

# MULTI-SCALE ATTENTION NETWORK FOR SINGLE IMAGE SUPER-RESOLUTION

Yan Wang\* Yusen Li Gang Wang Xiaoguang Liu

Nankai-Baidu Joint Lab, Nankai University, China

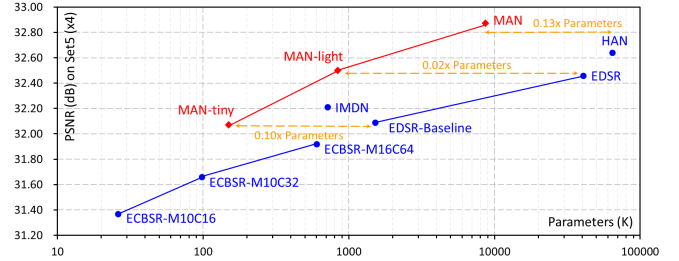
## ABSTRACT

By exploiting large kernel decomposition and attention mechanisms, convolutional neural networks (CNN) can compete with transformer-based methods in many high-level computer vision tasks. However, due to the advantage of long-range modeling, the transformers with self-attention still dominate the low-level vision, including the super-resolution task. In this paper, we propose a CNN-based multi-scale attention network (MAN), which consists of multi-scale large kernel attention (MLKA) and a gated spatial attention unit (GSAU), to improve the performance of convolutional SR networks. Within our MLKA, we rectify LKA with multi-scale and gate schemes to obtain the abundant attention map at various granularity levels, therefore jointly aggregating global and local information and avoiding the potential blocking artifacts. In GSAU, we integrate gate mechanism and spatial attention to remove the unnecessary linear layer and aggregate informative spatial context. To confirm the effectiveness of our designs, we evaluate MAN with multiple complexities by simply stacking different numbers of MLKA and GSAU. Experimental results illustrate that our MAN can achieve varied trade-offs between state-of-the-art performance and computations. Code is available at <https://github.com/icandle/MAN>.

**Index Terms**— Image Super-Resolution, Image Processing, Convolutional Neural Networks, Multi-scale Attention

## 1. INTRODUCTION

Image super-resolution (SR) is a widely concerned low-level computer vision task that focuses on rebuilding the missing high-frequency information from the low-quality (LQ) input [1, 2, 3, 4]. However, it is ill-posed that one low-resolution (LR) image corresponds to countless potential high-resolution (HR) images, leading to difficulties in learning the proper correlations between the LR and HR pixels. With the recently rapid development of deep neural networks, plenty of CNN- and transformer-based SR models [5, 6, 7, 8, 9] have been introduced to exploit prior and intra-image information to improve the reconstruction quality. Generally, they attempt to resolve the issues from the following two perspectives.



**Fig. 1.** Trade-off between performance and model size on Set5 [13] with  $\times 4$  SR scale. MAN family can achieve higher PSNR with fewer parameters.

The first and the simplest way is to enlarge the model capacity by training the network with larger datasets and better strategies. Specifically, IPT [10] used ImageNet [11] to accomplish a sophisticated pre-training. RCAN-it [12] leveraged reasonable training strategies to help RCAN [8] regain the state-of-the-art (SOTA) performance. Generally, these approaches are universal for all neural SR models but increase training and data processing consumption.

Another effective way is to activate more intra-image information via better network design. The primary idea is to enlarge the perceptive fields, which means a deeper and wider network topology. Following this, [4, 7, 8] continuously expanded the network depth to 200 layers. Nevertheless, the improvement brought by this strategy is capped due to the over-fitting and huge training cost. With a view to capturing more helpful information, the complex topology and attention mechanism have been employed in [14, 15, 16]. In detail, the RDN [17] utilized residual dense connection to fully exploit the hierarchical features. The MSRN [18] introduced the multi-scale mechanism, which can be regarded as a parallel topology, to detect the features at various scales. The NLSA [19] proposed the non-local sparse attention to attain long-range information with linear computational complexity. Recently, the transformer-based SR models [10, 9, 20] have shown the impressive representation ability of self-attention (SA). However, the quadratic complexity of calculating SA brings inconvenience in low-level tasks. Besides, SA ignores learning the essential texture from local features. To enable long-range dependency with simple convolutions, VAN [21] explored the kernel decomposition and proposed

\*Corresponding author. E-mail: wangy@nbjl.nankai.edu.cn.

large kernel attention (LKA), where a large kernel can be replaced by stacked depth-wise, depth-wise dilation, and point-wise convolution layers. Despite LKA’s efficiency, the dilation convolution causes blocking artifacts, which hurts the restored performance. Besides, a fixed-size LKA is inflexible to fully exploit the pending features because surrounding and remote pixels are equally important during reconstructing.

Motivated by these issues, we propose multi-scale large kernel attention (MLKA) that combines multi-scale mechanism and LKA to build various-range correlations with relatively few computations. Moreover, to avoid potential block artifacts aroused by dilation, we adopt the gate mechanism to recalibrate the generated attention maps adaptively. To maximize the benefits of MLKA, we place it on the MetaFormer [22]-style structure rather than RCAN [8]-style to construct a multi-attention block (MAB). Although transformer-style MAB can reach higher performance, the MLP feed-forward module is too heavy for large images. Inspired by simple-gate (SG) [23] and convolutional feed-forward network (CFF) [24], we propose a simplified gated spatial attention unit (GSAU), by applying spatial attention and gate mechanism to reduce calculations and include spatial information. Arming with MLKA and GSAU, the MABs are stacked to build the multi-scale attention network (MAN) for the SR task. In Fig. 1, we present the superior performance of our MAN.

To summarize, our contributions are as follows:

- We propose multi-scale large kernel attention (MLKA) to obtain long-range dependencies at various granularity levels, therefore significantly improving the model representation capability.
- We integrate the gate mechanism and spatial attention to build a simplified feed-forward network, GSAU, which can reduce parameters and calculations while maintaining the performance.
- We present multi-scale attention network (MAN) family that can achieve a trade-off between the performance and model complexity in both lightweight and performance-oriented SR tasks.

## 2. RELATED WORK

### 2.1. Single Image Super-Resolution

Numerous deep-learning models [4, 6, 25] have been proposed for SISR since the pioneering work SRCNN [5] introduced a 3-layers convolutional neural network (CNN) to map the correlation between LR and HR images. According to the model complexity, we can divide these solutions into two categories: the classical SR for desktop GPU, and the tiny and lightweight SR for mobile devices.

For the classical SR task, models are delicately designed for better reconstruction quality. Specifically, VDSR [4] and EDSR [7] were proposed to exploit deeper information by residual learning and increasing depth and width. RCAN [8] then developed EDSR by introducing channel attention (CA) and residual in residual (RIR) to further excavate intermediate features. Following RCAN, many works [26, 27, 16] inserted attention mechanism into the EDSR structure to further boost the performance. Very recently, vision transformers [10, 9] with self-attention (SA) are introduced to image restoration and refresh the SOTA performance.

For tiny and lightweight SR, the model size is constrained for mobile device deployment. The recursive learning was considered effective in decreasing the parameters in DRCN [28], DRRN [29], and LapSRN [30]. However, recursively using modules only reduces model size but maintains high computation costs. More recent works leverage productive operations, *e.g.*, channel splits and attention module, to exploit the hierarchical features. For example, IMDN [31] proposed information multi-distillation and contrast-aware channel attention.

### 2.2. Attentions in Super-Resolution

The attention mechanism can be viewed as a discriminative selection process that makes models focus on informative regions and ignore the irrelevant noise of pending features. Many SR networks apply attention modules to exploit latent correlations among the immediate features. Following RCAN [8] that first adopted channel attention, SAN [26] leveraged second-order channel attention to adapt the channel-wise features through second-order statistics. Several works introduced spatial attention to enrich the feature maps, *e.g.*, enhanced spatial attention in RFANet [14], and spatial-channel attention in HAN [27]. Additional CNN-based works have utilized and refined non-local attention (NLA) to obtain long-range correlations [19, 16] and achieved appreciable performance gain. Inspired by vision transformers [32, 24], self-attention has been employed in SR to capture long-term adaptability, *e.g.*, IPT [10] and SwinIR [9].

