

113-2 Embodied Vision

Diffusion Models and Flow Matching

Jia-Wei Liao

Ph.D. Candidate in Computer Science

National Taiwan University



Learning Resources for Diffusion Models

- [【漫士科普】人工智慧博士生告訴你 SORA 擴散模型究竟是怎麼產生影片的？](#)
- [Hung-Yi Lee YouTube](#)
- [Jia-Bin Huang YouTube](#)
- [Diffusion and Score-Based Generative Models \(Yang Song\)](#)
- [Evolution of Diffusion Models: From Birth to Enhanced Efficiency and Controllability \(Jesse\)](#)
- [Lil'Log What are Diffusion Models?](#)
- [生成擴散模型漫談 \(蘇劍林\)](#)

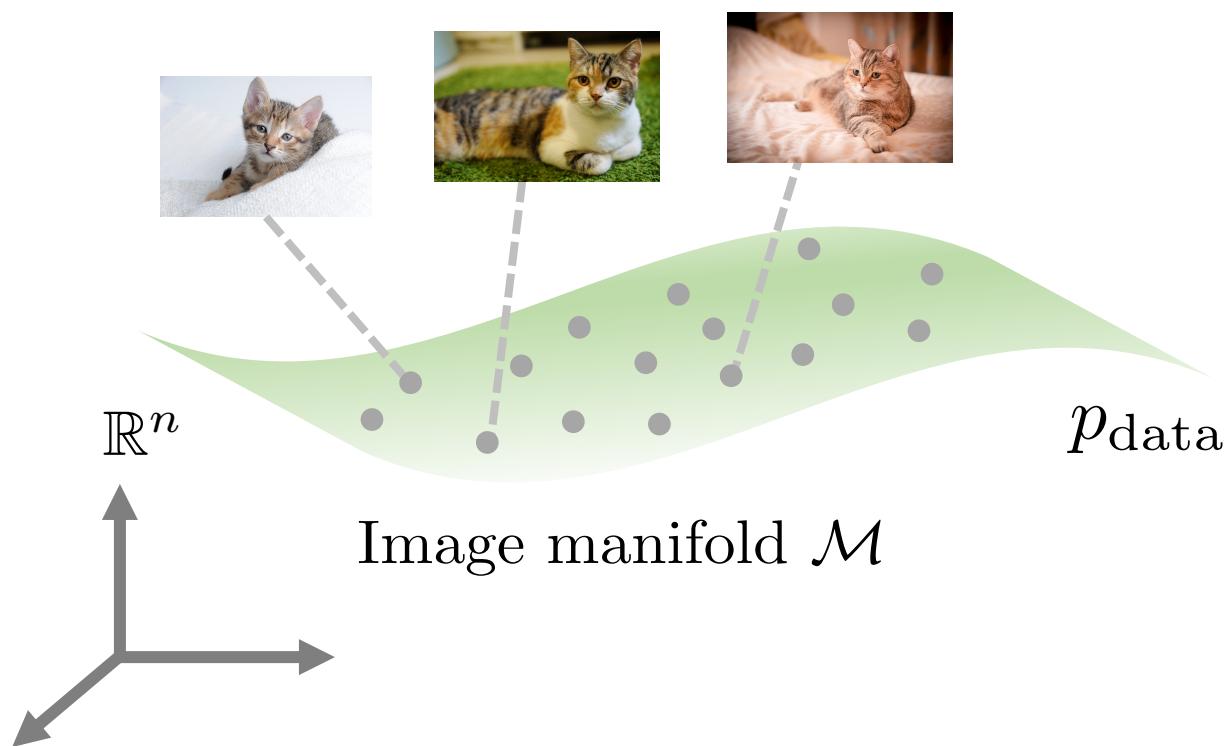
Outline

- Theoretical Foundations (ODE/SDE, Score Function)
- Training and Sampling
- Advanced Variants (DDIM, CM, CTM, Flow Matching, Rectified Flow)
- Controllability and Applications (CG, CFG, ControlNet)

What is Generative Model Learning?

Data Manifold Assumption

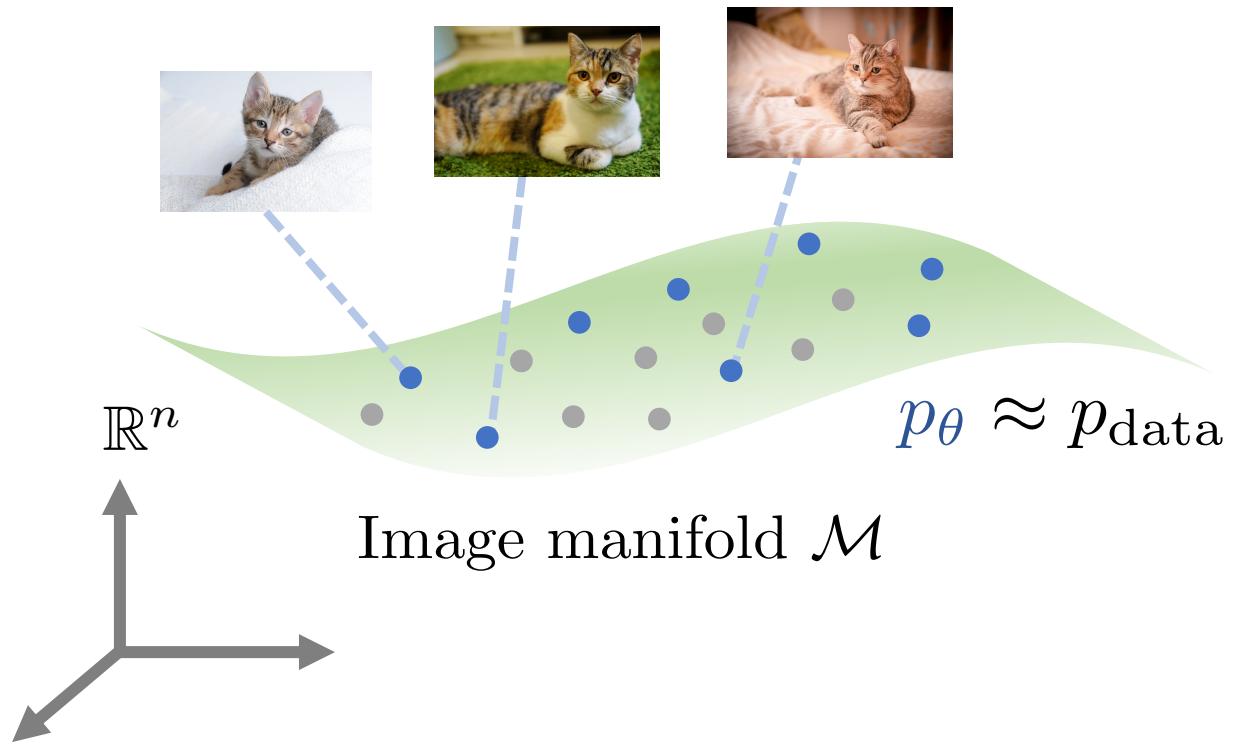
Natural high-dimensional data concentrate close to a non-linear low-dimensional manifold



What is Generative Model Learning?

Data Manifold Assumption

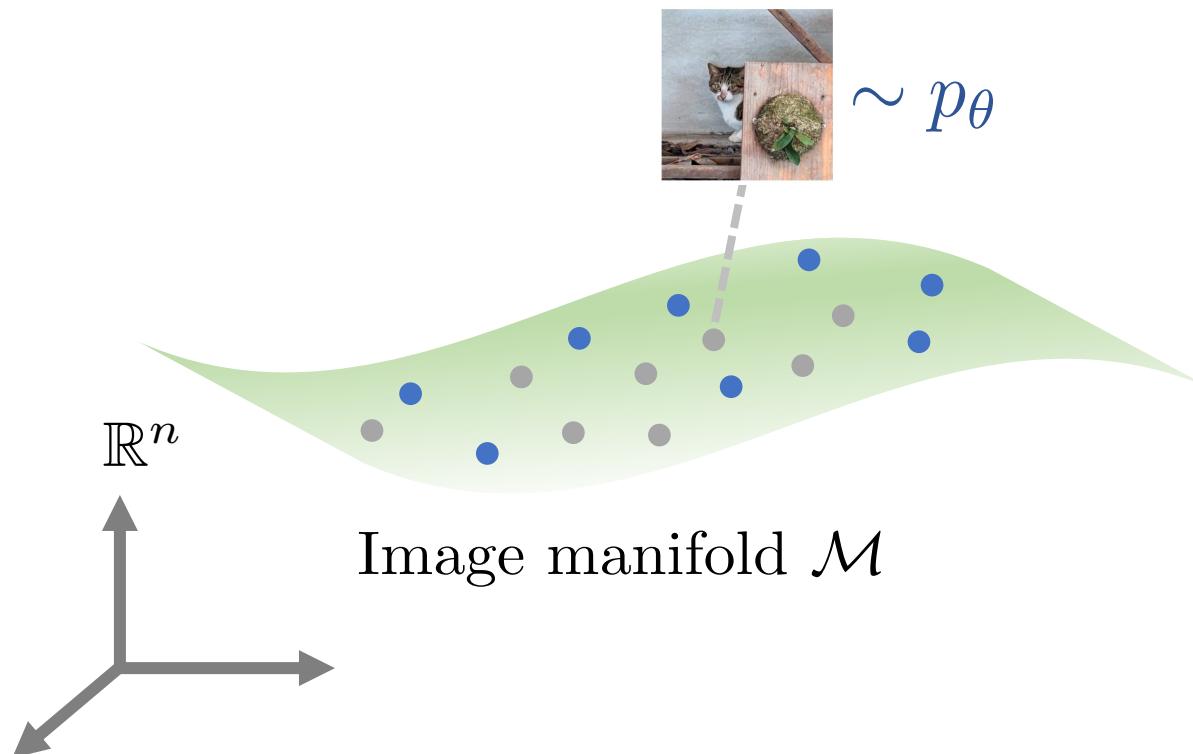
Natural high-dimensional data concentrate close to a non-linear low-dimensional manifold



What is Generative Model Learning?

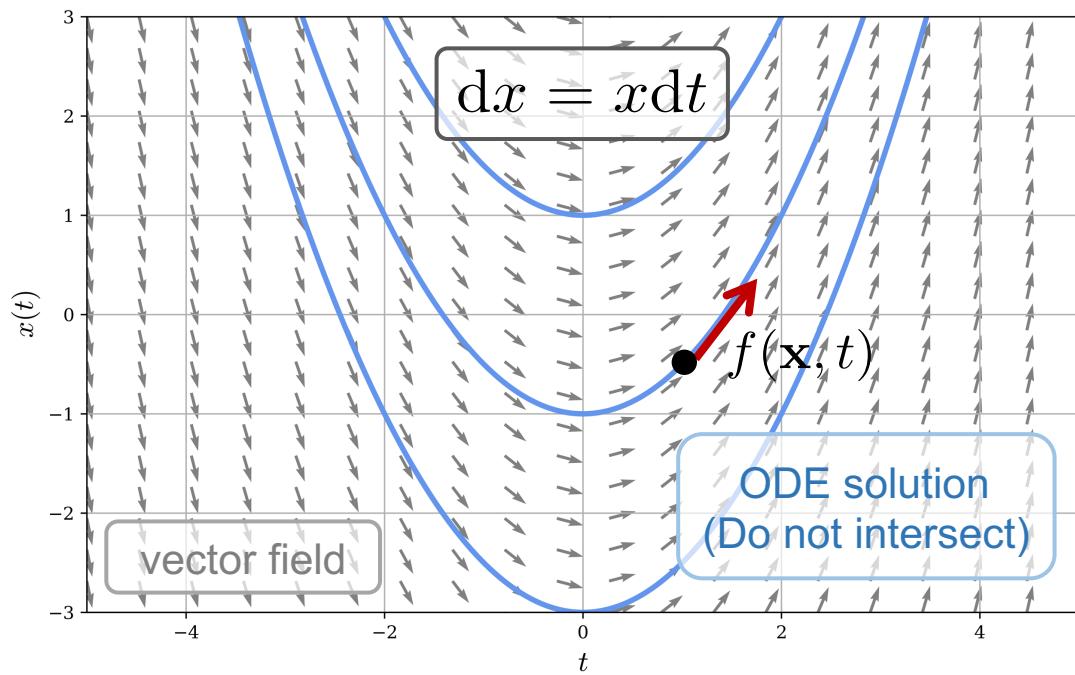
Data Manifold Assumption

Natural high-dimensional data concentrate close to a non-linear low-dimensional manifold

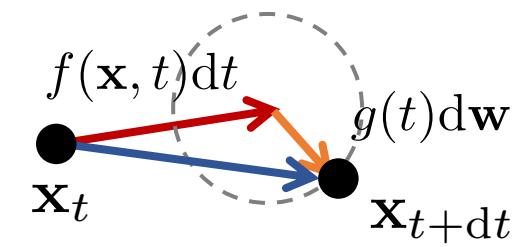


Recap ODE and SDE

$$d\mathbf{x} = f(\mathbf{x}, t)dt \quad \text{or} \quad \frac{d\mathbf{x}}{dt} = f(\mathbf{x}, t)$$

