

# End-to-end Alternating Optimization for Real-World Blind Super Resolution

Zhengxiong Luo<sup>1,2,3</sup>, Yan Huang<sup>2,3\*</sup>, Shang Li<sup>1,2</sup>, Liang Wang<sup>2,3,4</sup> and Tieniu Tan<sup>2,3,5</sup>

<sup>1</sup>Artificial Intelligence School, University of Chinese Academy of Sciences, Beijing, 100049, China.

<sup>2</sup>Institute of Automation, Chinese Academy of Sciences, Beijing, 100190, China.

<sup>3</sup>Center for Research on Intelligent Perception and Computing, Beijing, China.

<sup>4</sup>Center for Excellence in Brain Science and Intelligence Technology, Beijing, China.

<sup>5</sup> Nanjing University, Nanjing, China .

\*Corresponding author(s). E-mail(s): [yhuang@nlpr.ia.ac.cn](mailto:yhuang@nlpr.ia.ac.cn);

Contributing authors: [zhengxiong.luo@cripac.ia.ac.cn](mailto:zhengxiong.luo@cripac.ia.ac.cn); [lishang2018@ia.ac.cn](mailto:lishang2018@ia.ac.cn);  
[wangliang@nlpr.ia.ac.cn](mailto:wangliang@nlpr.ia.ac.cn); [tnt@nlpr.ia.ac.cn](mailto:tnt@nlpr.ia.ac.cn);

## Abstract

Blind Super-Resolution (SR) usually involves two sub-problems: 1) estimating the degradation of the given low-resolution (LR) image; 2) super-resolving the LR image to its high-resolution (HR) counterpart. Both problems are ill-posed due to the information loss in the degrading process. Most previous methods try to solve the two problems independently, but often fall into a dilemma: a good super-resolved HR result requires an accurate degradation estimation, which however, is difficult to be obtained without the help of original HR information. To address this issue, instead of considering these two problems independently, we adopt an alternating optimization algorithm, which can estimate the degradation and restore the SR image in a single model. Specifically, we design two convolutional neural modules, namely *Restorer* and *Estimator*. *Restorer* restores the SR image based on the estimated degradation, and *Estimator* estimates the degradation with the help of the restored SR image. We alternate these two modules repeatedly and unfold this process to form an end-to-end trainable network. In this way, both *Restorer* and *Estimator* could get benefited from the intermediate results of each other, and make each sub-problem easier. Moreover, *Restorer* and *Estimator* are optimized in an end-to-end manner, thus they could get more tolerant of the estimation deviations of each other and cooperate better to achieve more robust and accurate final results. Extensive experiments on both synthetic datasets and real-world images show that the proposed method can largely outperform state-of-the-art methods and produce more visually favorable results. The codes are released at <https://github.com/greatlog/RealDAN.git>.

**Keywords:** blind super resolution, degradation estimation, alternating optimization, *Restorer*, *Estimator*.

## 1 Introduction

Single image super-resolution (SR) aims to recover the high-resolution (HR) version of a given

degraded low-resolution (LR) image. It has wide applications in video enhancement, medical imaging, as well as security and surveillance imaging. Generally, the degradation process can be formulated as

$$\mathbf{y} = [(\mathbf{x} \otimes \mathbf{k}) \downarrow_s + \mathbf{n}]_q \quad (1)$$

where  $\mathbf{x}$  is the original HR image,  $\mathbf{y}$  is the degraded LR image,  $\otimes$  denotes the two-dimensional convolution of  $\mathbf{x}$  with blur kernel  $\mathbf{k}$ ,  $\mathbf{n}$  denotes the random noise,  $\downarrow_s$  denotes the standard  $s$ -fold downsampler (keeping only the upper-left pixel for each distinct  $s \times s$  patch) [1], and  $[\cdot]_q$  denotes JPEG compression with quality factor  $q$ . Then SR refers to the process of recovering  $\mathbf{x}$  from  $\mathbf{y}$ . In a blind case, not only the HR image  $\mathbf{x}$ , but also the degradation parameters ( $\mathbf{k}$ ,  $\mathbf{n}$  and  $q$ ) are unknown, which makes blind SR a quite challenging task [2].

Most previous methods try to decompose blind SR into two relatively easier steps: 1) *degradation estimation*, and 2) restoring the SR image with the estimated degradation, which is called *non-blind SR*. Following this framework, degradation estimation and non-blind SR have been independently studied for years, and many successful methods have been proposed in both research fields, such as Michaeli *et al.* [3], KernelGAN [4] for degradation estimation and DPSR [5], USRNet [1] for non-blind SR.

However, two problems may exist in this framework: 1) The degradation-estimation model and the SR model are independently optimized. It is likely that the two models could not well cooperate, *i.e.* small deviations of degradation estimation may lead to terrible SR results. 2) Degradation estimation is inherently ambiguous due to the information loss during downscaling [6]. It is extremely difficult to get an accurate degradation in the absence of original HR information. A straightforward idea to address the second problem is using SR results to help improve the accuracy of degradation estimation. However, as we have described in the first problem, a good SR result also requires an accurate degradation in the first place. As a result, previous methods often fall into a dilemma: how to simultaneously improve the accuracy of degradation estimation and SR performance?

To break this dilemma, instead of considering these two steps separately, we adopt an alternating optimization algorithm, which can estimate

the degradation and restore the SR image in a single model. In detail, we design two convolutional neural modules, namely *Restorer* and *Estimator*. *Restorer* restores the SR image based on the degradation estimated by *Estimator*, and the restored SR image is further used to help *Estimator* estimate a more accurate degradation. Once the degradation is manually initialized, the two modules can well cooperate with each other to form a closed loop, which can be iterated over and over. We fix the number of iterations and unfold the iterating process to form an end-to-end trainable network, which we call the deep alternating network (DAN). To ensure the convergence of the iteration, DAN is directly supervised at the last iteration during training. Thus, both *Restorer* and *Estimator* may learn to substantially improve their results during the iterating and finally converge to a stable point. In the framework of DAN, *Estimator* can utilize the information of intermediate SR results, which makes the degradation estimation easier. More importantly, as *Restorer* and *Estimator* are jointly optimized, they may get more tolerant of the deviations of each other and cooperate well to achieve a better final result.

We need to note that a preliminary version of this work has been presented as a conference paper [7], which is denoted as DAN-Pre in this paper. In the current version, we incorporate additional content in significant ways.

Firstly, DAN-Pre can only process blurry LR images (with mild additive white Gaussian noise (AWGN)). While in DAN, we consider much more complex degradations, including multiple blur, resize, noise, and JPEG compression. To deal with such complex degradations, we parameterize the whole degrading pipeline and each random degradation can be represented by a unique vector that is computable for both *Restorer* and *Estimator*.

Secondly, in DAN-Pre, the *Restorer* and *Estimator* are iterated in the image space and degradation space respectively, while in DAN, the two modules are iterated in the feature space. In this way, richer information can be passed between different iterations, which may make the training process more stable and get the whole network better optimized. Moreover, iterating in the feature space also saves computations as we do not need to compute the output in each iteration.

Thirdly, we reorganize the experiments and make more comprehensive comparisons with the

most recent state-of-the-art methods. We also add more experiments to better analyze the proposed method.

We summarize our contributions into following points:

1. We propose an alternating optimization algorithm that considers the degradation estimation and SR in a single network, in which way, both modules can utilize the intermediate results of each other and could get well compatible to produce better final results than previous two-step solutions.
2. To the best of our knowledge, the preliminary version of this work [7] proposes the first end-to-end network for blind SR, which largely simplifies and accelerates the training and inference of blind SR methods.
3. We parameterize the complex degradation process (including multiple blur, resize, noise, and compression) and make it computable for convolutional networks. To the best of our knowledge, the proposed DAN is the first network that can simultaneously estimate the complex degradations and restore SR images for real-world blind SR.
4. We design two convolutional neural modules, which can be alternated repeatedly to estimate the degradation and restore the SR image respectively.
5. Extensive experiments on synthetic and real-world images show that our model can largely outperform state-of-the-art blind-SR methods and produce more visually favorable results.

## 2 Related Work

In the recent decade, deep-learning (DL)-based SR methods [8–13] have made remarkable achievements and have shown great advantages against traditional methods [14–16]. Thus, we mainly discuss DL-based SR methods in this paper.

### 2.1 SR for Bicubically Downscaled Images

DL-based SR methods usually require a large number of paired HR-LR images as training samples. However, these paired samples are difficult to be collected in the real world. Consequently, synthetic data is usually used as an alternative. Early

researchers consider only the simplest case, *i.e.* the LR images are obtained by downscaling the HR images with bicubic interpolation [17]. In this way, a large number of training samples can be cheaply synthesized. In this case, as data is not the concern, most researchers concentrate on designing the structures of SR networks. In SRCNN [17], Dong *et al.* propose the first convolutional neural network (CNN) for SR, which has only three convolutional layers. In the following years, many CNN-based SR methods [18–22] have been proposed and strategies such as post-upsampling [23], residual learning [24], and pixel-shuffle operation [25] nearly become the default choices for building an SR network. After the proposal of RCAN [8], RRDB [9] and SAN [26], the CNN-based SR performance even starts to get saturate on common benchmark datasets. Recently, some transformer [27]-based methods, such as IPT [28] and SwinIR [10] further advance the performance for bicubically downsampled images.

However, although these methods perform well for super-resolving bicubically downsampled images, it is still difficult for them to get applied in real scenarios. The degradations in real scenarios are various and unknown, which are much more complex than the bicubic downscaling. Consequently, due to the domain gap between real and synthesized data, methods designed for bicubically downsampled images will suffer performance drop in real applications [29, 30]. To address this issue, researchers begin to work on more challenging cases where degradations of test images are unknown, *i.e.* blind SR.

