

ACDMSR: Accelerated Conditional Diffusion Models for Single Image Super-Resolution

Axi Niu, Pham Xuan Trung, Kang Zhang, Jinqiu Sun*, Yu Zhu, In So Kweon, *Member, IEEE*,
and Yanning Zhang, *Senior Member, IEEE*

Abstract—Diffusion models have gained significant popularity in the field of image-to-image translation. Previous efforts applying diffusion models to image super-resolution (SR) have demonstrated that iteratively refining pure Gaussian noise using a U-Net architecture trained on denoising at various noise levels can yield satisfactory high-resolution images from low-resolution inputs. However, this iterative refinement process comes with the drawback of low inference speed, which strongly limits its applications. To speed up inference and further enhance the performance, our research revisits diffusion models in image super-resolution and proposes a straightforward yet significant diffusion model-based super-resolution method called ACDMSR (accelerated conditional diffusion model for image super-resolution). Specifically, our method adapts the standard diffusion model to perform super-resolution through a deterministic iterative denoising process. Our study also highlights the effectiveness of using a pre-trained SR model to provide the conditional image of the given low-resolution (LR) image to achieve superior high-resolution results. We demonstrate that our method surpasses previous attempts in qualitative and quantitative results through extensive experiments conducted on benchmark datasets such as Set5, Set14, Urban100, BSD100, and Manga109. Moreover, our approach generates more visually realistic counterparts for low-resolution images, emphasizing its effectiveness in practical scenarios.

Index Terms—Diffusion Models, Image-to-Image Translation, Conditional Image Generation, Image Super-resolution.

I. INTRODUCTION

SINGLE IMAGE SUPER-RESOLUTION (SISR) has drawn active attention due to its wide applications in computer vision, such as object recognition, remote sensing and so on [1], [2], [3], [4], [5], [6], [7], [8]. SISR aims to obtain a high-resolution (HR) image containing great details and textures from a low-resolution (LR) image by a super-resolution method, which is a classic ill-posed inverse problem [9], [10]. To establish the mapping between HR and LR images, lots of CNN-based methods have emerged [11], [12], [13], [14], [15],

[16], [17], [18], [19]. These methods focus on designing novel architectures by adopting different network modules, such as residual blocks [20], [21], attention blocks [22], [23], non-local blocks [24], [25], transformer layers [26], [27], and contrastive learning [28], [29], [30]. For optimizing the training process, they prefer to use the MAE or MSE loss (e.g., L_1 or L_2) to optimize the architectures, which often leads to over-smooth results because the above losses provide a straightforward learning objective and optimize for the popular PSNR (peak signal-to-noise-ratio) metric [31], [32], [33], [34], [35].

With deep generative models of all kinds exhibiting high-quality samples in a wide variety of data modalities, approaches based on the deep generative model have become one of the mainstream, mainly including GAN-based methods [36], [37], [38] and flow-based methods [39], [40], [41], which have shown convincing image generation ability. GAN-based SISR methods [36], [37], [38] often introduce a generator and a discriminator in an adversarial way to push the generator to generate realistic images. The generator can generate an SR result for the input LR, and the discriminator aims to distinguish if the generated SR result is true. The training process is optimized by combining content loss and adversarial losses, which have strong learning abilities [38], [42], [43]. While GAN-based methods have an obvious drawback in that they easily fall into mode collapse, the training process is challenging to converge with complex optimization [44], [45], [46]. Furthermore, adversarial losses often introduce artifacts not present in the original clean image, leading to large distortion [47], [35]. Flow-based SR methods are another famous line based on the deep generative model. They directly account for the ill-posed problem with an invertible encoder [48], [49]. The flow-based operation transforms a Gaussian distribution into an HR image space instead of modeling one single output and inherently resolves the pathology of the original "one-to-many" SR problem. Optimized by a negative loglikelihood loss, these methods avoid training instability. Still, they suffer from enormous footprints and high training costs due to the strong architectural constraints to keep the bijection between latents and data [48].

Lately, the broad adoption of diffusion models has shown promising results in image generative tasks [50]. In SRDiff [51], the authors propose a two-stage SR framework. First, they design a super-resolution structure and pre-train it to obtain a conditional image for the diffusion process. Then they redesign the U-net structure in diffusion models. The training process of this method is relatively complicated, and it does not consider combining existing pre-trained SR

This work is funded in part by the Project of the National Natural Science Foundation of China under Grant 61871328, Natural Science Basic Research Program of Shaanxi under Grant 2021JCW-03, as well as the Joint Funds of the National Natural Science Foundation of China under Grant U19B2037.). (* Corresponding author: Jinqiu Sun.)

Axi Niu, Yu Zhu, and Yanning Zhang are with the School of Computer Science, Northwestern Polytechnical University, Xi'an, 710072, China, and also with the National Engineering Laboratory for Integrated Aero-Space-Ground-Ocean Big Data Application Technology, Xi'an, 710072, China (email: nax@mail.nwpu.edu.cn, yuzhu@mail.nwpu.edu.cn, ynzhang@nwpu.edu.cn).

Pham Xuan Trung, Kangzhang, and In So Kweon are with the School of Electrical Engineering, Korea Advanced Institute of Science and Technology. (email: trungpx@kaist.ac.kr, kangzhang@kaist.ac.kr, iskweon77@kaist.ac.kr)

Jinqiu Sun is with the School of Astronautics, Northwestern Polytechnical University, Xi'an 710072, China (email: sunjinqiu@nwpu.edu.cn)

arXiv:2307.00781v1 [cs.CV] 3 Jul 2023

models, such as EDSR [20], RCAN [12], and SwinIR [26]. Similarly, SR3 [48] directly applies the bicubic up-sampled LR image as the conditional image. Nevertheless, the stochastic sampling style in the inference phase makes the reconstruction process complex and slow. Unlike them, we propose a simple but non-trivial method for image super-resolution based on the conditional diffusion model, *i.e.*, ACDMSR (accelerated conditional diffusion model for image super-resolution). Our work shares some similarities with SRDiff, which first applies diffusion models to the SR tasks. Different from the existing technique [51], [48], our ACDMSR adopts the current pertained SR methods to provide the conditional image, which is more plausible than the one in [48], [51]. Moreover, it helps to significantly improve perceptual quality over existing state-of-the-art methods across multiple standard benchmarks. Furthermore, to accelerate the inference steps, we build a n -th order sampler that decreases the 1000-step- to 40-step inference and keeps the good quality. Compared with previous diffusion model-based methods SR3 and SRDiff, which need 1000 inference steps, ours significantly shortened the acquisition of final results. By simply concatenating a Gaussian noise and the conditional image with L_1 loss optimizing the diffusion model, our method makes the training process more concise compared with [48], [51]. The main contributions of this work are listed as follows:

- To the best of our knowledge, we are the first to combine diffusion models and the existing pre-trained SR models to conduct image super-resolution, which can also be taken as a post-process framework.
- Compared with existing diffusion model-based SR methods, our ACDMSR adopts a deterministic sampling way in the inference phase. It can effectively reduce the inference steps from 1000 to just 40, achieving an improved equilibrium between distortion and perceptual quality.
- Compared to existing SOTA SR methods, our ACDMSR achieves superior perceptive results and can generate more photo-realistic SR results on various benchmarks.

II. RELATED WORK

A. Single Image Super-resolution Methods

CNN-based methods. CNN-based methods are a trendy line for image super-resolution, and much great work is coming out. For example, [11] employs the ResNet architecture from [52] and solves the time and memory issues with good performance. Then [20] further optimizes it by analyzing and removing unnecessary modules to simplify the network architecture and produce better results. After them, RCAN [12] and MCAN [13], and EMASRN [53] adopt the attention mechanism [54] and design new residual dense networks. Then MLRN [17], SRNIF [18], and BSRT [55] proposed multi-scale fusion or internal and external features fusion architecture to solve the problem that the existing SISR could not make full use of the characteristic information of the middle network layer and internal features. In addition, SwinIR [26] and ESRT [27] apply transformer technology to improve the performance further. While these methods aim at pursuing higher PSNR (peak signal-to-noise-ratio) by

designing novel architectures and using the MSE or MAE loss (*e.g.*, L_1 or L_2) to optimize the architectures, which often leads to smooth results because the above losses provide a straightforward learning objective [31], [33], [34], [35].

Generative model-based methods.

Because deep generative models have recently exhibited promising results in generating images with rich details, it has become popular to adopt generative models to conduct image super-resolution, such as GAN-based methods [11], [36], [37], [38] and flow-based methods [39], [40], [41]. SRGAN [11] is the first GAN-based SISR method. It adopts the GAN technology to push the generator to produce results with better Visual effects. Compared with SRGAN, ESRAGN [36] trains the discriminator to predict the authenticity of the generated image instead of predicting if the generated image is valid. NatSR [37] proposes a Naturalness Loss based on a pre-trained natural manifold discriminator to improve the ability of the discriminator and achieve comparable results to recent CNNs. However, GAN-based methods have an obvious drawback that is jointly optimizing the whole training process by combining MAE or MSE makes the model easy to fall into mode collapse, and the training process is not easy to converge with complex optimization [44], [45]. Furthermore, adversarial losses often introduce artifacts not present in the original clean image, leading to large distortion [47], [35]. Flow-based SR methods are another famous line based on the deep generative model. They directly account for the ill-posed problem with an invertible encoder [48], [49]. The flow-based operation transforms a Gaussian distribution into an HR image space instead of modeling one single output and inherently resolves the pathology of the original "one-to-many" SR problem. Optimized by a negative loglikelihood loss, these methods avoid training instability. Still, they suffer from enormous footprints and high training costs due to the strong architectural constraints to keep the bijection between latents and data [48].

