

MULTI-SCALE CONDITIONAL GENERATIVE MODELING FOR MICROSCOPIC IMAGE RESTORATION

Anonymous authors

Paper under double-blind review

ABSTRACT

The advance of diffusion-based generative models in recent years has revolutionized state-of-the-art (SOTA) techniques in a wide variety of image analysis and synthesis tasks, whereas their adaptation on image restoration, particularly within computational microscopy remains theoretically and empirically underexplored. In this research, we introduce a multi-scale generative model that enhances conditional image restoration through a novel exploitation of the Brownian Bridge process within wavelet domain. By initiating the Brownian Bridge diffusion process specifically at the lowest-frequency subband and applying generative adversarial networks at subsequent multi-scale high-frequency subbands in the wavelet domain, our method provides significant acceleration during training and sampling while sustaining a high image generation quality and diversity on par with SOTA diffusion models. Experimental results on various computational microscopy and imaging tasks confirm our method’s robust performance and its considerable reduction in its sampling steps and time. This pioneering technique offers an efficient image restoration framework that harmonizes efficiency with quality, signifying a major stride in incorporating cutting-edge generative models into computational microscopy workflows.

1 INTRODUCTION

Within the last decade, the landscape of image synthesis has been radically transformed by the advent of generative models (GMs) (Song et al., 2020b; Ho et al., 2020; Song & Ermon, 2019). Among their broad success in various image synthesis applications, image restoration, including super-resolution, shadow removal, inpainting, etc, have caught much attention due to their importance in various practical scenarios. Image restoration aims to recover high-quality target image from low-quality images measured by an imaging system with assorted degradation effects, e.g., downsampling, aberration and noise. Numerous tasks in photography, sensing and microscopy can be formulated as image restoration problems, and therefore the importance of image restoration algorithms is self-evident in practical scenarios (Isola et al., 2017; Kupyn et al., 2018; Weigert et al., 2018; Wang et al., 2019).

Due to the ill-posedness of most image restoration problems, the application of generative learning becomes crucial for achieving high-quality image reconstruction. The wide applications of deep learning (DL)-based generative models in image restoration began with the success of generative adversarial networks (GANs) (Goodfellow et al., 2014). The emergence of conditional adversarial learning achieved unprecedented success and outperformed classical algorithms in various applications of computational imaging and microscopy (Zhu et al., 2017; Isola et al., 2017; Wang et al., 2018c; Kupyn et al., 2018; Nazeri et al., 2018; Karras et al., 2019; Weigert et al., 2018; Wang et al., 2019). Although GANs have achieved remarkable success in image restoration, they are also known to be prone to training instability and mode collapse. These issues significantly restrict the diversity of outputs produced by GAN models (Kodali et al., 2017; Gui et al., 2021). Recently, diffusion models (DMs) were intensively studied and outperformed GANs in various image generation tasks (Ho et al., 2020; Dhariwal & Nichol, 2021). Derived from the stochastic diffusion process, DMs employ neural networks to approximate the reverse diffusion process and sample target images from white noise with high quality and good mode coverage. Moreover, DMs have also been applied to image restoration, including super-resolution, dehazing, colorization, etc (Saharia et al., 2022b; Rombach et al., 2022; Li et al., 2023; Batzolis et al., 2021; Saharia et al., 2022a; Luo et al., 2023a). Most of the

existing methods regard the low-quality image as one additional condition in the reverse process and utilize it as one of the network arguments for inference. However, these image restoration models lack a clear modelling of the conditional image in the forward process and a theoretical guarantee that the terminal state of the diffusion is closely related to and dependent on the condition.

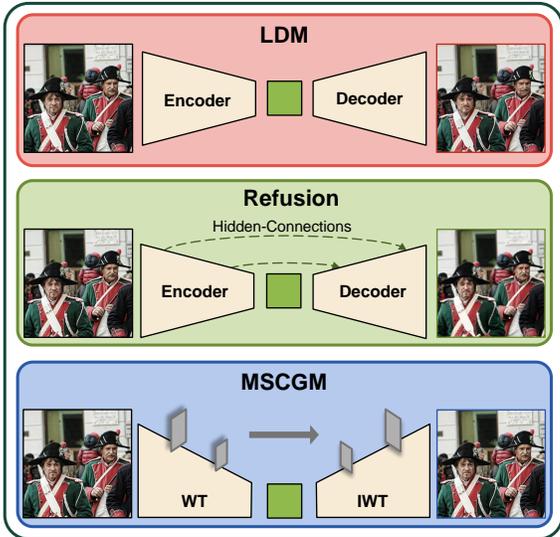


Figure 1: **Inherent information loss in autoencoder backbones.** The same high-quality image passes the three generative models without going through the diffusion process. Public pre-trained checkpoints are used for LDM (Rombach et al., 2022) and Refusion (Luo et al., 2023b). LDM and Refusion suffer from reconstruction loss due to lossy autoencoder backbones while the wavelet and inverse wavelet transform pairs are lossless.

There exist inevitable trade-offs between the compactness of the latent space and the reconstruction fidelity of the autoencoder, and between the normality of its latent distribution and its sample diversity (Saharia et al., 2022b; Li et al., 2022; Luo et al., 2023a). Some study, e.g., Refusion, has explored to improve the fidelity of autoencoder by adding hidden connections between the encoder and decoder (Luo et al., 2023b), whereas the inherent information loss of learned compression cannot be completely eliminated. Though these methods reported improved perceptual scores between their outputs and reference images, in image restoration tasks where pixel-wise structural correspondence is emphasized, the reconstruction quality is constrained by the fidelity of autoencoder backbones, producing noticeable distortions and hallucinations as shown in Figure 1.

On the other hand, microscopy image restoration problems are usually composed of complicated degradation models, sparse image distributions and strict constraints on structural correspondence of the outputs to the conditional images (Meiniet al., 2018; Solomon et al., 2018; Zhao et al., 2022). To preserve the signal sparsity of common biological samples and local feature correspondence to the conditional image, existing methods leveraged CNN-based models and trained on carefully registered datasets (Rivenson et al., 2017; Wang et al., 2019; Wu et al., 2019). Moreover, to enhance training and inference stability, conditional GAN were widely utilized to directly predict the corresponding ground truth image in a deterministic manner (Wang et al., 2019). However, the adaptability and effectiveness of advanced generative models like DMs to microscopy image restoration tasks are under-explored.

To address the aforementioned limitations, here we introduce a novel multi-scale conditional generative model (MSCGM) for image restoration based on multi-scale wavelet transform and Brownian bridge stochastic process. For one thing, multi-scale wavelet transform effectively and losslessly compresses the spatial dimensions of conditional images, eliminating the lossy encoding process of current autoencoder-based DMs. Notably, our method does not involve the pre-training of autoencoders in competitive methods on a sufficiently large and diverse dataset sampled from the image domain. For another, Brownian bridge stochastic process incorporates the modelling of low-

Besides, compared to GANs, DMs are further limited by time-consuming sampling and iterative refinement. Many efforts have been made to improve the sampling speed while ensuring the quality of image generation, among which the Latent Diffusion Model (LDM) (Rombach et al., 2022) stands out as a notable example. LDM projects the diffusion process into a low-dimensional latent space of a pre-trained autoencoder, and hence greatly reduces the computational resource required for high-resolution image generation and restoration. The success of LDM inspired a few studies leveraging the latent representations of pretrained networks to reconstruct high-quality images from low-quality images (Yin et al., 2022; Luo et al., 2023a). However, the performance of diffusion models in these approaches is highly dependent on the quality of latent representations of their pretrained backbones.

108 quality conditional image into both the forward and reverse diffusion process and better utilizes the
109 information of the conditional images. Besides, we theoretically analyze the distributions of low-
110 and high-frequency wavelet subbands and apply Brownian bridge diffusion process (BBDP) and
111 adversarial learning to the multi-scale generation of low- and high-frequency subbands, respectively.
112 In sum, the contributions of this work are three-fold:

- 113 1. We are the first to factorize the conditional image generation within multi-scale wavelet domains,
114 establishing a theoretical groundwork for a multi-scale conditional generative modeling;
- 115 2. Capitalizing on the unique distribution characteristics of wavelet subbands, we propose the
116 innovative MSCGM, which seamlessly integrates BBDP and adversarial learning;
- 117 3. We evaluate the MSCGM on various image restoration tasks, demonstrating its superior perfor-
118 mance in both sampling speed and image quality than competitive methods.
119

120 121 2 RELATED WORKS

122 123 2.1 MICROSCOPY IMAGE RESTORATION

124
125 As an important branch of computational imaging, computational microscopy springs up in recent
126 years and aims to restore high-quality, multi-dimensional images from low-quality, low-dimensional
127 measurements, usually with under-resourced equipment. Since the first work on microscopy image
128 super-resolution reported in 2017 Rivenson et al. (2017), DL has enabled a wide spectrum of novel
129 applications that were impossible with conventional optical technologies, e.g., microscopy image
130 super-resolution surpassing the physical resolution limit of microscopic imaging systems Wang et al.
131 (2019), volumetric imaging reconstructing 3D sample volumes from sparse 2D measurements Wu
132 et al. (2019), and virtual image labelling to match the contrast conventionally provided by chemical or
133 biological markers Rivenson et al. (2019). Compared to general image restoration in computational
134 imaging, microscopy image restoration mainly differs in two aspects: (1) The degeneration process,
135 including the transfer function, noise and aberration of the imaging system is generally complex,
136 unknown and hard to measure precisely; and such degeneration process could vary significantly in
137 real-world scenarios due to the variations of subjects, hardware and imaging protocols. (2) Strict
138 pixel-wise correspondence between output and ground truth images and consistency with physical
139 laws are generally emphasized Barbastathis et al. (2019).

140 141 2.2 GENERATIVE MODELS

142 GANs are well-known for generating high-quality, photorealistic samples rapidly Goodfellow et al.
143 (2014); Gui et al. (2021). Through training a discriminator that tells ground truth images apart from
144 fake images generated by the generator network, GANs outperformed traditional CNNs trained with
145 hand-crafted, pixel-based structural losses such as L_1 and L_2 by providing a high-level, learnable
146 perceptual objective. Conditional GANs such as Pix2Pix Isola et al. (2017), pix2pix HD Wang
147 et al. (2018c) and starGAN Choi et al. (2018) have been successfully applied in a wide spectrum of
148 image-to-image translation and image restoration tasks, including image colorization Nazari et al.
149 (2018), style transfer Karras et al. (2019), image deblurring Kupyn et al. (2018), etc. Unsupervised
150 image-to-image translation has also been extensively explored, such as cycleGAN Zhu et al. (2017),
151 UNIT Liu et al. (2017), DualGAN Yi et al. (2017), etc. In the fields of biomedical and microscopy
152 imaging, researchers have also explored the applications of GANs, e.g., reconstructing low-dose CT
153 and MRI images Yang et al. (2018); Hammernik et al. (2018), denoising microscopy images Weigert
154 et al. (2018), super-resolving diffraction-limited microscopy images Wang et al. (2019), among others.
155 Recently, WGSR Korkmaz et al. (2024) has also provided wavelet guided GAN model for image
156 super-resolution tasks.

156 Alternatively, transformer and related architectures have recently emerged and shown superior
157 performance over convolutional neural networks (CNNs)-based GANs. Swin transformer Liu et al.
158 (2021) and SwinIR Liang et al. (2021) have established a strong transformer-based baseline with
159 competitive performance to CNNs. TransGAN substituted CNNs in the common GAN framework
160 with transformers and improved the overall performance Jiang et al. (2021). More recently, diffusion
161 models (DMs) have been introduced and proved to be the state-of-the-art generative model in
162 various image generation benchmarks Ho et al. (2020); Nichol & Dhariwal (2021). Dhariwal et

162 al. reported diffusion models beat GANs on various image synthesis tasks Dhariwal & Nichol
 163 (2021). Nevertheless, the democratization of diffusion models was limited by its huge demand
 164 for computational resources in both training and sampling. Recent advancements in fast sampling
 165 diffusion models Xiao et al. (2021); Rombach et al. (2022); Song et al. (2020a); Phung et al. (2023);
 166 Kong & Ping (2021); Ho et al. (2022) have accelerated the sampling process in image generation tasks.
 167 DDGAN Xiao et al. (2021) effectively combines the strengths of GANs with diffusion principles for
 168 quicker outputs, while LDM utilizes a compressed latent space to speed up generation. DDIM Song
 169 et al. (2020a) optimizes denoising steps for efficiency, FastDPM Kong & Ping (2021) focuses on
 170 algorithmic improvements for rapid sampling, and CDM Ho et al. (2022) employs a staged approach
 171 to simplify the diffusion process. Furthermore, recent works like ?

172 3 METHODS

173 3.1 PRELIMINARIES

174 The outstanding success and wide applications of score-based diffusion models have been witnessed
 175 in the past years. Generally, for a Gaussian process $\{\mathbf{x}_t, t = 1, \dots, T\}$ defined as:

$$176 q(\mathbf{x}_t|\mathbf{x}_{t-1}) = N(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), t = 1, \dots, T, \quad (1)$$

177 the denoising DM attempts to solve the reverse process by parameterizing the conditional reverse
 178 distribution as:

$$179 p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = N\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right), \sigma_t\mathbf{I}\right). \quad (2)$$

180 Here $\alpha_t, \bar{\alpha}_t, \beta_t, \sigma_t$ are constants, and ϵ_θ is the network estimating the mean value of the reverse
 181 process. Despite the high image quality achieved by the denoising process, the total sampling steps T
 182 can be very large and make the sampling process time-consuming.

183 **Theorem 1.** For a given error ε between the generated distribution p_θ and the true distribution p ,
 184 the sampling steps needed can be expressed by Guth et al. (2022):

$$185 T = O(\varepsilon^{-2}\kappa^3), \quad (3)$$

186 where κ is the condition number of the covariance matrix of p .

