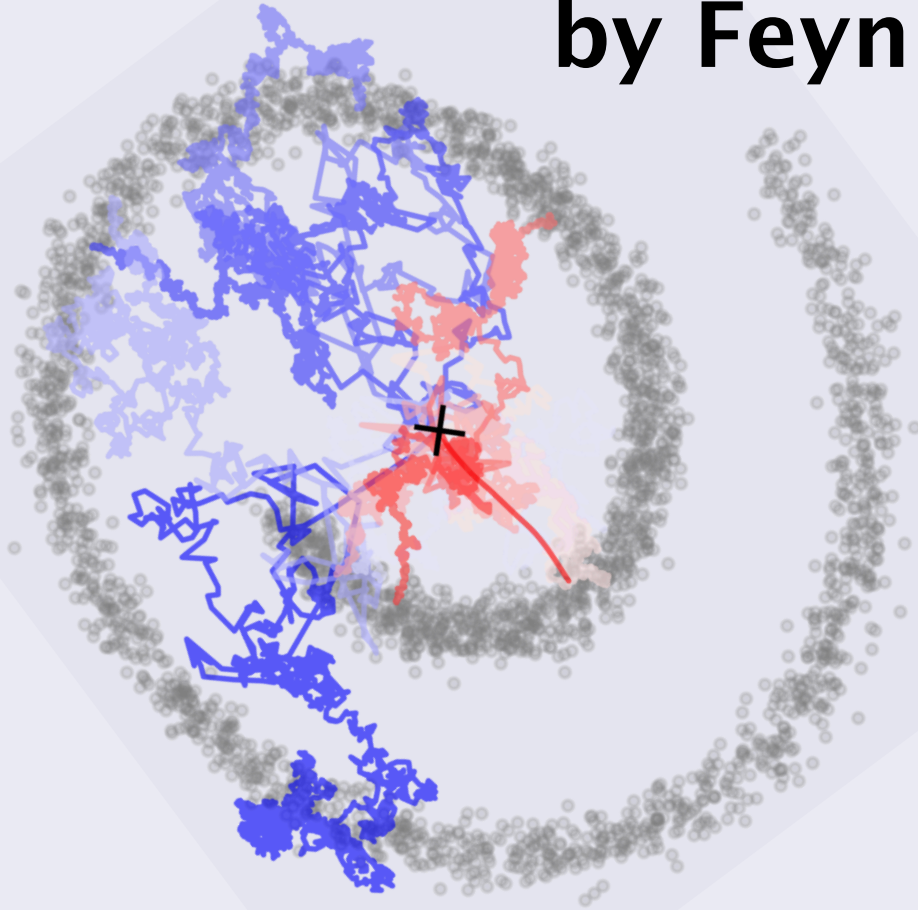# Understanding Diffusion Models by Feynman's Path Integral
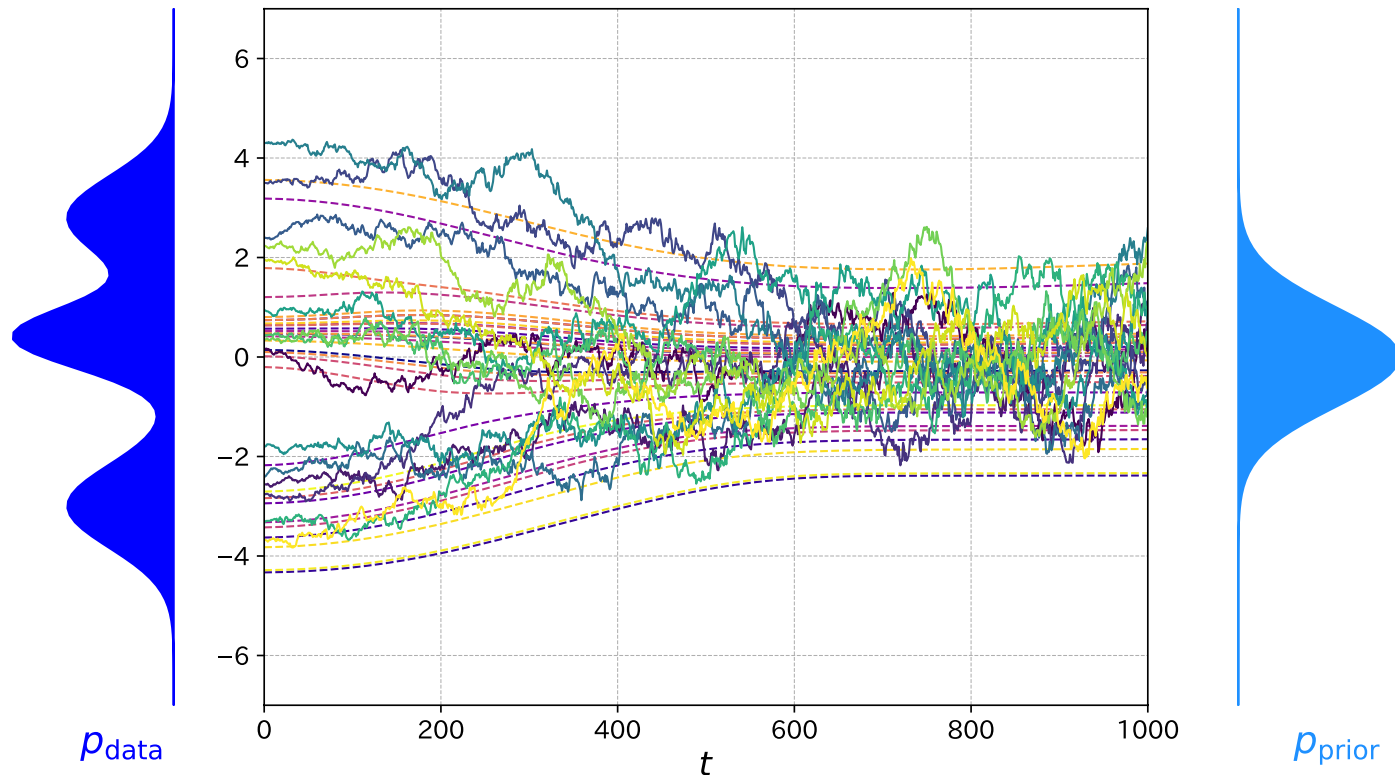
## Yuji Hirono (Univ. of Tsukuba)

