FORWARD-ONLY DIFFUSION PROBABILISTIC MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

This work presents a forward-only diffusion (FoD) approach for generative modelling. In contrast to traditional diffusion models that rely on a coupled forward-backward diffusion scheme, FoD directly learns data generation through a single forward diffusion process, yielding a simple yet efficient generative framework. The core of FoD is a state-dependent stochastic differential equation that involves a mean-reverting term in both the drift and diffusion functions. This mean-reversion property guarantees the convergence to clean data, naturally simulating a stochastic interpolation between source and target distributions. More importantly, FoD is analytically tractable and is trained using a simple stochastic flow matching objective, enabling a few-step non-Markov chain sampling during inference. The proposed FoD model—despite its simplicity—achieves state-of-the-art performance on various image restoration tasks. Its general applicability on image-conditioned generation is also demonstrated via qualitative results on image-to-image translation.

1 Introduction

The diffusion model has become a central theme in generative modelling (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021; Karras et al., 2022). A crucial feature of the diffusion model is the use of a forward process that gradually perturbs the data into noise, coupled with a backward process that learns to transform noise back to data (Sohl-Dickstein et al., 2015; Ho et al., 2020). Benefiting from this forward-backward framework, the diffusion models have achieved remarkable performance in producing high-quality results across a wide range of applications, including image synthesis (Dhariwal & Nichol, 2021; Saharia et al., 2022b; Rombach et al., 2022; Peebles & Xie, 2023; Podell et al., 2024), translation (Meng et al., 2022; Saharia et al., 2022a; Tumanyan et al., 2023; Kawar et al., 2023; Brooks et al., 2023), and restoration (Saharia et al., 2022c; Kawar et al., 2022c; Luo et al., 2023a; Wang et al., 2024; Lin et al., 2024).

Despite this success, the reliance on the coupled forward-backward construction substantially increases the algorithmic complexity at the same time as it leads to a challenging model training problem. In addition, the necessity of corrupting data to noise in diffusion models further imposes an undesirable constraint for image-conditioned generation (Kawar et al., 2022; Saharia et al., 2022a), where ideally the generative process should start with image conditions that are structurally more informative than noise (Luo et al., 2023a; Saharia et al., 2022c; Liu et al., 2023). This naturally leads to a fundamental question:

"Could a simpler, single diffusion process suffice for effective generative modelling?"

This paper answers this question affirmatively by introducing a probabilistic forward-only diffusion model (FoD). Our exploration starts from the mean-reverting stochastic differential equation (SDE) (Gillespie, 1996; Luo et al., 2023a), where the data is stochastically driven toward a specified state characterized by a fixed mean and variance. Notably, setting the mean to zero recovers the standard forward diffusion process (Song et al., 2021). Inspired by this, we propose a new form of the mean-reverting SDE, which adds mean-reversion to *both* the drift and diffusion functions as a state-dependent diffusion process. Here, we highlight the mean-reversion diffusion function as it guarantees the convergence to the noise-free mean state. By setting the mean to the target data, FoD naturally simulates the data transition between source and target distributions, without requiring a separate backward process.

We further demonstrate that the FoD process is analytically tractable and follows a multiplicative stochastic structure. Moreover, we show that the model can be learned by approximating the vector

field from each noisy state to the final clean data, a process we refer to as *stochastic flow matching*. The result is a simple, yet effective training process. Based on the tractable solution and the flow matching objective, FoD enables a few-step sampling strategy with both Markov and non-Markov chains, enabling more efficient data generation without compromising sample quality.

In addition, as a closely related work of our method, it is worth noting that flow matching (Lipman et al., 2022; Liu et al., 2022) can also eliminate the need for a separate data perturbation process by modelling a continuous flow from source distribution to target distribution. However, the noise injection, which has been shown crucial in generative models (Song & Ermon, 2019), is also eliminated due to the modelling of ordinary differential equations (ODEs). As a result, its performance drops significantly when handling image-conditioned generation tasks, such as image restoration (IR), which aims to recover high-quality images from their degraded low-quality counterparts (Albergo et al., 2023a; Martin et al., 2024; Ohayon et al., 2024). In contrast, FoD is a stochastic extension of flow matching that avoids the issue mentioned above by simulating SDEs with a state-dependent diffusion process, making it well-suited for image-conditioned generation.

Our experiments focus on image-conditioned generation, an active and fundamental direction of generative modelling with a wide range of real-world applications, including image restoration and image-to-image translation. Compared to existing diffusion and flow matching-based approaches, the proposed FoD achieves strong empirical performance across diverse tasks and datasets. Moreover, we provide a comprehensive analysis of efficient sampling using both Markov and non-Markov chains, and illustrate how the noise is injected and subsequently removed during the forward diffusion process, highlighting the importance of noise injection in image generation.

2 BACKGROUND

Given a source distribution p_{prior} and an unknown target data distribution p_{data} , our goal is to build a probability path $\{p(x_t)\}_{t=0}^T$ that transports between the source distribution $p(x_0) = p_{\text{prior}}$ and the target distribution $p(x_T) = p_{\text{data}}$. In this paper, the source can be either noise, for unconditional generation, or images, for image-conditioned generation, e.g., image restoration.

2.1 DIFFUSION MODELS

Given a target data point $x_T \sim p_{\text{data}}$, diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) define a Markov chain forward process to progressively perturb the data into noise $(x_T \to x_0)$ and then learn its reverse process to reconstruct the data $(x_0 \to x_T)$. This coupled forward-backward process can be defined by stochastic differential equations (SDEs) (Song et al., 2021), given by:

$$\underbrace{\mathrm{d}x_t = f(x_t, t)\,\mathrm{d}t + g(t)\,\mathrm{d}w_t}_{\text{Forward process}} \quad \text{and} \quad \underbrace{\mathrm{d}x_t = \left[f(x_t, t) - g(t)^2 \,\nabla \log p_t(x_t)\right]\mathrm{d}t + g(t)\,\mathrm{d}\bar{w}_t}_{\text{Backward process}}, \quad (1)$$

where f(x,t) is the *drift* function and g(t) is the *diffusion* function. Furthermore, w and \bar{w} are the standard Wiener process and its reverse process, respectively. We use $p_t(x_t)$ to denote the marginal probability density of x_t . The term $\nabla \log p_t(x_t)$, called the *score function*, is the sought-after objective in the backward (also called the reverse-time) SDE, which is often learned by a time-dependent neural network (Song et al., 2021) via score-matching. The training objective can also be converted to learn noise matching as in DDPMs (Ho et al., 2020). Moreover, the source distribution in diffusion models is often a Gaussian with a predefined mean and variance. Diffusion models typically require thousands of sampling steps to generate high-quality samples.

2.2 FLOW MATCHING GENERATIVE MODELS

Flow matching (Lipman et al., 2022; Albergo & Vanden-Eijnden, 2022; Liu et al., 2022) is a simple regression objective used for learning the velocity field $v(x_t,t)$ that transports a sample x_t from the source distribution to the target distribution along the probability path $p(x_t)$ (Lipman et al., 2024). More specifically, flow matching models aim to learn the ordinary differential equation (ODE): $\mathrm{d}x_t = v(x_t,t)\,\mathrm{d}t$, where $x_0 \sim p_{\mathrm{prior}}$ and the drift $v(x_t,t)$ transports samples from x_0 to $x_1 \sim p_{\mathrm{data}}$. Here, each latent variable x_t in the ODE path is drawn by linearly interpolating source and target data samples, i.e., $x_t = tx_1 + (1-t)x_0$. Then the training can be performed by uniformly sampling

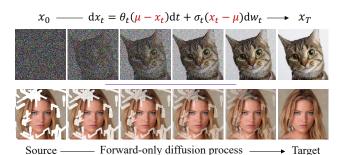


Figure 1: The proposed forward-only diffusion (FoD) probabilistic model. FoD introduces the mean reversion term (marked in red color) into *both* the drift and diffusion functions, enabling high-quality data samples with a single diffusion process. This method can be easily extended from unconditional generation (top row) to image-conditioned generation, such as the image restoration in the second row.

data pairs and timesteps and optimizing a flow matching objective, as:

$$L_{\text{FM}}(\phi) = \mathbb{E}_{x_0, x_1, t \sim \mathcal{U}(0, 1)} [\|(x_1 - x_0) - v_\phi(x_t, t)\|^2], \tag{2}$$

where $v_{\phi}(x_t,t)$ is a neural network approximating the true velocity field. Flow matching models eliminate the diffusion term from the generative process and thus lead to a simpler and more direct learning procedure based on ODE paths. However, in this paper, we observe that applying it to image-conditioned generation tasks, such as image restoration, leads to a significant performance drop due to the lack of stochastic noise injection (see Section 4.1 for more details). Moreover, it is worth noting that both diffusion models and flow matching models can be unified into the stochastic interpolants (Albergo et al., 2023a) framework.

3 FORWARD-ONLY DIFFUSION PROBABILISTIC MODELS

The key value of the forward-only diffusion (FoD) model lies in defining an analytically solvable, forward-only process that removes the need to approximate or learn a reverse SDE. This makes the generative process conceptually simple and stable, and easier to extend to image-conditioned generation, as illustrated in Figure 1.

3.1 Preliminaries: Mean-reverting SDE

Our exploration starts from a mean-reverting SDE (Gillespie, 1996; Luo et al., 2023a) where the data is stochastically driven towards a state characterized by a specified mean μ and variance λ^2 :

$$dx_t = \theta_t (\mu - x_t) dt + \sigma_t dw_t, \tag{3}$$

where $\{\theta_t\}_{t=0}^T$ and $\{\sigma_t\}_{t=0}^T$ are positive mean-reversion and diffusion schedules, respectively. By coupling the schedules as $\sigma_t^2/\theta_t=2~\lambda^2$ for all t, we obtain the following solution (Luo et al., 2023a)

$$x_t = \mu + (x_0 - \mu) e^{-\int_0^t \theta_z \, dz} + \int_0^t \sigma_z e^{-\int_s^t \theta_s \, ds} \, dw_z.$$
 (4)

As $t \to \infty$, the SDE converges to a stationary state $x_T \sim \mathcal{N}(x_t \mid \mu, \lambda^2)$. This property suggests constructing a process that transports samples from the source distribution p_{data} , by setting the mean μ to be a sample from p_{data} . However, as can be observed from Eq. (4), the resulting sample x_T is still noisy, with variance λ^2 , which works against our goal of generating high-quality clean data samples. In the following sections, we address this problem by introducing mean-reversion in both the drift and diffusion functions.

