

# HST: Hierarchical Swin Transformer for Compressed Image Super-resolution

Bingchen Li, Xin Li, Yiting Lu, Sen Liu, Ruoyu Feng, and Zhibo Chen\*

University of Science and Technology of China, Hefei 230027, China  
{1bc31415926, lixin666, luyt31415}@mail.ustc.edu.cn,  
elsen@iat.ustc.edu.cn, ustcfry@mail.ustc.edu.cn, chenzhibo@ustc.edu.cn

**Abstract.** Compressed Image Super-resolution has achieved great attention in recent years, where images are degraded with compression artifacts and low-resolution artifacts. Since the complex hybrid distortions, it is hard to restore the distorted image with the simple cooperation of super-resolution and compression artifacts removing. In this paper, we take a step forward to propose the Hierarchical Swin Transformer (HST) network to restore the low-resolution compressed image, which jointly captures the hierarchical feature representations and enhances each-scale representation with Swin transformer, respectively. Moreover, we find that the pretraining with Super-resolution (SR) task is vital in compressed image super-resolution. To explore the effects of different SR pretraining, we take the commonly-used SR tasks (*e.g.*, bicubic and different real super-resolution simulations) as our pretraining tasks, and reveal that SR plays an irreplaceable role in the compressed image super-resolution. With the cooperation of HST and pre-training, our HST achieves the fifth place in AIM 2022 challenge on the low-quality compressed image super-resolution track, with the PSNR of 23.51dB. Extensive experiments and ablation studies have validated the effectiveness of our proposed methods. The code and models are available at <https://github.com/USTC-IMCL/HST-for-Compressed-Image-SR>.

**Keywords:** Hierarchical network, Transformer, Compressed image super-resolution, Pretraining, AIM 2022 challenge

## 1 Introduction

Image super-resolution (SR) has achieved a quantum leap with the development of deep neural networks, which aims to restore the high-resolution (HR) images from their low-resolution counterparts. Existing SR can be roughly divided into three categories, simulated SR [12,30,19,57,47] (*e.g.*, bicubic downsampling), real-world SR [5,17,50,46,55,24,49] and blind SR [2,14,29,45,36], respectively. In particular, real-world and blind SR are greatly developed in recent years, of which the degradations are more consistent with unknown real-world distortions. However, not all images suffer from real-world degradation. In most cases,

---

\* corresponding author

the images are susceptible to various compression artifacts together with low-resolution, since the image compression [52,25,40,41,4], transmission and storage. This hybrid degradation poses a challenging image process task, *i.e.*, compressed image super-resolution.

As shown in Fig. 1, unlike general image SR and compression artifacts, the degradations of compressed low-resolution images are more severe, which composes of blurring, block artifacts, and noise, etc. Existing methods on image SR [30,57,12,47] and compression artifacts removing [11,43,6,56,15] cannot work well on such brand-new degradation, since the large distribution shift. As the pioneering works, a series of works [58,27,51] began to investigate the compressed video super-resolution. To further promote the development of compressed image/video super-resolution, AIM2022 [53] firstly holds the significant competition on compressed image super-resolution, where images are firstly down-sampled with the scale 1/4, and then, are compressed with JPEG using an extreme low-quality parameter  $Q = 10$ . A naïve and intuitive strategy to deal with it is exploiting a well-trained SR network and JPEG artifacts removing network to restore the distorted images in a sequential manner. However, the above strategy always fails since the distribution shift between hybrid distortions [23,32] and single distortion. Compressed image super-resolution requires that the restoration network have the strong representation capability to learn structure and texture jointly.

In this paper, we present Hierarchical Swin Transformer, namely HST, to tackle the compressed image super-resolution problem. Specifically, previous works [26,56] have shown superior advantages of hierarchical architecture on compression artifacts removing due to their great representation ability. Meanwhile, the variants of the transformer have been explored for image processing and quality assessment [8,31,35], *e.g.*, SwinIR [28], which achieve remarkable performance compared with their CNN counterparts, since their capability of global contextual representation. Inspired by these, we present the Hierarchical Swin Transformer (HST) by incorporating the individual advantage of the above two architectures. In particular, our HST consists of four modules: hierarchical feature extraction module, feature enhancement module, fusion module, and HR reconstruction module. The hierarchical feature extraction module uses multiple convolution layers with different strides to obtain hierarchical feature maps at different scales. Then, the residual swin transformer block (RSTB) from SwinIR [28] is used for the feature enhancement in each hierarchical branch. After getting the enhanced hierarchical features, we fuse them by concatenating the upsampling low-scale feature and high-scale feature, and then, input them into a convolution layer to obtain the fused feature. Lastly, we can get the super-resolved HR image with the HR reconstruction module, which is composed of convolution layers and pixelshuffle layers.

We also investigate the compressed image super-resolution from the perspective of pretraining. We observe that the pretraining with image super-resolution plays a vital role in the compressed image SR. Specifically, we systematically explore the effects of different image super-resolution tasks, including traditional

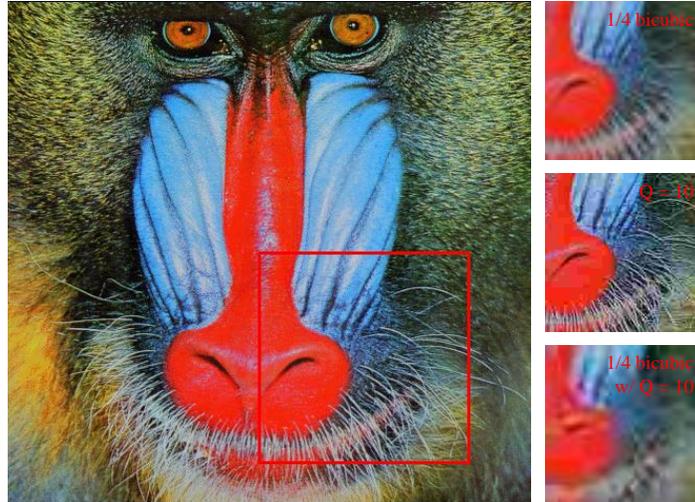


Fig. 1: A comparison between different degradations. The left image is the high-resolution reference image. The images from the top right to bottom right are 1/4 bicubic downsampling, JPEG compression with a quality factor of 10, and the combination of the above two distortions, respectively. Note that the bottom right image has the most severe degradation, thus requiring a stronger representation ability for the network to remove distortion.

