# WaveDM: Wavelet-Based Diffusion Models for Image Restoration

Yi Huang\*, Jiancheng Huang\*, Jianzhuang Liu, *Senior Member, IEEE,*
Mingfu Yan, Yu Dong, Jiaxi Lyu, Chaoqi Chen, Shifeng Chen[†]

*Abstract*—Latest diffusion-based methods for many image restoration tasks outperform traditional models, but they encounter the long-time inference problem. To tackle it, this paper proposes a Wavelet-Based Diffusion Model (WaveDM). WaveDM learns the distribution of clean images in the wavelet domain conditioned on the wavelet spectrum of degraded images after wavelet transform, which is more time-saving in each step of sampling than modeling in the spatial domain. To ensure restoration performance, a unique training strategy is proposed where the low-frequency and high-frequency spectrums are learned using distinct modules. In addition, an Efficient Conditional Sampling (ECS) strategy is developed from experiments, which reduces the number of total sampling steps to around 5. Evaluations on twelve benchmark datasets including image raindrop removal, rain steaks removal, dehazing, defocus deblurring, demoiréing, and denoising demonstrate that WaveDM achieves state-of-the-art performance with the efficiency that is comparable to traditional one-pass methods and over $100\times$ faster than existing image restoration methods using vanilla diffusion models. The code is available at https://github.com/stayalive16/WaveDM.

*Index Terms*—Diffusion models, image restoration, wavelet transform

## I. INTRODUCTION

IMAGE restoration, aiming to remove degradations (e.g., blur, raindrops, moiré, noise and so on) from a degraded image to generate a high-quality one, has raised great attention in computer vision research. Most previous methods depend on strong priors or estimate the degradation functions for specific tasks [21], [34], [67], [99], [103]. With the development of deep learning, deep neural network-driven methods have become the mainstay. These methodologies are primarily built on architectures like Convolutional Neural Networks (CNNs) [5], [16], [31], [68], [118], [119] and Transformers [26], [43], [47], [106], [117]. However, some of these deep learning models, generally relying on regression techniques, tend to yield results that are usually over-smoothing and lose subtle details. On the other hand, unsupervised methods [19], [30], [63], which are implemented without labeled data, promise impressive generalizability, especially in scenarios not seen during training. However, the absence of explicit guidance

sometimes results in outputs that may be over-enhanced in colors or contain amplified noise.

Another popular approach is through task-specific generative modeling, frequently leveraging Generative Adversarial Networks (GANs) [17], [45], [77], [120], [127]. These generative models aim to capture the latent data distribution of clean images and apply this prior to the degraded samples. While showing powerful generalization capabilities, GAN-based restoration techniques have their own drawbacks. The use of adversarial losses often induces artifacts that are absent in the original clean images, introducing distortions. Besides, the instability of GAN training further intensifies this challenge, and in certain scenarios, can even lead to mode collapse [90]. Another type, flow-based methods [59], [108], directly accounts for the ill-posed problem with an invertible encoder, which maps clean images to the flow-space latents conditioned on degraded inputs. However, the need for a strict bijection between latent and data spaces adds to their complexity.

Recently, diffusion models [13], [23], [72], [82], [85] have come into the spotlight. Their achievements span various computer vision tasks such as conditional image generation [53], [82], [89], image super-resolution [42], [86], image-to-image translation [84], [88], [104], [130], and face restoration [71], [75], [115]. These models have become popular because of many distinct benefits diffusion models possess. One is the outstanding generative capability as diffusion models can better capture the data distribution compared to other approaches such as GANs. Furthermore, diffusion models excel in countering diverse degradations, ranging from noise and blur to more complex corruptions due to their ability in modeling intricate data distributions. In addition, diffusion models are inherently resistant to mode collapse, ensuring comprehensive data distribution coverage [22], [90], [95]. This leads to more stable training, mitigating chances of unpredictable outputs and affirming their reliability in restoration. However, the biggest challenge among them is the heavy computational burden as diffusion models usually require many steps of sampling. Earlier works [13], [23] start from generating a low-resolution image and gradually upsample it through pretrained super-resolution models to reduce the processing time in each step. Rombach et al. [82] apply the diffusion models in the latent space of powerful pretrained autoencoders for high-resolution image synthesis. Some other works mainly focus on reducing the evaluation steps by accelerated deterministic implicit sampling [92], knowledge distillation [66], [87], changing the diffusion strategy [60] and reformulating the solution to the diffusion ordinary differential equations [56],

\*: Equal contribution. [†]: Corresponding author.

Yi Huang, Jiancheng Huang, Jianzhuang Liu, Mingfu Yan, Yu Dong, Jiaxi Lyu and Shifeng Chen are all with Shenzhen Key Lab of Computer Vision and Pattern Recognition, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen, 518055, China and also with University of Chinese Academy of Sciences, Beijing, 100039, China (e-mail: yi.huang; jc.huang; jz.liu; mf.yan; yu.dong; jx.lv1; shifeng.chen@siat.ac.cn).

Chaoqi Chen is with The University of Hong Kong (email: cqchen1994@gmail.com).

[57]. However, they are still restricted to practical applications of high-resolution image restoration.

Recently, Kawar et al. propose Denoising Diffusion Restoration Models (DDRM) [32] which takes advantage of a pre-trained diffusion model for solving linear inverse restoration problems without extra training, but it cannot handle images with nonlinear degradation. Some other works [9]–[11], [70], [84], [93], [105] seek to use diffusion models to address nonlinear inverse imaging problems, in which the forms and parameters of the degradation functions have to be known. However, the degradation models of most real-world restoration problems such as deraining cannot be obtained. Ozan et al. introduce an approach [73] to restore vision under adverse weather conditions with size-agnostic inputs by cutting images into multiple overlapping small patches. Although it can process high-resolution images with better performance than traditional one-pass methods, its computational complexity increases quadratically with the increase of image sizes. For example, one $2176 \times 1536$ image is cut into $12369$ $64 \times 64$ overlapping patches, requiring about 650 seconds for 25-step sampling on a regular GPU.

The first is to decrease the time of processing images in each step. Specifically, WaveDM learns the distribution of clean images in the wavelet domain, which is different from most of the current diffusion models that focus on the spatial domain. After wavelet transform for $n$ times, the spatial size of the original image is reduced by $1/4^n$, thus saving a lot of computation. Note that other popular transforms such as Fourier transform cannot achieve this because the size of the Fourier spectrum is the same as that of the image. Although some recent works [18], [27], [74] also introduce wavelet transform into diffusion models for image or 3D generative tasks, to the best of our knowledge, this attempt has not been explored in addressing the restoration problems. In our model, the input images are first decomposed into multiple frequency bands using wavelet transform. In the training phase, a diffusion model is utilized to learn the distribution of low-frequency bands of clean images by perturbing them with random noise at different moments of time. In addition, a lightweight high-frequency refinement module is constructed to provide the high-frequency bands, which also serve as the essential condition. The sampling starts from a random Gaussian noise to predict the low-frequency bands through a reverse diffusion process, which are then combined with the output from the high-frequency refinement module to generate a clean image through inverse wavelet transform.

The second scheme of acceleration is to reduce the total sampling steps, which is realized by an Efficient Conditional Sampling (ECS) strategy we obtain from experiments. ECS follows the same sampling procedure as the deterministic implicit sampling [92] during the initial sampling period and then stops at an intermediate moment to predict clean images directly instead of completing all the sampling steps. During this procedure, the degraded images serve as the essential conditions that provide strong priors such as global tone and spatial structure to remove the noise till the end. Due to its simple implementation, ECS can further reduce the sampling steps to as few as 4 without extra computation. Additionally,

experimental results on several datasets show that ECS is also capable of maintaining or even improving the restoration performance by setting the intermediate moment reasonably.

The main contributions of this work are summarized as follows:

- A wavelet-based diffusion model is proposed to learn the distribution of clean images in the wavelet domain, which dramatically reduces the computational expenses typically encountered in the spatial domain.
- A unique training strategy is proposed where the low-frequency and high-frequency spectrums are learned using distinct modules, which facilitates the effective restoration of degraded images.
- An efficient conditional sampling strategy is found based on our experiments to reduce the number of sampling steps to around 5 without extra computation, while maintaining the restoration performance compared to other diffusion-based methods using 25 steps or more.
- Comprehensive experiments conducted on twelve restoration benchmark datasets verify that our WaveDM achieves state-of-the-art performance with the efficiency that is comparable to traditional one-pass methods and over $100\times$ faster than existing diffusion-based models.

## II. RELATED WORK

### A. Image Restoration

Earlier restoration methods [21], [34], [67], [99] mainly rely on seeking strong priors for specific tasks. In recent years, deep neural works are widely used for general image restoration owing to their superb performance. These learning-based approaches usually require a specific model architecture constructed by CNN [5], [16], [50], [76], [110], [112], [118], [119], [121], [129], [131] or Transformer [43], [47], [106], [117]. Most convolutional encoder-decoder designs [2], [6], [7], [36], [44], [113] could be viewed as variants of a classical solution, U-Net [83], the effectiveness of which has been validated for their hierarchical representations while keeping computationally efficient. Extensively, spatial and channel attentions are also injected in it to capture some key information thus boosting the performance. The Vision Transformer [15], [55], first introduced for image classification, is capable of building strong relationships between image patches due to its self-attention mechanism. Naturally, a lot of transformer-based works are also studied for the low-level vision tasks like super-resolution [47], [49], [111], denoising [43], [106], [117], deraining [109], colorization [35], etc. Different from them, this paper aims to tackle the restoration problem from the view of generative modeling, implemented by a wavelet-based diffusion model.

### B. Diffusion Models

Diffusion models, a new type of generative models, are inspired by non-equilibrium thermodynamics. They learn to reverse the forward process of sequentially corrupting data samples with additive random noise following the Markov chain, until reconstructing the desired data that matches the

source data distribution from noise. Previous diffusion models can be roughly classified into diffusion based [90] and score-matching based [28], [102]. Following them, denoising diffusion probabilistic models [22], [72] and noise-conditioned score network [94], [95], [97] are proposed to synthesize high-quality images, respectively.

