

Diffusion Distillation



Nikita Starodubcev
Yandex Research

Lecture 4 | ODE-Based Diffusion Distillation

- 01** Motivation
- 02** Knowledge Distillation
- 03** Consistency Distillation
- 04** Multi-Boundary Consistency Distillation



1. Motivation

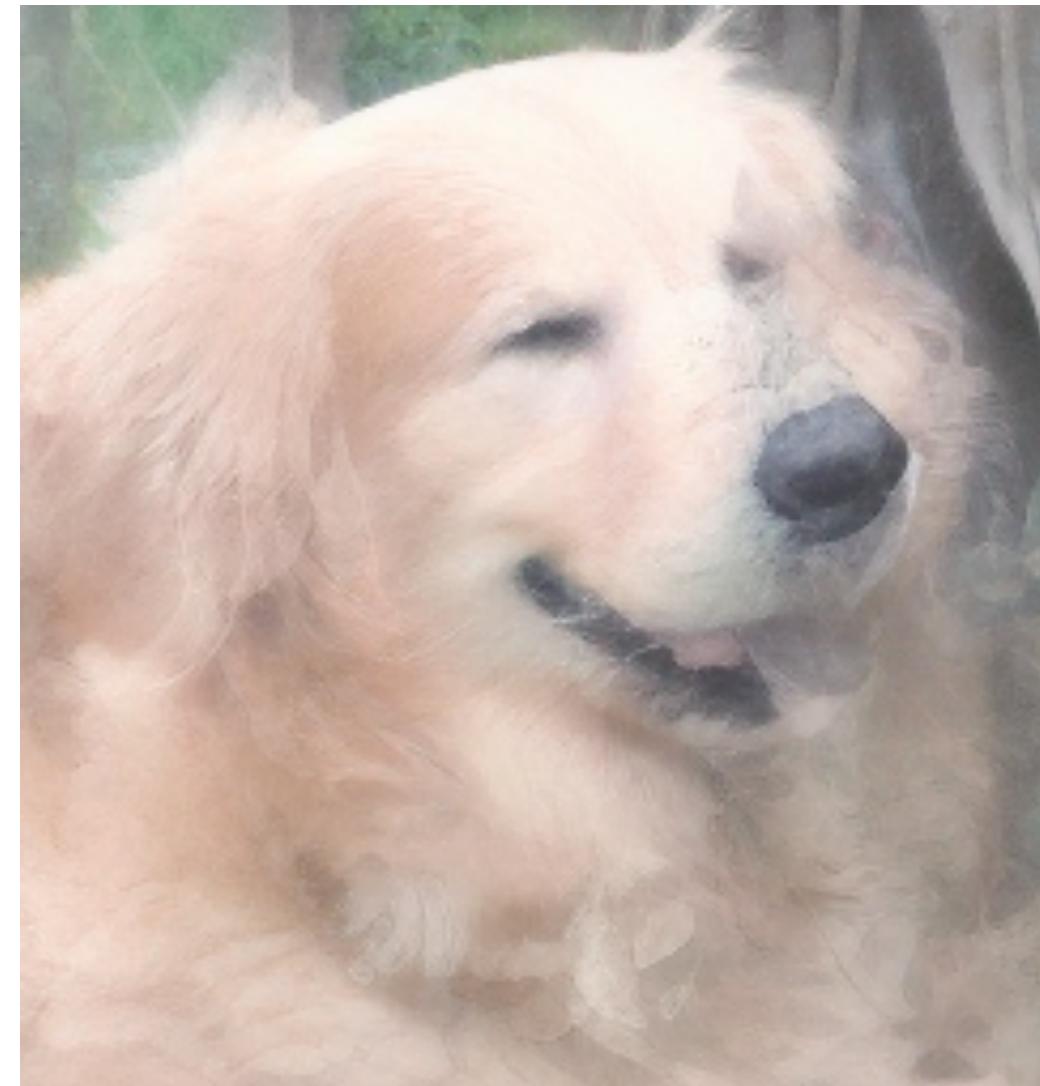
Diffusion Models are slow. Why?

30 steps

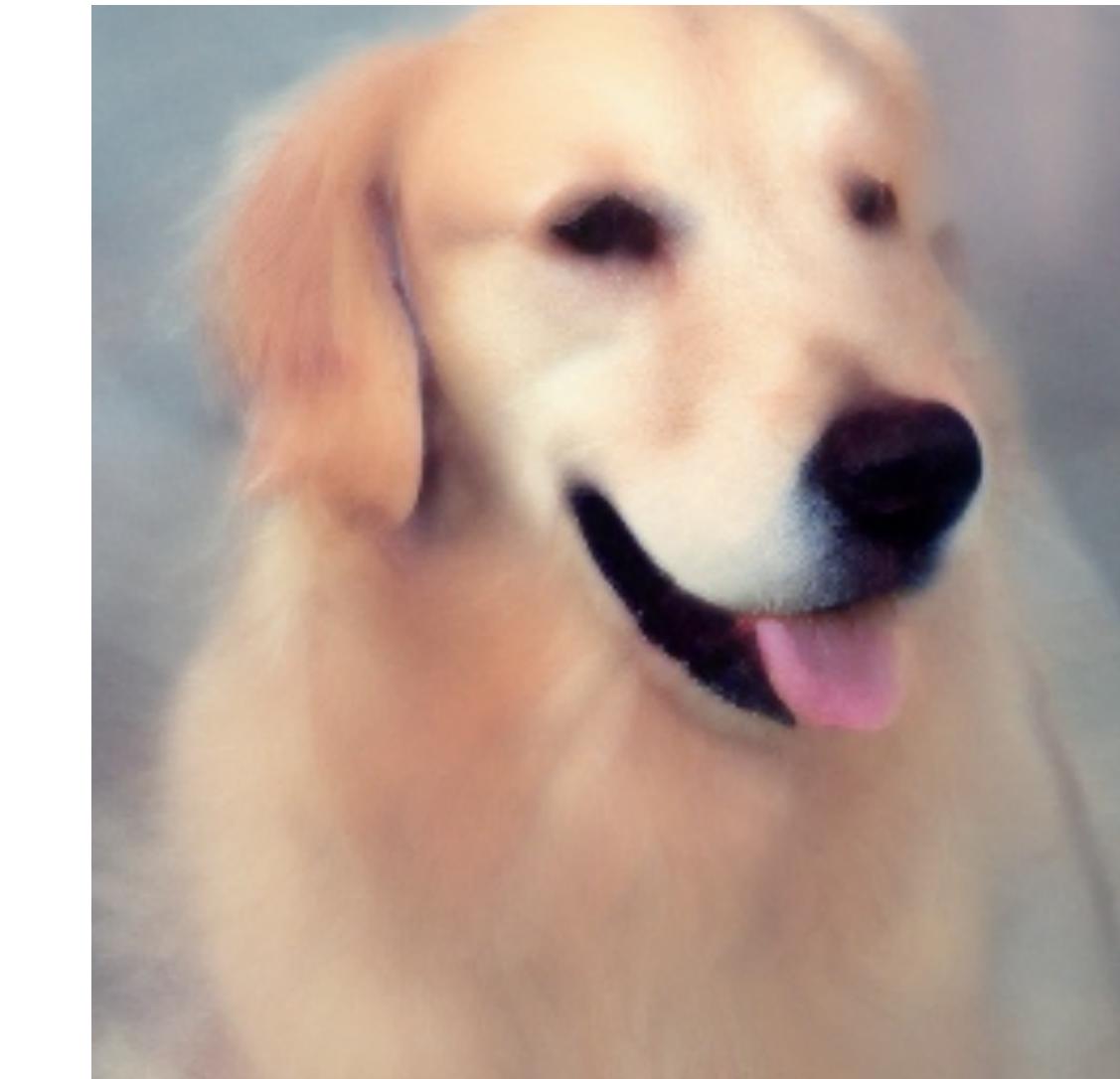


...

