# Activating More Pixels in Image Super-Resolution Transformer

Xiangyu Chen[1,2], Xintao Wang[3], Jiantao Zhou[1], and Chao Dong[2,4]

[1]University of Macau [2]Shenzhen Institute of Advanced Technology,
Chinese Academy of Sciences [3]ARC Lab, Tencent PCG [4]Shanghai AI Laboratory
{chxy95, xintao.alpha}@gmail.com, jtzhou@um.edu.mo, chao.dong@siat.ac.cn

**Abstract.** Transformer-based methods have shown impressive performance in low-level vision tasks, such as image super-resolution. 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 reconstruction, we propose a novel Hybrid Attention Transformer (HAT). It combines channel attention and 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 propose a same-task pre-training strategy to bring further improvement. Extensive experiments show the effectiveness of the proposed modules, and the overall method significantly outperforms the state-of-the-art methods by more than **1dB**. Codes and models will be available at https://github.com/chxy95/HAT.

## 1   Introduction

Single image super-resolution (SR) is a classic problem in computer vision and image processing. It aims to reconstruct a high-resolution image from a given low-resolution input. Since deep learning has been successfully applied to the SR task [9], numerous methods based on the convolutional neural network (CNN) have been proposed [10,11,19,30,65,67,7,41] and almost dominate this field in the past few years. Recently, due to the success in natural language processing, Transformer [50] has attracted the attention of the computer vision community. After making rapid progress on high-level vision tasks [13,36,51], Transformer-based methods are also developed for low-level vision tasks [5,54,62], as well as for SR [26,29]. Especially, a newly designed network, SwinIR [29], obtains a breakthrough improvement in this task.

Despite the success, "why Transformer is better than CNN" remains a mystery. An intuitive explanation is that this kind of network can benefit from the self-attention mechanism and utilize long-range information. However, we employ the attribution analysis method LAM [14] to examine the involved range of utilized information for reconstruction in SwinIR. Interestingly, we find that
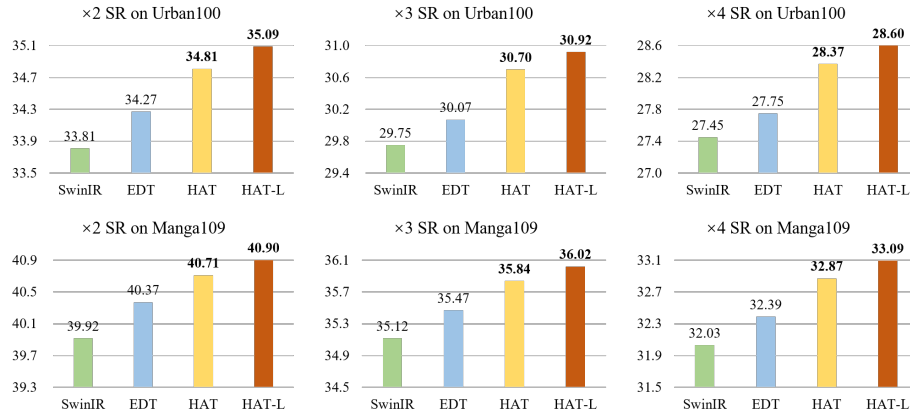
Fig. 1: Performance comparison of the proposed HAT with the state-of-the-art methods SwinIR [29] and EDT [26]. HAT-L represents a larger variant of HAT. Our approach can surpass the state-of-the-art methods by 0.3dB∼1.2dB.

SwinIR does NOT exploit more input pixels than CNN-based methods (*e.g.*, RCAN [65]) in super-resolution, as shown in Fig. 2(a). Besides, although SwinIR obtains higher quantitative performance, it produces inferior results to RCAN in some cases, due to the limited range of utilized information. These phenomena illustrate that Transformer has a stronger ability to model local information, but the range of its utilized information needs to be expanded.

To address the above-mentioned problem and further develop the potential of Transformer for SR, we propose a Hybrid Attention Transformer, namely HAT. Our HAT 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, to better aggregate the cross-window information, we also introduce an overlapping cross-attention module. Motivated by [5,26], we additionally explore the effect of pre-training on the SR task and provide a *same-task pre-training* strategy. Experimental results show that this strategy can perform better than multi-related-task pre-training [26]. Equipped with the above designs and improvements, our approach can surpass the state-of-the-art methods by a huge margin(0.3dB∼1.2dB), as shown in Fig. 1.

Overall, our contributions include three aspects: 1) We introduce channel attention to Transformer to utilize more input information. 2) We propose an overlapping cross-attention module to better aggregate the cross-window information. 3) We provide a same-task pre-training strategy to further activate the potential of the proposed network.

## 2   Related Work

### 2.1   Deep Networks for Image SR

Since SRCNN [9] first introduces deep convolution neural networks (CNNs) to the image SR task and obtains superior performance over conventional SR meth-

ods, numerous deep networks [10,11,45,19,23,30,67,65,7,41,40,29,26] 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 [24,30] and dense block [53,67], to enhance the model representation ability. Several works explore more different frameworks like recursive neural network [20,46] and graph neural network [69]. To improve perceptual quality, [24,53,64,52] introduce adversarial learning to generate more realistic results. By using attention mechanism, [65,7,66,33,41,40] achieve further improvement in terms of reconstruction fidelity. Recently, a series of Transformer-based networks [5,29,26] are proposed and constantly refresh the state-of-the-art of SR task, showing the powerful representation ability of Transformer.

To better understand the working mechanisms of SR networks, several works [14,35,59,22] are proposed to analyze and interpret the SR networks. LAM [14] adopts the integral gradient method to explore which input pixels contribute most to the final performance. DDR [35] reveals the deep semantic representations in SR networks based on deep feature dimensionality reduction and visualization. FAIG [59] is proposed to find discriminative filters for specific degradations in blind SR. [22] introduces channel saliency map to demonstrate that Dropout can help prevent co-adapting for real SR networks. In this work, we exploit LAM [14] to analyse and understand the behavior of SR networks.
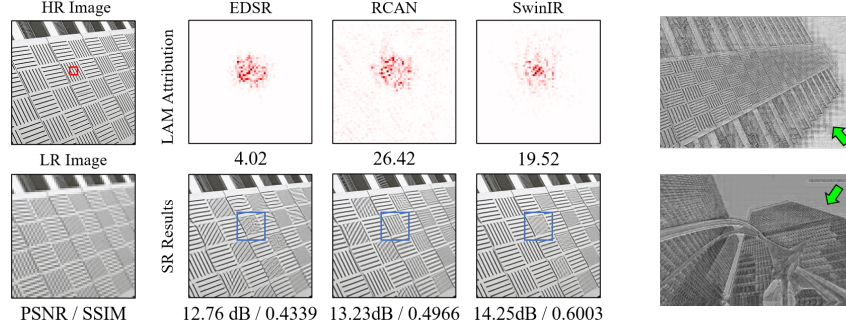
### 2.2   Vision Transformer

Recently, Transformer [50] 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 [13,27,56,51,36,18,12,6,57,60,25,42] have been developed for high-level vision tasks, including image classification [36,13,27,44,49], object detection [34,48,36,4,6], segmentation [55,51,16,2], *etc.* Although vision Transformer has shown its superiority on modeling long-range dependency [13,43], there are still many works demonstrating that the convolution can help Transformer achieve better visual representation [56,58,61,60,25]. Due to the impressive performance, Transformer has also been introduced for low-level vision tasks [5,54,37,29,3,62,28,26]. Specifically, [5] develops a ViT-style network and introduces multi-task pre-training for image processing. SwinIR [29] proposes an image restoration Transformer based on [36]. [3,28] introduce Transformer-based networks to video restoration. [26] 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 information for better reconstruction.

## 3   Method

### 3.1   Motivation

Swin Transformer [36] has already demonstrated excellent performance in image super-resolution [29]. Then we are eager to know what makes it work better

| HR Image | EDSR | RCAN | SwinIR | |
| LAM Attribution | | | | |
| LR Image | 4.02 | 26.42 | 19.52 | |
| SR Results | | | | |
| PSNR / SSIM | 12.76 dB / 0.4339 | 13.23dB / 0.4966 | 14.25dB / 0.6003 | |

(a) LAM [14] results for different networks.       (b) Blocking artifacts.

