

# DSRGAN: Detail Prior-Assisted Perceptual Single Image Super-Resolution via Generative Adversarial Networks

Ziyang Liu, Zhengguo Li, Xingming Wu, Zhong Liu, and Weihai Chen\*

**Abstract**—The generative adversarial network (GAN) is successfully applied to study the perceptual single image super-resolution (SISR). However, the GAN often tends to generate images with high frequency details being inconsistent with the real ones. Inspired by conventional detail enhancement algorithms, we propose a novel prior knowledge, the detail prior, to assist the GAN in alleviating this problem and restoring more realistic details. The proposed method, named DSRGAN, includes a well designed detail extraction algorithm to capture the most important high frequency information from images. Then, two discriminators are utilized for supervision on image-domain and detail-domain restorations, respectively. The DSRGAN merges the restored detail into the final output via a detail enhancement manner. The special design of DSRGAN takes advantages from both the model-based conventional algorithm and the data-driven deep learning network. Experimental results demonstrate that the DSRGAN outperforms the state-of-the-art SISR methods on perceptual metrics and achieves comparable results in terms of fidelity metrics simultaneously. Following the DSRGAN, it is feasible to incorporate other conventional image processing algorithms into a deep learning network to form a model-based deep SISR.

**Index Terms**—Single image super-resolution, generative adversarial networks, detail prior, model-based and data-driven

## I. INTRODUCTION

GENERATING a high-resolution (HR) image from one single low-resolution (LR) image, which is called single image super-resolution (SISR), is a classic and key problem in the field of image processing [1]. The generated image is referred to as the super-resolution (SR) image. Recently most SISR methods are deep learning-based. They can be divided into two categories according to their orientations to different image quality metrics: fidelity metrics such as the peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM), and perceptual metrics such as the perceptual index (PI) [2] and learned perceptual image patch similarity (LPIPS) [3].

The fidelity-oriented SISR methods are commonly formalized as the minimization of the mean squared error (MSE) or mean absolute error (MAE) between the recovered SR image and the real HR image [4]. Higher PSNR and SSIM have been achieved by these methods with the development of convolutional neural network (CNN) [5], [6], [7]. However, the minimization of the Euclidean distance between predicted

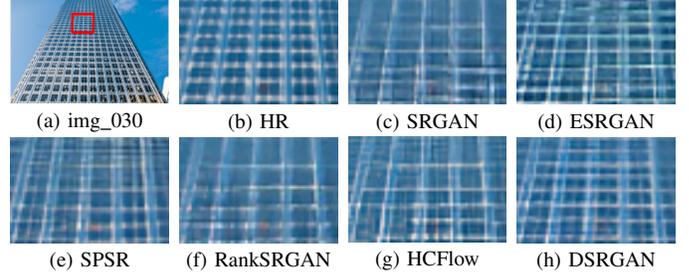


Fig. 1: SR results of different perception-oriented SISR methods. Our proposed DSRGAN preserves more sharp and geometric-consistent edges than existing methods including the SRGAN [8], ESRGAN [9], SPSR [10], RankSRGAN [11], and HCFlow [12].

and ground truth pixels will encourage the CNN to average all plausible outputs and produce blurry results, where edges and textures are smoothed.

The generative adversarial network (GAN) [13] has been developed to tackle this problem. The generator and discriminator in GAN are formalized to play a two-player min-max game. By playing the game, the generator network is capable of producing images with rich details, which are visual-pleasing and photo-realistic [8]. The GAN has been widely adopted to study the perception-oriented SISR [9], [10], [11]. Although these methods perform well on the PI and LPIPS, they often tend to produce distortions in SR images. Their performances on the PSNR and SSIM are unsatisfactory, which implies that the details generated by the generator are different from the real ones. As shown in Figures 1(c) and 1(d), a few edges restored by the SRGAN [8] and ESRGAN [9] are twisted. The reason is that the discriminator in GAN may introduce unstable factors to the optimization procedure [10]. In addition, the ill-posed nature of the SISR problem makes the generator difficult to restore details consistent with the real ones.

Since restoring the real high frequency details becomes the toughest problem in the SISR, a detail prior-assisted SISR paradigm via the GAN, named DSRGAN, is proposed to alleviate this issue. The proposal is inspired by a few conventional detail enhancement algorithms [14], [15], [16], [17], [18], [19], [20], [21], according to which a natural image can be decomposed into a base layer and a detail layer. The base layer is formed by homogeneous image regions with overly smoothed textures. The detail layer reveals the sharp edges of each local region and the rich textures within it. We pay more attention to the detail layer and exploit it to explicitly

\*Corresponding author: Weihai Chen, whchen@buaa.edu.cn.

Z. Liu, X. Wu, Z. Liu, and W. Chen are with School of Automation Science and Electrical Engineering, Beihang University, 100191, Beijing, China.

Z. Li is with SRO Department, Institute for Infocomm Research, Singapore, 138632, Singapore.

guide the perception-oriented SISR networks to enhance the SR performance.

Specifically, the DSRGAN includes a novel model-based detail extraction algorithm to extract the detail layer. The algorithm is able to capture the most important high frequency information of image, such as sharp edges and fine textures. It can also intelligently distinguish and exclude unwanted information such as noises, which obstruct the SISR reconstruction. Moreover, the algorithm contains a regularization term to ensure a sparse attribute of the detail layer, such that it can guide the network more efficiently. Once the detail layer is extracted from the LR image, it serves as a prior knowledge for the generator in DSRGAN. The DSRGAN includes two discriminators, one for supervision on image-domain LR-to-HR translation and the other for supervision on detail-domain LR-to-HR translation. Finally, the DSRGAN merges the restored detail into the SR image via a detail enhancement manner [14], [16]. The special design of DSRGAN takes advantages from both the model-based conventional algorithm and the data-driven deep learning network. As such, it is capable of restoring more realistic details, as shown in Figure 1(h). Notably, the DSRGAN is network-agnostic, which can be used for off-the-shelf SISR networks. Furthermore, the detail extraction algorithm is not required to be differentiable. Following the DSRGAN, it is also feasible to incorporate other conventional image processing algorithms into a deep network to form a model-based deep SISR, which will have a great potential research value. In summary, our main contributions are:

- A novel detail extraction algorithm is proposed to exploit the detail prior for explicitly guiding the perception-oriented SISR reconstruction.
- A detail prior-assisted SISR paradigm via the GAN, named DSRGAN, is proposed, which takes advantages from both the model-based conventional algorithm and the data-driven deep learning network.
- The proposed DSRGAN achieves the best on perceptual metrics compared with the most recent state-of-the-art SISR methods, meanwhile brings minor distortions to fidelity metrics.

The rest of this paper is organized as below. Relevant works on the SISR are presented in Section II. Details on the proposed DSRGAN are given in Section III. Extensive experimental results are provided in Section IV to verify the performance of DSRGAN. Finally, conclusion remarks are provided in Section V.

## II. RELATED WORK

In this section, existing fidelity-oriented and perception-oriented SISR methods are reviewed. The SISR methods assisted by different kinds of prior are also summarized.

### A. Fidelity-oriented SISR methods

In the field of deep learning-based SISR, Dong et al. firstly proposed a simple three-layer CNN to realize the SISR end-to-end [4]. Subsequently, networks of various architectures were proposed for a higher restoration accuracy [5], [22], [23], [24],

[6], [7]. These methods target at achieving high PSNR and SSIM by adopting L1 or L2 regularizations as loss functions. However, the fidelity-oriented SISR methods tend to average plausible outputs and produce blurry results, where edges and textures are smoothed.

