

Cross-SRN: Structure-Preserving Super-Resolution Network with Cross Convolution

Yuqing Liu, Qi Jia, Xin Fan, *Senior Member, IEEE*, Shanshe Wang, Siwei Ma, *Member, IEEE*,
and Wen Gao, *Fellow, IEEE*

Abstract—It is challenging to restore low-resolution (LR) images to super-resolution (SR) images with correct and clear details. Existing deep learning works almost neglect the inherent structural information of images, which acts as an important role for visual perception of SR results. In this paper, we design a hierarchical feature exploitation network to probe and preserve structural information in a multi-scale feature fusion manner. First, we propose a cross convolution upon traditional edge detectors to localize and represent edge features. Then, cross convolution blocks (CCBs) are designed with feature normalization and channel attention to consider the inherent correlations of features. Finally, we leverage multi-scale feature fusion group (MFFG) to embed the cross convolution blocks and develop the relations of structural features in different scales hierarchically, invoking a lightweight structure-preserving network named as Cross-SRN. Experimental results demonstrate the Cross-SRN achieves competitive or superior restoration performances against the state-of-the-art methods with accurate and clear structural details. Moreover, we set a criterion to select images with rich structural textures. The proposed Cross-SRN outperforms the state-of-the-art methods on the selected benchmark, which demonstrates that our network has a significant advantage in preserving edges.

Index Terms—Image super-resolution, cross convolution, multi-scale feature fusion, structure-preservation.

I. INTRODUCTION

GIVEN a low-resolution (LR) image, the task of super-resolution (SR) aims to find the corresponding high-resolution (HR) instance with refined details. Image SR has been considered in numerous computer vision tasks, such as recognition, person Re-Identification, semantic segmentation, and video compression.

As a highly ill-posed issue, image SR suffers from diverse degradation models, such as down-sampling, noise, and blur. Convolutional neural network (CNN) has shown its superior performance on complex information restoration [1], which is widely used by recent image SR works. SRCNN [2] is the first CNN-based image SR method with a three-layer network. Recently, some well-designed architectures, such as

DRN [3], HAN [4], and LatticeNet [5], achieve the state-of-the-art performances by building deeper or wider networks to explore features more effectively. However, they almost neglect the inherent structural information of images, such as the outer contour of objects, lines, and curves, which is a vital factor to evaluate the restoration quality.

Gradient information in fixed direction is crucial to detect structural information. Traditional edge detectors, such as Prewitt and Sobel [6], design filter-like templates to explore gradients maps in horizontal and vertical directions. Canny [7] obtains more accurate edge map by introducing Gaussian filter and double thresholds. These traditional methods employ templates with fix parameters, which can only detect edges with specified intensity. Meanwhile, traditional edge detectors are sensitive to image scale changes. Recently, CNN-based edge detectors achieve the state-of-the-art performances by learning filters on multi-scale images. BDCN [8] devises a bi-directional cascade network for perceptual edge detection. RCF [9] considers the multi-scale hierarchical edge detection by fusing the features from different stages derived from VGG-16 backbone.

According to the mechanism of human visual system (HVS), compared with other components, human eyes are most sensitive to the edge information on the image [10]. As such, edge information is a vital important characteristic for vision [11]. The visual quality of image is highly correlated with the edge information [12]. As one of the semantic visual information, the reconstruction on edge maps represents the capacity of the method on detail restoration [13]. A clear and accurate edge map shows that the high-quality restored image is with few artifacts [14], [15]. There are also works concentrating on edge-preserving image super-resolution for better visualization performance [14], [15], [16]. SeaNet [17] restores the HR images with a branch to recover the edges. DEGREE [18] introduces the loss between edges of LR and HR images for high-frequency information recovery. However, existing works only use edge maps as constraints while neglecting to construct specific filters or components to explore structure information directly.

Inspired by the edge detectors, this paper proposes a novel cross convolution to explore the structural information of features, which is composed of two factorized asymmetric filters. Two filters are applied simultaneously to increase the matrix rank and preserve more structural information. Based on the cross convolution, cross convolution block (CCB) is devised with feature normalization (F-Norm) [23] and CA [24] to consider the inherent correlations of features, which concentrate

Y. Liu is with the School of Software, Dalian University of Technology, Dalian 116620, China (e-mail: liuyuqing@mail.dlut.edu.cn).

Q. Jia and X. Fan are with International School of Information Science and Engineering, Dalian University of Technology, Dalian 116620, China (e-mail: jiaqi@dlut.edu.cn; xin.fan@dlut.edu.cn).

S. Wang, S. Ma, and W. Gao are with the School of Electronics Engineering and Computer Science, Institute of Digital Media, Peking University, Beijing 100871, China (e-mail: sswang@pku.edu.cn; swma@pku.edu.cn).

W. Gao is with the School of Electronics Engineering and Computer Science, Institute of Digital Media, Peking University, Beijing 100871, China, and Peng Cheng Laboratory, Shenzhen 518000, China (email: wgao@pku.edu.cn)

Corresponding author: Qi Jia

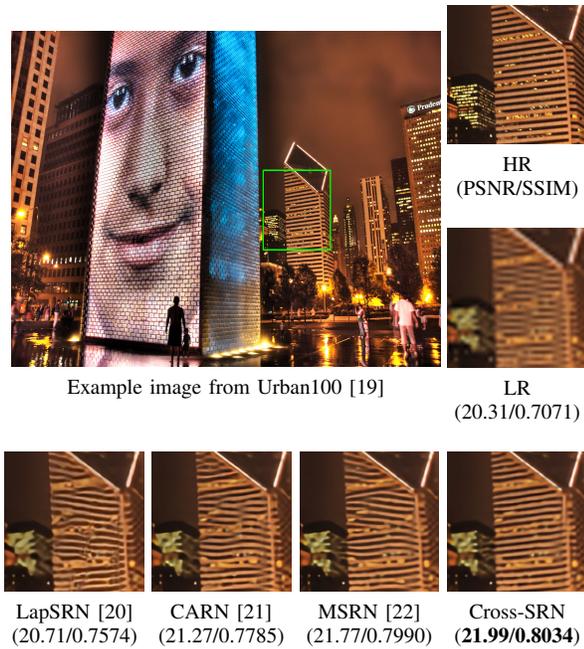


Fig. 1: Visual quality comparison for various image SR methods with scaling factor $\times 4$.

on the spatial and channel-wise information separately. CCBs are grouped in a multi-scale manner for feature exploration, termed as MFFG. MFFG considers the hierarchical edge information and progressively explore the structural information, where a padding structure is utilized to fully explore the features and residual connection is considered to keep the information. The MFFG modules are cascaded to constitute the final Cross-SRN. Experimental results show Cross-SRN achieves competitive or better performance than other works with more accurate structural information. Figure 1 shows the visual quality comparison for various image SR methods, where Cross-SRN restores more correct line and edge textures. To demonstrate the restoration performance quantitatively, we build a selected benchmark from numerical existing textured HR images with plentiful structural information, which shows our network can preserve the structural textures effectively.

Our contributions can be concluded as follows:

- Motivated by the edge detection methods, we design a cross convolution for effective structural information exploration.
- We devise a cross convolution block (CCB) to learn the relations of the edge features, and embedding the CCBs in a multi-scale feature fusion manner (MFFG) to explore the hierarchical features in different feature scales.
- The proposed Cross-SRN achieves competitive or better performance against the state-of-the-art methods with more accurate edge restoration. Especially, our network has a significant advantage over the selected benchmark with plentiful structure information.

II. RELATED WORKS

A. Deep Learning for Image Super-Resolution

CNN has proved to be an effective method for image restoration. SRCNN [2] is the first deep learning method for image SR, which contained three convolutional layers to present a sparse-coding like architecture. Then, VDSR [25] introduced a very deep network with residual learning for restoration. EDSR [26] removed batch normalization and built a large network with residual blocks. Recently, CNN-based works mainly focus on building an effective mapping from LR to HR with elaborate blocks and architectures. Inspired by Laplacian pyramid, LapSRN [20] and MS-LapSRN [27] built the networks to restore the images with different scaling factors simultaneously. Ahn *et al.* devised a cascading block in CARN [21] for fast and accurate image restoration. MSRN [22], RCAN [28], RDN [29], and other recent works also achieved state-of-the-art performances with well-designed blocks. However, these works seldom address the structural information exploration.

There are works considering the edge and gradient map as a prior for restoration. DEGREE [18] introduced the gradient loss between HR and SR images to constrain the structural information generation. Ma *et al.* investigated a GAN-based structure with gradient guidance in SPSR [30]. Fang *et al.* regarded edge as a prior in SeaNet [17] to restore the high-frequency information. These works consider the edge or gradient as a guidance or prior to restore the structural information, but almost neglect to explore edge information directly.

Information distillation provides an efficient way for feature exploration, which is usually implemented in a multi-scale feature fusion manner. As far as we know, IDN [31] proposed by Hui *et al.* is the first SR network with information distillation, which utilized channel separation to distill the important features. Hui *et al.* proposed IMDN with multi-distillation and channel attention for better performances. RFDN [32] was derived from IMDN and improved the network structure for fast and accurate restoration. However, these works seldom consider edge features in different scales upon limited computation costs.