Fig. 2: **(a)** The LAM attribution reflects the importance of each pixel in the input LR image when reconstructing the patch marked with a red box. DI [14] values are provided below the LAM results. (DI reflects the range of involved pixels. A higher DI represents a wider range of utilized pixels.) The results indicate that SwinIR utilize less information compared to RCAN. **(b)** The blocking artifacts in the intermediate features of SwinIR. The top is a feature map after the 1st RSTB in SwinIR [29] and the bottom is after the 3rd RSTB.

than CNN-based methods. To reveal its working mechanisms, we resort to a diagnostic tool – LAM [14], which is an attribution method designed for SR. With LAM, we could tell which input pixels contribute most to the selected region. As shown in Fig. 2(a), the red marked points are informative pixels that contribute to the reconstruction. Intuitively, the more information is utilized, the better performance can be obtained. This is true for CNN-based methods, as comparing RCAN [65] and EDSR [30]. However, for the Transformer-based method – SwinIR, its LAM does not show a larger range than RCAN. This is in contradiction with our common sense, but could also provide us with additional insights. First, it implies that SwinIR has a much stronger mapping capability than CNN, and thus could use less information to achieve better performance. Second, SwinIR still has improvement space if it could exploit more input pixels. As depicted in Fig. 2(a), the reconstructed pattern that is marked in blue box by SwinIR is inferior to RCAN. The channel attention scheme helps RCAN see more pixels, which may also be beneficial for Transformer.

Besides, we can observe obvious blocking artifacts in the intermediate feature maps of SwinIR, as presented in Fig. 2(b). These artifacts are caused by the window partition mechanism, and this phenomenon suggests that the shifted window mechanism is inefficient to build the cross-window connection. Some works for high-level vision tasks [12,18,57,42] also point out that enhancing the connection among windows can improve the window-based self-attention methods. Based on the above two points, we investigate channel attention in the Transformer-based model and propose an overlapping cross-attention module to better aggregate cross-window information for the window-based SR Transformer.
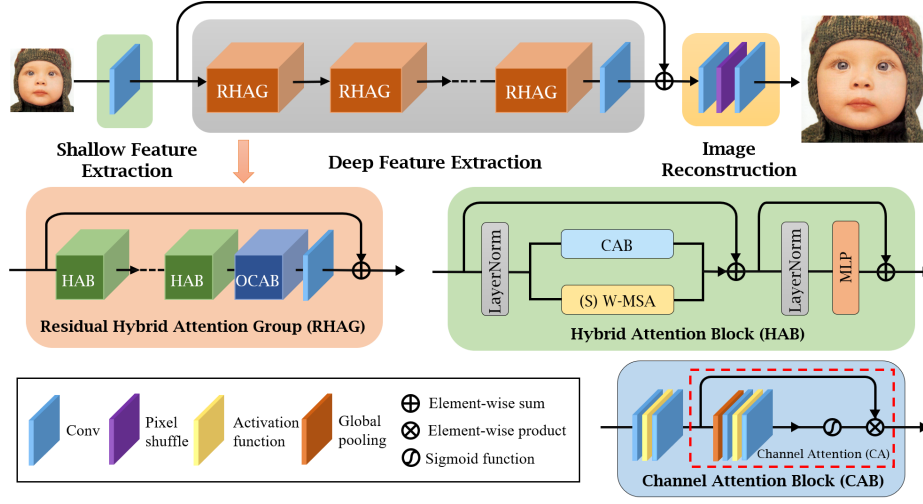
Fig. 3: The overall architecture of HAT and the structure of RHAG and HAB.

## 3.2   Network Architecture

**The Overall Structure.** As shown in Fig. 3, the overall network consists of three parts, including shallow feature extraction, deep feature extraction and image reconstruction. Specifically, for a given low-resolution (LR) input $I_{LR} \in \mathbb{R}^{H \times W \times C_{in}}$, we first use one convolutional layer $H_{SF}(\cdot)$ to extract the shallow feature $F_0 \in \mathbb{R}^{H \times W \times C}$ as

$$F_0 = H_{SF}(I_{LR}), \tag{1}$$

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 convolutional layer can help learn better visual representation [25] and lead to stable optimization [58]. 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

$$F_{DF} = H_{DF}(F_0), \tag{2}$$

where $H_{DF}(\cdot)$ consists of $N$ residual hybrid attention groups (RHAG) and one $3 \times 3$ convolutional layer $H_{Conv}(\cdot)$, which can progressively process the intermediate features as

$$F_i = H_{RHAG_i}(F_{i-1}), i = 1, 2, ..., N,$$
$$F_{DF} = H_{Conv}(F_N), \tag{3}$$

where $H_{RHAG_i}(\cdot)$ represents the $i$-th RHAG. Following [29], we also introduce a convolutional 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-resolution result via a reconstruction module as

$$I_{SR} = H_{Rec}(F_0 + F_D F), \tag{4}$$

where $H_{Rec}(\cdot)$ denotes the reconstruction module. Specifically, we adopt the pixel-shuffle method [45] to up-sample the fused feature. We simply use $L_1$ loss to optimize the parameters.

**Residual Hybrid Attention Group (RHAG).** As depicted in Fig. 3, each RHAG contains $M$ hybrid attention blocks (HAB), an overlapping cross-attention block (OCAB) and a $3\times3$ convolutional layer. To be specific, for the $i$-th RHAG, it can be formulated as

$$
\begin{aligned}
F_{i-1,0} &= F_{i-1}, \\
F_{i-1,j} &= H_{HAB_{i,j}}(F_{i-1,j-1}), j = 1, 2, ..., M, \\
F_i &= H_{Conv_i}(H_{OCAB_i}(F_{i-1,M})) + F_{i-1},
\end{aligned}
\tag{5}
$$

where $F_{i-1,0}$ indicates the input feature of $i$-th RHAG and $F_{i-1,j}$ represents the $j$-th output feature of $j$-th HAB in $i$-th RHAG. After the mapping of a series of HABs, we insert an OCAB to enlarge the receptive field for the window-based self-attention and better aggregate cross-window information. At the end of RHAG, we reserve the convolutional layer following [29]. A residual connection is also added to stabilize the training process [65].

**Hybrid Attention Block (HAB).** As shown in Fig. 2(a), more pixels are activated when channel attention is adopted, since global information is involved to calculate the channel attention weights. Besides, many works illustrate that convolution can help Transformer get better visual representation or achieve easier optimization [56,58,60,25,68]. Therefore, we incorporate a channel attention-based convolution block into the standard Transformer block to further enhance the representation ability of the network. As demonstrated in Fig. 3, a channel attention block (CAB) is inserted into the standard Swin Transformer block after the first LayerNorm (LN) layer in parallel with the window-based multi-head self-attention (W-MSA) module. Note that shifted window-based self-attention (SW-MSA) is adopted at intervals in consecutive HABs similar to [36,29]. To avoid the possible conflict of CAB and MSA on optimization and visual representation, a small constant $\alpha$ is multiplied to the output of CAB. For a given input feature $X$, the whole process of HAB is computed as

$$
\begin{aligned}
X_N &= \text{LN}(X), \\
X_M &= \text{(S)W-MSA}(X_N) + \alpha\text{CAB}(X_N) + X, \\
Y &= \text{MLP}(\text{LN}(X_M)) + X_M,
\end{aligned}
\tag{6}
$$

where $X_N$ and $X_M$ denote the intermediate features. $Y$ represents the output of HAB. Especially, we treat each pixel as a token for embedding (*i.e.*, set patch size as 1 for patch embedding following [29]). MLP denotes a multi-layer perceptron.