### B. Perception-oriented SISR methods

It is important to preserve the sharp edges and fine textures such that the SR image looks more realistic. To achieve this, Johnson et al. proposed a loss function defined on the feature-domain to measure the perceptual differences between SR and HR images [25]. Ledig et al. proposed a GAN framework, named SRGAN, to perform a deterministic LR-to-HR image translation [8]. Then, the feature-domain loss and GAN are widely studied for producing SR images with high visual quality as perceived by human observers. Sajjadi et al. proposed a texture loss inspired by the feature-domain loss [26]. Wang et al. introduced a relativistic GAN to better measure the distribution differences between SR and HR images [9]. Zhang et al. proposed a rank-content loss for optimizing network parameters in the direction of better no-reference perceptual metrics [11]. Shaham et al. proposed an internal learning scheme to conduct the SISR reconstruction by learning from a single LR image [27]. Shaham et al. disabled the feature-domain loss for preserving higher PSNR and SSIM. Fuoli et al. proposed a fourier space loss to encourage the network to recover high frequency information [28]. Lugmayr et al. presented the normalizing flow to predict the diverse photo-realistic HR images instead of using GAN [29], [12].

### C. Prior-assisted SISR methods

Although perception-oriented SISR methods are able to produce SR images of high perceptual quality, they sometimes generate unnatural artifacts and fake high frequency details. Recently, a few methods attempted to exploit extra image prior knowledge to alleviate the ill-posed SR problem and to restore high frequency information closer to the original HR images. The image gradient prior is the most widely used kind of prior. Yang et al. employed an off-the-shelf edge detector to extract the edge gradient, which was then used to guide the deep network to reconstruct SR images with more clear edges [30]. Ma et al. proposed a gradient-domain loss based on a differentiable gradient extraction module, to help the network concentrate more on geometric structures and suppress undesired structural distortions [10]. Semantic prior can also be a guide for the SR reconstruction. Wang et al. proposed to pre-compute the semantic information of the LR image, and then utilized the condition mechanism to generate textures and details that are more in line with their corresponding semantic categories [31]. Li et al. introduced the face semantic information to achieve better SR reconstruction for face images [32].

Different from these methods, we introduce a novel detail prior to guide the deep network to restore SR images with more realistic details, which can be of higher perceptual quality. Compared with the gradient prior, the detail prior preserves

not only edge profile information, but also finer texture information. To the best of our knowledge, no perception-oriented SISR methods have explicitly exploited detail information for the SR reconstruction.

### III. THE PROPOSED DSRGAN

The overall framework of DSRGAN is shown in Figure 2. It includes a novel detail extraction algorithm, a generator network, and two discriminator networks. In this section, we give details on each part of the DSRGAN and present how it combines the model-based conventional algorithm and the data-driven deep learning network.

#### A. Detail extraction

In the field of image processing, an image  $I$  is often required to be decomposed into a base layer  $I_B$  and a detail layer  $I_D$ .  $I_B$  captures the image structure information and is mainly formed by homogeneous regions. While  $I_D$  carries the high frequency information and is constituted of residual smaller scale details. We focus on  $I_D$  as restoring high frequency details is the most critical part of the SISR reconstruction. Inspired by our previous work [14], a detail extraction algorithm is proposed to acquire  $I_D$  for the SISR. The detail extraction is divided into two stages: 1) generation of a vector field; 2) extraction of fine details from the vector field.

The RGB image  $I$  is firstly converted into the YCbCr image and the luminance component  $Y$  is extracted. Then, the vector field  $(V_h, V_v)$  (including horizontal and vertical dimensions) is constructed by computing gradients in  $\log$  domain as:

$$V_h(i, j) = (1 + \alpha) \log_2 \left( \frac{Y(i+1, j) + 1}{Y(i, j) + 1} \right), \quad (1)$$

$$V_v(i, j) = (1 + \alpha) \log_2 \left( \frac{Y(i, j+1) + 1}{Y(i, j) + 1} \right), \quad (2)$$

where  $(i, j)$  represents the pixel coordinate.  $\alpha (\geq 0)$  is a hyper parameter and it amplifies the details such that they have more chances to be extracted. Notably, a large  $\alpha$  also amplifies the noises.

The detail layer  $I_D$  is extracted from  $(V_h, V_v)$  by solving the following quadratic optimization problem [14]:

$$\min_{I_D} \left\{ \|I_D\|^2 + \lambda \left[ \left( \frac{V_h - \frac{\partial I_D}{\partial h}}{\psi(V_h)} \right)^2 + \left( \frac{V_v - \frac{\partial I_D}{\partial v}}{\psi(V_v)} \right)^2 \right] \right\}, \quad (3)$$

where the function  $\psi(z)$  is defined as:

$$\psi(z) = \sqrt{|z|^\gamma + \epsilon}. \quad (4)$$

The objective (3) includes a regularization term and a fidelity term. The regularization term  $\|I_D\|^2$  encourages the detail layer to be sparse. The sparsity can help the network concentrate more on useful information. The fidelity term preserves the details by measuring the difference between  $\nabla I_D$  and  $V$ .  $\lambda$  is a hyper parameter that achieves a tradeoff between two energy terms. The function  $\psi$  is to prevent the noises of  $V$  from appearing in  $I_D$ .  $\gamma$  determines the sensitivity of the objective to  $V$ .  $\epsilon$  prevents division by zero and a larger value excludes more noises from  $I_D$ . Notably, small scale details can also be excluded if  $\epsilon$  is too large.

The optimal solution to the objective (3) can be obtained by solving the following linear equation:

$$(E + \lambda(D'_h \mathcal{A}(V_h) D_h + D'_v \mathcal{A}(V_v) D_v)) I_D = \lambda(D'_h \mathcal{A}(V_h) V_h + D'_v \mathcal{A}(V_v) V_v), \quad (5)$$

where  $D'_h$  represents the transpose of  $D_h$ . The function  $\mathcal{A}(z)$  is defined as:

$$\mathcal{A}(z) = \text{diag} \left( \frac{1}{\psi^2(z)} \right). \quad (6)$$

$E$  is the identity matrix.  $D_h$  and  $D_v$  are discrete differentiation operators along the horizontal and vertical dimensions, respectively.

Solving the matrix  $(E + \lambda(D'_h \mathcal{A}(V_h) D_h + D'_v \mathcal{A}(V_v) D_v))^{-1}$  is complicated. Different from [14], we apply a 1-Dimensional edge-preserving smoothing technique in [15] to approximate the matrix via several iterations:

$$(E + \lambda(D'_h \mathcal{A}(V_h) D_h + D'_v \mathcal{A}(V_v) D_v))^{-1} \approx \prod_{t=1}^T (E + \lambda_t(D'_h \mathcal{A}(V_h) D_h))^{-1} (E + \lambda_t(D'_v \mathcal{A}(V_v) D_v))^{-1}, \quad (7)$$

where  $T$  represents the number of iterations and  $\lambda_t$  is obtained as follows:

$$\lambda_t = \frac{3}{2} \frac{4^{T-t}}{4^T - 1} \lambda. \quad (8)$$

In other words, the objective (3) is decomposed into two 1-Dimensional optimization problems along the horizontal and vertical dimensions, i.e.,  $I_D^h$  and  $I_D^v$ , respectively:

$$\min_{I_D^h} \sum_{i=0}^{W-1} \left[ (I_D^h(i))^2 + \lambda_t \left( \frac{V_h(i) - \frac{\partial I_D^h(i)}{\partial i}}{\psi(V_h(i))} \right)^2 \right], \quad (9)$$

$$\min_{I_D^v} \sum_{j=0}^{H-1} \left[ (I_D^v(j))^2 + \lambda_t \left( \frac{V_v(j) - \frac{\partial I_D^v(j)}{\partial j}}{\psi(V_v(j))} \right)^2 \right], \quad (10)$$

where  $H$  and  $W$  correspond to the height and width of image, respectively. Clearly, solving the matrices  $(E + \lambda_t(D'_h \mathcal{A}(V_h) D_h))^{-1}$  and  $(E + \lambda_t(D'_v \mathcal{A}(V_v) D_v))^{-1}$  for the objectives (9) and (10) becomes much easier than solving the objective (3). The matrix have an exact optimal solution that can be obtained using the Gaussian elimination algorithm via a complexity of  $O(H)$  or  $O(W)$ . The overall computational complexity of the proposed detail extraction algorithm is  $O(T \times H \times W)$ , which is greatly reduced compared with [14].