SR, *i.e.*, bicubic downsampling, and two RealSR simulation methods from BSR-GAN [55] and DRTL [22]. Extensive experiments reveal that the pretraining with the RealSR simulation from DRTL [22] is better for compressed image SR.

The contributions of this paper can be summarized as:

1. We present the Hierarchical Swin Transformer (HST) for compressed image super-resolution, which incorporating the advantages of strong representation ability and global information utilization.
2. We investigate compressed image super-resolution from the pretraining perspective. Based on the observation, we find one proper pretraining scheme for compressed image super-resolution.
3. Extensive experiment results show that our HST achieve a remarkable result on compressed image super-resolution task under heavy distortion (compression quality  $Q = 10$  combined with 1/4 downsampling).

## 2 Related Works

### 2.1 Single Image Super-Resolution

Single Image Super-resolution (SISR) has been developed expeditiously with the advances of deep neural networks. SRCNN [12], as the pioneering work, firstly introduces the CNN to SISR and learns the network by minimizing the mean square error (MSE) between the generated images and their corresponding high-resolution (HR) images. Then, a series of works for SISR [30,57] are proposed by designing or modulating the network architecture. EDSR [30] revises the conventional residual module by removing the BatchNorm layers. RCAN [57] adds channel attention to the residual blocks, which focus on more informative channels. And SAN [10] introduces the second-order channel attention to utilize the second-order feature statistics for more discriminative representations.

However, the above works exhibit poor capability for subjective quality improving. To tackle the above challenge, SRGAN [21] firstly introduces Generative Adversarial Network (GAN) to SISR, and adopts the adversarial loss for approximating the natural image manifold. As an improved version of SRGAN, ESRGAN [47] exploits Relativistic average GAN [18] to enhance the discriminator and computes the VGG feature before the activation function to calculate the perceptual loss. To further improve the discriminator, FSMR [20] comes up with feature statistics mixing regularization, which encourages the discriminator’s prediction to remain invariant to the style of the input image.

Recently, real-world image super-resolution (RealSR) [5,17,50,55,46] and blind image super-resolution [14,29,2,45] have been proposed to solve more severe and unknown hybrid distortions existed in real-world low-resolution images. To tackle unseen distortions (*i.e.*, blind distortions), KernelGAN [2] train an internal-gan to estimate the degradation kernel contains in low-resolution images. IKC [14] ameliorates the estimation process into an iterative one, which can deal with more complex blind distortions. Different from the aforementioned works that predict the degradation kernel, RealSR directly trains networks on synthesized real-world distorted image pairs, such as BSRGAN [55] and ESRGAN [47]. In this paper, we focus on the compressed image super-resolution, which is more significant and valuable in the real world.

### 2.2 Compression Artifacts Removal

Compression artifacts removal aims to remove the distortions caused by image/video codecs. Early works of compression artifact removal mainly focus on the design of manual filters in the DCT domain. Due to the success of CNN in image denoising and image super-resolution, Yu et al. [11] propose ARCNN, the first CNN-based method for compression artifacts removal. Svoboda et al. [43] introduce residual learning to deepen the network under the assumption of “deeper is better”. However, ARCNN and its follow-up works only process artifacts in the pixel domain. DDCN [15], DMCNN [56] and  $D^3$  [48] utilize DCT domain prior on the basis of pixel domain. Based on the network of extracting the dual

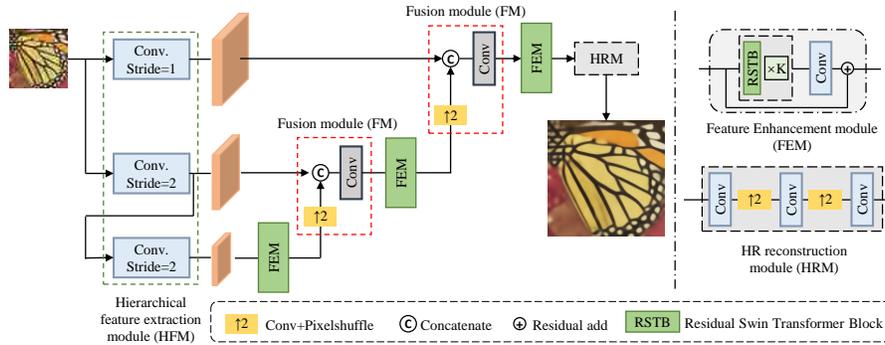


Fig. 2: The architecture of proposed HST

domain knowledge, Fu et al. [13] use dilated convolution for multi-scale feature extraction and added convolutional sparse coding to make the model more compact and explainable. Recently, there are a series of works that explore the hierarchical structures for compression artifacts removal. Lu et al. [34] prove that adding multi-scale priors to the image restoration network can effectively eliminate compression artifacts. Inspired by Lu et al. [34], Li et.al add a non-local attention module to fuse multi-scale features effectively and obtain the post-processing network MSGDN [26] of VCC Intra coding. Based on the above excellent works, we also introduce the hierarchical module to our compressed image super-resolution network.

