ShuffleMixer [21], CFSSR attains significant performance gains, enhancing PSNR by **0.19** and **0.17**, and SSIM by **0.0055** and **0.0014** on the Urban100 and Manga109 test datasets, respectively. Importantly, CFSSR achieves these results with reduced parameter counts and computational costs, underscoring its efficiency and effectiveness.

We believe that these results corroborate the potential of CFSSR as a resource-efficient and performance-oriented model for lightweight image super-resolution tasks.

**Qualitative evaluation.** We conducted a visual quality comparison of SR results between our proposed CFSSR and five representative models, including CARN [14], IMDN [15], RFDN [60], PAN [18], and ShuffleMixer [21]. The  $\times 4$  SR results are presented in Figure 5. When we take a closer look at the results in the second column, one can find that our CFSSR is able to recover the main structures with sharp textures. Moreover, results of samples 'img\_058' and 'img\_024' in Urban100 dataset showcase that CFSSR can obtain clearer edges while others fail.

Furthermore, we use LAM [54], which represents the range of attribution pixels, to visualize receptive fields. Visual results are presented in Figure 6, showing that our CFSSR can take advantage of a wider range of information than CARN, IMDN, PAN, ShuffleMixer, obtaining large receptive fields effectively. Let us take the 'Tiger' image in the second row of Figure 6 as the example. One can find that the proposed CFSSR can achieve richer textures with the larger receptive field compared to other advanced methods.

**Comparison on unknown degradation.** Given that the primary objective of image Super-Resolution (SR) is to address complex real-world degradations and generate visually appealing images, we conduct comprehensive evaluation on the RealSR [71] and unknown DIV2K datasets, known for their intricate degradation patterns. For a fair comparison, we retrain and evaluate the ShuffleMixer model and our CFSSR

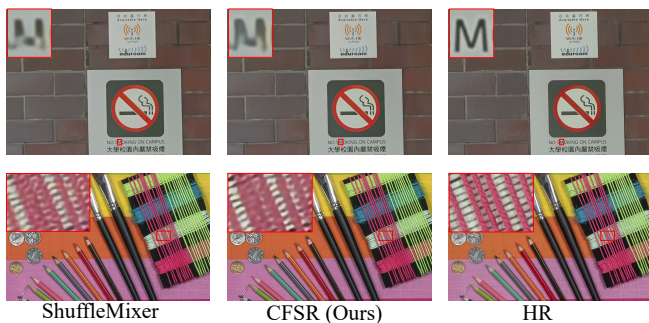


Fig. 7. Visual comparison on RealSR dataset. Super-resolved results of our CFSR can achieve more accurate textures even with complex degradation.

on these datasets, respectively. The results, as detailed in Table IV, clearly demonstrate that our CFSR model outperforms the advanced ShuffleMixer in handling complex degradation tasks.

Furthermore, in subjective image quality assessments, the CFSR model exhibits superior performance, and some results are presented in Figure 7. One can find that our CFSR is able to accurately reconstruct finer details and maintain a high level of clarity, even in areas of intricate patterns and textures. We believe these comparisons showcase the practical effectiveness of the CFSR model in real-world applications.

TABLE IV  
COMPARISON ON UNKNOWN DEGRADATION. WE EVALUATE THE PROPOSED CFSR AND SHUFFLEMIXER [21] ON REALSR TEST DATASET AND UNKNOWN DIV2K EVALUATION DATASET.

Method	Params.	FLOPs	RealSR	DIV2K unknown
ShuffleMixer	411K	28G	29.16/0.8261	29.17/0.8049
<b>CFSR</b>	<b>307K</b>	<b>17.5G</b>	<b>29.25/0.8266</b>	<b>29.39/0.8111</b>

### C. Ablation Studies

In this section, we conduct in-depth ablation studies on the core component of CFSR, the ConvFormer layer. The ConvFormer layer consists of two main elements: the Large Kernel Feature Mixer (LK Mixer) and the Edge-preserving Feed-forward Network (EFN). Each of these is ablated separately to elucidate their individual impact on the overall model performance. Specifically, we analyze the influence of various kernel sizes in the LK Mixer and examine the effect of our proposed Edge-preserving Depthwise Convolution (EDC) within the EFN.

TABLE V  
COMPARISON OF FFN WITH OR WITHOUT THE PROPOSED EDC FOR  $\times 4$  SR. WE REPORT THE RESULTS ON URBAN100 AND MANGA109 DATASETS.