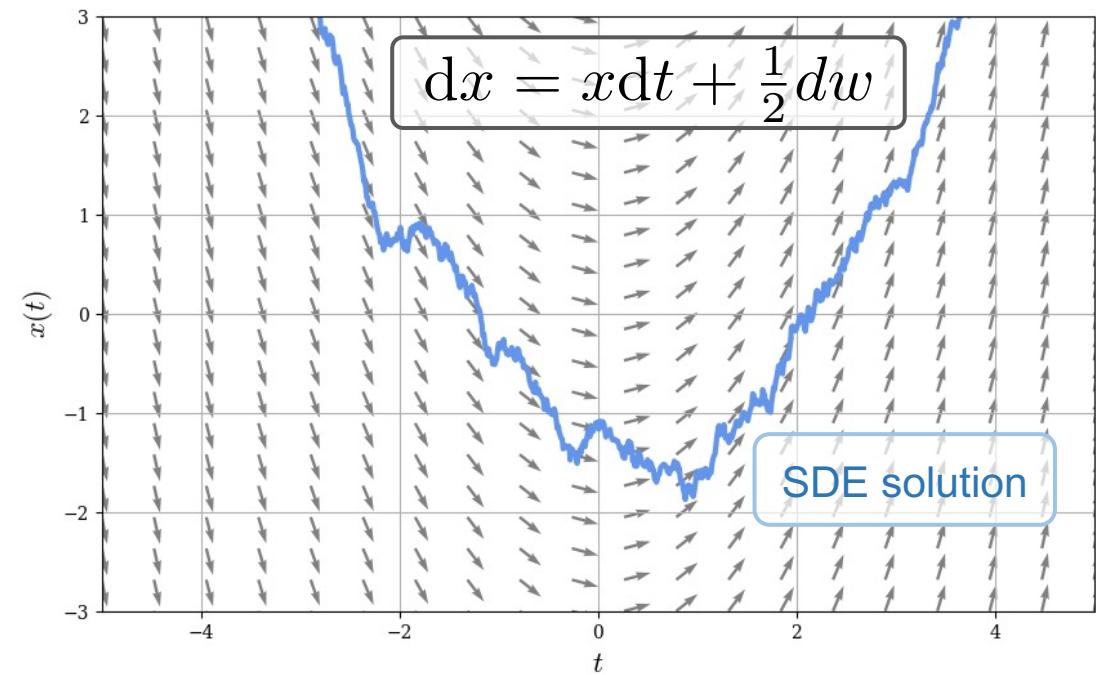


The arrow indicates the direction of movement at the moment of downward shift

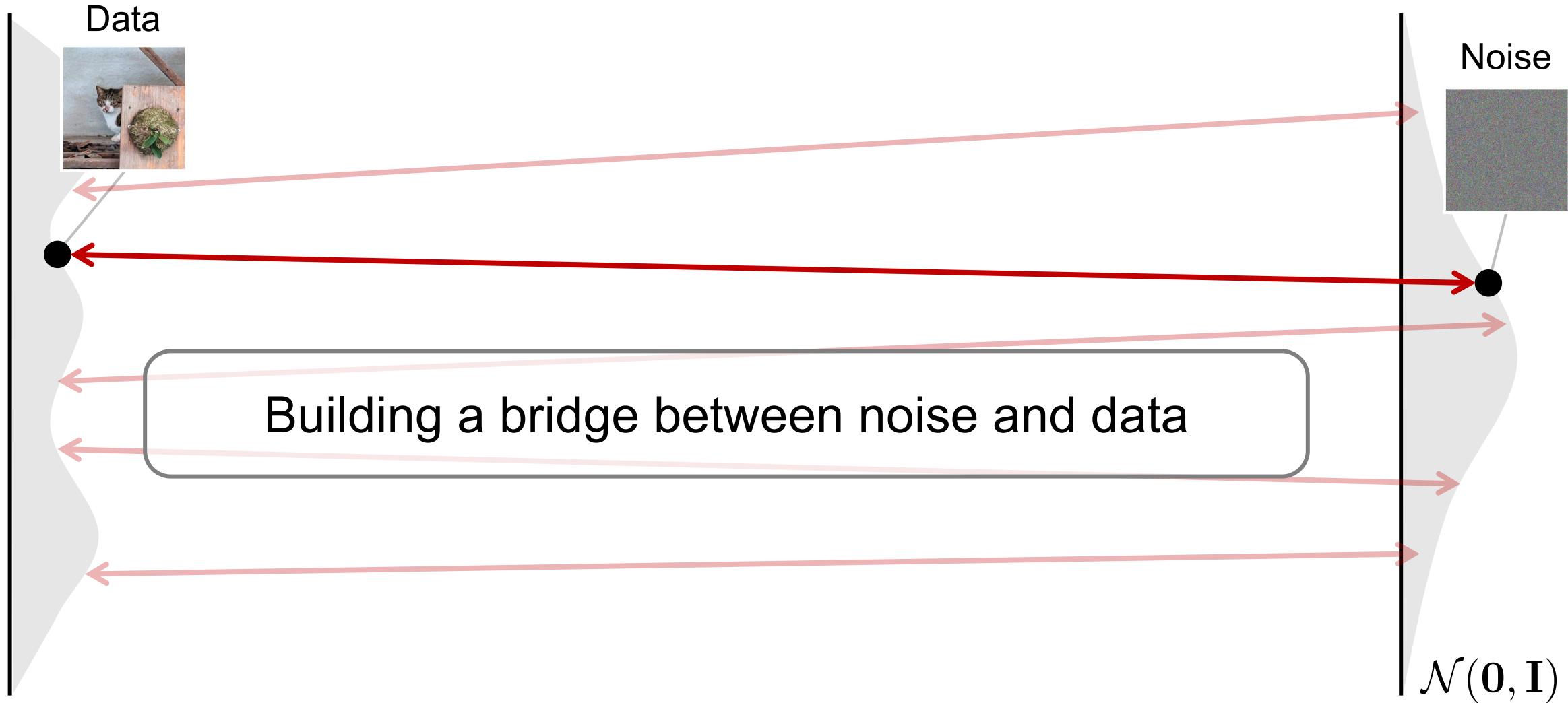


$$d\mathbf{x} = f(\mathbf{x}, t)dt + g(t)d\mathbf{w}$$

drift diffusion



The Goal of Diffusion Models



What is Diffusion Model?

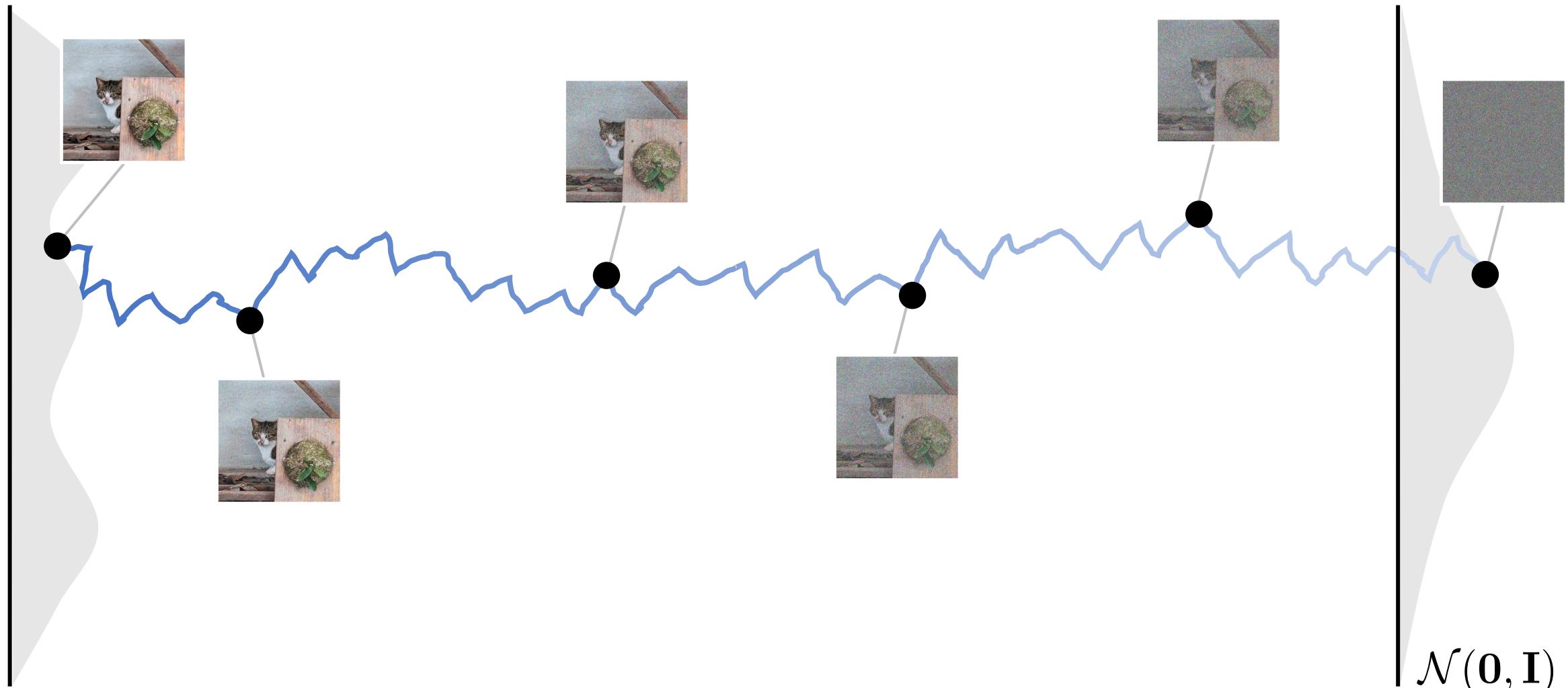
Forward Process: add noise step by step, from data to pure noise



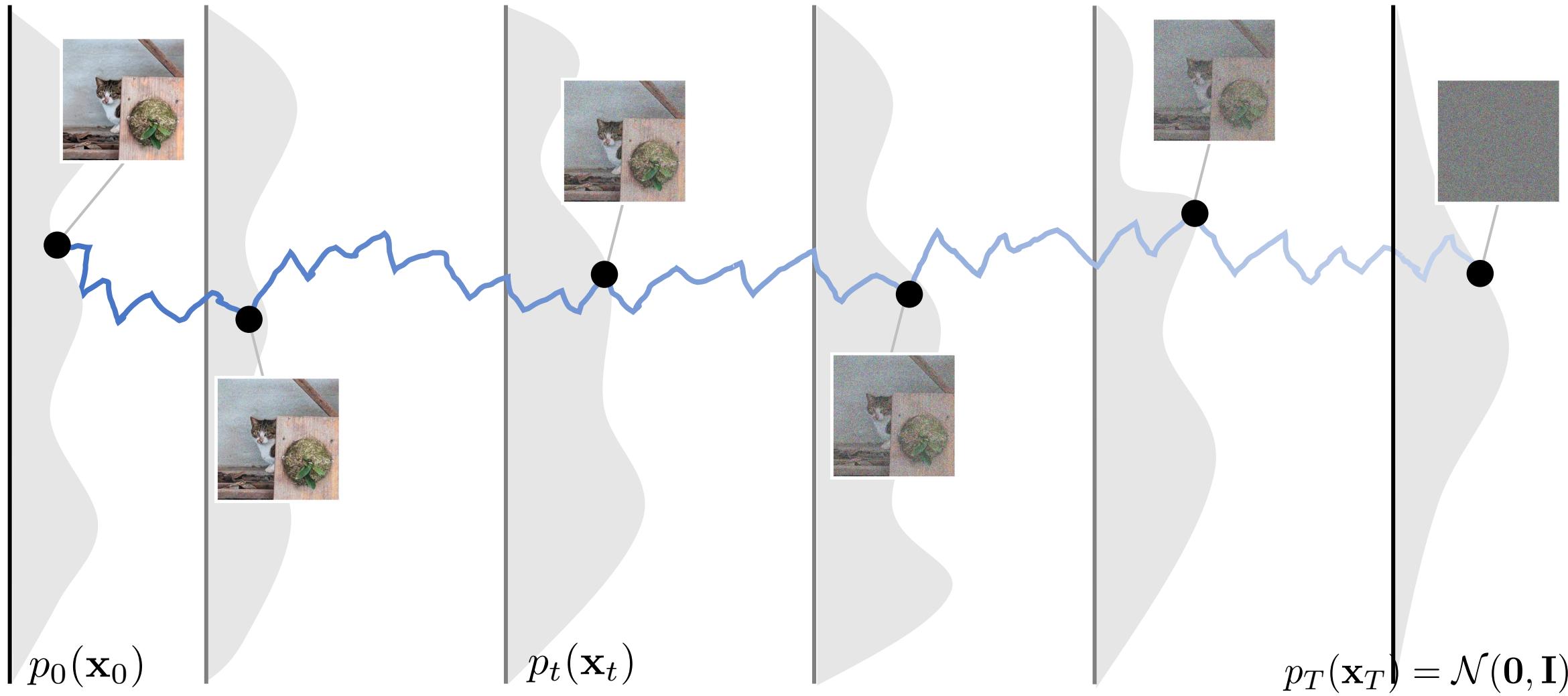
Reverse Process: generate data from pure noise by denoising

Creating noise from data is easy; creating data from noise is generative modeling – Yang Song

Diffusion Models



Diffusion Models



Diffusion Models

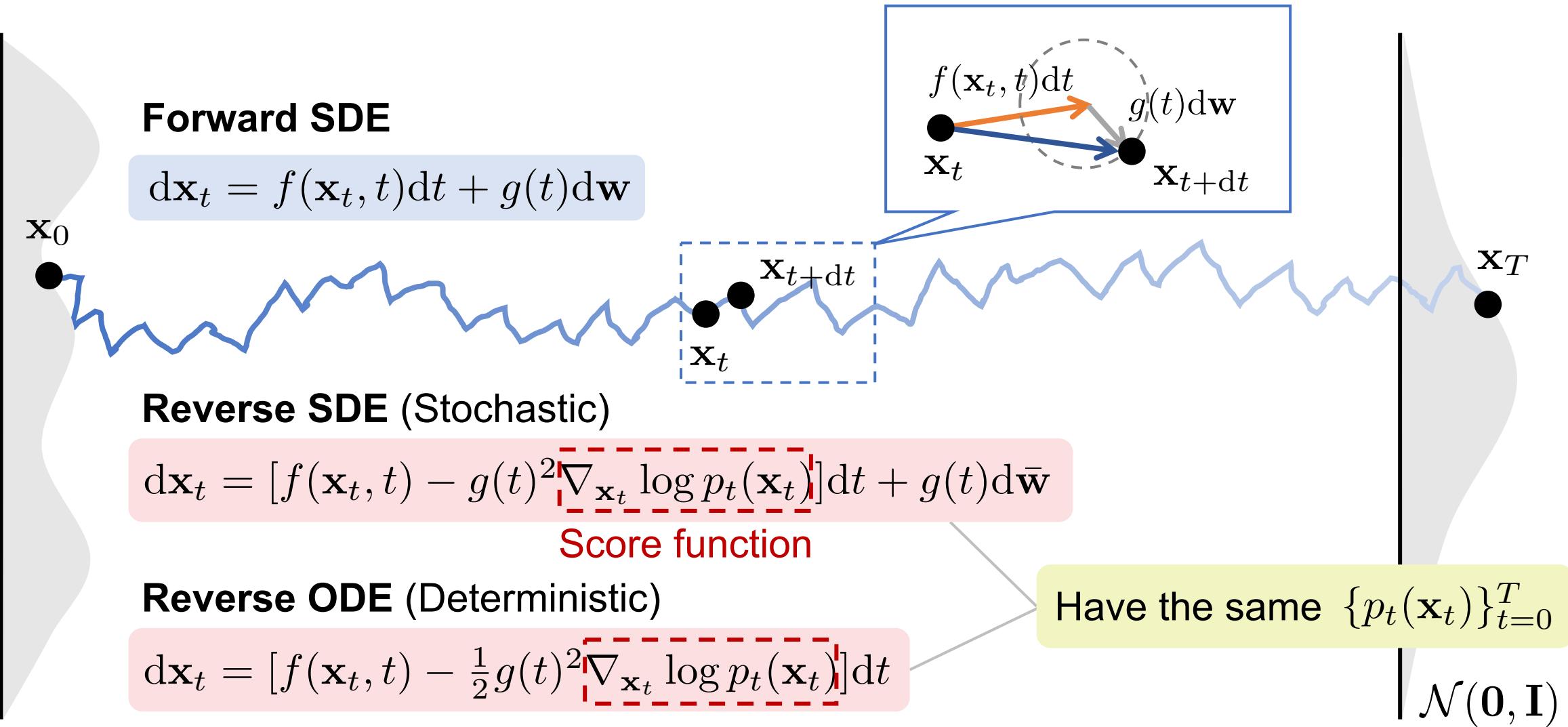
How to illustrate this process?