B. Attention Mechanism

Attention mechanism is well established to emphasise vital information, which acts in a weighting distribution manner. Channel attention [24] (CA), as an efficient design for image SR, has been applied in recent state-of-the-art SR works. As far as we know, SENet [24] is the first work which introduces CA into deep learning. Recently, numerous works with CA has shown state-of-the-art performances on image SR. RCAN [28] investigated residual-in-residual blocks with CA to improve the performance. Dai *et al.* considered both CA and non-local attention in SAN [33]. IMDN [34], RFDN [32], HAN [4] and DRN [3] also demonstrates superior performances with CA. However, existing CA estimates the information by global average pooling, neglecting the spatial relation between features.

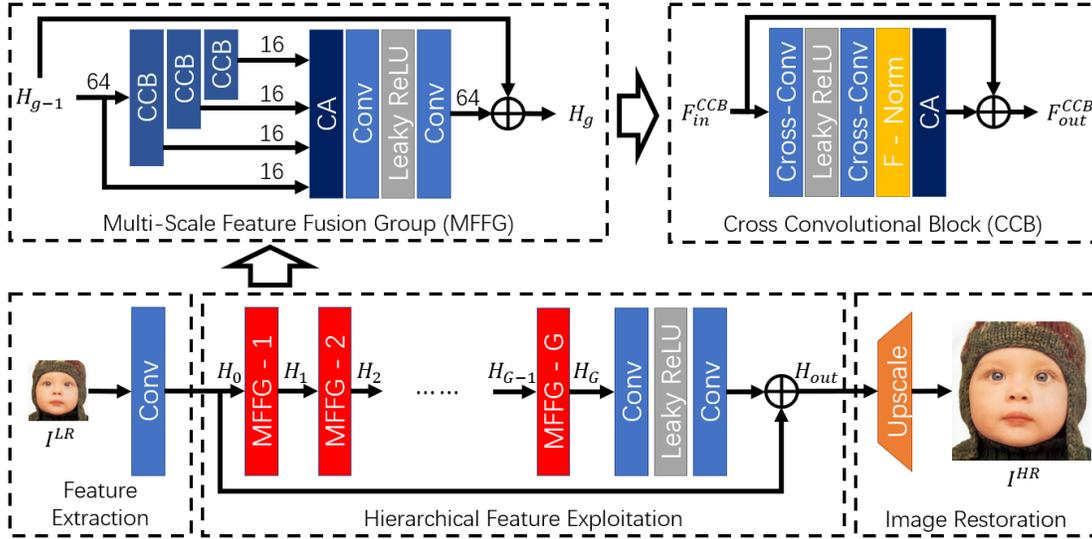


Fig. 2: Network structure of Cross-SRN. The holistic network Cross-SRN is illustrated in the image below, including three steps: feature extraction, hierarchical feature exploitation, and image restoration. The detail of multi-scale feature fusion group (MFFG) in red and the emended cross convolutional block (CCB) in dark blue are demonstrated in the upper left and upper right, respectively.

III. METHODOLOGY

In this section, we provide an overview of the proposed method. Then, we introduce the cross convolution and multi-scale feature fusion group (MFFG) modules in detail.

A. Network Design

An overview of the network design is shown in Figure 2. There are three steps in Cross-SRN, termed as feature extraction, hierarchical feature exploration, and image restoration, separated by different dashed boxes. Let I^{LR} and I^{HR} denote LR and HR instances, respectively. First, we use a convolutional layer to extract the original features from I^{LR} as

$$H_0 = f_{FE}(I^{LR}), \quad (1)$$

where $f_{FE}(\cdot)$ denotes the feature extraction step. The convolution expands the channel number of the input instance and maps the instance into a specific space containing more potential information than the RGB space.

Then, we try to retain valuable information in different scales of features, while emphasizing the edge features. Thus, we design G cascaded MFFGs for hierarchical feature exploitation with global residual learning. There is,

$$H_g = f_{MFFG}^g(H_{g-1}), \quad (2)$$

where $f_{MFFG}^g(\cdot)$ denotes the g -th MFFG, and H_{g-1} and H_g denote the input and output features of the MFFG module, respectively. The final output of the cascaded MFFGs H_G is fed into the residual module, which is composed of two convolutional layers with LeakyReLU layer. The padding structure with residual learning is devised as,

$$H_{out} = f_{PAD}(H_G) + H_0, \quad (3)$$

where $f_{PAD}(\cdot)$ denotes the padding.

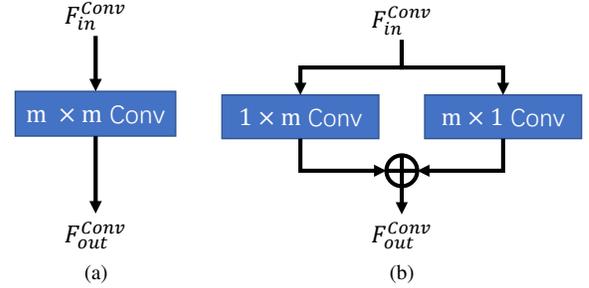


Fig. 3: Illustration of proposed cross convolution. (a) Vanilla convolution. (b) Cross convolution. The input feature is parallelly processed by the factorized asymmetric filters, whose addition is regarded as the output.

Finally, the HR images are restored by one convolution and a sup-pixel convolution, as,

$$I^{HR} = f_{IR}(H_{out}). \quad (4)$$

The final convolution decreases the channel number to restore the HR image and maps the features into RGB space.

B. Cross Convolution

As shown in Figure 3, different from vanilla convolutions, cross convolution is devised with two asymmetric perpendicular filters, which are denoted as $\mathbf{k}_{1 \times m}$ and $\mathbf{k}_{m \times 1}$ with receptive fields $1 \times m$ and $m \times 1$ separately. Let F_{in}^{Conv} , F_{out}^{Conv} be the input and output feature, then there is,

$$F_{out}^{Conv} = \mathbf{k}_{1 \times m} \otimes F_{in}^{Conv} + \mathbf{k}_{m \times 1} \otimes F_{in}^{Conv} + \mathbf{b}, \quad (5)$$

where \otimes denotes the convolution, \mathbf{b} is the bias term.

Cross convolution emphasizes the edge information by exploiting the vertical and horizontal gradient information

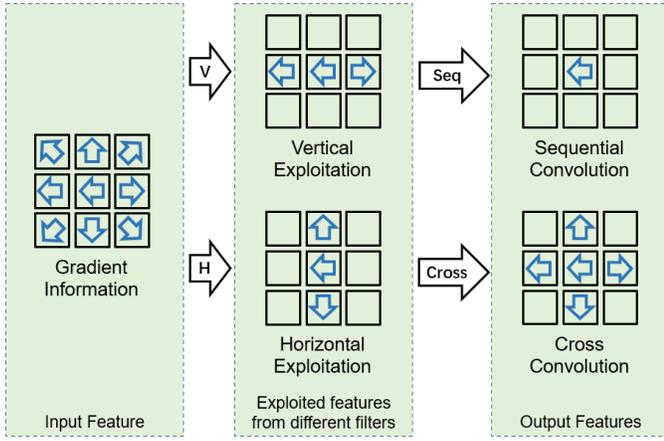


Fig. 4: An example of the difference on information preservation between sequential and cross convolutions. The blue arrows demonstrate the gradient directions of every pixel. V and H denote the vertical and horizontal exploitation separately.

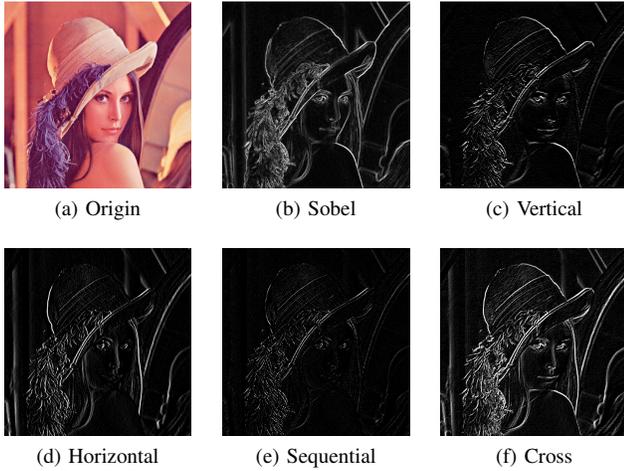


Fig. 5: An example of structural textures extracted by different filters. (a) Origin image. (b) Result of Sobel edge detector. (c) Edges extracted by vertical filter. (d) Edges extracted by horizontal filter. (e) Result of sequential convolution. (f) Result of cross convolution.

parallelly. The parallel design indicates fewer computation complexity and parameters than traditional filters with the same receptive field. Meanwhile, the parallel exploitation can also preserve more information than sequential design. Figure 4 demonstrates the difference on information preservation between sequential and cross convolutions. The input feature contains gradient information in various directions, and the blue arrows demonstrate the gradient directions of every pixel. After vertical and horizontal exploration by two asymmetric filters, gradients from different directions are exploited. Finally, the output features with sequential convolution only focuses on the main gradient direction of the feature, while the cross convolution can preserve more gradient directions.