*1) Diffusion Models in Low-Level Vision Tasks:* Recently, diffusion-based models show great potential in various computer vision tasks under conditions such as class-conditioned image synthesis with and without classifier guidance [13], [24], [33], image inpainting [58], super-resolution [40], [86], deblurring [107], and image-to-image translation (e.g., colorization and style transfer) [8], [37], [84], [104], [128]. Similarly, the applications following the score-based conditional modeling are also widely explored [12], [65]. Beyond synthesis, some works apply diffusion models for image restoration. Most of the restoration methods are trained either on large-scale datasets or with samples that come from some specific types (e.g., faces [71], [75]) to obtain high-quality generation performance. However, they may somewhat change the original spatial structure of conditional degraded images. Kawar et al. [32] propose DDRM to solve linear inverse image restoration problems, but it cannot be adapted to the inversion of nonlinear degradation. Some works [9]–[11], [70], [84], [93], [105] use diffusion models to address nonlinear inverse imaging problems, in which the forms and parameters of the degradation functions have be to known. Ozan et al. [73] propose patch-based diffusion models, which is the first diffusion-based work that achieves state-of-the-art performance on three real-world blind restoration tasks in terms of pixel-wise evaluation metrics such as PSNR. However, the main limitation of it is the much longer inference time than traditional one-pass methods due to a large amount of image patches and many sampling steps. This paper aims to solve this problem through the wavelet-based diffusion model with an efficient conditional sampling strategy, preserving the state-of-the-art restoration performance simultaneously.

*2) Accelerating Diffusion Models:* Though diffusion models are capable of generating high-quality images, their iterative sampling procedure usually results in long inference time. Song et al. [92] propose the deterministic implicit sampling that requires only 25 steps. Lu et al. [56] reformulate the exact solution to the diffusion ordinary differential equations (ODEs) and propose a fast dedicated high-order solver for diffusion ODE speedup using around 10 steps. Ma et al. [61] investigate this problem by viewing the diffusion sampling process as a Metropolis adjusted Langevin algorithm and introduce a model-agnostic preconditioned diffusion sampling that leverages matrix preconditioning, which accelerates vanilla diffusion models by up to $29\times$. Lyu et al. [60] start the reverse denoising process from a non-Gaussian distribution, which enables stopping the diffusion process early where only a few initial diffusion steps are considered. However, it requires an extra generative model (e.g., GAN or VAE) to approximate the real data distribution to start sampling. Different from them, our work focuses on accelerating the conditional image restoration diffusion model from two aspects: reducing the processing of each step implemented by a wavelet-based

diffusion model, and reducing the number of total sampling steps by an efficient conditional sampling strategy without extra training.

### C. Wavelet Transform-Based Methods

Wavelet transform has been widely explored in computer vision tasks, especially combined with deep neural networks. For example, Liu et al. [51] propose a multilevel Wavelet-CNN to enlarge receptive fields with a better trade-off between efficiency and restoration performance via multi-level wavelets. Liu et al. [50] design a wavelet-based dual-branch network with a spatial attention mechanism for image demoiréing. Xin et al. [110] first decompose the low-resolution image into a series of wavelet coefficients (WCs) and then use a CNN to predict the corresponding series of high-resolution WCs, which are then utilized to reconstruct the high-resolution image. Li et al. [43] propose an efficient wavelet transformer for image denoising. It is the first attempt to utilize Transformer in the wavelet domain, implemented by an efficient multi-level feature aggregation module, thus significantly reducing the device resource consumption of the conventional Transformer model. All the methods mentioned above combine wavelet transform with deep neural networks like CNNs and Transformers by designing task-specific network structures without using diffusion models, while our method combines a diffusion model with wavelet transform for various image restoration tasks by employing the general convolutional U-Net architecture. Our approach can achieve superior performance on multiple image restoration tasks while maintaining comparable processing efficiency.

In recent years, wavelet diffusion-based methods have emerged as a prominent approach, particularly in the realm of image generation. As evident from the works of Phung et al. [74] and Guth et al. [18], they focus on leveraging wavelet diffusion for image synthesis. Hui et al. [27] extend the framework to 3D shape generation. These approaches, while significant, primarily target generation tasks. Different from them, our WaveDM is architected with deliberate design, leveraging the diffusion principle innovatively for image restoration.

## III. PRELIMINARIES

### A. Denoising Diffusion Probabilistic Models

Denoising Diffusion Probabilistic Models (DDPMs) [22], [72] are a class of generative models that work by destroying training data through the successive addition of Gaussian noise, and then learning to recover the data by reversing this noising process. During training, the forward noising process follows the Markov chain that transforms a data sample from the real data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ into a sequence of noisy samples $\mathbf{x}_t$ in $T$ steps with a variance schedule $\beta_1, \ldots, \beta_T$:

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

Diffusion models learn to reverse the above process through a joint distribution $p_\theta(\mathbf{x}_{0:T})$ that follows the Markov chain

with parameters $\theta$, starting at a noisy sample from a standard Gaussian distribution $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \tag{2}$$

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)). \tag{3}$$

The parameters $\theta$ are usually optimized by a neural network that predicts $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ and $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)$ of Gaussian distributions, which is simplified by predicting noise vectors $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ with the following objective [22]:

$$\begin{aligned}
\mathbb{E}_{q(\mathbf{x}_0)}[-\log p_\theta(\mathbf{x}_0)] &\leq \mathbb{E}_q\left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \right] \\
&= \mathbb{E}_q\Big[ \underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_T \mid \mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} \underbrace{- \log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1)}_{L_0} \\
&\quad + \sum_{t>1} \underbrace{D_{\mathrm{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t))}_{L_{t-1}} \Big].
\end{aligned} \tag{4}$$

Obviously, the $L_{t-1}$ term actually trains the network to perform one reverse diffusion step. As reported by [22], the optimization of $L_{t-1}$ can be converted to training a network $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ that estimates the mean value of the posterior distribution $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$. Furthermore, the model can instead be trained to predict the noise vector $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ using an alternative reparameterization of the reverse process by:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right), \tag{5}$$

where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$. As a result, the training objective is transformed into a re-weighted simplified form given as:

$$L_{simple} = \mathbb{E}_{\mathbf{x}_0, t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})}\left[ \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right]. \tag{6}$$

Consequently, the sampling phase with the learned parameterized Gaussian transitions $p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$ can start from $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ by:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}, \tag{7}$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\alpha_t = 1 - \beta_t$, and $\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$.

### B. Deterministic Implicit Sampling

Denoising Diffusion Implicit Models (DDIMs) [92] generalize DDPMs to obtain the same training objective as Eq. 6 by defining a non-Markovian diffusion process:

$$q_\sigma(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \sigma_t^2 \mathbf{I}). \tag{8}$$

By setting $\sigma_t^2 = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t$, the forward process becomes Markovian and remains the same as DDPMs.

A deterministic implicit sampling (also called DDIM sampling) is implemented by setting $\sigma_t^2 = 0$, and thus the sampling process based on Eq. 8 can be accomplished by:

$$\begin{aligned}
\mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\left( \frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} \right) \\
&\quad + \sqrt{1-\bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t),
\end{aligned} \tag{9}$$

which enables a faster sampling procedure. Specifically, DDIMs replace the complete reverse sampling sequence $\mathbf{x}_T, \mathbf{x}_{T-1}, \dots, \mathbf{x}_1, \mathbf{x}_0$ with one of its sub-sequence $\mathbf{x}_T, \mathbf{x}_{\tau_S}, \mathbf{x}_{\tau_{S-1}}, \dots, \mathbf{x}_{\tau_1}$ which can be obtained by:

$$\tau_i = (i-1) \cdot T/S, \tag{10}$$

where $S$ denotes the total sampling steps for acceleration. Thus, the faster DDIM sampling procedure is formulated as:

$$\begin{aligned}
\mathbf{x}_{t-1} &= \sqrt{\bar{\alpha}_{t-1}}\left( \frac{\mathbf{x}_t - \sqrt{1-\bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}} \right) \\
&\quad + \sqrt{1-\bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t), \quad t = T, \tau_S, \dots, \tau_1.
\end{aligned} \tag{11}$$

### C. Wavelet Transform

*1) 2D Discrete Wavelet Transform (2D DWT):* Given an image $I \in \mathbb{R}^{H \times W \times C}$, where $M \times N$ is the spatial size, $C$ is the number of channesl, the 2D DWT decomposes the image into four sub-bands:

$$I_{LL}, I_{LH}, I_{HL}, I_{HH} = \mathrm{DWT}_{2D}(I). \tag{12}$$

The sub-band $I_{LL}$ represents the approximation coefficients and has a size of $\frac{M}{2} \times \frac{N}{2} \times C$. The other sub-bands $I_{LH}$, $I_{HL}$, and $I_{HH}$ correspond to the horizontal, vertical, and diagonal detail coefficients, respectively, and each of them possesses a size of $\frac{M}{2} \times \frac{N}{2} \times C$.

For multi-level wavelet decomposition, the DWT is recursively applied to the $I_{LL}$ sub-band from the previous level. After $k$ decompositions, the size of the $I_{LL}$ sub-band reduces to $\frac{M}{2^k} \times \frac{N}{2^k} \times C$.

The wavelet used for decomposition can be of various types, such as the Haar wavelet, which provides a simple and effective basis for image decomposition.

*2) 2D Inverse Discrete Wavelet Transform (2D IDWT):* Starting with the four sub-bands, the original image is reconstructed:

$$I' = \mathrm{IDWT}_{2D}(I_{LL}, I_{LH}, I_{HL}, I_{HH}), \tag{13}$$

where $I' \in \mathbb{R}^{H \times W \times C}$. The process of multi-level reconstruction begins from the deepest decomposition level and sequentially moves towards the first level, eventually yielding a reconstructed image $I'$ of the original size.

*3) 2D Full Wavelet Packet Transform (2D FWPT):* Unlike the 2D DWT, which only recursively decomposes the $I_{LL}$ sub-band, the 2D Full Wavelet Packet Transform (2D FWPT) exhaustively decomposes every sub-band at each level.

For a single level decomposition of an image $I \in \mathbb{R}^{H \times W \times C}$, the 2D FWPT yields 4 sub-bands:

$$\{I_{i,j}\}_{i,j \in \{L,H\}} = \mathrm{FWPT}_{2D}(I), \tag{14}$$

where each sub-band $I_{i,j} \in \mathbb{R}^{\frac{M}{2} \times \frac{N}{2} \times C}$.

For a 2-level FWPT, each of the initial sub-bands is further decomposed, leading to a total of 16 sub-bands, each of size $\frac{M}{4} \times \frac{N}{4} \times C$.