$$p_0(\mathbf{x}_0)$$

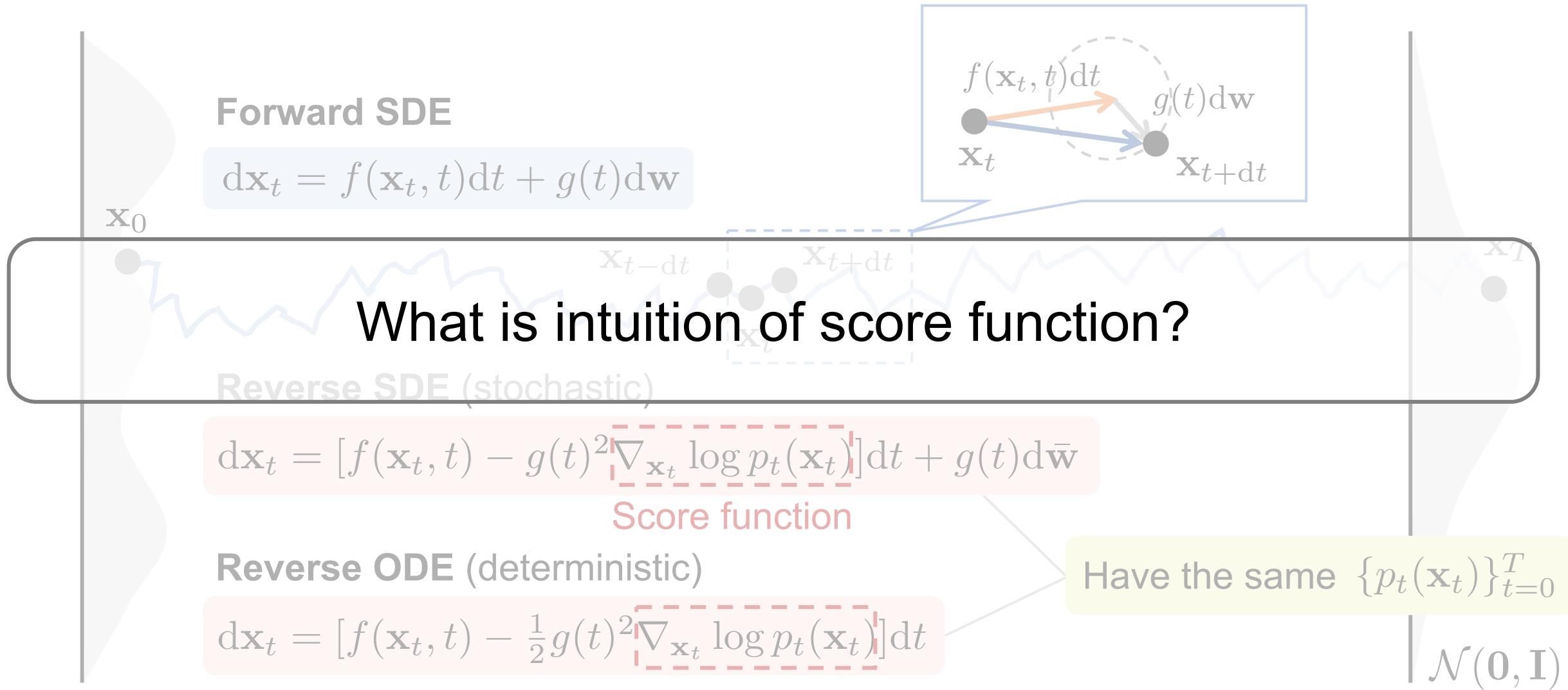
$$p_t(\mathbf{x}_t)$$

$$p_T(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Score-based Diffusion Models [Song+ ICLR'21]

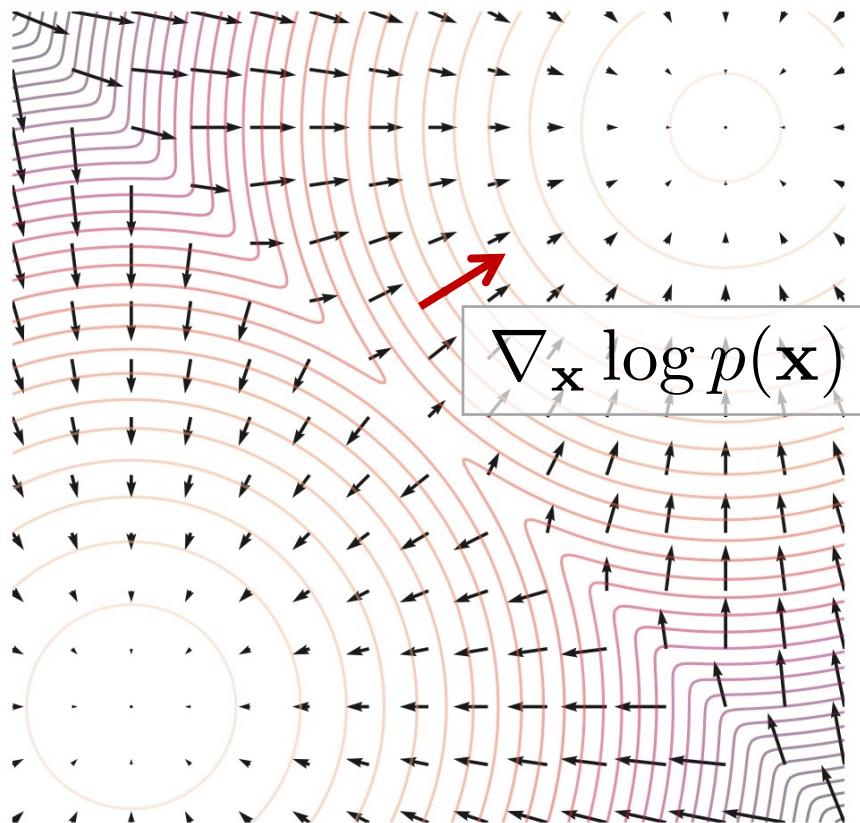


Score-based Diffusion Models [Song+ ICLR'21]



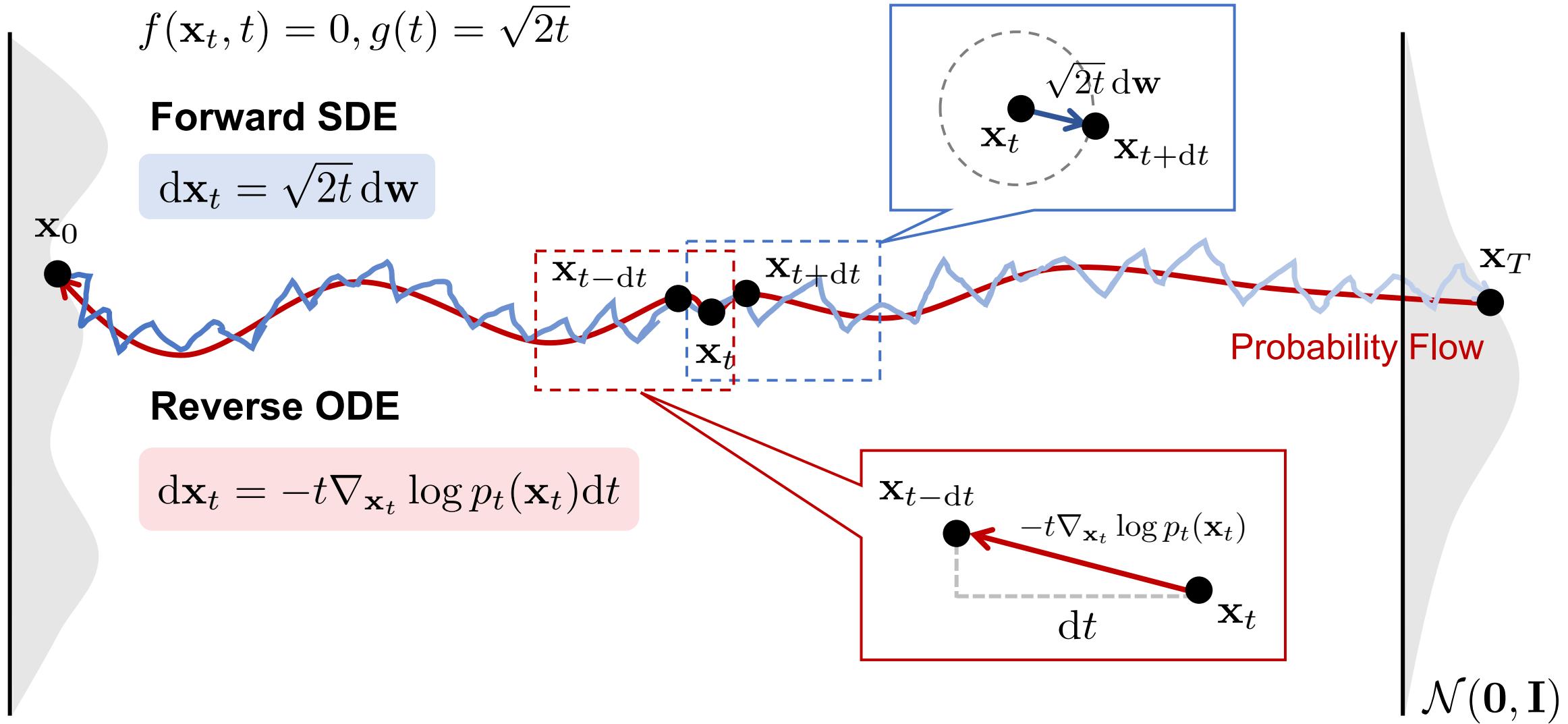
Score Function

Probability Density Function $p(\mathbf{x})$



Score Function: The direction in which the probability density increases most rapidly

Score-based Diffusion Models [Song+ ICLR'21]

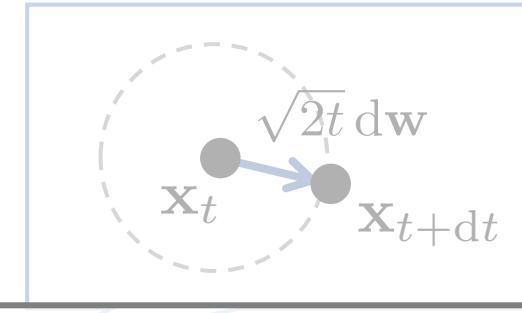


Score-based Diffusion Models [Song+ ICLR'21]

$$f(\mathbf{x}_t, t) = 0, g(t) = \sqrt{2t}$$

Forward SDE

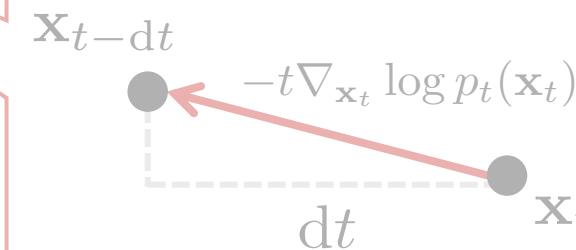
$$d\mathbf{x}_t = \sqrt{2t} dw$$



Can we compute $\mathbf{x}_0 = \mathbf{x}_T + \int_0^T -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$?

Very difficulty!

$$d\mathbf{x}_t = -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$$



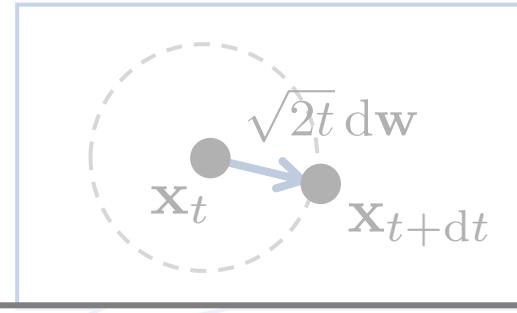
$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

Score-based Diffusion Models [Song+ ICLR'21]

$$f(\mathbf{x}_t, t) = 0, g(t) = \sqrt{2t}$$

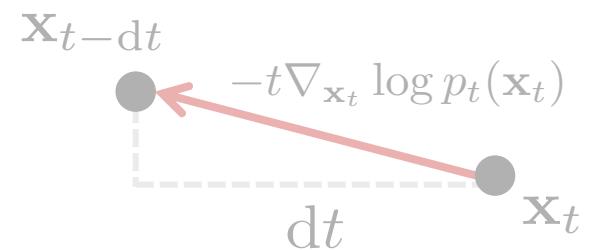
Forward SDE

$$d\mathbf{x}_t = \sqrt{2t} dw$$



1. How do we obtain the score function during sampling?
2. How to solve the reverse ODE?

$$d\mathbf{x}_t = -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$$



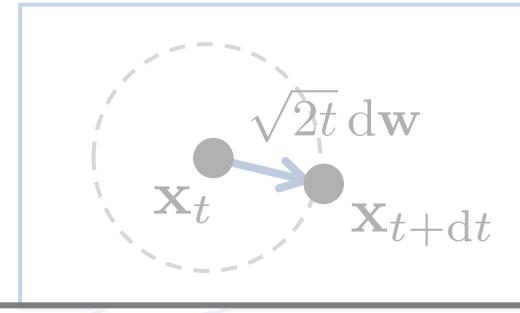
$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

Score-based Diffusion Models [Song+ ICLR'21]

$$f(\mathbf{x}_t, t) = 0, g(t) = \sqrt{2t}$$

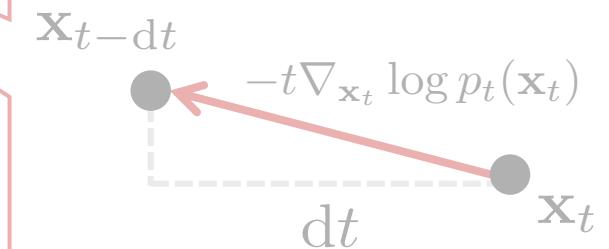
Forward SDE

$$d\mathbf{x}_t = \sqrt{2t} dw$$



1. How do we obtain the score function during sampling?
2. How to solve the reverse ODE?

$$d\mathbf{x}_t = -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$$



$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

Training Score-based Diffusion Models

We train the model to approximate the score function

$$\mathbf{s}_\theta(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$$

During sampling, we use the model to predict the score and plug it into SDE/ODE

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

The training objective can be expressed as

$$\mathcal{L}_{\text{SM}}(\mathbf{x}_0; \theta) = \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \quad (\text{Score Matching})$$

Training Score-based Diffusion Models

We train the model to approximate the score function

$$\mathbf{s}_\theta(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$$

During sampling, we use the model to predict the score and plug it into SDE (ODE)

Can we easily compute the score function during training?