## 3. METHODOLOGY

### 3.1. Network Architecture

As illustrated in Fig. 2, the proposed MAN is constituted of three components: the shallow feature extraction module (SF), the deep feature extraction module (DF) based on multiple multi-scale attention blocks (MAB), and the high-quality image reconstruction module. Given an input LR image  $I_{LR} \in \mathbb{R}^{3 \times H \times W}$ , the SF module is first utilized to extract the primitive feature  $F_p \in \mathbb{R}^{C \times H \times W}$  by a single  $3 \times 3$  convolution function  $f_{SF}(\cdot)$  as:

$$F_p = f_{SF}(I_{LR}). \quad (1)$$

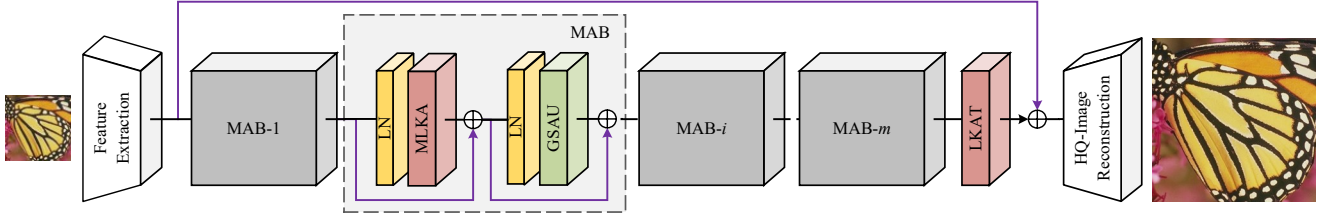


Fig. 2. Overview of our multi-scale attention network (MAN).

The  $F_p$  is then sent to cascading MABs for further extraction, termed as  $f_{DF}(\cdot)$ , which can be formulated as:

$$F_r = f_{DF}(F_p), \quad (2)$$

where the  $F_r$  is the estimated high-frequency feature for final restoration. By adding the long residual feature, the final reconstruction component restores the HQ images  $I_{SR} \in \mathbb{R}^{3 \times H \times W}$  by:

$$I_{SR} = f_{RC}(F_p + F_r), \quad (3)$$

where  $f_{RC}(\cdot)$  represents reconstruction module implemented by a  $3 \times 3$  convolution and a pixel-shuffle layer for efficiency.

In terms of optimization, we utilize the commonly used  $\ell_1$  loss for fair comparison with state-of-the-art methods [7, 8, 27]. Specifically, supposing an input batch of  $N$  images, i.e.  $\{I_i^{LR}, I_i^{HR}\}_{i=1}^N$ , training MAN is to minimize the  $\ell_1$ :

$$\ell_1(\Theta) = \frac{1}{N} \sum_{i=1}^N \|f_{MAN}(I_i^{LR}) - I_i^{HR}\|_1 \quad (4)$$

where  $f_{MAN}(\cdot)$  is the proposed network and  $\Theta$  denotes its trainable parameters.

### 3.2. Multi-scale Attention Block (MAB)

Inspired by recent breakthroughs on transformers, we reconsider the basic convolutional block for feature extraction in SISR task. Different to many RCAN-style [8] blocks, the proposed MAB absorbs MetaFormer [22] style design to obtain a promising extraction capability. As shown in Fig. 3, MAB consists of two components: the multi-scale large kernel attention (MLKA) module and the gate spatial attention unit (GSAU).

Given the input feature  $X$ , the whole process of MAB is formulated as:

$$\begin{aligned} N &= LN(X), \\ X &= X + \lambda_1 f_3(MLKA(f_1(N)) \otimes f_2(N)), \\ N &= LN(X), \\ X &= X + \lambda_2 f_6(GSAU(f_4(N), f_5(N))), \end{aligned} \quad (5)$$

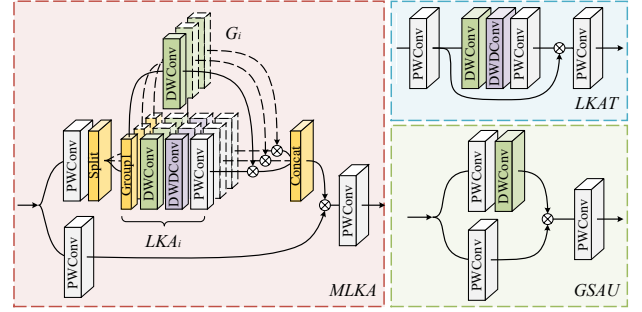


Fig. 3. Details of proposed modules.

where  $LN(\cdot)$  and  $\lambda$  are layer normalization and learnable scaling factors, separately.  $MLKA(\cdot)$  and  $GSAU(\cdot)$  are proposed MLKA and GSAU module introduced in following sections.  $\otimes$  and  $f_i(\cdot)$  represent element-wise multiplication and  $i$ -th point-wise convolution that keeps the dimensions.

### 3.3. Multi-scale Large Kernel Attention (MLKA)

Attention mechanisms can force networks to focus on crucial information and ignore irrelevant ones. Previous SR models adopt a series of attention mechanisms, including channel attention (CA) and self-attention (SA), to obtain more informative features. However, these methods fail to simultaneously uptake local information and long-range dependence, and they often consider the attention maps at a fixed reception field. Enlightened by the latest visual attention researches [21], we propose multi-scale large kernel attention (MLKA) to resolve these problems by combining large kernel decomposition and multi-scale learning. Specifically, the MLKA consists of three main functions, large kernel attentions (LKA) for establishing interdependence, multi-scale mechanism for obtaining heterogeneous-scale correlation, and the gated aggregation for dynamic recalibration.

**Large kernel attention.** Given the input feature maps  $X \in \mathbb{R}^{C \times H \times W}$ , the LKA adaptively builds the long-range relationship by decomposing a  $K \times K$  convolution into three components: a  $(2d - 1) \times (2d - 1)$  depth-wise convolution  $f_{DW}(\cdot)$ , a  $\lceil \frac{K}{d} \times \frac{K}{d} \rceil$  depth-wise  $d$ -dilation convolution

$f_{DWD}(\cdot)$ , and a point-wise convolution  $f_{PW}(\cdot)$ , which can be formulated as:

$$LKA(X) = f_{PW}(f_{DWD}(f_{DW}(X))). \quad (6)$$

**Multi-scale mechanism.** To learn the attention maps with omni-scale information, we enhance the LKA with group-wise multi-scale mechanism. Supposing the input feature maps  $X \in \mathbb{R}^{C \times H \times W}$ , the module first splits it into  $n$ -pieces  $X_1, X_2, \dots, X_n$  of  $\lfloor \frac{C}{n} \rfloor \times H \times W$ . For  $i$ -th group of features  $X_i$ , a LKA decomposed by  $\{K_i, d_i\}$  is utilized to generate a homogeneous scale attention map  $LKA_i$ . In detail, we leverage three groups of LKA:  $\{7, 2\}$  implemented by 3-5-1,  $\{21, 3\}$  by 5-7-1, and  $\{35, 4\}$  by 7-9-1, where  $a$ - $b$ -1 means cascading  $a \times a$  depth-wise convolution,  $b \times b$  depth-wise-dilated convolution, and point-wise convolution.

**Gated aggregation.** Different to many high-level computer vision tasks, the SR task has a worse tolerance for dilation and partition. As shown in the Fig. 4, although the larger LKA captures wider responses of pixels, the blocking artifacts appear in the generated attention maps of larger LKA. For  $i$ -th group input  $X_i$ , to avoid the block effect, as well as to learn more local information, we leverage spatial gate to dynamically adapt  $LKA_i(\cdot)$  into  $MLKA_i(\cdot)$  by:

$$MLKA_i(X_i) = G_i(X_i) \otimes LKA_i(X_i), \quad (7)$$

where  $G_i(\cdot)$  is the  $i$ -th gate generated by  $a_i \times a_i$  depth-wise convolution, and  $LKA_i(\cdot)$  is the LKA decomposed by  $a_i$ - $b_i$ -1. In Fig. 4, we provide the visual results of the gated aggregation. It can be observed that the block effects are removed from the attention maps and the  $MLKA_i$ s are more reasonable. In particular, the LKA with larger receptive fields react more on long-range dependence while the smaller LKA tends to retain local texture.

