

Blind Image Super-Resolution: A Survey and Beyond

Anran Liu, Yihao Liu, Jinjin Gu, Yu Qiao, and Chao Dong

Abstract—Blind image super-resolution (SR), aiming to super-resolve low-resolution images with unknown degradation, has attracted increasing attention due to its significance in promoting real-world applications. Many novel and effective solutions have been proposed recently, especially with the powerful deep learning techniques. Despite years of efforts, it still remains as a challenging research problem. This paper serves as a systematic review on recent progress in blind image SR, and proposes a taxonomy to categorize existing methods into three different classes according to their ways of degradation modelling and the data used for solving the SR model. This taxonomy helps summarize and distinguish among existing methods. We hope to provide insights into current research states, as well as to reveal novel research directions worth exploring. In addition, we make a summary on commonly used datasets and previous competitions related to blind image SR. Last but not least, a comparison among different methods is provided with detailed analysis on their merits and demerits using both synthetic and real testing images.

Index Terms—Image Super-Resolution, Deep Learning, Degradation Modelling.



1 INTRODUCTION

SINGLE-IMAGE super-resolution (SISR) has long been a fundamental problem in low-level vision, aiming to recover a high-resolution (HR) image from an observed low-resolution (LR) input. Years of efforts from the research community have brought about remarkable progress in this field, especially with the booming deep learning techniques [1], [2], [3], [4], [5]. However, most existing methods assume a pre-defined degradation process (e.g., bicubic downsampling) from an HR image to an LR one, which can hardly hold true for real-world images with complex degradation types. Towards filling this gap, growing attention has been paid in recent years to approaches for unknown degradations, namely *blind SR*. Despite many exciting improvements, these proposed methods tend to fail in many real-world scenarios, as their performance is usually limited to certain kinds of inputs and will drop dramatically in other cases. The main reason is that they still make some assumptions on the degradation types related to the input LR. Readers can see Fig.1(a) for an illustration, which shows four different LR inputs with assumed degradation types of some state-of-the-art methods but targeting at the same HR. Therefore, when given an arbitrary input deviating from their assumed data distributions, these methods inevitably produce much less pleasing results. Fig.1(b) demonstrates different SR results for a real-world image cropped from the famous film *Forrest Gump*, which are generated by four state-of-the-art methods. We may find none of these methods have lived up to our expectation for a good viewing experience, since this real-world image does not strictly follow

their assumptions on inputs. In fact, it is not rare that we feel confused about which method to choose for a certain image at hand, or whether we can really get a high-quality result using existing methods.

In this paper, we try to relieve this confusion through a systematic survey on recent progress in blind SR with our own insight. What's more, it is highly necessary that we look back and reflect on the proposed methods to have a clear understanding of the current research state and remaining gaps. As stated above, we often have difficulty in selecting a proper method when facing so many ones: KernelGAN [6] for a single image looks cool, but how about IKC [7] with iterative scheme or CinCGAN [8] with unpaired training data? Also, even if every single blind SR method is claimed to work well for real images, we may still struggle to obtain a satisfactory output for our own image, just like the case in Fig.1. At this stage of development, it is time to ask: To what extent have we solved the problem? What is holding us back and where should we go for future endeavour?

Hence, this paper aims to serve much more than a list of recent progress. Specifically, we propose a taxonomy to effectively categorize existing approaches, which clearly distinguishes among different methods and naturally reveals some research gaps. Based on this taxonomy, our goal is to let each method have its own position within a broad picture composed of existing work. This picture can provide a guideline on reasonable and fair comparison between different kinds of methods in future work. In addition, we will make a summary on the application scopes along with limitations of each kind of approaches, helping readers to efficiently select appropriate methods for various scenarios. Note that this paper focuses on SISR for general natural images, not including domain-specific topics like face SR or depth map SR.

Our contributions are mainly three-fold: 1) We present a systematic survey on recent progress in blind image super-resolution, including the improvements and limitations of

- A. Liu is with The University of Hong Kong, Hong Kong SAR, China. E-mail: liuar616@connect.hku.hk
- J. Gu is with the School of Electrical and Information Engineering, The University of Sydney. E-mail: jinjin.gu@sydney.edu.au
- Y. Liu, Y. Qiao and C. Dong are with Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China. E-mail: liuyihao14@mails.ucas.ac.cn, {yu.qiao, chao.dong}@siat.ac.cn

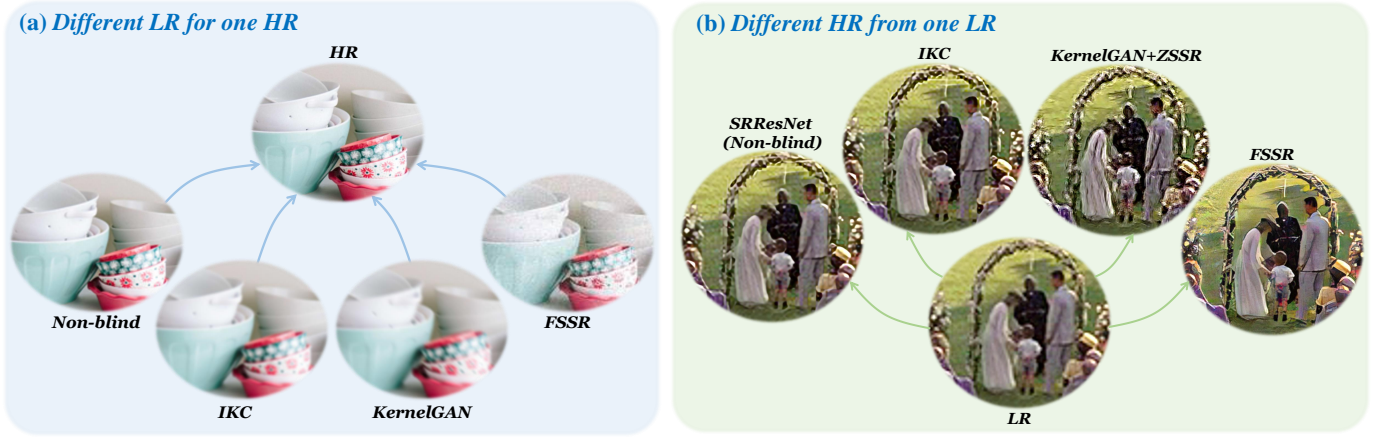


Fig. 1: Left: an HR image and its different LR versions with assumed degradation types of four SISR methods. Right: Different HR images generated by each method from an LR input crop of *Forrest Gump*.

different approaches. 2) We propose a taxonomy to effectively categorize existing methods and reveal some research gaps. 3) We provide our deep insight into current research state and promising future directions.

In the following sections, we first introduce the mathematical formulations of some commonly used SR models in Sec.2, and discuss about the challenges from real-world images that we face when addressing blind SR in Sec.3. Then we put forward our proposed taxonomy in Sec.4. A quick review on non-blind SISR is provided in Sec.5, since the research state in non-blind setting has set up the foundation for blind SR. Then we elaborate on methods of each category in Sec.6 and Sec.7, followed by a summary on commonly used datasets and previous competitions in the field of blind SR in Sec.8. Quantitative and qualitative comparison among some representative methods is included in Sec.9. Finally, we draw a conclusion on our insight through this survey and our perspective on future directions in Sec.10.

2 PROBLEM FORMULATION

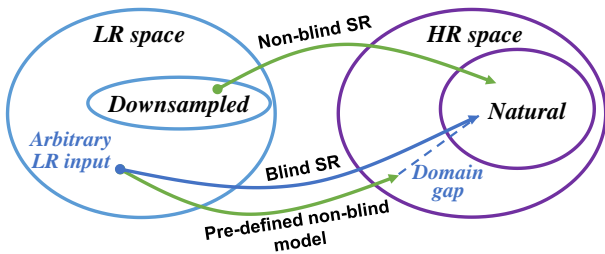


Fig. 2: Domain interpretation of differences between non-blind and blind SR. There exists a large domain gap between the SR result and desired high-quality HR, which is caused by applying a pre-trained non-blind model to LR input with degradation deviating from the assumed one (e.g., downsampling).

In this section, we introduce some mathematical formulations of the SISR problem. Specifically, SISR refers to the task of reconstructing an HR image from a given LR input, especially the high-frequency contents in HR. The underlying degradation process from HR to LR can be generally expressed with the following equation:

$$y = f(x; s), \tag{1}$$

where x, y denote HR image and LR image respectively, f is the degradation function with a scale factor s . Therefore, SR problem is equivalent to modelling and solving the inverse function f^{-1} . In the background of *non-blind SR*, f is usually assumed to be bicubic downsampling:

$$y = x \downarrow_s^{bic}, \tag{2}$$

or the combination of downsampling and a fixed Gaussian blur with kernel k_g :

$$y = (x \otimes k_g) \downarrow_s, \tag{3}$$

where \otimes denotes convolutional operation. Under either assumption, the corresponding SR model is only able to handle LR inputs with this specific kind of degradation. For other LR images with different degradation types, the mismatch between SR model and intrinsic degradation of inputs may cause severe artifacts in SR results [7], [12]. Fig.2 gives an illustration on this mismatch from the perspective of image domain adaptation: if an SR model corresponding to a pre-defined degradation is applied to an arbitrary LR input, there will be a large domain gap between the SR output and desired image samples from the target *Natural HR* domain, thus leading to a poor-quality result.

Hence, the topic of blind SR for unknown degradation is proposed in an attempt to bridge this gap. Up till now, there have been two different ways of modelling the degradation process for blind SR: explicit modelling based on an extension of Equation (3), and implicit modelling through inherent distribution within external dataset.

To be specific, explicit modelling usually employs a so-called classical degradation model, which is a more general form of Equation (3):

$$y = (x \otimes k) \downarrow_s + n, \tag{4}$$

where the SR blur kernel k and additive noise n are two main factors involved in the degradation process, and parameters related to these two factors will be unknown for an arbitrary LR input. Fig.3(a) shows several image examples with different k and n , which are much more degraded than their bicubic-downsampled counterpart. Some approaches utilize external dataset to learn an SR model well adapted to a large set of various k or n , such as IKC [7] and SRMD

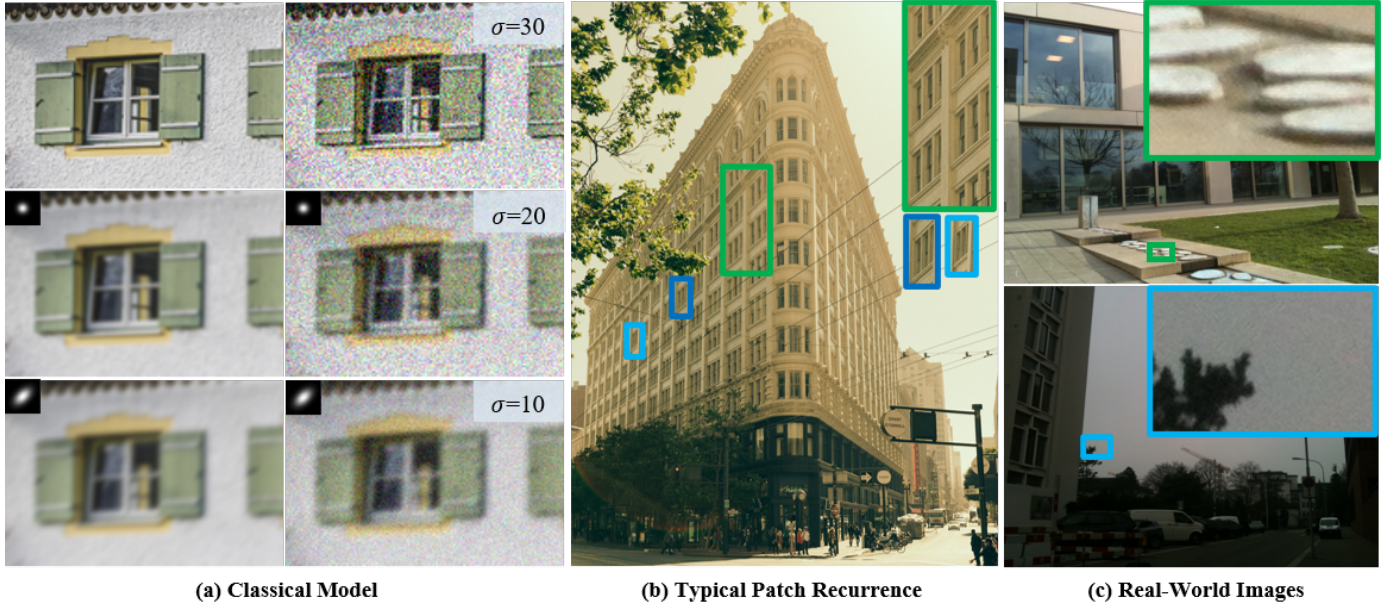


Fig. 3: Examples of degraded LR image. (a) Clean LR (top left, generated with bicubic downsampling) and its blurry or noisy versions with different k and n . The 2^{nd} row is with isotropic Gaussian kernels while the 3^{rd} row with anisotropic Gaussian, and σ is variance of additive Gaussian noise. The example image is from DIV8K [9]; (b) Image with typical patch recurrence, both within and across different scales of the same image. The image is from Urban100 [10]; (c) Real-world smartphone images with complex unknown degradations. Images are from DPED dataset [11].