**Collaborators: A. Tanaka (RIKEN), K. Fukushima(U. Tokyo)**

# Generative diffusion models

- **Diffusion models**: Generative AI for images, movies, and texts
  - Text-to-image generation, inpainting, hyper-resolution, Large Language Model, …
- We reformulate diffusion models via the path-integral method
  - Understanding ML models via physics methods

[Hirono–Tanaka–Fukushima, ICML'24]



$p_{\text{data}}$      $t$      $p_{\text{prior}}$

# Examples of generated images

"A family of lions in a cozy ramen shop"

"A set of sushi that look like dogs"

# Examples of generated movies

"There is a family of lions in a cozy ramen shop. They eat noodles with chopsticks"

# Examples of generated movies

"A family of **real majestic** lions in a cozy ramen shop. They use chopsticks to eat noodles. "

# Understanding diffusion models via path integral

- **How diffusion models work**

- **Path-integral formulation of diffusion models**
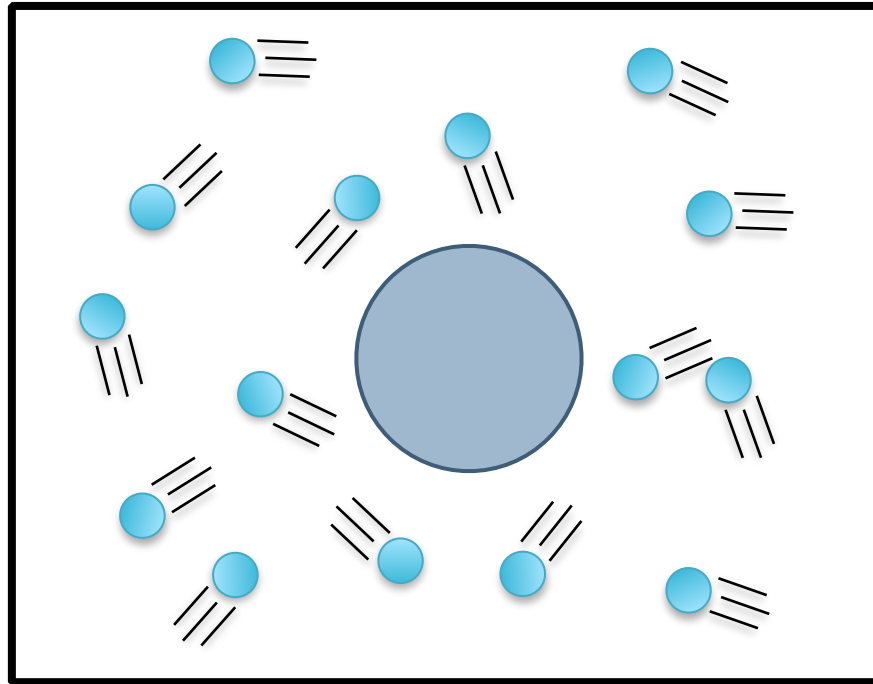
- **"Classical limit" and beyond**