### 3 Method

In this section, we will explain our HST and clarify our pretraining strategy in detail. As shown in Fig. 2, our HST is composed of four main components, respectively as hierarchical feature extraction module (HFM), feature enhancement module (FEM), fusion module (FM) and HR reconstruction module (HRM).

#### 3.1 Hierarchical Feature Extraction

Previous works [9,39] have revealed that extracting hierarchical features at different scales from images and processing them in a divide-and-conquer manner, can provide the network a strong representation ability. And thus, it can deal with more severe and complex image degradation effectively. To achieve a trade-off between network parameters and performance, we choose a three-branches hierarchical architecture as our backbone. More specifically, the input LR image is gradually passed through three different convolution layers with different kernel

sizes and strides, to extract the hierarchical representations with three scales. Following the implementation in [39], we design upper branch convolution by  $k7n60s1p3$ , where  $k, n, s, p$  stands for kernel size, number of channels, stride and padding respectively. For other two branches, we use convolution  $k5n60s2p2$  to obtain the middle-scale feature map, and convolution  $k3n60s2p1$  to obtain the low branch’s feature map from the aforementioned feature map. The whole process can be formed as Eq. 1.

$$\begin{aligned} F_h &= \text{Conv}_{k7n60s1p3}(I_l) \\ F_m &= \text{Conv}_{k5n60s2p2}(I_l) \\ F_l &= \text{Conv}_{k3n60s2p1}(F_m) \end{aligned} \quad (1)$$

where  $I_l \in \mathbb{R}^{H \times W \times C}$  is the compressed low resolution image and  $H, W, C$  refer to its height, width and color channel, respectively.  $F_h, F_m, F_l$  represent the features of three branches.

Through this process, we obtain the hierarchical features  $\{F_h, F_m, F_l\}$  at different scales. Then, we will input them into the feature enhancement module to process in a divide-and-conquer manner.

### 3.2 Feature Enhancement and Fusion

The feature enhancement module and feature fusion module are the important components of HST. We will clarify them carefully in this section.

**Feature enhancement module** Different from previous hierarchical networks [33,26,56,39] for image restoration and super-resolution, where convolutional neural network (CNN) is used as the feature enhancement module for each branch, we use swin transformer architecture as ours. As proved by [28], swin transformer-based architecture can model long-range dependency enabled by the shifted window mechanism. Therefore, this architecture is more suitable for difficult degradation removal tasks, *e.g.*, compressed image super-resolution. Moreover, with the help of swin transformer, we can get better performance with less parameters.

Specifically, we directly apply multiple residual swin transformer blocks (RSTB) from [28], as our feature enhancement module. The architecture of RSTB is shown in Fig. 3. Each RSTB is composed of several swin transformer layers (STL), a convolution layer and a residual skip connection. This process can be formulated as Eq. 2

$$\begin{aligned} F_0 &= F_{in} \\ F_i &= \text{STL}(F_{i-1}), \quad i = 1, 2, \dots, K \\ F_{out} &= \text{Conv}(F_K) + F_0 \end{aligned} \quad (2)$$

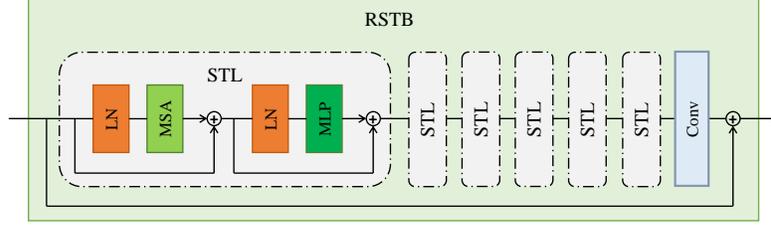


Fig. 3: Network structure of residual swin transformer block (RSTB)

where  $F_{in}$  is the input feature of one STL layer, and  $STL(\cdot)$  means each STL layer inside RSTB, which can be formulated as Eq. 3

$$\begin{aligned} X &= MSA(LN(X)) \\ X &= MLP(LN(X)) \end{aligned} \quad (3)$$

where  $MSA(\cdot)$  stands for multi-head self-attention,  $MLP(\cdot)$  stands for a multi-layer perceptron with two fully-connected layers and GELU as activation, and  $LN(\cdot)$  stands for LayerNorm. Since STL is not our contribution and previous works have already proved the effectiveness of this module, we directly utilize the same architecture as is presented in [28].

**Feature fusion module** After getting the enhanced features  $F_l^*$  from the low-branch enhancement module, composed of several RSTB blocks. We will integrate it into the higher feature with the fusion module, which aims to bring the contextual information from low-scale to high-scale. To demonstrate the fusion process clearly, we take the fusion of the low-branch feature and the middle-branch feature as an example. As described in Eq. 4, the low-branch feature  $F_l$  is enhanced to  $F_l^*$  with low-branch feature enhancement module  $FEM_l$ . Then we concatenate the super-resolved low-branch feature  $F_l^*_{\uparrow 2}$  and middle-branch feature  $F_m$ , and exploit the convolution layer to fuse these two components. Finally, we can obtain the enhanced middle-branch feature  $F_m^*$  by passing the fused feature  $F_m$  into the middle-branch enhancement module  $FEM_m$ . It is worthy to notice that, the up-sampling operation is implemented with Pixelshuffle [42], which can bring more stable results. The fusion of middle branch and high-branch features are implemented in the same way in Eq. 2.

$$\begin{aligned} F_l^* &= FEM_l(F_l), \\ F_m &= Conv(F_m \textcircled{C} F_l^*_{\uparrow 2}) \\ F_m^* &= FEM_m(F_m) \end{aligned} \quad (4)$$

### 3.3 HR Reconstruction Module

Since the compressed image super-resolution in the competition requires the resolution of network output to be  $4\times$  higher than their input image, the HR reconstruction module aims to produce the final three-channel RGB high-resolution clean image with the enhanced high-branch feature  $F_h^*$ .