### 2.2 SR for Blurry Downsampled LR Images

Early blind SR methods walk a step further in practicability: they also consider cases where LR images are blurred by various kernels (sometimes AWGN is also considered) before downscaling. Compared with SR for bicubically downsampled images, apart from the HR image  $\mathbf{x}$ , there is one more unknown variable, *i.e.* the blur kernel  $\mathbf{k}$ . Thus, this problem usually involves two steps: kernel estimation and SR with the given kernel. Some methods focus on only one of them and some try to solve them simultaneously.

**Kernel Estimation.** Some methods only focus on kernel estimation. As this is an ill-posed problem [31, 32], some priors are usually needed to get it properly solved. In [3], a non-parametric method is used by utilizing the patch recurrence between the test image and its downsampled version. A similar idea is also adopted in [4], but powered with neural networks and adversarial training [33]. Another widely used prior is the extreme channel priors. In [34, 35], Pan *et al.* firstly propose the dark channel prior, *i.e.* the dark channel in a natural image is usually sparse, which can be used for solving the blur kernel from a blurred image. In [36, 37], the bright channel prior is further proposed and the idea is augmented to extreme channel priors.

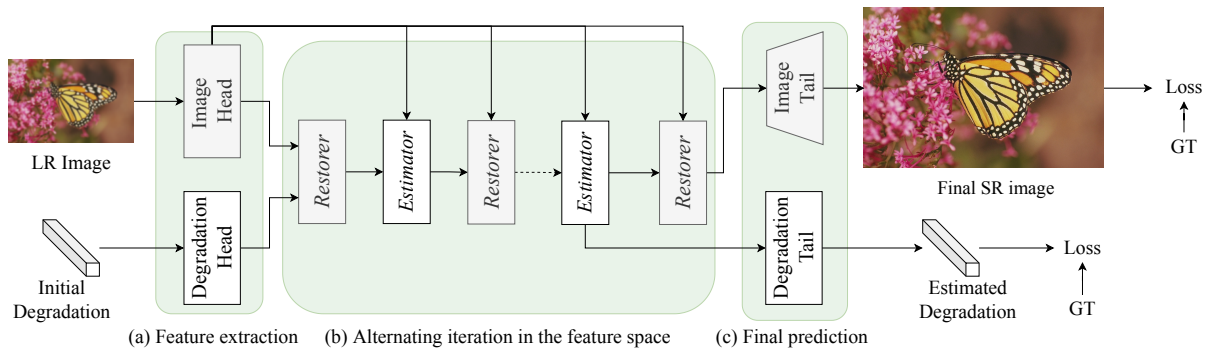
**SR with given kernels.** Some methods assume that the kernel estimation has been solved and only focus on SR with given kernels. In [38–40], the blur kernel is used to down sample images and synthesize training samples, which can be used to train a specific model for a given kernel and LR image. In SRMD [41], the kernel and LR image are directly concatenated at the first layer of a CNN. Thus, the SR result can be closely correlated to both LR image and blur kernel. In DPSR [5], Zhang *et al.* propose a method based on the ADMM algorithm. They interpret this problem as MAP optimization and solve the data term and prior term alternately. A similar idea is adopted in USRNet [1]. These methods can achieve remarkable performance as long as the ground-truth blur kernel is known. However, in real applications, the blur kernels are predicted by kernel-estimation methods and are likely to be biased from the ground-truth ones. This bias will cause a serious performance drop when the kernel-estimation step and the SR step are combined.

**Integrated SR with Kernel Estimation.** To avoid the performance drop when the two steps are combined, recent methods try to integrate them into a single model, which consists of a kernel-estimation module and a SR module. In [42], Gu *et al.* propose to finetune the kernel-estimation module together with the SR module to get them more compatible. Its improvement is evident, but the training and inference are complicated and time-consuming. In the preliminary version of this paper [7], we adopt an alternating optimization algorithm and propose the first end-to-end

method for integrated SR with kernel estimation, which further improves the compatibility of the two modules. And since then the end-to-end framework becomes the prevalent choice for blind SR methods [43, 44]. In DASR [45], the kernel representations are firstly extracted via contrastive learning and then input to a SR module. In DCLS [6], the kernel is firstly estimated by the dynamic deep linear module and then input into a deep-constrained-least-squares module to help restore the SR image. However, the kernel-estimation module in these methods utilizes only information from the LR images, which may limit the accuracy of kernel estimation and the final SR performance. While in DAN, the accuracy of degradation estimation may be improved with the help of intermediate SR results.

## 2.3 SR for Real Images

Although the blurry LR image is a better assumption than the bicubically downsampled one, it may still be far from real images [46–48]. The domain gap between training and testing data will largely destroy the practicability of SR methods. To address this issue, some researchers try to manually collected paired HR-LR samples [49, 50]. However, it is expensive and time-consuming, and there may still have a domain gap between images collected by different cameras. Recently, in BSRGAN [47] and Real-ESRGAN [46], researchers propose to synthesize training samples by much more complex degradations, including multiple blur, resize, noise, and compression. In this way, the synthesized data may be diverse enough to include most cases in real scenarios. And the SR model trained with these samples may be practical enough in applications. However, the two methods are degradation-unaware, *i.e.* directly super-resolving the LR image regardless of its degradation. Consequently, they may fail to exploit a more general relationship between SR under various degradations and could only achieve sub-optimal results. While in the current version of DAN, we re-parameterize the complex degradation process and adopt an alternating optimization algorithm to simultaneously do SR and degradation estimation, which may lead the SR module to be more degradation-specific and achieve better performance.



**Fig. 1:** The overview of the deep alternating network (DAN).

### 3 End-to-End Alternating Optimization

#### 3.1 Formulation

Generally, the degradation process shown in Equation 1 can be formulated as:

$$\mathbf{y} = \mathcal{D}(\mathbf{x}; \theta), \quad (2)$$

where  $\mathcal{D}$  is the degradation function, and  $\theta$  is the involved parameter. Then the blind SR problem can be mathematically expressed an optimization problem in the Maximum A Posterior (MAP) framework [51]:

$$\mathbf{x}, \theta = \arg \min_{\mathbf{x}, \theta} \|\mathbf{y} - \mathcal{D}(\mathbf{x}; \theta)\|_1. \quad (3)$$

However, as there are too many unknown variables, this optimization problem is still ill-posed and has an infinite number of solutions [2]. To get it properly solved, some prior terms are usually added [52, 53]:

$$\mathbf{x}, \theta = \arg \min_{\mathbf{x}, \theta} \|\mathbf{y} - \mathcal{D}(\mathbf{x}; \theta)\|_1 + \phi(\mathbf{x}) + \psi(\theta), \quad (4)$$

where  $\phi(\mathbf{x})$  denotes the prior for HR image, and  $\psi(\theta)$  represents the prior for degradation parameter.

#### 3.2 Two-step Solution

In the context of SR for blurry downsampled images,  $\mathcal{D}$  represents blur and downscaling. And  $\theta$  represents the blur kernel  $\mathbf{k}$  and downscaling factor  $s$ . In many methods,  $\mathbf{k}$  is assumed to be a Gaussian

kernel [42, 45]. Such prior makes it easier to solve the degradation parameter, which can be further used to solve the HR image. In this case, previous usually adopt a two-step solution:

$$\begin{cases} \theta = \mathcal{K}(\mathbf{y}) \\ \mathbf{x} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathcal{D}(\mathbf{x}; \theta)\|_1 + \phi(\mathbf{x}) \end{cases} \quad (5)$$

where  $\mathcal{K}(\cdot)$  denotes the function that estimates  $\theta$ , *i.e.* blur kernel  $\mathbf{k}$  in this case, from  $\mathbf{y}$ . And the second step is usually solved by a non-blind SR method such as DPSR [5], USRNet [1], *etc.*

As we have mentioned in Sec 2.2, the two steps are independent research fields in most cases. Both of them only consider the performance under their own given conditions, while ignoring the overall performance. This two-step solution has its drawbacks threefold. Firstly, this algorithm usually requires the training of two or even more models, which is rather complicated. Secondly,  $\mathcal{K}(\cdot)$  can only utilize information from  $\mathbf{y}$ . However, this is also an ill-posed problem:  $\theta$  could not be properly solved without information from  $\mathbf{x}$ . At last, the non-blind SR model for the second step is trained with ground-truth degradations. While during testing, it only has access to degradations estimated in the first step. The difference between ground-truth and estimated degradations will usually cause a serve performance drop of the non-blind SR model [42]. Moreover, in cases of more complex degradations, the deviations of estimated degradations are likely to be larger, which may further destroy the SR performance of the second step.

### 3.3 Unfolding the Alternating Optimization

Towards the drawbacks of two-step solution, we propose an end-to-end network that can largely alleviate these issues. Specifically, we still split it into two subproblems. But instead of solving them in sequential, we adopt an alternating optimization algorithm, which restores the SR image and estimates the degradation alternately. The mathematical expression is

$$\begin{cases} \theta_{i+1} = \arg \min_{\theta} \|\mathbf{y} - \mathcal{D}(\mathbf{x}_i; \theta)\|_1 + \psi(\theta) \\ \mathbf{x}_{i+1} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathcal{D}(\mathbf{x}; \theta_i)\|_1 + \phi(\mathbf{x}). \end{cases} \quad (6)$$

We define two solvers, namely *Estimator* and *Restorer* for the two subproblems respectively. And The blind SR problem may be solved by repeatedly iterating the two solvers. Since the prior terms, *i.e.*  $\psi(\theta)$  and  $\phi(\mathbf{x})$ , are difficult to be mathematically expressed (unless in some simple cases), it is also hard to find the analytic solutions for both solvers. Given the remarkable achievements of DL networks, we try to construct both solvers with convolutional neural modules. We hope that once the two modules are trained, they could automatically solve the sub-problems respectively.

