

DR2: Diffusion-based Robust Degradation Remover for Blind Face Restoration

Zhixin Wang¹ Ziying Zhang¹ Xiaoyun Zhang^{1†} Huangjie Zheng²
 Mingyuan Zhou² Ya Zhang^{1,3} Yanfeng Wang^{1,3†}

¹CMIC*, Shanghai Jiao Tong University, ³Shanghai AI Lab, ²The University of Texas at Austin

{dedsec.z, zyzzhang2000, xiaoyun.zhang, ya_zhang, wangyanfeng} @sjtu.edu.cn

huangjie.zheng@utexas.edu, mingyuan.zhou@mcombs.utexas.edu

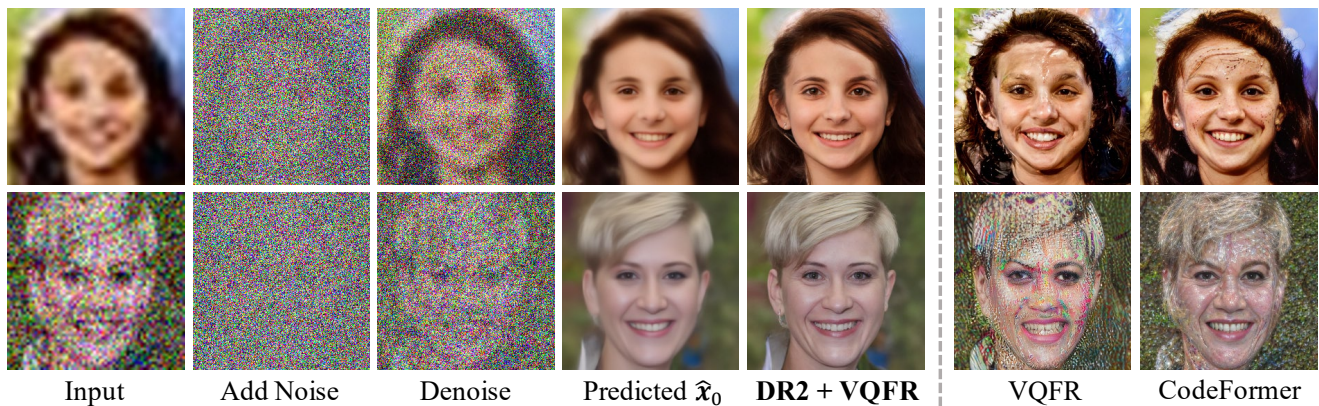


Figure 1. **DR2 uses Denoising Diffusion Probabilistic Models to remove degradation.** The generative process is conditioned on the low-quality input after being diffused into a noisy status. As a result, DR2 predicts coarse faces \hat{x}_0 regardless of the degradation type. On severely degraded images, our final restoration results achieve high quality with fewer artifacts than previous arts [14, 48].

Abstract

Blind face restoration usually synthesizes degraded low-quality data with a pre-defined degradation model for training, while more complex cases could happen in the real world. This gap between the assumed and actual degradation hurts the restoration performance where artifacts are often observed in the output. However, it is expensive and infeasible to include every type of degradation to cover real-world cases in the training data. To tackle this robustness issue, we propose **Diffusion-based Robust Degradation Remover (DR2)** to first transform the degraded image to a coarse but degradation-invariant prediction, then employ an enhancement module to restore the coarse prediction to a high-quality image. By leveraging a well-performing denoising diffusion probabilistic model, our DR2 diffuses input images to a noisy status where various types of degradation give way to Gaussian noise, and then captures semantic information through iterative denoising steps. As a result, DR2 is robust against common degradation (e.g. blur, resize, noise and compression) and compatible with different

designs of enhancement modules. Experiments in various settings show that our framework outperforms state-of-the-art methods on heavily degraded synthetic and real-world datasets.

1. Introduction

Blind face restoration aims to restore high-quality face images from their low-quality counterparts suffering from unknown degradation, such as low-resolution [5, 11, 27], blur [45], noise [23, 36], compression [10], etc. Great improvement in restoration quality has been witnessed over the past few years with the exploitation of various facial priors. Geometric priors such as facial landmarks [5], parsing maps [4, 5], and heatmaps [43] are pivotal to recovering the shapes of facial components. Reference priors [9, 25, 26] of high-quality images are used as guidance to improve details. Recent research investigates generative priors [39, 42] and high-quality dictionaries [14, 24, 48], which help to generate photo-realistic details and textures.

Despite the great progress in visual quality, these methods lack a robust mechanism to handle degraded inputs besides relying on pre-defined degradation to synthesize the

[†]Corresponding author.

*Cooperative Medianet Innovation Center.

training data. When applying them to images of severe or unseen degradation, undesired results with obvious artifacts can be observed. As shown in Fig. 1, artifacts typically appear when 1) the input image lacks high-frequency information due to downsampling or blur (1st row), in which case restoration networks can not generate adequate information, or 2) the input image bears corrupted high-frequency information due to noise or other degradation (2nd row), and restoration networks mistakenly use the corrupted information for restoration. The primary cause of this inadaptability is the inconsistency between the synthetic degradation of training data and the actual degradation in the real world.

Expanding the synthetic degradation model for training would improve the models’ adaptability but it is apparently difficult and expensive to simulate every possible degradation in the real world. To alleviate the dependency on synthetic degradation, we leverage a well-performing denoising diffusion probabilistic model (DDPM) [16, 37] to remove the degradation from inputs. DDPM generates images through a stochastic iterative denoising process and Gaussian noisy images can provide guidance to the generative process [6, 29]. As shown in Fig. 2, noisy images are **degradation-irrelevant** conditions for DDPM generative process. Adding extra Gaussian noise (right) makes different degradation less distinguishable compared with the original distribution (left), while DDPM can still capture the semantic information within this noise status and recover clean face images. This property of pretrained DDPM makes it a robust degradation removal module though only high-quality face images are used for training the DDPM.

Our overall blind face restoration framework DR2E consists of the **D**iffusion-based **R**obust **D**egradation **R**emover (DR2) and an **E**nhancement module. In the first stage, DR2 first transforms the degraded images into coarse, smooth, and visually clean intermediate results, which fall into a degradation-invariant distribution (4th column in Fig. 1). In the second stage, the degradation-invariant images are further processed by the enhancement module for high-quality details. By this design, the enhancement module is compatible with various designs of restoration methods in seeking the best restoration quality, ensuring our DR2E achieves both strong robustness and high quality.