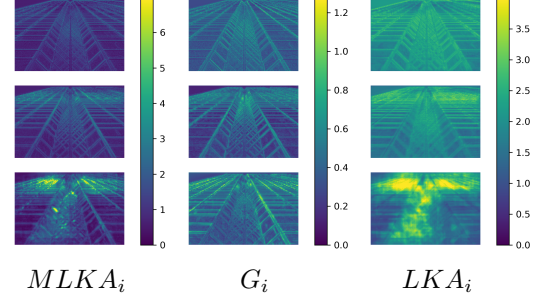
### 3.4. Gated Spatial Attention Unit (GSAU)

In transformer blocks, the feed-forward network (FFN) is a significant component to enhance the feature representation. However, the MLP with wide intermediate channels is too heavy for SR, especially for large image inputs. Inspired by [33, 23, 34, 24], we integrate simple spatial attention (SSA) and gated linear unit (GLU) into the proposed GSAU to enable an adaptive gating mechanism and reduce the parameters and calculations.

To capture spatial information more efficiently, we adopt a single layer depth-wise convolution to weight the feature map. Given the dense-transformed  $X$  and  $Y$ , the key process of GSAU can be represented as:

$$GSAU(X, Y) = f_{DW}(X) \otimes Y, \quad (8)$$

where  $f_{DW}(\cdot)$  and  $\otimes$  indicate depth-wise convolution and element-wise multiplication, respectively. By applying spatial gate, the GSAU can remove the nonlinear layer and capture local continuity under considerate complexity.



**Fig. 4.** Visual activation maps of Eq. (8) in the 16-th layer of MAN-light. From top to bottom are the results of 3-5-1, 5-7-1, and 7-9-1, respectively.

### 3.5. Large Kernel Attention Tail (LKAT)

In previous SR networks [7, 8, 26, 27, 9], the  $3 \times 3$  convolution layer is widely used as the tail of the deep extraction backbone. However, it has a flaws in establish long-range connections, therefor limiting the representative capability of the final reconstruction feature. In order to summarize more reasonable information from the stacked MABs, we introduce the 7-9-1 LKA in the tail module. Concretely, the LKA is wrapped by two  $1 \times 1$  convolutions as depicted in Fig. 3. The experimental results show that the LKAT module can efficiently aggregate the useful information and significantly improve the reconstruction quality.

## 4. EXPERIMENTS

### 4.1. Datasets and Metrics

Following latest works [35, 9, 12], we utilize DIV2K [36] and Flickr2K [7], which contain 800 and 2650 training images, to train our models. For testing, we valuate our method on five commonly used datasets: Set5 [13], Set14 [37], BSD100 [38], Urban100 [39], and Manga109 [40]. In addition, two standard evaluation metrics, peak-signal-to-noise-ratio (PSNR) and the structural similarity index (SSIM) [41], are applied in  $Y$  channel of the YCbCr images to measure the quality of restoration.

### 4.2. Implementation Details

For more comprehensive evaluation of the proposed methods, we train three different versions of MAN: tiny, light, and classical, to resolve the classic SR tasks under different complexity. Following [9], we stack 5/24/36 MABs and set the channel width to 48/60/180 in the corresponding tiny/light/classical MAN. Three multi-scale decomposition modes are utilized in MLKA, listed as 3-5-1, 5-7-1, and 7-9-1. The  $7 \times 7$  depth-wise convolution is used in the GSAU.

In the training stage, bicubic interpolation is used to generate LR-HR image pairs. The training pairs are further augmented by horizontal flip and random rotations of  $90^\circ$ ,  $180^\circ$ ,



**Table 1.** Ablation studies on components of MAN. The impact of LKAT, Multi-scale mechanism, and GSAU are shown upon MAN-tiny/light ( $\times 2$ ). In detail, we replace LKAT with convolution layer, Multi-scale with LKA(5-7-1), and GSAU with MLP.

Method	LKAT	Multi-Scale	GSAU	#Params	#Mult-Adds	Set5		Set14		BSD100		Urban100	
						PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
MAN-tiny	✓			108K	24.2G	37.71	0.9594	33.24	0.9148	31.97	0.8973	31.08	0.9178
	✓			121K	27.2G	37.75	0.9595	33.27	0.9154	32.01	0.8979	31.25	0.9199
	✓	✓		133K	29.9G	37.77	0.9596	33.30	0.9153	32.01	0.8979	31.30	0.9202
	✓	✓	✓	134K	29.9G	37.79	0.9598	33.31	0.9155	32.02	0.8980	31.33	0.9206
MAN-light	✓			737K	165.8G	38.01	0.9605	33.55	0.9179	32.23	0.9005	32.14	0.9287
	✓			756K	170.0G	38.05	0.9607	33.60	0.9182	32.25	0.9007	32.23	0.9297
	✓	✓		835K	187.6G	38.07	0.9607	33.62	0.9181	32.26	0.9009	32.42	0.9308
	✓	✓	✓	820K	184.0G	38.07	0.9608	33.69	0.9188	32.29	0.9012	32.43	0.9316

**Table 2.** Ablation study on block structure. The performance of RCAN-style and transformer-style MABs are shown upon MAN-light on  $\times 4$  SR task.

Method	#Params	#Mult-Adds	Set5	
			PSNR	SSIM
RCAN-style	924K	53.0G	32.16	0.8945
Metaformer-style	840K	47.1G	32.33	0.8967

270°. The {patch size, batch size} is set to  $\{48 \times 48, 32\}$  and  $\{64 \times 64, 16\}$  in the training-from-scratch and fine-tuning stage, respectively. The  $\ell_1$  loss is adopted to discriminate the pixel-wise restoration quality. All models are trained using the Adam optimizer [42] with  $\beta_1=0.9$ ,  $\beta_2=0.99$ . The learning rate is initialized as  $5e^{-4}$  and scheduled by cosine annealing learning for 160K iterations in training anew, while setting as  $1e^{-4}$  for 80K in fine-tuning. For ablation studies, we train all models in 20K iterations. All experiments are conducted by Pytorch [43] framework on 4 Nvidia RTX 3090 GPUs.

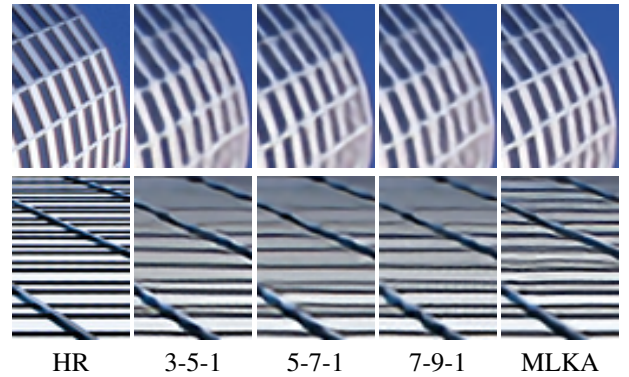
### 4.3. Ablation Studies

In this section, we validate the effectiveness of the proposed components from coarse to fine. In detail, we first investigate the combination of all proposed module and then investigate each of them individually.

**Overall study on components of MAN.** In the Table 1, we present the results of deploying proposed components on our tiny and light networks. Comparing the base models without proposed modules, we can find that employing any of them can improve performance. Specifically, 0.25 dB and 0.29 dB promoting on Urban100 [39] can be observed in MAN-tiny and MAN-light, while the parameters and calculations increase negligibly. Among these components, the LKAT module and multi-scale mechanism are more significant ones to enhance quality. Without any of them, the PSNR will drop by 0.09 dB. The GSAU is an economical replacement for MLP. It reduces 15K parameter and 3.6G calculations while bringing significant improvements across all datasets.

**Table 3.** Ablation study on GSAU. The MLP [44], Simple-Gate (SG) [23], CFF [24], and GSAU are tested on MAN-light for  $\times 4$  SR.

Method	#Params	#MAAdd	Set5	Set14	B100	U100
			PSNR	PSNR	PSNR	PSNR
MLP	854K	48.0G	32.31	28.73	27.65	26.26
SG	768K	43.1G	32.28	28.74	27.66	26.28
CFF	1140K	64.3G	32.35	28.76	27.67	26.34
GSAU	840K	47.1G	32.33	28.76	27.67	26.31



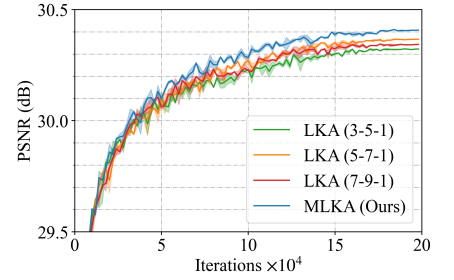
**Fig. 5.** Visual comparison between LKA and MLKA.