3.2 FORWARD-ONLY DIFFUSION PROCESS

We begin by designing an SDE with mean-reversion terms in both the drift and diffusion functions, as

$$dx_t = \theta_t (\mu - x_t) dt + \sigma_t (x_t - \mu) dw_t.$$
 (5)

This is a state-dependent linear SDE with multiplicative noise, where the diffusion volatility increases in the beginning steps and then decreases to zero when x_t converges to μ . We typically use $x_t - \mu$ in the diffusion function such that this SDE simulates a reverse Wiener process as in diffusion models (Song et al., 2021). For image generation, the noise $\mathrm{d}w_t$ is added independently at each pixel, meaning that this SDE is applied for image transitions pixel-by-pixel, under the Itô interpretation.

We refer to this SDE as the *forward-only diffusion* (FoD) process, and present its solution as follows:

Proposition 3.1. Given an initial state x_s at time s < t, the unique solution to the SDE (5) is

$$x_t = \left(x_s - \mu\right) e^{-\int_s^t \left(\theta_z + \frac{1}{2}\sigma_z^2\right) dz + \int_s^t \sigma_z dw_z} + \mu,\tag{6}$$

where the stochastic integral is interpreted in the Itô sense and can be reparameterised as $\bar{\sigma}_{s:t} \epsilon$, where $\bar{\sigma}_{s:t} = \sqrt{\int_s^t \sigma_z^2 \, \mathrm{d}z}$, and $\epsilon \sim \mathcal{N}(0, I)$ is a standard Gaussian noise.

The proof is provided in Appendix A.1. In addition, the solution in Eq. (6) shows that the stochastic flow field $\mu - x_t$ forms a Geometric Brownian motion (Ross, 2014) and yields the following corollary:

Corollary 3.2. Under the same assumptions as in Proposition 3.1, the stochastic flow field $\mu - x_t$ satisfies the multiplicative stochastic structure. More precisely, it is log-normally distributed by

$$\log(\mu - x_t) \sim \mathcal{N}\left(\log(\mu - x_s) - \int_s^t \left(\theta_z + \frac{1}{2}\sigma_z^2\right) dz, \int_s^t \sigma_z^2 dz I\right). \tag{7}$$

This follows directly from Proposition 3.1, by rearranging $\mu - x_t$ to the left of Eq. (6) and applying the logarithm to both sides (see Appendix A.1). The subtractive form of the logarithm reflects that the flow field decays multiplicatively from its initial value with a stochastic exponential scaling.

Notational Clarifications: Although the sign of $\mu - x_t$ in general can be either positive or negative, it remains consistent across all times t for a given sample; therefore, we choose to omit absolute values inside the logarithmic terms in Eq. (7) for notational convenience. In addition, we further let $\bar{m}_{s:t} = -\int_s^t (\theta_z + \frac{1}{2}\sigma_z^2) \, \mathrm{d}z$ and $\bar{m}_t = \bar{m}_{0:t}$ in the rest of the paper to simplify the notation.

3.3 STOCHASTIC FLOW MATCHING

Let us now explain how we can learn this FoD process, i.e., transforming data from a known source distribution p_{prior} to an unknown target distribution p_{data} . Following DDPMs (Ho et al., 2020), we define the FoD model as $p_{\phi}(x_{0:T})$, a joint distribution with learnable transitions starting at x_0 , as

$$p_{\phi}(x_{0:T}) = p_{\text{prior}}(x_0) \prod_{t=0}^{T-1} p_{\phi}(x_{t+1} \mid x_t), \qquad x_0 \sim p_{\text{prior}}.$$
 (8)

We propose to set the transition kernel $p_{\phi}(x_{t+1}|x_t)$ to be in the same log-Gaussian form as Eq. (6). The training can then be performed by minimizing the negative log-likelihood of $p_{\phi}(x_T)$, which is equivalent to optimizing the following objective:

$$\mathbb{E}_{p} \Big[\sum_{t=0}^{T-1} D_{KL}(p(x_{t+1} \mid x_{t}, x_{T}) \parallel p_{\phi}(x_{t+1} \mid x_{t})) \Big]. \tag{9}$$

The proof is provided in Appendix A.2. During training, we set x_T equal to μ such that the SDE (5) converges to data μ exactly. Then, the conditional distribution $p(x_{t+1}|x_t,x_T=\mu)$ is tractable as shown in Eq. (6). By letting the functions $f_{\mu}(x_t) = \mu - x_t$ and $f_{\phi}(x_t,t) = \hat{\mu}_{\phi} - x_t$ denote the ground truth and the model prediction of the stochastic flow field, respectively, we transform the distributions in Eq. (9) from SDE states to stochastic flow fields. Note that this transformation, i.e., from $p(\cdot|x_t,\mu)$ to $p(\cdot|f_{\mu}(x_t))$, holds because its Jacobian determinant equals one. Instead of Eq. (9), we can therefore minimize the KL divergence between two stochastic flow distributions:

$$\mathbb{E}_{p} \Big[\sum_{t=0}^{T-1} D_{KL}(p(f_{\mu}(x_{t+1}) \mid f_{\mu}(x_{t}) \parallel p(f_{\phi}(x_{t+1}, t) \mid f_{\phi}(x_{t}, t)) \Big]. \tag{10}$$

Combining this with Corollary 3.2, we obtain the final objective:

$$L_{\text{SFM}}(\phi) := \mathbb{E}_{\mu \sim p_{\text{data}}, x_t \sim p(x_t | x_0, \mu)} \Big[\| \log(\mu - x_t) - \log f_{\phi}(x_t, t) \|^2 \Big]$$

$$\approx \mathbb{E}_{\mu \sim p_{\text{data}}, x_t \sim p(x_t | x_0, \mu)} \Big[\| (\mu - x_t) - f_{\phi}(x_t, t) \|^2 \Big],$$
(11)

where the approximation follows from a first-order Taylor expansion close to the optimum. Please refer to Appendix A.3 for more details. This objective is referred to as *stochastic flow matching* and it is entirely linear, which leads to a simple and numerically stable training process.

```
216
              Algorithm 1 FoD Training
                                                                                             Algorithm 2 FoD Sampling
217
218
                                                                                             Require: p_{\text{prior}}, time interval \Delta t, model f_{\phi}
              Require: p_{\text{prior}}, p_{\text{data}}, \text{ model } f_{\phi}
                                                                                               1: x_0 \sim p_{\text{prior}}
               1: repeat
219
                                                                                               2: for t = 0, ..., T - 1 do
               2: x_0 \sim p_{\text{prior}}, \mu \sim p_{\text{data}}
220
               3: \epsilon \sim \mathcal{N}(0, I), t \sim \text{Uniform}(\{1, \dots, T\})
                                                                                               3: \epsilon \sim \mathcal{N}(0, I)
221
                   x_t = (x_0 - \mu) e^{\bar{m}_t + \sigma_t \epsilon} + \mu
                                                                                              4: \Delta x = \theta_t f_{\phi}(x_t, t) \cdot \Delta t - \sigma_t f_{\phi}(x_t, t) \cdot \sqrt{\Delta t} \epsilon
222
                    Take gradient descent step on
                                                                                               5: x_{t+1} = x_t + \Delta x
                                  \nabla_{\phi} \| (\mu - x_t) - f_{\phi}(x_t, t) \|^2
                                                                                               6: end for
224
               6: until converged
                                                                                               7: return x_T
225
226
              Algorithm 3 Markov Chain Sampling
                                                                                             Algorithm 4 Non-Markov Chain Sampling
227
              Require: p_{\text{prior}}, step size k, model f_{\phi}
                                                                                             Require: p_{\text{prior}}, step size k, model f_{\phi}
228
               1: x_0 \sim p_{\text{prior}}
                                                                                               1: x_0 \sim p_{\text{prior}}
229
               2: for t = 0, k, 2k, ..., T do
                                                                                               2: for t = 0, k, 2k, ..., T do
230
               3: \epsilon \sim \mathcal{N}(0, I)
                                                                                               3: \epsilon \sim \mathcal{N}(0, I)
231
               4: \hat{\mu} = x_t + f_{\phi}(x_t, t)
                                                                                               4: \hat{\mu} = x_t + f_{\phi}(x_t, t)
               5: x_{t+k} = (x_t - \hat{\mu}) e^{\bar{m}_{t:t+k} + \epsilon \cdot \bar{\sigma}_{t:t+k}} + \hat{\mu}
                                                                                              5: x_{t+k} = (x_0 - \hat{\mu}) e^{\bar{m}_{t+k} + \epsilon \cdot \bar{\sigma}_{t+k}} + \hat{\mu}
232
233
               6: end for
                                                                                               6: end for
234
               7: return x_T
                                                                                               7: return x_T
```

The standard training and sampling (via the Euler–Maruyama method) procedures are provided in Algorithm 1 and Algorithm 2, respectively. In addition, the target data estimate $\hat{\mu}$ is given by

$$\hat{\mu} = x_t + f_\phi(x_t, t),\tag{12}$$

which can be applied to the forward transition (6) for fast data sampling.

Fast Sampling with Markov and non-Markov Chains While the generation can be performed by iteratively solving the SDE (5) with numerical schemes such as the Euler-Maruyama method, it often requires hundreds of sampling steps. Fortunately, the tractable solution of FoD naturally enables fast sampling during inference, by choosing times discretely with a larger step size k, as $t = [0, k, 2k, 3k, \ldots, T]$, where T is the total number of timesteps. Since our prediction at each step is the sought-after target data $\hat{\mu}$, the next state x_{t+k} can be sampled following Eq. (6) with either Markov or non-Markov chains. This is done by setting the transition to $x_t \to x_{t+k}$ or $x_0 \to x_{t+k}$, as illustrated in Algorithm 3 and Algorithm 4, respectively. A further discussion is provided in Section 5.