As shown in Fig. 2, this part is composed of two sub-pixel convolution layers, including two convolution layers and two PixelShuffle layers [42]. Following the previous SR works [30,57,28], we utilize two sub-pixel convolution layers for the  $4\times$  upsampling. Finally, a convolution layer is used to generate the output HR image.

### 3.4 Pretraining with SR

To further boost the capability of the network, We also explore one simple but effective pertaining strategy for compressed image super-resolution. It is noteworthy that pretraining with more relevant distortions can bring better knowledge transfer. Particularly, we select three SR tasks as the pretraining schemes, *i.e.*, traditional SR, two RealSR simulation methods from BSRGAN [55] and DRTL [22], and explore their effectiveness for compressed image super-resolution. The relevant experimental analyses are shown in Sec. 4.3, which demonstrates pretraining with RealSR simulations leads to promising results, especially with the simulation in DRTL [22].

### 3.5 Loss Functions

In order to enable our HST to be competent for the task of compressed image super-resolution, we first pretrain our network on the  $\times 4$  super-resolution task, and then finetune it for compressed image super-resolution. For the  $\times 4$  super-resolution pretraining, we optimize network parameters by minimizing the  $L_1$  pixel loss as:

$$\mathcal{L} = \|I_{SR} - I_{HR}\|_1, \quad (5)$$

where  $I_{SR}$  is obtained by passing low-resolution images through the network, and  $I_{HR}$  is the corresponding ground-truth HR image. For compressed image super-resolution, we optimize network parameters by minimizing the Charbonnier loss [7].

$$\mathcal{L} = \sqrt{\|I_{SR} - I_{HR}\|^2 + \epsilon}, \quad (6)$$

where  $\epsilon$  is set as default value  $10^{-9}$ .

## 4 Experiments

### 4.1 Datasets

We produce the experimental results in our paper with two training datasets, DIV2K [1] (including 800 high-resolution images) and Flickr2K [44] (including 2650 high-resolution images). In the competition AIM2022 [53], we also collect extra 746 high-resolution images from CLIC 2021 official website<sup>1</sup> as the additional training data, which is only used for the competition results in Sec. 4.6. For the testing stage in this paper, we adopt Set5 [3], Set14 [54], BSD100 [37], Urban100 [16], Manga109 [38] and DIV2K [1] validation as our testing datasets.

### 4.2 Implementation Details

We use a three-branch HST for our experiments. The channel numbers of three feature enhancement modules  $FEM_h$ ,  $FEM_m$ ,  $FEM_l$  are set to 60, 60, 60, respectively. The spatial resolution of the high branch is  $64 \times 64$ , and halved for each downscale branch. Following [28], we set the number of swin transformer layers (STLs) as 6 for all residual swin transformer blocks (RSTBs) in HST. We use 2, 4, and 6 RSTBs for low branch, middle branch and high branch, respectively. The window size is set to 8 throughout the experiment.

We train our HST using four NVIDIA 2080Ti GPUs, with a batch size of 16. We offline generate training image pairs by the MATLAB bicubic kernel, then add JPEG compression with specified quality factor through the OpenCV function. We randomly crop LR into  $64 \times 64$  patches for training. For data augmentation, we leverage random flipping and random rotation simultaneously. In the stage of pretraining, the total training iterations are set to 400K. We adopt Adam optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ , the initial learning rate is set to  $2e-4$  and reduced by half at [100K, 250K]. In the stage of finetuning, we load network parameters from the pretraining stage. We conduct experiments on four different compression levels, with quality factors at 40, 30, 20 and 10, respectively. The training is first finished on quality factor at 40, with the initial learning rate and total iterations as  $1e-4$  and 200K. And the learning rate is halved after 100K iterations. The rest of tasks are finetuned based on the first task (*i.e.*, quality factor at 40), with the learning rate as  $8e-5$  and total iterations as 100K.

### 4.3 Effects of different pretraining schemes

As discussed in Sec. 3.4, pretraining is crucial for compressed image super-resolution task. To find out the optimal pretraining scheme, we conduct an ablation study on four different strategies, including: without pretraining, pure bicubic  $\times 4$  pretraining, pretraining with RealSR simulation from BSRGAN [55], and pretraining with RealSR simulation from DRTL [22]. BSRGAN [55] uses

<sup>1</sup> <http://clic.compression.cc/2021/tasks/index.html>

Table 1: Quantitative comparison for ablation study of network pretraining scheme. Results are tested on  $\times 4$  with compression quality 10 on Urban100 [16] dataset in terms of PSNR/SSIM. Best performance are in red.