The benefit of this exhaustive decomposition is that all the sub-bands at each level have the same spatial dimension, allowing easier concatenation and analysis of frequency details in a structured manner.
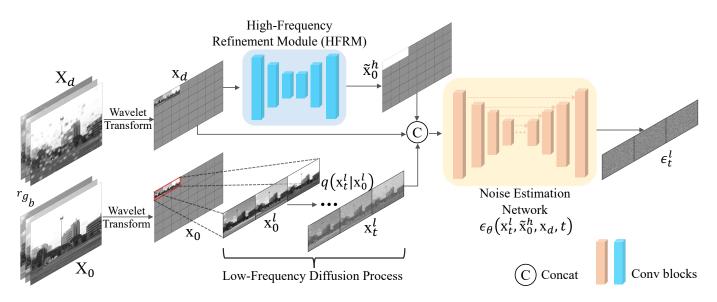
Fig. 1. Training of the wavelet-based diffusion model (WaveDM) for image restoration, where $\mathbf{X}_d$ and $\mathbf{X}_0$ stand for a pair of RGB degraded and clean images. $\mathbf{x}_d$ and $\mathbf{x}_0$ are the wavelet spectrum of $\mathbf{X}_d$ and $\mathbf{X}_0$ after the Haar wavelet transform, respectively. $\mathbf{x}_t^l$ is the diffusion result of the low-frequency spectrum $\mathbf{x}_0^l$ extracted from the first three bands of $\mathbf{x}_0$. $\tilde{\mathbf{x}}_0^h$ denotes the high-frequency spectrum of the clean image based on $\mathbf{x}_d$ with the HFRM. $\mathbf{x}_d$, $\tilde{\mathbf{x}}_0^h$ and $\mathbf{x}_t^l$ are concatenated together as input to the noise estimation network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, \tilde{\mathbf{x}}_0^h, \mathbf{x}_d, t)$ to predict the noise $\boldsymbol{\epsilon}_t^l$ at all time moments.

*4) 2D Inverse Full Wavelet Packet Transform (2D IFWPT):* Given the sub-bands obtained from the 2D FWPT, the 2D IFWPT reconstructs the original image. For a 1-level FWPT:

$$I' = \text{IFWPT}_{2D}(\{I_{i,j}\}_{i,j\in\{L,H\}}), \tag{15}$$

where the reconstructed image $I' \in \mathbb{R}^{H\times W\times 1}$. The reconstruction process, similar to 2D IDWT, starts from the deepest decomposition level and works its way up to the first level, combining all sub-bands together to form the original image.

## IV. METHOD

### A. Overview

Recently, diffusion models are increasingly favored over alternatives like GANs in image restoration due to their better ability to capture complex data distributions and their inherently stable training processes. However, current methods, such as [73] and [32], apply diffusion models directly in the spatial domain, resulting in long inference time. To mitigate this computational challenge, we leverage the wavelet transform's capability for image size reduction with no information loss and frequency sub-band separation. Consequently, we propose a wavelet-based diffusion model (WaveDM) that learns the distribution of clean images in the wavelet domain, where the low-frequency and high-frequency spectrums are learned using distinct modules for restoration quality.

Fig. 1 depicts the WaveDM training procedure. The pipeline consists of three primary parts: the High-Frequency Refinement Module (HFRM), the low-frequency diffusion process, and the noise estimation network. Firstly, both degraded and clean image pairs are transformed to the wavelet domain using the 2D FWPT, with the high and low-frequency details extracted and resolution reduced. The full wavelet spectrum of the degraded image is taken as input to HFRM to estimate the clean image's high-frequency spectrum. Concurrently, we

add Gaussian noise to the low-frequency spectrum of the clean image. The noisy low-frequency spectrum, concatenated with the input and output of HFRM, is then sent to the noise estimation network for noise prediction. Comprehensive training details are described in Section IV-B.

Fig. 2 describes the sampling process of WaveDM. First, 2D FWPT captures the wavelet spectrum of a degraded image, serving as HFRM's input. The sampling starts from the concatenation of HFRM's input and output with a Gaussian noise, which is then sent to the noise estimation network, yielding a noisy low-frequency wavelet spectrum at the first step. This process iterates using the efficient conditional sampling strategy, described in Section IV-B, to produce a clean low-frequency wavelet spectrum at the end of sampling. Then the final clean RGB image is obtained from the concatenation of this spectrum with HFRM's output, followed by 2D IFWPT.

### B. Training of WaveDM

As shown in Fig. 1, given a degraded image $\mathbf{X}_d \in \mathbb{R}^{H\times W\times 3}$ and its corresponding ground truth $\mathbf{X}_0 \in \mathbb{R}^{H\times W\times 3}$, we employ a 2-level 2D FWPT using the Haar wavelet. The Haar wavelet transform iteratively applies low-pass and high-pass decomposition filters, coupled with downsampling, to compute the wavelet coefficients. Specifically, the low-pass filter, with coefficients $(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$, captures the average information, while the high-pass filter, with coefficients $(\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}})$, focuses on the details or transitions in the image. The transformation process begins by applying these filters to each row of the image, resulting in two intermediate forms. These forms are then subjected to the same filter application along their columns, decomposing the original image into four distinct sub-bands: LL (averaged information), LH (details along columns), HL (details along rows), and HH (details in both rows and columns). For the 2-level 2D FWPT, this decomposition process is recursively applied to all sub-bands. As a

result, each image is transformed into the wavelet spectrum $\mathbf{x}_d, \mathbf{x}_0 \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 48}$, consisting of 48 bands with the same spatial dimension, which can be represented as:

$$\begin{aligned} \mathbf{x}_d &= \text{FWPT}_{2D}(\mathbf{X}_d), \\ \mathbf{x}_0 &= \text{FWPT}_{2D}(\mathbf{X}_0). \end{aligned} \quad (16)$$

Instead of adopting a naive diffusion approach in the wavelet domain, which directly corrupts all wavelet bands of the clean image using additive Gaussian noise and then reversing the process during sampling, we introduce an optimized approach. Essential experiments, discussed in Section V-B3, demonstrated the ineffectiveness of the naive method. Drawing from experimental insights, only the low-frequency spectrum $\mathbf{x}_0^l \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3}$ which is derived from the first three bands of the clean image wavelet spectrum $\mathbf{x}_0$, is corrupted with Gaussian random noise. This corruption follows a forward diffusion process defined as: $q(\mathbf{x}_t^l \mid \mathbf{x}_0^l) = \mathcal{N}(\mathbf{x}_t^l; \sqrt{\bar{\alpha}_t} \mathbf{x}_0^l, (1 - \bar{\alpha}_t) \mathbf{I}), t = 1, 2, \ldots, T$.

Additionally, recognizing the importance of high-frequency information that remains unmodeled in the low-frequency spectrum, we design a High Frequency Refinement Module (HFRM). This lightweight module estimates the high-frequency spectrum $\tilde{\mathbf{x}}_0^h \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 45}$ of the clean image $\mathbf{x}_0$ from $\mathbf{x}_d$ in a single pass. This can be presented as:

$$\tilde{\mathbf{x}}_0^h = \text{HFRM}(\mathbf{x}_d). \quad (17)$$

During each step, the degraded image's wavelet spectrum $\mathbf{x}_d$ and the estimated high-frequency spectrum $\tilde{\mathbf{x}}_0^h$ serve as conditions to model the low-frequency spectrum distribution of clean images. Specifically, the concatenated diffusion result $\mathbf{x}_t^l, t = 1, 2, \ldots, T$, $\tilde{\mathbf{x}}_0^h$, and $\mathbf{x}_d$ across channels feed into the noise estimation network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, \tilde{\mathbf{x}}_0^h, \mathbf{x}_d, t)$. By transitioning the diffusion model from the spatial domain to the wavelet domain using 2D FWPT, we achieve a spatial size reduction of $1/16$ for input images, leading to a substantial speedup in processing.

For training, we employ a combined objective function to optimize the diffusion process in the wavelet domain and refine the high-frequency bands. Specifically, the primary objective $L_{simple}$, as defined in Eq. 6, is utilized to optimize $\boldsymbol{\epsilon}_\theta$. The HFRM, which is independent of the variable $t$, is trained using the objective $L_1 = ||\tilde{\mathbf{x}}_0^h - \mathbf{x}_0^h||_1$, where $\mathbf{x}_0^h$ denotes the high-frequency bands of $\mathbf{x}_0$. The total training loss is given by:

$$L_{total} = L_{simple} + \lambda L_1. \quad (18)$$

where $\lambda$ acts as a weighting hyperparameter.

### C. Sampling of WaveDM

The WaveDM framework, after training, adopts a sequential inference approach in processing the wavelet bands. Firstly, high-frequency wavelet bands are predicted. Subsequently, the low-frequency wavelet bands are sampled. These combined bands are then utilized to generate a clean RGB image using 2D IFWPT. This entire operation is represented in Fig. 2.

For a degraded image denoted by $\mathbf{X}_d \in \mathbb{R}^{H \times W \times 3}$, we apply a 2-level 2D FWPT implemented by the Haar wavelet. This transformation uses the same filters as in the training of
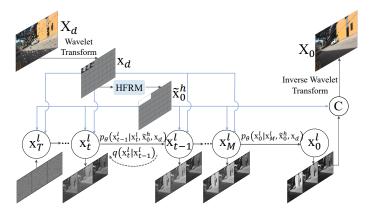


Fig. 2. Overview of WaveDM with ECS. $q(\mathbf{x}_t^l \mid \mathbf{x}_{t-1}^l)$ stands for the forward diffusion (dashed line). The sampling process $p_\theta(\mathbf{x}_{t-1}^l \mid \mathbf{x}_t^l, \tilde{\mathbf{x}}_0^h, \mathbf{x}_d)$ (solid lines) starts from a standard Gaussian noise $\mathbf{x}_T^l \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to generate the low-frequency spectrum of the clean image, where $\mathbf{x}_d$ and $\tilde{\mathbf{x}}_0^h$ serve as conditions (blue solid lines) from step $T$ to step $M$. Then the intermediate result $\mathbf{x}_M^l$ is utilized to predict the low-frequency spectrum $\mathbf{x}_0^l$ of the clean image directly, followed by inverse wavelet transform that turns the concatenation of $\tilde{\mathbf{x}}_0^h$ and $\mathbf{x}_0^l$ into a clean RGB image $\mathbf{X}_0$.

WaveDM. The output of this operation is the wavelet spectrum $\mathbf{x}_d \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 48}$ as outlined in Eq. 16. Then the spectrum $\mathbf{x}_d$, when fed into HFRM, produces the predicted high-frequency bands $\tilde{\mathbf{x}}_0^h$ of the restored image, which is described in Eq. 17.

