Xin Li (ID), Weisheng Dong (ID), Jinjian Wu (ID),
Leida Li (ID), and Guangming Shi (ID)

# Superresolution Image Reconstruction

*Selective milestones and open problems*



©SHUTTERSTOCK.COM/TRIFF

In multidimensional signal processing, such as image and video processing, superresolution (SR) imaging is a classical problem. Over the past 25 years, academia and industry have been interested in reconstructing high-resolution (HR) images from their low-resolution (LR) counterparts. We review the development of SR technology in this tutorial article based on the evolution of key insights associated with the prior knowledge or regularization method from analytical representations to data-driven deep models. The coevolution of SR with other technical fields, such as autoregressive modeling, sparse coding, and deep learning, will be highlighted in both model-based and learning-based approaches. Model-based SR includes geometry-driven, sparsity-based, and gradient-profile priors; learning-based SR covers three types of neural network (NN) architectures, namely residual networks (ResNet), generative adversarial networks (GANs), and pretrained models (PTMs). Both model-based and learning-based SR are united by highlighting their limitations from the perspective of model-data mismatch. Our new perspective allows us to maintain a healthy skepticism about current practice and advocate for a hybrid approach that combines the strengths of model-based and learning-based SR. We will also discuss several open challenges, including arbitrary-ratio, reference-based, and domain-specific SR.

## Introduction

In image processing, *SR* refers to techniques that increase image resolution. The use of SR imaging can be implemented on a hardware basis (e.g., optical solutions) or on a software basis (e.g., digital zooming or image scaling). Software-based (as well as computational) SR imaging approaches can be classified in several ways according to the assumptions about the relationship between LR images and HR images: single image versus multiframe, nonblind versus blind, fixed versus arbitrary scaling ratios, etc. In the past quarter-century, SR techniques have evolved into two categories: model based (1998 to present) [2], [7], [12], [20], [31], [37]

and learning based (2014 to present) [4], [6], [14], [18], [19], [22], [27], [32], [33], [39]. Model-based approaches rely on mathematical models to connect LR and HR data; the main difference is in how the LR observation and HR image prior are characterized. Learning nonlinear mapping between LR and HR image data can be greatly facilitated by the simple idea of skip connections (i.e., ResNet) in learning-based approaches. Recently, researchers have focused on developing novel network architectures [e.g., Generative Latent Bank (GLEAN) [3] and nonlocal blocks [25], [34]] and applying them to realistic scenarios (e.g., locally discriminative learning (LDL) [21]).

These four perspectives can be used to justify the importance of studying SR. SR imaging has a wide range of applications, ranging from nanometer- to light-year scale (for example, SR microscopy won the Nobel Prize in Chemistry in 2014). Watson and Crick's discovery about DNA's double-helix structure could become trivial if SR microscopy technology reveals DNA's detailed structure on a nanometer scale. In terms of technology, SR imaging shows how expensive hardware (i.e., optical zoom) can be traded for more cost-effective software (i.e., SR algorithms). Single-lens reflex cameras are phasing out as SR technology advances, resulting in smartphone photography. SR imaging has also been applied to a variety of engineering systems, including Mars Curiosity and NASA's James Webb Space Telescope. Last but not least, SR image reconstruction is a class of inverse problems that have been extensively studied by mathematicians. SR image reconstruction solutions often have profound implications for inverse problems, such as blind image deconvolution and medical image reconstruction.

There are two main motivations behind this tutorial article. Instead of mathematically approximating LR and HR images with nonlinear mappings $f : X_{LR} \rightarrow X_{HR}$, SR has evolved to data-driven or learning-based methods of determining surrogate models. During the past seven years, extensive research has been conducted along the following two lines. First, skip connections and squeeze and excitation modules have been introduced into ResNet-like NN architectures to alleviate the vanishing gradient problem. Second, model-based approaches can be leveraged to provide new insights, such as the importance of exploiting higher order attention and nonlocal dependency [4]. Model-based SR can also be unfolded directly into deep NNs (DNNs) [27]. In contrast, learning-based SR has coevolved with other fields in computer vision and machine learning. Using a discriminative model, SRGAN intelligently separates the truth of the ground (real HR) from the result of the SR reconstruction (fake HR) as a result of the invention of the GAN. The attention mechanism has sparked interest in transformer-based models, which have been successfully applied to SR (e.g., [8]). Recent advances in blind image restoration have renewed interest in solving the blind real-world SR problem with an LDL approach [21]. By simultaneously estimating the blur kernel and HR image, a network is unfolded to solve the joint optimization problem.

A systematic review of SR's evolution over the last 25 years is presented in this tutorial. The purpose of this article is not to provide a comprehensive review of image SR; interested readers are referred to three recent survey articles [1], [23], [36]. We aim to highlight the rich connections between image processing and other technical fields rather than focusing on a wide range of topics. SR has evolved with Wiener filtering, compressed sensing, and NN design since 1998 (the 50th anniversary of the IEEE Signal Processing Society). SR is a class of inverse problems extensively studied in the literature from a mathematical perspective. As a scientific concept, SR is related to the Rayleigh criterion, a limit for diffraction in optical imaging systems. Engineering applications of SR range from biomedical imaging to consumer electronics. Smartphones, high-definition television (HDTV), remote sensing, and smart health are examples of SR technology in our daily lives. Following are the key new insights offered in this tutorial in addition to the scientific challenges and key milestones of SR:

> **SR imaging has also been applied to a variety of engineering systems, including Mars Curiosity and NASA's James Webb Space Telescope.**

- The first is a selective review of SR milestones in the past 25 years with an emphasis on theoretical insights: i.e., how can missing high-frequency information be approximated or recovered?
- The second is a healthy skepticism toward well-cited SR algorithms. To illustrate progress in both model-based and learning-based approaches, we will highlight failure examples.
- Three open challenges have been selected in the field of SR image reconstruction: arbitrary-ratio, reference-based, and domain-specific SR. We will discuss the current state of the art and future directions for each challenge.

## Problem formulation

### Observation model

Generally speaking, the problem of single-image SR (SISR) refers to the reconstruction of an HR image from its corresponding LR observation [refer to Figure 1(a)]. For a layperson, SISR is widely known as *digital zoom*, which is in contrast to optical zoom. Digital zoom and optical zoom represent software- and hardware-based approaches to enhance the resolution of digital images; the latter is often conceived as the upper bound for the former when optical imaging systems operate within the diffraction limit. From a computational imaging perspective, SISR or digital zoom represents a cost-effective approximation of optical zoom. Closing the gap between software-based and hardware-based approaches has been the holy grail of SR technology in the past 25 years.

We note that the SISR problem formulation has made several simplified assumptions to make it technically more tractable. Depending on the assumption, with the LR observation model, we can formulate the SISR into image interpolation where LR $Y$ is simply a decimated version of HR $X$, as shown in Figure 1(b), or SR image reconstruction where LR is obtained from HR by several operators (e.g., warping, blur, and downsampling), as shown in Figure 1(c). When the degradation is unknown (i.e., the so-called blind or real-world scenario), the problem of SISR is more challenging than its nonblind formulation (i.e., with complete knowledge of the LR observation model). Blind or real-world SISR [23] is one of the frontiers of SR research these days.

In the framework of Bayesian inference, a maximum a posteriori (MAP) estimation of an HR image $X$ from its version of the LR observation $Y$ can be formulated as $\arg\max P(X \mid Y) \approx \arg\max P(Y \mid X) P(X)$ using the Bayesian formula. The LR observation model deals with the likelihood term $P(Y \mid X)$ that characterizes the degradation process of the LR image acquisition. For example, one might start with a parametric observation model $Y = DHX + n$, where $D$ and $H$ denote downsampling/blurring operators, respectively, and $n$ is additive noise. Note that the image interpolation problem is a special case with $H, n$ being skipped; spatially invariant blur $H$ is already an oversimplified abstraction of image degradation in the real world. In the meantime, the source of additive noise $n$ can be sensor related (e.g., shot noise in a raw color filter array) or transmission related (e.g., image compression artifacts). In the formulation of blind problems, the blurring kernel $H$ is unknown and even spatially varying; therefore, we have to address the problem of estimating the blurring kernel and reconstructing the HR image simultaneously.

> From a computational imaging perspective, SISR or digital zoom represents a cost-effective approximation of optical zoom.