We summarize the contributions as follows. (1) We propose DR2 that leverages a pretrained diffusion model to remove degradation, achieving robustness against complex degradation without using synthetic degradation for training. (2) Together with an enhancement module, we employ DR2 in a two-stage blind face restoration framework. The enhancement module has great flexibility in incorporating a variety of restoration methods to achieve high restoration quality. (3) Comprehensive experiments show that our framework outperforms state-of-the-art methods on heavily degraded synthetic and real-world datasets.

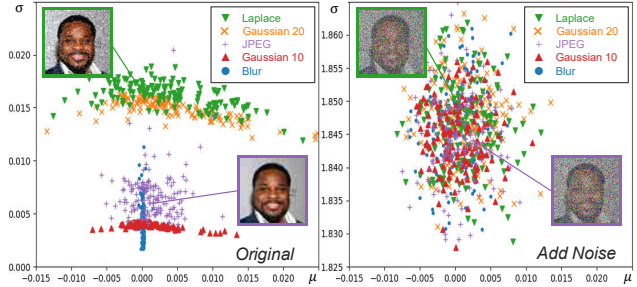


Figure 2. **Mean and standard variation of pixel-wise error distribution.** (Left) the error between original degraded input y and its ground truth low-resolution image \hat{y} (only bicubically down-sampled); (Right) the error between $q(y_{500}|y)$ and $q(\hat{y}_{500}|\hat{y})$ sampled by Eq. (2), with extra Gaussian noise added by the diffusion function.

2. Related Work

Blind Face Restoration Based on face hallucination or face super-resolution [1, 18, 40, 44], blind face restoration aims to restore high-quality faces from low-quality images with unknown and complex degradation. Many facial priors are exploited to alleviate dependency on degraded inputs. Geometry priors, including facial landmarks [5, 22, 49], parsing maps [4, 5, 36], and facial component heatmaps [43] help to recover accurate shapes but contain no information on details in themselves. Reference priors [9, 25, 26] of high-quality images are used to recover details or preserve identity. To further boost restoration quality, generative priors like pretrained StyleGAN [20, 21] are used to provide vivid textures and details. PULSE [30] uses latent optimization to find latent code of high-quality face, while more efficiently, GPEN [42], GFP-GAN [39], and GLEAN [2] embed generative priors into the encoder-decoder structure. Another category of methods utilizes pretrained Vector-Quantize [12, 32, 38] codebooks. DFDNet [24] suggests constructing dictionaries of each component (e.g. eyes, mouth), while recent VQFR [14] and CodeFormer [48] pretrain high-quality dictionaries on entire faces, acquiring rich expressiveness.

Diffusion Models Denoising Diffusion Probabilistic Models (DDPM) [16, 37] are a fast-developing class of generative models in unconditional image generation rivaling Generative Adversarial Networks (GAN) [13, 19, 31]. Recent research utilizes it for super-resolution. SR3 [35] modifies DDPM to be conditioned on low-resolution images through channel-wise concatenation. However, it fixes the degradation to simple downsampling and does not apply to other degradation settings. Latent Diffusion [33] performs super-resolution in a similar concatenation manner but in a low-dimensional latent space. ILVR [6] proposes a conditioning method to control the generative process of pretrained DDPM for image-translation tasks. Diffusion-based methods face a common problem of slow sampling speed, while

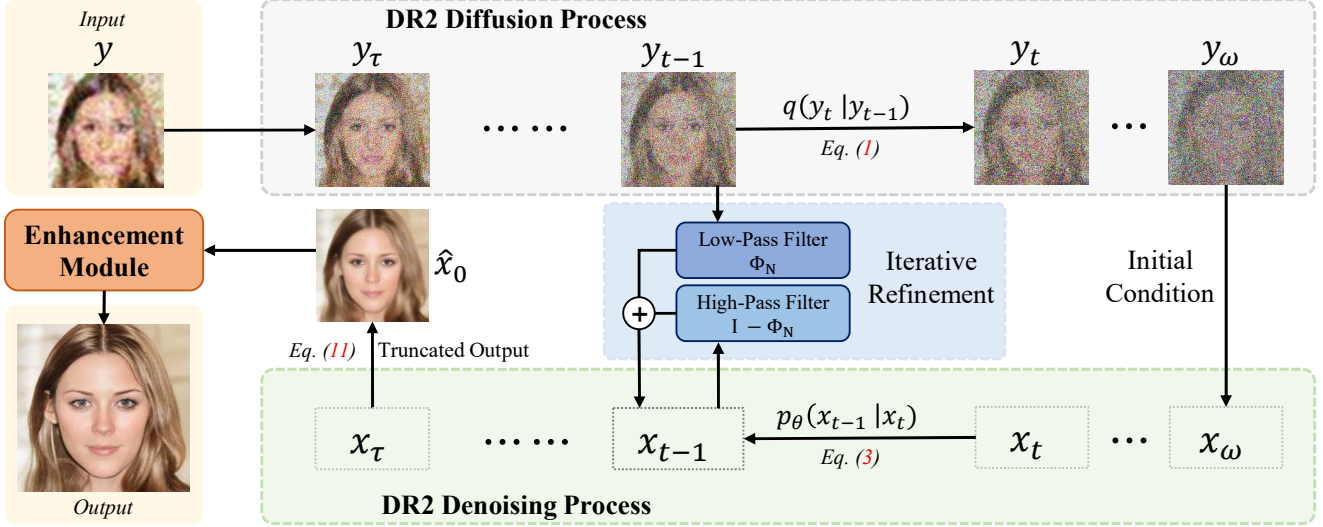


Figure 3. **Overall DR2E framework.** It consists of DR2 as the degradation removal module and an enhancement module. During inference, we sample $\mathbf{y}_\tau, \mathbf{y}_{\tau+1}, \dots, \mathbf{y}_\omega$ through diffusion process and use them as guidance. We use \mathbf{y}_ω as \mathbf{x}_ω and start the denoising process from step ω to τ . After each transition from \mathbf{x}_t to \mathbf{x}_{t-1} , we combine the low-frequency of \mathbf{y}_{t-1} and high-frequency of \mathbf{x}_{t-1} . At step τ , we predict $\hat{\mathbf{x}}_0$ based on the estimated noise. Then the enhancement module produces high-quality output from $\hat{\mathbf{x}}_0$.

our DR2E adopts a hybrid architecture like [47] to speed up the sampling process.