The advantage of the proposed cross convolution can be validated by the amount of information it holds. For a vanilla

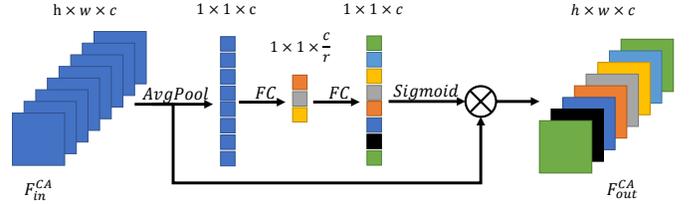


Fig. 6: Illustration of channel attention (CA).

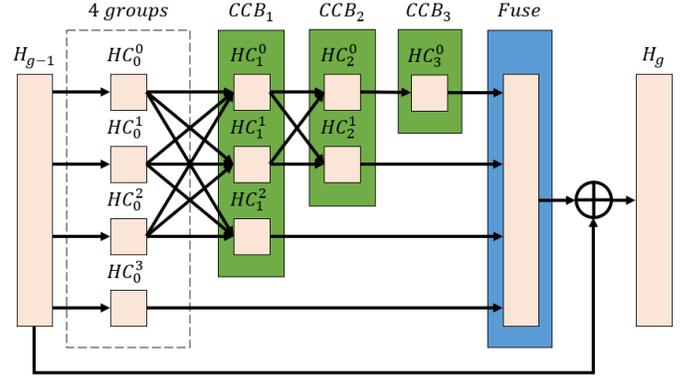


Fig. 7: Illustration of multi-scale feature fusion group (MFFG). The fuse structure is composed of two convolutional layers, one CA layer, and a Leaky ReLU.

convolution filter $\mathbf{k}_{m \times m}$, the rank is $\text{rank}(\mathbf{k}_{m \times m}) \leq m$. For the sequential combination of $\mathbf{k}_{1 \times m}$ and $\mathbf{k}_{m \times 1}$, the rank is $\text{rank}(\mathbf{k}_{1 \times m} \cdot \mathbf{k}_{m \times 1}) = 1$. The filter with lower rank preserves less information than the filter with higher rank. The rank of cross convolution $\mathbf{k}_{\text{cross}}$ is $\text{rank}(\mathbf{k}_{\text{cross}}) \leq 2$, which can preserve more latent information than the sequential one.

In fact, the cross convolution holds a similar formulation with the traditional edge detectors. Figure 5 shows an example of structural textures extracted by different filters. We compare the cross convolution with Sobel operator, which detects the edges from both vertical and horizontal directions. Figure 5 (b) and (f) are the results of Sobel and cross convolution, respectively. The comparison demonstrates that cross convolution can preserve most of the edge information with explicit and sharp edges, validating the capacity of structure texture exploration. Specially, we also compare the extracted edge maps with vertical, horizontal filters and sequential convolution, in figures (c), (d), and (e), which demonstrate that cross convolution can extract more edge information than other methods.

It is worth noting that Inception [35] also has similar filters to Cross convolution, but their method is inapplicable to the challenging issues that we are addressing. Cross convolution is specially designed to explore edge information for image SR, while filters in Inception aim to save parameters for image classification, without explicit edge features. Cross convolution explores and overlays edge features in two perpendicular directions, invoking higher matrix rank of filters to preserve more information than channel aggregation of Inception. More importantly, Cross convolution is capable of any position of the networks for edge exploration, while Inception blocks can't be used in the beginning of the network as stated in their paper.

Furthermore, The sequentially stacked filters in Inception are discussed and termed as “**Seq**” in Figure 4 and Section IV-B1. The contrast test in Table I demonstrates Cross convolution preserves more information than **Seq**.

The cross convolution is designed with two factorized filters and concentrates on the vertical and horizontal structural information exploration. Similar to Canny, Sobel, and other traditional edge detectors, the factorized filters explore the edges in two orthogonal directions and derive the final structural information by summarization. As such, the cross convolution is inherently suitable for structural information exploration.

Separable convolution [36] organizes the asymmetric filters in a sequential manner, while the proposed cross convolution designs the filters in a parallel fashion and achieves better quantitative performance with the same number of parameters. The motivation of spatial separable convolution is to convert the regular convolution into a parameter efficient design with less matrix multiplications and keep the receptive field. Compared with spatial separable convolution, the proposed cross convolution for image super-resolution only requires one extra addition operation with insignificant computational cost. The “**Seq**” in Table I denotes the same design as the separable convolution. In the table, our cross convolution achieves better restoration performance with similar computational complexity. The cross convolution concentrates more on the structural information, which holds a similar formulation to the edge detectors with higher matrix rank and more potential information than the separable convolution.

C. Multi-Scale Feature Fusion Group

In order to obtain accurate edge information, we build the cross convolution block (CCB) based on the basic cross convolutions, as shown in Figure 2. Two cross convolutions with a Leaky ReLU activation are utilized to explore the structural information. Besides convolutions, F-Norm [23] and CA [24] are considered to emphasize important spatial and channel-wise information separately. Figure 6 shows the operation of CA, where global average pooling is applied to squeeze the information, and two full connection layers with ReLU activation explore the non-linear attention for every channel. F-Norm concentrates on the diversity of spatial information, which can be formulated as,

$$F_{out}^{(i)} = (F_{in}^{(i)} \otimes \mathbf{k}^{(i)} + \mathbf{b}^{(i)}) + F_{in}^{(i)}, \quad (6)$$

where $F_{in}^{(i)}$, $F_{out}^{(i)}$ are the input and output features of the i -th channel, $\mathbf{k}^{(i)}$ and $\mathbf{b}^{(i)}$ are the filter and bias for the i -th channel.

CCB in the network organizes the cross convolution in a residual block design. Residual block design has been widely considered in advanced networks for boosting the performance. Besides the residual connection, channel-wise attention and feature normalization are utilized to improve the exploration performance, which prove to be effective components for image SR.

As edge information is sensitive to scale changes, CCBs are grouped in a multi-scale feature fusion manner in MFFG

to explore features in different scales. As shown in Figure 7, for the g -th MFFG, the input feature H_{g-1} is divided into several groups with the same number of channels in average. As demonstrated in Figure 7, let $f_{CCB}^j(\cdot)$ be the j -th CCB in MFFG, and the input feature H_{g-1} is divided into four groups from HC_0^0 to HC_3^0 . The multi-scale feature fusion can be demonstrated as,

$$\begin{aligned} [HC_1^0, HC_1^1, HC_1^2] &= f_{CCB}^1([HC_0^0, HC_0^1, HC_0^2]), \\ [HC_2^0, HC_2^1] &= f_{CCB}^2([HC_1^0, HC_1^1]), \\ [HC_3^0] &= f_{CCB}^3([HC_2^0]), \end{aligned} \quad (7)$$

where HC_j^k denotes the k -th group after the j -th CCB, and $[\cdot]$ denotes the group combination.

The hierarchical features are aggregated by a residual block structure, which is composed of two convolutional layers with Leaky ReLU and one CA layer. Finally, the output of MFFG is,

$$H_g = f_{Fuse}([HC_3^0, HC_2^1, HC_1^2, HC_0^3]) + H_{g-1}, \quad (8)$$

where $f_{Fuse}(\cdot)$ is the fusion structure.

Herein, MFFG keeps the original information and emphasizes the structural information hierarchically. As shown in the upper left of the Figure 2, the input feature H_{g-1} are separated into four groups. Three groups are sequentially processed by CCBs for hierarchical structural information exploitation. In order to preserve the potential information lost during edges exploring, the last group keeps the identical original features. After the CCBs, a vanilla residual block structure with channel attention is utilized for effective feature exploration and gradient transmission.

IV. EXPERIMENTS AND ANALYSIS

In this section, we introduce the experiment settings of our Cross-SRN in the beginning. Then, we provide the ablation study on the cross convolution and MFFG to show the effectiveness of the proposed model. Finally, we compare our Cross-SRN with the state-of-the-art works.

A. Experiment Settings

In Cross-SRN, there are $G = 10$ MFFGs for non-linear exploitation. Filter numbers of convolution layers are set as $c = 64$ except for the image restoration module, and the kernel sizes of filters are set as $m = 3$. Since the largest scaling factor in this paper is set as $\times 4$, the image restoration module can be regarded as a down-sampling step from the feature space to RGB space. We train the network with DIV2K [37] dataset. DIV2K is firstly proposed in New Trends in Image Restoration and Enhancement (NTIRE) 2017 competition and has been widely used for image SR tasks. We choose 895 images for training, and 5 images for validation. We leverage bicubic [38] (**BI**) to obtain the degraded images by different scaling factors. For training sets, images are cropped with patch size 48×48 , and then randomly flipped and rotated for augmentation. Cross-SRN is updated by an Adam [39] optimizer with learning rate $lr = 10^{-4}$. We train the network for 1000 epochs, and halve the learning rate for every 200

TABLE I: Comparison of three different convolution designs on PSNR/SSIM with **BI** $\times 4$ degradation.