### Image prior

The key technical challenge of SISR lies in the construction of an image prior $P(X)$ (as well as the regularization functional in the literature of image restoration and inverse problems) [28]. During the past 25 years, the effort devoted to image prior construction can be classified into two paradigms: model based (1998 to present) and learning based (2014 to present). In the paradigm of model-based SR, the unifying theme is to construct mathematical models (e.g., geometry driven, sparsity based, or gradient domain) for the class of HR images. In the paradigm of learning-based SR, the common objective is to learn a nonlinear mapping [e.g., NNs consisting of several building blocks such as convolution and max-pooling layers, rectified linear unit (ReLU), and batch normalization modules] from the space of LR images to that of HR images. The paradigm shift from model based to learning based is catalyzed by rapid advances in data science (e.g., the large-scale collection of training data such as ImageNet) and deep learning (i.e., the replacement of Moore's law for CPU acceleration by Huang's law for GPU acceleration). (Huang's law is an observation in computer science and engineering that advancements in GPUs are growing at a rate much faster than with traditional CPUs.)

Image prior/regularizer construction or learning represents the state of the art in developing analytical or numerical representations to explain intensity distributions in images,
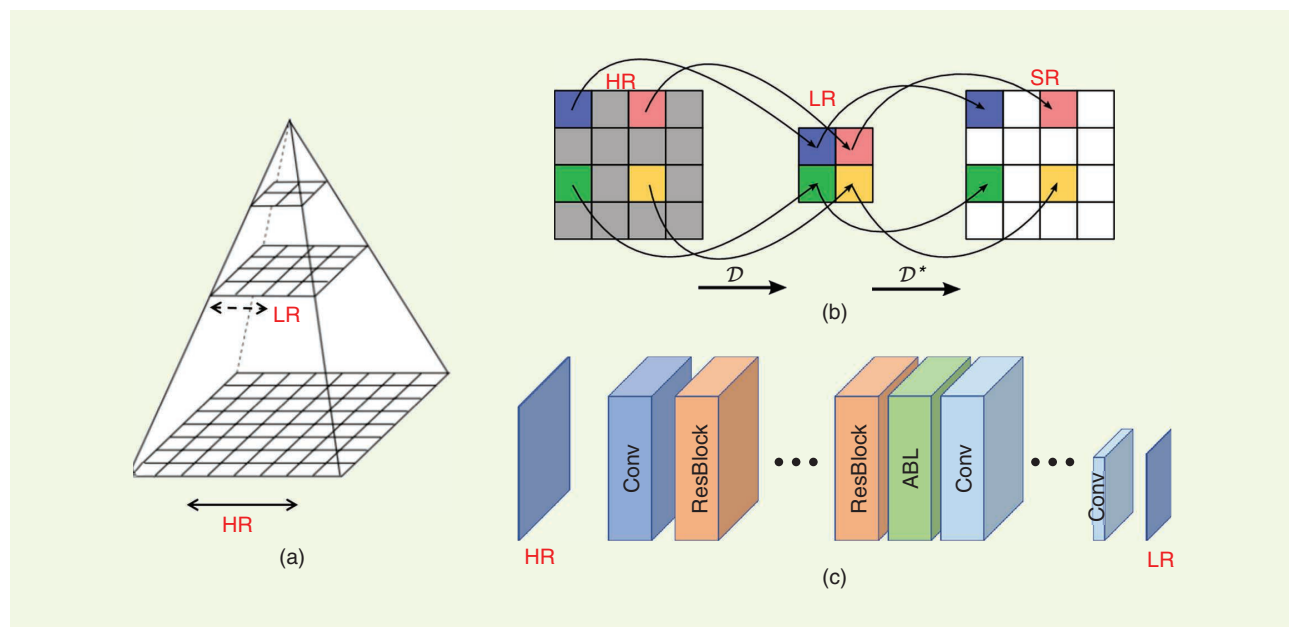


**FIGURE 1.** The problem formulation of SISR. (a) The abstract relationship between LR and HR from the pinhole imaging model. (b) A simplified SISR formulation (model-based image interpolation), where LR is a decimated version of HR. (c) Degradation modeling for more accurate characterization of LR observation from HR (learning-based image SR).

regardless of the model- or learning-based paradigm. Wavelet, partial differential equations, Markov random fields, and NNs serve only as tools to communicate ideas abstracted from the physical properties of visual information. Such abstraction, generative or discriminative, allows us to handle a wide range of image data regardless of their semantic contents (e.g., biometrics or surveillance), acquisition conditions (e.g., camera distance and illumination), or physical origins (e.g., optical sensors versus MRI scanners).

## Model-based SR: From edge directed to sparsity based

In this section, we review model-based SR based on geometry-driven, sparsity-based, and gradient-profile priors that were developed during the first decade of the new millennium. They are constructed from varying insights about the prior knowledge of unknown HR images.

### Adaptive image interpolation via geometric invariance

In the simplified image interpolation situation, LR pixels correspond directly to the decimated version of HR, as shown in Figure 2. For a scaling factor of two, the task of image interpolation boils down to guessing the missing pixels that occupy three-quarters of the spatial locations. The new edge-directed interpolation (NEDI) [20] extends the classic Wiener filtering [mathematically equivalent to least-square (LS) estimation] from prediction to interpolation. As shown in Figure 2, missing pixels as unknown HR sampling locations are denoted by yellow dots. Each yellow pixel (labeled "0") must be predicted from the linear combination of its four surrounding black pixels (labeled as "1" to "4"). Wiener filtering or LS-based estimation of weighting coefficients requires the calculation of local covariance at the HR (marked by solid lines with different colors), which is infeasible due to the missing yellow pixels. Based on the observation that edge orientation is scale invariant, NEDI calculates the local covariances at the LR (marked by dashed lines with different colors) and uses them as the surrogate covariance to drive the derivation of LS-based estimation at the HR.

The effectiveness of NEDI can be interpreted from the following perspectives. First, local geometric information on the direction of the edge can be viewed as being implicitly embedded in the four linear weights in the LS formula. Such an implicit exploitation of the geometry-related prior (i.e., the scale-invariant property of edge orientation) makes the NEDI model a good fit for an arbitrarily oriented edge. Second, there is an elegant duality between step 1 and step 2 of NEDI implementation—they are geometrically isomorphic (up to a rotation by 45° clockwise). Note that the pixels interpolated from step 1 will be treated the same as the given LR (i.e., yellow pixels in step 1 become black ones in step 2). Such geometric duality demonstrates the potential of quincunx sampling as an improved strategy to hierarchically organize visual information compared to conventional raster sampling. (Quincunx

is a geometric pattern consisting of five points arranged in a cross, with four of them forming a square or rectangle and a fifth at its center.)

The limitations of NEDI are summarized next. First, NEDI is a localized model that ignores the nonlocal dependency within the natural images. Second, the geometry-related prior exploited by NEDI matches only a certain class of image structures. For example, edge-directed insight is not applicable to irregular texture patterns whose local statistics are more sophisticated and violate the scale-invariant assumption. Third, the two-step implementation of NEDI is open loop, ignoring the issue of possible inconsistency between adjacent windows. A closed-loop optimization of LS-based autoregressive models was later studied in the literature (e.g., [38]).

### Image SR via sparse coding

The birth of compressed sensing theory around 2006 has inspired many novel applications of sparse representations, including SISR. A key observation obtained from the theory of sparse coding or compressed sensing is that image patches can be decomposed into a sparse linear combination of elements from an overcomplete dictionary. Such observations suggest that the sparse representation can be faithfully recovered from the downsampled signals under mild conditions (the theoretical foundation for SR image reconstruction). In [37], a sparse coding-based approach to SR is developed by jointly training two dictionaries for the LR and HR image patches. Unlike geometry-driven NEDI, the new insight is to enforce the similarity of sparse representations between the LR and HR image patch pairs with respect to their own dictionaries. Along this line of reasoning, the
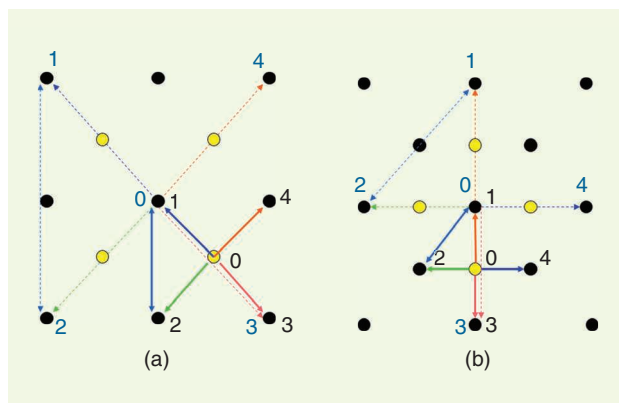
> **Blind or real-world SISR is one of the frontiers of SR research these days.**

**FIGURE 2.** A geometric invariance property exploited by model-based SR, such as NEDI [20]. On the basis of the observation that edge orientation is scale invariant, we can replace the fine-scale correlation (marked by solid lines) with their coarse-scale counterparts (marked by dashed lines). In other words, a multiscale extension of classic Wiener filtering was at the core of NEDI to adapt the interpolation based on the local covariance estimates (implicitly conveying the information about local edge direction). Note that the correspondence between the LR and HR pixel pairs is marked by different colors, and step 2 is isomorphic to step 1 (up to the rotation of 45° clockwise). (a) Step 1 of NEDI. (b) Step 2 of NEDI.

sparse representation of an LR image patch can be used as a surrogate model for the HR image patch dictionary to generate an HR image patch. As shown in Figure 3, the learned dictionary pair is a more compact representation of the image patch pairs, substantially reducing the computational cost. Further development along this line of research includes an adaptive selection of the sparse domain and nonlocal extension, as presented in [7].