B. Diffusion Models

Diffusion models have achieved promising results in image generation [50], [56], [57]. It aims to use a Markov chain to transform latent variables in simple distributions (*e.g.*, Gaussian) to data in complex distributions. The core technology for the success of diffusion models is their iterative sampling process. It progressively removes noise from a random noise vector. This iterative refinement procedure repetitively evaluates the diffusion model, allowing for the trade-off of compute for sample quality: by using extra compute for more iterations, a small-sized model can unroll into a larger computational graph and generate higher quality samples [58], [59], [60]. Inspired by the above works, some researchers apply diffusion models in low-level vision tasks [35], [51], [48]. In [35], authors propose a novel framework for blind image deblurring based on conditional diffusion models, which employs a stochastic sampler to refine the output of a deterministic predictor and produces a diverse set of plausible reconstructions for a given input, leading to a significant improvement in perceptual quality over existing state-of-the-art methods.

SRdiff [51] also discuss the drawbacks of current generative models-based SR methods. It designs a novel single-image super-resolution model based on diffusion models, which can provide diverse and realistic super-resolution predictions while avoiding issues with over-smoothing, mode collapse, or large model footprints. At the same time, it has to combine the counterpart output from the pre-trained SR model for the LR input, which makes the whole training process and forward diffusion process very complex and may struggle with images that contain complex textures or patterns. Unlike SRdiff, SR3 [48] presents a straightforward style to introduce diffusion models to help image super-resolution. It just takes the bicubic low-resolution image as the conditional image and uses denoising diffusion probabilistic models to perform stochastic denoising and achieve super-resolution through iterative refinement using a U-Net model trained on denoising at various noise levels, achieving strong performance on super-resolution tasks on faces and natural images, as well as effective cascaded image generation. Though these methods have achieved plausible visual quality, they have an obvious drawback: the sampling speed needs to be improved in the inference time.

III. PERLIMINARIES: OVERVIEW OF DIFFUSION MODELS

In diffusion models, a Markov chain of diffusion steps generates data by progressively perturbing the data with Gaussian noise. Subsequently, these models aim to learn how to reverse the diffusion process and reconstruct desired data samples from the noise. This section begins by revisiting the standard denoising diffusion probabilistic model (DDPM) [50] to provide a basic understanding. A typical probabilistic diffusion model consists of four main components: the forward process, the reverse process, the optimization of the diffusion model, and the inference stage. We will now introduce each of these components in the following sections:

A. Forward process

Suppose we have a real data distribution $\mathbf{x}_0 \sim q(\mathbf{x})$. The forward process gradually adds noise into a sampled image \mathbf{x}_0 using a variance (noise) schedule β_1, \dots, β_T ($\beta_t \in (0, 1), 1 \leq t \leq T$) to generate noised versions $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ from the original image \mathbf{x}_0 . This process can be defined with a Markovian structure:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad 1 \leq t \leq T. \quad (1)$$

By leveraging the properties of the Gaussian distribution and marginalizing the intermediate steps, we can sample \mathbf{x}_t at any given time-step t using the following formulation:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(x_t; \sqrt{\hat{\alpha}_t}\mathbf{x}_0, (1 - \hat{\alpha}_t)\mathbf{I}), \quad (2)$$

where $\alpha_t = 1 - \beta_t$ and $\hat{\alpha}_t = \prod_{s=1}^t \alpha_s$. This formulation allows us to express \mathbf{x}_t using the reparameterization trick:

$$\mathbf{x}_t(\mathbf{x}_0, \epsilon) = \sqrt{\hat{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \hat{\alpha}_t}\epsilon, \quad (3)$$

where ϵ is a Gaussian noise vector with $\epsilon \sim \mathcal{N}(0, \mathbf{I})$.

B. Reverse process

In order to acquire a real sample \mathbf{x}_0 from a Gaussian noise input $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$, the reversal of the preceding forward process is required. This involves the construction of the inverse of Eq. 1 and the iterative reversal using $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$. It is worth highlighting that if β_t is sufficiently small, $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ will also follow a Gaussian distribution. However, estimating $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ presents a challenge as it requires the utilization of the complete dataset. Furthermore, when conditioned on \mathbf{x}_0 it becomes tractable:

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \mu(\mathbf{x}_t, \mathbf{x}_0), \sigma(\mathbf{x}_t, \mathbf{x}_0)), \quad (4)$$

where $\mu(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\hat{\alpha}_{t-1}\beta_t}}{1 - \hat{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\hat{\alpha}_t(1 - \hat{\alpha}_{t-1})}}{1 - \hat{\alpha}_t}\mathbf{x}_t$ and $\sigma(\mathbf{x}_t, \mathbf{x}_0) := \frac{1 - \hat{\alpha}_{t-1}}{1 - \hat{\alpha}_t}\beta_t$. By substitution Eq. 3, $\mathbf{x}_0 = (\mathbf{x}_t - \sqrt{1 - \hat{\alpha}_t}\epsilon)/\sqrt{\hat{\alpha}_t}$, into the $\mu(\mathbf{x}_t, \mathbf{x}_0)$, we can have

$$\mu(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \hat{\alpha}_t}}\epsilon). \quad (5)$$

Following the choice of [50], if we train a model $q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \sum_\theta(\mathbf{x}_t, t)\mathbf{I})$ to learn the above reverse process, $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$, and set $\sum_\theta(\mathbf{x}_t, t)$ as $\sum_\theta(\mathbf{x}_t, t) = \sigma(\mathbf{x}_t, \mathbf{x}_0)$, we can use network f_θ to predict the noise $\epsilon \approx f_\theta(\mathbf{x}_t, t)$ so that the reverse process becomes learnable:

$$q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}\left(\frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \hat{\alpha}_t}}f_\theta(\mathbf{x}_t, t)), \frac{1 - \hat{\alpha}_{t-1}}{1 - \hat{\alpha}_t}\beta_t\mathbf{I}\right). \quad (6)$$

C. Optimize the diffusion model

In [50], it has been demonstrated that reweighted evidence lower bound proves to be an effective loss function in practical applications:

$$L(\theta) = \mathbb{E}_{t, \mathbf{x}, \epsilon} \|f_\theta(\mathbf{x}_t, t) - \epsilon\|^2, \quad (7)$$

where the model learns to predict the added noise ϵ . The pseudocode for the training is shown in the training part of Algorithm 1.

D. Inference

After training, the inference becomes trivial now, since given the start point \mathbf{x}_T , we can get the formulation of next step image \mathbf{x}_{t-1} with the reparametrization trick for Equation 6 as follows:

$$\mathbf{x}_{t-1} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \hat{\alpha}_t}}f_\theta(\mathbf{x}_t, t)) + \sqrt{1 - \hat{\alpha}_t}\epsilon_t, \quad (8)$$

where $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ is the random noise added in each denoise step. We can sample the final image \mathbf{x}_0 by iteratively applying the above equation. The pseudocode for the inference is shown in the inference part of Algorithm 1.

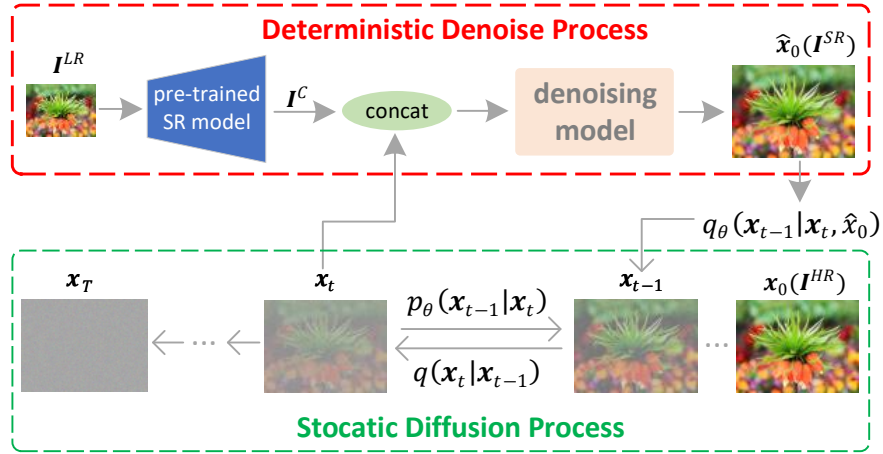


Fig. 1: Illustration of our method. The model contains a stochastic forward diffusion process, gradually adding noise to an I^{HR} image. And a deterministic denoise process is applied to recover high-resolution and realistic images I^{SR} corresponding to I^{LR} images.

Algorithm 1 DDPM

Input: Dataset D , noise predictor f_θ , noise schedule $\hat{\alpha}_t$, total timestep T

Training: train f_θ

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim D$
- 3: $t \sim [1, \dots, T]$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take a gradient descent step on $\|\epsilon - f_\theta(\sqrt{\hat{\alpha}_t}\mathbf{x}_0 + \sqrt{(1-\hat{\alpha}_t)}\epsilon, t)\|^2$
- 7: **until** converged

Inference: sampling \mathbf{x}_0

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: **for** $t = T, \dots, 1$ **do**
 - 3: $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, if $t > 1$, else $\epsilon_t = \mathbf{0}$
 - 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\hat{\alpha}_t}}(x_t - \frac{1-\hat{\alpha}_t}{\sqrt{1-\hat{\alpha}_t}}f_\theta(\mathbf{x}_t, t)) + \beta_t\epsilon_t$
 - 5: **end for.**
-

IV. METHODOLOGY

Our method can be seen as a post-process for single image super-resolution (SISR). As shown in Fig. 1, our ACDMSR consists of a stochastic diffusion process forward procedure that gradually adds noise to an image until a fully normal Gaussian noise and a deterministic denoising reverse process that conditions on I^C to reconstruct the image from noise. Algorithm 2 shows the whole process of our ACDMSR. The following section introduces our method in detail.

A. Stochastic Diffusion Process

Given a SISR dataset $(I^{HR}, I^{LR}) \sim D$, we adopt the diffusion model [50], [57] to map a normal Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{1})$ to a high-resolution image $\mathbf{x}_0 = I^{HR}$ with a corresponding conditional image $\mathbf{x}^C = I^{LR}$. We will talk about the choice of the conditional image later. The diffusion model contains latent variables $\mathbf{x} = \{\mathbf{x}_t | t = 0, 1, \dots, T\}$, where $\mathbf{x}_0 = I^{HR}$, $\mathbf{x}_T = \mathcal{N}(0, \mathbf{1})$. The same noise schedule with [50] is used for our method, β_1, \dots, β_T where $1 \leq t \leq T$.