# How diffusion models work

# Diffusion



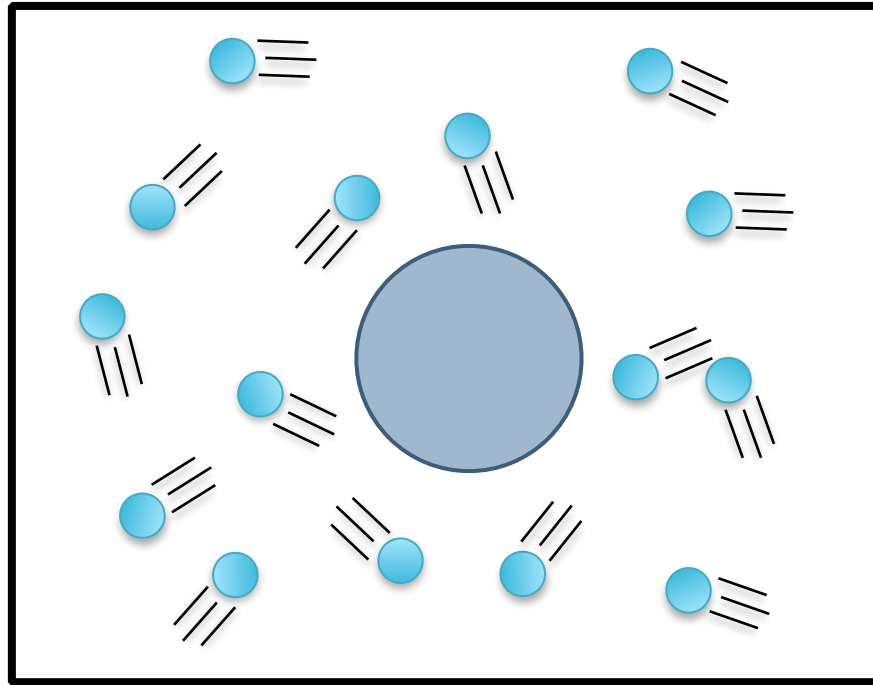https://www.youtube.com/watch?v=_Owb7Nbhhkg

# Brownian motion



- EOM of an ink molecule in a fluid

$$\mathbf{x}_{t+\mathrm{d}t} = \mathbf{x}_t + \underline{\mathrm{d}\mathbf{w}_t}$$