The performance of SR via sparse representations is tightly coupled with the quality of the training dataset used for dictionary learning. The selection of the patch size and the optimal dictionary size for SR image reconstruction remain open issues to address. For example, a special dictionary was learned for face hallucination in [37]; can we generalize such a result to other specific application domains? Algorithm 1 [Figure 3(b)] uses an initial SR $X_0$ as the stepping stone; can a related reference image refine such an estimate? Furthermore, the observation model in problem formulation assumes a fixed scaling factor. A different dictionary needs to be trained for a different scaling factor. These weaknesses will be addressed in the three open problems later.

> **The birth of compressed sensing theory around 2006 has inspired many novel applications of sparse representations, including SISR.**

### Image SR via gradient profile prior

Gradient-domain image processing, also known as *Poisson image editing*, deals with image gradients rather than original intensity values. The mathematical foundation of gradient-domain image processing is the numerical solution to the Poisson equation. Conceptually, the horizontal and vertical gradient fields can be viewed as a redundant representation of the original image (each pixel is associated with a pair of gradients). Image reconstruction from gradient profiles can be interpreted as a nontrivial back-projection operation from the gradient space to the image space. In the context of gradient-domain image processing, we can address the problem of SISR by prioritizing gradient profiles instead of intensity values, as shown in Figure 4.

The key observation behind the gradient profile prior (GPP) is that the sharpness of natural images can be characterized by a parametric model such as a generalized exponential distribution. To impose such an image prior, it is possible to design a gradient-field transformation, as shown on the right of Figure 4. The role of the gradient transform is to match the distribution of gradient fields between the target and the observed images. The transformed gradient field is then used to reconstruct the enhanced images. In this way, the objective of the SR image reconstruction is achieved in the gradient domain. Similar to other geometry-driven priors (e.g., total-variation models), the performance of the GPP often degrades for the class of texture images.

### Learning-based SR: Evolution of NN architectures

A rise in deep learning can be seen in 2015. SR via convolutional NN (SRCNN) [5] represented a pioneering work in deep learning-based SR, as shown in Figure 5. Since then, there has been an explosion of literature related to learning-based SR. Due to space limitations, we have to selectively review the most representative work from the perspective of the evolution of network architectures.



**Algorithm 1 (SR Via Sparse Representation).**

1: **Input:** training dictionaries $D_h$ and $D_l$, a low-resolution image $Y$.
2: **For** each $3 \times 3$ patch $y$ of $Y$, taken starting from the upper-left corner with 1 pixel overlap in each direction,
  • Compute the mean pixel value $m$ of patch $y$.
  • Solve the optimization problem with $\tilde{D}$ and $\tilde{y}$. defined in (8): $\min_{\alpha} || \tilde{D}\alpha - \tilde{y} ||_2^2 + \lambda || \alpha ||_1$.
  • Generate the high-resolution patch $x = D_h \alpha^*$. Put the patch $x + m$ into a high-resolution image $X_0$.
3: **End**
4: Using gradient descent, find the closest image to $X_0$ which satisfies the reconstruction constraint

$$X^* = \arg \min_{X} || SHX - Y ||_2^2 + c || X - X_0 ||_2^2.$$
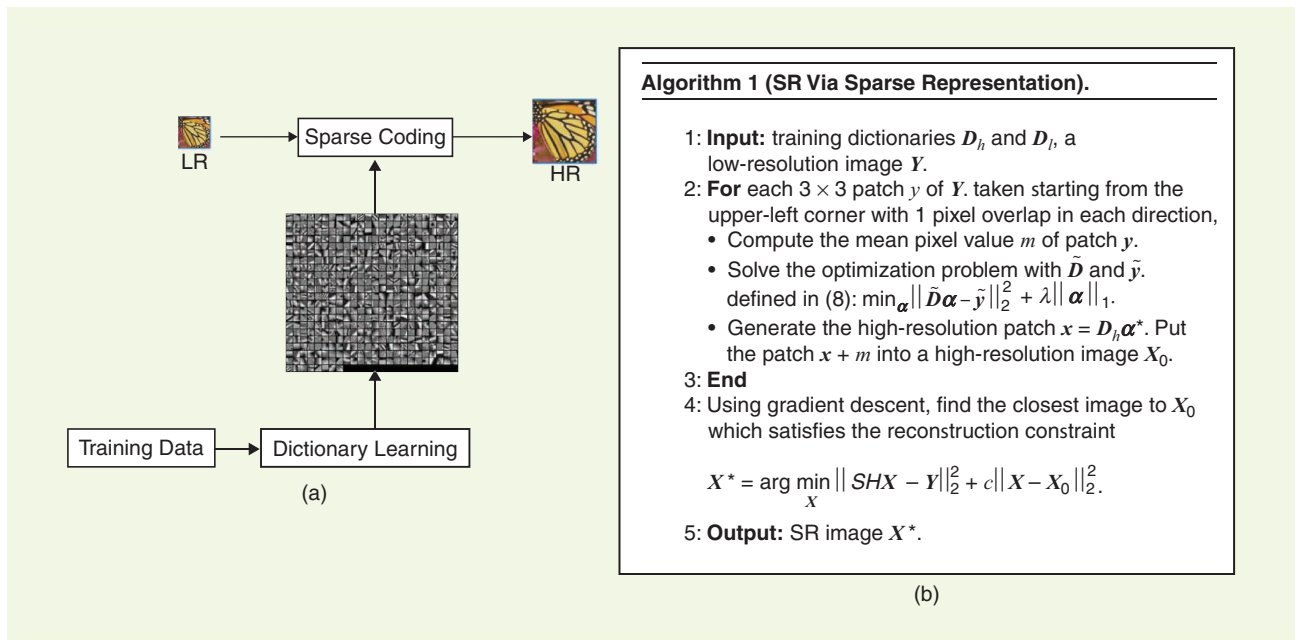
5: **Output:** SR image $X^*$.

**FIGURE 3.** (a) and (b) SISR via sparse representation [37]. The key idea is to enforce the similarity of sparse representations between the LR and HR image patches with respect to their own dictionaries.

## SR image reconstruction via residue refinement

The first category is inspired by the celebrated ResNet architecture. The simple yet elegant idea behind ResNet is to alleviate the notorious vanishing gradient problem via skip connections (mathematically equivalent to predictive coding). This idea is naturally consistent with the objective of SR image reconstruction because missing high-frequency information can be interpreted as residue signals, the target of nonlinear mapping, as shown in Figure 5(a). If we make an analogy between traffic flow (from source to destination) and information flow (from input to output), the construction of the network architecture for SISR shares an objective similar to that of the transportation network. The common objective is to maximize the flow capacity of a transportation network or the amount of residue information in an image reconstruction network.

Many well-cited papers have been published under the framework mentioned previously, as shown in Figure 5(b). Early work such as the deep recursive convolutional network (DRCN) [17], the deep recursive residual network (DRRN) [32], the enhanced deep SR network (EDSR) [22], and the Laplacian pyramid SR network (LapSRN) [18] focused on the construction of network architectures to facilitate the prediction of high-frequency residuals (e.g., via recurrent layers [17], [32] and multiscale decomposition [18], [22]). This line of research was further enhanced by the introduction of the squeeze and excitation module in residual channel attention networks (RCANs) [39] and residual dense networks (RDNs) [40]. Other improvements include considering the error feedback mechanism in deep back-projection networks (DBPNs) [14] and higher order attention mechanisms such as the second-order attention network (SAN) [4].

There are two open questions related to the construction of ResNet-inspired SR networks. First, what is the fundamental limit of this residual refinement strategy? An improved theoretical understanding of what can be learned (i.e., what missing high-frequency information can be recovered?) will offer valuable guidance to the design of a high-order attention mechanism in DNNs. The latest work on iterative refinement with denoising diffusion probabilistic models [15],

[29] contains some promising results. The second is related to the interpretability of the network design. From a practical perspective, a transparent design is expected to help strike an improved tradeoff between cost and performance. In our recent work [27], we have presented a model-guided deep unfolding network (MoG-DUN) implementation, which achieves an improved tradeoff between the performance of the SR reconstruction [measured by the peak signal-to-noise ratio (PSNR) values] and the cost (measured by the number of parameters).