As shown in Fig 1 (a), we firstly initialize the degradation as  $\theta_0$  (set as learnable and initialized as null-vector in our experiments) and then encode the LR image and initial degradation to the feature space with two head modules:

$$\mathbf{f}_0^x = H^x(\mathbf{y}), \quad \mathbf{f}_0^\theta = H^\theta(\theta_0), \quad (7)$$

where  $H^x$  and  $H^\theta$  are the head modules for images and degradations respectively,  $\mathbf{f}_0^x$  and  $\mathbf{f}_0^\theta$  are the initial features of images and degradations respectively. As shown in Fig 1 (b), we then iterate the *Restorer*  $R$  and the *Estimator*  $E$  in the feature space to alternatively solve the features for images and degradations:

$$\begin{cases} \mathbf{f}_{i+1}^x = R(\mathbf{y}, \mathbf{f}_i^\theta) \\ \mathbf{f}_{i+1}^\theta = E(\mathbf{y}, \mathbf{f}_i^x). \end{cases} \quad (8)$$

After  $T$  iterations, the image features are super-resolved to the HR image, and the degradation

features are regressed to the estimated degradation (shown in Fig 1 (c)):

$$\mathbf{x} = T^x(\mathbf{f}_T^x), \quad \theta = T^\theta(\mathbf{f}_T^\theta), \quad (9)$$

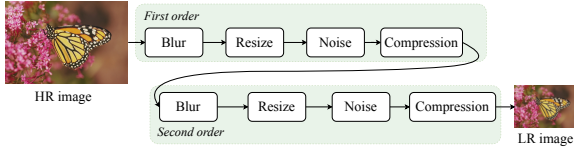
where  $T^x$  and  $T^\theta$  are the tail modules for the images and degradations respectively. We fix the number of iterations and unfold the iterating process. Then the whole pipeline can form an end-to-end trainable network, which is called the deep alternating network (DAN). Since we hope that the final iteration results could converge to a stable point, we supervise DAN at the end by the ground-truth degradation and HR image.

### 3.4 Discussion

Compared with previous two-step methods, DAN has its benefits threefold: 1) DAN can be trained and tested in an end-to-end manner, which largely simplifies and accelerate blind SR methods; 2) the *Estimator* can utilize the information from intermediate SR results, which makes the degradation estimation easier; 3) *Estimator* and *Restorer* are jointly optimized, and thus they may get more compatible with each other and cooperate better to achieve preferable final results.

It should be noted that despite both DAN and IKC utilizing an iterative approach, DAN and IKC both adopt an iterative strategy, their optimization algorithm differs significantly. In IKC, the SR module is initially trained as a non-blind SR model and then kept fixed during the training of the kernel correction module. This sequential training process in IKC can be cumbersome. Additionally, the SR module in IKC only has access to ground-truth degradations during its training phase, rather than estimated degradations. Consequently, the compatibility between the SR module and kernel-correction module may be compromised, leading to sub-optimal performance. In contrast, in DAN, the *Restorer* and *Estimator* are jointly optimized in an unfolded network. Thus both modules are trained in an end-to-end manner, simplifying the training process. Moreover, the SR module (*i.e.* the *Restorer*) in DAN is trained with estimated degradations, which may lead to better compatibility with the *Estimator* and potentially achieve superior results.

After the proposal of the preliminary version of this work, many other end-to-end blind



**Fig. 2:** The multi-order degradation model.

SR methods [6, 45] are also proposed. However, their degradation-estimation modules usually utilize only the limited information from the LR image, which may restrict their performance.

Compared with the preliminary version of this work, the current framework has two main differences. Firstly, the initial degradation is manually set in the preliminary version while it is learnable in the current version. In this way, it is more likely that DAN could adaptively find a good initial point for the alternating optimization, which may help improve the final results. Secondly, in the preliminary version, the information of different iterations can only be passed through the SR image (degradation). While in the current version, *Estimator* and *Restorer* are iterated in the feature space and predict the final results only at the last iteration. In this way, richer information could be passed through intermediate features. And this strategy could also reduce the redundant computations of final results.

### 3.5 Reparamization of the Degradation

Previous methods [7, 42] consider only blurry LR images, in which case the degradation is parameterized as the blur kernel. While in real applications, degradations are much more complex and are difficult to be parameterized and calculated. In Real-ESRGAN [46], Wang *et al.* propose a multi-order degradation model (shown in Fig 2), which includes multiple blur, resize, noise, and compression. Although each single degradation operation is simple, they could be integrated to well simulate the complex real degradations. We adopt the same degradation model in work and we try to reparameterize this model by parameterizing each degradation operation.

#### 3.5.1 Blur

In the multi-order degradation model, the original HR image is firstly blurred. In [46], three kinds of

blur kernels are considered, including generalized Gaussian kernels [54], generalized plateau kernels, and what they call *sinc* kernels. The generalized kernel can be mathematically expressed as

$$\mathbf{k}(i, j) = \exp\left(-\frac{1}{2}([i \ j] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} i \\ j \end{bmatrix})^\beta\right), \quad (10)$$

where  $i, j$  is coordinate of  $\mathbf{k}$ ,  $\boldsymbol{\Sigma}$  is the covariance matrix, and  $\beta$  is the shape parameter. If  $\beta = 1$ ,  $\mathbf{k}$  becomes a common Gaussian kernel. Similarly, the plateau kernel can be mathematically expressed as

$$\mathbf{k}(i, j) = \frac{1}{1 + ([i \ j] \boldsymbol{\Sigma}^{-1} \begin{bmatrix} i \\ j \end{bmatrix})^\beta}. \quad (11)$$

For both Gaussian and plateau kernels, the covariance matrix can be expressed as:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_x & 0 \\ 0 & \sigma_y \end{bmatrix} \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}. \quad (12)$$

Thus, either the Gaussian kernel or plateau kernel can be parameterized as four parameters, *i.e.*  $\sigma_x$ ,  $\sigma_y$ ,  $\theta$ , and  $\beta$ . The *sinc* kernel is expressed as

$$\mathbf{k}(i, j) = \frac{\omega_c}{2\pi\sqrt{i^2 + j^2}} J_1(\omega_c\sqrt{i^2 + j^2}), \quad (13)$$

where  $\omega_c$  is the cutoff frequency, and  $J_1$  is the first order Bessel function of the first kind. Thus, the *sinc* kernel can be parameterized as one parameter  $\omega_c$ .

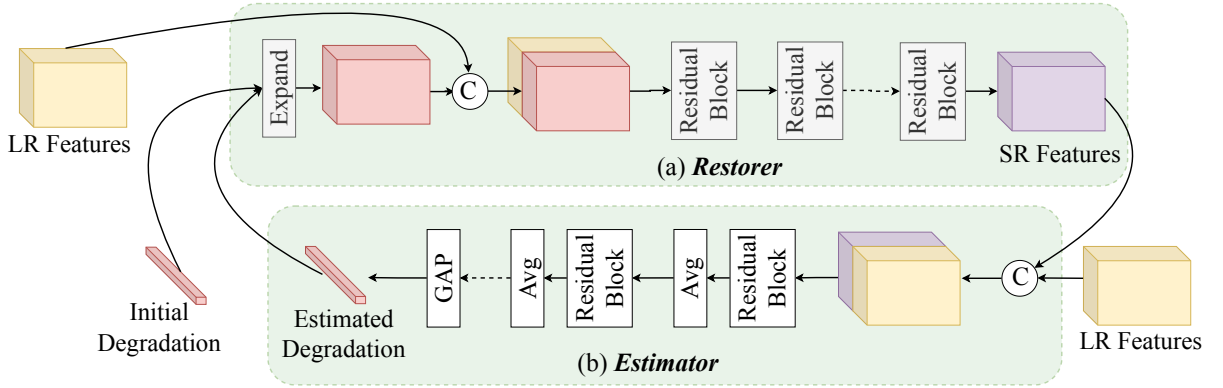
We further define a variable  $k_s$  to indicate the kernel size and a vector  $[k_g, k_c]$  to indicate the type of the blur kernel. Specifically,  $k_g = 0, k_c = 0$  indicates a plateau kernel,  $k_g = 1, k_c = 0$  indicates a Gaussian kernel, and  $k_g = 0, k_c = 1$  indicates a *sinc* kernel. Then, the blur kernel can be parameterized as

$$\mathbf{B} = [k_g, k_c, k_s, \sigma_x, \sigma_y, \theta, \beta, \omega_c]. \quad (14)$$

We need to note when the  $\mathbf{b}$  represents a *sinc* kernel, we set  $\sigma_x = \sigma_y = \theta = \beta = 0$ , since a *sinc* kernel dose not have those parameters.

#### 3.5.2 Resize

After being blurred, the image is then resized. In [46], there are three different modes of resizing,



**Fig. 3:** The frameworks of (a) *Restorer* and (b) *Estimator*. © denotes concatenating in the channel dimension. And ‘Expand’ means expanding the degradation features in spatial dimensions. ‘Avg’ denotes average pooling, and ‘GAP’ denotes global average pooling.

namely *area*, *bilinear*, and *bicubic*. Each image is resized via the mode randomly chosen from them. We use a one-hot vector  $[r_{area}, r_{bil}, r_{bic}]$  to indicate the resizing mode. In detail,  $[1, 0, 0]$  indicates the *area* mode,  $[0, 1, 0]$  indicates the *bilinear* mode, and  $[0, 0, 1]$  indicates the *bicubic* mode. We further denote the scale factor as  $s$ . And the resizing operation can be parameterized as:

$$\mathbf{R} = [r_{area}, r_{bil}, r_{bic}, s]. \quad (15)$$

### 3.5.3 Noise

In both Real-ESRGAN and BSRGAN, two kinds of noises are considered, *i.e.* the Gaussian noise and the Poisson noise. The Gaussian noise  $\mathbf{n}_g \sim \mathcal{N}(0, \sigma_g)$  is parameterized by the standard deviation  $\sigma_g$ . And the Poisson noise  $\mathbf{n}_p$  is parameterized by the noise level  $\lambda$  [46, 55]. We define a variable  $n_t$  to indicate the type of the noise. Specifically,  $n_t = 1$  indicates Gaussian noise, and  $n_t = 0$  indicates Poisson noise. There are also two kinds of color for the noise, *i.e.* gray and RGB. Thus, we further define a variable  $n_c$  to indicate the noise color (1 for RGB and 0 for gray). Finally, the noise  $\mathbf{n}$  can be parameterized as

$$\mathbf{N} = [n_t, n_c, \sigma_g, \lambda]. \quad (16)$$

We need to note when  $n$  represents the Gaussian noise, we have  $\lambda = 0$  since there is no Poisson noise. Similarly, when  $n$  represents the Poisson noise, we also have  $\sigma_g = 0$ .