Notably, the detail extraction algorithm is different from [19], [20], [21], whereas  $I_B$  is filtered from  $I$  by edge-preserving smoothing algorithms and  $I_D$  is obtained by subtracting  $I_B$  from  $I$ . As a result,  $I_D$  tend to hold negative values and possible noises that filtered out by smoothing. While in the proposed detail extraction algorithm, a quadratic objective is directly defined on  $I$  in  $\log$  domain. After solving the objective, an exponential operation is required for the solution to obtain the final  $I_D$  [14]. Thus,  $I_D$  is strictly positive. We define a function  $\mathcal{E}$  to denote the entire detail extraction algorithm in this paper, i.e.,  $I_D = \mathcal{E}(I)$ .

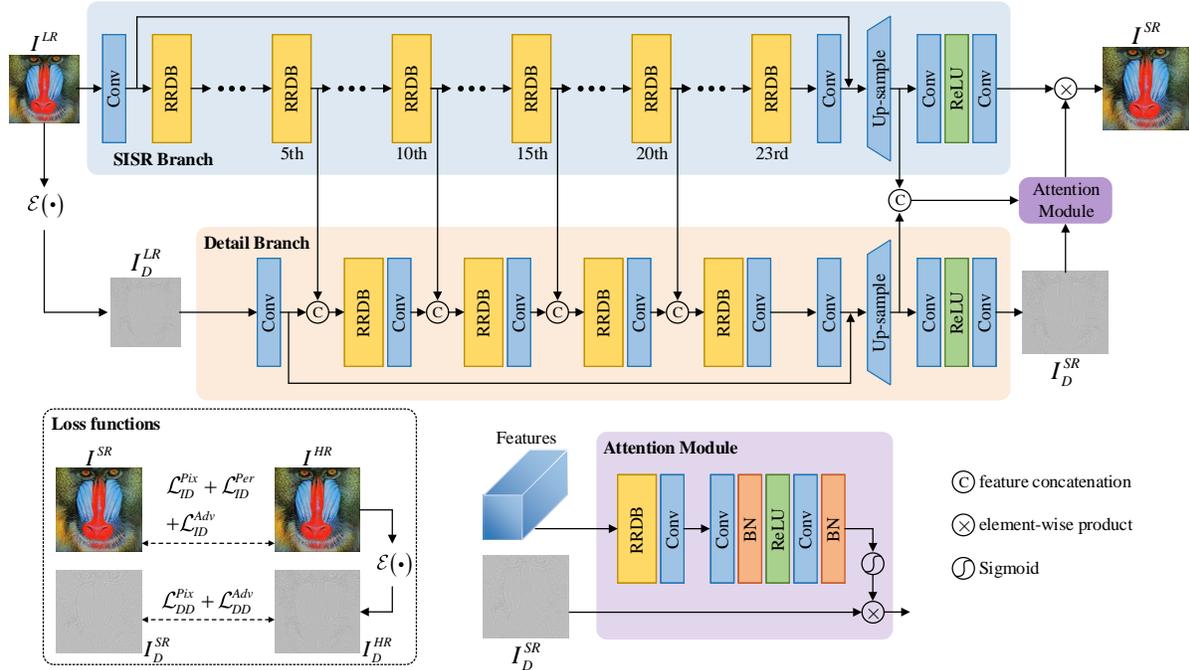


Fig. 2: The overall framework of the proposed DSRGAN. It includes a novel model-based detail extraction algorithm, a generator, an image-domain discriminator, and a detail-domain discriminator. The generator network is composed of a main SISR branch, an auxiliary detail branch, and an attention module.

### B. The generator network of DSRGAN

The DSRGAN learns the high frequency details in an explicit way with the help of  $\mathcal{E}$ . The generator of DSRGAN includes three parts: a main SISR branch, an auxiliary detail branch, and an attention module.

The SISR branch is composed of 23 residual-in-residual dense blocks (RRDBs) [9]. The nearest neighbor interpolation is performed at the end of the branch to resize features into the same size of the HR image  $I^{HR}$ . Finally, the features are mapped to the image domain via the convolution-ReLU-convolution.

The detail branch takes the detail layer of  $I^{LR}$  as input, which is obtained by  $I_D^{LR} = \mathcal{E}(I^{LR})$ , and restores its HR counterpart  $I_D^{SR}$ . The ultimate goal of detail branch is to learn high frequency details explicitly and to feedback the details to the main SISR branch. Following [10], the detail branch is composed of only 4 RRDBs considering the increment of network parameters. However, a lightweight detail branch is struggling to restore the seriously lost details. Therefore, we incorporate the features from the SISR branch to enhance the representation ability of the detail branch. Specifically, the features from the 5th, 10th, 15th, and 20th RRDBs in the SISR branch are concatenated with the features from the 1st, 2nd, 3rd, and 4th RRDBs in the detail branch, as shown in Figure 2. Each RRDB in the detail branch is followed by a convolution layer to reduce the feature channels. Finally, the upsampled features are mapped to the detail domain via the convolution-ReLU-convolution.

Once the  $I_D^{SR}$  is generated, it is multiplied back to the image domain to generate the final SR image  $I^{SR}$ . The multiplication is inspired by the conventional detail enhancement algorithm [14]. It is also shown in [16] that such detail enhancement

can improve both the subjective and the objective qualities of image.

Furthermore, we apply an attention on the  $I_D^{SR}$  to prevent artifacts caused by inappropriate detail enhancement. The attention module takes features from both the SISR and detail branches and computes the attention map through a few convolutions, activations, and normalizations. Owing to this design, the restored detail layer can enhance the SR image in a more adaptive way.

We define a function  $\mathcal{G}$  to denote the generator of DSRGAN and we get:

$$I^{SR}, I_D^{SR} = \mathcal{G}(I^{LR}, I_D^{LR} | \theta_G), \quad (11)$$

where  $\theta_G$  corresponds to the network parameters.

### C. Loss functions

The loss functions can be divided into image-domain losses and detail-domain losses, respectively. The DSRGAN includes an image-domain discriminator  $\mathcal{D}_{ID}$  and a detail-domain discriminator  $\mathcal{D}_{DD}$ , for supervision on restoring  $I^{SR}$  and  $I_D^{SR}$ , respectively.

1) *Image-domain loss functions*: Commonly used image-domain loss functions are constituted of a pixel loss, a perceptual loss, and an adversarial loss [8], [9], [10]. The pixel loss measures the pixel-wise differences between  $I^{SR}$  and ground truth  $I^{HR}$ :

$$\mathcal{L}_{ID}^{Pix} = \mathbb{E}_{I^{SR}} \|I^{SR} - I^{HR}\|_1, \quad (12)$$

The perceptual loss measures the differences between the deep-level features of SR and HR images, which are extracted from a pre-trained VGG model [33]:

$$\mathcal{L}_{ID}^{Per} = \mathbb{E}_{I^{SR}} \|\mathcal{V}(I^{SR}) - \mathcal{V}(I^{HR})\|_1, \quad (13)$$

where  $\mathcal{V}$  represents the VGG model.

The adversarial loss  $\mathcal{L}_{ID}^{Adv}$  includes an energy function for the  $\mathcal{G}$  and  $\mathcal{D}_{ID}$ . They are optimized by playing a two-player game:

$$\begin{aligned} \mathcal{L}_{ID}^{Adv}(\theta_{\mathcal{D}_{ID}}) = & -\mathbb{E}_{I^{SR}}[\log(1 - \mathcal{D}_{ID}(I^{SR}))] \\ & -\mathbb{E}_{I^{HR}}[\log(\mathcal{D}_{ID}(I^{HR}))], \end{aligned} \quad (14)$$

$$\mathcal{L}_{ID}^{Adv}(\theta_{\mathcal{G}}) = -\mathbb{E}_{I^{SR}}[\log(\mathcal{D}_{ID}(I^{SR}))], \quad (15)$$

where  $\mathcal{D}_{ID}$  is implemented by a encoder network similar to [8], [9], [10]. We also employ a relativistic average GAN proposed in [9] to boost the performance of DSRGAN.