Random force

**Langevin equation**

# Brownian motion and discovery of atoms

Estimation of the Avogadro constant $N_{\mathrm{A}}$

$$\langle (\mathbf{x}_t - \mathbf{x}_0)^2 \rangle = 2Dt$$

Mean squared displacement

Diffusion constant

$$D = \frac{RT}{N_{\mathrm{A}}\,\gamma}$$

$$\gamma = 6\pi a\eta$$

$R$ : Gas constant
$T$ : Temperature

$\eta$ : Viscosity
$a$ : Radius

# Brownian motion and discovery of atoms

Time: 0.00 s

Estimation of the Avogadro constant $N_\mathrm{A}$

$$\langle (\mathbf{x}_t - \mathbf{x}_0)^2 \rangle = 2Dt$$

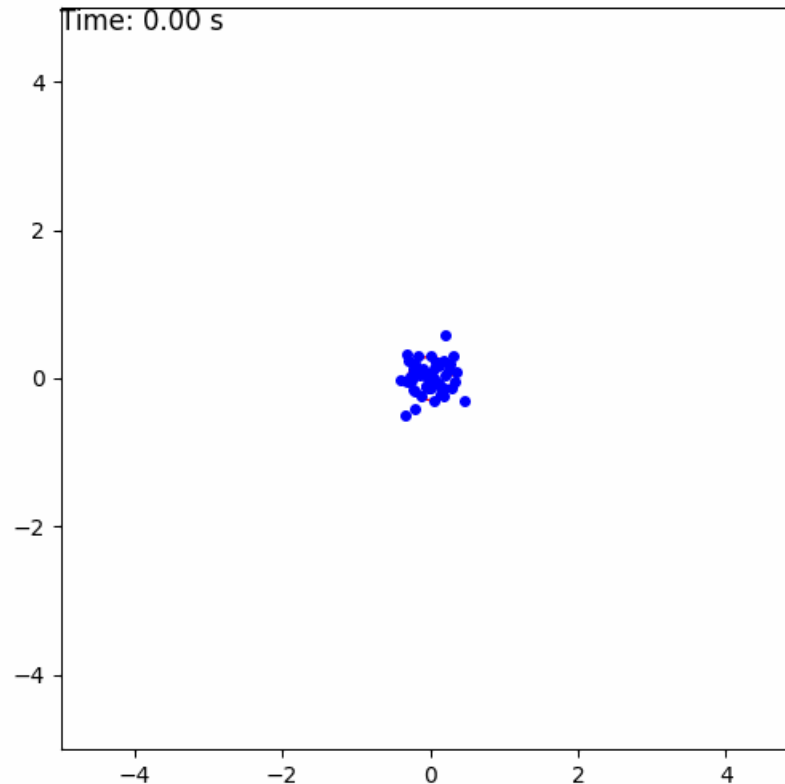Mean squared displacement

Diffusion constant

$$D = \frac{RT}{N_\mathrm{A}\,\gamma}$$

$$\gamma = 6\pi a\eta$$

$R$ : Gas constant
$T$ : Temperature

$\eta$ : Viscosity
$a$ : Radius

# Stochastic Differential Equations

- We consider the stochastic differential equation (SDE) of the form

$$\mathrm{d}\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}_t \qquad \text{where } \mathbf{w}_t \text{ is a Wiener process}$$

- This process is equivalent to the following Fokker–Planck equation

$$\partial_t p_t(\mathbf{x}) = -\nabla \cdot \left[ \mathbf{f}(\mathbf{x}, t)p_t(\mathbf{x}) - \frac{g^2(t)}{2}\nabla p_t(\mathbf{x}) \right]$$

# Diffusion models via SDEs [Song et. al., ICML'21]

Training data $\sim$ $p_{\text{data}}(\mathbf{x})$

- We'd like to sample from $p_{\text{data}}$

  - $p_{\text{data}}$ is unknown

  - Even if we know $p_{\text{data}}$, sampling via Markov Chain Monte-Carlo (MCMC) is inefficient

# Diffusion models via SDEs [Song et. al., ICML'21]

Forward: adding noise



Training data

$\sim$

$t = 0$

$p_{\text{data}}(\mathbf{x}) = p_0(\mathbf{x})$

$p_t$

$t = T$

$p_T = p_{\text{prior}}$

Backward: denoising for sampling

- For the forward process, one can employ, for example,

$$\partial_t p_t(\mathbf{x}) = -\nabla \cdot \left[ \mathbf{f}(\mathbf{x}, t) p_t(\mathbf{x}) - \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) \right] \quad \text{with} \quad \mathbf{f}(\mathbf{x}, t) = -\beta \mathbf{x}, \quad g(t) = \sqrt{\beta}$$
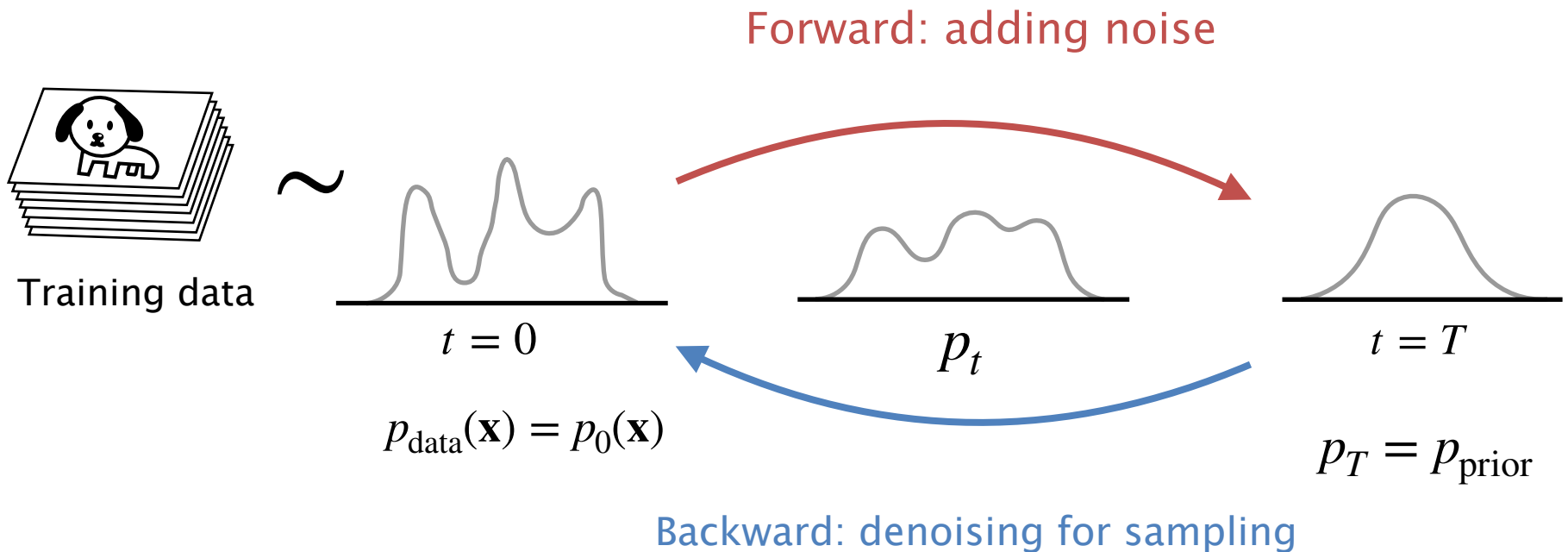
The stationary distribution is $\quad p_{\text{ss}}(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{1})$