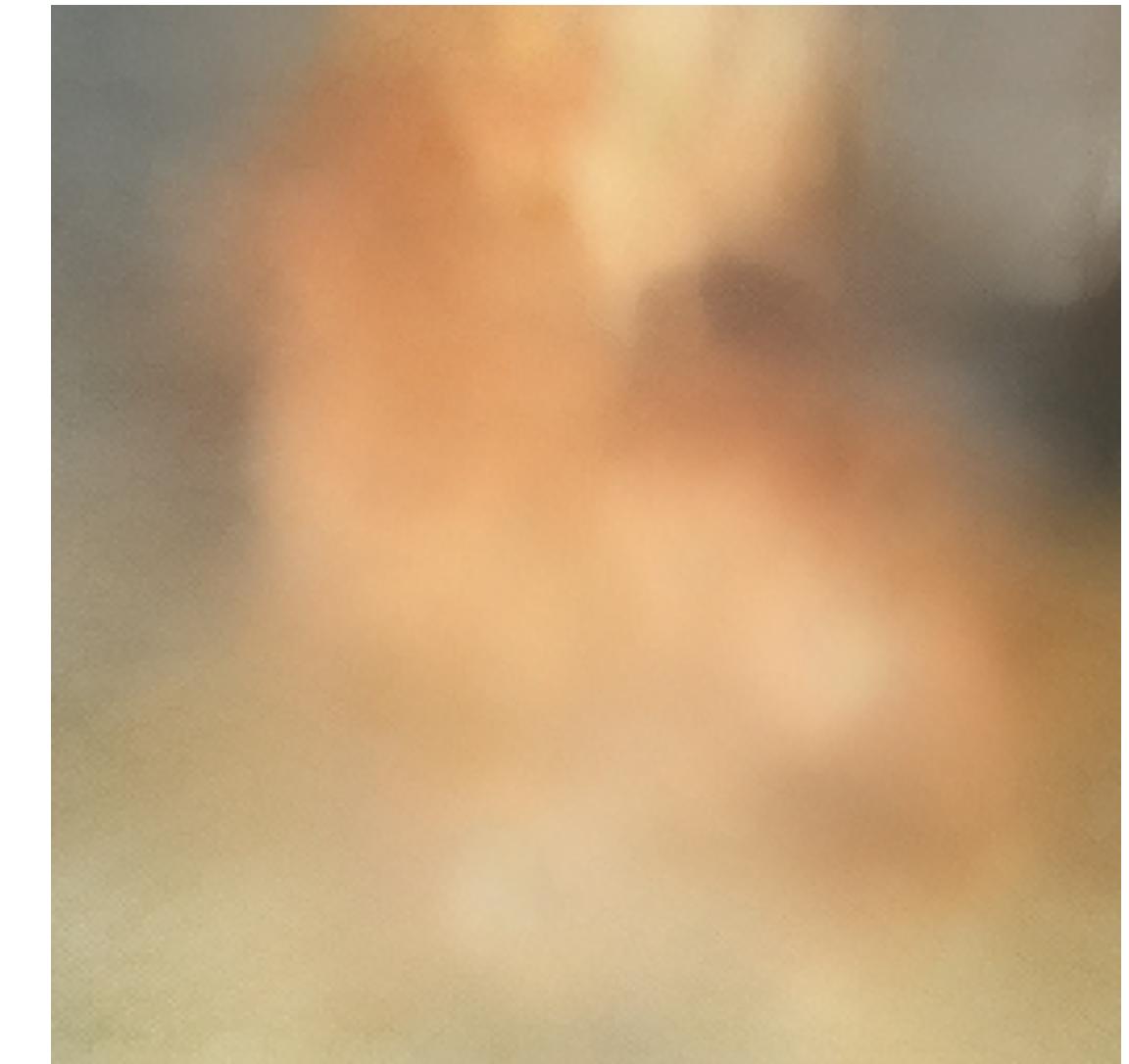
7 steps



5 steps



3 steps



1. PF-ODE, VE case

Remind the reverse PF-ODE equation

$$dx = \left(f(x, t) - \frac{1}{2} g(t)^2 \nabla_x \log p_t(x) \right) dt, t \in [0, 1]$$

Variance Exploding case

$$f(x, t) = 0, \quad g(t) = \sqrt{\frac{d\sigma_t^2}{dt}}$$

$$dx = -\sigma_t \nabla_x \log p_t(x) d\sigma_t$$

$$\nabla_x \log p_t(x) = -\frac{1}{\sigma_t^2} (x - \mathbb{E} x_0 | x)$$

1. Diffusion loss

Remind the diffusion loss function

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(x_0)} \mathbb{E}_{q_t(x|x_0)} \|x_0 - f_\theta(x, t)\|_2^2 \rightarrow \min_{\theta}$$

$\mathcal{U}(0,1)$ – Uniform distribution of time

$p_{data}(x_0)$ – Data distribution (e.g., images)

$q_t(x|x_0) = \mathcal{N}(x|x_0, \sigma_t^2)$ – Transition kernels (how to noise the data)

1. Why Diffusion Models are slow?

Potential sources of low speed?

1. Low capacity of a neural network that approximates $\mathbb{E} [x_0 | x]$?
2. Maybe we need a solver with lower discretization error?

