

Swin2SR: SwinV2 Transformer for Compressed Image Super-Resolution and Restoration

Marcos V. Conde¹, Ui-Jin Choi², Maxime Burchi¹, and Radu Timofte¹

¹ Computer Vision Lab, CAIDAS, University of Würzburg, Germany
{marcos.conde-osorio,radu.timofte}@uni-wuerzburg.de
² MegaStudyEdu, South Korea

Abstract. Compression plays an important role on the efficient transmission and storage of images and videos through band-limited systems such as streaming services, virtual reality or videogames. However, compression unavoidably leads to artifacts and the loss of the original information, which may severely degrade the visual quality. For these reasons, quality enhancement of compressed images has become a popular research topic. While most state-of-the-art image restoration methods are based on convolutional neural networks, other transformers-based methods such as SwinIR, show impressive performance on these tasks. In this paper, we explore the novel Swin Transformer V2, to improve SwinIR for image super-resolution, and in particular, the compressed input scenario. Using this method we can tackle the major issues in training transformer vision models, such as training instability, resolution gaps between pre-training and fine-tuning, and hunger on data. We conduct experiments on three representative tasks: JPEG compression artifacts removal, image super-resolution (classical and lightweight), and compressed image super-resolution. Experimental results demonstrate that our method, Swin2SR, can improve the training convergence and performance of SwinIR, and is a top-5 solution at the “AIM 2022 Challenge on Super-Resolution of Compressed Image and Video”. Our code can be found at <https://github.com/mv-lab/swin2sr>.

Keywords: Super-Resolution, Image Compression, Transformer, JPEG

1 Introduction

Compression plays an important role on the efficient transmission and storage of images and videos through band-limited systems such as streaming services, virtual reality, cloud storage for images, videoconferences or videogames. However, compression leads to artifacts and the loss of the original information, which may severely degrade the visual quality of the image. For these reasons, quality enhancement and restoration of compressed images has become a popular research topic. Image restoration techniques, such as image super-resolution (SR) and JPEG compression artifact reduction, aim to reconstruct the high-quality clean image from its low-quality degraded (or compressed) counterpart. During the past decade, several revolutionary works were proposed for single image

super-resolution, most of them are CNN-based methods [17, 21, 29, 32, 55, 62–68]. We can also find plenty of proposed methods for the reduction of JPEG artifacts [19, 28, 46]. Recently, the blind super-resolution [23, 57, 63] methods have been proposed. They are able to use one model to jointly handle the tasks of super-resolution, deblurring, JPEG artifacts reduction, etc. Although the performance of these deep learning methods significantly improved compared with traditional methods [49], they generally suffer from two basic problems that arise from the basic convolution layer receptive field: (i) the interactions between images and kernels are content-independent, therefore, using the same kernel to restore different image regions may not be the best. (ii) Under the principle of locality, convolution is not effective for long-range dependency modelling [33].

As an alternative to CNNs, Transformer [53] designs a self-attention mechanism to capture global interactions between contexts and has shown promising performance in several vision problems [6, 18, 37, 51]. Recently, Swin Transformer [37] has shown great promise as it leverages the advantages of both CNN and Transformers (*i.e.* CNN to process image with large size due to the local attention mechanism, and transformer to model long-range dependency with the shifted window scheme). Compared with classical CNN-based image restoration models, Transformer-based methods have several benefits: (i) content-based interactions between image content and attention weights, which can be interpreted as spatially varying convolution [52]. (ii) long-range dependency modelling are enabled by the shifted window mechanism. (iii) in some cases, better performance with less parameters. In this context, Liang *et al.* SwinIR [33], based on Swin Transformer [37], represents the state-of-the-art of transformer-based models for image restoration.

AIM 2022 challenge on Super-Resolution of Compressed Image and Video This challenge is a step forward for establishing a new benchmark for the super-resolution of JPEG images and videos. The methods proposed in this challenge also have the potential to solve various super-resolution tasks. The challenge utilizes the famous DIV2K [1] dataset for evaluating methods. Other related challenges such as “NTIRE 2022 challenge on super-resolution and quality enhancement of compressed video” [58, 60] and “NTIRE 2020 challenge on real-world image SR” [38] also represent the SOTA in this field.

In this paper, we propose Swin2SR, a SwinV2 Transformer-based model [36, 37] for Compressed Image Super-Resolution and Restoration. This model represents a possible improvement or update of SwinIR [33] for these particular tasks. SwinV2 [36] (CVPR ’22) allows us to tackle the major issues in training large transformer-based vision models, including training instability and duration, and resolution gaps between pre-training and fine-tuning [33]. We are the first work to explore successfully other transformer blocks beyond Swin Transformer [37] for image super-resolution and restoration. In some scenarios, our model can achieve similar results as SwinIR [33], yet training 33% less.

We also provide extensive comparisons with state-of-the-art methods, and achieve competitive results at the related AIM 2022 Challenge.

2 Related Work

2.1 Image Restoration

Image restoration is split in a large number of sub-problems, for instance image denoising, image deblurring, super-resolution and compression artifacts removal among others. Traditional model-based methods for image restoration were usually defined by hand-crafted priors that narrowed the ill-posed nature of the problems by reducing the set of plausible solutions [12, 48, 49]. Learning-based methods based on CNNs have recently gained great popularity for image restoration, and they represent current state-of-the-art in most low-level vision tasks (*i.e.* denoising, deblurring, compression artifacts removal). The first remarkable work on denoising with deep learning is probably Zhang *et al.* [64] DnCNN. Other pioneering works include Dong *et al.* SRCNN [17] for image super-resolution and ARCNN [16] for JPEG compression artifact removal. Since research has moved towards deep learning, multiple CNN-based approaches have been proposed to improve the learned representations using more more complex neural network architectures, such as residual blocks, dense residual blocks, and laplacian operators [7, 29, 30, 62, 70, 71]. Other solutions attempt to exploit the attention mechanism in CNNs, such as channel attention and spatial attention [15, 34, 42, 43, 68].