**Study on block structures.** Within MAB, we choose the Metaformer-style rather than the RCAN-style structure to deploy MLKA. To fully explore their effectiveness, we implement and compare two versions MABs in the Table 2. The experimental results indicate that the transformer-style MAB surpasses the RCAN-style one at a large margin. On Set5 [13], the PSNR is increased from 32.15 dB to 32.33 dB by employing transformer structure. The results show the transformer-style MAB is a more efficient choice to balance the performance and computations.

**Study on MLKA.** To justify our design of MLKA, we conduct ablation experiments on multi-scale and kernel decomposition. Specifically, we take three LKA and three MLKA implementations into consideration in Table 4. We

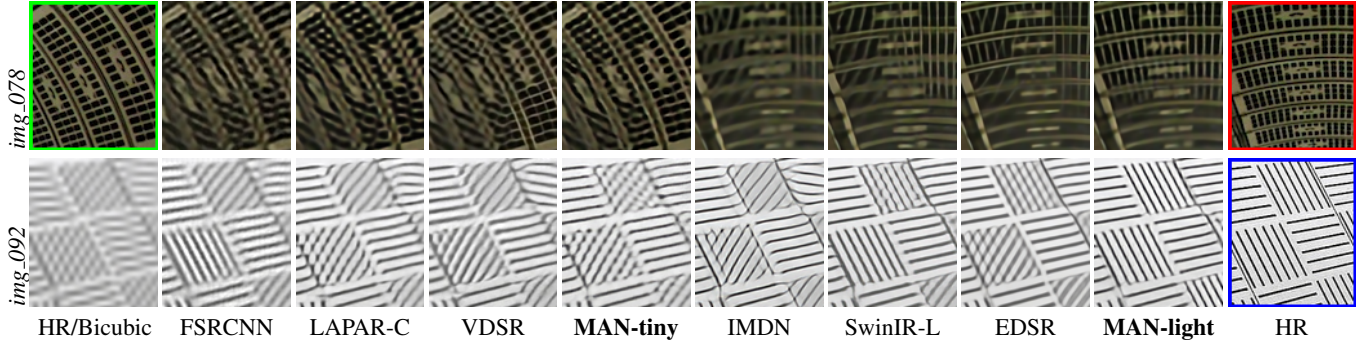
**Table 4.** Ablation studies of multi-scale and decomposition type (LKA/MLKA). The results are tested on MAN-light for  $\times 4$ .

Method	Decomposition			#Params	#Mult-Adds	Set5		Set14	
	3-5-1	5-7-1	7-9-1			PSNR	SSIM	PSNR	SSIM
LKA	✓			703K	39.4G	32.23	0.8956	28.70	0.7842
		✓		761K	42.7G	32.27	0.8963	28.72	0.7846
			✓	841K	47.4G	32.25	0.8958	28.71	0.7845
MLKA	✓	✓		803K	45.0G	32.32	0.8968	28.72	0.7848
		✓	✓	900K	50.6G	32.33	0.8968	28.74	0.7852
	✓	✓	✓	840K	47.1G	32.33	0.8967	28.76	0.7856



**Table 5.** Quantitative comparison (average PSNR/SSIM) with state-of-the-art approaches for **tiny/light image SR** on benchmark datasets ( $\times 4$ ). Best and second best performance are **highlighted** and underlined, respectively.

Method	Scale	#Params	#Mult-Adds	Set5		Set14		BSD100		Urban100		Manga109	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
FSRCNN	$\times 4$	12K	4.6G	30.71	0.8657	27.59	0.7535	26.98	0.7150	24.62	0.7280	27.90	0.8517
LAPAR-C	$\times 4$	115K	25.0G	<u>31.72</u>	0.8884	28.31	0.7740	27.40	0.7292	25.49	0.7651	<u>29.50</u>	<u>0.8951</u>
ECBSR-M10C32	$\times 4$	98K	5.7G	31.66	<u>0.8911</u>	<u>28.15</u>	<u>0.7776</u>	27.34	<u>0.7363</u>	25.41	<u>0.7653</u>	-	-
MAN-tiny (Ours)	$\times 4$	150K	8.4G	<b>32.07</b>	<b>0.8930</b>	<b>28.53</b>	<b>0.7801</b>	<b>27.51</b>	<b>0.7345</b>	<b>25.84</b>	<b>0.7786</b>	<b>30.18</b>	<b>0.9047</b>
EDSR-baseline	$\times 4$	1518K	114G	32.09	0.8938	28.58	0.7813	27.57	0.7357	26.04	0.7849	30.35	0.9067
IMDN	$\times 4$	715K	40.9G	32.21	0.8948	28.58	0.7811	27.56	0.7353	26.04	0.7838	30.45	0.9075
LatticeNet	$\times 4$	777K	43.6G	32.30	0.8962	28.68	0.7830	27.62	0.7367	26.25	0.7873	-	-
SwinIR-light	$\times 4$	897K	49.6G	<u>32.44</u>	<u>0.8976</u>	<u>28.77</u>	<u>0.7858</u>	<u>27.69</u>	<u>0.7406</u>	<u>26.47</u>	<u>0.7980</u>	<u>30.92</u>	<u>0.9151</u>
MAN-light (Ours)	$\times 4$	840K	47.1G	<b>32.50</b>	<b>0.8988</b>	<b>28.87</b>	<b>0.7885</b>	<b>27.77</b>	<b>0.7429</b>	<b>26.70</b>	<b>0.8052</b>	<b>31.25</b>	<b>0.9170</b>
EDSR	$\times 4$	43090K	2895G	32.46	0.8968	28.80	0.7876	27.71	0.7420	26.64	0.8033	31.02	0.9148



**Fig. 6.** Visual comparison for tiny/lightweight SR models with an upscaling factor  $\times 4$ .

first investigate the influence of kernel size upon LKA. When we increase the kernel size, the PSNR decreases after an initial increase, which is inconsistent with high-level tasks [21]. This is due to that both long-range correlation and local textural information are indispensable in image restoration task. And up until this point, we introduce MLKA to refine feature at comprehensive scales. In the Table 4, we also illustrate the training evaluation results of LKAs and proposed MLKA. The MLKA outperforms other LKAs throughout the training phase. For the visual comparison shown in Fig. 5, MLKA also helps to recover more details on both images from Urban100. In addition, we also discuss MLKA of different combinations. The results suggest the MLKA with all

three decomposition types can trade off between parameters, computations, and performance.

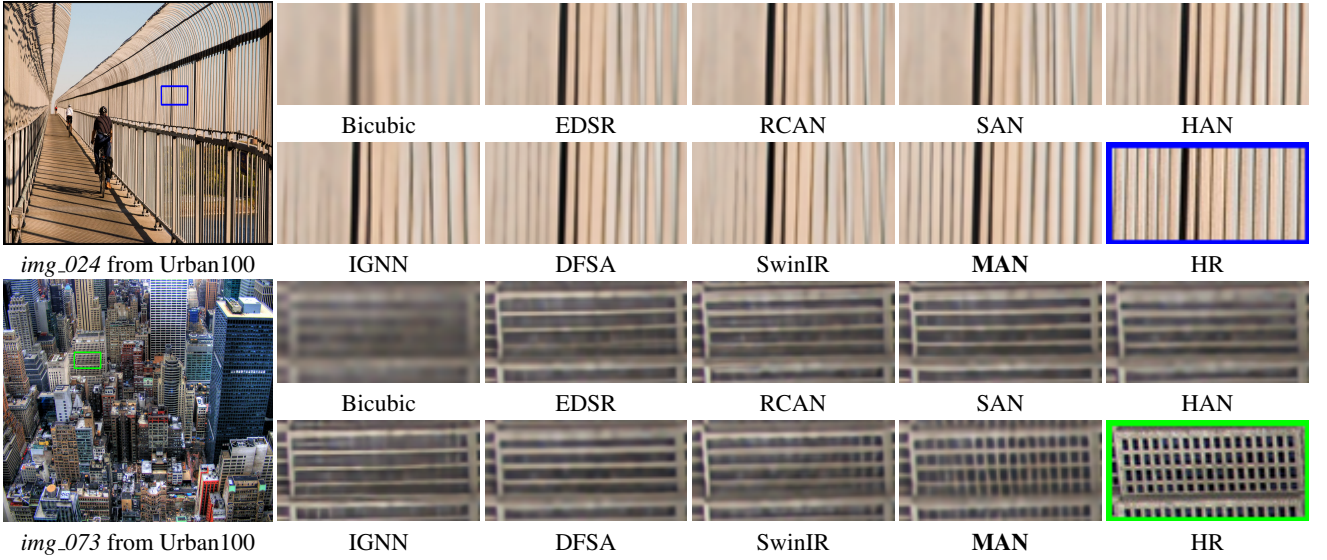
**Study on GSAU.** To further confirm the efficiency of the proposed GSAU, we compare it with some other FFNs. In Table 3, we validate four advanced designs: MLP, SG, CFF, and our GSAU. The evaluation results suggest that the GSAU achieves comparable performance with powerful CFF while occupying 73% of parameters and calculations.

