

Are We Learning the Score Function?

Reinterpreting Diffusion Models Through Wasserstein Gradient Flow Matching

An Vuong¹, Michael T. McCann², Javier E. Santos², Yen Ting Lin^{2,*}

¹Oregon State University ²Los Alamos National Laboratory

Standard Interpretation of Diffusion Models

Current understanding of the mechanism of diffusion model relies on the notion of reverse-time process. Take Ornstein-Uhlenbeck process (as in DDPM or as in Variance-Preserving SDE¹) for example, the **forward process**

$$dX_t = -X_t dt + \sqrt{2} dW_t \quad (1)$$

is used as the data-corruption process from $t: 0 \rightarrow \infty$ to transport $\mu_0 \rightarrow \mu_\infty$, and the **reverse-time** process from $\tau: -\infty \rightarrow 0$ to revert the process

$$dX_\tau = [X_\tau + 2\nabla_x \log p(X_\tau, -\tau)] d\tau + \sqrt{2} dW_\tau, \quad X_{-\infty} \sim \mu_\infty. \quad (2)$$

is utilized for inference/sampling.

The key idea is to use the samples generated by the forward process to learn the **score-function** $S(x, t) := \nabla_x \log \rho(x, t)$, where $\rho(x, t) := (d/dx)\mathbb{P}\{X_t > x\}$ is the probability density of the forward process.

¹ Santos & Lin, *Using Ornstein-Uhlenbeck Process to understand Denoising Diffusion Probabilistic Model and its Noise Schedules*, 2023

Otto Calculus and Wasserstein Gradient Flow

Jordan, Kinderlehrer, and Otto² and Otto³ established the equivalence of the marginal distribution induced by the forward diffusion process (1) and that by a flow described by an ordinary differential equation

$$\frac{d}{dt}x(t) = v_{\text{WGF}}(x(t)) := -x(t) - \nabla_x \log \rho(x(t), t). \quad (\text{WGF})$$

The equivalence can be established by observing the same forward Chapman-Komogorov equations (CKE); for diffusion process, the CK equation is the Fokker-Planck Equation:

$$\partial_t \rho(\xi, t|\zeta, 0) = \nabla_\xi [\xi \rho(\xi, t|\zeta, 0)] + \nabla_\xi^2 \rho(\xi, t|\zeta, 0),$$

whereas for WGF, CKE is a Liouville Equation:

$$\partial_t \rho(\xi, t|\zeta, 0) = \nabla_\xi [(\xi + \nabla_\xi \log \rho(\xi, t|\zeta, 0)) \rho(\xi, t|\zeta, 0)].$$

² Jordan, Kinderlehrer, and Otto, *The Variational Formulation of the Fokker-Planck Equation*, 1998

³ Otto, *The geometry of dissipative evolution equations: The porous medium equation*, 2001

Are We Really Learning the Score Function?

Because the score function $S(x, t)$ is a gradient of a scalar field $\log \rho(\cdot, t)$, it is a **conservative field**. Specifically, at any given time t , it must satisfy the **integral constraint**

$$\oint_{\mathcal{C}} S(x, t) dx = 0, \quad \text{for any closed path } \mathcal{C}, \quad (\text{IC})$$

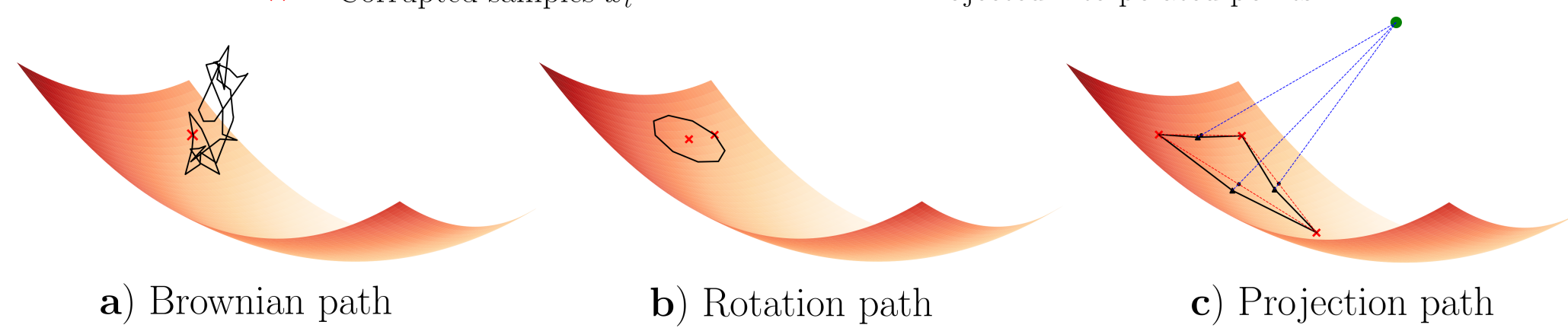
as well as the **differential constraint**