2.2 Vision Transformer

The Transformer architecture [53] has recently gained much popularity in the computer vision community. Originally designed for neural machine translation, the Transformer architecture has successfully been applied to image classification [13, 14, 18, 37, 52], object detection [6, 51], object segmentation [4] and perceptual quality assessment (IQA) [10, 22]. The attention mechanism learns complex global interactions by attending to important regions in the image. Due to its impressive performance, transformers have also been introduced to image restoration [5, 8, 56]. More recently, Chen *et al.* [8] proposed IPT, a general backbone model for multiple image restoration tasks based on the standard Transformer [53]. This model shows promising performance on several tasks, however, it relies on a large number of parameters and heavy computation (over 115.5M parameters), and a large-scale dataset like ImageNet (over 1M images). VSR-Transformer proposed by Cao *et al.* [5] combines the self-attention mechanism and CNN-based feature extraction to fuse better features in video super-resolution. Note that many transformer-based approaches such as IPT [8] and VSR-Transformer [5] use patch-wise attention, which may not be optimal for image restoration. Liang *et al.* proposed SwinIR [33] based Swin Transformer [37], which represents the state-of-the-art in many restoration tasks.

In this context, the Swin Transformer [37] improved the Vision Transformer architecture by using shifted window based self-attention with progressive image downsampling like CNNs. Window self-attention is computed for non-overlapped image patches reducing attention computational complexity from Eq. 1 to Eq. 2:

$$O(MSA) = 4hwC^2 + 2(hw)^2C \quad (1)$$

$$O(WMSA) = 4hwC^2 + 2M^2hwC \quad (2)$$

for an image of size $h \times w$ and patches of size $M \times M$. The former quadratic computational complexity is replaced by a linear complexity when M is fixed. Learned relative positional bias are also added to include position information while computing similarities for each head.

The Swin Transformer V2 [36] modified the Swin Attention [37] module to better scale model capacity and window resolution. They first replace the *pre-norm* by a *post-norm* configuration, use a *scaled cosine attention* instead of the *dot product attention* and use a *log-spaced continuous* relative position bias approach to replace the previous *parameterized* approach. The attention output is:

$$\text{Attention}(Q, K, V) = \text{Softmax}(\cos(Q, K)/\tau + S)V \quad (3)$$

Where $Q, K, V \in \mathbb{R}^{M^2 \times d}$ are the query, key and value matrices. $S \in \mathbb{R}^{M^2 \times M^2}$ are the relative to absolute positional embeddings obtained by projecting the position bias after re-indexing. τ is a learnable scalar, non-shared across heads and layers. This block is illustrated in Figure 1.

3 Our Method

Our method Swin2SR is illustrated in Figure 1. We propose some modifications of SwinIR [33], which is based on Swin Transformer [37], that enhance the model’s capabilities for Super-Resolution, and in particular, for Compressed Input SR. We update the original Residual Transformer Block (RSTB) by using the new SwinV2 transformer [36] (CVPR’22) layers and attention to scale up capacity and resolution [36]. Our method has a classical upscaling branch which uses a bicubic interpolation, as shown in the AIM 2022 Challenge Leaderboard [59] and our results (5), this alone can recover basic structural information. For this reason, the output of our model is added to the basic upscaled image, to enhance it. We also explore different loss functions to make our model more robust to JPEG compression artifacts, being able to recover high-frequency details from the compressed LR image, and therefore, achieve better performance.

Advantages of updating to SwinV2 The SwinV2 architecture modifies the shifted window self-attention module to better scale model capacity and window resolution. The use of *post normalization* instead of *pre normalization* reduce the average feature variance of deeper layers and increase numerical stability during training. This allows to scale the SwinV2 Transformer up to 3 billion parameters without training instabilities [36]. The use of *scaled cosine attention* instead of *dot product* between queries and keys reduce the dominance of some attention heads for a few pixel pairs. In some tasks, our Swin2SR model achieved the same results as SwinIR [33], yet training 33% less iterations. Finally, the use of *log-spaced continuous* relative position bias allows us to generalize to higher input resolution at inference time.

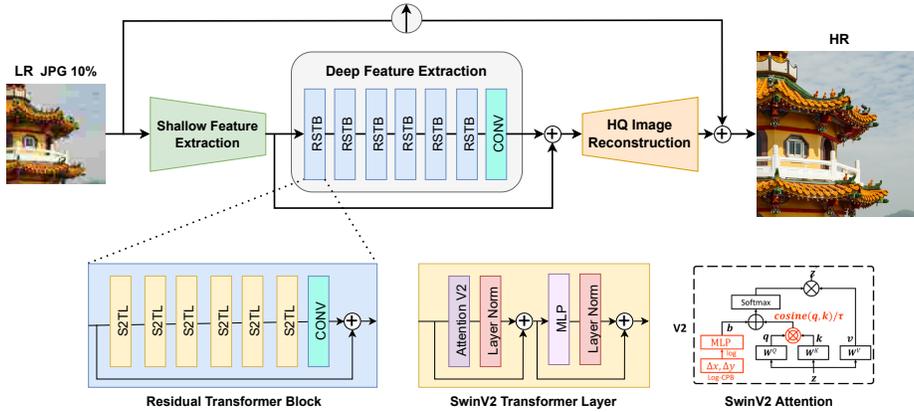


Fig. 1: The architecture of the proposed Swin2SR [11]. In this case, we show our method applied to Super-Resolution of Compressed Image [59].

3.1 Experimental Setup

For a fair comparison and ensure reproducibility, we follow the same experimental setup as SwinIR [33] and other state-of-the-art methods [63, 70].

We evaluate our model on three tasks: JPEG compression artifacts removal (Section 4.1), classical and lightweight image super-resolution (Section 4.2) and compressed image super-resolution (Section 4.4). We mainly use the DIV2K dataset for training and validation [1], and following the tradition of image SR, we report PSNR and SSIM on the Y channel of the YCbCr space [33, 63, 70].

Our model Swin2SR has the following elements, similar to SwinIR [33]: shallow feature extraction, deep feature extraction and high-quality image reconstruction modules. The **shallow feature extraction** module uses a convolution layer to extract features, which are directly transmitted to the reconstruction module to preserve low-frequency information [33, 64]. The **Deep feature extraction** module is mainly composed of Residual SwinV2 Transformer blocks (RSTB), each of which utilizes several SwinV2 Transformer [36] layers (S2TL) for local attention and cross-window interaction. Finally, both shallow and deep features are fused in the reconstruction module for high-quality image reconstruction. To upscale the image, we use standard a pixel shuffle operation.

The hyper-parameters of the architecture are as follows: the RSTB number, S2TL number, window size, channel number and attention head number are generally set to 6, 6, 8, 180 and 6, respectively. For lightweight image SR, we explain the details in Section 4.2.