To estimate the low-frequency wavelet bands of the restored image, both $\mathbf{x}_d$ and $\tilde{\mathbf{x}}_0^h$ are employed. This operation conventionally begins with a random Gaussian noise generation, denoted as $\mathbf{x}_T^l \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 3} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ at timestep $T$. Usually, this noise acts as a starting point for the DDIM sampling which samples a clean low-frequency spectrum $\mathbf{x}_0^l$. The procedure using DDIM sampling is methodically detailed in Eq. 19:

$$\begin{aligned} \mathbf{x}_{t-1}^l &= \sqrt{\bar{\alpha}_{t-1}} \left( \frac{\mathbf{x}_t^l - \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, t)}{\sqrt{\bar{\alpha}_t}} \right) \\ &+ \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, t), \quad t = T, \tau_S, \ldots, \tau_1, \end{aligned} \quad (19)$$

where the number of total sampling steps is $S$.

However, in our experimental observations, using DDIM sampling directly requires more than 20 steps for achieving the desired restoration performance. This inefficiency in DDIM sampling drives us to explore alternative strategies. After conducting extensive experiments, we find and develop the Efficient Conditional Sampling (ECS) strategy, the effectiveness and efficiency of which are demonstrated in Section V-B5. Not only does ECS significantly reduce the sampling steps to around 5, but also brings an enhancement in the restoration quality compared to the conventional DDIM sampling.

In the ECS methodology, instead of allowing DDIM sampling to run its full process, we strategically interrupt it at a specific intermediate step denoted as $M$. At this moment, rather than continuing with the usual diffusion sampling iterations, we leverage the information contained in the noisy spectrum $\mathbf{x}_M^l$. With this information, we compute the desired $\mathbf{x}_0^l$ directly by a portion of the DDIM equation (Eq. 19), effectively simplifying the process and mitigating the need for additional iterative steps. This ECS procedure is represented

in Eq. 20, in which the number of total sampling steps is $S(T - M)/T + 1$. It is noteworthy that for $\epsilon_\theta$, the input variables $\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h$ are concatenated channel-wise.

$$\begin{cases} \mathbf{x}_{t-1}^l = \sqrt{\bar{\alpha}_{t-1}} \left( \dfrac{\mathbf{x}_t^l - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_\theta(\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, t)}{\sqrt{\bar{\alpha}_t}} \right) \\ \qquad + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \epsilon_\theta(\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, t), \ t = T, \tau_S, \dots, M + \dfrac{T}{S}, \\ \\ \hat{\mathbf{x}}_0^l = \dfrac{\mathbf{x}_M^l - \sqrt{1 - \bar{\alpha}_M} \cdot \epsilon_\theta(\mathbf{x}_M^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, M)}{\sqrt{\bar{\alpha}_M}}. \end{cases}$$
$$(20)$$

Upon acquiring the clean low-frequency wavelet spectrum $\hat{\mathbf{x}}_0^l$ through Eq. 20, the restored clean RGB image $\mathbf{X}_0$ is obtained using 2D IFWPT. This is expressed as:

$$\mathbf{X}_0 = \text{IFWPT}_{2D}(\hat{\mathbf{x}}_0^l, \tilde{\mathbf{x}}_0^h), \qquad (21)$$

where $\tilde{\mathbf{x}}_0^h$ is the high-frequency wavelet spectrum predcited by HFRM.

## V. EXPERIMENTS

### A. Datasets and Settings

We evaluate WaveDM on twelve benchmark datasets for several image restoration tasks: (i) RainDrop [77] (861 training images and 58 testing images of size $720 \times 480$) for image raindrop removal, (ii) Outdoor-rain [45] (9000 training images and 750 testing images of size $720 \times 480$) for image rain steaks removal, (iii) SOTS-Outdoor [41] (72135 training images and 500 testing images with about $600 \times 400$ resolution) for image dehazing, (iv) DPDD [2] (350 training images of size $1680 \times 1120$ and 76 testing images of size $1664 \times 1120$) for both single-pixel and dual-pixel defocus deblurring, (v) London's Buildings [50] (561 training images and 53 testing images with about $2200 \times 1600$ resolution) for image demoiréing, (vi) SIDD [1] (about 30000 training images and 1280 testing images of size $256 \times 256$) for real image denoising, and (vii) DFWB for training with 6 benchmark datasets for testing Gaussian image denoising. Specifically, DFWB denotes the combination of DIV2K [4] (800 images), Flickr2K [48] (2650 images), WED [62] (4744 images), and BSD500 [64] (400 images).

The framework of the Patch-based Diffusion Models [73] (PatchDM) is adopted as the baseline, with which we share the same training settings. (e.g., 1000 diffusion steps with linear noise corruption strategy, sinusoidal positional encoding [101] to encode time embeddings for $t$, 2000000 training iterations, Adam optimized with a fixed learning rate of $4 \times e^{-4}$ without weight decay, and exponential moving average with a weight of 0.999 to facilitate more stable training). A similar U-Net architecture based on WideResNet [116] is used as the backbone of the noise estimation network with minor revision to adapt to the input size. As for HFRM, we use the same architecture with fewer residual blocks and reduced number of feature channels. The whole training is implemented on eight NVIDIA Tesla V100 GPUs. We use Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), Frechet Inception Distance (FID), and inference time on a single NVIDIA Tesla V100 GPU as the main evaluation

TABLE I
COMPARISON OF DIFFERENT SETTINGS FOR LEARNING THE DISTRIBUTIONS OF CLEAN IMAGES ON RAINDROP. $Comp.$ : DIFFUSION COMPONENTS. $Cond.$ : CONDITIONAL COMPONENTS. ✔: USED. ✗: NOT USED.

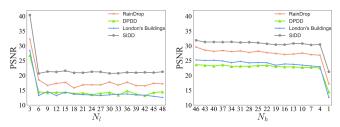| Method | HFRM | $Comp.$ | $Cond.$ | PSNR | SSIM | Time |
|--------|------|---------|---------|------|------|------|
| PatchDM | ✗ | $\mathbf{X}_t$ | $\mathbf{X}_d$ | 32.08 | 0.937 | 61.27s |
| WaveDM₁ | ✗ | $\mathbf{x}_t$ | $\mathbf{x}_d$ | 17.16 | 0.391 | 1.12s |
| WaveDM₂ | ✗ | $\mathbf{x}_t^l$ | $\mathbf{x}_d^l$ | 29.80 | 0.924 | **0.75s** |
| WaveDM₃ | ✔ | $\mathbf{x}_t^l$ | $\mathbf{x}_d, \tilde{\mathbf{x}}_0^h$ | **32.23** | **0.944** | 0.97s |

Performance over Number of Wavelet Bands



Fig. 3. Different numbers of the wavelet bands for diffusion. $N_l$ indicates using the 1st to the $N_l$-th bands. $N_h$ indicates using the 48-th to the $N_h$-th bands.

metrics. Besides, we also test the number of model parameters and memory consumption for reference.

TABLE II
MODEL CONFIGURATIONS AND PARAMETER CHOICES.

| Network | Setting | Time |
|---------|---------|------|
| Noise Estimation Network | Base channels | 128 |
| | Channel multipliers | {1, 1, 2, 2, 4, 4} |
| | Residual blocks per resolution | 2 |
| | Attention resolutions | $h/4$ ($h$: input height) |
| | Time step embedding length | 512 |
| HFRM | Base channels | 32 |
| | Channel multipliers | {1, 2, 4, 8, 16} |
| | Residual blocks per resolution | 1 |
| | Attention resolutions | $h/4$ ($h$: input height) |

TABLE III
EVALUATION OF TWO MODULES IN TERMS OF PSNR ON THE RAINDROP DATASET.

| NEN | HFRM | Description | PSNR (↑) |
|-----|------|-------------|----------|
| Default | Default | Base channels (NEN): 128, Multipliers (NEN): {1, 1, 2, 2, 4, 4} Base channels (HFRM): 32, Multipliers (HFRM): {1, 2, 4, 8, 16} | 32.25dB |
| Variant 1 | Default | Base channels (NEN): 128, Multipliers (NEN): {1, 1, 2, 2, 4, 6} Base channels (HFRM): 32, Multipliers (HFRM): {1, 2, 4, 8, 16} | 32.37dB |
| Variant 2 | Default | Base channels (NEN): 256, Multipliers (NEN): {1, 1, 2, 2, 4, 6} Base channels (HFRM): 32, Multipliers (HFRM): {1, 2, 4, 8, 16} | 32.39dB |
| Default | Variant 1 | Base channels (NEN): 128, Multipliers (NEN): {1, 1, 2, 2, 4, 4} Base channels (HFRM): 32 | 32.22dB |

### B. Ablation Studies

*1) Input Conditions:* In this section, we explore how the varieties of conditions influence the restoration performance. Several choices for them with corresponding quantitative results are shown in Table I, in which all methods use 25
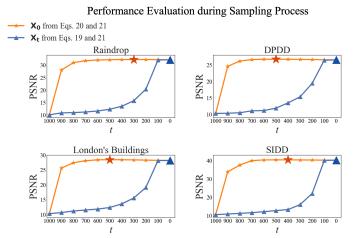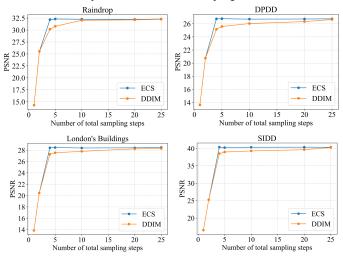
Fig. 4. Performance evaluation during sampling process using 10 steps of stride 100 on four datasets. The ★ and ▲ denote the best PSNR values of the obtained $\mathbf{X}_0$ from Eq. 20 and 21 and $\mathbf{X}_t$ from Eq. 19 and 21, respectively. $t$ represents the current time moment of sampling.



Fig. 5. Restoration performance comparison between DDIM sampling and ECS under multiple sampling step settings on four datasets. $t$ represents the current time moment of sampling.

TABLE IV
PERFORMANCE COMPARISON BETWEEN DDIM SAMPLING AND ECS UNDER TWO SAMPLING STEP SETTINGS ON FOUR DATASETS. THE PSNR VALUES ARE COMPUTED AT $t = 0$ AVERAGED ON EACH DATASET.

| Method | Step | Sampling Trajectory | PSNR | | | |
|--------|------|---------------------|---------|--------|--------------------|---------|
| | | | Raindrop | DPDD | London's Buildings | SIDD |
| ECS | 4 | $1000{\to}870{\to}730$ $\to\mathbf{600}{\to}0$ | 32.19dB | 26.75dB | 28.42dB | 40.38dB |
| DDIM | | $1000{\to}750{\to}500$ $\to250{\to}0$ | 30.17dB | 25.15dB | 26.95dB | 38.52dB |
| ECS | 5 | $1000{\to}900{\to}800$ $\to700{\to}\mathbf{600}{\to}0$ | 32.21dB | 26.77dB | 28.47dB | 40.24dB |
| DDIM | | $1000{\to}800{\to}600$ $\to400{\to}200{\to}0$ | 30.79dB | 25.59dB | 27.29dB | 38.99dB |

to predict the low-frequency bands. The PSNR/SSIM gain by WaveDM$_3$ verifies that HFRM is effective with little extra computation (one-pass for $\tilde{\mathbf{x}}_0^h$).