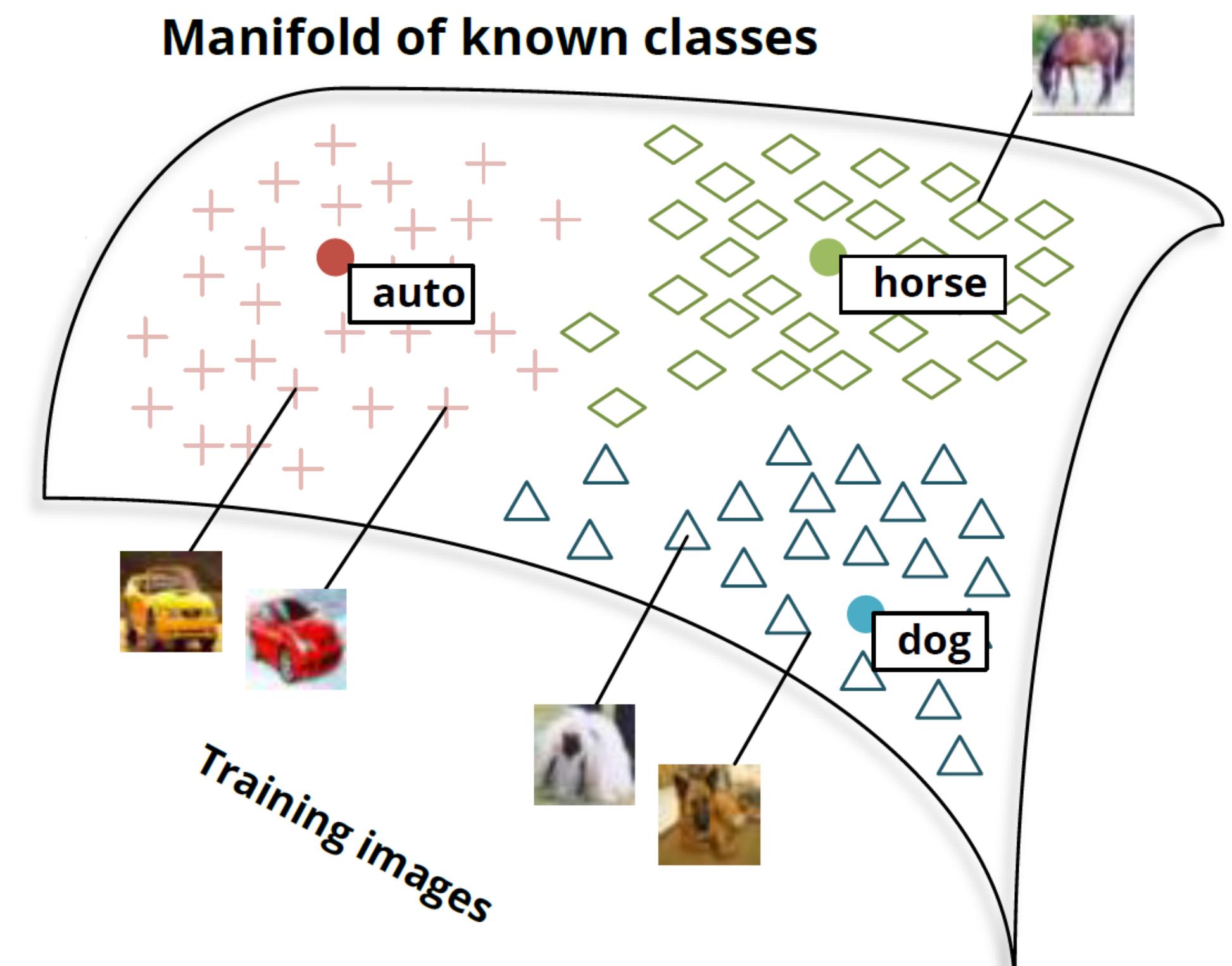
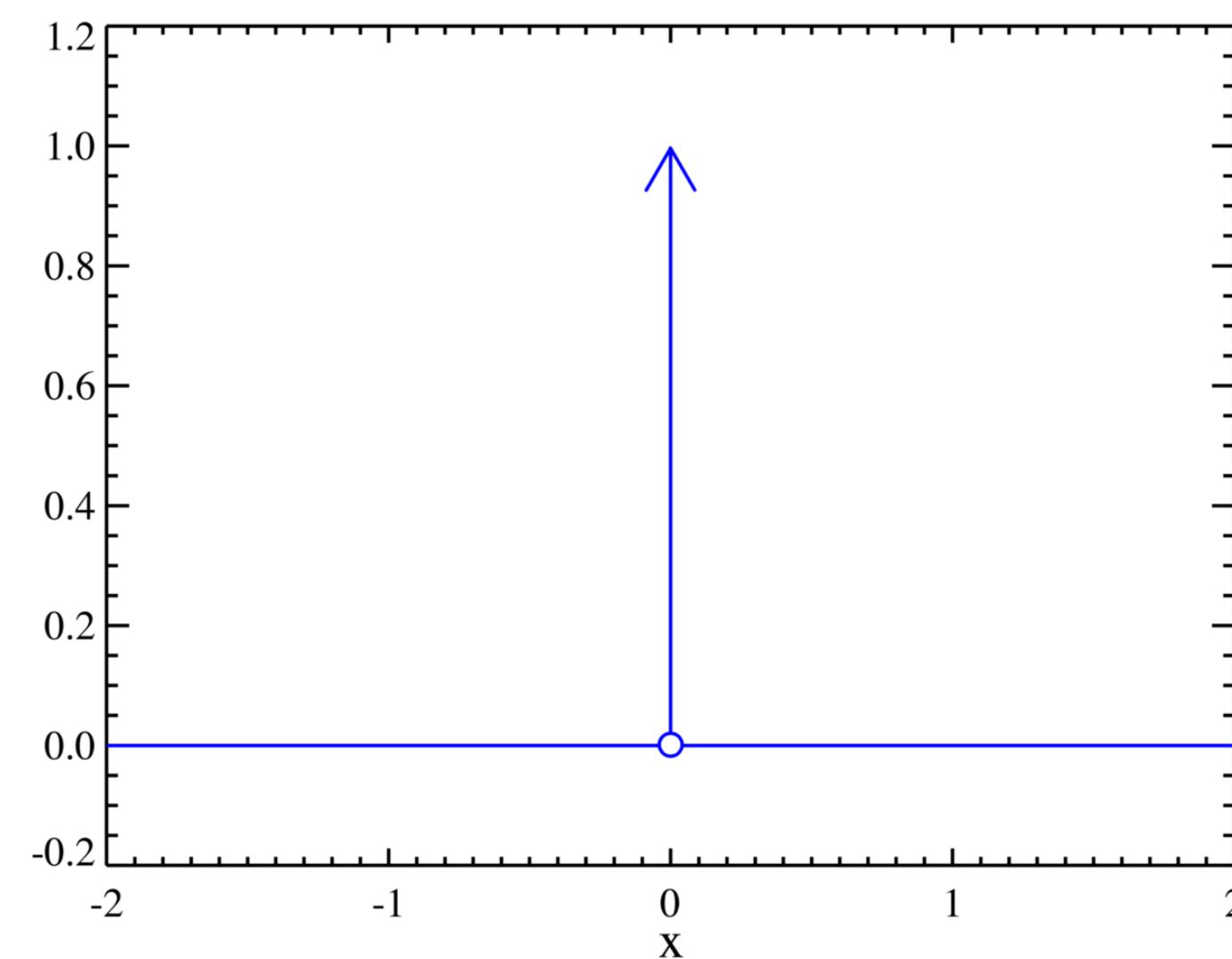
1. Towards ideal approximation

What if our neural network perfectly approximates $\mathbb{E} [x_0 | x]$?

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(x_0)} \mathbb{E}_{q_t(x|x_0)} \|x_0 - f_\theta(x, t)\|_2^2 - \text{optimum?}$$

Assume that

$$p_{data}(x_0) = \frac{1}{N} \sum_{j=1}^N \delta(x_0 - x_0^j); \{x_0^j\}_{j=1}^N - \text{our dataset}$$



Expand the expectations:

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(x_0)} \mathbb{E}_{q_t(x|x_0)} \|x_0 - f_\theta(x, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t | 0, 1)$$

Expand the expectations:

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_t(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t | 0, 1) \underbrace{\int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j)}$$

Expand the expectations:

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_t(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t | 0, 1) \int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{x}_0 - \mathbf{x}_0^j) \int_{\mathbb{R}^n} d\mathbf{x} \mathcal{N}(\mathbf{x} | \mathbf{x}_0, \sigma_t^2) \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2 =$$

Expand the expectations:

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_t(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t | 0, 1) \frac{1}{N} \sum_{j=1}^N \int_{\mathbb{R}^n} d\mathbf{x}_0 \delta(\mathbf{x}_0 - \mathbf{x}_0^j) \left[\int_{\mathbb{R}^n} d\mathbf{x} \mathcal{N}(\mathbf{x} | \mathbf{x}_0, \sigma_t^2) \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2 \right] =$$

We can use this property $\int_{\mathbb{R}^n} dx \delta(x - y) g(x) = g(y)$

Expand the expectations:

$$L(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(x_0)} \mathbb{E}_{q_t(x|x_0)} \|x_0 - f_\theta(x, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t|0,1) \frac{1}{N} \sum_{j=1}^N \int_{\mathbb{R}^n} dx_0 \delta(x_0 - x_0^j) \left[\int_{\mathbb{R}^n} dx \mathcal{N}(x|x_0, \sigma_t^2) \|x_0 - f_\theta(x, t)\|_2^2 \right] =$$