2) *Detail-domain loss functions*: In typical deep learning-based SISR methods, details are implicitly learned by residual learning, where the input and output of network are skipped connected [5], [6], [7]. However, the deep learning network works in a black-box way, which does not guarantee that the residuals are of high frequency. Different from these works, the DSRGAN learns high frequency information in a more interpretable way. To achieve this, the detail-domain loss functions are proposed for supervision of the detail layer. The loss functions constituted of a pixel loss and an adversarial loss, respectively.

Similar to  $\mathcal{L}_{ID}^{Pix}$ , the pixel loss defined on the detail domain measures the pixel-wise differences between  $I_D^{SR}$  and ground truth  $I_D^{HR}$ :

$$\mathcal{L}_{DD}^{Pix} = \mathbb{E}_{I_D^{SR}} \|I_D^{SR} - I_D^{HR}\|_1, \quad (16)$$

where  $I_D^{HR} = \mathcal{E}(I^{HR})$ .

A simple L1 regularization can lead the  $I_D^{SR}$  to be smooth and it is too weak to capture high frequency differences. We further employ a powerful GAN to solve this problem. Similar to  $\mathcal{L}_{ID}^{Adv}$ , a detail-domain discriminator  $\mathcal{D}_{DD}$  is introduced and  $\mathcal{G}$  and  $\mathcal{D}_{DD}$  are optimized by playing another two-player game:

$$\begin{aligned} \mathcal{L}_{DD}^{Adv}(\theta_{\mathcal{D}_{DD}}) = & -\mathbb{E}_{I_D^{SR}}[\log(1 - \mathcal{D}_{DD}(I_D^{SR}))] \\ & -\mathbb{E}_{I_D^{HR}}[\log(\mathcal{D}_{DD}(I_D^{HR}))], \end{aligned} \quad (17)$$

$$\mathcal{L}_{DD}^{Adv}(\theta_{\mathcal{G}}) = -\mathbb{E}_{I_D^{SR}}[\log(\mathcal{D}_{DD}(I_D^{SR}))], \quad (18)$$

where  $\mathcal{D}_{DD}$  shares the same network architecture with  $\mathcal{D}_{ID}$ . In practice, a relativistic average GAN is utilized. Notably, the extracted  $I_D^{LR}$  and  $I_D^{HR}$  are sparse, which is very friendly to the optimization procedure of  $\mathcal{G}$  and  $\mathcal{D}_{DD}$ .

The overall loss function is defined as follows:

$$\begin{aligned} \mathcal{L} = & \eta_{ID}^{Pix} \mathcal{L}_{ID}^{Pix} + \eta_{ID}^{Per} \mathcal{L}_{ID}^{Per} + \eta_{ID}^{Adv} \mathcal{L}_{ID}^{Adv} \\ & + \eta_{DD}^{Pix} \mathcal{L}_{DD}^{Pix} + \eta_{DD}^{Adv} \mathcal{L}_{DD}^{Adv}. \end{aligned} \quad (19)$$

The hyper-parameter  $\eta$  leverages the penalties on different losses.

## D. Differences between the DSRGAN and SPSR

The proposed DSRGAN shares a similar network architecture to the SPSR in [10]. However, the DSRGAN is fundamentally different from the SPSR: a) The SPSR includes a gradient branch to learn image gradients, while the DSRGAN focus on image details. Compared with the gradients, the details are able to record fine textures in addition to sharp edges. Thus, the DSRGAN can pay more attention on recovering these high frequency information; b) In the DSRGAN, the fusion of SISR branch and detail branch does not rely on feature concatenation as in the SPSR. The fusion is realized by a more interpretable way, where the detail layer is multiplied back to the image domain. Such a fusion manner adheres to the conventional detail enhancement algorithm [14] and brings a more powerful ability to the network in generating high perceptual quality images; c) The gradient extraction has to be differentiable for training the SPSR. While in the DSRGAN, the detail extraction  $\mathcal{E}$  is not required to be differentiable for an end-to-end training. Therefore, following the DSRGAN, it is feasible to blend other model-based conventional algorithms with recent data-driven deep learning methods and take advantages from both of them.

Notably, the SISR and detail branches in the DSRGAN are agnostic to the specific selection of network architectures, and any other deep learning models are feasible. Therefore, the DSRGAN is able to benefit from off-the-shelf networks that have powerful representation abilities.

## IV. EXPERIMENTS

Extensive experiments are conducted to demonstrate the performance of our proposed DSRGAN. For qualitative comparison, readers are invited to view the electronic version of figures and zoom in them so as to better check differences among all images.

### A. Datasets and Metrics

1) *Datasets*: The training split of DIV2K [34] is used for training. It contains 800 2k-resolution images. For test, 5 commonly used datasets including Set5, Set14, BSD100, Urban100, and General100 are used. For both training and test, the LR images are stimulated by downsampling the HR images with the scaling factor as 4 via Bicubic kernel.

2) *Metrics*: For evaluating the quality of SR images, two metrics on measuring fidelity and two metrics on measuring perceptual quality are used, respectively. The fidelity metrics include PSNR and SSIM, which compute the dissimilarities between SR and HR images. One of the perceptual metrics, LPIPS [3], computes the dissimilarities on feature-domain. The other one, PI [2], is a no-reference metric which does not need the ground truth HR images for evaluation. For the PSNR and SSIM, higher is better. For the PI and LPIPS, lower is better.

### B. Implementation details

We use Pytorch based on Python to implement all methods. Experiments are conducted on the 4029GP-TRT server carried

TABLE I: Ablation study. The best and second best results are shown in red and blue.

| No. | $\mathcal{L}_{DD}^{Pix}$ | $\mathcal{L}_{DD}^{Adv}$ | Attention | Set5   |        |        | Set14  |        |        | B100   |        |        |
|-----|--------------------------|--------------------------|-----------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
|     |                          |                          |           | PI     | LPIPS  | PSNR   | PI     | LPIPS  | PSNR   | PI     | LPIPS  | PSNR   |
| 1   | X                        | X                        | X         | 3.4177 | 0.0656 | 30.232 | 2.9016 | 0.1330 | 26.566 | 2.3994 | 0.1637 | 25.396 |
| 2   | ✓                        | X                        | X         | 3.2641 | 0.0685 | 30.161 | 2.8334 | 0.1369 | 26.453 | 2.3975 | 0.1615 | 25.336 |
| 3   | ✓                        | ✓                        | X         | 3.4642 | 0.0640 | 30.301 | 2.7814 | 0.1324 | 26.504 | 2.3543 | 0.1596 | 25.452 |
| 4   | ✓                        | ✓                        | ✓         | 3.2019 | 0.0641 | 30.172 | 2.7387 | 0.1303 | 26.554 | 2.3525 | 0.1590 | 25.645 |

with NVIDIA TITAN RTX GPU of 24GB memory. All the metrics of SR results are computed using MATLAB 2020a. We give more details about the method implementation and network training.

1) *Method details*: For the proposed detail extraction algorithm,  $\alpha$  in Equations (1) and (2) is selected as 4 to amplify the details.  $\lambda$  in Equation (3) is selected as 1 for balancing the regularization and fidelity terms, such that the extracted details can be sparse.  $\gamma$  and  $\epsilon$  in Equation (4) are selected as 0.75 and 2 for simultaneously preserving details and excluding noises. The iteration times  $T$  in Equation (7) is selected as 4. For the generator network, the number of feature channels of the RRDB in the SISR branch is selected as 64 [9]. While the number of feature channels of the RRDB in the detail branch is selected as 128 [10].