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

The training objective can be expressed as

$$\mathcal{L}_{SM}(\mathbf{x}_0; \theta) = \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \text{ (Score Matching)}$$

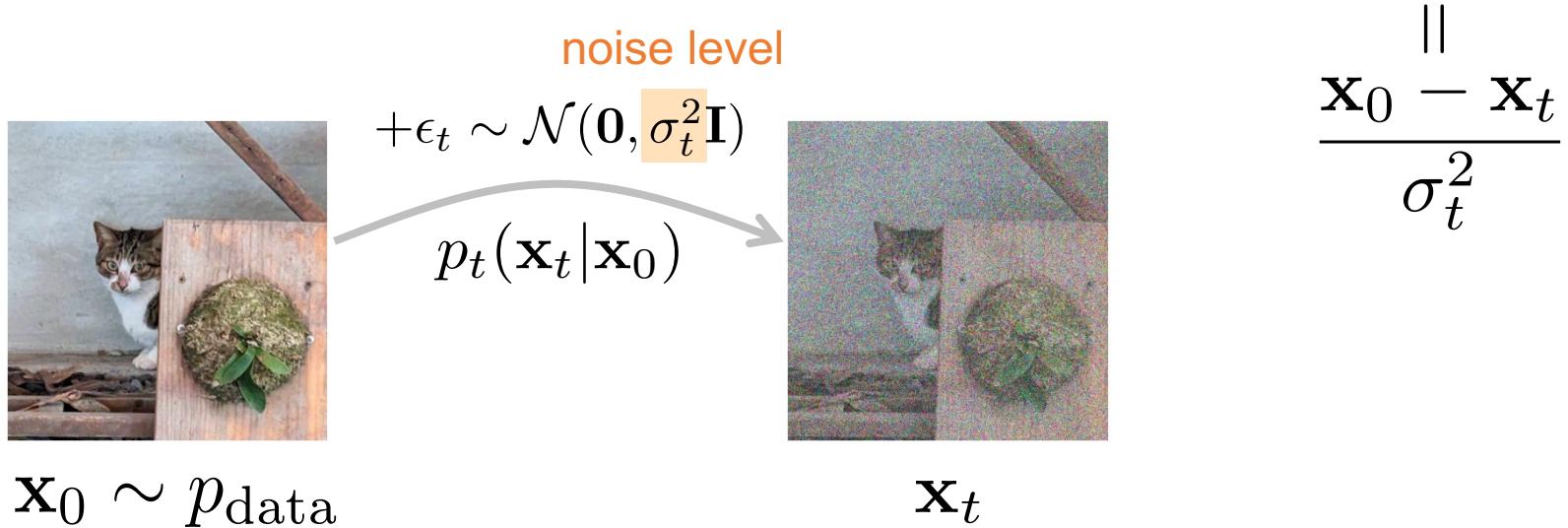
Training Score-based Diffusion Models

Score Matching (SM) [Hyvarinen JMLR'05]:

$$\mathcal{L}_{\text{SM}}(\mathbf{x}_0; \theta) = \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2$$

Denoising Score Matching (DSM) [Vincent Neural Comput.'11]:

$$\mathcal{L}_{\text{DSM}}(\mathbf{x}_0; \theta) = \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \boxed{\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{x}_0)}\|_2^2$$

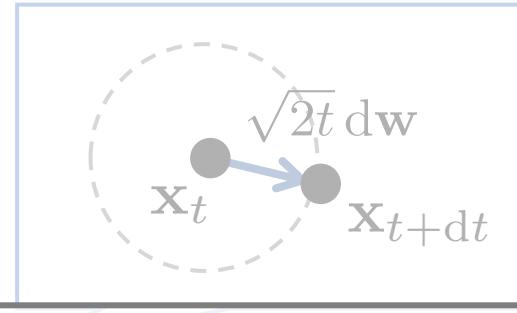


Score-based Diffusion Models [Song+ ICLR'21]

$$f(\mathbf{x}_t, t) = 0, g(t) = \sqrt{2t}$$

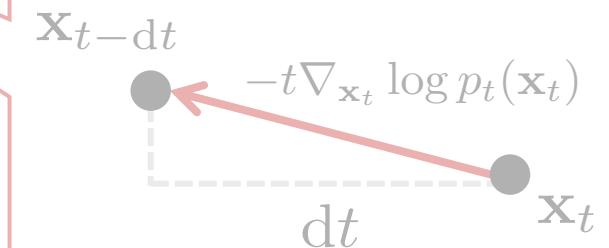
Forward SDE

$$d\mathbf{x}_t = \sqrt{2t} dw$$



1. How do we obtain the score function during sampling?
2. How to solve the reverse ODE?

$$d\mathbf{x}_t = -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$$



$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

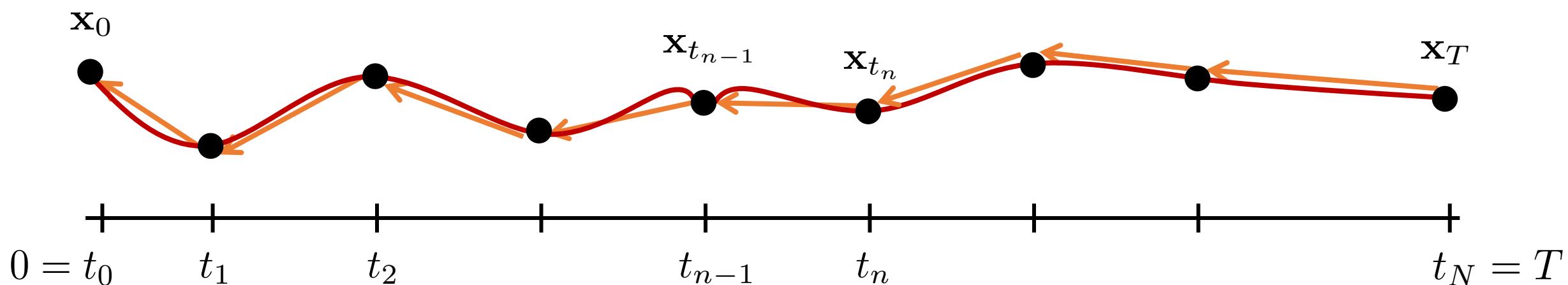
How to Solve the Reverse ODE?

Euler's Method (1st order):

$$\frac{d\mathbf{x}_t}{dt} = -ts\theta(\mathbf{x}_t, t)$$

||

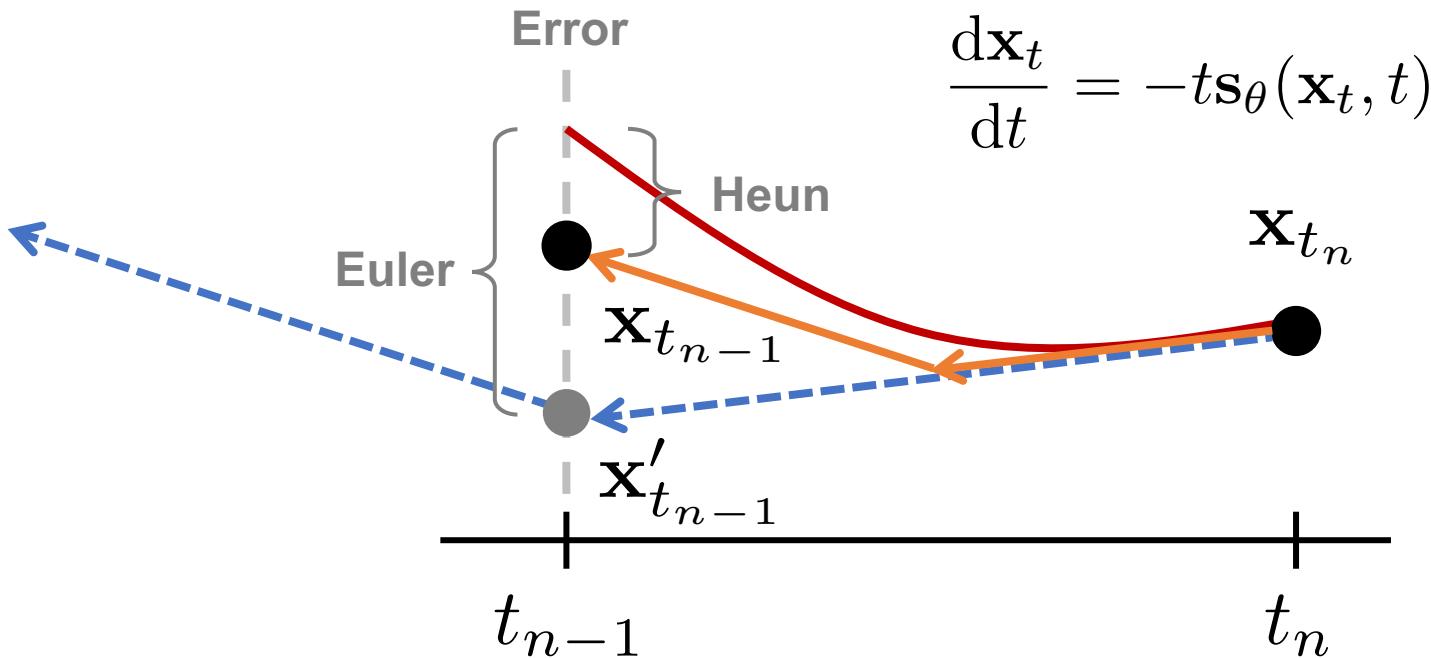
$$\frac{\mathbf{x}_{t_{n-1}} - \mathbf{x}_{t_n}}{t_{n-1} - t_n} + O(n)$$



$$\mathbf{x}_{t_{n-1}} = \mathbf{x}_{t_n} - (t_{n-1} - t_n) t_n \mathbf{s}_\theta(\mathbf{x}_{t_n}, t_n)$$

How to Solve the Reverse ODE?

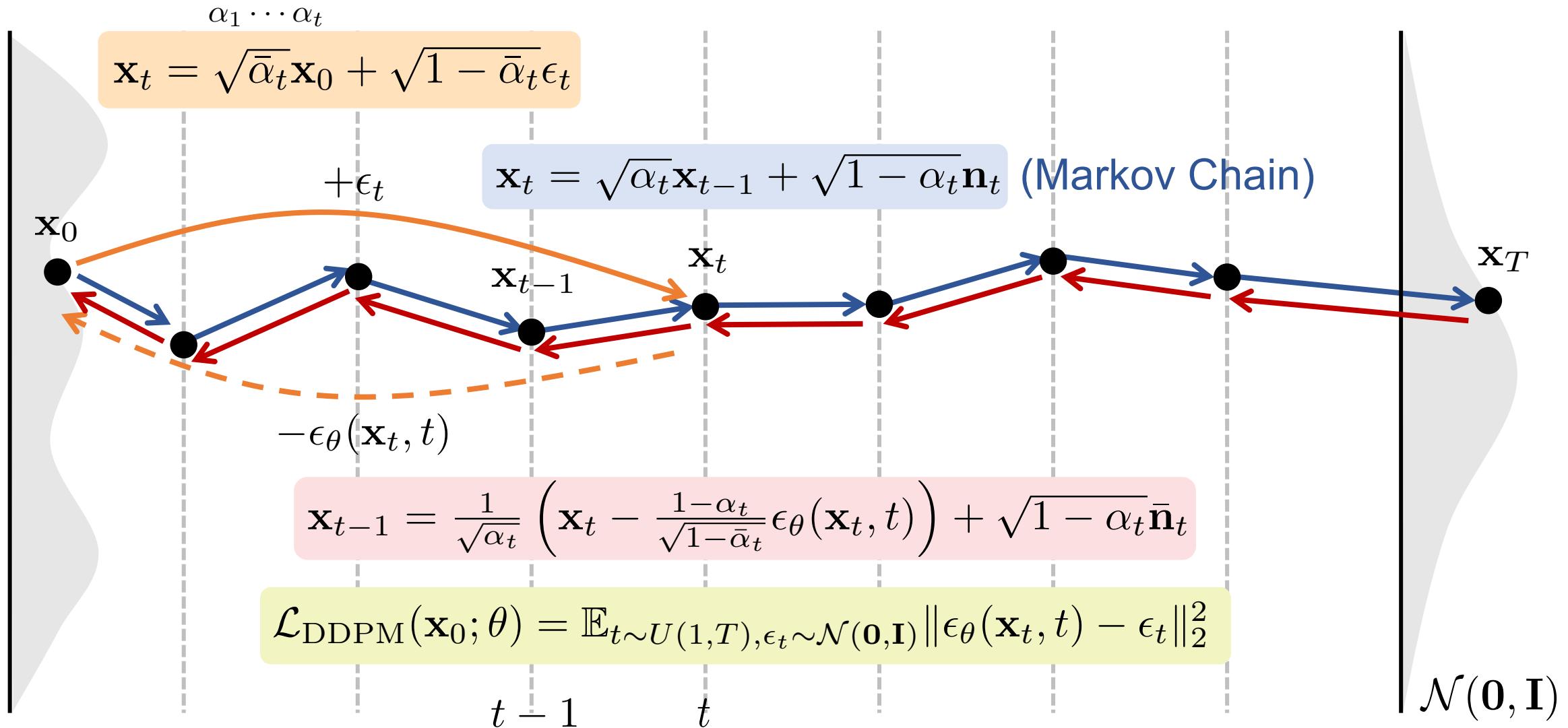
Heun's Method (2nd order):



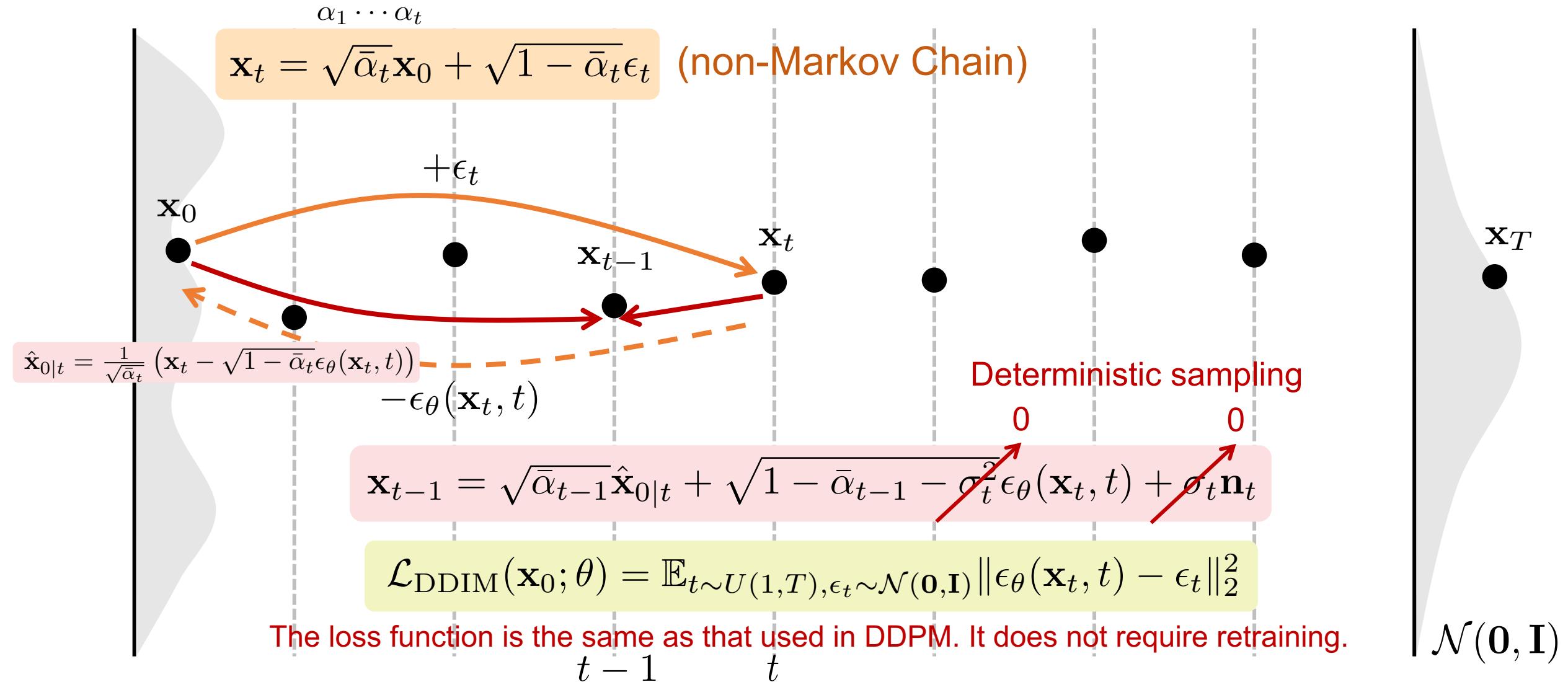
$$\mathbf{x}'_{t_{n-1}} = \mathbf{x}_{t_n} - (t_{n-1} - t_n) t_n \mathbf{s}_\theta(\mathbf{x}_{t_n}, t_n)$$

$$\mathbf{x}_{t_{n-1}} = \mathbf{x}_{t_n} - \frac{1}{2}(t_{n-1} - t_n)[t_n \mathbf{s}_\theta(\mathbf{x}_{t_n}, t_n) + t_{n-1} \mathbf{s}_\theta(\mathbf{x}'_{t_{n-1}}, t_{n-1})]$$

