

HIPA: Hierarchical Patch Transformer for Single Image Super Resolution

Qing Cai, *Member, IEEE* Yiming Qian, Jinxing Li, Jun Lyu, Yee-Hong Yang, *Senior Member, IEEE*
Feng Wu, *Fellow, IEEE* and David Zhang, *Life Fellow, IEEE*

Abstract—Transformer-based architectures start to emerge in single image super resolution (SISR) and have achieved promising performance. However, most existing vision Transformer-based SISR methods still have two shortcomings: (1) they divide images into the same number of patches with a *fixed* size, which may not be optimal for restoring patches with different levels of texture richness; and (2) their position encodings treat all input tokens equally and hence, neglect the dependencies among them. This paper presents a HIPA, which stands for a novel Transformer architecture that progressively recovers the high resolution image using a hierarchical patch partition. Specifically, we build a cascaded model that processes an input image in multiple stages, where we start with tokens with small patch sizes and gradually merge them to form the full resolution. Such a hierarchical patch mechanism not only explicitly enables feature aggregation at multiple resolutions but also adaptively learns patch-aware features for different image regions, e.g., using a smaller patch for areas with fine details and a larger patch for textureless regions. Meanwhile, a new attention-based position encoding scheme for Transformer is proposed to let the network focus on which tokens should be paid more attention by assigning different weights to different tokens, which is the first time to our best knowledge. Furthermore, we also propose a multi-receptive field attention module to enlarge the convolution receptive field from different branches. The experimental results on several public datasets demonstrate the superior performance of the proposed HIPA over previous methods quantitatively and qualitatively. We will share our code and models when the paper is accepted.

Index Terms—Image restoration, single image super-resolution,

This work was supported in part by the National Science Foundation of China under Grant 62102338, Grant 61906162, and Grant 62172347; in part by the Natural Science Foundation of Shandong Province under Grant ZR2020QF031; in part by the Qingdao Postdoctoral Innovation Project under Grant QDBSH20230101001; in part by the CUHK(SZ)- Linklogis Joint Laboratory of Computer Vision and Artificial Intelligence; in part by the Shenzhen Institute of Artificial Intelligence and Robotics for Society; in part by the Shenzhen Research Institute of Big Data; and in part by the Natural Sciences and Engineering Research Council of Canada and the University of Alberta. (*Corresponding author: David Zhang, Jun Lyu*)

Qing Cai is with the Faculty of Information Science and Engineering, Ocean University of China, Qingdao, Shandong, 266100, China.

Yiming Qian is with the Department of Computer Science, University of Manitoba, Winnipeg, Manitoba, R3T 2N2, Canada.

Jinxing Li is with the School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen, Guangdong, 518055, China.

Jun Lyu is with the School of Nursing, The Hong Kong Polytechnic University, Hong Kong (e-mail:ljdream0710@pku.edu.cn).

Yee-Hong Yang is with the Department of Computing Science, University of Alberta, Edmonton, Alberta T6G 2E9, Canada.

Feng Wu is with the School of Information Science and Technology, University of Science and Technology of China, Hefei, Anhui, 230026, China.

David Zhang is with the School of Data Science, The Chinese University of Hong Kong, Shenzhen, Guangdong 518172, China, also with the Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen, Guangdong 518000, China, and also with the CUHK(SZ)- Linklogis Joint Laboratory of Computer Vision and Artificial Intelligence, Shenzhen, Guangdong 518172, China (e-mail: davidzhang@cuhk.edu.cn).



Fig. 1. Visual comparisons of $\times 4$ SISR on “img_063” from Urban 100. It can be seen that our method obtains better visual quality and recovers more textures and details compared with that of other state-of-the-art methods. The colors red and blue represent the best and the second best methods.

hierarchical patch Transformer, attention-based position embedding

I. INTRODUCTION

Single Image Super-Resolution (SISR), aiming to recover a high-resolution (HR) image from its corresponding degraded low-resolution (LR) version, plays an important and fundamental role in computer vision and image processing due to its wide range of real-world applications, such as medical imaging [1], surveillance [2] and remote sensing [3], amongst others. SISR is a very challenging and ill-posed problem because there is no unique solution for any given LR input [4, 5].

Deep convolutional neural networks (CNNs) have achieved remarkable success in SISR and various architectures have been presented so far, for example, residual learning [6–8], dense connections [9, 10], UNet-like architectures with skip connections [11, 12], dilated convolutions [13, 14], generative models [15–18] and other kinds of CNNs [19–21]. However, the convolution in CNN uses a sliding window to extract local features and hence, is weak in capturing long-range or non-local dependencies, which are important for SISR. In particular, for some regions with fine textures, faithful reconstruction depends not only on local relationships but also on long-range dependencies [22, 23]. To alleviate this issue, many attention mechanisms have been proposed and introduced into SISR, such as global attention mechanism [24–28] and non-local attention mechanism [29–32]. As shown in

Fig. 1, although some state-of-the-art (SOTA) methods such as NLSN [31] could recover some amounts of high-frequency details, the reconstructed slanted line structures exhibit fuzzy and blurry boundaries, which are faithfully restored using our hierarchical patch Transformer.

Inspired by the significant success of Transformer in natural language processing [33] for its advantages in modeling long-range context, vision Transformer is also introduced into the field of SISR [34–38] and has obtained superior results than many SOTA CNN-based methods due to the multi-head self-attention mechanism that is capable of modeling long-distance dependencies [39]. Very recently, hybrid architectures combining CNN and Transformer start to emerge in the community [37] to fully utilize the advantage of CNN in extracting local features and the advantage of Transformer in establishing long-range dependencies. Although existing Transformer-based SISR models have achieved superior results, the recovered results still exhibit blurry boundaries as shown in the result of SwinIR [37] in Fig. 1. **The main reasons may lie in two shortcomings of existing vision Transformer-based SISR methods.** First, almost all of them partition all input images into the same number of fixed-size patches, which may not be ideal considering different images on image regions have their own characteristics [40]. Second, the position encoding of most vision Transformer-based SISR methods treats all input tokens equally. However, the low-resolution input tokens contain abundant information for SISR, which are treated equally across tokens and hence, the representation ability of Transformer is limited.

In order to compensate the above two shortcomings, in this paper, we propose a **Hierarchical Patch (HIPA) Transformer** by partitioning an input image into a hierarchy of patches with different sizes. In particular, a multi-stage architecture is first developed by alternately stacking CNN and Transformer to boost their benefits in feature extraction. Then, to achieve different size patch input for the Transformer and to let the Transformer establish global dependencies from different numbers of tokens, the LR image is first partitioned into a hierarchy of subblocks, which are used as inputs to the Transformer by starting from the small-size blocks and gradually merging them in the next stage. In addition, we design a novel attention-based position encoding scheme for the Transformer based on dilated channel attention to model the position information with a continuous dynamical model. Besides, a multi-receptive field attention module is proposed based on dilated convolution with different dilation factors to enlarge the convolution receptive field from different branches. As shown in Fig. 1, our HIPA obtains better visual quality compared with that of other state-of-the-art SISR methods.

Briefly, the contributions of this paper mainly include:

- A novel hierarchical patch Transformer has been designed to achieve multi-size patches for Transformers. This approach is more effective than treating all samples with the same number of fixed-size patches because the hierarchical patch Transformer allows patches with different texture richness to adopt different sizes, rather than a single size patch;

- A new attention-based position encoding scheme is proposed for Transformer that allows the network to focus on which tokens should be paid more attention, which is the first time to our best knowledge;
- A multi-receptive field dilated attention module is designed to enlarge the convolution receptive field from different branches, which achieves relatively smaller increase of the computational complexity compared to the one by increasing the depth and the filter size of a CNN to enlarge the receptive field.

The rest of the paper is organized as follows: Section II briefly overviews related works. Section III presents the proposed HIPA Transformer and discusses its advantages and differences with existing methods. Section IV presents the experimental results and analysis of the proposed method by comparing it with state-of-the-art models. Finally, the paper concludes in Section V.

II. RELATED WORK

CNN-based Models: The SRCNN model proposed by Dong *et al.* [41] is a pioneering work to apply CNN to single image super-resolution, which has achieved superior performance than traditional methods [42–44] by using only a three-layer CNN to represent the mapping between LF and HR images. Based on the SRCNN, many deeper and wider CNN based SISR models have been proposed to achieve better restoration performance. However, blindly increasing the depth of a network does not necessarily improve the performance but may introduce many new issues for training, for example, the vanishing or exploding gradient [45]. Later, residual learning is introduced into SISR to ease the training difficulty of deeper networks. For example, by introducing residual learning into a deeper network, Kim *et al.* can stack more convolutional layers and propose the VDSR [46]. However, all of the above models need to first pre-process the LR input to obtain the desired image size using interpolation, which is not only time consuming but also often introduces noise and blurriness in the input image. To address the above issues, Dong *et al.* [47] introduce a deconvolution layer as the last layer and achieve end-to-end training for SISR. Such a deconvolution layer is then substituted by a more efficient sub-pixel convolution layer [48] proposed by Shi *et al.*, which is also adopted by our method similar to the EDSR [45] and the RCAN [24]. However, all of these models treat the LR features equally across channels, which inevitably limits the restoration capability of CNNs. Even worse, the convolution kernel usually has a limited receptive field and cannot sufficiently extract long-range or non-local features. As a result, for some regions with fine details, these methods yield poor performance.