> **SR via convolutional NN represented a pioneering work in deep learning-based SR.**

## Perceptual optimization via adversarial learning

The second category is inspired by the influential GAN architecture. In the pioneering work of SRGAN [19], the objective of perceptual optimization was achieved by introducing an adversarial loss, which pushes the superresolved image closer to the manifold of natural images. In SRGAN, as shown in Figure 6(a), a dedicated discriminator network is trained to differentiate
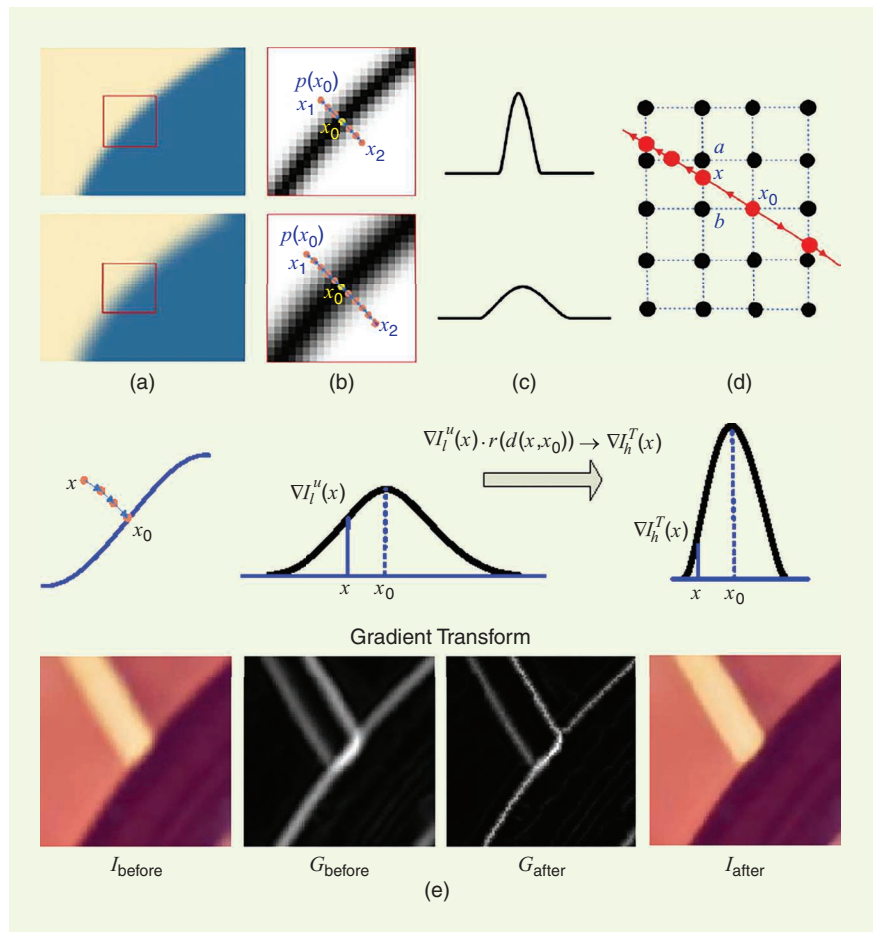


**FIGURE 4.** Image SR via a gradient profile prior (GPP) [31]. (a)–(d) Gradient profile. (a) Two edges with different sharpness; (b) gradient maps; (c) 1D gradient profiles; and (d) tracing the curve of the gradient profile requires subpixel interpolation. (e) Gradient transformation (before and after corresponds to before and after imposing the GPP). By imposing an image prior to gradient transformation, the GPP-based image SR achieves the objective of sharpening edges in the reconstructed HR images.

between superresolved images (the fake sample) and original photorealistic images (the real sample). Note that ideas inspired by ResNet, such as residual-in-residual dense block, have also been incorporated into SRGAN, further improving the performance of adversarial learning. Other ideas, such as relativistic GAN and improved perceptual loss, have also shown impressive performance improvements in enhanced SRGAN (ESRGAN) [35]).

In the context of perceptual optimization, the most controversial issue will be the compromise between improving the details of the image and avoiding the generation of artifacts [21]. Differentiating between texture-like signals and artifact-like noise requires the sophisticated modeling of visual perception by a human vision system. LDL [21] represents the first step toward explicitly discriminating visual artifacts from realistic details. However, LDL requires an HR image as a reference during the discrimination process; how to relax such a requirement (e.g., using a related HR image as a reference) is an interesting topic worthy of further study.

### Large-factor SISR via PTMs

More recently, PTMs such as StyleGAN have been proposed as a latent bank to improve the restoration quality of large-factor image SR (e.g., PULSE [26]). Unlike existing SISR approaches that attempt to generate realistic textures through adversarial learning, Generative Latent Bank (GLEAN) [3] made a significant departure by directly leveraging rich and diverse priors encapsulated in a PTM. GLEAN also differs from the prevalent GAN inversion methods that require expensive image-specific optimization at runtime because it needs only a single forward pass to generate the SR image.
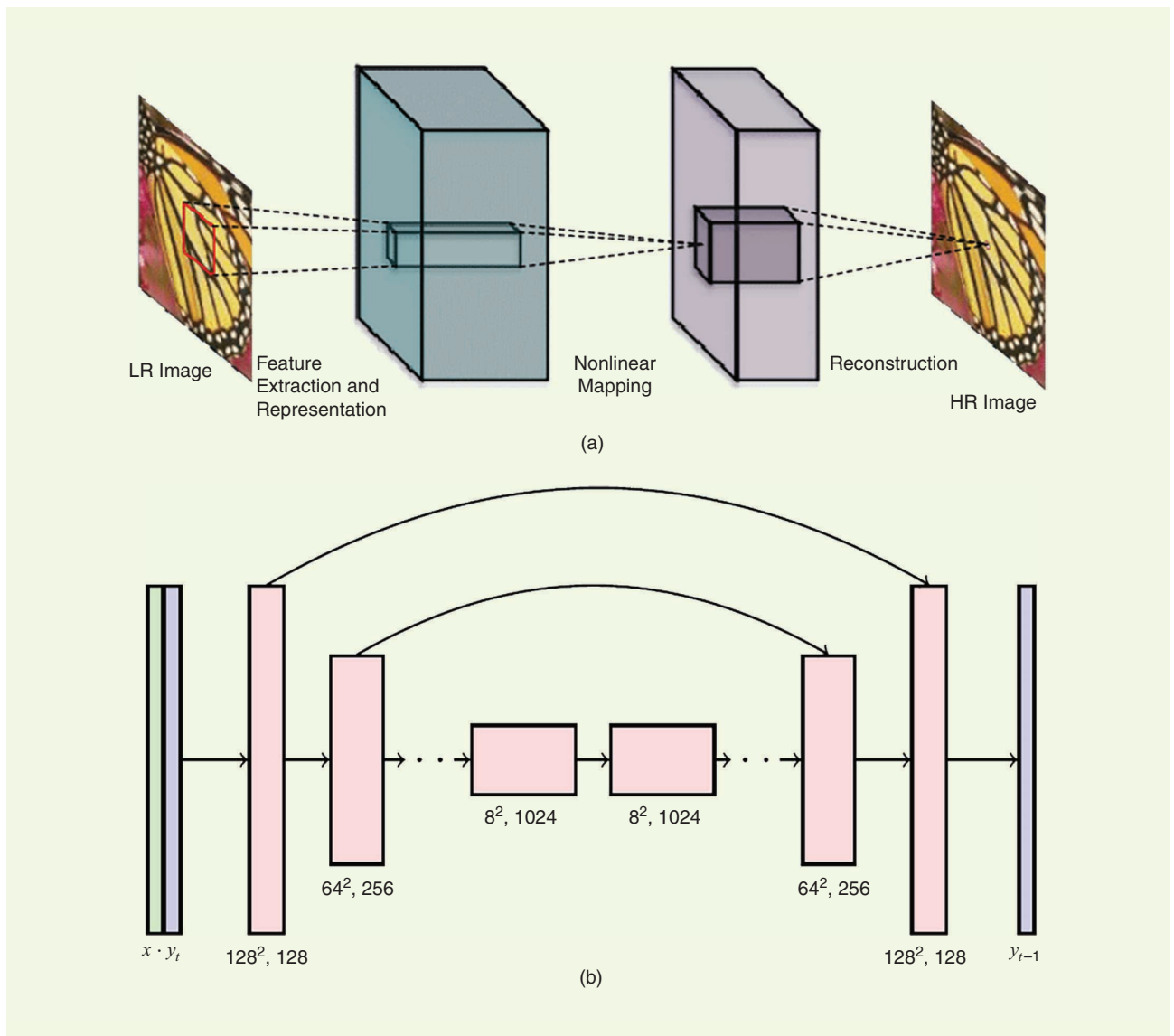


**FIGURE 5.** Learning-based SR. (a) An early attempt, such as SRCNN [5], learns a nonlinear mapping from the space of LR images to HR images. (b) The latest advances achieve SR via U-Net-based iterative residue refinement using denoising diffusion probabilistic models [29].