$$\partial_{x_j} S_i(x, t) = \partial_{x_i} S_j(x, t). \quad (\text{DC})$$

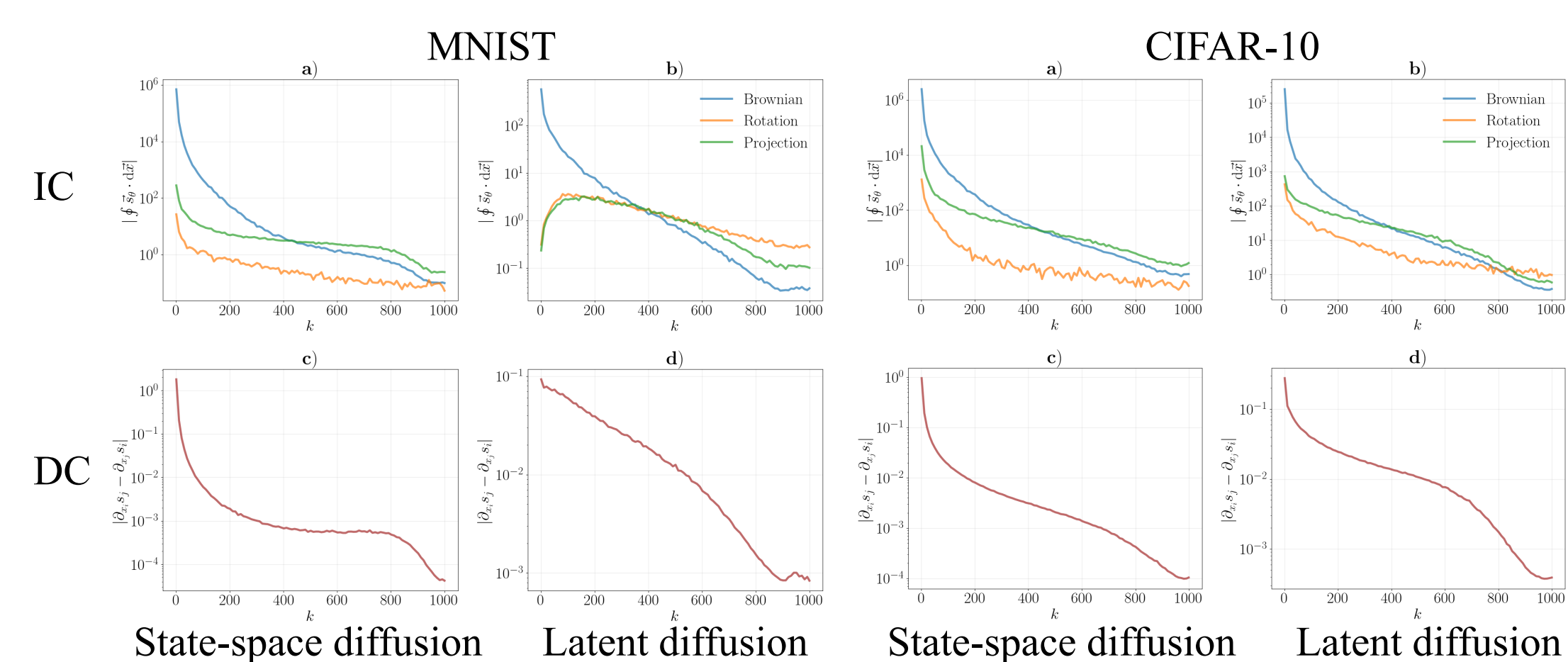
We are interested in checking if the trained diffusion models respect IC and DC.

