

HAT: Hybrid Attention Transformer for Image Restoration

Xiangyu Chen, Xiantao Wang, Wenlong Zhang, Xiangtao Kong,
 Yu Qiao, *Senior Member, IEEE*, Jiantao Zhou, *Senior Member, IEEE*, Chao Dong

Abstract—Transformer-based methods have shown impressive performance in image restoration tasks, such as image super-resolution and denoising. However, we find that these networks can only utilize a limited spatial range of input information through attribution analysis. This implies that the potential of Transformer is still not fully exploited in existing networks. In order to activate more input pixels for better restoration, we propose a new Hybrid Attention Transformer (HAT). It combines both channel attention and window-based self-attention schemes, thus making use of their complementary advantages. Moreover, to better aggregate the cross-window information, we introduce an overlapping cross-attention module to enhance the interaction between neighboring window features. In the training stage, we additionally adopt a same-task pre-training strategy to further exploit the potential of the model for further improvement. Extensive experiments have demonstrated the effectiveness of the proposed modules. We further scale up the model to show that the performance of the SR task can be greatly improved. Besides, we extend HAT to more image restoration applications, including real-world image super-resolution, Gaussian image denoising and image compression artifacts reduction. Experiments on benchmark and real-world datasets demonstrate that our HAT achieves state-of-the-art performance both quantitatively and qualitatively. Codes and models are publicly available at <https://github.com/XPixelGroup/HAT>.

Index Terms—Image restoration, image super-resolution, image denoising, Transformer

I. INTRODUCTION

IMAGE restoration (IR) is a classic problem in computer vision. It aims to reconstruct a high-quality (HQ) image from a given low-quality (LQ) input. Classic IR tasks encompass image super-resolution, image denoising, compression artifacts reduction, and etc. Image restoration plays an important role in computer vision and has widespread application in areas such as AI photography [1], surveillance imaging [2], medical imaging [3], and image generation [4]. Since deep learning has been successfully applied to IR tasks [5]–[7], numerous methods based on the convolutional neural network (CNN) have been proposed [8]–[13] and almost dominate this field in the past few years. Recently, due to the success in natural language processing, Transformer [14] has attracted

Xiangyu Chen and Jiantao Zhou are with State Key Laboratory of Internet of Things for Smart City, University of Macau.

Xiangyu Chen, Yu Qiao and Chao Dong are with Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China.

Xiangyu Chen, Wenlong Zhang, Xiangtao Kong, Chao Dong and Yu Qiao are with Shanghai Artificial Intelligence Laboratory, Shanghai, China.

Xiantao Wang is with the ARC Lab, Tencent PCG, Shenzhen, China.

Corresponding Authors: Jiantao Zhou (jtzhou@um.edu.mo) and Chao Dong (chao.dong@siat.ac.cn).

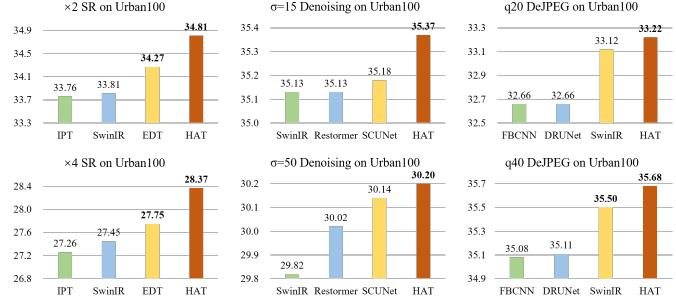


Fig. 1. Performance comparison of the proposed HAT on various image restoration tasks with the state-of-the-art methods.

increasing attention in the computer vision community. After making rapid progress on high-level vision tasks [15]–[17], Transformer-based methods are also developed for low-level vision tasks [18]–[22]. A successful example is SwinIR [22], which obtains a breakthrough improvement on IR tasks.

Despite its success, existing work has rarely discussed why Transformer outperforms CNN. An intuitive explanation provided in prior study is that Transformer benefits from the self-attention mechanism, allowing it to leverage long-range information [22]. To verify whether this is indeed the case for image restoration, we take image super-resolution (SR) as an example task and employ an attribution analysis method — Local Attribution Map (LAM) [23] to examine the range of information used in SwinIR. Interestingly, we find that although SwinIR achieves higher average quantitative performance, it does NOT utilize more input pixels than CNN-based methods (e.g., RCAN [10]), as shown in Fig. 2. This contradicts the conclusion in LAM [23] that there is a positive correlation between the range of information a network uses and its reconstruction performance. Since the aforementioned conclusion is primarily derived from networks of the same type (i.e., CNNs), we believe that the superior performance of SwinIR can be attributed to its stronger ability to model local information compared to CNN. However, it is also limited by the restricted range of information it utilizes, leading to inferior results on samples where broader contextual information could produce better outcomes. Additionally, we observe block artifacts in the intermediate features of SwinIR, as shown in Fig. 4, suggesting that the shifted window mechanism does not fully achieve cross-window information interaction. This may be one of the reasons why SwinIR does not achieve better long-range information utilization.

To address the above-mentioned limitations of the existing IR Transformer and further develop the potential of such

networks, we propose a Hybrid Attention Transformer, namely HAT. It combines channel attention and self-attention schemes, in order to take advantage of the former's capability in using global information and the powerful representative ability of the latter. Besides, we introduce an overlapping cross-attention module to achieve more direct interaction of adjacent window features. Benefiting from these designs, our model can activate more pixels for reconstruction and thus obtains significant performance improvement. Since Transformers do not have an inductive bias like CNNs, large-scale data pre-training is important to unlock the potential of such models. In this paper, we provide an effective *same-task pre-training* strategy. Different from IPT [18] using multiple restoration tasks for pre-training and EDT [21] utilizing multiple degradation levels of a specific task for pre-training, we directly perform pre-training using large-scale dataset on the same task. We believe that large-scale data is what really matters for pre-training. Experimental results show the superiority of our strategy. Equipped with the above designs, HAT surpasses the state-of-the-art methods by a large margin on SR, as well as several other image restoration tasks, as shown in Fig. 1.

Overall, our main contributions are four-fold:

- We design a Hybrid Attention Transformer (HAT) that combines self-attention, channel attention and a new overlapping cross-attention for high-quality image restoration.
- We propose an effective same-task pre-training strategy to further exploit the potential of SR Transformer and show the importance of large-scale data pre-training.
- Our method significantly outperforms existing state-of-the-art methods on the SR task. By further scaling up HAT to build a large model, we greatly extend the performance upper bound of the SR task.
- Our method also achieves state-of-the-art performance on image denoising and compression artifacts reduction, showing its superiority on various image restoration tasks.

A preliminary version of this work was presented at CVPR2023 [24]. The present work expands upon the initial version in several significant ways. Firstly, we provide a theoretic illustration of LAM and augment the analysis with CEM results. This can facilitate readers' understanding of the motivation behind our method and the rationality of its design. Secondly, we investigate a flexible plain architecture of HAT for application to various IR tasks, allowing us to further explore the potential breadth of this method. Thirdly, we extend HAT to real-world image super-resolution based on practical degradation models. The promising results show the potential of HAT for real-world applications. Additionally, we further extend HAT for image denoising and compression artifacts reduction. Extensive experiments show that our method achieves state-of-the-art performance on several IR tasks.

II. RELATED WORK

A. Image Super-Resolution

Since SRCNN [5] first introduces deep convolution neural networks (CNNs) to the image SR task and obtains superior performance over conventional SR methods, numerous deep networks [8], [10], [21], [22], [25]–[32] have been proposed

for SR to further improve the reconstruction quality. For instance, many methods apply more elaborate convolution module designs, such as residual block [28], [33] and dense block [12], [29], to enhance the model representation ability. Several works explore more different frameworks like recursive neural network [34], [35] and graph neural network [36]. To improve perceptual quality, [12], [33], [37], [38] introduce adversarial learning to generate more realistic results. By using attention mechanism, [10], [11], [30]–[32], [39] achieve further improvement in terms of reconstruction fidelity. Recently, a series of Transformer-based networks [18], [21], [22] are proposed and constantly refresh the state-of-the-art of SR task, showing the powerful representation ability of Transformer.

To deepen the understanding of SR networks, several studies [23], [40]–[44] are conducted to analyze and interpret their working mechanism. LAM [23] adopts the integral gradient method to explore which input pixels contribute to the final performance. DDR [40] reveals the deep semantic representations in SR networks based on deep feature dimensionality reduction and visualization. FAIG [41] is proposed to find discriminative filters for specific degradations in blind SR. [42] introduces channel saliency map to show that Dropout can help prevent co-adapting for real SR networks. SRGA [43] aims to evaluate the generalization ability of SR methods. CEM [44] interprets the low-level vision models based on causal effect theory. In this work, we exploit LAM [23] and CEM [44] to analyze and understand the behavior of different networks.

B. Vision Transformer

Recently, Transformer [14] has attracted the attention of computer vision community due to its success in the field of natural language processing. A series of Transformer-based methods [15]–[17], [45]–[53] have been developed for high-level vision tasks, including image classification [15], [16], [45], [54], [55], object detection [16], [49], [56]–[58], segmentation [17], [59]–[61], etc. Although vision Transformer has shown its superiority on modeling long-range dependency [15], [62], there are still many works demonstrating that the convolution can help Transformer achieve better visual representation [46], [51], [52], [63], [64]. Due to the impressive performance, Transformer has also been introduced for low-level vision tasks [18]–[22], [65]–[67]. Specifically, IPT [18] develops a ViT-style network and introduces multi-task pre-training for image processing. SwinIR [22] proposes an image restoration Transformer based on [16]. VRT [67] introduces Transformer-based networks to video restoration. EDT [21] adopts self-attention mechanism and multi-related-task pre-training strategy to further refresh the state-of-the-art of SR. However, existing works still cannot fully exploit the potential of Transformer, while our method can activate more input pixels for better reconstruction.

C. Deep Networks for Image Restoration

Image restoration, which aims to recover high-quality images from degraded inputs, has seen significant progress with the rise of deep learning. Early successes are achieved in tasks like image super-resolution [25], image denoising [7],

and compression artifact reduction [6]. Numerous CNN-based networks have since been proposed for image restoration [9], [19]–[22], [66], [68]–[74]. Before the advent of Transformers in low-level vision tasks, CNNs dominated the field. For example, ARCNN [6] employs stacked convolutional layers to address JPEG compression artifacts, and DnCNN [7] combines convolution with batch normalization for image denoising. RDN [69] introduces a residual dense CNN architecture, excelling in various restoration tasks. A comprehensive investigation of CNN-based methods for SR can be found in [75]. As Transformers have gained prominence in computer vision, Transformer-based image restoration methods have emerged. SwinIR [22], built on the Swin Transformer [16], demonstrates excellent performance on image super-resolution, denoising, and JPEG artifact reduction. Uformer [19] introduces a U-Net-style Transformer for diverse restoration tasks, while Restormer [20] innovates with transposed self-attention to achieve state-of-the-art results. SCUNet [74] combines CNNs and Transformers to create a highly effective denoising network. Transformer-based networks have demonstrated superior performance compared to previous CNN-based methods. In this paper, we introduce a hybrid attention mechanism, further improving the performance of image restoration Transformer.

III. MOTIVATION

Swin Transformer [16] has already demonstrated excellent performance in image restoration tasks [22]. We are thus eager to understand what makes it superior to CNN-based methods and what its potential shortcomings are that could be improved. To explore these questions, we seek to derive insights using interpretability tools and visualization analysis. Given that existing IR Transformer shows remarkable advancements on SR, and that the analytical tools developed for SR are more mature, we focus our analysis primarily on SR. In this section, we present the motivation behind our approach. We begin by reviewing the LAM method, an attribution analysis tool for SR networks. Next, we apply LAM to several classical SR networks and augment the results with causal analysis using the Causal Effect Map (CEM) [44]. Finally, we present feature visualization results that provide additional insights.

A. An overview of LAM

Local Attribution Map (LAM) [23] is an attribution analysis method designed for SR networks. It builds on the integrated gradient method [76], but introduces modifications tailored to the SR task. By constructing a baseline input and a path function appropriate for SR, LAM provides insights into how SR networks utilize spatial information from the input image.

Formally, Let $F : \mathbb{R}^n \mapsto [0, 1]$ represent a deep classification network, $I \in \mathbb{R}^n$ be the input, and $I' \in \mathbb{R}^n$ denote the baseline input (e.g., a black image), which means the absence of important features. The integrated gradient along the i^{th} dimension, $\text{IntegratedGrads}_i(I)$, measures the contribution of each input dimension to the output and is defined as:

$$\text{IntegratedGrads}_i(I) := (I_i - I'_i) \times \int_{\theta=0}^1 \frac{\partial F(I'_i + \theta \times (I - I'_i))}{\partial I_i} d\theta. \quad (1)$$

Let $\lambda(\theta) : [0, 1] \mapsto \mathbb{R}^{h \times w}$ be a smooth function that describes a path from the baseline I' to the input I (i.e., $\lambda(0) = I'$ and $\lambda(1) = I$). The path-integrated gradient along the i^{th} dimension for input I is given by:

$$\text{PathIntegratedGrads}_i^\lambda(I) := \int_0^1 \frac{\partial F(\lambda(\theta))}{\partial \lambda_i(\theta)} \times \frac{\partial \lambda_i(\theta)}{\partial \theta} d\theta, \quad (2)$$

where $\frac{\partial F(I)}{\partial I_i}$ is the gradient along the i^{th} dimension at I .