3.2 Implementation details

The method was implemented in Pytorch using as baseline <https://github.com/cszn/KAIR> and the official repository for SwinIR [33]. We initially train

Swin2SR from scratch using the basic \mathcal{L}_1 loss for reconstruction. While training, we randomly crop HR images using 192px patch size and crop correspondingly the LR image generated offline using MATLAB, we also use standard augmentations that include all variations of flipping and rotations [50]. We use mainly the DIV2K [1]. In some experiments, to explore the potential benefits of more training data, we also use the Flickr2K dataset (2650 images).

In the particular scenario of **Compressed Input Super-Resolution** [59] (Section 4.4), we explore different loss functions to improve the performance and robustness of our method; these are represented in Figure 2.

First, we add an Auxiliary Loss that minimizes the \mathcal{L}_1 distance between the downsampled prediction \hat{y} and the downsampled reference y .png, as follows:

$$\mathcal{L}_{aux} = \|D(y) - D(\hat{y})\|_1 \quad (4)$$

where x is the low-resolution degraded image, y is the high-resolution clean image, $f(x) = \hat{y}$ is the restored image using our model f , and $D(\cdot)$ is a down-sampling operator (*i.e.* $\times 4$ bicubic kernel). This helps to ensure consistency also at lower-resolution. In order to minimize Eq. 4 the restored image at a lower resolution should not have artifacts (*i.e.* the prediction at lower resolution should be close to the downsampled reference .png without artifacts).

Second, we extract the high-frequency (HF) information from the High-Resolution images. This loss is formulated as follows:

$$\mathcal{L}_{hf} = \|(y - (y * b)) - (\hat{y} - (\hat{y} * b))\|_1 = \|HF(y) - HF(\hat{y})\|_1 \quad (5)$$

where $HR(\cdot)$ denotes the high-frequency information of an image. To obtain this, we convolve a simple 5×5 kernel b as a gaussian blur operation. This term enforces the prediction to have the same high-frequency details as the reference, and therefore, it helps to improve the sharpness and quality of the results.

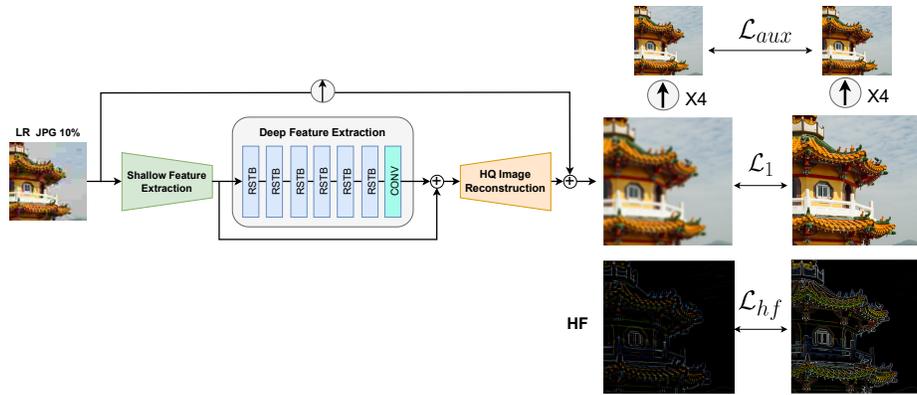


Fig. 2: Swin2SR training with additional regularization.

4 Experimental Results

4.1 JPEG Compression Artifacts Removal

Table 1 shows the comparison of Swin2SR with *state-of-the-art* JPEG compression artifact reduction methods: ARCNN [16], DnCNN-3 [64], QGAC [19], RNAN [69], and MWCNN [35]. All of compared methods are CNN-based models trained specifically for each quality type (*i.e.* four models per dataset). Due to our limited resources, and seeking for a more flexible approach, we train a single model able to deal with the four different quality factors. For this reason, we do not compare directly with DRUNet [62], as we consider it an unfair comparison. Moreover, Swin2SR only has 12M parameters, while DRUNet [62], is a large model that has 32.7M parameters. Note that we perform these comparisons using the same setup as [33]. Following [33, 62, 71], we test different methods on two benchmark datasets: (i) Classic5 [20] and (ii) LIVE1 [45]; using JPEG quality factors (q) 10, 20, 30 and 40. As we can see in Table 1, our Swin2SR achieves state-of-the-art results in compression artifacts removal.

Table 1: Quantitative comparison (average PSNR/SSIM) with state-of-the-art methods for **JPEG compression artifact reduction** on benchmark datasets. Best and second best performance are in red and blue colors, respectively. Note that Swin2SR is a single model that generalizes to different qualities, meanwhile, some methods are trained for each specific quality. Some numbers are from [28].

Dataset	q	ARCNN [16]	DnCNN [64]	QGAC [19]	RNAN [69]	MWCNN [35]	SwinIR [33]	Swin2SR
Classic5 [20]	10	29.03/0.79	29.40/0.80	29.84/0.83	29.96/0.81	30.01/0.82	30.27/0.82	30.02/0.81
	20	31.15/0.85	31.63/0.86	31.98/0.88	32.11/0.86	32.16/ 0.87	31.32/0.85	32.26/0.87
	30	32.51/0.88	32.91/0.88	33.22/0.90	33.38/0.89	33.43/0.89	31.39/0.853	33.51/0.89
	40	33.32/0.89	33.77/0.90	-	34.27/0.90	34.27/0.90	31.38/0.85	34.33/0.90
LIVE1 [45]	10	28.96/0.80	29.19/0.81	29.53/0.84	29.63/0.82	29.69/0.82	29.86/0.82	29.67/ 0.82
	20	31.29/0.87	31.59/0.88	31.86/0.90	32.03/0.88	32.04/0.89	31.00/0.86	32.07/0.89
	30	32.67/0.90	32.98/0.90	33.23/0.92	33.45/0.91	33.45/0.91	31.08/0.86	33.49/0.91
	40	33.63/0.91	33.96/0.92	-	34.47/0.92	34.45/0.93	31.05/0.86	34.49/0.92

In the case of SwinIR [33], which is also *state-of-the-art* for JPEG artifacts reduction, authors train one model per quality factor (*i.e.* four models) for 1600K iterations, and $q = 10/20/30$ models are fine-tuned using the $q = 40$ model as general baseline. We train a single model using the same setup [33], only for 800k iterations (*i.e.* $\times 2$ less training than SwinIR [33]), and JPEG compression as an augmentation. For this reason in Table 1 we compare with SwinIR trained for the most challenging $q = 10$. We also compare with MWCNN [35], IDCN [72] and FBCNN-C [28] using RGB color images. Attending to Tables 1 and 2, we consider our model a more general and flexible approach for grayscale or color compression artifacts removal, since it can be trained faster and generalizes to different compression quality factors. We also provide **qualitative results** in Figure 3. Swin2SR can restore compressed images and generate high-quality results. We provide additional results in the supplementary material.