Method	Param	MACs	Selective	Set5	Set14	B100	Urban100	Manga109
Seq	1,296K	74.2G	24.39/0.8527	32.20/0.8949	28.61/0.7817	27.57/0.7361	26.15/0.7875	30.49/0.9077
Conv	1,509K	86.5G	24.44/0.8549	32.20/0.8946	28.60/0.7823	27.58/0.7367	26.17/0.7889	30.54/0.9092
Cross	1,296K	74.2G	24.46/0.8539	32.24/0.8954	28.59/0.7817	27.58/0.7364	26.16/0.7881	30.53/0.9081

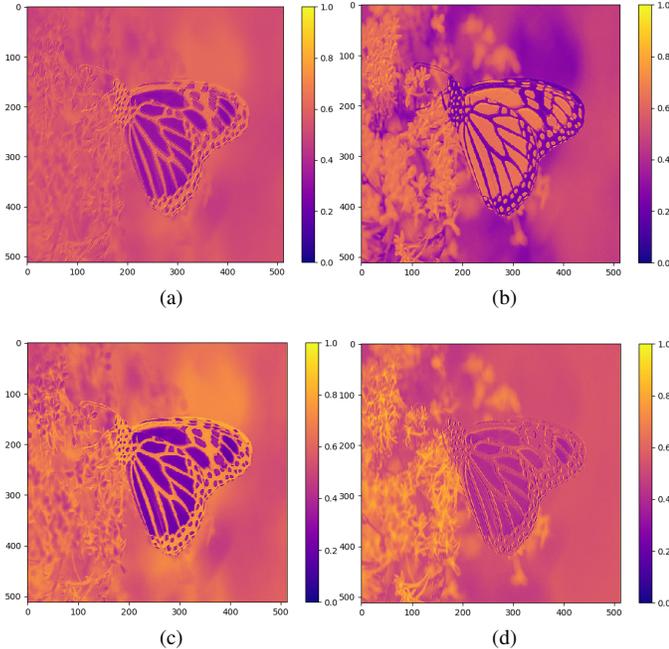


Fig. 8: Visualized feature maps processed by different convolution designs. (a) Input feature. (b) Feature processed by the vertical filter. (c) Feature processed by the horizontal filter. (d) Output feature of cross convolution. The values are calculated by averaging the feature maps and normalized in range 0 to 1.

TABLE II: Investigation on multi-scale feature fusion on PSNR/SSIM with **BI** $\times 4$ degradation.

Method	Param	MACs	Set5	Set14	B100
w/o MF	2,366K	135.2G	32.24/0.8957	28.65/0.7837	27.62/0.7379
w/o CCB	847K	48.5G	32.04/0.8926	28.48/0.7792	27.50/0.7338
Ours	1,296K	74.2G	32.24/0.8954	28.59/0.7817	27.58/0.7364

epochs. The evaluation indicators are chosen as peak signal-to-noise ratio (PSNR), structural similarity (SSIM) and multiply-accumulate operations (MACs). The higher PSNR/SSIM result means better performance, and the lower MACs means faster speed. MACs is calculated by restoring a 1280×720 (720P) image with scaling factor $\times 4$.

Five testing benchmarks are used for performance comparison: Set5 [40], Set14 [41], B100 [42], Urban100 [19], and Manga109 [43]. Urban100 is a testing benchmark with real-world HR images. B100 is also a benchmark from real-world with abundant complex textures. Manga109 is composed of comic covers. All of the three benchmarks contain plentiful high-frequency and structural information. Furthermore, we build a selective benchmark to compare the capacities of structural information recovery, which is composed of images

TABLE III: PSNR/SSIM, parameters and MACs comparisons between sequential convolution and cross convolution on Set5 with scaling factor $\times 4$.

Receptive Field	3×3		5×5		7×7	
	Seq	Cross	Seq	Cross	Seq	Cross
PSNR	31.99	32.07	32.10	32.14	32.08	32.08
SSIM	0.8924	0.8926	0.8931	0.8936	0.8932	0.8931
Params (M)	1.29	1.29	1.58	1.58	1.87	1.87
MACs (G)	74.22	74.22	90.74	90.74	107.25	107.25

TABLE IV: Investigation on channel attention (CA) and feature normalization (F-Norm).

CA	FN	Set5	Urban100	Manga109
✓	✗	32.20/0.8952	26.15/0.7875	30.47/0.9079
✗	✓	32.13/0.8948	26.11/0.7872	30.50/0.9086
✓	✓	32.24/0.8954	26.16/0.7881	30.53/0.9081

from Urban100 and Manga109 with numerical edges and lines. The benchmark is built according to the following steps. Firstly, we blur the images with 7×7 Gaussian kernel for denoising. After filtering, Sobel operator is utilized to extract the edges from the images. We use threshold $t_e = 128$ to remove the weak responses from edge maps, and calculate the average response of image. The images with responses higher than $t_r = 12$ are included in the benchmark. There are 38 images in total, including 18 images from Urban100 and 20 images from Manga109.

t_e and t_r are used to select images with more edge information. t_e is used to explore edges with high response and t_r is used to evaluate the average intensity of edge information in images. The edge map extracted by Sobel represents the responding of edge information in every pixel with range 0-255. t_e is set as 128, an average value of the range. Pixels with a value less than t_e is set as 0. Then, the average response of each pixel is calculated to estimate the intensity of edge information in each image. If the response is larger than an empirical value $t_r=12$, the image is included in the selective benchmark.

B. Ablation Study

Herein, we investigate the performance gained from the cross convolution, the multi-scale feature fusion and the network structure separately.

1) *Investigation on Cross Convolution*: To demonstrate the performance of cross convolution, we design three convolution structures for comparison: **Cross** denotes the proposed cross convolution, **Seq** denotes the sequential convolution, and **Conv** denotes the vanilla convolution. To investigate the effectiveness of **Cross** on structural information preservation, we test the three convolution structures on the selective benchmarks and other widely used benchmarks. Table I demonstrates the comparison on PSNR/SSIM. We can see **Cross** achieves competitive or better PSNR/SSIM performance on the testing

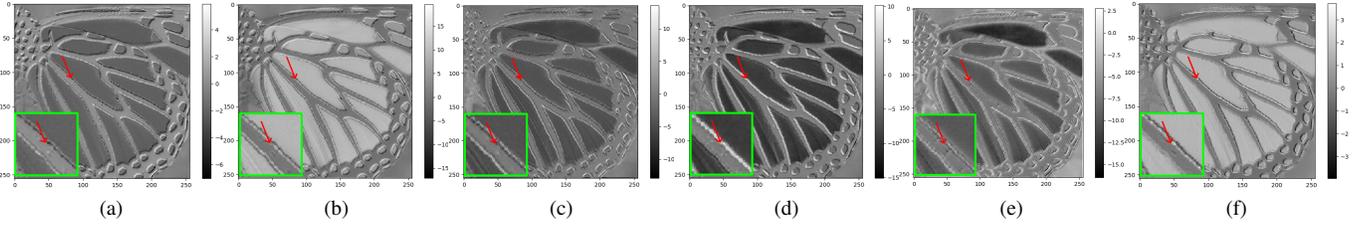


Fig. 9: Visualized feature maps of multi-scale feature fusion. (a): input feature. (b)-(e): different groups after CCBs, which denote from HC_0^3 to HC_3^0 in Eq. 7. (f): output feature. All the images are illustrated by the same gray-scale strategy.

benchmarks with fewer parameters and MACs, which proves our hypothesis on matrix rank. Specially, on the selective benchmark with plentiful edge structures, **Cross** achieves higher PSNR than **Conv** and **Seq**, which means the proposed **Cross** achieves better performance on structural information preservation with only 76.8% MACs and 85.8% parameters of **Conv** or **Seq**. Thus, **Cross** makes a good balance on edge feature exploration and information loss.

The visualized analysis of three convolution structures is demonstrated in Figure 8. In Figure 8 (b) and (c), the features processed by two factorized filters response to different kinds of areas. Figures (b) and (c) are processed by different directions of filters, which have high responses to two directional margins, respectively. On the contrary, the proposed cross convolution extracts edge features parallelly from the same input. As shown in Figure 8 (d), the edges are enhanced after cross convolution, which also demonstrates our capacity on structural texture restoration.

The cross convolution is designed to substitute the ordinary convolution for addressing the structural information with restricted parameters and MACs. A vanilla $m \times m$ convolutional filter can be regarded as a matrix with shape $m \times m$, and the rank of the filter $\text{rank}(\mathbf{k}_{m \times m}) \leq m$. In other words, the upper bound of the filter's rank is m . The cross convolution is designed as the combination of $\mathbf{k}_{m \times 1}$ and $\mathbf{k}_{1 \times m}$, whose upper bound is 2. As such, the cross convolution can preserve the most information of ordinary convolution when m is small. Furthermore, the cross convolution follows a similar formulation to the traditional edge detectors, which can effectively explore the structural information. The experimental results show cross convolution achieves competitive PSNR/SSIM performance than the ordinary convolution with fewer parameters and MACs.

The sequential situation **Seq** holds a same receptive field as the vanilla convolution, but the sequential exploration losses much information. The rank of sequential convolution $\text{rank}(\mathbf{k}_{1 \times m} \cdot \mathbf{k}_{m \times 1}) = 1$, which means the sequential convolution potentially losses half of the information than the cross convolution.