Attention-based Models: To address the above issues, attention mechanism [24, 26, 28, 30] is introduced into SISR to guide the deep neural network to selectively pay more attention on features where there is more information. For example, by integrating channel attention and residual blocks, Zhang *et al.* propose the RCAN [24], which markedly improves the representational performance of the CNN. Dai *et al.* propose the SAN [28] using a novel trainable second-order channel at-

tention. However, the channel attention treats different convolution layers independently and neglects the correlation among them. To alleviate this issue, Niu *et al.* propose the HAN [26] by integrating a layer attention module and a channel-spatial attention module into the residual blocks. More recently, non-local attention modules [29–32] are proposed to address the inherent issue of CNNs in establishing long range or non-local dependencies among exacted features. For example, Zhang *et al.* [29], propose the RNAN by mixing a local masked branch and a non-local attention mechanism, which are, respectively, in charge of concentrating on extracting more local structures and considering more long-range dependencies in the extracted features. Mei *et al.* [30] propose the CSNLN by integrating a Cross-Scale Non-Local prior with local and in-scale non-local priors using a recurrent neural network, which can efficiently explore the existing cross-scale feature similarities in images. Xia *et al.* [32] propose an efficient non-local attention module by using the kernel function of approximation and the associative law of matrix multiplication, which successfully achieves comparable performance compared to that of the previous non-local attention module while requires only linear computation and space complexity with respect to the LR size. However, these models are still incapable of adequately and comprehensively compensate for the shortcomings of CNNs in establishing long-range dependencies.

Transformer-based Models: Inspired by the significant performance of the vision Transformer [33, 49, 50], it has also been applied to the SISR field [34–38]. For example, Chen *et al.* propose the image processing Transformer (IPT) [35] model for various image restoration tasks based on a pre-trained standard Transformer [33]. Recently, to capture local relationships, researchers begin to introduce convolutions to Transformers by integrating the vision Transformer module with convolution [37, 51–54]. For example, Liang *et al.* [37] propose the Swin Transformer-based image resolution model (SwinIR) by combining CNN and Transformer and achieves superior performance while maintaining computational efficiency. Huang *et al.* propose DGSM-Swin [52] by introducing a learned Gaussian Scale Mixture (GSM) prior into the Swin Transformer. In addition to classic performance-oriented SISR, hybrid architectures have also emerged in the field of light-weight SISR [53–55]. For example, Lu *et al.* propose a novel Efficient Super-Resolution Transformer (ESRT) [53], which integrates a lightweight CNN backbone and a lightweight Transformer backbone to achieve a small GPU memory footprint using an efficient multi-head attention. Fang *et al.* proposed a Hybrid Network of CNN and Transformer (HNCT) [54] for lightweight image SISR, which can exploit both local and non-local priors by integrating CNN and Transformer. Although Transformer-based SISR have achieved impressive results, most existing Transformers divide images into the same number of fixed-size patches, which may not be ideal for restoring patches with different levels of texture richness. Besides, the position encodings used in most existing Transformer are predefined and treat the positional information of different tokens equally.

III. METHODOLOGY

A. Issues and Motivations

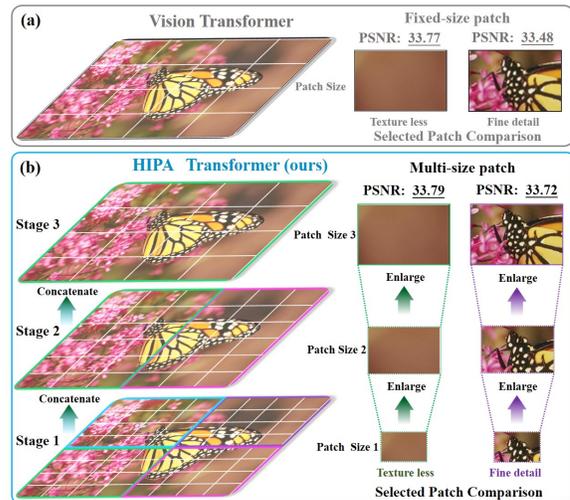


Fig. 2. Illustration of comparison between (a) the existing vision Transformer with fixed-size patch and (b) our HIPA Transformer with multi-size patch on regions with different texture richness.

As discussed in the contribution part of Section I and the Transformer-based models part of Section II, using fixed-size patches with varying level of texture richness is suboptimal and can limit the recovery performance of many existing Transformer-based SISR models. To further explain this, we provide an example shown in Fig. 2 that demonstrates the impact of patch size on the “monarch” image from the Set14 dataset. Fig. 2(a) shows the fixed-size split of vision Transformer-based methods (left column), in which the input image is split into the same number of fixed-size tokens, a selected textureless background region (middle column) and a butterfly head region with fine detail (right column). Fig. 2(b) shows the multi-size split of our HIPA Transformer (left column), in which the input image is partitioned into multiple stages, where tokens with small patch sizes are used first and gradually merged with larger patches to form the full resolution, two selected regions (last two columns) the same as that in Fig. 2(a). From the visual and quantitative comparison of the selected background region (middle column) between using the existing fixed-size patch and our multi-size patch, it can be observed that their visual quality and PSNR values are very similar, which suggests that using a large patch size for textureless region is enough for the network to finish the final restoration. However, from the recovery performance comparison of the selected butterfly head region (right column) between using the existing fixed-size patch and using our multi-size patch, it can be found that their visual quality and PSNR values have a certain gap, which suggests that using a large patch size for region with fine detail is not optimal. In contrast, in this case, using a smaller patch size is more helpful to recover fine details, which is demonstrated by the improved PSNR value using our multi-size patch. From the above discussion, we can summarize that: (1) using fixed-size patches with different texture richness in the whole restoration process is inappropriate, which is the reason that many existing

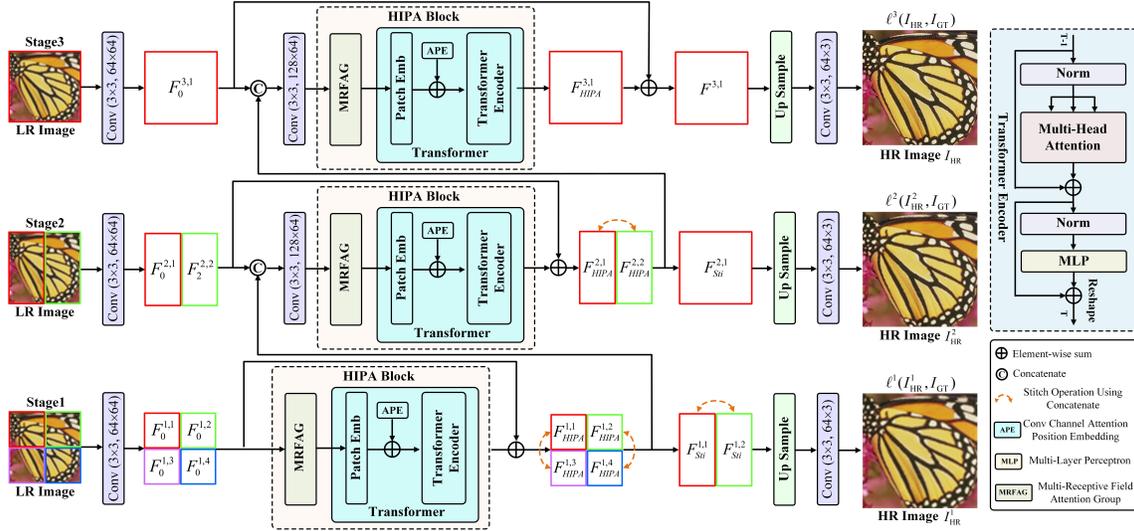


Fig. 3. Overall framework of our HIPA method for progressive SISR. Each stage in the proposed model is constructed based on the proposed HIPA block, which consists of two main modules. The first module is a Transformer designed to learn global dependencies between contexts (See Section III-B for details). The second module is a multi-receptive field attention group used to exhaustively mine local features contained in the original LR image (See Section III-D for details). The Transformers of the first two stages mainly learn the broad contextual information, while the last stage focuses more on learning the desired details. Finally, the training loss is defined based on the summation over all outputs $\ell^1(\cdot)$, $\ell^2(\cdot)$, $\ell^3(\cdot)$ of different stages to optimize our HIPA.

Transformer-based models still exhibit blurry boundaries; (2) using multi-size patches with different richness is helpful to improve the restoration results, which motivates the proposed HIPA Transformer.

B. Hierarchical Patch Transformer

As shown in Fig. 3, the proposed HIPA consists of three stages to progressively recover the high-resolution (HR) image from its low-resolution (LR) input. The Transformer of the first two stages mainly learn broad contextual information, while the last stage focuses on learning the desired details. Each stage is constructed based on the proposed HIPA block, which mainly consists of two modules: a multi-receptive field attention group and a designed Transformer. To achieve multi-size patch input for the HIPA, we adopt the hierarchical patch partition on the input LR image. Specifically, we first split the LR image into different non-overlapping patches for different stages: four for the first stage, two for the second stage, and the entire LR image for the last stage, and then, gradually integrate intermediate results in the next stage.