We can use this property $\int_{\mathbb{R}^n} dx \delta(x - y)g(x) = g(y)$

$$= \int_0^1 dt \mathcal{U}(t|0,1) \frac{1}{N} \sum_{j=1}^N \int_{\mathbb{R}^n} dx \mathcal{N}(x|x_0^j, \sigma_t^2) \|x_0^j - f_\theta(x, t)\|_2^2 =$$

$$= \int_0^1 dt \mathcal{U}(t|0,1) \frac{1}{N} \sum_{j=1}^N \int_{\mathbb{R}^n} dx \mathcal{N}\left(x|x_0^j, \sigma_t^2\right) \|x_0^j - f_\theta(x, t)\|_2^2 =$$

$$= \int_0^1 dt \int_{\mathbb{R}^n} \frac{1}{N} \sum_{j=1}^N \mathcal{N}\left(x|x_0^j, \sigma_t^2\right) \|x_0^j - f_\theta(x, t)\|_2^2 dx =$$

$$= \int_0^1 dt \int_{\mathbb{R}^n} \underline{\mathbf{J}(f_\theta, x, t)} dx$$

We can minimize $L(f_\theta)$ by minimizing $\mathbf{J}(f_\theta, x, t)$ independently for each x, t

$$\nabla_{f_\theta} J(f_\theta) = 0$$

$$\nabla_{f_\theta} J(f_\theta) = 0$$

$$\nabla_{f_\theta} \left(\frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x \mid x_0^j, \sigma_t^2 \right) \|x_0^j - f_\theta(x, t)\|_2^2 \right) =$$

$$\nabla_{f_\theta} J(f_\theta) = 0$$

$$\nabla_{f_\theta} \left(\frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x \mid x_0^j, \sigma_t^2 \right) \|x_0^j - f_\theta(x, t)\|_2^2 \right) = \frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x \mid x_0^j, \sigma_t^2 \right) \nabla_{f_\theta} \|x_0^j - f_\theta(x, t)\|_2^2 =$$

$$\nabla_{f_\theta} J(f_\theta) = 0$$

$$\begin{aligned} \nabla_{f_\theta} \left(\frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) \|x_0^j - f_\theta(x, t)\|_2^2 \right) &= \frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) \nabla_{f_\theta} \|x_0^j - f_\theta(x, t)\|_2^2 = \\ &= -\frac{2}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) (x_0^j - f_\theta(x, t)) = 0 \end{aligned}$$

$$\nabla_{f_\theta} J(f_\theta) = 0$$

$$\nabla_{f_\theta} \left(\frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) \|x_0^j - f_\theta(x, t)\|_2^2 \right) = \frac{1}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) \nabla_{f_\theta} \|x_0^j - f_\theta(x, t)\|_2^2 =$$

$$= -\frac{2}{N} \sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) (x_0^j - f_\theta(x, t)) = 0$$

$$f_\theta(x, t) = \frac{\sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right) x_0^j}{\sum_{j=1}^N \mathcal{N} \left(x | x_0^j, \sigma_t^2 \right)}$$

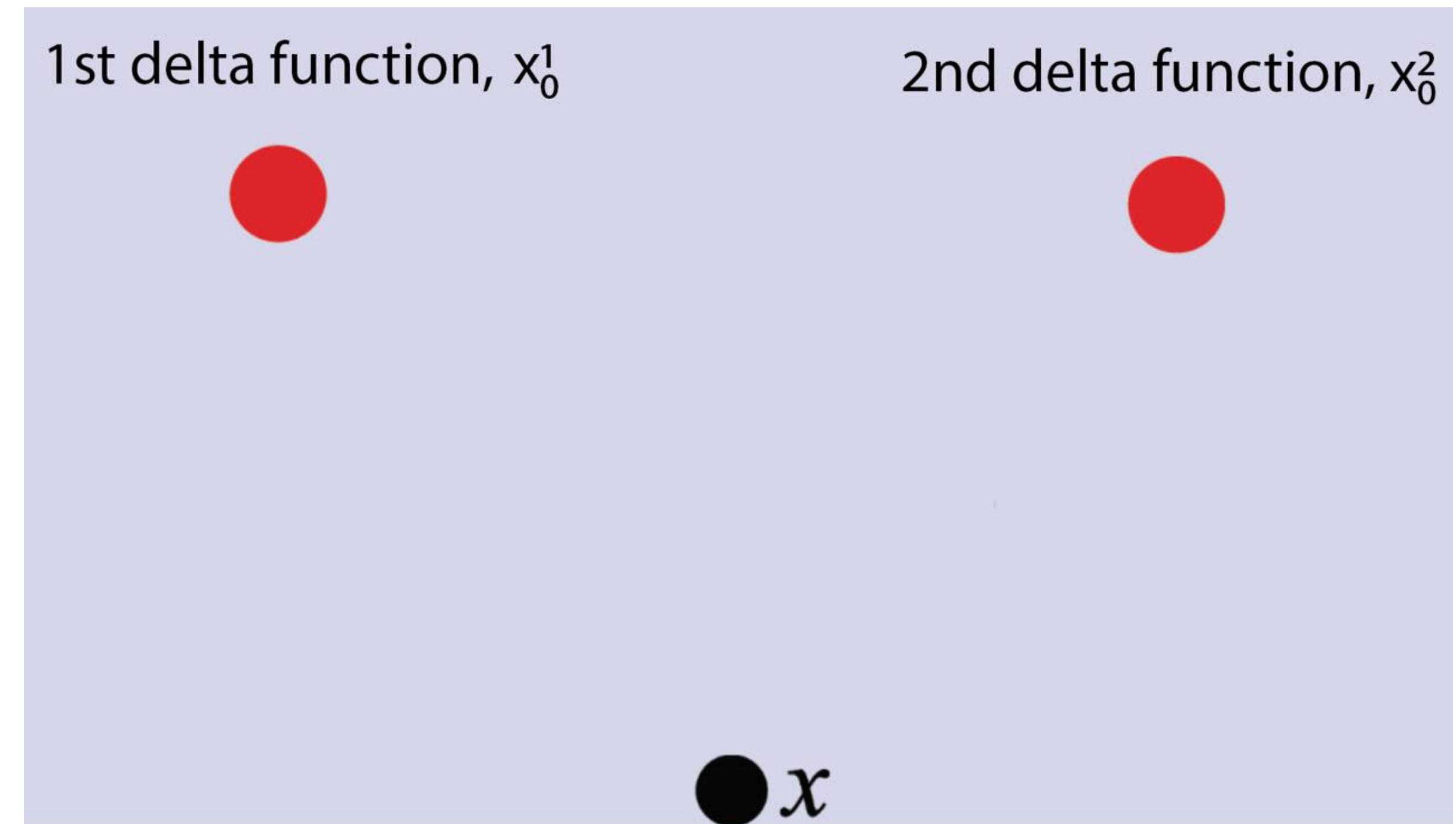
1. Ideal diffusion model

The form of the conditional expectation

$$f_{\theta}(x, t) = \frac{\sum_{j=1}^N \mathcal{N}(x | x_0^j, \sigma_t^2) x_0^j}{\sum_{j=1}^N \mathcal{N}(x | x_0^j, \sigma_t^2)}$$