To better illustrate the potential information loss, we perform the ablation studies to compare the PSNR/SSIM. The results are shown in Table. III

Table III shows the PSNR/SSIM, parameters and MACs comparisons between sequential convolution and cross convolution on Set5 benchmark with scaling factor $\times 4$. For a

fair comparison, all models are re-trained for 200 epochs. In the table, we can find that the cross convolution achieves better PSNR/SSIM performance than the sequential design on all receptive fields. When the receptive field is 3×3 , there is 0.8 dB PSNR improvement, which is significant. This observation is in accordance with our conclusion on the rank of filters which shows the cross convolution can preserve more information than the sequential design. Furthermore, we can find that the performances of cross convolution are similar when the receptive fields are different. As we discussed on the filter's rank, the upper bound of cross convolution's rank is 2. So, choosing filter size $m = 3$ is the most efficient way for restoring the HR image.

There is another reason for choosing $m = 3$. According to the linearity of convolution operation, the filter with larger receptive fields can be equivalently substituted by a combination of filters with smaller receptive fields [35]. The receptive field 3×3 is the economic choice to build a deeper neural network, which has been widely considered in recent SR works [26], [29], [28].

2) *Investigation on Multi-Scale Feature Fusion*: To investigate the effectiveness of the multi-scale feature fusion module, we compare the proposed MFFG design with the other two optional designs, one is MFFG without channel division and multi-scale feature fusion, which is denoted as **w/o MFF**. The other is MFFG without CCBs, which is termed as **w/o CCB**. The **w/o CCB** means we remove all the CCBs from the Cross-SRN. In other words, the three CCBs in MFFG are omitted while other components are not modified. The ablation study of **w/o CCB** aims to investigate the effectiveness of the backbone of Cross-SRN. In the **w/o MFF** version, the three CCBs all take all the feature channels as input, and all the CCBs are processed sequentially. In other words, **w/o MFF** means all the CCBs in MFFG are modified with channel number as 64, and no channel separation is considered. The MFFG without multi-scale feature fusion can be regarded as a residual-in-residual-like architecture by stacking three CCBs, two convolutional layers, one LeakyReLU and one CA. The comparison results on PSNR/SSIM are illustrated in Table II. We can see **Ours** achieves competitive performance than **w/o MFF** with around half parameters and MACs. **Ours** drops less than 0.01 dB on Set5 and 0.05db on B100. The model without MFF holds similar PSNR/SSIM to MFFG but the Params and MACs are much higher. The results demonstrate that multi-scale feature fusion preserves most edge information

by hierarchical exploration. The model without CCB drops near 0.2 dB on Set5 when compared with MFFG, which means CCB is crucial for restoration.

According to Figure 7, we visualize the feature maps in four groups before fusing and the final feature map of the MFFG in Figure 9. The input feature of MFFG is demonstrated in Figure 9 (a). Figures 9 (b)-(e) are the feature maps of groups from different CCBs, indicating features from HC_0^3 to HC_3^0 in Eq. 7, respectively. Figure 9 (f) is the output of MFFG. All the images are produced by the same strategy. For each group, the feature maps are averaged and normalized in the range 0 to 1. We can obtain more and more refined edges and textures from (b) to (e), as the structural information is enhanced with the increase of CCBs. We take a zoomed-in green rectangle area as an example, which is guided by the red arrow. The small edge is clearer with enhanced contrast from (b) to (e). Thus, the experiment validates that the edges and structural textures are sharper in the output feature after the CCBs processing, which is obvious in the comparison between (a) and (f).

3) *Analysis on Network Structure*: We analysis and compare the performance of different network structure settings, including channel attention, and the different feature normalization mechanisms.

We demonstrate the effectiveness of the channel attention (CA) and feature normalization by testing the performance with or without two structures. The results are demonstrated in Table IV. We can see our model with CA and F-Norm achieves the best performance. The results show that model without F-Norm drops 0.06 dB on Manga109 benchmark, and the model without CA drops 0.05 dB on Urban100. In this point of view, CA and F-Norm are effective components for boosting the network performance.

C. Comparison with State-of-the-Art Methods

We compare Cross-SRN with three kinds of representative image SR works, including classical image SR works (SRCNN [2], VDSR [25], LapSRN [27], and CARN [21]), multi-scale works (MRFN [44] and IMDN [34]), and structure-preserving works with edge map prior (DEGREE [18] and SeaNet-baseline [17]). The PSNR/SSIM results are shown in Table V. All the methods for comparison follow the same training and testing protocol, and the performances are provided by the published papers. The last column of Table V demonstrates the proposed Cross-SRN achieves competitive or superior performances, and 90% of the results achieve the best (in bold) or second best (underline) performance. Especially, when compared with the classical and multi-scale works, Cross-SRN has near 0.1db PSNR improvement on Urban100 and Manga109, and achieves better performance on B100. The superior performance demonstrates that our network can recover structural textures more efficiently.

The proposed Cross-SRN also has obvious advantage over edge map based methods DEGREE [18] and SeaNet-baseline [17]. In Table V, Cross-SRN achieves superior performances to DEGREE on all testing benchmarks. Compared with SeaNet, Cross-SRN achieves 90% best or second best results over all the scale factors and benchmarks, while

the percentage for SeaNet is only 43%. On Urban100 and Manga109, the PSNR for Cross-SRN is 0.1-0.3 dB higher than SeaNet. It is worth noting that SeaNet-baseline contains 4.1M parameters, while Cross-SRN only contains around 1.3M. In this point of view, Cross-SRN can restore the structural information more effectively than other edge map based works with less parameters.

We also compare our Cross-SRN with other works on the selective benchmark to show the effectiveness of structure information preservation. As shown in Table VI, the selected 18 images from Urban100 are denoted as Urban100 (S), 20 images from Manga109 are denoted as Manga109 (S), and the total 38 images are denoted as Selected. The experiment result demonstrated Cross-SRN achieves superior performance to other works, which shows that our network has higher capacity on preserving structural information.

To demonstrate the restoration performance, visualization comparisons are shown in Figure 10, Figure 11 and Figure 12. The results on Set14 dataset demonstrate the effectiveness on structural information recovery. Urban100 benchmark is composed of high-resolution real world images with plentiful complex structural textures. Two representative instances of building are chosen to represent the restoration capacity. From Figure 11, Cross-SRN recover the grids and lines from tall buildings more efficiently. Compared with IMDN, Cross-SRN can prevent more texture mixtures, and regain more correct structural textures. It should be noted that IMDN is another image SR network with multi-scale fusion design. From this point of view, Cross-SRN has a convincing restoration capacity on structural information restoration.

Besides Urban100, we select two representative instances from Manga109 for comparison. Manga109 is a benchmark of comics with sharply defined areas and much high-frequency information. We compare Cross-SRN with LR, LapSRN and IMDN. The results are shown in Figure 12. From the comparisons, Cross-SRN can restore the high-frequency information more accurately. Specially, the tiny textures can be recovered more effectively by Cross-SRN, such as areas of the word and the eye. Meanwhile, the serried line textures, which are mixed by down-sampling, can be well recovered by Cross-SRN. From the visualization comparison, Cross-SRN gains superior restoration performance than other lightweight works.

It is worth noting that the proposed Cross-SRN achieves superior performance with fewer parameters and lower computation complexity than other methods. Figure 13 demonstrates an intuitive comparison on the PSNR and corresponding parameters and MACs on Urban100 dataset. Our method is labeled by red star, while other methods are labeled by blue points. The proposed method achieves the best performance with the red star over all the blue points with less parameters and lower MACs. Thus, Cross-SRN proves to be an efficient design for structure-preserving image restoration.

Besides the objective and subjective comparisons, we also compare the speed of different methods. We calculate the time cost by restoring a random 720P image with $\mathbf{BI} \times 4$ degradation for 100 times, and choose the average value for comparison. The results are shown in Table VII. In the table, our network achieves the best PSNR/SSIM performance with



Fig. 10: Visualization comparison on Set14 dataset.