#### 4.4. Comparisons with tiny and light-weight SR models

To verify the efficiency and scalability of our MAN, we compare MAN-tiny and MAN-light to some state-of-the-art tiny [6, 45, 46] and lightweight [31, 47, 9] SR mod-

**Table 6.** Quantitative comparison (average PSNR/SSIM) with CNN-based state-of-the-art approaches for **classical image SR**. Best and second best performance are **highlighted** and underlined, respectively. “†” and “+” indicate using pre-training and self-ensemble strategy, respectively.

Method	Scale	#Params	#FLOPs	Set5		Set14		BSD100		Urban100		Manga109	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
RCAN	$\times 2$	15.4M	3.5T	38.27	0.9614	34.12	0.9216	32.41	0.9027	33.34	0.9384	39.44	0.9786
SAN	$\times 2$	15.9M	3.1T	38.31	0.9620	34.07	0.9213	32.42	0.9028	33.10	0.9370	39.32	0.9792
HAN	$\times 2$	63.6M	14.6T	38.27	0.9614	34.16	0.9217	32.41	0.9027	33.35	0.9385	39.46	0.9785
IGNN	$\times 2$	49.5M	-	38.24	0.9613	34.07	0.9217	32.41	0.9025	33.23	0.9383	39.35	0.9786
NLSA	$\times 2$	41.8M	9.6T	38.34	0.9618	34.08	0.9231	32.43	0.9027	33.42	0.9394	39.59	0.9789
DFSA+	$\times 2$	-	-	38.38	0.9620	34.33	0.9232	32.50	0.9036	33.66	0.9412	39.98	0.9798
MAN (Ours)	$\times 2$	8.7M	1.7T	<u>38.42</u>	<u>0.9622</u>	<u>34.40</u>	<u>0.9242</u>	<u>32.53</u>	<u>0.9043</u>	<u>33.73</u>	<u>0.9422</u>	<u>40.02</u>	<u>0.9801</u>
MAN+ (Ours)	$\times 2$	8.7M	-	<b>38.44</b>	<b>0.9623</b>	<b>34.49</b>	<b>0.9248</b>	<b>32.55</b>	<b>0.9045</b>	<b>33.86</b>	<b>0.9430</b>	<b>40.13</b>	<b>0.9804</b>
SwinIR <sup>†</sup>	$\times 2$	11.8M	2.3T	38.42	0.9623	34.46	0.9250	32.53	0.9041	33.81	0.9427	39.92	0.9797
RCAN	$\times 3$	15.6M	1.6T	34.74	0.9299	30.65	0.8482	29.32	0.8111	29.09	0.8702	34.44	0.9499
SAN	$\times 3$	15.9M	1.6T	34.75	0.9300	30.59	0.8476	29.33	0.8112	28.93	0.8671	34.30	0.9494
HAN	$\times 3$	64.3M	6.5T	34.75	0.9299	30.67	0.8483	29.32	0.8110	29.10	0.8705	34.48	0.9500
IGNN	$\times 3$	49.5M	-	34.72	0.9298	30.66	0.8484	29.31	0.8105	29.03	0.8696	34.39	0.9496
NLSA	$\times 3$	44.7M	4.6T	34.85	0.9306	30.70	0.8485	29.34	0.8117	29.25	0.8726	34.57	0.9508
DFSA+	$\times 3$	-	-	34.92	<u>0.9312</u>	30.83	0.8507	29.42	0.8128	29.44	0.8761	<u>35.07</u>	0.9525
MAN (Ours)	$\times 3$	8.7M	0.8T	<u>34.91</u>	<u>0.9312</u>	<u>30.88</u>	<u>0.8514</u>	<u>29.43</u>	<u>0.8138</u>	<u>29.52</u>	<u>0.8782</u>	<u>35.06</u>	<u>0.9526</u>
MAN+ (Ours)	$\times 3$	8.7M	-	<b>34.97</b>	<b>0.9315</b>	<b>30.91</b>	<b>0.8522</b>	<b>29.47</b>	<b>0.8144</b>	<b>29.65</b>	<b>0.8799</b>	<b>35.21</b>	<b>0.9533</b>
SwinIR <sup>†</sup>	$\times 3$	11.9M	1.0T	34.97	0.9318	30.93	0.8534	29.46	0.8145	29.75	0.8826	35.12	0.9537
RCAN	$\times 4$	15.6M	0.9T	32.63	0.9002	28.87	0.7889	27.77	0.7436	26.82	0.8087	31.22	0.9173
SAN	$\times 4$	15.9M	0.9T	32.64	0.9003	28.92	0.7888	27.78	0.7436	26.79	0.8068	31.18	0.9169
HAN	$\times 4$	64.2M	3.8T	32.64	0.9002	28.90	0.7890	27.80	0.7442	26.85	0.8094	31.42	0.9177
IGNN	$\times 4$	49.5M	-	32.57	0.8998	28.85	0.7891	27.77	0.7434	26.84	0.8090	31.28	0.9182
NLSA	$\times 4$	44.2M	3.0T	32.59	0.9000	28.87	0.7891	27.78	0.7444	26.96	0.8109	31.27	0.9184
DFSA+	$\times 4$	-	-	32.79	0.9019	29.06	0.7922	27.87	0.7458	27.17	0.8163	31.88	<b>0.9266</b>
MAN (Ours)	$\times 4$	8.7M	0.4T	<u>32.81</u>	<u>0.9024</u>	<u>29.07</u>	<u>0.7934</u>	<u>27.90</u>	<u>0.7477</u>	<u>27.26</u>	<u>0.8197</u>	<u>31.92</u>	0.9230
MAN+ (Ours)	$\times 4$	8.7M	-	<b>32.87</b>	<b>0.9030</b>	<b>29.12</b>	<b>0.7941</b>	<b>27.93</b>	<b>0.7483</b>	<b>27.39</b>	<b>0.8223</b>	<b>32.13</b>	<u>0.9248</u>
SwinIR <sup>†</sup>	$\times 4$	11.9M	0.6T	32.92	0.9044	29.09	0.7950	27.92	0.7489	27.45	0.8254	32.03	0.9260



**Fig. 7.** Visual comparison for classical SR models with an upscaling factor  $\times 4$ .



els. In Table 5, we present the numerical results that our MAN-tiny/light surpasses all other tiny/lightweight methods. Specifically, MAN-tiny exceeds second place by about 0.2dB on Set5, Urban100, and Manga109, and around 0.07dB on Set14 and BSD100. We also list EDSR-baseline [7] for reference. It can be observed that our tiny model has less than 150K parameters but achieves a similar restoration quality with EDSR-baseline, which is  $10\times$  larger than ours. Similarly, our MAN-light surpasses both CNN-based and transformer-based SR models. Compared to IMDN (CNN) and SwinIR-light (transformer), our model leads a 0.66 dB/0.23 dB improvements on Urban100 ( $\times 4$ ) benchmark. Moreover, our MAN-light is excelled to some traditional normal SR models such as EDSR. In detail, the proposed model takes only 2% parameters and computations of EDSR while getting better results on all datasets.

In Fig. 6, we also exhibit the visual results of several tiny and lightweight models on Urban100 ( $\times 4$ ). For *img\_078*, the tiny and light models are tested with the patches framed by green and red boxed, respectively. Generally, our MANs can rescale the correct texture of the air vent well than other methods, including EDSR. Similar results can be observed in *img\_092*, where only MAN-light can rightly recover all line ornamentation. The experimental results confirm that even on the light-weight framework, our multi-scale attention design can efficiently capture long-term dependency and greatly improve the SR quality.