While the above formulation works for classification networks, it is not directly applicable to the SR task. SR networks focus on reconstructing high-frequency details, such as textures and edges, and thus require a different baseline input and path function. Therefore, instead of using a black image, LAM selects the blurred version of the LR image as the baseline input I' , which is defined as:

$$I' = \omega(\sigma) \otimes I, \quad (3)$$

where \otimes denotes the convolution operation and $\omega(\sigma)$ is a Gaussian blur kernel parameterized by the kernel width σ . This choice of baseline effectively removes the high-frequency components (e.g., textures) that are crucial for SR. To define a reasonable integration path, LAM introduces a progressive blurring path function, denoted by λ_{pb} , which smoothly transitions from the baseline I' to the input I by progressively reducing the blur. The path function is defined as:

$$\lambda_{pb}(\theta) = \omega(\sigma - \theta\sigma) \otimes I, \quad (4)$$

where $\lambda_{pb}(0) = I'$ and $\lambda_{pb}(1) = I$. This path ensures a gradual recovery of high-frequency details as we move from the blurred baseline to the original input image.

Unlike classification, where the output $F(I)$ (e.g., the Softmax output) directly reflects the network's response, the output of SR networks is a reconstructed image. In SR, the key interest is whether the reconstructed image contains clear edges and textures, rather than pixel intensity alone. To capture this, LAM employs a gradient detector D_{xy} (e.g., Gabor filter [77]) to quantify the presence of edges and textures in a specific $l \times l$ patch located at (x, y) , which is denoted as:

$$D_{xy}(I) = \sum_{m \in [x, x+l], n \in [y, y+l]} \nabla_{mn} I. \quad (5)$$

where $\nabla_{mn} I$ represents the gradient at location (m, n) , indicating the presence of edges or textures.

Overall, LAM interprets the SR network F by combining the gradient detector D with the progressive blurring path function λ_{pb} . The contribution of the i^{th} input pixel to the reconstruction of edges and textures in the output is given by:

$$\text{LAM}_{F,D}(\lambda_{pb})_i := \int_0^1 \frac{\partial D(F(\lambda_{pb}(\theta)))}{\partial \lambda_{pb}(\theta)_i} \times \frac{\partial \lambda_{pb}(\theta)_i}{\partial \theta} d\theta. \quad (6)$$

As shown in Fig. 2, LAM visualizations depict how input pixels influence the SR reconstruction. The more red points in the LAM results, the more pixels are involved in the reconstruction. Darker red points indicate a stronger influence of the corresponding pixel on the reconstruction result.

To quantify the range of pixels involved in the reconstruction, LAM introduces the Diffusion Index (DI), a metric based

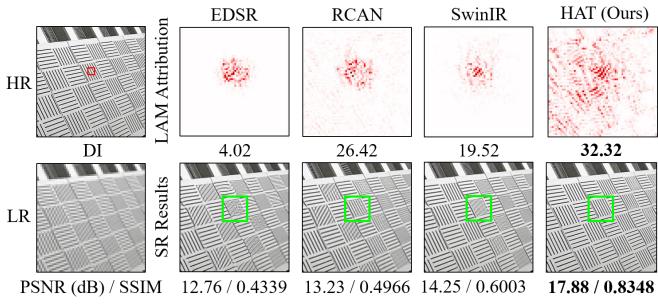


Fig. 2. LAM [23] results of different networks. SwinIR utilizes less information compared to RCAN, while HAT uses the most pixels for reconstruction.

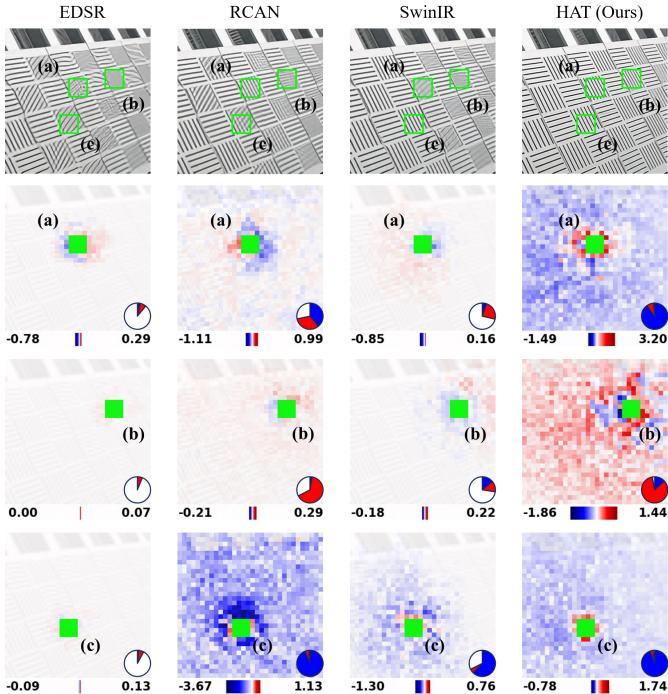


Fig. 3. CEM [44] results of different networks. Activating more input information for Transformer is crucial to the reconstruction performance.

on the Gini coefficient [78]. A larger DI indicates that a wider range of pixels contributes to the reconstruction process, as demonstrated in Fig. 2. SR networks with a higher DI tend to involve more pixels in the SR task, which often correlates with better reconstruction quality [23].

B. Interpretability analysis

We first employ LAM to perform attribution analysis on several classic SR networks, as shown in Fig. 2. Intuitively, SR networks that utilize more input information achieve superior reconstruction performance. This relationship is clearly observed in the comparison between EDSR [28] and RCAN [10]. However, this conclusion is the opposite in the comparison between RCAN and SwinIR. SwinIR achieves better reconstruction results despite utilizing significantly less input information. First, this LAM observation contradicts the intuition in existing literature [22], which suggests that Transformers

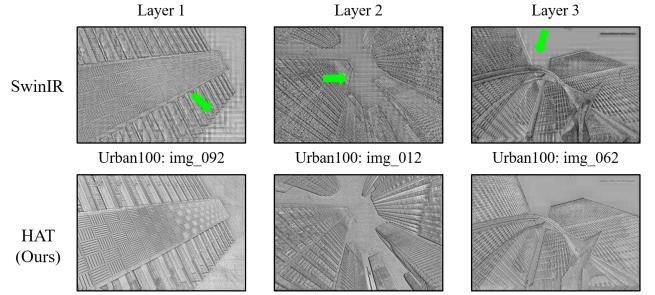


Fig. 4. Intermediate features visualization. “Layer N” means the intermediate features after the N_{th} layer (*i.e.*, RSTB in SwinIR and RHAG in HAT).

perform better by more effectively modeling long-range dependency. Second, it means that SwinIR, which employs a window-based self-attention mechanism, excels at capturing local information and can achieve superior performance with less input information. Additionally, we observe that SwinIR produces incorrect texture reconstruction in case where RCAN successfully restores the texture, which may be attributed to SwinIR’s limited information utilization.

To further investigate the network behaviors, we apply a causal analysis method designed for low-level vision tasks, Causal Effect Map (CEM) [44]. As shown in Fig. 3, CEM results reflect the causal effect of each patch within the input image on the network’s reconstruction of the ROI, (*i.e.* the green marked area). Patches with positive or negative causal effects are indicated in red and blue. The pie chart records the percentage of patches with different causal effects, and the color bar expresses the effect range. For example (a), RCAN activates a large area of input information, with many patches exerting a positive effect on reconstructing the ROI. Conversely, SwinIR utilizes much less input information, leading to wrong texture. Our HAT, on the other hand, activates information over a near-global area but primarily relies on local information around the ROI to achieve accurate reconstruction. Example (b) presents a more challenging scenario: neither RCAN nor SwinIR successfully restores the correct texture. However, RCAN utilizes a wider range of input information, resulting in a clearer texture than SwinIR. HAT leverages both local and global information more effectively, producing superior results. Example (c) further highlights an interesting case: RCAN utilizes almost all input information, but much of it yields a negative effect on ROI reconstruction, leading to incorrect textures in an otherwise relatively simple scenario. In contrast, SwinIR, despite using limited input, reconstructs accurate textures. In this case, HAT again demonstrates superior performance by effectively utilizing neighborhood information to produce both accurate textures and sharper edges.

In conclusion, we posit that the performance of SR networks primarily depends on their ability to utilize information from the region to be reconstructed with its neighborhood. While the use of global information does not always yield beneficial outcomes, it can be crucial in certain cases. Therefore, enhancing the effective range of information utilized by Transformer, as well as their ability to capture global information, is instrumental in developing more advanced networks.

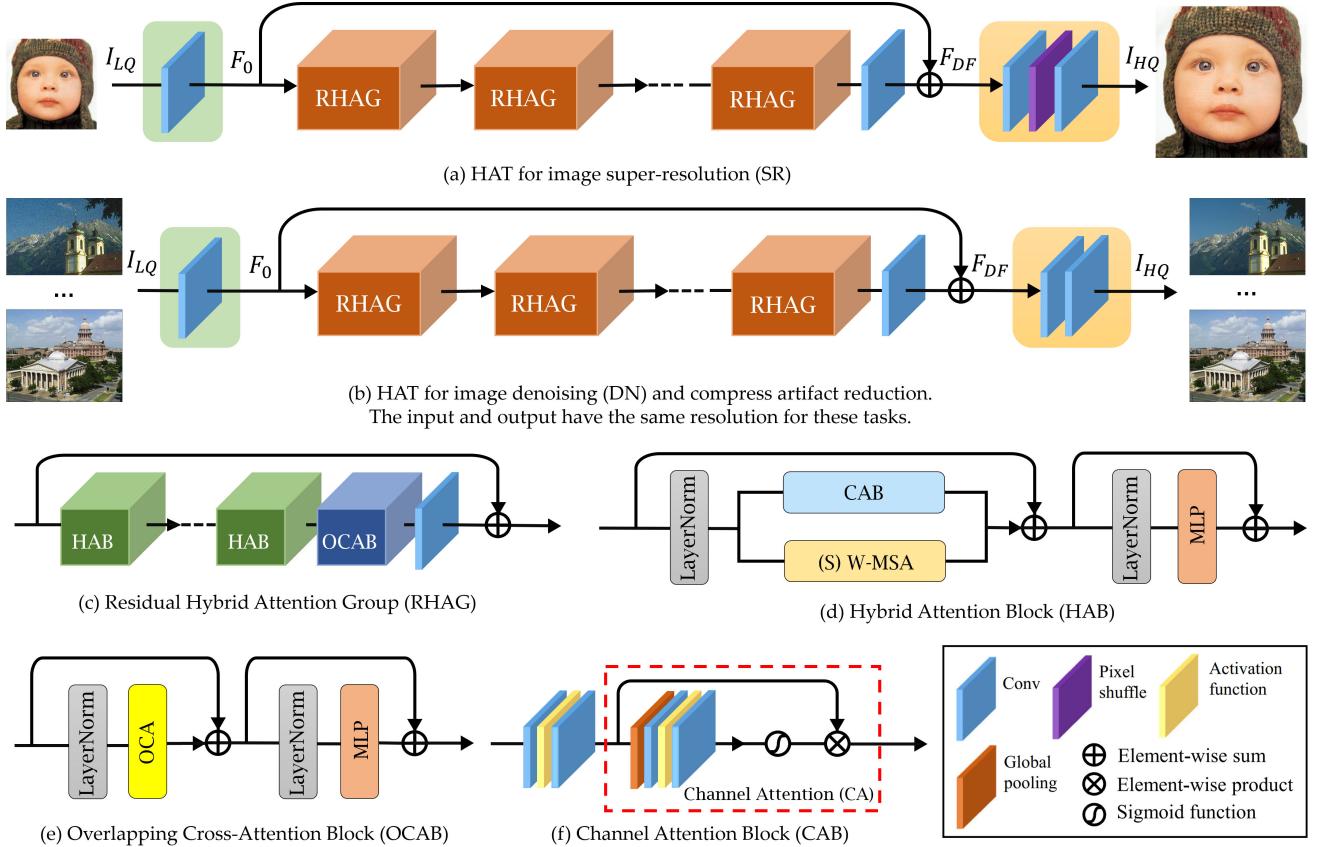


Fig. 5. The overall architecture of HAT and the structure of RHAG and HAB.

C. Feature visualization

SwinIR, as a new architecture distinct from traditional CNN designs, motivates us to examine its intermediate features to gain further insights. As shown in Figure 4, we observe noticeable block artifacts in SwinIR. Interestingly, the size of these blocks coincides with the window size, suggesting that these artifacts are likely caused by the window partitioning mechanism. This indicates that the shifted window approach may be insufficient for effectively integrating information across windows. This limitation could be one of the reasons why SwinIR fails to utilize more pixels for reconstruction, as evidenced in Figures 2 and 3. Several studies on high-level vision tasks have also pointed out that enhancing connections between windows can improve window-based self-attention mechanisms [47], [48], [50], [53]. Consequently, we enhance the interaction of information across windows in our method. We can see that the block artifacts in the intermediate features of our HAT are significantly alleviated.

IV. METHODOLOGY

Based on the above analysis, we aim to design a better image restoration network by enhancing the ability of the existing Transformer model to efficiently utilize more input information, integrating global information, and improving the cross-window interaction. In this section, we provide a detailed introduction to our approach, HAT, including the overall architecture, key module designs, training strategy, implementation details, as well as discussions with other methods.