For simplicity, in the notation that follows, I_{LR} and I_{HR} denote the original LR input and the final HR output of the HIPA, respectively. $I_{LR}^{i,j}$ denotes the j -th patch at Stage i . For example, $I_{LR}^{1,1}$ denotes the 1-st patch at Stage 1, i.e., the upper left corner patch of Stage 1 input shown in Fig. 3.

Following [24, 45], we also use one convolution layer to extract the shallow feature (SF) $F_0^{1,j}$ from the original LR image. For Stage 1:

$$F_0^{1,j} = H_{SF}(I_{LR}^{1,j}) \quad (j = 1, 2, 3, 4), \quad (1)$$

where H_{SF} denotes the convolution operation. Then, the extracted shallow feature is input to the proposed HIPA block to further extract deep features:

$$F_{HIPA}^{1,j} = H_{HIPA}(F_0^{1,j}) \quad (j = 1, 2, 3, 4), \quad (2)$$

where H_{HIPA} denotes the proposed HIPA block. After stitching $F_{HIPA}^{1,1}$ with $F_{HIPA}^{1,3}$ and stitching $F_{HIPA}^{1,2}$ with $F_{HIPA}^{1,4}$ using concatenate operation, dubbed vertical stitching, we obtain the output features of Stage 1, which are then concatenated with the shallow features of Stage 2 as shown in Fig. 3:

$$\begin{aligned} F_{Sti}^{1,1} &= H_{Sti}(F_{HIPA}^{1,1}, F_{HIPA}^{1,3}) + H_{Sti}(F_0^{1,1}, F_0^{1,3}), \\ F_{Sti}^{1,2} &= H_{Sti}(F_{HIPA}^{1,2}, F_{HIPA}^{1,4}) + H_{Sti}(F_0^{1,2}, F_0^{1,4}), \end{aligned} \quad (3)$$

where H_{Sti} denotes the stitch using the concatenate operation. We utilize vertical stitching for sub-patches rather than horizontal stitching, i.e., stitching $F_{HIPA}^{1,1}$ with $F_{HIPA}^{1,3}$ and stitching $F_{HIPA}^{1,2}$ with $F_{HIPA}^{1,4}$ using the concatenate operation. Although we also investigated horizontal stitching, it did not yield significant differences. Finally, the recovered HR image of Stage 1: I_{HR}^1 , is obtained by further stitching $F_{Sti}^{1,1}$ and $F_{Sti}^{1,2}$ using concatenate operation, and then successively input the stitched result into an upscale module and a reconstruction module (i.e., one convolution layer) as follows:

$$I_{HR}^1 = H_{Rec}(H_{UP}(H_{Sti}(F_{Sti}^{1,1}, F_{Sti}^{1,2}))), \quad (4)$$

where H_{UP} and H_{Rec} denote the upscale and reconstruction module, respectively.

For Stage 2 and Stage 3, the extracted shallow features $F_0^{2,j}$ ($j = 1, 2$) and $F_0^{3,1}$ need to be first concatenated with the output features of the upper stage, which is then input into the next operation similar to Stage 1. Finally, the recovered HR images of Stage 2: I_{HR}^2 and Stage 3: I_{HR} can be obtained. As shown in Fig. 3, the predictions of the three stages are gradually improved. For example, the prediction of Stage 2 is the refinement of Stage 1. With the multi-stage refinement, image regions with high spatial frequency are gradually recovered.

Finally, the proposed HIPA is trained using a training loss, which is the sum over all the outputs of I_{HR}^1 (Stage 1), I_{HR}^2

(Stage 2) and I_{HR} (Stage 3):

$$L(\Theta) = \ell^1(I_{HR}^1, I_{GT}) + \ell^2(I_{HR}^2, I_{GT-T}) + \ell^3(I_{HR}, I_{GT}), \quad (5)$$

where Θ denotes the parameter set of the proposed network. $\ell^1(\cdot)$, $\ell^2(\cdot)$ and $\ell^3(\cdot)$, respectively, stand for the loss of Stage 1, Stage 2 and Stage 3. This work also uses the L_1 loss following previous work for the sake of fairness. I_{GT} denotes the ground-truth HR image.

C. Attention-based Position Encoding

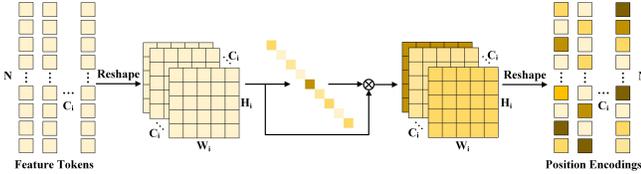


Fig. 4. Illustration of attention-based position encoding (APE). $N = \frac{H_i \times W_i}{P_i^2}$ denotes the number of tokens. W , H and C , respectively, denote the width, height and channel number of the feature map.

To incorporate the order of the token sequence, position encodings are usually adopted in Transformers [56]. However, the original position embedding of ViT is pre-defined and independent of input tokens. When an input LR image with a new size is input, the number of patches will be different from that before and the learned position embedding will be mismatched with the new size. So, the input image with a new size has to be first interpolated to the desired size, which not only reduces the overall performance of the ViT but also seriously limits its application. To address the above issue, Chu *et al.* [56] propose a conditional position encoding (CPE) by introducing a 2-D convolution to embed position encoding, which can easily generalize to an input LR image with a new input size. However, the CPE treats all input tokens equally and may neglect the dependencies among them. To address this, we propose a new attention-based position encoding (APE) method by introducing attention into position embedding to let the Transformer focus on important tokens. Specifically, the patch embedding module first reshapes the extracted feature $F_{MRFAG}^{i,j} \in R^{H_i \times W_i \times C_i}$ into a number of flattened 2D patches $x_p \in R^{\frac{H_i \times W_i}{P_i^2} \times P_i^2 \times C_i}$ by partitioning the input into non-overlapping $P_i \times P_i$ patches where (H_i, W_i) , C_i , $\frac{H_i \times W_i}{P_i^2}$ and P_i , respectively, denote the resolution of Stage i input, the number of channel, the number of patches and the patch size. Then, as shown in Fig. 4, the flattened feature tokens are reshaped to the 2D image space. In the 2D image space, a convolution and a channel attention are applied to produce the final position encoding. With the help of attention, the final position encoding can let the network focus on the important tokens.

D. Multi-Receptive Field Attention Group

We now show our MRFAG, which mainly consists of G multi-receptive field attention modules (MRFAMs) as shown in Fig. 5. Each MRFAM consists of three dilated convolution based channel attention connected in parallel, a fusion module

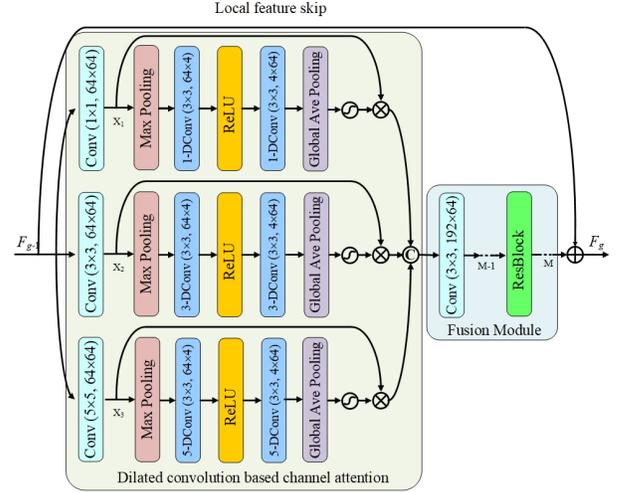


Fig. 5. The architecture of the multi-receptive field attention module (MRFAM). Note the numbers ($kernel \times kernel$, $input \times output$) in Conv and Dconv denote kernel size and input and output feature map number. i -DCConv denotes dilated convolution with dilation factor i .

and a local feature skip (LFS). It has been verified that although increasing the depth and the filter size of the CNN can, respectively, enlarge the receptive field and extract more information contained in the low-quality image, it not only introduces more parameters but also increases the computational complexity [14]. Thus, we propose the dilated convolution based channel attention to enlarge the receptive field of the networks, which is the most significant difference between ours and the Squeeze-and-Excitation network (SE) [57].

Specifically, for each dilated convolution based channel attention shown in Fig. 5, denote $X_i = [x_{i,1}, \dots, x_{i,c}, \dots, x_{i,C}]$ to be the input, which contains C 2D feature map $x_{i,c} \in R^{H \times W}$, where H and W , respectively, are the height and width of the feature map. Firstly, by shrinking the extracted features using max pooling, the output feature $Z_i = [z_{i,1}, \dots, z_{i,c}, \dots, z_{i,C}]$ of each branch can be obtained, where $z_{i,c} \in R^{\frac{H}{Stride} \times \frac{W}{Stride} \times C_i}$ denotes the output feature. Then, two dilated convolution layers and an activation function are applied to fully exploit feature dependencies from the aggregated information. Finally, the sigmoid function is adopted as the activation function:

$$s_{i,c} = f(H_{GPL}(W_U \delta(W_D z_{i,c}))), \quad (6)$$

where $f(\cdot)$, $H_{GPL}(\cdot)$ and $\delta(\cdot)$, respectively, stand for the sigmoid function, the global average pooling and the ReLU function. W_D is the weight set of the first dilated convolution layer in the channel attention shown in Fig. 5, which plays the role of downscaling with a reduction ratio γ (we set $\gamma = 16$). After the ReLU function, the low-dimension feature is then upsampled with ratio γ by the second dilated convolution layer. W_U denotes its weight set. The channel statistics s can be obtained to rescale the input $x_{i,c}$:

$$\hat{x}_{i,c} = s_{i,c} \cdot x_{i,c}, \quad (7)$$

where $s_{i,c}$ and $x_{i,c}$ denote, respectively, the scaling factor and feature maps of the c -th channel.