Task	Methods(PSNR/SSIM)			
	w/o	bicubic $\times 4$	BSRGAN [55]	DRTL [22]
$\times 4, Q=10$	19.70/0.5181	19.92/0.5301	20.04/0.5375	<b>20.06/0.5383</b>

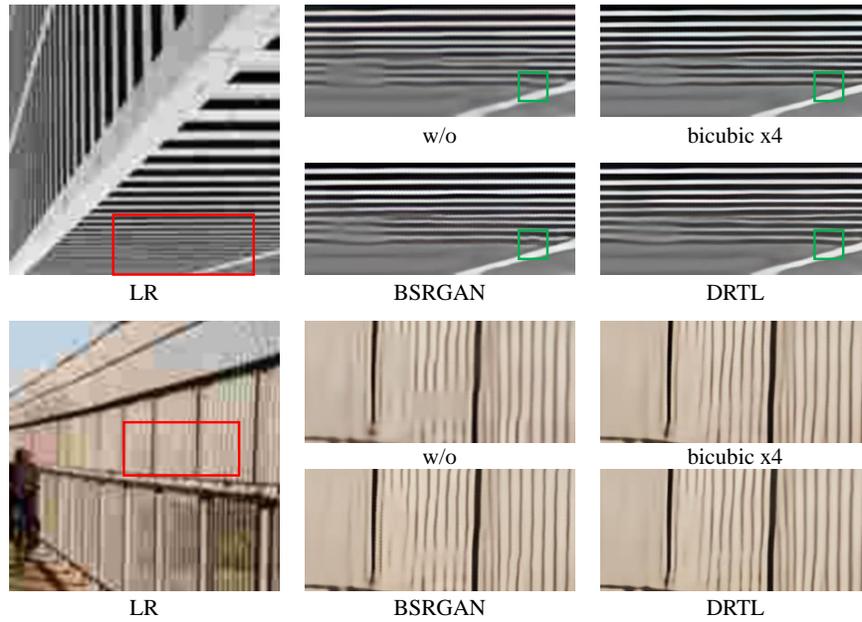


Fig. 4: Qualitative comparison for different pretraining schemes on  $\times 4$  image super-resolution with compression quality 10. Testing images are “011” and “024” from Urban100 [16] respectively.

a practical complex degradation simulation process, which demonstrates its effectiveness on real-world distortion removal. DRTL [22] proposes a multi-task degradation training scheme, to simulate distortion in real-world scenarios, and works well on few-shot real-world image super-resolution problems. For fast convergence and convincing results, we use SwinIR-s [28] as a training model, and test network performance on Urban100 [16].

As shown in Table 1, quantitative results show that the pretraining with RealSR leads to a gain of 0.36dB/0.0202 on the test dataset, which reveals that the pretraining is vital for compressed image super-resolution. Another observation is that pretraining with RealSR can achieve a better performance compared with simple bicubic downsampling, especially with the simulation in

DRTL [22]. The reason for this might be that RealSR simulations contain lots of hybrid distortions, which are more complex and the knowledge is more likely to be transferred to the severely compressed image super-resolution task.

#### 4.4 Effects of hierarchical architecture

To explore the advantage of introducing a hierarchical network structure, we set network branches from 1 to 3 and observe their performances on  $\times 4$  super-resolution with compression quality 40. Note that, one branch framework is almost the same as SwinIR-M [28]. As shown in Table 2, more branches lead to higher performance. However, it also brings an increase of computational complexity. In this paper, we choose a three-branch HST to achieve the best performance. In addition, benefits from structure and texture information compensation from lower branches, three-branch HST can generate images with clearer lines and more structural components, as shown in Fig. 5

Table 2: Quantitative comparison for ablation study of network scales. The number of parameters is listed in the bracket. Results are tested on  $\times 4$  with compression quality 40 in terms of PSNR/SSIM. Best performance are in red.

Methods	Q	Datasets(PSNR/SSIM)				
		Set5	Set14	BSD100	Urban100	Manga109
HST-1(11.90M)	40	25.28/0.726	23.78/0.613	23.82/0.583	22.21/0.652	23.69/0.767
HST-2(12.98M)		25.35/0.727	23.82/0.614	23.84/0.584	22.21/0.651	23.78/0.769
HST-3(16.58M)		25.39/0.728	23.84/0.614	23.87/0.584	22.23/0.651	23.85/0.768

#### 4.5 Comparison with other frameworks

We compare our HST with two other state-of-the-art models in image super-resolution, and one real-SR method [46] for qualitative comparison. Among them, RRDB [47] uses residual in residual dense blocks to deepen the network structure, thus having the ability to better aggregate image structure and texture information from multi-levels. SwinIR [28] introduces transformer into image restoration tasks and outperforms previous CNN-based models. The performances are tested on Set5 [3], Set14 [54], BSD100 [37], urban100 [16], Manga109 [38], respectively, with PSNR and SSIM in RGB channels. Moreover, we also test three models' performance on AIM2022 [53] official validation dataset, which includes 100 images from the DIV2K validation dataset. Quantitative and qualitative results are shown in Table 3,4 and Fig. 6, respectively. We denote the model using a self-ensemble strategy [30] with \*.

Extensive experiments show that our HST outperforms other methods by 0.25dB at most on compressed image super-resolution tasks. Even without self-ensemble, HST can still achieve an increase of 0.16dB at most. As shown in Fig.