3. Methodology

Our proposed DR2E framework is depicted in Fig. 3, which consists of the degradation remover DR2 and an enhancement module. Given an input image \mathbf{y} suffering from unknown degradation, diffused low-quality information \mathbf{y}_{t-1} is provided to refine the generative process. As a result, DR2 recovers a coarse result $\hat{\mathbf{x}}_0$ that is semantically close to \mathbf{y} and degradation-invariant. Then the enhancement module maps $\hat{\mathbf{x}}_0$ to the final output with higher resolution and high-quality details.

3.1. Preliminary

Denosing Diffusion Probabilistic Models (DDPM) [16, 37] are a class of generative models that first pre-defines a variance schedule $\{\beta_1, \beta_2, \dots, \beta_T\}$ to progressively corrupt an image \mathbf{x}_0 to a noisy status through forward (diffusion) process:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (1)$$

Moreover, based on the property of the Markov chain, for any intermediate timestep $t \in \{1, 2, \dots, T\}$, the corresponding noisy distribution has an analytic form:

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} \end{aligned} \quad (2)$$

where $\bar{\alpha}_t := \prod_{s=1}^t (1 - \beta_s)$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if T is big enough, usually $T = 1000$.

The model progressively generates images by reversing the forward process. The generative process is also a Gaussian transition with the learned mean $\boldsymbol{\mu}_\theta$:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}) \quad (3)$$

where σ_t is usually a pre-defined constant related to the variance schedule, and $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ is usually parameterized by a denoising U-Net $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ [34] with the following equivalence:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)) \quad (4)$$

3.2. Framework Overview

Suppose the low-quality image \mathbf{y} is degraded from the high-quality ground truth $\mathbf{x} \sim \mathcal{X}(\mathbf{x})$ as $\mathbf{y} = \mathcal{T}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} describes the degradation model. Previous studies constructs the inverse function $\mathcal{T}^{-1}(\cdot, \mathbf{z})$ by modeling $p(\mathbf{x} | \mathbf{y}, \mathbf{z})$ with a pre-defined \mathbf{z} [5, 11, 27]. It meets the adaptation problem when actual degradation \mathbf{z}' in the real world is far from \mathbf{z} .

To overcome this challenge, we propose to model $p(\mathbf{x} | \mathbf{y})$ without a known \mathbf{z} by a two-stage framework: it first removes degradation from inputs and get $\hat{\mathbf{x}}_0$, then maps degradation-invariant $\hat{\mathbf{x}}_0$ to high-quality outputs. Our target is to maximize the likelihood:

$$\begin{aligned} p_{\psi, \phi}(\mathbf{x} | \mathbf{y}) &= \int p_\psi(\mathbf{x} | \hat{\mathbf{x}}_0) p_\phi(\hat{\mathbf{x}}_0 | \mathbf{y}) d\hat{\mathbf{x}}_0 \\ &= \mathbb{E}_{\hat{\mathbf{x}}_0 \sim p_\phi(\hat{\mathbf{x}}_0 | \mathbf{y})} [p_\psi(\mathbf{x} | \hat{\mathbf{x}}_0)], \end{aligned} \quad (5)$$

$p_\phi(\hat{x}_0|\mathbf{y})$ corresponds to the degradation removal module, and $p_\psi(\mathbf{x}|\hat{x}_0)$ corresponds to the enhancement module. For the first stage, instead of directly learning the mapping from \mathbf{y} to \hat{x}_0 which usually involves a pre-defined degradation model z , we come up with an important assumption and propose a diffusion-based method to remove degradation.

Assumption. For the diffusion process defined in Eq. (2), (1) there exists an intermediate timestep τ such that for $t > \tau$, the distance between $q(\mathbf{x}_t|\mathbf{x})$ and $q(\mathbf{y}_t|\mathbf{y})$ is close especially in the low-frequency part; (2) there exists $\omega > \tau$ such that the distance between $q(\mathbf{x}_\omega|\mathbf{x})$ and $q(\mathbf{y}_\omega|\mathbf{y})$ is eventually small enough, satisfying $q(\mathbf{x}_\omega|\mathbf{x}) \approx q(\mathbf{y}_\omega|\mathbf{y})$.

Note this assumption is not strong, as paired \mathbf{x} and \mathbf{y} would share similar low-frequency contents, and for sufficiently large $t \approx T$, $q(\mathbf{x}_t|\mathbf{x})$ and $q(\mathbf{y}_t|\mathbf{y})$ are naturally close to the standard $\mathcal{N}(\mathbf{0}, \mathbf{I})$. This assumption is also qualitatively justified in Fig. 2. Intuitively, if \mathbf{x} and \mathbf{y} are close in distribution (implying mild degradation), we can find ω and τ in a relatively small value and vice versa.

Then we rewrite the objective of the degradation removal module by applying the assumption $q(\mathbf{x}_\omega|\mathbf{x}) \approx q(\mathbf{y}_\omega|\mathbf{y})$:

$$p_\phi(\hat{x}_0|\mathbf{y}) = \int p(\hat{x}_0|\mathbf{x}_\tau)p_\theta(\mathbf{x}_\tau|\mathbf{y}_\omega)q(\mathbf{y}_\omega|\mathbf{y})d\mathbf{x}_\tau d\mathbf{y}_\omega \quad (6)$$

$$\approx \int p(\hat{x}_0|\mathbf{x}_\tau)p_\theta(\mathbf{x}_\tau|\mathbf{x}_\omega)q(\mathbf{x}_\omega|\mathbf{x})d\mathbf{x}_\tau d\mathbf{x}_\omega \quad (7)$$

$$p_\theta(\mathbf{x}_\tau|\mathbf{x}_\omega) = \prod_{t=\tau+1}^{\omega} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (8)$$

By replacing variable from \mathbf{y}_ω to \mathbf{x}_ω , Eq. (7) and Eq. (8) naturally yields a DDPM model that denoises \mathbf{x}_ω back to \mathbf{x}_τ , and we can further predict \hat{x}_0 by the reverse of Eq. (2). \hat{x}_0 would maintain semantics with \mathbf{x} if proper conditioning methods like [6] is adopted. So by leveraging a DDPM, we propose Diffusion-based Robust Degradation Remover (DR2) according to Eq. (6).

3.3. Diffusion-based Robust Degradation Remover

Consider a pretrained DDPM $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ (Eq. (3)) with a denoising U-Net $\epsilon_\theta(\mathbf{x}_t, t)$ pretrained on high-quality face dataset. We respectively implement $q(\mathbf{y}_\omega|\mathbf{y})$, $p_\theta(\mathbf{x}_\tau|\mathbf{y}_\omega)$ and $p(\hat{x}_0|\mathbf{x}_\tau)$ in Eq. (6) by three steps in below.