#### 4.5. Comparisons with classical SR models

To validate the effectiveness of our MAN, we compare our normal model to several SOTA CNN-based classical [8, 26, 48, 27, 19, 35, 12] SR models. We also add SwinIR [9] for reference. In Table 6, the quantitative results show that our MAN exceeds other CNN-based methods to a large extent. The maximum improvement on PSNR reaches 0.69 dB for  $\times 2$ , 0.77 dB for  $\times 3$ , and 0.81 dB for  $\times 4$ . Moreover, we compare our MAN with SwinIR. For  $\times 2$ , our MAN achieves competitive or even better performance than SwinIR. The PSNR value on Manga109 is promoted from 39.92 dB to 40.02 dB. For  $\times 4$ , MAN is slightly behind SwinIR because the latter uses  $\times 2$  model as the pre-trained model. More important, our MAN is significantly smaller than existing methods. Compared to HAN, it only takes 13% parameters and 11% calculations while advancing 0.2 dB-0.7 dB for  $\times 4$  SR.

In Fig. 7, we also visualize the qualitative results of several models on Urban100 ( $\times 4$ ) benchmark dataset. For *img\_024*, compared with other models generating the distorted fence, our MAN rebuilds a clear structure from the blurred input. Similarly, in *img\_073*, MAN is the only model that restores the windows of the building. The quantitative and qualitative results demonstrate the advantages of MAN in capturing the local and global structural information.

## 5. CONCLUSION

In this paper, we propose a multi-scale attention network (MAN) for rescaling super-resolution images under multiple environments. MAN adopts transformer-style block to obtain better modeling representation capability. In order to effectively and flexibly establish long-range correlations among various regions, we develop multi-scale large kernel attention (MLKA) that combines large kernel decomposition and multi-scale mechanisms. Furthermore, we propose a simplified feed-forward network that integrates gate mechanisms and spatial attention to activate local information and reduce model complexity. Extensive experiments show that our CNN-based MAN can achieve better performance than previous SOTA models in a more efficient manner.

## 6. REFERENCES

- [1] Jian Sun, Zongben Xu, and Heung-Yeung Shum, “Image super-resolution using gradient profile prior,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Anchorage, USA, 2008, pp. 1–8. 1
- [2] Yu-Wing Tai, Shuaicheng Liu, Michael S Brown, and Stephen Lin, “Super resolution using edge prior and single image detail synthesis,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, San Francisco, USA, 2010, pp. 2400–2407. 1
- [3] He He and Wan-Chi Siu, “Single image super-resolution using gaussian process regression,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Colorado Springs, USA, 2011, pp. 449–456. 1
- [4] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee, “Accurate image super-resolution using very deep convolutional networks,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA*. 2016, pp. 1646–1654, IEEE Computer Society. 1, 2
- [5] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang, “Image super-resolution using deep convolutional networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, 2016. 1, 2
- [6] Chao Dong, Chen Change Loy, and Xiaoou Tang, “Accelerating the super-resolution convolutional neural network,” in *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, Eds. 2016, vol. 9906, pp. 391–407, Springer. 1, 2, 6



- [7] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee, “Enhanced deep residual networks for single image super-resolution,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2017, pp. 1132–1140, IEEE Computer Society. 1, 2, 3, 4, 8
- [8] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu, “Image super-resolution using very deep residual channel attention networks,” in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany*. 2018, vol. 11211, pp. 294–310, Springer. 1, 2, 3, 4, 8
- [9] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte, “Swinir: Image restoration using swin transformer,” in *IEEE/CVF International Conference on Computer Vision Workshops, Montreal, Canada*. 2021, pp. 1833–1844, IEEE. 1, 2, 4, 6, 8
- [10] 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 *IEEE Conference on Computer Vision and Pattern Recognition, virtual*. 2021, pp. 12299–12310, Computer Vision Foundation / IEEE. 1, 2
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Miami, USA*. 2009, pp. 248–255, IEEE Computer Society. 1
- [12] Zudi Lin, Prateek Garg, Atmadeep Banerjee, Salma Abdel Magid, Deqing Sun, Yulun Zhang, Luc Van Gool, Donglai Wei, and Hanspeter Pfister, “Revisiting rcnn: Improved training for image super-resolution,” *arXiv preprint arXiv:2201.11279*, 2022. 1, 4, 8, 11
- [13] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie-Line Alberi-Morel, “Low-complexity single-image super-resolution based on nonnegative neighbor embedding,” in *British Machine Vision Conference, BMVC 2012, Surrey, UK*. 2012, pp. 1–10, BMVA Press. 1, 4, 5
- [14] Jie Liu, Wenjie Zhang, Yuting Tang, Jie Tang, and Gangshan Wu, “Residual feature aggregation network for image super-resolution,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Seattle, USA*. 2020, pp. 2356–2365, Computer Vision Foundation / IEEE. 1, 2
- [15] Kui Jiang, Zhongyuan Wang, Peng Yi, and Junjun Jiang, “Hierarchical dense recursive network for image super-resolution,” *Pattern Recognit.*, vol. 107, pp. 107475, 2020. 1
- [16] Bin Xia, Yucheng Hang, Yapeng Tian, Wenming Yang, Qingmin Liao, and Jie Zhou, “Efficient non-local contrastive attention for image super-resolution,” *arXiv preprint arXiv:2201.03794*, 2022. 1, 2
- [17] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu, “Residual dense network for image super-resolution,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, USA*. 2018, pp. 2472–2481, Computer Vision Foundation / IEEE Computer Society. 1
- [18] Juncheng Li, Faming Fang, Kangfu Mei, and Guixu Zhang, “Multi-scale residual network for image super-resolution,” in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany*. 2018, vol. 11212, pp. 527–542, Springer. 1
- [19] Yiqun Mei, Yuchen Fan, and Yuqian Zhou, “Image super-resolution with non-local sparse attention,” in *IEEE Conference on Computer Vision and Pattern Recognition, virtual*. 2021, pp. 3517–3526, Computer Vision Foundation / IEEE. 1, 2, 8
- [20] Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang, “Efficient long-range attention network for image super-resolution,” *arXiv preprint arXiv:2203.06697*, 2022. 1, 11
- [21] Meng-Hao Guo, Cheng-Ze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shi-Min Hu, “Visual attention network,” *arXiv preprint arXiv:2202.09741*, 2022. 1, 3, 6
- [22] 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,” *arXiv preprint arXiv:2111.11418*, 2021. 2, 3
- [23] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun, “Simple baselines for image restoration,” *arXiv preprint arXiv:2204.04676*, 2022. 2, 4, 5
- [24] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao, “PVT v2: Improved baselines with pyramid vision transformer,” *Comput. Vis. Media*, vol. 8, no. 3, pp. 415–424, 2022. 2, 4, 5
- [25] Yan Wang, “Edge-enhanced feature distillation network for efficient super-resolution,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2022, pp. 777–785. 2

- [26] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang, "Second-order attention network for single image super-resolution," in *IEEE Conference on Computer Vision and Pattern Recognition, Long Beach, USA*. 2019, pp. 11065–11074, Computer Vision Foundation / IEEE. 2, 4, 8
- [27] 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 *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK*. 2020, vol. 12357, pp. 191–207, Springer. 2, 3, 4, 8
- [28] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee, "Deeply-recursive convolutional network for image super-resolution," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, USA*. 2016, pp. 1637–1645, IEEE Computer Society. 2
- [29] Ying Tai, Jian Yang, and Xiaoming Liu, "Image super-resolution via deep recursive residual network," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA*. 2017, pp. 2790–2798, IEEE Computer Society. 2
- [30] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang, "Deep laplacian pyramid networks for fast and accurate super-resolution," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, USA*. 2017, pp. 5835–5843, IEEE Computer Society. 2
- [31] Zheng Hui, Xinbo Gao, Yunchu Yang, and Xiumei Wang, "Lightweight image super-resolution with information multi-distillation network," in *Proceedings of the 27th ACM International Conference on Multimedia, Nice, France*. 2019, pp. 2024–2032, ACM. 2, 6
- [32] 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 *2021 IEEE/CVF International Conference on Computer Vision, Montreal, Canada*. 2021, pp. 9992–10002, IEEE. 2
- [33] Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier, "Language modeling with gated convolutional networks," in *Proceedings of the 34th International Conference on Machine Learning, Sydney, Australia*. 2017, vol. 70 of *Proceedings of Machine Learning Research*, pp. 933–941, PMLR. 4
- [34] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc V Le, "Transformer quality in linear time," *arXiv preprint arXiv:2202.10447*, 2022. 4
- [35] Salma Abdel Magid, Yulun Zhang, Donglai Wei, Wondong Jang, Zudi Lin, Yun Fu, and Hanspeter Pfister, "Dynamic high-pass filtering and multi-spectral attention for image super-resolution," in *2021 IEEE/CVF International Conference on Computer Vision, Montreal, Canada*. 2021, pp. 4268–4277, IEEE. 4, 8
- [36] Eirikur Agustsson and Radu Timofte, "NTIRE 2017 challenge on single image super-resolution: Dataset and study," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, Honolulu, USA*. 2017, pp. 1122–1131, IEEE Computer Society. 4
- [37] Roman Zeyde, Michael Elad, and Matan Protter, "On single image scale-up using sparse-representations," in *Curves and Surfaces - 7th International Conference, Avignon, France*. 2010, vol. 6920, pp. 711–730, Springer. 4
- [38] David R. Martin, Charless 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 Eighth International Conference On Computer Vision, Vancouver, Canada*. 2001, pp. 416–425, IEEE Computer Society. 4
- [39] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja, "Single image super-resolution from transformed self-exemplars," in *IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA*. 2015, pp. 5197–5206, IEEE Computer Society. 4, 5
- [40] 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.*, vol. 76, no. 20, pp. 21811–21838, 2017. 4
- [41] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004. 4
- [42] Diederik P. Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, San Diego, USA*, Yoshua Bengio and Yann LeCun, Eds., 2015. 5
- [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al., "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32, Vancouver, Canada*, 2019, pp. 8024–8035. 5