As shown in Figure 7, GLEAN can be easily incorporated into a simple encoder-bank-decoder architecture with multi-resolution skip connections, making it versatile with images from various categories.

Despite the impressive performance achieved by GLEAN (e.g., as much as 5 dB of PSNR improvement over PULSE on certain classes of images), it still suffers from two fundamental limitations. First, the performance of GLEAN on real-world LR images has remained poor due to a strong assumption with paired training data. In a real-world scenario, SISR is blind because only unpaired LR and HR images are available for training. Degradation learning plays an equally important role in prior learning. How to jointly optimize the interacting components of degradation and prior learning in a Bayesian
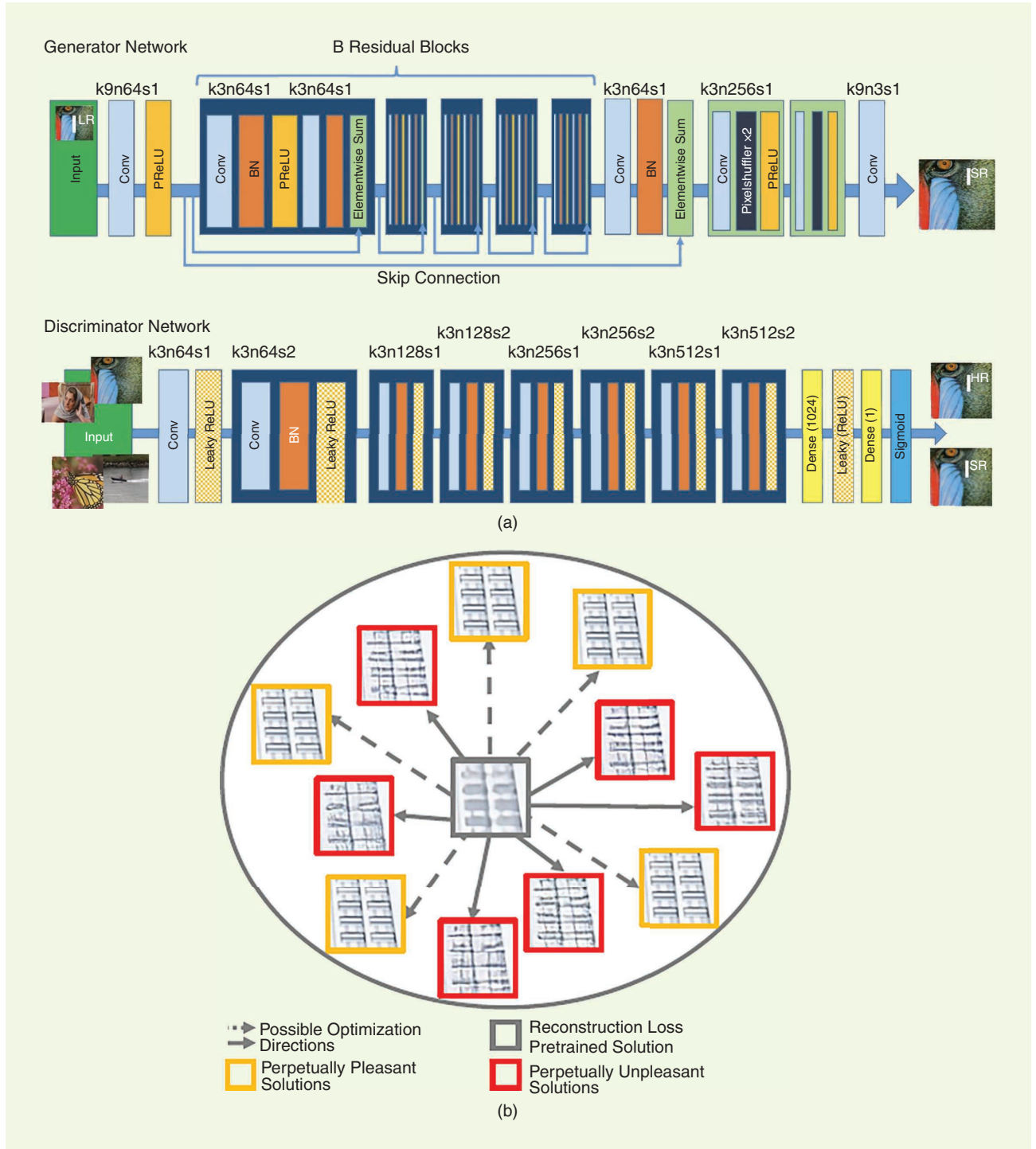


**FIGURE 6.** Adversarial learning-based SR. (a) SRGAN [19] uses an HR image as the reference (real) to distinguish it from a reconstructed SR image (fake). (b) Perceptual optimization of the GAN-based SR model [21]. PReLU: parametric rectified linear unit; BN: batch normalization.

framework is the key to the next milestone in realistic SR. Recent work has reported some initial success along this line of research [23]. Second, the generality of the GLEAN model made it suboptimal for a specific class of images (e.g., human faces). It is possible to design a more powerful and optimized generative prior for face images alone (e.g., generative face prior GFP-GAN).

Most recently, denoising diffusion probabilistic models [15] have been successfully applied to perform SR through a stochastic iterative denoising process in SR3 [29]. The key idea behind SR3 is to iteratively refine the reconstructed HR images by a U-Net architecture trained on denoising at various noise levels and conditioned on the LR input image. SR3 has demonstrated strong SR performance at varying magnification factors and diverse image contents. Additional latest advances are the extension of SISR to blind SR through a joint MAP formulation in KULNet [9] and deep constrained least squares (DCLS) [24]. To estimate the unknown kernel and HR image simultaneously, KULNet introduces uncertainty learning in the latent space to facilitate the estimation of the blur kernel. The joint MAP estimator was unfolded into a deep CNN-based implementation with a learned Laplacian scale mixture prior and the estimated kernel. DCLS reformulates the degradation model so that the deblurring kernel estimation can be transferred into the space of LR images. The reconstructed feature and the LR image feature are jointly fed into a dual-path structured SR network and restore the final HR image.

## Open problems: Arbitrary-ratio SR, reference-based SR, and domain-specific SR

Despite the rapid progress of SR in the last 25 years, there are still many open problems in the field. From the signal processing perspective, we have handpicked the three most significant challenges based on their potential impact in real-world applications. In this section, we will discuss why they are important and what the promising attacks are.

### Arbitrary-ratio SR: Beyond integer factors

Most articles published in the literature on SISR have considered only integer factors (e.g., $\times 2, \times 3, \times 4$). Such integer-factor constraints are simplified situations that make it easier to develop SISR algorithms. In practice, digital zooming often requires noninteger scenarios, e.g., upsampling a $640 \times 480$ image to $1,024 \times 768$ will require a fractional factor of 8/5. Meta-SR [16] is one of few methods that can deal with an arbitrary scaling ratio $r$ (Figure 8). In this method, local projection, weighted prediction, and feature mapping are jointly exploited to implement the noninteger meta-upscale module $r$. Note that such meta-upscale modules with fractional ratios offer an intellectually appealing alternative to integer-factor upscaling, e.g., $2 = 8/7 \times 7/6 \times 6/5 \times 5/4$. Therefore, a particularly promising direction to work with small fractional factors $r > 1$ is the exploitation of local self-similarity (LSS), as advocated in [11].