Table 2: Quantitative comparison on **color** JPEG images with **single** compression. We report average PSNR/SSIM on benchmark datasets. Our model outperforms networks designed for this particular task (although we recognise that training with more data). Some numbers are from [28].

Dataset	q	JPEG	ARCNN [16]	QGAC [19]	MWCNN [35]	IDCN [72]	FBCNN-C [28]	Swin2SR
LIVE1 [45]	10	25.69/0.74	26.91/0.79	27.62/0.80	27.45/0.80	27.63/0.81	27.77/0.80	27.98/0.82
	40	30.28/0.88	-	32.05/0.91	-	-	32.34/0.91	32.53/0.92
ICB [44]	10	29.44/0.75	30.06/0.77	32.06/0.81	30.76/0.77	31.71/0.80	32.18/0.81	32.46/0.81
	40	33.95/0.84	-	32.25/0.91	-	-	36.02/0.86	36.25/0.86

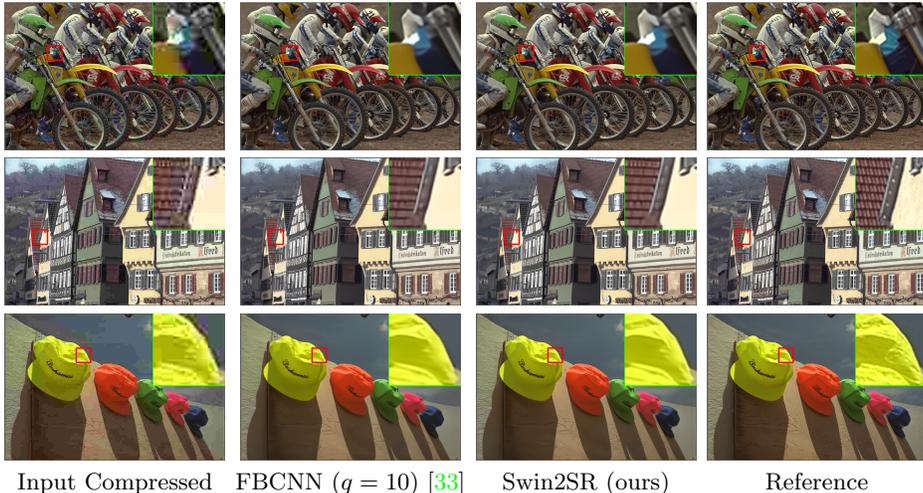


Fig. 3: Qualitative samples of JPEG Compression Artifacts Removal. We show the JPEG compressed image at quality $q = 10$. All images have the same resolution. Images from Classic5 [20] and LIVE1 [45]. Best viewed by zooming.

4.2 Classical Image Super-Resolution

For classical and lightweight image SR, following [33,62,63], we train Swin2SR on 800 training images of DIV2K and 2650 images from Flickr2K. For fair comparison with SwinIR [33], we use 64×64 LQ image patches, and the HQ-LQ image pairs are obtained by the MATLAB bicubic kernel. We train our model from scratch during 500k iterations, and fine-tune it for the $\times 4$ task. Table 3 shows the quantitative comparisons between Swin2SR and *state-of-the-art methods*: DBPN [24], RCAN [68], RRDB [55], SAN [15], IGNN [73], HAN [43], NLSA [42], IPT [8] and SwinIR [33]. All the CNN-based methods perform worse than the studied transformer-based methods, IPT [8], SwinIR [33] and Swin2SR. Moreover, Swin2SR was trained using only DIV2K+Flickr2K and achieves better performance than IPT [8], even though IPT [8] utilizes ImageNet (more than 1.3M images) in training and has huge number of parameters (115.5M). In con-

Table 3: Quantitative comparison (average PSNR/SSIM) with state-of-the-art methods for **classical image SR** on benchmark datasets. Best and second best performance are in **red** and **blue** colors, respectively.

Method	Scale	Training Dataset	Set5 [3]		Set14 [61]		BSD100 [40]		Urban100 [25]		Manga109 [41]	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
RCAN [68]	$\times 2$	DIV2K	38.27	0.9614	34.12	0.9216	32.41	0.9027	33.34	0.9384	39.44	0.9786
SAN [15]	$\times 2$	DIV2K	38.31	0.9620	34.07	0.9213	32.42	0.9028	33.10	0.9370	39.32	0.9792
IGNN [73]	$\times 2$	DIV2K	38.24	0.9613	34.07	0.9217	32.41	0.9025	33.23	0.9383	39.35	0.9786
HAN [43]	$\times 2$	DIV2K	38.27	0.9614	34.16	0.9217	32.41	0.9027	33.35	0.9385	39.46	0.9785
NLSA [42]	$\times 2$	DIV2K	38.34	0.9618	34.08	0.9231	32.43	0.9027	33.42	0.9394	39.59	0.9789
DBPN [24]	$\times 2$	DIV2K+Flickr2K	38.09	0.9600	33.85	0.9190	32.27	0.9000	32.55	0.9324	38.89	0.9775
IPT [8]	$\times 2$	ImageNet	38.37	-	34.43	-	32.48	-	33.76	-	-	-
SwinIR [33]	$\times 2$	DIV2K+Flickr2K	38.42	0.9623	34.46	0.9250	32.53	0.9041	33.81	0.9427	39.92	0.9797
Swin2SR	$\times 2$	DIV2K+Flickr2K	38.43	0.9623	34.48	0.9256	32.54	0.905	33.89	0.9431	39.88	0.9798
Swin2SR-D	$\times 2$	DIV2K+Flickr2K	38.06	-	33.81	-	32.32	-	32.6	-	38.98	-
RCAN [68]	$\times 4$	DIV2K	32.63	0.9002	28.87	0.7889	27.77	0.7436	26.82	0.8087	31.22	0.9173
SAN [15]	$\times 4$	DIV2K	32.64	0.9003	28.92	0.7888	27.78	0.7436	26.79	0.8068	31.18	0.9169
IGNN [73]	$\times 4$	DIV2K	32.57	0.8998	28.85	0.7891	27.77	0.7434	26.84	0.8090	31.28	0.9182
HAN [43]	$\times 4$	DIV2K	32.64	0.9002	28.90	0.7890	27.80	0.7442	26.85	0.8094	31.42	0.9177
NLSA [42]	$\times 4$	DIV2K	32.59	0.9000	28.87	0.7891	27.78	0.7444	26.96	0.8109	31.27	0.9184
DBPN [24]	$\times 4$	DIV2K+Flickr2K	32.47	0.8980	28.82	0.7860	27.72	0.7400	26.38	0.7946	30.91	0.9137
IPT [8]	$\times 4$	ImageNet	32.64	-	29.01	-	27.82	-	27.26	-	-	-
RRDB [55]	$\times 4$	DIV2K+Flickr2K	32.73	0.9011	28.99	0.7917	27.85	0.7455	27.03	0.8153	31.66	0.9196
SwinIR [33]	$\times 4$	DIV2K+Flickr2K	32.92	0.9044	29.09	0.7950	27.92	0.7489	27.45	0.8254	32.03	0.9260
Swin2SR	$\times 4$	DIV2K+Flickr2K	32.92	0.9039	29.06	0.7946	27.92	0.7505	27.51	0.8271	31.03	0.9256
Swin2SR-D	$\times 4$	DIV2K+Flickr2K	32.41	-	28.75	-	27.69	-	26.4	-	30.96	-