Method	Params	FLOPs	Urban100 PSNR/SSIM	Manga109 PSNR/SSIM
w/o	307K	17.5G	26.13/0.7875	30.60/0.9098
w	307K	17.5G	<b>26.21/0.7897</b>	<b>30.72/0.9111</b>

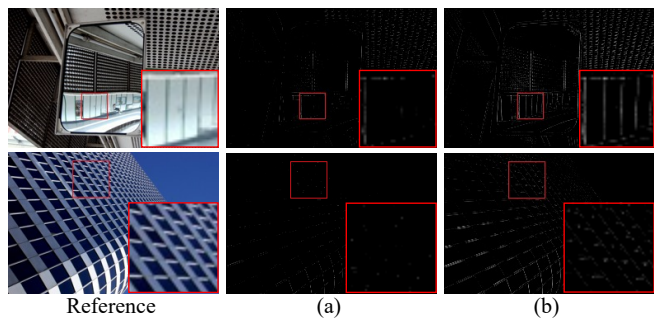


Fig. 8. Visual comparisons between latent features learned with or without the EDC (correspond to the ablation results in V). The first column is the reference image sampled from Urban100 dataset. (a) Latent features learned by the vanilla DwConv. (b) Feature maps extracted by our EFN with the EDC layer. One can find that EFN can substantially obtain clear and robust high-frequency information in latent feature maps.

**Impact of the kernel size.** Our in-depth analysis, as presented in Table III, explores the implications of diverse kernel sizes on model performance. The results explicitly underscore that model performance escalates with increasing kernel size. Meanwhile, we present corresponding visualization of the activating pixels by LAM [54] in Figure 4, intuitively showcase the superiority of larger kernels in effectively extending the receptive field.

Let us take a closer look at the results in Table III. Despite larger kernel sizes yielding superior performance, this comes with an associated surge in computational demands. Therefore, we limited our investigation to kernel sizes not exceeding 11, owing to both an evident saturation effect and computational constraints. In detail, a progressive enhancement in PSNR/SSIM values of over 0.05dB/0.001 was recorded as the kernel size grew from 5 to 9 across both Urban100 and Manga109 datasets. This consistent gain emphasizes the potency of large kernels in improving the performance. However, the advancement from a 9 to 11 kernel size revealed a more subdued gain of only 0.03/0.0005 in PSNR/SSIM. Given this negligible enhancement and the saturation observed in escalating kernel size benefits, we determined 9 as the default size for the CFSR. Furthermore, LAM results of different kernel sizes in LK Mixer are presented in Figure 4. One can find that the proposed LK Mixer in CFSR can effectively provide large receptive fields.

**Impact of EDC.** The efficacy of EDC within the Edge-preserving Feed-forward Network (EFN) is evident from the results presented in Table V. The effectiveness of the EDC within the EFN is demonstrated in Table V. Importantly, the employment of re-parameterization ensures that the model with EDC does not incur additional complexity during the inference phase. To provide a more comprehensive understanding, we intensify our analysis by visualizing the latent feature maps, as presented in Figure 8. Fascinatingly, the incorporation of EDC yields a pronounced improvement in the extraction of high-frequency information within intermediate features. From Figure 8, one can find that there are more edges and textures in latent feature map of our EFN. This observation significantly demonstrates the effectiveness of our enhanced FFN by EDC.

## V. CONCLUSION

The field of Single-Image Super-Resolution (SISR) has undergone significant advancements in recent years, largely attributed to the deep learning methods. In this paper we propose a Transformer-like, convolutional network for lightweight SR tasks, named CFSR, which achieves an advanced performance. CFSR employs a large kernel feature mixer (LK Mixer) as an efficient substitute for the complicate self-attention module, enabling the modeling of extensive receptive fields while significantly reducing computational overhead. Moreover, we propose an edge-preserving feed-forward network (EFN) to extract local features while preserving high-frequency information. To enhance the edge-preserving capabilities of our network, we propose the edge-preserving depth-wise convolution (EDC), which provides richer high-frequency information without introducing extra complexity during inference by a re-parameterizing strategy. To understand the influence of these components, detailed ablation studies are provided. Comprehensive experiments highlight the superior performance of CFSR, showcasing its ability to outclass existing state-of-the-art methods while maintaining a lean profile in terms of parameter count and computational complexity. Overall, our proposed method holds great potential for advancing the field of SISR and enabling the deployment of super-resolution algorithms on resource-constrained devices.