### 3.5.4 JPEG Compression

In real applications, some LR images will be compressed by the JPEG compression [56]. We define a variable  $j$  to indicate whether the image is compressed (1 for yes and 0 for no). For the compressed images, the JPEG compression can be parameterized by the quality factor  $q$  (sampled in  $[1, 100]$ ). Thus, the compression operation  $\mathbf{J}$  can be parameterized as

$$\mathbf{J} = [j, q]. \quad (17)$$

We need to note when  $j = 0$ , we have  $q = 100$ , which indicates an uncompressed image.

### 3.5.5 Multi-order Degradation

As shown in Fig 2, in the multi-order degradation model, the HR image will go through the *blur-resize-noise-compression* pipeline for twice before it is degraded to a LR image. As we have discussed above, this pipeline can be parameterized as

$$\mathbf{D} = [\mathbf{B}, \mathbf{R}, \mathbf{N}, \mathbf{J}], \quad (18)$$

which is a 1-dimension vector. Thus, the whole degradation model can be parameterized as

$$\theta = [\mathbf{D}_1, \mathbf{D}_2], \quad (19)$$

where  $\mathbf{D}_1$  and  $\mathbf{D}_2$  represent the parameters of the first and the second order degradation respectively, and  $\theta$  is also a 1-dimension vector. In this



way, each degradation can be uniquely represented and calculated.

### 3.6 Instantiation

**Head modules.** As we have described in Sec 3.3 and Fig 1, the features of the LR image and degradation are firstly extracted by two *head* modules respectively. The image head in Fig 1 is constructed by one convolutional layer. While the degradation head is constructed by one fully connected layer since the degradation is a one-dimension vector. The extracted shallow features are then input to the *Restorer* and *Estimator*.

***Restorer* and *Estimator*.** Both of the *Restorer* and *Estimator* need to cope with two inputs. In the preliminary version [7], we propose a conditional residual block to help the two modules exploit the inter information of their two inputs and ensure that the output of the *Restorer* (*Estimator*) is closely related to both of its inputs. However, we later experimentally find that the same effect can be achieved by simple concatenation at the beginning of *Restorer* (*Estimator*).

For the *Restorer*, as shown in Fig 3 (a), we simply concatenate the LR feature and the degradation feature at the beginning and then restore them to the SR features via the body of the *Restorer*. We need to note that the degradation feature is a one-dimension vector (Eq 19). Thus we need to expand it spatially before it is concatenated with the LR feature. In this way, the body of *Restorer* can take the benefits of the architectures in recent state-of-the-art SR models, such as EDSR [57], RCAN [8], and RRDB [9]. For simplicity, we use 16 residual blocks proposed in EDSR to construct the body of *Restorer*.

For the *Estimator*, as shown in Fig 3 (b), the two inputs, *i.e.* the LR features and the SR features, are also concatenated at the beginning. The body of *Estimator* is constructed by residual blocks and average pooling. Since the degradation estimation usually requires global information, especially for estimation of noise level and quality factor of jpeg compression. Thus each residual block is followed by an average pooling (AvgP) layer to enlarge the receptive fields of *Estimator*. And at the end of *Estimator*, an global average pooling (GAP) layer is used to predict the degradation features.

Compared with the DAN-Pre, the architectures of *Restorer* and *Estimator* are largely simplified, which further accelerates the training and testing of DAN (which will be discussed in Sec 4.4). While experiments in Sec 4.2 show that the DAN can achieve even better results than DAN-Pre.

**Tail modules.** As shown in Fig 1, the SR feature and the degradation feature are input into two *tail* modules. The image tail upscales the SR feature and reconstructs it to the SR image. While the degradation tail projects the degradation feature to the final estimated degradation. The image tail consists of a PixelShuffle [25] layer and several convolutional layers. And the degradation tail consists of only fully connected layers.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets** Following previous blind SR methods, DAN is also trained on datasets synthesized from the HR images in the training set of DIV2K [59] and Flickr2K [60]. The two datasets contain 3450 HR images (800 from DIV2K and 2650 from Flickr2K). Based on the discussion in Sec 3.5, DAN can deal with LR images with various degradations. However, most previous SR methods can only handle blurry LR images. Thus, to make comprehensive comparisons with other methods, we train DAN with two different settings.

Under the first setting, DAN is trained and evaluated on blurry LR images. Following the setting in [4, 6, 45], The training dataset is synthesized by various random anisotropic Gaussian kernels. For scale factor  $\times 2$  and  $\times 4$ , the kernel size set as  $11 \times 11$  and  $31 \times 31$  respectively. The  $\beta$  in Eq 10 is set as 1,  $\sigma_x$  and  $\sigma_y$  are uniformly sampled in (0.6, 5.0), and  $\theta$  is uniformly sampled in  $[-\pi, \pi]$ . We also apply uniform multiplicative noise (up to 25% of each pixel value of the kernel) and normalize it to sum to one. Under this setting, DAN is evaluated on the DIV2KRK [4] dataset.

Under the second setting, DAN is trained on the dataset synthesized by more complex degradations. We use the same degradation model in Real-ESRGAN [46] to synthesize HR-LR training pairs. The degradation is parameterized as we described in Sec 3.5. In this case, DAN is evaluated on real images from the track2 of NTIRE2020 [61]

**Table 1:** Quantitative comparison with state-of-the-art methods for blurry LR images. Average PSNR and SSIM results on DIV2KRK [4] are reported.  $\times 2$  and  $\times 4$  denote the scale factors.  $\uparrow$  denotes the higher the better. The best two results are indicated in bold and underlined respectively.

Methods	$\times 2$		$\times 4$	
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
<i>Two-step methods</i>				
KernelGAN [4]+ZSSR [39]	30.36	0.8669	26.81	0.7316
KernelGAN [4]+SRMD [41]	29.57	0.8564	27.51	0.7265
KernelGAN [4]+USRNet [1]	-	-	20.06	0.5359
Michaeli [3]+SRMD [41]	25.51	0.8083	23.34	0.6530
Michaeli [3]+ZSSR [39]	29.37	0.8370	26.09	0.7138
<i>End-to-end methods</i>				
Bicubic	28.73	0.8040	25.33	0.6795
EDSR [57]	32.42	0.9034	28.68	0.7883
IKC [42]	31.20	0.8767	27.69	0.7657
AdaTarget [58]	-	-	28.42	0.7854
KOALAnet [43]	31.89	0.8852	27.77	0.7637
DCLS [6]	<u>32.75</u>	<u>0.9094</u>	<b>28.99</b>	<u>0.7946</u>
DAN-Pre [7]	32.56	0.8997	27.55	0.7582
DAN	<b>32.96</b>	<b>0.9114</b>	<u>28.90</u>	<b>0.7961</b>

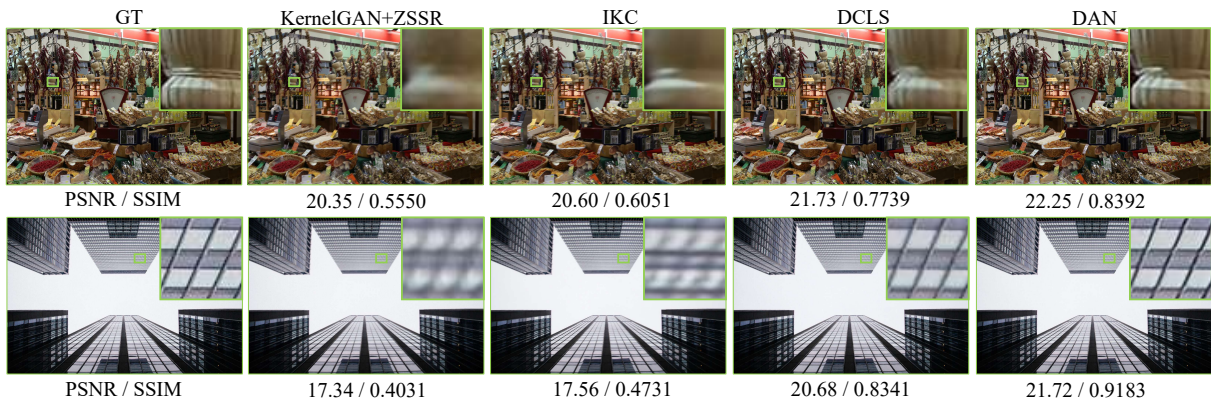


**Fig. 4:** Visual comparisons between different methods on DIV2KRK for scale factor  $\times 2$ . Best viewed in color.

(denoted as 2020Track2), and the RealSRSet collected in [47]. The 2020Track2 contains 100 real images taken by iPhone, and RealSRSet contains 20 real images either downloaded from the internet or directly chosen from existing testing datasets [62–65]. However, those real images do not have corresponding ground-truth HR images,

which makes it difficult to quantitatively measure the performance of different methods. Thus, for more comprehensive comparisons, we also synthesize a validation dataset, namely DIV2K-Real, with the same degradation model in [46].