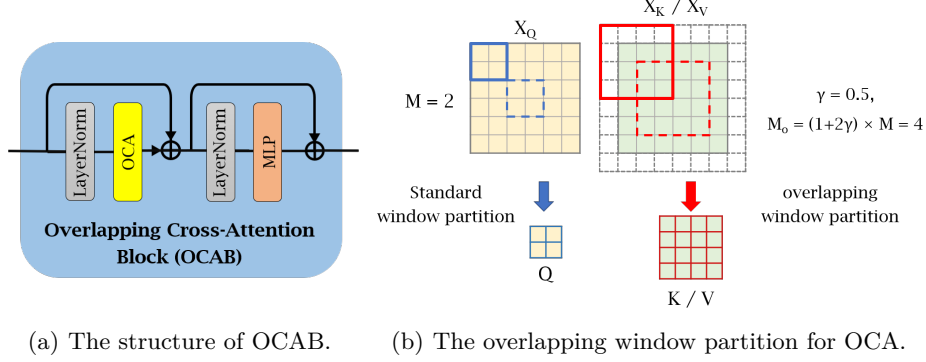(a) The structure of OCAB.        (b) The overlapping window partition for OCA.

Fig. 4: The structure of OCAB is similar to the standard Swin Transformer block. The difference is that self-attention in OCAB is calculated based on the overlapping window partition, which generates *key/value* from a larger cross-window feature than *query*.

For specific 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

$$\text{Attention}(Q, K, V) = \text{SoftMax}(QK^T/\sqrt{d} + B)V, \tag{7}$$

where $d$ represents the dimension of *query/key*. $B$ denotes the relative position encoding and is calculated as [50]. Besides, to build the connections between neighboring non-overlapping windows, we also utilize the shifted window partitioning approach [36] and set the shift size to half of the window size.

A CAB consists of two standard convolution layers with a GELU activation function [15] between them and a channel attention (CA) module, as shown in Fig. 3. Since the Transformer-based structure often requires a large number of channels for token embedding, directly using convolutions with constant width incurs a large computation cost. Thus, we compress the channel number by a constant $\beta$ between the two convolutional layers. For an input feature with $C$ channels, the channel number of the output feature after the first convolutional layer is squeezed to $\frac{C}{\beta}$, then the feature is expanded to $C$ channels through the second layer. Next, a standard CA module [65] is exploited to adaptively rescale channel-wise features. The whole process is formulated as

$$X_{out} = \text{CA}(Conv_2(\text{GELU}(Conv_1(X_{in})))), \tag{8}$$

where $X_{in}$, $X_{out}$, $Conv_1$, $Conv_2$ indicate the input feature, the output feature, the first convolutional layer and the second convolutional layer, respectively.

**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 [36]. But for OCA, as depicted in Fig. 4, we use different window sizes to partition the projected features. Concretely, 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

$$M_o = (1 + 2\gamma) \times M, \tag{9}$$

where $\gamma$ 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 $\gamma M$ is used to ensure the size consistency of overlapping windows. The attention matrix is calculated as Equ. 7, 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 [42] performs similar overlapping window partition, our OCA is different from MOA, since MOA calculates global attention using window features as tokens while OCA computes cross-attention inside each window feature using pixel token.

### 3.3   Pre-training on ImageNet

Recent works [5,26] demonstrate that pre-training plays an important role in low-level tasks. IPT [5] emphasizes the use of various low-level tasks, such as denoising, deraining, super-resolution and *etc.*, while EDT [26] utilizes different degradation levels of a specific task to do pre-training. These works aim to explore 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 [8]) based on the same task. For example, when we want to train a model for ×4 SR, we first train a ×4 SR model on ImageNet, then fine-tune it on the specific dataset, such as DF2K. As presented in Sec. 4.5, 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 believe 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.

# 4    Experiments

## 4.1    Experimental Setup

We use DF2K dataset (DIV2K [31]+Flicker2K [47]) as the original training dataset, following the latest publications [29,32]. When utilizing pre-training, we adopt ImageNet [8] following [5,26]. For the structure of HAT, the RHAG number and HAB number are both 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 ($\alpha$), the squeeze factor between two convolution layers in CAB ($\beta$), and the overlapping ratio of OCA ($\gamma$) as 0.01, 3 and 0.5, respectively. For the larger variant HAT-L, we directly double the depth of HAT by increasing the number of RHAG from 6 to 12. To evaluate the quantitative performance, PSNR and SSIM (calculated on the Y channel) are reported. More training details can refer to the appendix.

## 4.2    Effects of different window sizes

As discussed in Sec. 3.1, activating more input pixels for SR tends to achieve better performance. Enlarging window size for the calculation of self-attention is an intuitive way. In [26], the authors investigate the effects of different window sizes. However, they conduct experiments based on the shifted cross local attention and only explore the window size up to 12×12. We further explore how the window size of self-attention influences the representation ability. To eliminate the influence of our newly-introduced blocks, we conduct the following experiments directly on SwinIR. As shown in Tab. 1, the model with a large window size of 16 obtains better performance, especially on the Urban100 [17]. We also provide the qualitative comparison as depicted in Fig. 5. For the marked patch in the red box, the model with window size of 16 utilizes much more input pixels than the model with window size of 8. The quantitative performance of the reconstructed results also demonstrates the effectiveness of large window size. Based on this conclusion, we directly use window size 16 as our default setting.

Table 1: Quantitative results of models with different window sizes for ×4 SR.

| Window size | Set5 | | Set14 | | BSD100 | | Urban100 | | Manga109 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| (8,8) | 32.88 | 0.9033 | 29.09 | 0.7946 | 27.92 | 0.7489 | 27.45 | 0.8254 | 32.03 | 0.9260 |
| (16,16) | 32.97 | 0.9049 | 29.12 | 0.7958 | 27.95 | 0.7504 | 27.81 | 0.8336 | 32.15 | 0.9274 |

## 4.3    Ablation Study

**Effectiveness of OCAB and CAB.** We conduct experiments to demonstrate the effectiveness of the proposed CAB and OCAB. The quantitative performance tested on the Urban100 dataset for ×4 SR is shown in Tab. 2. Compared with the baseline, both the proposed OCAB and CAB bring a performance gain of 0.1dB. Benefiting from the two modules, the model obtains a further performance
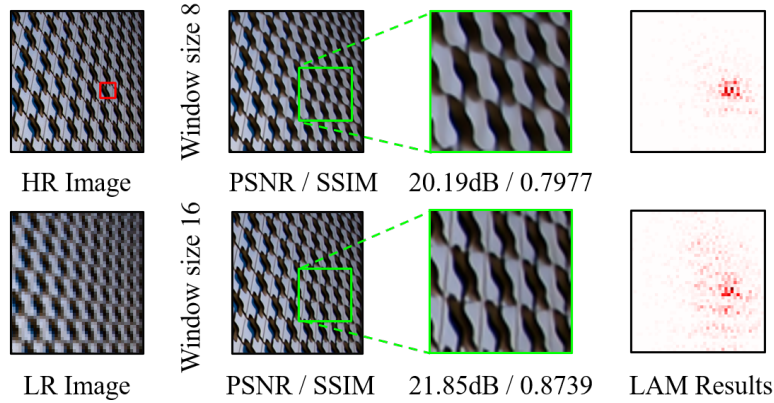
Fig. 5: Qualitative comparison on models with different window sizes.

improvement of 0.16dB. We also provide qualitative comparison to further illustrate the influence of OCAB and CAB, as depicted in Fig. 6. 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 the full image. Moreover, the result of our method with OCAB and CAB obtains the highest DI[14], which means our method utilizes the most input pixels. Although it obtains a little lower performance than w/OCAB, our method gets the highest SSIM and reconstructs the clearest textures.

Table 2: Ablation study on the proposed OCAB and CAB.

| | Baseline | | | |
|---|---|---|---|---|
| OCAB | X | ✓ | X | ✓ |
| CAB | X | X | ✓ | ✓ |
| PSNR/SSIM | 27.81/0.8336 | 27.91/0.8352 | 27.91/0.8355 | 27.97/0.8366 |

**Effects of the overlapping size.** As presented in Sec. 4, we set a constant $\gamma$ to control the overlapping size of OCAB. To explore the effects of different overlapping sizes for the method, we set a group of $\gamma$ from 0 to 0.75 to examine the performance change, as shown in Tab. 3. 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 $\gamma$ 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.