There are several open questions related to the development of meta-upscale modules. First, the training dataset is obtained by bicubic resampling of popular DIV2K images (available at https://data.vision.ee.ethz.ch/cvl/DIV2K/). It is desirable to collect a more realistic training dataset by varying the focal length of a digital camera. We believe that the ultimate objective of arbitrary-ratio SR is to provide a cost-effective solution to optical zoom. Second, the design of local projection, weighted prediction, and feature mapping can be optimized end to end. For example, if we consider the dual operator of meta-upscale $f(m/n)$, namely meta-downscale $g = f(n/m)$, the concatenation of $f$ and $g$ should become an identity operator [16]. Third, natural images are characterized by LSS, as shown in [11]. Such an LSS is best preserved for factors close to unity $((m/n) \rightarrow 1)$.
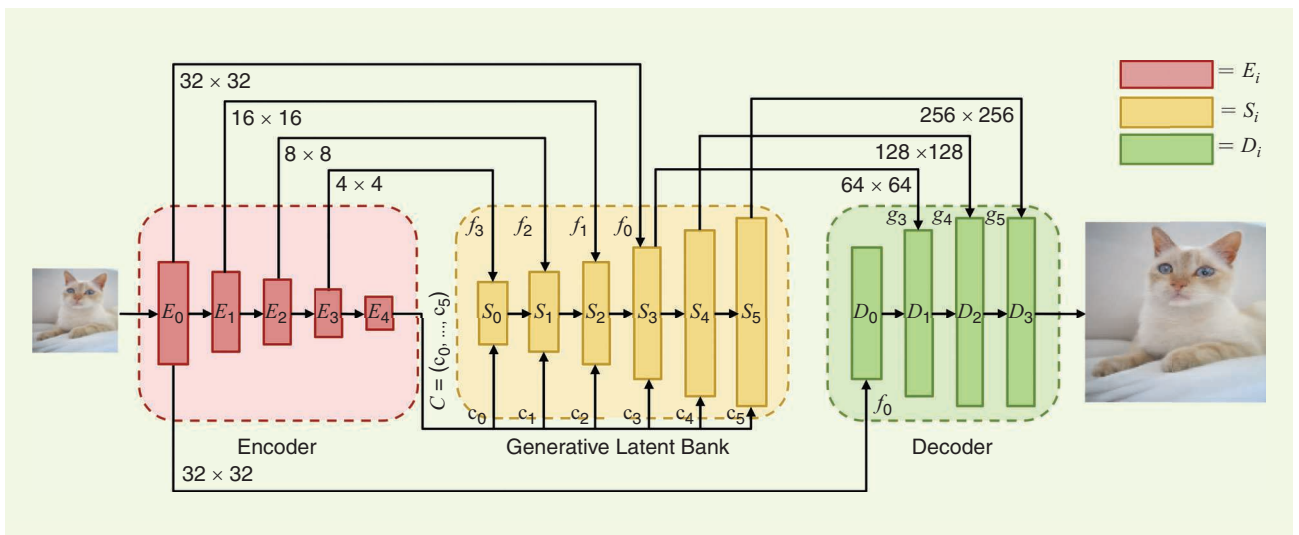


**FIGURE 7.** SISR via pretrained GAN. GLEAN [3] uses the generator as a dictionary and conditions it on the convolutional features provided by the encoder. With a pretrained GAN that captures the natural image prior to processing, GLEAN is capable of achieving a large-scale SISR (a scaling factor of $\times 8$ is shown).

The question of how to exploit LSS using nonlocal NNs [34] is a fascinating topic.

## Reference-based SR via knowledge distillation

Since SISR is an ill-posed inverse problem, it is generally challenging to accurately reconstruct the missing high-frequency details of the unknown HR images from LR observation. A more plausible approach to recover missing high-frequency details is to "borrow" them from a reference HR image with similar content. With additional help from the reference image, this class of reference-based SR (RefSR), as well as guided image SR [42], has the potential to overcome the fundamental limitations of SISR. A different perspective is to view RefSR as a constrained formulation of example-based SR; instead of working with a whole dataset, we aim at utilizing the most relevant reference (containing similar content) to generate rich textures. The key technical challenge is how to pass on the missing high-frequency details from the teacher (reference HR image) to the student (reconstructed SR image).

Similarity Search and Extraction Network (SSEN) [30] represents an example solution to RefSR based on knowledge distillation. As shown in Figure 9, SSEN uses a Siamese network with shared parameters as the backbone of the



**FIGURE 8.** SISR with an arbitrary scaling ratio [16]. Note that in real-world applications, the magnification ratio is often not an integer but some fractional numbers.
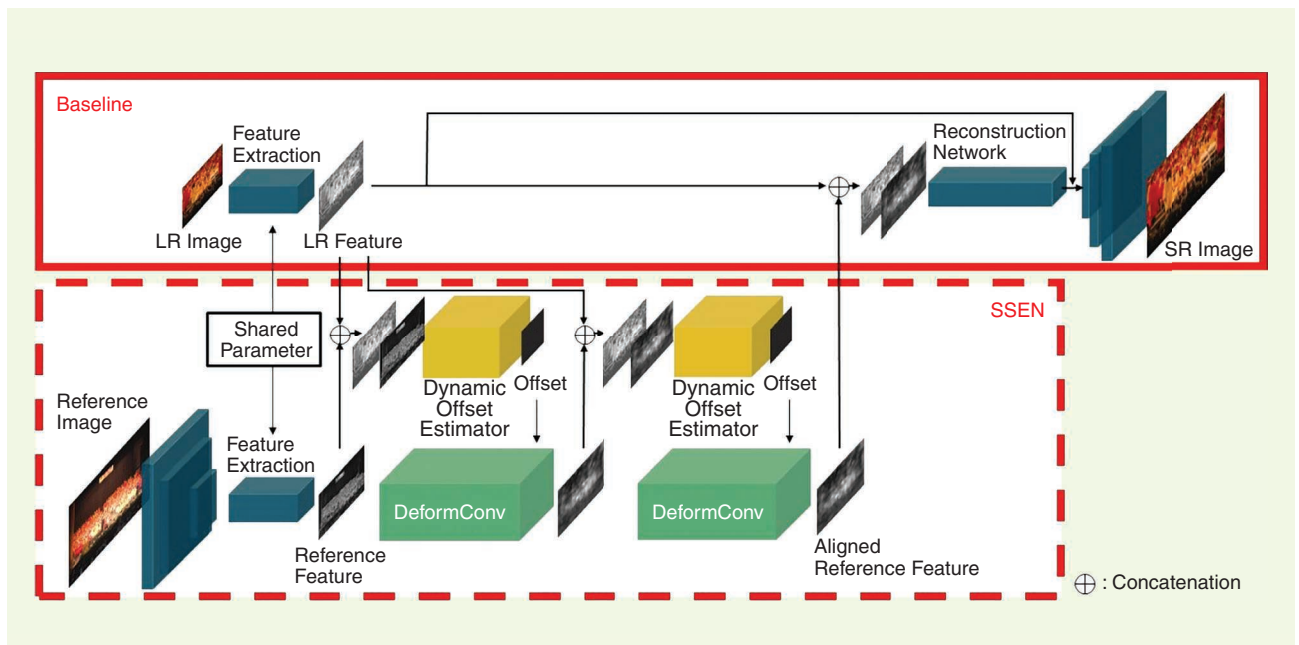


**FIGURE 9.** RefSR. The availability of a reference image provides relevant side information to facilitate the reconstruction of the SR image [30].

teacher-student network. Inspired by the feature alignment capability of deformable convolution, RefSR can be formulated as an integrative reconstruction process of matching similar contents between input and reference features and extracting the reference features (distilled knowledge) in aligned form. Since similar patches can occur at different scales, a multiscale search with a progressively larger receptive field can be achieved by stacking deformable convolution layers. The combination of a multiscale structure with nonlocal blocks makes it convenient to estimate the offsets for deformable convolution kernels. Note that SSEN can also take the input LR image as a self-reference, which is conceptually similar to the idea of self-similarity-based SISR [13].

One of the open challenges in RefSR is the selection of a suitable reference image. Dual-camera zoom in modern smartphone design offers a natural choice in that a pair of images, with different zoomed observations, can be acquired simultaneously. The one with more zoom (telephoto) can serve as a reference for the other with less zoom (short focus). Such a problem formulation of RefSR admits a self-supervised learning-based solution because telephoto, with proper alignment, serves as a self-supervision reference for a digital zoom of short focus. Another closely related extension of RefSR is from image based to video based. With adjacent frames available, video-based SR faces the challenge of fusing relevant information from multiple reference images. How to jointly optimize the interaction between image alignment and SR reconstruction has remained an under-researched topic.

## Domain-specific SR: Connecting domain knowledge with network architecture

The last category for which SR is likely to attract increasing attention is computational imaging in physical and biological sciences. SR imaging is the key to enhancing mankind's vision capability at extreme scales (e.g., nanometers and light years) by breaking the barrier in the physical world. From microscopy to astronomical imaging, domain-specific SR includes a class of customized design challenges where SR imaging has

to be jointly optimized with the imaging modality and for specific applications. The central question is how to incorporate domain knowledge (related to the object of interest or the imaging modality itself) into domain-specific SR algorithms.