3.4 Connection to Prior Work

235236

237

238

239240

241242

243

244

245

246

247

248

249250251

252

253

254

255

256257

258

259

260

261262

263

264

265

266267

268

269

In this section, we establish the theoretical connections between FoD and two closely related prior works: stochastic interpolants (SI) (Albergo et al., 2023a) and flow matching (FM) (Lipman et al., 2022). SI provides a unified stochastic framework to bridge two arbitrary distributions, while FM formulates a deterministic transport map between two distributions via an ODE.

Stochastic Interpolants Let us recall the solution in Eq. (6) of the FoD process. By setting the initial state to x_0 and rearranging the equation, we obtain a stochastic process in the interpolant form:

$$x_t = I(t, x_0, \mu) = x_0 \alpha_t + \mu (1 - \alpha_t), \quad \alpha_t = e^{-\int_0^t (\theta_z + \frac{1}{2}\sigma_z^2) dz + \int_0^t \sigma_z dw_z}.$$
 (13)

Here, $I(t,x_0,\mu)$ satisfies the boundary conditions of a stochastic interpolant, with randomness introduced via $\mathrm{d}w$. FoD can thus be viewed as a powerful instantiation of SI, distinguished by two key properties: multiplicative log-normal interpolation and a state-dependent stochastic path from x_0 to μ . This formulation allows noise to be gradually added and subsequently removed within a single forward process. This perspective helps unify FoD with a broader class of generative frameworks.

Flow Matching We consider a deterministic version of the FoD process in Eq. (5), i.e., omitting the diffusion term or setting $\sigma_t = 0$ for all times. This gives a mean-reverting ODE that bridges two distributions without noise injection, as $\mathrm{d}x_t = \theta_t \, (\mu - x_t) \, \mathrm{d}t$ with solution $x_t = (x_s - \mu) \, \mathrm{e}^{-\int_s^t \theta_z \, \mathrm{d}z} + \frac{1}{2} \, \mathrm{d}t$

Table 1: Quantitative comparison of our method with other diffusion and flow matching approaches on four different image restoration datasets, evaluated using both distortion and perceptual metrics.

| Method | Disto | rtion | Perceptual | | |
|----------------|-------|-------|------------|-------|--|
| Wiemou | PSNR↑ | SSIM↑ | LPIPS↓ | FID↓ | |
| U-Net baseline | 29.12 | 0.882 | 0.153 | 57.55 | |
| IR-SDE | 31.65 | 0.904 | 0.047 | 18.64 | |
| GOUB | 31.96 | 0.903 | 0.046 | 18.14 | |
| ReFlow | 28.36 | 0.871 | 0.152 | 64.81 | |
| PMRF | 29.01 | 0.857 | 0.173 | 69.25 | |
| FoD (Ours) | 32.56 | 0.925 | 0.038 | 14.10 | |

| Method | Disto | ortion | Percep | Perceptual | | |
|----------------|-------|--------|--------|------------|--|--|
| Wediod | PSNR↑ | SSIM↑ | LPIPS↓ | FID↓ | | |
| U-Net baseline | 20.51 | 0.808 | 0.162 | 75.84 | | |
| IR-SDE | 20.45 | 0.787 | 0.129 | 47.28 | | |
| GOUB | 19.29 | 0.775 | 0.148 | 50.44 | | |
| ReFlow | 19.62 | 0.767 | 0.221 | 91.93 | | |
| PMRF | 19.32 | 0.753 | 0.189 | 81.59 | | |
| FoD (Ours) | 21.61 | 0.819 | 0.105 | 41.31 | | |

(a) Deraining results on the Rain100H dataset.

(b) Low-light enhancement on the LOL dataset.

| Method | Disto | ortion | Perceptual | | |
|----------------|-------|--------|------------|-------|--|
| Wiemod | PSNR↑ | SSIM↑ | LPIPS↓ | FID↓ | |
| U-Net baseline | 22.88 | 0.906 | 0.065 | 15.65 | |
| IR-SDE | 25.25 | 0.906 | 0.060 | 8.33 | |
| GOUB | 25.31 | 0.908 | 0.048 | 8.21 | |
| ReFlow | 20.84 | 0.864 | 0.081 | 23.53 | |
| PMRF | 22.45 | 0.868 | 0.092 | 24.09 | |
| FoD (Ours) | 26.57 | 0.932 | 0.033 | 8.14 | |

| Disto | rtion | Perceptual | | |
|-------|--|--|---|--|
| PSNR↑ | SSIM↑ | LPIPS↓ | FID↓ | |
| 27.97 | 0.889 | 0.097 | 58.78 | |
| 29.83 | 0.904 | 0.045 | 26.30 | |
| 29.81 | 0.916 | 0.039 | 23.39 | |
| 29.84 | 0.912 | 0.065 | 38.65 | |
| 30.45 | 0.901 | 0.082 | 55.40 | |
| 30.28 | 0.923 | 0.029 | 16.12 | |
| | PSNR↑ 27.97 29.83 29.81 29.84 30.45 | 27.97 0.889 29.83 0.904 29.81 0.916 29.84 0.912 30.45 0.901 | PSNR↑ SSIM↑ LPIPS↓ 27.97 0.889 0.097 29.83 0.904 0.045 29.81 0.916 0.039 29.84 0.912 0.065 30.45 0.901 0.082 | |

(c) Dehazing results on the RESIDE-6k dataset. (d) Inpainting results on the CelebA-HQ dataset.

 μ . Setting s to 0 and rewriting this solution yields an interpolation between x_0 and μ :

$$x_t = x_0 \alpha_t + \mu (1 - \alpha_t), \quad \alpha_t = e^{-\int_0^t \theta_z dz}, \tag{14}$$

which forms a similar transportation path as in flow matching but with a special velocity field given by θ_t ($\mu - x_t$). We can then learn the drift, resulting in a conditional flow matching objective:

$$L_{\text{CFM}} := \mathbb{E}_{\mu, x_t} \Big[\|(\mu - x_t) - f_{\phi}(x_t, t)\|^2 \Big] = \mathbb{E}_{\mu, x_t} \Big[\|\alpha_t(\mu - x_0) - f_{\phi}(x_t, t)\|^2 \Big], \tag{15}$$

which is a deterministic form of the stochastic flow matching in Eq. (11). In practice, we can learn the target displacement $\mu - x_0$ directly and define the α schedule to be linear, i.e., decreasing from 1 to 0, in which case this mean-reverting ODE becomes flow matching with a straight-line path (Lipman et al., 2022; Liu et al., 2022) exactly. In other words, our primary FoD model can also be regarded as a stochastic extension of flow matching models.

EXPERIMENTS

Our experiments mainly focus on image restoration (IR), a fundamental problem in computer vision which aims to accurately recover high-quality images from their degraded low-quality counterparts. The general applicability of our FoD model on image-conditioned generation is further demonstrated via qualitative results on diverse image-to-image translation tasks¹.

Implementation and Setup We use a U-Net (Ronneberger et al., 2015) architecture similar to DDPM (Ho et al., 2020) for flow prediction in all tasks. Attention layers are removed for efficient training and testing, similar to IR-SDE (Luo et al., 2023a;b). We choose the commonly used cosine and linear schedules (Nichol & Dhariwal, 2021) for θ_t and σ_t , respectively, and normalize σ_t^2 to sum to 1 to ensure numerical stability under multiplicative noise perturbation. The number of sampling steps is fixed to 100 for all tasks. We use the AdamW (Loshchilov & Hutter, 2017) optimizer with parameters $\beta_1 = 0.9$ and $\beta_2 = 0.99$. The training requires 500 000 iterations with a learning rate of 10^{-4} . All our models are trained on an A100 GPU with 40 GB of memory for about 1.5 days.

The Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al., 2018) and Fréchet Inception Distance (FID) (Heusel et al., 2017) are reported to evaluate the perceptual fidelity and overall visual quality. Additionally, Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM) (Wang et al., 2004) are also included to evaluate pixel-level and structural similarity.

¹We provide more training details, datasets, and unconditional generation results in the Appendix.



Figure 2: Comparison of FoD with other approaches on four IR tasks. Here, we add one column for the results of 'FoD w/o noise' to illustrate the importance of noise injection in image restoration.

4.1 IMAGE RESTORATION

We evaluate our method on four IR tasks: 1) image deraining on the Rain100H dataset (Yang et al., 2017), 2) dehazing on the RESIDE-6k dataset (Qin et al., 2020), 3) low-light enhancement on LOL (Wei et al., 2018), and 4) face inpainting on CelebA-HQ (Karras et al., 2017).

In our experiments, we select IR-SDE (Luo et al., 2023a) as the main comparison method to evaluate the performance gap between forward-backward and forward-only schemes for diffusion-based restoration. We also compare with a diffusion bridge model GOUB (Yue et al., 2024) and other flow-based approaches, including Rectified flow (Liu et al., 2022; Liu, 2022) and posterior-mean rectified flow (PMRF) (Ohayon et al., 2024), that learn ODEs and also allow the model to generate images with a single forward process. In addition, a U-Net model, using the same architecture as our FoD, is trained with the ℓ_1 loss as a CNN baseline on all tasks for reference.

The quantitative comparisons on four IR tasks are reported in Table 1. The proposed FoD achieves the best results across all datasets in comparison to other diffusion-based and flow-based approaches. Compared to the U-Net baseline, IR-SDE and GOUB successfully improve the results on perceptual metrics (LPIPS and FID) across all tasks, proving the effectiveness of the forward-backward based diffusion IR schemes. We observe that flow-based approaches, such as ReFlow and PMRF, perform inferiorly on all distortion and perceptual metrics. While PMRF improves the PSNR results for flow matching-based IR, the performance gain potentially comes from the two-stage training strategy and the small noise injection in the initial state of rectified flow.

We also provide visual comparisons in Figure 2, showing that our FoD produces the most realistic and high-fidelity results. In particular, while all deterministic approaches without noise injection (ReFlow, PMRF and 'FoD w/o noise') tend to generate overly smooth outputs (see e.g. the left eye area in the face inpainting case), the proposed FoD model consistently produces sharper and more detailed images. Further discussion on the role of noise injection is provided in Section 5.