For differential constraints, we rely on reverse-mode automatic differentiation. For integral constraints, we perform three different closed path-generation mechanisms:

- Corrupted distribution $\rho(\cdot|x_0)$
- Integration path
- Corrupted samples x_t
- Mean of the corrupted distribution $x_0 e^{-t}$
- Linearly interpolated points
- Projected interpolated points



We examined the violation of integral and differential constraints of trained diffusion models, including plain and latent diffusion, on different datasets. The results indicated that both IC and DC are violated by a significant margin:



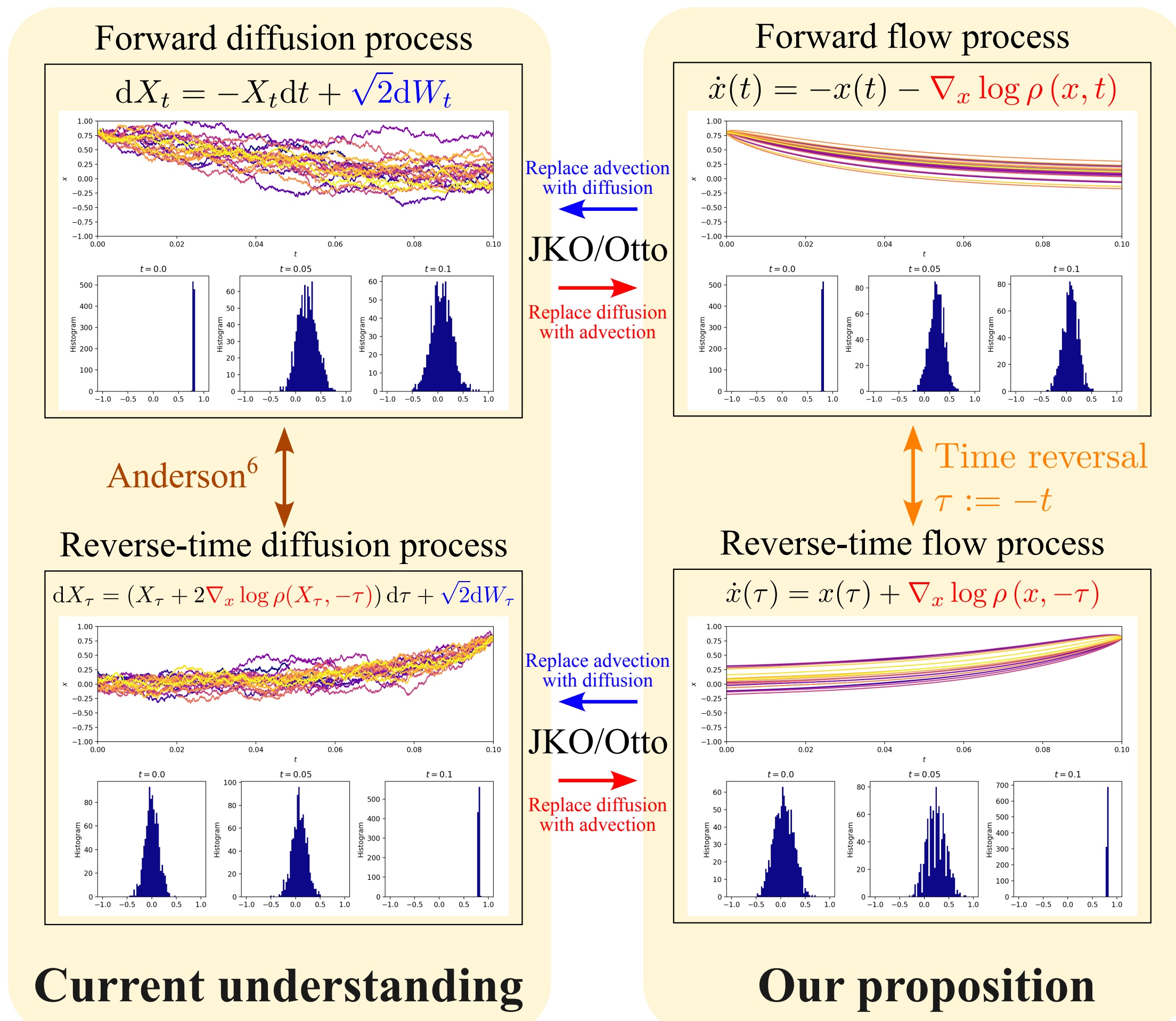
The violations raised an interesting question: Clearly, the neural network is not capable of learning the conservative score function needed for the actual reverse-time process. Why is existing computational procedure remains effective for generative modeling?

Understanding Diffusion Models as Wasserstein Gradient Flow Matching

We propose a hypothesis leveraging WGF theory to clarify what diffusion models actually learn:

Existing diffusion modeling is better understood as modeling a normalizing flow⁴, through performing flow matching⁵ to the WGF velocity, rather than learning the reverse stochastic differential equation established by Anderson⁶.