Besides, we introduce the LFS connection to ensure stability in training the network and to bypass redundant features in the low-quality image. The final output of MRFAG is obtained as

$$F_{MRFAG} = F_0^{i,j} + \omega_{LFS}(F_{MRFAM_G}), \quad (8)$$

where ω_{LFS} denotes the weight of the convolution layer at the tail of MRFAG. F_{MRFAG} and F_{MRFAM_G} , respectively, denote the output of MRFAG and the G -th output of MRFAM.

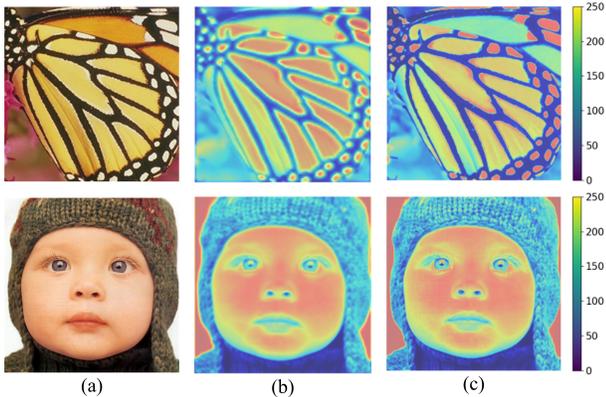


Fig. 6. Class activate map (CAM) comparison. (a) Input images. (b) CAM before using MRFAG. (c) CAM after using MRFAG.

Fig. 6 shows a comparison of the class active map (CAM) before and after using the MRFAG. It can be found that the CAM after using MRFAG becomes sharper than before, validating that the MRFAG has a better capability of recovering high-frequency signals. This also enables the network to focus more on recovering textures and details.

E. Discussions

Below, we discuss the significant differences between our HIPA Transformer and two most closely related Transformers: IPT [35] and SwinIR [37].

Differences to the IPT model: Based on a pre-trained standard Transformer [33], Chen *et al.* propose the image processing Transformer (IPT) [35] model for image restoration tasks, which achieves superior performance than most CNN-based SISR methods. Although it is a Transformer-based SISR method similar to our method, there are three significant differences between the IPT and our HIPA Transformer: (i) the IPT model uses a pre-trained Transformer, which means that when we apply it for image super-resolution, we need to first pre-train the Transformer using large labeled datasets and then fine-tune the whole network. As a result, its performance is limited by the lack of sufficient labeled samples for training. In contrast, our HIPA Transformer is an end-to-end network, which successfully avoids the tedious pre-training and fine-tuning process; (ii) IPT uses fixed-size patches for all input tokens with different richness, which is not optimal and limits its performance as discussed in Section III-A, while our HIPA Transformer uses multi-size patches for tokens with different richness, e.g., using a smaller patch for areas with fine details and a large patch for textureless regions; (iii) The IPT uses Transformer to extract features and to construct long-range dependencies, while our HIPA Transformer is a hybrid

architecture combining CNN and Transformer, which can fully utilize the advantage of CNN in local feature extraction and the advantage of Transformer in establishing long-range dependencies. More comparisons of experimental results are shown in Section IV-B.

Differences to the SwinIR model: Liang *et al.* propose the SwinIR [37] by combining CNN with the Swin Transformer [64] into one network, which is also a hybrid architecture similar to our HIPA Transformer. Our key distinctions from it are summarized as follows: (i) The SwinIR uses a plain concatenation of CNN and Transformer with a fixed patch size, while we use a multi-stage model that divides the input into different blocks and aggregates them from small to large by alternating CNN and Transformer, which not only explicitly enables feature aggregation at multiple resolution but also adaptively learns patch-aware features for different image regions; (ii) The local-resolution input tokens contain abundant information for SISR, however, the SwinIR treats all the input tokens equally and hence, limits its representation ability. In contrast, we design a novel attention-based encoding method to focus on the important tokens and to improve its performance for regions with fine details; (iii) The CNN used in the SwinIR model detects local image features using the same scale, treats all LR image features equally, and neglects the dependencies among them. In contrast, our HIPA proposes a multi-receptive field attention module, which lets the proposed network know where to pay more attention and to sufficiently extract local features of LR images. The experimental results, which demonstrate the advantages of our method, are shown in Section IV-B.

IV. EXPERIMENTS

A. Settings

Datasets: Following previous works [24, 26, 28, 30], we also choose DIV2K [65] as our training dataset, which contains 800 training images and 100 validation images. For testing, we select the standard public datasets: Set5 [66], Set14 [67], B100 [68], Urban100 [69], and Manga109 [70] as our test datasets. All degraded datasets are obtained by the bicubic interpolation model.

Evaluation Metrics: To quantitatively compare the recovered HR results of the proposed model with that of the state-of-the-art models, PSNR and SSIM are used, which are calculated based on the luminance channel of the YCbCr space of the recovered RGB results.

Training Settings: We set the number of MRFAMs as $G = 5, 5, 20$ in the MRFAG structure for Stage 1, Stage 2 and Stage 3, respectively. In each MRFAM, we set the number of residual blocks as $M = 5$. All the convolution layers have $C = 64$ filters except for those in the dilated convolution layer as shown in Fig. 5, where the convolution layer has $C = 4$ filters. We use 3×3 as the filter size for all convolution layers except for those in the dilated convolution based channel attention where the kernel sizes are 1×1 , 3×3 and 5×5 , which are shown in Fig. 5. Following previous works [24, 28, 45], we adopt the sub-pixel convolution [48] to upsample the LR features to HR. During training, we also augment the training

TABLE I
 QUANTITATIVE COMPARISONS WITH STATE-OF-THE-ART PERFORMANCE-ORIENTED SISR METHODS ON FIVE BENCHMARK DATASETS FOR SCALE FACTOR $\times 2$, $\times 3$ AND $\times 4$. THE BEST RESULTS ARE HIGHLIGHTED IN RED AND THE SECOND BEST IN BLUE.