2) *Training details*: Before training the DSRGAN, we first use the pre-trained fidelity-oriented model in [9] to initialize the network parameters of the main SISR branch. During training, the HR and corresponding LR images are randomly cropped into patches with the size of  $128 \times 128$  and  $32 \times 32$  on the fly. The batch size is set as 20. The random cropping, flipping, and rotating are applied for data augmentation. In the loss function (19),  $\eta_{ID}^{Pix}$ ,  $\eta_{ID}^{Per}$ , and  $\eta_{DD}^{Pix}$  are set as 0.2, 1, and 0.5, respectively. The  $\eta_{ID}^{Adv}$  for the generator and discriminator are set as 0.005 and 1, respectively, the same for  $\eta_{DD}^{Adv}$ . The optimizer is Adam [35] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The initial learning rates for generator and two discriminators are all set as  $1e-4$ , and they are decreased by a factor of 0.5 at 50k, 100k, 200k, and 300k iterations. The number of total training iterations is up to 400k and it takes about 4-5 days to finish the training.

### C. Ablation study

To investigate the effects of the model-based detail extraction algorithm on the data-driven deep learning network, the ablation study is conducted on Set5, Set14, and B100 datasets. The improvements are gradually added to the baseline method to see how they contribute to the SR results. Different ablation methods are listed on the left of Table I, and each of them is explained as follows:

- The No. 1 ablation represents the baseline method, where the detail extraction algorithm, the detail branch, and the detail-domain discriminator are disabled. In practice, it is implemented by an enhanced ESRGAN, whose generator is composed of 27 RRDBs instead of 23 for a fair comparison with the rest of ablation methods (i.e., for roughly equivalent number of network parameters).
- The No. 2 ablation represents the DSRGAN casting off  $\mathcal{L}_{DD}^{Adv}$  and the attention module, where the detail layer

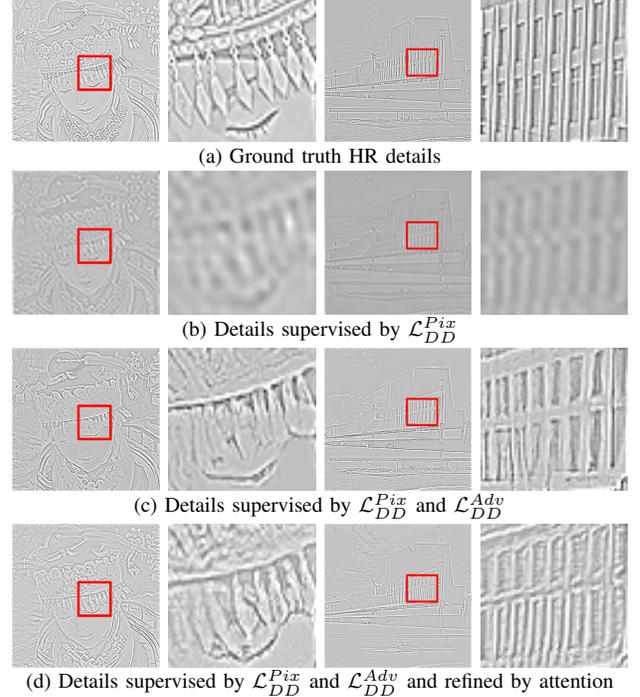


Fig. 3: The visualization of details. They are normalized into the range from 0 to 1. The details are exhibited on the left column and their corresponding parts in the red rectangular box are enlarged on the right. The examples are ‘comic’ from Set14 and ‘78004’ from B100.

$I_D^{LR}$  generated by the detail branch is only supervised by  $\mathcal{L}_{DD}^{Pix}$  with penalty  $\eta_{DD}^{Pix}$ .

- The No. 3 ablation represents the DSRGAN casting off the attention module, where  $I_D^{LR}$  is supervised by  $\mathcal{L}_{DD}^{Pix}$  with penalty  $\eta_{DD}^{Pix}$  and by  $\mathcal{L}_{DD}^{Adv}$  with penalty  $\eta_{DD}^{Adv}$ .
- The No. 4 ablation represents the complete DSRGAN, where  $I_D^{LR}$  is refined by the attention module before being multiplied back to the image domain.

According to Table I, the No. 1 and 2 ablation methods have their own strengths on three datasets in terms of LPIPS. The improvement is not significant as the  $I_D^{SR}$  is restored badly by training with the  $\mathcal{L}_{DD}^{Pix}$  only. As shown in Figures 3(a) and 3(b), edges and textures are seriously smoothed in the restored  $I_D^{SR}$ . This problem is the same with that met by fidelity-oriented SISR methods, where the L1 regularization is hard to capture high frequency differences.

The No. 3 ablation method improves the PI and LPIPS by a large margin from the No. 1 and 2. As shown in Figure 3(c), the high frequency information is preserved well in  $I_D^{SR}$  by adding the adversarial loss  $\mathcal{L}_{DD}^{Adv}$ . The experimental results demonstrate that an explicit learning of detail prior plays a very positive effect on the perceptual SISR reconstruction.

TABLE II: Comparison with state-of-the-art perception-oriented SISR methods. The best and second best results are shown in red and blue.

| Method                    | Set5          |               | Set14         |               | B100          |               | Urban100      |               | General100    |               |
|---------------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|                           | PSNR          | SSIM          |
| Bicubic                   | 28.386        | 0.8249        | 26.097        | 0.7853        | 25.961        | 0.6675        | 23.145        | 0.9011        | 28.018        | 0.8281        |
| SRGAN [8] (CVPR'17)       | 29.879        | <b>0.8694</b> | 26.552        | 0.7891        | 25.498        | 0.6520        | 24.387        | 0.9376        | 29.358        | 0.8531        |
| EnhanceNet [26] (ICCV'17) | 28.548        | 0.8374        | 25.763        | 0.7737        | 24.931        | 0.6259        | 23.543        | 0.9357        | 28.069        | 0.8308        |
| SFTGAN [31] (CVPR'18)     | 29.913        | 0.8672        | 26.226        | 0.7860        | 25.505        | 0.6549        | 24.015        | 0.9365        | 29.034        | 0.8512        |
| ESRGAN [9] (ECCV'18)      | <b>30.422</b> | <b>0.8683</b> | 26.272        | 0.7837        | 25.317        | 0.6506        | 24.360        | 0.9452        | 29.412        | 0.8546        |
| SinGAN [27] (ICCV'19)     | 26.419        | 0.7920        | 23.134        | 0.6807        | 24.330        | 0.5950        | 21.300        | 0.8203        | 25.392        | 0.7664        |
| SPSR [10] (CVPR'20)       | <b>30.382</b> | 0.8635        | <b>26.635</b> | <b>0.7939</b> | 25.505        | 0.6576        | <b>24.799</b> | 0.9481        | 29.414        | 0.8536        |
| SRFlow [29] (ECCV'20)     | 30.184        | 0.8609        | <b>26.880</b> | <b>0.8006</b> | <b>26.055</b> | <b>0.6721</b> | <b>25.268</b> | <b>0.9535</b> | <b>29.770</b> | <b>0.8591</b> |
| RankSRGAN [11] (TPAMI'21) | 29.640        | 0.8587        | 26.446        | 0.7864        | 25.453        | 0.6484        | 24.472        | 0.9413        | 29.098        | 0.8480        |
| Fourier [28] (ICCV'21)    | -             | -             | -             | -             | 25.660        | 0.6560        | 24.690        | 0.7230        | -             | -             |
| HCFlow [12] (ICCV'21)     | 30.360        | 0.8636        | 26.146        | 0.7809        | 25.537        | 0.6608        | 24.510        | 0.9471        | <b>29.646</b> | <b>0.8557</b> |
| DSRGAN (Ours)             | 30.172        | 0.8607        | 26.554        | 0.7932        | <b>25.645</b> | <b>0.6614</b> | 24.686        | <b>0.9482</b> | 29.476        | 0.8548        |
|                           | <b>PI</b>     | <b>LPIPS</b>  |
| Bicubic                   | 7.3903        | 0.3440        | 7.0207        | 0.4412        | 7.0010        | 0.5249        | 6.9413        | 0.4726        | 7.9351        | 0.3528        |
| SRGAN [8] (CVPR'17)       | 3.5684        | 0.0814        | 2.9234        | 0.1452        | 2.3816        | 0.1779        | 3.4848        | 0.1425        | 4.2173        | 0.0961        |
| EnhanceNet [26] (ICCV'17) | <b>2.9600</b> | 0.1015        | 2.9842        | 0.1629        | 2.9071        | 0.2013        | <b>3.4649</b> | 0.1632        | 4.1251        | 0.1327        |
| SFTGAN [31] (CVPR'18)     | 3.7726        | 0.0892        | 2.9027        | 0.1489        | 2.3775        | 0.1769        | 3.6141        | 0.1433        | 4.2843        | 0.1028        |
| ESRGAN [9] (ECCV'18)      | 3.7760        | 0.0758        | 2.8780        | 0.1332        | 2.4806        | 0.1614        | 3.7628        | 0.1229        | 4.3188        | 0.0879        |
| SinGAN [27] (ICCV'19)     | 4.1939        | 0.1901        | 2.9864        | 0.3034        | 2.6490        | 0.2916        | 3.7266        | 0.2657        | 4.1062        | 0.2208        |
| SPSR [10] (CVPR'20)       | 3.3124        | <b>0.0647</b> | 2.8770        | 0.1327        | <b>2.3513</b> | <b>0.1610</b> | 3.5437        | <b>0.1184</b> | 4.0969        | 0.0863        |
| SRFlow [29] (ECCV'20)     | 4.2089        | 0.0818        | 3.0433        | <b>0.1316</b> | 2.5971        | 0.1828        | 3.7927        | 0.1269        | 4.6271        | 0.0956        |
| RankSRGAN [11] (TPAMI'21) | <b>3.1266</b> | 0.0734        | <b>2.5123</b> | 0.1371        | <b>2.1234</b> | 0.1745        | <b>3.3304</b> | 0.1380        | <b>3.8916</b> | 0.0955        |
| Fourier [28] (ICCV'21)    | -             | -             | -             | -             | -             | 0.1720        | -             | 0.1320        | -             | -             |
| HCFlow [12] (ICCV'21)     | 4.0203        | 0.0752        | 2.9530        | 0.1325        | 2.5951        | 0.1638        | 3.6320        | 0.1242        | 4.3690        | <b>0.0859</b> |
| DSRGAN (Ours)             | 3.2019        | <b>0.0641</b> | <b>2.7387</b> | <b>0.1303</b> | 2.3525        | <b>0.1590</b> | 3.4918        | <b>0.1175</b> | <b>4.0319</b> | <b>0.0847</b> |