*2) Model Configurations:* In this part, we explore the significance of structures of the Noise Estimation Network (NEN) and HFRM within our WaveDM, particularly focusing on restoration performance. The default configurations for these modules are detailed in Table II. Both modules have the U-Net architecture similar to PatchDM [73]. Specifically, the setting of NEN is the same as PatchDM's. Meanwhile, HFRM uses fewer blocks and a reduced number of feature channels for a lightweight design. To further understand the influence of these modules, we devise two alternative configurations of NEN and one for HFRM. These variants are then evaluated on the Raindrop dataset. The comparative results, presented in Table V-A, indicate that the performance slightly drops when NEN scales down (compare rows 1–3) and the complexity of HFRM has little impact on the restoration quality (compare rows 1 and 4).

*3) Wavelet Bands:* To further explore what wavelet bands should be used in the diffusion model, we select $N$ wavelet bands for diffusion, and the other $48 - N$ bands in $\mathbf{x}_d$ serve as the input to HFRM. Experimental comparison on the four datasets is shown in Fig. 3, from which we can observe that the restoration performance reaches the best when modeling the first three low-frequency bands for diffusion. Therefore, this setting is used in all the following experiments.

*4) Wavelet Levels:* We also conduct an experiment to explore the effect of different wavelet transform levels on the restoration performance. For levels 1, 2 and 3, the values of PSNR/time of WaveDM with the 25-step DDIM sampling on London's Buildings are 28.12dB/126.16s, 28.39dB/5.21s and 24.14dB/0.72s. After the 1-level Haar wavelet transform, the wavelet bands still have a large spatial size and also need to be cut into patches for processing, which is time-consuming. However, in the 2-level wavelet transform, the inference time can be reduced to 5.21s with a similar PSNR. When the level is further increased (3 or higher), the performance is harmed because too many details are lost in the low-frequency bands for diffusion.

*5) Efficient Conditional Sampling:* The development of the ECS is based on a series of exploratory experiments. Fig. 4 provides a visual representation of a subset of these experiments. The finding from these experiments, serves as a foundation of the ECS formulation. Upon examining the re-

steps of the DDIM sampling. According to them, we can see that although PatchDM performs well in the spatial domain, the processing of a large number of small-size patches is extremely time-consuming. Instead, when switching to the wavelet domain, the total sampling time is reduced by around $1/60$ due to the small spatial size after wavelet transform. However, when modeling the distribution on all 48 bands ($\mathbf{x}_t$) of clean images with all-frequency components ($\mathbf{x}_d$) of degraded images as the condition (WaveDM$_1$), the results are the worst. WaveDM$_2$ models the distribution of clean images on the first three low-frequency bands ($\mathbf{x}_t^l$) with the corresponding bands ($\mathbf{x}_d^l$) of degraded images as the condition, of which the performance is much better than WaveDM$_1$. WaveDM$_3$ is our full WaveDM, where an additional HFRM is added to estimate the high-frequency bands ($\tilde{\mathbf{x}}_0^h$) of clean images, which serves as not only the essential high-frequency bands for inverse wavelet transform but also an extra condition

| $\mathbf{X}_t$ | | | | | | |
| $\mathbf{X}_0$ | | | | | | |
| $t$ | 1000 | 800 | 600 | 400 | 200 | 0 |

Fig. 6. Visual results of the $\mathbf{X}_t$ from Eq. 19 and 21, and $\mathbf{X}_0$ from Eq. 20 and 21 during the sampling process for image raindrop removal.



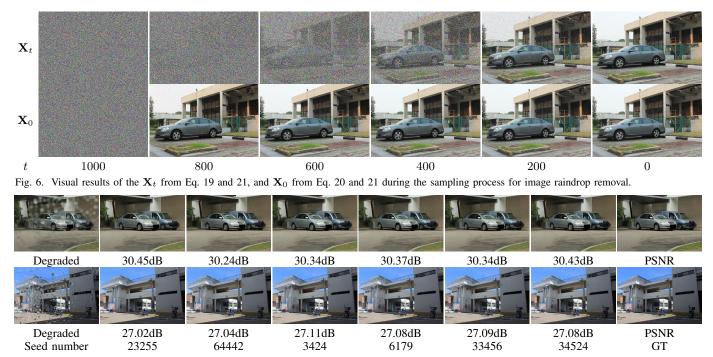| Degraded | 30.45dB | 30.24dB | 30.34dB | 30.37dB | 30.34dB | 30.43dB | PSNR |
| Degraded | 27.02dB | 27.04dB | 27.11dB | 27.08dB | 27.09dB | 27.08dB | PSNR |
| Seed number | 23255 | 64442 | 3424 | 6179 | 33456 | 34524 | GT |

Fig. 7. Visual results of the generated samples with different seeds for image raindrop removal. Each column is generated from the same random seed.
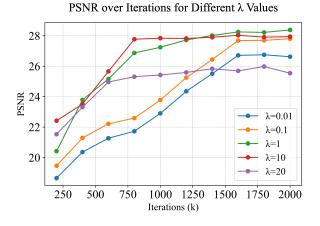


Fig. 8. Performance comparison in terms of PSNR across different training iterations for varying values of $\lambda$ on the London's Buildings dataset.

sults, as depicted in Fig. 4, we observe that when the sampling time reaches the moment around $t = 600$, the PSNR between the predicted $\mathbf{X}_0$ from $\mathbf{X}_{600}$ and the ground truth (GT) reaches a significant value, even if not the highest. After this moment ($t < 600$), the PSNR remains relatively stable with negligible fluctuations. Based on these experiments, we select $M = 600$ as the default setting for ECS for all experiments. Besides, the PSNR of predicted $\mathbf{X}_0$ decreases when the sampling continues after ★, the reason of which comes from the fact that the second term $\sqrt{1 - \bar{\alpha}_{t-1}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^l, \mathbf{x}_d, \tilde{\mathbf{x}}_0^h, t)$ in Eq. 20 introduces extra noise to the results. Due to the small value of the weight term $\sqrt{1 - \bar{\alpha}_{t-1}}$, the decline in PSNR is slight. We also present an example of the denoising process for synthesizing clean images in Fig. 6. To compare ECS with the conventional DDIM sampling method, we show PSNR values for both strategies with the same number of total sampling steps but

with different trajectories. Table IV highlights this comparison for two exemplary settings. An extended comparison with more settings of sampling steps is presented in Fig. 5. From the results, it is observable that as the sampling steps increase, the difference in PSNR decreases. However, when the number of total sampling steps is set below 10, ECS shows better performance, demonstrating its superior efficiency over the traditional DDIM sampling.

*6) Weight of Loss Functions:* In the process of optimizing our model, the weighting parameter $\lambda$ is introduced in the loss function defined in Eq. 18. This parameter serves as a trade-off between the losses associated with the low-frequency wavelet bands, denoted as $L_{\text{simple}}$, and the high-frequency bands, represented as $L_1$. Fig. 8 visually presents the experimental results, portraying the performance of our model across different iterations for various choices of $\lambda$. This evaluation is conducted on the London's Buildings dataset [50] for image demoiréing. From the results, we can obtain the following observations. First, the choice of $\lambda$ impacts the model's convergence rate. A smaller $\lambda$ (e.g., 0.01) leans towards slower convergence due to the focus on low-frequency bands. Conversely, higher values of $\lambda$ like 1 or 10 accelerate the convergence. Second, despite the variance in convergence, the final restoration quality, measured in PSNR, remains almost consistent across a large range of $\lambda$ values (e.g., [0.1, 10]), indicating that our model is not very sensitive to $\lambda$. However, too small or too large $\lambda$ would harm the performance, which reveals the intricate relationship between wavelet spectrum learning and the balance of high- and low-frequency representations. An inappropriate weighting might lead the model to overly focus on either the detailed high frequencies or the coarse low frequencies. In all the following experiments, $\lambda$ is set to 1.

TABLE V
QUANTITATIVE COMPARISON WITH SOTA METHODS ON VARIOUS RESTORATION TASKS.