## REFERENCES

- [1] Wenming Yang, Xuechen Zhang, Yapeng Tian, Wei Wang, Jing-Hao Xue, and Qingmin Liao. Deep learning for single image super-resolution: A brief review. *IEEE Trans. Multim.*, 21(12):3106–3121, 2019. 1
- [2] Juncheng Li, Zehua Pei, and Tiejong Zeng. From beginner to master: A survey for deep learning-based single-image super-resolution. *arXiv preprint arXiv:2109.14335*, 2021. 1, 2
- [3] Zhihao Wang, Jian Chen, and Steven C. H. Hoi. Deep learning for image super-resolution: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(10):3365–3387, 2021. 1
- [4] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2016. 1, 6, 8
- [5] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 6, 8
- [6] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017. 1, 2, 6
- [7] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [8] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Deep back-projection networks for super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [9] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, 2018. 1, 2
- [10] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [11] Ben Niu, Weilei Wen, Wenqi Ren, Xiangde Zhang, Lianping Yang, Shuzhen Wang, Kaihao Zhang, Xiaochun Cao, and Haifeng Shen. Single image super-resolution via a holistic attention network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2
- [12] Yiqun Mei, Yuchen Fan, and Yuqian Zhou. Image super-resolution with non-local sparse attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2
- [13] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 1
- [14] Namhyuk Ahn, Byungkon Kang, and Kyung-Ah Sohn. Fast, accurate, and lightweight super-resolution with cascading residual network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 1, 3, 6, 8, 9
- [15] Zheng Hui, Xinbo Gao, Yunchu Yang, and Xiumei Wang. Lightweight image super-resolution with information multi-distillation network. In *Proceedings of the ACM International Conference on Multimedia (ACM MM)*, 2019. 1, 3, 6, 8, 9
- [16] Wenbo Li, Kun Zhou, Lu Qi, Nianjuan Jiang, Jiangbo Lu, and Jiaya Jia. LAPAR: linearly-assembled pixel-adaptive regression network for single image super-resolution and beyond. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1, 3, 6, 8
- [17] Longguang Wang, Xiaoyu Dong, Yingqian Wang, Xinyi Ying, Zaiping Lin, Wei An, and Yulan Guo. Exploring sparsity in image super-resolution for efficient inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 3, 6, 8
- [18] Hengyuan Zhao, Xiangtao Kong, Jingwen He, Yu Qiao, and Chao Dong. Efficient image super-resolution using pixel attention. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2020. 1, 4, 6, 8, 9
- [19] Xindong Zhang, Hui Zeng, and Lei Zhang. Edge-oriented convolution block for real-time super resolution on mobile devices. In *Proceedings of the ACM International Conference on Multimedia (ACM MM)*, 2021. 1, 2, 3, 4, 6, 8
- [20] Guangwei Gao, Wenjie Li, Juncheng Li, Fei Wu, Huimin Lu, and Yi Yu. Feature distillation interaction weighting network for lightweight image super-resolution. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2022. 1, 3
- [21] Long Sun, Jinshan Pan, and Jinhui Tang. Shufflemixer: An efficient convnet for image super-resolution. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:17314–17326, 2022. 1, 3, 6, 8, 9, 10
- [22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. 2, 3, 5
- [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2, 3, 5
- [24] Wenbo Li, Xin Lu, Jiangbo Lu, Xiangyu Zhang, and Jiaya Jia. On efficient transformer and image pre-training for low-level vision. *CoRR*, abs/2112.10175, 2021. 2, 3
- [25] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. SwinIR: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2021. 2, 3, 6, 8, 9
- [26] Zheng Chen, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. Cross aggregation transformer for image restoration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2
- [27] Yawei Li, Yuchen Fan, Xiaoyu Xiang, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. Efficient and explicit modelling of image hierarchies for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 3
- [28] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10809–10819, 2022. 2
- [29] Qibin Hou, Chengze Lu, Ming-Ming Cheng, and Jiashi Feng. Conv2former: A simple transformer-style convnet for visual recognition. *CoRR*, abs/2211.11943, 2022. 2, 3, 4
- [30] Meng-Hao Guo, Chengze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shimin Hu. Visual attention network. *CoRR*, abs/2202.09741, 2022. 2
- [31] Xiaohan Ding, Xiangyu Zhang, Yizhuang Zhou, Jungong Han, Guiguang Ding, and Jian Sun. Scaling up your kernels to 31x31: Revisiting large