# Diffusion models via SDEs <span>[Song et. al., ICML'21]</span>
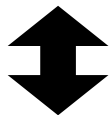
- The forward process is governed by

$$\partial_t p_t(\mathbf{x}) = -\nabla \cdot \left[ \mathbf{f}(\mathbf{x}, t) p_t(\mathbf{x}) - \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) \right] = \frac{g^2(t)}{2} \nabla^2 p_t(\mathbf{x}) + \cdots$$

- Since $p_{\text{data}}$ is unknown, this process is performed for samples by

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t) dt + g(t) d\mathbf{w}_t$$

- What is the SDE corresponding to the time−reversed FP equation?

$$\partial_t p_t(\mathbf{x}) = -\nabla \cdot \left[ \mathbf{f}(\mathbf{x}, t) p_t(\mathbf{x}) - \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) \ -\frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) + \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) \right]$$

$$= -\nabla \cdot \left[ \left( \mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla \log p_t(\mathbf{x}) \right) p_t(\mathbf{x}) + \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) \right]$$
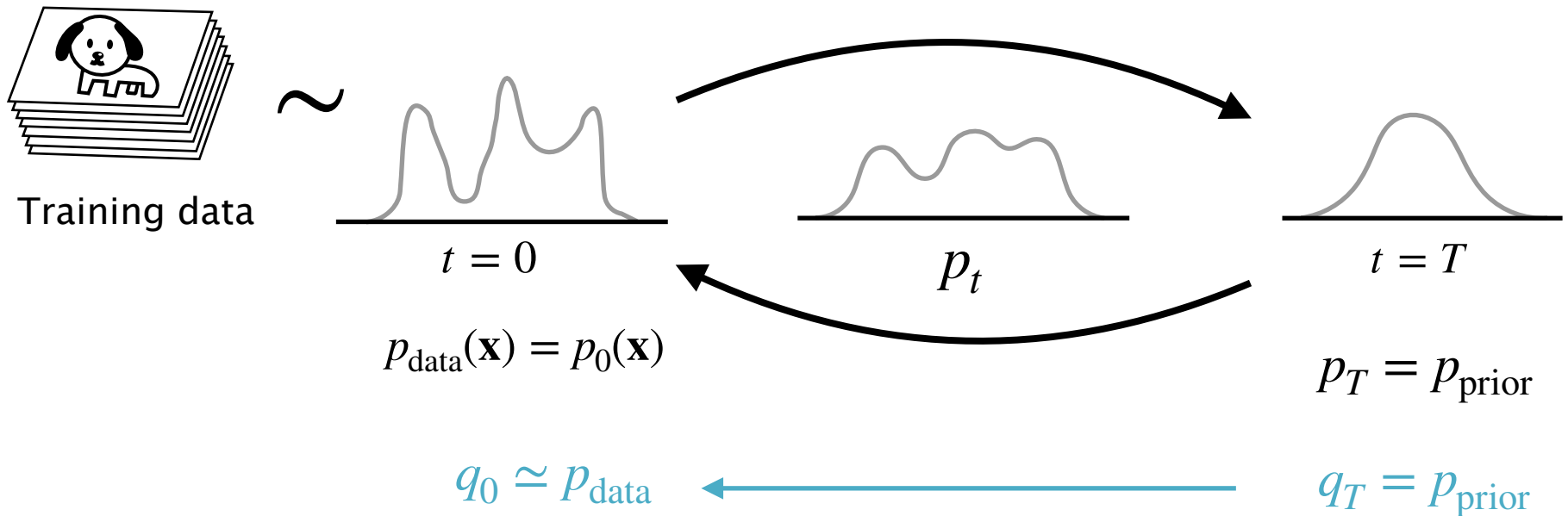
$$d\mathbf{x}_t = \left[ \mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \nabla \log p_t(\mathbf{x}_t) \right] dt + g(t) d\tilde{\mathbf{w}}_t$$

**SDE for the backward process**

<span>15</span>

# Diffusion models via SDEs [Song et. al., ICML'21]

Forward: noising $\qquad \mathrm{d}\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}_t$