**Evaluation metrics.** For validation sets that have ground-truth HR images, we use PSNR, SSIM [66], and LPIPS [67] (the network is set

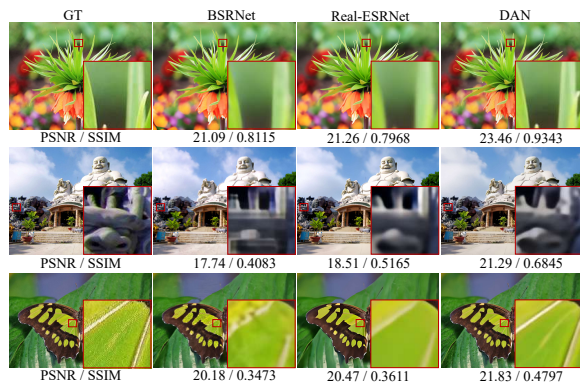


**Fig. 5:** Visual comparisons between different methods on DIV2KRC for scale factor  $\times 4$ . Best viewed in color.

as AlexNet [68]) to evaluate the performance of different methods. We need to note that both PSNR and SSIM are calculated on the Y channel (*i.e.* luminance) of transformed YCbCr space. For those real images that have no ground-truth references, the non-reference metrics, *i.e.* NIQE [69], NRQM [70], and PI [71] are used.

**Loss functions.** As shown in Fig 1, DAN is supervised by the ground-truth HR image and degradation. The supervision of HR image is applied via L1 loss, while the supervision of degradation is applied via L2 loss. In this paper, we also study the GAN [33, 72] version of DAN, which is denoted as DAN-GAN. Following the setting in RealESRGAN [46], DAN-GAN is also supervised by additional perceptual loss and adversarial loss. The perceptual loss is calculated in the feature space of VGG19 [73]. The discriminator that is used to calculate the adversarial loss adopts the same U-Shape architecture in [46].

**Training details.** During training, the LR images are cropped into patches of  $48 \times 48$ . For scale factor  $\times 4$ , the HR images are cropped into  $192 \times 192$ , and for scale factor  $\times 2$ , the HR images are cropped into  $96 \times 96$ . The batch size is set as 64. The model is trained for  $6 \times 10^5$  steps. The learning rate is initialized as  $2 \times 10^{-4}$  and is decayed by half every  $2 \times 10^5$  steps. We use Adam [74] as the optimizer. And all models are trained on 4 RTX 3090 GPUs.



**Fig. 6:** Visual comparisons between different methods on DIV2K-Real for scale factor  $\times 4$ . Best viewed in color.

## 4.2 Comparisons on Blurry LR images

In this subsection, we explore the performance of DAN for blurry LR images. We mainly compare two types of methods: the two-step methods and end-to-end methods. A two-step method is usually the combination of a degradation estimation method and a non-blind SR method. Thus, we denote the two-step methods as 'A+B', where A is the degradation estimation method, and B is the non-blind SR method. In this section, we mainly compare with five kinds of combinations, including KernelGAN [4]+ZSSR [39], KernelGAN+SRMD [41], KernelGAN+USRNet [1], Michaeli *et al.* [3]+SRMD, and Michaeli *et al.* +ZSSR. For the end-to-end methods, we also mainly compare with five methods, *i.e.* EDSR [57],

**Table 2:** Quantitative comparison with state-of-the-art methods for real-world LR images. Average results for  $\times 4$  models are reported on DIV2K-Real, the track2 of NTIRE2020, and RealSet [47].  $\uparrow$  denotes the higher the better, and  $\downarrow$  denotes the lower the better. The best two results are indicated in bold.

Methods	DIV2K-Real			2020Track2			RealSRSet		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	NIQE $\downarrow$	NRQM $\uparrow$	PI $\downarrow$	NIQE $\downarrow$	NRQM $\uparrow$	PI $\downarrow$
BSRNet [47]	25.12	0.6774	0.5113	8.475	3.447	7.514	8.830	<b>4.416</b>	6.807
Real-ESRNet [46]	24.75	0.6834	<b>0.4917</b>	8.083	3.653	7.215	8.082	4.300	6.892
DAN	<b>25.47</b>	<b>0.6886</b>	0.5036	<b>7.891</b>	<b>3.673</b>	<b>7.109</b>	<b>7.367</b>	3.958	<b>6.704</b>
BSRGAN [47]	24.28	0.6306	0.3861	<b>4.761</b>	6.164	4.300	5.468	6.254	4.607
Real-ESRGAN [46]	23.45	0.6291	<b>0.3527</b>	5.114	5.600	4.758	5.615	5.988	4.813
Ji <i>et al.</i> (DPED) [75]	20.96	0.4102	0.6374	4.866	<b>6.742</b>	<b>4.062</b>	<b>4.567</b>	<b>6.473</b>	<b>4.047</b>
DAN-GAN	<b>24.91</b>	<b>0.6591</b>	0.4069	5.871	5.288	5.291	5.019	6.300	4.359

**Table 3:** Comparisons of complexity and speed of different methods. The best results are denoted in bold.

Methods	# Params (M)	Multi-Adds (G)	Speed (s/image)
KernelGAN [4]+ZSSR [39]	0.30	-	120.27
IKC [42]	5.29	2178.72	1.14
DCLS [6]	13.63	368.15	0.18
RealESRGAN [46]	16.70	871.24	0.17
DAN-Pre [7]	4.33	926.72	0.21
DAN	1.95	<b>221.61</b>	<b>0.05</b>

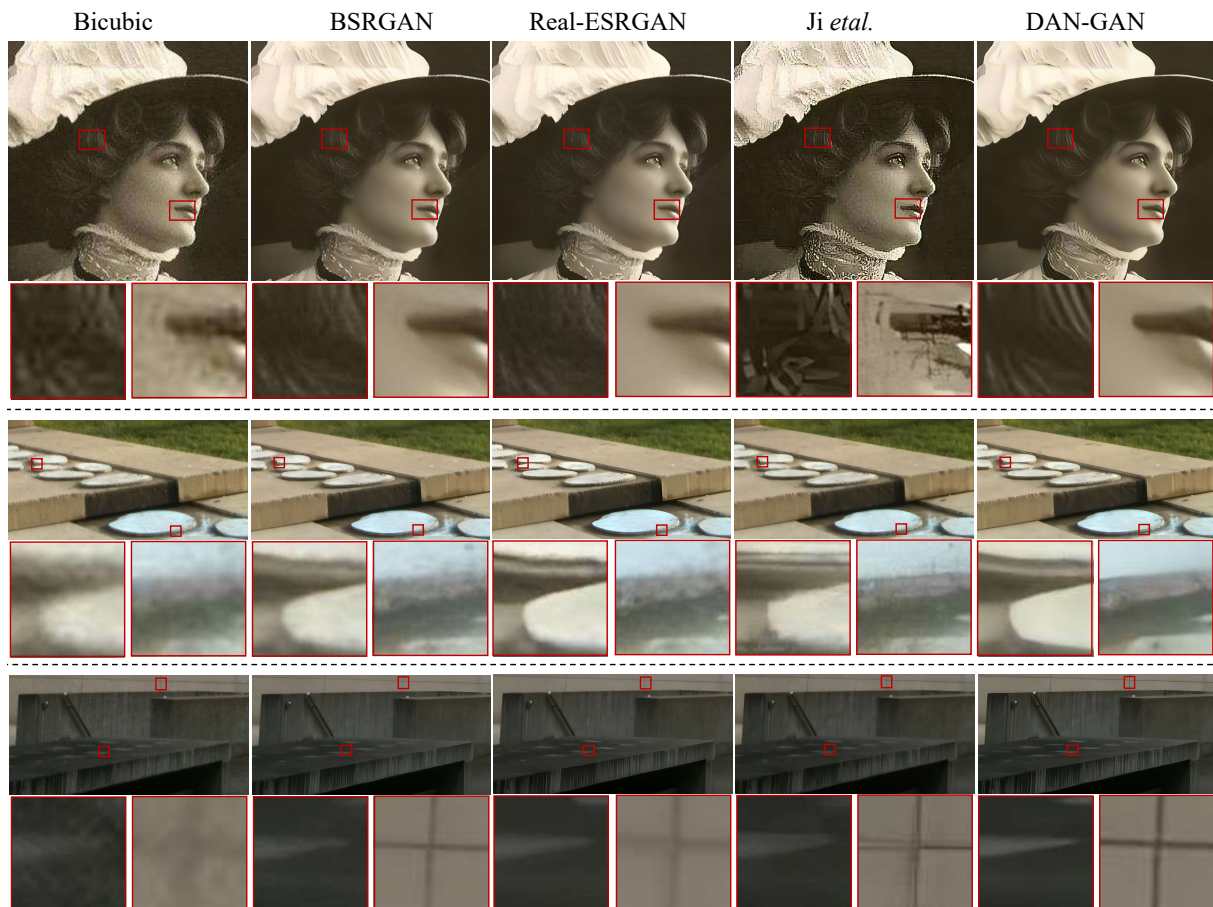
IKC [42], AdaTarget [58], KOALANet [43], and DCLS [6]. The result of bicubic interpolation is also added as a reference. We also include a comparison with the preliminary version of DAN, which is denoted as DAN-Pre. We need to note that the original EDSR is trained with bicubically synthesized samples. To make a comprehensive comparison, we retrain EDSR in the same way that we train DAN here. As we have discussed in Sec 3.6, the *Restorer* body in DAN can be chosen freely. To make the comparison with EDSR fair, we directly use the body of EDSR as the *Restorer* body in DAN.