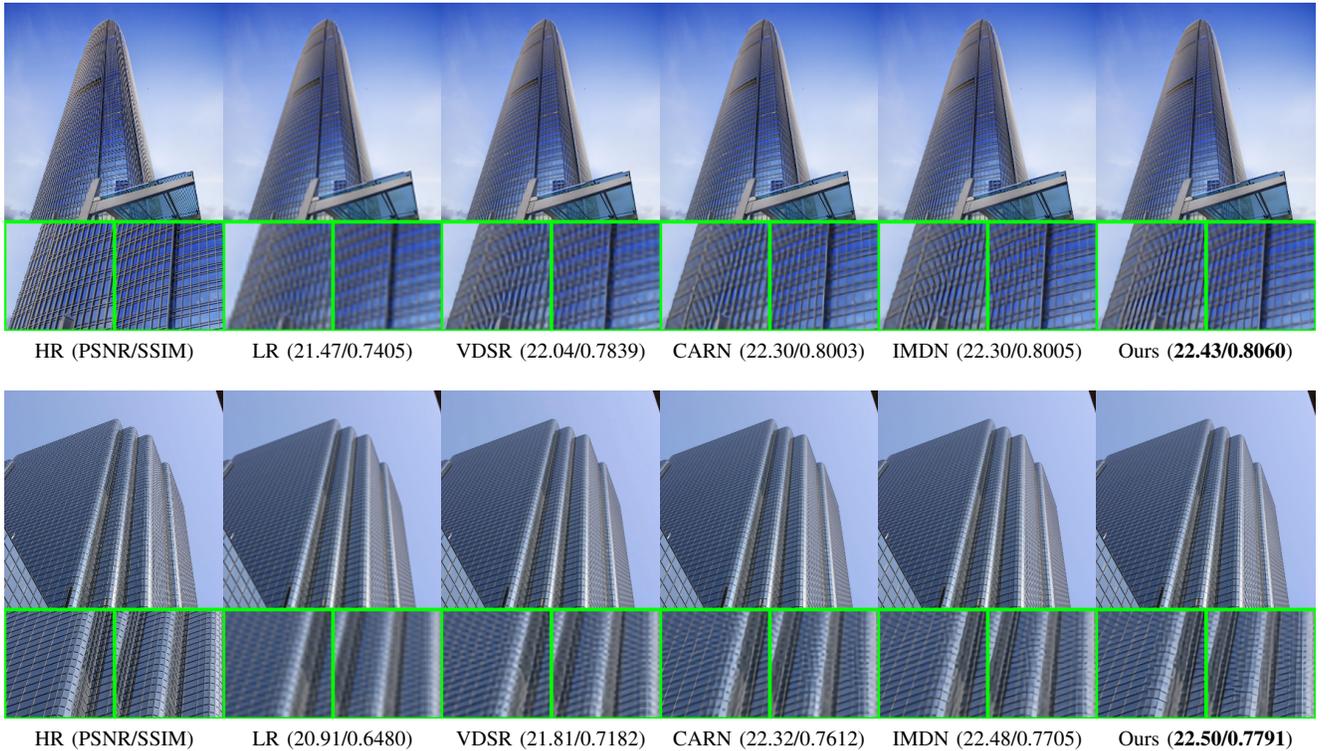


Fig. 11: Visualization comparison on Urban100 dataset.

a fast speed. Compared with CARN, we achieve 0.1 dB PSNR improvement and only require 20 ms more time cost.

V. CONCLUSION

In this paper, we propose a network termed as Cross-SRN for edge-preserving image super-resolution. Inspired by edge detection methods, a novel cross convolution operation is designed to exploit the structural information more effectively. The cross convolution leverages two perpendicular factorized filters parallelly to increase the matrix rank and preserve more information. Based on the cross convolution, CCBs are

investigated with CA and F-Norm to concern the inherent correlation of channel-wise and spatial features separately. MFFG groups CCBs in a multi-scale feature fusion manner for efficient hierarchical feature exploration. Experimental result shows that Cross-SRN has demonstrated superior restoration capacity than other lightweight works on both quantitative and qualitative comparisons.

ACKNOWLEDGEMENT

This work is partially supported by the NSF of China under grand Nos. 62088102, 62031013 and High-performance

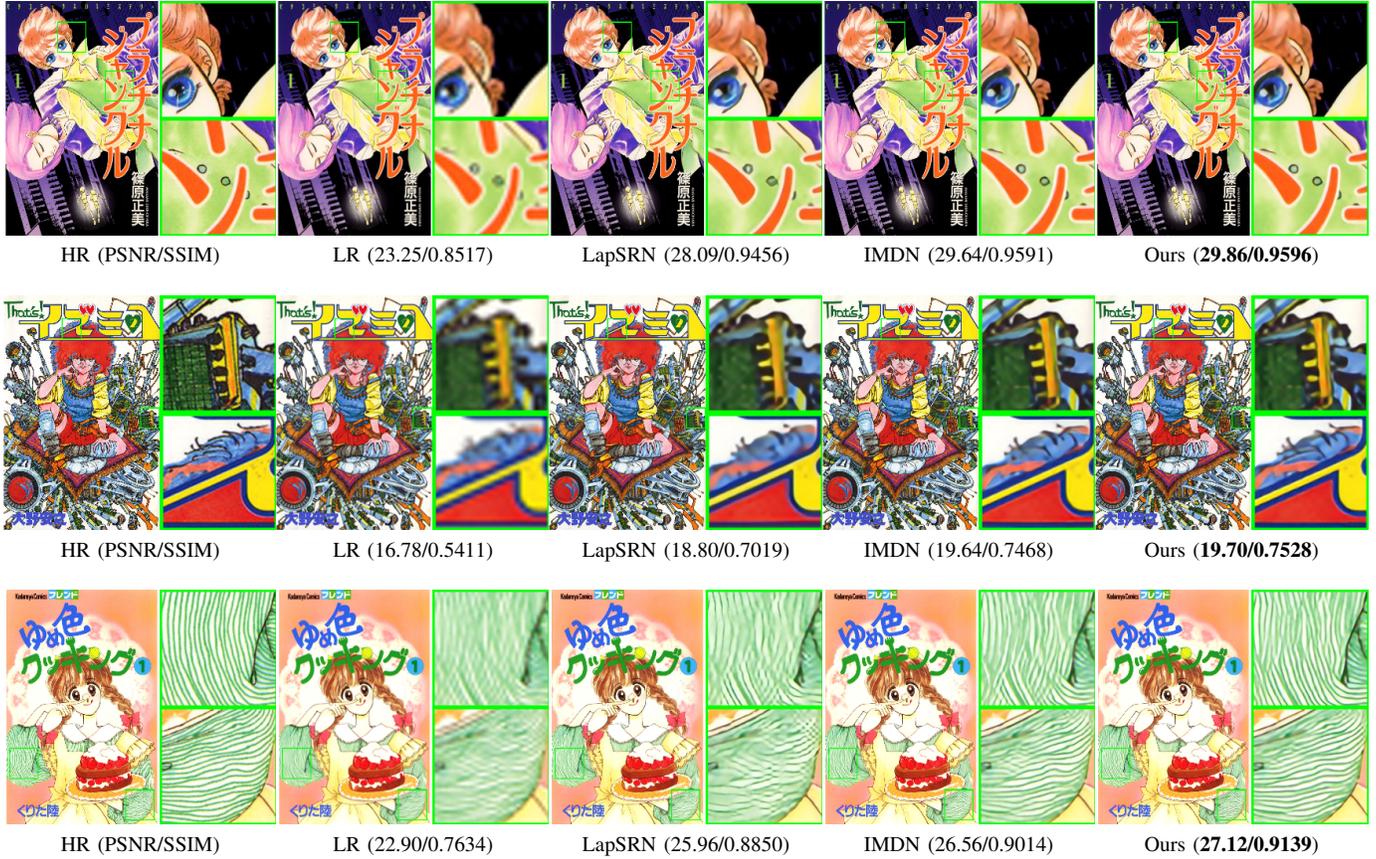


Fig. 12: Visualization comparison on Manga109 dataset.

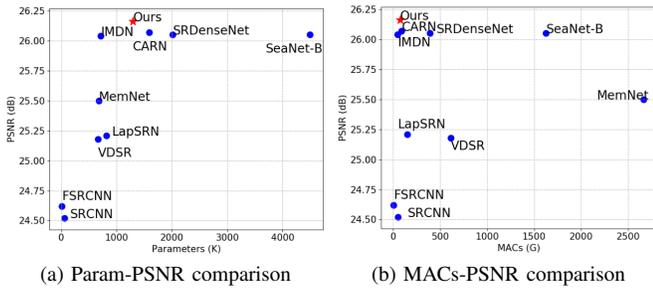


Fig. 13: Comparisons of parameters, MACs and performances on Urban100 with $BI \times 4$ degradation. (a) Parameter comparison. (b) MACs comparison.

Computing Platform of Peking University, which are gratefully acknowledged. This work is also partially supported by the NSF of China under grant Nos. 61876030, 61733002.

REFERENCES

[1] W. Gao, S. Ma, L. Duan, Y. Tian, P. Xing, Y. Wang, S. Wang, H. Jia, and T. Huang, "Digital retina: A way to make the city brain more efficient by visual coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 11, pp. 4147–4161, 2021.

[2] 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, 2016.

[3] Y. Guo, J. Chen, J. Wang, Q. Chen, J. Cao, Z. Deng, Y. Xu, and M. Tan, "Closed-loop matters: Dual regression networks for single image super-resolution," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5406–5415, 2020.

[4] B. Niu, W. Wen, W. Ren, X. Zhang, L. Yang, S. Wang, K. Zhang, X. Cao, and H. Shen, "Single image super-resolution via a holistic attention network," in *European Conference on Computer Vision*, vol. 12357, pp. 191–207, 2020.

[5] X. Luo, Y. Xie, Y. Zhang, Y. Qu, C. Li, and Y. Fu, "Latticenet: Towards lightweight image super-resolution with lattice block," in *European Conference on Computer Vision*, Springer, 2020.

[6] Y. Chien, "Pattern classification and scene analysis," *IEEE Transactions on Automatic Control*, vol. 19, no. 4, pp. 462–463, 1974.

[7] J. Canny, "A computational approach to edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PAMI-8, no. 6, pp. 679–698, 1986.