187 The proof of Theorem 1 is detailed in Appendix B.

188 Therefore, for highly non-Gaussian distributed images, e.g., microscopy images (please refer to
 189 Appendix C, where we theoretically and statistically demonstrate the high degree of non-Gaussianity
 190 of microscopy datasets), which tends to have high sparsity, resolution and contrast, standard diffusion
 191 models may not be practical due to their slow sampling speed and large sampling steps required for
 192 such distributions.

193 To overcome this limitation, we turned to multi-scale wavelet transform (as detailed in Appendix D),
 194 which offers an excellent latent space (i.e., wavelet domain) for generative modeling. The wavelet
 195 domain not only facilitates lossless compression but also provides low-frequency coefficients with
 196 near-Gaussian distributions. Consequently, we transform the diffusion process into the wavelet
 197 domain, thereby introducing a novel multi-scale wavelet-based generative model. It can be shown
 198 that the diffusion in wavelet domain is a dual problem to the original diffusion in the spatial domain.
 199 For a detailed explanation of this duality, please see Appendix E. Due to the distinct characteristics
 200 of wavelet coefficients in low- and high-frequency subbands (as detailed in Appendix C), we adopt
 201 different generative modeling approaches for the low- and high-frequency coefficients, marking a
 202 key innovation in our work. Specifically, while the low-frequency wavelet coefficients exhibit a
 203 Gaussian tendency, the high-frequency coefficients are sparse and non-Gaussian. Therefore, for the
 204 low-frequency coefficients, we employed the Brownian Bridge Diffusion Process (BBDP) Li et al.
 205 (2023), and for the high-frequency coefficients, we utilized a Generative Adversarial Network (GAN)
 206 based generative method.

207 3.2 BROWNIAN BRIDGE DIFFUSION PROCESS

208 We leverage BBDP to better model the conditional diffusion process and apply it to image restoration.
 209 Image restoration tasks focus on the generation of the target image $\mathbf{x}_0 \in \mathbb{R}^{H \times W \times C}$ from a conditional

image $\mathbf{y} \in \mathbb{R}^{H \times W \times C}$. Most existing DMs tackle the conditional image \mathbf{y} as an additional input argument to ϵ_θ , without integrating the conditional probability distribution on \mathbf{y} into the diffusion theory. Different from the standard diffusion process, we adapt the Brownian bridge process and derive a conditional diffusion process for image-to-image translation, termed as BBDP.

Definition 2 (Forward Brownian bridge process). The forward Brownian bridge with initial state \mathbf{x}_0 and terminal state \mathbf{y} is defined as

$$q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_T = \mathbf{y}) = N\left(\left(1 - \frac{t}{T}\right)\mathbf{x}_0 + \frac{t}{T}\mathbf{y}, \frac{t(T-t)}{T^2}\mathbf{I}\right), \quad (4)$$

By denoting $m_t = \frac{t}{T}$ and $\delta_t = \frac{t(T-t)}{T^2}$, we can reparameterize the distribution of \mathbf{x}_t as:

$$\mathbf{x}_t = \mathbf{x}_0 + m_t(\mathbf{y} - \mathbf{x}_0) + \sqrt{\delta_t}\epsilon_t, \epsilon_t \sim N(\mathbf{0}, \mathbf{I}) \quad (5)$$

Following existing works on diffusion models, a network is trained to estimate \mathbf{x}_0 from \mathbf{x}_t, \mathbf{y} ; in other words, a network ϵ_θ is trained to estimate $m_t(\mathbf{y} - \mathbf{x}_0) + \sqrt{\delta_t}\epsilon_t$.

The loss function of ϵ_θ is defined as below where γ_t is the weight for each t :

$$L = \sum_t \gamma_t \mathbb{E}_{(\mathbf{x}_0, \mathbf{y}), \epsilon_t} \|m_t(\mathbf{y} - \mathbf{x}_0) + \sqrt{\delta_t}\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)\|_2^2, \quad (6)$$

Theorem 3. The reverse process can be shown to be a Gaussian process with mean $\mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})$ and variance $\delta'_t \mathbf{I}$

$$p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) = N(\mathbf{x}_{t-1}; \mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}), \delta'_t \mathbf{I}), \quad (7)$$

$$\mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) = c_{xt}\mathbf{x}_t + c_{yt}\mathbf{y} - c_{et}\epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t) \quad (8)$$

$$\delta'_t = \frac{\delta_{t|t-1}\delta_{t-1}}{\delta_t}. \quad (9)$$

Equation 8 indicates the posterior sampling process of BBDP and the training process of BBDP is detailed here. The input is the low-frequency wavelet coefficients of the original image and the output is the low-frequency wavelet coefficients of the restored image.

Algorithm 1 Training of BBDP

- 1: **repeat**
 - 2: $\mathbf{x}_0, \mathbf{y} \sim q(\mathbf{x}_0, \mathbf{y})$
 - 3: $t \sim \text{Uniform}([1, \dots, T])$
 - 4: $\epsilon_t \sim N(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
 $\nabla_\theta \|m_t(\mathbf{y} - \mathbf{x}_0) + \sqrt{\delta_t}\epsilon_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)\|_2^2$
 - 6: **until** converged
-

The proof of the expression of c_{xt}, c_{yt}, c_{et} is provided as below:

Proof of Theorem 3. Brownian bridge process can prove to be Markovian and hence defined alternatively by its one-step forward process. The one-step forward process $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y})$ can be derived as:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}) = N\left(\frac{1 - m_t}{1 - m_{t-1}}\mathbf{x}_{t-1} + \left(m_t - \frac{1 - m_t}{1 - m_{t-1}}m_{t-1}\right)\mathbf{y}, \delta_{t|t-1}\mathbf{I}\right). \quad (10)$$

Here $\delta_{t|t-1} = \delta_t - \delta_{t-1} \frac{(1-m_t)^2}{(1-m_{t-1})^2}$ and $m_t = \frac{t}{T}$.

The reverse process can be derived by Bayesian formula and shown to be a Gaussian process with mean $\mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})$ and variance $\delta'_t \mathbf{I}$:

$$p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) = \frac{p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y})p(\mathbf{x}_{t-1} | \mathbf{x}_0, \mathbf{y})}{p(\mathbf{x}_t | \mathbf{x}_0, \mathbf{y})} = N(\mathbf{x}_{t-1}; \mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}), \delta'_t \mathbf{I}), \quad (11)$$

$$\mu'_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) = \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}}\mathbf{x}_t + (1 - m_{t-1}) \frac{\delta_{t|t-1}}{\delta_t}\mathbf{x}_0 + (m_{t-1} - m_t) \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}}\mathbf{y}, \quad (12)$$

$$\delta'_t = \frac{\delta_{t|t-1}\delta_{t-1}}{\delta_t}. \quad (13)$$

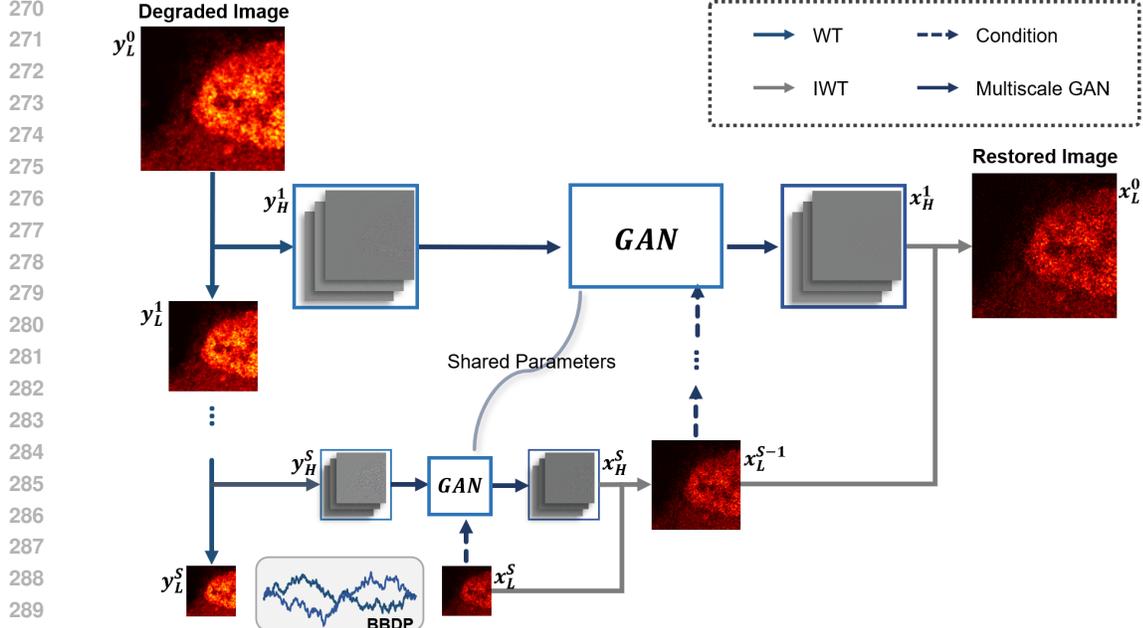


Figure 2: **Schematic diagram of MSCGM.** The conditional image y_L^0 is first decomposed by multi-scale wavelet transform (WT). In the coarsest wavelet layer, a BBDP transforms the low-frequency subband of conditional image to the low-frequency subband of the target image. A multi-scale adversarial learning process transforms subsequent high-frequency subbands of conditional image to the high-frequency subbands of the target image and recovers the full-resolution image using inverse wavelet transform (IWT). y_L^i and y_H^i represent the low- and high-frequency wavelet coefficients of the conditional image at the i^{th} level of wavelet transform, respectively. Similarly, x_L^i and x_H^i denote the low- and high-frequency wavelet coefficients of the target image at the i^{th} level of wavelet transform.

Utilizing the estimate of x_0 by ϵ_θ in Eq. 5, we could eliminate x_0 in Eq. 12 and rewrite Eq. 12 as Li et al. (2023):

$$\begin{aligned} \mu'_t(\mathbf{x}_t, \mathbf{y}) &= \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}} \mathbf{x}_t + (m_{t-1} - m_t) \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}} \mathbf{y} \\ &+ (1 - m_{t-1}) \frac{\delta_{t|t-1}}{\delta_t} (\mathbf{x}_t - \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)) = c_{xt} \mathbf{x}_t + c_{yt} \mathbf{y} - c_{et} \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t), \end{aligned} \quad (14)$$

$$\begin{aligned} c_{xt} &= \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}} + (1 - m_{t-1}) \frac{\delta_{t|t-1}}{\delta_t} \\ c_{yt} &= m_{t-1} - m_t \frac{\delta_{t-1}}{\delta_t} \frac{1 - m_t}{1 - m_{t-1}} \\ c_{et} &= (1 - m_{t-1}) \frac{\delta_{t|t-1}}{\delta_t}. \end{aligned} \quad (15)$$

3.3 MULTI-SCALE CONDITIONAL GENERATIVE MODEL

Wavelet transform, with its theoretical details outlined in Appendix D, is characterized by an orthogonal transform matrix $\mathbf{A} \in \mathbb{R}^{N^2 \times N^2}$. The wavelet transform decomposes an image $\mathbf{x} \in \mathbb{R}^{N^2}$ to one low-frequency (LL) subband $\mathbf{x}_L^1 \in \mathbb{R}^{\frac{N^2}{4}}$ and remaining high-frequency subbands $\mathbf{x}_H^1 \in \mathbb{R}^{\frac{3N^2}{4}}$.

Definition 4 (Multi-scale wavelet decomposition of conditional image generation). With multi-scale wavelet transformation, we can reformulate the conditional probability distribution of \mathbf{x}_0 on \mathbf{y} as

$$p(\mathbf{x}_0 | \mathbf{y}) = \prod_{k=1}^S p(\mathbf{x}_H^k | \mathbf{x}_L^k, \mathbf{y}_H^k) p(\mathbf{x}_L^S | \mathbf{y}_L^S), \quad (16)$$

where S denotes the maximum scale and

$$(\mathbf{x}_H^1, \mathbf{x}_L^1)^T = \mathbf{A} \mathbf{x}_0, (\mathbf{x}_H^{k+1}, \mathbf{x}_L^{k+1})^T = \mathbf{A} \mathbf{x}_L^k, k = 1, \dots \quad (17)$$

Different from existing approaches, our method leverages *BBDP* and *adversarial learning* process inspired by GANs to handle low- and high-frequency subbands at various scales respectively, and the schematic diagram of our model is illustrated in Fig. 2. For the coarsest level low-frequency subband \mathbf{x}_L^S , due to the whitening effect of the low-frequency subband after wavelet transform, DMs can effectively and efficiently approximate $p(\mathbf{x}_L^S | \mathbf{y}_L^S)$ with fewer sampling steps, while generating diverse and photorealistic images. For another, though the conditional distribution of high-frequency subbands deviates from unimodal Gaussian distributions considerably, the multi-scale adversarial learning process is able to approximate their multi-modal distribution and sample the full-resolution images rapidly in a coarse-to-fine style. Since the BBDP at the coarsest level produces samples with good diversity and fidelity, the possibility of mode collapse commonly observed in pure GAN models can be minimized.