A. Network Structure of HAT

The overall network structure of HAT follows the classic Residual in Residual (RIR) architecture similar to [10], [22]. As shown in Fig. 5, HAT consists of three parts, including shallow feature extraction, deep feature extraction and image reconstruction. Concretely, for a given low-quality (LQ) input image $I_{LQ} \in \mathbb{R}^{H \times W \times C_{in}}$, we use one 3×3 convolution layer $H_{Conv}(\cdot)$ to extract the shallow feature $F_0 \in \mathbb{R}^{H \times W \times C}$ as:

$$F_0 = H_{SF}(I_{LQ}), \quad (7)$$

where C_{in} and C denote the channel number of the input and the intermediate feature, respectively. The shallow feature extraction can simply map the input from low-dimensional space to high-dimensional space, while achieving the high-dimensional embedding for each pixel token. Moreover, the early convolution layer can help learn better visual representation [52] and lead to stable optimization [63]. We then perform deep feature extraction $H_{DF}(\cdot)$ to further obtain the deep feature $F_{DF} \in \mathbb{R}^{H \times W \times C}$ as:

$$F_{DF} = H_{DF}(F_0), \quad (8)$$

where $H_{DF}(\cdot)$ consists of N_1 residual hybrid attention groups (RHAG) and one 3×3 convolution layer $H_{Conv}(\cdot)$. These RHAGs progressively process the intermediate features as:

$$\begin{aligned} F_i &= H_{RHAG_i}(F_{i-1}), i = 1, 2, \dots, N, \\ F_{DF} &= H_{Conv}(F_N), \end{aligned} \quad (9)$$

where $H_{RHAG_i}(\cdot)$ represents the i -th RHAG. Following [22], we also introduce a convolution layer at the tail of this part to better aggregate information of deep features. After that, we add a global residual connection to fuse shallow features and deep features, and then reconstruct the high-quality (HQ) result via a reconstruction module as:

$$I_{HQ} = H_{Rec}(F_0 + F_{DF}), \quad (10)$$

where $H_{Rec}(\cdot)$ denotes the reconstruction module. We adopt the pixel-shuffle method [26] to up-sample the fused feature for the SR task shown in Fig. 5(a), and use two convolutions for the tasks where have the input and output have the same resolution shown in Fig. 5(b). The key component RHAG consists of N_2 hybrid attention blocks (HAB), one overlapping cross-attention block (OCAB), one 3×3 convolution layer, with a residual connection, as presented in Fig. 5(c).

B. Hybrid Attention Block

In this section, we detail our proposed Hybrid Attention Block (HAB), illustrated in Fig. 5(d). HAB adopts a structure similar to the standard Swin Transformer block, preserving the window-based self-attention mechanism. However, we enhance the representative ability of self-attention and introduce a channel attention block to capture global information. As discussed in Sec.III, we aim to activate more input pixels for Transformer to achieve stronger reconstruction capability. Unlike convolution that expands the receptive field by stacking layers, self-attention possesses a global receptive field within its scope. Therefore, a natural approach to expand the range of information utilized by window-based self-attention is to enlarge the window size. Previous work [22] limits the window size used for self-attention calculations to a small range (i.e., 7 or 8), relying on a shifted window mechanism to gradually expand the receptive field. While this method reduces computational cost, it compromises the effectiveness of self-attention. As discussed in Sec. V-A, we find that window size is a crucial factor influencing the ability of window-based self-attention to exploit information. Appropriately increasing the window size can significantly improve the Transformer's performance. Therefore, in HAB, we adopt a larger window size (i.e., 16).

In addition to window-based self-attention, global information can also be captured by incorporating channel attention. We think that global information may help for cases where many similar textures are present, as shown in Fig.3. Moreover, several studies [20], [72] have demonstrated that channel-wise dynamic mapping is beneficial for low-level vision tasks. Therefore, we introduce the channel attention mechanism into our network. Given the evidence that convolution can improve the visual representation of Transformer models and facilitate easier optimization [46], [51], [52], [63], [79], we incorporate a channel attention-based convolution block, referred to as the channel attention block (CAB), into the standard Transformer block to construct our HAB (see Fig. 5(f)). In addition, global information can be utilized when channel attention is adopted, as it is involved to calculate the channel attention weights. We think that the global information may help for cases where many similar textures exist, such as the examples given in

Fig. 2. Besides, several works [20], [72] have demonstrated that the channel-wise dynamic mapping is beneficial for low-level vision tasks. Therefore, we want to introduce the channel attention mechanism to our network. Since many works have shown that convolution can help Transformer get better visual representation or achieve easier optimization [46], [51], [52], [63], [79], we incorporate a channel attention-based convolution block, i.e., the channel attention block (CAB) in Fig. 5(f), into the standard Transformer block to build our HAB.

To avoid the possible conflict of CAB and MSA on optimization and visual representation, we combine them in parallel and set a small constant α to control the weight of the CAB output. Overall, for a given input feature X , the whole process of HAB is computed as:

$$\begin{aligned} X_N &= \text{LN}(X), \\ X_M &= (\text{S})\text{W-MSA}(X_N) + \alpha \text{CAB}(X_N) + X, \\ Y &= \text{MLP}(\text{LN}(X_M)) + X_M, \end{aligned} \quad (11)$$

where X_N and X_M denote the intermediate features. Y represents the output of HAB. LN represents the layer normalization operation and MLP denotes a multi-layer perceptron. (S)W-MSA means the standard and shifted window multihead self-attention modules. Especially, we treat each pixel as a token for embedding (i.e., set patch size as 1 for patch embedding following [22]). For calculation of the self-attention module, given an input feature of size $H \times W \times C$, it is first partitioned into $\frac{HW}{M^2}$ local windows of size $M \times M$, then self-attention is calculated inside each window. For a local window feature $X_W \in \mathbb{R}^{M^2 \times C}$, the *query*, *key* and *value* matrices are computed by linear mappings as Q , K and V . Then the window-based self-attention is formulated as:

$$\text{Attention}(Q, K, V) = \text{SoftMax}(QK^T / \sqrt{d} + B)V, \quad (12)$$

where d represents the dimension of *query/key*. B denotes the relative position encoding and is calculated as [14]. Besides, to build the connections between neighboring non-overlapping windows, we also use the shifted window partitioning approach [16], with the shift size set to half of the window size.

A CAB consists of two standard convolution layers with GELU activation [80] and a channel attention (CA) module, as shown in Fig.5(f). Since Transformer-based structures often require a large number of channels for token embedding, directly using convolutions with constant width would result in high computational costs. To address this, we compress the number of channels in the two convolution layers by a constant factor β . For an input feature with C channels, the number of channels is reduced to $\frac{C}{\beta}$ after the first convolution layer, and then expanded back to C channels through the second layer. Finally, a standard CA module [10] is employed to adaptively rescale channel-wise features.

C. Overlapping Cross-Attention Block (OCAB)

We introduce OCAB to directly establish cross-window connections and enhance the representative ability for the window self-attention. Our OCAB consists of an overlapping cross-attention (OCA) layer and an MLP layer similar to the standard Swin Transformer block [16]. But for OCA, as depicted in

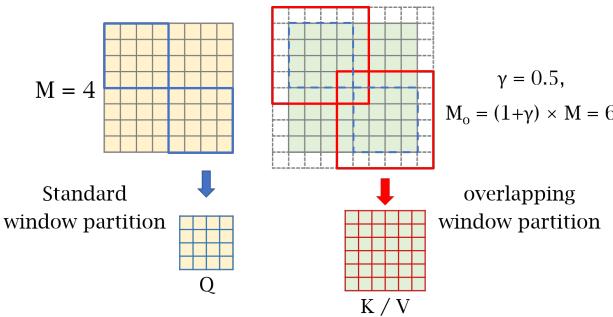


Fig. 6. The overlapping window partition for OCA.

Fig. 6, we use different window sizes to partition the projected features. Specifically, for the $X_Q, X_K, X_V \in \mathbb{R}^{H \times W \times C}$ of the input feature X , X_Q is partitioned into $\frac{HW}{M^2}$ non-overlapping windows of size $M \times M$, while X_K, X_V are unfolded to $\frac{HW}{M^2}$ overlapping windows of size $M_o \times M_o$. It is calculated as

$$M_o = (1 + \gamma) \times M, \quad (13)$$

where γ is a constant to control the overlapping size. To better understand this operation, the standard window partition can be considered as a sliding partition with the kernel size and the stride both equal to the window size M . In contrast, the overlapping window partition can be viewed as a sliding partition with the kernel size equal to M_o , while the stride is equal to M . Zero-padding with size $\frac{\gamma M}{2}$ is used to ensure the size consistency of overlapping windows. The attention matrix is calculated as Eq. (12), and the relative position bias $B \in \mathbb{R}^{M \times M_o}$ is also adopted. Unlike WSA whose *query*, *key* and *value* are calculated from the same window feature, OCA computes *key/value* from a larger field where more useful information can be utilized for the *query*. Note that although Multi-resolution Overlapped Attention (MOA) module in [53] performs similar overlapping window partition, our OCA is fundamentally different from MOA. MOA calculates global attention using window features as tokens, while OCA computes cross-attention inside each window using pixel tokens.

D. The Same-task Pre-training

Pre-training is proven effective on many high-level vision tasks [15], [81], [82]. Recent works [18], [21] also demonstrate that pre-training is beneficial to low-level vision tasks. IPT [18] emphasizes the use of various low-level tasks, such as denoising, deraining, SR and *etc.*, while EDT [21] utilizes different degradation levels of a specific task to do pre-training. These works focus on investigating the effect of multi-task pre-training for a target task. In contrast, we directly perform pre-training on a larger-scale dataset (*i.e.*, ImageNet [83]) based on the same task, showing that the effectiveness of pre-training depends more on the scale and diversity of data. For example, when we want to train a model for $\times 4$ SR, we first train a $\times 4$ SR model on ImageNet, then fine-tune it on the specific dataset, such as DF2K. The proposed strategy, namely *same-task pre-training*, is simpler while bringing more performance improvements. It is worth mentioning that sufficient training iterations for pre-training and an appropriate small learning rate for fine-tuning are very important for the effectiveness of the

pre-training strategy. We think that it is because Transformer requires more data and iterations to learn general knowledge for the task, but needs a small learning rate for fine-tuning to avoid overfitting to the specific dataset.

E. Discussions

In this part, we analyze the distinctions of our HAT and several relevant works, including SwinIR [22], EDT [21], SCUNet [74] and HaloNet [55].

Difference to SwinIR. SwinIR [22] is the first work to successfully use Swin Transformer [16] for low-level vision tasks. It builds an image restoration network by using the original Swin Transformer block. Our HAT is inspired by SwinIR and retains the core design of window-based self-attention. However, we address the problem of limited range of utilized information in SwinIR by enlarging the window size and introducing channel attention. At the same time, we introduce a newly designed OCA to further enhance the ability to implement cross-window interaction. This work aims to design a more powerful backbone for image restoration tasks.

Difference to EDT. EDT [21] builds an image restoration Transformer based on the shifted crossed local attention, which also calculates self-attention in the windows of fixed size. For a given feature map, it splits the feature into two parts and performs self-attention in an either horizontal or vertical rectangle window. In contrast, our HAT adopts the vanilla window-based self-attention and shifted window mechanism similar to Swin Transformer [16]. EDT also studies the pre-training strategy and emphasizes the advantages of multi-related-task pre-training (*i.e.*, performing pre-training on a specific task with multiple degradation levels). However, HAT shows that training on the same task but using a large-scale dataset is the key factor in the effectiveness of pre-training.

Difference to SCUNet. SCUNet [74] is also an image restoration network that integrates the strengths of Transformers and CNN. It utilizes the Swin Transformer block alongside a classic convolution block within its U-Net architecture, forming a Swin-Conv Block, and achieves excellent performance for denoising. Unlike our approach that originates from the SR task, SCUNet is primarily designed for denoising, with a focus on capturing multi-scale information. In contrast, our method emphasizes the benefits of window self-attention for local information fitting and addresses its limitations in cross-window interaction and global information acquisition. Therefore, the two methods differ significantly in motivation, overall architecture and the design details of key modules.

Difference to HaloNet. HaloNet [55] incorporates a similar window partition mechanism to our OCA, enabling the calculation of self-attention within overlapping window features. HaloNet employs this overlapping self-attention as the fundamental module to build the network, inevitably leading to a large computational cost. This design could impose a substantial computational burden and is not friendly to image restoration tasks. On the contrary, our HAT leverages only a limited number of OCA modules to augment the interaction between adjacent windows. This approach can effectively enhance the image restoration Transformer without incurring excessive computational costs.

TABLE I
QUANTITATIVE RESULTS ON PSNR(DB) FOR DIFFERENT WINDOW SIZES.

Size	Set5	Set14	BSD100	Urban100	Manga109
(8,8)	32.88	29.09	27.92	27.45	32.03
(16,16)	32.97	29.12	27.95	27.81	32.15

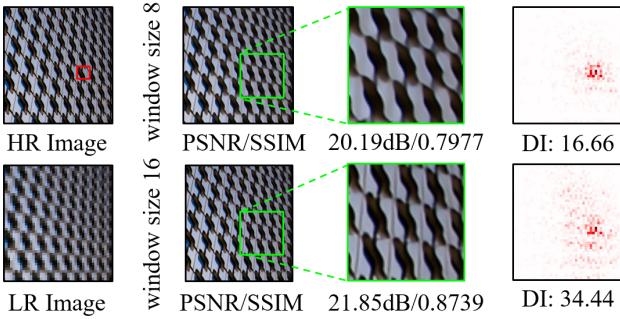


Fig. 7. Qualitative comparison of different window sizes.