SR microscopy is best known for winning the 2014 Nobel Prize in Chemistry. The development of superresolved fluorescence microscopy overcomes the barrier of diffraction limits and brings optical microscopy into the nanodimension. SR microscopy has become an indispensable tool for understanding biological functions at the molecular level in the biomedical research community. One can imagine that the great discovery (the double-helix structure of DNA) made by Watson and Crick indirectly using an XRD image would have been almost straightforward if we could directly observe the DNA structure at nanometer scales. From the signal processing perspective, one of the emerging opportunities is multiframe SR image reconstruction [10]. To break the diffraction limit, one can utilize fluorescent probes that switch between active and inactive states so that only a small optically resolvable fraction of the fluorophores is detected in every snapshot. Such a stochastic excitation strategy ensures that the positions of the active sites can be determined with high precision from the center positions of the fluorescent spots. With multiple snapshots of the sample, each capturing a random subset of the object, a final SR image can be reconstructed from the accumulated positions.

Astronomical imaging is another promising domain on the other scale of physics (distances measured by light years) where SR has great potential in practice. In 2019, for the first time, mankind obtained a photo of a black hole [see Figure 10(a)] that was captured by the Event Horizon Telescope. Due to the far distance, it is not trivial to peek at the supermassive black hole in the M87 galaxy, which is 6.5 billion times larger than our sun. The new launch of the James Webb Telescope has equipped humans with unprecedented capabilities to probe deep space. However, SR techniques, if cleverly combined with optical hardware, can further break the fundamental limit of physical laws (conceptually similar to the microscopic world). Computational imaging techniques such as SR still have much

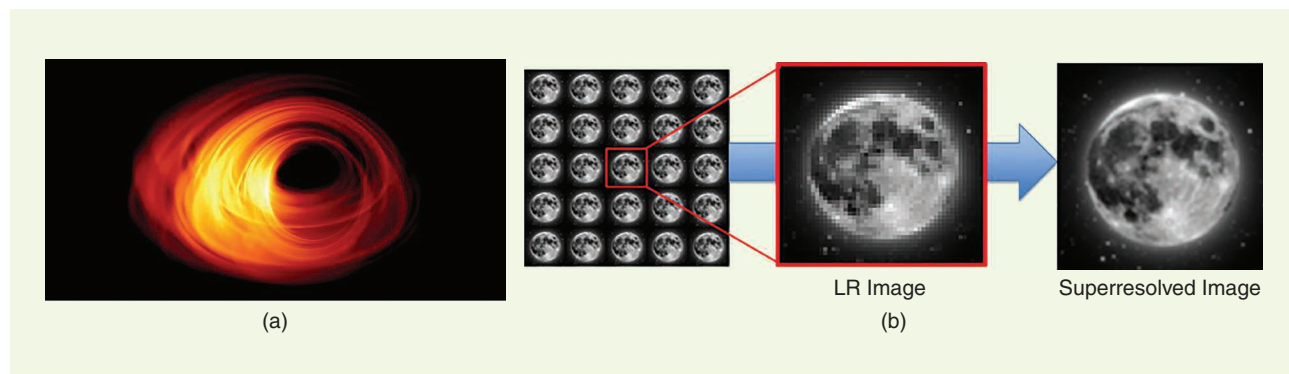> **SR microscopy is best known for winning the 2014 Nobel Prize in Chemistry.**



**FIGURE 10.** Domain-specific SR for astronomical imaging. (a) The first photo of a black hole captured by the Event Horizon Telescope. (b) SR image reconstruction from burst imaging (joint optimization of registration and reconstruction will be needed to suppress undesirable artifacts).

to offer to transform the practice of observing deep space from Earth. Figure 10(b) illustrates an emerging concept called *burst imaging*. By trading space with time, one can improve the spatial resolution of an image by combining the information acquired from multiple timings.

## Conclusion

In this article, we review the evolution of SR technology in the last 25 years from model-based to learning-based SISR. A priori knowledge about HR images, representing the abstraction of 2D data, can be incorporated into the regularization functional in analytical models or loss functions in NNs. Model-based approaches enjoy the benefit of excellent interpretability but suffer from the limitation of a potential mismatch with real-world data. As G. Box once said, "All models are wrong; some of them are useful." On the contrary, learning-based approaches are conceptually closer to the data (but there is still a potential mismatch between training and test data) at the sacrifice of transparency. Perhaps a hybrid approach, combining the strengths of model-based and learning-based paradigms (e.g., [41]), can achieve both good generalization and interpretability.

Looking ahead, what will we see in the next 25 years? Bayesian deep learning provides a new theoretical framework for quantifying various uncertainty factors in deep learning models. By unfolding Bayesian iterative optimization into a DNN-based implementation, we can achieve a principled approach to model and estimate uncertainty for learning-based SR. On the application end, we can foresee that SR technology will reach a higher impact in computational imaging by finding novel applications from the two extreme scales, nanometers and light years. Because most of the information processed by the human brain is visual, we expect that SR imaging will continue to be a key enabling technology in human adventures to the unexplored territories in biological and physical sciences.

## Acknowledgment

## Authors

*Xin Li* received his B.S. degree (highest honors) in electronic engineering and information science from the University of Science and Technology of China, Hefei, China, in 1996 and his Ph.D. degree in electrical engineering from Princeton University, Princeton, NJ, USA, in 2000. Since January 2003, he has been a faculty member with the Lane Department of Computer Science and Electrical Engineering, West Virginia University, Morgantown, WV 26506-6109 USA. He was a member of the technical staff with Sharp Laboratories of America, Camas, WA, USA, from August 2000 to December 2002. He is a Fellow of IEEE.

*Weisheng Dong* received his B.S. degree in electronic engineering from the Huazhong University of Science and Technology, Wuhan, China, in 2004 and his Ph.D. degree in circuits and systems from Xidian University, Xi'an, China, in 2010. In 2010, he joined the School of Electronic Engineering, Xidian University, Xi'an 710071, China, as a lecturer, where he has been a professor since 2016. He was a visiting student at Microsoft Research Asia, Beijing, China, in 2006. From 2009 to 2010, he was a research assistant with the Department of Computing, Hong Kong Polytechnic University, Hong Kong. His research interests include inverse problems in image processing, sparse signal representation, and image compression. He was a recipient of the Best Paper Award at SPIE Visual Communication and Image Processing (VCIP) in 2010. He is currently serving as an associate editor of *IEEE Transactions on Image Processing*. He is a Member of IEEE.

*Jinjian Wu* received his B.Sc. and Ph.D. degrees from Xidian University, Xi'an, China, in 2008 and 2013, respectively. Since July 2015, he has been an associate professor with the School of Electronic Engineering, Xidian University, Xi'an 710071, China. From September 2011 to March 2013, he was a research assistant at Nanyang Technological University, Singapore. From August 2013 to August 2014, he was a postdoctoral research fellow at Nanyang Technological University. From July 2013 to June 2015, he was a lecturer at Xidian University. His research interests include visual perceptual modeling, saliency estimation, quality evaluation, and just noticeable difference estimation. He has served as the special section chair for IEEE Visual Communications and Image Processing (VCIP) 2017 and section chair/organizer/TPC member for ICME2014-2015, PCM2015-2016, ICIP2015, and QoMEX2016. He was awarded the Best Student Paper at ISCAS 2013. He is a Member of IEEE.

*Leida Li* received his B.S. and Ph.D. degrees from Xidian University, Xi'an, China, in 2004 and 2009, respectively. He is currently a professor with the School of Artificial Intelligence, Xidian University, Xi'an 710071, China. From 2014 to 2015, he was a visiting research fellow with the Rapid-Rich Object Search (ROSE) Laboratory, School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, where he was a senior research fellow from 2016 to 2017. His research interests include multimedia quality assessment, affective computing, information hiding, and image forensics. He was the senior program committee member for IJCAI 2019-2020, session chair for ICMR 2019 and PCM 2015, and TPC for AAAI 2019, ACM MM 2019-2020, ACM MM-Asia 2019, ACII 2019, and PCM 2016. He is an associate editor for the *Journal of Visual Communication and Image Representation* and the *EURASIP Journal on Image and Video Processing*. He is a Member of IEEE.

*Guangming Shi* received his B.S. degree in automatic control in 1985, his M.S. degree in computer control, and his Ph.D. degree in electronic information technology, all from Xidian University in 1988 and 2002, respectively. Presently, he is the deputy director of the School of Electronic Engineering, Xidian University, Xi'an 710071, China, and the academic

leader in the subject of circuits and systems. He joined the School of Electronic Engineering, Xidian University, in 1988. Since 2003, he has been a professor in the School of Electronic Engineering at Xidian University, and in 2004, he became the head of the National Instruction Base of Electrician and Electronic (NIBEE). From 1994 to 1996, as a research assistant, he cooperated with the Department of Electronic Engineering at the University of Hong Kong. From June to December 2004, he studied in the Department of Electronic Engineering, University of Illinois at Urbana Champaign.