Forward stochastic diffusion process. We define the forward process $q(\mathbf{x}_t | I^{HR}) := q(\mathbf{x}_t | \mathbf{x}_0)$ of diffusion model with a Gaussian process by the Markovian structure:

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}_{t-1}) &= \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \\ q(\mathbf{x}_t | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\hat{\alpha}_t}\mathbf{x}_0, (1 - \hat{\alpha}_t)\mathbf{I}). \end{aligned} \quad (9)$$

Same with DDPM [50], the forward process gradually adds noise into an image \mathbf{x}_0 to generate latent variables $\mathbf{x}_1, \dots, \mathbf{x}_T$ for the original image \mathbf{x}_0 . With the Gaussian distribution reparameterization trick, we can write the latent variable \mathbf{x}_t as Eq. 3, $\mathbf{x}_t(\mathbf{x}_0, \epsilon) = \sqrt{\hat{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \hat{\alpha}_t}\epsilon$.

Model training. According to [61], our findings demonstrate that predicting the image, rather than focusing on the noise, yields superior outcomes when applied in super-resolution tasks. We have proved it in Sec. V-C. Therefore, the optimization target of our diffusion model is denoising $\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)$ to get estimated $\hat{\mathbf{x}}_0$ with a U-Net $f_\theta(\mathbf{x}_t, t, \mathbf{x}^C) := \hat{\mathbf{x}}_0 \approx \mathbf{x}_0$. We use the following loss function to train the model:

$$L := \mathbb{E}_{t, (\mathbf{x}_0, I^C), \epsilon} [\|\mathbf{x}_0 - f_\theta(\hat{\alpha}_t\mathbf{x}_0 + \sigma_t\epsilon, t, \mathbf{x}^C)\|^2], \quad (10)$$

where t is uniformly sampled between 1 and T . With Eq. 3, $\epsilon = (\mathbf{x}_t - \sqrt{\hat{\alpha}_t}\hat{\mathbf{x}}_0)/\sqrt{1 - \hat{\alpha}_t}$, we can easily predict the added noise to the image \mathbf{x}_t .

Here, different with [50], we add an additional input \mathbf{x}^C as the conditional image to guide the model f_θ to keep the same content with \mathbf{x}^C during the denoising process.

B. Conditional image choice

To get realistic super-resolution images, [51], [48] also introduced diffusion models with conditional denoising on a pre-trained feature extractor or a bicubic upsampled image on a low-resolution image. In this work, we leverage the power of the current development of SISR to provide a better conditional image. Specifically, given a low-resolution image I^{LR} and a pre-trained super-resolution model ϕ_θ , we generate

Algorithm 2 ACDMSR

Training : train denoising model f_θ

Input: Dataset D , schedule α_t, σ_t , timesteps T , pre-trained super-resolution model ϕ_θ

- 1: **repeat**
- 2: $(\mathbf{I}^{HR}, \mathbf{I}^{LR}) \sim D, t \sim \text{Uniform}(\{1, \dots, T\}),$
 $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 3: $\mathbf{x}_0 = \mathbf{I}^{HR}, \mathbf{x}^C = \phi_\theta(\mathbf{I}^{LR})$
- 4: $\mathbf{x}_t = \alpha_t \mathbf{I}^{HR} + \sigma_t \epsilon$
- 5: Take a gradient descent step on
 $\nabla_\theta \|\mathbf{I}^{HR} - f_\theta(\mathbf{x}_t, t, \mathbf{x}^C)\|^2$
- 6: **until** converged

Inference: super resolve \mathbf{I}^{LR}

Input: trained denoising model f_θ , pre-trained super-resolution model ϕ_θ , diffusion sampler $\mathcal{F}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C)$

- 1: $\mathbf{x}^C = \phi_\theta(\mathbf{I}^{LR})$
 - 2: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 3: **for** $t = T, \dots, 1$ **do**
 - 4: $\mathbf{x}_{t-1} = \mathcal{F}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C)$
 - 5: **end for**
-

our conditional image by $\mathbf{x}^C = \phi_\theta(\mathbf{I}^{LR})$, which has been proved to be more plausible for obtaining results with better perceptual quality in ablation study V-C.

C. Sampling Process

Diffusion models are known to be slow and need thousands of forward evaluation steps to achieve the generated image with good quality. Similarly, the diffusion model-based super-resolution method inherited this drawback. To remedy this issue, we propose a n -th order sampler that adapts two accelerating sampling strategies from existing work, *i.e.*, DDIM [56] and DPM-solver [62]. These two sampling strategy has been shown to help the diffusion model achieve good image quality and keep a short sampling time. In this section, we first define what a sampler is. Then we describe the proposed super-resolution sampler in detail.

Iterative super-resolution sampler. Given a pre-trained model f_θ with objective Eq 10, and a low-resolution conditional image \mathbf{x}^C , we define a iterative super-resolution sampler from $t = T$ to $t = 0$ as:

$$\mathbf{x}_{t-1} = \mathcal{F}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C) \quad (11)$$

where \mathbf{x}_t is the ancestor of \mathbf{x}_{t-1} .

First order deterministic sampling. Different from SR3 and SRdiff sampling via a stochastic way, we use a deterministic sampling method to conduct the iterative reverse process $\mathbf{x}_{t-1} = \mathcal{F}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C)$ in a DDIM-like manner which has been shown achieve a high-quality image in limited inference steps. Given the image \mathbf{x}_t at step t , we can write the generation process of \mathbf{x}_{t-1} as follows:

$$\begin{aligned} \mathbf{x}_{t-1} &= \mathcal{F}_{1st}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C) \\ &= \sqrt{\hat{\alpha}_{t-1}} \hat{\mathbf{x}}_0 + \sqrt{1 - \hat{\alpha}_{t-1}} \frac{\mathbf{x}_t - \sqrt{\hat{\alpha}_t} \hat{\mathbf{x}}_0}{\sqrt{1 - \hat{\alpha}_t}}, \end{aligned} \quad (12)$$

where $\hat{\mathbf{x}}_0$ is predicted with trained denoising model $\hat{\mathbf{x}}_0 = f_\theta(\mathbf{x}_t, t, \mathbf{x}^C)$. Compared to the DDPM sampling process in Eq. 8, the above sampling does not add noise in each step,

making it a deterministic method. Since, in each step, we need only one forward model evaluation, we call this method a first-order method.

Second order deterministic sampling. [63] view the diffusion model as a stochastic differential equation (SDE), which has the same transition distribution $q(\mathbf{x}_t|\mathbf{x}_0)$ as in Eq 2 for any $t \in [0, T]$:

$$d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\mathbf{w}_t, \quad \mathbf{x}_0 \sim q_0(\mathbf{x}_0), \quad (13)$$

where $\mathbf{w}_t \in \mathbb{R}^D$ is the standard Wiener process, and

$$f(t) = \frac{d \log \sqrt{\hat{\alpha}_t}}{dt}, g(t) = \frac{1 - \hat{\alpha}_t}{dt} - 2f(t)\sqrt{\hat{\alpha}_t}. \quad (14)$$

With some regularity, [64] shows that the above forward SDE Eq.13 has an equivalent reverse process starting from the marginal distribution $q(\mathbf{x}_T)$ at time T to time step 0:

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log q(\mathbf{x}_t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (15)$$

where score function $\nabla_{\mathbf{x}} \log q(\mathbf{x}_t)$ can be replaced with the noise prediction of a model $\epsilon_\theta(\mathbf{x}_t, t)$, such that:

$$\frac{d\mathbf{x}_t}{dt} = f(t)\mathbf{x}_t - \frac{g^2(t)}{2\sqrt{1 - \hat{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t), \quad \mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (16)$$

This probability flows ordinary differential equation (ODE) has the same marginal distribution at each time t as that of the SED in Eq. 13. Sampling can be done by solving the integral of the above ODE from T to 0. [62] identifies the integral of the above ODE Eq. 16 has a linear part $f(t)\mathbf{x}_t$ which can be solved exactly and a nonlinear part $\frac{g^2(t)}{2\sqrt{1 - \hat{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t)$ which needs a black-box ODE solver to approximate. Compared to solving the whole ODE using a black-box solver, this semilinear property enables the elimination of the approximation error of the linear part. We build our second-order deterministic sampler in a DPM-Solver [62] way. To this end, define $\lambda_t = \lambda(t) = \log \sqrt{(\hat{\alpha}_t(1 - \hat{\alpha}_t))}$ and its inverse function $t_\lambda(\cdot)$ such that $t = t_\lambda(\lambda(t))$, we formulate our second-order sampling method on image \mathbf{x}_t as follows:

$$\begin{aligned} s &= t_\lambda\left(\frac{\lambda_t + \lambda_{t-1}}{2}\right), \\ \mathbf{u} &= \mathcal{F}_{1st}(f_\theta, \mathbf{x}_t, s, \mathbf{x}^C), \\ \mathbf{x}_{t-1} &= \mathcal{F}_{1st}(f_\theta, \mathbf{u}, t, \mathbf{x}^C). \end{aligned} \quad (17)$$

Since there uses first order two times, we call the above iterative sampler a second order deterministic sampler and denote it as $\mathcal{F}_{2ed}(f_\theta, \mathbf{x}_t, t, \mathbf{x}^C)$.

We conducted experiments to compare these three sampling methods. As shown in Fig. 2, the PSNR of the original DDPM sampling method is below 20dB in 500 forward steps, which is due to the nature of the stochastic reverse process. DDPM requires many steps to remove the randomness added during each step during the reverse process. First-order and second-order deterministic sampling methods perform much better in small sampling steps. As shown in Fig. 2, the first-order and second-order methods demonstrate a trade-off between visual quality and image distortion in the low sampling steps region. As the number of sampling steps increases, the PSNR decreases while the NIQE visual quality measure improves.