F. Implementation Details

For the structure of HAT, both the RHAG number and HAB number are set to 6. The channel number of the whole network is set to 180. The attention head number and window size are set to 6 and 16 for both (S)W-MSA and OCA. For the specific hyper-parameters of the proposed modules, we set the weighting factor of CAB output (α), the squeeze factor between two convolution layers in CAB (β), and the overlapping ratio of OCA (γ) as 0.01, 3 and 0.5, respectively. For the large variant HAT-L, we double the depth of HAT by increasing the RHAG number from 6 to 12. We also provide a small version HAT-S with fewer parameters and similar computation to SwinIR. For HAT-S, we set the channel number to 144 and set β to 24 in CAB. When implementing the pre-training strategy, we adopt ImageNet [83] as the pre-training dataset following [18], [21]. We conduct the main experiments and ablation study on image SR. Therefore, we use the DF2K dataset (DIV2K [84]+Flicker2K [85]) as the training dataset, following [22], [86]. PSNR/SSIM calculated on the Y channel is reported for the quantitative metrics.

V. NETWORK INVESTIGATION

A. Effects of different window sizes

As discussed in Sec. III, activating more input pixels for the restoration tends to achieve better performance. Enlarging window size for the window-based self-attention is an intuitive way to realize the goal. In this section, we explore how the window size of self-attention influences the performance of Transformer on image SR. To eliminate the influence of our newly-introduced blocks, we conduct the following experiments directly on the preliminary version of SwinIR. As shown in Tab. I, the model with a large window size of 16×16 obtains better performance, especially on Urban100. We also provide the qualitative comparison in Fig. 7. The result produced by the model with a larger window size has much clearer textures. For LAM results, the model with window size of 16 utilizes much more input pixels than the model with window size of 8. These experiments show that enlarging the window size can effectively improve the performance of SR Transformer.

TABLE II
ABLATION STUDY ON THE PROPOSED OCAB AND CAB.

	Baseline			
OCAB	X	✓	X	✓
CAB	X	X	✓	✓
Set5	32.99dB	33.02dB	33.00dB	33.04dB
Set14	29.13dB	29.19dB	29.16dB	29.23dB
Urban100	27.81dB	27.91dB	27.91dB	27.97dB

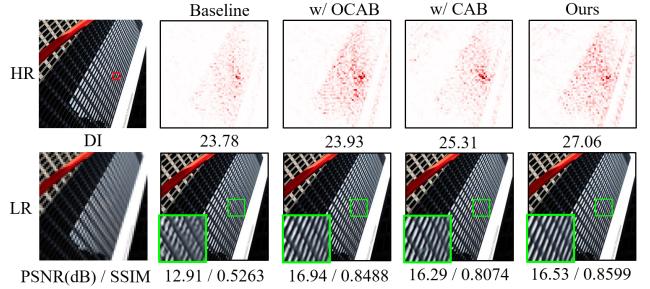


Fig. 8. Ablation study on the proposed OCAB and CAB.

B. Ablation Study

Effectiveness of OCAB and CAB. We conduct experiments to demonstrate the effectiveness of the proposed CAB and OCAB. The quantitative performance reported on three benchmark datasets for $\times 4$ SR is shown in Tab. II. On Urban100, compared with the baseline results, both OCAB and CAB bring the performance gain of 0.1dB. Benefiting from the two modules, the model obtains a further performance gain of 0.16dB. On Set5 and Set14, the proposed OCAB and CAB can also bring considerable performance improvement. We think that the performance improvement comes from two aspects. On the one hand, improving the window interaction by OCAB and utilizing the global statistics by CAB both help the model better deal with the long-term patterns (e.g., self-similarity of repeated textures). On the other hand, the two modules enrich and enhance the model ability by introducing cross-attention and convolution blocks. We also provide qualitative comparison to further illustrate the influence of OCAB and CAB, as presented in Fig. 8. We can observe that the model with OCAB has a larger scope of the utilized pixels and generate better-reconstructed results. When CAB is adopted, the used pixels even expand to almost the full image. Moreover, the result of our method with OCAB and CAB obtains the highest DI [23], which means our method utilizes the most input pixels.

Effects of different designs of CAB. We conduct experiments to explore the effects of different designs of CAB. First, we investigate the influence of channel attention. As shown in Tab. III, the model using CA achieves a performance gain of 0.05dB compared to the model without CA. It demonstrates the effectiveness of the channel attention in our network. We also conduct experiments to explore the effects of the weighting factor α of CAB. As presented in Sec. IV-B, α is used to control the weight of CAB features for feature fusion. A larger α means a larger weight of features extracted by CAB and $\alpha = 0$ represents CAB is not used. As shown in Tab. IV, the model with α of 0.01 obtains the best performance. It indicates that CAB and self-attention may have potential issue in optimization, while a small weighting factor for the CAB

TABLE III
EFFECTS OF THE CHANNEL ATTENTION (CA) MODULE IN CAB.

Structure	w/o CA	w/ CA
PSNR / SSIM	27.92dB / 0.8362	27.97dB / 0.8367

TABLE IV
EFFECTS OF THE WEIGHTING FACTOR α IN CAB.

α	0	1	0.1	0.01
PSNR	27.81dB	27.86dB	27.90dB	27.97dB

TABLE V
ABLATION STUDY ON THE OVERLAPPING RATIO OF OCAB.

γ	0	0.25	0.5	0.75
PSNR	27.85dB	27.81dB	27.91dB	27.86dB

branch can suppress this issue for the better combination.

Effects of the overlapping ratio. In OCAB, we set a constant γ to control the overlapping size for the overlapping cross-attention, as illustrated in Sec IV-C. To explore the effects of different overlapping ratios, we set a group of γ from 0 to 0.75 to examine the performance change, as shown in Tab. V. Note that $\gamma = 0$ means a standard Transformer block. It can be found that the model with $\gamma = 0.5$ performs best. In contrast, when γ is set to 0.25 or 0.75, the model has no obvious performance gain or even has a performance drop. It illustrates that inappropriate overlapping size cannot benefit the interaction of neighboring windows.

C. Analysis of Model Complexity

We first conduct experiments to analyze the computational complexity of our method from three aspects: window size for the calculation of self-attention, ablation of OCAB and CAB, and the different designs of CAB. We evaluate the performance based on the results of $\times 4$ SR on Urban100 and the number of Multiply-Add operations is counted at the input size of 64×64 . Note that pre-training is **Not** used for all models, and all methods are fairly compared under the same training settings.

First, we use the standard Swin Transformer block [22] as the backbone to explore the influence on different window sizes. As shown in Tab. VI, enlarging window size can bring a large performance gain (+0.36dB) with a little increase in parameters and $\sim\!19\%$ increase in Multi-Adds.

Then, we investigate the computational complexity of OCAB and CAB. As illustrated in Tab. VII, our OCAB obtains a performance gain with a limited increase of parameters and Multi-Adds. It demonstrates that the effectiveness and efficiency of the proposed OCAB. Adding CAB to the baseline model will increase the computational complexity, but it can achieve considerable performance improvement.

Since CAB seems to be computationally expensive, we further explore the influence on CAB sizes by modulating the squeeze factor β (mentioned in Sec IV-B). As shown in Tab. VIII, adding a small CAB whose β equals 6 can bring considerable performance improvement. When we continuously reduce β , the performance increases but with larger model sizes. To balance the performance and computations, we set β to 3 as the default setting.

TABLE VI
MODEL COMPLEXITY FOR DIFFERENT WINDOW SIZES.

window size	#Params.	#Multi-Adds.	PSNR
(8, 8)	11.9M	53.6G	27.45dB
(16, 16)	12.1M	63.8G	27.81dB

TABLE VII
MODEL COMPLEXITY FOR OCAB AND CAB.

Method	#Params.	#Multi-Adds.	PSNR
Baseline	12.1M	63.8G	27.81dB
w/ OCAB	13.7M	74.7G	27.91dB
w/ CAB	19.2M	92.8G	27.91dB
Ours	20.8M	103.7G	27.97dB

TABLE VIII
MODEL COMPLEXITY FOR DIFFERENT CAB SIZES.

β in CAB	#Params.	#Multi-Adds.	PSNR
1	33.2M	150.1G	27.97dB
2	22.7M	107.1G	27.92dB
3 (default)	19.2M	92.8G	27.91dB
6	15.7M	78.5G	27.88dB
w/o CAB	12.1M	63.8G	27.81dB

TABLE IX
MODEL COMPLEXITY COMPARISON OF SWINIR AND HAT.

Method	#Params.	#Multi-Adds.	PSNR
SwinIR	11.9M	53.6G	27.45dB
HAT-S (ours)	9.6M	54.9G	27.80dB
SwinIR-L1	24.0M	104.4G	27.53dB
SwinIR-L2	23.1M	102.4G	27.58dB
HAT (ours)	20.8M	103.7G	27.97dB

Furthermore, we compare HAT and SwinIR with the similar numbers of parameters and Multi-Adds in two settings, as presented in Tab. IX. First, we compare HAT-S with the original version of SwinIR. With less parameters and comparable computations, HAT-S significantly outperforms SwinIR. Second, we enlarge SwinIR by increasing the width and depth to achieve similar computations to HAT, denoted as SwinIR-L1 and SwinIR-L2. HAT achieves the best performance at the lowest computational cost. This demonstrates that HAT outperforms SwinIR in performance and computational efficiency.

Overall, we find that enlarging the window size for the calculation of self-attention is a very cost-effective way to improve the Transformer model. Moreover, the proposed OCAB can bring an obvious performance gain with limited increase of computations. Although CAB is not as efficient as above two schemes, it can also bring stable and considerable performance improvement. Benefiting from the three designs, HAT can substantially outperforms the state-of-the-art method SwinIR with comparable computations.

D. Study on the pre-training strategy

From Tab. XI, we can see that HAT can benefit greatly from the pre-training strategy, by comparing the performance of HAT and HAT^\dagger . To show the superiority of the proposed same-task pre-training, we also apply the multi-related-task pre-training [21] to HAT for comparison using full ImageNet, under the same training settings as [21]. As depicted as Tab. X, the same-task pre-training performs better, not only

TABLE X
QUANTITATIVE RESULTS ON PSNR(DB) OF HAT USING TWO KINDS OF PRE-TRAINING STRATEGIES ON $\times 4$ SR UNDER THE SAME TRAINING SETTING.

Strategy	Stage	Set5	Set14	Urban100
Multi-related-task pre-training	pre-training	32.94	29.17	28.05
	fine-tuning	33.06	29.33	28.21
Same-task pre-training(ours)	pre-training	33.02	29.20	28.11
	fine-tuning	33.07	29.34	28.28

in the pre-training stage but also in the fine-tuning process. From this perspective, multi-task pre-training probably impairs the restoration performance of the network on a specific degradation, while the same-task pre-training can maximize the performance gain brought by large-scale data. To further investigate the influences of our pre-training strategy for different networks, we apply our pre-training to four networks: SRResNet (1.5M), RRDBNet (16.7M), SwinIR (11.9M) and HAT (20.8M), as shown in Fig. 9. First, we can see that all four networks can benefit from pre-training, showing the effectiveness of the proposed same-task pre-training strategy. Second, for the same type of network (*i.e.*, CNN or Transformer), the larger the network capacity, the more performance gain from pre-training. Third, although with less parameters, SwinIR obtains greater performance improvement from the pre-training compared to RRDBNet. It suggests that Transformer needs more data to exploit the potential of the model. HAT obtains the largest gain from pre-training, indicating the necessity of the pre-training strategy for such large models. Equipped with big models and large-scale data, we show that the performance upper bound of this task can be significantly extended.

VI. EXPERIMENTAL RESULTS

A. Training Settings

For classic image super-resolution, we use DF2K (DIV2K + Flickr2K) with 3450 images as the training dataset when training from scratch. The low-resolution images are generated from the ground truth images by the “bicubic” down-sampling in MATLAB. We set the input patch size to 64×64 and use random rotation and horizontally flipping for data augmentation. The mini-batch size is set to 32 and total training iterations are set to 500K. The learning rate is initialized as 2e-4 and reduced by half at [250K,400K,450K,475K]. For $\times 4$ SR, we initialize the model with pre-trained $\times 2$ SR weights and halve the iterations for each learning rate decay as well as total iterations. We adopt Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ to train the model. When using the same-task pre-training, we exploit the full ImageNet dataset with 1.28 million images to pre-train the model for 800K iterations. The initial learning rate is also set to 2e-4 but reduced by half at [300K,500K,650K,700K,750k]. Then, we adopt DF2K dataset to fine-tune the pre-trained model. For fine-tuning, we set the initial learning rate to 1e-5 and halve it at [125K,200K,230K,240K] for total 250K iterations.

For real-world image super-resolution, we train HAT models based on two simulated real-world degradation models, *i.e.*, BSRGAN [92] and Real-ESRGAN [38]. The total batch size is set to 32 and the input patch size is set to 64×64 . The network

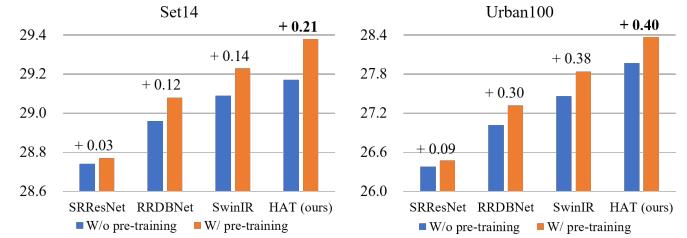


Fig. 9. Quantitative comparison on PSNR(dB) of four different networks without and with the same-task pre-training on $\times 4$ SR.

structure is the same as the basic version of HAT for classic image super-resolution. Following Real-ESRGAN [38], we first train the MSE-based model and introduce the generative adversarial training to fine-tune the GAN-based model. More training and degradation details can refer to [92] and [38].