Methods	Scale	Year	Set5		Set14		B100		Urban100		Manga109	
			PSNR	SSIM								
SRMD [58]	$\times 2$	2018	37.79	0.9601	33.32	0.9159	32.05	0.8985	31.33	0.9204	38.07	0.9761
DBPN [59]	$\times 2$	2018	38.09	0.9600	33.85	0.9190	32.27	0.9000	32.55	0.9324	38.89	0.9775
RDN [9]	$\times 2$	2018	38.24	0.9614	34.01	0.9212	32.34	0.9017	32.89	0.9353	39.18	0.9780
MSRN [60]	$\times 2$	2018	38.08	0.9605	33.74	0.9170	32.23	0.9013	32.22	0.9326	38.82	0.9768
RCAN [24]	$\times 2$	2018	38.27	0.9614	34.12	0.9216	32.41	0.9027	33.34	0.9384	39.44	0.9786
SRFBN [61]	$\times 2$	2019	38.11	0.9609	33.82	0.9196	32.29	0.9010	32.62	0.9328	39.08	0.9779
SAN [28]	$\times 2$	2019	38.31	0.9620	34.07	0.9213	32.42	0.9028	33.10	0.9370	39.32	0.9792
CSNLTN [30]	$\times 2$	2020	38.28	0.9616	34.12	0.9223	32.40	0.9024	33.25	0.9386	39.37	0.9785
HAN [26]	$\times 2$	2020	38.27	0.9614	34.16	0.9217	32.41	0.9027	33.35	0.9385	39.46	0.9785
NSR [62]	$\times 2$	2020	38.23	0.9614	33.94	0.9203	32.34	0.9020	33.02	0.9367	39.31	0.9782
IGNN [63]	$\times 2$	2020	38.24	0.9613	34.07	0.9217	32.41	0.9025	33.23	0.9386	39.35	0.9786
RFANet [27]	$\times 2$	2020	38.26	0.9615	34.16	0.9220	32.41	0.9026	33.33	0.9389	39.44	0.9783
NLSN [31]	$\times 2$	2021	38.34	0.9618	34.08	0.9231	32.43	0.9027	33.42	0.9394	39.59	0.9789
SwinIR [37]	$\times 2$	2021	38.35	0.9620	34.14	0.9227	32.44	0.9030	33.40	0.9393	39.60	0.9792
TDPN [5]	$\times 2$	2022	38.31	0.9621	34.16	0.9225	32.52	0.9045	33.36	0.9386	39.57	0.9795
ELAN [51]	$\times 2$	2022	38.36	0.9620	34.20	0.9228	32.45	0.9030	33.44	0.9391	39.62	0.9793
DGSM-Swin [52]	$\times 2$	2023	38.24	0.9615	33.93	0.9217	32.36	0.9019	32.95	0.9442	39.31	0.9783
HIPA(ours)	$\times 2$	2023	38.38	0.9621	34.25	0.9235	32.48	0.9033	33.50	0.9400	39.75	0.9794
HIPA+(ours)	$\times 2$	2023	38.41	0.9623	34.30	0.9238	32.51	0.9036	33.57	0.9409	39.81	0.9795
SRMD [58]	$\times 3$	2018	34.12	0.9254	30.04	0.8382	28.97	0.8025	27.57	0.8398	33.00	0.9403
RDN [9]	$\times 3$	2018	34.71	0.9296	30.57	0.8468	29.26	0.8093	28.80	0.8653	34.13	0.9484
MSRN [60]	$\times 3$	2018	34.38	0.9262	30.34	0.8395	29.08	0.8041	28.08	0.8554	33.44	0.9427
RCAN [24]	$\times 3$	2018	34.74	0.9299	30.65	0.8482	29.32	0.8111	29.09	0.8702	34.44	0.9499
SRFBN [61]	$\times 3$	2019	34.70	0.9292	30.51	0.8461	29.24	0.8084	28.73	0.8641	34.18	0.9481
SAN [28]	$\times 3$	2019	34.75	0.9300	30.59	0.8476	29.33	0.8112	28.93	0.8671	34.30	0.9494
CSNLTN [30]	$\times 3$	2020	34.74	0.9300	30.66	0.8482	29.33	0.8105	29.13	0.8712	34.45	0.9502
HAN [26]	$\times 3$	2020	34.75	0.9299	30.67	0.8483	29.32	0.8110	29.10	0.8705	34.48	0.9500
NSR [62]	$\times 3$	2020	34.62	0.9289	30.57	0.8475	29.26	0.8100	28.83	0.8663	34.27	0.9484
IGNN [63]	$\times 3$	2020	34.72	0.9298	30.66	0.8484	29.31	0.8105	29.03	0.8696	34.39	0.9496
RFANet [27]	$\times 3$	2020	34.79	0.9300	30.67	0.8487	29.34	0.8115	29.15	0.8720	34.59	0.9506
NLSN [31]	$\times 3$	2021	34.85	0.9306	30.70	0.8485	29.34	0.8117	29.25	0.8726	34.57	0.9508
SwinIR [37]	$\times 3$	2021	34.89	0.9312	30.77	0.8503	29.37	0.8124	29.29	0.8744	34.74	0.9518
TDPN [5]	$\times 3$	2022	34.86	0.9312	30.79	0.8501	29.45	0.8126	29.26	0.8724	34.48	0.9508
ELAN [51]	$\times 3$	2022	34.90	0.9313	30.80	0.8504	29.38	0.8124	29.32	0.8745	34.73	0.9517
DGSM-Swin [52]	$\times 3$	2023	34.77	0.9300	30.65	0.8490	29.29	0.8109	28.93	0.8684	34.30	0.9498
HIPA(ours)	$\times 3$	2023	34.95	0.9318	30.84	0.8515	29.45	0.8140	29.41	0.8760	34.88	0.9521
HIPA+(ours)	$\times 3$	2023	35.01	0.9320	30.90	0.8530	29.49	0.8151	29.50	0.8784	34.96	0.9528
SRMD [58]	$\times 4$	2018	31.96	0.8925	28.35	0.7787	27.49	0.7337	25.68	0.7731	30.09	0.9024
DBPN [59]	$\times 4$	2018	32.47	0.8980	28.82	0.7860	27.72	0.7400	26.38	0.7946	30.91	0.9137
RDN [9]	$\times 4$	2018	32.47	0.8990	28.81	0.7871	27.72	0.7419	26.61	0.8028	31.00	0.9151
MSRN [60]	$\times 4$	2018	32.07	0.8903	28.60	0.7751	27.52	0.7273	26.04	0.7896	30.17	0.9034
RCAN [24]	$\times 4$	2018	32.63	0.9002	28.87	0.7889	27.77	0.7436	26.83	0.8087	31.22	0.9173
SRFBN [61]	$\times 4$	2019	32.47	0.8983	28.81	0.7868	27.72	0.7409	26.60	0.8015	31.15	0.9160
SAN [28]	$\times 4$	2019	32.64	0.9003	28.92	0.7888	27.78	0.7436	26.79	0.8068	31.18	0.9169
CSNLTN [30]	$\times 4$	2020	32.68	0.9004	28.95	0.7888	27.80	0.7439	27.22	0.8168	31.43	0.9201
HAN [26]	$\times 4$	2020	32.64	0.9002	28.90	0.7890	27.80	0.7442	26.85	0.8094	31.42	0.9177
NSR [62]	$\times 4$	2020	32.55	0.8987	28.79	0.7876	27.72	0.7414	26.61	0.8025	31.10	0.9145
IGNN [63]	$\times 4$	2020	32.57	0.8998	28.85	0.7891	27.77	0.7434	26.84	0.8090	31.28	0.9182
RFANet [27]	$\times 4$	2020	32.66	0.9004	28.88	0.7894	27.79	0.7442	26.92	0.8112	31.41	0.9187
NLSN [31]	$\times 4$	2021	32.59	0.9000	28.87	0.7891	27.78	0.7444	26.96	0.8109	31.27	0.9184
SwinIR [37]	$\times 4$	2021	32.72	0.9021	28.94	0.7914	27.83	0.7459	27.07	0.8164	31.67	0.9226
TDPN [5]	$\times 4$	2022	32.69	0.9005	29.01	0.7943	27.93	0.7460	27.24	0.8171	31.58	0.9218
ELAN [51]	$\times 4$	2022	32.75	0.9022	28.96	0.7914	27.83	0.7459	27.13	0.8167	31.68	0.9226
DGSM-Swin [52]	$\times 4$	2023	32.61	0.9005	28.91	0.7903	27.78	0.7445	26.73	0.8068	31.25	0.9193
HIPA(ours)	$\times 4$	2023	32.78	0.9025	29.07	0.7935	27.90	0.7479	27.27	0.8191	31.83	0.9365
HIPA+(ours)	$\times 4$	2023	32.84	0.9034	29.14	0.7955	27.98	0.7490	27.31	0.8214	31.91	0.9438

dataset by randomly rotating by 90° , 180° , 270° and flipping horizontally [24, 28, 45]. In each training batch, LR images with patch size 48×48 are cropped as inputs. The proposed model is trained by the ADAM optimizer with a fixed initial learning rate of 10^{-4} . The whole process is implemented in the PyTorch platform with 4 Nvidia TITAN TRX GPUs, each with 24GB of memory.

B. Comparisons with State-of-the-arts

In this section, we compare our HIPA with 17 state-of-the-art SISR methods: SRMD [58] DBPN [59], RDN [9], MSRN [60], RCAN [24], SRFBN [61], SAN [28], CSNLTN [30], HAN [26], NSR [62], IGNN [63], RFANet [27], NLSN [31], SwinIR [37], TDPN [5], ELAN [51] and DGSM-Swin [52]. Following previous works [24, 26, 28, 37], we also perform self-ensemble on our HIPA to further improve its performance and dub it HIPA+.

Quantitative Comparison: Table I reports the quantitative comparisons between our method and 17 state-of-the-art SISR methods on five benchmark datasets for scale factor $2\times$, $3\times$ and $4\times$. The best results are highlighted in red and the second best in blue. All the reported methods are proposed in recent 5 years and have achieved competitive results. Compared with these methods, our HIPA+ achieves the best results on multiple benchmarks for all scaling factors and surpasses most state-of-the-art methods in terms of PSNR and SSIM. Without using self-ensemble our network HIPA still achieves the best results on multiple benchmarks for all scale factors. It is noteworthy that our proposed HIPA is superior to SwinIR [37] and DGSM-Swin [52], both of which are all hybrid architecture similar to HIPA. Specifically, the values of PSNR on the Urban100 dataset for scale factor $\times 4$ are improved by **0.2 dB** and **0.54 dB**, respectively, compared to SwinIR and DGSM-Swin. The main reasons may lie in that i) the designed multi-stage