4.2 IMAGE-TO-IMAGE TRANSLATION

We also perform qualitative experiments on diverse image-to-image translation tasks, to further demonstrate the general applicability of the proposed FoD method in image-conditioned generation. These tasks include edges to handbags and shoes (Isola et al., 2017), facades to labels (Tyleček & Šára, 2013), aerial photos to maps (Isola et al., 2017), and night to day (Laffont et al., 2014). We adopt the same implementation as in the image restoration setting, except that all images are resized to 64×64 resolution. Qualitative results across these different tasks are shown in Figure 3. One can observe that FoD is capable of handling complex image-to-image translation problems, even when the source and target domains differ significantly, such as in edges to photos or photos to labels/maps.

Table 2: Results of different sampling approaches using the same trained FoD model. Here, 'FoD w/ EM' denotes the *100-step* Euler–Maruyama sampling method, while 'FoD w/ MC' and 'FoD w/ NMC' denote *10-step* Markov and non-Markov chain fast sampling, respectively.

| Method | Deraining | | Deraining Low-light enhance | | Dehazing | | | Inpainting | | | | |
|------------|-----------|--------|-----------------------------|-------|----------|-------|-------|------------|-------|-------|--------|-------|
| | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ |
| FoD w/ EM | 32.56 | 0.038 | 14.10 | 21.61 | 0.105 | 41.31 | 26.57 | 0.033 | 8.14 | 30.28 | 0.029 | 16.12 |
| FoD w/ MC | 33.27 | 0.039 | 15.14 | 23.12 | 0.093 | 32.37 | 26.76 | 0.031 | 10.07 | 31.02 | 0.031 | 18.06 |
| FoD w/ NMC | 33.63 | 0.041 | 15.64 | 23.05 | 0.098 | 47.87 | 26.77 | 0.032 | 10.31 | 31.32 | 0.038 | 23.28 |



Figure 3: Qualitative results of our FoD method on diverse image-to-image translation tasks.

Notably, the night to day dataset involves many-to-many mappings due to temporal variability, yet FoD still produces satisfactory results, demonstrating its strong generative capability.

5 DISCUSSION

Fast Sampling As discussed in Section 3.3, FoD naturally supports fast sampling using both Markov and non-Markov chains. Table 2 demonstrates that with only 10 sampling steps, the non-Markov variant achieves even better results than the standard Euler-Maruyama solver in terms of PSNR, without significantly compromising perceptual quality. To investigate the impact of sampling steps, we provide a detailed analysis on image restoration in Figure 4. For comparison, the Euler-Maruyama solver with 100 steps is included as a baseline. We observe that both sampling strategies yield improved distortion metrics as the number of steps decreases, especially in the lowstep areas (5–20), where they substantially outperform the baseline with a small perceptual performance drop. Similar trends are observed for LPIPS and FID, although both metrics show a modest decline in performance when reducing the number of steps from 20 to 5. More comparisons and an illustration of the forward process are provided in

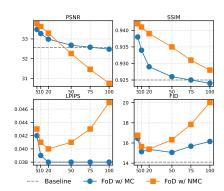


Figure 4: Comparison of fast sampling with Markov chain (MC) vs. with non-Markov chain (NMC) on deraining task, using the same pretrained model.

Appendix C.2. The overall results suggest that using 10 sampling steps can serve as a practical rule of thumb for fast sampling with both Markov and non-Markov chains.

Effectiveness of Noise Injection This section explores the importance of noise injection in image-conditioned generation. To this end, we conduct an additional experiment where a noise-free variant of FoD is trained on four image restoration tasks, by setting $\sigma_t = 0$ for all t. This effectively reduces FoD to a flow matching model, as described in Section 3.4. We keep the θ schedule the same as FoD for a fair comparison. Quantitative results in Table 3 show a significant drop in performance across all tasks and metrics. The corresponding visual results are also provided in the second-to-last column of Figure 2, where the produced images are blurry and unclear compared to those of our original FoD. Moreover, the training curves of FoD with and without noise injection on different image restoration tasks are provided in Figure 5. FoD consistently outperforms its noise-free variant on all tasks, further demonstrating the effectiveness and importance of noise injection in image-conditioned generation.

Table 3: Quantitative results of FoD and its noise-free variant on four image restoration tasks.

| Method | | Deraining | | Low | -light enha | nce | Dehazing | | | Inpainting | | |
|---------------|-------|-----------|-------|-------|-------------|-------|----------|--------|-------|------------|--------|-------|
| | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ | PSNR↑ | LPIPS↓ | FID↓ |
| FoD w/o noise | 29.44 | 0.132 | 59.54 | 19.96 | 0.193 | 86.21 | 24.18 | 0.048 | 13.80 | 29.94 | 0.065 | 38.78 |
| FoD (Ours) | 32.56 | 0.038 | 14.10 | 21.61 | 0.105 | 41.31 | 26.57 | 0.033 | 8.14 | 30.28 | 0.029 | 16.12 |

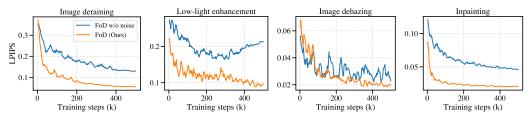


Figure 5: Training curves of FoD and its noise-free variant on four image restoration tasks.

Limitations and Future Work Although FoD performs well for image-conditioned generation, its multiplicative structure in the forward process poses a challenge for unconditional generation, where the source distribution is typically Gaussian. Specifically, injecting log-Gaussian noise (as defined in Eq. (6)) into a source sample $x_0 \sim \mathcal{N}(0, I)$ complicates the learning process and leads to a decline in sample quality (See Appendix C.3 for additional details). In future work, we plan to explore more advanced strategies, such as log-space transformations and optimal transport-based drift paths, to further improve the unconditional generation capabilities of FoD. Additionally, this work adopts the commonly used cosine and linear noise schedules and finds them effective for most tasks. Alternatives will be explored in future work, as we believe that FoD offers a more flexible framework for exploring different schedules due to its forward-only formulation.

6 RELATED WORK

Denoising diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021; Karras et al., 2022) and flow matching models (Lipman et al., 2022; Liu et al., 2022; Lipman et al., 2024; Albergo et al., 2023a; Gat et al., 2024) are two popular frameworks in generative modelling and have been widely applied to various applications including image generation (Dhariwal & Nichol, 2021; Ho et al., 2022; Peebles & Xie, 2023; Ho & Salimans, 2022; Ma et al., 2024), text-to-image generation (Rombach et al., 2022; Ruiz et al., 2023; Podell et al., 2024; Saharia et al., 2022b), image translation (Meng et al., 2022; Lugmayr et al., 2022; Saharia et al., 2022a; Su et al., 2022; Xia et al., 2024b; Liu et al., 2023; Ben-Hamu et al., 2024; Li et al., 2023; Xia et al., 2024a; Zheng et al., 2024), etc. Inspired by their success in producing photo-realistic images conforming to human preference, these models have recently been applied to image restoration for advanced performance (Wang et al., 2024; Saharia et al., 2022c; Yue et al., 2023; Kawar et al., 2022; Yue et al., 2024; Luo et al., 2024a;b; Liu et al., 2024; Shi et al., 2024). In addition, Albergo et al. (2023a) unify diffusion and flow matching models through stochastic interpolants. This formulation has also been applied to image restoration (Albergo et al., 2023b), where stochastic flows guided by corrupted observations recover clean images, effectively serving as a stochastic extension of flow matching for inverse problems. Subsequent works enhance this approach using pretrained flow matching models (Ben-Hamu et al., 2024) or refined training pipelines (Ohayon et al., 2024). These works are closely related to ours, but they all adopt noise-free generation processes. In contrast, FoD involves a state-dependent diffusion process for image generation, which is naturally well-suited for image-conditioned generation, and which degrades to flow matching when the diffusion term vanishes.

7 CONCLUSION

This paper presents a new framework, named FoD, for generative modelling with a single forward diffusion process. We show that FoD is analytically tractable and can be trained using a simple flow matching objective. Our model is evaluated on various image-conditioned generation tasks, including image restoration and image-to-image translation. FoD achieves strong performance compared to other diffusion models and flow matching approaches, demonstrating its effectiveness and efficiency in generative modelling, particularly for image restoration.

REFERENCES

- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. arXiv preprint arXiv:2209.15571, 2022. 2
- Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. arXiv preprint arXiv:2303.08797, 2023a. 2, 3, 5, 9
- Michael S Albergo, Mark Goldstein, Nicholas M Boffi, Rajesh Ranganath, and Eric Vanden-Eijnden. Stochastic interpolants with data-dependent couplings. arXiv preprint arXiv:2310.03725, 2023b.
- Heli Ben-Hamu, Omri Puny, Itai Gat, Brian Karrer, Uriel Singer, and Yaron Lipman. D-flow: Differentiating through flows for controlled generation. arXiv preprint arXiv:2402.14017, 2024. 9
- Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern</u> Recognition (CVPR), pp. 18392–18402, 2023. 1
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. <u>Advances in Neural Information Processing Systems (NeurIPS)</u>, 34:8780–8794, 2021. 1, 9
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. Advances in Neural Information Processing Systems (NeurIPS), 37:133345–133385, 2024. 9
- Daniel T Gillespie. Exact numerical simulation of the Ornstein-Uhlenbeck process and its integral. Physical review E, 54(2):2084, 1996. 1, 3
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In <u>Proceedings of Advances in Neural Information Processing Systems (NeurIPS)</u>, volume 30, 2017. 6
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. <u>arXiv preprint arXiv:2207.12598</u>, 2022. 9
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems (NeurIPS), 33:6840–6851, 2020. 1, 2, 4, 6, 9, 16, 18
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. <u>Journal of Machine Learning</u> Research, 23(47):1–33, 2022. 9
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1125–1134, 2017. 7, 17
- Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. arXiv preprint arXiv:1710.10196, 2017. 7, 25
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. Advances in Neural Information Processing Systems (NeurIPS), 35: 26565–26577, 2022. 1, 9
- Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. <u>Advances in Neural Information Processing Systems (NeurIPS)</u>, 35:23593–23606, 2022. 1, 9
- Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 6007–6017, 2023. 1
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Toronto, ON, Canada, 2009. 18