[13]. Besides blurring and noise, more complex and realistic degradation types can also be involved into the formulation, like JPEG compression with a quality factor q [14]:

$$y = ((x \otimes k) \downarrow_s + n)_{JPEG_q}, \quad (5)$$

Another group of methods leverage the internal statistics within a single image derived from the classical degradation model, thus requiring no external dataset for training, like ZSSR [12] and DGDML-SR [15]. In fact, internal statistics just reflects the patch recurrence property of an image, and readers can refer to Fig.3(b) for an illustration.

Nevertheless, real-world degradations are usually too complex to be modelled with an explicit combination of multiple degradation types, as shown in Fig.3(c). Therefore, implicit modelling attempts to circumvent the explicit modelling function. Instead, it defines the degradation process f implicitly through data distribution, and all the existing approaches with implicit modelling require an external dataset for training. Typically, these methods utilize data distribution learning with Generative Adversarial Network (GAN) [16] to grasp the implicit degradation model possessed within training dataset, like CinCGAN [8].

Although so many models have been put forward in blind SR, there is still a long way to go since we have only tackled a small set of real-world images. Existing methods often claim to focus on real-world settings, but they actually assume a certain scene, like images taken by some digital cameras [17], [18]. In fact, real-world images are greatly different in their underlying degradation types, and an SR model designed for a specific type can easily fail for another. In the next section, we will give a brief discussion on different kinds of real-world images, which have posed severe challenges to the field of blind SR.

3 CHALLENGES FROM REAL-WORLD IMAGES

With the development of modern imaging devices, we are now embracing the world with a surge of visual data. Such a variety of image sources also pose more challenges, especially in terms of degradations types. Generally speaking, there are three main factors causing different degradations:

- (1) Different imaging devices. This era of technology has given birth to a dazzling array of digital cameras, not to mention smartphones with advanced camera systems. However, these devices differ greatly in the characteristics of the photos taken [11]. For example, a DSLR (digital single-lens reflex) camera is able to capture high-quality images with a sense of stereoscopy by adjusting its focal length, while a smartphone camera is nowhere near DSLR-quality, tending to produce a “flattened” and noisy scene due to its physical limitations in sensor size and lenses. Another type of low-quality imaging is surveillance video, which often suffers greatly from loss of focus. Readers can see Fig.4 for some image examples. Images captured with different devices can thus have distinct degradations from one another.
- (2) Image processing algorithms. This problem is mostly related to digital and smartphone cameras, since it is an image signal processor on the chip that actually processes digital signals into images. The processing pipeline usually involves multiple steps, such as pixel correction, white balance correction, denoising and sharpening. During this process, complex unknown degradations can be introduced [21], which are unpredictable and varying among different devices. A typical pipeline is shown in Fig.5.
- (3) Degradations rising from storage. To reduce the resource consumption for transmitting and storing



Fig. 4: Real-world images (or video frames) captured with different devices. Images are from DPED [11], VIRAT [19], UCF-Crime [20] datasets.

data, images and videos are always compressed. Accompanying compressed images are compression artifacts, which will lead to degradations like blurring and blocky effects. In addition, time itself can gradually deteriorate images, especially for old photos and movies recorded on films. Such degradations are mainly caused by poor imaging equipments or erosion in the air, including film grain, sepia effect and color fading [22]. Some example images are presented in Fig.6. This kind of degradations can hardly be expressed with explicit functions or covered by a few external datasets, thus demanding more efforts in designing restoration algorithms.

The real-world images discussed above all bear their own degradations and challenges. Nonetheless, previous work usually focus on a single type of real images, like those taken by smartphones, which greatly limits their performance in diverse scenes. We expect to see more explorations on different types of real-world images in the future. Specifically, effective solutions for each distinct type, if not a general solution to all, should be the ultimate goal of our research community.

4 TAXONOMY

In this section, we will elaborate on our proposed taxonomy to serve as the guideline for our review and analysis. According to Sec.2, there have been two ways of modelling the degradation process involved in blind SR: explicit modelling based on the classical degradation model or its variants, and implicit modelling using data distribution among external dataset. The basic idea of explicit modelling is to learn an

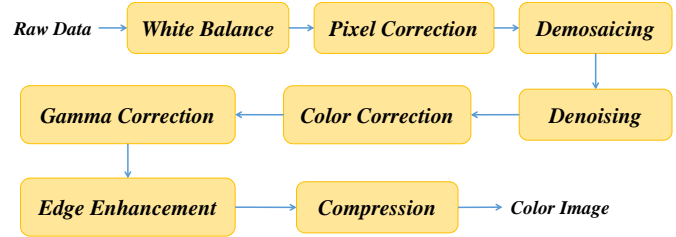


Fig. 5: Image signal processing pipeline.

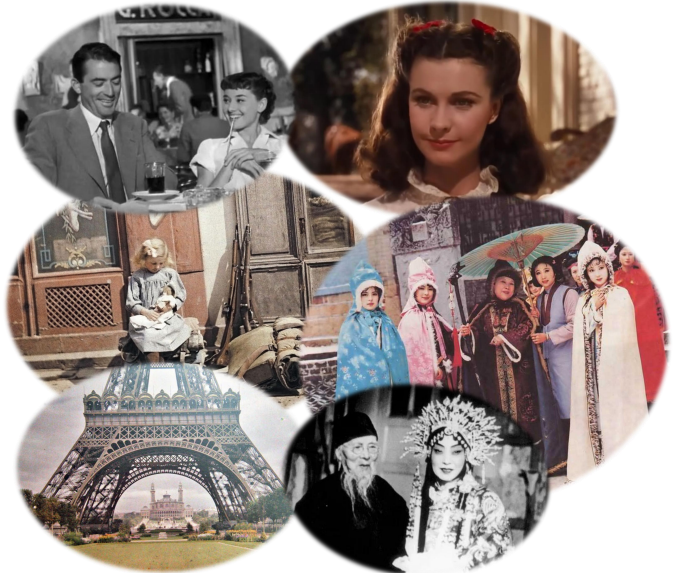


Fig. 6: Old photos with degradation rising from storage.

SR model with external training data covering a large set of degradations, which are usually parameterized with k and n in Equation (4). Representative approaches include SRMD [13], IKC [7] and KMSR [23]. Another group of methods propose to exploit internal statistics of patch recurrence, like KernelGAN [6] and ZSSR [12], which can directly work on a single input image. This kind of modelling is primarily based on the classical degradation model. On the other hand, methods with implicit modelling do not depend on any explicit parameterization, and they usually learn the underlying SR model implicitly through data distribution within external dataset. Among such methods are CinGAN [8] and FSSR [24].

Therefore, we propose a taxonomy to effectively classify existing approaches according to their ways of degradation modelling and the used data for solving the SR model: explicit modelling or implicit modelling, external dataset or a single input image, as shown in Fig.7. Our reasons for adopting this categorization are three-fold: first, distinguishing between explicit and implicit modelling helps us to understand the assumption of a certain method, i.e., what kind of degradations this method aims to deal with; second, whether using external dataset or a single input image indicates different strategies of image-specific adaptation with explicit modelling; finally, after categorizing existing approaches into these classes, one remaining research gap naturally reveals itself - *implicit modelling with a single image*. We argue that this direction is promising in terms of ad-

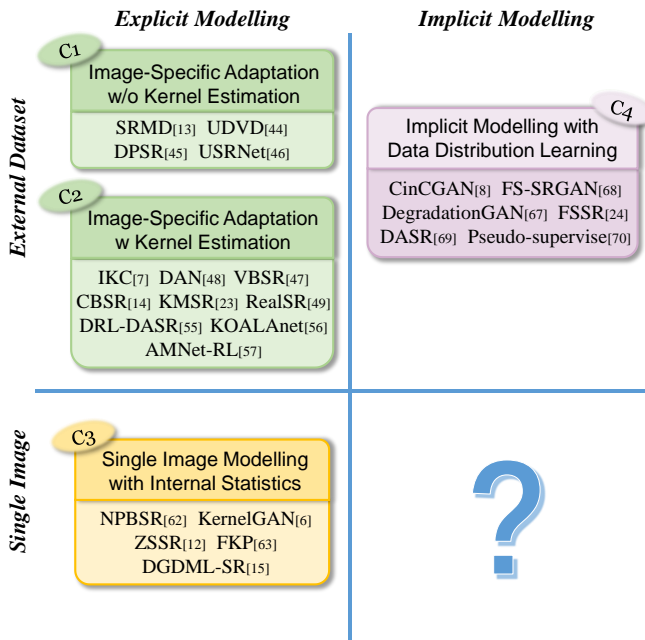


Fig. 7: Our proposed taxonomy and the corresponding representative methods. Our taxonomy distinguishes the ways of degradation modelling and data used for solving SR models, which also naturally reveals the remaining research gap.

addressing general real-world images with diverse contents, and we will also try to propose feasible suggestions for new solutions in this direction.

In the next sections, we first give a quick overview on non-blind SISR, which sets the basis for blind SR methods. Then methods with explicit modelling are introduced in Sec.6, and those using implicit modelling are discussed in Sec.7. For each type of methods, we will unfold the review along its *course of development*, and make an analysis on their limitations to inspire future work.

5 OVERVIEW OF NON-BLIND SINGLE-IMAGE SUPER-RESOLUTION

As explained in Sec.2, non-blind SR assumes a fixed known degradation process for solving HR outputs. Before the development of deep learning techniques, many traditional techniques are example-based. [25], [26], [27], [28] learn the mapping function from LR to HR with external HR-LR exemplar pairs, where the mapping learning is usually based on a compact dictionary or manifold space. Some others [29], [30] utilize the property of internal self-similarity within a single image without employing external dataset. In 2014, the pioneering work of SRCNN [31] opened a new era of deploying convolutional neural network (CNN) to tackle this task, and it also set up the basic framework for later works, as shown in Fig.8.

The commonly adopted CNN framework for SISR task includes three main modules: shallow feature extraction to convert an input LR image into feature maps, deep feature extraction or mapping based on extracted shallow features, and finally SR output reconstruction. Residual learning has also been widely adopted to ease the training process, either in image-level [33] or feature-level [34]. Recent years have

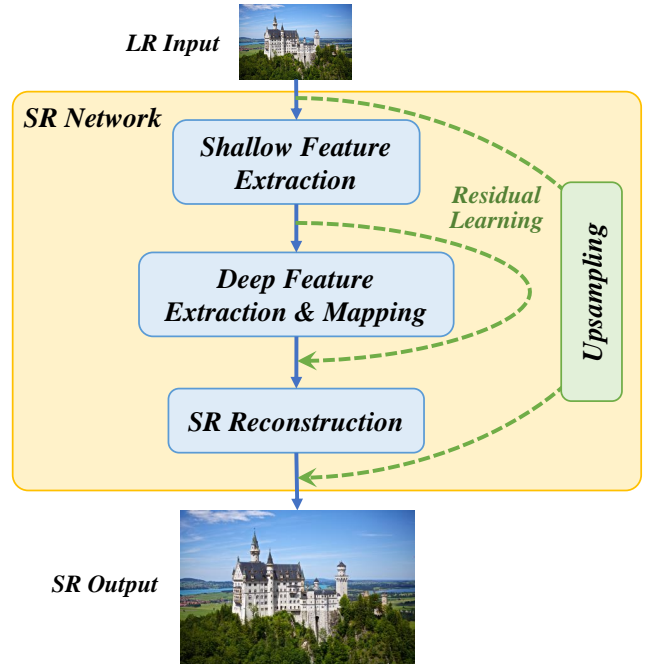


Fig. 8: Common CNN framework of non-blind SISR.

witnessed many improvements on deep feature extraction and SR reconstruction modules, such as introducing residual blocks [34], [35], [36], recursive or recurrent structure [37], [38], attention mechanism [39], [40], sub-pixel convolution [41], etc. In addition, multiple loss functions are also proposed for better perceptual quality of SR results [32], [42], [43]. These techniques bring about remarkable progress in terms of both reconstruction accuracy and efficiency, and non-blind SISR with bicubic-downsampling assumption actually reaches maturity.