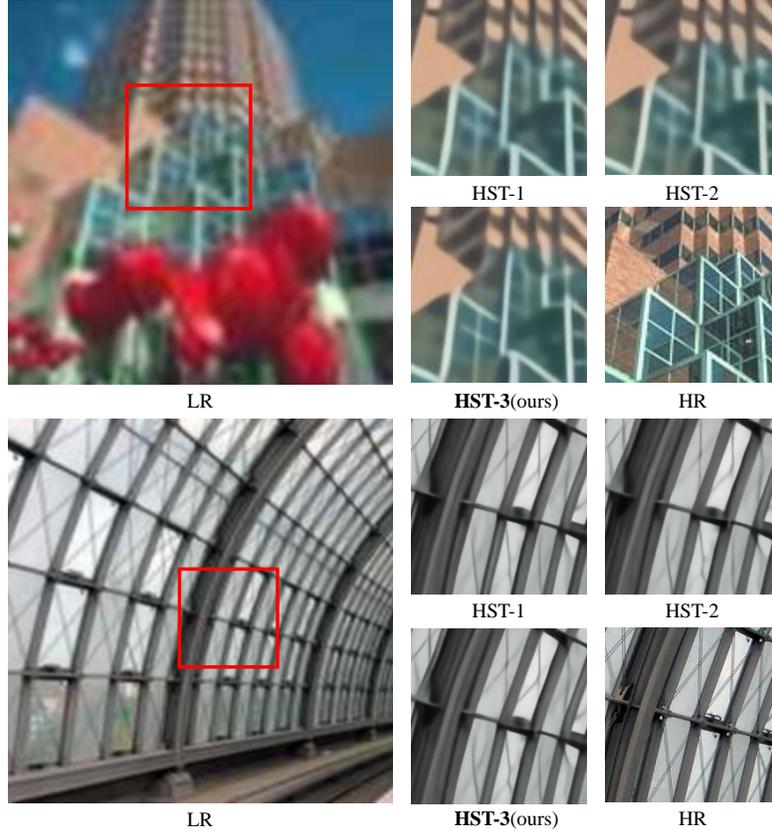


Fig. 5: Qualitative comparison for different network scales on  $\times 4$  image super-resolution with compression quality 40. Testing images are “095” from BSD100 [37] and “002” from Urban100 [16] respectively.

6, Real-ESRGAN [46] generates unnatural textures although its degradation process includes JPEG compression. Compared with other methods, our HST can generate SR with fewer artifacts. Moreover, HST performs better in rich texture areas, resulting in pleasant perception. All these benefit from a hierarchical network structure, which captures features at different scales and enhances the network’s representation ability.

#### 4.6 AIM2022 challenge

To further explore the performance of our HST, we follow the training process we used in AIM2022 [53] competition to train our HST. More specifically, we use

Table 3: Quantitative comparison for compressed image super-resolution on benchmark datasets. Results are tested on  $\times 4$  with different compression qualities in terms of PSNR/SSIM. Best performance are in **red**.

Methods	Q	Datasets(PSNR/SSIM)				
		Set5	Set14	BSD100	Urban100	Manga109
RRDB [47]	10	22.36/0.629	21.75/0.538	22.13/0.514	20.24/0.553	20.66/0.677
SwinIR [28]		22.45/0.636	21.79/0.541	22.16/0.517	20.35/0.561	20.81/0.685
<b>HST</b>		22.49/0.637	21.84/0.542	22.18/0.517	20.38/0.559	20.88/0.684
<b>HST*</b>		22.51/0.637	21.86/0.542	22.20/0.518	20.43/0.561	20.94/0.686
RRDB [47]	20	23.73/0.674	22.81/0.575	23.06/0.550	21.17/0.599	22.17/0.722
SwinIR [28]		23.81/0.682	22.87/0.577	23.09/0.551	21.32/0.608	22.35/0.729
<b>HST</b>		23.91/0.683	22.93/0.578	23.11/0.551	21.33/0.607	22.41/0.728
<b>HST*</b>		23.96/0.684	22.95/0.579	23.13/0.551	21.38/0.607	22.48/0.729
RRDB [47]	30	24.74/0.708	23.42/0.599	23.53/0.569	21.77/0.630	23.09/0.750
SwinIR [28]		24.83/0.713	23.43/0.600	23.53/0.571	21.85/0.636	23.20/0.755
<b>HST</b>		24.89/0.713	23.49/0.600	23.57/0.571	21.91/0.635	23.30/0.754
<b>HST*</b>		24.94/0.714	23.52/0.601	23.59/0.571	21.96/0.636	23.39/0.756
RRDB [47]	40	25.05/0.717	23.67/0.609	23.78/0.581	21.93/0.638	23.37/0.756
SwinIR [28]		25.28/0.726	23.78/0.613	23.82/0.583	22.21/0.652	23.69/0.767
<b>HST</b>		25.39/0.728	23.84/0.614	23.87/0.584	22.23/0.651	23.85/0.768
<b>HST*</b>		25.43/0.729	23.87/0.614	23.89/0.585	22.29/0.653	23.94/0.771

Table 4: Quantitative comparison for compressed image super-resolution on DIV2K [1] validation datasets. Results are tested on  $\times 4$  with different compression qualities in terms of PSNR/SSIM. Best performance are in **red**.

Datasets	Q	Methods(PSNR/SSIM)			
		RRDB	SwinIR	<b>HST</b>	<b>HST*</b>
DIV2K [1]	10	23.52/0.6400	23.57/0.6436	23.62/0.6436	23.65/0.6443
	20	24.68/0.6746	24.73/0.6771	24.77/0.6769	24.80/0.6777
	30	25.31/0.6949	25.32/0.6966	25.38/0.6963	25.41/0.6971
	40	25.58/0.7038	25.67/0.7077	25.74/0.7085	25.78/0.7093

all three training datasets described in Sec. 4.1 to finetune the network. After 100k iterations’ finetuning with Charbonnier Loss [7], we further use MSE Loss to optimize the network until convergence. The result on the official validation dataset shows that, with hierarchical network architecture, HST outperforms the one-branch network we used in the competition by 0.05dB, with a final PSNR of 23.80dB.