- Pierre-Yves Laffont, Zhile Ren, Xiaofeng Tao, Chao Qian, and James Hays. Transient attributes for high-level understanding and editing of outdoor scenes. <u>ACM Transactions on graphics (TOG)</u>, 33(4):1–11, 2014. 7
 - Bo Li, Kaitao Xue, Bin Liu, and Yu-Kun Lai. Bbdm: Image-to-image translation with brownian bridge diffusion models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 1952–1961, 2023. 9
 - Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Yu Qiao, Wanli Ouyang, and Chao Dong. Diffbir: Toward blind image restoration with generative diffusion prior. In European Conference on Computer Vision, pp. 430–448. Springer, 2024. 1
 - Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. arXiv preprint arXiv:2210.02747, 2022. 2, 5, 6, 9, 19
 - Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. <u>arXiv:2412.06264</u>, 2024. 2, 9
 - Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos Theodorou, Weili Nie, and Anima Anandkumar. I²SB: Image-to-image Schrödinger bridge. In <u>International Conference on Machine Learning (ICML)</u>, pp. 22042–22062. PMLR, 2023. 1, 9
 - Jiawei Liu, Qiang Wang, Huijie Fan, Yinong Wang, Yandong Tang, and Liangqiong Qu. Residual denoising diffusion models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 2773–2783, 2024. 9
 - Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. <u>arXiv preprint</u> arXiv:2209.14577, 2022. 7
 - Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. arXiv preprint arXiv:2209.03003, 2022. 2, 6, 7, 9, 19, 20
 - Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. <u>arXiv preprint</u> <u>arXiv:1711.05101, 2017.</u> 6
 - Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 11461–11471, 2022. 9, 17
 - Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Image restoration with mean-reverting stochastic differential equations. In International Conference on Machine Learning (ICML), pp. 23045–23066. PMLR, 2023a. 1, 3, 6, 7
 - Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Refusion: Enabling large-size realistic image restoration with latent-space diffusion models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 1680–1691, 2023b. 6
 - Ziwei Luo, Fredrik K. Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B. Schön. Controlling vision-language models for multi-task image restoration. In <u>The Twelfth International Conference on Learning Representations</u>, 2024a. URL https://openreview.net/forum?id=t3vnnLeajU.9
 - Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Photo-realistic image restoration in the wild with controlled vision-language models. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 6641–6651, 2024b. 9
 - Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In European Conference on Computer Vision, pp. 23–40. Springer, 2024. 9

- Ségolène Martin, Anne Gagneux, Paul Hagemann, and Gabriele Steidl. PnP-Flow: Plug-and-play image restoration with flow matching. arXiv preprint arXiv:2410.02423, 2024. 2
 - 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 <u>International</u> Conference on Learning Representations, 2022. 1, 9, 19
 - Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In International Conference on Machine Learning (ICML), pp. 8162–8171. PMLR, 2021. 6
 - Guy Ohayon, Tomer Michaeli, and Michael Elad. Posterior-mean rectified flow: Towards minimum MSE photo-realistic image restoration. arXiv preprint arXiv:2410.00418, 2024. 2, 7, 9
 - William Peebles and Saining Xie. Scalable diffusion models with transformers. In <u>Proceedings of</u> the IEEE/CVF International Conference on Computer Vision, pp. 4195–4205, 2023. 1, 9
 - Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In The Twelfth International Conference on Learning Representations, 2024. 1, 9
 - Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. FFA-Net: Feature fusion attention network for single image dehazing. In Proceedings of the AAAI Conference on Artificial Intelligence, pp. 11908–11915, 2020. 7, 17, 24
 - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In <u>Proceedings of the IEEE/CVF</u> Conference on Computer Vision and Pattern Recognition (CVPR), pp. 10684–10695, 2022. 1, 9
 - Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18, pp. 234–241. Springer, 2015. 6
 - Sheldon M Ross. Introduction to probability models. Academic press, 2014. 4, 14
 - Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 22500–22510, 2023. 9
 - Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In <u>ACM SIGGRAPH 2022</u> conference proceedings, pp. 1–10, 2022a. 1, 9
 - 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. <u>Advances in Neural Information</u> Processing Systems (NeurIPS), 35:36479–36494, 2022b. 1, 9
 - Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. <u>IEEE transactions on pattern analysis and machine</u> intelligence, 45(4):4713–4726, 2022c. 1, 9
 - Zhenning Shi, Chen Xu, Changsheng Dong, Bin Pan, Along He, Tao Li, Huazhu Fu, et al. Resfusion: Denoising diffusion probabilistic models for image restoration based on prior residual noise. Advances in Neural Information Processing Systems (NeurIPS), 37:130664–130693, 2024. 9
 - Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In <u>International Conference on Machine Learning (ICML)</u>, pp. 2256–2265. pmlr, 2015. 1, 2, 9
 - Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Advances in Neural Information Processing Systems (NeurIPS), 32, 2019. 2, 18

- Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. Advances in Neural Information Processing Systems (NeurIPS), 33:12438–12448, 2020. 18
 - Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. <u>International</u> Conference on Learning Representations, 2021. 1, 2, 3, 9, 16, 18
 - Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual diffusion implicit bridges for image-to-image translation. arXiv preprint arXiv:2203.08382, 2022. 9
 - Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1921–1930, 2023. 1
 - Radim Tyleček and Radim Šára. Spatial pattern templates for recognition of objects with regular structure. In German conference on pattern recognition, pp. 364–374. Springer, 2013. 7
 - Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. <u>International Journal of Computer Vision</u>, 132(12):5929–5949, 2024. 1, 9
 - Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. <u>IEEE Transactions on Image Processing</u>, 13(4):600–612, 2004. 6
 - Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. arXiv preprint arXiv:1808.04560, 2018. 7, 17, 24
 - Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang, Radu Timofte, and Luc Van Gool. Diffi2i: Efficient diffusion model for image-to-image translation. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2024a. 9
 - Mengfei Xia, Yu Zhou, Ran Yi, Yong-Jin Liu, and Wenping Wang. A diffusion model translator for efficient image-to-image translation. <u>IEEE Transactions on Pattern Analysis and Machine Intelligence</u>, 2024b. 9
 - Wenhan Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan. Deep joint rain detection and removal from a single image. In <u>Proceedings of the IEEE conference on computer vision and pattern recognition</u>, pp. 1357–1366, 2017. 7, 17, 23
 - Conghan Yue, Zhengwei Peng, Junlong Ma, Shiyan Du, Pengxu Wei, and Dongyu Zhang. Image restoration through generalized ornstein-uhlenbeck bridge. In <u>International Conference on Machine Learning (ICML)</u>, pp. 58068–58089. PMLR, 2024. 7, 9
 - Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. <u>Advances in Neural Information Processing Systems</u> (NeurIPS), 36:13294–13307, 2023. 9
 - Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <u>Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</u>, pp. 586–595, 2018. 6
 - Dian Zheng, Xiao-Ming Wu, Shuzhou Yang, Jian Zhang, Jian-Fang Hu, and Wei-Shi Zheng. Selective hourglass mapping for universal image restoration based on diffusion model. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 25445–25455, 2024. 9

A PROOFS

A.1 Proof to Proposition 3.1 and Corollary 3.2

Proposition 3.1. Given an initial state x_s at time s < t, the unique solution to the SDE (5) is

$$x_t = (x_s - \mu) e^{-\int_s^t \left(\theta_z + \frac{1}{2}\sigma_z^2\right) dz + \int_s^t \sigma_z dw_z + \mu, \tag{16}$$

where the stochastic integral is interpreted in the Itô sense and can be reparameterised as $\bar{\sigma}_{s:t} \epsilon$, where $\bar{\sigma}_{s:t} = \sqrt{\int_s^t \sigma_z^2 \, \mathrm{d}z}$, and $\epsilon \sim \mathcal{N}(0, I)$ is a standard Gaussian noise.

Corollary 3.2. Under the same assumptions as in Proposition 3.1, the stochastic flow field $\mu - x_t$ satisfies the multiplicative stochastic structure. More precisely, it is log-normally distributed by

$$\log(\mu - x_t) \sim \mathcal{N}\left(\log(\mu - x_s) - \int_s^t \left(\theta_z + \frac{1}{2}\sigma_z^2\right) dz, \int_s^t \sigma_z^2 dz I\right). \tag{17}$$

Proof. Recall the FoD process from Eq. (5):

$$dx_t = \theta_t (\mu - x_t) dt + \sigma_t (x_t - \mu) dw_t.$$
(18)

To solve this SDE, we introduce a new variable $y_t = x_t - \mu$ and replace it into Eq. (18), which typically yields a Geometric Brownian motion (Ross, 2014) on y_t , given by

$$dy_t = -\theta_t y_t dt + \sigma_t y_t dw_t. (19)$$

To simplify the notation, we use y rather than y_t in all the following equations. This equation can be solved by applying Itô's formula:

$$d\psi(y,t) = \frac{\partial \psi}{\partial t}(y,t) dt + \frac{\partial \psi}{\partial y}(y,t) f(y,t) dt$$

$$+ \frac{1}{2} \frac{\partial^2 \psi}{\partial y^2}(y,t) g(t)^2 dt$$

$$+ \frac{\partial \psi}{\partial y}(y,t) g(t) dw,$$
(20)

where $\psi(y,t) = \ln |y|$ is a surrogate differentiable function. By substituting f(y,t) and g(t) with the drift and the diffusion functions in (19), we obtain

$$d\psi(y,t) = -(\theta_t + \frac{\sigma_t^2}{2}) dt + \sigma_t dw.$$
(21)

Then we can solve y_t conditioned on y_s , by integrating both sides:

$$\ln|y_t| - \ln|y_s| = -\int_s^t (\theta_z + \frac{\sigma_z^2}{2}) \,\mathrm{d}z + \int_s^t \sigma_z \mathrm{d}w(z)$$
 (22)

where the stochastic interaction follow a Gaussian distribution, i.e., $\int_s^t \sigma_z \, \mathrm{d}w(z) \sim \mathcal{N} \left(0, \int_s^t \sigma_z^2 \mathrm{d}z\right)$, then we can rewrite:

$$\ln|y_t| = \ln|y_s| - (\bar{\theta}_{s:t} + \frac{\bar{\sigma}_{s:t}^2}{2}) + \bar{\sigma}_{s:t}\epsilon_{s\to t}, \quad \epsilon_{s\to t} \sim \mathcal{N}(0, I), \tag{23}$$

where $\bar{\theta}_{s:t} = \int_s^t \theta_z \, dz$, $\bar{\sigma}_{s:t}^2 = \int_s^t \sigma_z^2 \, dz$, and $\bar{\sigma}_{s:t} = \sqrt{\bar{\sigma}_{s:t}^2}$. By replacing y_t with the original $x_t - \mu$, we have the following

$$\ln|x_t - \mu| = \ln|x_s - \mu| - (\bar{\theta}_{s:t} + \frac{\bar{\sigma}_{s:t}^2}{2}) + \bar{\sigma}_{s:t}\epsilon_{s \to t}.$$
 (24)

Note that the sign of $x_t - \mu$ remains consistent across all times t, therefore, we can safely omit the absolute value inside the logarithm, which leads to a log-normal distribution:

$$\log(\mu - x_t) \sim \mathcal{N}\left(\log(\mu - x_s) - \int_s^t \left(\theta_z + \frac{1}{2}\sigma_z^2\right) dz, \int_s^t \sigma_z^2 dz I\right), \tag{25}$$

which gives the Corollary 3.2. In addition, applying the exponential function to both sides yields

$$(x_t - \mu) = (x_s - \mu)e^{-(\bar{\theta}_{s:t} + \frac{\bar{\sigma}_{s:t}^2}{2}) + \bar{\sigma}_{s:t}\epsilon_{s \to t}}$$

$$(26)$$

which is the solution to the SDE, and thus we complete the proof.