However, these non-blind models usually struggle to generalize to input images with more complex degradations deviating from their assumed one. Some failure cases of a non-blind SR network are shown in Fig.9, where the network performs well on bicubically downsampled clean input in accordance with the assumed degradation model, but cannot handle blurry or noisy input images. Hence, it is demanding to propose methods for blind SR setting, which are the main focus of this survey and will be explored in detail in the following two sections.

6 EXPLICIT DEGRADATION MODELLING

This section covers recently proposed blind SR methods with explicit modelling of degradation process, usually based on the classical degradation model shown by Equation (4). What's more, these approaches can be further classified into two sub-classes according to whether they employ external dataset or rely on a single input image to solve the SR problem.

6.1 Classical Degradation Model with External Dataset

This kind of approaches utilize external dataset to train an SR model well adapted to variant SR blur kernels k and noises n , especially the former. Typically, the SR model is



Fig. 9: Failure cases of non-blind SR network, e.g. SRResNet [32]. Compared with results generated by bicubic interpolation, SRResNet recovers little sharp texture for the blurry input, but also unfavourably keeps the noises for the noisy input.

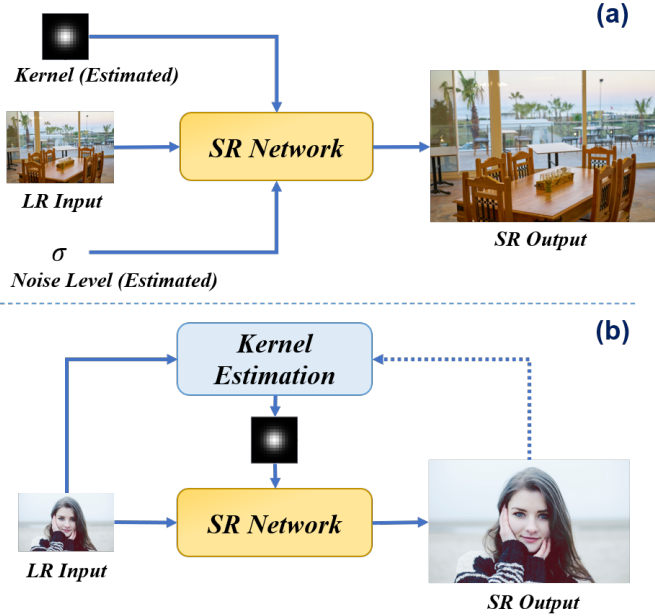


Fig. 10: Two types of overall frameworks for methods with explicit modelling and external dataset. (a) Image-specific adaptation *without* kernel estimation; (b) Image-specific adaptation *with* kernel estimation. The connection with dotted line indicates that it is optional.

parameterized with a convolutional neural network (CNN), and an estimation on k or n for a specific LR image is used as conditional input to the SR model for feature adaptation purpose. After the training process, the model will be able to produce satisfactory results for LR inputs with degradation types covered in the training dataset. According to whether a certain approach includes degradation estimation in its proposed framework, we further divide these approaches into two types: image-specific adaptation *without* kernel estimation, and image-specific adaptation *with* kernel estimation. To be more specific, the first type receives estimated degradation information as additional inputs and is focused on how to utilize the estimation input for image-specific adaptation, while the second one pays special attention to kernel estimation along with the SR process. An illustration of their overall frameworks is presented in Fig.10.

6.1.1 Image-Specific Adaptation without Kernel Estimation

Super-resolution for multiple degradations (SRMD) [13] proposes to directly concatenate an LR input image with its degradation map as a unified input to the SR model,

thus allowing feature adaptation according to the specific degradation and covering multiple degradation types in a single model. In order to generate degradation map with the same dimension as the LR image, it introduces a strategy called dimensionality stretching. Specifically, an SR blur kernel with size $r \times r$ is flattened to an r^2 -length vector and reduced to t -dim with principal component analysis (PCA) to get the kernel coding. After concatenating with the estimated noise level σ related to n , the $(t+1)$ -dim vector is stretched both vertically and horizontally to get the final $H \times W \times (t+1)$ -dim degradation map, where H and W are height and width of the LR image. This strategy can be easily extended to non-uniform maps for spatially variant degradations. The SR reconstruction network of SRMD is similar to those commonly adopted in non-blind SR. The whole pipeline is presented in Fig.11(a).

Following SRMD, UDVD [44] also uses the degradation map as an additional input for SR reconstruction, yet it makes one step forward by employing per-pixel dynamic convolution to more effectively deal with variational degradations across images. Specifically, a refinement network composed of several dynamic blocks with dynamic convolution is cascaded to feature extraction module, and each of these blocks refines SR output in an iterative way based on the result of its previous block. In addition, an improvement on the kernel coding operation is proposed by variant blind SR [47] to replace the PCA technique with a shallow neural network, which can potentially learn a kernel mapping more fitted to the specific SR model.

Although SRMD extends the generalization capacity of an SR model to variant SR kernels and noise levels, it still has very limited scope since it is usually non-trivial to effectively encode an arbitrary kernel and handle it with a single model, especially for those with irregular patterns like motion blur. Hence, another group of methods have been proposed based on an MAP framework, which requires no kernel coding for degradation map generation. Specifically, deep plug-and-play super-resolution (DPSR) [45] incorporates SR network into an MAP-based iterative optimization scheme. It primarily solves the HR image by minimizing the following objective function, which consists of a data term D and a prior term P regularized by a parameter λ :

$$E(x) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{x} \downarrow_s \otimes \mathbf{k}\|^2 + \lambda\Phi(x) = D + \lambda P, \quad (6)$$

whose corresponding degradation model is a modified version of Equation (4), decoupling the downsampling process

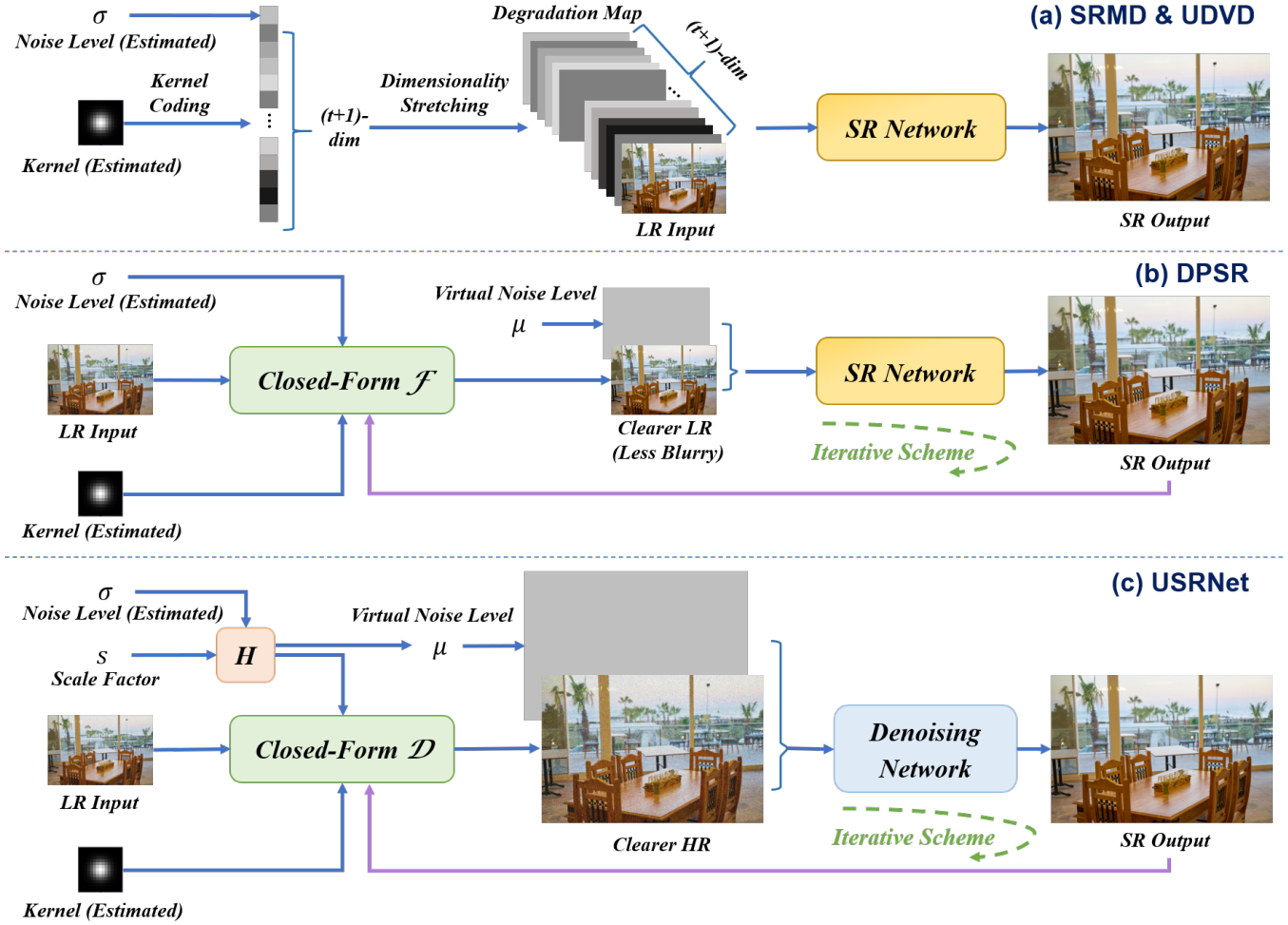


Fig. 11: Pipeline of methods in *Image-Specific Adaptation without Kernel Estimation*. (a) SRMD [13] and UDVD [44]; (b) DPSR [45]; (c) USRNet [46].

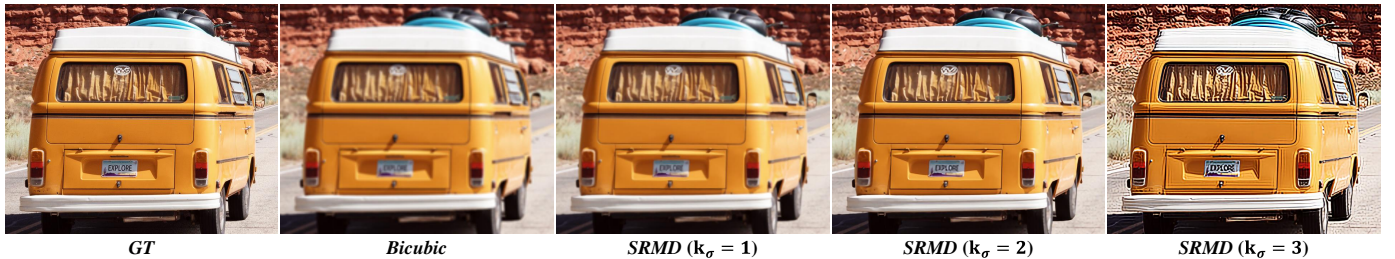


Fig. 12: Illustration on the limitations of *Image-Specific Adaptation without Kernel Estimation*. Input image is blurred and downsampled with Gaussian kernel $k_{\sigma}=2$. An incorrect kernel estimation input either leads to blurry output or unnatural ringing artifacts with over-enhanced textures.

from blurring operation:

$$y = (x \downarrow_s \otimes k) + n. \quad (7)$$

The objective function shown by Equation (6) can be split into two sub-problems using the half-quadratic splitting (HQS) algorithm: one addresses deblurring task and is related to data term D with parameter k , while the other aims to super-resolve a bicubic downsampled image with some virtual noise level μ and is related to prior term P . Fortunately, the first sub-problem can be solved in a closed form with Fast Fourier Transform without any kernel coding, thus allowing the model to cope with more complex

kernels. Moreover, thanks to the decoupling of blurring and downsampling operations, the second sub-problem can be modelled by a non-blind SR network capable of dealing with additive noise, and this network can be directly adapted from SRMD framework with a single noise map as additional input. Unfolding super-resolution network (USRNet) [46] also adopts the MAP framework but is based on the original degradation model in Equation (4), and the corresponding two sub-problems become super-resolving an LR image blurred by kernel k and denoising an HR image with a virtual noise level μ . It enhances the solution framework by unfolding the iterative optimization process