**Quantitative comparisons.** We evaluate the reference methods on DIV2KRRK, and the results are shown in Table 1. USRNet [1] and SRMD [41] are the leading non-blind SR methods. They perform excellently as long as the ground-truth blur kernels are provided. However, in real applications, the blur kernels are estimated, which may deviate from the ground-truth ones. Consequently, as shown in the table, KernelGAN+USRNet and Michaeli *et al.* +SRMD can only achieve results even worse than simple bicubic interpolation

(20.06 dB and 23.34 dB V.S. 25.33 dB for scale factor  $\times 4$ ). This may be because the two parts of the two-step method are separately optimized and thus they could not cooperate well with each other. As one can see from the table, most end-to-end methods perform better than all of the listed two-step methods, which strongly indicates the benefits of joint optimization.

Compared with KOALANet and DCLS, which also adopt the end-to-end framework (proposed after DAN-Pre), DAN can still maintain its superiority (although DCLS achieves comparable results with DAN for scale factor  $\times 4$ , it has a much larger model size than DAN, which will be discussed in Sec 4.4). It may be because the alternating optimization algorithm of DAN enables *Restorer* and *Estimator* to utilize the intermediate results of each other. This comparison also demonstrates the advantages of alternating optimization.

We also note that EDSR, which has a very simple network architecture (consists of only residual blocks), can perform better than most end-to-end methods. It indicates that the elaborated complex network architecture may not help improve



**Fig. 7:** Visual comparisons between different methods on real images for scale factor  $\times 4$ . Best viewed in color.

the performance of blind SR methods. Inspired by the performance of EDSR, we also largely simplify the architecture of DAN-Pre. As shown in the table, the current version of DAN can perform much better than the preliminary version. Additionally, the body of *Restorer* in DAN has the same architecture as that of EDSR. While DAN can largely outperform EDSR, which suggests that the *Estimator* in DAN plays an important role in improving the performance of blind SR. EDSR super-resolves LR images regardless of their degradations. While *Estimator* enables DAN to be degradation-aware and perform better over various degradations.

**Qualitative comparisons.** We visualize the  $\times 2$  SR results of *img 823* and *img 872* in Fig 4. As one can see, KerneGAN+ZSSR and IKC fail to remove the blur and can only produce

over-smoothed results. The SR images produced by DCLS are much sharper but contain some unpleasant artifacts, such as twisted lines. While the SR results of DAN are clearer, sharper, and contain fewer artifacts. The same comparisons are also shown in Fig 5, which is the  $\times 4$  SR results of *img 837* and *img 845*.

### 4.3 Comparisons on Real-World LR Images

In this section, we explore the performance of DAN for real-world LR images. As we have described above, we use the degradation model in Real-ESRGAN [46] to synthesize training samples, which are then used to train DAN for real-world LR images. In this case, we mainly compare with three methods, *i.e.* Ji *et al.* [75], BSRGAN [47],

and Real-ESRGAN [46]. These methods have two versions: the PSNR-oriented version and the GAN version. As the two versions have different behaviors, methods of different versions are compared independently. And we also provide the results of DAN and DAN-GAN.

**Comparison results.** We evaluate different methods on three datasets: the DIV2K-Real that we synthesized via the degradation model in [46], 2020Track2 [61], and the RealSRSet collected in [47]. Since there are no ground-truth for 2020Track2 and RealSRSet, we use the non-reference metrics for reference. As shown in Table 2, DAN and DAN-GAN achieve the best PSNR and SSIM results on the synthetic dataset among the PSNR-oriented methods and perceptual-oriented methods respectively. Fig 6 shows the visual comparisons on the synthetic dataset. As one can see, compared with BSR-Net and Real-ESRNet, DAN produce can produce SR images with sharper and clearer textures. On 2020Track2 and RealSRSet, DAN-GAN fail to achieve promising quantitative results. However, as discussed in [47, 76], these metrics may fail to measure the visual quality of SR images. As shown in Fig 7, although Ji *et al.* achieves the best quantitative results, its SR images contain serious artifacts, which are the worst in the reference methods. BSRGAN and Real-ESRGAN perform better, but the edges in their produced images are not clear enough. This may be because these methods are degradation-unaware, and are likely to produce over-smoothed images. As a comparison, DAN-GAN is degradation-specific, *i.e.* performing SR according to the estimated degradation. As a result, DAN-GAN can produce sharper and clearer SR images.

#### 4.4 Comparisons on Complexity and Speed

Compared with other blind SR methods, our end-to-end model also has superiority in model complexity inference speed. To make a quantitative comparison, we evaluate the average speed of different methods on the same platform with an RTX 3090 GPU. We choose the  $\times 4$  models of KernelGAN [4] + ZSSR [39], IKC [42], and RealESRGAN [46] as the comparison methods. The model complexity is measured by the number of parameters and multi-adds. And the speed is

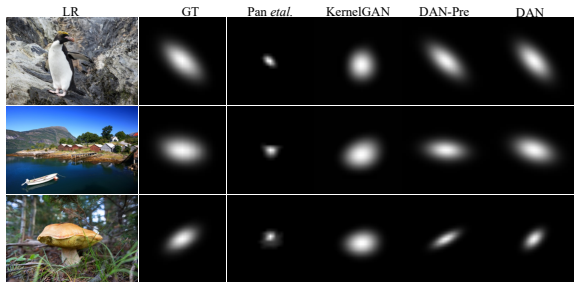
evaluated by the average time of processing 1000 images. The multi-adds and speed are calculated when the size of their output images is  $180 \times 270 \times 3$ . The comparisons are shown in Table 3. The number of multi-adds of KernelGAN+ZSSR is left out because it re-trains a different model for each test image. In that case, multi-adds can not indicate the model complexity.

As one can see, the average speed of DAN-Pre is 0.21 seconds per image, nearly 554 times faster than KernelGAN + ZSSR, and 5 times faster than IKC, which demonstrates the speed superiority of our end-to-end framework. Moreover, in the current version, DAN is further simplified. Compared with DAN-Pre, DAN has 55% fewer parameters and 77% fewer multi-adds. The speed of DAN is also 4 times faster than DAN-Pre. Compared with RealESRGAN, the current version of DAN also has 89% fewer parameters, 74% fewer multi-adds, and 3.4 times faster speed. While as we have discussed in Sec 2, DAN also performs better.

#### 4.5 Study of Estimated Degradations

**Accuracy.** Previous degradation estimation methods mainly focus on estimating the blur kernel. Thus, to make better comparisons with other methods, we also evaluate the accuracy of estimated kernels in the case of SR for blurry LR images. We use two metrics to measure the accuracy: 1) we calculate the mean squared error (MSE) between the ground-truth kernel and the predicted kernel; 2) we degrade the original HR image with our predicted kernel and calculate the PSNR between the original LR image and our generated LR image.

We choose four reference methods, *i.e.* KernelGAN [4], CorrFilter [77], DCLS [6], and DAN-Pre [7]. The quantitative results are shown in Table 4. It should be noted that the original version of DAN-Pre, as presented in our conference paper, can only estimate the PCA feature of the kernel instead of the whole kernel, whereas in this experiment, we have modified DAN-Pre to predict the entire kernel. As one can see, in terms of both Kernel-MSE and LP-PSNR, the estimation accuracy of DAN-Pre is much better than KernelGAN and CorrFilter. This superiority could be attributed to two reasons: 1) the *Estimator* of DAN can utilize the information of intermediate



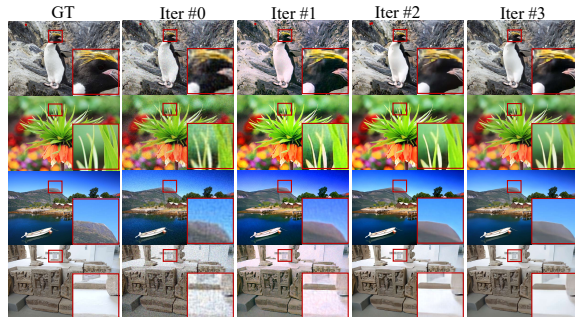
**Fig. 8:** Visualization of estimated blur kernels.

SR results, which makes it easier for *Estimator* to predict accurate degradation; 2) the *Estimator* and *Restorer* are optimized in an end-to-end network and are likely to get better compatible with each other. DCLS, which is proposed after DAN-Pre and also adopts an end-to-end framework, achieves slightly better kernel-estimation performance than DAN-Pre. While the current version DAN surpasses it by a large margin. This may be because the *Estimator* in DAN predicts the parameters ( $\sigma_x$ ,  $\sigma_y$ ,  $\theta$ , etc.) instead of the whole kernel, in which way the possible space the predicted kernel can be largely reduced, and the kernel estimation also becomes easier.

**Visualization.** we also visualize some estimated kernels on DIV2K<sub>RR</sub> (scale factor  $\times 4$ ) to qualitatively measure the performance of *Estimator*. We use the results of KernelGAN [4], Pan *et al.* [78], DAN-Pre [7] as comparisons. As shown in Fig 8, the kernels estimated by KernelGAN are likely to be isotropic and look very different from the ground-truth kernels. Instead, both DAN-Pre and DAN can estimate the kernel much more accurately, even if the ground-truth kernels are highly anisotropic.

## 4.6 Non-blind Setting

To explore the influence of the degradation-estimation deviation, we replace the estimated degradations with the ground-truth ones and observe the performance of the *Restorer*. The comparisons are shown in Table 5. As one can see, when the ground-truth degradations are provided, the results of DAN almost keep unchanged. It indicates that *Restorer* is not sensitive to estimation deviation of the degradation. This is because the *Estimator* and the *Restorer* are jointly optimized and they can be more tolerant of the estimation



**Fig. 9:** Visualization of the intermediate SR images during the alternating process.

deviations of each other. The superiority of DAN also partially comes from the good cooperation between its *Estimator* and *Restorer*.