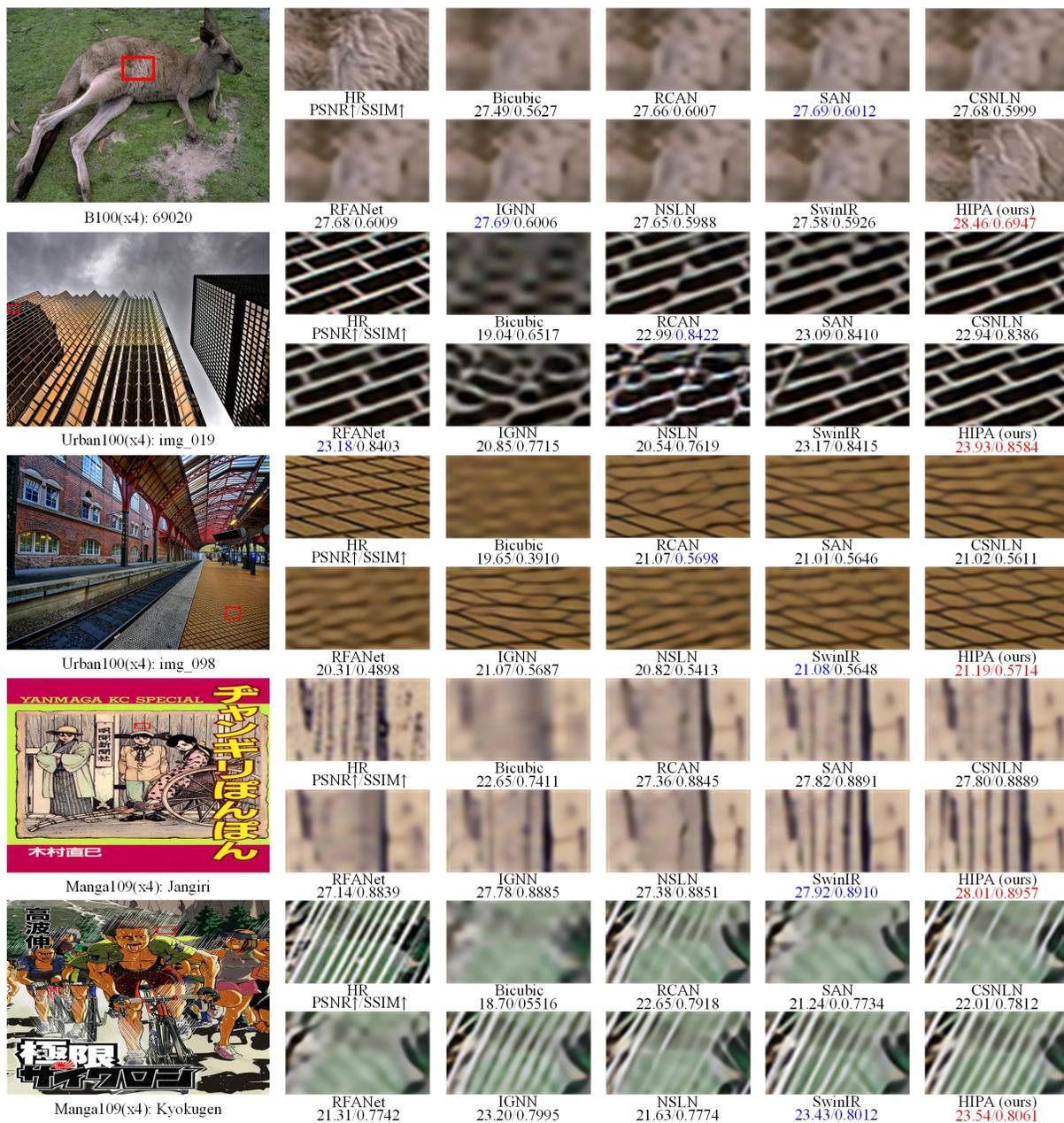


Fig. 7. Visual comparisons with state-of-the-art SISR methods for 4× SR on the B100, the Urban100 and the Manga109 datasets. Best viewed on screen.

progressive model not only can exploit features from different size patches but also can gradually recover the HR image from coarse to fine; and ii) the proposed MRFAG can let the network exhaustively mine local features contained in the original LR image from different receptive fields based on dilated convolution with different dilation factors.

Qualitative Comparison: In Fig. 7, we also visually illustrate the zoomed in comparison results with some state-of-the-art methods on several images from the test datasets. From the results, we find that our proposed HIPA can always obtain sharper results and recover more high-frequency textures and details, while most competing SISR models suffer from some unpleasant blurring artifacts. Take “img_109” in Manga109 shown in Fig. 7 as an example, existing methods obtain

heavy blurring artifacts. The early proposed Bicubic fails to generate the clear structures. Although more recent methods, e.g. RCAN [24], SAN [28], CSNLN [30], RFANet [27], IGNN [63], NSLN [31] and SwinIR [37] can recover the main outlines, they fail to recover textures and details, and even generate some distorted and deformed textures. In contrast, our method effectively recovers textures through using the proposed HIPA and MRFAG.

Further Comparison: Table II compares the number of parameters, computational complexity and average running time comparisons for various SISR methods. Except for ESRT [53], HNCT [54] and Swin2SR-s [55], which are light weight SISR, the other methods (including our HIPA) are classic performance-oriented SISR. EDSR [45], RDN [9] and

TABLE II
TOTAL NUMBER OF PARAMETERS, COMPUTATIONAL COMPLEXITY,
RUNNING TIME AND PSNR COMPARISON ON URBAN100 DATASET FOR
SCALE FACTOR $\times 4$ OF DIFFERENT MODELS.

Model	Params(M)	FLOPs(G)	Time(s)	PSNR(dB)
EDSR [45]	43	2875	0.5437	26.64
RDN [9]	22.3	1305	0.3127	26.61
RCAN [24]	16	912	0.3891	26.83
IPT [35]	114	1480	1.3550	27.26
SwinIR [37]	11.8	978	0.4331	27.07
ESRT [53]	0.68	65.2	0.0106	26.39
HNCT [54]	0.37	78.8	0.0154	26.20
Swin2SR-s [55]	1	146.5	0.0223	26.58
HIPa(ours)	11.3	764	0.2615	27.27

RCAN [24] are CNN-based SISr methods, while IPT [35], SwinIR [37], ESRT [53], HNCT [54] and Swin2SR-s [55] are state-of-the-art Transformer-based SISrs.

Compared to EDSR [45], RDN [9] and RCAN [24], HIPA not only has fewer parameters but also achieves a better PSNR value. Compared to two similar performance-oriented methods: IPT [35] and SwinIR [37], our method is more efficient in both computational time and memory usage. Although classic performance-oriented SISrs have more parameters and high computational complexity than light weight SISrs, they have better PSNR values because they focus more on performance.

C. Ablation Study

TABLE III
ABLATION STUDY OF THE DESIGNED MULTI-SIZE PATCH INPUT FOR THE
PROPOSED METHOD. ALL THE EXPERIMENTS ARE CONDUCTED WITH THE
SAME EXPERIMENTAL CONDITIONS EXCEPT THAT FOR THE PATCH
EMBEDDING USED IN HIPA TRANSFORMER, AND TESTED ON THE SET14
AND URBAN100 DATASETS FOR SCALE FACTOR $\times 4$.

Design	Set14			Urban 100		
	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$
Fixed-size patch	34.21	30.76	28.95	33.43	29.30	27.18
Our multi-size patch	34.25	30.84	29.07	33.50	29.41	27.27

Ablation Study of HIPA Transformer: In Table III, we report the quantitative comparisons between the proposed HIPA Transformer with fixed-size patches and with multi-size patches by letting the network with and without partitioning the input LR image into a hierarchy of subblocks for scale factor $\times 2$, $\times 3$ and $\times 4$ on the Set14 and Urban100 datasets. From the PSNR results, we find that the HIPA Transformer using patches of different sizes outperforms that using fixed-size patches by a maximum of 0.12dB. The main reason is that the hierarchy of subblocks let the network learn one LR image from different sizes and improves the overall performance of the final results. Our result not only validates the effectiveness of the proposed multi-size patch but also further validates the effectiveness of the proposed hierarchical multi-stage structure.

Besides, in Table IV, we show the effects of the HIPA Transformer size on model performance. It can be found that the PSNR is positively correlated with the HIPA Transformer size. Even though the performance keeps increasing, the total number of parameters of the proposed HIPA Transformer grows also. To balance the performance and model size, we

TABLE IV
IMPACT OF HIPA TRANSFORMER SIZE FOR THE PROPOSED METHOD.
HIPA_S, HIPA_M AND HIPA_L DENOTE SMALL, MEDIUM AND LARGE
VERSION OF HIPA TRANSFORMER, RESPECTIVELY. PATS, HEADN AND
LAYERN, RESPECTIVELY, DENOTE THE PATCH SIZE, THE HEAD NUMBER
AND THE LAYER NUMBER OF HIPA TRANSFORMER. ALL THE
EXPERIMENTS ARE CONDUCTED WITH THE SAME EXPERIMENTAL
CONDITIONS, AND TESTED ON THE SET14 DATASET FOR SCALE FACTOR
 $\times 4$.

HIPa Index	PatS	HeadN	LayerN	Params	PSNR
HIPa_S	4	4	4	8.9 M	28.99
HIPa_M	8	8	8	11.3 M	29.07
HIPa_L	16	16	16	16.1 M	29.12

TABLE V
ABLATION STUDY ON THE HIPA TRANSFORMER USING ‘PE’, ‘CPE’ AND
THE PROPOSED ‘APE’ FOR SCALE $\times 2$, $\times 3$ AND $\times 4$ ON THE MANGA109
DATASET.

Different PE	$\times 2$	$\times 3$	$\times 4$
PE [71]	39.69	34.80	31.74
CPE [56]	34.71	34.83	31.77
APE (Ours)	39.75	34.88	31.83

choose HIPA_M (PatS = 8, HeadNr = 8 and LayerN = 8) in the rest of the experiments.

Ablation Study of the Proposed APE: To validate the effectiveness of the proposed attention position encoding (APE), a comparison experiment between the proposed method using the previous position embedding (PE) [71], the condition position encoding (CPE) [56] and the proposed APE is conducted for scale $\times 2$, $\times 3$ and $\times 4$ on the Set14 and Urban100 datasets. From the PSNR results shown in Table V, we find that the HIPA Transformer using the proposed APE obtains superior performance than that using the previous PE and CPE for all scales on the two datasets, which validates the effectiveness of the proposed APE.