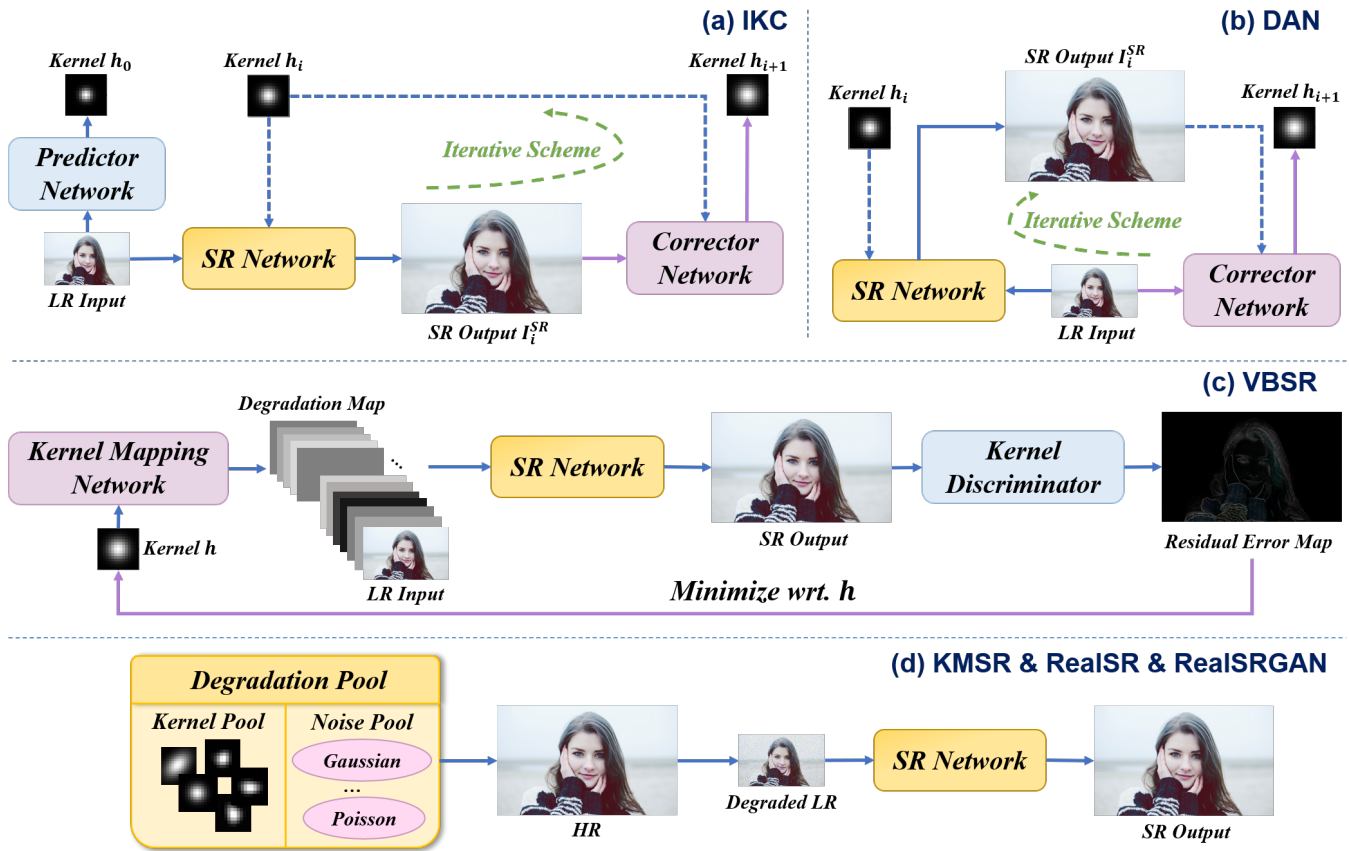


Fig. 13: Detailed frameworks of methods from *Image-Specific Adaptation with Kernel Estimation*. (a) IKC [7]; (b) DAN [48]; (c) VBSR [47]; (d) KMSR [23], RealSR [49] and RealSRGAN [50]. A connection with dotted line denotes a conditional input.



Fig. 14: Illustration on the limitations of *Image-Specific Adaptation with Kernel Estimation*. For input images with degradation types not covered in the SR model, like JPEG compression and Gaussian noise for IKC [7], the obtained SR results deteriorate dramatically. Best viewed on screen.

of DPSR into an end-to-end trainable network with iterative scheme, enabling joint optimization between the two sub-problems. A comparison between solution frameworks of DPSR and USRNet is depicted by Fig.11(b) and (c). Besides, some other methods exploiting plug-and-play technique include [51], [52], [53].

Limitation: In spite of the aforementioned progress, this kind of methods have one obvious drawback: they all rely on an additional input of degradation estimation, especially the SR kernel k . However, estimating the correct kernel from an arbitrary LR image is not an easy task, and an inaccurate estimation input will cause kernel mismatch and greatly undermine the SR performance [7], [12]. Fig.12 shows the comparison between SR results with correct and incorrect kernels based on the method SRMD. Therefore, only if

one has some method at hand for reliable degradation estimation can he quickly obtain a satisfactory SR output, otherwise he may come to the tedious work of manually choosing a proper estimation input for better result. Hence, we introduce another kind of approaches in the next part, which incorporate kernel estimation into SR framework for more robust performance.

6.1.2 Image-Specific Adaptation with Kernel Estimation

Iterative kernel correction (IKC) [7] proposes to correct kernel estimation in an iterative way to gradually approach a satisfactory result. The highlight of this method is to take advantage of the intermediate SR results, since the artifacts within an SR image caused by kernel mismatch tend to have regular patterns. Specifically, a corrector network is used to estimate the kernel correcting residual given an SR image

conditioned on the current kernel. Then the updated kernel is used to generate a new SR result with fewer artifacts. The SR network includes spatial feature transform [54] layers in each residual block, and the current kernel is used to generate transforming parameters for feature adaptation, which can be more effective than direct concatenation of inputs as proposed by SRMD. In addition, a predictor network is applied for kernel initialization based on the input LR image alone, and dimensionality stretching is adopted for kernel coding. A more recent work, deep alternating network (DAN) [48], further enhances the IKC framework. It unifies the corrector and SR network into an end-to-end trainable one instead of training each sub-network separately as IKC does. This joint training strategy can make the two networks more compatible to each other. Moreover, the corrector uses original LR input for kernel estimation conditioned on an intermediate SR result, which is beneficial to more robust kernel estimation performance. The overall frameworks of IKC and DAN are illustrated in Fig.13(a) and (b). The idea of making use of SR artifacts for kernel estimation is also employed in variant blind SR (VBSR) [47], yet it trains a kernel discriminator to estimate the error map of an SR output instead of the kernel itself, and finds the optimal kernel by minimizing the error of SR output during the inference stage, as shown by Fig.13(c). In addition to the SR kernel, estimation on more degradation types has also been studied. CBSR [14] combines two sub-networks for noise and kernel estimation with a non-blind SR network, thus forming a unified cascaded architecture for blind SR.

In fact, the iterative scheme adopted by IKC and DAN can be interpreted well from the perspective of domain adaptation: instead of producing the final SR output in a single stroke like SRMD, it chooses several intermediate SR results as interchange stations during the long trip from input LR to the target *Natural HR* domain, passing across the domain gap in Fig.2 step-by-step. These two methods can have more robust performance than SRMD framework depending on the accuracy of kernel estimation input.

Nevertheless, such an iterative scheme usually consumes more inference time and requires human intervention to choose the optimal number of iterations. To tackle these issues, some recent works propose non-iterative frameworks by introducing more accurate degradation estimation or more efficient feature adaptation strategies. Unsupervised degradation representation learning for blind SR (DRL-DASR) [55] tries to estimate the degradation information with a trainable encoder in the latent feature space, and the degradation encoder is trained with contrastive learning in an unsupervised manner. Specifically, LR samples with the same degradation as the query input are considered as the positive exemplars while those with different degradations are taken as negative ones. Then the mutual information among all samples is maximized in the latent space, leading to content-invariant degradation representations. Moreover, the estimated degradation representation is used to generate the corresponding convolutional kernels and modulation coefficients in SR network. Such a framework can achieve satisfactory SR results with a single forward pass. Kernel-oriented adaptive local adjustment (KOALANet) [56] also utilizes a similar dynamic kernel strategy that adapts the SR network to a specific degradation, and it further

extends the non-iterative framework to spatially-variant degradation with a downsampling network for local kernel estimation. Another work, adaptive modulation network with reinforcement learning (AMNet-RL) [57], proposes a modified version of adaptive instance norm (AdaIN) [58] to incorporate kernel estimation into the SR network, and it also pioneered in optimizing the blind SR model with in-differentiable perceptual metrics (e.g., NIQE [59]) under reinforcement learning framework.

There are also some other approaches proposing to learn a blind SR model by merely covering more degradations in the training dataset, especially more realistic kernels estimated from real images. For instance, kernel modelling super-resolution (KMSR) [23] builds a large kernel pool with data distribution learning based on some realistic SR kernels estimated from real LR images. Kernels from this pool are then used to synthesize HR-LR training pairs according to the classical degradation model, and the training process just follows non-blind setting with supervised learning. Usually, a more general training dataset enables the SR model to implicitly distinguish and adaptively deal with LR inputs with different degradations. In other words, the SR model will be implicitly endowed with more capacity for kernel estimation in the training process, thus avoiding explicit kernel estimation in the framework. However, such a direct way may not lead to top performance, as have been argued in [13]. A similar strategy is employed in RealSR [49] and RealSRGAN [50] to build more generic training dataset with more kinds of realistic kernels. This process is presented in Fig.13(d). Besides these methods, a correction filter [60] is designed for modifying an LR input to match the SR model with a pre-defined degradation, which is primarily based on the kernel estimation from LR.

Limitation: Compared with approaches without kernel estimation, these methods practically save us from efforts in searching for kernel estimation algorithms, especially during inference stage, and have demonstrated impressive performance. Yet, they still cannot avoid the inherent disadvantage of explicit modelling: they cannot give out satisfactory results for images with degradations not covered in their model. As presented in Fig.14, for an SR model focusing on the degradation caused by kernel k , like IKC, it can hardly deal with LR inputs with degradations out of its modelling scope. This limitation is really too tough for complex real-world images. Even if we are willing to retrain the model with more degradation types, it is impractical for us to explicitly model the degradation in an arbitrary LR and gather enough external training data, as stated in Sec.3. Next, let us step into another type of methods, which utilize a single input image alone for image-specific SR modelling.

6.2 Single Image Modelling with Internal Statistics

SR modelling with a single image is based on the internal statistics of natural images: patches of a single image tend to recur within and across different scales of this image. The internal statistics has been quantified and proved to have more predictive power than external statistics for many natural images [61]. A theoretical formulation is given by [62]. Specifically, one may assume that an HR image \mathbf{h} and



Fig. 15: Network structure of DGDML-SR [15]. Two training paths are included, which are based on HR and LR crops from different image depths, respectively.

its LR counterpart l are taken by the same camera, but with an s -scale zoom-in for the latter:

$$\mathbf{h}[n] = \int I(x)b_H\left(\frac{n}{s} - x\right)dx, \quad (8)$$

$$\mathbf{l}[n] = \int I(x)b_L(n - x)dx, \quad (9)$$

where $I(x)$ is the continuous-space image, b is the point spread function (PSF) of the camera, and b_H should be a downsampled version of b_L in the case of optical zoom:

$$b_H(x) = sb_L(sx). \quad (10)$$

Also, using the classical degradation modelling in Equation (4) without noise \mathbf{n} , the relationship between \mathbf{h} and \mathbf{l} expressed in discrete form is:

$$\mathbf{l}[n] = \sum_m \mathbf{h}[m]\mathbf{k}[sn - m]. \quad (11)$$

Now, for a given LR image, let q and r be two local patches with recurring pattern $P(x)$ in the continuous scene, where r is s -times larger than q . Then there will be:

$$\mathbf{r}[n] = \int P\left(\frac{x}{s}\right)b_L(n - x)dx = \int sP(x)b_L(n - sx)dx, \quad (12)$$

$$\mathbf{q}[n] = \int P(x)b_L(n - x)dx, \quad (13)$$

based on Equation (10), one can finally arrive:

$$\mathbf{q}[n] = \sum_m \mathbf{r}[m]\mathbf{k}[sn - m], \quad (14)$$

which means that the relationship between q and r in the same LR is equivalent to two patches from an HR and its LR version related by kernel \mathbf{k} . This property can be used to estimate \mathbf{k} and solve the unknown HR.

Glaser et al. [30] did the pioneering work in 2009 to introduce internal statistics into solving SISR problem from a single image. Latter, nonparametric blind SR (NPBSR) [62] further extends this framework to blind SR setting. Specifically, it proposes an MAP framework to estimate the SR blur kernel, based on the observation from Equation (14) that the optimal kernel \mathbf{k} is the one that maximizes the similarity among recurring patches across different scales. In addition, NPBSR proves that the optimal \mathbf{k} is not PSF of the camera but should be one with a smaller width, contrary to the common sense in its era.