Table 3: Ablation study on the overlapping ratio of OCAB.

| $\gamma$ | 0 | 0.25 | 0.5 | 0.75 |
|---|---|---|---|---|
| PSNR/SSIM | 27.85/0.8341 | 27.81/0.8338 | 27.91/0.8352 | 27.86/0.8347 |

**Effects of different designs of CAB.** We conduct experiments to explore the effects of different designs of CAB. First, we investigate the influence of convolution design and channel attention. As shown in Tab.4, using depth-wise
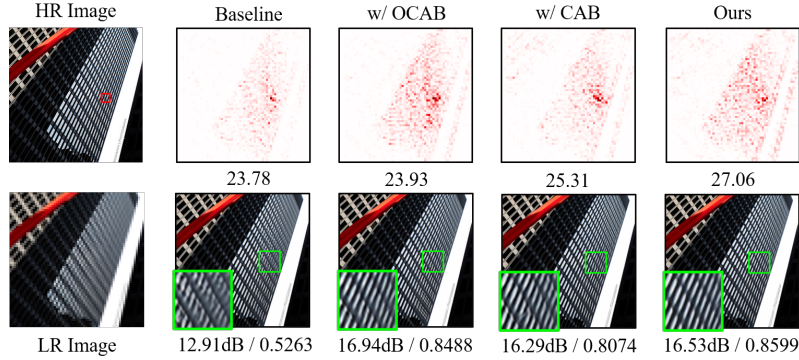
Fig. 6: Qualitative comparisons for the proposed OCAB and CAB.

convolution causes a severe performance drop, which means depth-wise convolution weakens the representative ability of CAB. Besides, we can observe that channel attention brings obvious performance improvement. It demonstrates the effectiveness of channel attention. We also conduct experiments to explore the effects of the weighting factor $\alpha$ of CAB. As presented in Sec. 3.2, $\alpha$ is used to control the weight of CAB features for feature fusion. A larger $\alpha$ means a larger weight of features extracted by CAB and $\alpha = 0$ represents CAB is not used. As shown in Tab. 5, the model with $\alpha$ 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 can suppress this issue for the better combination.

Table 4: Effects of different structures of CAB. "DWCAB" means depth-wise convolution is adopted in CAB. "CA" represents channel attention.

| Structure | DWCAB | w/o CA | w/ CA |
|---|---|---|---|
| PSNR/SSIM | 27.86/0.8329 | 27.92/0.8362 | 27.97/0.8367 |

Table 5: Effects of the weighting factor $\alpha$.

| $\alpha$ | 0 | 1 | 0.1 | 0.01 |
|---|---|---|---|---|
| PSNR/SSIM | 27.81/0.8336 | 27.86/0.8347 | 27.90/0.8358 | 27.97/0.8367 |

### 4.4   Comparison with State-of-the-Art Methods

**Quantitative results.** Tab. 6 shows the quantitative comparison of our approach and the state-of-the-art methods: EDSR [30], RCAN [65], SAN [7], IGNN [69], HAN [41], NLSN [40], RCAN-it [32], as well as approaches using ImageNet pre-training, *i.e.*, IPT [5] and EDT [26]. We can see that our approach outperforms the other methods significantly on all benchmark datasets. Especially, HAT surpasses SwinIR by 0.48dB∼0.64dB on the Urban100 dataset and 0.34dB∼0.45dB on the Manga109 dataset. When compared with the approaches using pre-training, HAT also has large performance gains of more than 0.5dB against EDT on the Urban100 dataset for all three scales. Besides, HAT with pre-training outperforms SwinIR by a huge margin of up to 1dB on the Urban100 dataset for ×2 SR. Moreover, the large model HAT-L can even bring further improvement. It is noteworthy that the performance gaps are much larger on the

Table 6: Quantitative comparison with state-of-the-art methods on benchmark datasets. The top three results are marked in red, blue and green. "†" indicates that methods adopt pre-training strategy on ImageNet.