Denoising Diffusion Probabilistic Models (DDPM) [Ho+ NeurIPS'20]

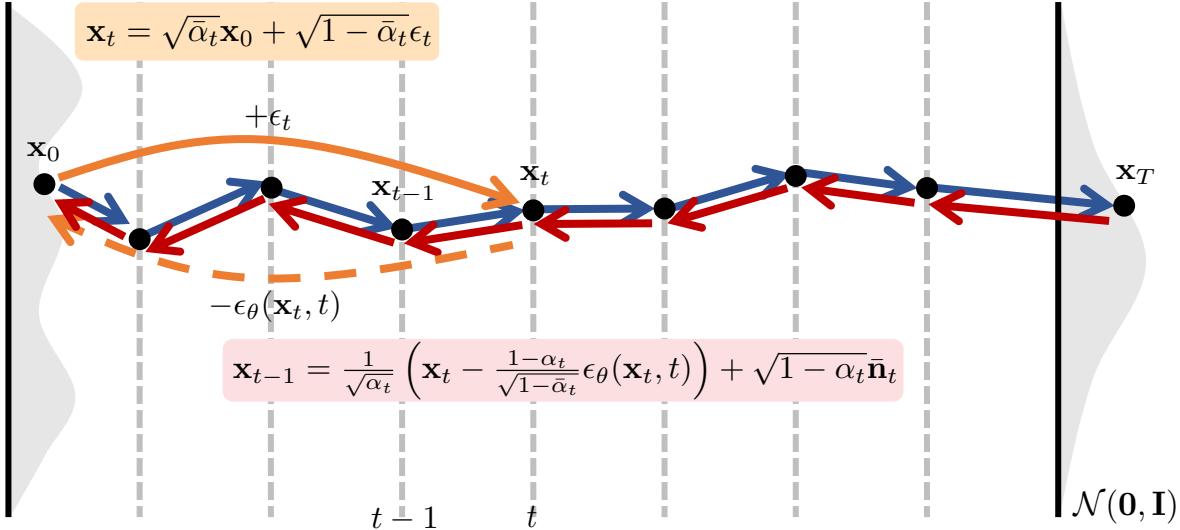


Denoising Diffusion Implicit Models (DDIM) [Song+ ICLR'21]



Summary: Discrete vs Continuous

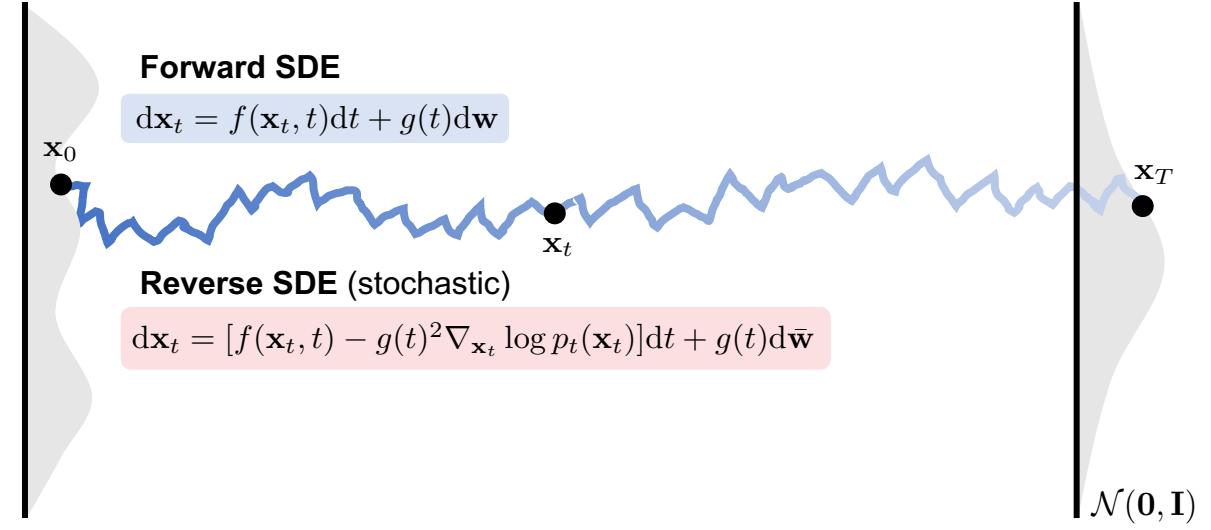
DDPM



$$\mathcal{L}_{\text{DDPM}}(\mathbf{x}_0; \theta) = \mathbb{E}_{t \sim U(1, T), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \|\epsilon_\theta(\mathbf{x}_t, t) - \epsilon_t\|_2^2$$

Discretized in the training phase

Score-based SDE

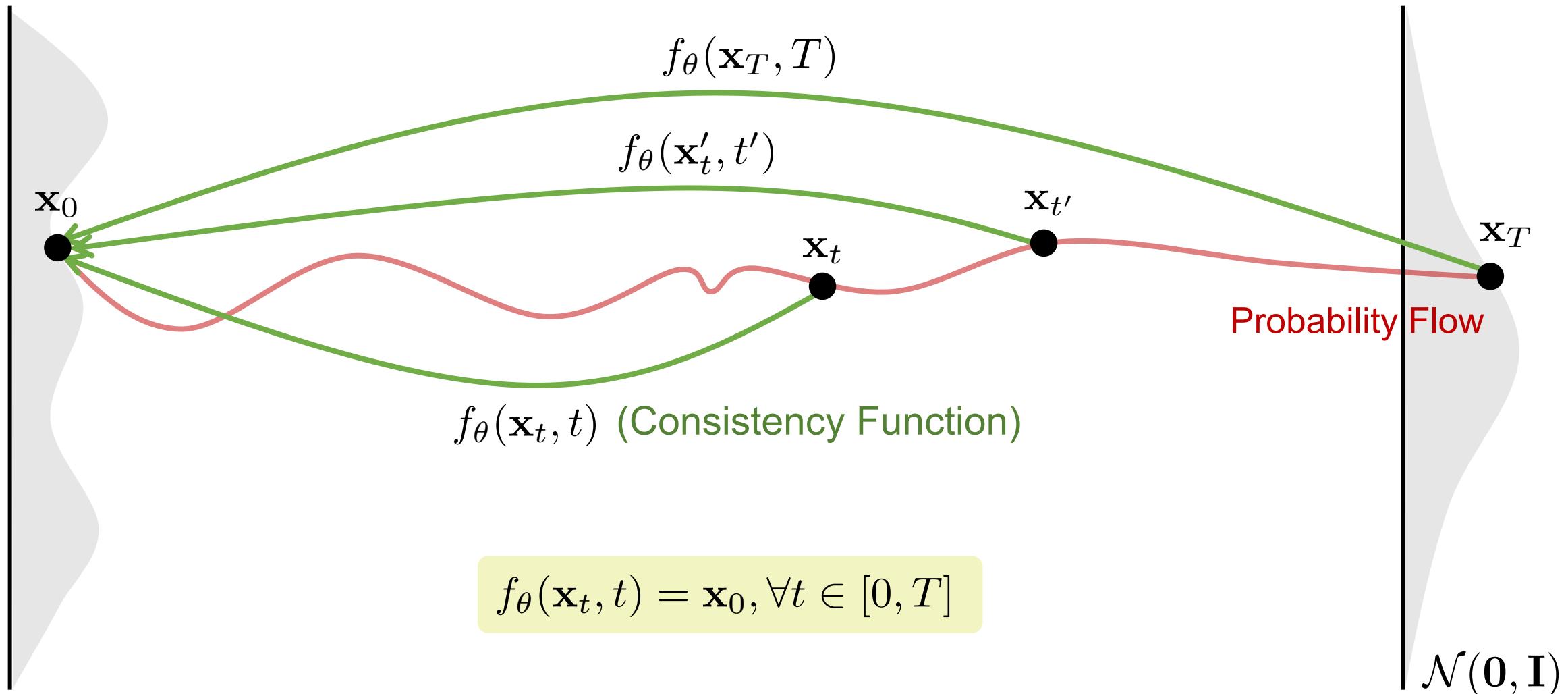


$$\mathcal{L}_{\text{DSM}}(\mathbf{x}_0; \theta) = \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{x}_0)\|_2^2$$

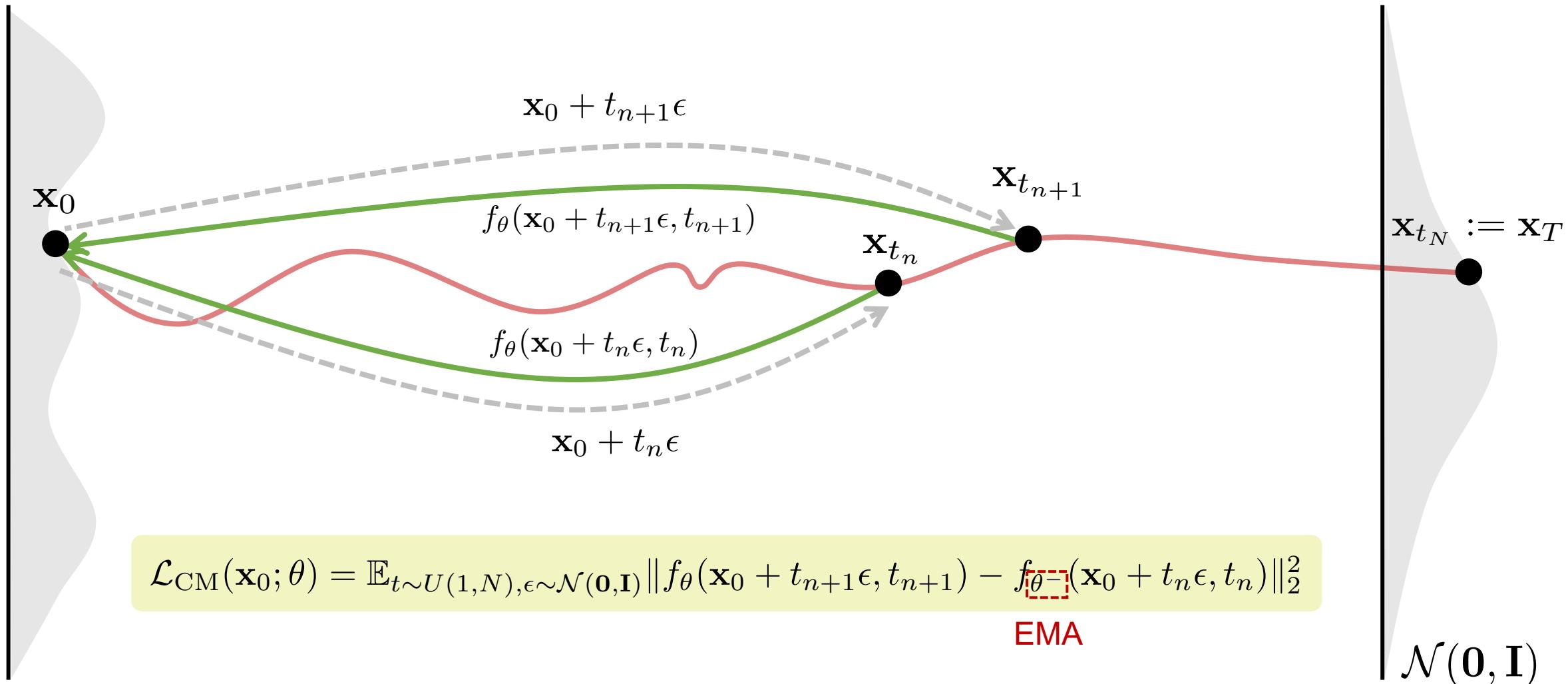
Discretized during sampling
The error is only affected by the order of the solver

Consistency Models (CM) [Song+ ICML'23]