trast, Swin2SR has only 12M parameters, which is competitive even compared with state-of-the-art CNN-based models (15.4~44.3M). Note that our models achieve essentially the same performance as SwinIR [33], yet trained for 400k iterations from scratch, without fine-tuning or pre-training, in comparison with SwinIR [33] models trained during 500k, and in the case of $\times 4$ fine-tuned using the $\times 2$ model. We provide visual comparisons in Figures 5. Swin2SR can remove artifacts and recover structural information and high-frequency details.

Dynamic Super-Resolution Likewise Section 4.1, we explore the performance of a single super-resolution model to upscale directly using any arbitrary \times factor. We call this a Dynamic Super-Resolution model, referred as Swin2SR-D.

In SwinIR [33] we can find an upsampling layer designed to upscale images using s particular factor (*i.e.* $\times 2$). This layer cannot be adjusted to a different factor on-line, therefore, SwinIR [33] trains one model for each different factor. To deal with this problem, we implemented a Dynamic upsampling layer, which initially can super-resolve the images using $\times 2$, $\times 3$, and $\times 4$ factors on-line in the same module. We show in Table 3 the potential of this method, as this single model can perform $\times 2$ and $\times 4$ super-resolution indistinctly.

Lightweight image SR. We also provide comparison of Swin2SR-s with *state-of-the-art methods* lightweight image SR methods: CARN [2], FALSAR-A [9], IMDN [26], LAPAR-A [31], LatticeNet [39] and SwinIR (small) [33].

Our lightweight model is designed as SwinIR (small) [33], we decrease the number of Residual Swin Transformer Blocks (RSTB) and convolution channels to 4 and 60, respectively. However, the number of Swin Transformer Layers (STL) in each RSTB, window size and attention head number still set to 6, 8 and 6, respectively (as in Swin2SR base model).

In addition to PSNR and SSIM, we also report the total numbers of parameters and multiply-accumulate operations for different methods [33]. These MACs are calculated using a 1280×720 image. As shown in Table 4, Swin2SR outperforms competitive methods [2, 9, 26, 31] on different benchmark datasets, with similar total numbers of parameters and multiply-accumulate operations. In our experiments, Swin2SR can achieve the same results as SwinIR (small) [33], yet, training almost 33% less iterations.

Table 4: Quantitative comparison (average PSNR/SSIM) with state-of-the-art methods for **lightweight image SR $\times 2$** on benchmark datasets. Best and second best performance are in red and blue colors, respectively. In our experiments, Swin2SR-s converges faster than SwinIR (small) [33].

Method	# Params	# Mult-Adds	Set5 [3]		Set14 [61]		BSD100 [40]		Urban100 [25]		Manga109 [41]	
			PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
CARN [2]	1,592K	222.8G	37.76	0.9590	33.52	0.9166	32.09	0.8978	31.92	0.9256	38.36	0.9765
FALSR-A [9]	1,021K	234.7G	37.82	0.959	33.55	0.9168	32.1	0.8987	31.93	0.9256	-	-
IMDN [26]	694K	158.8G	38.00	0.9605	33.63	0.9177	32.19	0.8996	32.17	0.9283	38.88	0.9774
LAPAR-A [31]	548K	171.0G	38.01	0.9605	33.62	0.9183	32.19	0.8999	32.10	0.9283	38.67	0.9772
LatticeNet [39]	756K	169.5G	38.15	0.9610	33.78	0.9193	32.25	0.9005	32.43	0.9302	-	-
SwinIR [33]	878K	195.6G	38.14	0.9611	33.86	0.9206	32.31	0.9012	32.76	0.9340	39.12	0.9783
Swin2SR-s	1000K	199.0G	38.17	0.9613	33.95	0.9216	32.35	0.9024	32.85	0.9349	39.32	0.9787

4.3 Real-world Image Super-Resolution

We also test our approach using real-world images and prove the generalization capabilities of Swin2SR. We use the same setup as SwinIR [33] for training and testing our methods to exploit the full potential of these transformer-based approaches. Since there is no ground-truth high-quality images, we only provide visual comparison with representative bicubic model in Figure 4. Our model produces detailed images without artifacts. Due to the limitations of space and visualization in this document, we include the comparison with ESRGAN [55] and state-of-the-art real-world image SR models such as RealSR [27], BSRGAN [63], Real-ESRGAN [54] and SwinIR [33] in the supplementary material.