| Task | Type | Method | Step | PSNR↑ | SSIM↑ | FID↓ | Time↓ | Parameters↓ | Memory↓ |
|------|------|--------|------|-------|-------|------|-------|-------------|---------|
| Raindrop Removal | One-pass | DuRN [54] | | 31.24dB | 0.926 | 30.63 | **0.09s** | 10.2M | 4108MB |
| | | CCN [79] | | 31.44dB | 0.947 | 28.94 | 0.80s | 12.4M | 3738MB |
| | | RainAttn [80] | 1 | 31.44dB | 0.926 | 28.47 | 0.41s | **6.24M** | 8335MB |
| | | AttnGAN [77] | | 31.59dB | 0.917 | 27.84 | 0.63s | 7.08M | 8797MB |
| | | IDT [109] | | 31.87dB | 0.931 | 25.51 | 0.39s | 16.4M | **3660MB** |
| | Iterative | PatchDM$_{64}$ [73] | 10 | 32.13dB | 0.939 | 25.12 | 24.36s | 110M | 8758MB |
| | | PatchDM$_{128}$ [73] | 50 | **32.31dB** | 0.946 | **20.57** | 301.35s | 110M | 22313MB |
| | | Ours | 8 | 32.25dB | **0.948** | 23.53 | 0.30s | 124M | 4965MB |
| Rain steaks Removal | One-pass | HRGAN [45] | | 21.56dB | 0.855 | 69.25 | 0.80s | 50.4M | 3663MB |
| | | PCNet [29] | | 26.19dB | 0.901 | 44.57 | **0.06s** | **0.63M** | **1845MB** |
| | | MPRNet [119] | 1 | 28.03dB | 0.919 | 30.61 | 0.12s | 20.1M | 6942MB |
| | | All-in-One [46] | | 24.71dB | 0.898 | \ | \ | \ | \ |
| | | TransWeather [100] | | 28.83dB | 0.900 | 22.52 | 0.16s | 38.1M | 4734MB |
| | Iterative | PatchDM$_{64}$ [73] | 25 | 28.38dB | 0.932 | 17.36 | 59.97s | 110M | 8759MB |
| | | Ours | 4 | **31.39dB** | **0.943** | **11.42** | 0.16s | 124M | 4965MB |
| Dehazing | One-pass | DCP [21] | | 19.13dB | 0.815 | 20.03 | **0.05s** | \ | 2388MB |
| | | GridDehazeNet [52] | | 30.86dB | 0.982 | 4.76 | 0.20s | **0.96M** | 1956MB |
| | | MSBDN [14] | 1 | 33.48dB | 0.982 | 5.59 | 0.16s | 31.4M | **1756MB** |
| | | FFA-Net [78] | | 33.57dB | 0.984 | 6.43 | 0.36s | 4.46M | 2246MB |
| | | DehazeFormer-B [96] | | 34.95dB | 0.984 | 4.58 | 0.14s | 2.51M | 1760MB |
| | Iterative | PatchDM$_{64}$ [73] | 25 | 35.52dB | 0.989 | 5.75 | 19.31s | 110M | 8759MB |
| | | Ours | 4 | **37.00dB** | **0.994** | **2.80** | 0.15s | 124M | 6336MB |
| Single-pixel Defocus Debluring | One-pass | DMENet [38] | | 23.41dB | 0.714 | 54.51 | 1.79s | 26.9M | 8954MB |
| | | DPDNet [2] | | 24.34dB | 0.747 | 55.21 | 0.32s | 32.3M | 11747MB |
| | | KPAC [91] | 1 | 25.22dB | 0.774 | 46.49 | 0.33s | **2.64M** | 12575MB |
| | | IFAN [39] | | 25.37dB | 0.789 | 46.47 | **0.20s** | 10.5M | 19273MB |
| | | Restormer [117] | | 25.98dB | 0.811 | **43.13** | 3.22s | 25.5M | 26256MB |
| | Iterative | PatchDM$_{64}$ [73] | 25 | 26.49dB | 0.812 | 47.92 | 365.20s | 110M | **8759MB** |
| | | Ours | 4 | **26.75dB** | **0.822** | 45.43 | 0.47s | 124M | 12190MB |
| Dual-pixel Defocus Debluring | One-pass | DPDNet [2] | | 25.13dB | 0.786 | 45.52 | 0.32s | 32.3M | **12265MB** |
| | | RDPD [3] | | 25.39dB | 0.772 | 39.71 | 0.29s | 24.3M | 18492MB |
| | | IFAN [39] | 1 | 25.99dB | 0.804 | 36.87 | **0.20s** | **10.5M** | 20127MB |
| | | Restormer [117] | | 26.66dB | 0.833 | 34.49 | 3.22s | 25.5M | 28214MB |
| | Iterative | Ours | 4 | **27.49dB** | **0.855** | **31.28** | 0.48s | 124M | 12326MB |
| Demoiréing | One-pass | MultiscaleNet [98] | | 23.64dB | 0.791 | 71.39 | 0.59s | **0.65M** | 15486MB |
| | | WDNet [50] | 1 | 24.12dB | 0.847 | 51.65 | **0.18s** | 3.92M | 21472MB |
| | | FHDe$^2$Net [20] | | 24.31dB | 0.799 | 41.38 | 2.03s | 13.6M | 27686MB |
| | | ESDNet [114] | | 25.67dB | 0.871 | 58.92 | 0.23s | 5.93M | 28432MB |
| | Iterative | PatchDM$_{64}$ [73] | 25 | 28.09dB | 0.934 | 33.51 | 656.75s | 110M | **8759MB** |
| | | Ours | 4 | **28.42dB** | **0.942** | **23.14** | 1.01s | 124M | 13052MB |
| Real Denoising | One-pass | MIRNet [118] | | 39.72dB | 0.959 | 47.71 | 0.090s | 31.8M | 5805MB |
| | | MPRNet [119] | 1 | 39.71dB | 0.958 | 49.54 | 0.055s | **20.1M** | **2861MB** |
| | | Uformer [106] | | 39.77dB | 0.959 | 47.17 | **0.031s** | 50.9M | 9157MB |
| | | Restormer [117] | | 40.02dB | 0.960 | 47.28 | 0.114s | 25.5M | 7702MB |
| | Iterative | PatchDM$_{64}$ [73] | 25 | 39.86dB | 0.959 | 47.59 | 9.332s | 110M | 8759MB |
| | | Ours | 4 | **40.38dB** | **0.962** | **47.01** | 0.062s | 124M | 3430MB |

*7) Conditional Sampling Variability:* To delve deeper into WaveDM's sampling capability from the conditional distribution, we conduct an experiment where we vary the seed to produce different samples while maintaining a constant degraded image as the condition. Fig. 7 exhibits two groups of generated images for six distinct seeds. To offer a quantitative measure, we also compute the PSNR between each of these samples and the ground truth of the degraded image. The visual inspection of these images along with the quantitive comparison brings an observation that the differences between the samples, even though generated using different seeds, are extremely subtle both visually and quantitively. On the contrary, with different conditions (i.e., different degraded images), the results are generated differently and guided by the conditions. This clearly shows WaveDM's consistent ability to produce high-quality restorations, regardless of the minor variabilities introduced by different seeds.

| Degraded Image | 25.12dB Raindrop | 24.99dB AttnGAN [77] | 27.71dB DuRN [54] | 30.10dB RainAttn [80] | 30.12dB PatchDM [73] | 30.52dB Ours | PSNR GT |

Fig. 9. Visual comparison on image raindrop removal. The PSNR values are computed on the whole images.



| Degraded Image | 25.12dB Raindrop | 20.72dB HRGAN [45] | 28.31dB MPRNet [119] | 29.03dB TransWeather [54] | 31.09dB PatchDM [73] | 34.77dB Ours | PSNR GT |

Fig. 10. Visual comparison on image rain steaks removal. The PSNR values are computed on the whole images.



| Degraded Image | 28.58dB DCP [21] | 27.88dB GridDehazeNet [52] | 38.77dB MSBDN [14] | 29.99dB FFA-Net [78] | 37.58dB DehazeFormer [96] | 38.99dB Ours | PSNR GT |

Fig. 11. Visual comparison on image dehazing.



| Degraded Image | 20.68dB Defocus | 22.03dB DPDNet [2] | 21.90dB IFAN [39] | 21.48dB Restormer [117] | 21.67dB PatchDM [73] | 22.62dB Ours | PSNR GT |

Fig. 12. Visual comparison on image defocus deblurring. The PSNR values are computed on the whole images.



| Degraded Image | 21.69dB Moiré | 24.73dB FHDe2Net [20] | 25.22dB WDNet [50] | 25.75dB ESDNet [114] | 27.83dB PatchDM [73] | 28.12dB Ours | PSNR GT |

Fig. 13. Visual comparison on image demoiréing. The PSNR values are computed on the whole images.



| 26.79dB Noisy Image | 40.27dB MIRNet [118] | 40.42dB MPRNet [119] | 40.64dB Uformer [106] | 40.67dB Restormer [117] | 41.21dB PatchDM [73] | 42.47dB Ours | PSNR GT |

Fig. 14. Visual comparison on real image denoising.

## C. Comparison with State-of-the-Art Methods

We evaluate our WaveDM with other state-of-the-art (SOTA) methods on twelve benchmark datasets.

All results are obtained either by copying from their papers or retraining and testing with their official code and released pretrained models. We re-implement PatchDM for this task as the baseline. The best and second best values are indicated in **bold** and <u>underlined</u>, respectively. Since the default PatchDM cuts images into multiple patches of size $64 \times 64$, denoted

as $PatchDM_{64}$, its memory used keeps unchanged across different image sizes. PatchDM also provides another version $PatchDM_{128}$, which cuts images into $128 \times 128$ patches. Besides, as our noise estimation network keeps the same as the baseline PatchDM for a fair comparison, the extra parameters only come from HFRM. The FID, processing time, parameter count, and the memory usage of method All-in-One [46] cannot be obtained due to unavailability of its code.

*1) Image Raindrop Removal:* In our evaluation on the Rain-Drop dataset [77], various methods are analyzed for their rain-

TABLE VI
QUANTITATIVE COMPARISON WITH SOTA METHODS FOR GAUSSIAN GRAYSCALE IMAGE DENOISING ON THREE COMMON BENCHMARKS.

| Methods | Set12 [123] | | | BSD68 [64] | | | Urban100 [25] | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ |
| MWCNN [51] | 33.15 | 30.79 | 27.74 | 31.86 | 29.41 | 26.53 | 33.17 | 30.66 | 27.42 |
| DeamNet [81] | 33.19 | 30.81 | 27.74 | 31.91 | 29.44 | 26.54 | 33.37 | 30.85 | 27.53 |
| DAGL [69] | 33.28 | 30.93 | 27.81 | 31.93 | 29.46 | 26.51 | 33.79 | 31.39 | 27.97 |
| SwinIR [47] | 33.36 | 31.01 | 27.91 | **31.97** | 29.50 | 26.58 | 33.70 | 31.30 | 27.98 |
| Restormer [117] | _33.42_ | _31.08_ | _28.00_ | _31.96_ | _29.52_ | **26.62** | _33.79_ | _31.46_ | **28.29** |
| Ours | **33.75** | **31.47** | **28.44** | 31.95 | **29.58** | _26.60_ | **33.92** | **31.86** | _28.21_ |

TABLE VII
QUANTITATIVE COMPARISON WITH SOTA METHODS FOR GAUSSIAN COLOR IMAGE DENOISING ON THREE COMMON BENCHMARKS.

| Methods | CBSD68 [64] | | | Kodak24 [126] | | | Urban100 [25] | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ | $\sigma = 15$ | $\sigma = 25$ | $\sigma = 50$ |
| IRCNN [124] | 33.86 | 31.16 | 27.86 | 34.69 | 32.18 | 28.93 | 33.78 | 31.20 | 27.70 |
| DnCNN [123] | 33.90 | 31.24 | 27.95 | 34.60 | 32.14 | 28.95 | 32.98 | 30.81 | 27.59 |
| FFDNet [125] | 33.87 | 31.21 | 27.96 | 34.63 | 32.13 | 28.98 | 33.83 | 31.40 | 28.05 |
| DRUNet [122] | 34.30 | 31.69 | 28.51 | 35.31 | 32.89 | 29.86 | 34.81 | 32.60 | 29.61 |
| Restormer [117] | _34.39_ | _31.78_ | _28.59_ | **35.44** | **33.02** | **30.00** | _35.06_ | **32.91** | _30.02_ |
| Ours | **34.85** | **31.81** | **28.78** | _35.41_ | _32.97_ | **30.23** | 35.31 | _32.89_ | **30.22** |

drop removal efficiency, with the results tabulated in Table V. While PatchDM [73] slightly edges out in terms of PSNR (a marginal 0.06 dB advantage over WaveDM), its practicality is limited due to the extensive inference time involved in patch processing. WaveDM, in contrast, demonstrates comparable performance in a fraction of the time, proving its efficiency in dealing with challenging conditions like heavy raindrop obstruction, as substantiated in Fig. 9.