The recent development of GAN gives birth to a new realization of using patch recurrence for blind kernel estimation. KernelGAN [6] interprets the maximization of patch recurrence within a single image as a data distribution learning problem. It assumes that the downsampled version of an LR image generated by the optimal \mathbf{k} should share the same patch distribution with the original LR. Under GAN framework, a deep linear network is used as generator to parameterize the underlying SR kernel, and a discriminator distinguishes generated patches from those in original LR image. Once the training finishes, one can explicitly obtain the kernel estimation by convolving together all convolutional filters in generator. It is worth noting that the training process relies merely on the input LR without any external dataset, which can be seen as self-supervised learning. Flow-based kernel prior (FKP) [63] develops a more effective approach for kernel optimization, where a kernel prior in latent space is learned with normalizing flow (NF) [64], [65] technique. Thanks to the invertible mapping between latent and pixel spaces enabled by NF, the search for the optimal \mathbf{k} can be conducted in the learned kernel manifold. This process can be more efficient than directly optimizing a randomly initialized deep network, thus leading to more robust kernel estimation results.

The idea of self-supervision based on patch recurrence property can also be directly applied to performing SR. Zero-shot super-resolution (ZSSR) [12], developed by authors from the same group as NPBSR and KernelGAN, made the very first attempt to train an image-specific CNN for super-resolving each input LR without any pre-training step. The training is conducted with HR-LR pairs generated from a single LR input \mathbf{y} , where \mathbf{y} is regarded as HR and coarser LR images are produced by downsampling with kernel \mathbf{k} . Data augmentation is utilized to make full use of information from the input image alone. The CNN trained with these image pairs will be capable of inferring specific relationships across different scales of \mathbf{y} , which is then used to super-resolve \mathbf{y} . In addition, ZSSR can be more robust to distracting artifacts (e.g., Gaussian noises, JPEG artifacts) by adding some noise to LR training samples, since it argues that only correlated image contents tend to recur across scales rather than noises.

In fact, ZSSR is still not well designed for blind setting: it requires estimated SR blur kernel \mathbf{k} as input to guide the generation of coarser LR images for training. A unified self-

supervision framework is thus proposed in DGDML-SR [15] - depth guided degradation model for learning-based SR. It combines a degradation network and an SR network into a single architecture, where the former is trained to simulate the degradation process, similar to the function of KernelGAN, and the latter aims to perform SR task like ZSSR does. This joint framework allows directly using generated LR as input to SR network without explicit extraction of SR kernel. In addition, DGDML-SR proposes to sample HR and LR patches in an unpaired way according to the depth map of input image, assuming that patches with smaller depth is equivalent to HR and those with larger depth to LR. A two-cycle training scheme similar to CycleGAN [66] structure is employed to simultaneously train the two networks (see Fig.15), where the unpaired HR and LR patches are used as real samples for data distribution learning.

Limitation: The idea of self-supervision with internal statistics seems attractive for solving SR images from LR with variant degradation types, since it requires no effort in gathering large external training dataset. Nevertheless, its basic assumption may easily fail, especially for natural images with diverse contents (e.g., animals) or monotonous scenes (e.g., sky), since it is hard to exploit recurring information across scales to robustly perform SR with this kind of input images. Hence, these approaches can only produce favourable SR outputs for a very limited set of images with frequently recurring contents across scales, and new methods for single-image modelling are waiting to be explored for more general natural images.

So far, we have had an overview on approaches with explicit degradation modelling, as well as their merits and demerits. Explicit modelling of degradation process is clear and straightforward, yet it can be too simple to model more complex degradations other than blurring and additive noise, such as real-world degradations originating from camera sensors. In fact, real-world images usually include multiple degradations, and we can hardly express these entangled factors with an explicit well-defined function. Hence, another group of methods propose to implicitly model the degradations through data distribution learning. To the best of our knowledge, so far there have only been approaches based on external dataset for implicit modelling, and we will talk about them in the following section.

7 IMPLICIT DEGRADATION MODELLING

7.1 Learning Data Distribution within External Dataset

This kind of approaches aim to implicitly grasp the underlying degradation model through learning with external dataset. For dataset with paired HR-LR images, supervised learning with cautious design of SR network may be already enough to achieve satisfactory results, just like the top solutions proposed in NTIRE 2018 [71] and AIM 2020 [18] challenges. A more difficult setting is learning with unpaired data, where ground truth of LR images with realistic degradations are unavailable. Existing approaches usually exploit data distribution learning with GAN framework [72], and one or more discriminators are used to distinguish generated image samples from real ones, pushing the generator towards the appropriate modelling direction. In most

cases, two datasets are used to train the model, including HR and unpaired LR respectively, and one may regard the two datasets as representing target and source domains for learning the domain adaptation.

Among the earliest attempts on implicit modelling with unpaired data is CinCGAN [8]. It proposes to first transform an LR input to the *Bicubic LR* domain before performing SR with a pre-trained non-blind model. The corresponding adaptation process is illustrated in Fig.16(a). The *Bicubic LR* domain is also regarded as *Clean LR*, because its samples are generated with bicubic downsampling from HR images and assumed to be without any noise. Two CycleGAN structures [66] are respectively applied to transformation from *LR* to *Clean LR* and to target *HR*, helping to maintain cycle consistency in the transformation process. In this way no paired data is required during training, thus forming an unsupervised training scheme. However, unsupervised domain adaptation is not an easy task, in that it is usually hard for a single 0/1 plane modelled by a discriminator to separate the right target domain.

To leverage both supervised learning and data distribution learning with GAN, some approaches focus on learning the degradation process from HR to LR, and the generated LR samples with realistic degradations are used to train the SR model in a paired manner. This process is depicted in Fig.16(b). Degradation GAN [67] adopts this scheme and combines a High-to-Low degradation generator with a Low-to-High SR one in a single framework. Specifically, the High-to-Low generator simulates the degradation from *HR* to *LR* domain with an *LR* discriminator, and the degraded LR image is used as input to Low-to-High generator for SR training. Adversarial loss is designed to dominate the training process, different from SRGAN [32] and ESRGAN [34] which mainly use pixel-wise loss for supervised training. Such learning strategy - unpaired degradation and paired SR - can also be found in some later works, including frequency separation for real-world SR (FSSR) [24] and frequency separation SRGAN (FS-SRGAN) [68]. These two methods both apply adversarial loss only to high-frequency contents in order to relieve the difficulty in adversarial training. FS-SRGAN further introduces a color attention module into its High-to-Low generator to alleviate color shifting problem during domain adaptation.

However, these approaches still cannot avoid the problems related to GAN framework. The generated LR images may have a large domain gap from real LR samples, thus undermining SR training performance. Towards addressing this issue, DASR [69] proposes a domain-gap aware training strategy, where both generated LR images and real LR samples from dataset are used to train the SR model (see Fig.16(c)). Another strategy called domain-distance weighted supervision is also employed to assign different loss weights to the LR inputs according to their domain distance to real *LR*, helping to reduce negative influence caused by generated far-away LR samples. Specifically, predictions from *LR* discriminator are used to quantitatively measure the domain distance for each LR. Another work, pseudo-supervision [70], combines the forward SR reconstruction path in CinCGAN architecture with degradation learning to deal with problems caused by domain gap. This approach keeps the forward adaptation route of CinCGAN from

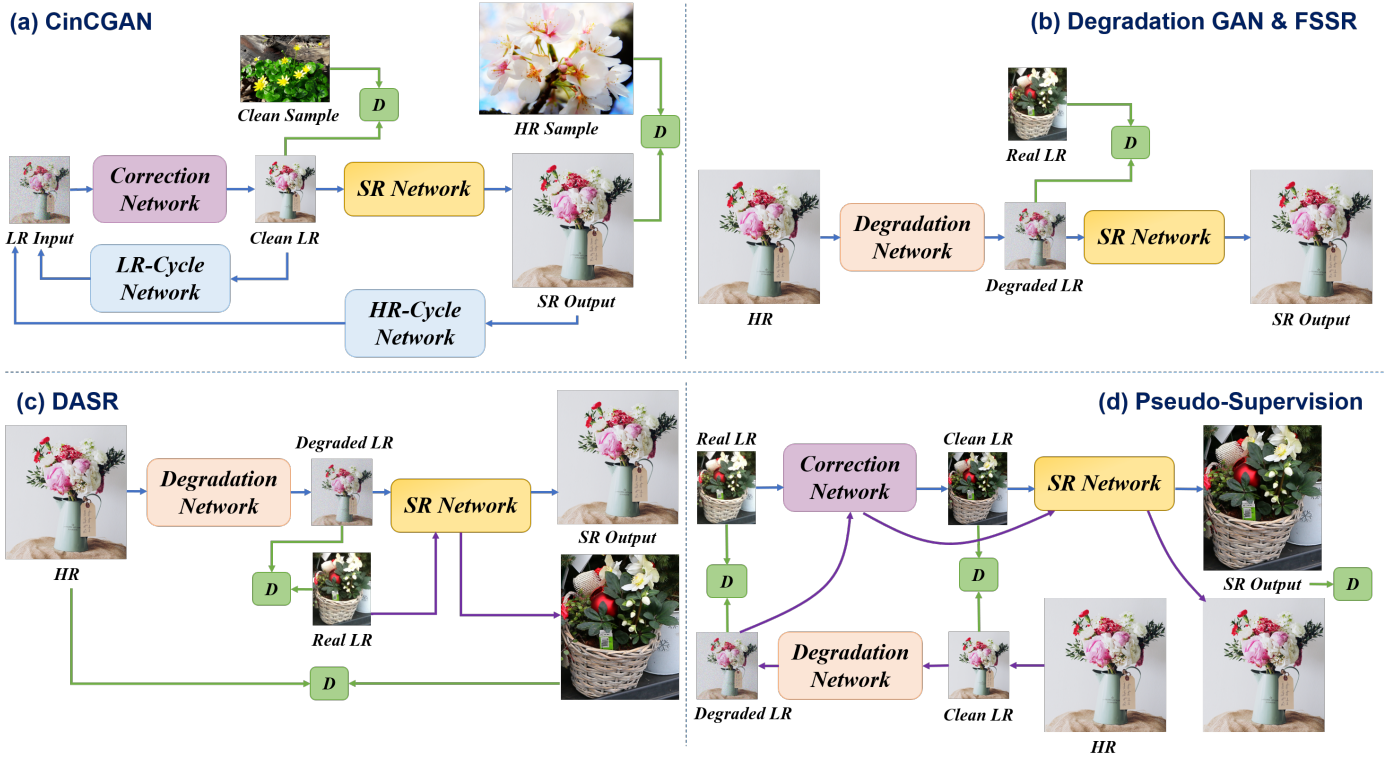


Fig. 16: Overall architectures of methods with implicit modelling for data distribution learning. (a)CinCGAN [8]; (b)Degradation GAN [67] and FSSR [24], FS-SRGAN [68]; (c)DASR [69]; (d)pseudo-supervision [70]. "D" represents discriminator network with GAN framework.

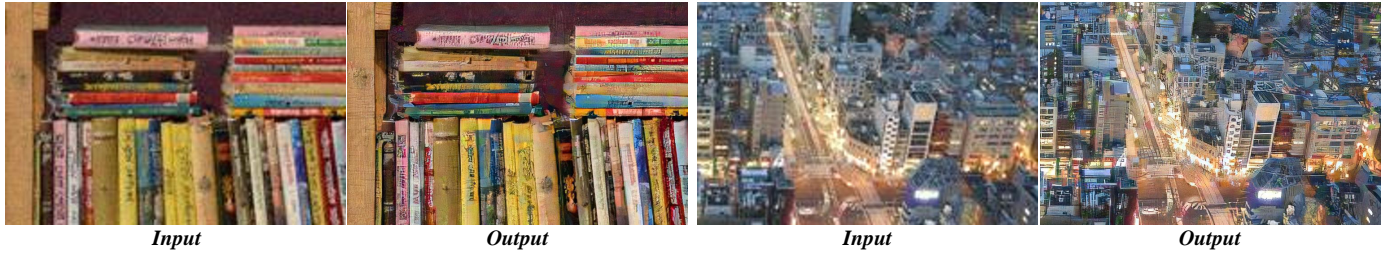


Fig. 17: Illustration on the limitations of *Learning Data Distribution within External Dataset*. For GAN-based framework, the SR results usually include severe artifacts and fake textures. Best viewed on screen.

degraded LR input to SR output, and adds a supervised path to it, as shown by Fig.16(d). This supervised path starts from HR to degraded LR, then goes back to HR across the intermediate *Clean LR* domain. This idea is in fact somewhat similar to the domain-gap aware training in DASR, if the CinCGAN route with *Real LR* is regarded as the additional path of real LR images besides the supervised one.

Limitation: Though seemingly flexible and powerful, this kind of methods is still far from a cure-all in blind SR. On the one hand, these methods must rely on large external datasets to learn the SR model through implicit data distribution, but this data-hungry manner is not suitable for certain tasks, including old photo restoration. On the other hand, most of them exploit GAN framework for unsupervised data distribution learning. The GAN-based framework may be difficult to train, and it will frequently produce severe artifacts in SR results. These artifacts are harmful to many real-world applications, such as high-definition display and old photo/film restoration. Readers

can refer to Fig.17 for illustration on the performance of a GAN-based method, FSSR.