(1) Initial Condition at ω . We first “forward” the degraded image \mathbf{y} to an initial condition \mathbf{y}_ω by sampling from Eq. (2) and use it as \mathbf{x}_ω :

$$\mathbf{x}_\omega := \mathbf{y}_\omega = \sqrt{\bar{\alpha}_\omega}\mathbf{y} + \sqrt{1 - \bar{\alpha}_\omega}\boldsymbol{\epsilon}, \quad (9)$$

$\omega \in \{1, 2, \dots, T\}$. This corresponds to $q(\mathbf{x}_\omega|\mathbf{y})$ in Eq. (6). Then the DR2 denoising process starts at step ω . This reduces the samplings steps and helps to speed up as well.

(2) Iterative Refinement. After each transition from \mathbf{x}_t to \mathbf{x}_{t-1} ($\tau + 1 \leq t \leq \omega$), we sample \mathbf{y}_{t-1} from \mathbf{y} through

Eq. (2). Based on Assumption (1), we replace the low-frequency part of \mathbf{x}_{t-1} with that of \mathbf{y}_{t-1} because they are close in distribution, which is formulated as:

$$\mathbf{x}_{t-1} := \Phi_N(\mathbf{y}_{t-1}) + (\mathbf{I} - \Phi_N)(\mathbf{x}_{t-1}) \quad (10)$$

where $\Phi_N(\cdot)$ denotes a low-pass filter implemented by downsampling and upsampling the image with a sharing scale factor N . We drop the high-frequency part of \mathbf{y} for it contains little information due to degradation. Unfiltered degradation that remained in the low-frequency part would be covered by the added noise. These conditional denoising steps correspond to $p_\theta(\mathbf{x}_\tau|\mathbf{y}_\omega)$ in Eq. (6), which ensure the result shares basic semantics with \mathbf{y} .

Iterative refinement is pivotal for preserving the low-frequency information of the input images. With the iterative refinement, the choice of ω and the randomness of Gaussian noise affect little to the result. We present ablation study in the supplementary for illustration.

(3) Truncated Output at τ . As t gets smaller, the noise level gets milder and the distance between $q(\mathbf{x}_t|\mathbf{x})$ and $q(\mathbf{y}_t|\mathbf{y})$ gets larger. For small t , the original degradation is more dominating in $q(\mathbf{y}_t|\mathbf{y})$ than the added Gaussian noise. So the denoising process is truncated before t is too small. We use predicted noise at step τ ($0 < \tau < \omega$) to estimate the generation result as follows:

$$\hat{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_\tau}}(\mathbf{x}_\tau - \sqrt{1 - \bar{\alpha}_\tau}\epsilon_\theta(\mathbf{x}_\tau, \tau)) \quad (11)$$

This corresponds to $p(\hat{x}_0|\mathbf{x}_\tau)$ in Eq. (6). \hat{x}_0 is the output of DR2, which maintains the basic semantics of \mathbf{y} and is removed from various degradation.

Selection of N and τ . Downsampling factor N and output step τ have significant effects on the fidelity and “cleaness” of \hat{x}_0 . We conduct ablation studies in Sec. 4.4 to show the effects of these two hyper-parameters. The best choices of N and τ are data-dependent. Generally speaking, big N and τ are more effective to remove the degradation but lead to lower fidelity. On the contrary, small N and τ leads to high fidelity, but may keep the degradation in the outputs. While ω is empirically fixed to $\tau + 0.25T$.

3.4. Enhancement Module

With outputs of DR2, restoring the high-quality details only requires training an enhancement module $p_\psi(\mathbf{x}|\hat{x}_0)$ (Eq. (5)). Here we do not hypothesize about the specific method or architecture of this module. Any neural network that can be trained to map a low-quality image to its high-quality counterpart can be plugged in our framework. And the enhancement module is independently trained with its proposed loss functions.

Backbones. In practice, without loss of generality, we choose SPARNetHD [3] that utilized no facial priors, and VQFR [14] that pretrain a high-quality VQ codebook [12,

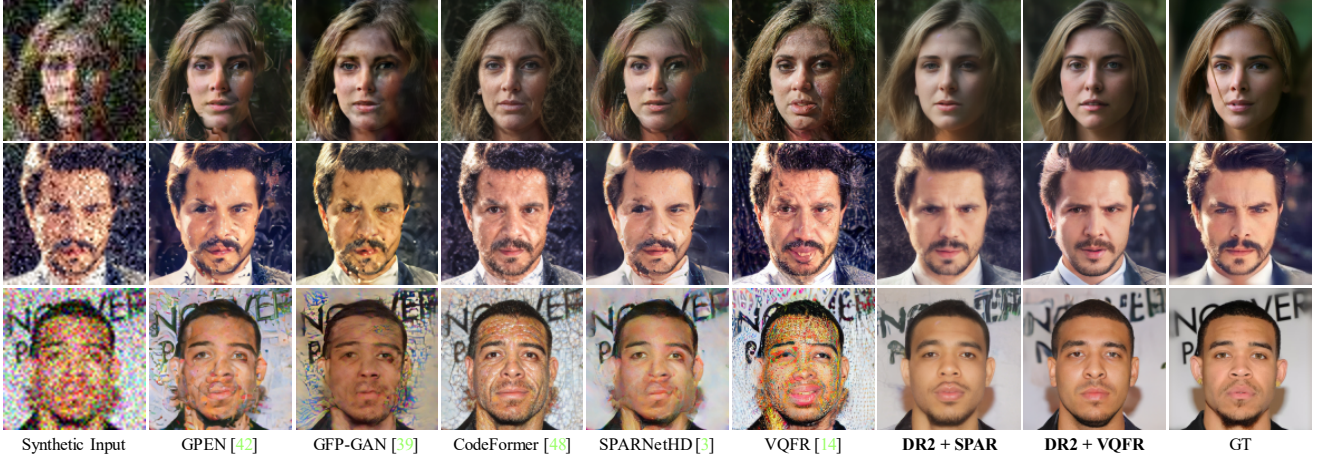


Figure 4. **Qualitative comparison on the CelebA-Test dataset.** Our method with different enhancement module backbones achieve higher restoration quality with fewer artifacts despite the heavy degradation in inputs.

[32, 38] as two alternative backbones for our enhancement module to justify that it can be compatible with a broad choice of existing methods. We denote them as DR2 + SPAR and DR2 + VQFR respectively.