*2) Image Rain Streaks Removal:* Rain streaks present a different challenge compared to raindrops. Fig. 10 exhibits the visual results on the Outdoor-rain dataset [45], which validate that WaveDM shows an excellent capability in removing heavy rain streaks without compromising image details. Besides, quantitative evaluations are presented in Table V, also demonstrating WaveDM's better performance against others' with a competitive processing time.

*3) Image Dehazing:* In addition to the rainy scenarios, we also apply our method to another adverse weather condition, haze. We select 6000 images from the training set SOTS [41] that contains over 70000 images for training, and evaluate our model on the SOTS-Outdoor benchmark. The quantitive results shown in Table V demonstrate that WaveDM achieves the best PSNR and SSIM. Besides, WaveDM obtains 2.8 in terms of FID, showing high fidelity of restored samples. The visual samples presented in Fig. 11 evidence WaveDM's dehazing performance, including the image's clarity, sharpness, and color preservation.

*4) Image Defocus Deblurring:* Our experiments on the DPDD dataset [2] for both single and dual-pixel defocus deblurring, as captured in Table V and Fig. 12, reveal WaveDM's superior performance over other SOTA methodologies with a comparable inference time to the one-pass methods.

*5) Image Demoiréing:* In dealing with moiré patterns in images, especially those from the London's Buildings dataset [50], WaveDM proves to be a strong competitor. While diffusion models generally perform better than one-pass methods, WaveDM distinguishes itself by obtaining superb performance with quick inference, matching the speed of one-pass systems,

TABLE VIII
COMPARISON WITH LATENT DIFFUSION IMPLEMENTATION.

| Implementation | Testing | | | Upper Bound | |
|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | Time↓ | PSNR↑ | SSIM↑ |
| WaveDM | 31.39dB | 0.943 | 0.81s | 56.38dB | 0.999 |
| LDM | 24.31dB | 0.873 | 1.05s | 25.82dB | 0.895 |

as evidant in Table V and Fig. 13.

*6) Real Image Denoising:* On the SIDD dataset [1] for real image denoising tasks, WaveDM's effectiveness is further confirmed. It not only outperforms other one-pass methods but also matches the performance of diffusion-based methods such as PatchDM. The images in Fig. 14 and numerical evaluations in Table V clearly showcase its strengths.

*7) Gaussian Image Denoising:* In addition to real image denoising, we also apply WaveDM to Gaussian image denoising. For grayscale image denoising, we employ three widely-used benchmark datasets: Set12 [123], BSD68 [64], and Urban100 [25]. The quantitative results, presented in Table VI, show that our method overall outperforms other competing methods across different noise levels. Specifically, for noise level $\sigma = 50$, our method achieves a remarkable PSNR of 28.44dB on Set12, which is notably higher than other methods. Similarly, for color image denoising, our experiments span across datasets CBSD68 [64], Kodak24 [126], and Urban100 [25]. The results, detailed in Table VII, further solidify our method's superior performance. For instance, on the Urban100 dataset at noise level $\sigma = 50$, our model achieves an impressive PSNR of 30.22dB, surpassing all competitors. These experimental results offer solid evidence that WaveDM is not only robust to varying degrees of Gaussian noise but also consistently outperforms SOTA methods, emphasizing its effectiveness and adaptability.

### D. Comparison with Latent Diffusion Implementation

To further demonstrate the efficiency and effectiveness of WaveDM, which employs wavelet transform for image size reduction and diffusion modeling in the wavelet domain, we conduct a comparative experiment on the Outdoor-rain dataset [45] with the Latent Diffusion Model (LDM) [82] implementation, a method that can also reduce image size using VAE-based subsampling. Specifically, the images, sized $720 \times 480 \times 3$, are processed by a 4-downsampled pretrained VAE from [82] to be transformed into a latent space. These transformed images are then used as input for the latent diffusion model (LDM), which keeps its architecture the same as WaveDM's. The results of this comparative evaluation are presented in Table VIII, in which "Upper Bound" gives the maximum results WaveDM and LDM can achieve, where the PSNR and SSIM values are computed by directly applying either the wavelet transform (and its inverse) or the VAE's encoder-decoder on clean images (ground truth), without any diffusion processing. From the table, we observe that the VAE's reconstruction restricts LDM's restoration capability. Additionally, LDM tends to be slightly slower than WaveDM, since the wavelet transformation is inherently more efficient

than VAE processing. In conclusion, WaveDM demonstrates superior restoration and efficiency compared to the LDM alternative.

## VI. CONCLUSION AND LIMITATION

This paper proposes a wavelet-based diffusion model (WaveDM) to reduce the inference time of diffusion-based models for image restoration. WaveDM learns the distribution in the wavelet domain of clean images, which saves a lot of time in each step of sampling. In addition, an efficient conditional sampling technique is developed from experiments to reduce the total sampling steps to around 5. Experiments on twelve image datasets validate that our WaveDM achieves SOTA performance with the efficiency that is over $100\times$ faster than the previous diffusion-based SOTA PatchDM and is also comparable to traditional one-pass methods.

The major limitation is that WaveDM requires millions of training iterations for several days, especially for large-scale datasets, which is left to deal with in future work.

## REFERENCES

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, 2018.

[2] Abdullah Abuolaim and Michael S Brown. Defocus deblurring using dual-pixel data. In *ECCV*, 2020.

[3] Abdullah Abuolaim, Mauricio Delbracio, Damien Kelly, Michael S Brown, and Peyman Milanfar. Learning to reduce defocus blur by realistically modeling dual-pixel data. In *ICCV*, 2021.

[4] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *CVPR*, 2017.

[5] Saeed Anwar and Nick Barnes. Densely residual laplacian super-resolution. *IEEE TPAMI*, 44(3):1192–1204, 2020.

[6] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, 2022.

[7] Sung-Jin Cho, Seo-Won Ji, Jun-Pyo Hong, Seung-Won Jung, and Sung-Jea Ko. Rethinking coarse-to-fine approach in single image deblurring. In *ICCV*, 2021.

[8] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *ICCV*, 2021.

[9] Hyungjin Chung, Jeongsol Kim, Sehui Kim, and Jong Chul Ye. Parallel diffusion models of operator and image for blind inverse problems. In *CVPR*, 2023.

[10] Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. In *ICLR*, 2023.

[11] Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. *NeurIPS*, 2022.

[12] Hyungjin Chung, Byeongsu Sim, and Jong Chul Ye. Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction. In *CVPR*, 2022.

[13] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *NeurIPS*, 2021.

[14] Hang Dong, Jinshan Pan, Lei Xiang, Zhe Hu, Xinyi Zhang, Fei Wang, and Ming-Hsuan Yang. Multi-scale boosted dehazing network with dense feature fusion. In *CVPR*, 2020.

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.

[16] Akshay Dudhane, Syed Waqas Zamir, Salman Khan, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Burst image restoration and enhancement. In *CVPR*, 2022.

[17] Guangwei Gao, Lei Tang, Fei Wu, Huimin Lu, and Jian Yang. Jdsr-gan: Constructing an efficient joint learning network for masked face super-resolution. *IEEE TMM*, 25:1505–1512, 2023.

[18] Florentin Guth, Simon Coste, Valentin De Bortoli, and Stephane Mallat. Wavelet score-based generative modeling. *NeurIPS*, 2022.

[19] Kewen Han and Xinguang Xiang. Decomposed cyclegan for single image deraining with unpaired data. In *ICASSP*, 2020.

[20] Bin He, Ce Wang, Boxin Shi, and Ling-Yu Duan. Fhde 2 net: Full high definition demoireing network. In *ECCV*, 2020.

[21] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE TPAMI*, 33(12):2341–2353, 2010.

[22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.

[23] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *JMLR*, 23:47–1, 2022.

[24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop*, 2021.

[25] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed selfexemplars. In *CVPR*, 2015.

[26] Yi Huang, Yu Dong, He Zhang, Jiancheng Huang, and Shifeng Chen. Learning image-adaptive lookup tables with spatial awareness for image harmonization. *IEEE Transactions on Consumer Electronics*, 2023.

[27] Ka-Hei Hui, Ruihui Li, Jingyu Hu, and Chi-Wing Fu. Neural wavelet-domain diffusion for 3d shape generation. In *SIGGRAPH Asia*, 2022.

[28] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *JMLR*, 6(4), 2005.

[29] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Zheng Wang, and Chia-Wen Lin. Pcnet: progressive coupled network for real-time image deraining. In *ICIP*, 2021.

[30] Xin Jin, Zhibo Chen, Jianxin Lin, Zhikai Chen, and Wei Zhou. Unsupervised single image deraining with self-supervised constraints. In *ICIP*, 2019.

[31] Zhi Jin, Muhammad Zafar Iqbal, Dmytro Bobkov, Wenbin Zou, Xia Li, and Eckehard Steinbach. A flexible deep cnn framework for image restoration. *IEEE TMM*, 22(4):1055–1068, 2019.

[32] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *NeurIPS*, 2022.

[33] Bahjat Kawar, Roy Ganz, and Michael Elad. Enhancing diffusion-based image synthesis with robust classifier guidance. *TMLR*, 2022.

[34] Johannes Kopf, Boris Neubert, Billy Chen, Michael Cohen, Daniel Cohen-Or, Oliver Deussen, Matt Uyttendaele, and Dani Lischinski. Deep photo: Model-based photograph enhancement and viewing. *ACM TOG*, 27(5):1–10, 2008.

[35] Manoj Kumar, Dirk Weissenborn, and Nal Kalchbrenner. Colorization transformer. In *ICLR*, 2021.

[36] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *ICCV*, 2019.

[37] Gihyun Kwon and Jong Chul Ye. Diffusion-based image translation using disentangled style and content representation. In *ICLR*, 2023.

[38] Junyong Lee, Sungkil Lee, Sunghyun Cho, and Seungyong Lee. Deep defocus map estimation using domain adaptation. In *CVPR*, 2019.

[39] Junyong Lee, Hyeongseok Son, Jaesung Rim, Sunghyun Cho, and Seungyong Lee. Iterative filter adaptive network for single image defocus deblurring. In *CVPR*, 2021.