7.2 Implicit Modelling with a Single Image: a Future Direction

The idea of implicit modelling seems promising to deal with complex real-world degradations, as long as the source LR and target HR datasets are provided. However, there is still a long way to go since existing approaches mainly rely on GAN framework for data distribution learning, and the artifacts caused by GAN are harmful to many real-world applications. Besides exploring for more robust generative models, another direction is also worth noting, which has never been proposed so far: implicit modelling with a single LR image, as revealed in Fig.7.

As stated in previous sections, existing approaches all have their limitations, especially for the outliers of internal statistics with complex real-world degradations. Examples of such kind include surveillance video, old photos and

films, as demonstrated in Sec.3. These images are commonly seen in our daily life, and they pose great challenges to existing methods: they not only lack patch redundancy across scales to serve as hint for explicit modelling, but also cannot be covered with a few external datasets due to the complexity of unpredictable degradations. We argue that this research gap can be possibly filled by methods based on implicit modelling with a single input image. Up till now, there have been no related work in this field, but we believe this is a worthwhile direction for future research.

The main difficulty lying in this direction is the lack of effective SR prior, and one possible solution is to apply human intervention as additional information. For example, image restoration network with modulation [73] can be used to manipulate the SR outputs with a proper controlling coefficient, or to manually choose a clear image with similar contents to the LR input as an SR reference. The key point of these proposals is to increase the amount of useful information to make blind SR possible.

8 DATASETS AND COMPETITIONS

8.1 Datasets

A large portion of methods covered in this paper, especially those with explicit degradation modelling and external dataset, require HR-LR image pairs for solving and evaluating SR models. However, due to the difficulty of obtaining real paired data, so far there have been only a few real-world datasets and most methods still synthesize LR inputs from HR images. Sec.8.1.1 discusses the common ways of building synthetic dataset, and Sec.8.1.2 gives an introduction on a few available real-world datasets.

8.1.1 Synthetic Dataset

For methods with explicit degradation modelling, the process of synthesizing degraded LR images from HR ground truth usually follows Equation (4), where a kernel k or noise with level σ is sequentially applied to the HR image together with the downscaling operation. For kernel k , Gaussian kernels have been the most widely adopted kernel type. A typical practice is introduced by SRMD [13]. An isotropic Gaussian kernel can be generated with a kernel width uniformly sampled from a pre-defined range, while an anisotropic one is characterized with a covariance matrix Σ , where the rotation angle of its eigenvectors and the corresponding eigenvalues determine the kernel shape. As for noise n , additive Gaussian noise is mostly used to simulate real-world noises, and the level σ can also be sampled from a specific range. Examples of Gaussian kernels and different noise levels are included in Fig.3(a). In addition, one can also enlarge the dataset with more realistic kernel or noise types, just as done by KMSR [23] and RealSR [49].

The most popular HR datasets in the non-blind setting are also employed in blind SR. For example, DIV2K [74] and Flickr2K [75] are often used for training, while Set5 [76], Set14 [77], BSD100 [78] and Urban100 [10] are usually for testing. Specifically, a blind SR benchmark *DIV2KRRK* is proposed in KernelGAN [6], where each image in DIV2K validation set is blurred and downsampled by a randomly generated anisotropic Gaussian kernel with some multiplicative noise to simulate more complex degradation. Other

synthetic datasets with unknown degradation, e.g. DIV2K wild [71], are also used by some methods with implicit modelling like CinCGAN [8].

8.1.2 Real-World Dataset

Up till now, there have been several real-world datasets with paired HR-LR images, and these datasets are built with carefully designed techniques and advanced digital devices. Representatives are City100 [79], DRealSR [80] and RealSR [81]. Among them, DRealSR is the largest one with around 800 image pairs for each scale factor. Usually, an HR image and its corresponding LR observations are captured by adjusting the focal lengths of imaging devices, then HR-LR pairs are accurately aligned with image registration and color rectification. Some other real-world datasets without HR ground truth have also been used as source domain images for SR network inputs, like DPED [11] in NTIRE 2020 Real-World Image SR challenge [17]. Compared with synthetic data, these real-world datasets serve as an important benchmark for investigating blind SR in real setting. However, building such a dataset is time-consuming and expensive, and also cannot cover all scenarios due to complicated variations among different imaging systems.

8.2 Competitions

In order to gauge and promote the development of superior solutions, some competitions have been held in the field of SISR, including some tracks with unknown degradation or blind SR setting. Here, we make a summary on previous competitions related to blind SR in Table 1, hoping to reveal some research trends from another perspective.

As the first competition related to blind setting, NTIRE 2017 challenge on SISR goes one step ahead from bicubic downscaling and assumes unknown blur and decimation operations to get LR images. More complex degradation types are then introduced by succeeding competitions, including real-world degradations arising from digital cameras. These competitions help to probe the state-of-the-art in this field. However, since paired training dataset is provided, their top solutions were largely focused on network structure enhancement as for non-blind SR. AIM 2019 challenge made the first attempt to address the real-world setting where paired training data is unavailable, hoping to stimulate research endeavours towards unsupervised learning. It can be seen that rising attention has been paid to this emerging task, especially from its follower NTIRE 2020, and the research community is expected to propose more novel solutions to fit into real-world problems.

9 QUANTITATIVE COMPARISON

This section includes detailed analysis on the performance and limitations of different methods with some testing examples. In fact, it is not an easy task to provide a comprehensive and fair comparison here, and we will state several problems hindering this practice in Sec.9.1. Despite the difficulty, we still present some testing results based on officially released pre-trained models in Sec.9.2, in order to provide readers with some useful information on the performance of each kind of methods.

TABLE 1: Details of challenges on blind image super-resolution. "T" denotes "track", "Synthe" means synthetic data generated from HR images, and — indicates the same with previous row.

Competition Name	Task & Degradation (Degrad.)	Train Set	Teams	Champion Solution
NTIRE 2017 SISR [75]	T2: Unknown downscaling, only blur	DIV2K Synthe. (paired)	17	EDSR [82]: modified SRResNet + a downsampling network to learn the degradation for new training data generation with Flickr2K
NTIRE 2018 SISR [71]	T2: Realistic mild condition, motion blur & Poisson noise, same degradation level for all images	—	18	WDSR [83]: modified EDSR + pre-alignment to reduce random shift effects between HR-LR pairs
	T3: Realistic difficult condition, stronger degradation than T2	—	17	
	T4: Realistic wild condition, different degradation levels for each image	—	14	PDN [84]: PolyDense block motivated by PolyNet and DenseNet
NTIRE 2019 Real Image SR [85]	Real-world SISR proposed in [81], real-world HR-LR image pairs captured with different focal lens of DSLR cameras	RealSR (paired)	36	UDSR [86]: U-shaped network with U-Net structure + three-stage cascaded refinement framework with coarse to fine supervision
AIM 2019 Real-World Image SR [87]	T1: Same domain, synthetic degradation in source domain, SR results preserve low-level image characteristics, i.e. keep source domain quality	Flickr2K Synthe., no target (unpaired)	7	FSSR [24]: DSGAN to learn to generate LR images with degradations in source domain + ESRGAN trained with generated paired data + frequency-separation loss
	T2: Target domain, synthetic degradation in source domain, SR results should be of high-quality defined with target domain	Flickr2K Synthe., DIV2K (unpaired)	4	
NTIRE 2020 Real World Image SR [17]	T1: Image processing artifacts, unknown synthetic degradation generated from image processing methods, source (input) and target domains are unpaired for training	—	16	RealSR [49]: build more generic paired training dataset by covering more realistic kernels and noise types
	T2: Smartphone images, real images captured with a low-quality smartphone camera	DPED [11] iPhone3, DIV2K (unpaired)	14	
AIM 2020 Real Image SR [18]	Real image SISR proposed in [80], real natural images captured by five DSLR cameras with aligned HR-LR pairs	DRealSR [80] (paired)	24	Proposed by Baidu team [88]: apply Neural Architecture Search among networks composed of dense residual blocks

9.1 Problem: Difficulty of Fair Comparison

- *Inaccessible Code*: Among the methods covered in our survey, some have released their full set of codes for training and testing, such as SRMD, IKC, RealSR, KernelGAN with ZSSR, and FSSR. But there are also some others which are not able to make their codes publicly available, and it is also non-trivial to reproduce these methods due to the lack of some important implementation details in original paper, especially for GAN-based approaches.
- *Different Training Data*: Even though the official implementations and pre-trained models are at hand, we still cannot sit back and relax for fair comparison since most of them adopt different datasets or different degradation types during training. These variables can greatly affect their performance and

generalization capacity, especially for the blind setting with a large variety of distinct degradations.

Hence, it is necessary and demanding to set up a benchmark for existing approaches. This benchmark should provide a fair and comprehensive comparison based on a uniform platform, and helps to gauge and push forward the state-of-the-art for blind SR. In this survey, we only show some quantitative results from published papers and a few testing examples produced by released pre-trained models.

9.2 Comparison and Analysis Based on Pre-Trained Models

This section includes an overview on the quantitative performance of some representative approaches. Note that all the quantitative results shown in this section are copied from related papers. In order to be as fair as possible, every

single table in this section is from a single paper, so results in the same table are ensured with a fair comparison. Also, each table mainly corresponds to a specific category in our proposed taxonomy. These results can provide readers with a straightforward comparison between different methods in each category.

For approaches with explicit modelling, we first show in Table 2 the comparison between two kernel estimation methods (NPBSR [62] and KernelGAN [6] in Sec.6.2) combined with two SR models with kernel estimation input (SRMD [13] in Sec.6.1.1, ZSSR [12] in Sec.6.2). These results are reported in KernelGAN paper on the testing dataset DIV2K, which consists of 100 DIV2K validation images blurred and downsampled using randomly generated anisotropic Gaussian kernels. These results indicate that methods using degradation information as additional input can well be fitted into blind setting if combined with an appropriate kernel estimation algorithm. However, there still remains a considerable performance gap between using ground truth kernel and applying estimation methods, since it is non-trivial to accurately estimate degradation information from an arbitrary image. Also, these methods can hardly outperform those combining kernel estimation and SR into a single framework, such as DAN, showing the advantage of joint optimization.

The results of another two representative methods with explicit modelling and incorporated kernel estimation, namely IKC [7] and DAN [48] in Sec.6.1.2, are shown in Table 3. These results are copied from paper of DAN [48]. The testing data includes two classic datasets widely used in SR task, namely Set14 [77] and BSD100 [78], and HR images are blurred and downsampled using eight selected isotropic Gaussian kernels for each scale factor to synthesize degraded LR. According to the comparison, DAN outperforms IKC as well as ZSSR for all the three scales.

For methods with implicit modelling, we show the results from DASR [69] in Table 4, including comparison among three representative data distribution learning models: CinCGAN [8], FSSR [24] and DASR [69]. Besides PSNR and SSIM, LPIPS [89] is used to better evaluate the perceptual quality of SR images generated from GAN. Two dataset are used for testing: AIM [87] is released by AIM Challenge on Real World SR in ICCV 2019, while RealSR [81] is a real-world dataset composed of HR-LR pairs captured with different focal lengths of the camera. We can see that DASR demonstrates the best performance, especially in terms of visual quality, owing to its better training strategies for narrowing the domain gap.

Qualitative comparison of some testing examples is shown in Fig.18~20, and we use pre-trained models provided by the authors of the corresponding papers. The first two degraded LR images are synthesized with isotropic Gaussian blur or additive Gaussian noise, and the third one is a real-world image from NTIRE real-world SR challenge [17] without ground truth. Note that for SRMD and USRNet, we directly use the ground truth kernel or noise level as additional inputs in order to validate the efficacy of the SR model alone. Also, for RealSR and FSSR, while there have been multiple pre-trained models based on different training datasets, we choose the one trained with real images taken by mobile devices from DPED dataset [11]. In

TABLE 2: Quantitative comparison (PSNR(dB)/SSIM) of image-specific adaptation without kernel estimation (SRMD [13], ZSSR [12]) plus blind kernel estimation (NPBSR [62], KernelGAN [6]).