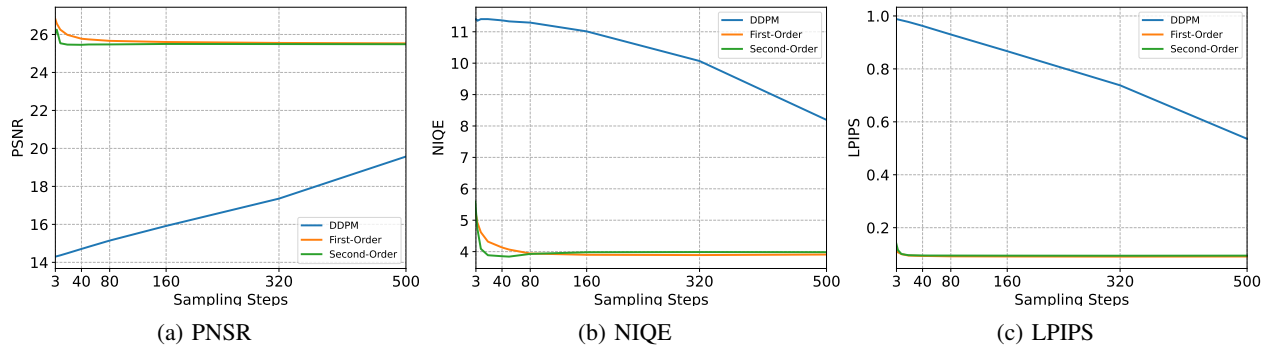


Fig. 2: Sampling steps comparison for DDPM sampling method, first-order deterministic, and second-order deterministic sampling method. (Conducted on Urban100 under $\times 4$ scale.)

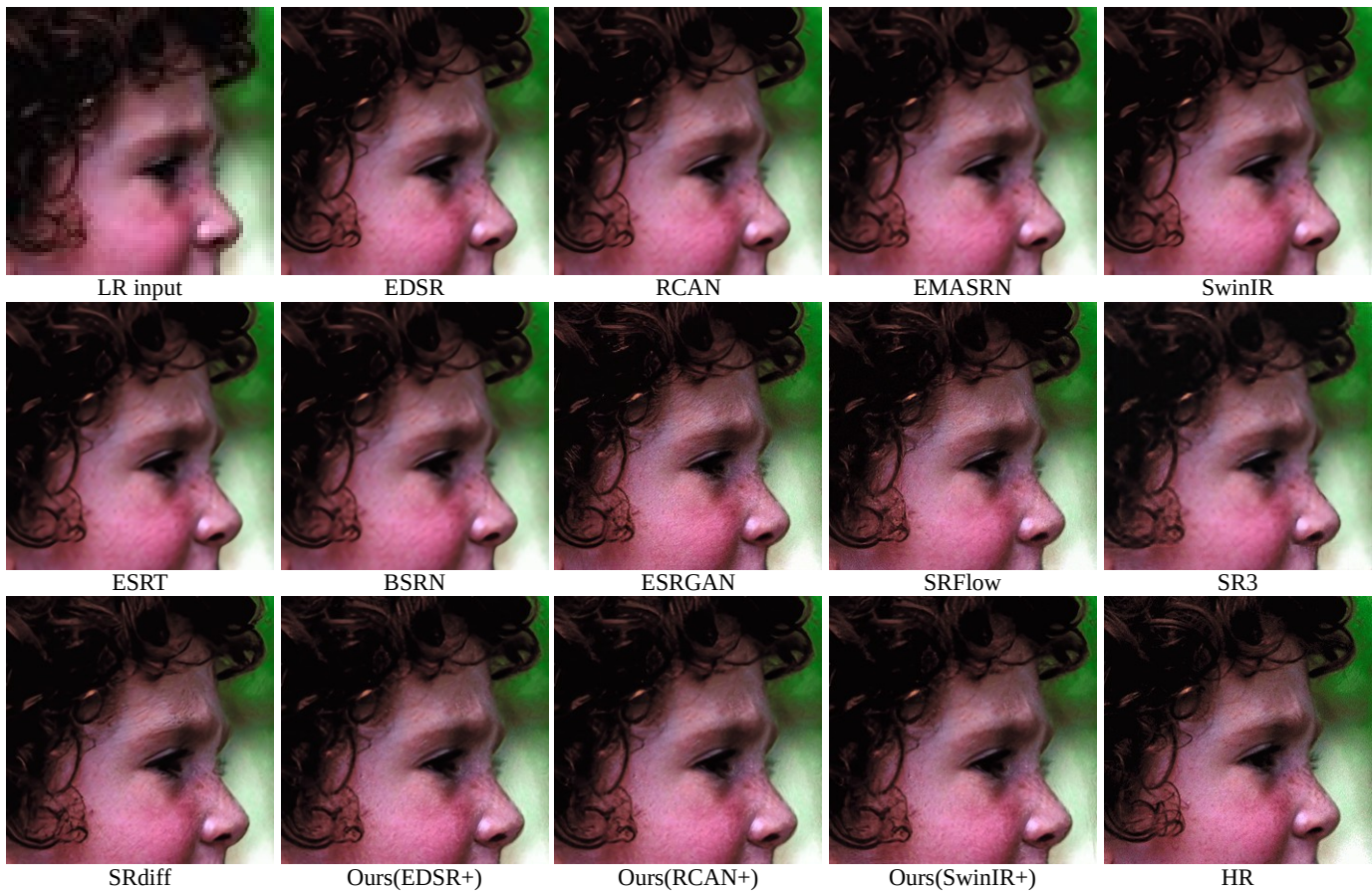


Fig. 3: Qualitative comparison with SOTAs performed on image ‘head’ from Set5 ($\times 4$ scale, best view in zoomed-in.)

Fig. 2 shows that the second-order sampler can achieve the lowest NIQE scores in just 40 steps. Therefore, we chose second-order sampling as our sampling method and set $T = 40$ during inference because we can achieve good perception quality from 40 feedforward steps.

V. EXPERIMENTS

A. Experimental Settings

Dataset. We use 800 image pairs in DIV2K as the training set. We take public benchmark datasets, *i.e.*, Set5, Set14,

Urban100, BSD100, and Manga109 as the test set to compare with other methods.

Setups. We set $T = 1000$ for training and $T = 40$ during the inference time for the diffusion model. We take the pre-trained super-resolution models (EDSR [20], and RCAN [12], SwinIR [26]) to provide the initial super-resolution image, *i.e.*the conditional image. The conditional diffusion model is trained with Adam optimizer and batch size 16, with a learning rate of 1×10^{-4} for 400k steps. The architecture of the model is the same as that in [48].

Metrics. The previous study has shown that distortion and

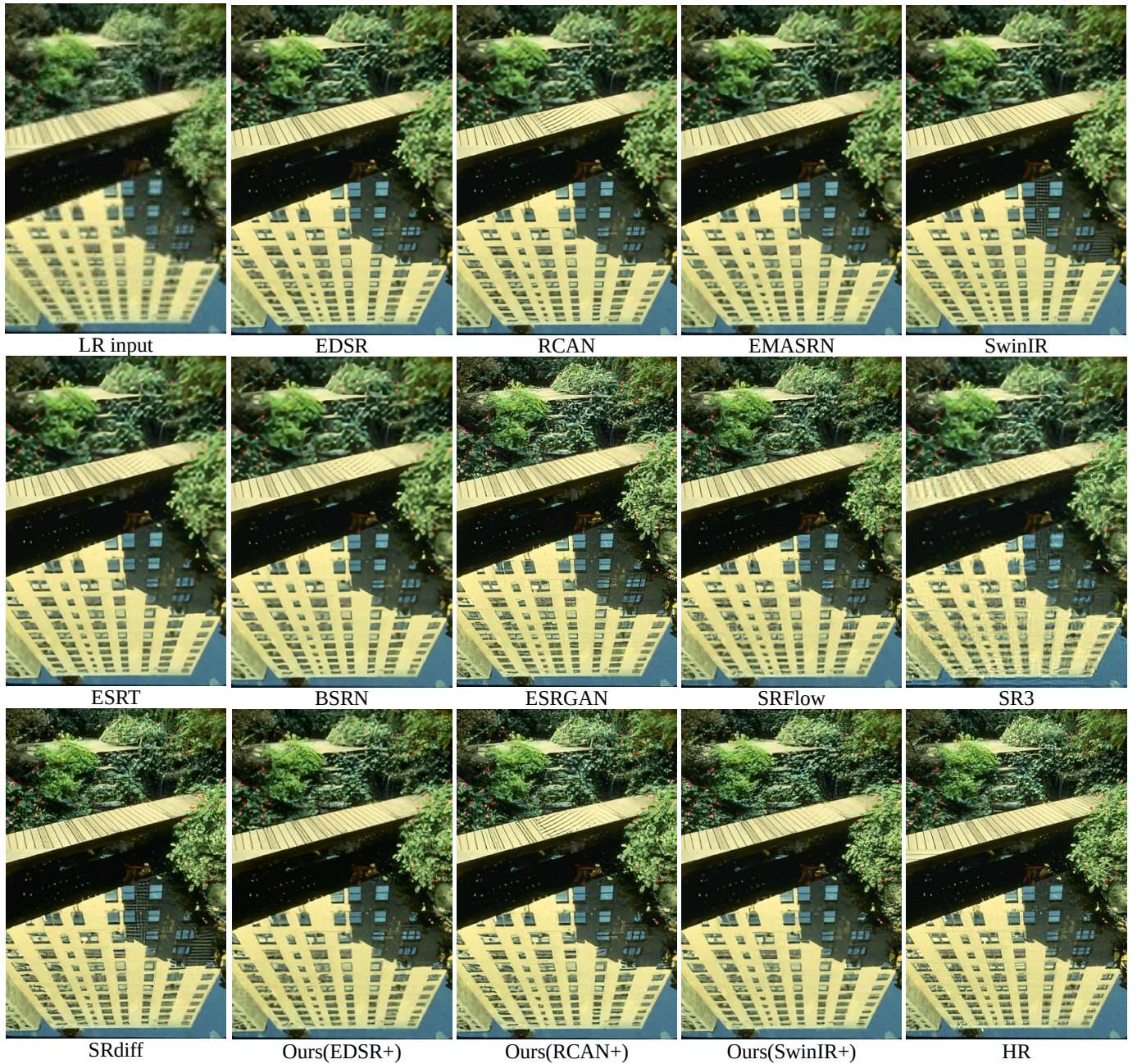


Fig. 4: Qualitative comparison with SOTAs performed on image ‘184026’ from BSD100 ($\times 4$ scale. Cropped and zoomed in for a better view.)

perceptual quality are at odds with each other, and there is a trade-off between them [31]. Since our work focuses on the perceptual quality, except the distortion metrics: PSNR and SSIM, we also provide perceptual metrics: LPIPS [65] and NIQE [66] to show that our method can generate better perceptual results than other methods. LPIPS is recently introduced as a reference-based image quality evaluation metric, which computes the perceptual similarity between the ground truth and the SR image. NIQE is a no-reference image quality score built on a “quality aware” collection of statistical features based on a simple and successful space domain natural scene statistic model.