~

Training data

$t = 0$

$p_{\mathrm{data}}(\mathbf{x}) = p_0(\mathbf{x})$

$p_t$

$t = T$

$p_T = p_{\mathrm{prior}}$

$q_0 \simeq p_{\mathrm{data}}$ $\qquad\longleftarrow\qquad$ $q_T = p_{\mathrm{prior}}$

Score function

Backward: sampling $\quad \mathrm{d}\mathbf{x}_t = \left[ \mathbf{f}(t, \mathbf{x}_t) - g(t)^2 \, \nabla \log p_t(\mathbf{x}_t) \right] \mathrm{d}t + g(t)\mathrm{d}\tilde{\mathbf{w}}_t$

$$\simeq \mathbf{s}_\theta(\mathbf{x}_t, t)$$

# Training objective

- $\mathbf{s}_\theta(\mathbf{x}, t)$ is learned to minimize the following loss function

$$\mathscr{L}(\theta) = \int_0^T \frac{g(t)^2}{2} \mathbb{E}_{\mathbf{x} \sim p_t} \left[ \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla \log p_t(\mathbf{x}_t) \right\|^2 \right] \mathrm{d}t$$

- This loss function gives the upper bound of $D_{\mathrm{KL}}(p_0 \| q_0)$,

$$D_{\mathrm{KL}}(p_0 \| q_0) \leq D_{\mathrm{KL}}(p_T \| q_T) + \mathscr{L}(\theta)$$

# Path-integral formulation of diffusion models

# Path integral in quantum mechanics

- Path integral: a formulation of quantum mechanics / QFTs

[Feynman, Rev. Mod. Phys. 20 (1948)]

- Expectation value of an observable $\mathscr{O}(\{\mathbf{x}_t\})$ is represented as

$$\langle \mathscr{O}(\mathbf{x}_t) \rangle = N \sum_{\text{paths}} e^{i\mathscr{A}[\{\mathbf{x}_t\}]/\hbar} \mathscr{O}(\mathbf{x}_t)$$

$\mathscr{A}[\{\mathbf{x}_t\}_{t \in [0,T]}]$: "Action"



- Classical mechanics: a path with least action $\delta\mathscr{A}[\{\mathbf{x}_t\}] = 0$

# Path integral formulation of diffusion models

- The expectation value of $\mathcal{O}(\{\mathbf{x}_t\})$ where $\mathbf{x}_t$ obeys the SDE

$$\mathrm{d}\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}_t$$

can be represented in the path-integral form as

$$\mathbb{E}\left[\mathcal{O}(\{\mathbf{x}_t\})\right] = \int [D\mathbf{x}_t]\,\mathcal{O}(\{\mathbf{x}_t\})\underline{p_0(\mathbf{x}_0)e^{-\mathscr{A}}}$$

$$=: P(\{\mathbf{x}_t\}_{t\in[0,T]}) \quad \text{Path probability}$$

"Action" $\qquad \mathscr{A} := \displaystyle\int_0^T \frac{1}{2g(t)^2}\,\left\| \dot{\mathbf{x}}_t - \mathbf{f}(\mathbf{x}_t, t) \right\|^2 \mathrm{d}t$

Known as Onsager–Machlup function

[Onsager–Machlup, Phys. Rev., 1953]

# Backward process in path integral

- Backward SDE:    $d\mathbf{x}_t = \tilde{\mathbf{f}}(\mathbf{x}, t)dt + g(t)d\tilde{\mathbf{w}}_t$

  where    $\tilde{\mathbf{f}}(\mathbf{x}, t) := \mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla \log p_t(\mathbf{x})$

- Backward SDE can be derived naturally in path integral:

  The path probability can be written as

  $$P(\{\mathbf{x}_t\}_{t \in [0,T]}) = p_0(\mathbf{x}_0)\, e^{-\mathscr{A}} = e^{-\tilde{\mathscr{A}}} p_T(\mathbf{x}_T)$$

  Action for backward process    $\tilde{\mathscr{A}} = \int_0^T \frac{1}{2g(t)^2} \left\| \dot{\mathbf{x}}_t - \tilde{\mathbf{f}}(\mathbf{x}_t, t) \right\|^2 dt$

# "Classical limit" and beyond

# Sampling processes of diffusion models

- Stochastic: $\mathrm{d}\mathbf{x}_t = \left[ \mathbf{f}(\mathbf{x}, t) - g(t)^2 \mathbf{s}_\theta(\mathbf{x}, t) \right] \mathrm{d}t + g(t)\mathrm{d}\tilde{\mathbf{w}}_t$

- Deterministic: **Probability Flow (PF) ODE**    [Song et. al., ICML'21]

$$\mathrm{d}\mathbf{x}_t = \left[ \mathbf{f}(\mathbf{x}, t) - \frac{1}{2} g(t)^2 \mathbf{s}_\theta(\mathbf{x}, t) \right] \mathrm{d}t$$

- These sampling methods are equivalent if the learned score is perfect, $\mathbf{s}_\theta(\mathbf{x}, t) = \nabla \log p_t(\mathbf{x})$.