Method	Scale	DIV2K
Bicubic		28.73 / 0.8040
Bicubic kernel + ZSSR		29.10 / 0.8215
NPBSR + SRMD		25.51 / 0.8083
NPBSR + ZSSR		29.37 / 0.8370
KernelGAN + SRMD	2	29.57 / 0.8564
KernelGAN + ZSSR		30.36 / 0.8669
Ground-truth kernel + SRMD		31.96 / 0.8955
Ground-truth kernel + ZSSR		32.44 / 0.8992
DAN		32.56 / 0.8997
Bicubic		25.33 / 0.6795
Bicubic kernel + ZSSR		25.61 / 0.6911
NPBSR + SRMD		23.34 / 0.6530
NPBSR + ZSSR		26.08 / 0.7138
KernelGAN + SRMD	4	25.71 / 0.7265
KernelGAN + ZSSR		26.81 / 0.7316
Ground-truth kernel + SRMD		27.38 / 0.7655
Ground-truth kernel + ZSSR		27.53 / 0.7446
DAN		27.55 / 0.7582

TABLE 3: Quantitative comparison (PSNR(dB)/SSIM) of image-specific adaptation with kernel estimation (IKC [7], DAN [48]).

Method	Scale	Set14	BSD100
Bicubic		26.02 / 0.7634	25.92 / 0.7310
ZSSR	2	28.35 / 0.7933	27.92 / 0.7632
IKC		32.82 / 0.8999	31.36 / 0.9097
DAN		33.07 / 0.9068	31.76 / 0.9213
Bicubic		24.01 / 0.6662	24.25 / 0.6356
ZSSR	3	26.11 / 0.6942	26.06 / 0.6633
IKC		29.46 / 0.8229	28.56 / 0.8493
DAN		30.09 / 0.8287	28.94 / 0.8606
Bicubic		22.79 / 0.6032	23.29 / 0.5786
ZSSR	4	24.78 / 0.6268	24.97 / 0.5989
IKC		28.26 / 0.7688	27.29 / 0.8014
DAN		28.43 / 0.7693	27.51 / 0.8078

addition, we include two renowned non-blind SR models, SRRResNet [4] and ESRGAN [34], into the comparison list for readers' reference. Based on these testing examples, some important observations can be drawn as following:

- (1) For methods exploiting external dataset, their generalization largely depends on the coverage of degradation modelling or training data distribution. For example, approaches with explicit modelling can only handle noisy inputs if noise is directly covered as a degradation factor in the SR modelling, like SRMD and USRNet. Otherwise, they will not be capable of noise removal and consequently bring unfavourable artifacts in SR results, like SRRResNet and IKC. On the other hand, models trained on real image dataset hardly give out visually pleasing results on synthetic data, such as RealSR and FSSR.

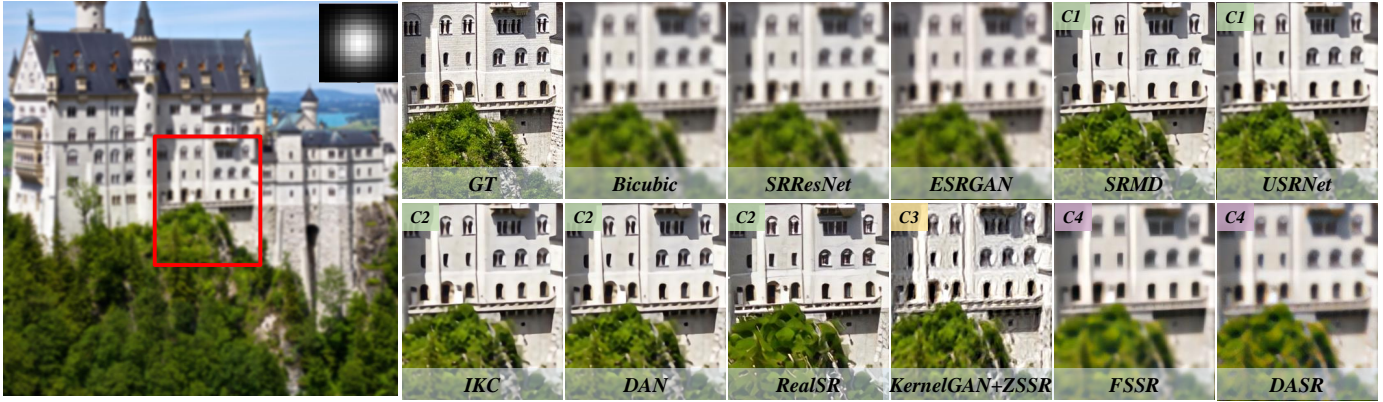


Fig. 18: Qualitative comparison of synthetic testing example with Gaussian blur.



Fig. 19: Qualitative comparison of synthetic testing example with Gaussian blur and noise.



Fig. 20: Qualitative comparison of real testing image captured with iPhone.

TABLE 4: Quantitative comparison of implicit modelling with data distribution learning (CinCGAN [8], FSSR [24], DASR [69]). Note that ESRGAN [34] is trained with paired data.

Method	PSNR(dB) / SSIM / LPIPS↓	
	AIM	RealSR
ESRGAN	- / - / -	25.70 / 0.7487 / 0.20
ZSSR	22.33 / 0.6022 / 0.63	26.01 / 0.7482 / 0.39
CinCGAN	21.60 / 0.6129 / 0.46	25.09 / 0.7459 / 0.41
FSSR	20.82 / 0.5103 / 0.39	25.99 / 0.7388 / 0.27
DASR	21.60 / 0.5640 / 0.34	26.78 / 0.7822 / 0.23

(2) Real-world images *do* include more complex degradations, which deviate a lot from synthetic data dis-

tribution. Models with explicit modelling generally perform well on the synthetic images within their degradation coverage, like SRMD and IKC. However, none of the methods generate satisfactory result for the third real image, including those with implicit modelling. Specifically, SRResNet and IKC tend to keep the noise texture and cause artifacts. SRMD and USRNet alleviate the noisy artifacts to some extent but in the sacrifice of high-frequency details, and it also needs some effort to estimate kernel and noise level for a real image with unknown degradations. RealSR and FSSR do have better results in terms of removing noise and preserving sharp textures, yet they still generate fake textures or artifacts, leading to the unnatural-looking of SR images.

9.3 Suggestions on Fair Comparison

We would like to suggest that our readers make better use of our proposed taxonomy in their future work, especially for effective and fair comparison among different methods in their paper. One may first try to place their own method into the corresponding category, and then pay special attention to previous methods belonging to the same category for evaluation and comparison. As for methods out of the specific categorial scope, since they have employed different degradation modelling or data sources (external or internal), direct comparison may become unfair and unnecessary, sometimes causing difficulty and even confusion. Hence, we recommend here that future work may as well follow our taxonomy to make comparison among methods from the same category, both for proposing new solutions and benchmark setting.

10 CONCLUSION

In this paper, we present a systematic survey on recent progress in blind image SR. In order to effectively classify and summarize existing methods, we propose a taxonomy according to their ways of degradation modelling and the data used for solving the SR model: explicit or implicit modelling with external dataset or a single LR image. Except implicit modelling with a single image, the other three categories all have representative existing approaches, and we make a conclusion on them as following:

- **Explicit modelling with external dataset:** representatives are SRMD and IKC, which utilize the classical degradation model or its variants for image-specific adaptation based on degradation information. These methods perform well on degradation types covered in its modelling, but their performance will severely deteriorate on other degradations.
- **Explicit modelling with a single LR image:** including NPBSR and KernelGAN for blind kernel estimation, as well as ZSSR and DGDML-SR for SR. They leverage the internal statistics of natural images - patch recurrence across scales, which can also be theoretically derived from classical degradation model. However, these methods may fail for more general natural images with diverse or monotonous scenes, i.e., those without enough recurring clues for SR task.
- **Implicit modelling with external dataset:** such as CinCGAN, FSSR and DASR. These methods assume that real-world degradations are too complex to be explicitly modelled, but can be implicitly learned with data distribution learning under GAN framework. However, domain adaptation is a non-trivial task due to the large space of the natural image domain, making it hard to train a GAN-based model with good performance and generalization capacity.

From our perspective, implicit modelling with a single image, which has not been proposed yet, is a direction worth exploring in future research, especially for general natural images with complex degradations and without strong internal statistics. One possible solution to this problem is

utilizing human intervention to provide additional information as SR prior, and restoration network with modulation or manually choosing an HR reference image may be of help. We hope this paper can inspire some new ideas for future work and make contributions to the prosperity of blind image SR techniques.

REFERENCES

- [1] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *ECCV*, 2014, pp. 184–199.
- [2] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *ECCV*, 2016, pp. 391–407.
- [3] 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 *CVPR*, 2016, pp. 1874–1883.
- [4] 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 *CVPR*, 2017, pp. 4681–4690.
- [5] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," in *ECCV*, 2018, pp. 286–301.
- [6] S. Bell-Kligler, A. Shocher, and M. Irani, "Blind super-resolution kernel estimation using an internal-gan," in *NeurIPS*, 2019, pp. 284–293.
- [7] J. Gu, H. Lu, W. Zuo, and C. Dong, "Blind super-resolution with iterative kernel correction," in *CVPR*, 2019, pp. 1604–1613.
- [8] Y. Yuan, S. Liu, J. Zhang, Y. Zhang, C. Dong, and L. Lin, "Un-supervised image super-resolution using cycle-in-cycle generative adversarial networks," in *CVPRW*, 2018, pp. 701–710.
- [9] S. Gu, A. Lugmayr, M. Danelljan, M. Fritsche, J. Lamour, and R. Timofte, "DIV8K: diverse 8k resolution image dataset," in *2019 IEEE/CVF International Conference on Computer Vision Workshops, ICCV Workshops 2019*. IEEE, 2019, pp. 3512–3516.
- [10] J. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *CVPR*, 2015, pp. 5197–5206.
- [11] A. Ignatov, N. Kobyshev, R. Timofte, K. Vanhoey, and L. V. Gool, "Dslr-quality photos on mobile devices with deep convolutional networks," in *ICCV*, 2017, pp. 3297–3305.
- [12] A. Shocher, N. Cohen, and M. Irani, "zero-shot" super-resolution using deep internal learning," in *CVPR*, 2018, pp. 3118–3126.
- [13] K. Zhang, W. Zuo, and L. Zhang, "Learning a single convolutional super-resolution network for multiple degradations," in *CVPR*, 2018, pp. 3262–3271.
- [14] P. Liu, H. Zhang, Y. Cao, S. Liu, D. Ren, and W. Zuo, "Learning cascaded convolutional networks for blind single image super-resolution," *Neurocomputing*, vol. 417, pp. 371–383, 2020.
- [15] X. Cheng, Z. Fu, and J. Yang, "Zero-shot image super-resolution with depth guided internal degradation learning," in *ECCV*, 2020.
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014, pp. 2672–2680.
- [17] A. Lugmayr, M. Danelljan, and R. Timofte, "Ntire 2020 challenge on real-world image super-resolution: Methods and results," in *CVPRW*, 2020, pp. 494–495.
- [18] P. Wei, H. Lu, R. Timofte, L. Lin *et al.*, "AIM 2020 challenge on real image super-resolution: Methods and results," in *Computer Vision - ECCV 2020 Workshops*, A. Bartoli and A. Fusiello, Eds., vol. 12537. Springer, 2020, pp. 392–422.
- [19] S. Oh, A. Hoogs, A. G. A. Perera, N. P. Cuntoor, and *et al.*, "A large-scale benchmark dataset for event recognition in surveillance video," in *The 24th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011*. IEEE Computer Society, 2011, pp. 3153–3160.
- [20] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018*. IEEE Computer Society, 2018, pp. 6479–6488.
- [21] G. Qian, J. Gu, J. S. Ren, C. Dong, F. Zhao, and J. Lin, "Trinity of pixel enhancement: a joint solution for demosaicking, denoising and super-resolution," *arXiv preprint arXiv:1905.02538*, 2019.