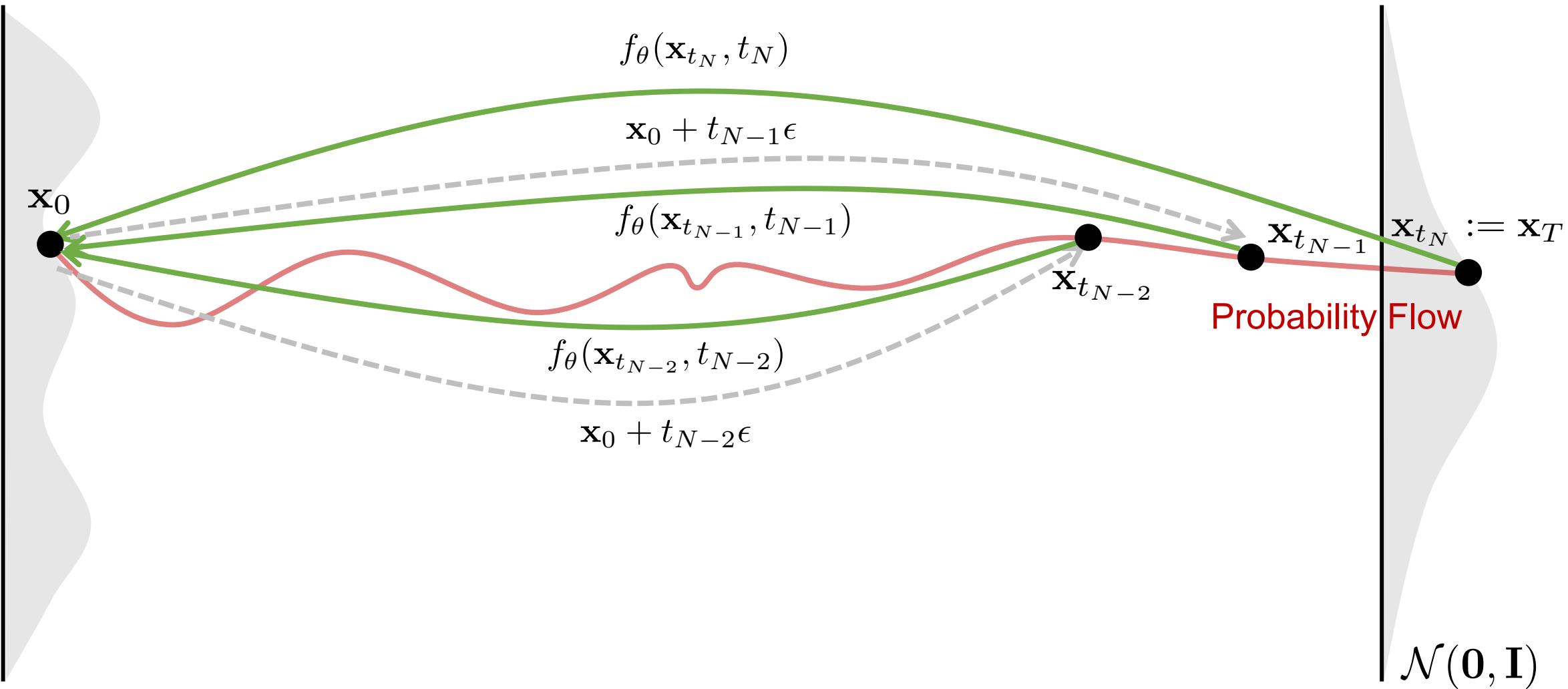
$$d\mathbf{x}_t = -t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) dt$$



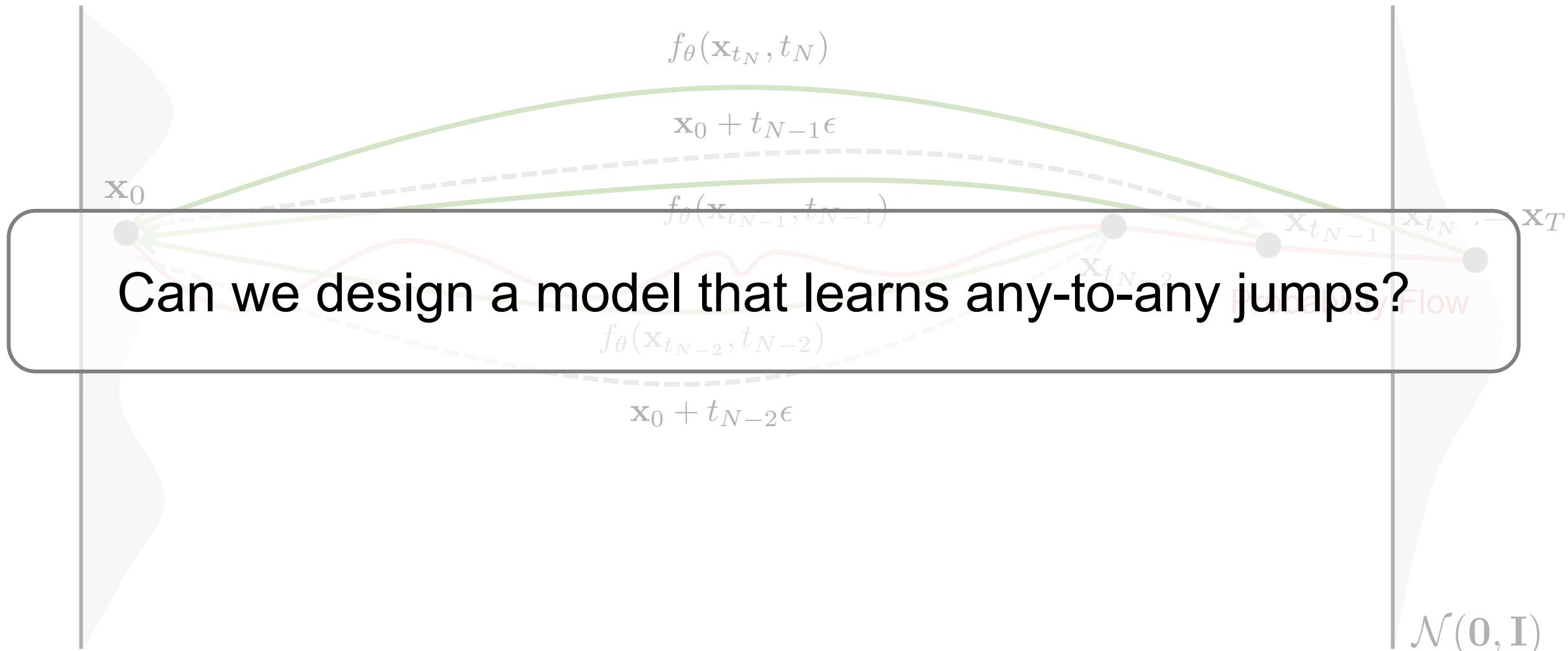
Training CM



Sampling with CM

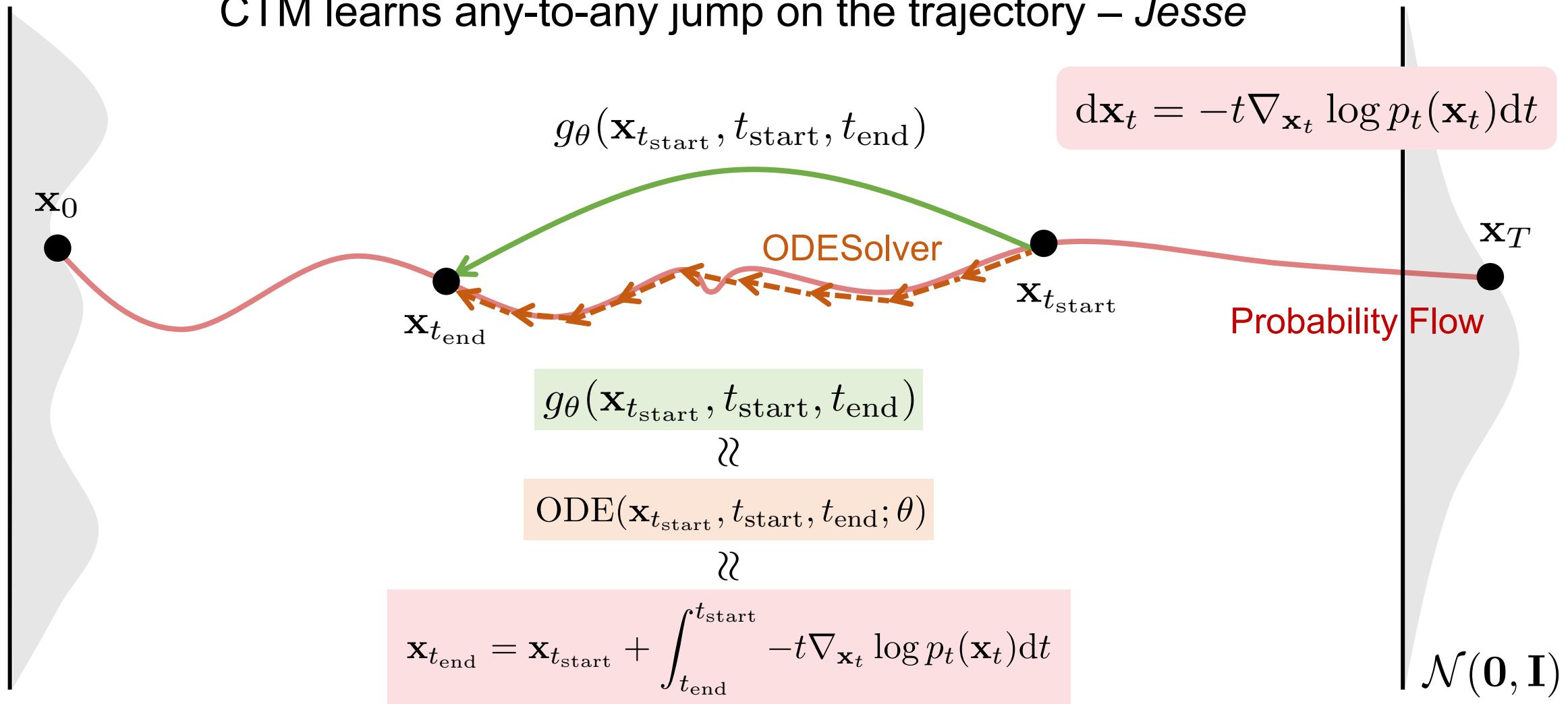


Sampling with CM

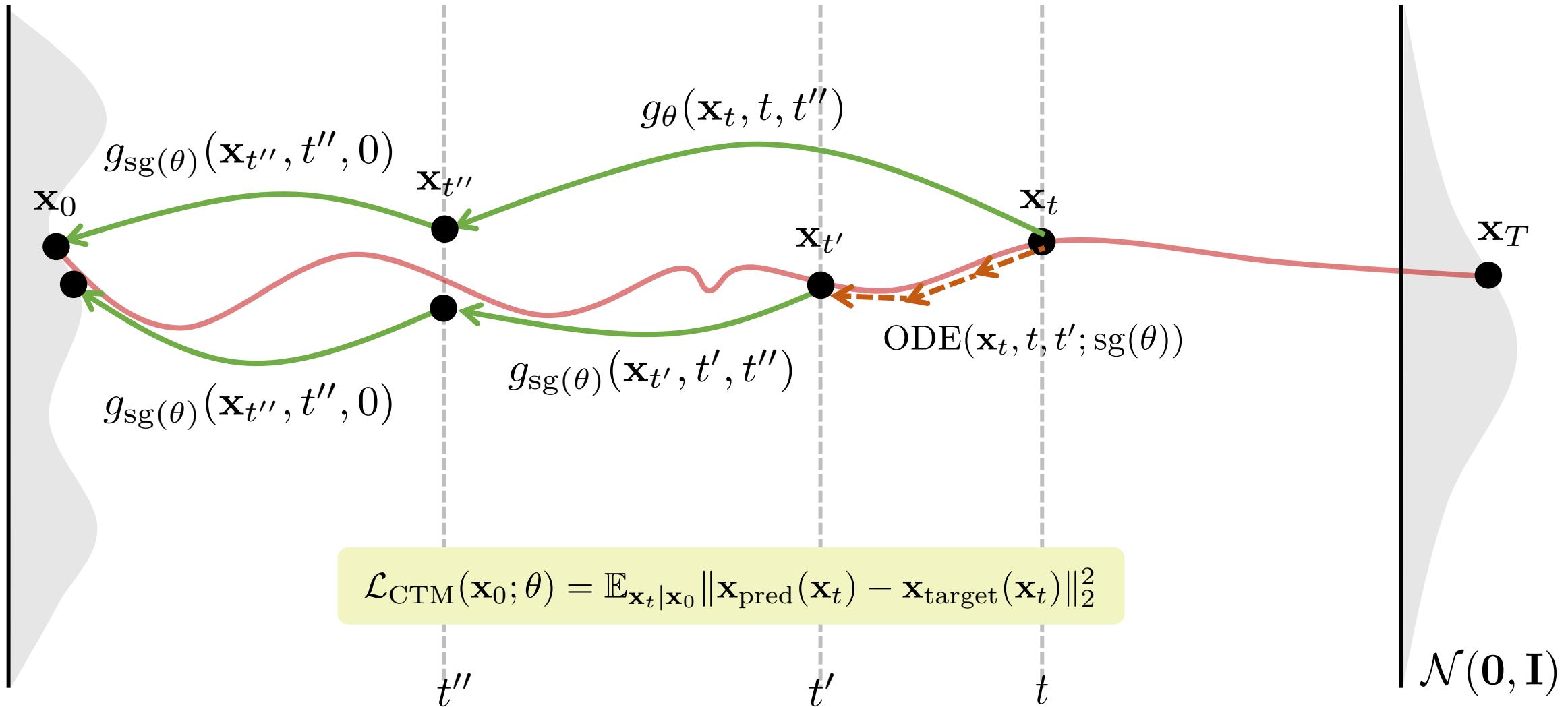


Consistency Trajectory Models (CTM) [Kim+ ICLR'24]

CTM learns any-to-any jump on the trajectory – Jesse



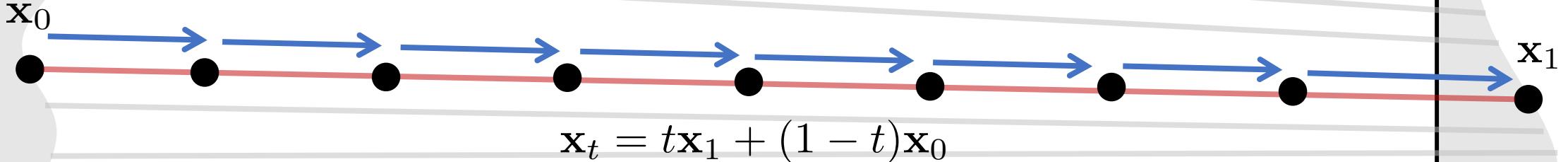
Training CTM



Flow Matching [Lipman+ ICLR'23][Liu+ ICLR'23]

Forward ODE

$$\begin{aligned} d\mathbf{x}_t &= (\mathbf{x}_1 - \mathbf{x}_0) dt \\ &\Downarrow \\ \mathbf{v}_\theta(\mathbf{x}_t, t) \end{aligned}$$



Reverse ODE

$$d\mathbf{x}_t = -\mathbf{v}_\theta(\mathbf{x}_t, t) dt$$

$$\mathcal{L}_{\text{FM}}(\mathbf{x}_0, \mathbf{x}_1; \theta) = \mathbb{E}_{t \sim U(0,1)} \|\mathbf{v}_\theta(t\mathbf{x}_1 + (1-t)\mathbf{x}_0, t) - (\mathbf{x}_1 - \mathbf{x}_0)\|_2^2$$

$\mathcal{N}(\mathbf{0}, \mathbf{I})$

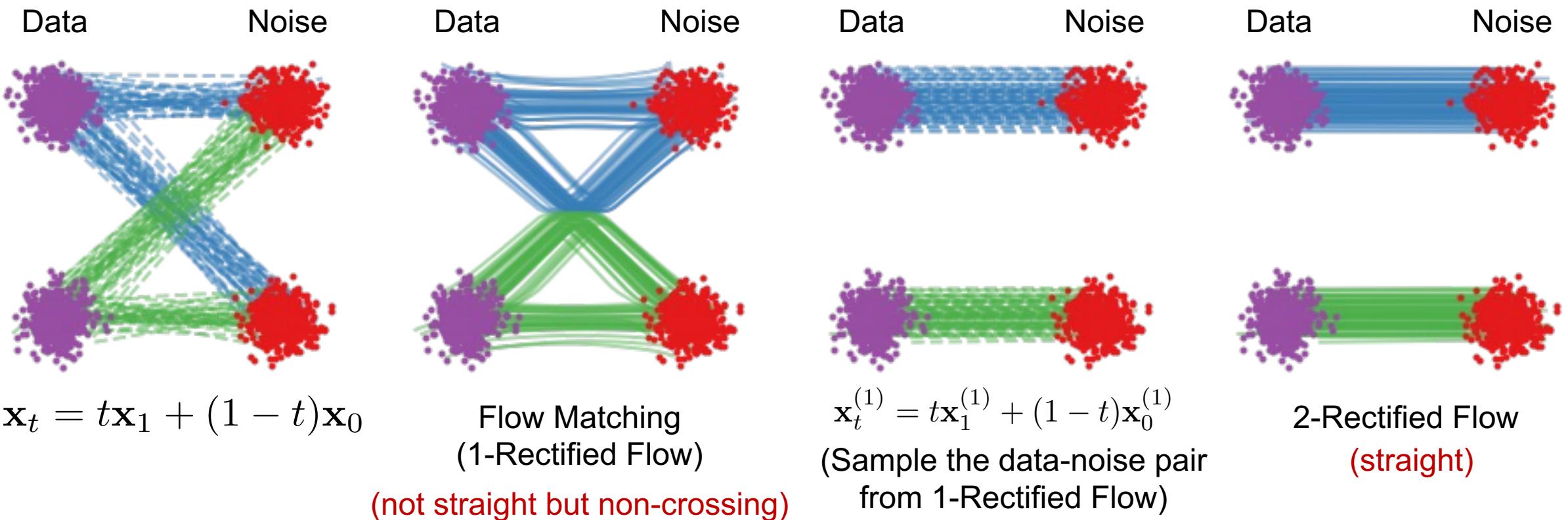
Flow Matching

[Lipman+ ICLR'23][Liu+ ICLR'23]

Is it possible that the paths learned by the model intersect?