Our proposition can be tersely expressed as the following diagram:



Reinterpreting the existing algorithm

Training:

- Flow-matching: $\min_{\theta} \mathbb{E}_{t \sim \text{Unif}(0, T)} \mathbb{E}_{x \sim \rho(\cdot, t)} \|v_{\theta}(\cdot, t) - S(\cdot, t)\|_2^2$

Sampling/Inference: $\tau := -t$, $v_{\text{reverse}} := x - v_{\theta}$

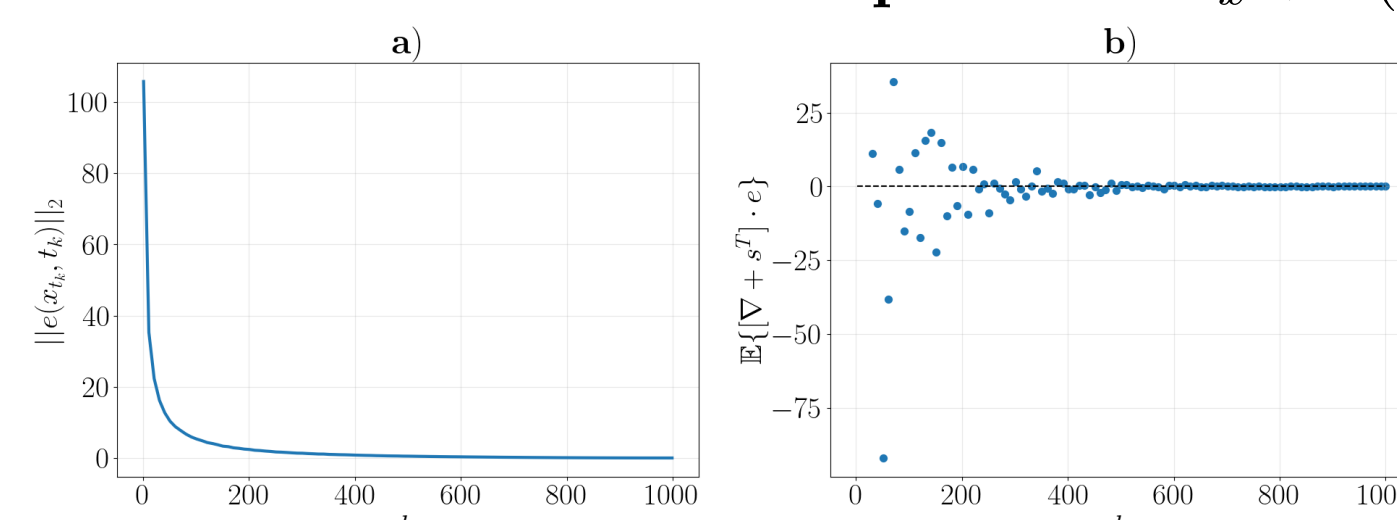
- Time-reversed Wasserstein Gradient Flow \equiv Probability Flow⁷
- WGF $\xrightarrow{\text{Reverse JKO/Otto}}$ Diffusion \equiv reversed-time SDE⁶

The advantage of the reinterpretation includes:

- Natural formulation and no need to explain the missing coefficients in DDPM (Cf. Sohl-Dickstein et al.⁹) and no need to reinvent the so-called “Probability Flow⁷”;
- Involving the existence of reversed-time SDE⁶ is not necessary: generative modeling only needs the **marginal distribution** (in contrast, reversed-time SDE⁶ ensures the **joint distribution** and the **path measure**);
- An error/null-kernel analysis could elucidate why existing models are robust even failing to learn the conservative $S(x, t)$, showing that generative modeling can be effective if

$$0 = [\nabla_x + S(x, t)] \cdot [\text{NN}_{\theta}(x, t) - S(x, t)].$$

A natural connection to the **Stein operator**⁸ $\nabla_x + S(x, t)$ emerged.



⁴ Chen et al. *Neural Ordinary Differential Equations*, 2019

⁵ Lipman et al. *Flow Matching for Generative Modeling*, 2022

⁶ Anderson, *Reverse-time diffusion equation models*, 1982

⁷ Song et al. *Score-based generative modeling through stochastic differential equations*, 2021

⁸ Liu & Wang, *Stein variational gradient descent*, 2016

⁹ Sohl-Dickstein et al. *Deep Unsupervised Learning using Non-Eq. Thermodynamics*, 2015