# Sampling processes of diffusion models

- Introducing a parameter $\hbar$, FP equation can be written as

$$\partial_t p_t(\mathbf{x}) = -\nabla \cdot \left[ \mathbf{f}(\mathbf{x}, t) p_t(\mathbf{x}) - \frac{g^2(t)}{2} \nabla p_t(\mathbf{x}) - \frac{\hbar g^2(t)}{2} \nabla p_t(\mathbf{x}) + \frac{\hbar g^2(t)}{2} \nabla p_t(\mathbf{x}) \right]$$

$$= -\nabla \cdot \left[ \left( \mathbf{f}(\mathbf{x}, t) - \frac{1 + \hbar}{2} g(t)^2 \nabla \log p_t(\mathbf{x}) \right) p_t(\mathbf{x}) + \frac{\hbar g^2(t)}{2} \nabla p_t(\mathbf{x}) \right]$$

$$\mathrm{d}\mathbf{x}_t = \left[ \mathbf{f}(\mathbf{x}_t, t) - \frac{1 + \hbar}{2} g(t)^2 \nabla \log p_t(\mathbf{x}) \right] \mathrm{d}t + \sqrt{\hbar} g(t) \mathrm{d}\tilde{\mathbf{w}}_t$$

$\hbar = 1$     Stochastic sampling

$\hbar = 0$     Deterministic sampling (PF–ODE)

# Sampling processes of diffusion models

- Stochastic:    $d\mathbf{x}_t = \left[ \mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \mathbf{s}_\theta(\mathbf{x}_t, t) \right] dt + g(t) d\tilde{\mathbf{w}}_t$

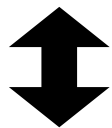- Deterministic: **Probability Flow (PF) ODE**    [Song et. al., ICML'21]

$$dx_t = \left[ \mathbf{f}(\mathbf{x}_t, t) - \frac{1}{2} g(t)^2 \mathbf{s}_\theta(\mathbf{x}_t, t) \right] dt$$

- Pro: one–to–one correspondence:

  - Faster sampling

  - The log–likelihood of an image $\mathbf{x}_0$ can be evaluated by

$$\log p_0(\mathbf{x}_0) = \log p_T(\mathbf{x}_T) + \int_0^T \nabla \cdot \left[ \mathbf{f}(\mathbf{x}_t, t) - \frac{1}{2} g(t)^2 \mathbf{s}_\theta(\mathbf{x}_t, t) \right] dt$$

    using the solution $\{\mathbf{x}_t\}_{t \in [0, T]}$ of PF ODE

- Con: worse sample quality

  When the score is imperfect, the two sampling processes are inequivalent

# "Classical limit" and beyond

- Stochastic & deterministic samplings are continuously interpolated:

$$\mathrm{d}\mathbf{x}_t = \tilde{\mathbf{F}}_{\theta,\mathfrak{h}}(\mathbf{x}_t, t)\mathrm{d}t + \sqrt{\mathfrak{h}}\, g(t)\mathrm{d}\tilde{\mathbf{w}}_t \qquad \mathfrak{h} = 1 \quad \text{Stochastic}$$

$$\tilde{\mathbf{F}}_{\theta,\mathfrak{h}}(\mathbf{x}, t) := \mathbf{f}(\mathbf{x}, t) - \frac{1+\mathfrak{h}}{2}g(t)^2 \mathbf{s}_\theta(\mathbf{x}, t) \qquad \mathfrak{h} = 0 \quad \text{Deterministic}$$

- Path probability of a model is

$$Q_{\mathfrak{h}}(\{\mathbf{x}_t\}_{t\in[0,T]}) = e^{-\frac{1}{\mathfrak{h}}\tilde{\mathscr{A}}_{\theta,\mathfrak{h}}} q_T(\mathbf{x}_T) \quad \longleftrightarrow \quad e^{i\frac{1}{\hbar}\mathscr{A}} \quad \text{in QM}$$

$$\tilde{\mathscr{A}}_{\theta,\mathfrak{h}} := \int_0^T \frac{1}{2g(t)^2} \left\| \dot{\mathbf{x}}_t - \tilde{\mathbf{F}}_{\theta,\mathfrak{h}}(\mathbf{x}, t) \right\|^2 \mathrm{d}t$$