$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

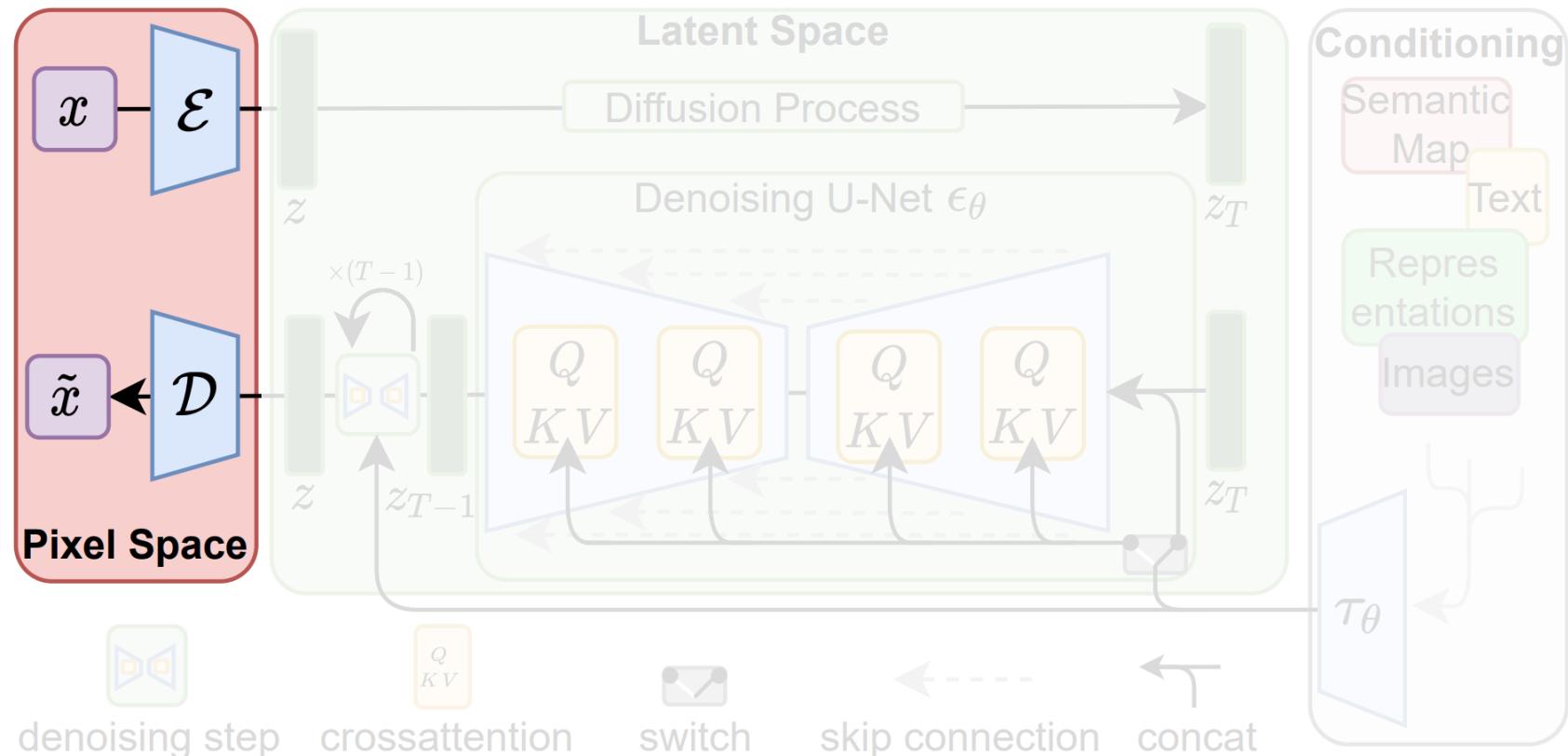
k -Rectified Flow [Liu+ ICLR'23]



How to reduce the computational cost of diffusion models?

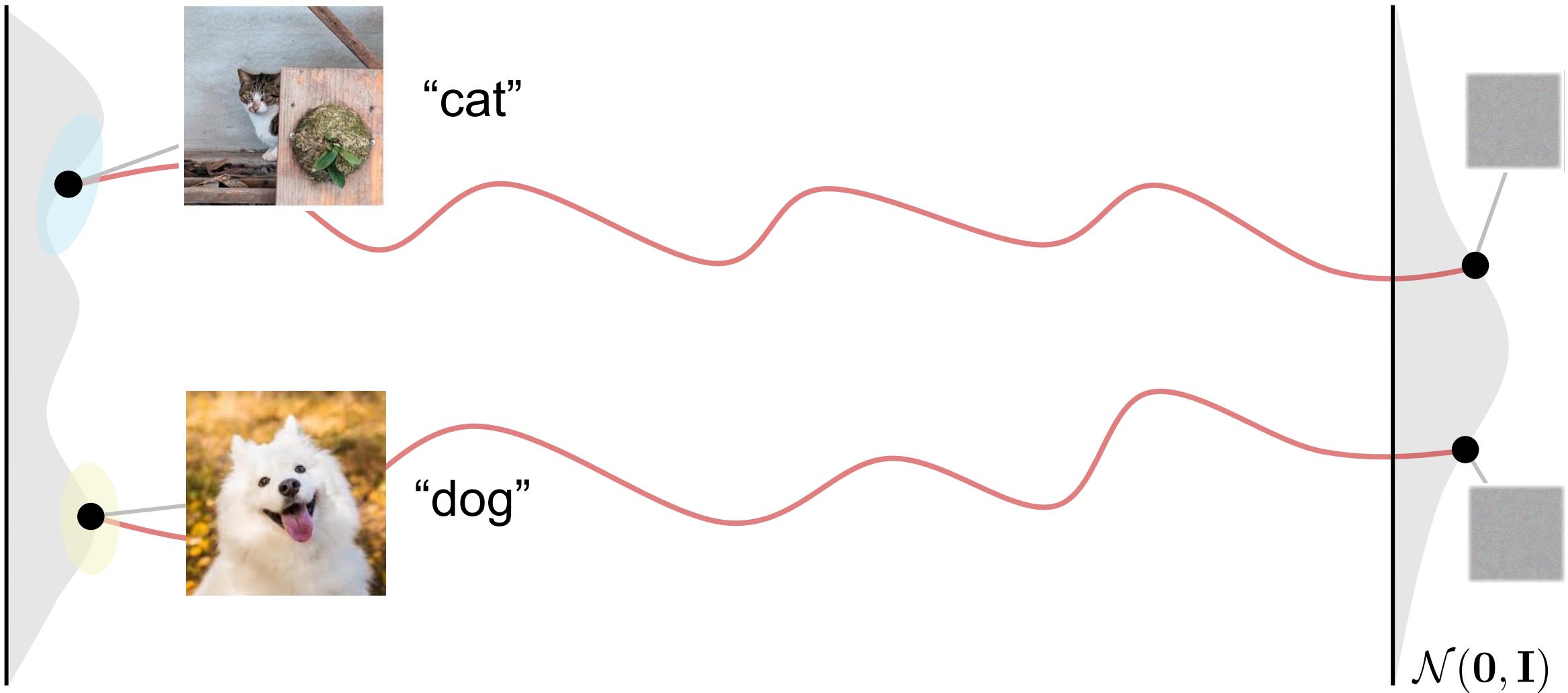
Latent Diffusion Model (LDM) [Rombach+ CVPR'22]

Use the pretrained-VAE to compress the image to latent reducing computation time

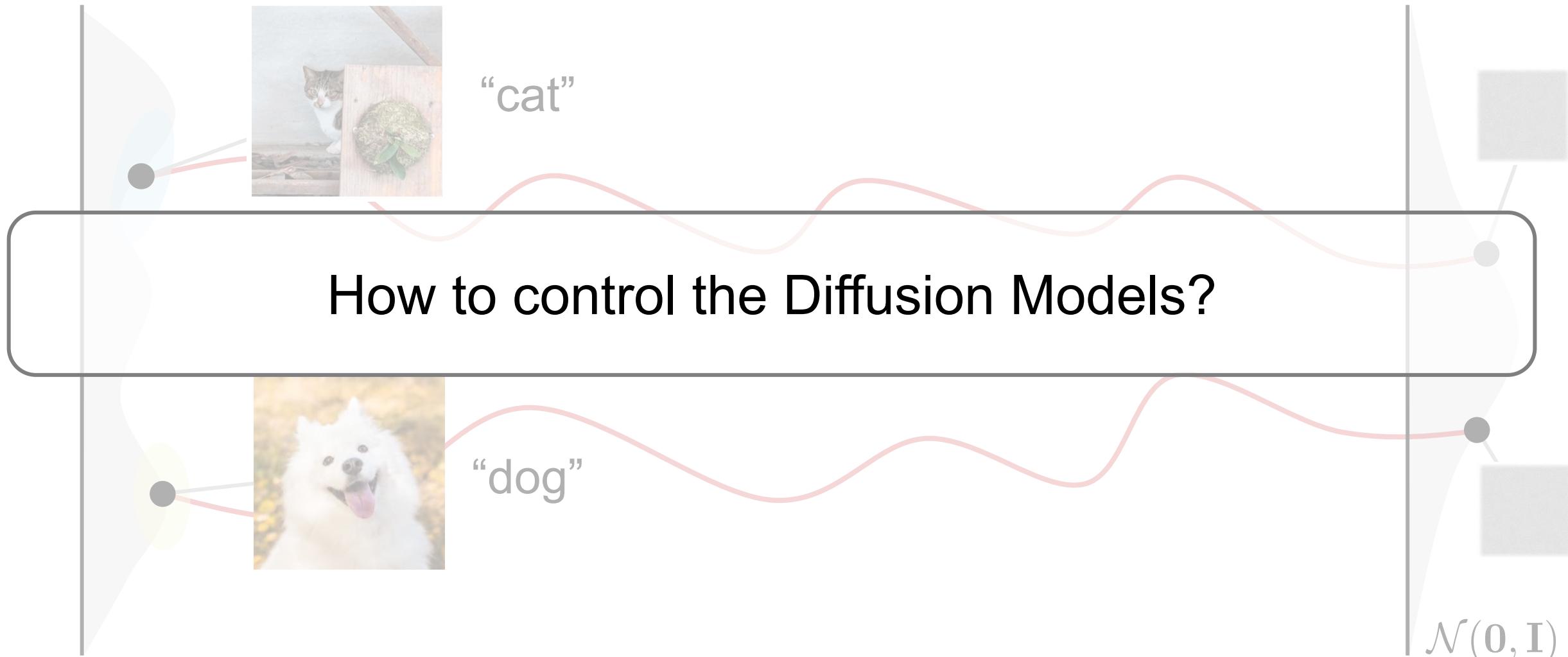


Can we control the Diffusion Models?