B. Quantitative and Qualitative Results

To verify the effectiveness of our ACDMSR, we select some SOTA generative methods to conduct the comparative experiments, including ESGAN [36], SRFlow [39], SRDiff [51], SR3 [48]. We selected EDSR [20], RCAN [12], and SwinIR [26] to provide the conditional image, respectively. Therefore, we report three cases for our cDPMASR, *i.e.*, **EDSR+**, **RCAN+**, and **SwinIR+**. In addition, we also compare our method with some SOTA tradition CNN-based SR methods to verify further the effectiveness of our ACDMSR, including EDSR [20], RCAN [12], EMASRN [53], SwinIR [26], ESRT [27], and BSRN [55]. All the results are obtained from

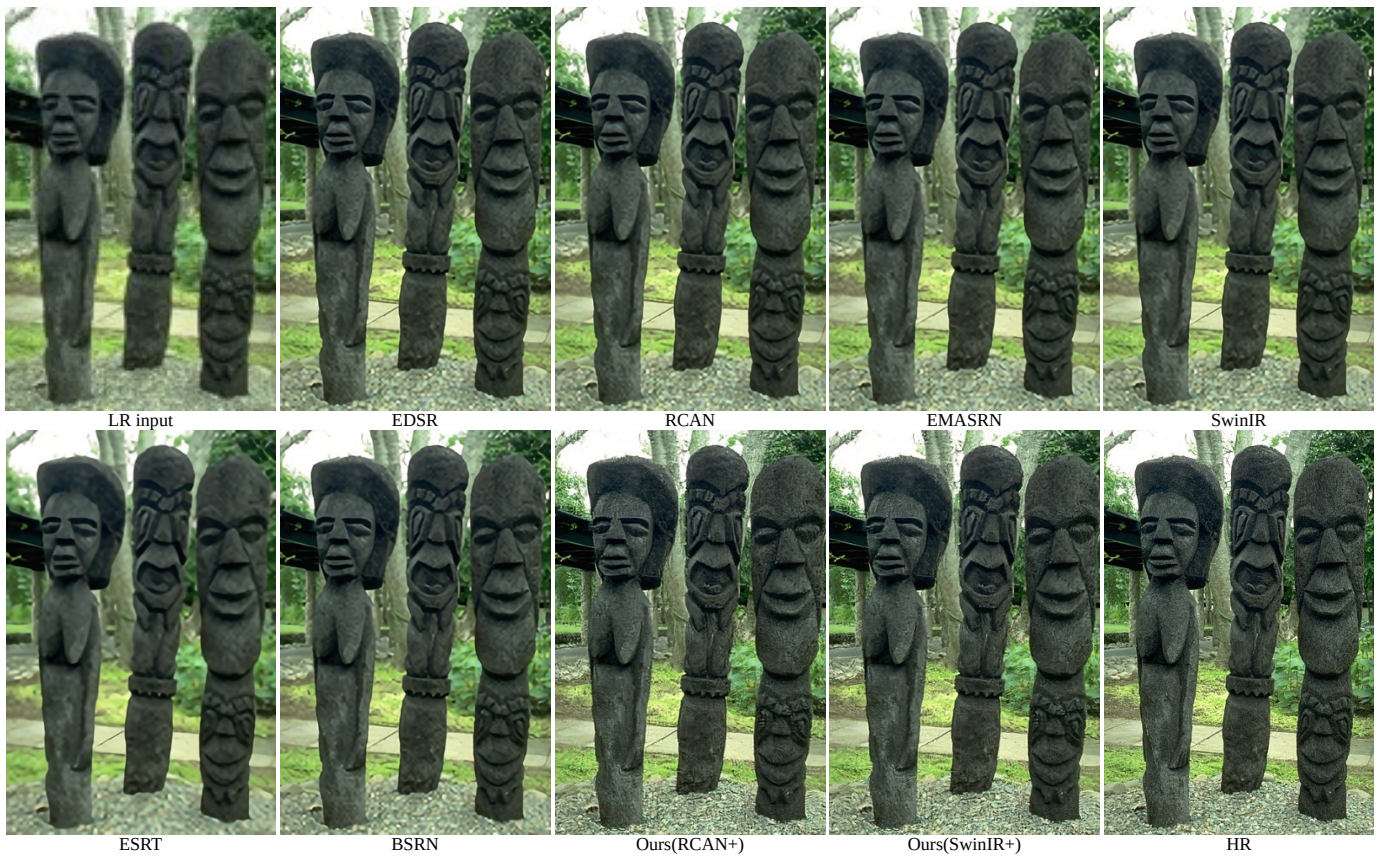


Fig. 5: Qualitative comparison with SOTAs performed on image ‘42012’ from BSD100 ($\times 3$ scale, best view in zoomed-in.)



Fig. 6: Qualitative comparison with SOTAs performed on image ‘AkkeraKanjinchou’ from Manga109 ($\times 8$ scale, best view in zoomed-in.)

TABLE I: Results on Set5, Set14, BSD100, Urban100, and Manga109. The best and the second-best results are highlighted in red and green.

Method			ESRGAN	SRFlow	SRDiff	SR3	Ours		
							EDSR+	RCAN+	SwinIR+
Set5	×4	LPIPS↓	0.0596	0.0767	0.0770	0.1084	0.0619	0.0606	0.0564
		PSNR↑	30.459	28.35	30.938	27.314	30.830	30.838	31.031
		SSIM↑	0.8516	0.8138	0.8738	0.7844	0.8684	0.8645	0.8676
LPIPS↓		0.0867	0.1318	0.1009	0.1284	0.0883	0.0865	0.0827	
PSNR↑		26.282	24.97	27.230	25.475	26.996	27.015	27.039	
SSIM↑		0.6980	0.6908	0.7432	0.6889	0.7316	0.7257	0.7341	
BSD100		LPIPS↓	0.0834	0.1831	0.1041	0.1392	0.0935	0.0953	0.0834
		PSNR↑	25.288	24.654	25.948	25.208	25.743	25.865	25.947
		SSIM↑	0.6495	0.6573	0.6833	0.6498	0.6681	0.6706	0.6743
Urban100	LPIPS↓	0.0944	0.1279	0.1077	0.1993	0.0997	0.0997	0.0934	
	PSNR↑	24.349	23.652	25.340	22.489	25.452	25.587	25.852	
	SSIM↑	0.7327	0.7312	0.7661	0.6336	0.7649	0.7681	0.7796	
Manga109	LPIPS↓	0.0420	0.0660	0.0473	0.1100	0.0409	0.0396	0.0374	
	PSNR↑	28.476	27.14	28.668	24.691	29.072	29.385	29.601	
Set5	×8	SSIM↑	0.8595	0.8244	0.8851	0.7568	0.8791	0.8816	0.8874
		LPIPS↓	0.2626	0.2304	0.3129	0.1872	0.2988	0.1813	0.1669
		PSNR↑	24.830	22.604	20.074	25.394	24.955	27.026	27.165
SSIM↑		0.6843	0.6062	0.5852	0.6897	0.6897	0.7834	0.7842	
Set14		LPIPS↓	0.2536	0.2965	0.2929	0.1983	0.2740	0.2083	0.2108
		PSNR↑	23.496	21.261	21.149	23.959	23.648	25.136	25.133
		SSIM↑	0.5854	0.4894	0.5393	0.5850	0.5864	0.6487	0.6485
BSD100		LPIPS↓	0.2503	0.3236	0.2879	0.1975	0.2734	0.2202	0.2219
		PSNR↑	23.868	21.619	19.162	23.918	24.015	24.983	24.944
	SSIM↑	0.5518	0.4634	0.4560	0.5345	0.5572	0.6054	0.6052	
Urban100	LPIPS↓	0.3062	0.2968	0.3494	0.2722	0.3384	0.2177	0.2245	
	PSNR↑	20.977	19.383	19.659	21.538	21.164	22.947	22.932	
	SSIM↑	0.5359	0.4999	0.5155	0.5546	0.5390	0.6428	0.6413	
Manga109	LPIPS↓	0.2364	0.2185	0.2660	0.1820	0.2564	0.1239	0.1322	
	PSNR↑	22.166	20.480	18.757	23.003	22.145	25.097	25.019	
		SSIM↑	0.6858	0.6434	0.6311	0.7126	0.6790	0.7984	0.7931

the provided codes or publicized papers.