We adopt the Wasserstein distances between fake and real images Arjovsky et al. (2017) to optimize the generator G and discriminator D . We adopted a pixel-wise L2 loss and a structural similarity index loss to penalize local and global mismatch, respectively. The training loss for multi-scale adversarial learning process is:

$$L_G = \sum_{k=1}^S [\lambda(G(\mathbf{x}_L^k, \mathbf{z}^k) - \mathbf{x}_H^k)^2 + \nu(1 - \text{SSIM}(G(\mathbf{x}_L^k, \mathbf{z}^k), \mathbf{x}_H^k)) - \alpha D(G(\mathbf{x}_L^k, \mathbf{z}^k))] \quad (18)$$

$$L_D = \sum_{k=1}^S (D(G(\mathbf{x}_L^k, \mathbf{z}^k)) - D(\mathbf{x}_H^k)) \quad (19)$$

Here \mathbf{z}^k refers to random white noise at scale k , $\text{SSIM}(\cdot, \cdot)$ is the structural similarity index measure Wang et al. (2004). The complete sampling process of our model is detailed here:

Algorithm 2 Sampling

```

1: Sample  $\mathbf{y} \sim q(\mathbf{y})$ 
2: Wavelet transform  $S$  times to get  $\{\mathbf{y}_L^S, \mathbf{y}_H^S, \dots, \mathbf{y}_H^1\}$ 
3: for  $t = T, T - 1, \dots, 1$  do
4:    $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
5:    $\mathbf{x}_{t-1,L}^S = c_{xt}\mathbf{x}_{t,L}^S + c_{yt}\mathbf{y}_L^S - c_{et}\epsilon_{\theta}(\mathbf{x}_{t,L}^S, \mathbf{y}_L^S, t) + \sqrt{\delta_t}\mathbf{z}$ 
6: end for
7:  $\mathbf{x}_L^S = \mathbf{x}_{0,L}^S$ 
8: for  $k = S, S - 1, \dots, 1$  do
9:    $\mathbf{x}_H^k = G(\mathbf{x}_L^k, \mathbf{y}_H^k, \mathbf{z}^k)$ 
10:   $\mathbf{x}_L^{k-1} = A^T(\mathbf{x}_H^k, \mathbf{x}_L^k)^T$ 
11: end for
12: Return  $\mathbf{x}_0 = \mathbf{x}_L^0$ 

```

4 EXPERIMENTS

In this section, we first elucidate the design and training details of our method, as well as the preparation of training and testing datasets. Then, we evaluate our method on various image restoration tasks in computational and microscopy imaging, and compare it with baseline methods.

4.1 EXPERIMENTAL SETUP AND IMPLEMENTATION DETAILS

For the BBDP at the coarsest wavelet scale, we adapt the UNet architecture with multi-head attention layers as practiced in Nichol & Dhariwal (2021). The number of sampling steps is set as 1000 for training. The Brownian bridge diffusion model (BBDM) baseline Li et al. (2023) is implemented at the full resolution scale without wavelet decomposition, and the same 1000 discretization step is used for training. The training of IR-SDE (image restoration stochastic differential equation) Luo et al. (2023a) and Refusion Luo et al. (2023a) follow the training setups as their original setups.

Inspired by Chen et al. (2022), the generator adopt a similar architecture to NAFNet. The generator contains 36 NAFBlocks distributed at 4 scales. A 2×2 convolutional layer with stride 2 doubling the

Methods	Trainable Params.↓	Sampling Time (s)↓	PSNR (dB)↑			SSIM↑		
			DIV2K	Set5	Set14	DIV2K	Set5	Set14
IR-SDE	135.3M	19.51	23.54	26.73	22.71	0.56	0.72	0.54
ReFusion	131.4M	17.25	21.39	22.82	22.13	0.43	0.52	0.49
BBDM	124.7M	32.29	31.50	31.60	30.39	0.66	0.76	0.68
MSCGM	192.5M	2.28	31.66	32.33	30.79	0.72	0.85	0.71

Table 1: **Comparison of IR-SDE, ReFusion, BBDM and our method (MSCGM) on $4\times$ super resolution experiment.** Sample steps are set as 1000 for all methods. Metrics are calculated on 256×256 center-cropped patches of DIV2K validation set, Set 5 and Set 14. Entries in bold indicate the best performance achieved among the compared methods. Sampling time is measured at the same resource cost.

channels connects adjacent scales in the encoding (downsampling) path, and a 1×1 convolutional layer with pixel shuffle layer connects adjacent scales in the decoding (upsampling)superre path. The number of channels in the first NAFBlock is 64. The discriminator consists of 5 convolutional blocks and 2 dense layers at the end, and each convolutional block halves the spatial dimension but doubles the number of channels. The number of channels for the first convolutional block in the discriminator is 64. The details about training and dataset are elucidated in Appendix F.

4.2 EVALUATION METRICS

Peak Signal-to-Noise Ratio (PSNR) is commonly used to measure the quality of reconstruction in generated images, with higher values indicating better image quality. Structural Similarity Index Measure (SSIM) Wang et al. (2004) assesses the high-level quality of images by focusing on changes in structural information, luminance, and contrast. Fréchet Inception Distance (FID) score Heusel et al. (2017) is used in generative models like GANs to compare the distribution of generated images against real ones, where lower FID values imply images more similar to real ones, indicating higher quality.

4.3 RESULTS AND COMPARISON

In this section, we first evaluate and compare our method with the current state-of-the-art model (cross-modality super-resolution (CMSR)Wang et al. (2019)) on two microscopy image restoration tasks with different samples, and the experimental results demonstrate that our method achieves the **best performance** on these tasks. To further validate the model’s generalization capability on standard nature image datasets, we then conduct comparative evaluations on three different natural image restoration tasks against multiple baselines, achieving similarly excellent restoration results.

First, we apply our method to microscopy images, where the degradation process is complex and unknown. Given the pronounced contrast and sparsity inherent in microscopy images, it is crucial to use generative models capable of handling multi-modal distributions to adapt effectively to complex microscopy datasets. We utilize our method to perform super-resolution on microscopy images of nano-beads and HeLa cells, transforming diffraction-limited confocal images to achieve resolution beyond the optical diffraction limit and match the image quality of STED microscopy. Next, we assess the adaptability of our method to various image restorations tasks of natural images, including a $4\times$ super-resolution task on DIV2K dataset, a shadow removal task on natural images (ISTD dataset) and on a low-light image enhancement task on natural images (LOL dataset). Through comparison against competitive methods on various testbeds, we demonstrate the superior effectiveness and versatility of our method for image restoration.

4.3.1 MICROSCOPY IMAGE RESTORATION

Microscopy Image Super-resolution: We evaluate our method on microscopy image super-resolution tasks and compare it with existing generative models in this field. Unlike natural image super-resolution, the LR images are not downsampled but sampled at the same spatial frequency as the HR images. However, the LR images are limited by the optical diffraction limit, which is equivalent to a convolution operation on the HR images with a low-pass point spread function (PSF). We apply our method to confocal (LR) images of fluorescence nanobeads to evaluate its capability to overcome the

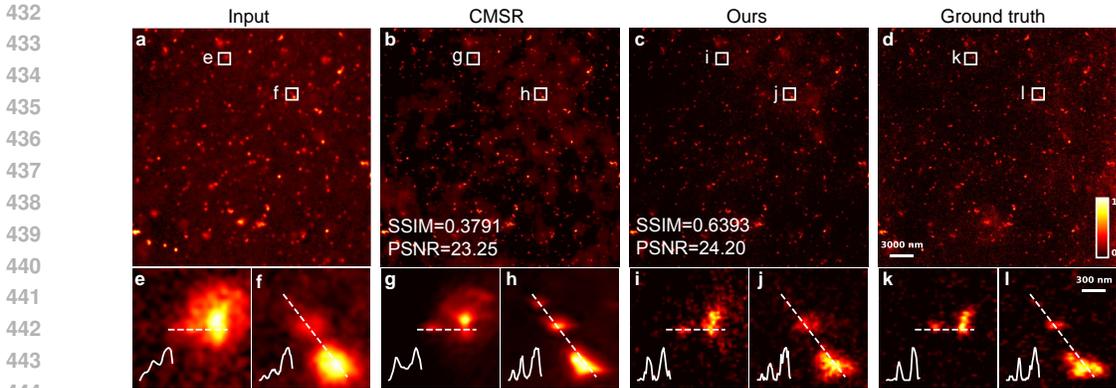


Figure 3: **Comparison of CMSR Wang et al. (2019) and our method on microscopy image super-resolution of nanobeads.** (a) Input images captured by a confocal microscope, (b, c) SR outputs of CMSR and our method, and (d) ground truths captured by an STED microscope of the same FOV. (e-l) Zoom-in regions marked by the corresponding white boxes in (a-c). Cross-section intensity values along the dashed line are plotted.

optical diffraction limit (see Methods for sample and dataset details). Figure 3 illustrates one typical field-of-view (FOV) of nanobead samples (~ 20 -nm) captured using confocal and STED microscopy. The dimensions of nanobeads are considerably smaller than the optical diffraction limit (~ 250 -nm) and nearby beads cannot be distinguished in confocal images. MSCGM and CMSR model are trained on the same training data and learn to transform LR confocal images to match HR STED images.

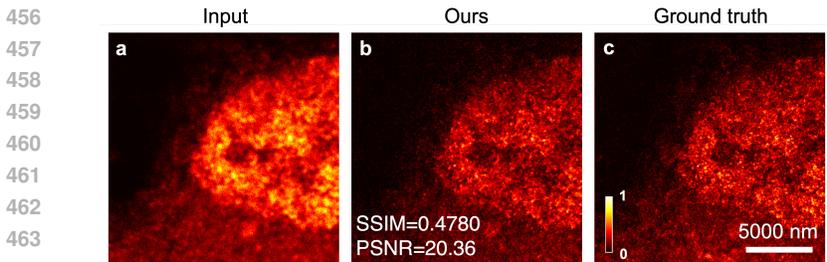


Figure 4: **Microscopy image super-resolution of our method on HeLa cells.** (a) LR confocal image, (b) SR output image and (c) HR STED image of the same FOV.

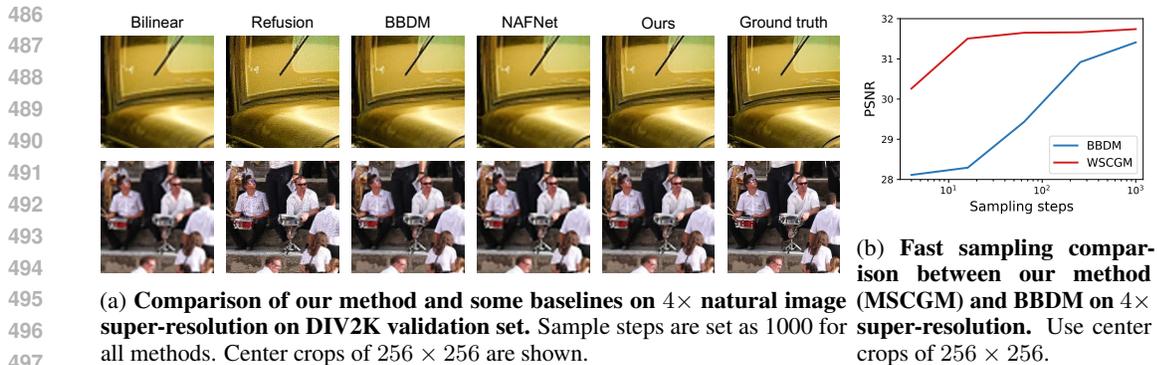
ages. As another demonstration, we apply our method to fluorescence imaging of HeLa cells and present Fig. 4, Appendix Fig. 10, 11 to further support its success in microscopy image restoration.

4.3.2 EVALUATION ON STANDARD IMAGE DATASETS

To further verify the generalization and scalability of our MSCGM, we conducted more experiments on three different standard natural image datasets.

Natural Image Super-resolution: We train the three methods on the DIV2K training dataset and then test them on the DIV2K validation set, Set5, and Set14. Table 1 quantifies the super-resolution performance in terms of PSNR and SSIM on the three test sets. The reported MSCGM scores indicate better restoration quality over competitive methods. Table 1 also presents the sampling time of each method under the same resource cost (GPU memory). Our method (MSCGM) achieves superior image generation results and exhibits a 10x improvement in average single-image generation time compared to IR-SDE, an 8x improvement compared to ReFusion, and a 16x improvement compared to BBDM. Additional visualization results of the three methods on the DIV2K validation set are displayed in Fig. 5a, showing better reconstruction fidelity and perceptual quality.

Moreover, we implement fast sampling with fewer sampling steps on the same super-resolution tasks, from 4 to 1000 steps, and depict the results in Fig. 5b. Additional comparative samples, showcasing



495
496
497
498
499
Figure 5: Visualization and fast sampling results on natural image super-resolution

500 the performance of our model against BBDM at fewer sampling steps, are provided in Appendix Figs.
501 12 and 13.

502 **Natural Image Shadow Removal:** For image shadow removal task, we train our paradigm on
503 the ISTD training dataset Wang et al. (2018a) of 1330 image triplets (shadow image, mask and
504 clean image) and evaluate it on the ISTD test set of 540 triplets. Appendix Table 4 summarizes the
505 performance of our method against competitive methods, including DC-ShadowNet Jin et al. (2023),
506 ST-CGAN Wang et al. (2018b), DSC Hu et al. (2019), DHAN Cun et al. (2020), BMNet Zhu et al.
507 (2022), ShadowFormer Guo et al. (2023) and BBDM Li et al. (2023). Appendix Fig. 9 further show
508 the visualization results of our method and competitive methods. In summary, our method effectively
509 restores high-quality images from shadowed and low-light images while minimizing the artifacts and
510 inconsistency in the output images.



522
523
524
525
526
527
528
529
530
Figure 6: Image enhancement results on LOL dataset. Center-cropped 256×256 patches shown.

Low-light Natural Image Enhancement: In this task, we train and evaluate our paradigm on
LOL dataset Wei et al. (2018), which contains 485 training and 15 testing image pairs. Figure 6
qualitatively showcases the results of our method in comparison with RetinexFormer Cai et al. (2023)
and BBDM. Our method shows better image enhancement performance, and the restored images have
more natural colors and less noise details. This effectiveness of our method on image enhancement is
further confirmed by the quantitative evaluation results against RetinexFormer, HWMNet Fan et al.
(2022) and BBDM, as shown in Appendix Table 3.

531 CONCLUSION

532
533
534
535
536
537
538
539
To address the limitations of existing diffusion models in conditional image restoration, we demon-
strate a novel generative model for image restoration based on Brownian bridge process and multi-
scale wavelet transform. By factorizing the image restoration process in the multi-scale wavelet
domains, we utilize Brownian bridge diffusion process and generative adversarial networks to recover
different wavelet subbands according to their distribution properties, consequently accelerate the
sampling speed significantly and achieve high sample quality and diversity competitive to diffusion
model baselines. Future implementation could integrate standard acceleration techniques for diffusion
models.