| Method | Scale | Training Dataset | Set5 [1] PSNR | Set5 [1] SSIM | Set14 [63] PSNR | Set14 [63] SSIM | BSD100 [38] PSNR | BSD100 [38] SSIM | Urban100 [17] PSNR | Urban100 [17] SSIM | Manga109 [39] PSNR | Manga109 [39] SSIM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EDSR [30] | ×2 | DIV2K | 38.11 | 0.9602 | 33.92 | 0.9195 | 32.32 | 0.9013 | 32.93 | 0.9351 | 39.10 | 0.9773 |
| RCAN [65] | ×2 | DIV2K | 38.27 | 0.9614 | 34.12 | 0.9216 | 32.41 | 0.9027 | 33.34 | 0.9384 | 39.44 | 0.9786 |
| SAN [7] | ×2 | DIV2K | 38.31 | 0.9620 | 34.07 | 0.9213 | 32.42 | 0.9028 | 33.10 | 0.9370 | 39.32 | 0.9792 |
| IGNN [69] | ×2 | DIV2K | 38.24 | 0.9613 | 34.07 | 0.9217 | 32.41 | 0.9025 | 33.23 | 0.9383 | 39.35 | 0.9786 |
| HAN [41] | ×2 | DIV2K | 38.27 | 0.9614 | 34.16 | 0.9217 | 32.41 | 0.9027 | 33.35 | 0.9385 | 39.46 | 0.9785 |
| NLSN [40] | ×2 | DIV2K | 38.34 | 0.9618 | 34.08 | 0.9231 | 32.43 | 0.9027 | 33.42 | 0.9394 | 39.59 | 0.9789 |
| RCAN-it [32] | ×2 | DF2K | 38.37 | 0.9620 | 34.49 | 0.9250 | 32.48 | 0.9034 | 33.62 | 0.9410 | 39.88 | 0.9799 |
| SwinIR [29] | ×2 | DF2K | 38.42 | 0.9623 | 34.46 | 0.9250 | 32.53 | 0.9041 | 33.81 | 0.9427 | 39.92 | 0.9797 |
| EDT [26] | ×2 | DF2K | 38.45 | 0.9624 | 34.57 | 0.9258 | 32.52 | 0.9041 | 33.80 | 0.9425 | 39.93 | 0.9800 |
| **HAT** (ours) | ×2 | DF2K | 38.63 | 0.9630 | 34.86 | 0.9274 | 32.62 | 0.9053 | 34.45 | 0.9466 | 40.26 | 0.9809 |
| IPT† [5] | ×2 | ImageNet | 38.37 | - | 34.43 | - | 32.48 | - | 33.76 | - | - | - |
| EDT† [26] | ×2 | DF2K | 38.63 | 0.9632 | 34.80 | 0.9273 | 32.62 | 0.9052 | 34.27 | 0.9456 | 40.37 | 0.9811 |
| **HAT†** (ours) | ×2 | DF2K | 38.73 | 0.9637 | 35.13 | 0.9282 | 32.69 | 0.9060 | 34.81 | 0.9489 | 40.71 | 0.9819 |
| **HAT-L†** (ours) | ×2 | DF2K | 38.91 | 0.9646 | 35.29 | 0.9293 | 32.74 | 0.9066 | 35.09 | 0.9505 | 41.01 | 0.9831 |
| EDSR [30] | ×3 | DIV2K | 34.65 | 0.9280 | 30.52 | 0.8462 | 29.25 | 0.8093 | 28.80 | 0.8653 | 34.17 | 0.9476 |
| RCAN [65] | ×3 | DIV2K | 34.74 | 0.9299 | 30.65 | 0.8482 | 29.32 | 0.8111 | 29.09 | 0.8702 | 34.44 | 0.9499 |
| SAN [7] | ×3 | DIV2K | 34.75 | 0.9300 | 30.59 | 0.8476 | 29.33 | 0.8112 | 28.93 | 0.8671 | 34.30 | 0.9494 |
| IGNN [69] | ×3 | DIV2K | 34.72 | 0.9298 | 30.66 | 0.8484 | 29.31 | 0.8105 | 29.03 | 0.8696 | 34.39 | 0.9496 |
| HAN [41] | ×3 | DIV2K | 34.75 | 0.9299 | 30.67 | 0.8483 | 29.32 | 0.8110 | 29.10 | 0.8705 | 34.48 | 0.9500 |
| NLSN [40] | ×3 | DIV2K | 34.85 | 0.9306 | 30.70 | 0.8485 | 29.34 | 0.8117 | 29.25 | 0.8726 | 34.57 | 0.9508 |
| RCAN-it [32] | ×3 | DF2K | 34.86 | 0.9308 | 30.76 | 0.8505 | 29.39 | 0.8125 | 29.38 | 0.8755 | 34.92 | 0.9520 |
| SwinIR [29] | ×3 | DF2K | 34.97 | 0.9318 | 30.93 | 0.8534 | 29.46 | 0.8145 | 29.75 | 0.8826 | 35.12 | 0.9537 |
| EDT [26] | ×3 | DF2K | 34.97 | 0.9316 | 30.89 | 0.8527 | 29.44 | 0.8142 | 29.72 | 0.8814 | 35.13 | 0.9534 |
| **HAT** (ours) | ×3 | DF2K | 35.06 | 0.9329 | 31.08 | 0.8555 | 29.54 | 0.8167 | 30.23 | 0.8896 | 35.53 | 0.9552 |
| IPT† [5] | ×3 | ImageNet | 34.81 | - | 30.85 | - | 29.38 | - | 29.49 | - | - | - |
| EDT† [26] | ×3 | DF2K | 35.13 | 0.9328 | 31.09 | 0.8553 | 29.53 | 0.8165 | 30.07 | 0.8863 | 35.47 | 0.9550 |
| **HAT†** (ours) | ×3 | DF2K | 35.16 | 0.9335 | 31.33 | 0.8576 | 29.59 | 0.8177 | 30.70 | 0.8949 | 35.84 | 0.9567 |
| **HAT-L†** (ours) | ×3 | DF2K | 35.28 | 0.9345 | 31.47 | 0.8584 | 29.63 | 0.8191 | 30.92 | 0.8981 | 36.02 | 0.9576 |
| EDSR [30] | ×4 | DIV2K | 32.46 | 0.8968 | 28.80 | 0.7876 | 27.71 | 0.7420 | 26.64 | 0.8033 | 31.02 | 0.9148 |
| RCAN [65] | ×4 | DIV2K | 32.63 | 0.9002 | 28.87 | 0.7889 | 27.77 | 0.7436 | 26.82 | 0.8087 | 31.22 | 0.9173 |
| SAN [7] | ×4 | DIV2K | 32.64 | 0.9003 | 28.92 | 0.7888 | 27.78 | 0.7436 | 26.79 | 0.8068 | 31.18 | 0.9169 |
| IGNN [69] | ×4 | DIV2K | 32.57 | 0.8998 | 28.85 | 0.7891 | 27.77 | 0.7434 | 26.84 | 0.8090 | 31.28 | 0.9182 |
| HAN [41] | ×4 | DIV2K | 32.64 | 0.9002 | 28.90 | 0.7890 | 27.80 | 0.7442 | 26.85 | 0.8094 | 31.42 | 0.9177 |
| NLSN [40] | ×4 | DIV2K | 32.59 | 0.9000 | 28.87 | 0.7891 | 27.78 | 0.7444 | 26.96 | 0.8109 | 31.27 | 0.9184 |
| RRDB [53] | ×4 | DF2K | 32.73 | 0.9011 | 28.99 | 0.7917 | 27.85 | 0.7455 | 27.03 | 0.8153 | 31.66 | 0.9196 |
| RCAN-it [32] | ×4 | DF2K | 32.69 | 0.9007 | 28.99 | 0.7922 | 27.87 | 0.7459 | 27.16 | 0.8168 | 31.78 | 0.9217 |
| SwinIR [29] | ×4 | DF2K | 32.92 | 0.9044 | 29.09 | 0.7950 | 27.92 | 0.7489 | 27.45 | 0.8254 | 32.03 | 0.9260 |
| EDT [26] | ×4 | DF2K | 32.82 | 0.9031 | 29.09 | 0.7939 | 27.91 | 0.7483 | 27.46 | 0.8246 | 32.05 | 0.9254 |
| **HAT** (ours) | ×4 | DF2K | 33.04 | 0.9056 | 29.23 | 0.7973 | 28.00 | 0.7517 | 27.97 | 0.8368 | 32.48 | 0.9292 |
| IPT† [5] | ×4 | ImageNet | 32.64 | - | 29.01 | - | 27.82 | - | 27.26 | - | - | - |
| EDT† [26] | ×4 | DF2K | 33.06 | 0.9055 | 29.23 | 0.7971 | 27.99 | 0.7510 | 27.75 | 0.8317 | 32.39 | 0.9283 |
| **HAT†** (ours) | ×4 | DF2K | 33.18 | 0.9073 | 29.38 | 0.8001 | 28.05 | 0.7534 | 28.37 | 0.8447 | 32.87 | 0.9319 |
| **HAT-L†** (ours) | ×4 | DF2K | 33.30 | 0.9083 | 29.47 | 0.8015 | 28.09 | 0.7551 | 28.60 | 0.8498 | 33.09 | 0.9335 |

Urban100 dataset, as it contains more structured and self-repeated patterns. This is beneficial for our method to utilize more useful information. All these results have demonstrated the effectiveness of our method.

**Visual comparison.** We also provide the visual comparison as shown in Fig. 7. For the images "img_011", "img_044" and "img_073" in Urban100 dataset, our method can successfully recover the clear lattice content. In contrast, the other approaches all suffer from severe blurry effects.
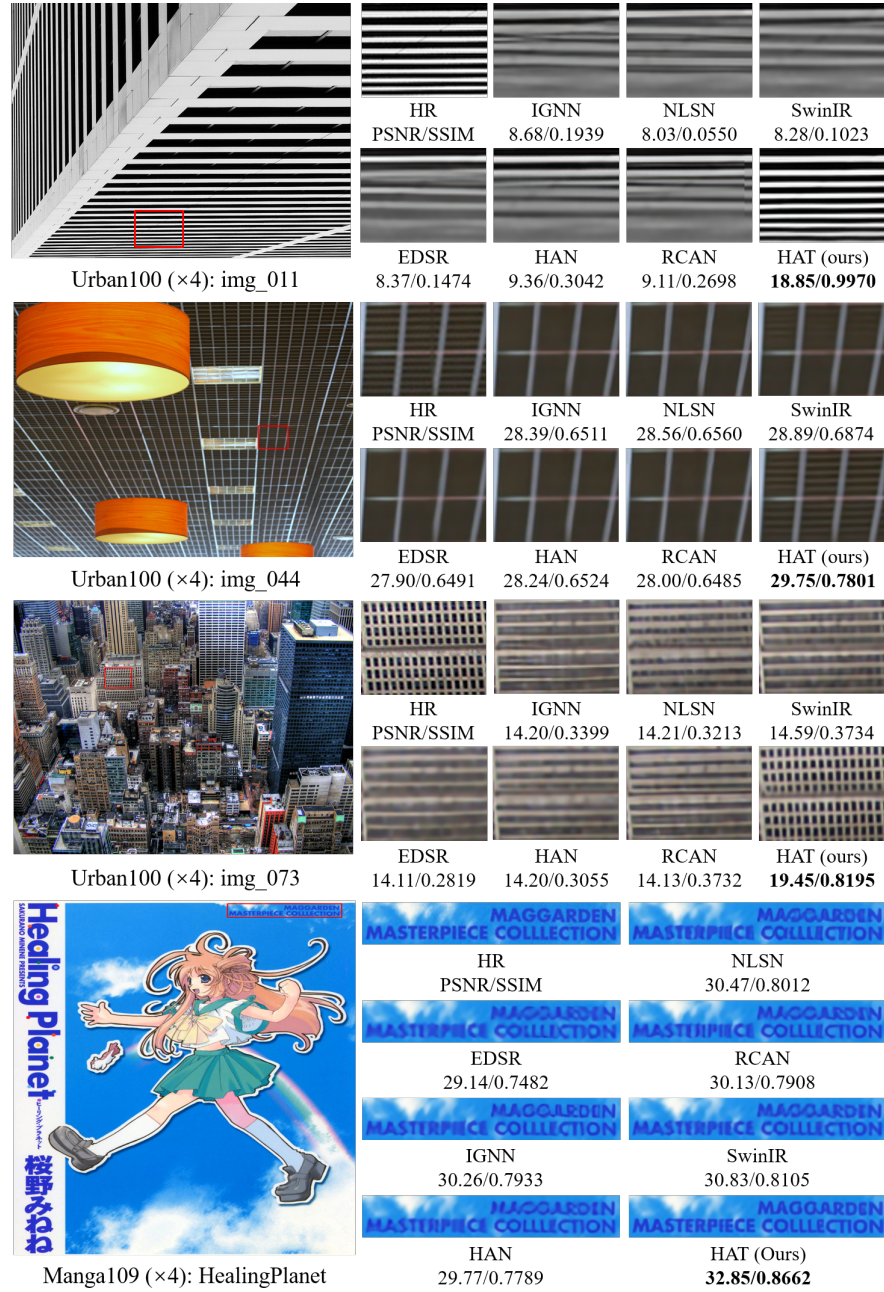
| | | | |
|---|---|---|---|
| HR<br>PSNR/SSIM | IGNN<br>8.68/0.1939 | NLSN<br>8.03/0.0550 | SwinIR<br>8.28/0.1023 |
| EDSR<br>8.37/0.1474 | HAN<br>9.36/0.3042 | RCAN<br>9.11/0.2698 | HAT (ours)<br>**18.85/0.9970** |