A.2 PROOF OF THE KL DIVERGENCE

The KL divergence of (9) can be derived as follows:

$$\begin{split} \tilde{L} &:= -\log p_{\phi}(x_T) \\ &= -\log \int p_{\phi}(x_{0:T}) \, \mathrm{d}x_{0:T-1} \\ &= -\log \int \frac{p_{\phi}(x_{0:T}) \, \mathrm{d}x_{0:T-1} \, | \, x_T)}{q(x_{0:T-1} \, | \, x_T)} \, \mathrm{d}x_{0:T-1} \\ &= -\log \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[\frac{p_{\phi}(x_{0:T})}{q(x_{0:T-1} | \, x_T)} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log \frac{p_{\phi}(x_{0:T})}{q(x_{0:T-1} | \, x_T)} \right] & \text{(Jensen's Inequality)} \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log \frac{p(x_0) \prod_{t=1}^T p_{\phi}(x_t | \, x_{t-1})}{\prod_{t=1}^T q(x_{t-1} | \, x_t)} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log p(x_0) - \sum_{t=1}^T \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_{t-1} | \, x_t)} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log p(x_0) - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} \cdot \frac{q(x_t | \, x_T)}{q(x_{t-1} | \, x_T)} - \log \frac{p_{\phi}(x_T | \, x_{T-1})}{q(x_{T-1} | \, x_T)} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log p(x_0) - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log \frac{q(x_{T-1} | \, x_T)}{q(x_{T-1} | \, x_T)} - \log \frac{p_{\phi}(x_T | \, x_{T-1})}{q(x_{T-1} | \, x_T)} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log p(x_0) - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log \frac{q(x_{T-1} | \, x_T)}{q(x_0 | \, x_T)} - \log \frac{p_{\phi}(x_T | \, x_{T-1})}{q(x_T | \, x_{T-1})} \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log \frac{p(x_0)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log p_{\phi}(x_T | \, x_{T-1}) \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log \frac{p(x_0)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log p_{\phi}(x_T | \, x_{T-1}) \right] \\ &= \mathbb{E}_{q(x_{0:T-1} | x_T)} \left[-\log \frac{p(x_0)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log p_{\phi}(x_T | \, x_{T-1}) \right] \\ &= \mathbb{E}_{q(x_0 | \, x_T)} \left[-\log \frac{p(x_0)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log \frac{p_{\phi}(x_t | \, x_{t-1})}{q(x_t | \, x_{t-1}, x_T)} - \log p_{\phi}(x_T | \, x_{T-1}) \right] \\ &= \mathbb{E}_{q(x_0 | \, x_T)} \left[-\log \frac{p(x_0 | \, x_T)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log \frac{p(x_0 | \, x_T)}{q(x_0 | \, x_T)} \right] \\ &= \mathbb{E}_{q(x_0 | \, x_T)} \left[-\log \frac{p(x_0 | \, x_T)}{q(x_0 | \, x_T)} - \sum_{t=1}^{T-1} \log$$

 $+ \sum_{t=1}^{T-1} D_{KL}(q(x_t \mid x_{t-1}, x_T) \mid\mid p_{\phi}(x_t \mid x_{t-1})) - \mathbb{E}_{q(x_{T-1} \mid x_T)} \Big[\log p_{\phi}(x_T \mid x_{T-1}) \Big],$

where the first term can be ignored since it doesn't have trainable parameters, and the third term can be merged to the final stochastic flow matching objective.

A.3 DERIVATION OF STOCHASTIC FLOW MATCHING

Remark. For two log-normal distributions with $\log p_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ and $\log p_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, the KL divergence between them is given by

$$D_{KL}(p_1||p_2) = \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2} + \frac{\sigma_1^2}{2\sigma_2^2} + \ln\frac{\sigma_2}{\sigma_1} - \frac{1}{2}.$$
 (28)

This result is for scalars but can be naturally extended to high-dimensional cases. In the following, we will use it to derive our stochastic flow matching objective.

More specifically, given the KL divergence

$$\mathbb{E}_{p}\left[\sum_{t=0}^{T-1} D_{KL}(p(f_{\mu}(x_{t+1}) \mid f_{\mu}(x_{t}) \parallel p(f_{\phi}(x_{t+1}, t) \mid f_{\phi}(x_{t}, t)))\right]$$
(29)

and the truth that $f_{\phi}(x_t, t)$ approximates $\mu - x_t$. Since the two transitions are log-normally distributed and share the same parameters θ_t and σ_t as in Eq. (7), we can initially obtain a log space loss based

on the Remark (28), as

$$L := \mathbb{E}_{\mu \sim p_{\text{data}}, x_t \sim q(x_t | x_0, \mu)} \left[\frac{1}{2\sigma_{t+1}^2} \|\log f_{\mu}(x_t, t) - \log f_{\phi}(x_t, t)\|^2 \right]. \tag{30}$$

However, this would run into issues when $\mu - x_t \le 0$. Though this can be alleviated using absolute values and adding a small additive term ϵ , it complicates the learning and is still unstable.

To address it, we assume that, after some training, $f_{\phi}(x_t,t)/f_{\mu}(x_t,t)=1+\delta$ where $|\delta|\ll 1$. The first-order Taylor approximation is then

$$\log f_{\phi}(x_t, t) - \log f_{\mu}(x_t, t) = \log(1 + \delta) \approx \frac{f_{\phi}(x_t, t) - f_{\mu}(x_t, t)}{f_{\mu}(x_t, t)}.$$
 (31)

Since the denominator does not depend on the parameters, we obtain our simplified loss function by omitting all non-trainable weights:

$$L := \mathbb{E}_{\mu \sim p_{\text{data}}, x_t \sim q(x_t | x_0, \mu)} \left[\| f_{\mu}(x_t, t) - f_{\phi}(x_t, t) \|^2 \right], \tag{32}$$

which is the proposed stochastic flow matching objective.

Note that the first-order Taylor approximation is not valid during the early stages of training, when the predicted flow is typically far from the ground truth. In such cases, Eq. (11) no longer corresponds to an exact KL objective, but instead serves as a surrogate loss for directly learning the flow. Nevertheless, we emphasize that this surrogate objective can empirically achieve strong performance and simplify the optimization. This is conceptually analogous to denoising objectives and those used in score matching (not exact KL objectives, but have still proven effective in practice). Table 4 below tracks both the approximation error $(\log(1+\delta)-\delta)$ and magnitude of the higher-order terms $(-\frac{1}{2}\delta^2)$. As one can see, the approximation is loose during the early stages of training but becomes valid (≈ 0.1) after $\sim 50,000$ iterations.

Table 4: Tracking the approximation error and magnitude of the higher-order terms.

| Training steps | Approximation error | Magnitude |
|----------------|---------------------|-----------|
| 1,000 | 0.404 | 0.333 |
| 10,000 | 0.183 | 0.224 |
| 50,000 | 0.105 | 0.108 |
| 100,000 | 0.095 | 0.086 |

B MAXIMUM LIKELIHOOD ESTIMATION

Following DDPMs (Ho et al., 2020), both our training and sampling are implemented with discrete times, which can of course be converted into continuous times, as in Score SDEs (Song et al., 2021), but that requires θ and σ schedules to be integrable. Below, we further show that the solution to FoD also allows us to compute the maximum likelihood:

Given the clean data μ , assuming that there exists an optimal forward transition from z_t to z_{t+1} , where $z_t = |\mu - x_t|$. In other words, we want to maximize the likelihood of $p(z_{t+1} \mid z_t)$, which is a log-normal distribution as illustrated in Corollary 3.2, and its density is given by

$$p(z_{t+1} \mid z_t) = \frac{1}{z_{t+1}\sigma_{t+1}\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_{t+1}^2} \left[\ln z_{t+1} - \ln z_t + \left(\theta_{t+1} + \frac{\sigma_{t+1}^2}{2}\right)\right]^2\right).$$
(33)

Then, we can minimize the negative log-likelihood:

$$-\ln p(z_{t+1} \mid z_t) = \ln z_{t+1} + \frac{1}{2\sigma_{t+1}^2} \left[\ln z_{t+1} - \ln z_t + \left(\theta_{t+1} + \frac{\sigma_{t+1}^2}{2}\right) \right]^2 + \underline{\ln(\sigma_{t+1}\sqrt{2\pi})}, \quad (34)$$

which can be solved by setting the gradient to 0, as

$$\nabla_{z_{t+1}} - \ln p(z_{t+1} \mid z_t) = \frac{1}{z_{t+1}} \left[1 + \frac{1}{\sigma_{t+1}^2} \left(\ln z_{t+1} - \ln z_t + \left(\theta_{t+1} + \frac{\sigma_{t+1}^2}{2} \right) \right) \right] = 0.$$
 (35)

Since z_{t+1} is not zero, the optimal solution z_{t+1}^* is obtained according to

$$z_{t+1} = z_t e^{-(\theta_{t+1} + \frac{\sigma_{t+1}^2}{2}) - \sigma_{t+1}^2}.$$
(36)

Recall that $\mu - x_t$ has the same sign for all times t. Replacing z_t with $|\mu - x_t|$ gives the following:

$$(\mu - x_{t+1})^* = (\mu - x_t)e^{-(\theta_{t+1} + \frac{\sigma_{t+1}^2}{2}) - \sigma_{t+1}^2},$$
(37)

which is the optimal forward flow from x_{t+1} to μ .