For image denoising and JPEG compression artifacts reduction, we directly use the combination of DIV2K, Flickr2K, BSD500 [89] and WED images [93] datasets to train the models¹, following [13], [22], [69]. The network is the same as the basic version for classic image super-resolution without up-sampling. Noisy images are generated by adding additive white Gaussian noises with noise level σ and compressed images are obtained by the MATLAB JPEG encoder with JPEG level q . To speed up the training, we first train models with the batch size of 32 and the patch size of 64×64 for 800 iterations. We then proceed to fine-tune models with the batch size of 8 and the patch size of 128×128 for 500 iterations.

B. Classic Image Super-Resolution

Quantitative results. Tab. XI shows the quantitative comparison of our approach with the state-of-the-art methods: EDSR [28], RCAN [10], SAN [30], IGNN [36], HAN [31], NLSN [32], RCAN-it [86], as well as approaches using ImageNet pre-training, *i.e.*, IPT [18] and EDT [21]. We can see that our method significantly outperforms the other methods on all five benchmark datasets. Concretely, with the same depth and width, HAT surpasses SwinIR by 0.48dB~0.64dB on Urban100 and 0.34dB~0.45dB on Manga109. When compared with the approaches using pre-training, HAT[†] also has large performance gains of more than 0.5dB against EDT on Urban100 for all three scales. Besides, HAT equipped with pre-training outperforms SwinIR by a huge margin of up to 1dB on Urban100 for $\times 2$ SR. Moreover, the large model HAT-L can even bring further improvement and greatly expands the performance upper bound of this task. HAT-S with fewer parameters and similar computation can also significantly outperforms the state-of-the-art method SwinIR. (Detailed computational complexity comparison can be found in Sec. V-C.) Note that the performance gaps are much larger on Urban100, as it contains more structured and self-repeated patterns that can provide more useful pixels for reconstruction when the utilized range of information is enlarged. All these results show the effectiveness of our method.

¹ImageNet pre-training has limited impact on the performance of these two tasks. Similar conclusions are also mentioned in [21].

TABLE XI

QUANTITATIVE COMPARISON WITH STATE-OF-THE-ART METHODS FOR CLASSIC IMAGE SUPER-RESOLUTION ON BENCHMARK DATASETS. THE BEST AND SECOND BEST RESULTS ARE MARKED IN **BOLD** AND UNDERLINE. “ \dagger ” INDICATES THAT METHODS ADOPT PRE-TRAINING STRATEGY ON IMAGENET.

Method	Scale	Training Dataset	Set5 [87]		Set14 [88]		BSD100 [89]		Urban100 [90]		Manga109 [91]	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
EDSR [28]	$\times 2$	DIV2K	38.11	0.9602	33.92	0.9195	32.32	0.9013	32.93	0.9351	39.10	0.9773
RCAN [10]	$\times 2$	DIV2K	38.27	0.9614	34.12	0.9216	32.41	0.9027	33.34	0.9384	39.44	0.9786
SAN [30]	$\times 2$	DIV2K	38.31	0.9620	34.07	0.9213	32.42	0.9028	33.10	0.9370	39.32	0.9792
IGNN [36]	$\times 2$	DIV2K	38.24	0.9613	34.07	0.9217	32.41	0.9025	33.23	0.9383	39.35	0.9786
HAN [31]	$\times 2$	DIV2K	38.27	0.9614	34.16	0.9217	32.41	0.9027	33.35	0.9385	39.46	0.9785
NLSN [32]	$\times 2$	DIV2K	38.34	0.9618	34.08	0.9231	32.43	0.9027	33.42	0.9394	39.59	0.9789
RCAN-it [86]	$\times 2$	DF2K	38.37	0.9620	34.49	0.9250	32.48	0.9034	33.62	0.9410	39.88	0.9799
SwinIR [22]	$\times 2$	DF2K	38.42	0.9623	34.46	0.9250	32.53	0.9041	33.81	0.9427	39.92	0.9797
EDT [21]	$\times 2$	DF2K	38.45	0.9624	34.57	0.9258	32.52	0.9041	33.80	0.9425	39.93	0.9800
HAT-S (ours)	$\times 2$	DF2K	38.58	0.9628	34.70	0.9261	32.59	0.9050	34.31	0.9459	40.14	0.9805
HAT (ours)	$\times 2$	DF2K	38.63	0.9630	34.86	0.9274	32.62	0.9053	34.45	0.9466	40.26	0.9809
IPT \dagger [18]	$\times 2$	ImageNet	38.37	-	34.43	-	32.48	-	33.76	-	-	-
EDT \dagger [21]	$\times 2$	DF2K	38.63	0.9632	34.80	0.9273	32.62	0.9052	34.27	0.9456	40.37	0.9811
HAT \dagger (ours)	$\times 2$	DF2K	38.73	<u>0.9637</u>	<u>35.13</u>	<u>0.9282</u>	<u>32.69</u>	<u>0.9060</u>	<u>34.81</u>	<u>0.9489</u>	<u>40.71</u>	<u>0.9819</u>
HAT-L \dagger (ours)	$\times 2$	DF2K	38.91	0.9646	35.29	0.9293	32.74	0.9066	35.09	0.9505	41.01	0.9831
EDSR [28]	$\times 3$	DIV2K	34.65	0.9280	30.52	0.8462	29.25	0.8093	28.80	0.8653	34.17	0.9476
RCAN [10]	$\times 3$	DIV2K	34.74	0.9299	30.65	0.8482	29.32	0.8111	29.09	0.8702	34.44	0.9499
SAN [30]	$\times 3$	DIV2K	34.75	0.9300	30.59	0.8476	29.33	0.8112	28.93	0.8671	34.30	0.9494
IGNN [36]	$\times 3$	DIV2K	34.72	0.9298	30.66	0.8484	29.31	0.8105	29.03	0.8696	34.39	0.9496
HAN [31]	$\times 3$	DIV2K	34.75	0.9299	30.67	0.8483	29.32	0.8110	29.10	0.8705	34.48	0.9500
NLSN [32]	$\times 3$	DIV2K	34.85	0.9306	30.70	0.8485	29.34	0.8117	29.25	0.8726	34.57	0.9508
RCAN-it [86]	$\times 3$	DF2K	34.86	0.9308	30.76	0.8505	29.39	0.8125	29.38	0.8755	34.92	0.9520
SwinIR [22]	$\times 3$	DF2K	34.97	0.9318	30.93	0.8534	29.46	0.8145	29.75	0.8826	35.12	0.9537
EDT [21]	$\times 3$	DF2K	34.97	0.9316	30.89	0.8527	29.44	0.8142	29.72	0.8814	35.13	0.9534
HAT-S (ours)	$\times 3$	DF2K	35.01	0.9325	31.05	0.8550	29.50	0.8158	30.15	0.8879	35.40	0.9547
HAT (ours)	$\times 3$	DF2K	35.07	0.9329	31.08	0.8555	29.54	0.8167	30.23	0.8896	35.53	0.9552
IPT \dagger [18]	$\times 3$	ImageNet	34.81	-	30.85	-	29.38	-	29.49	-	-	-
EDT \dagger [21]	$\times 3$	DF2K	35.13	0.9328	31.09	0.8553	29.53	0.8165	30.07	0.8863	35.47	0.9550
HAT \dagger (ours)	$\times 3$	DF2K	<u>35.16</u>	<u>0.9335</u>	<u>31.33</u>	<u>0.8576</u>	<u>29.59</u>	<u>0.8177</u>	<u>30.70</u>	<u>0.8949</u>	<u>35.84</u>	<u>0.9567</u>
HAT-L \dagger (ours)	$\times 3$	DF2K	35.28	0.9345	31.47	0.8584	29.63	0.8191	30.92	0.8981	36.02	0.9576
EDSR [28]	$\times 4$	DIV2K	32.46	0.8968	28.80	0.7876	27.71	0.7420	26.64	0.8033	31.02	0.9148
RCAN [10]	$\times 4$	DIV2K	32.63	0.9002	28.87	0.7889	27.77	0.7436	26.82	0.8087	31.22	0.9173
SAN [30]	$\times 4$	DIV2K	32.64	0.9003	28.92	0.7888	27.78	0.7436	26.79	0.8068	31.18	0.9169
IGNN [36]	$\times 4$	DIV2K	32.57	0.8998	28.85	0.7891	27.77	0.7434	26.84	0.8090	31.28	0.9182
HAN [31]	$\times 4$	DIV2K	32.64	0.9002	28.90	0.7890	27.80	0.7442	26.85	0.8094	31.42	0.9177
NLSN [32]	$\times 4$	DIV2K	32.59	0.9000	28.87	0.7891	27.78	0.7444	26.96	0.8109	31.27	0.9184
RRDB [12]	$\times 4$	DF2K	32.73	0.9011	28.99	0.7917	27.85	0.7455	27.03	0.8153	31.66	0.9196
RCAN-it [86]	$\times 4$	DF2K	32.69	0.9007	28.99	0.7922	27.87	0.7459	27.16	0.8168	31.78	0.9217
SwinIR [22]	$\times 4$	DF2K	32.92	0.9044	29.09	0.7950	27.92	0.7489	27.45	0.8254	32.03	0.9260
EDT [21]	$\times 4$	DF2K	32.82	0.9031	29.09	0.7939	27.91	0.7483	27.46	0.8246	32.05	0.9254
HAT-S (ours)	$\times 4$	DF2K	32.92	0.9047	29.15	0.7958	27.97	0.7505	27.87	0.8346	32.35	0.9283
HAT (ours)	$\times 4$	DF2K	33.04	0.9056	29.23	0.7973	28.00	0.7517	27.97	0.8368	32.48	0.9292
IPT \dagger [18]	$\times 4$	ImageNet	32.64	-	29.01	-	27.82	-	27.26	-	-	-
EDT \dagger [21]	$\times 4$	DF2K	33.06	<u>0.9055</u>	29.23	<u>0.7971</u>	27.99	<u>0.7510</u>	27.75	<u>0.8317</u>	32.39	0.9283
HAT \dagger (ours)	$\times 4$	DF2K	<u>33.18</u>	<u>0.9073</u>	<u>29.38</u>	<u>0.8001</u>	<u>28.05</u>	<u>0.7534</u>	<u>28.37</u>	<u>0.8447</u>	<u>32.87</u>	<u>0.9319</u>
HAT-L \dagger (ours)	$\times 4$	DF2K	33.30	0.9083	29.47	0.8015	28.09	0.7551	28.60	0.8498	33.09	0.9335

Visual comparison. We provide the visual comparison of different approaches. As shown in Fig. 10, HAT successfully recovers the clear lattice content for the images “img_002”, “img_011”, “img_030”, “img_044” and “img_073” in the Urban100 dataset. In contrast, the other approaches cannot restore correct textures or suffer from severe blurry effects. We can also observe similar behaviors on “PrayerHaNemurenai” in the Manga109 dataset. When recovering the characters in the image, HAT obtains much clearer textures than the other methods. The visual results also demonstrate the superiority of our approach on classic image super-resolution.

LAM comparison. We provide more visual results with LAM to compare the state-of-the method SwinIR and our HAT. As shown in Fig. 11, the utilized pixels for reconstruction of HAT expands to the almost full image, while that of SwinIR only gathers in a limited range. For the quantitative metric, HAT also obtains a much higher DI value than SwinIR. These results demonstrate that our method can activate more pixels to reconstruct the low-resolution input image. As a result, SR results generated by our method have higher PSNR/SSIM and better visual quality. We can observe that HAT restores much clearer textures and edges than SwinIR.