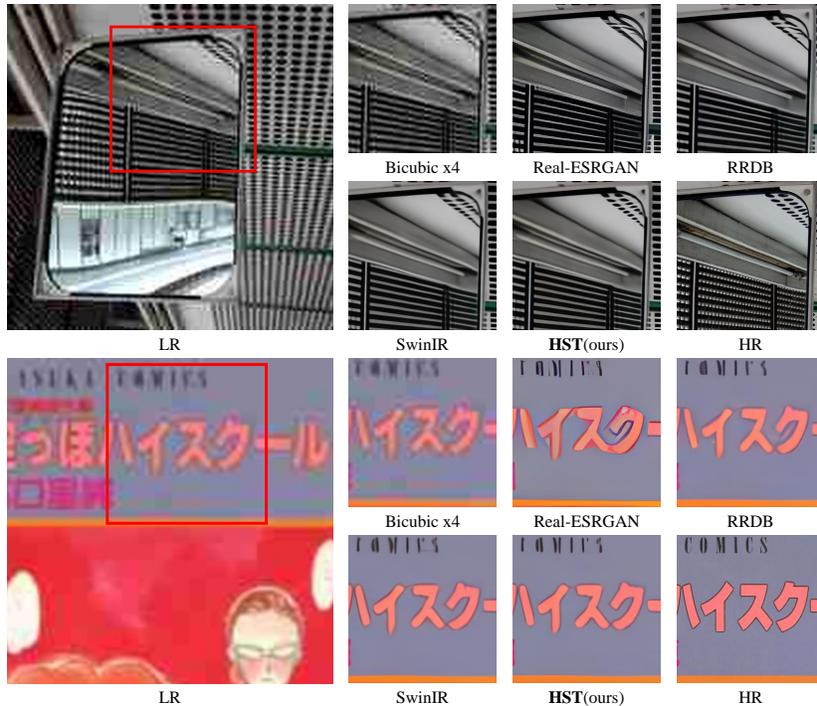


Fig. 6: Qualitative comparison for  $\times 4$  image super-resolution with compression quality 20. Testing images are “004” from Urban100 [16] and “KarappoHigh-school” from Manga109 [38], respectively.

## 5 Conclusion

In this paper, we propose the Hierarchical Swin Transformer for compressed image super-resolution, which incorporates the advantages of the hierarchical structure and Swin Transformer. Moreover, we find that pretraining with SR is vital and effective for compressed image super-resolution. Particularly, we explore three pretraining tasks, *i.e.*, traditional SR, and two RealSR simulations from BSRGAN and DRTL, respectively, of which the experimental results show that pretraining with RealSR simulations can bring better performance, especially with the simulation in DRTL [22]. Extensive experiments demonstrate that, with a pretraining and hierarchical network structure, our HST achieves the best performance on compressed image super-resolution tasks. In addition, our model achieves the fifth place in the AIM2022 challenge, with a PSNR of 23.51dB.

## References

1. Agustsson, E., Timofte, R.: Ntire 2017 challenge on single image super-resolution: Dataset and study. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 126–135 (2017)
2. Bell-Kligler, S., Shocher, A., Irani, M.: Blind super-resolution kernel estimation using an internal-gan. *Advances in Neural Information Processing Systems* **32** (2019)
3. Bevilacqua, M., Roumy, A., Guillemot, C., Alberi-Morel, M.L.: Low-complexity single-image super-resolution based on nonnegative neighbor embedding (2012)
4. Bross, B., Chen, J., Ohm, J.R., Sullivan, G.J., Wang, Y.K.: Developments in international video coding standardization after avc, with an overview of versatile video coding (vvc). *Proceedings of the IEEE* **109**(9), 1463–1493 (2021)
5. Cai, J., Zeng, H., Yong, H., Cao, Z., Zhang, L.: 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. pp. 3086–3095 (2019)
6. Cavigelli, L., Hager, P., Benini, L.: Cas-cnn: A deep convolutional neural network for image compression artifact suppression. In: 2017 International Joint Conference on Neural Networks (IJCNN). pp. 752–759. *IEEE* (2017)
7. Charbonnier, P., Blanc-Feraud, L., Aubert, G., Barlaud, M.: Two deterministic half-quadratic regularization algorithms for computed imaging. In: Proceedings of 1st International Conference on Image Processing. vol. 2, pp. 168–172. *IEEE* (1994)
8. 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)
9. Chen, Y., Fan, H., Xu, B., Yan, Z., Kalantidis, Y., Rohrbach, M., Yan, S., Feng, J.: Drop an octave: Reducing spatial redundancy in convolutional neural networks with octave convolution. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3435–3444 (2019)
10. 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)
11. Dong, C., Deng, Y., Loy, C.C., Tang, X.: Compression artifacts reduction by a deep convolutional network. In: Proceedings of the IEEE international conference on computer vision. pp. 576–584 (2015)
12. 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)
13. Fu, X., Zha, Z.J., Wu, F., Ding, X., Paisley, J.: Jpeg artifacts reduction via deep convolutional sparse coding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2501–2510 (2019)
14. Gu, J., Lu, H., Zuo, W., Dong, C.: Blind super-resolution with iterative kernel correction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1604–1613 (2019)
15. Guo, J., Chao, H.: Building dual-domain representations for compression artifacts reduction. In: European Conference on Computer Vision. pp. 628–644. Springer (2016)
16. 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)