Tab. I reports the PSNR, SSIM, and LPIPS values for those generative methods. Our method achieves superior performance under these quantitative metrics in terms of both distortion and perceptual quality across multiple standard datasets. ESRGAN is a typical GAN-based SR method, which includes an SR image generator and an SR image discriminator to push the generator to generate more realistic images. It achieves better LPIPSs, lower PSNRs, and lower SSIMs on different datasets under different scales compared with SRDiff. It seems the results generated by ESRGAN in Fig. 3, Fig. 4 and Fig.6 include more details than other methods, but it introduces too many false artifacts compared to the ground truth. SRFlow adopts the flow model to obtain reasonable high-resolution images by learning a conditional distribution when given low-resolution images. But the flow model needs invertible parameterized transformations with a tractable Jacobian determinant, which limits their expressiveness [48] and obtains worse LPIPSs, lower PSNRs, and lower SSIM compared with SRDiff and our method. And the results of SRflow seem noisy. To our knowledge, SRDiff and SR3 are state-of-the-art SR methods based on the diffusion model. SRDiff employs a two-stage structure, first pre-training an SR model and then optimizing the diffusion model. SR3 proposes an intuitive SR diffusion model based on the standard diffusion model in [50]. Our method is similar to these two methods. However, we use

existing SR methods to provide the conditional image instead of pretraining a new conditional-provided model and adjusting the optimization method by predicting the original image instead of the noise, which is more suitable for the SR task. With better conditional image, our method exhibits superior performance on both quantitative and qualitative results than SR3 [48]. Though SRdiff obtains some comparable numeric results in Tab.I, the visual results of our ACDMSR are closer to ground truths (Especially the forehead in Fig.3, the plants and the building in Fig.4). In Sec.V-C, we have further conducted ablation studies to prove that a better conditional image indeed helps improve the SR performance of the diffusion model.

Tab. II reports the PSNR, SSIM, LPIPS, and NIQE values for those traditional CNN-based SR methods. Because these methods are PSNR-directed and they all focus on obtaining results with good distortion [31], they can perform well on PSNR and SSIM, which are well-known to only partially correspond to human perception and can lead to algorithms with visibly lower quality in the reconstructed images [48]. The SR results of these PSNR-oriented methods are obviously so over-smooth that some details are missing. Though the PSNR and SSIM numbers of our method are slightly lower than theirs, it performs better when considering the metrics more in line with the human visual system.

In addition, we present Fig. 3, Fig. 4, Fig. 5, and Fig. 6 to illustrate the SR visual results on different datasets with

TABLE II: Results of different scales on Set5, Set14, BSD100, Urban100, and Manga109. The **bold** represents the best result.

Method			EDSR	RCAN	EMASRN	SwinIR	ESRT	BSRN	Ours		
									EDSR+	RCAN+	SwinIR+
Set5		LPIPS↓	0.0322	0.0321	-	0.0316	0.0609	0.0611	0.0121	0.0121	0.0125
		NIQE↓	5.3005	5.2721	-	5.3325	5.3226	5.3824	4.3374	4.4035	4.4151
		PSNR↑	38.193	38.271	-	38.357	38.088	38.072	36.241	36.255	36.462
		SSIM↑	0.9609	0.9614	-	0.9620	0.9598	0.9597	0.9404	0.9399	0.9431
Set14		LPIPS↓	0.0458	0.0446	-	0.0433	0.0968	0.0937	0.0275	0.0262	0.0252
		NIQE↓	5.0109	4.9893	-	4.9729	5.2071	5.1882	3.8002	3.8680	3.8216
		PSNR↑	33.948	34.126	-	34.141	33.690	33.642	32.073	32.315	32.277
		SSIM↑	0.9202	0.9216	-	0.9227	0.9183	0.9186	0.8834	0.8858	0.8863
BSD100	×2	LPIPS↓	0.0623	0.0615	-	0.0608	0.1463	0.1458	0.0281	0.0278	0.0283
		NIQE↓	4.9536	4.9673	-	4.9238	5.1657	5.1902	3.3873	3.5461	3.4376
		PSNR↑	32.352	32.389	-	32.448	32.272	32.221	30.176	30.293	30.346
		SSIM↑	0.9019	0.9024	-	0.9030	0.8993	0.8987	0.8544	0.8572	0.8588
Urban100		LPIPS↓	0.0359	0.0346	-	0.0333	0.0619	0.0619	0.0273	0.0269	0.0272
		NIQE↓	4.5070	4.4983	-	4.4880	4.6086	4.5908	3.9641	3.9850	3.9371
		PSNR↑	32.967	33.175	-	33.404	32.602	32.324	31.386	31.538	31.721
		SSIM↑	0.9359	0.9371	-	0.9394	0.9320	0.9296	0.9124	0.9112	0.9152
Manga109		LPIPS↓	0.0106	0.0102	-	0.0100	0.0228	0.0226	0.0072	0.0070	0.0067
		NIQE↓	4.5104	4.5217	-	4.4956	4.6483	4.6622	3.8334	3.9864	3.8553
		PSNR↑	39.193	39.438	-	39.586	39.073	38.992	37.046	37.518	37.607
		SSIM↑	0.9782	0.9787	-	0.9791	0.9773	0.9771	0.9650	0.9660	0.9669
Set5		LPIPS↓	0.0758	0.0747	0.1356	0.0734	0.1363	0.1378	0.0365	0.0354	0.0363
		NIQE↓	6.4616	6.4571	6.5556	6.6240	6.6755	6.8924	5.0188	4.8185	4.8930
		PSNR↑	34.680	34.758	34.361	34.878	34.612	34.499	32.618	32.715	33.001
		SSIM↑	0.9294	0.9300	0.9264	0.9312	0.9271	0.9262	0.8989	0.9002	0.9059
Set14		LPIPS↓	0.1002	0.1001	0.2175	0.0976	0.2288	0.2092	0.0630	0.0607	0.0637
		NIQE↓	5.5798	5.6797	5.9351	5.6477	5.9953	5.9011	3.8107	3.8364	3.7388
		PSNR↑	30.533	30.627	28.571	30.771	30.583	30.379	28.423	28.578	28.762
		SSIM↑	0.8465	0.8476	0.7809	0.8502	0.8341	0.8435	0.7861	0.7898	0.7953
BSD100	×3	LPIPS↓	0.1163	0.1150	0.2967	0.1124	0.2968	0.2944	0.0641	0.0648	0.0634
		NIQE↓	5.7653	5.8292	6.0468	5.7018	6.2079	6.0124	3.4016	3.4324	3.3014
		PSNR↑	29.263	29.301	29.053	29.367	29.224	29.181	26.938	27.139	27.163
		SSIM↑	0.8096	0.8106	0.8035	0.8124	0.8049	0.8035	0.7355	0.7415	0.7417
Urban100		LPIPS↓	0.0863	0.0830	0.1675	0.0798	0.1674	0.1581	0.0661	0.0654	0.0647
		NIQE↓	5.0547	5.1298	5.2835	5.0891	5.3741	5.2855	4.0667	4.0781	4.0287
		PSNR↑	28.812	29.009	28.042	29.288	28.469	28.389	27.424	27.722	27.889
		SSIM↑	0.8659	0.8685	0.8493	0.8744	0.8578	0.8558	0.8309	0.8361	0.8412
Manga109		LPIPS↓	0.0328	0.0320	0.0662	0.0307	0.0669	0.0638	0.0231	0.0215	0.0220
		NIQE↓	4.8532	4.9141	4.9435	4.8789	5.0512	4.9802	3.9195	3.8586	3.6899
		PSNR↑	34.200	34.429	33.433	34.749	34.109	33.982	32.207	32.338	32.428
		SSIM↑	0.9486	0.9498	0.9433	0.9517	0.9454	0.9450	0.9245	0.9236	0.9259
Set5		LPIPS↓	0.1098	0.1096	0.1820	0.1087	0.1889	0.1865	0.0619	0.0606	0.0564
		NIQE↓	7.2500	7.1562	7.2289	7.0368	6.9859	7.2315	5.5288	5.6325	4.9999
		PSNR↑	32.426	32.638	32.173	32.722	32.442	32.387	30.830	30.838	31.031
		SSIM↑	0.8985	0.9002	0.8948	0.9021	0.8960	0.8949	0.8684	0.8645	0.8676
Set14		LPIPS↓	0.1415	0.1387	0.2886	0.1369	0.2911	0.2871	0.0883	0.0865	0.0827
		NIQE↓	6.0475	6.1797	6.3646	6.2370	6.3369	6.2940	3.8035	3.8188	3.7958
		PSNR↑	28.679	28.851	28.572	28.937	28.614	28.534	26.996	27.015	27.039
		SSIM↑	0.7883	0.7885	0.7809	0.7914	0.7845	0.7837	0.7316	0.7257	0.7341
BSD100	×4	LPIPS↓	0.1551	0.1536	0.3847	0.1542	0.3881	0.3829	0.0935	0.0953	0.0834
		NIQE↓	6.3351	6.3104	6.5912	6.3638	6.6465	6.5235	3.3751	3.4251	3.4096
		PSNR↑	27.734	27.743	27.552	27.841	27.725	27.675	25.743	25.865	25.947
		SSIM↑	0.7425	0.7430	0.7351	0.7461	0.7369	0.7353	0.6681	0.6706	0.6743
Urban100		LPIPS↓	0.1220	0.1220	0.2342	0.1200	0.2396	0.2315	0.0997	0.0997	0.0934
		NIQE↓	5.4302	5.4886	5.6733	5.4203	5.9356	5.7585	4.0943	4.0521	4.1578
		PSNR↑	26.645	26.745	26.012	27.075	26.522	26.278	25.452	25.587	25.852
		SSIM↑	0.8039	0.8066	0.7837	0.8165	0.7965	0.7903	0.7649	0.7681	0.7796
Manga109		LPIPS↓	0.0562	0.0544	0.1066	0.1033	0.1109	0.1035	0.0409	0.0396	0.0374
		NIQE↓	5.1480	5.2272	5.2393	5.1456	5.4343	5.3175	3.7661	3.7599	3.7985
		PSNR↑	31.057	31.197	30.413	31.668	30.979	30.837	29.072	29.385	29.601
		SSIM↑	0.9160	0.9170	0.9076	0.9226	0.9107	0.9097	0.8791	0.8816	0.8874

varying scales. Our methods perform well on a variety of content, including humans, plants, text, and animals. These results further demonstrate the effectiveness of our approach in achieving both metric and perceptual quality.