Urban100 (×4): img_011

| | | | |
|---|---|---|---|
| HR<br>PSNR/SSIM | IGNN<br>28.39/0.6511 | NLSN<br>28.56/0.6560 | SwinIR<br>28.89/0.6874 |
| EDSR<br>27.90/0.6491 | HAN<br>28.24/0.6524 | RCAN<br>28.00/0.6485 | HAT (ours)<br>**29.75/0.7801** |

Urban100 (×4): img_044

| | | | |
|---|---|---|---|
| HR<br>PSNR/SSIM | IGNN<br>14.20/0.3399 | NLSN<br>14.21/0.3213 | SwinIR<br>14.59/0.3734 |
| EDSR<br>14.11/0.2819 | HAN<br>14.20/0.3055 | RCAN<br>14.13/0.3732 | HAT (ours)<br>**19.45/0.8195** |

Urban100 (×4): img_073

| | |
|---|---|
| HR<br>PSNR/SSIM | NLSN<br>30.47/0.8012 |
| EDSR<br>29.14/0.7482 | RCAN<br>30.13/0.7908 |
| IGNN<br>30.26/0.7933 | SwinIR<br>30.83/0.8105 |
| HAN<br>29.77/0.7789 | HAT (Ours)<br>**32.85/0.8662** |

Manga109 (×4): HealingPlanet

Fig. 7: Visual comparison for ×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.

We can observe similar behaviors on "HealingPlanet" in Manga109 dataset. When recovering the characters in the image, HAT obtains clearer textures than the other methods. Combining the quantitative comparison results, we demonstrate the superiority of our method.

### 4.5  Effectiveness of the same-task pre-training

As shown in Tab. 6, models using the same-task pre-training strategy significantly outperform models without pre-training. EDT [26] also explores the effects of different pre-training strategies for the SR task. It demonstrates that pre-training on ImageNet based on multi-related-tasks (*i.e.*, pre-train the SR model on ×2, ×3, ×4 SR) is the most effective strategy compared to single-task pre-training (*e.g.*, pre-training on ×2 SR setup for ×4 SR task) and multi-unrelated-task pre-training similar to [5]. To demonstrate the effectiveness and superiority of our strategy, we also apply the same strategy as EDT in our network. For a fair comparison, both strategies adopt the full ImageNet dataset with 1.28 million images. We provide the quantitative results of the pre-trained models as well as the fine-tuned models on ×4 SR, as shown in Tab. 7. As one can see that the same-task pre-training performs better, not only in the pre-training stage but also in the fine-tuning process. Compared to pre-training on the specific task, the multi-task pre-training seems to weaken the performance. From this perspective, we tend to believe that the reason "why pre-training works" is attributed to the diversity of data instead of the correlation between tasks.

Table 7: Quantitative results of HAT using two kinds of pre-training strategies on ×4 SR under the same training setting. The full ImageNet dataset is adopted to perform pre-training and DF2K dataset is used for fine-tuning.

| Pre-training Strategy | Stage | Set5 | | Set14 | | BSD100 | | Urban100 | | Manga109 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| Multi-related-task pre-training [26] | pre-training | 32.94 | 0.9057 | 29.17 | 0.7970 | 27.99 | 0.7524 | 28.05 | 0.8405 | 31.66 | 0.9266 |
| | fine-tuning | 33.06 | 0.9063 | 29.34 | 0.7988 | 28.03 | 0.7525 | 28.21 | 0.8414 | 32.71 | 0.9309 |
| Same-task pre-training (ours) | pre-training | 33.02 | 0.9063 | 29.20 | 0.7975 | 27.99 | 0.7528 | 28.11 | 0.8422 | 31.72 | 0.9273 |
| | fine-tuning | 33.07 | 0.9064 | 29.33 | 0.7991 | 28.04 | 0.7528 | 28.28 | 0.8432 | 32.79 | 0.9314 |

## 5  Conclusion

In this paper, we propose a novel Hybrid Attention Transformer, HAT, for image super-resolution. Our model combines channel attention and self-attention to activate more pixels for reconstructing high-resolution results. Besides, we propose an overlapping cross-attention module that calculates attention between features with different window sizes to better aggregate the cross-window information. Moreover, we introduce a same-task pre-training strategy to further activate the potential of the proposed model. Extensive experiments show the effectiveness of the proposed modules, and our HAT significantly outperforms the state-of-the-art methods.

# References

1. Bevilacqua, M., Roumy, A., Guillemot, C., Alberi-Morel, M.L.: Low-complexity single-image super-resolution based on nonnegative neighbor embedding (2012)
2. Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-unet: Unet-like pure transformer for medical image segmentation. arXiv preprint arXiv:2105.05537 (2021)
3. Cao, J., Li, Y., Zhang, K., Van Gool, L.: Video super-resolution transformer. arXiv preprint arXiv:2106.06847 (2021)
4. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV, 2020
5. Chen, H., Wang, Y., Guo, T., Xu, C., Deng, Y., Liu, Z., Ma, S., Xu, C., Xu, C., Gao, W.: Pre-trained image processing transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12299–12310 (2021)
6. Chu, X., Tian, Z., Wang, Y., Zhang, B., Ren, H., Wei, X., Xia, H., Shen, C.: Twins: Revisiting the design of spatial attention in vision transformers. Advances in Neural Information Processing Systems **34** (2021)
7. Dai, T., Cai, J., Zhang, Y., Xia, S.T., Zhang, L.: Second-order attention network for single image super-resolution. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11065–11074 (2019)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
9. Dong, C., Loy, C.C., He, K., Tang, X.: Learning a deep convolutional network for image super-resolution. In: European conference on computer vision. pp. 184–199. Springer (2014)
10. Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. IEEE transactions on pattern analysis and machine intelligence **38**(2), 295–307 (2015)
11. Dong, C., Loy, C.C., Tang, X.: Accelerating the super-resolution convolutional neural network. In: European conference on computer vision. pp. 391–407. Springer (2016)
12. Dong, X., Bao, J., Chen, D., Zhang, W., Yu, N., Yuan, L., Chen, D., Guo, B.: Cswin transformer: A general vision transformer backbone with cross-shaped windows. arXiv preprint arXiv:2107.00652 (2021)
13. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020
14. Gu, J., Dong, C.: Interpreting super-resolution networks with local attribution maps. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9199–9208 (2021)
15. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 (2016)
16. Huang, G., Wang, Y., Lv, K., Jiang, H., Huang, W., Qi, P., Song, S.: Glance and focus networks for dynamic visual recognition. arXiv preprint arXiv:2201.03014 (2022)
17. Huang, J.B., Singh, A., Ahuja, N.: Single image super-resolution from transformed self-exemplars. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5197–5206 (2015)