# References

[1] S. Anwar, S. Khan, and N. Barnes, "A deep journey into super-resolution: A survey," *ACM Comput. Surv.*, vol. 53, no. 3, pp. 1–34, May 2020, doi: 10.1145/3390462.

[2] S. Baker and T. Kanade, "Limits on super-resolution and how to break them," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 9, pp. 1167–1183, Sep. 2002, doi: 10.1109/TPAMI.2002.1033210.

[3] K. C. K. Chan, X. Wang, X. Xu, J. Gu, and C. C. Loy, "GLEAN: Generative latent bank for large-factor image super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 14,245–14,254, doi: 10.1109/CVPR46437.2021.01402.

[4] T. Dai, J. Cai, Y. Zhang, S.-T. Xia, and L. Zhang, "Second-order attention network for single image super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 11,065–11,074, doi: 10.1109/CVPR.2019.01132.

[5] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, Feb. 2016, doi: 10.1109/TPAMI.2015.2439281.

[6] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *Proc. Eur. Conf. Comput. Vis.*, Cham, Switzerland: Springer-Verlag, 2016, pp. 391–407, doi: 10.1007/978-3-319-46475-6_25.

[7] W. Dong, L. Zhang, G. Shi, and X. Wu, "Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization," *IEEE Trans. Image Process.*, vol. 20, no. 7, pp. 1838–1857, Jul. 2011, doi: 10.1109/TIP.2011.2108306.

[8] P. Esser, R. Rombach, and B. Ommer, "Taming transformers for high-resolution image synthesis," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 12,873–12,883, doi: 10.1109/CVPR46437.2021.01268.

[9] Z. Fang, W. Dong, X. Li, J. Wu, L. Li, and G. Shi, "Uncertainty learning in kernel estimation for multi-stage blind image super-resolution," in *Proc. 17th Eur. Conf. Comput. Vis.*, Tel Aviv, Israel: Springer-Verlag, Oct. 2022, pp. 144–161, doi: 10.1007/978-3-031-19797-0_9.

[10] S. Farsiu, M. D. Robinson, M. Elad, and P. Milanfar, "Fast and robust multiframe super resolution," *IEEE Trans. Image Process.*, vol. 13, no. 10, pp. 1327–1344, Oct. 2004, doi: 10.1109/TIP.2004.834669.

[11] G. Freedman and R. Fattal, "Image and video upscaling from local self-examples," *ACM Trans. Graph.*, vol. 30, no. 2, pp. 1–11, Apr. 2011, doi: 10.1145/1944846.1944852.

[12] W. T. Freeman, T. R. Jones, and E. C. Pasztor, "Example-based super-resolution," *IEEE Comput. Graph. Appl.*, vol. 22, no. 2, pp. 56–65, Mar./Apr. 2002, doi: 10.1109/38.988747.

[13] D. Glasner, S. Bagon, and M. Irani, "Super-resolution from a single image," in *Proc. IEEE 12th Int. Conf. Comput. Vis.*, 2009, pp. 349–356, doi: 10.1109/ICCV.2009.5459271.

[14] M. Haris, G. Shakhnarovich, and N. Ukita, "Deep back-projection networks for super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 1664–1673, doi: 10.1109/CVPR.2018.00179.

[15] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *Proc. Adv. Neural Inf. Process. Syst.*, 2020, vol. 33, pp. 6840–6851.

[16] X. Hu, H. Mu, X. Zhang, Z. Wang, T. Tan, and J. Sun, "Meta-SR: A magnification-arbitrary network for super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1575–1584, doi: 10.1109/CVPR.2019.00167.

[17] J. Kim, J. K. Lee, and K. M. Lee, "Deeply-recursive convolutional network for image super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 1637–1645, doi: 10.1109/CVPR.2016.181.

[18] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep Laplacian pyramid networks for fast and accurate super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 624–632, doi: 10.1109/CVPR.2017.618.

[19] C. Ledig et al., "Photo-realistic single image super-resolution using a generative adversarial network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 4681–4690.

[20] X. Li and M. T. Orchard, "New edge-directed interpolation," *IEEE Trans. Image Process.*, vol. 10, no. 10, pp. 1521–1527, Oct. 2001, doi: 10.1109/83.951537.

[21] J. Liang, H. Zeng, and L. Zhang, "Details or artifacts: A locally discriminative learning approach to realistic image super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2022, pp. 5657–5666.

[22] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, 2017, pp. 136–144, doi: 10.1109/CVPRW.2017.151.

[23] A. Liu, Y. Liu, J. Gu, Y. Qiao, and C. Dong, "Blind image super-resolution: A survey and beyond," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 5, pp. 5461–5480, May 2022, doi: 10.1109/TPAMI.2022.3203009.

[24] Z. Luo, H. Huang, L. Yu, Y. Li, H. Fan, and S. Liu, "Deep constrained least squares for blind image super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 17,642–17,652 , doi: 10.1109/CVPR52688.2022.01712.

[25] Y. Mei, Y. Fan, and Y. Zhou, "Image super-resolution with non-local sparse attention," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 3517–3526, doi: 10.1109/CVPR46437.2021.00352.

[26] S. Menon, A. Damian, S. Hu, N. Ravi, and C. Rudin, "PULSE: Self-supervised photo upsampling via latent space exploration of generative models," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 2437–2445, doi: 10.1109/CVPR42600.2020.00251.

[27] Q. Ning, W. Dong, G. Shi, L. Li, and X. Li, "Accurate and lightweight image super-resolution with model-guided deep unfolding network," *IEEE J. Sel. Topics Signal Process.*, vol. 15, no. 2, pp. 240–252, Feb. 2021, doi: 10.1109/JSTSP.2020.3037516.

[28] S. C. Park, M. K. Park, and M. G. Kang, "Super-resolution image reconstruction: A technical overview," *IEEE Signal Process. Mag.*, vol. 20, no. 3, pp. 21–36, May 2003, doi: 10.1109/MSP.2003.1203207.

[29] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, "Image super-resolution via iterative refinement," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 4, pp. 4713–4726, Apr. 2023, doi: 10.1109/TPAMI.2022.3204461.

[30] G. Shim, J. Park, and I. S. Kweon, "Robust reference-based super-resolution with similarity-aware deformable convolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 8425–8434, doi: 10.1109/CVPR42600.2020.00845.

[31] J. Sun, Z. Xu, and H.-Y. Shum, "Image super-resolution using gradient profile prior," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2008, pp. 1–8, doi: 10.1109/CVPR.2008.4587659.

[32] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 3147–3155, doi: 10.1109/CVPR.2017.298.

[33] R. Timofte, V. De Smet, and L. Van Gool, "A+: Adjusted anchored neighborhood regression for fast super-resolution," in *Proc. Asian Conf. Comput. Vis.*, Cham, Switzerland: Springer-Verlag, 2015, pp. 111–126.

[34] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 7794–7803, doi: 10.1109/CVPR.2018.00813.

[35] X. Wang et al., "ESRGAN: Enhanced super-resolution generative adversarial networks," in *Proc. Eur. Conf. Comput. Vis. Workshops*, 2018, pp. 63–79, doi: 10.1007/978-3-030-11021-5_5.

[36] Z. Wang, J. Chen, and S. C. Hoi, "Deep learning for image super-resolution: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 10, pp. 3365–3387, Oct. 2021, doi: 10.1109/TPAMI.2020.2982166.

[37] 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, Nov. 2010, doi: 10.1109/TIP.2010.2050625.

[38] X. Zhang and X. Wu, "Image interpolation by adaptive 2-D autoregressive modeling and soft-decision estimation," *IEEE Trans. Image Process.*, vol. 17, no. 6, pp. 887–896, Jun. 2008, doi: 10.1109/TIP.2008.924279.

[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 *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2018, pp. 286–301.

[40] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2472–2481, doi: 10.1109/CVPR.2018.00262.

[41] Z. Zhang, S. Yu, W. Qin, X. Liang, Y. Xie, and G. Cao, "Self-supervised CT super-resolution with hybrid model," *Comput. Biol. Med.*, vol. 138, Nov. 2021, Art. no. 104775, doi: 10.1016/j.compbiomed.2021.104775.

[42] M. Zhou, K. Yan, J. Pan, W. Ren, Q. Xie, and X. Cao, "Memory-augmented deep unfolding network for guided image super-resolution," *Int. J. Comput. Vis.*, vol. 131, no. 1, pp. 215–242, Jan. 2023, doi: 10.1007/s11263-022-01699-1.

**SP**