The deep learning network is strengthened by the model-based conventional algorithm.

The refinement by attention module further highlights the edges and textures in  $I_D^{SR}$ , as shown in Figure 3(d). The attention can also adaptively controls the degree of detail enhancement for per pixel. According to Table I, it improves the perceptual quality of SR images, though slightly. Notably, the PSNR is not affected by the detail enhancement and even gets better on a few datasets.

#### D. Comparison with state-of-the-art methods

The complete DSRGAN is compared with several state-of-the-art perception-oriented SISR methods on benchmark SR datasets. Notably, for a fair comparison, we implement all these methods (except for Fourier [28]) by using the official source codes and trained model weights provided by their authors instead of retraining them from scratch. The experimental results of Fourier is directly quoted from its paper. All the results are summarized in Table II.

For perceptual metrics, the proposed DSRGAN achieves the best LPIPS on all 5 datasets. The DSRGAN is also very competitive on the PI and achieves the second best on Set14 and General100. The RankSRGAN [11] performs the best PI on most datasets as it is trained in the direction of PI via a rank-content loss. However, its LPIPS is worse than most recent methods. Compared with the most relevant method SPSR [10], the DSRGAN improves both PI and LPIPS significantly on all datasets. For the fidelity metrics, the DSRGAN still has an advantage over most perception-oriented SISR methods. It achieves the second best SSIM on B100 and Urban100. The SRFlow [29] gets the highest PSNR and SSIM

on most datasets. However, its performances on PI and LPIPS are very inferior to the DSRGAN. Since the SRFlow is a flow-based method and doesn't include the adversarial loss, it is more like a fidelity-oriented SISR method. The overall results demonstrate the superior ability of DSRGAN in producing nature and photo-realistic SR images with minor fidelity distortions simultaneously. In addition, the transcendence of the DSRGAN over the SFTGAN [31] and SPSR [10] indicates that the detail prior is a more powerful assistance for the perceptual SISR than the semantic or gradient priors. The detail prior partially alleviates the ill-posed problem of the SISR reconstruction.

A few restored SR images are shown in Figure 4 for qualitative comparison. Readers are invited to view the electronic version of figures and zoom in them so as to better check differences among all images. For the 'baboon' from Set14, the DSRGAN restores the beard as realistic as ground truth HR image. This is benefiting from the explicit texture extraction by model-based conventional algorithm. For the '24077' from B100, enabled by taking multiplication as the fusion manner, the DSRGAN generates even more clear and natural pillars than ground truth. For the 'img\_044' and 'img\_054' from Urban100, the SR images produced by DSRGAN have the most similar sharp edges to HR counterparts among all the methods. This shows that the detail prior is able to preserve sharp edges and structures, with the same as the gradient prior [10]. Notably, the detail prior additionally includes the fine texture information that the gradient prior doesn't have. Moreover, owing to the concept of detail enhancement [14], the image restored by DSRGAN is even of higher perceptual quality than the original image. For instance, in the 'img\_044',