REFERENCES

- 540
541
542 Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution:
543 Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*
544 *Workshops*, July 2017.
- 545 Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017. URL [http://](http://arxiv.org/abs/1701.07875)
546 arxiv.org/abs/1701.07875. cite arxiv:1701.07875.
- 547 Robert B Ash and Catherine A Doléans-Dade. *Probability and measure theory*. Academic press,
548 2000.
- 550 George Barbastathis, Aydogan Ozcan, and Guohai Situ. On the use of deep learning for computational
551 imaging. *Optica*, 6(8):921–943, 2019.
- 552 Georgios Batzolis, Jan Stanczuk, Carola-Bibiane Schönlieb, and Christian Etmann. Conditional
553 image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*, 2021.
- 554 Yuanhao Cai, Hao Bian, Jing Lin, Haoqian Wang, Radu Timofte, and Yulun Zhang. Retinex-
555 former: One-stage retinex-based transformer for low-light image enhancement. *arXiv preprint*
556 *arXiv:2303.06705*, 2023.
- 557
558 Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration.
559 In *European Conference on Computer Vision*, pp. 17–33. Springer, 2022.
- 560
561 Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Star-
562 gan: Unified generative adversarial networks for multi-domain image-to-image translation. In
563 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- 564
565 Xiaodong Cun, Chi-Man Pun, and Cheng Shi. Towards ghost-free shadow removal via dual hierar-
566 chical aggregation network and shadow matting gan. In *Proceedings of the AAAI Conference on*
567 *Artificial Intelligence*, volume 34, pp. 10680–10687, 2020.
- 568 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
569 *in neural information processing systems*, 34:8780–8794, 2021.
- 570
571 Magnus Ekström. A general central limit theorem for strong mixing sequences. *Statistics &*
572 *Probability Letters*, 94:236–238, 2014.
- 573
574 Chi-Mao Fan, Tsung-Jung Liu, and Kuan-Hsien Liu. Half wavelet attention on m-net+ for low-light
575 image enhancement. In *2022 IEEE International Conference on Image Processing (ICIP)*, pp.
576 3878–3882. IEEE, 2022.
- 577
578 Douglas F Fraser, James F Gilliam, Michael J Daley, An N Le, and Garrick T Skalski. Explaining
579 leptokurtic movement distributions: intrapopulation variation in boldness and exploration. *The*
American Naturalist, 158(2):124–135, 2001.
- 580
581 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
582 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural informa-*
tion processing systems, 2014.
- 583
584 Richard A Groeneveld and Glen Meeden. Measuring skewness and kurtosis. *Journal of the Royal*
585 *Statistical Society Series D: The Statistician*, 33(4):391–399, 1984.
- 586
587 Jie Gui, Zhenan Sun, Yonggang Wen, Dacheng Tao, and Jieping Ye. A review on generative
588 adversarial networks: Algorithms, theory, and applications. *IEEE transactions on knowledge and*
data engineering, 35(4):3313–3332, 2021.
- 589
590 Lanqing Guo, Siyu Huang, Ding Liu, Hao Cheng, and Bihan Wen. Shadowformer: global context
591 helps shadow removal. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
592 pp. 710–718, 2023.
- 593
Gaurav Gupta, Xiongye Xiao, and Paul Bogdan. Multiwavelet-based operator learning for differential
equations. *Advances in neural information processing systems*, 34:24048–24062, 2021.

- 594 Gaurav Gupta, Xiongye Xiao, Radu Balan, and Paul Bogdan. Non-linear operator approximations
595 for initial value problems. In *International Conference on Learning Representations (ICLR)*, 2022.
596
- 597 Florentin Guth, Simon Coste, Valentin De Bortoli, and Stephane Mallat. Wavelet score-based
598 generative modeling. *Advances in Neural Information Processing Systems*, 35:478–491, 2022.
599
- 600 Kerstin Hammernik, Teresa Klatzer, Erich Kobler, Michael P Recht, Daniel K Sodickson, Thomas
601 Pock, and Florian Knoll. Learning a variational network for reconstruction of accelerated mri data.
602 *Magnetic resonance in medicine*, 79(6):3055–3071, 2018.
- 603 Christopher E Heil and David F Walnut. Continuous and discrete wavelet transforms. *SIAM review*,
604 31(4):628–666, 1989.
605
- 606 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans
607 trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural
608 information processing systems*, 30, 2017.
- 609 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in
610 neural information processing systems*, 33:6840–6851, 2020.
611
- 612 Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans.
613 Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning
614 Research*, 23(1):2249–2281, 2022.
- 615 Xiaowei Hu, Chi-Wing Fu, Lei Zhu, Jing Qin, and Pheng-Ann Heng. Direction-aware spatial context
616 features for shadow detection and removal. *IEEE transactions on pattern analysis and machine
617 intelligence*, 42(11):2795–2808, 2019.
618
- 619 Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with
620 conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and
621 Pattern Recognition (CVPR)*, 2017.
- 622 Björn Jawerth and Wim Sweldens. An overview of wavelet based multiresolution analyses. *SIAM
623 review*, 36(3):377–412, 1994.
624
- 625 Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two transformers can make one strong
626 gan. *arXiv preprint arXiv:2102.07074*, 1(3), 2021.
627
- 628 Yeying Jin, Aashish Sharma, and Robby T. Tan. Dc-shadownet: Single-image hard and soft shadow
629 removal using unsupervised domain-classifier guided network, 2023.
- 630 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative
631 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
632 recognition*, pp. 4401–4410, 2019.
633
- 634 Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. On convergence and stability of gans.
635 *arXiv preprint arXiv:1705.07215*, 2017.
- 636 Zhifeng Kong and Wei Ping. On fast sampling of diffusion probabilistic models. *arXiv preprint
637 arXiv:2106.00132*, 2021.
638
- 639 Cansu Korkmaz, A Murat Tekalp, and Zafer Dogan. Training generative image super-resolution
640 models by wavelet-domain losses enables better control of artifacts. In *Proceedings of the
641 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5926–5936, 2024.
- 642 Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan:
643 Blind motion deblurring using conditional adversarial networks. In *Proceedings of the IEEE
644 conference on computer vision and pattern recognition*, pp. 8183–8192, 2018.
645
- 646 Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. Bbdm: Image-to-image translation with brownian
647 bridge diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
Pattern Recognition*, pp. 1952–1961, 2023.

- 648 Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting
649 Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*,
650 479:47–59, 2022.
- 651
652 Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Im-
653 age restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference*
654 *on computer vision*, pp. 1833–1844, 2021.
- 655 Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks.
656 *Advances in neural information processing systems*, 30, 2017.
- 657
658 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
659 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*
660 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 661 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
662 *arXiv:1711.05101*, 2017.
- 663
664 Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Image restoration
665 with mean-reverting stochastic differential equations. *arXiv preprint arXiv:2301.11699*, 2023a.
- 666
667 Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Refusion:
668 Enabling large-size realistic image restoration with latent-space diffusion models. In *Proceedings*
669 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1680–1691, 2023b.
- 670
671 William Meinel, Jean-Christophe Olivo-Marin, and Elsa D. Angelini. Denoising of microscopy
672 images: A review of the state-of-the-art, and a new sparsity-based method. *IEEE Transactions on*
673 *Image Processing*, 27(8):3842–3856, 2018. doi: 10.1109/TIP.2018.2819821.
- 674
675 Kamyar Nazeri, Eric Ng, and Mehran Ebrahimi. Image colorization using generative adversarial
676 networks. In *Articulated Motion and Deformable Objects: 10th International Conference, AMDO*
677 *2018, Palma de Mallorca, Spain, July 12-13, 2018, Proceedings 10*, pp. 85–94. Springer, 2018.
- 678
679 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
680 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 681
682 Hao Phung, Quan Dao, and Anh Tran. Wavelet diffusion models are fast and scalable image generators.
683 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
684 10199–10208, 2023.
- 685
686 Mikko Ranta. *Wavelet multiresolution analysis of financial time series*. Vaasan yliopisto, 2010.
- 687
688 Yair Rivenson, Zoltán Göröcs, Harun Günaydin, Yibo Zhang, Hongda Wang, and Aydogan Ozcan.
689 Deep learning microscopy. *Optica*, 4(11):1437–1443, 2017.
- 690
691 Yair Rivenson, Hongda Wang, Zhensong Wei, Kevin de Haan, Yibo Zhang, Yichen Wu, Harun
692 Günaydin, Jonathan E Zuckerman, Thomas Chong, Anthony E Sisk, et al. Virtual histological
693 staining of unlabelled tissue-autofluorescence images via deep learning. *Nature biomedical*
694 *engineering*, 3(6):466–477, 2019.
- 695
696 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
697 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
698 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 699
700 Murray Rosenblatt. A central limit theorem and a strong mixing condition. *Proceedings of the*
701 *national Academy of Sciences*, 42(1):43–47, 1956.
- 702
703 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet,
704 and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022*
705 *Conference Proceedings*, pp. 1–10, 2022a.
- 706
707 Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi.
708 Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and*
709 *Machine Intelligence*, 45(4):4713–4726, 2022b.

- 702 Mark J Shensa et al. The discrete wavelet transform: wedding the a trous and mallat algorithms.
703 *IEEE Transactions on signal processing*, 40(10):2464–2482, 1992.
704
- 705 Oren Solomon, Maor Mutzafi, Mordechai Segev, and Yonina C. Eldar. Sparsity-based super-resolution
706 microscopy from correlation information. *Opt. Express*, 26(14):18238–18269, Jul 2018. doi:
707 10.1364/OE.26.018238. URL [https://opg.optica.org/oe/abstract.cfm?URI=](https://opg.optica.org/oe/abstract.cfm?URI=oe-26-14-18238)
708 [oe-26-14-18238](https://opg.optica.org/oe/abstract.cfm?URI=oe-26-14-18238).
- 709 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
710 *preprint arXiv:2010.02502*, 2020a.
711
- 712 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
713 *Advances in neural information processing systems*, 32, 2019.
714
- 715 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
716 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
717 *arXiv:2011.13456*, 2020b.
- 718 Radomir S Stanković and Bogdan J Falkowski. The haar wavelet transform: its status and achieve-
719 ments. *Computers & Electrical Engineering*, 29(1):25–44, 2003.
720
- 721 Hongda Wang, Yair Rivenson, Yiyin Jin, Zhensong Wei, Ronald Gao, Harun Günaydın, Laurent A
722 Bentolila, Comert Kural, and Aydogan Ozcan. Deep learning enables cross-modality super-
723 resolution in fluorescence microscopy. *Nature methods*, 16(1):103–110, 2019.
724
- 725 Jifeng Wang, Xiang Li, and Jian Yang. Stacked conditional generative adversarial networks for jointly
726 learning shadow detection and shadow removal. In *CVPR*, 2018a.
- 727 Jifeng Wang, Xiang Li, and Jian Yang. Stacked conditional generative adversarial networks for
728 jointly learning shadow detection and shadow removal. In *Proceedings of the IEEE conference on*
729 *computer vision and pattern recognition*, pp. 1788–1797, 2018b.
730
- 731 Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-
732 resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of*
733 *the IEEE Conference on Computer Vision and Pattern Recognition*, 2018c.
- 734 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from
735 error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612,
736 2004.
737
- 738 Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light
739 enhancement. In *British Machine Vision Conference*, 2018.
- 740 Martin Weigert, Uwe Schmidt, Tobias Boothe, Andreas Müller, Alexandr Dibrov, Akanksha Jain,
741 Benjamin Wilhelm, Deborah Schmidt, Coleman Broaddus, Siân Culley, et al. Content-aware image
742 restoration: pushing the limits of fluorescence microscopy. *Nature methods*, 15(12):1090–1097,
743 2018.
744
- 745 Christopher Stroude Withers. Central limit theorems for dependent variables. i. *Zeitschrift für*
746 *Wahrscheinlichkeitstheorie und verwandte Gebiete*, 57(4):509–534, 1981.
747
- 748 Yichen Wu, Yair Rivenson, Hongda Wang, Yilin Luo, Eyal Ben-David, Laurent A Bentolila, Christian
749 Pritz, and Aydogan Ozcan. Three-dimensional virtual refocusing of fluorescence microscopy
750 images using deep learning. *Nature methods*, 16(12):1323–1331, 2019.
- 751 Xiongye Xiao, Defu Cao, Ruochen Yang, Gaurav Gupta, Gengshuo Liu, Chenzhong Yin, Radu Balan,
752 and Paul Bogdan. Coupled multiwavelet neural operator learning for coupled partial differential
753 equations. *arXiv preprint arXiv:2303.02304*, 2023.
754
- 755 Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with
denoising diffusion gans. *arXiv preprint arXiv:2112.07804*, 2021.

756 Qingsong Yang, Pingkun Yan, Yanbo Zhang, Hengyong Yu, Yongyi Shi, Xuanqin Mou, Mannudeep K
757 Kalra, Yi Zhang, Ling Sun, and Ge Wang. Low-dose ct image denoising using a generative
758 adversarial network with wasserstein distance and perceptual loss. *IEEE transactions on medical*
759 *imaging*, 37(6):1348–1357, 2018.

760 Zili Yi, Hao Zhang, Ping Tan, and Minglun Gong. Dualgan: Unsupervised dual learning for image-
761 to-image translation. In *Proceedings of the IEEE international conference on computer vision*, pp.
762 2849–2857, 2017.

764 Yueqin Yin, Lianghua Huang, Yu Liu, and Kaiqi Huang. Diffgar: Model-agnostic restoration from
765 generative artifacts using image-to-image diffusion models. In *Proceedings of the 2022 6th*
766 *International Conference on Computer Science and Artificial Intelligence*, pp. 55–62, 2022.

767 Weisong Zhao, Shiqun Zhao, Liuju Li, Xiaoshuai Huang, Shijia Xing, Yulin Zhang, Guohua Qiu,
768 Zhenqian Han, Yingxu Shang, De-en Sun, et al. Sparse deconvolution improves the resolution of
769 live-cell super-resolution fluorescence microscopy. *Nature biotechnology*, 40(4):606–617, 2022.