4.4 Compressed Image Super-Resolution

The ‘‘AIM 2022 Challenge on Super-Resolution of Compressed Image’’ [59] is a step forward for establishing a benchmark of the super-resolution of JPEG images. In this challenge, we use the popular dataset DIV2K [1] as the training, validation and test sets. JPEG is the most commonly used image compression

standard. We target the $\times 4$ super-resolution of the images compressed with JPEG with the quality factor of 10. Figure 1 illustrates this process. We propose two solutions for this problem based on previous Sections 4.1 and 4.2:

1. **Swin2SR-CI** An end-to-end model for JPEG artifacts removal and super-resolution (*i.e.* Figure 1).
2. A 2-stage approach where first we remove JPEG compression artifacts in the LR input image using **Swin2SR-DJPEG**, and second, we upscale using **Swin2SRx4** (*i.e.* the model trained for Classical SR, Section 4.2). We refer to this experiment as “Swin2SR-CI2”.

As we show in Table 5⁽³⁾, our method is a top solution at the challenge. We trained Swin2SR using only DIV2K [1] and Flickr2K [47] datasets, in comparison with other teams like CASIA LCVG, which trained using 1 million images. Our average testing time of Swin2SR model is 1.41s using single GPU A100.

In Figure 5 we show extensive qualitative results of compressed input super-resolution [59]. Our model can recover information from the low-quality low-resolution input image, and generates high-resolution high-quality images. Among the limitations of our model, we can appreciate a clear blur effect, nevertheless, we find SwinIR [33] (and other *state-of-the-art* methods) to have the same issues.

Table 5: Results of AIM 2022 Challenge on Super-Resolution of Compressed Image. Our solutions are placed among the top teams, while our methods can process a single image in under a second (w/o self-ensemble).

Team	Test PSNR (dB)	Runtime (s)	Hardware
VUE	23.6677	120	Tesla V100
BSR	23.5731	63.96	Tesla A100
CASIA LCVG	23.5597	78.09	Tesla A100
USTC-IR	23.5085	19.2	2080ti
Swin2SR-CI2	23.4946	24	Tesla A100
MSDRSR	23.4545	7.94	Tesla V100
Giantpandacv	23.4249	0.248	RTX 3090
Swin2SR-CI	23.4033	9.39	Tesla A100
MVideo	23.3250	1.7	RTX 3090
UESTC+XJU CV	23.2911	3.0	RTX 3090
cvlab	23.2828	6.0	1080 Ti
Bicubic $\times 4$	22.2420	-	-

Ensembles and fusion strategies. We use classical self-ensemble techniques where the input image is flipped and rotated several times, and the resultant images are averaged [38, 50]. We only use this technique in the related

³ online leaderboard <https://codalab.lisn.upsaclay.fr/competitions/5076>

AIM 2022 Challenge (Section 4.4 and Table 5), and the marginal improvement of this technique was approximately 0.02dB PSNR.

In Table 6 we show our ablation studies using the challenge DIV2K [1] validation set. The use of additional loss functions helped the model to converge faster, however after certain number of iterations (*i.e.* 250k) the model converges. As previously mentioned, among the **limitations** of our model, we can appreciate a clear blur effect in the qualitative samples in Figure 5, indicating that our model is struggling to recover fine details and sharpness. Nevertheless, we find SwinIR [33] (and other *state-of-the-art* methods) to have the same issues to recover the high-frequency details. However, the overall results look very impressive considering the level of degradation of the input image (downsampled and compressed using JPEG at quality $q = 10$). We also provide additional results and samples for DIV2K [1] in the supplementary material.

Table 6: Ablation study of our experiments in the AIM 2022 Compressed Image Super-Resolution Challenge. The additional loss functions, and our new design Swin2SR help to converge faster and produce competitive results. Note that we compare with SwinIR pre-trained model while we trained using only the challenge DIV2K [1] data.

Exp.	Method	PSNR
1	Bicubic	22.350
2	RDN [70]	23.320
3	SwinIR [33]	23.546
4	Swin2SR (Ours)	23.580
5	Swin2SR + AuxLoss	23.585
6	Swin2SR + AuxLoss + HFLoss	23.590
7	Self-ensemble Exp6	23.616

5 Conclusion

In this paper we propose Swin2SR, a SwinV2 Transformer-based model for super-resolution and restoration of compressed images. This model is a possible improvement of SwinIR (based on Swin Transformer), allowing faster training and convergence, and bigger capacity and resolution. Extensive experiments show that Swin2SR achieves state-of-the-art performance on: JPEG compression artifacts removal, image super-resolution (classical and lightweight), and compressed image super-resolution. Our method also achieves competitive results at the “AIM 2022 Challenge on Super-Resolution of Compressed Image and Video”, being ranked among the top-5, and therefore, it helps to advance the state-of-the-art in super-resolution of compressed inputs, which will play an essential role in industries like streaming services, virtual reality or video games.

Acknowledgments This work was partly supported by The Alexander von Humboldt Foundation (AvH).



Fig. 4: Qualitative results on **real-world** SR datasets (RealSRSet, 5images). Our model can recover textures, remove noise and produce pleasant results.