18. Huang, Z., Ben, Y., Luo, G., Cheng, P., Yu, G., Fu, B.: Shuffle transformer: Rethinking spatial shuffle for vision transformer. arXiv preprint arXiv:2106.03650 (2021)
19. Kim, J., Lee, J.K., Lee, K.M.: Accurate image super-resolution using very deep convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1646–1654 (2016)
20. Kim, J., Lee, J.K., Lee, K.M.: Deeply-recursive convolutional network for image super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1637–1645 (2016)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
22. Kong, X., Liu, X., Gu, J., Qiao, Y., Dong, C.: Reflash dropout in image super-resolution. arXiv preprint arXiv:2112.12089 (2021)
23. Lai, W.S., Huang, J.B., Ahuja, N., Yang, M.H.: Deep laplacian pyramid networks for fast and accurate super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 624–632 (2017)
24. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., 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. pp. 4681–4690 (2017)
25. Li, K., Wang, Y., Zhang, J., Gao, P., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unifying convolution and self-attention for visual recognition. arXiv preprint arXiv:2201.09450 (2022)
26. Li, W., Lu, X., Lu, J., Zhang, X., Jia, J.: On efficient transformer and image pre-training for low-level vision. arXiv preprint arXiv:2112.10175 (2021)
27. Li, Y., Zhang, K., Cao, J., Timofte, R., Van Gool, L.: Localvit: Bringing locality to vision transformers. arXiv preprint arXiv:2104.05707 (2021)
28. Liang, J., Cao, J., Fan, Y., Zhang, K., Ranjan, R., Li, Y., Timofte, R., Van Gool, L.: Vrt: A video restoration transformer. arXiv preprint arXiv:2201.12288 (2022)
29. Liang, J., Cao, J., Sun, G., Zhang, K., Van Gool, L., Timofte, R.: Swinir: Image restoration using swin transformer. In: ICCVW, 2021
30. Lim, B., Son, S., Kim, H., Nah, S., Mu Lee, K.: Enhanced deep residual networks for single image super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 136–144 (2017)
31. Lim, B., Son, S., Kim, H., Nah, S., Mu Lee, K.: Enhanced deep residual networks for single image super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 136–144 (2017)
32. Lin, Z., Garg, P., Banerjee, A., Magid, S.A., Sun, D., Zhang, Y., Van Gool, L., Wei, D., Pfister, H.: Revisiting rcan: Improved training for image super-resolution. arXiv preprint arXiv:2201.11279 (2022)
33. Liu, D., Wen, B., Fan, Y., Loy, C.C., Huang, T.S.: Non-local recurrent network for image restoration. Advances in neural information processing systems **31** (2018)
34. Liu, L., Ouyang, W., Wang, X., Fieguth, P., Chen, J., Liu, X., Pietikäinen, M.: Deep learning for generic object detection: A survey. International journal of computer vision **128**(2), 261–318 (2020)
35. Liu, Y., Liu, A., Gu, J., Zhang, Z., Wu, W., Qiao, Y., Dong, C.: Discovering" semantics" in super-resolution networks. arXiv preprint arXiv:2108.00406 (2021)
36. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. arXiv preprint arXiv:2103.14030 (2021)

37. Lu, Z., Liu, H., Li, J., Zhang, L.: Efficient transformer for single image super-resolution. arXiv preprint arXiv:2108.11084 (2021)
38. Martin, D., Fowlkes, C., Tal, D., Malik, J.: 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, pp. 416–423. IEEE (2001)
39. Matsui, Y., Ito, K., Aramaki, Y., Fujimoto, A., Ogawa, T., Yamasaki, T., Aizawa, K.: Sketch-based manga retrieval using manga109 dataset. Multimedia Tools and Applications **76**(20), 21811–21838 (2017)
40. Mei, Y., Fan, Y., Zhou, Y.: Image super-resolution with non-local sparse attention. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3517–3526 (2021)
41. Niu, B., Wen, W., Ren, W., Zhang, X., Yang, L., Wang, S., Zhang, K., Cao, X., Shen, H.: Single image super-resolution via a holistic attention network. In: European conference on computer vision. pp. 191–207. Springer (2020)
42. Patel, K., Bur, A.M., Li, F., Wang, G.: Aggregating global features into local vision transformer. arXiv preprint arXiv:2201.12903 (2022)
43. Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., Dosovitskiy, A.: Do vision transformers see like convolutional neural networks? Advances in Neural Information Processing Systems **34** (2021)
44. Ramachandran, P., Parmar, N., Vaswani, A., Bello, I., Levskaya, A., Shlens, J.: Studying stand-alone self-attention in vision models (2019)
45. Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: 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. pp. 1874–1883 (2016)
46. Tai, Y., Yang, J., Liu, X.: Image super-resolution via deep recursive residual network. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3147–3155 (2017)
47. Timofte, R., Agustsson, E., Van Gool, L., Yang, M.H., Zhang, L.: Ntire 2017 challenge on single image super-resolution: Methods and results. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 114–125 (2017)
48. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. pp. 10347–10357. PMLR (2021)
49. Vaswani, A., Ramachandran, P., Srinivas, A., Parmar, N., Hechtman, B., Shlens, J.: Scaling local self-attention for parameter efficient visual backbones. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12894–12904 (2021)
50. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
51. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 568–578 (2021)
52. Wang, X., Xie, L., Dong, C., Shan, Y.: Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1905–1914 (2021)

53. Wang, X., Yu, K., Wu, S., Gu, J., Liu, Y., Dong, C., Qiao, Y., Change Loy, C.: Esrgan: Enhanced super-resolution generative adversarial networks. In: Proceedings of the European conference on computer vision (ECCV) workshops. pp. 0–0 (2018)

54. Wang, Z., Cun, X., Bao, J., Liu, J.: Uformer: A general u-shaped transformer for image restoration. arXiv preprint arXiv:2106.03106 (2021)

55. Wu, B., Xu, C., Dai, X., Wan, A., Zhang, P., Yan, Z., Tomizuka, M., Gonzalez, J., Keutzer, K., Vajda, P.: Visual transformers: Token-based image representation and processing for computer vision. arXiv preprint arXiv:2006.03677 (2020)

56. Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., Zhang, L.: Cvt: Introducing convolutions to vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 22–31 (2021)

57. Wu, S., Wu, T., Tan, H., Guo, G.: Pale transformer: A general vision transformer backbone with pale-shaped attention. arXiv preprint arXiv:2112.14000 (2021)

58. Xiao, T., Dollar, P., Singh, M., Mintun, E., Darrell, T., Girshick, R.: Early convolutions help transformers see better. Advances in Neural Information Processing Systems **34** (2021)

59. Xie, L., Wang, X., Dong, C., Qi, Z., Shan, Y.: Finding discriminative filters for specific degradations in blind super-resolution. Advances in Neural Information Processing Systems **34** (2021)

60. Yuan, K., Guo, S., Liu, Z., Zhou, A., Yu, F., Wu, W.: Incorporating convolution designs into visual transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 579–588 (2021)

61. Yuan, Y., Fu, R., Huang, L., Lin, W., Zhang, C., Chen, X., Wang, J.: Hrformer: High-resolution transformer for dense prediction. arXiv preprint arXiv:2110.09408 (2021)

62. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H.: Restormer: Efficient transformer for high-resolution image restoration. arXiv preprint arXiv:2111.09881 (2021)

63. Zeyde, R., Elad, M., Protter, M.: On single image scale-up using sparse-representations. In: International conference on curves and surfaces. pp. 711–730. Springer (2010)

64. Zhang, W., Liu, Y., Dong, C., Qiao, Y.: Ranksrgan: Generative adversarial networks with ranker for image super-resolution. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3096–3105 (2019)