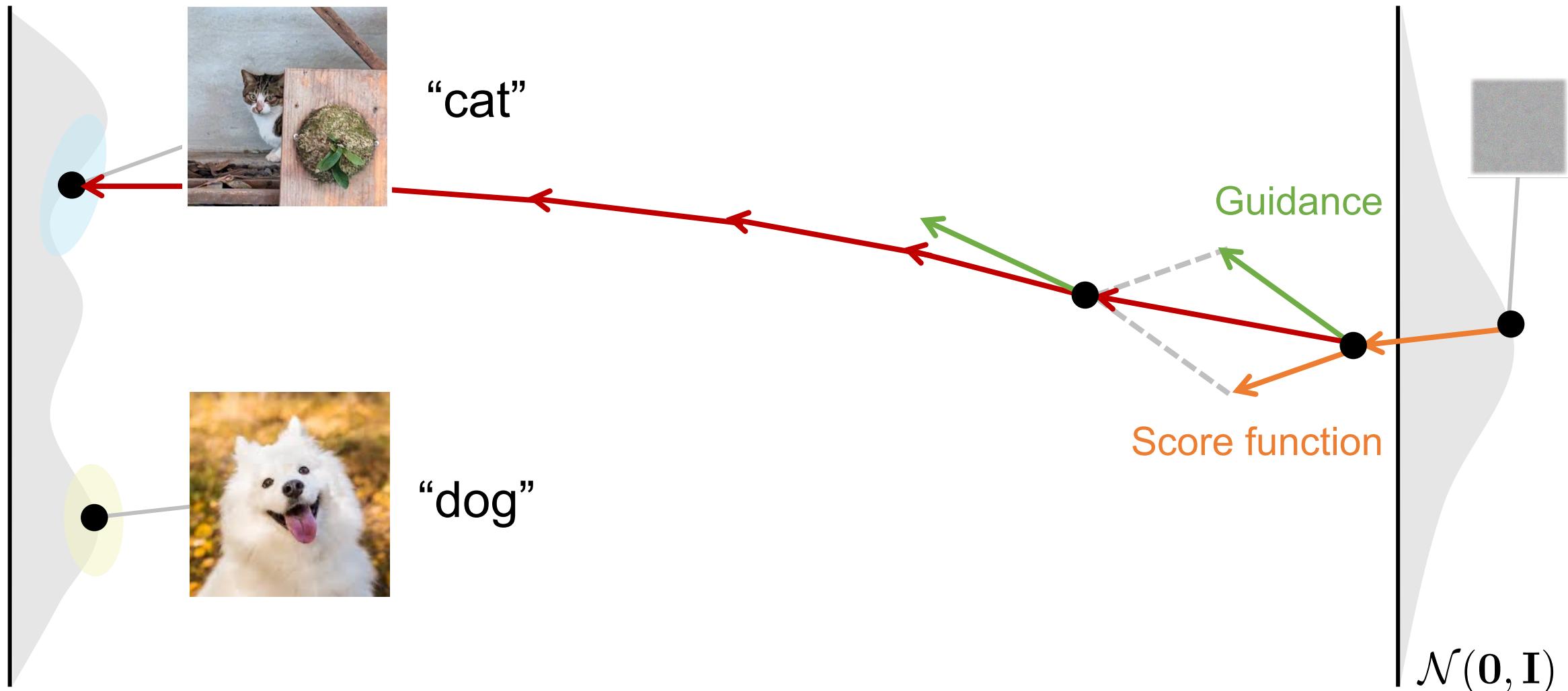
Conditional Diffusion Models



Conditional Diffusion Models



Conditional Diffusion Models



Classifier Guidance (CG) [Dhariwal+ NeurIPS'21]

- Bayes' Rule:

$$P\left(\mathbf{x}_t = \begin{array}{|c|} \hline \text{cat} \\ \hline \text{ball} \\ \hline \end{array} \mid \mathbf{y} = \text{"cat"}\right) = \frac{p(\mathbf{x}_t)p(\mathbf{y}|\mathbf{x}_t)}{p(\mathbf{y})}$$

- Bayes' Rule for Score Function:

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{y})^0 \\ &= \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)} + \nabla_{\mathbf{x}_t} \log \boxed{p(\mathbf{y}|\mathbf{x}_t)}\end{aligned}$$

Unconditional score Classifier (need to additional training)

Classifier-Free Guidance (CFG) [Ho+ NeurIPS'21]

From CG, we have

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

Reweight the coefficient between unconditional score and classifier score

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log \tilde{p}(\mathbf{x}_t|\mathbf{y}) &:= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \boxed{\gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)} \\ &= \gamma \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y})} + (1 - \gamma) \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}\end{aligned}$$

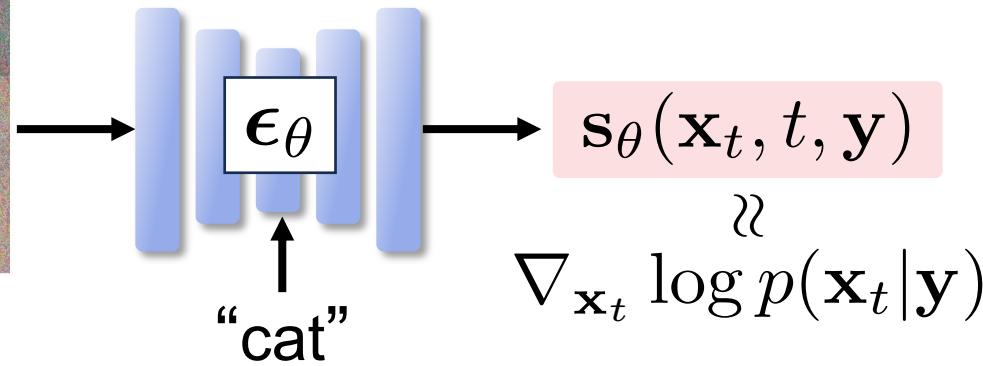
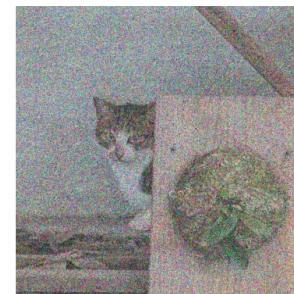
Conditional score Unconditional score

Training CFG

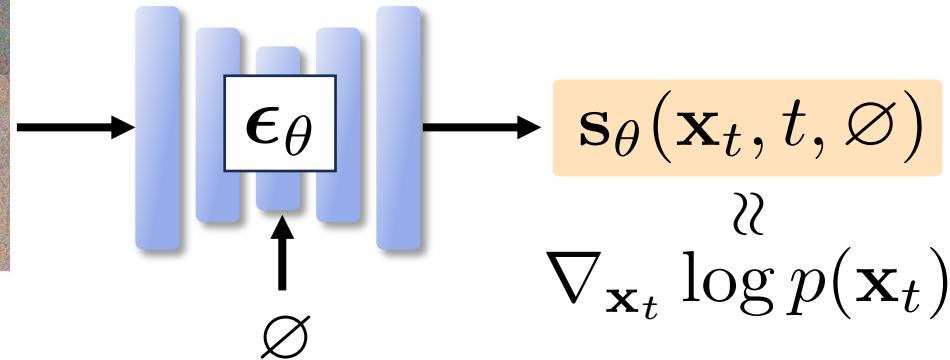
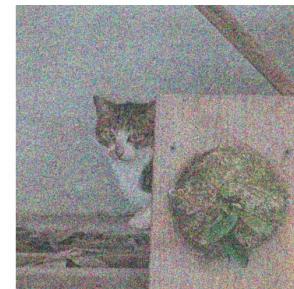
$$\nabla_{\mathbf{x}_t} \log \tilde{p}(\mathbf{x}_t | \mathbf{y}) = \gamma \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})} + (1 - \gamma) \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}$$

Conditional score Unconditional score

Predict
Conditional Score



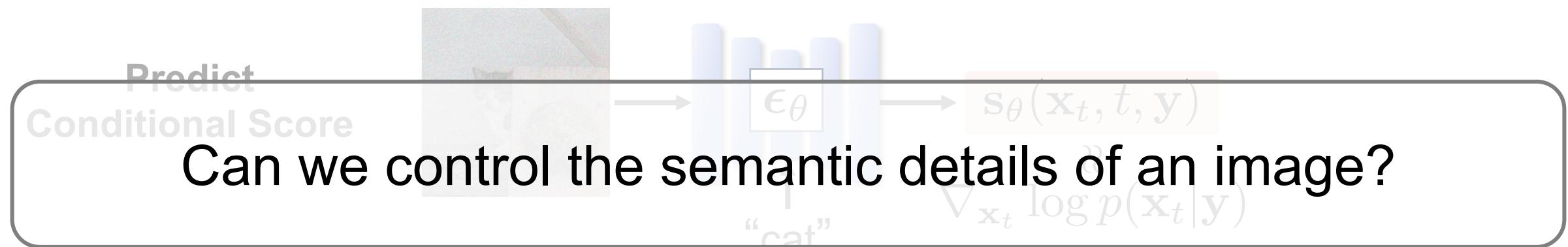
Predict
Unconditional Score



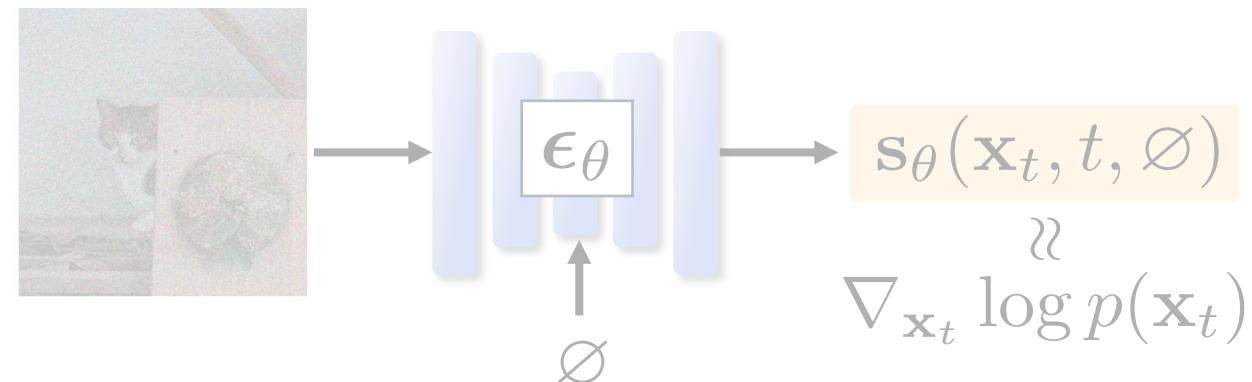
Training CFG

$$\nabla_{\mathbf{x}_t} \log \tilde{p}(\mathbf{x}_t | \mathbf{y}) = \gamma \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y})} + (1 - \gamma) \boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}$$

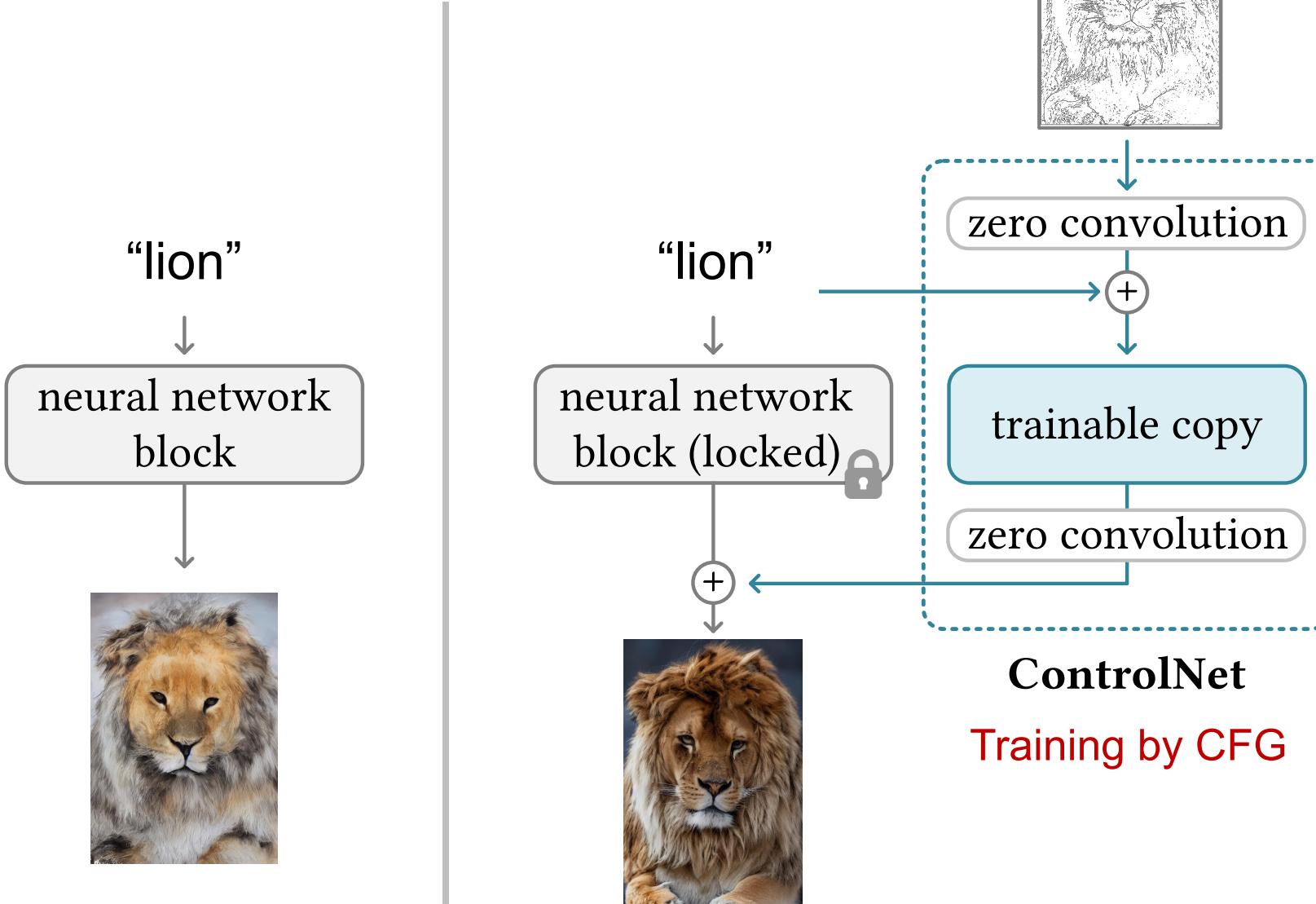
Conditional score Unconditional score



Predict
Unconditional Score



ControlNet [Zhang+ ICCV'23]



Applications of Diffusion Models in 3D Vision



Zero-1-to-3
2023 / 5



Magic123
2023 / 6



One-2-3-45++
2023 / 11

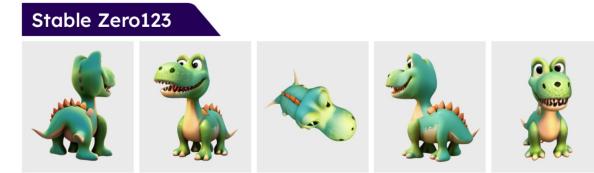
DreamFusion
2022 / 9



One-2-3-45
2023 / 6



Zero123++
2023 / 10



Stable Zero123
2023 / 12

Takeaway

- Diffusion models build a bridge between noise and data, forming a powerful generative modeling framework.
- Score-based models leverage SDE/ODE formulations and score functions to guide the reverse process.
- Variants like DDIM, CM, CTM, and Flow Matching offer trade-offs in speed, quality, and control.
- Conditional methods like CFG and ControlNet significantly improve controllability and application scope.
- Diffusion models are already widely applied in 2D/3D image synthesis, speech, and multi-modal generation.

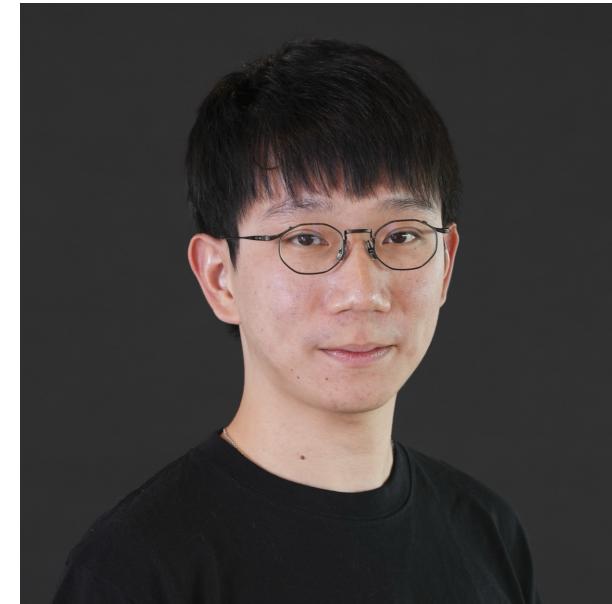
Acknowledgement

Some concepts and insights in these slides are inspired by Yang Song and Jesse



[Yang Song](#)

Research Scientist at OpenAI



[Chieh-Hsin \(Jesse\) Lai](#)

Research Scientist at Sony AI

Thank you