C. Ablation Study

In this section, we conduct ablation studies to verify the influence of different conditional images on our ACDMSR. In addition, we also investigate how stochastic sampling and deterministic sampling influence the reconstruction results.

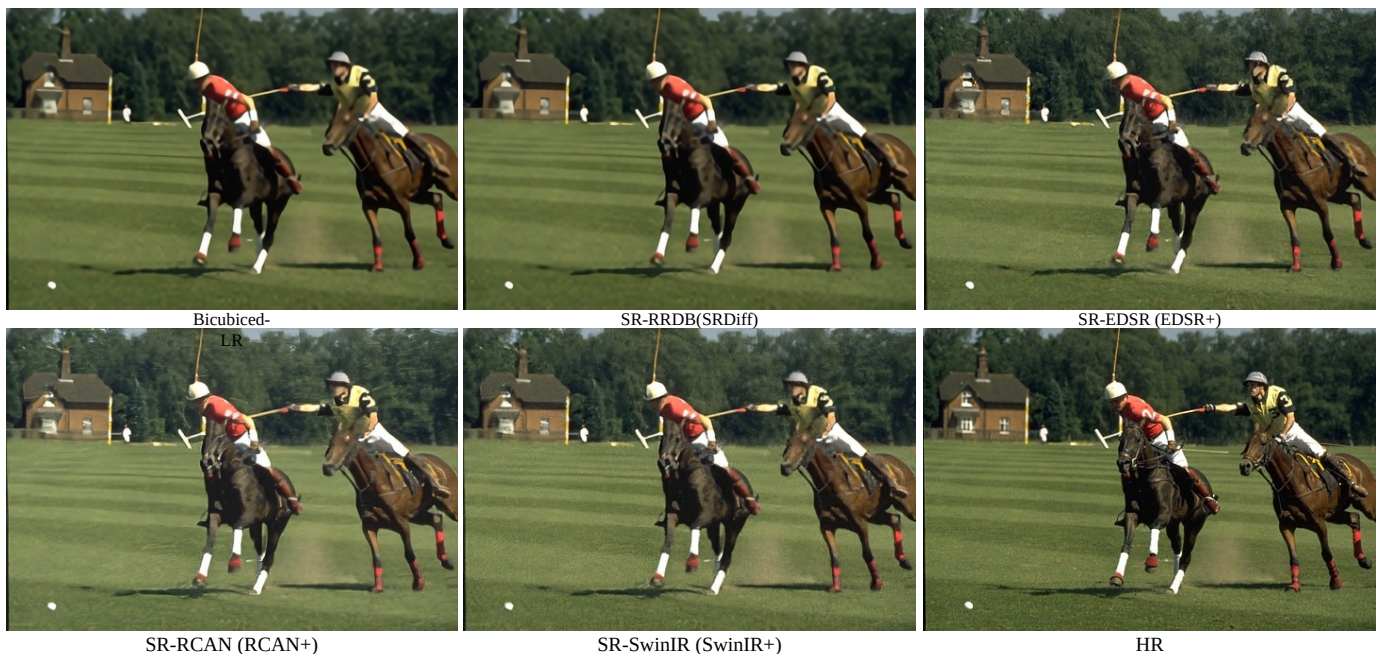


Fig. 7: Visual results on '361010' from BSD100 under different conditional images ($\times 4$ scale. Best view with zoomed-in.)

Furthermore, we conduct experiments to verify the effectiveness of different loss functions.

Different conditional images. Here, we conduct experiments to verify how different conditional images influence performance. We adopt LR, SRs generated by EDSR, RCAN, SwinIR, and the RRDB trained in SRDiff [51] as the conditional images to perform experiments, respectively. As shown in Tab.III and Fig.7, without any pre-processing, the result under LR conditional performs worst in both quantitative and qualitative. After being pre-trained by RRDB, EDSR, RCAN, and SwinIR, the conditional images can restore more details, pushing our ACDMSR model to perform better.

TABLE III: Results of ablation study for different conditional images on BSD100. ($4 \times SR$)

Method	LR	SR-RRDB	ours		
			EDSR+	RCAN+	SwinIR+
LPIPS↓	0.1412	0.1096	0.0935	0.0953	0.0834
PSNR↑	24.353	25.208	25.743	25.865	25.947
SSIM↑	0.6402	0.6589	0.6681	0.6706	0.6743

Noise-predicted Loss VS. Image-predicted Loss. We conduct an experiment on the Urban100 dataset with scale factor 4 to verify whether training the model to predict noise or

TABLE IV: Results of ablation study for different loss with $4 \times SR$ on Urban100.

Method		LPIPS	PSNR	SSIM
EDSR+	Image-predicted	0.0997	25.452	0.7649
	Noise-predicted	0.1106	24.866	0.7805
RCAN+	Image-predicted	0.0997	25.587	0.7681
	Noise-predicted	0.1051	25.035	0.7820
SwinIR+	Image-predicted	0.0934	25.852	0.7796
	Noise-predicted	0.1018	25.130	0.7915

images can achieve better performance. As shown in Tab. IV, the image prediction model can achieve both better distortion metric (PSNR, SSIM) and perceptual quality (LPIPS), compared with SR3 [48] and SRdiff [51], whose model predicts the added noise. It is because the image-predict model is more likely to learn the distribution of image information, which helps obtain good results for super-resolution reconstruction.

VI. CONCLUSION

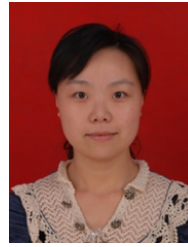
Our work revisits diffusion models in super-resolution and reveals that taking a pre-super-resolved version for the given LR image as the conditional image can help to achieve a better high-resolution image. Based on this, we propose a simple but non-trivial DPM-based super-resolution post-process framework, *i.e.*, ACDMSR. By taking a pre-super-resolved version of the given LR image and adapting the standard diffusion models to perform super-resolution, our ACDMSR improves both qualitative and quantitative results and can generate more photo-realistic counterparts for the low-resolution images on benchmark datasets (Set5, Set14, Urban100, BSD100, Manga109). In the future, we will extend our ACDMSR to images with more complex degradation.

Although our method achieves impressive results in generating high-quality images in single image super-resolution, however, it inherits the natural issue of the diffusion models that require multiple feedforwards to achieve the final output. The recent progress in the research community attempts to resolve this drawback of the diffusion model to shorten it to a single step with promising results [59], [67], [60], which can be beneficial for the SISR framework proposed by our method. In future work, we will focus on accelerating the inference process of diffusion models for image super-resolution.