Based on it, we can get the maximum likelihood learning objective:

$$L = \|(\mu - x_{t+1})^* - \mathbb{E}[\mu_{\phi} - x_{t+1}]\|^2, \tag{38}$$

$$L = \|x_{t+1}^* - \mathbb{E}_{\phi}[x_{t+1} \mid x_t]\|^2, \tag{39}$$

where $\mathbb{E}_{\phi}[x_{t+1} \mid x_t]$ is the expectation given x_t in discrete time: $\mathbb{E}_{\phi}[x_{t+1} \mid x_t] = x_t + dx_t$.

$$\mathbb{E}\left[\mu_{\phi} - x_{t+1}\right] = (\mu_{\phi} - x_{t})e^{-(\theta_{t+1} + \frac{\sigma_{t+1}^{2}}{2}) + \frac{\sigma_{t+1}^{2}}{2}}$$

$$= (\mu_{\phi} - x_{t})e^{-\theta_{t+1}}.$$
(40)

Combining (37) and (40) predicts x_{t+1} in the optimal path. This maximum likelihood-based loss function performs similarly to stochastic flow matching and can be potentially used in future work.

C MORE EXPERIMENTAL DETAILS

C.1 IMPLEMENTATION AND DATASETS

We set the number of diffusion steps to 100 for all tasks. Practically, we choose to use a discrete time implementation for our method, where we let $\bar{\theta}_t = \int_0^t \theta_z \, \mathrm{d}z \approx \sum \theta_t \Delta t$ and $\bar{\sigma}_t^2 = \int_s^t \sigma_z^2 \, \mathrm{d}z \approx \sum \sigma_t^2 \Delta t$. To ensure that FoD converges to the clean data μ , we let the deterministic exponential term at the terminal state be a smaller value, i.e. $\mathrm{e}^{-\int_0^t (\theta_s + \frac{1}{2}\sigma_s^2) \, \mathrm{d}s} = \delta = 0.001$. Solving it leads to an updated time interval $\Delta t = \frac{\log \delta}{\int_0^t (\theta_s + \frac{1}{2}\sigma_s^2) \, \mathrm{d}s}^2$. For image-to-image translation, most datasets are the same as in Pix2Pix (Isola et al., 2017), but with all images resized to 64×64 to improve training and testing efficiency. All these image-to-image translation experiments share the same settings, as the goal is to illustrate the general applicability of FoD rather than to optimize performance for each task. In addition, the details of image restoration datasets are listed below:

- Deraining: collected from the Rain100H (Yang et al., 2017) dataset containing 1800 images for training and 100 images for testing.
- *Dehazing*: collected from the RESIDE-6k (Qin et al., 2020) dataset which has mixed indoor and outdoor images with 6000 images for training and 1000 images for testing.
- Low-light enhancement: collected from the LOL (Wei et al., 2018) dataset containing 485 images for training and 15 images for testing.
- Face inpainting: we use CelebaHQ as the training dataset and divide 100 images with 100 thin masks from RePaint (Lugmayr et al., 2022) for testing.

C.2 Additional Discussions

Fast Sampling As mentioned in Section 3.3, we provide two fast sampling strategies with Markov and non-Markov Chains. Their comparison on image deraining is illustrated in Section 5 of the main paper. Here, we give more comparison results on four image restoration tasks, reporting PSNR, SSIM, LPIPS, and FID values in Figure 9. The results with standard FoD sampling (See Algorithm 2) are also reported as the baseline. It can be observed that, in most tasks, decreasing the number of steps in fast sampling approaches leads to better performance, particularly in terms of PSNR and SSIM. Their perceptual performance also decreases when the number of steps decreases from 20 to 5,

²Our code is provided in the supplementary material.



Figure 6: Visualization of the diffusion process using trained FoD models on various tasks, including deraining, dehazing, low-light enhancement, face inpainting, and unconditional generation. In each case, FoD gradually injects noise into the degraded regions and subsequently denoises these intermediate states, restoring images with enhanced and corrected details.

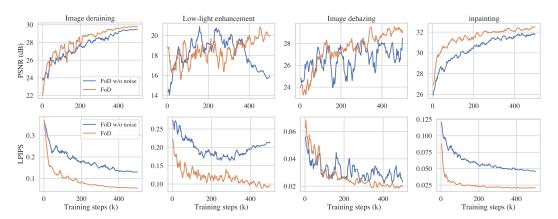


Figure 7: Training curves of our FoD and its noise-free variant on different tasks in terms of PSNR and LPIPS. All these results demonstrate the effectiveness of noise injection in image restoration.

which suggests we set the sample step to 10 as a practical rule of thumb for efficient sampling. Note that the key difference between faster sampling approaches and the standard FoD sampling is that the latter's iterative process highly depends on the μ prediction, while the former gradually refines the state x_t rather than μ . Also, we would recommend using non-Markov chain sampling for small degradation problems (which means μ prediction is easier), such as image deraining and dehazing. The visual comparisons are illustrated in Figure 10.

Illustration of the Diffusion Process To clearly show the noise injection process with our model, we apply the trained FoD on various tasks, including both image restoration and unconditional image generation, and illustrate their intermediate states along timesteps as shown in Figure 6. It can be observed that, for image restoration tasks, the noise level for each task and even for each image is different. This is given by the term $\mu - x_0$ in the solution, where areas with large difference (between LQ and HQ images) tend to produce large noise. It also means our model focuses more on restoring the degradations instead of reconstructing the whole image, yielding a more efficient solution to the image restoration problem. We also find that the noise is only injected into the degraded areas, such as the masked regions in deraining and face inpainting. Figure 7 provides additional training curves of our FoD and its noise-free variant on different tasks to illustrate the significance of noise injection in image restoration. More examples for the FoD diffusion process are provided in Figure 11.

C.3 UNCONDITIONAL IMAGE GENERATION

In this section, we evaluate the unconditional generation performance of our model on the CIFAR-10 dataset (Krizhevsky et al., 2009), using the same architecture as in image-conditioned generation (i.e., attention is not used). Specifically, we showcase results from both our FoD model and its ODE-based variant, FoD-ODE (see Section 3.4). Both models share the same θ schedule and are sampled with 100 steps. We compare against several baselines, including 1) diffusion models with forward-backward schemes such as DDPM (Ho et al., 2020) and Score SDE (Song et al., 2021), as well as 2) forward-only frameworks: score-based generative models (SGMs) like NSCN (Song & Ermon, 2019) and NSCNv2 (Song & Ermon, 2020), flow matching generative models using diffusion

| - | |
|-------------------------------|-------|
| Method | FID↓ |
| DDPM | 3.17 |
| DDPM* | 4.36 |
| Score SDE | 2.38 |
| forward-only scheme | |
| NCSN | 25.32 |
| NCSNv2 | 10.87 |
| Flow Matching w/ diff path | 10.31 |
| Flow Matching w/ OT path | 6.96 |
| Rectified Flow | 2.58 |
| FoD-SDE (T=100) | 7.89 |
| FoD-ODE (T=100) | 5.01 |
| FoD-ODE (T=1000) | 4.60 |
| FoD-ODE w/o α (T=1000) | 4.33 |

Figure 8: Left: Visual results of unconditional generation on the CIFAR-10 dataset by two variants of our FoD model. **Right:** Quantitative comparison of our methods with other approaches on CIFAR-10. Here, '*' means a re-implementation using our U-Net architecture and hyperparameter setting.

and optimal transport (OT) paths (Lipman et al., 2022), and Rectified Flow (Liu et al., 2022) which also adopts the OT path for data transformation.

Results We present the generated image samples and quantitative comparisons in Figure 8, with visual examples on the left and numerical results summarized in the table on the right. Under the forward-only framework, our FoD model with SDE sampling achieves a competitive FID of 7.89, outperforming other score-based generative models as well as the flow matching approach using diffusion paths. The ODE-based variant of FoD improves the performance to a FID of 5.01, surpassing the flow matching model based on the optimal transport path. Similar to standard diffusion models, increasing the sampling steps to 1000 improves the performance. Moreover, adopting the rectified flow objective (i.e., ignoring α_t in (15)) further decreases the FID to 4.33. Nonetheless, our overall performance on CIFAR-10 remains inferior to Rectified Flow and conventional diffusion models using a forward-backward scheme. We partly attribute this performance gap to differences in architectural choices, hyperparameter tuning, etc. However, we also note that the main aim of this paper is to construct an effective single diffusion process model, which we expect to be particularly well-suited for image-conditioned generation tasks such as image restoration.

Table 5: FID results of our FoD and its noise-free variant on four image-to-image translation tasks.

| Method | Edges to handbags | Facades to labels | Photos to maps | Night to day |
|---------------|-------------------|-------------------|----------------|--------------|
| FoD w/o noise | 25.29 | 41.33 | 17.78 | 78.52 |
| FoD (Ours) | 8.45 | 7.95 | 0.93 | 52.11 |

Table 6: Comparison of our method with other approaches on the edges to handbags dataset.

| Method | $MSE\!\!\downarrow$ | LPIPS↓ | $\text{FID}{\downarrow}$ |
|-----------------------------------|---------------------|--------|--------------------------|
| SDEdit (Meng et al., 2022) | 0.510 | 0.271 | 26.5 |
| Rectified Flow (Liu et al., 2022) | 0.088 | 0.241 | 25.3 |
| FoD (Ours) | 0.025 | 0.198 | 8.45 |

ADDITIONAL RESULTS

In this section, we provide more results for four image restoration tasks including image deraining, low-light enhancement, image dehazing, and image inpainting in Figure 12, Figure 13, Figure 14, and Figure 15. In most tasks, the results produced by our method are sharper and more realistic. For image-to-image translation, we report the FID results of our FoD and its noise-free variant on the edges to handbags, facades to labels, photos to maps, and night to day datasets in Table 5. The results further prove the importance of noise injection in image-conditioned generation. In addition, Table 6 illustrates the comparison of our FoD with other approaches (SDEdit (Meng et al., 2022) and

Rectified Flow (Liu et al., 2022)) on the edges to handbags dataset. For unconditional generation, more results of our model with SDE and ODE on CIFAR-10 are provided in Figure 16 and Figure 17, respectively. Please zoom in for the best view.