[8] J. He, S. Zhang, M. Yang, Y. Shan, and T. Huang, "Bdcn: Bi-directional cascade network for perceptual edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.

[9] Y. Liu, M. Cheng, X. Hu, J. Bian, L. Zhang, X. Bai, and J. Tang, "Richer convolutional features for edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 8, pp. 1939–1946, 2019.

[10] X. Ran and N. Farvardin, "A perceptually motivated three-component image model-part i: description of the model," *IEEE Transactions on Image Processing*, vol. 4, no. 4, pp. 401–415, 1995.

[11] H. Wang, X. Hu, X. Zhao, and Y. Zhang, "Wide weighted attention multi-scale network for accurate mr image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2021.

[12] S. Mandal and A. K. Sao, "Edge preserving single image super resolution in sparse environment," in *2013 IEEE International Conference on Image Processing*, pp. 967–971, 2013.

[13] Y. Li, J. Liu, W. Yang, and Z. Guo, "Neighborhood regression for edge-preserving image super-resolution," in *2015 IEEE International*

TABLE V: Average PSNR/SSIM with degradation model **BI** $\times 2$, $\times 3$, and $\times 4$ on five benchmarks. The best and second performances are shown in **bold** and underline.

Scale	Model	Set5 [40] PSNR/SSIM	Set14 [41] PSNR/SSIM	B100 [42] PSNR/SSIM	Urban100 [19] PSNR/SSIM	Manga109 [43] PSNR/SSIM
$\times 2$	SRCNN [2]	36.66 / 0.9542	32.42 / 0.9063	31.36 / 0.8879	29.50 / 0.8946	35.74 / 0.9661
	FSRCNN [45]	37.00 / 0.9558	32.63 / 0.9088	31.53 / 0.8920	29.88 / 0.9020	36.67 / 0.9694
	VDSR [25]	37.53 / 0.9587	33.03 / 0.9124	31.90 / 0.8960	30.76 / 0.9140	37.22 / 0.9729
	DRCN [46]	37.63 / 0.9588	33.04 / 0.9118	31.85 / 0.8942	30.75 / 0.9133	37.63 / 0.9723
	CNF [47]	37.66 / 0.9590	33.38 / 0.9136	31.91 / 0.8962	-	-
	LapSRN [20]	37.52 / 0.9590	33.08 / 0.9130	31.80 / 0.8950	30.41 / 0.9100	37.27 / 0.9740
	DRRN [48]	37.74 / 0.9591	33.23 / 0.9136	32.05 / 0.8973	31.23 / 0.9188	37.92 / 0.9760
	BTSRN [49]	37.75 / -	33.20 / -	32.05 / -	31.63 / -	-
	MemNet [50]	37.78 / 0.9597	33.28 / 0.9142	32.08 / 0.8978	31.31 / 0.9195	37.72 / 0.9740
	SelNet [51]	37.89 / 0.9598	33.61 / 0.9160	32.08 / 0.8984	-	-
	CARN [21]	37.76 / 0.9590	33.52 / 0.9166	32.09 / 0.8978	31.92 / 0.9256	38.36 / 0.9765
	IMDN [34]	<u>38.00 / 0.9605</u>	<u>33.63 / 0.9177</u>	<u>32.19 / 0.8996</u>	<u>32.17 / 0.9283</u>	38.88 / 0.9774
	RAN [52]	37.58 / 0.9592	33.10 / 0.9133	31.92 / 0.8963	-	-
	DNCL [53]	37.65 / 0.9599	33.18 / 0.9141	31.97 / 0.8971	30.89 / 0.9158	-
	FilterNet [54]	37.86 / 0.9610	33.34 / 0.9150	32.09 / 0.8990	31.24 / 0.9200	-
	MRFN [44]	37.98 / 0.9611	33.41 / 0.9159	32.14 / 0.8997	31.45 / 0.9221	38.29 / 0.9759
SeaNet-baseline [17]	37.99 / 0.9607	33.60 / 0.9174	32.18 / 0.8995	32.08 / 0.9276	38.48 / 0.9768	
DEGREE [18]	37.58 / 0.9587	33.06 / 0.9123	31.80 / 0.8974	-	-	
Cross-SRN (Ours)	38.03 / 0.9606	33.62 / 0.9180	32.19 / 0.8997	32.28 / 0.9290	38.75 / 0.9773	
$\times 3$	SRCNN [2]	32.75 / 0.9090	29.28 / 0.8209	28.41 / 0.7863	26.24 / 0.7989	30.59 / 0.9107
	FSRCNN [45]	33.16 / 0.9140	29.43 / 0.8242	28.53 / 0.7910	26.43 / 0.8080	30.98 / 0.9212
	VDSR [25]	33.66 / 0.9213	29.77 / 0.8314	28.82 / 0.7976	27.14 / 0.8279	32.01 / 0.9310
	DRCN [46]	33.82 / 0.9226	29.76 / 0.8311	28.80 / 0.7963	27.15 / 0.8276	32.31 / 0.9328
	CNF [47]	33.74 / 0.9226	29.90 / 0.8322	28.82 / 0.7980	-	-
	DRRN [46]	34.03 / 0.9244	29.96 / 0.8349	28.95 / 0.8004	27.53 / 0.8378	32.74 / 0.9390
	BTSRN [49]	34.03 / -	29.90 / -	28.97 / -	27.75 / -	-
	MemNet [50]	34.09 / 0.9248	30.00 / 0.8350	28.96 / 0.8001	27.56 / 0.8376	32.51 / 0.9369
	SelNet [51]	34.27 / 0.9257	30.30 / 0.8399	28.97 / 0.8025	-	-
	CARN [21]	34.29 / 0.9255	30.29 / 0.8407	29.06 / 0.8034	28.06 / 0.8493	33.50 / 0.9440
	IMDN [34]	34.36 / 0.9270	30.32 / 0.8417	<u>29.09 / 0.8046</u>	28.17 / 0.8519	<u>33.61 / 0.9445</u>
	RAN [52]	33.71 / 0.9223	29.84 / 0.8326	28.84 / 0.7981	-	-
	DNCL [53]	33.95 / 0.9232	29.93 / 0.8340	28.91 / 0.7995	27.27 / 0.8326	-
	FilterNet [54]	34.08 / 0.9250	30.03 / 0.8370	28.95 / 0.8030	27.55 / 0.8380	-
	MRFN [44]	34.21 / 0.9267	30.03 / 0.8363	28.99 / 0.8029	27.53 / 0.8389	32.82 / 0.9396
	SeaNet-baseline [17]	34.36 / 0.9280	30.34 / 0.8428	29.09 / 0.8053	28.17 / 0.8527	33.40 / 0.9444
DEGREE [18]	33.76 / 0.9211	29.82 / 0.8326	28.74 / 0.7950	-	-	
Cross-SRN (Ours)	34.43 / 0.9275	30.33 / 0.8417	29.09 / 0.8050	28.23 / 0.8535	33.65 / 0.9448	
$\times 4$	SRCNN [2]	30.48 / 0.8628	27.49 / 0.7503	26.90 / 0.7101	24.52 / 0.7221	27.66 / 0.8505
	FSRCNN [45]	30.71 / 0.8657	27.59 / 0.7535	26.98 / 0.7150	24.62 / 0.7280	27.90 / 0.8517
	VDSR [25]	31.35 / 0.8838	28.01 / 0.7674	27.29 / 0.7251	25.18 / 0.7524	28.83 / 0.8809
	DRCN [46]	31.53 / 0.8854	28.02 / 0.7670	27.23 / 0.7233	25.14 / 0.7510	28.98 / 0.8816
	CNF [47]	31.55 / 0.8856	28.15 / 0.7680	27.32 / 0.7253	-	-
	LapSRN [20]	31.54 / 0.8850	28.19 / 0.7720	27.32 / 0.7280	25.21 / 0.7560	29.09 / 0.8845
	DRRN [46]	31.68 / 0.8888	28.21 / 0.7720	27.38 / 0.7284	25.44 / 0.7638	29.46 / 0.8960
	BTSRN [49]	31.85 / -	28.20 / -	27.47 / -	25.74 / -	-
	MemNet [50]	31.74 / 0.8893	28.26 / 0.7723	27.40 / 0.7281	25.50 / 0.7630	29.42 / 0.8942
	SelNet [51]	32.00 / 0.8931	28.49 / 0.7783	27.44 / 0.7325	-	-
	SRDenseNet [55]	32.02 / 0.8934	28.50 / 0.7782	27.53 / 0.7337	26.05 / 0.7819	-
	CARN [21]	32.13 / 0.8937	28.60 / 0.7806	27.58 / 0.7349	26.07 / 0.7837	30.47 / 0.9084
	IMDN [34]	<u>32.21 / 0.8948</u>	<u>28.58 / 0.7811</u>	27.56 / 0.7353	<u>26.04 / 0.7838</u>	30.45 / 0.9075
	RAN [52]	31.43 / 0.8847	28.09 / 0.7691	27.31 / 0.7260	-	-
	DNCL [53]	31.66 / 0.8871	28.23 / 0.7717	27.39 / 0.7282	25.36 / 0.7606	-
	FilterNet [54]	31.74 / 0.8900	28.27 / 0.7730	27.39 / 0.7290	25.53 / 0.7680	-
MRFN [44]	31.90 / 0.8916	28.31 / 0.7746	27.43 / 0.7309	25.46 / 0.7654	29.57 / 0.8962	
SeaNet-baseline [17]	32.18 / 0.8948	28.61 / 0.7822	27.57 / 0.7359	26.05 / 0.7896	30.44 / 0.9088	
DEGREE [18]	31.47 / 0.8837	28.10 / 0.7669	27.20 / 0.7216	-	-	
Cross-SRN (Ours)	32.24 / 0.8954	28.59 / 0.7817	27.58 / 0.7364	26.16 / 0.7881	30.53 / 0.9081	

Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1201–1205, 2015.