- **Parameter $\mathfrak{h}$ is the counterpart of Planck's constant $\hbar$**

# "Classical limit" and beyond

- Deterministic sampling appears as the "classical limit" $\hbar \to 0$

$$Q_\hbar(\{\mathbf{x}_t\}_{t\in[0,T]}) = e^{-\frac{1}{\hbar}\tilde{\mathscr{A}}_{\theta,\hbar}} q_T(\mathbf{x}_T) \qquad \tilde{\mathscr{A}}_{\theta,\hbar} := \int_0^T \frac{1}{2g(t)^2} \left\| \dot{\mathbf{x}}_t - \tilde{\mathbf{F}}_{\theta,\hbar}(\mathbf{x},t) \right\|^2 dt$$

$$\stackrel{\hbar \to 0}{\to} \prod_{t\in[0,T]} \delta\left(\dot{\mathbf{x}}_t - \tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,t)\right) q_T(\mathbf{x}_T) \qquad \tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x},t) := \mathbf{f}(\mathbf{x},t) - \frac{1}{2}g(t)^2 \mathbf{s}_\theta(\mathbf{x},t)$$

- Likelihood computation for $\hbar \neq 0$ via **WKB expansion**

  - To the first order in $\hbar$,

$$\log q_0^\hbar(\mathbf{x}_0) = \log q_0^{\hbar=0}(\mathbf{x}_0) + \hbar\left[\delta\mathbf{x}_T \cdot \nabla \log q_T^{\hbar=0}(\mathbf{x}_T) + \int_0^T \nabla \cdot \delta\tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,\delta\mathbf{x}_t,t)dt\right]$$

First–order correction to log–likelihood

where $\{\mathbf{x}_t, \delta\mathbf{x}_t\}$ are the solution of the following ODE with $\mathbf{x}_{t=0} = \mathbf{x}_0$, $\delta\mathbf{x}_{t=0} = \mathbf{0}$

$$\begin{cases} \dot{\mathbf{x}}_t = \tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,t) \\ \dot{\delta\mathbf{x}}_t = \delta\tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,\delta\mathbf{x}_t,t) \end{cases} \qquad \delta\tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,\delta\mathbf{x}_t,t) := (\delta\mathbf{x}_t \cdot \nabla)\tilde{\mathbf{f}}_\theta^{\mathrm{PF}}(\mathbf{x}_t,t) - \frac{\mathfrak{g}(t)^2}{2}[\mathbf{s}_\theta(\mathbf{x}_t,t) - \nabla \log q_t^{\hbar=0}(\mathbf{x}_t)]$$

# Noise improves sample quality

**SWISS-ROLL**

| SDE (NLL) | tol | 1ST-CORR | ERRORS |
|---|---|---|---|
| SIMPLE $(1.39\pm 0.05)$ | 1e-3 1e-5 | -0.31$\pm$0.21 -0.44$\pm$0.38 | 0.13$\pm$0.00 0.13$\pm$0.00 |
| COSINE $(1.42\pm 0.02)$ | 1e-3 1e-5 | -1.59$\pm$0.57 -3.27$\pm$1.11 | 0.35$\pm$0.00 0.37$\pm$0.02 |

**25-GAUSSIAN**

| SDE (NLL) | tol | 1ST-CORR | ERRORS |
|---|---|---|---|
| SIMPLE $(-1.22\pm 0.01)$ | 1e-3 1e-5 | -3.64$\pm$0.49 -3.61$\pm$0.64 | 0.31$\pm$0.00 0.32$\pm$0.01 |
| COSINE $(-1.71\pm 0.02)$ | 1e-3 1e-5 | -17.57$\pm$5.56 -19.65$\pm$17.46 | 0.70$\pm$0.01 0.67$\pm$0.03 |

$O(\mathfrak{h}^1)$ correction to Negative Log Likelihood (NLL) $\mathbb{E}\left[-\log q_0\right]$

- Negative correction $\rightarrow$ noise improves sample quality

# Path integral formulation of diffusion models

- Useful for physicists in understanding various aspects of diffusion models: backward process, training objective

- Deterministic sampling by PF ODE appears as "classical limit"

$$Q_{\hbar}(\{\mathbf{x}_t\}_{t\in[0,T]}) = e^{-\frac{1}{\hbar}\tilde{\mathscr{A}}_{\theta,\hbar}} \, q_T(\mathbf{x}_T) \quad \longleftrightarrow \quad e^{i\frac{1}{\hbar}\mathscr{A}} \quad \text{in QM}$$

- Likelihood evaluation via WKB expansion

- Physics methods for analyzing generative AI models

# Quantum × AI in the next 100 years?

- Discrete diffusion models are used to build Large Language models

  - Mercury Coder (Feb. 2025 – )
    Gemini Diffusion (May. 2025 –)

- Based on "classical" stochastic dynamics

- A possible new mechanism of LLMs based on quantum dynamics?