- [44] 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 *9th International Conference on Learning Representations, Virtual Event, Austria*, 2021. 5
- [45] 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 33, virtual*, 2020. 6
- [46] Xindong Zhang, Hui Zeng, and Lei Zhang, “Edge-oriented convolution block for real-time super resolution on mobile devices,” in *MM ’21: ACM Multimedia Conference, Virtual Event, China*. 2021, pp. 4034–4043, ACM. 6
- [47] Xiaotong Luo, Yuan Xie, Yulun Zhang, Yanyun Qu, Cuihua Li, and Yun Fu, “Latticenet: Towards lightweight image super-resolution with lattice block,” in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK*. 2020, vol. 12367, pp. 272–289, Springer. 6
- [48] Shangchen Zhou, Jiawei Zhang, Wangmeng Zuo, and Chen Change Loy, “Cross-scale internal graph neural network for image super-resolution,” in *Advances in Neural Information Processing Systems 33, virtual*, 2020. 8
- [49] Kartikeya Bhardwaj, Milos Milosavljevic, Liam O’Neil, Dibakar Gope, Ramon Matas Navarro, Alex Chalfin, Naveen Suda, Lingchuan Meng, and Danny Loh, “Collapsible linear blocks for super-efficient super resolution,” in *Proceedings of Machine Learning and Systems 2022, Santa Clara, USA*, 2022. 11
- [50] Long Sun, Jinshan Pan, and Jinhui Tang, “Shufflemixer: An efficient convnet for image super-resolution,” *arXiv preprint arXiv:2205.15175*, 2022. 11
- [51] Zhisheng Lu, Hong Liu, Juncheng Li, and Linlin Zhang, “Efficient transformer for single image super-resolution,” *arXiv preprint arXiv:2108.11084*, 2021. 11
- [52] Yulun Zhang, Donglai Wei, Can Qin, Huan Wang, Hanspeter Pfister, and Yun Fu, “Context reasoning attention network for image super-resolution,” in *2021 IEEE/CVF International Conference on Computer Vision, Montreal, Canada*. 2021, pp. 4258–4267, IEEE. 11

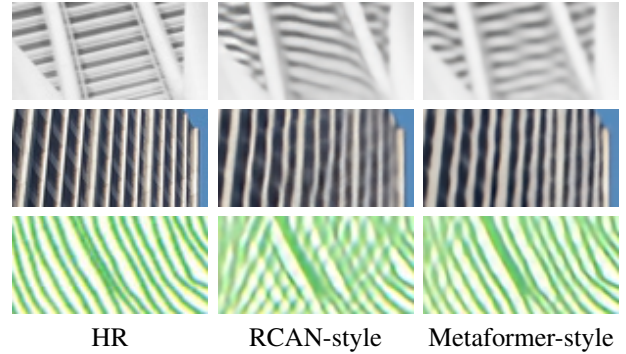
## 7. APPENDIX

### 7.1. More Quantitative Comparison

We show more results of comparison with SOTA methods on Set5, Set14, BSD100, Urban100, and Manga109 with  $\times 2$ ,  $\times 3$ , and  $\times 4$  upscaling factors. In Table 7, we compare our tiny version with other methods, including additional latest works, SESR [49] and ShuffleMixer [50]. For lightweight and classical SR, we also add ESRT [51], ELAN [20], CRAN [52], and RCAN-it [12] to perform supplementary comparison in Table 8 and Table 6, respectively.

### 7.2. More Qualitative Comparison

We also present more visual ablation results of the proposed designs. In Fig. 8, we exhibit the Urban100 and Manga109 results generated by MAN with two differently designed MABa. Specifically, Metaformer-style MAB obtains better restoration on the textures.



**Fig. 8.** Visual comparison between RCAN- and Metaformer-style MAB.

We further compare EDSR family with MAN family in Fig. 9 for further visual comparison. The model size and computation are correspondingly listed for reference. The results reflect that our MAN can compete EDSR with larger size.

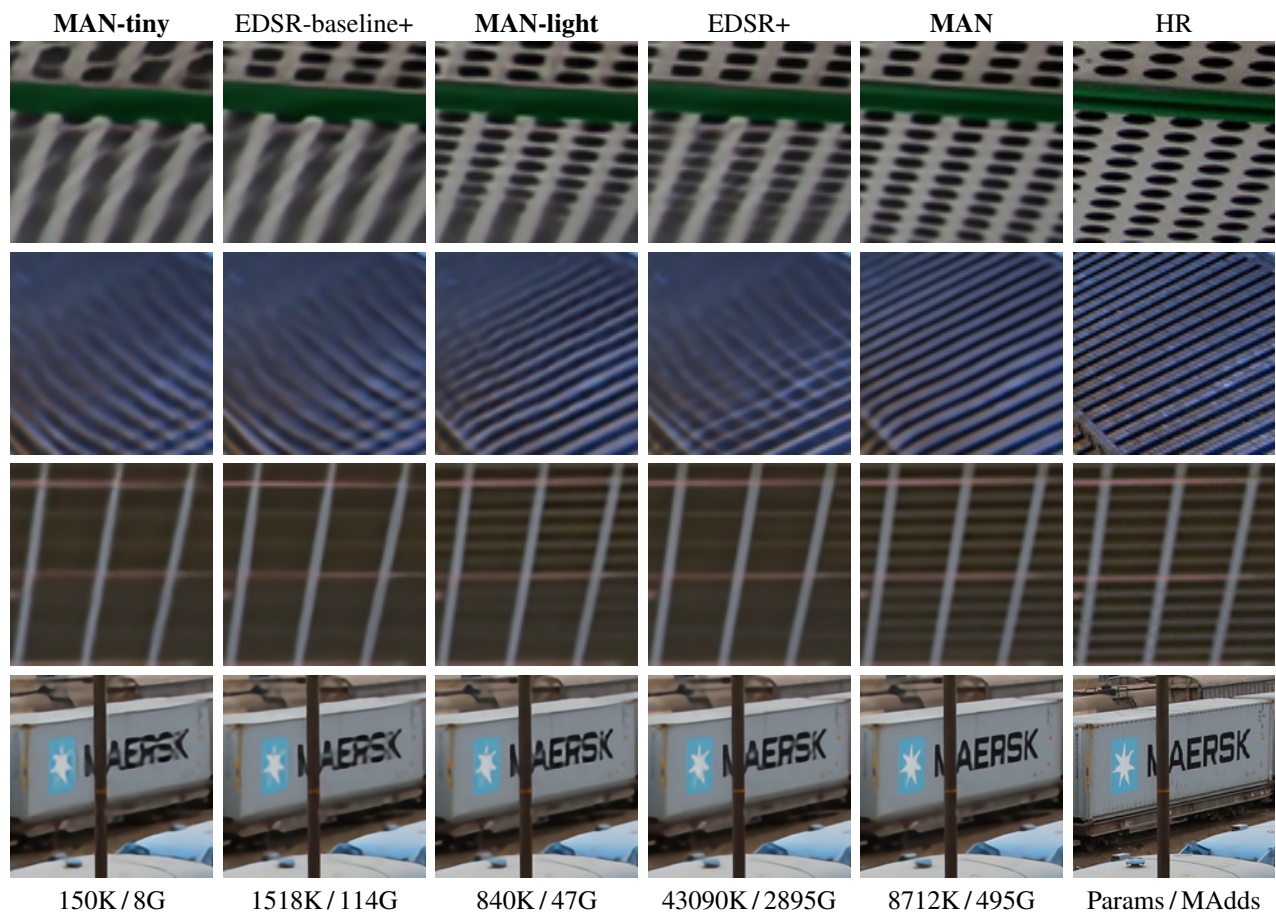
**Table 7.** Quantitative comparison (average PSNR/SSIM) with state-of-the-art approaches for tiny image SR. Best and second best results are **highlighted** and underlined, respectively.