770 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation
771 using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference*
772 *on computer vision*, pp. 2223–2232, 2017.

774 Yurui Zhu, Jie Huang, Xueyang Fu, Feng Zhao, Qibin Sun, and Zheng-Jun Zha. Bijective mapping
775 network for shadow removal. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
776 *and Pattern Recognition*, pp. 5627–5636, 2022.

777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810 A CODE AND DATA AVAILABILITY

811
812 The codes of our reported method is available at <https://anonymous.4open.science/r/MSCGM-E114>. The DIV2K dataset Agustsson & Timofte (2017) is obtained from <https://data.vision.ee.ethz.ch/cvl/DIV2K/>, and the ISTD dataset Wang et al. (2018a) is obtained
814 from <https://github.com/DeepInsight-PCALab/ST-CGAN>. The microscopy image
815 datasets (nanobeads and HeLa cells) are requested from Wang et al. Wang et al. (2019) and partial
816 demo images are uploaded to <https://anonymous.4open.science/r/MSCGM-E114>.
817

818 B SCORE REGULARITY FOR DISCRETIZATION

819 **Theorem 5.** Suppose the Gaussian distribution $p = N(0, \Sigma)$ and distribution \tilde{p}_0 from time reversed
820 SDE, the Kullback-Leibler divergence between p and \tilde{p}_0 relates to the covariance matrix Σ as:
821 $KL(p \parallel \tilde{p}_0) \leq \Psi_T + \Psi_{\Delta t} + \Psi_{T, \Delta t}$, with:

$$822 \Psi_T = f(e^{-4T} |\text{Tr}((\Sigma - \text{Id})\Sigma)|), \quad (20)$$

$$823 \Psi_{\Delta t} = f(\Delta t |\text{Tr}(\Sigma^{-1} - \Sigma(\Sigma - \text{Id})^{-1} \log(\Sigma)/2 + (\text{Id} - \Sigma^{-1})/3)|), \quad (21)$$

$$824 \Psi_{T, \Delta t} = o(\Delta t + e^{-4T}), \quad \Delta t \rightarrow 0, T \rightarrow +\infty \quad (22)$$

825 where $f(t) = t - \log(1 + t)$ and d is the dimension of Σ , $\text{Tr}(\Sigma) = d$.

826 **Proposition 1.** For any $\epsilon > 0$, there exists $T, \Delta t \geq 0$ such that:

$$827 (1/d)(\Psi_T + \Psi_{\Delta t}) \leq \epsilon, \quad (23)$$

$$828 T/\Delta t \leq C\epsilon^{-2}\kappa^3, \quad (24)$$

829 where $C \geq 0$ is a universal constant, and κ is the condition number of Σ .

830 Guth et al. (2022) provides the proof outline for Theorem 5, based on the following Theorem 6,

831 **Theorem 6.** Let $N \in \mathbb{N}$, $\Delta t > 0$, and $T = N\Delta t$. Then, we have that $\bar{x}_t^N \sim N(\hat{\mu}_N, \Sigma^{\hat{N}})$ with

$$832 \Sigma^{\hat{N}} = \Sigma + \exp(-4T)\Sigma^{\hat{T}} + \Delta t\Psi^{\hat{T}} + (\Delta t)^2 R^{\hat{T}, \Delta t}, \quad (25)$$

$$833 \hat{\mu}_N = \mu + \exp(-2T)\hat{\mu}_T + \Delta t e^{\hat{T}} + \frac{(\Delta t)^2}{2} r^{T, \Delta t}, \quad (26)$$

834 where $\Sigma^{\hat{T}}, \Psi^{\hat{T}}, R^{\hat{T}, \Delta t} \in \mathbb{R}^{d \times d}$, $\hat{\mu}_T, e^{\hat{T}}, r^{T, \Delta t} \in \mathbb{R}^d$, and $\|R^{\hat{T}, \Delta t}\| + \|r^{T, \Delta t}\| \leq R$, not dependent
835 on $T \geq 0$ and $\Delta t > 0$. We have that

$$836 \Sigma^{\hat{T}} = -(\Sigma - \text{Id})(\Sigma\Sigma^{-1})^2, \quad (27)$$

$$837 \Psi^{\hat{T}} = \text{Id} - \frac{1}{2}\Sigma^2(\Sigma - \text{Id})^{-1} \log(\Sigma) + \exp(-2T)\Psi^{\tilde{T}}. \quad (28)$$

838 In addition, we have

$$839 \hat{\mu}_T = -\Sigma^{-1}T\Sigma\mu, \quad (29)$$

$$840 e^{\hat{T}} = \left\{ -2\Sigma^{-1} - \frac{1}{4}\Sigma(\Sigma - \text{Id})^{-1} \log(\Sigma) \right\} \mu + \exp(-2T)\tilde{\mu}_T, \quad (30)$$

841 with $\Psi^{\tilde{T}}, \tilde{\mu}_T$ bounded and not dependent on T .

842 **Theorem 7.** Suppose that $\nabla \log p_t(x)$ is φ^2 in both t and x such that:

$$843 \sup_{x, t} \|\nabla^2 \log p_t(x)\| \leq K, \quad \|\partial_t \nabla \log p_t(x)\| \leq Me^{-\alpha t} \|x\| \quad (31)$$

844 for some $K, M, \alpha > 0$. Then, $\|p - \tilde{p}_0\|_{TV} \leq \Psi_T + \Psi_{\Delta t} + \Psi_{T, \Delta t}$, where:

$$845 \Psi_T = \sqrt{2}e^{-T} \text{KL}(p \parallel N(0, \text{Id}))^{1/2} \quad (32)$$

$$846 \Psi_{\Delta t} = 6\sqrt{\Delta t} \left[1 + \mathbb{E}_p(\|x\|^4)^{1/4} \right] \left[1 + K + M(1 + 1/2\alpha)^{1/2} \right] \quad (33)$$

$$847 \Psi_{T, \Delta t} = o\left(\sqrt{\Delta t} + e^{-T}\right) \quad \Delta t \rightarrow 0, T \rightarrow +\infty \quad (34)$$

848 Theorem 7 generalizes Theorem 5 to non-Gaussian processes. Please refer to Guth et al. (2022) for
849 the complete proof.

C CHARACTERISTICS OF HIGH AND LOW FREQUENCY COEFFICIENTS IN THE WAVELET DOMAIN

C.1 GAUSSIAN TENDENCY OF LOW-FREQUENCY COEFFICIENTS IN HIGHER SCALES

In an image, pixel intensities are represented as random variables, with adjacent pixels exhibiting correlation due to their spatial proximity. This correlation often follows a power-law decay:

$$C(d) = \frac{1}{(1 + \alpha d)^\beta}, \tag{35}$$

where $C(d)$ is the correlation between pixels separated by distance d , and α and β characterize the rate of decay.

The wavelet transform (i.e., Haar wavelet transform), particularly its down-sampling step, increases the effective distance d among pixels, thereby reducing their original spatial correlation. This reduction is crucial for applying the generalized Central Limit Theorem Rosenblatt (1956); Withers (1981); Ekström (2014); Ash & Doléans-Dade (2000), which requires that the individual variables (pixels, in this case) are not strongly correlated.

At scale k in the wavelet decomposition, the low-frequency coefficients, \bar{X}_k , representing the average intensity over n_k pixels, are calculated as:

$$\bar{X}_k = \frac{1}{n_k}(X_1 + X_2 + \dots + X_{n_k}), \tag{36}$$

where n_k is the number of pixels in each group at scale k .

As the scale increases, the effect of averaging over larger groups of pixels, combined with the reduced correlation due to down-sampling, leads to a scenario where the generalized Central Limit Theorem can be applied. Consequently, the distribution of \bar{X}_k tends towards a Gaussian distribution:

$$\bar{X}_k \xrightarrow{d} N(\mu_k, \frac{\sigma_k^2}{n_k}), \tag{37}$$

where μ_k and σ_k^2 are the mean and variance of the averaged intensities at scale k , respectively. This Gaussian tendency becomes more pronounced at higher scales due to the combination of reduced pixel correlation and the averaging process.

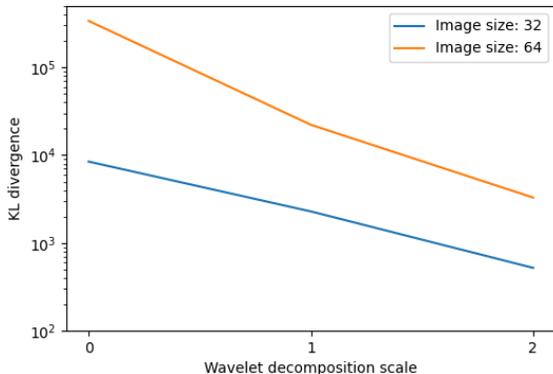


Figure 7: KL divergence between the standard normal distribution and normalized sample distribution with respect to the wavelet scale. Images were sampled from DIV2K dataset by 32 × 32 and 64 × 64 patches.

In Fig. 7, we quantify the Gaussianity of low-frequency wavelet subbands at different scales using the Kullback-Leibler (KL) divergence, which measures the distance between the standard normal distribution and normalized sample distribution of the low-frequency wavelet coefficients. We

918 sampled random patches of different resolutions (32×32 and 64×64) from DIV2K dataset to
 919 calculate KL divergence. With the increasing of the scale, KL divergence decreases, validating
 920 the Gaussian tendency of low-frequency subbands after multi-scale wavelet transforms. Figure 8
 921 further validates this tendency by plotting the kurtosis of sample distribution of microscopy images
 922 of nanobeads with respect to the scale.

924 C.2 SPARSITY AND NON-GAUSSIANITY OF HIGH-FREQUENCY COEFFICIENTS

925 High-frequency coefficients, when analyzed through wavelet transform, exhibit a distinct property
 926 of sparsity, characterized by a majority of wavelet coefficients being near or at zero, with only a
 927 sparse representation of significant non-zero coefficients. This sparsity highlights the efficiency of
 928 wavelet transforms in encoding signal details and abrupt changes. Furthermore, these high-frequency
 929 components often deviate from Gaussian distributions, tending towards leptokurtic distributions Fraser
 930 et al. (2001) with higher peaks and heavier tails. This non-Gaussian nature suggests a concentration
 931 of energy in fewer coefficients and is crucial in applications like signal denoising and compression,
 932 where recognizing and preserving these vital characteristics is paramount.

933 In the following proposition, we theoretically show that the conditional distribution of \mathbf{x}_H^k on \mathbf{x}_L^k
 934 exhibits highly non-Gaussian properties and yields sparse samples. For a given image \mathbf{x} and threshold
 935 t , the sparsity of its high-frequency coefficients at k -scale is defined as:

$$936 s(\mathbf{x}_H^k) = \frac{\|\mathbf{1}\{\mathbf{x}_H^k \leq t\}\|}{L^2}, \quad k = 1, 2, \dots \quad (38)$$

939 Here $\|\cdot\|$ is the norm counting the number of 1s in the vector. In this way, we could estimate the
 940 expected sparsity of the true marginal distribution $p(\mathbf{x}_H^k)$. Considering that the LL coefficients with
 941 approximate Gaussian distribution given the whitening effect of wavelet decomposition, we have the
 942 following proposition.

943 **Proposition 2.** *For a sufficiently large k , if the expected sparsity of \mathbf{x}_H^k has a lower bound α*

$$944 \mathbb{E}(s(\mathbf{x}_H^k)) \geq \alpha, \quad (39)$$

945 where $\alpha \in [0, 1]$. Then the conditional expected sparsity of \mathbf{x}_H^k on \mathbf{x}_L^k is bounded by

$$946 \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) \geq \alpha - \varepsilon, \quad (40)$$

947 where $\varepsilon > 0$ is a small positive number determined by k .

950 *Proof.* According to Eq. 37, for a sufficiently large k we could assume that

$$951 \int |p(\mathbf{x}_L^k) - f_k(\mathbf{x}_L^k)| d\mathbf{x}_L^k \leq \varepsilon, \quad (41)$$

952 where $f_k(\mathbf{x}_L^k)$ is the PDF of standard Gaussian distribution. Notice that

$$953 \mathbb{E}(s(\mathbf{x}_H^k)) = \iint s(\mathbf{x}_H^k) p(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) p(\mathbf{x}_L^k) d\mathbf{x}_L^k ds \quad (42)$$

$$954 = \int \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) p(\mathbf{x}_L^k) d\mathbf{x}_L^k \geq \alpha \quad (43)$$

955 Since s is a bounded function in $[0, 1]$, $\mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k)$ has an uniform lower bound with respect to all
 956 \mathbf{x}_L^k , denoted as α' . In other words, there exists $\alpha' \in [0, 1]$ such that

$$957 \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) \geq \alpha', \quad \forall \mathbf{x}_L^k \quad (44)$$

958 We can get

$$959 \int \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) p(\mathbf{x}_L^k) d\mathbf{x}_L^k \\
 960 = \int \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) f_k(\mathbf{x}_L^k) d\mathbf{x}_L^k \\
 961 + \int \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) (p(\mathbf{x}_L^k) - f_k(\mathbf{x}_L^k)) d\mathbf{x}_L^k \geq \alpha' \quad (45)$$

Similarly, it is easy to see that 1 is a trivial uniform upper bound for $\mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k)$. Thus,

$$\mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k) \geq \alpha' = \int \mathbb{E}(s(\mathbf{x}_H^k)|\mathbf{x}_L^k)p(\mathbf{x}_L^k)d\mathbf{x}_L^k - \varepsilon \geq \alpha - \varepsilon. \quad (46)$$

□

C.3 QUANTIFYING NON-GAUSSIANITY OF DATASETS

Third- and fourth-order sample cumulants, i.e., skewness and kurtosis, to quantify the non-Gaussianity of certain sample distributions Groeneveld & Meeden (1984). The non-Gaussianity of high-frequency subbands can be evidenced by the kurtosis plot with respect to wavelet scales in Fig. 8. The kurtosis of high-frequency subbands of microscopy images increases with the wavelet scales, showing the high non-Gaussianity of the distribution of high-frequency coefficients.