Training Data. Any pretrained blind face restoration models can be directly plugged-in without further finetuning, but in order to help the enhancement module adapt better and faster to DR2 outputs, we suggest constructing training data for the enhancement module using DR2 as follows:

$$\mathbf{y} = DR2(\mathbf{x}; N, \tau) \otimes k_\sigma \quad (12)$$

Given a high-quality image \mathbf{x} , we first use DR2 to reconstruct itself with controlling parameters (N, τ) then convolve it with an Gaussian blur kernel k_σ . This helps the enhancement module adapt better and faster to DR2 outputs, which is recommended but not compulsory. Noting that beside this augmentation, **no** other degradation model is required in the training process as what previous works [3, 14, 39, 42, 48] do by using Eq. (13).

4. Experiments

4.1. Datasets and Implementation

Implementation. DR2 and the enhancement module are independently trained on FFHQ dataset [20], which contains 70,000 high-quality face images. We use pretrained DDPM proposed by [6] for our DR2. As introduced in Sec. 3.4, we choose SPARNetHD [3] and VQFR [14] as two alternative architectures for the enhancement module. We train SPARNetHD backbone from scratch with training data constructed by Eq. (12). We set $N = 4$ and randomly sample τ, σ from $\{50, 100, 150, 200\}, \{1 : 7\}$, respectively. As for VQFR backbone, we use its official pretrained model.

Testing Datasets. We construct one synthetic dataset and four real-world datasets for testing. A brief introduction of

each is as followed:

- *CelebA-Test.* Following previous works [3, 14, 39, 42, 48], we adopt a commonly used degradation model as follows to synthesize testing data from CelebA-HQ [28]:

$$\mathbf{y} = [(\mathbf{x} \otimes k_\sigma) \downarrow_r + n_\delta]_{JPEG_q} \quad (13)$$

A high-quality image \mathbf{x} is first convolved with a Gaussian blur kernel k_σ , then bicubically downsampled with a scale factor r . n_δ represents additive noise and is randomly chosen from Gaussian, Laplace, and Poisson. Finally, JPEG compression with quality q is applied. We use $r = 16, 8$, and 4 to form three restoration tasks denoted as $16\times$, $8\times$, and $4\times$. For each upsampling factor, we generate three splits with different levels of degradation and each split contains 1,000 images. The *mild* split randomly samples σ, δ and q from $\{3 : 5\}, \{5 : 20\}, \{60 : 80\}$, respectively. The *medium* from $\{5 : 7\}, \{15 : 40\}, \{40 : 60\}$. And the *severe* split from $\{7 : 9\}, \{25 : 50\}, \{30 : 40\}$.

- *WIDER-Normal* and *WIDER-Critical.* We select 400 critical cases suffering from heavy degradation (mainly low-resolution) from WIDER-face dataset [41] to form the WIDER-Critical dataset and another 400 regular cases for WIDER-Normal dataset.

- *CelebChild* contains 180 child faces of celebrities collected from the Internet. Most of them are only mildly degraded.

- *LFW-Test.* LFW [17] contains low-quality images with mild degradation from the Internet. We choose 1,000 testing images of different identities.

During testing, we conduct grid search for best controlling parameters (N, τ) of DR2 for each dataset. Detailed parameter settings are presented in the supplementary.

4.2. Comparisons with State-of-the-art Methods

We compare our method with several state-of-the-art face restoration methods: DFDNet [24], SPARNetHD [3],

Methods	$\times 16$				$\times 8$				$\times 4$			
	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FID \downarrow	PSNR \uparrow	SSIM \uparrow
DFDNet* [24]	0.5511	109.41	20.80	0.4929	0.5033	120.13	21.75	0.4758	0.4405	98.10	23.81	0.5357
GPEN [42]	0.4313	81.57	21.77	0.5916	0.3745	64.00	24.02	0.6398	0.2934	53.56	26.38	0.7057
GFP-GAN [39]	0.5430	139.13	18.35	0.4578	0.3233	56.88	23.36	0.6695	0.2720	58.78	24.94	0.7244
CodeFormer [48]	0.5176	117.17	19.70	0.4553	0.3465	71.22	23.04	0.5950	0.2587	61.41	26.33	0.7065
SPARNetHD [3]	0.4289	77.02	22.28	0.6114	0.3361	59.66	24.71	0.6743	0.2638	53.20	26.59	0.7255
VQFR [14]	0.6312	152.56	17.73	0.3381	0.4214	66.54	21.83	0.5345	0.3094	52.39	23.52	0.6335
DR2 + SPAR(ours)	0.3908	53.22	22.29	0.6587	0.3218	56.29	24.78	0.6966	0.2635	51.44	26.28	0.7263
DR2 + VQFR(ours)	0.3893	47.29	21.29	0.6222	0.3167	53.82	23.40	0.6802	0.2902	51.41	24.04	0.6844

Table 1. Quantitative comparisons on **CelebA-Test** dataset. **Red** and **blue** indicates the best and the second best performance. ‘*’ denotes using ground-truth facial landmarks as input.

GFP-GAN [39], GPEN [42], VQFR [14], and Codeformer [48]. We adopt their *official* codes and pretrained models.

For evaluation, we adopt pixel-wise metrics (PSNR and SSIM) and the perceptual metric (LPIPS [46]) for the CelebA-Test with ground truth. We also employ the widely-used non-reference perceptual metric FID [15].

Synthetic CelebA-Test. For each upsampling factor, we calculate evaluation metrics on three splits and present the average in Tab. 1. For $16\times$ and $8\times$ upsampling tasks where degradation is severe due to low resolution, DR2 + VQFR and DR2 + SPAR achieve the best and the second-best LPIPS and FID scores, indicating our results are perceptually close to the ground truth. Noting that DR2 + VQFR is better at perceptual metrics (LPIPS and FID) thanks to the pretrained high-quality codebook, and DR2 + SPAR is better at pixel-wise metrics (PSNR and SSIM) because without facial priors, the outputs have higher fidelity to the inputs. For $4\times$ upsampling task where degradation is relatively milder, previous methods trained on similar synthetic degradation manage to produce high-quality images without obvious artifacts. But our methods still obtain superior FID scores, showing our outputs have closer distribution to ground truth on different settings.

Qualitative comparisons from are presented in Fig. 4. Our methods produce fewer artifacts on severely degraded inputs compared with previous methods.

Real-World Datasets. We evaluate FID scores on different real-world datasets and present quantitative results in Tab. 2. On severely degraded dataset WIDER-Critical, our DR2 + VQFR and DR2 + SPAR achieve the best and the second best FID. On other datasets with only mild degradation, the restoration quality rather than robustness becomes the bottleneck, so DR2 + SPAR with no facial priors struggles to stand out, while DR2 + VQFR still achieves the best performance.