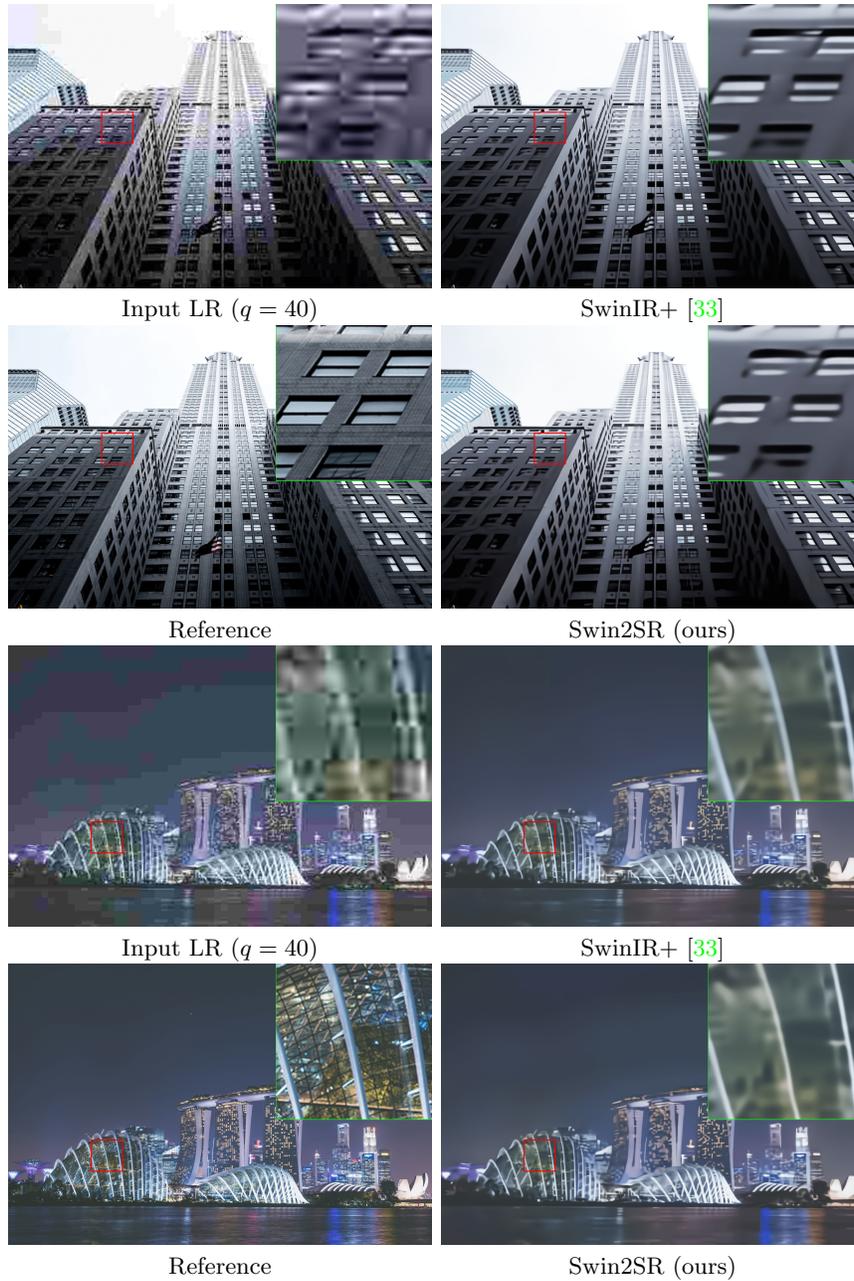


Fig. 5: Qualitative samples from the AIM 2022 Challenge on Super-Resolution of Compressed Image. Validation images from the DIV2K [1].

References

1. Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. [2](#), [5](#), [6](#), [10](#), [11](#), [12](#), [14](#)
2. Namhyuk Ahn, Byungkong Kang, and Kyung-Ah Sohn. Fast, accurate, and lightweight super-resolution with cascading residual network. In *European Conference on Computer Vision*, pages 252–268, 2018. [9](#), [10](#)
3. Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie-Line Alberi Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. In *British Machine Vision Conference*, pages 135.1–135.10, 2012. [9](#), [10](#)
4. Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. *arXiv preprint arXiv:2105.05537*, 2021. [3](#)
5. Jie Zhang Cao, Yawei Li, Kai Zhang, and Luc Van Gool. Video super-resolution transformer. *arXiv preprint arXiv:2106.06847*, 2021. [3](#)
6. Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. [2](#), [3](#)
7. Lukas Cavigelli, Pascal Hager, and Luca Benini. Cas-cnn: A deep convolutional neural network for image compression artifact suppression. In *2017 International Joint Conference on Neural Networks*, pages 752–759, 2017. [3](#)
8. Hanqing Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 12299–12310, 2021. [3](#), [8](#), [9](#)
9. Xiangxiang Chu, Bo Zhang, Hailong Ma, Ruijun Xu, and Qingyuan Li. Fast, accurate and lightweight super-resolution with neural architecture search. In *International Conference on Pattern Recognition*, pages 59–64. IEEE, 2020. [9](#), [10](#)
10. Marcos V Conde, Maxime Burchi, and Radu Timofte. Conformer and blind noisy students for improved image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 940–950, 2022. [3](#)
11. Marcos V Conde, Ui-Jin Choi, Maxime Burchi, and Radu Timofte. Swin2SR: Swinv2 transformer for compressed image super-resolution and restoration. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2022. [5](#)
12. Marcos V. Conde, Steven McDonagh, Matteo Maggioni, Ales Leonardis, and Eduardo Pérez-Pellitero. Model-based image signal processors via learnable dictionaries. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(1):481–489, Jun. 2022. [3](#)
13. Marcos V Conde and Kerem Turgutlu. Clip-art: contrastive pre-training for fine-grained art classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3956–3960, 2021. [3](#)
14. Marcos V Conde and Kerem Turgutlu. Exploring vision transformers for fine-grained classification. *arXiv preprint arXiv:2106.10587*, 2021. [3](#)
15. Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11065–11074, 2019. [3](#), [8](#), [9](#)

16. Chao Dong, Yubin Deng, Chen Change Loy, and Xiaoou Tang. Compression artifacts reduction by a deep convolutional network. In *IEEE International Conference on Computer Vision*, pages 576–584, 2015. [3](#), [7](#), [8](#)
17. Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European Conference on Computer Vision*, pages 184–199, 2014. [2](#), [3](#)
18. Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [2](#), [3](#)
19. Max Ehrlich, Larry Davis, Ser-Nam Lim, and Abhinav Shrivastava. Quantization guided jpeg artifact correction. In *European Conference on Computer Vision*, pages 293–309, 2020. [2](#), [7](#), [8](#)
20. Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Pointwise shape-adaptive dct for high-quality denoising and deblocking of grayscale and color images. *IEEE Transactions on Image Processing*, 16(5):1395–1411, 2007. [7](#), [8](#)
21. Manuel Fritsche, Shuhang Gu, and Radu Timofte. Frequency separation for real-world super-resolution. In *IEEE Conference on International Conference on Computer Vision Workshops*, pages 3599–3608, 2019. [2](#)
22. Jinjin Gu, Haoming Cai, Chao Dong, Jimmy S Ren, Radu Timofte, Yuan Gong, Shanshan Lao, Shuwei Shi, Jiahao Wang, Sidi Yang, et al. Ntire 2022 challenge on perceptual image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 951–967, 2022. [3](#)
23. Jinjin Gu, Hannan Lu, Wangmeng Zuo, and Chao Gu. Blind super-resolution with iterative kernel correction. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1604–1613, 2019. [2](#)
24. Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Deep back-projection networks for super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1664–1673, 2018. [8](#), [9](#)
25. Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5197–5206, 2015. [9](#), [10](#)
26. Zheng Hui, Xinbo Gao, Yunchu Yang, and Xiumei Wang. Lightweight image super-resolution with information multi-distillation network. In *ACM International Conference on Multimedia*, pages 2024–2032, 2019. [9](#), [10](#)
27. Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 466–467, 2020. [10](#)
28. Jiayi Jiang, Kai Zhang, and Radu Timofte. Towards flexible blind jpeg artifacts removal. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4997–5006, 2021. [2](#), [7](#), [8](#)
29. Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1646–1654, 2016. [2](#), [3](#)
30. Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 624–632, 2017. [3](#)
31. Wenbo Li, Kun Zhou, Lu Qi, Nianjuan Jiang, Jiangbo Lu, and Jiaya Jia. Lapa: Linearly-assembled pixel-adaptive regression network for single image super-resolution and beyond. *arXiv preprint arXiv:2105.10422*, 2021. [9](#), [10](#)