What if we generate data using this score? Let's consider two delta functions

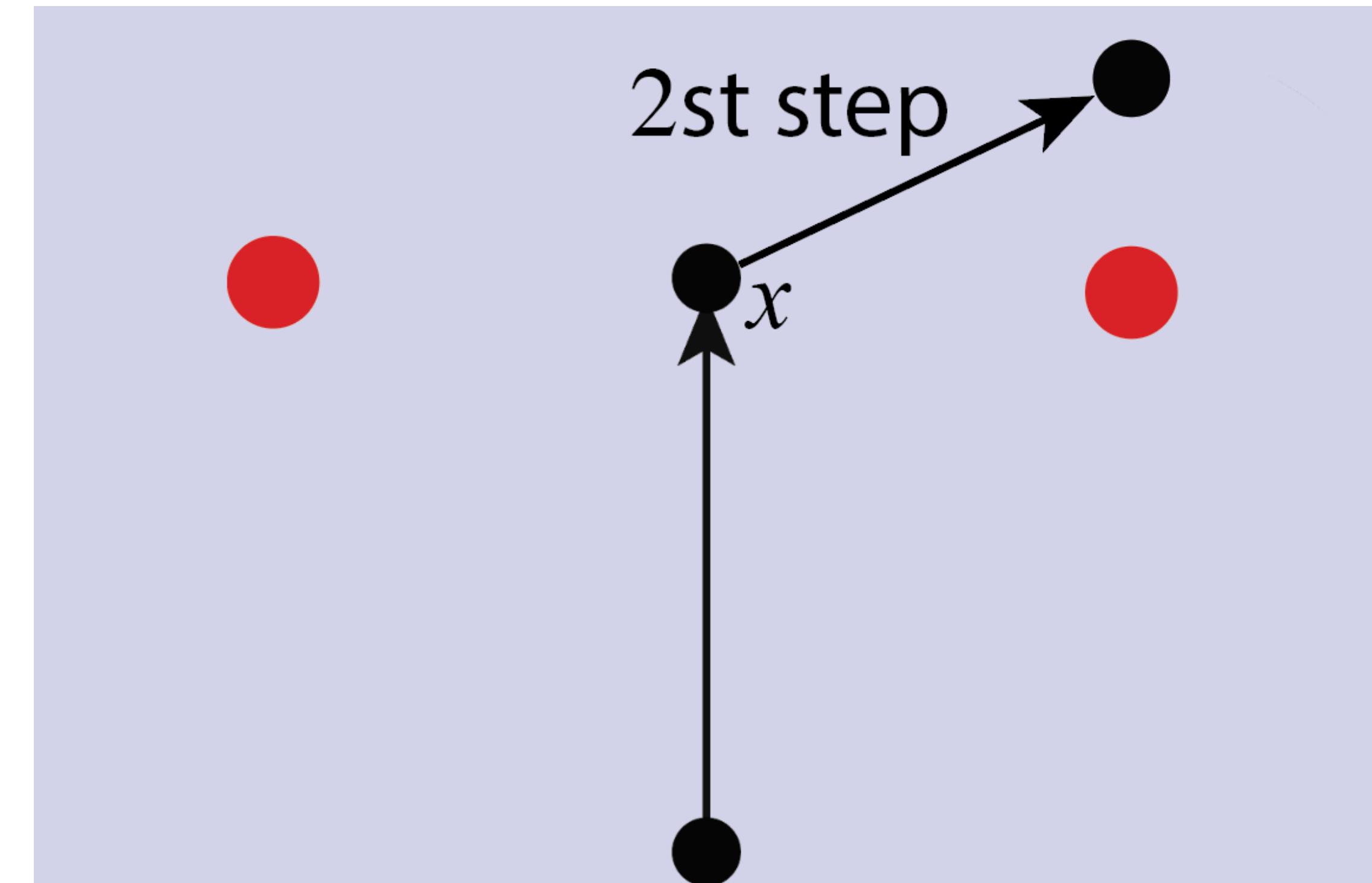
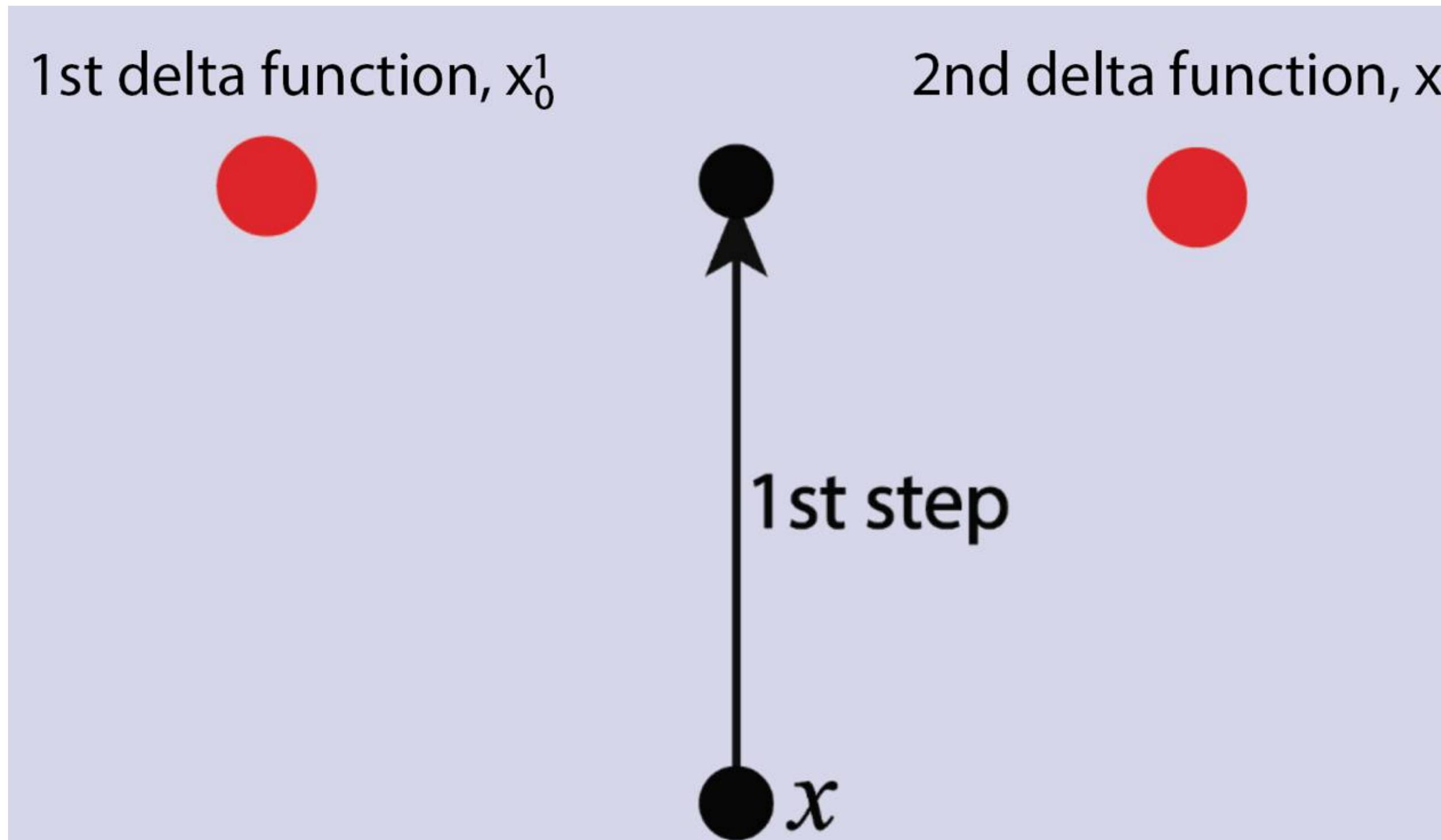
$$f_{\theta}(x, t) = \frac{\mathcal{N}(x | x_0^1, \sigma_t^2) x_0^1 + \mathcal{N}(x | x_0^2, \sigma_t^2) x_0^2}{\mathcal{N}(x | x_0^1, \sigma_t^2) + \mathcal{N}(x | x_0^2, \sigma_t^2)}$$



1. Sampling from ideal diffusion model

$$f_{\theta}(x) = \frac{0.5x_0^1 + 0.5x_0^2}{0.5 + 0.5} = \frac{x_0^1 + x_0^2}{2}$$

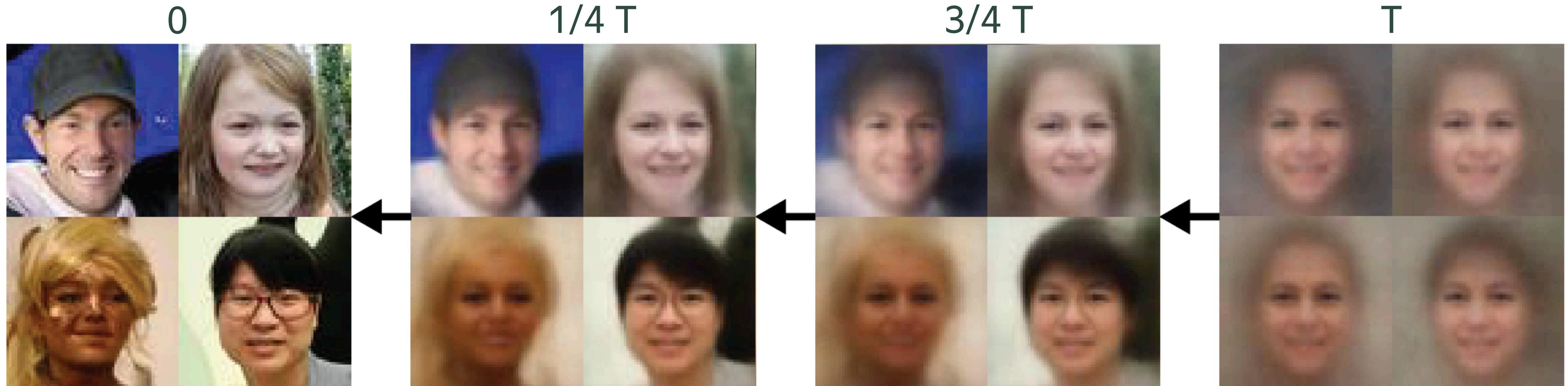
$$f_{\theta}(x) = 0.3x_0^1 + 0.7x_0^2$$



Even in the simplest case, DM is not suitable for obtaining samples in a few steps