Qualitative results on WIDER-Critical are shown in Fig. 5. When input images’ resolutions are very low, previous methods fail to complement adequate information for

Datasets Methods	W-Cr FID \downarrow	W-Nm FID \downarrow	Celeb-C FID \downarrow	LFW FID \downarrow
DFDNet [24]	78.87	73.12	107.18	64.89
GPEN [42]	65.06	67.85	107.27	55.77
GFP-GAN [39]	64.14	59.20	111.79	54.84
CodeFormer [48]	66.84	60.10	114.34	56.15
SPARNetHD [3]	69.79	61.34	110.30	52.28
VQFR [14]	81.37	60.84	104.39	51.81
DR2 + SPAR(ours)	61.66	63.69	107.00	52.27
DR2 + VQFR(ours)	60.06	58.78	103.91	50.98

Table 2. Quantitative comparisons on **WIDER-Critical** (W-Cr), **WIDER-Normal** (W-Nm), **CelebChild** (Celeb-C) and **LFW-Test** (LFW). **Red** and **blue** indicates the best and the second best performance.

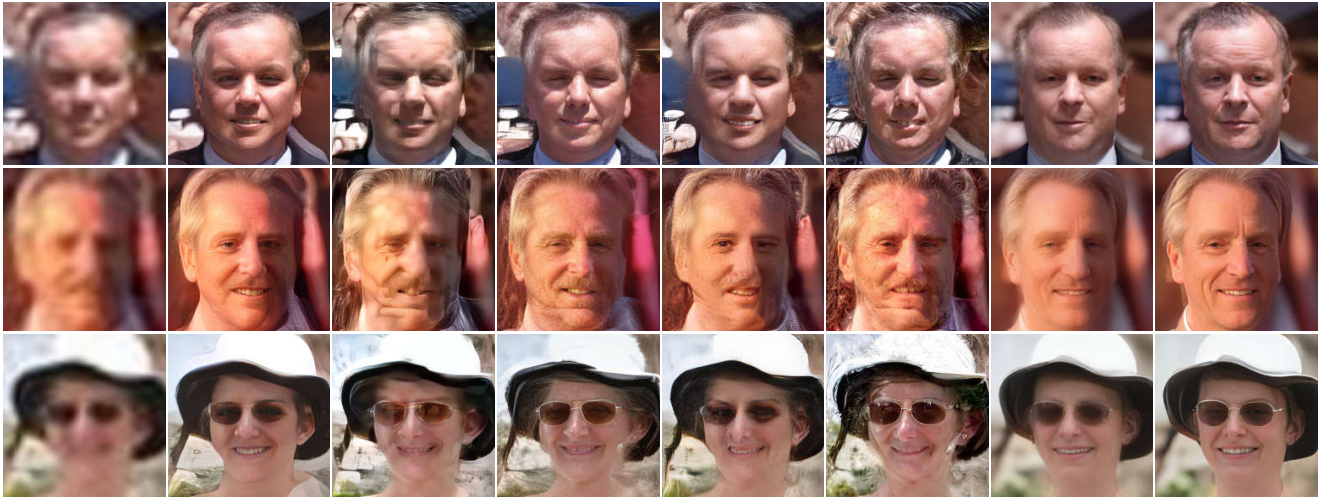
pleasant faces, while our outputs are visually more pleasant thanks to the generative ability of DDPM.

4.3. Comparisons with Diffusion-based Methods

Diffusion-based super-resolution methods can be grouped into two categories by whether feeding auxiliary input to the denoising U-Net.

SR3 [35] typically uses the concatenation of low-resolution images and x_t as the input of the denoising U-Net. But SR3 fixes degradation to bicubic downsampling during training, which makes it highly degradation-sensitive. For visual comparisons, we re-implement the concatenation-based method based on [8]. As shown in Fig. 6, minor noise in the second input evidently harm the performance of this concatenation-based method. Eventually, this type of method would rely on synthetic degradation to improve robustness like [33], while our DR2 have good robustness against different degradation without training on specifically degraded data.

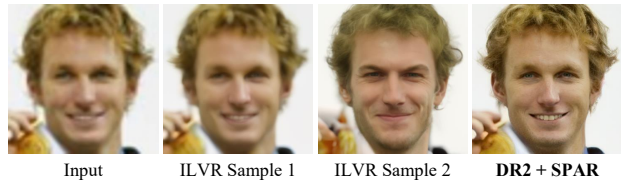
Another category of methods is training-free, exploiting pretrained diffusion methods like ILVR [6]. It shows the



Real Input GPEN [43] GFP-GAN [40] CodeFormer [49] SPARNetHD [3] VQFR [14] DR2 + SPAR DR2 + VQFR
 Figure 5. Qualitative comparisons on **WIDER-critical** dataset. Thanks to the generative ability of DR2, our methods produce more visually pleasant results when inputs are of low-resolution.



Input Concatenation DR2 Input Concatenation DR2
 Figure 6. Comparisons with diffusion-based super-resolution method based on concatenation of input images. This method is highly degradation-sensitive.



Input ILVR Sample 1 ILVR Sample 2 DR2 + SPAR
 Figure 7. Trade-off problem of ILVR for blind face restoration. ILVR Sample 1 is sampled with wide conditioning step range [6] and Sample 2 with narrow conditioning range.

ability to transform both clean and degraded low-resolution images into high-resolution outputs. However, relying solely on ILVR for blind face restoration faces the trade-off problem between fidelity and quality (realness). As shown in Fig. 7, ILVR Sample 1 has high fidelity to input but low visual quality because the conditioning information is over-used. On the contrary, under-use of conditions leads to high quality but low fidelity as ILVR Sample 2. In our framework, fidelity is controlled by DR2 and high-quality details are restored by the enhancement module, thus alleviating the trade-off problem.

4.4. Effect of Different N and τ

In this section, we explore the property of DR2 output in terms of the controlling parameter (N, τ) so that we can have a better intuitions for choosing appropriate parameters

for variant input data. To avoid the influence of the enhancement modules varying in structures, embedded facial priors, and training strategies, we only evaluate DR2 outputs with no enhancement.

In Fig. 8, DR2 outputs are generated with different combinations of N and τ . Bigger N and τ are effective to remove degradation but tent to make results deviant from the input. On the contrary, small N and τ lead to high fidelity, but may keep the degradation in outputs.