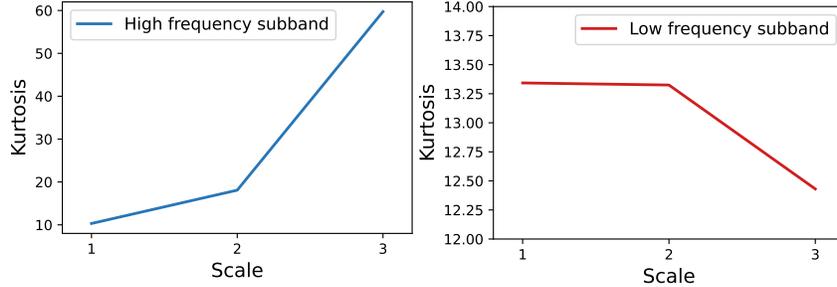


Figure 8: Kurtosis of (left) low-frequency and (right) high-frequency subband coefficients with respect to the wavelet scale. Metrics were calculated on microscopy images of nanobeads.

In the following subsections, we first introduce the definitions of the two metrics and then evaluate and compare the non-Gaussianity of DIV2K and microscopy nanobead datasets.

C.3.1 SKEWNESS (γ_1)

Skewness is a measure of the asymmetry of the probability distribution of a real-valued random variable. It quantifies how much the distribution deviates from a normal distribution in terms of asymmetry. The skewness value can be positive, zero, negative, or undefined. In a perfectly symmetrical distribution, skewness is zero. Positive skewness indicates a distribution with an extended tail on the right side, while negative skewness shows an extended tail on the left side.

The mathematical formula for skewness is given by:

$$\gamma_1 = E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] \quad (47)$$

where X is the random variable, μ is the mean of X , σ is the standard deviation of X , and E denotes the expected value.

The greater the absolute value of the skewness, the higher the degree of non-Gaussianity in the distribution.

C.3.2 KURTOSIS (β_2)

Kurtosis is a measure of the "tailedness" of the probability distribution of a real-valued random variable. It provides insights into the shape of the distribution's tails and peak. High kurtosis in a data set suggests a distribution with heavy tails and a sharper peak (leptokurtic), while low kurtosis indicates a distribution with lighter tails and a more flattened peak (platykurtic). Kurtosis is often compared to the normal distribution, which has a kurtosis of 3 (excess kurtosis of 0).

The formula for kurtosis is:

$$\beta_2 = E \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] - 3 \quad (48)$$

where the variables represent the same as in the skewness formula.

These statistical measures, skewness (γ_1) and kurtosis (β_2), are crucial for quantifying and analyzing the non-Gaussianity in image data. They provide valuable insights into the distribution characteristics of image pixel intensities, particularly in highlighting deviations from the normal distribution.

The higher the kurtosis, the greater the degree of non-Gaussianity in the distribution, indicating a distribution with heavier tails than a normal distribution.

C.3.3 NON-GAUSSIANITY OF DATASETS

Here we examine the non-Gaussianity of the distribution of DIV2K training dataset and the microscopy nanobead dataset used in this work. Table 2 summarizes their skewnesses and kurtoses. We observe that microscopy images tend to have larger absolute values of skewness and kurtosis, confirming their highly non-Gaussian distribution and the high condition number for microscopy image restoration problem. As a result, standard DMs with the assumption that the distribution of target images is close to normal does not hold for microscopy images and may lead to performance degradation and excessive sampling time.

Dataset	Skewness	Kurtosis
DIV2K	0.2193	-1.0373
Microscopy nanobeads	2.5852	13.3429

Table 2: Quantitative measure of non-Gaussianity of DIV2K and microscopy image datasets.

D WAVELET TRANSFORM

Wavelet transforms are derived from a single prototype function known as the ‘mother wavelet’. This function undergoes various scaling and translation processes to generate a family of wavelets. The general form of a wavelet function $\psi(x)$, derived from the mother wavelet, is expressed as:

$$\psi_{a,b}(x) = \frac{1}{\sqrt{|a|}} \psi\left(\frac{x-b}{a}\right), \quad (49)$$

where $a, b \in \mathbb{R}$, $a \neq 0$, and $x \in D$, with a and b representing scaling and translation parameters, respectively, and D being the domain of the wavelets.

D.1 MULTIREOLUTION ANALYSIS (MRA)

MRA is crucial in the Wavelet Transform Jawerth & Sweldens (1994); Ranta (2010). It involves constructing a basic functional basis in a subspace V_0 within $L^2(\mathbb{R})$ and then expanding this basis through scaling and translation to cover the entire space $L^2(\mathbb{R})$ for multiscale analysis.

For each $k \in \mathbb{Z}$ and $k \in \mathbb{N}$, we define the scale function space as $\mathbf{V}_k = \{f | f \text{ is restricted to } (2^{-k}l, 2^{-k}(l+1)) \text{ for all } l = 0, 1, \dots, 2^k - 1 \text{ and vanishes elsewhere}\}$. Each space \mathbf{V}_k encompasses a dimension of 2^k , with each subspace \mathbf{V}_i nested within \mathbf{V}_{i+1} :

$$\mathbf{V}_0 \subseteq \mathbf{V}_1 \subseteq \mathbf{V}_2 \subseteq \dots \subseteq \mathbf{V}_n \subseteq \dots \quad (50)$$

Using a base function $\varphi(x)$ in \mathbf{V}_0 , we can span \mathbf{V}_k with 2^k functions derived from $\varphi(x)$ through scaling and translation:

$$\varphi_l^k(x) = 2^{k/2} \varphi(2^k x - l), \quad l = 0, 1, \dots, 2^k - 1. \quad (51)$$

These functions, $\varphi_l(x)$, are known as scaling functions and they project any function into the approximation space \mathbf{V}_0 . The orthogonal complement of \mathbf{V}_k in \mathbf{V}_{k+1} is the wavelet subspace \mathbf{W}_k , satisfying:

$$\mathbf{V}_k \oplus \mathbf{W}_k = \mathbf{V}_{k+1}, \quad \mathbf{V}_k \perp \mathbf{W}_k. \quad (52)$$

Thus, we can decompose \mathbf{V}_k as:

$$\mathbf{V}_k = \mathbf{V}_0 \oplus \mathbf{W}_0 \oplus \mathbf{W}_1 \oplus \dots \oplus \mathbf{W}_{k-1}. \quad (53)$$

To form the orthogonal basis for \mathbf{W}_k , we construct it within $L^2(\mathbb{R})$. This basis is derived from the wavelet function $\psi(x)$, orthogonal to the scaling function $\varphi(x)$. The wavelet function is defined as:

$$\psi_l^k(x) = 2^{k/2}\psi(2^k x - l), \quad l = 0, 1, \dots, 2^k - 1. \quad (54)$$

A key property of these wavelet functions is their orthogonality to polynomials of lower order, exemplified by the vanishing moments criterion for first-order polynomials:

$$\int_{-\infty}^{\infty} x\psi_j(x)dx = 0. \quad (55)$$

Orthogonality and vanishing moments are central to wavelets, enabling efficient data representation and feature extraction by capturing unique data characteristics without redundancy. This efficiency is particularly useful in areas like signal processing and image compression.

D.2 WAVELET DECOMPOSITION AND RECONSTRUCTION

The Discrete Wavelet Transform (DWT) provides a multi-resolution analysis of signals, useful in various applications Gupta et al. (2021; 2022); Xiao et al. (2023). For a discrete signal $f[n]$, DWT decomposes it into approximation coefficients cA and detail coefficients cD using the scaling function $\varphi(x)$ and the wavelet function $\psi(x)$, respectively:

$$cA[k] = \sum_n h[n - 2k]f[n], \quad (56)$$

$$cD[k] = \sum_n g[n - 2k]f[n], \quad (57)$$

where $h[n]$ and $g[n]$ are the low-pass and high-pass filters, respectively. This decomposition process can be recursively applied for deeper multi-level analysis.

Reconstruction of the signal $f'[n]$ from its wavelet coefficients uses inverse transformations, employing synthesis filters $h_0[n]$ and $g_0[n]$:

$$f'[n] = \sum_k cA[k] \cdot h_0[n - 2k] + cD[k] \cdot g_0[n - 2k], \quad (58)$$

where the synthesis filters are typically the time-reversed counterparts of the decomposition filters. In multi-level decompositions, reconstruction is a stepwise process, beginning with the coarsest approximation and progressively incorporating higher-level details until the original signal is reconstructed.

D.3 HAAR WAVELET TRANSFORM

The Haar wavelet Stanković & Falkowski (2003), known for its simplicity and orthogonality, is a fundamental tool in digital signal processing. Its straightforward nature makes it an ideal choice for a variety of applications, which is why we have incorporated it into our project. The discrete wavelet transform (DWT) using Haar wavelets allows for the efficient decomposition of an image into a coarse approximation of its main features and detailed components representing high-frequency aspects. This process, enhanced by multi-resolution analysis (MRA), facilitates the examination of the image at various scales, thereby uncovering more intricate details.

The mathematical representation of the Haar wavelet and its scaling function is as follows:

$$\psi(t) = \begin{cases} 1 & \text{if } 0 \leq t < 0.5, \\ -1 & \text{if } 0.5 \leq t < 1, \\ 0 & \text{otherwise,} \end{cases} \quad (59)$$

$$\phi(t) = \begin{cases} 1 & \text{if } 0 \leq t < 1, \\ 0 & \text{otherwise,} \end{cases} \quad (60)$$

where $\psi(t)$ denotes the Haar wavelet function and $\phi(t)$ is the corresponding scaling function.

The application of Haar wavelet transform extends to two-dimensional spaces, particularly in image processing. This extension utilizes the same decomposition approach as for one-dimensional signals but applies it to both rows and columns of the image. The filter coefficients for Haar wavelets are calculated through these inner product evaluations:

- Low-pass filter coefficients (h):

$$\begin{aligned} h_0 &= \int_0^1 \phi(t) \cdot \phi(2t) dt, \\ h_1 &= \int_0^1 \phi(t) \cdot \phi(2t - 1) dt. \end{aligned} \quad (61)$$

- High-pass filter coefficients (g):

$$\begin{aligned} g_0 &= \int_0^1 \psi(t) \cdot \psi(2t) dt, \\ g_1 &= - \int_0^1 \psi(t) \cdot \psi(2t - 1) dt. \end{aligned} \quad (62)$$

To normalize these coefficients (ensuring their L2 norm is 1), we find $h_0 = \frac{1}{\sqrt{2}}$ and similar values for the other coefficients. Thus, the Haar filter coefficients are:

$$h = \left[\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right], \quad g = \left[\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right]. \quad (63)$$

For two-dimensional DWT in image processing, these filter coefficients are matrix-operated on the image. The horizontal operation uses the outer product of the filter vector with a column vector, and the vertical operation uses the outer product with a row vector. The resulting filter matrices are:

$$\begin{aligned} H &= h^T \otimes h = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}, \\ G &= g^T \otimes g = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}. \end{aligned} \quad (64)$$

The matrix H corresponds to low-pass filtering (approximation), while G captures the high-frequency details. In image transformation using Haar wavelets, these matrices help derive various coefficients representing different aspects of the image, such as approximation(LL), horizontal(LH), vertical(HL), and diagonal detail(HH) components.

E DUALITY PROOF

E.1 GENERATIVE MODELING IN SPATIAL DOMAIN

For the Score-based Generative Model (SGM) Song et al. (2020b); Ho et al. (2020); Song & Ermon (2019), the forward/noising process is mathematically formulated as the Ornstein-Uhlenbeck (OU) process. The general time-rescaled OU process is expressed as:

$$d\mathbf{X}_t = -g(t)^2 \mathbf{X}_t dt + \sqrt{2}g(t)d\mathbf{B}_t. \quad (65)$$

Here, $(\mathbf{X}_t)_{t \in [0, T]}$ represents the noising process starting with \mathbf{X}_0 sampled from the data distribution. $(\mathbf{B}_t)_{t \in [0, T]}$ denotes a standard d -dimensional Brownian motion. The reverse process, $\mathbf{X}_t^{\leftarrow}$, is defined such that $(\mathbf{X}_t^{\leftarrow})_{t \in [0, T]} = (\mathbf{X}_{T-t})_{t \in [0, T]}$. Assuming $g(t) = 1$ in standard diffusion models, the reverse process is:

$$d\mathbf{X}_t^{\leftarrow} = (\mathbf{X}_t^{\leftarrow} + 2\nabla \log p_{T-t}(\mathbf{X}_t^{\leftarrow})) dt + \sqrt{2}d\mathbf{B}_t. \quad (66)$$

Here, p_t is the marginal density of \mathbf{X}_t , and $\nabla \log p_t$ is the score. To revert \mathbf{X}_0 from \mathbf{X}_T via the time-reversed SDE, accurate estimation of the score $\nabla \log p_t$ at each time t is essential, alongside minimal error introduction during SDE discretization.

The reverse process approximation, as specified in Eq. 66, involves time discretization and score approximation $\nabla \log p_t$ by s_t , forming a Markov chain approximation of the time-reversed SDE. The chain starts with $\tilde{\mathbf{x}}_T \sim N(0, I_d)$, evolving over uniform time intervals Δt , from $t_N = T$ to $t_0 = 0$. The discretized process is detailed as:

$$\tilde{\mathbf{x}}_{t-1} = \tilde{\mathbf{x}}_t + \Delta t (\tilde{\mathbf{x}}_t + 2s_t(\tilde{\mathbf{x}}_t)) + \sqrt{2\Delta t}\mathbf{z}_t, \quad (67)$$

where \mathbf{z}_t are instances of Brownian motion \mathbf{B}_t .