Method	Scale	#Params	#Mult-Adds	Set5		Set14		BSD100		Urban100		Manga109	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
FSRCNN	×2	12K	6.0G	37.00	0.9558	32.63	0.9088	31.53	0.8920	29.88	0.9020	36.67	0.9710
LAPAR-C	×2	87K	35.0G	37.65	0.9593	33.20	0.9141	31.95	0.8969	31.10	0.9178	37.75	0.9752
SESR-XL	×2	105K	24.3G	37.77	0.9601	33.24	0.9145	31.99	0.8976	31.16	0.9184	38.01	0.9759
ECBSR-M10C32	×2	94K	21.8G	37.76	<b>0.9609</b>	33.26	0.9146	32.04	<u>0.8986</u>	<u>31.25</u>	0.9190	-	-
ShuffleMixer-tiny	×2	108K	25.0G	<u>37.85</u>	0.9600	<u>33.33</u>	<u>0.9153</u>	31.99	0.8972	31.22	0.9183	<u>38.25</u>	<u>0.9761</u>
<b>MAN-tiny (Ours)</b>	×2	134K	29.9G	<b>37.91</b>	<u>0.9603</u>	<b>33.47</b>	<b>0.9170</b>	<b>32.13</b>	<b>0.8993</b>	<b>31.75</b>	<b>0.9250</b>	<b>38.63</b>	<b>0.9771</b>
EDSR-baseline	×2	1370K	316G	37.99	0.9604	33.57	0.9175	32.16	0.8994	31.98	0.9272	38.54	0.9769
FSRCNN	×4	12K	4.6G	30.71	0.8657	27.59	0.7535	26.98	0.7150	24.62	0.7280	27.90	0.8517
LAPAR-C	×4	115K	25.0G	31.72	0.8884	28.31	0.7740	27.40	0.7292	25.49	0.7651	29.50	0.8951
SESR-XL	×4	115K	6.6G	31.54	0.8866	28.12	0.7712	27.31	0.7277	25.31	0.7604	29.04	0.8901
ECBSR-M10C32	×4	98K	5.7G	31.66	0.8911	28.15	0.7776	27.34	<u>0.7363</u>	25.41	0.7653	-	-
ShuffleMixer-tiny	×4	113K	8.0G	31.88	0.8912	28.46	0.7779	27.45	0.7313	25.66	0.7690	29.96	0.9006
<b>MAN-tiny (Ours)</b>	×4	150K	8.4G	<b>32.07</b>	<b>0.8930</b>	<b>28.53</b>	<b>0.7801</b>	<b>27.51</b>	<b>0.7345</b>	<b>25.84</b>	<b>0.7786</b>	<b>30.18</b>	<b>0.9047</b>
EDSR-baseline	×4	1518K	114G	32.09	0.8938	28.58	0.7813	27.57	0.7357	26.04	0.7849	30.35	0.9067

**Table 8.** Quantitative comparison (average PSNR/SSIM) with state-of-the-art approaches for lightweight image SR. Best and second best results are **highlighted** and underlined, respectively.

Method	Scale	#Params	#Mult-Adds	Set5		Set14		BSD100		Urban100		Manga109	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
IMDN	×2	694K	158.8G	38.00	0.9605	33.63	0.9177	32.19	0.8996	32.17	0.9283	38.88	0.9774
LatticeNet	×2	756K	169.5G	38.15	0.9610	33.78	0.9193	32.25	0.9005	32.43	0.9302	-	-
ESRT	×2	677K	-	38.03	0.9600	33.75	0.9184	32.25	0.9001	32.58	0.9318	<u>39.12</u>	0.9774
SwinIR-light	×2	878K	195.6G	38.14	<u>0.9611</u>	33.86	0.9206	<u>32.31</u>	<u>0.9012</u>	32.76	0.9340	<u>39.12</u>	<u>0.9783</u>
ELAN-light	×2	582K	168.4G	<u>38.17</u>	<u>0.9611</u>	<b>33.94</b>	0.9207	32.30	<u>0.9012</u>	<u>32.76</u>	0.9340	39.11	0.9782
<b>MAN-light (Ours)</b>	×2	820K	184.0G	<b>38.18</b>	<b>0.9612</b>	<u>33.93</u>	<b>0.9213</b>	<b>32.36</b>	<b>0.9022</b>	<b>32.92</b>	<b>0.9364</b>	<b>39.44</b>	<b>0.9786</b>
EDSR	×2	40730K	9387G	38.11	0.9602	33.92	0.9195	32.32	0.9013	32.93	0.9351	39.10	0.9773
IMDN	×3	703K	71.5G	34.36	0.9270	30.32	0.8417	29.09	0.8046	28.17	0.8519	<u>33.61</u>	<u>0.9445</u>
LatticeNet	×3	765K	76.3G	34.53	0.9281	30.39	0.8424	29.15	0.8059	28.33	0.8538	-	-
ESRT	×3	770K	-	34.42	0.9268	30.43	0.8433	29.15	0.8063	28.46	0.8574	33.95	0.9455
SwinIR-light	×3	886K	87.2G	<u>34.62</u>	<u>0.9289</u>	30.54	<u>0.8463</u>	29.20	<u>0.8082</u>	28.66	<u>0.8624</u>	33.98	<u>0.9478</u>
ELAN-light	×3	590K	75.7G	34.61	0.9288	30.55	0.8463	29.21	0.8081	28.69	<u>0.8624</u>	34.00	0.9478
<b>MAN-light (Ours)</b>	×3	829K	82.5G	<b>34.65</b>	<b>0.9292</b>	<b>30.60</b>	<b>0.8476</b>	<b>29.29</b>	<b>0.8101</b>	<b>28.87</b>	<b>0.8671</b>	<b>34.40</b>	<b>0.9493</b>
EDSR	×3	43680K	4470G	34.65	0.9280	30.52	0.8462	29.25	0.8093	28.80	0.8653	34.17	0.9476
IMDN	×4	715K	40.9G	32.21	0.8948	28.58	0.7811	27.56	0.7353	26.04	0.7838	30.45	0.9075
LatticeNet	×4	777K	43.6G	32.30	0.8962	28.68	0.7830	27.62	0.7367	26.25	0.7873	-	-
ESRT	×4	751K	-	32.19	0.8947	28.69	0.7833	27.69	0.7379	26.39	0.7962	30.75	0.9100
SwinIR-light	×4	897K	49.6G	<u>32.44</u>	<u>0.8976</u>	28.77	<u>0.7858</u>	<u>27.69</u>	<u>0.7406</u>	26.47	0.7980	<u>30.92</u>	<u>0.9151</u>
ELAN-light	×4	601K	43.2G	32.43	0.8975	<u>28.78</u>	<u>0.7858</u>	<u>27.69</u>	<u>0.7406</u>	<u>26.54</u>	<u>0.7982</u>	<u>30.92</u>	0.9150
<b>MAN-light (Ours)</b>	×4	840K	47.1G	<b>32.50</b>	<b>0.8988</b>	<b>28.87</b>	<b>0.7885</b>	<b>27.77</b>	<b>0.7429</b>	<b>26.70</b>	<b>0.8052</b>	<b>31.25</b>	<b>0.9170</b>
EDSR	×4	43090K	2895G	32.46	0.8968	28.80	0.7876	27.71	0.7420	26.64	0.8033	31.02	0.9148





**Fig. 9.** Visual comparison between EDSR series and MAN series on Urban100 and DIV2K ( $\times 4$ ). “+” indicates self-ensemble strategy. The Params and MAdds represent parameters and multiply-add operations counted under  $1280 \times 720$  outputs.