We provide quantitative evaluations on CelebA-Test ($8\times$, medium split) dataset in Fig. 9. With bicubically downsampled low-resolution images used as ground truth, we adopt pixel-wise metric (PSNR \uparrow) and identity distance (Deg \downarrow) based on the embedding angle of ArcFace [7] for evaluating the quality and fidelity of DR2 outputs. For scale $N = 4, 8, \text{ and } 16$, PSNR first goes up and Deg goes down because degradation is gradually removed as τ increases. Then they hit the optimal point at the same time before the outputs begin to deviate from the input as τ continues to grow. Optimal τ is bigger for smaller N . For $N = 2$, PSNR stops to increase before Deg reaches the optimality because Gaussian noise starts to appear in the output (like results sampled with $(N, \tau) = (2, 350)$ in Fig. 8). This cause of the appearance of Gaussian noise is that \mathbf{y}_t sampled by Eq. (2) contains heavy Gaussian noise when t ($t > \tau$) is big and most part of \mathbf{y}_t is utilized by Eq. (10) when N is small.

4.5. Discussion and Limitations

Our DR2 is built on a pretrained DDPM, so it would face the problem of slow sampling speed even we only perform $0.25T$ steps in total. But DR2 can be combined with diffusion acceleration methods like inference every 10 steps. And keep the output resolution of DR2 relatively low (256^2

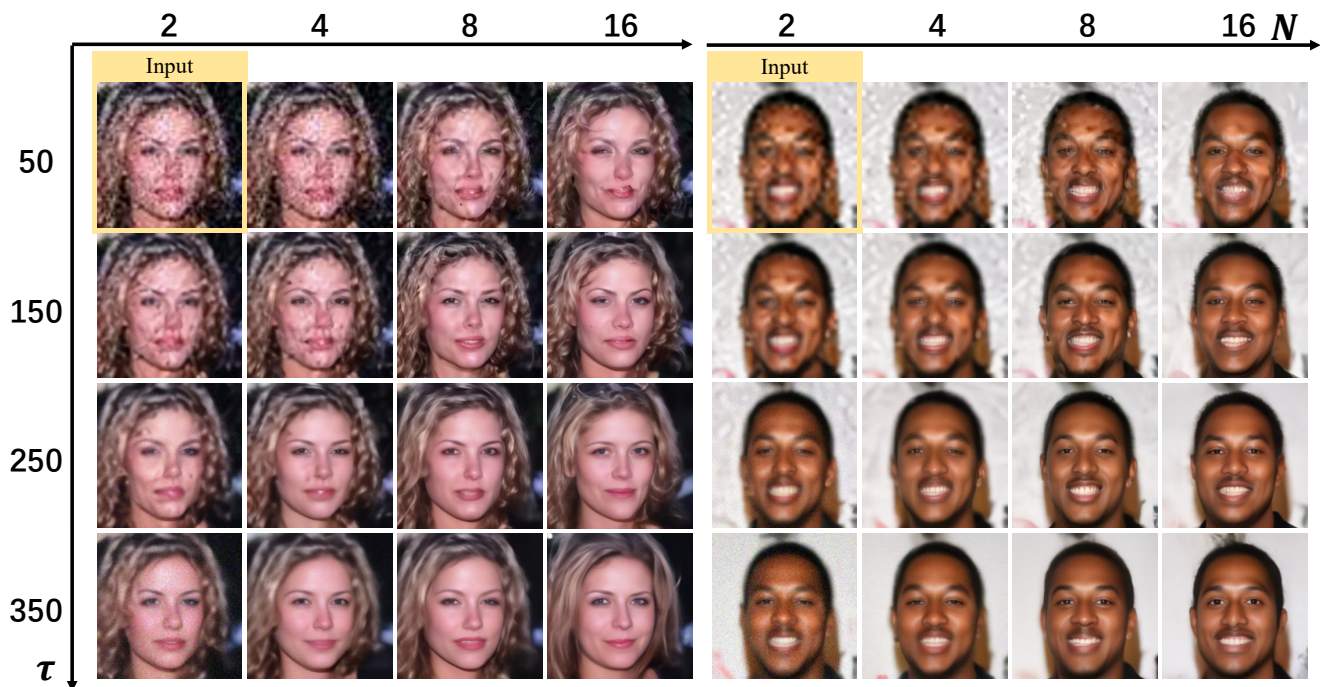


Figure 8. **DR2 outputs with various controlling parameters.** Input images are denoted with yellow boxes. Samples obtained from bigger N and τ have lower fidelity but contain fewer degradation. Whiles samples from smaller N and τ are more similar to the input but still contain degradation.

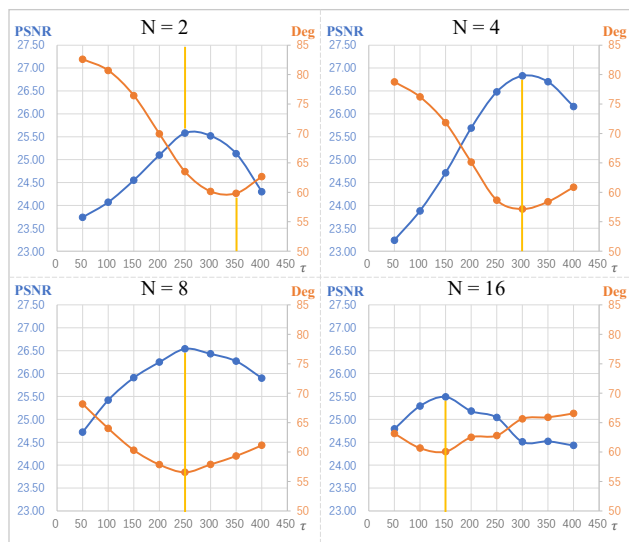


Figure 9. **Quantitative evaluation of DR2 output with various controlling parameters.** On different settings, both PSNR \uparrow and Deg \downarrow first get better then worse as τ increase. The golden line indicates the best result for each metric.

in our practice) and leave the upsampling for enhancement module for faster speed.

Another major limitation of our proposed DR2 is the manual choosing for controlling parameters N and τ . As a future work, we are exploring whether image quality assessment scores (like NIQE) can be used to develop an au-

tomatic search algorithms for N and τ .

Furthermore, for inputs with slight degradation, DR2 is less necessary because previous methods can also be effective and faster. And in extreme cases where input images contains very slight degradation or even no degradation, DR2 transformation may remove details in the inputs, but that is not common cases for blind face restoration.