[40] Sangyun Lee, Hyungjin Chung, Jaehyeon Kim, and Jong Chul Ye. Progressive deblurring of diffusion models for coarse-to-fine image synthesis. In *NeurIPS Workshop*, 2022.

[41] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *IEEE TIP*, 28(1):492–505, 2018.

[42] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022.

[43] Juncheng Li, Bodong Cheng, Ying Chen, Guangwei Gao, and Tieyong Zeng. Ewt: Efficient wavelet-transformer for single image denoising. *arXiv preprint arXiv:2304.06274*, 2023.

[44] Juncheng Li, Faming Fang, Kangfu Mei, and Guixu Zhang. Multi-scale residual network for image super-resolution. In *ECCV*, 2018.

[45] Ruoteng Li, Loong-Fah Cheong, and Robby T Tan. Heavy rain image restoration: Integrating physics model and conditional adversarial learning. In *CVPR*, 2019.

[46] Ruoteng Li, Robby T Tan, and Loong-Fah Cheong. All in one bad weather removal using architectural search. In *CVPR*, 2020.

[47] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *ICCV*, 2021.

[48] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *CVPRW*, 2017.

[49] Jianxin Lin, Lianying Yin, and Yijun Wang. Steformer: Efficient stereo image super-resolution with transformer. *IEEE TMM*, 2023.

[50] Lin Liu, Jianzhuang Liu, Shanxin Yuan, Gregory Slabaugh, Aleš

Leonardis, Wengang Zhou, and Qi Tian. Wavelet-based dual-branch network for image demoiréing. In *ECCV*, 2020.

[51] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *CVPRW*, 2018.

[52] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Griddehazenet: Attention-based multi-scale network for image dehazing. In *ICCV*, 2019.

[53] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. In *WACV*, 2023.

[54] Xing Liu, Masanori Suganuma, Zhun Sun, and Takayuki Okatani. Dual residual networks leveraging the potential of paired operations for image restoration. In *CVPR*, 2019.

[55] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.

[56] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *NeurIPS*, 2022.

[57] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.

[58] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, 2022.

[59] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflow: Learning the super-resolution space with normalizing flow. In *ECCV*, 2020.

[60] Zhaoyang Lyu, Xudong Xu, Ceyuan Yang, Dahua Lin, and Bo Dai. Accelerating diffusion models via early stop of the diffusion process. *arXiv preprint arXiv:2205.12524*, 2022.

[61] Hengyuan Ma, Li Zhang, Xiatian Zhu, and Jianfeng Feng. Accelerating score-based generative models with preconditioned diffusion sampling. In *ECCV*, 2022.

[62] Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei Zhang. Waterloo exploration database: New challenges for image quality assessment models. *IEEE TIP*, 26(2):1004–1016, 2016.

[63] Long Ma, Tengyu Ma, Risheng Liu, Xin Fan, and Zhongxuan Luo. Toward fast, flexible, and robust low-light image enhancement. In *CVPR*, 2022.

[64] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*. IEEE, 2001.

[65] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2022.

[66] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *CVPR*, 2023.

[67] Tomer Michaeli and Michal Irani. Nonparametric blind super-resolution. In *ICCV*, 2013.

[68] Chong Mou, Jian Zhang, Xiaopeng Fan, Hangfan Liu, and Ronggang Wang. Cola-net: Collaborative attention network for image restoration. *IEEE TMM*, 24:1366–1377, 2021.

[69] Chong Mou, Jian Zhang, and Zhuoyuan Wu. Dynamic attentive graph learning for image restoration. In *ICCV*, 2021.

[70] Naoki Murata, Koichi Saito, Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Yuki Mitsufuji, and Stefano Ermon. Gibbsddrm: A partially collapsed gibbs sampler for solving blind inverse problems with denoising diffusion restoration. In *ICML*, 2023.

[71] Nithin Gopalakrishnan Nair, Kangfu Mei, and Vishal M Patel. At-ddpm: Restoring faces degraded by atmospheric turbulence using denoising diffusion probabilistic models. In *WACV*, 2022.

[72] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, 2021.

[73] Ozan Özdenizci and Robert Legenstein. Restoring vision in adverse weather conditions with patch-based denoising diffusion models. *IEEE TPAMI*, 2023.

[74] Hao Phung, Quan Dao, and Anh Tran. Wavelet diffusion models are fast and scalable image generators. In *CVPR*, 2023.

[75] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. In *CVPR*, 2022.

[76] Densen Puthussery, Hrishikesh Panikkasseril Sethumadhavan, Melvin Kuriakose, and Jiji Charangatt Victor. Wdrn: A wavelet decomposed relightnet for image relighting. In *ECCV*, 2020.

[77] Rui Qian, Robby T Tan, Wenhan Yang, Jiajun Su, and Jiaying Liu.

[78] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. Ffa-net: Feature fusion attention network for single image dehazing. In *AAAI*, 2020.

[79] Ruijie Quan, Xin Yu, Yuanzhi Liang, and Yi Yang. Removing raindrops and rain streaks in one go. In *CVPR*, 2021.

[80] Yuhui Quan, Shijie Deng, Yixin Chen, and Hui Ji. Deep learning for seeing through window with raindrops. In *ICCV*, 2019.

[81] Chao Ren, Xiaohai He, Chuncheng Wang, and Zhibo Zhao. Adaptive consistency prior based deep network for image denoising. In *CVPR*, 2021.

[82] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.

[83] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.

[84] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *SIGGRAPH*, 2022.

[85] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 2022.

[86] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE TPAMI*, 2022.

[87] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.

[88] Hiroshi Sasaki, Chris G Willcocks, and Toby P Breckon. Unit-ddpm: Unpaired image translation with denoising diffusion probabilistic models. *arXiv preprint arXiv:2104.05358*, 2021.

[89] Jie Shi, Chenfei Wu, Jian Liang, Xiang Liu, and Nan Duan. Divae: Photorealistic images synthesis with denoising diffusion decoder. *arXiv preprint arXiv:2206.00386*, 2022.

[90] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015.

[91] Hyeongseok Son, Junyong Lee, Sunghyun Cho, and Seungyong Lee. Single image defocus deblurring using kernel-sharing parallel atrous convolutions. In *ICCV*, 2021.

[92] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.

[93] Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. In *ICLR*, 2023.

[94] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 2019.

[95] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *NeurIPS*, 2020.

[96] Yuda Song, Zhuqing He, Hui Qian, and Xin Du. Vision transformers for single image dehazing. *IEEE TIP*, 32:1927–1941, 2023.

[97] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *ICLR*, 2021.

[98] Yujing Sun, Yizhou Yu, and Wenping Wang. Moiré photo restoration using multiresolution convolutional neural networks. *IEEE TIP*, 27(8):4160–4172, 2018.

[99] Radu Timofte, Vincent De Smet, and Luc Van Gool. Anchored neighborhood regression for fast example-based super-resolution. In *ICCV*, 2013.

[100] Jeya Maria Jose Valanarasu, Rajeev Yasarla, and Vishal M Patel. Transweather: Transformer-based restoration of images degraded by adverse weather conditions. In *CVPR*, 2022.

[101] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017.

[102] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

[103] Lifeng Wang, Zhouchen Lin, Tian Fang, Xu Yang, Xuan Yu, and Sing Bing Kang. Real-time rendering of realistic rain. In *SIGGRAPH Sketches*. 2006.

[104] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022.

[105] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. In *ICLR*, 2023.

[106] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped

transformer for image restoration. In *CVPR*, 2022.

[107] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *CVPR*, 2022.

[108] Jay Whang, Erik Lindgren, and Alex Dimakis. Composing normalizing flows for inverse problems. In *ICML*, 2021.

[109] Jie Xiao, Xueyang Fu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. Image de-raining transformer. *IEEE TPAMI*, 2022.

[110] Jingwei Xin, Jie Li, Xinrui Jiang, Nannan Wang, Heng Huang, and Xinbo Gao. Wavelet-based dual recursive network for image super-resolution. *IEEE TNNLS*, 33(2):707–720, 2020.

[111] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Baining Guo. Learning texture transformer network for image super-resolution. In *CVPR*, 2020.

[112] Qiaosi Yi, Juncheng Li, Qinyan Dai, Faming Fang, Guixu Zhang, and Tieyong Zeng. Structure-preserving deraining with residue channel prior guidance. In *ICCV*, 2021.

[113] Qiaosi Yi, Juncheng Li, Faming Fang, Aiwen Jiang, and Guixu Zhang. Efficient and accurate multi-scale topological network for single image dehazing. *IEEE TMM*, 24:3114–3128, 2021.

[114] Xin Yu, Peng Dai, Wenbo Li, Lan Ma, Jiajun Shen, Jia Li, and Xiaojuan Qi. Towards efficient and scale-robust ultra-high-definition image demoiréing. In *ECCV*, 2022.

[115] Zongsheng Yue and Chen Change Loy. Difface: Blind face restoration with diffused error contraction. *arXiv preprint arXiv:2212.06512*, 2022.

[116] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *BMVC*, 2016.

[117] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, 2022.

[118] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, 2020.

[119] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, 2021.

[120] He Zhang, Vishwanath Sindagi, and Vishal M Patel. Image de-raining using a conditional generative adversarial network. *IEEE TCSVT*, 30(11):3943–3956, 2019.

[121] Kaihao Zhang, Dongxu Li, Wenhan Luo, and Wenqi Ren. Dual attention-in-attention model for joint rain streak and raindrop removal. *IEEE TIP*, 30:7608–7619, 2021.

[122] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE TPAMI*, 44(10):6360–6376, 2021.

[123] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE TIP*, 26(7):3142–3155, 2017.

[124] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep cnn denoiser prior for image restoration. In *CVPR*, 2017.

[125] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE TIP*, 27(9):4608–4622, 2018.

[126] Lei Zhang, Xiaolin Wu, Antoni Buades, and Xin Li. Color demosaicking by local directional interpolation and nonlocal adaptive thresholding. *Journal of Electronic imaging*, 20(2):023016–023016, 2011.

[127] Menglei Zhang and Qiang Ling. Supervised pixel-wise gan for face super-resolution. *IEEE TMM*, 23:1938–1950, 2020.

[128] Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. Inversion-based style transfer with diffusion models. In *CVPR*, 2023.

[129] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *ECCV*, 2018.

[130] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. *NeurIPS*, 2022.

[131] Bolun Zheng, Shanxin Yuan, Chenggang Yan, Xiang Tian, Jiyong Zhang, Yaoqi Sun, Lin Liu, Aleš Leonardis, and Gregory Slabaugh. Learning frequency domain priors for image demoireing. *IEEE TPAMI*, 44(11):7705–7717, 2021.