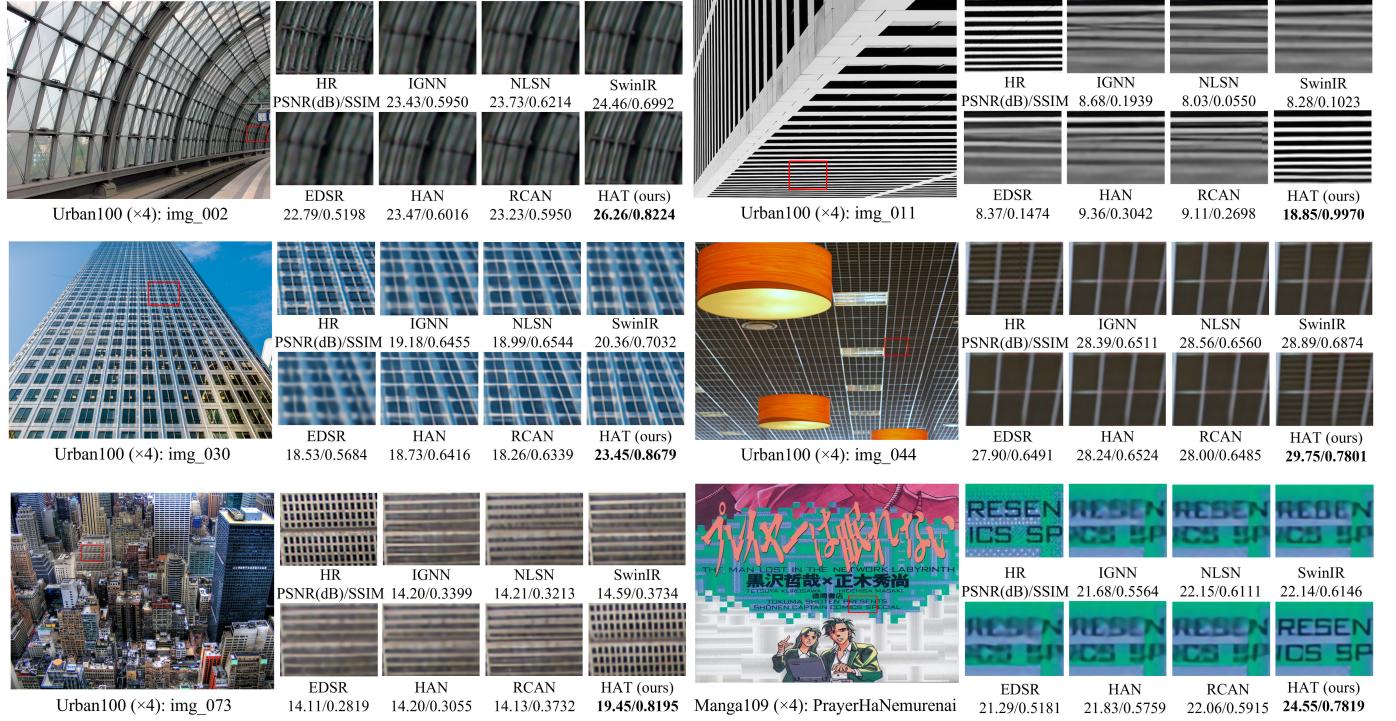


Fig. 10. Visual comparison for $\times 4$ SR. The patches for comparison are marked with red boxes in the original images. PSNR/SSIM is calculated based on the patches to better reflect the performance difference.

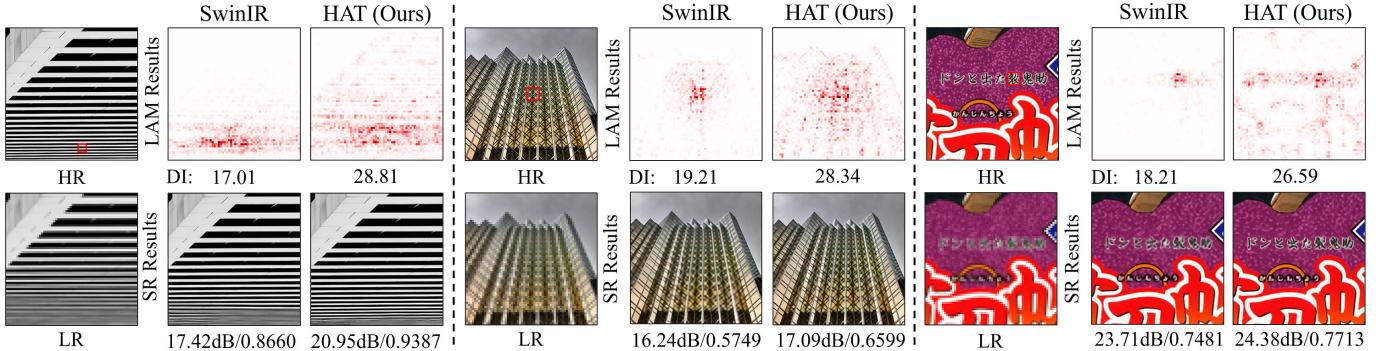


Fig. 11. Comparison of LAM results between SwinIR and HAT.

C. Real-world Image Super-Resolution

Quantitative results. Table XII presents a quantitative comparison of our method with state-of-the-art approaches: ESRGAN [12], BSRGAN [92], Real-ESRGAN [38], DASR [94], and SwinIR [22]. All models are GAN-based, and our HAT models are trained using the BSRGAN degradation model [92] (i.e., HAT-1) and Real-ESRGAN degradation model [38] (i.e., HAT-2). For compared methods, we utilize the officially released models and evaluate them on four datasets. Real-SR-cano [95] and Real-Nikon [95] feature real-world data pairs from specific cameras. AIM2019-val, from the AIM2019 challenge [96], is a synthetic dataset with realistic, unknown degradations². Additionally, we construct a synthetic dataset

using 100 DIV2K_valid [84] images with a high-order degradation model [38], following DASR [94]. We use PSNR and LPIPS [97] as metrics, with PSNR measuring fidelity and LPIPS assessing perceptual quality. Since all of the methods are GAN-based, the PSNR performance of different models on various datasets is not consistent. Nevertheless, similar PSNR values suggest comparable fidelity among the methods. Notably, HAT-1 achieves the best balance between PSNR and LPIPS among the three methods. HAT-2 obtains the best performance on all four datasets, indicating that it generates the results with the best perceptual quality.

Visual comparison. We show the visual results from different methods on real-world low-resolution images in Fig. 12. We adopt the RealSRSet+5images [92] as the test set, which is commonly used for evaluating real-world SR models. As different degradation models from BSRGAN [92] and Real-

²We use paired data of the validation set from this challenge. The competition officials do not release the degradation model.

TABLE XII

QUANTITATIVE COMPARISON WITH STATE-OF-THE-ART METHODS FOR **REAL-WORLD IMAGE SUPER-RESOLUTION** ON BENCHMARK DATASETS. THE BEST AND SECOND BEST RESULTS ARE MARKED IN **BOLD** AND UNDERLINE.

Method	RealSR-cano [95]		RealSR-Nikon [95]		AIM2019-val [96]		DIV2K-SysReal	
	PSNR (\uparrow)	LPIPS (\downarrow)	PSNR (\uparrow)	LPIPS (\downarrow)	PSNR (\uparrow)	LPIPS (\downarrow)	PSNR (\uparrow)	LPIPS (\downarrow)
ESRGAN [12]	27.67	0.412	27.46	0.425	23.16	0.550	23.48	0.627
DASR [94]	<u>27.40</u>	0.393	26.35	0.401	23.76	0.421	23.78	0.473
BSRGAN [92]	26.91	0.371	25.56	0.391	24.20	0.400	<u>23.83</u>	0.469
SwinIR [22]	26.64	<u>0.357</u>	25.76	<u>0.364</u>	23.89	0.387	23.31	0.449
HAT-1 (ours)	27.17	0.360	<u>26.52</u>	0.376	24.09	<u>0.380</u>	23.62	<u>0.439</u>
Real-ESRGAN [38]	26.14	0.378	25.49	0.388	23.89	0.396	23.58	0.446
HAT-2 (ours)	26.68	0.342	25.85	0.358	<u>24.19</u>	0.370	23.98	0.423

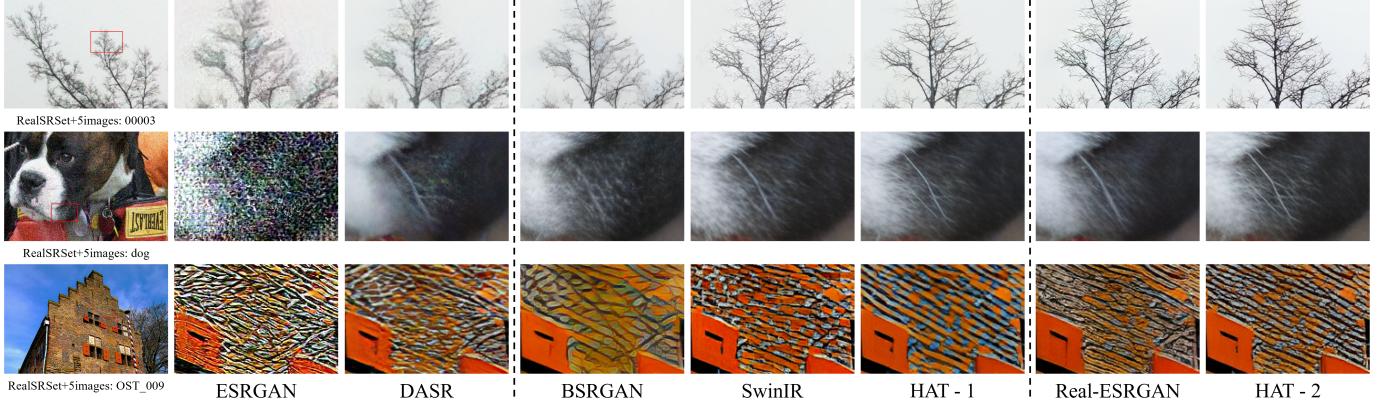


Fig. 12. Real-world image super-resolution results on SR $\times 4$. HAT-1 uses BSRGAN degradation model and HAT-2 adopts Real-ESRGAN degradation model.

ESRGAN [95] may produce varied visual properties, we present results for both models, denoted as HAT-1 and HAT-2. In the first and second rows of visual comparisons, our HAT produces much clearer branches and whiskers than other methods. In the third row, results from different degradation models exhibit significant variation. BSRGAN-based models show notable differences in color and texture. The details of BSRGAN results are relatively smooth. SwinIR generates clear textures but with obvious color deviations. In contrast, HAT achieves a relatively good balance. It handles the details well, and its larger receptive field also enhances the processing accuracy of low-frequency information. Using the degradation model from Real-ESRGAN, we can see that the result of HAT appear neater and more brick-like. Overall, our method produces visually appealing results with sharp, clear edges, demonstrating its potential for real-world applications.

D. Image Denoising

1) *Grayscale image denoising*: Tab. XIII shows the quantitative comparison on grayscale image denoising of our approach with the state-of-the-art method: DnCNN [7], IRCNN [9], FFDNet [68], RNAN [11], RDN [69], IPT [18], DRUNet [13], SwinIR [22], Restormer [20] and SCUNet [74]. Compared to SwinIR [22], HAT achieves better performance on all datasets with multiple noise levels. On Urban100, HAT achieves the largest performance gain by up to 0.64dB for $\sigma = 50$. Compared to current state-of-the-art methods Restormer [20] and SCUNet [74], HAT can still outperform the former and obtains comparable performance with the latter.

We provide the visual results of different methods in Fig. 13. For “09” in Set12, our HAT restores clear lines while other approaches suffers from severe blurs. For “test021” in BSD68 and “img_061” in Urban100, HAT reconstructs much sharper edges than other methods. For “img_076” in Urban100, the results produced by HAT have the clearest textures. Overall, HAT obtains the best visual quality among all methods.

2) *Color image denoising*: Tab. XIV shows the quantitative comparison results on color image denoising of our approach with the state-of-the-art methods: DnCNN [7], IRCNN [9], FFDNet [68], RNAN [11], RDN [69], IPT [18], DRUNet [13], SwinIR [22], Restormer [20] and SCUNet [74]. We can observe that HAT achieves the best performance on almost all four benchmark datasets. Specifically, HAT outperforms SwinIR from 0.24dB to 0.38dB and surpasses SCUNet by a large margin of 0.19dB on Urban100 with $\sigma = 15$.

We present the visual results of different methods in Fig. 14. For images “img_039”, “img_042” and “img_074” on Urban100, HAT reconstructs complete and clear edges, whereas other methods cannot produce complete lines or suffer from severe blurring effects. For the image “img_085” on Urban100, our method successfully restores the correct shape and clear texture, while other methods all fail. All these results demonstrate the superiority of the proposed HAT on image denoising.

E. JPEG compression artifacts Reduction

Tab. XV shows the quantitative comparison results on JPEG compression artifacts reduction of our approach with the state-of-the-art methods: ARCNN [6], DnCNN [7], RNAN [11],

TABLE XIII

QUANTITATIVE COMPARISON (AVERAGE PSNR) WITH STATE-OF-THE-ART METHODS FOR GRayscale IMAGE DENOISING ON BENCHMARK DATASETS. THE BEST AND SECOND BEST RESULTS ARE MARKED IN **BOLD** AND UNDERLINE, RESPECTIVELY.

Dataset	σ	DnCNN [7]	IRCNN [9]	FFDNet [68]	NLRN [39]	RNAN [11]	MWCNN [98]	DRUNet [13]	SwinIR [22]	Restormer [20]	SCUNet [74]	HAT (ours)
Set12 [7]	15	32.86	32.76	32.75	33.16	-	33.15	33.25	33.36	33.42	<u>33.43</u>	33.49
	25	30.44	30.37	30.43	30.80	-	30.79	30.94	31.01	31.08	<u>31.09</u>	31.13
	50	27.18	27.12	27.32	27.64	27.70	27.74	27.90	27.91	28.00	<u>28.04</u>	28.07
BSD68 [89]	15	31.73	31.63	31.63	31.88	-	31.86	31.91	<u>31.97</u>	31.96	31.99	31.99
	25	29.23	29.15	29.19	29.41	-	29.41	29.48	29.50	<u>29.52</u>	29.55	<u>29.52</u>
	50	26.23	26.19	26.29	26.47	26.48	26.53	26.59	26.58	<u>26.62</u>	26.67	26.60
Urban100 [90]	15	32.64	32.46	32.40	33.45	-	33.17	33.44	33.70	<u>33.79</u>	<u>33.88</u>	33.99
	25	29.95	29.80	29.90	30.94	-	30.66	31.11	31.30	31.46	<u>31.58</u>	31.67
	50	26.26	26.22	26.50	27.49	27.65	27.42	27.96	27.98	28.29	<u>28.56</u>	28.62

TABLE XIV

QUANTITATIVE COMPARISON (AVERAGE PSNR) WITH STATE-OF-THE-ART METHODS FOR COLOR IMAGE DENOISING ON BENCHMARK DATASETS. THE BEST, SECOND BEST RESULTS ARE MARKED IN **BOLD** AND UNDERLINE, RESPECTIVELY.