E.2 GENERATIVE MODELING IN WAVELET DOMAIN

Consider \mathbf{X} as the image vector in the spatial domain. The discrete wavelet transform (DWT) Shensa et al. (1992); Heil & Walnut (1989) of \mathbf{X} can be expressed as:

$$\widehat{\mathbf{X}} = \mathbf{A}\mathbf{X}, \quad \mathbf{X} \in \mathbb{R}^d.$$

Here, \mathbf{A} represents the discrete wavelet transform matrix, which is orthogonal, satisfying $\mathbf{A}\mathbf{A}^\top = \mathbf{I}$. Various forms of \mathbf{A} are utilized in practice, such as Haar wavelets.

In the context of score-based generative modeling, we consider the forward (or noising) process, which can be formulated by the Ornstein–Uhlenbeck (OU) process. This is mathematically described as:

$$d\mathbf{X}_t = -\gamma(t)^2\mathbf{X}_t dt + \sqrt{2\gamma(t)}d\mathbf{B}_t, \quad (68)$$

where \mathbf{B}_t is a standard d -dimensional Brownian motion. Upon applying DWT to \mathbf{X}_t , the transformed $\widehat{\mathbf{X}}_t$ also follows the OU process:

$$d\widehat{\mathbf{X}}_t = -\gamma(t)^2\widehat{\mathbf{X}}_t dt + \sqrt{2\gamma(t)}\mathbf{A}d\mathbf{B}_t, \quad \widehat{\mathbf{X}}_0 = \mathbf{A}\mathbf{X}_0. \quad (69)$$

Defining $\widehat{\mathbf{B}}_t = \mathbf{A}\mathbf{B}_t$, which also behaves as a standard Brownian motion. If \mathbf{X}_0 is sampled from a distribution p , then $\widehat{\mathbf{X}}_0$ originates from the distribution

$$\hat{q} = \mathcal{T}_{\mathbf{A}}\#p, \quad (70)$$

where $\mathcal{T}_{\mathbf{A}}$ denotes the linear transformation operation by \mathbf{A} and $\#$ represents the pushforward operation. Consequently, we have

$$\hat{q}(\mathbf{x}) = p(\mathbf{A}^\top\mathbf{x}). \quad (71)$$

Let p_t be the density distribution of \mathbf{X}_t and \hat{q}_t that of $\widehat{\mathbf{X}}_t$. Then,

$$\hat{q}_t = \mathcal{T}_{\mathbf{A}}\#p_t, \quad \hat{q}_t(\mathbf{x}) = p_t(\mathbf{A}^\top\mathbf{x}). \quad (72)$$

Define the score functions for both processes as:

$$\mathbf{s}_t = \nabla \log p_t, \quad \mathbf{r}_t = \nabla \log \hat{q}_t. \quad (73)$$

These functions are related by:

$$\mathbf{r}_t(\mathbf{x}) = \frac{\nabla \hat{q}_t(\mathbf{x})}{\hat{q}_t(\mathbf{x})} = \frac{\mathbf{A}\nabla p_t(\mathbf{A}^\top\mathbf{x})}{p_t(\mathbf{A}^\top\mathbf{x})} = \mathbf{A}\mathbf{s}_t(\mathbf{A}^\top\mathbf{x}). \quad (74)$$

For the reverse processes denoted as \mathbf{X}_t^\leftarrow and $\widehat{\mathbf{X}}_t^\leftarrow$, assuming $\gamma(t) = 1$, they are given by:

$$\begin{aligned} d\mathbf{X}_t^\leftarrow &= (\mathbf{X}_t^\leftarrow + 2\mathbf{s}_{T-t}(\mathbf{X}_t^\leftarrow)) dt + \sqrt{2}d\mathbf{B}_t, \\ d\widehat{\mathbf{X}}_t^\leftarrow &= (\widehat{\mathbf{X}}_t^\leftarrow + 2\mathbf{r}_{T-t}(\widehat{\mathbf{X}}_t^\leftarrow)) dt + \sqrt{2}d\widehat{\mathbf{B}}_t. \end{aligned} \quad (75)$$

Here, $\widehat{\mathbf{B}}_t = \mathbf{A}\mathbf{B}_t$. Exploring the second SDE:

$$\begin{aligned} d\widehat{\mathbf{X}}_t^\leftarrow &= (\widehat{\mathbf{X}}_t^\leftarrow + 2\mathbf{r}_{T-t}(\widehat{\mathbf{X}}_t^\leftarrow)) dt + \sqrt{2}d\widehat{\mathbf{B}}_t \\ &= (\widehat{\mathbf{X}}_t^\leftarrow + 2\mathbf{A}\mathbf{s}_{T-t}(\mathbf{A}^\top\widehat{\mathbf{X}}_t^\leftarrow)) dt + \sqrt{2}d\widehat{\mathbf{B}}_t \\ \mathbf{A}^\top d\widehat{\mathbf{X}}_t^\leftarrow &= (\mathbf{A}^\top\widehat{\mathbf{X}}_t^\leftarrow + 2\mathbf{s}_{T-t}(\mathbf{A}^\top\widehat{\mathbf{X}}_t^\leftarrow)) dt + \sqrt{2}\mathbf{A}^\top d\widehat{\mathbf{B}}_t. \end{aligned} \quad (76)$$

Substituting $\mathbf{A}^\top\widehat{\mathbf{X}}_t^\leftarrow$ with \mathbf{X}_t^\leftarrow brings us back to the first equation. The training processes for $s_\theta, \mathbf{r}_\theta$ with $\mathbf{X}_t^{(i)}, \widehat{\mathbf{X}}_t^{(i)}$ also follow the standard denoising score matching loss function:

$$\begin{aligned} &\mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{X}_0} \mathbb{E}_{\mathbf{X}_t | \mathbf{X}_0} \left[\|\mathbf{s}_\theta(\mathbf{X}_t, t) - \nabla_{\mathbf{X}_t} \log p_{0t}(\mathbf{X}_t | \mathbf{X}_0)\|^2 \right] \right\} \\ &\mathbb{E}_t \left\{ \hat{\lambda}(t) \mathbb{E}_{\widehat{\mathbf{X}}_0} \mathbb{E}_{\widehat{\mathbf{X}}_t | \widehat{\mathbf{X}}_0} \left[\|\mathbf{r}_\theta(\widehat{\mathbf{X}}_t, t) - \nabla_{\widehat{\mathbf{X}}_t} \log \hat{q}_{0t}(\widehat{\mathbf{X}}_t | \widehat{\mathbf{X}}_0)\|^2 \right] \right\}. \end{aligned} \quad (77)$$

The forward and reverse probability distribution functions p_{0t} and \hat{q}_{0t} are defined as per the standard SGM model:

$$\begin{aligned} p_{0t}(\mathbf{X}_t|\mathbf{X}_0) &= N(\mathbf{X}_t; \sqrt{\bar{\alpha}_t}\mathbf{X}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \\ \hat{q}_{0t}(\hat{\mathbf{X}}_t|\hat{\mathbf{X}}_0) &= N(\hat{\mathbf{X}}_t; \sqrt{\bar{\alpha}_t}\hat{\mathbf{X}}_0, (1 - \bar{\alpha}_t)\mathbf{I}). \end{aligned} \quad (78)$$

F DATASETS AND IMPLEMENTATION DETAILS

For super-resolution experiments on DIV2K, the LR images were bicubically downsampled to 64×64 center-cropped patches and bilinearly upsampled to 256×256 as the input images, and the HR images were center-cropped as 256×256 . For shadow removal on the ISTD dataset, we used a similar method, pairing original shadow images with their shadow-free counterparts, both consistently sized to maintain uniformity in processing.

For super-resolution experiments involving HeLa cell nuclei and nano-beads, we utilized a Leica TCS SP8 stimulated emission depletion (STED) confocal microscope, equipped with a Leica HC PL APO 100 \times /1.40-NA Oil STED White objective. The samples were prepared on high-performance coverslips ($170\mu\text{m} \pm 10\mu\text{m}$, Carl Zeiss Microscopy), facilitating precise imaging. In the staining process, HeLa cell nuclei were treated with Rabbit anti-Histone H3 and Atto-647N Goat anti-rabbit IgG antibodies. The nano-beads, used for the STED experiments, were processed with methanol and ProLong Diamond antifade reagents. Both samples were excited with a 633-nm wavelength laser. Emission detection was carried out using a HyD SMD photodetector (Leica Microsystems) through a 645–752-nm bandpass filter.

For super-resolution experiments on microscopy images of nanobeads, the low-resolution (LR) images were acquired with specific settings (16 times line average and 30 times frame average for nano-beads; 8 times line average and 6 times frame average for cell nuclei), ensuring the acquisition of high-quality images necessary for our computer vision analysis. The scanning step size, i.e., the effective pixel size is 30 nm. Here, all LR and HR images were center cropped as 256×256 patches to match the setting in other experiments.

AdamW Loshchilov & Hutter (2017) optimizers were used in all experiments with an initial learning rate of $1e-4$. Models with exponential moving averaged parameters with a rate of 0.999 was saved and evaluated. The BBDM models at the coarsest wavelet scale were trained for 100000 steps with a batch size of 6.

In our GAN training, we used a batch size of 3, with images cropped at 256×256 resolution. We employed AdamW optimizers, setting the learning rate at $1e-4$ for the generator and $1e-5$ for the discriminator. The training loss function is a weighted sum of diverse loss terms: L1 loss with a weight of 20.0, adversarial loss at 0.1, and structural similarity index measure (SSIM) Wang et al. (2004) loss weighted at 0.5. The training was conducted for 200 epochs.

All models were implemented on NVIDIA V100 graphic cards using PyTorch. For speed test, a batch size of 16 was used for BBDM, IR-SDE and Refusion models, and a batch size of 64 was used for MSCGM so that they consume equivalent computational resources.

The complete training process of WSCGM is described below:

Algorithm 3 WSCGM Training

- 1: Sample $(\mathbf{x}_0, \mathbf{y}) \sim q(\mathbf{x}_0, \mathbf{y})$
 - 2: Wavelet transform \mathbf{x}_0, \mathbf{y} by S times to get $\{\mathbf{x}_L^S, \mathbf{x}_H^S, \dots, \mathbf{x}_H^1, \mathbf{y}_L^S, \mathbf{y}_H^S, \dots, \mathbf{y}_H^1\}$
 - 3: Sample $t \sim \text{Uniform}([1, \dots, T])$, $\epsilon_t \sim N(\mathbf{0}, \mathbf{I})$
 - 4: Take gradient step on $\nabla_{\theta} \|m_t(\mathbf{y}_L^S - \mathbf{x}_L^S) + \sqrt{\delta}\epsilon_t - \epsilon_{\theta}(\mathbf{x}_{t,L}^S, \mathbf{y}, t)\|^2$
 - 5: Take gradient step on L_G and L_D
-

G ADDITIONAL RESULTS

G.1 RESULTS ON IMAGE ENHANCEMENT

Here we provide the metrics comparison results of image enhancement.

Methods	PSNR (dB) \uparrow	SSIM \uparrow	FID \downarrow
BBDM	28.66	0.860	98.49
HWMNet	29.40	0.902	82.30
RetinexFormer	29.55	0.897	109.75
MSCGM	29.27	0.863	115.96

Table 3: **Comparison of HWMNetFan et al. (2022), RetinexFormerCai et al. (2023) and our method (MSCGM) on image enhancement experiment.** Metrics calculated on 256×256 center-cropped images of LOL dataset. Entries in bold indicate the best performance achieved among the compared method.

G.2 RESULTS ON SHADOW REMOVAL TASK

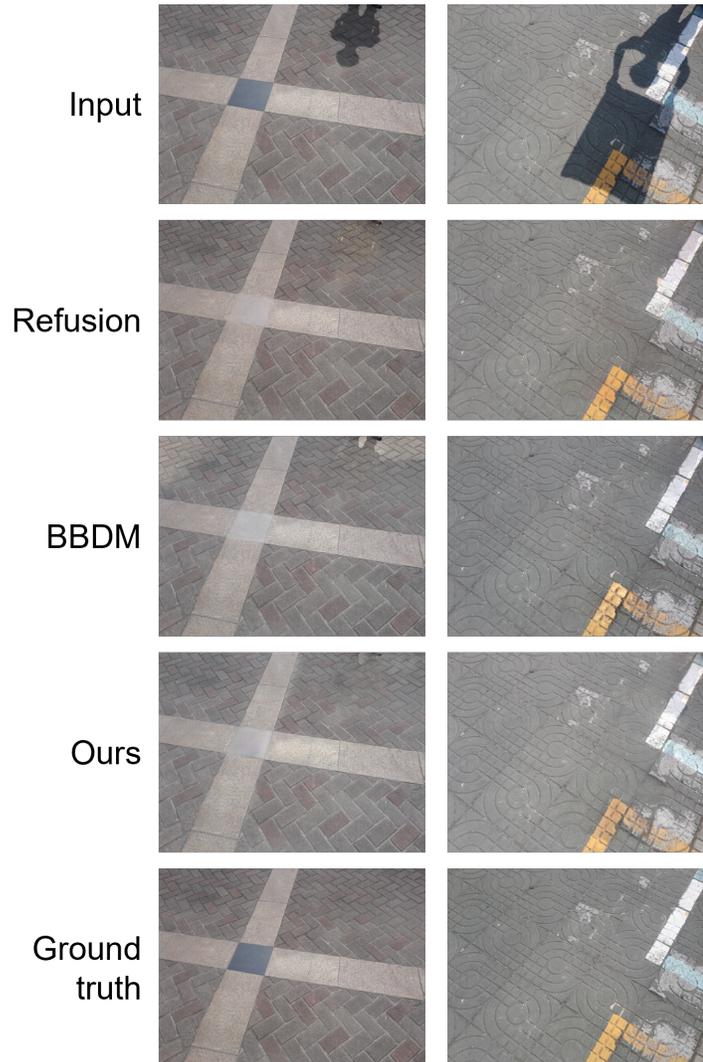


Figure 9: Comparison of our method and baselines on the natural image shadow removal task on ISTD test set. Sample steps were set as 1000 for all methods.