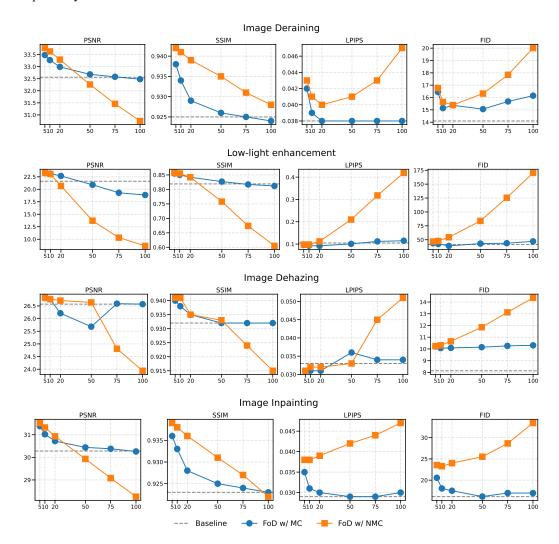


Figure 9: Comparison of different sampling approaches with pretrained FoD models on four image restoration tasks. The baseline is the Euler-Maruyama method with 100 sampling steps.

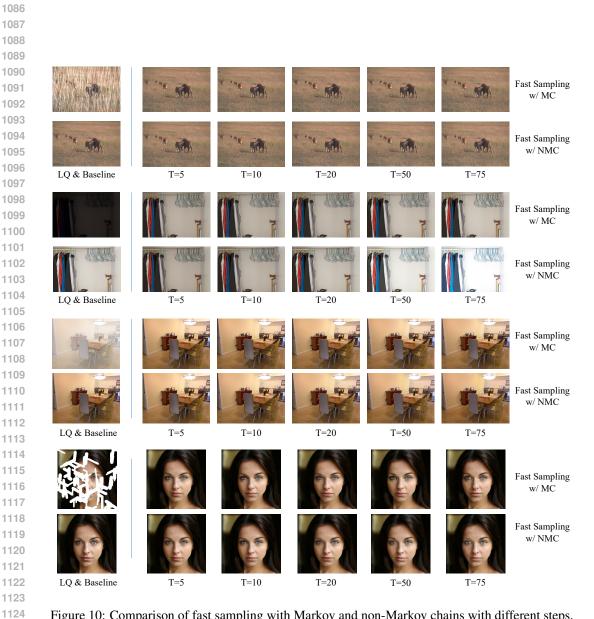


Figure 10: Comparison of fast sampling with Markov and non-Markov chains with different steps.

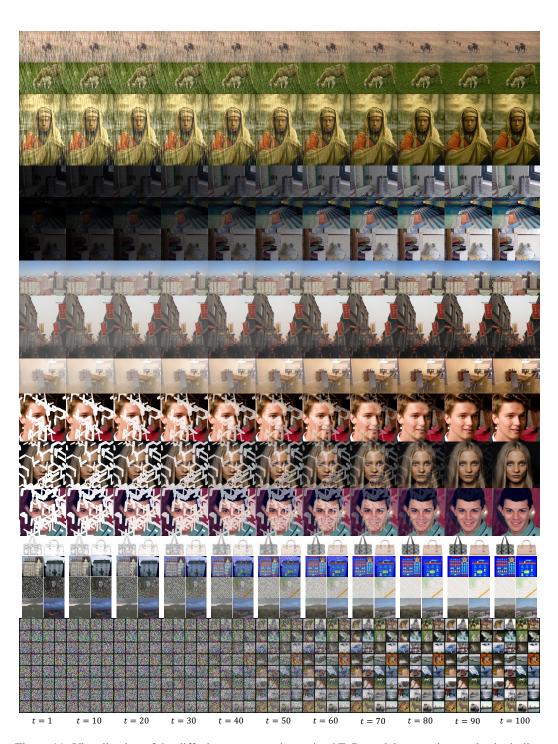


Figure 11: Visualization of the diffusion process using trained FoD models on various tasks, including deraining, dehazing, low-light enhancement, face inpainting, and unconditional generation. In each case, FoD gradually injects noise into the degraded regions and subsequently denoises these intermediate states, restoring images with enhanced and corrected details.

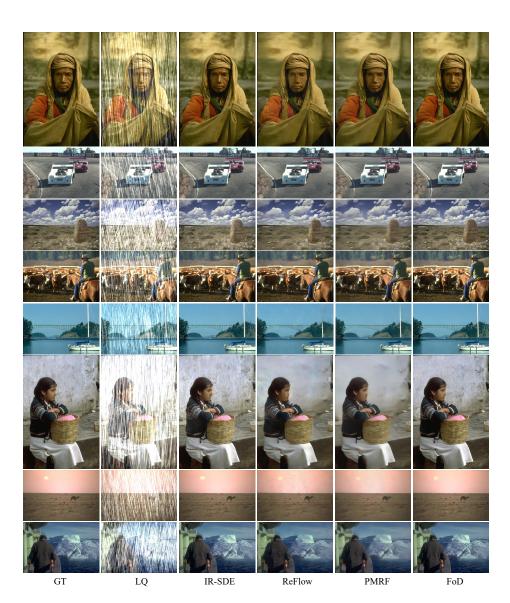


Figure 12: Visual results of image deraining on Rain100H (Yang et al., 2017) dataset.

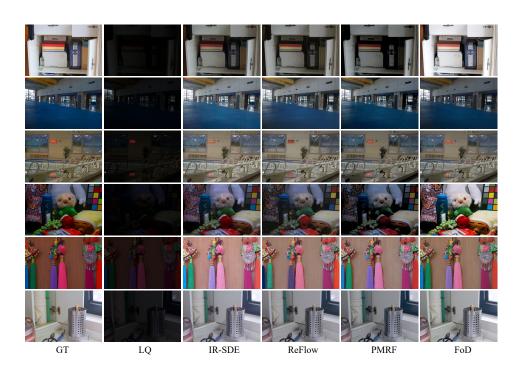


Figure 13: Visual results of image low-light enhancement on LOL (Wei et al., 2018) dataset.

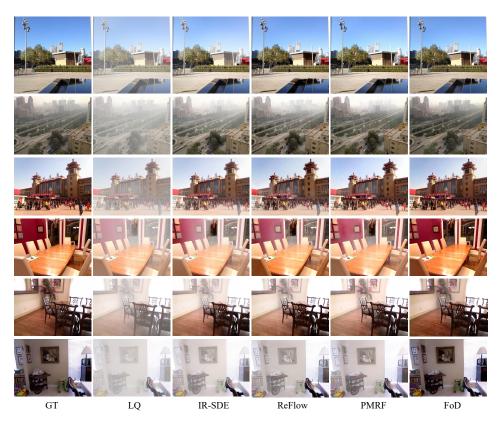


Figure 14: Visual results of image dehazing on RESIDE-6k (Qin et al., 2020) dataset.

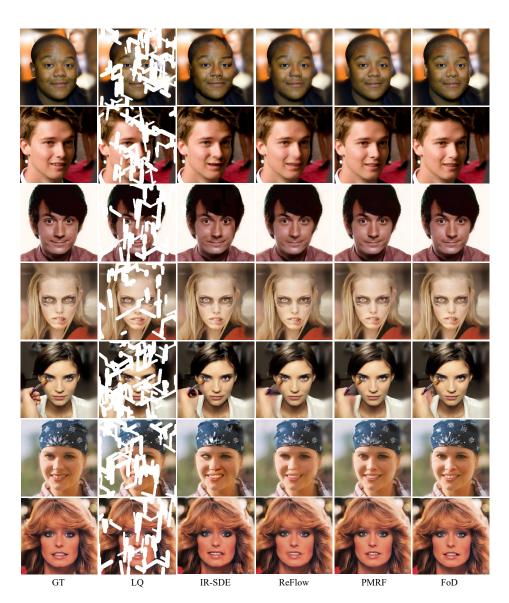


Figure 15: Visual results of image inpainting on CelebA-HQ (Karras et al., 2017) dataset.



Figure 16: Unconditional generation by FoD (SDE sampler).

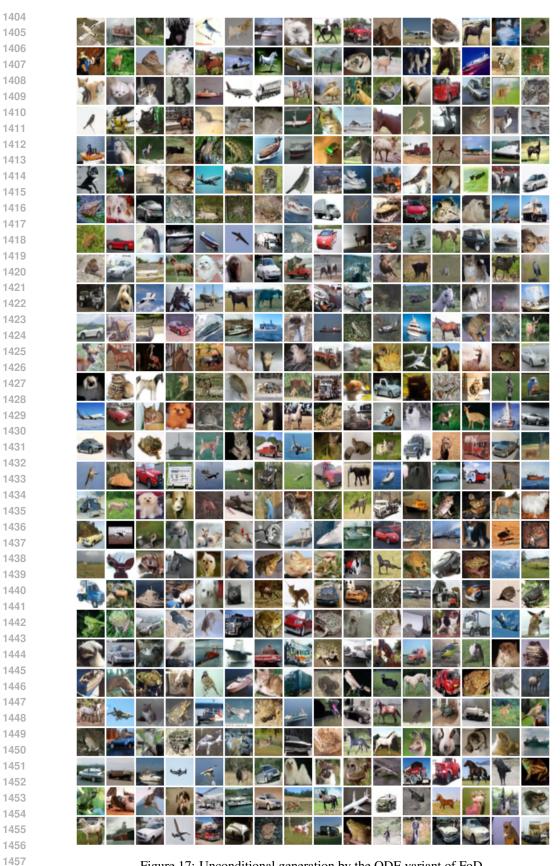


Figure 17: Unconditional generation by the ODE variant of FoD.