REFERENCES

- [1] C. Yan, Y. Hao, L. Li, J. Yin, A. Liu, Z. Mao, Z. Chen, and X. Gao, "Task-adaptive attention for image captioning," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [2] Z. Chen, K.-Y. Lin, and W.-S. Zheng, "Consistent intra-video contrastive learning with asynchronous long-term memory bank," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [3] T. X. Pham, A. Niu, Z. Kang, S. R. Madjid, J. W. Hong, D. Kim, J. T. J. Tee, and C. D. Yoo, "Self-supervised visual representation learning via residual momentum," *arXiv preprint arXiv:2211.09861*, 2022.
- [4] A. Niu, K. Zhang, C. Zhang, C. Zhang, I. S. Kweon, C. D. Yoo, and Y. Zhang, "Fast adversarial training with noise augmentation: A unified perspective on randstart and gradalign," *arXiv preprint arXiv:2202.05488*, 2022.
- [5] F. Pan, I. Shin, F. Rameau, S. Lee, and I. S. Kweon, "Unsupervised intra-domain adaptation for semantic segmentation through self-supervision," in *CVPR*, 2020.
- [6] F. Pan, S. Hur, S. Lee, J. Kim, and I. S. Kweon, "MI-bpm: Multi-teacher learning with bidirectional photometric mixing for open compound domain adaptation in semantic segmentation," in *ECCV*, 2022.
- [7] F. Pan, F. Rameau, and I. S. Kweon, "Labeling where adapting fails: Cross-domain semantic segmentation with point supervision via active selection," *arXiv preprint arXiv:2206.00181*, 2022.
- [8] K. Chang, H. Li, Y. Tan, P. L. K. Ding, and B. Li, "A two-stage convolutional neural network for joint demosaicking and super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [9] A. Niu, Y. Zhu, C. Zhang, J. Sun, P. Wang, I. S. Kweon, and Y. Zhang, "Ms2net: Multi-scale and multi-stage feature fusion for blurred image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [10] A. Niu, K. Zhang, T. X. Pham, J. Sun, Y. Zhu, I. S. Kweon, and Y. Zhang, "Cdpmr: Conditional diffusion probabilistic models for single image super-resolution," *arXiv preprint arXiv:2302.12831*, 2023.
- [11] 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.
- [12] 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.
- [13] H. Ma, X. Chu, B. Zhang, and S. Wan, "A matrix-in-matrix neural network for image super resolution," *arXiv preprint arXiv:1903.07949*, 2019.
- [14] 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.
- [15] Y. Zuo, Q. Wu, Y. Fang, P. An, L. Huang, and Z. Chen, "Multi-scale frequency reconstruction for guided depth map super-resolution via deep residual network," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [16] Z.-S. Liu, W.-C. Siu, and Y.-L. Chan, "Photo-realistic image super-resolution via variational autoencoders," *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [17] J. Lyn, "Multi-level feature fusion mechanism for single image super-resolution," *arXiv preprint arXiv:2002.05962*, 2020.
- [18] X. Li and Z. Chen, "Single image super-resolution reconstruction based on fusion of internal and external features," *Multimedia Tools and Applications*, 2021.
- [19] H. Zhang, Y. Zhu, J. Sun, and Y. Zhang, "Real-world image super-resolution via kernel augmentation and stochastic variation," in *ICIP*, 2022.
- [20] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *CVPR workshops*, 2017.
- [21] Y. Hu, J. Li, Y. Huang, and X. Gao, "Channel-wise and spatial feature modulation network for single image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [22] F. Li, H. Bai, and Y. Zhao, "Filternet: Adaptive information filtering network for accurate and fast image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [23] A. Niu, P. Wang, Y. Zhu, J. Sun, Q. Yan, and Y. Zhang, "Gran: Ghost residual attention network for single image super resolution," *Multimedia Tools and Applications*, 2023.
- [24] L.-J. Deng, W. Guo, and T.-Z. Huang, "Single-image super-resolution via an iterative reproducing kernel hilbert space method," *IEEE Transactions on Circuits and Systems for Video Technology*, 2015.
- [25] Y. Mei, Y. Fan, and Y. Zhou, "Image super-resolution with non-local sparse attention," in *CVPR*, 2021.
- [26] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "Swinir: Image restoration using swin transformer," in *ICCV*, 2021.
- [27] Z. Lu, J. Li, H. Liu, C. Huang, L. Zhang, and T. Zeng, "Transformer for single image super-resolution," in *CVPR workshop*, 2022.
- [28] G. Wu, J. Jiang, X. Liu, and J. Ma, "A practical contrastive learning framework for single image super-resolution," *arXiv preprint arXiv:2111.13924*, 2021.
- [29] Y. Zhu, H. Shuai, G. Liu, and Q. Liu, "Self-supervised video representation learning using improved instance-wise contrastive learning and deep clustering," *IEEE Transactions on Circuits and Systems for Video Technology*, 2022.
- [30] A. Niu, K. Zhang, T. X. Pham, P. Wang, J. Sun, I. S. Kweon, and Y. Zhang, "Learning from multi-perception features for real-word image super-resolution," *arXiv preprint arXiv:2305.18547*, 2023.
- [31] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in *CVPR*, 2018.
- [32] Z. He, S. Tang, J. Yang, Y. Cao, M. Y. Yang, and Y. Cao, "Cascaded deep networks with multiple receptive fields for infrared image super-resolution," *IEEE transactions on circuits and systems for video technology*, 2018.
- [33] M. Delbracio, H. Talebi, and P. Milanfar, "Projected distribution loss for image enhancement," *arXiv preprint arXiv:2012.09289*, 2020.
- [34] D. Freirich, T. Michaeli, and R. Meir, "A theory of the distortion-perception tradeoff in wasserstein space," *NeurIPS*, 2021.
- [35] J. Whang, M. Delbracio, H. Talebi, C. Saharia, A. G. Dimakis, and P. Milanfar, "Deblurring via stochastic refinement," in *CVPR*, 2022.
- [36] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *ECCV workshops*, 2018.
- [37] J. W. Soh, G. Y. Park, J. Jo, and N. I. Cho, "Natural and realistic single image super-resolution with explicit natural manifold discrimination," in *CVPR*, 2019.
- [38] C. Tian, X. Zhang, J. C.-W. Lin, W. Zuo, Y. Zhang, and C.-W. Lin, "Generative adversarial networks for image super-resolution: A survey," *arXiv preprint arXiv:2204.13620*, 2022.
- [39] A. Lugmayr, M. Danelljan, L. V. Gool, and R. Timofte, "Srfflow: Learning the super-resolution space with normalizing flow," in *ECCV*, 2020.
- [40] J. Liang, A. Lugmayr, K. Zhang, M. Danelljan, L. Van Gool, and R. Timofte, "Hierarchical conditional flow: A unified framework for image super-resolution and image rescaling," in *ICCV*, 2021.
- [41] V. Wolf, A. Lugmayr, M. Danelljan, L. Van Gool, and R. Timofte, "Deflow: Learning complex image degradations from unpaired data with conditional flows," in *CVPR*, 2021.
- [42] S. Bell-Kligler, A. Shocher, and M. Irani, "Blind super-resolution kernel estimation using an internal-gan," *NeurIPS*, 2019.
- [43] M. Emad, M. Peemen, and H. Corporaal, "Dualsr: Zero-shot dual learning for real-world super-resolution," in *WACV*, 2021.
- [44] L. Metz, B. Poole, D. Pfau, and J. Sohl-Dickstein, "Unrolled generative adversarial networks," *arXiv preprint arXiv:1611.02163*, 2016.
- [45] S. Ravuri and O. Vinyals, "Classification accuracy score for conditional generative models," *NeurIPS*, 2019.
- [46] L. Huang and Y. Xia, "Fast blind image super resolution using matrix-variable optimization," *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [47] A. Lugmayr, M. Danelljan, and R. Timofte, "Ntire 2021 learning the super-resolution space challenge," in *CVPR*, 2021.
- [48] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, and M. Norouzi, "Image super-resolution via iterative refinement," *TPAMI*, 2022.
- [49] C. Laroche and M. Tassano, "Bridging the domain gap in real world super-resolution," in *ICIP*, 2022.
- [50] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *NeurIPS*, 2020.
- [51] H. Li, Y. Yang, M. Chang, S. Chen, H. Feng, Z. Xu, Q. Li, and Y. Chen, "Srdiff: Single image super-resolution with diffusion probabilistic models," *Neurocomputing*, 2022.
- [52] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [53] X. Zhu, K. Guo, S. Ren, B. Hu, M. Hu, and H. Fang, "Lightweight image super-resolution with expectation-maximization attention mechanism," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [54] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *ECCV*, 2018.

- [55] Z. Li, Y. Liu, X. Chen, H. Cai, J. Gu, Y. Qiao, and C. Dong, "Blueprint separable residual network for efficient image super-resolution," in *CVPR*, 2022.
- [56] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.
- [57] A. Bansal, E. Borgnia, H.-M. Chu, J. S. Li, H. Kazemi, F. Huang, M. Goldblum, J. Geiping, and T. Goldstein, "Cold diffusion: Inverting arbitrary image transforms without noise," *arXiv preprint arXiv:2208.09392*, 2022.
- [58] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," *arXiv preprint arXiv:2206.00364*, 2022.
- [59] T. Salimans and J. Ho, "Progressive distillation for fast sampling of diffusion models," *arXiv preprint arXiv:2202.00512*, 2022.
- [60] Y. Song, P. Dhariwal, M. Chen, and I. Sutskever, "Consistency models," *arXiv preprint arXiv:2303.01469*, 2023.
- [61] Y. Nikankin, N. Haim, and M. Irani, "Sinfusion: Training diffusion models on a single image or video," *arXiv preprint arXiv:2211.11743*, 2022.
- [62] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, "Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps," *arXiv preprint arXiv:2206.00927*, 2022.
- [63] D. Kingma, T. Salimans, B. Poole, and J. Ho, "Variational diffusion models," *NeurIPS*, 2021.
- [64] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.
- [65] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *CVPR*, 2018.
- [66] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal processing letters*, 2012.
- [67] H. Zheng, W. Nie, A. Vahdat, K. Azizzadenesheli, and A. Anandkumar, "Fast sampling of diffusion models via operator learning," *arXiv preprint arXiv:2211.13449*, 2022.



Jinqiu Sun received her B.S., M.S., and Ph.D. degrees from Northwestern Polytechnical University in 1999, 2004, and 2005, respectively. She is presently a Professor of the School of Astronomy at Northwestern Polytechnical University. Her research work focuses on signal and image processing, computer vision, and pattern recognition.



Yu Zhu received the B.S., M.S., and M.S. degrees from Northwestern Polytechnical University, Xi'an, China, in 2008, 2011, and 2017, respectively. He is presently an associate researcher at the School of Computer Science, Northwestern Polytechnical University. His current research interests include image processing and image super-resolution.



In So Kweon received the B.S. and the M.S. degrees in Mechanical Design and Production Engineering from Seoul National University, Korea, in 1981 and 1983, respectively, and the Ph.D. degree in Robotics from the Robotics Institute at Carnegie Mellon University in 1990. He is currently a Professor of electrical engineering (EE) and the director of the National Core Research Center – P3 DigiCar Center at KAIST. He served as the department head of Automation and Design Engineering (ADE) at KAIST in 1995-1998. His research interests include computer vision and robotics. He has co-authored several books, including "Metric Invariants for Camera Calibration," and more than 300 technical papers. He served as a Founding Associate-Editor-in-Chief for "The International Journal of Computer Vision and Applications", and has been an Editorial Board Member for "The International Journal of Computer Vision" since 2005. He is a member of many computer vision and robotics conference program committees and has been a program co-chair for several conferences and workshops. Most recently, he has been a general co-chair of the 2012 Asian Conference on Computer Vision (ACCV) Conference. He received several awards from international conferences, including "The Best Student Paper Runnerup Award in the IEEE-CVPR'2009" and "The Student Paper Award in the ICCAS'2008". He also earned several honors at KAIST, including the 2002 Best Teaching Award in EE. He is a member of KROS, ICROS, and IEEE.



Axi Niu received her B.S. and M.S. degrees from Henan University, Kaifeng, China, in 2014 and 2017. She is currently pursuing the Ph.D. degree with the School of Computer Science, Northwestern Polytechnical University, Xi'an, China. Her research interests include image processing and computer vision.



Kang Zhang received his B.S. degree from Harbin Institute of Technology, 2020. He is currently pursuing the Ph.D. degree at Korea Advanced Institute of Science & Technology. His research work focuses on Deep Learning, Self-Supervised Learning, and Adversarial Machine Learning.



Pham Xuan Trung received his B.S. degree in the School of Electronics and Telecommunications (SET) at Hanoi University of Science and Technology (HUST) in 2014. He is currently working toward his Ph.D. at KAIST under the supervision of Prof. Chang D. Yoo. His doctoral research interests include Speech Processing, Self-Supervised Learning, and Computer Vision.



Yanning Zhang received her B.S. degree from the Dalian University of Science and Engineering in 1988, M.S. and Ph.D. Degree from Northwestern Polytechnical University in 1993 and 1996, respectively. She is presently a Professor of School of Computer Science and Technology, Northwestern Polytechnical University. She is also the organization chair of ACCV2009 and the publicity chair of ICME2012. Her research focuses on signal and image processing, computer vision, and pattern recognition. She has published over 200 papers in these fields, including the ICCV2011 best student paper. She is a member of IEEE.