TABLE VI
ABLATION STUDY OF THE MRFAG ON THE B100, URBAN100 AND
MANGA109 DATASETS FOR SCALE FACTOR $\times 3$.

Module	B100	Urban100	Manga109
w/o MRFAG	29.41	29.35	34.80
w/ RCAB [24]	29.43	29.37	34.83
w/ MRFAG	29.45	29.41	34.88

Ablation Study of the Proposed MRFAG: To validate the effectiveness of the MRFAG, as shown in Table VI, we conduct a comparison experiment between the proposed method without using the proposed MRFAG, with using a classic and effective feature extraction module RCAB [24] and with the proposed MRFAG. Note that, for a fair comparison, the total number of parameters of our method with RCAB modules is 11.6M, which is close to the number of parameters for our method with MRFAG, which is 11.3M. It is found that the improvement of the proposed method with using the proposed MRFAG is greater than that of the proposed method using the RCAB, which validates the effectiveness of the proposed MRFAG.

In addition, as shown in Table VII, another comparison between the proposed MRFAG using SE attention and using the proposed dilated convolution based attention to validate the effectiveness of the proposed dilated convolution based

TABLE VII
PSNR COMPARISON BETWEEN THE PROPOSED METHOD USING SE ATTENTION AND USING THE PROPOSED CONVOLUTION BASED ATTENTION ON THE B100, URBAN100 AND MANGA109 DATASETS FOR SCALE FACTOR $\times 4$.

Module	B100	Urban100	Manga109
Standard SE attention	27.87	27.22	31.78
Dilated convolution attention	27.90	27.27	31.83

attention. It can be found that the proposed dilated convolution based attention can improve PSNR value by a mean 0.043 dB on the B100, Urban100 and Manga109 datasets for scale factor $\times 4$ than that of the standard SE attention.

V. CONCLUSION

In this paper, we propose the Hierarchical Patch Transformer (HIPA) for accurate single image super resolution, which progressively recovers the high resolution image by partitioning the input into a hierarchy of patches. Specifically, a multi-stage progressive model is employed where the earlier stages use smaller patches as tokens and the final stage operates at full resolution. Our architecture is a cascade CNNs and Transformers for feature aggregation across multiple stages. In addition, we develop a novel attention-based position encoding scheme that allows the Transformer focus on the important tokens and easily process an input low resolution images with varying sizes. Besides, the proposed multi-receptive field attention module can enlarge the convolution receptive field from different branches. The quantitative and qualitative evaluations on different benchmark datasets demonstrate the effectiveness of the hierarchical patch partition over using fixed-size patches, as well as the superior performance of the proposed HIPA over most state-of-the-art methods in PSNR, SSIM and visual quality.

REFERENCES

- [1] Y. Li, Y. Iwamoto, L. Lin, R. Xu, R. Tong, and Y.-W. Chen, "Volumenet: a lightweight parallel network for super-resolution of mr and ct volumetric data," *IEEE Transactions on Image Processing*, vol. 30, pp. 4840–4854, 2021.
- [2] W. Wen, W. Ren, Y. Shi, Y. Nie, J. Zhang, and X. Cao, "Video super-resolution via a spatio-temporal alignment network," *IEEE Transactions on Image Processing*, vol. 31, pp. 1761–1773, 2022.
- [3] M. R. Arefin, V. Michalski, P.-L. St-Charles, A. Kalaitzis, S. Kim, S. E. Kahou, and Y. Bengio, "Multi-image super-resolution for remote sensing using deep recurrent networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [4] K. Zhang, L. V. Gool, and R. Timofte, "Deep unfolding network for image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3217–3226.
- [5] Q. Cai, J. Li, H. Li, Y.-H. Yang, F. Wu, and D. Zhang, "Tdpn: Texture and detail-preserving network for single image super-resolution," *IEEE Transactions on Image Processing*, vol. 31, pp. 2375–2389, 2022.
- [6] S. Anwar and N. Barnes, "Real image denoising with feature attention," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 3155–3164.
- [7] S. Nie, C. Ma, D. Chen, S. Yin, H. Wang, L. Jiao, and F. Liu, "A dual residual network with channel attention for image

- restoration," in *Proceedings of the European Conference on Computer Vision*, 2020, pp. 352–363.
- [8] Q. Wang, Q. Gao, L. Wu, G. Sun, and L. Jiao, "Adversarial multi-path residual network for image super-resolution," *IEEE Transactions on Image Processing*, vol. 30, pp. 6648–6658, 2021.
- [9] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2472–2481.
- [10] D. Song, C. Xu, X. Jia, Y. Chen, C. Xu, and Y. Wang, "Efficient residual dense block search for image super-resolution," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 007–12 014.
- [11] X. Hu, M. A. Naei, A. Wong, M. Lamm, and P. Fieguth, "Runet: A robust unet architecture for image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [12] K. Prajapati, V. Chudasama, H. Patel, A. Sarvaiya, K. P. Upla, K. Raja, R. Ramachandra, and C. Busch, "Channel split convolutional neural network (chasnet) for thermal image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4368–4377.
- [13] W. Yang, R. T. Tan, J. Feng, J. Liu, Z. Guo, and S. Yan, "Deep joint rain detection and removal from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1357–1366.
- [14] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep cnn denoiser prior for image restoration," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3929–3938.
- [15] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4681–4690.
- [16] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *Proceedings of the European Conference on Computer Vision Workshops*, 2018, pp. 0–0.
- [17] K. Prajapati, V. Chudasama, H. Patel, K. Upla, K. Raja, R. Ramachandra, and C. Busch, "Direct unsupervised super-resolution using generative adversarial network (dus-gan) for real-world data," *IEEE Transactions on Image Processing*, vol. 30, pp. 8251–8264, 2021.
- [18] Z. Wenlong, L. Yihao, C. Dong, and Y. Qiao, "Ranksrgan: Generative adversarial networks with ranker for image super-resolution," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [19] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5407–5416.
- [20] L. Lu, W. Li, X. Tao, J. Lu, and J. Jia, "Masa-sr: Matching acceleration and spatial adaptation for reference-based image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6368–6377.
- [21] Y. Guo, X. Wu, and X. Shu, "Data acquisition and preparation for dual-reference deep learning of image super-resolution," *IEEE Transactions on Image Processing*, 2022.
- [22] Y. Zhang, D. Wei, C. Qin, H. Wang, H. Pfister, and Y. Fu, "Context reasoning attention network for image super-resolution," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 4278–4287.
- [23] G. Li, J. Lv, Y. Tian, Q. Dou, C. Wang, C. Xu, and J. Qin,

- “Transformer-empowered multi-scale contextual matching and aggregation for multi-contrast mri super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 20 636–20 645.
- [24] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, “Image super-resolution using very deep residual channel attention networks,” in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 286–301.
- [25] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, “Cycleisp: Real image restoration via improved data synthesis,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2696–2705.
- [26] 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 *Proceedings of the European Conference on Computer Vision*, 2020, pp. 191–207.
- [27] J. Liu, W. Zhang, Y. Tang, J. Tang, and G. Wu, “Residual feature aggregation network for image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2359–2368.
- [28] T. Dai, J. Cai, Y. Zhang, S.-T. Xia, and L. Zhang, “Second-order attention network for single image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 065–11 074.
- [29] Y. Zhang, K. Li, K. Li, B. Zhong, and Y. Fu, “Residual non-local attention networks for image restoration,” in *International Conference on Learning Representations*, 2019.
- [30] Y. Mei, Y. Fan, Y. Zhou, L. Huang, T. S. Huang, and H. Shi, “Image super-resolution with cross-scale non-local attention and exhaustive self-exemplars mining,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5690–5699.
- [31] Y. Mei, Y. Fan, and Y. Zhou, “Image super-resolution with non-local sparse attention,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3517–3526.
- [32] B. Xia, Y. Hang, Y. Tian, W. Yang, Q. Liao, and J. Zhou, “Efficient non-local contrastive attention for image super-resolution,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017.
- [34] F. Yang, H. Yang, J. Fu, H. Lu, and B. Guo, “Learning texture transformer network for image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5791–5800.
- [35] H. Chen, Y. Wang, T. Guo, C. Xu, Y. Deng, Z. Liu, S. Ma, C. Xu, C. Xu, and W. Gao, “Pre-trained image processing transformer,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 299–12 310.
- [36] J. Cao, Y. Li, K. Zhang, and L. Van Gool, “Video super-resolution transformer,” *arXiv preprint arXiv:2106.06847*, 2021.
- [37] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, “Swinir: Image restoration using swin transformer,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2021.
- [38] Z. Wang, X. Cun, J. Bao, W. Zhou, J. Liu, and H. Li, “Uformer: A general u-shaped transformer for image restoration,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17 683–17 693.
- [39] X. Chu, Z. Tian, Y. Wang, B. Zhang, H. Ren, X. Wei, H. Xia, and C. Shen, “Twins: Revisiting the design of spatial attention in vision transformers,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 9355–9366, 2021.
- [40] Y. Wang, R. Huang, S. Song, Z. Huang, and G. Huang, “Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 11 960–11 973, 2021.
- [41] C. Dong, C. C. Loy, K. He, and X. Tang, “Learning a deep convolutional network for image super-resolution,” in *Proceedings of the European Conference on Computer Vision*, 2014, pp. 184–199.
- [42] L. Zhang and X. Wu, “An edge-guided image interpolation algorithm via directional filtering and data fusion,” *IEEE Transactions on Image Processing*, vol. 15, no. 8, pp. 2226–2238, 2006.
- [43] K. Zhang, X. Gao, D. Tao, and X. Li, “Single image super-resolution with non-local means and steering kernel regression,” *IEEE Transactions on Image Processing*, vol. 21, no. 11, pp. 4544–4556, 2012.
- [44] C. L. P. Chen, L. Liu, L. Chen, Y. Y. Tang, and Y. Zhou, “Weighted couple sparse representation with classified regularization for impulse noise removal,” *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 4014–4026, 2015.
- [45] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, “Enhanced deep residual networks for single image super-resolution,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 136–144.
- [46] J. Kim, J. K. Lee, and K. M. Lee, “Accurate image super-resolution using very deep convolutional networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2016, pp. 1646–1654.
- [47] C. Dong, C. C. Loy, and X. Tang, “Accelerating the super-resolution convolutional neural network,” in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 391–407.
- [48] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1874–1883.
- [49] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [50] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, “Glue: A multi-task benchmark and analysis platform for natural language understanding,” in *International Conference on Learning Representations*, 2019.
- [51] X. Zhang, H. Zeng, S. Guo, and L. Zhang, “Efficient long-range attention network for image super-resolution,” in *Proceedings of the European Conference on Computer Vision*, 2022, pp. 649–667.
- [52] T. Huang, X. Yuan, W. Dong, J. Wu, and G. Shi, “Deep gaussian scale mixture prior for image reconstruction,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [53] Z. Lu, J. Li, H. Liu, C. Huang, L. Zhang, and T. Zeng, “Transformer for single image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2022, pp. 457–466.
- [54] J. Fang, H. Lin, X. Chen, and K. Zeng, “A hybrid network of cnn and transformer for lightweight image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2022, pp. 1103–1112.
- [55] M. V. Conde, U.-J. Choi, M. Burchi, and R. Timofte, “Swin2sr: Swin2 transformer for compressed image super-resolution and restoration,” in *Proceedings of the European Conference on Computer Vision Workshops*. Springer, 2023, pp. 669–687.
- [56] X. Chu, Z. Tian, B. Zhang, X. Wang, X. Wei, H. Xia, and C. Shen, “Conditional positional encodings for vision transformers,” *arXiv preprint arXiv:2102.10882*, 2021.
- [57] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.
- [58] K. Zhang, W. Zuo, and L. Zhang, “Learning a single convolutional super-resolution network for multiple degradations,” in

Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 3262–3271.

- [59] M. Haris, G. Shakhnarovich, and N. Ukita, “Deep back-projection networks for super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1664–1673.
- [60] J. Li, F. Fang, K. Mei, and G. Zhang, “Multi-scale residual network for image super-resolution,” in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 517–532.
- [61] Z. Li, J. Yang, Z. Liu, X. Yang, G. Jeon, and W. Wu, “Feedback network for image super-resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3867–3876.
- [62] Y. Fan, J. Yu, Y. Mei, Y. Zhang, Y. Fu, D. Liu, and T. S. Huang, “Neural sparse representation for image restoration,” in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 15 394–15 404.
- [63] S. Zhou, J. Zhang, W. Zuo, and C. C. Loy, “Cross-scale internal graph neural network for image super-resolution,” in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 3499–3509.
- [64] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 10 012–10 022.
- [65] E. Agustsson and R. Timofte, “Ntire 2017 challenge on single image super-resolution: Dataset and study,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 126–135.
- [66] M. Bevilacqua, A. Roumy, C. Guillemot, and M.-L. A. Morel, “Low-complexity single-image super-resolution based on non-negative neighbor embedding,” in *British Machine Vision Conference*, 2012.
- [67] R. Zeyde, M. Elad, and M. Protter, “On single image scale-up using sparse-representations,” in *International Conference on Curves and Surfaces*, 2010, pp. 711–730.
- [68] 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 *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2, 2001, pp. 416–423.
- [69] J. Huang, A. Singh, and N. Ahuja, “Single image super-resolution from transformed self-exemplars,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5197–5206.
- [70] 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. 21 811–21 838, 2017.
- [71] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2021.



Yiming Qian received the B.Sc. degree from University of Science and Technology of China, Hefei, China, in 2012, the M.Sc. degree from Memorial University of Newfoundland, St. John’s, Canada, in 2014, and the Ph.D. degree from the Department of Computing Science, University of Alberta, Edmonton, Canada, in 2019. He is currently an Assistant Professor with the University of Manitoba, Canada. His research interests include computer vision and computer graphics, while he recently focuses on 3D modeling.



Jinxing Li received the B.Sc. degree from the Department of Automation, Hangzhou Dianzi University, Hangzhou, China, in 2012, the M.Sc. degree from the Department of Automation, Chongqing University, Chongqing, China, in 2015, and the Ph.D. degree from the Department of Computing, Hong Kong Polytechnic University, Hong Kong, in 2018. He is currently an Associate Professor with the Harbin Institute of Technology at Shenzhen. His research interests are pattern recognition, deep learning, medical biometrics, and machine learning.



Jun Lyu received the B.Sc. degree in intelligence science and technology from Xidian University, in 2013, and the Ph.D. degree in biomechanics and medical engineering from Peking University, in 2018. From October 2017 to April 2018, she was a Visiting Ph.D. Student in the David Geffen School of Medicine at University of California, Los Angeles. She is currently a Postdoctoral Fellow with the School of Nursing, The Hong Kong Polytechnic University. Her research interests include deep learning, and medical image processing and analysis.



Yee-Hong Yang (Senior Member, IEEE) received the B.Sc. (first honors) from the University of Hong Kong, the M.Sc. from Simon Fraser University, and the Ph.D. from the University of Pittsburgh. He was a faculty member in the Department of Computer Science at the University of Saskatchewan from 1983 to 2001 and served as Graduate Chair from 1999 to 2001. While there, in addition to department level committees, he also served on many college and university level committees. Since July 2001, he has been a Professor in the Department of Computing Science at the University of Alberta. He served as Associate Chair (Graduate Studies) in the same department from 2003 to 2005. His research interests cover a wide range of topics from computer graphics to computer vision, which include physically based animation of Newtonian and non-Newtonian fluids, texture analysis and synthesis, human body motion analysis and synthesis, computational photography, stereo and multiple view computer vision, and underwater imaging. He has published over 100 papers in international journals and conference proceedings in the areas of computer vision and graphics. He is a Senior Member of the IEEE and serves on the Editorial Board of the journal *Pattern Recognition*. In addition to serving as a reviewer to numerous international journals, conferences, and funding agencies, he has served on the program committees of many national and international conferences. In 2007, he was invited to serve on the expert review panel to evaluate computer science research in Finland.



Qing Cai received the M.Sc. and Ph.D. degree from the Department of Automation, Northwestern Polytechnical University, Xi’an, China, in 2016 and 2019, respectively. From 2017 to 2018, he was a Visiting Ph.D. Student in the Department of Computing Science at University of Alberta. From 2020 to 2022, he was a Postdoctoral Fellow with The Chinese University of Hong Kong at Shenzhen, and also with the University of Science and Technology of China. He is currently an Associate Professor with the Faculty of Information Science and Engineering,

at Ocean University of China. His research interests include machine learning, deep learning, and computer vision, with a focus on image restoration, image segmentation, medical image processing and visual tracking.



Feng Wu (Fellow, IEEE) received the B.Sc. degree in electrical engineering from Xidian University, Xi’an, China, in 1992, and the M.Sc. and Ph.D. degrees in computer science from the Harbin Institute of Technology, Harbin, China, in 1996 and 1999, respectively. He was a Principle Researcher and a Research Manager with Microsoft Research Asia, Beijing, China. He is currently a Professor with the University of Science and Technology of China, Hefei, China, where he is also the Dean of the School of Information Science and Technology.

His research interests include image and video compression, media communication, and media analysis and synthesis.



David Zhang (Life Fellow, IEEE) graduated in Computer Science from Peking University. He received his MSc in 1982 and his PhD in 1985 in both Computer Science from the Harbin Institute of Technology (HIT), respectively. From 1986 to 1988 he was a Postdoctoral Fellow at Tsinghua University and then an Associate Professor at the Academia Sinica, Beijing. In 1994 he received his second PhD in Electrical and Computer Engineering from the University of Waterloo, Ontario, Canada. He has been a Chair Professor at the Hong Kong

Polytechnic University where he is the Founding Director of Biometrics Research Centre (UGC/CRC) supported by the Hong Kong SAR Government since 1998. Currently he is Presidential Chair Professor in Chinese University of Hong Kong (Shenzhen). So far, he has published over 20 monographs, 500+ international journal papers and 40+ patents from USA/Japan/HK/China. He has been continuously 8 years listed as a Global Highly Cited Researchers in Engineering by Clarivate Analytics during 2014-2021. He is also ranked about 85 with H-Index 123 at Top 1000 Scientists for international Computer Science and Electronics. Professor Zhang is a Croucher Senior Research Fellow, Distinguished Speaker of the IEEE Computer Society, and Fellows of both Royal Society of Canada and Canadian Academy of Engineering, as well as IEEE Life Fellow and IAPR/AAIA Fellow.