Figure 9 visualizes the shadow removal performance of our approach and competitive methods. Overall, our method generates output images with better consistency and milder artifacts. Refer to Table 4 for quantitative assessments on these results.

Methods	PSNR (dB) \uparrow	SSIM \uparrow
DC-ShadowNet Jin et al. (2023)	26.38	0.922
ST-CGAN Wang et al. (2018b)	27.44	0.929
DSC Hu et al. (2019)	30.64	0.843
DHANCun et al. (2020)	27.88	0.921
BMNet Zhu et al. (2022)	30.28	0.927
ShadowFormer Guo et al. (2023)	30.47	0.935
BBDM Li et al. (2023)	30.54	0.910
MSCGM(Ours)	31.08	0.915

Table 4: **Quantitative evaluation metrics of our method (MSCGM) and competitive methods on ISTD dataset.** Metrics calculated on full-resolution images.

G.3 ADDITIONAL RESULTS ON MICROSCOPY IMAGE SUPER-RESOLUTION

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

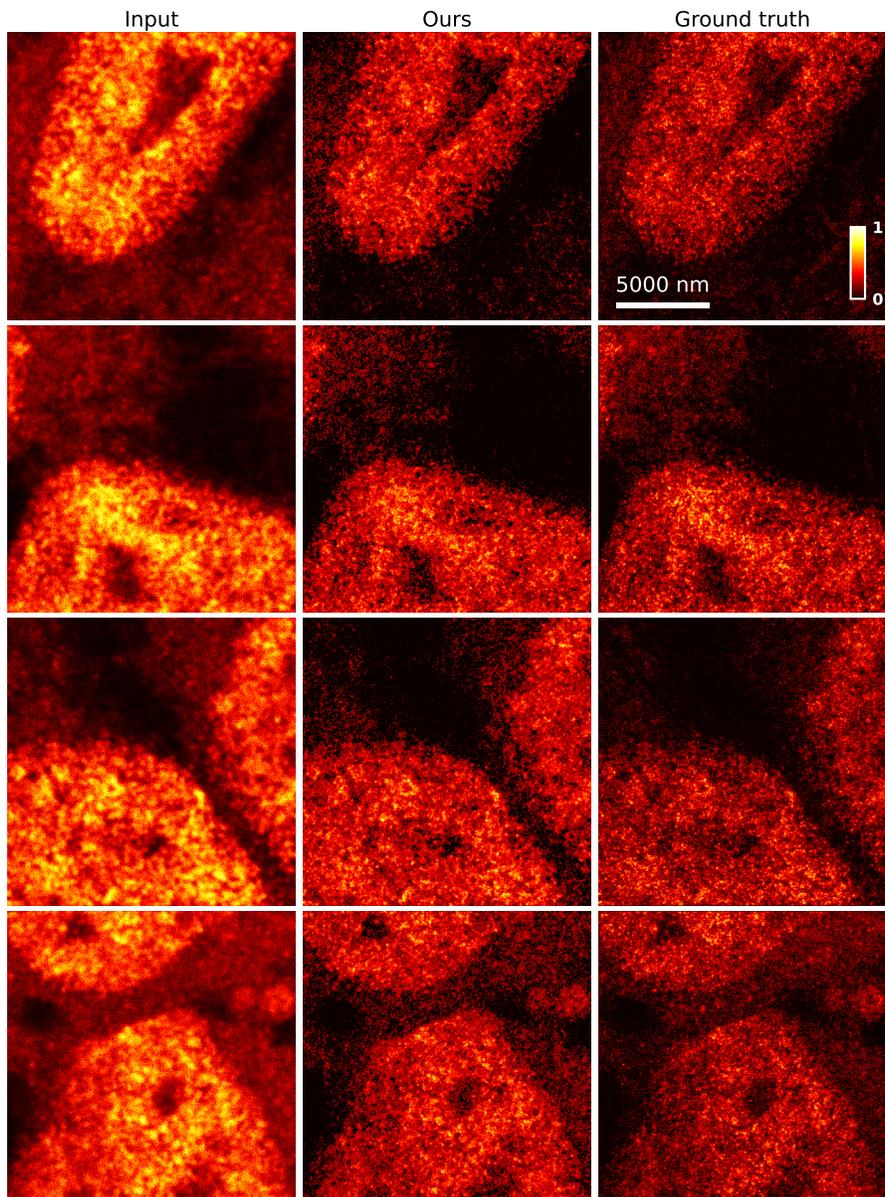


Figure 10: Additional microscopy image super-resolution results of HeLa cells.

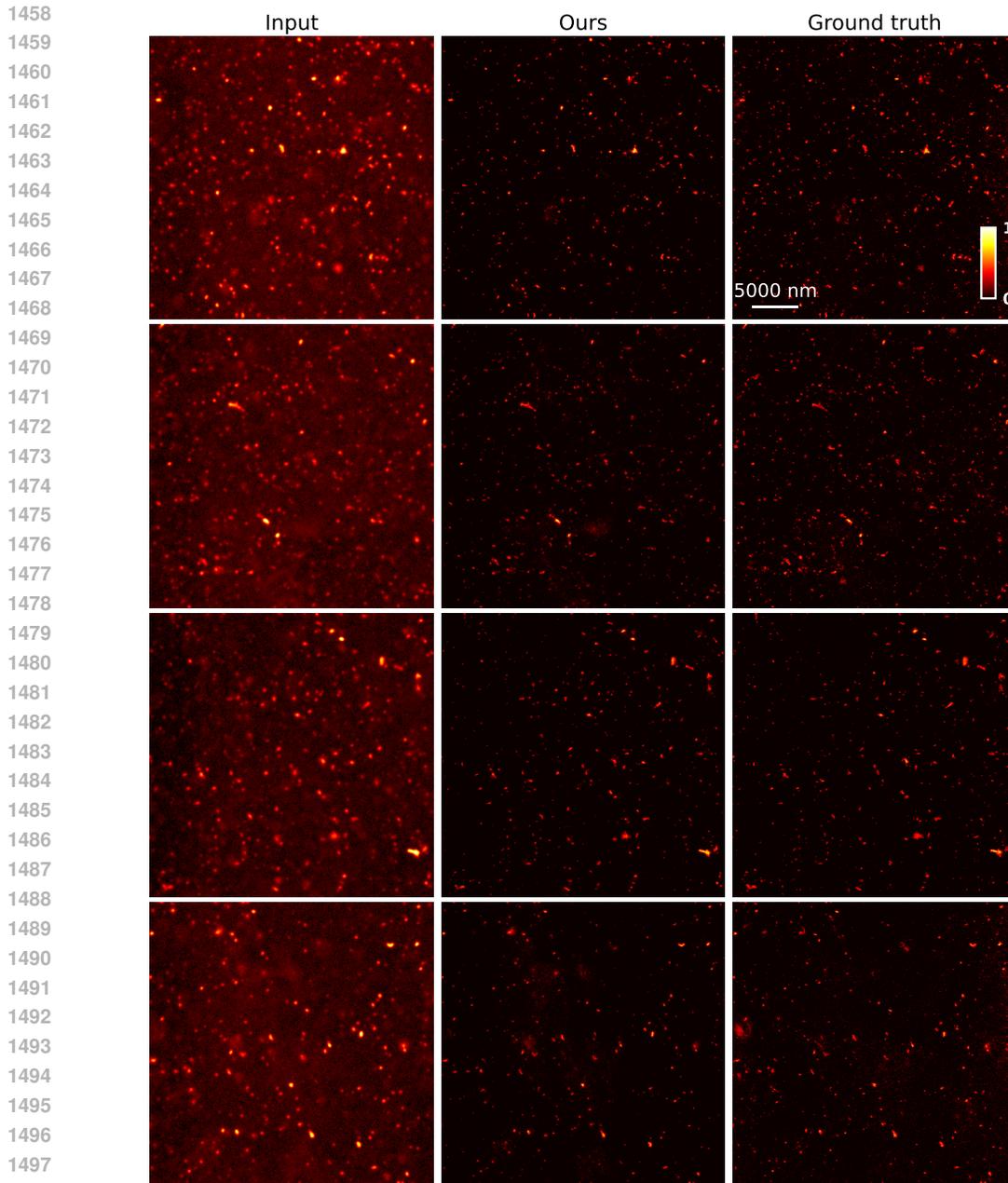


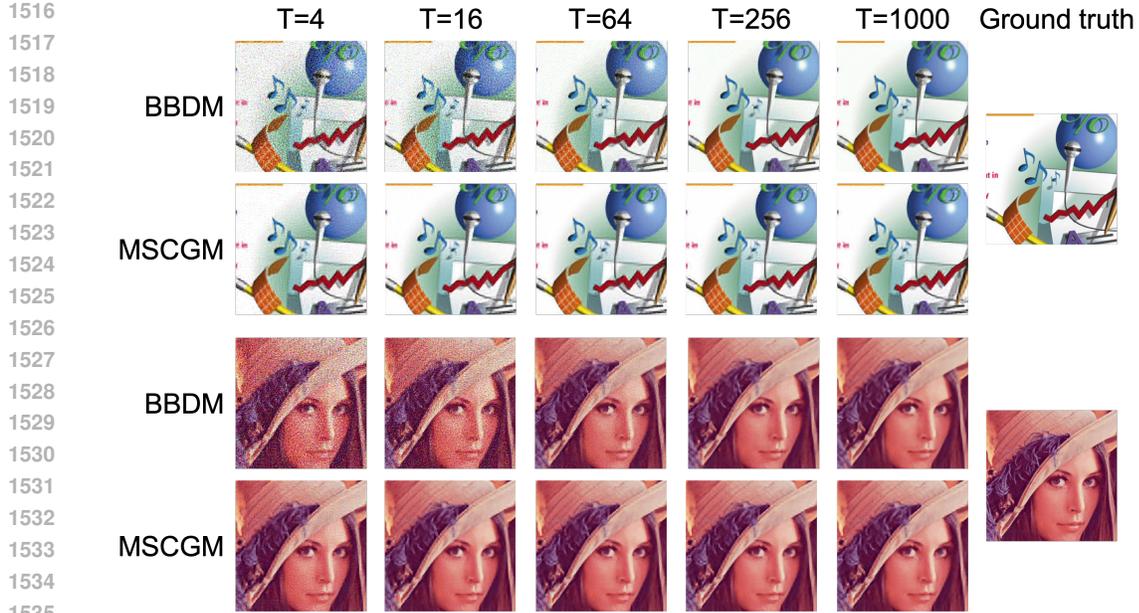
Figure 11: Additional microscopy image super-resolution results of nanobeads.

Figure 4 and 10 presents additional results on microscopy image super-resolution of HeLa cells and Fig. 11 showcases additional results on microscopy image super-resolution of nanobeads. In correspondence to the results shown in the main text, SR images generated by our method match with the ground truths very well.

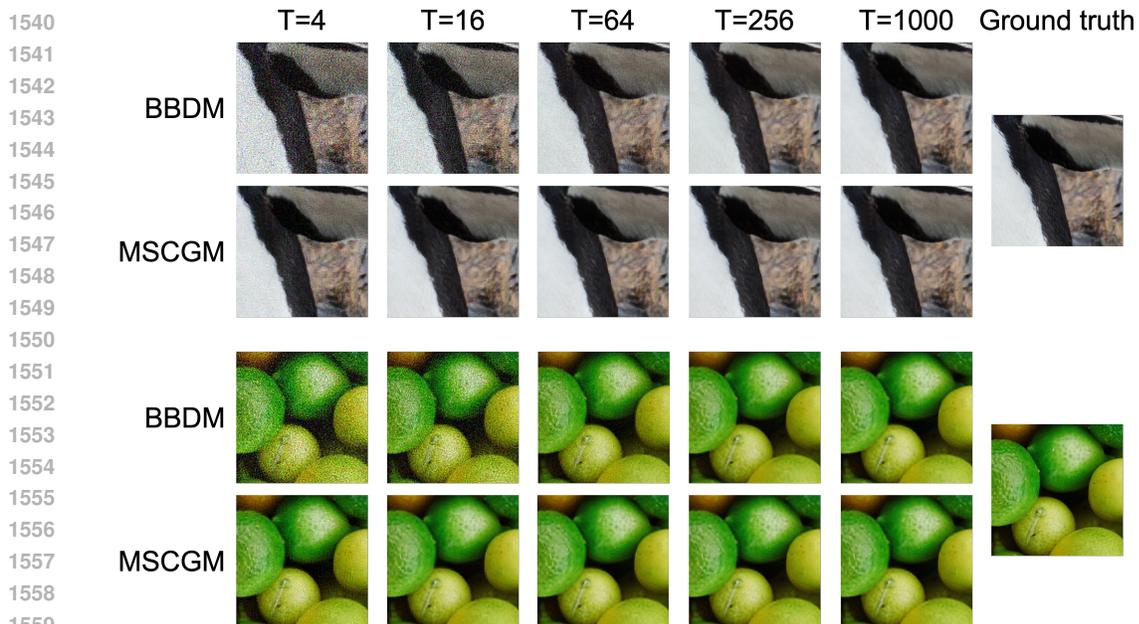
G.4 FAST SAMPLING RESULT

Here we present the sampling results of our method(MSCGM) and BBDM of each sampling steps (4, 16, 64, 256) on the Set14 and DIV2k validation dataset. Figure 12 show the sampling results of our method and BBDM with various sampling steps on Set14 dataset, and Fig. 13 illustrates the improvement of sampling quality with respect to the number of sampling steps from 4 to 1000 on

1512 DIV2K validation dataset. Remarkably, our method recovers high-quality image considerably faster in
 1513 fewer sampling steps, confirming its superiority in sampling speed compared to competitive diffusion
 1514 models.
 1515



1536 Figure 12: Sampling results of our method (MSCGM) and BBDM with various sampling steps from
 1537 4 to 1000. Images sampled from 64×64 LR images in Set14.
 1538



1560 Figure 13: Sampling results of our method (MSCGM) and BBDM with various sampling steps from
 1561 4 to 1000. Images sampled from 64×64 LR images in DIV2K validation dataset.
 1562

1563
 1564
 1565