Dataset	σ	DnCNN [7]	IRCNN [9]	FFDNet [68]	RNAN [11]	RDN [29]	IPT [18]	DRUNet [13]	SwinIR [22]	Restormer [20]	SCUNet [74]	HAT (ours)
CBSD68 [89]	15	33.90	33.86	33.87	-	-	-	34.30	34.42	<u>34.40</u>	<u>34.40</u>	34.42
	25	31.24	31.16	31.21	-	-	-	31.69	<u>31.78</u>	31.79	31.79	31.79
	50	27.95	27.86	27.96	28.27	28.31	28.39	28.51	28.56	<u>28.60</u>	28.61	28.58
Kodak24 [99]	15	34.60	34.69	34.63	-	-	-	35.31	35.34	<u>35.47</u>	35.34	35.50
	25	32.14	32.18	32.13	-	-	-	32.89	32.89	<u>33.04</u>	32.92	33.08
	50	28.95	28.93	28.98	29.58	29.66	29.64	29.86	29.79	<u>30.01</u>	29.87	30.02
McMaster [100]	15	33.45	34.58	34.66	-	-	-	35.40	<u>35.61</u>	<u>35.61</u>	35.60	35.64
	25	31.52	32.18	32.35	-	-	-	33.14	33.20	<u>33.34</u>	<u>33.34</u>	33.37
	50	28.62	28.91	29.18	29.72	-	29.98	30.08	30.22	30.30	30.29	30.26
Urban100 [90]	15	32.98	33.78	33.83	-	-	-	34.81	35.13	35.13	<u>35.18</u>	35.37
	25	30.81	31.20	31.40	-	-	-	32.60	32.90	32.96	<u>33.03</u>	33.14
	50	27.59	27.70	28.05	29.08	29.38	29.71	29.61	29.82	30.02	<u>30.14</u>	30.20

TABLE XV

QUANTITATIVE COMPARISON (AVERAGE PSNR/SSIM) WITH STATE-OF-THE-ART METHODS FOR JPEG COMPRESSION ARTIFACTS REDUCTION ON BENCHMARK DATASETS. THE BEST AND SECOND BEST RESULTS ARE MARKED IN **BOLD** AND UNDERLINE, RESPECTIVELY.

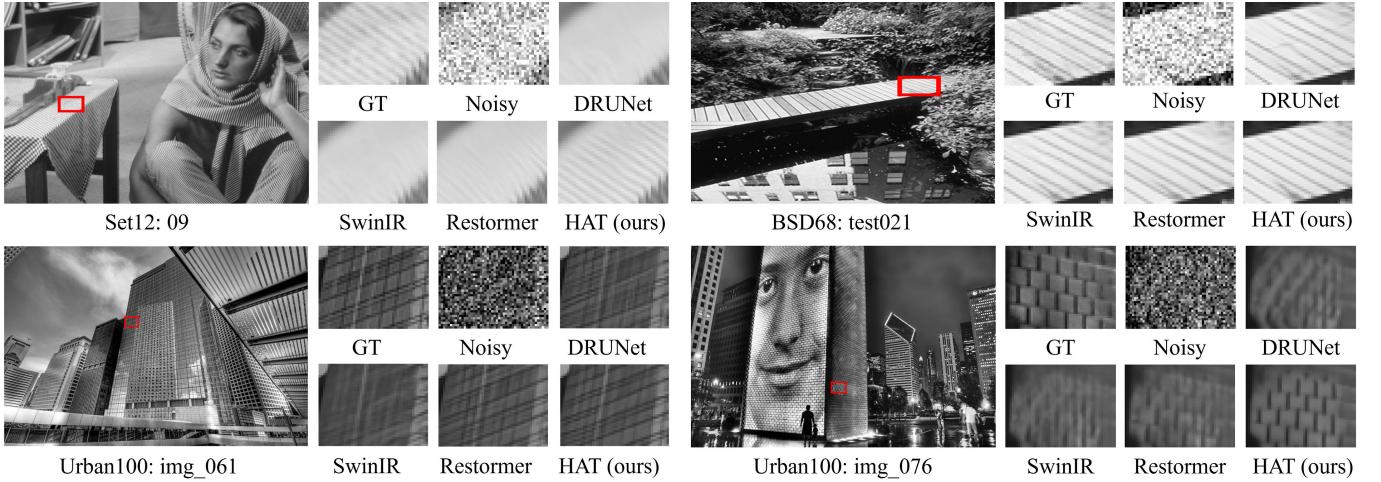
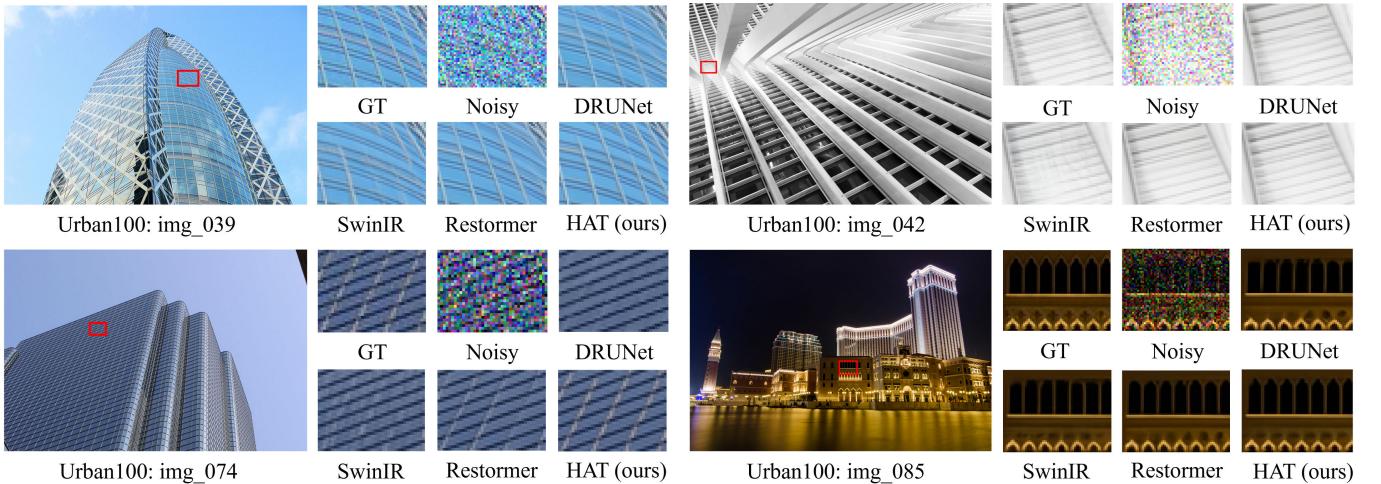
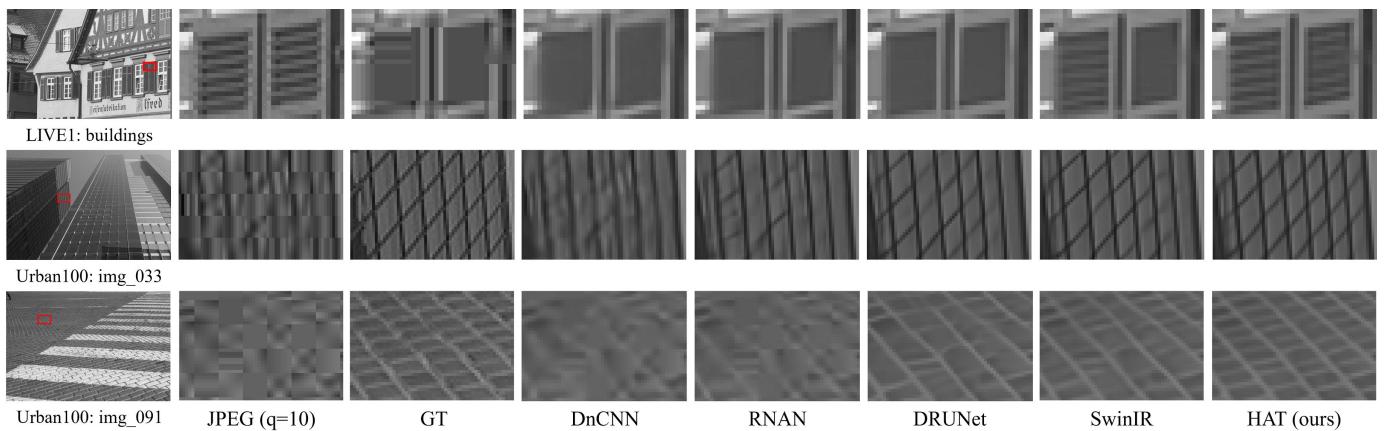
Dataset	q	ARCNN [6]	DnCNN [7]	RNAN [11]	FBCNN [101]	DRUNet [13]	SwinIR [22]	HAT (ours)
Classic5 [102]	10	29.03/0.7929	29.40/0.8026	29.96/0.8178	30.12/0.8223	<u>30.16/0.8234</u>	30.27/0.8249	<u>30.27/0.8246</u>
	20	31.15/0.8517	31.63/0.8610	32.11/0.8693	32.31/0.8724	32.39/0.8734	<u>32.52/0.8748</u>	32.54/0.8749
	30	32.51/0.8806	32.91/0.8861	33.38/0.8924	33.54/0.8943	33.59/0.8949	33.73/0.8961	<u>33.72/0.8958</u>
	40	33.32/0.8953	33.77/0.9003	34.27/0.9061	34.35/0.9070	34.41/0.9075	<u>34.52/0.9082</u>	34.58/0.9086
LIVE1 [103]	10	28.96/0.8076	29.19/0.8123	29.63/0.8239	29.75/0.8268	29.79/0.8278	29.86/0.8287	29.84/0.8283
	20	31.29/0.8733	31.59/0.8802	32.03/0.8877	32.13/0.8893	32.17/0.8899	<u>32.25/0.8909</u>	32.26/0.8907
	30	32.67/0.9043	32.98/0.9090	33.45/0.9149	33.54/0.9161	33.59/0.9166	33.69/0.9174	33.66/0.9170
	40	33.63/0.9198	33.96/0.9247	34.47/0.9299	34.53/0.9307	34.58/0.9312	<u>34.67/0.9317</u>	34.69/0.9317
Urban100 [90]	10	-	28.54/0.8487	29.76/0.8723	30.15/0.8795	30.05/0.8772	<u>30.55/0.8842</u>	30.60/0.8845
	20	-	31.01/0.9022	32.33/0.9179	32.66/0.9219	32.66/0.9216	<u>33.12/0.9254</u>	33.22/0.9260
	30	-	32.47/0.9248	33.83/0.9365	34.09/0.9392	34.13/0.9392	34.58/0.9417	34.58/0.9419
	40	-	33.49/0.9376	34.95/0.9476	35.08/0.9490	35.11/0.9491	<u>35.50/0.9509</u>	35.68/0.9517

DRUNet [13], FBCNN [101], SwinIR [22]. On Classic5 and LIVE1, HAT only achieves comparable performance to SwinIR. We consider this is because the demand for the model fitting ability of this task has approached saturation, particularly for low-resolution images. We further provide the performance comparison on Urban100. Then we can see that HAT achieves considerable performance gains over SwinIR, up to 0.18dB for JPEG quality $q = 40$. This can be attributed to the presence of a large number of regular textures and repeating patterns in the images of the Urban100 dataset. HAT is capable of activating more pixels for restoration. With a larger receptive field, it can restore sharper edges and textures.

We also provide the visual results of different approaches in Fig. 15. For the image ‘buildings’ in LIVE1, HAT obtains much clearer textures than other methods. For the self-repeated textures that appear in the images ‘img_033’ and ‘img_091’ of Urban100, our HAT successfully restores the correct results. All the quantitative and visual results demonstrate the superiority of our method on compression artifacts reduction.

VII. CONCLUSION

In this work, we propose a new Hybrid Attention Transformer, HAT, for image restoration. Our model combines channel attention and self-attention to activate more pixels for

Fig. 13. Grayscale image denoising results with noise level $\sigma = 50$.Fig. 14. Color image denoising results with noise level $\sigma = 50$.Fig. 15. Image compression artifacts reduction results with JPEG quality $q = 10$.

high-resolution reconstruction. Besides, we propose an overlapping cross-attention module to enhance the cross-window interaction. Moreover, we introduce a same-task pre-training

strategy for image super-resolution. Extensive benchmark and real-world evaluations demonstrate that HAT outperforms the state-of-the-art methods for several image restoration tasks.