65. Zhang, Y., Li, K., Li, K., Wang, L., Zhong, B., Fu, Y.: Image super-resolution using very deep residual channel attention networks. In: Proceedings of the European conference on computer vision (ECCV). pp. 286–301 (2018)

66. Zhang, Y., Li, K., Li, K., Zhong, B., Fu, Y.: Residual non-local attention networks for image restoration. arXiv preprint arXiv:1903.10082 (2019)

67. Zhang, Y., Tian, Y., Kong, Y., Zhong, B., Fu, Y.: Residual dense network for image super-resolution. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2472–2481 (2018)

68. Zhao, Y., Wang, G., Tang, C., Luo, C., Zeng, W., Zha, Z.J.: A battle of network structures: An empirical study of cnn, transformer, and mlp. arXiv preprint arXiv:2108.13002 (2021)

69. Zhou, S., Zhang, J., Zuo, W., Loy, C.C.: Cross-scale internal graph neural network for image super-resolution. Advances in neural information processing systems **33**, 3499–3509 (2020)

# A   Appendix

## A.1   Training Details

We use DF2K (DIV2K [31]+Flicker2K [47]) with 3360 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 $64 \times 64$ 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 [21] optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.99$ to train the model. For the same-task pre-training, the full ImageNet dataset [8] with 1.28 million images is first exploited 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 also 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 of 250K training iterations.

## A.2   Analysis of Model Complexity

We conduct experiments to analyze the computational complexity of our method from three aspects: large window size, overlapping cross-attention block (OCAB) and channel attention block (CAB) in the hybrid attention block. We evaluate the performance based on the results of $\times 4$ SR on the Urban100 [17] and the number of Multiply-Add operations is counted at the input size of $64 \times 64$. Note that basic models without using pre-training are used for experiments here.

First, we use the standard Swin Transformer block [29] as the backbone to explore the influence on different window sizes. As shown in Tab. 8, enlarging window size can bring significant performance improvement (+0.36dB) with a little increase in parameters and $\sim \%19$ increase in Multiply-Add operations.

We use window size 16 as the baseline to investigate the computational complexity of the proposed OCAB and CAB. As illustrated in Tab. 9, our OCAB obtains a performance gain with a limited increase of parameters and Multi-Adds. It demonstrates that the proposed OCAB is effective and efficient. Besides, Adding CAB to the baseline model also achieves better performance.

Since CAB seems to be computationally expensive, we further explore the influence on different sizes of CAB by modulating the squeeze factor $\beta$ in CAB (mentioned in Sec. 3.2 in the main paper). As shown in Tab. 10, adding a small CAB whose $\beta$ equals 6 can bring performance improvement. When we continuously reduce $\beta$, the performance increases but with larger model sizes. To balance the performance and computational cost, we set $\beta$ to 3 as the default setting.

Furthermore, we compare HAT and SwinIR [29] with the similar numbers of parameters and Multi-Adds in two different settings. For SwinIR-1 and SwinIR-2, we increase the width and depth of the original SwinIR to achieve similar

computational complexity as HAT. As shown in Tab.11, HAT obtains the best performance with the lowest computational cost.

Table 8: Comparison of computational complexity on different window sizes.

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

Table 9: Comparison of computational complexity on OCAB and CAB.

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

Table 10: Comparison of computational complexity on different sizes of CAB.

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

Table 11: Comparison of computational complexity between SwinIR and HAT.

| Method | #Params. (M) | #Multi-Adds. (G) | PSNR (dB) |
|---|---|---|---|
| SwinIR-1 | 24.0 | 104.4 | 27.53 |
| SwinIR-2 | 23.1 | 102.4 | 27.58 |
| HAT | 20.8 | 103.7 | 27.97 |

### A.3   More LAM Results

We provide more qualitative and quantitative comparisons of the LAM results between SwinIR and our method. The red points in LAM results represent the used pixels for reconstructing the patch marked with a red box in the HR image, and DI is computed to reflect the range of involved pixels. The more pixels are utilized to recover the specific input patch, the wider the distribution of red points is in LAM and the higher Diffusion Index (DI) [14] is. As shown in Fig. 8, the LAM attribution 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. All these results demonstrate that our method activates 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.
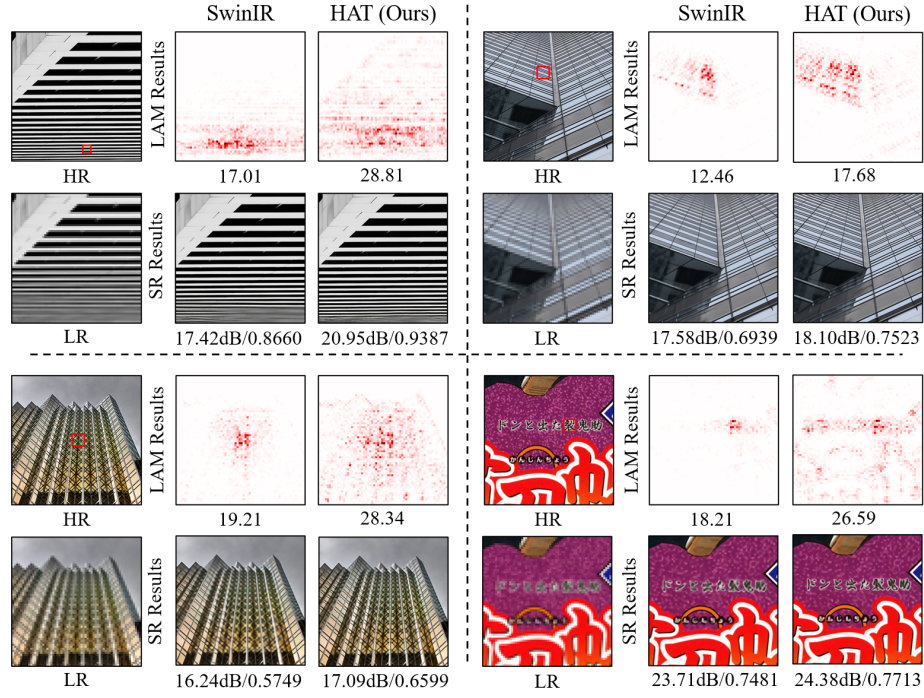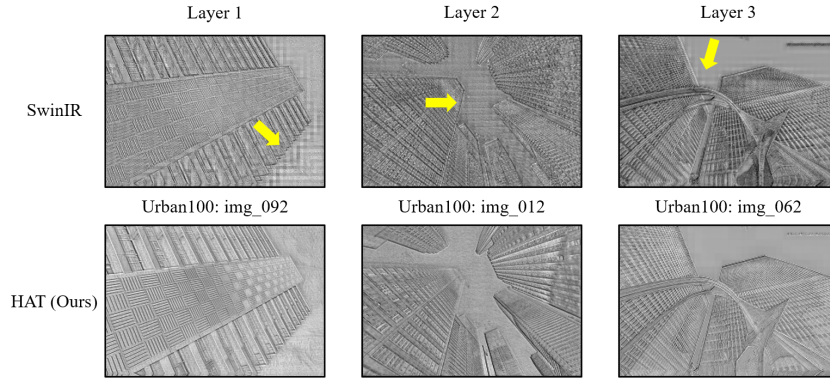
Fig. 8: More LAM comparisons between SwinIR and HAT.



Fig. 9: Visual comparisons of intermediate features between SwinIR and our method. "Layer $N$" represents the intermediate features after the $N_{th}$ layer (*i.e.*, RSTB in SwinIR and RHAG in our HAT, respectively).

## A.4    Comparison of Intermediate Features

We also provide the visual comparison of intermediate features. Since features generated by different models cannot be exactly aligned, to make an alignment, we select the features in the same layers as that in SwinIR by minimizing the $L1$ distance. As depicted in Fig. 9, severe blocking artifacts can be observed in the features generated by SwinIR, while our HAT obtains cleaner textures without blocking artifacts. Note that "Layer $N$" means the features are extracted after the $N_{th}$ layer (*i.e.*, RSTB in SwinIR and RHAG in HAT, respectively). The results demonstrate that our method can better aggregate cross-window information and alleviate the blocking artifacts in the intermediate features.