## 4.7 Study of Iterations

In this section, we experimentally explore the influence of the iterations of the alternating optimization. We gradually change the iterations from 1 to 4 and train DAN with different iterations. The quantitative results on DIV2K-Real are shown in Table 6. As one can see, as the number of iterations increases, the performance of DAN also grows monotonically. However, more iterations also require a longer inference time. We finally set the number of alternating iterations in DAN as 3 for the balance of cost and effectiveness.

We also change the iterations of a trained DAN to explore its iterating process. It should be noted that the tail modules in DAN are only trained to process the features at the last iterations. If we directly change the iterations in a trained DAN, the tail modules need to process features at intermediate iterations, which may be a different domain of features at the last one. Thus, to comprehensively explore the performance of intermediate iterations, we train tail modules for each iteration respectively, while keeping the weights of DAN fixed, in which case, the intermediate performance can be better evaluated. As shown in Fig 9, the SR results exhibit improved visual quality with cleaner images and richer details as the iterations progress. We also show the quantitative results in Table 7.

**Table 4:** Comparisons on the accuracy of estimated kernel. Results are calculated on the DIV2K<sub>RRK</sub> (scale factor  $\times 4$ ) dataset.  $\uparrow$  denotes the higher the better, and  $\downarrow$  denotes the lower the better. The best result is denoted in bold.

Methods	KernelGAN	CorrFilter	DCLS	DAN-Pre	DAN
Kernel-MSE $\downarrow$	0.1518	0.1392	0.0574	0.0817	<b><math>8.104 \times 10^{-7}</math></b>
LR-PSNR $\uparrow$	41.28	41.35	45.27	45.06	<b>47.35</b>

**Table 5:** Average results of DAN on the DIV2K-Real dataset.  $\uparrow$  denotes the higher the better, and  $\downarrow$  denotes the lower the better.

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
w/ GT	25.47	0.6885	0.5008
w/o GT	25.47	0.6886	0.5036

**Table 6:** The quantitative results of DAN trained with different number of alternating iterations. Average results are reported on the DIV2K-Real dataset.  $\uparrow$  denotes the higher the better, and  $\downarrow$  denotes the lower the better.

# Iters	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1	25.20	0.6822	0.5104
2	25.33	0.6840	0.5088
3	25.47	0.6886	0.5036
4	25.53	0.6904	0.5021

**Table 7:** The quantitative results of different alternating iterations in a trained DAN.  $\uparrow$  indicates the higher the better, and  $\downarrow$  indicates the lower the better.

# Iters	0	1	2	3
PSNR $\uparrow$	21.60	24.14	25.10	25.47
SSIM $\uparrow$	0.6115	0.6347	0.6737	0.6886
LPIPS $\downarrow$	0.7138	0.7085	0.5391	0.5036

## 4.8 Ablation Studies about the Network Architecture

In this section, we perform ablation studies to validate the the advantages of current version against the preliminary version [7], *i.e.* DAN-Pre. There are three main modifications, including A) concatenating the two inputs of *Estimator* and *Restorer* at the beginning; B) iterating in the

feature space; and C) setting the initial degradation learnable. We use the preliminary version as the baseline method. Then we apply the three modifications in turn to the baseline method and compare the results. As shown in Table 8, if the modification A is applied, cite the number of parameters can be largely reduced, while the performance becomes better instead. It suggests that the conditional residual block in DAN-Pre may have the only limited representing capacity and consume too many unnecessary parameters. If modification B is further applied, the performance of can be further improved (the PSNR result increases from 24.98 dB to 25.24 dB). It may be attributed to that iterating in the feature space can get the information better transferred among different iterations. And the modification C, *i.e.* setting the initial degradation learnable, improves the PSNR result from 25.24 dB to 25.41 dB. It suggests that a learnable initial state may make it easier for DAN to converge to a better point.

**Table 8:** Ablation studies on the network architectures. Average results are reported on the DIV2K-Real dataset.  $\uparrow$  denotes the higher the better, and  $\downarrow$  denotes the lower the better. The best two results are indicated in bold and underlined respectively.

A	B	C	# Params (M)	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
			4.33	24.33	0.6712	0.5251
$\checkmark$			1.94	24.98	0.6821	0.5121
$\checkmark$	$\checkmark$		1.94	25.24	0.6830	0.5090
$\checkmark$	$\checkmark$	$\checkmark$	1.95	25.47	0.6886	0.5036

## 4.9 Study of the Loss Function for *Estimator*

In this paper, the estimation of degradation parameters is formulated as a regression task and is supervised using L2 loss. While some degradation parameters ( *e.g.*  $k_g$ ,  $k_c$ ,  $n_t$  *etc.*) are



**Table 9:** The comparison between different loss types for *Estimator*.  $\uparrow$  indicates the higher the better, and  $\downarrow$  indicates the lower the better.

Loss Types	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Mixture of CE and L2 loss	25.47	0.6887	0.5021
Only L2 loss	25.47	0.6886	0.5036

binary-valued. The estimation of these parameters actually can be formulated as a classification task and supervised using cross-entropy (CE) loss. In this section, we partition the degradation parameters into two categories: discrete-values (*e.g.*, kernel types, noise types, scaling modes, *etc.*) and continuous-values (*e.g.*, kernel radius, noise level, quality factor, *etc.*). The discrete-value part is supervised using CE loss, while the continuous-value part is still supervised using L2 loss. As shown in Table 9, the performance of DAN is comparable with both types of losses. However, when a mixture of CE loss and L2 loss is used, slightly better SSIM and LPIPS scores are obtained. Nevertheless, implementing L2 loss alone is significantly more straightforward compared to a mixed loss. As a result, we have decided to use only L2 loss.

## 5 Conclusion

In this paper, we have proposed an end-to-end algorithm that can simultaneously estimate the complex degradations and restore the SR images in blind SR. This algorithm is based on alternating optimization and consists of two parts, namely *Restorer* and *Estimator*. We implement the two parts by convolutional modules and unfold the alternating process to form an end-to-end trainable network. In this way, *Estimator* can utilize information from intermediate SR images, which makes it easier to estimate the degradation. More importantly, *Restorer* is trained with the degradations estimated by *Estimator*, instead of the ground-truth ones. Thus *Restorer* could be more tolerant to the estimation error of *Estimator*. Experiments show that the well-compatibility of the two modules can largely improve the accuracy of blind SR, which demonstrates the importance of an end-to-end pipeline. Moreover, compared with those degradation-unaware methods, the proposed

method performs SR according to the degradations of LR images and is likely to produce sharper and clearer SR images.

## References

- [1] Zhang, K., Gool, L.V., Timofte, R.: Deep unfolding network for image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3217–3226 (2020)
- [2] Baker, S., Kanade, T.: Limits on super-resolution and how to break them. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **24**(9), 1167–1183 (2002)
- [3] Michaeli, T., Irani, M.: Nonparametric blind super-resolution. *IEEE International Conference on Computer Vision*, 945–952 (2013)
- [4] Bell-Kligler, S., Shocher, A., Irani, M.: Blind super-resolution kernel estimation using an internal-gan. In: *Advances in Neural Information Processing Systems* (2019)
- [5] Kai Zhang and Wangmeng Zuo and Lei Zhang: Deep plug-and-play super-resolution for arbitrary blur kernels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1671–1681 (2019)
- [6] Luo, Z., Huang, H., Yu, L., Li, Y., Fan, H., Liu, S.: Deep constrained least squares for blind image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
- [7] Luo, Z., Huang, Y., Li, S., Wang, L., Tan, T.: Unfolding the alternating optimization for blind super resolution. *Advances in Neural Information Processing Systems* **33** (2020)
- [8] Zhang, Y., Li, K., Li, K., Wang, L., Zhong, B., Fu, Y.: Image super-resolution using very deep residual channel attention networks. In: Proceedings of the European Conference on Computer Vision, pp. 286–301 (2018)
- [9] Wang, X., Yu, K., Wu, S., Gu, J., Liu, Y.,

- Dong, C., Qiao, Y., Change Loy, C.: Esrgan: Enhanced super-resolution generative adversarial networks. In: Proceedings of the European Conference on Computer Vision Workshops, pp. 0–0 (2018)
- [10] Liang, J., Cao, J., Sun, G., Zhang, K., Van Gool, L., Timofte, R.: Swinir: Image restoration using swin transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1833–1844 (2021)
- [11] Wang, Z., Chen, J., Hoi, S.C.: Deep learning for image super-resolution: A survey. *IEEE transactions on pattern analysis and machine intelligence* (2020)
- [12] Zhang, H., Li, Y., Chen, H., Gong, C., Bai, Z., Shen, C.: Memory-efficient hierarchical neural architecture search for image restoration. *International Journal of Computer Vision* **130**(1), 157–178 (2022)
- [13] Zhou, M., Yan, K., Pan, J., Ren, W., Xie, Q., Cao, X.: Memory-augmented deep unfolding network for guided image super-resolution. *International Journal of Computer Vision* (2022)
- [14] Yang, J., Wright, J., Huang, T.S., Ma, Y.: Image super-resolution via sparse representation. *IEEE transactions on image processing* **19**(11), 2861–2873 (2010)
- [15] Kim, K.I., Kwon, Y.: Single-image super-resolution using sparse regression and natural image prior. *IEEE transactions on pattern analysis and machine intelligence* **32**(6), 1127–1133 (2010)
- [16] Timofte, R., De Smet, V., Van Gool, L.: Anchored neighborhood regression for fast example-based super-resolution. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1920–1927 (2013)
- [17] Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **38**(2), 295–307 (2015)
- [18] Haris, M., Shakhnarovich, G., Ukita, N.: Deep back-projection networks for super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1664–1673 (2018)
- [19] Hu, X., Mu, H., Zhang, X., Wang, Z., Tan, T., Sun, J.: Meta-sr: A magnification-arbitrary network for super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1575–1584 (2019)
- [20] Hui, Z., Gao, X., Yang, Y., Wang, X.: Lightweight image super-resolution with information multi-distillation network. In: Proceedings of the 27th ACM International Conference on Multimedia, pp. 2024–2032 (2019)
- [21] Ahn, N., Kang, B., Sohn, K.-A.: Fast, accurate, and lightweight super-resolution with cascading residual network. In: Proceedings of the European Conference on Computer Vision, pp. 252–268 (2018)
- [22] Hui, Z., Wang, X., Gao, X.: Fast and accurate single image super-resolution via information distillation network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 723–731 (2018)
- [23] Dong, C., Loy, C.C., Tang, X.: Accelerating the super-resolution convolutional neural network. In: Proceedings of the European Conference on Computer Vision, pp. 391–407 (2016). Springer
- [24] Kim, J., Kwon Lee, J., Mu Lee, K.: Accurate image super-resolution using very deep convolutional networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1646–1654 (2016)
- [25] Shi, W., Caballero, J., Huszár, F., Totz, J., Aitken, A.P., Bishop, R., Rueckert, D., Wang, Z.: Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp.