5. Conclusion

We propose the DR2E, a two-stage blind face restoration framework that leverages a pretrained diffusion model that generates high-quality faces to remove degradation from inputs, and an enhancement module for detail restoration and up sampling. In the first stage, DR2 removes degradation by using diffused low-quality information as conditions to guide DR2's generative process. The entire framework requires no synthetically degraded data for training. Extensive comparisons demonstrate the strong robustness and high restoration quality of our DR2E framework.

Acknowledgements

This work is supported by National Natural Science Foundation of China (62271308), STCSM (22511105700, 22DZ2229005), 111 plan (BP0719010), and State Key Laboratory of UHD Video and Audio Production and Presentation.

References

- [1] Qingxing Cao, Liang Lin, Yukai Shi, Xiaodan Liang, and Guanbin Li. Attention-aware face hallucination via deep reinforcement learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 690–698, 2017. [2](#)
- [2] Kelvin CK Chan, Xintao Wang, Xiangyu Xu, Jinwei Gu, and Chen Change Loy. Glean: Generative latent bank for large-factor image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14245–14254, 2021. [2](#)
- [3] Chaofeng Chen, Dihong Gong, Hao Wang, Zhifeng Li, and Kwan-Yee K Wong. Learning spatial attention for face super-resolution. *IEEE Transactions on Image Processing*, 30:1219–1231, 2020. [4](#), [5](#), [6](#)
- [4] Chaofeng Chen, Xiaoming Li, Lingbo Yang, Xianhui Lin, Lei Zhang, and Kwan-Yee K Wong. Progressive semantic-aware style transformation for blind face restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11896–11905, 2021. [1](#), [2](#)
- [5] Yu Chen, Ying Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2492–2501, 2018. [1](#), [2](#), [3](#)
- [6] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14347–14356. IEEE, 2021. [2](#), [4](#), [5](#), [6](#), [7](#)
- [7] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. [7](#)
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. [6](#)
- [9] Berk Dogan, Shuhang Gu, and Radu Timofte. Exemplar guided face image super-resolution without facial landmarks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. [1](#), [2](#)
- [10] Chao Dong, Yubin Deng, Chen Change Loy, and Xiaoou Tang. Compression artifacts reduction by a deep convolutional network. In *Proceedings of the IEEE international conference on computer vision*, pages 576–584, 2015. [1](#)
- [11] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. [1](#), [3](#)
- [12] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. [2](#), [4](#)
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [2](#)
- [14] Yuchao Gu, Xintao Wang, Liangbin Xie, Chao Dong, Gen Li, Ying Shan, and Ming-Ming Cheng. Vqfr: Blind face restoration with vector-quantized dictionary and parallel decoder. In *ECCV*, 2022. [1](#), [2](#), [4](#), [5](#), [6](#)
- [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. [6](#)
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. [2](#), [3](#)
- [17] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in 'Real-Life' Images: detection, alignment, and recognition*, 2008. [5](#)
- [18] Huaibo Huang, Ran He, Zhenan Sun, and Tieniu Tan. Wavelet-srnet: A wavelet-based cnn for multi-scale face super resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1689–1697, 2017. [2](#)
- [19] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. [2](#)
- [20] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. [2](#), [5](#)
- [21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. [2](#)
- [22] Deokyun Kim, Minseon Kim, Gihyun Kwon, and Dae-Shik Kim. Progressive face super-resolution via attention to facial landmark. *arXiv preprint arXiv:1908.08239*, 2019. [2](#)
- [23] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8183–8192, 2018. [1](#)
- [24] Xiaoming Li, Chaofeng Chen, Shangchen Zhou, Xianhui Lin, Wangmeng Zuo, and Lei Zhang. Blind face restoration via deep multi-scale component dictionaries. In *European Conference on Computer Vision*, pages 399–415. Springer, 2020. [1](#), [2](#), [5](#), [6](#)
- [25] Xiaoming Li, Wenyu Li, Dongwei Ren, Hongzhi Zhang, Meng Wang, and Wangmeng Zuo. Enhanced blind face restoration with multi-exemplar images and adaptive spatial feature fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2706–2715, 2020. [1](#), [2](#)
- [26] Xiaoming Li, Ming Liu, Yuting Ye, Wangmeng Zuo, Liang Lin, and Ruigang Yang. Learning warped guidance for blind face restoration. In *Proceedings of the European conference on computer vision (ECCV)*, pages 272–289, 2018. [1](#), [2](#)

- [27] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 1, 3
- [28] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 5
- [29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2
- [30] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2437–2445, 2020. 2
- [31] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2
- [32] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019. 2, 4
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2, 6
- [34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 3
- [35] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2, 6
- [36] Ziyi Shen, Wei-Sheng Lai, Tingfa Xu, Jan Kautz, and Ming-Hsuan Yang. Deep semantic face deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8260–8269, 2018. 1, 2
- [37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 2, 3
- [38] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 2, 4
- [39] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9168–9178, 2021. 1, 2, 5, 6
- [40] Xiangyu Xu, Deqing Sun, Jinshan Pan, Yujin Zhang, Hanspeter Pfister, and Ming-Hsuan Yang. Learning to super-resolve blurry face and text images. In *Proceedings of the IEEE international conference on computer vision*, pages 251–260, 2017. 2
- [41] Shuo Yang, Ping Luo, Chen-Change Loy, and Xiaoou Tang. Wider face: A face detection benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5525–5533, 2016. 5
- [42] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 672–681, 2021. 1, 2, 5, 6
- [43] Xin Yu, Basura Fernando, Bernard Ghanem, Fatih Porikli, and Richard Hartley. Face super-resolution guided by facial component heatmaps. In *Proceedings of the European conference on computer vision (ECCV)*, pages 217–233, 2018. 1, 2
- [44] Xin Yu, Basura Fernando, Richard Hartley, and Fatih Porikli. Super-resolving very low-resolution face images with supplementary attributes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 908–917, 2018. 2
- [45] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155, 2017. 1
- [46] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6
- [47] Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Truncated diffusion probabilistic models and diffusion-based adversarial auto-encoders. *arXiv preprint arXiv:2202.09671*, 2022. 3
- [48] Shangchen Zhou, Kelvin C.K. Chan, Chongyi Li, and Chen Change Loy. Towards robust blind face restoration with codebook lookup transformer. In *NeurIPS*, 2022. 1, 2, 5, 6
- [49] Shizhan Zhu, Sifei Liu, Chen Change Loy, and Xiaoou Tang. Deep cascaded bi-network for face hallucination. In *European conference on computer vision*, pages 614–630. Springer, 2016. 2