1. Sampling from ideal diffusion model

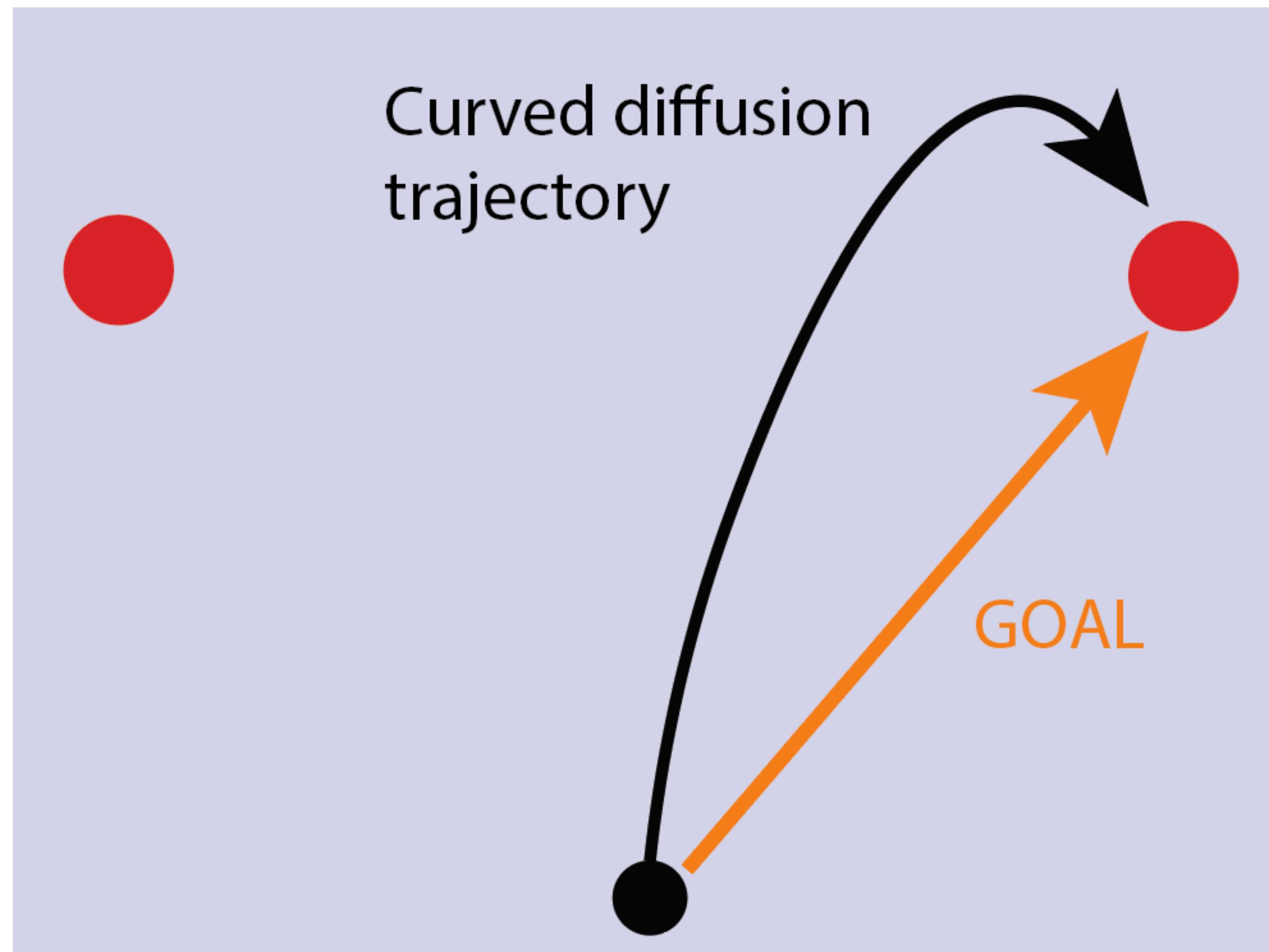
What about images?



Images become less and less averaged

1. DMs speed are limited!

Diffusion has curved trajectories



2. Moving towards fast DMs

Why DMs are slow? → The optimal solution is an averaged sample → Curved trajectories

No fast sampling with DMs :(But we want it!

1. Flow based approaches

Generalized PF-ODEs to obtain
more straightened trajectories

- Flow matching
- Rectified flows
- ...

2. Distillation (diffusion → few-step model)

Reformulate the problem to obtain

$$f_\phi(\mathbf{x}, t) = \mathbf{x}_0$$

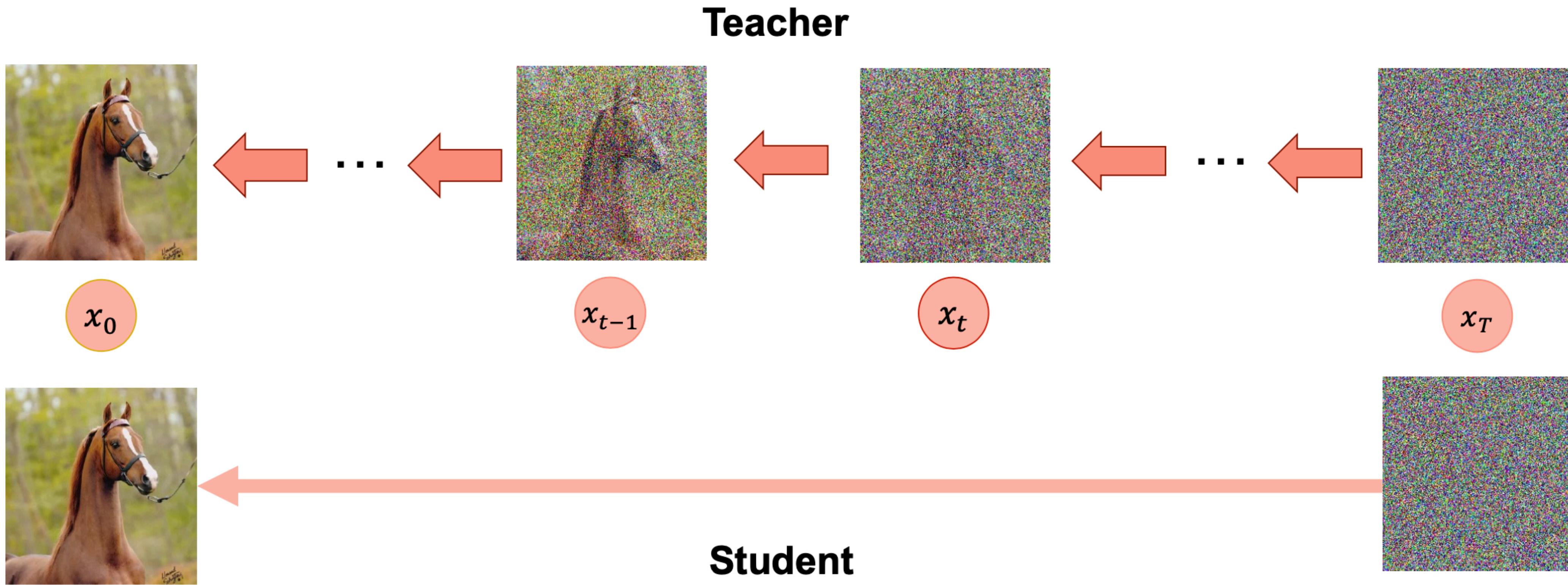
ODE-based

- Knowledge
- Consistency
- Progressive
- ...