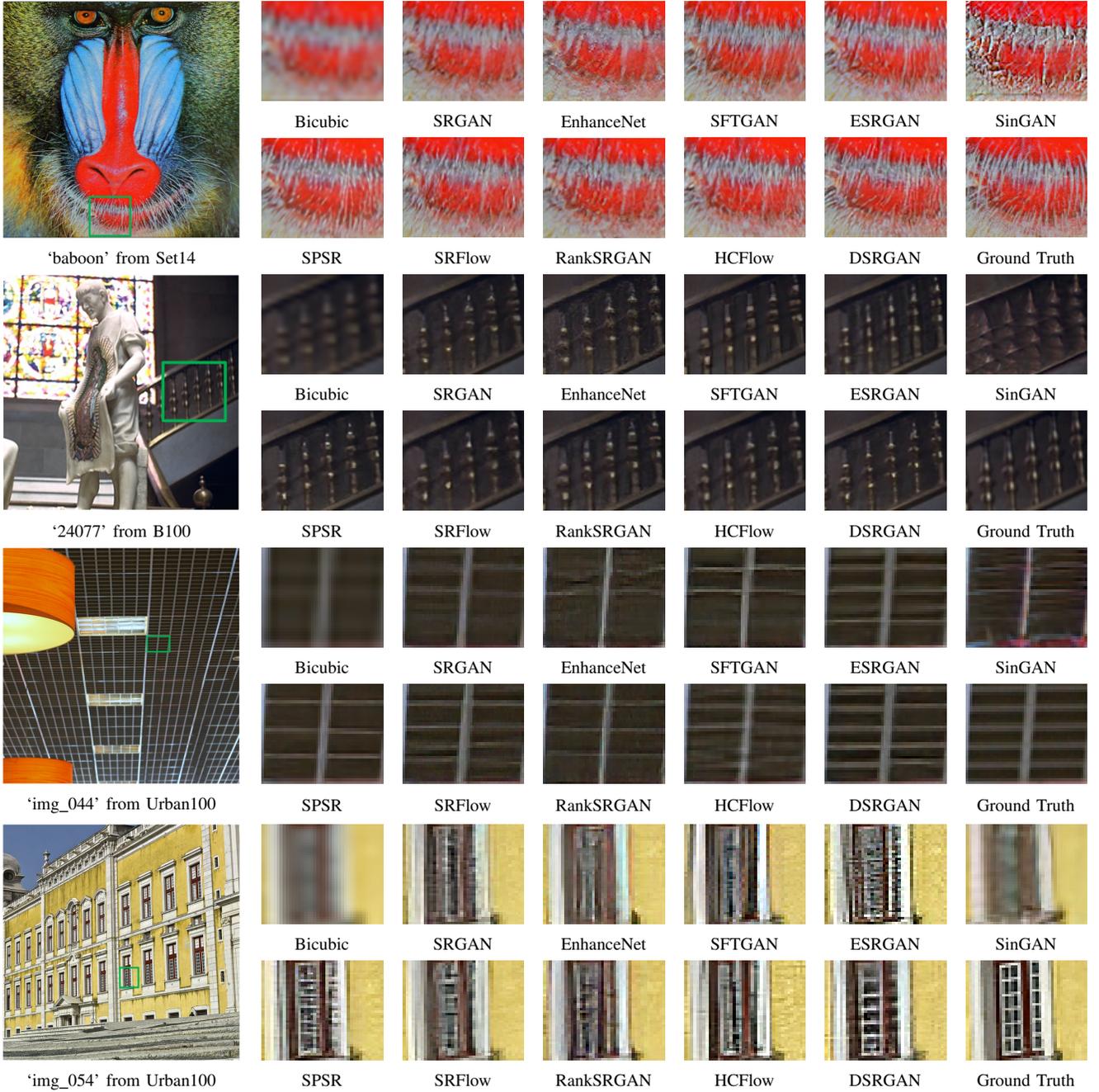


Fig. 4: The qualitative comparison of SR images restored by state-of-the-art perception-oriented SISR methods. The original images are exhibited on the left column and their SR counterparts in the green rectangular box are enlarged on the right.

the luminance of edges produced by DSRGAN is even brighter than ground truth.

#### E. Analysis of the detail prior

In this subsection, we replace  $\mathcal{E}$  in the DSRGAN with other commonly used detail extraction algorithms to see how they affect the SISR reconstruction. In the field of image processing,  $I_D$  is commonly extracted by subtracting  $I_B$  from  $I$ , where  $I_B$  is smoothed from  $I$  by image smoothing algorithm. Two classic smoothing algorithms, guided image filtering (GIF) [20] and multi-scale detail manipulation (MSDM) [19],

are utilized to replace the proposed detail extraction algorithm to obtain  $I_D^{LR}$  and  $I_D^{HR}$  in the complete DSRGAN.

For the GIF, the  $r_{GIF}$  and  $\epsilon_{GIF}$  in [20] are selected as 2 and 0.3. For the MSDM, the  $\lambda_{MSDM}$  and  $\alpha_{MSDM}$  in [19] are selected as 1 and 1.2. We denote the new DSRGAN based on GIF and MSDM as DSRGAN-GIF and DSRGAN-MSDM, respectively. Notably, in the DSRGAN-GIF and DSRGAN-MSDM, the restored detail  $I_D^{SR}$  is element-wise summed back to the image domain instead of multiplication, because the ground truth detail  $I_D^{HR}$  is subtracted from  $I^{HR}$ .

The performances of DSRGAN, DSRGAN-GIF, and DSRGAN-MSDM are compared on Set14 and B100, as shown

TABLE III: Comparison among DSRGANs assisted by different kinds of detail. The best results is shown in red.

| Method      | Set14         |               |               |               |
|-------------|---------------|---------------|---------------|---------------|
|             | PI            | LPIPS         | PSNR          | SSIM          |
| DSRGAN-GIF  | 2.8215        | 0.1372        | 26.314        | 0.7871        |
| DSRGAN-MSDM | 2.8484        | 0.1366        | 26.085        | 0.7884        |
| DSRGAN      | <b>2.7387</b> | <b>0.1303</b> | <b>26.554</b> | <b>0.7932</b> |
| Method      | B100          |               |               |               |
|             | PI            | LPIPS         | PSNR          | SSIM          |
| DSRGAN-GIF  | <b>2.3154</b> | 0.1602        | 25.543        | 0.6604        |
| DSRGAN-MSDM | 2.4765        | 0.1606        | 25.066        | 0.6415        |
| DSRGAN      | 2.3525        | <b>0.1590</b> | <b>25.645</b> | <b>0.6614</b> |

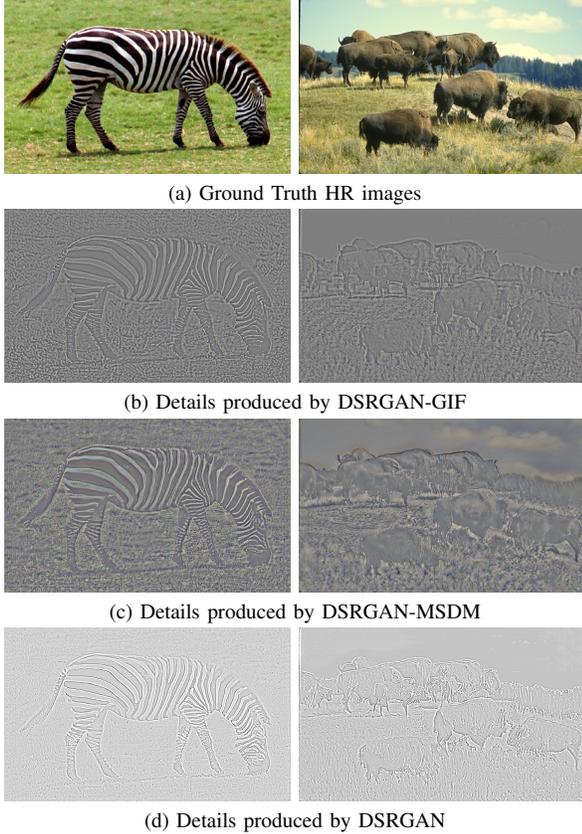


Fig. 5: The visualization of details. They are normalized into the range from 0 to 1. Notably, the details of DSRGAN have only one channel, while the details of DSRGAN-GIF and DSRGAN-MSDM have three channels. The examples are ‘zebra’ from Set14 and ‘38092’ from B100.

in Table III. We can see that these methods achieve close PI and LPIPS, demonstrating that conventional detail extraction algorithms do have a positive effect on deep learning-based SISR. Therefore, combining model-based conventional algorithms with deep learning networks for the SISR reconstruction has a great potential research value. The combination can be achieved via the DSRGAN, where the conventional algorithms are not required to be differentiable.

The DSRGAN still has a slight advantage over the DSRGAN-GIF and DSRGAN-MSDM on both perceptual metrics and fidelity metrics. The reason might be that the detail extracted by the proposed algorithm is more sparse than those extracted by the GIF and MSDM. The sparse attribute originates from the regularization term in Equation

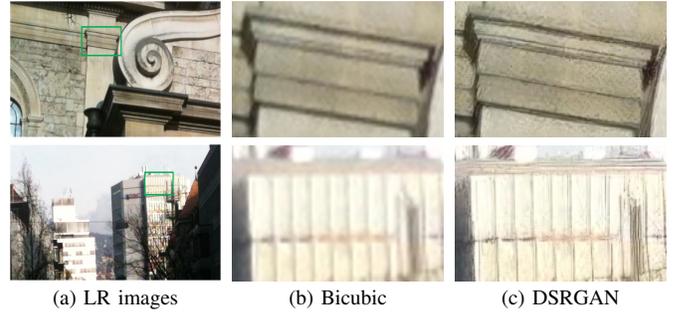


Fig. 6: The performance of DSRGAN on real-world image SR. The original LR images are exhibited on the left column and their SR counterparts in the green rectangular box are enlarged on the right. The examples are ‘00047’ and ‘00049’.