- kernel design in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 5
- [32] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 4
- [33] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17662–17672, 2022. 2, 5
- [34] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5718–5729, 2022. 2, 5
- [35] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13733–13742, 2021. 2, 4, 6
- [36] Xiaohan Ding, Yuchen Guo, Guiguang Ding, and Jungong Han. Acnet: Strengthening the kernel skeletons for powerful CNN via asymmetric convolution blocks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1911–1920, 2019. 2
- [37] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021. 2
- [38] Saeed Anwar, Salman H. Khan, and Nick Barnes. A deep journey into super-resolution: A survey. *ACM Comput. Surv.*, 2020. 2
- [39] Ying Tai, Jian Yang, and Xiaoming Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [40] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [41] Zheng Hui, Xiumei Wang, and Xinbo Gao. Fast and accurate single image super-resolution via information distillation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3, 6, 8
- [42] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [43] Xiangyu Chen, Xintao Wang, Jiantao Zhou, and Chao Dong. Activating more pixels in image super-resolution transformer. *arXiv preprint arXiv:2205.04437*, 2022. 3
- [44] Qing Cai, Yiming Qian, Jinxing Li, Jun Lyu, Yee-Hong Yang, Feng Wu, and David Zhang. Hipa: Hierarchical patch transformer for single image super resolution. *IEEE Transactions on Image Processing*, 32:3226–3237, 2023. 3
- [45] Gang Wu, Junjun Jiang, Yuanhao Bai, and Xianming Liu. Incorporating transformer designs into convolutions for lightweight image super-resolution. *CoRR*, abs/2303.14324, 2023. 3
- [46] Hang Wang, Xuanhong Chen, Bingbing Ni, Yutian Liu, and Jinfan Liu. Omni aggregation networks for lightweight image super-resolution. *CoRR*, abs/2304.10244, 2023. 3
- [47] Haram Choi, Jeongmin Lee, and Jihoon Yang. N-gram in swin transformers for efficient lightweight image super-resolution. *CoRR*, abs/2211.11436, 2022. 3
- [48] Shiwei Liu, Tianlong Chen, Xiaohan Chen, Xuxi Chen, Qiao Xiao, Boqian Wu, Mykola Pechenizkiy, Decebal Constantin Mocanu, and Zhangyang Wang. More convnets in the 2020s: Scaling up kernels beyond 51x51 using sparsity. *CoRR*, abs/2207.03620, 2022. 3
- [49] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, pages 22–31. IEEE, 2021. 3, 4
- [50] Xintao Wang, Chao Dong, and Ying Shan. Repsr: Training efficient vgg-style super-resolution networks with structural re-parameterization and batch normalization. In *ACM Multimedia*, pages 2556–2564. ACM, 2022. 4
- [51] Yan Wang, Tongtong Su, Yusen Li, Jiuwen Cao, Gang Wang, and Xiaoguang Liu. Ddistill-sr: Reparameterized dynamic distillation network for lightweight image super-resolution. *IEEE Transactions on Multimedia*, pages 1–13, 2022. 4
- [52] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *CVPR*, pages 11027–11036. Computer Vision Foundation / IEEE, 2020. 4
- [53] Faming Fang, Juncheng Li, and Tiejiong Zeng. Soft-edge assisted network for single image super-resolution. *IEEE Transactions on Image Processing*, 29:4656–4668, 2020. 5
- [54] Jinjin Gu and Chao Dong. Interpreting super-resolution networks with local attribution maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 6, 9, 10
- [55] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep Laplacian pyramid networks for fast and accurate super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 8
- [56] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew P. Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 8
- [57] Biao Li, Bo Wang, Jiabin Liu, Zhiquan Qi, and Yong Shi. s-lwrs: Super lightweight super-resolution network. *IEEE Transactions on Image Processing*, 29:8368–8380, 2020. 6, 8
- [58] Karam Park, Jae Woong Soh, and Nam Ik Cho. A dynamic residual self-attention network for lightweight single image super-resolution. *IEEE Transactions on Multimedia*, 25:907–918, 2023. 6, 8
- [59] Yan Wang, Tongtong Su, Yusen Li, Jiuwen Cao, Gang Wang, and Xiaoguang Liu. Ddistill-sr: Reparameterized dynamic distillation network for lightweight image super-resolution. *IEEE Transactions on Multimedia*, pages 1–13, 2022. 6, 8, 9
- [60] Jie Liu, Jie Tang, and Gangshan Wu. Residual feature distillation network for lightweight image super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2020. 6, 8, 9
- [61] Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang. Efficient long-range attention network for image super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, volume 13677 of *Lecture Notes in Computer Science*, pages 649–667. Springer, 2022. 6, 8, 9
- [62] Xiaotong Luo, Yuan Xie, Yulun Zhang, Yanyun Qu, Cuihua Li, and Yun Fu. Latticenet: Towards lightweight image super-resolution with lattice block. In *Proceedings of the European Conference on Computer Vision (ECCV)*, volume 12367 of *Lecture Notes in Computer Science*, pages 272–289. Springer, 2020. 6, 8
- [63] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, July 2017. 6
- [64] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie-Line Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2012. 6
- [65] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *Curves and Surfaces*, 2010. 6
- [66] David R. Martin, Charles C. Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2001. 6
- [67] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 6
- [68] Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multim. Tools Appl.*, 2017. 6
- [69] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 6
- [70] Adam Paszke, Sam Gross, Francisco Massa, et al. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8024–8035, 2019. 6
- [71] Jianrui Cai, Huiyu Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3086–3095, 2019. 9