REFERENCES

- [1] N. Fatima, "Ai in photography: Scrutinizing implementation of super-resolution techniques in photo-editors," in *2020 35th International Conference on Image and Vision Computing New Zealand (IVCNZ)*. IEEE, 2020, pp. 1–6.
- [2] W. W. Zou and P. C. Yuen, "Very low resolution face recognition problem," *IEEE Transactions on image processing*, vol. 21, no. 1, pp. 327–340, 2011.
- [3] W. Shi, J. Caballero, C. Ledig, X. Zhuang, W. Bai, K. Bhatia, A. M. S. M. de Marvao, T. Dawes, D. O'Regan, and D. Rueckert, "Cardiac image super-resolution with global correspondence using multi-atlas patchmatch," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2013: 16th International Conference, Nagoya, Japan, September 22–26, 2013, Proceedings, Part III 16*. Springer, 2013, pp. 9–16.
- [4] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.
- [5] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *European conference on computer vision*. Springer, 2014, pp. 184–199.
- [6] C. Dong, Y. Deng, C. C. Loy, and X. Tang, "Compression artifacts reduction by a deep convolutional network," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 576–584.
- [7] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE transactions on image processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [8] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *European conference on computer vision*. Springer, 2016, pp. 391–407.
- [9] 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.
- [10] 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 (ECCV)*, 2018, pp. 286–301.
- [11] 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*, 2018.
- [12] 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 (ECCV) workshops*, 2018, pp. 0–0.
- [13] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. Van Gool, and R. Timofte, "Plug-and-play image restoration with deep denoiser prior," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6360–6376, 2021.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [15] 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*, 2020.
- [16] 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/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [17] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 568–578.
- [18] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12299–12310.
- [19] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17683–17693.
- [20] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5728–5739.
- [21] W. Li, X. Lu, J. Lu, X. Zhang, and J. Jia, "On efficient transformer and image pre-training for low-level vision," 2021.
- [22] 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/CVF International Conference on Computer Vision*, 2021, pp. 1833–1844.
- [23] J. Gu and C. Dong, "Interpreting super-resolution networks with local attribution maps," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9199–9208.
- [24] X. Chen, X. Wang, J. Zhou, Y. Qiao, and C. Dong, "Activating more pixels in image super-resolution transformer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22367–22377.
- [25] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 295–307, 2015.
- [26] 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.
- [27] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1646–1654.
- [28] 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 conference on computer vision and pattern recognition workshops*, 2017, pp. 136–144.
- [29] 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.
- [30] 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/CVF conference on computer vision and pattern recognition*, 2019, pp. 11065–11074.
- [31] B. Niu, W. Wen, W. Ren, X. Zhang, L. Yang, S. Wang, K. Zhang, X. Cao, and H. Shen, "Single image super-resolution via a holistic attention network," in *European conference on computer vision*. Springer, 2020, pp. 191–207.
- [32] Y. Mei, Y. Fan, and Y. Zhou, "Image super-resolution with non-local sparse attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3517–3526.
- [33] 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.
- [34] J. Kim, J. K. Lee, and K. M. Lee, "Deeply-recursive convolutional network for image super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1637–1645.
- [35] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3147–3155.
- [36] S. Zhou, J. Zhang, W. Zuo, and C. C. Loy, "Cross-scale internal graph neural network for image super-resolution," *Advances in neural information processing systems*, vol. 33, pp. 3499–3509, 2020.
- [37] W. Zhang, Y. Liu, C. Dong, and Y. Qiao, "Ranksrgan: Generative adversarial networks with ranker for image super-resolution," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3096–3105.
- [38] X. Wang, L. Xie, C. Dong, and Y. Shan, "Real-esrgan: Training real-world blind super-resolution with pure synthetic data," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1905–1914.
- [39] D. Liu, B. Wen, Y. Fan, C. C. Loy, and T. S. Huang, "Non-local recurrent network for image restoration," *Advances in neural information processing systems*, vol. 31, 2018.
- [40] Y. Liu, A. Liu, J. Gu, Z. Zhang, W. Wu, Y. Qiao, and C. Dong, "Discovering semantics" in super-resolution networks," 2021.

- [41] L. Xie, X. Wang, C. Dong, Z. Qi, and Y. Shan, "Finding discriminative filters for specific degradations in blind super-resolution," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [42] X. Kong, X. Liu, J. Gu, Y. Qiao, and C. Dong, "Reflash dropout in image super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6002–6012.
- [43] Y. Liu, H. Zhao, J. Gu, Y. Qiao, and C. Dong, "Evaluating the generalization ability of super-resolution networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [44] J. Hu, J. Gu, S. Yu, F. Yu, Z. Li, Z. You, C. Lu, and C. Dong, "Interpreting low-level vision models with causal effect maps," *arXiv preprint arXiv:2407.19789*, 2024.
- [45] Y. Li, K. Zhang, J. Cao, R. Timofte, and L. Van Gool, "Localvit: Bringing locality to vision transformers," 2021.
- [46] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, and L. Zhang, "Cvt: Introducing convolutions to vision transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 22–31.
- [47] Z. Huang, Y. Ben, G. Luo, P. Cheng, G. Yu, and B. Fu, "Shuffle transformer: Rethinking spatial shuffle for vision transformer," 2021.
- [48] X. Dong, J. Bao, D. Chen, W. Zhang, N. Yu, L. Yuan, D. Chen, and B. Guo, "Cswin transformer: A general vision transformer backbone with cross-shaped windows," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 124–12 134.
- [49] 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, 2021.
- [50] S. Wu, T. Wu, H. Tan, and G. Guo, "Pale transformer: A general vision transformer backbone with pale-shaped attention," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 2731–2739.
- [51] K. Yuan, S. Guo, Z. Liu, A. Zhou, F. Yu, and W. Wu, "Incorporating convolution designs into visual transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 579–588.
- [52] K. Li, Y. Wang, J. Zhang, P. Gao, G. Song, Y. Liu, H. Li, and Y. Qiao, "Uniformer: Unifying convolution and self-attention for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [53] K. Patel, A. M. Bur, F. Li, and G. Wang, "Aggregating global features into local vision transformer," 2022.
- [54] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens, "Stand-alone self-attention in vision models," *Advances in neural information processing systems*, vol. 32, 2019.
- [55] A. Vaswani, P. Ramachandran, A. Srinivas, N. Parmar, B. Hechtman, and J. Shlens, "Scaling local self-attention for parameter efficient visual backbones," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 894–12 904.
- [56] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *International journal of computer vision*, vol. 128, no. 2, pp. 261–318, 2020.
- [57] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training data-efficient image transformers & distillation through attention," in *International Conference on Machine Learning*. PMLR, 2021, pp. 10 347–10 357.
- [58] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.
- [59] B. Wu, C. Xu, X. Dai, A. Wan, P. Zhang, Z. Yan, M. Tomizuka, J. Gonzalez, K. Keutzer, and P. Vajda, "Visual transformers: Token-based image representation and processing for computer vision," *arXiv preprint arXiv:2006.03677*, 2020.
- [60] G. Huang, Y. Wang, K. Lv, H. Jiang, W. Huang, P. Qi, and S. Song, "Glance and focus networks for dynamic visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4605–4621, 2022.
- [61] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," in *European conference on computer vision*. Springer, 2022, pp. 205–218.
- [62] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy, "Do vision transformers see like convolutional neural networks?" *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [63] T. Xiao, P. Dollar, M. Singh, E. Mintun, T. Darrell, and R. Girshick, "Early convolutions help transformers see better," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [64] Y. Yuan, R. Fu, L. Huang, W. Lin, C. Zhang, X. Chen, and J. Wang, "Hrformer: High-resolution vision transformer for dense predict," *Advances in Neural Information Processing Systems*, vol. 34, pp. 7281–7293, 2021.
- [65] J. Cao, Y. Li, K. Zhang, J. Liang, and L. Van Gool, "Video super-resolution transformer," *arXiv preprint arXiv:2106.06847*, 2021.
- [66] Z. Tu, H. Talebi, H. Zhang, F. Yang, P. Milanfar, A. Bovik, and Y. Li, "Maxim: Multi-axis mlp for image processing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5769–5780.
- [67] J. Liang, J. Cao, Y. Fan, K. Zhang, R. Ranjan, Y. Li, R. Timofte, and L. Van Gool, "Vrt: A video restoration transformer," *arXiv preprint arXiv:2201.12288*, 2022.
- [68] K. Zhang, W. Zuo, and L. Zhang, "Ffdnet: Toward a fast and flexible solution for cnn-based image denoising," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4608–4622, 2018.
- [69] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image restoration," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 7, pp. 2480–2495, 2020.
- [70] S. Nah, T. Hyun Kim, and K. Mu Lee, "Deep multi-scale convolutional neural network for dynamic scene deblurring," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3883–3891.
- [71] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, "Multi-stage progressive image restoration," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 14 821–14 831.
- [72] L. Chen, X. Chu, X. Zhang, and J. Sun, "Simple baselines for image restoration," in *European Conference on Computer Vision*. Springer, 2022, pp. 17–33.
- [73] L. Liu, L. Xie, X. Zhang, S. Yuan, X. Chen, W. Zhou, H. Li, and Q. Tian, "Tape: Task-agnostic prior embedding for image restoration," in *European Conference on Computer Vision*. Springer, 2022, pp. 447–464.
- [74] K. Zhang, Y. Li, J. Liang, J. Cao, Y. Zhang, H. Tang, D.-P. Fan, R. Timofte, and L. V. Gool, "Practical blind image denoising via swin-conv-unet and data synthesis," *Machine Intelligence Research*, vol. 20, no. 6, pp. 822–836, 2023.
- [75] S. Anwar, S. Khan, and N. Barnes, "A deep journey into super-resolution: A survey," *ACM Computing Surveys (CSUR)*, vol. 53, no. 3, pp. 1–34, 2020.
- [76] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International Conference on Machine Learning*, 2017, pp. 3319–3328.
- [77] G. H. Granlund, "In search of a general picture processing operator," *Computer Graphics and Image Processing*, vol. 8, no. 2, pp. 155–173, 1978.
- [78] S. Yitzhaki, "Relative deprivation and the gini coefficient," *The quarterly journal of economics*, pp. 321–324, 1979.
- [79] Y. Zhao, G. Wang, C. Tang, C. Luo, W. Zeng, and Z.-J. Zha, "A battle of network structures: An empirical study of cnn, transformer, and mlp," *arXiv preprint arXiv:2108.13002*, 2021.
- [80] D. Hendrycks and K. Gimpel, "Gaussian error linear units (gelus)," *arXiv preprint arXiv:1606.08415*, 2016.
- [81] H. Bao, L. Dong, S. Piao, and F. Wei, "Beit: Bert pre-training of image transformers," in *International Conference on Learning Representations*, 2021.
- [82] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 000–16 009.
- [83] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, 2009, pp. 248–255.
- [84] 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 conference on computer vision and pattern recognition workshops*, 2017, pp. 136–144.
- [85] R. Timofte, E. Agustsson, L. Van Gool, M.-H. Yang, and L. Zhang, "Ntire 2017 challenge on single image super-resolution: Methods and results," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 114–125.

- [86] Z. Lin, P. Garg, A. Banerjee, S. A. Magid, D. Sun, Y. Zhang, L. Van Gool, D. Wei, and H. Pfister, “Revisiting rcan: Improved training for image super-resolution,” *arXiv preprint arXiv:2201.11279*, 2022.
- [87] M. Bevilacqua, A. Roumy, C. Guillemot, and M.-L. A. Morel, “Low-complexity single-image super-resolution based on nonnegative neighbor embedding,” in *British Machine Vision Conference (BMVC)*, 2012.
- [88] R. Zeyde, M. Elad, and M. Protter, “On single image scale-up using sparse-representations,” in *International conference on curves and surfaces*. Springer, 2010, pp. 711–730.
- [89] 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 Eighth IEEE International Conference on Computer Vision. ICCV 2001*, vol. 2. IEEE, 2001, pp. 416–423.
- [90] J.-B. 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.
- [91] 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.
- [92] K. Zhang, J. Liang, L. Van Gool, and R. Timofte, “Designing a practical degradation model for deep blind image super-resolution,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4791–4800.
- [93] K. Ma, Z. Duanmu, Q. Wu, Z. Wang, H. Yong, H. Li, and L. Zhang, “Waterloo exploration database: New challenges for image quality assessment models,” *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 1004–1016, 2016.
- [94] J. Liang, H. Zeng, and L. Zhang, “Efficient and degradation-adaptive network for real-world image super-resolution,” in *European Conference on Computer Vision*. Springer, 2022, pp. 574–591.
- [95] J. Cai, H. Zeng, H. Yong, Z. Cao, and L. Zhang, “Toward real-world single image super-resolution: A new benchmark and a new model,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3086–3095.
- [96] A. Lugmayr, M. Danelljan, R. Timofte, M. Fritzsche, S. Gu, K. Purohit, P. Kandula, M. Suin, A. Rajagoapalan, N. H. Joon *et al.*, “Aim 2019 challenge on real-world image super-resolution: Methods and results,” in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2019, pp. 3575–3583.
- [97] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [98] P. Liu, H. Zhang, K. Zhang, L. Lin, and W. Zuo, “Multi-level wavelet-cnn for image restoration,” in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018, pp. 773–782.
- [99] R. Franzen, “Kodak lossless true color image suite,” *source: <http://r0k.us/graphics/kodak>*, vol. 4, no. 2, 1999.
- [100] L. Zhang, X. Wu, A. Buades, and X. Li, “Color demosaicking by local directional interpolation and nonlocal adaptive thresholding,” *Journal of Electronic Imaging*, vol. 20, no. 2, pp. 023 016–023 016, 2011.
- [101] J. Jiang, K. Zhang, and R. Timofte, “Towards flexible blind jpeg artifacts removal,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4997–5006.
- [102] A. Foi, V. Katkovnik, and K. Egiazarian, “Pointwise shape-adaptive dct for high-quality denoising and deblocking of grayscale and color images,” *IEEE Transactions on Image Processing*, vol. 16, no. 5, pp. 1395–1411, 2007.
- [103] H. Sheikh, “Live image quality assessment database release 2,” *<http://live.ece.utexas.edu/research/quality>*, 2005.