(3). As such, the burden of detail branch for restoring the  $I_D^{SR}$  is greatly reduced. Another reason might be that image noises are filtered out from the  $I_B$  by GIF and MSDM, but reappear in the  $I_D$ . The noises can be interferences to the image. We visualize the details restored by DSRGAN-GIF, DSRGAN-MSDM, and DSRGAN in Figure 5. The detail from DSRGAN do have fewer noises compared to the DSRGAN-GIF and DSRGAN-MSDM, and it is much more sparse. The details from DSRGAN-GIF and DSRGAN-MSDM not only preserve high frequency information, but also include unwanted redundant information such as colors and noises.

#### F. Limitation of the proposed DSRGAN

To investigate the generalization ability of the proposed DSRGAN, we conduct experiments on real-world image SR. The track 2 of NTIRE 2020 Challenge contains a set of LR images captured by smart phones, where ground truth HR images are non-existent [36]. The SR images restored by DSRGAN are shown in Figure 6. It can be observed that the DSRGAN preserves edges clear for the real-world image SR. Nevertheless, on the second row, the SR image is carried with a few artifacts and noises, which demonstrates that there still exists a performance gap of DSRGAN between the benchmark SR and the real-world SR. The DSRGAN can be extended to the real-world image SR by taking advantages of real-world-oriented SR methods, such as [37]. This will be our future work.

## V. CONCLUSION

Inspired by conventional detail enhancement algorithms, a novel detail prior-assisted SISR paradigm via the GAN, named DSRGAN, is proposed in this study. Extensive experiments demonstrate that the proposed detail prior is a powerful assistance for the deep learning network in producing SR images with realistic details. The DSRGAN benefits from advantages of both conventional model-based algorithms and data-driven deep learning networks. It is also network-agnostic, which can be used for off-the-shelf SISR networks. Moreover, the detail extraction algorithm is not required to be differentiable. Following the DSRGAN, it is feasible to incorporate other conventional image processing algorithms into a deep learning

network to form a model-based deep SISR, which will have a great potential research value.

One more interesting problem is to extend the DSRGAN to the video coding. Instead of coding a video directly, the video can be down-sampled. A super-resolution friendly rate control can be extended from existing rate control algorithms [38], [39] to code the down-sampled video. Then, the video can be super-resolved by the extended DSRGAN. These problems will be studied in our future research.

## REFERENCES

- [1] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE transactions on image processing*, vol. 19, no. 11, pp. 2861–2873, 2010.
- [2] Y. Blau, R. Mechrez, R. Timofte, T. Michaeli, and L. Zelnik-Manor, "The 2018 pirm challenge on perceptual image super-resolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 0–0, 2018.
- [3] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- [4] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 295–307, 2015.
- [5] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1646–1654, 2016.
- [6] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 136–144, 2017.
- [7] J. Yu, Y. Fan, J. Yang, N. Xu, Z. Wang, X. Wang, and T. Huang, "Wide activation for efficient and accurate image super-resolution," *arXiv preprint arXiv:1808.08718*, 2018.
- [8] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al., "Photo-realistic single image super-resolution using a generative adversarial network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4681–4690, 2017.
- [9] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 0–0, 2018.
- [10] C. Ma, Y. Rao, Y. Cheng, C. Chen, J. Lu, and J. Zhou, "Structure-preserving super resolution with gradient guidance," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7769–7778, 2020.
- [11] Z. Wenlong, L. Yihao, C. Dong, and Y. Qiao, "Ranksrgan: Generative adversarial networks with ranker for image super-resolution," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [12] J. Liang, A. Lugmayr, K. Zhang, M. Danelljan, L. Van Gool, and R. Timofte, "Hierarchical conditional flow: A unified framework for image super-resolution and image rescaling," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4076–4085, October 2021.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- [14] Z. G. Li, J. H. Zheng, and S. Rahardja, "Detail-enhanced exposure fusion," *IEEE Transactions on Image Processing*, vol. 21, no. 11, pp. 4672–4676, 2012.
- [15] D. Min, S. Choi, J. Lu, B. Ham, K. Sohn, and M. N. Do, "Fast global image smoothing based on weighted least squares," *IEEE Transactions on Image Processing*, vol. 23, no. 12, pp. 5638–5653, 2014.
- [16] Q. Wang, W. Chen, X. Wu, and Z. Li, "Detail-enhanced multi-scale exposure fusion in yuv color space," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 8, pp. 2418–2429, 2019.
- [17] F. Kou, Z. Wei, W. Chen, X. Wu, C. Wen, and Z. Li, "Intelligent detail enhancement for exposure fusion," *IEEE Transactions on Multimedia*, vol. 20, no. 2, pp. 484–495, 2017.
- [18] Z. Wei, C. Wen, and Z. Li, "Local inverse tone mapping for scalable high dynamic range image coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, no. 2, pp. 550–555, 2016.
- [19] Z. Farbman, R. Fattal, D. Lischinski, and R. Szeliski, "Edge-preserving decompositions for multi-scale tone and detail manipulation," *ACM Transactions on Graphics (TOG)*, vol. 27, no. 3, pp. 1–10, 2008.
- [20] K. He, J. Sun, and X. Tang, "Guided image filtering," in *European conference on computer vision*, pp. 1–14, Springer, 2010.
- [21] Z. Li, J. Zheng, Z. Zhu, W. Yao, and S. Wu, "Weighted guided image filtering," *IEEE Transactions on Image processing*, vol. 24, no. 1, pp. 120–129, Jan. 2015.
- [22] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *European conference on computer vision*, pp. 391–407, Springer, 2016.
- [23] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1874–1883, 2016.
- [24] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3147–3155, 2017.
- [25] H. Zhao, O. Gallo, I. Frosio, and J. Kautz, "Loss functions for image restoration with neural networks," *IEEE Transactions on Computational Imaging*, vol. 3, no. 1, pp. 47–57, 2016.
- [26] M. S. Sajjadi, B. Scholkopf, and M. Hirsch, "Enhancenet: Single image super-resolution through automated texture synthesis," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4491–4500, 2017.
- [27] T. R. Shaham, T. Dekel, and T. Michaeli, "Singan: Learning a generative model from a single natural image," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4570–4580, 2019.
- [28] D. Fuoli, L. Van Gool, and R. Timofte, "Fourier space losses for efficient perceptual image super-resolution," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2360–2369, October 2021.
- [29] A. Lugmayr, M. Danelljan, L. Van Gool, and R. Timofte, "SrfLOW: Learning the super-resolution space with normalizing flow," in *European Conference on Computer Vision*, pp. 715–732, Springer, 2020.
- [30] W. Yang, J. Feng, J. Yang, F. Zhao, J. Liu, Z. Guo, and S. Yan, "Deep edge guided recurrent residual learning for image super-resolution," *IEEE Transactions on Image Processing*, vol. 26, no. 12, pp. 5895–5907, 2017.
- [31] X. Wang, K. Yu, C. Dong, and C. Change Loy, "Recovering realistic texture in image super-resolution by deep spatial feature transform," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 606–615, 2018.
- [32] M. Li, Z. Zhang, J. Yu, and C. W. Chen, "Learning face image super-resolution through facial semantic attribute transformation and self-attentive structure enhancement," *IEEE Transactions on Multimedia*, vol. 23, pp. 468–483, 2020.
- [33] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*, pp. 694–711, Springer, 2016.
- [34] R. Timofte, E. Agustsson, L. Van Gool, M.-H. Yang, and L. Zhang, "Ntire 2017 challenge on single image super-resolution: Methods and results," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 114–125, 2017.
- [35] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [36] A. Lugmayr, M. Danelljan, and R. 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*, pp. 494–495, 2020.
- [37] X. Ji, Y. Cao, Y. Tai, C. Wang, J. Li, and F. Huang, "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.
- [38] Z. Li, "Adaptive basic unit layer rate control for jvt," in *JVT 7th Meeting, Pattaya, Mar2003*, 2003.
- [39] Y. Liu, Z. G. Li, and Y. C. Soh, "A novel rate control scheme for low delay video communication of h. 264/avc standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 17, no. 1, pp. 68–78, 2006.