- 1874–1883 (2016)
- [26] Dai, T., Cai, J., Zhang, Y., Xia, S.-T., Zhang, L.: Second-order attention network for single image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11065–11074 (2019)
- [27] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
- [28] Chen, H., Wang, Y., Guo, T., Xu, C., Deng, Y., Liu, Z., Ma, S., Xu, C., Xu, C., Gao, W.: Pre-trained image processing transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12299–12310 (2021)
- [29] Köhler, T., Bätz, M., Naderi, F., Kaup, A., Maier, A., Riess, C.: Toward bridging the simulated-to-real gap: Benchmarking super-resolution on real data. *IEEE transactions on pattern analysis and machine intelligence* **42**(11), 2944–2959 (2019)
- [30] Chen, S., Han, Z., Dai, E., Jia, X., Liu, Z., Xing, L., Zou, X., Xu, C., Liu, J., Tian, Q.: Unsupervised image super-resolution with an indirect supervised path. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 468–469 (2020)
- [31] Levin, A., Weiss, Y., Durand, F., Freeman, W.T.: Understanding and evaluating blind deconvolution algorithms. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1964–1971 (2009). IEEE
- [32] Levin, Anat and Weiss, Yair and Durand, Fredo and Freeman, William T: Efficient marginal likelihood optimization in blind deconvolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2657–2664 (2011). IEEE
- [33] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A.C., Bengio, Y.: Generative adversarial nets. In: *Advances in Neural Information Processing Systems* (2014)
- [34] Pan, J., Sun, D., Pfister, H., Yang, M.-H.: Deblurring images via dark channel prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **40**, 2315–2328 (2018)
- [35] Pan, Jinshan and Sun, Deqing and Pfister, Hanspeter and Yang, Ming-Hsuan: Blind image deblurring using dark channel prior. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1628–1636 (2016)
- [36] Yan, Y., Ren, W., Guo, Y., Wang, R., Cao, X.: Image deblurring via extreme channels prior. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4003–4011 (2017)
- [37] Cai, J., Zuo, W., Zhang, L.: Dark and bright channel prior embedded network for dynamic scene deblurring. *IEEE Transactions on Image Processing* **29**, 6885–6897 (2020)
- [38] Glasner, D., Bagon, S., Irani, M.: Super-resolution from a single image. 2009 IEEE 12th International Conference on Computer Vision, 349–356 (2009)
- [39] Shocher, A., Cohen, N., Irani, M.: ”zero-shot” super-resolution using deep internal learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2018)
- [40] Soh, J.W., Cho, S., Cho, N.I.: Meta-transfer learning for zero-shot super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3516–3525 (2020)
- [41] Zhang, K., Zuo, W., Zhang, L.: Learning a single convolutional super-resolution network for multiple degradations. In: Proceedings of the IEEE/CVF Conference on Computer

- Vision and Pattern Recognition, pp. 3262–3271 (2018)
- [42] Gu, J., Lu, H., Zuo, W., Dong, C.: Blind super-resolution with iterative kernel correction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1604–1613 (2019)
- [43] Kim, S.Y., Sim, H., Kim, M.: Koalanet: Blind super-resolution using kernel-oriented adaptive local adjustment. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10611–10620 (2021)
- [44] Liang, J., Zhang, K., Gu, S., Van Gool, L., Timofte, R.: Flow-based kernel prior with application to blind super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10601–10610 (2021)
- [45] Wei, Y., Gu, S., Li, Y., Timofte, R., Jin, L., Song, H.: Unsupervised real-world image super resolution via domain-distance aware training. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13385–13394 (2021)
- [46] Wang, X., Xie, L., Dong, C., Shan, Y.: Real-esrgan: Training real-world blind super-resolution with pure synthetic data. arXiv preprint arXiv:2107.10833 (2021)
- [47] Zhang, K., Liang, J., Van Gool, L., Timofte, R.: Designing a practical degradation model for deep blind image super-resolution. arXiv preprint arXiv:2103.14006 (2021)
- [48] Xie, J., Zhan, X., Liu, Z., Ong, Y.-S., Loy, C.C.: Delving into inter-image invariance for unsupervised visual representations. *International Journal of Computer Vision*, 1–20 (2022)
- [49] Wei, P., Xie, Z., Lu, H., Zhan, Z., Ye, Q., Zuo, W., Lin, L.: Component divide-and-conquer for real-world image super-resolution. In: European Conference on Computer Vision, pp. 101–117 (2020). Springer
- [50] Cai, J., Zeng, H., Yong, H., Cao, Z., Zhang, L.: Toward real-world single image super-resolution: A new benchmark and a new model. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3086–3095 (2019)
- [51] Zhang, K., Zuo, W., Gu, S., Zhang, L.: Learning deep cnn denoiser prior for image restoration. In: Proceedings of the IEEE/CVF Conference on Computer Cision and Pattern Recognition, pp. 3929–3938 (2017)
- [52] Ren, D., Zhang, K., Wang, Q., Hu, Q., Zuo, W.: Neural blind deconvolution using deep priors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3341–3350 (2020)
- [53] Kaufman, A., Fattal, R.: Deblurring using analysis-synthesis networks pair. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5811–5820 (2020)
- [54] Pascal, F., Bombrun, L., Tourneret, J.-Y., Berthoumieu, Y.: Parameter estimation for multivariate generalized gaussian distributions. *IEEE Transactions on Signal Processing* **61**(23), 5960–5971 (2013)
- [55] Brooks, T., Mildenhall, B., Xue, T., Chen, J., Sharlet, D., Barron, J.T.: Unprocessing images for learned raw denoising. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11036–11045 (2019)
- [56] Shin, R., Song, D.: Jpeg-resistant adversarial images. In: NIPS 2017 Workshop on Machine Learning and Computer Security, vol. 1 (2017)
- [57] Lim, B., Son, S., Kim, H., Nah, S., Mu Lee, K.: Enhanced deep residual networks for single image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 136–144 (2017)
- [58] Jo, Y., Oh, S.W., Vajda, P., Kim, S.J.: Tackling the ill-posedness of super-resolution

- through adaptive target generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16236–16245 (2021)
- [59] Agustsson, E., Timofte, R.: Ntire 2017 challenge on single image super-resolution: Dataset and study, 1122–1131 (2017)
- [60] et al., R.T.: Ntire 2017 challenge on single image super-resolution: Methods and results, 1110–1121 (2017)
- [61] Lugmayr, A., Danelljan, M., Timofte, R.: 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)
- [62] Ignatov, A., Kobyshev, N., Timofte, R., Vanhoey, K., Van Gool, L.: Dslr-quality photos on mobile devices with deep convolutional networks. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 3277–3285 (2017)
- [63] Martin, D., Fowlkes, C., Tal, D., Malik, J.: A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In: Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001, vol. 2, pp. 416–423 (2001). IEEE
- [64] Matsui, Y., Ito, K., Aramaki, Y., Fujimoto, A., Ogawa, T., Yamasaki, T., Aizawa, K.: Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications* **76**, 21811–21838 (2016)
- [65] Zhang, K., Zuo, W., Zhang, L.: Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Transactions on Image Processing* **27**(9), 4608–4622 (2018)
- [66] Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
- [67] Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2018)
- [68] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems* **25**, 1097–1105 (2012)
- [69] Mittal, A., Soundararajan, R., Bovik, A.C.: Making a “completely blind” image quality analyzer. *IEEE Signal Processing Letters* **20**(3), 209–212 (2012)
- [70] Ma, C., Yang, C.-Y., Yang, X., Yang, M.-H.: Learning a no-reference quality metric for single-image super-resolution. *Computer Vision and Image Understanding* **158**, 1–16 (2017)
- [71] Blau, Y., Mechrez, R., Timofte, R., Michaeli, T., Zelnik-Manor, L.: The 2018 pirm challenge on perceptual image super-resolution. In: Proceedings of the European Conference on Computer Vision (ECCV) Workshops, pp. 0–0 (2018)
- [72] Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., *et al.*: Photo-realistic single image super-resolution using a generative adversarial network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4681–4690 (2017)
- [73] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
- [74] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
- [75] Ji, X., Cao, Y., Tai, Y., Wang, C., Li, J., Huang, F.: Real-world super-resolution via

kernel estimation and noise injection. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 466–467 (2020)

- [76] Luo, Z., Huang, Y., , Li, S., Wang, L., Tan, T.: Learning the degradation distribution for blind image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2022)
- [77] Hussein, S.A., Tirer, T., Giryes, R.: Correction filter for single image super-resolution: Robustifying off-the-shelf deep super-resolvers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1428–1437 (2020)
- [78] Pan, J., Hu, Z., Su, Z., Yang, M.-H.: Deblurring text images via l0-regularized intensity and gradient prior. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2901–2908 (2014)