32. Zhen Li, Jinglei Yang, Zheng Liu, Xiaomin Yang, Gwanggil Jeon, and Wei Wu. Feedback network for image super-resolution. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3867–3876, 2019. 2
33. Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1833–1844, 2021. 2, 3, 4, 5, 7, 8, 9, 10, 11, 12, 14
34. Ding Liu, Bihan Wen, Yuchen Fan, Chen Change Loy, and Thomas S Huang. Non-local recurrent network for image restoration. *arXiv preprint arXiv:1806.02919*, 2018. 3
35. Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 773–782, 2018. 7, 8
36. Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12009–12019, 2022. 2, 4, 5
37. Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 2, 3, 4
38. Andreas Lugmayr, Martin Danelljan, and Radu Timofte. Ntire 2020 challenge on real-world image super-resolution: Methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 494–495, 2020. 2, 11
39. Xiaotong Luo, Yuan Xie, Yulun Zhang, Yanyun Qu, Cuihua Li, and Yun Fu. Latticenet: Towards lightweight image super-resolution with lattice block. In *European Conference on Computer Vision*, pages 272–289, 2020. 9, 10
40. David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *IEEE Conference on International Conference on Computer Vision*, pages 416–423, 2001. 9, 10
41. Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76(20):21811–21838, 2017. 9, 10
42. Yiqun Mei, Yuchen Fan, and Yuqian Zhou. Image super-resolution with non-local sparse attention. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3517–3526, 2021. 3, 8, 9
43. Ben Niu, Weilei Wen, Wenqi Ren, Xiangde Zhang, Lianping Yang, Shuzhen Wang, Kaihao Zhang, Xiaochun Cao, and Haifeng Shen. Single image super-resolution via a holistic attention network. In *European Conference on Computer Vision*, pages 191–207, 2020. 3, 8, 9
44. Rawzor. Image compression benchmark. 8
45. HR Sheikh. Live image quality assessment database release 2. <http://live.ece.utexas.edu/research/quality>, 2005. 7, 8
46. Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *IEEE International Conference on Computer Vision*, pages 4539–4547, 2017. 2
47. Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. Ntire 2017 challenge on single image super-resolution: Methods and re-

- sults. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 114–125, 2017. [11](#)
48. Radu Timofte, Vincent De Smet, and Luc Van Gool. Anchored neighborhood regression for fast example-based super-resolution. In *IEEE Conference on International Conference on Computer Vision*, pages 1920–1927, 2013. [3](#)
 49. Radu Timofte, Vincent De Smet, and Luc Van Gool. A+: Adjusted anchored neighborhood regression for fast super-resolution. In *Asian Conference on Computer Vision*, pages 111–126, 2014. [2](#), [3](#)
 50. Radu Timofte, Rasmus Rothe, and Luc Van Gool. Seven ways to improve example-based single image super resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1865–1873, 2016. [6](#), [11](#)
 51. Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. [2](#), [3](#)
 52. Ashish Vaswani, Prajit Ramachandran, Aravind Srinivas, Niki Parmar, Blake Hechtman, and Jonathon Shlens. Scaling local self-attention for parameter efficient visual backbones. *arXiv preprint arXiv:2103.12731*, 2021. [2](#), [3](#)
 53. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. [2](#), [3](#)
 54. Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. *arXiv preprint arXiv:2107.10833*, 2021. [10](#)
 55. Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *European Conference on Computer Vision Workshops*, pages 701–710, 2018. [2](#), [8](#), [9](#), [10](#)
 56. Zhendong Wang, Xiaodong Cun, Jianmin Bao, and Jianzhuang Liu. Uformer: A general u-shaped transformer for image restoration. *arXiv preprint arXiv:2106.03106*, 2021. [3](#)
 57. Mehmet Yamac, Baran Ataman, and Aakif Nawaz. Kernelnet: A blind super-resolution kernel estimation network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 453–462, 2021. [2](#)
 58. Ren Yang, Radu Timofte, et al. NTIRE 2021 challenge on quality enhancement of compressed video: Methods and results. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2021. [2](#)
 59. Ren Yang, Radu Timofte, 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 Workshops (ECCVW)*, 2022. [4](#), [5](#), [6](#), [10](#), [11](#)
 60. Ren Yang, Radu Timofte, et al. NTIRE 2022 challenge on super-resolution and quality enhancement of compressed video: Dataset, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2022. [2](#)
 61. Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *International Conference on Curves and Surfaces*, pages 711–730, 2010. [9](#), [10](#)
 62. Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [2](#), [3](#), [7](#), [8](#)

63. Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *IEEE Conference on International Conference on Computer Vision*, 2021. 2, 5, 8, 10
64. Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155, 2017. 2, 3, 5, 7
65. Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep cnn denoiser prior for image restoration. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3929–3938, 2017. 2
66. Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Transactions on Image Processing*, 27(9):4608–4622, 2018. 2
67. Kai Zhang, Wangmeng Zuo, and Lei Zhang. Learning a single convolutional super-resolution network for multiple degradations. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3262–3271, 2018. 2
68. Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *European Conference on Computer Vision*, pages 286–301, 2018. 2, 3, 8, 9
69. Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. *arXiv preprint arXiv:1903.10082*, 2019. 7
70. Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2472–2481, 2018. 3, 5, 12
71. Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image restoration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(7):2480–2495, 2020. 3, 7
72. Bolun Zheng, Yaowu Chen, Xiang Tian, Fan Zhou, and Xuesong Liu. Implicit dual-domain convolutional network for robust color image compression artifact reduction. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(11):3982–3994, 2019. 7, 8
73. Shangchen Zhou, Jiawei Zhang, Wangmeng Zuo, and Chen Change Loy. Cross-scale internal graph neural network for image super-resolution. *arXiv preprint arXiv:2006.16673*, 2020. 8, 9