ODE-free

- DMD
- Adversarial
- ...

2. Knowledge distillation



The main idea is simple:

- collect the data using diffusion, $\{(x_T^1, x_0^1), \dots, (x_T^N, x_0^N)\}$
- fine-tune the network, $f_\phi(x, t)$, on these data, $f_\phi(x_T) \rightarrow x_0$

2. Knowledge distillation

$$L_{KD}(f_\phi) = \mathbb{E}_{\mathcal{N}(\mathbf{x}|0,\sigma^2)} \|\hat{\mathbf{x}}_0(\mathbf{x}, \theta) - f_\phi(\mathbf{x}, t)\|_2^2$$

$\hat{\mathbf{x}}_0(\mathbf{x}, \theta)$ – sample from the pretrained DM using solver

Looks quite similar to the diffusion loss

$$L_{diffusion}(f_\theta) = \mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_t(\mathbf{x}|\mathbf{x}_0)} \|\mathbf{x}_0 - f_\theta(\mathbf{x}, t)\|_2^2$$

But the optimum is different!

2. Knowledge distillation. Optimal solution

$$\nabla_{f_\phi} \mathcal{L}_{KD}(f_\phi) = - ?$$

$$\mathcal{L}_{KD}(f_\phi) = \mathbb{E}_{\mathcal{N}(x|0,\sigma^2)} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 =$$

2. Knowledge distillation. Optimal solution

$$\nabla_{f_\phi} \mathcal{L}_{KD}(f_\phi) = - ?$$

$$\mathcal{L}_{KD}(f_\phi) = \mathbb{E}_{\mathcal{N}(x|0,\sigma^2)} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 = \int_{\mathbb{R}^n} \mathcal{N}(x | 0, \sigma^2) \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 dx,$$

2. Knowledge distillation. Optimal solution

$$\nabla_{f_\phi} \text{L}_{KD}(f_\phi) = - ?$$

$$\text{L}_{KD}(f_\phi) = \mathbb{E}_{\mathcal{N}(x|0,\sigma^2)} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 = \int_{\mathbb{R}^n} \mathcal{N}(x|0,\sigma^2) \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 dx,$$

$$\int_{\mathbb{R}^n} \mathcal{N}(x|0,\sigma^2) \nabla_{f_\phi} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 dx = 0 \rightarrow \nabla_{f_\phi} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 = 0$$

2. Knowledge distillation. Optimal solution

$$\nabla_{f_\phi} \text{L}_{KD}(f_\phi) = - ?$$

$$\text{L}_{KD}(f_\phi) = \mathbb{E}_{\mathcal{N}(x|0,\sigma^2)} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 = \int_{\mathbb{R}^n} \mathcal{N}(x|0,\sigma^2) \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 dx,$$

$$\int_{\mathbb{R}^n} \mathcal{N}(x|0,\sigma^2) \nabla_{f_\phi} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 dx = 0 \rightarrow \nabla_{f_\phi} \|\hat{x}_0(x, \theta) - f_\phi(x, t)\|_2^2 = 0$$

$$f_\phi(x, t) = \hat{x}_0(x, \theta)$$

$$f_\theta(x, t) = \frac{\sum_{j=1}^N \mathcal{N}(x|x_0^j, \sigma_t^2) x_0^j}{\sum_{j=1}^N \mathcal{N}(x|x_0^j, \sigma_t^2)}$$

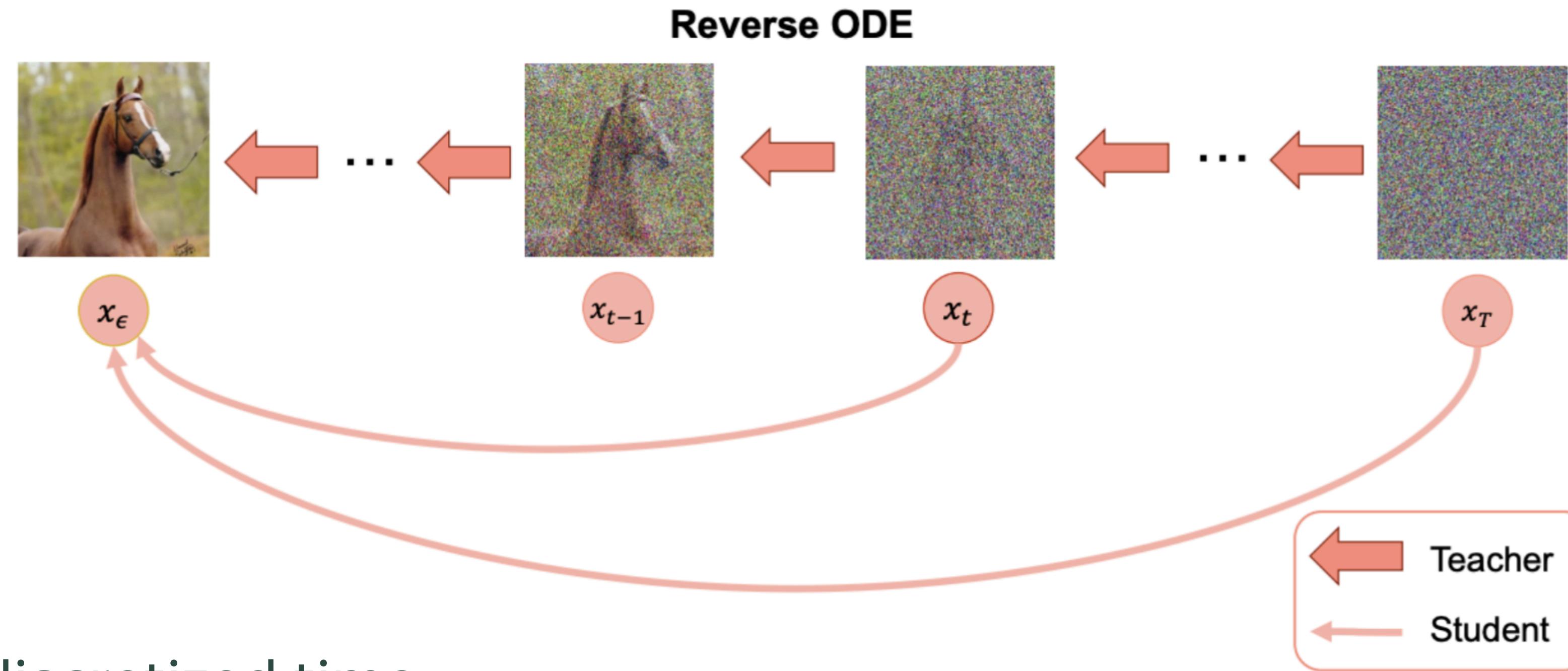
Distillation general principle

The form of different distillations may be different, but the principle remains the same.

The network should predict x_0 instead of an averaged sample!

3. Consistency Distillation. Main idea

We have to sample from DM in the KD. Consistency Distillation aims to avoid this.