17. Ji, X., Cao, Y., Tai, Y., Wang, C., Li, J., Huang, F.: Real-world super-resolution via kernel estimation and noise injection. In: proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops. pp. 466–467 (2020)
18. Jolicoeur-Martineau, A.: The relativistic discriminator: a key element missing from standard gan. arXiv preprint arXiv:1807.00734 (2018)
19. 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)
20. Kim, J., Choi, Y., Uh, Y.: Feature statistics mixing regularization for generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11294–11303 (2022)
21. 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)
22. Li, X., Jin, X., Fu, J., Yu, X., Tong, B., Chen, Z.: Few-shot real image restoration via distortion-relation guided transfer learning. arXiv preprint arXiv:2111.13078 (2021)
23. Li, X., Jin, X., Lin, J., Liu, S., Wu, Y., Yu, T., Zhou, W., Chen, Z.: Learning disentangled feature representation for hybrid-distorted image restoration. In: European Conference on Computer Vision. pp. 313–329. Springer (2020)
24. Li, X., Jin, X., Yu, T., Sun, S., Pang, Y., Zhang, Z., Chen, Z.: Learning omni-frequency region-adaptive representations for real image super-resolution. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 1975–1983 (2021)
25. Li, X., Shi, J., Chen, Z.: Task-driven semantic coding via reinforcement learning. *IEEE Transactions on Image Processing* **30**, 6307–6320 (2021)
26. Li, X., Sun, S., Zhang, Z., Chen, Z.: Multi-scale grouped dense network for vvc intra coding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 158–159 (2020)
27. Li, Y., Jin, P., Yang, F., Liu, C., Yang, M.H., Milanfar, P.: Comisr: Compression-informed video super-resolution. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2543–2552 (2021)
28. Liang, J., Cao, J., Sun, G., Zhang, K., Van Gool, L., Timofte, R.: Swinir: Image restoration using swin transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1833–1844 (2021)
29. Liang, J., Zhang, K., Gu, S., Van Gool, L., Timofte, R.: Flow-based kernel prior with application to blind super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10601–10610 (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. Liu, J., Li, X., Peng, Y., Yu, T., Chen, Z.: Swiniqa: Learned swin distance for compressed image quality assessment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1795–1799 (2022)
32. Liu, J., Lin, J., Li, X., Zhou, W., Liu, S., Chen, Z.: Lira: Lifelong image restoration from unknown blended distortions. In: European Conference on Computer Vision. pp. 616–632. Springer (2020)
33. Liu, P., Zhang, H., Lian, W., Zuo, W.: Multi-level wavelet convolutional neural networks. *IEEE Access* **7**, 74973–74985 (2019)

34. Lu, M., Chen, T., Liu, H., Ma, Z.: Learned image restoration for vvc intra coding. In: CVPR Workshops. p. 0 (2019)
35. Lu, Y., Fu, J., Li, X., Zhou, W., Liu, S., Zhang, X., Jia, C., Liu, Y., Chen, Z.: Rtn: Reinforced transformer network for coronary ct angiography vessel-level image quality assessment. arXiv preprint arXiv:2207.06177 (2022)
36. Luo, Z., Huang, H., Yu, L., Li, Y., Fan, H., Liu, S.: Deep constrained least squares for blind image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 17642–17652 (2022)
37. 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)
38. 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)
39. Pang, Y., Li, X., Jin, X., Wu, Y., Liu, J., Liu, S., Chen, Z.: Fan: Frequency aggregation network for real image super-resolution. In: European Conference on Computer Vision. pp. 468–483. Springer (2020)
40. Pennebaker, W.B., Mitchell, J.L.: JPEG: Still image data compression standard. Springer Science & Business Media (1992)
41. Rabbani, M., Joshi, R.: An overview of the jpeg 2000 still image compression standard. *Signal processing: Image communication* **17**(1), 3–48 (2002)
42. 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)
43. Svoboda, P., Hradis, M., Barina, D., Zemcik, P.: Compression artifacts removal using convolutional neural networks. arXiv preprint arXiv:1605.00366 (2016)
44. 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)
45. Wang, L., Wang, Y., Dong, X., Xu, Q., Yang, J., An, W., Guo, Y.: Unsupervised degradation representation learning for blind super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10581–10590 (2021)
46. 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)
47. 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)
48. Wang, Z., Liu, D., Chang, S., Ling, Q., Yang, Y., Huang, T.S.: D3: Deep dual-domain based fast restoration of jpeg-compressed images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2764–2772 (2016)
49. Wei, P., Lu, H., Timofte, R., Lin, L., Zuo, W., Pan, Z., Li, B., Xi, T., Fan, Y., Zhang, G., et al.: Aim 2020 challenge on real image super-resolution: Methods and results. In: European Conference on Computer Vision. pp. 392–422. Springer (2020)

50. Wei, P., Xie, Z., Lu, H., Zhan, Z., Ye, Q., Zuo, W., Lin, L.: Component divide-and-conquer for real-world image super-resolution. In: European Conference on Computer Vision. pp. 101–117. Springer (2020)
51. Wu, Y., Wang, X., Li, G., Shan, Y.: Animesr: Learning real-world super-resolution models for animation videos. arXiv preprint arXiv:2206.07038 (2022)
52. Wu, Y., Li, X., Zhang, Z., Jin, X., Chen, Z.: Learned block-based hybrid image compression. *IEEE Transactions on Circuits and Systems for Video Technology* (2021)
53. Yang, R., Timofte, R., et al.: AIM 2022 challenge on super-resolution of compressed image and video: Dataset, methods and results. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops (2022)
54. 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)
55. Zhang, K., Liang, J., Van Gool, L., Timofte, R.: Designing a practical degradation model for deep blind image super-resolution. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4791–4800 (2021)
56. Zhang, X., Yang, W., Hu, Y., Liu, J.: Dmccnn: Dual-domain multi-scale convolutional neural network for compression artifacts removal. In: 2018 25th IEEE International Conference on Image Processing (ICIP). pp. 390–394. IEEE (2018)
57. 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)
58. Zheng, M., Xing, Q., Qiao, M., Xu, M., Jiang, L., Liu, H., Chen, Y.: Progressive training of a two-stage framework for video restoration. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1024–1031 (2022)