- [22] Z. Wan, B. Zhang, D. Chen, P. Zhang, D. Chen, J. Liao, and F. Wen, "Old photo restoration via deep latent space translation," *CoRR*, vol. abs/2009.07047, 2020.
- [23] R. Zhou and S. Susstrunk, "Kernel modeling super-resolution on real low-resolution images," in *ICCV*, 2019, pp. 2433–2443.
- [24] M. Fritsche, S. Gu, and R. Timofte, "Frequency separation for real-world super-resolution," in *ICCV Workshop*, 2019, pp. 3599–3608.
- [25] H. Chang, D. Yeung, and Y. Xiong, "Super-resolution through neighbor embedding," in *2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2004)*. IEEE Computer Society, 2004, pp. 275–282.
- [26] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution as sparse representation of raw image patches," in *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008)*. IEEE Computer Society, 2008.
- [27] R. Timofte, V. D. Smet, and L. V. Gool, "Anchored neighborhood regression for fast example-based super-resolution," in *IEEE International Conference on Computer Vision, ICCV 2013*. IEEE Computer Society, 2013, pp. 1920–1927.
- [28] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE Trans. Image Process.*, vol. 19, no. 11, pp. 2861–2873, 2010.
- [29] G. Freedman and R. Fattal, "Image and video upscaling from local self-examples," *ACM Trans. Graph.*, vol. 30, no. 2, pp. 12:1–12:11, 2011.
- [30] D. Glasner, S. Bagon, and M. Irani, "Super-resolution from a single image," in *IEEE 12th International Conference on Computer Vision, ICCV 2009*. IEEE Computer Society, 2009, pp. 349–356.
- [31] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," in *Computer Vision - ECCV 2014*, D. J. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds., vol. 8692. Springer, 2014, pp. 184–199.
- [32] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, "Photo-realistic single image super-resolution using a generative adversarial network," in *CVPR*, 2017, pp. 105–114.
- [33] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*. IEEE Computer Society, 2016, pp. 1646–1654.
- [34] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "ESRGAN: enhanced super-resolution generative adversarial networks," in *ECCV Workshop*, ser. Lecture Notes in Computer Science, vol. 11133, 2018, pp. 63–79.
- [35] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017*. IEEE Computer Society, 2017, pp. 1132–1140.
- [36] T. Tong, G. Li, X. Liu, and Q. Gao, "Image super-resolution using dense skip connections," in *IEEE International Conference on Computer Vision, ICCV 2017*. IEEE Computer Society, 2017, pp. 4809–4817.
- [37] J. Kim, J. K. Lee, and K. M. Lee, "Deeply-recursive convolutional network for image super-resolution," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*. IEEE Computer Society, 2016, pp. 1637–1645.
- [38] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*. IEEE Computer Society, 2017, pp. 2790–2798.
- [39] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," in *Computer Vision - ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11211. Springer, 2018, pp. 294–310.
- [40] T. Dai, J. Cai, Y. Zhang, S. Xia, and L. Zhang, "Second-order attention network for single image super-resolution," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019*. Computer Vision Foundation / IEEE, 2019, pp. 11 065–11 074.
- [41] W. Shi, J. Caballero, F. Huszar, 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 *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*. IEEE Computer Society, 2016, pp. 1874–1883.
- [42] M. S. M. Sajjadi, B. Schölkopf, and M. Hirsch, "Enhancenet: Single image super-resolution through automated texture synthesis," in *IEEE International Conference on Computer Vision, ICCV 2017*. IEEE Computer Society, 2017, pp. 4501–4510.
- [43] W. Zhang, Y. Liu, C. Dong, and Y. Qiao, "Ranksrgan: Generative adversarial networks with ranker for image super-resolution," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019*. IEEE, 2019, pp. 3096–3105.
- [44] Y.-S. Xu, S.-Y. R. Tseng, Y. Tseng, H.-K. Kuo, and Y.-M. Tsai, "Unified dynamic convolutional network for super-resolution with variational degradations," in *CVPR*, 2020, pp. 12 496–12 505.
- [45] K. Zhang, W. Zuo, and L. Zhang, "Deep plug-and-play super-resolution for arbitrary blur kernels," in *CVPR*, June 2019.
- [46] K. Zhang, L. V. Gool, and R. Timofte, "Deep unfolding network for image super-resolution," in *CVPR*, 2020, pp. 3214–3223.
- [47] V. Cornillere, A. Djelouah, W. Yifan, O. Sorkine-Hornung, and C. Schroers, "Blind image super-resolution with spatially variant degradations," *ACM Trans. Graph.*, vol. 38, no. 6, pp. 1–13, 2019.
- [48] Y. Huang, S. Li, L. Wang, T. Tan *et al.*, "Unfolding the alternating optimization for blind super resolution," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [49] X. Ji, Y. Cao, Y. Tai, C. Wang, J. Li, and F. Huang, "Real-world super-resolution via kernel estimation and noise injection," in *CVPRW*, June 2020.
- [50] H. Ren, A. Kheradmand, M. El-Khamy, S. Wang, D. Bai, and J. Lee, "Real-world super-resolution using generative adversarial networks," in *CVPRW*, 2020, pp. 436–437.
- [51] A. Brifman, Y. Romano, and M. Elad, "Turning a denoiser into a super-resolver using plug and play priors," in *2016 IEEE International Conference on Image Processing, ICIP 2016*. IEEE, 2016, pp. 1404–1408.
- [52] S. H. Chan, X. Wang, and O. A. Elgandy, "Plug-and-play ADMM for image restoration: Fixed-point convergence and applications," *IEEE Trans. Computational Imaging*, vol. 3, no. 1, pp. 84–98, 2017.
- [53] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep CNN denoiser prior for image restoration," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*. IEEE Computer Society, 2017, pp. 2808–2817.
- [54] X. Wang, K. Yu, C. Dong, and C. Change Loy, "Recovering realistic texture in image super-resolution by deep spatial feature transform," in *CVPR*, 2018, pp. 606–615.
- [55] L. Wang, Y. Wang, X. Dong, Q. Xu, J. Yang, W. An, and Y. Guo, "Unsupervised degradation representation learning for blind super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10 581–10 590.
- [56] S. Y. Kim, H. Sim, and M. Kim, "Koalernet: Blind super-resolution using kernel-oriented adaptive local adjustment," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10 611–10 620.
- [57] Z. Hui, J. Li, X. Wang, and X. Gao, "Learning the non-differentiable optimization for blind super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 2093–2102.
- [58] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1501–1510.
- [59] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.
- [60] S. A. Hussein, T. Tirer, and R. Giryes, "Correction filter for single image super-resolution: Robustifying off-the-shelf deep super-resolvers," in *CVPR*, 2020, pp. 1428–1437.
- [61] M. Zontak and M. Irani, "Internal statistics of a single natural image," in *CVPR*, 2011, pp. 977–984.
- [62] T. Michaeli and M. Irani, "Nonparametric blind super-resolution," in *ICCV*, December 2013.
- [63] J. Liang, K. Zhang, S. Gu, L. Van Gool, and R. Timofte, "Flow-based kernel prior with application to blind super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 10 601–10 610.
- [64] L. Dinh, D. Krueger, and Y. Bengio, "NICE: non-linear independent components estimation," in *3rd International Conference on Learning Representations, ICLR 2015, Workshop Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.
- [65] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real NVP," in *5th International Conference on Learning Representations, ICLR 2017, Conference Track Proceedings*, 2017.

- [66] J. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017, pp. 2242–2251.
- [67] A. Bulat, J. Yang, and G. Tzimiropoulos, "To learn image super-resolution, use a gan to learn how to do image degradation first," in *ECCV*, September 2018.
- [68] Y. Zhou, W. Deng, T. Tong, and Q. Gao, "Guided frequency separation network for real-world super-resolution," in *CVPRW*, 2020, pp. 428–429.
- [69] Y. Wei, S. Gu, Y. Li, R. Timofte, L. Jin, and H. Song, "Unsupervised real-world image super resolution via domain-distance aware training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 13 385–13 394.
- [70] S. Maeda, "Unpaired image super-resolution using pseudo-supervision," in *CVPR*, 2020, pp. 291–300.
- [71] R. Timofte, S. Gu, J. Wu, L. Van Gool, L. Zhang, M.-H. Yang, M. Haris *et al.*, "Ntire 2018 challenge on single image super-resolution: Methods and results," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2018.
- [72] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial networks," *CoRR*, vol. abs/1406.2661, 2014.
- [73] J. He, C. Dong, and Y. Qiao, "Modulating image restoration with continual levels via adaptive feature modification layers," in *CVPR*, 2019, pp. 11 056–11 064.
- [74] E. Agustsson and R. Timofte, "Ntire 2017 challenge on single image super-resolution: Dataset and study," in *CVPRW*, 2017, pp. 126–135.
- [75] R. Timofte, E. Agustsson, L. V. Gool, M. Yang, L. Zhang, B. Lim, S. Son *et al.*, "NTIRE 2017 challenge on single image super-resolution: Methods and results," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017*. IEEE Computer Society, 2017, pp. 1110–1121.
- [76] M. Bevilacqua, A. Roumy, C. Guillemot, and M. Alberi-Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," in *British Machine Vision Conference, BMVC 2012*, R. Bowden, J. P. Collomosse, and K. Mikolajczyk, Eds. BMVA Press, 2012, pp. 1–10.
- [77] R. Zeyde, M. Elad, and M. Protter, "On single image scale-up using sparse-representations," in *International conference on curves and surfaces*. Springer, 2010, pp. 711–730.
- [78] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *CVPR*, 2015, pp. 5197–5206.
- [79] C. Chen, Z. Xiong, X. Tian, Z. Zha, and F. Wu, "Camera lens super-resolution," in *CVPR*, 2019, pp. 1652–1660.
- [80] P. Wei, Z. Xie, H. Lu, Z. Zhan, Q. Ye, W. Zuo, and L. Lin, "Component divide-and-conquer for real-world image super-resolution," in *Computer Vision - ECCV 2020*, vol. 12353. Springer, 2020, pp. 101–117.
- [81] J. Cai, H. Zeng, H. Yong, Z. Cao, and L. Zhang, "Toward real-world single image super-resolution: A new benchmark and a new model," in *ICCV*, 2019, pp. 3086–3095.
- [82] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *CVPRW*, 2017, pp. 136–144.
- [83] J. Yu, Y. Fan, J. Yang, N. Xu, Z. Wang, X. Wang, and T. S. Huang, "Wide activation for efficient and accurate image super-resolution," *CoRR*, vol. abs/1808.08718, 2018.
- [84] X. Wang, K. Yu, T.-W. Hui, C. Dong, L. Lin, and L. C. Change, "Deep poly-dense network for image superresolution," 2018.
- [85] J. Cai, S. Gu, R. Timofte, L. Zhang *et al.*, "NTIRE 2019 challenge on real image super-resolution: Methods and results," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2019*, 2019, pp. 2211–2223.
- [86] X. Liu, Y. Ding, D. He, C. Li, Y. Fu, and S. Wen, "Cascaded u-shaped deep super-resolution," 2018.
- [87] A. Lugmayr, N. H. Joon, Y. S. Won, G. Kim *et al.*, "AIM 2019 challenge on real-world image super-resolution: Methods and results," in *ICCV Workshop*, 2019, pp. 3575–3583.
- [88] Z. Pan, B. Li, T. Xi, Y. Fan, G. Zhang, J. Liu, J. Han, and E. Ding, "Real image super resolution via heterogeneous model ensemble using GP-NAS," in *Computer Vision - ECCV 2020 Workshops*, ser. Lecture Notes in Computer Science, A. Bartoli and A. Fusiello, Eds., vol. 12537. Springer, 2020, pp. 423–436.
- [89] 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*, 2018, pp. 586–595.



Anran Liu received the B.Eng. degree in 2019 from Tsinghua University, Beijing, China. She is now pursuing a Ph.D. degree at the Department of Computer Science, the University of Hong Kong. Her research interests include computer vision and deep learning, particularly focusing on image processing and super-resolution.



Yihao Liu received the B.S. degree from University of Chinese Academy of Sciences, Beijing, in 2018. He is now working towards the Ph.D. degree in Multimedia Laboratory, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences. He is supervised by Prof. Yu Qiao and Prof. Chao Dong. His research interests include computer vision and image/video enhancement.



Jinjin Gu received his B.Eng. degree in computer science and engineering from the Chinese University of Hong Kong, Shenzhen, in 2020. He is currently pursuing a Ph.D. degree in Engineering and IT with the University of Sydney. His research interests include computer vision, image processing, interpretability of deep learning algorithms and the applications of machine learning in industrial.



Yu Qiao (Senior Member, IEEE) is currently a Professor with the Shenzhen Institutes of Advanced Technology (SIAT), Chinese Academy of Science, and the Director of the Institute of Advanced Computing and Digital Engineering. He has published more than 180 articles in international journals and conferences, including T-PAMI, IJCV, T-IP, T-SP, CVPR, and ICCV. His research interests include computer vision, deep learning, and bioinformation. He received the First Prize of the Guangdong Technological Invention Award, and the Jiayi Lv Young Researcher Award from the Chinese Academy of Sciences. His group achieved the first runner-up at the ImageNet Large Scale Visual Recognition Challenge 2015 in scene recognition, and the Winner at the ActivityNet Large Scale Activity Recognition Challenge 2016 in video classification. He has served as the Program Chair of the IEEE ICIST 2014.



Chao Dong is currently an associate professor in Shenzhen Institute of Advanced Technology, Chinese Academy of Science. He received his Ph.D. degree from The Chinese University of Hong Kong in 2016. In 2014, he first introduced deep learning method – SRCNN into the super-resolution field. This seminal work was chosen as one of the top ten “Most Popular Articles” of TPAMI in 2016. His team has won several championships in international challenges –NTIRE2018, PIRM2018, NTIRE2019,

NTIRE2020 and AIM2020. He worked in SenseTime from 2016 to 2018, as the team leader of Super-Resolution Group. His current research interest focuses on low-level vision problems, such as image/video super-resolution, denoising and enhancement.