- [14] S. Huang, J. Sun, Y. Yang, Y. Fang, P. Lin, and Y. Que, “Robust single-image super-resolution based on adaptive edge-preserving smoothing regularization,” *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 2650–2663, 2018.
- [15] S. Pelletier and J. R. Cooperstock, “Preconditioning for edge-preserving image super resolution,” *IEEE Transactions on Image Processing*, vol. 21, no. 1, pp. 67–79, 2012.
- [16] L. Wang, S. Xiang, G. Meng, H. Wu, and C. Pan, “Edge-directed single-image super-resolution via adaptive gradient magnitude self-interpolation,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 8, pp. 1289–1299, 2013.
- [17] F. Fang, J. Li, and T. Zeng, “Soft-edge assisted network for single image super-resolution,” *IEEE Transactions on Image Processing*, vol. 29, pp. 4656–4668, 2020.
- [18] W. Yang, J. Feng, J. Yang, F. Zhao, J. Liu, Z. Guo, and S. Yan, “Deep edge guided recurrent residual learning for image super-resolution,” *IEEE Transactions on Image Processing*, vol. 26, no. 12, pp. 5895–5907, 2017.
- [19] J. Huang, A. Singh, and N. Ahuja, “Single image super-resolution from transformed self-exemplars,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5197–5206, 2015.
- [20] W. Lai, J. Huang, N. Ahuja, and M. Yang, “Deep laplacian pyramid

TABLE VI: Average PSNR/SSIM on selected benchmark with $\mathbf{BI}\times 4$ degradation.

	LR	MS-LapSRN [27]	VDSR [25]	IMDN [34]	Cross-SRN
Urban100 (S)	18.59 / 0.6356	21.80 / 0.7916	21.37 / 0.7696	22.44 / 0.8074	22.71 / 0.8145
Mangal09 (S)	20.39 / 0.7314	25.11 / 0.8795	-	25.96 / 0.8913	26.03 / 0.8890
Selected	19.54 / 0.6861	23.55 / 0.8379	-	24.30 / 0.8515	24.46 / 0.8537

TABLE VII: Running time and PSNR/SSIM comparisons on Set5 with $\mathbf{BI}\times 4$ degradation.

Method	LapSRN [27]	SRDenseNet [55]	SeaNet (Baseline) [17]	CARN [21]	Ours
PSNR	31.54	32.02	32.18	32.13	32.24
SSIM	0.8850	0.8934	0.8948	0.8937	0.8954
Time Cost (ms)	62.6	304.8	717.5	41.0	62.4

networks for fast and accurate super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5835–5843, 2017.

- [21] N. Ahn, B. Kang, and K. Sohn, “Fast, accurate, and lightweight super-resolution with cascading residual network,” in *European Conference on Computer Vision*, vol. 11214, pp. 256–272, 2018.
- [22] J. Li, F. Fang, K. Mei, and G. Zhang, “Multi-scale residual network for image super-resolution,” in *European Conference on Computer Vision*, vol. 11212, pp. 527–542, 2018.
- [23] Y. Liu, S. Wang, J. Zhang, S. Wang, S. Ma, and W. Gao, “Iterative network for image super-resolution,” *IEEE Transactions on Multimedia*, pp. 1–1, 2021.
- [24] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, “Squeeze-and-excitation networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 8, pp. 2011–2023, 2020.
- [25] J. Kim, J. K. Lee, and K. M. Lee, “Accurate image super-resolution using very deep convolutional networks,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1646–1654, 2016.
- [26] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, “Enhanced deep residual networks for single image super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1132–1140, 2017.
- [27] W. Lai, J. Huang, N. Ahuja, and M. Yang, “Fast and accurate image super-resolution with deep laplacian pyramid networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 11, pp. 2599–2613, 2019.
- [28] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, “Image super-resolution using very deep residual channel attention networks,” in *European Conference on Computer Vision*, vol. 11211, pp. 294–310, 2018.
- [29] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, “Residual dense network for image restoration,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.
- [30] C. Ma, Y. Rao, Y. Cheng, C. Chen, J. Lu, and J. Zhou, “Structure-preserving super resolution with gradient guidance,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7766–7775, 2020.
- [31] Z. Hui, X. Wang, and X. Gao, “Fast and accurate single image super-resolution via information distillation network,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 723–731, 2018.
- [32] J. Liu, J. Tang, and G. Wu, “Residual feature distillation network for lightweight image super-resolution,” in *European Conference on Computer Vision*, vol. 12537, pp. 41–55, 2020.
- [33] T. Dai, J. Cai, Y. Zhang, S. Xia, and L. Zhang, “Second-order attention network for single image super-resolution,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11057–11066, 2019.
- [34] Z. Hui, X. Gao, Y. Yang, and X. Wang, “Lightweight image super-resolution with information multi-distillation network,” in *ACM International Conference on Multimedia*, p. 2024–2032, 2019.
- [35] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2818–2826, 2016.
- [36] A. Gholami, K. Kwon, B. Wu, Z. Tai, X. Yue, P. Jin, S. Zhao, and K. Keutzer, “SqueezeNext: Hardware-aware neural network design,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1719–171909, 2018.
- [37] E. Agustsson and R. Timofte, “Ntire 2017 challenge on single image super-resolution: Dataset and study,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1122–1131, 2017.
- [38] 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.
- [39] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations*, 2015.
- [40] M. Bevilacqua, A. Roumy, C. Guillemot, and M. line Alberi Morel, “Low-complexity single-image super-resolution based on nonnegative neighbor embedding,” in *British Machine Vision Conference*, pp. 135.1–135.10, 2012.
- [41] R. Zeyde, M. Elad, and M. Protter, “On single image scale-up using sparse-representations,” in *Curves and Surfaces*, (Berlin, Heidelberg), pp. 711–730, Springer Berlin Heidelberg, 2012.
- [42] D. Martin, C. Fowlkes, D. Tal, and J. Malik, “A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics,” in *IEEE International Conference on Computer Vision*, vol. 2, pp. 416–423 vol.2, 2001.
- [43] Y. Matsui, K. Ito, Y. Aramaki, A. Fujimoto, T. Ogawa, T. Yamasaki, and K. Aizawa, “Sketch-based manga retrieval using manga109 dataset,” *Multimedia Tools and Applications*, vol. 76, no. 20, pp. 21811–21838, 2017.
- [44] Z. He, Y. Cao, L. Du, B. Xu, J. Yang, Y. Cao, S. Tang, and Y. Zhuang, “MRFN: multi-receptive-field network for fast and accurate single image super-resolution,” *IEEE Transactions on Multimedia*, vol. 22, no. 4, pp. 1042–1054, 2020.
- [45] C. Dong, C. C. Loy, and X. Tang, “Accelerating the super-resolution convolutional neural network,” in *European Conference on Computer Vision*, pp. 391–407, 2016.
- [46] J. Kim, J. K. Lee, and K. M. Lee, “Deeply-recursive convolutional network for image super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1637–1645, 2016.
- [47] H. Ren, M. El-Khamy, and J. Lee, “Image super resolution based on fusing multiple convolution neural networks,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1050–1057, 2017.
- [48] Y. Tai, J. Yang, and X. Liu, “Image super-resolution via deep recursive residual network,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2790–2798, 2017.
- [49] Y. Fan, H. Shi, J. Yu, D. Liu, W. Han, H. Yu, Z. Wang, X. Wang, and T. S. Huang, “Balanced two-stage residual networks for image super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1157–1164, 2017.
- [50] Y. Tai, J. Yang, X. Liu, and C. Xu, “Memnet: A persistent memory network for image restoration,” in *IEEE International Conference on Computer Vision*, pp. 4549–4557, 2017.
- [51] J. Choi and M. Kim, “A deep convolutional neural network with selection units for super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1150–1156, 2017.
- [52] Y. Wang, L. Wang, H. Wang, and P. Li, “Resolution-aware network for image super-resolution,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 5, pp. 1259–1269, 2019.
- [53] C. Xie, W. Zeng, and X. Lu, “Fast single-image super-resolution via deep network with component learning,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 12, pp. 3473–3486, 2019.
- [54] F. Li, H. Bai, and Y. Zhao, “Filternet: Adaptive information filtering network for accurate and fast image super-resolution,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 6, pp. 1511–1523, 2020.
- [55] T. Tong, G. Li, X. Liu, and Q. Gao, “Image super-resolution using dense skip connections,” in *IEEE International Conference on Computer Vision*, pp. 4809–4817, 2017.