$\{t_1, \dots, t_N\}$ – discretized time

$f_\phi(x_t, t) \rightarrow x_0, \forall x_t$ from the same ODE trajectory

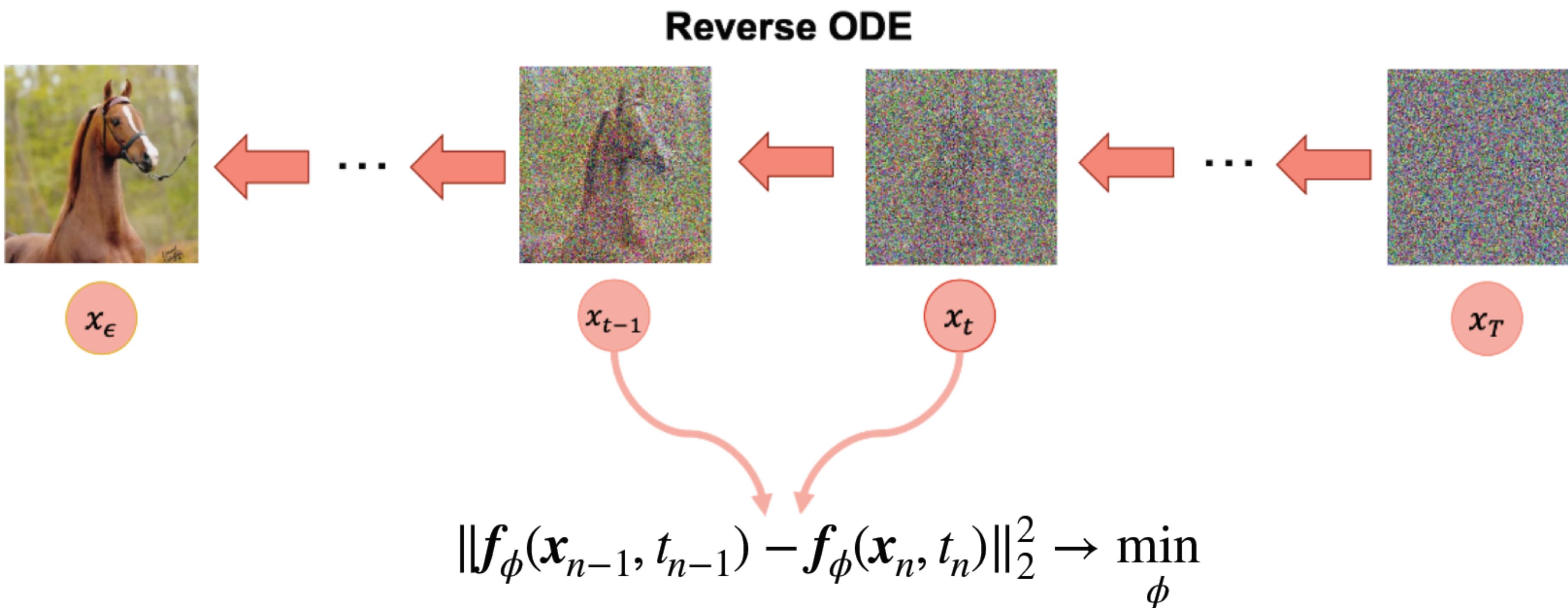
$$f_\phi(x_1, t_1) = f_\phi(x_2, t_2) = \dots = f_\phi(x_N, t_N) = x_0$$

3. Consistency Distillation. Main idea

To obtain such a model, enforce the self-consistency property

$$\|f_{\phi}(x_{n-1}, t_{n-1}) - f_{\phi}(x_n, t_n)\|_2^2 \rightarrow \min_{\phi}$$

where x_{n-1}, x_n – adjacent points on the ODE trajectory



3. Consistency Distillation. Loss function

How to obtain the adjacent points? $x_n \sim p_{t_n}(x_n | x_0); \hat{x}_{n-1} = \text{Solver}(x_n, t_n, t_{n-1} | \theta)$

$$\mathbb{E}_{n \sim \mathcal{U}(2,N)} \mathbb{E}_{p_{data}(x_0)} \mathbb{E}_{q_{t_n}(x_n | x_0)} \|f_\phi(\hat{x}_{n-1}, t_{n-1}) - f_\phi(x_n, t_n)\|_2^2 \rightarrow \min_{\phi}$$

$\mathcal{U}(2,N)$ – Discrete uniform distribution of time points

$p_{data}(x_0)$ – Data distribution (e.g., images)

$q_{t_n}(x | x_0) = \mathcal{N}(x | x_0, \sigma_{t_n}^2)$ – Transition kernels (how to noise the data)

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\hat{\mathbf{x}}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2,N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

$$\sum_{n=2}^N \frac{1}{N-2}$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{\underline{p_{data}(\mathbf{x}_0)}} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

$$\sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j)$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

$$\sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j) \int_{\mathbb{R}^n} d\mathbf{x}_n \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0, \sigma_t^2) \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

$$\sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j) \int_{\mathbb{R}^n} d\mathbf{x}_n \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0, \sigma_t^2) \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 =$$

$$= \sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_n \frac{1}{M} \sum_{j=1}^M \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0^j, \sigma_{t_n}^2) \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{f_\phi}$$

$$p_{data}(\mathbf{x}_0) = \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j); \{\mathbf{x}_0^j\}_{j=1}^M - \text{our dataset}$$

$$\sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_0 \frac{1}{M} \sum_{j=1}^M \delta(\mathbf{x}_0 - \mathbf{x}_0^j) \int_{\mathbb{R}^n} d\mathbf{x}_n \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0, \sigma_t^2) \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 =$$

$$= \sum_{n=2}^N \frac{1}{N-2} \int_{\mathbb{R}^n} d\mathbf{x}_n \frac{1}{M} \sum_{j=1}^M \mathcal{N}(\mathbf{x}_n | \mathbf{x}_0^j, \sigma_{t_n}^2) \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2$$

$$\mathbf{J}(f_\phi, \mathbf{x}_n, n)$$

3. What about optimal solution?

$$\nabla_{f_\phi} J(f_\phi, x_n, n) = 0$$

$$f_\phi(x_n, t_n) = \frac{\sum_{j=1}^M \left(\mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2) f_\phi(\hat{x}_{n-1}, t_{n-1}) \right)}{\sum_{j=1}^M \mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2)}$$

Hm, again a mixture? NO!

3. What about optimal solution?

$$\nabla_{f_\phi} J(f_\phi, x_n, n) = 0$$

$$f_\phi(x_n, t_n) = \frac{\sum_{j=1}^M \left(\mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2) f_\phi(\hat{x}_{n-1}, t_{n-1}) \right)}{\sum_{j=1}^M \mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2)}$$

Hm, again a mixture? NO!

$x_n \sim p_{t_n}(x_n | x_0^k)$; $\hat{x}_{n-1} = \text{Solver}(x_n, t_n, t_{n-1} | \theta) \rightarrow \hat{x}_{n-1}$ depends only on x_0^k

3. What about optimal solution?

$$\nabla_{f_\phi} J(f_\phi, x_n, n) = 0$$

$$f_\phi(x_n, t_n) = \frac{\sum_{j=1}^M \left(\mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2) f_\phi(\hat{x}_{n-1}, t_{n-1}) \right)}{\sum_{j=1}^M \mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2)}$$

Hm, again a mixture? NO!

$x_n \sim p_{t_n}(x_n | x_0^k)$; $\hat{x}_{n-1} = \text{Solver}(x_n, t_n, t_{n-1} | \theta) \rightarrow \hat{x}_{n-1}$ depends only on x_0^k

$$f_\phi(x_n, t_n) = \frac{f_\phi(\hat{x}_{n-1}, t_{n-1}) \sum_{j=1}^M \mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2)}{\sum_{j=1}^M \mathcal{N}(x_n | x_0^j, \sigma_{t_n}^2)} = f_\phi(\hat{x}_{n-1}, t_{n-1})$$

3. What about optimal solution?

$$\mathbb{E}_{n \sim \mathcal{U}(2, N)} \mathbb{E}_{p_{data}(\mathbf{x}_0)} \mathbb{E}_{q_{t_n}(\mathbf{x}_n | \mathbf{x}_0)} \|f_\phi(\hat{\mathbf{x}}_{n-1}, t_{n-1}) - f_\phi(\mathbf{x}_n, t_n)\|_2^2 \rightarrow \min_{\phi}$$

$$f_\phi(\mathbf{x}_N, t_N) = f_\phi(\mathbf{x}_{N-1}, t_{N-1}) = \dots = f_\phi(\mathbf{x}_1, t_1)$$

But where \mathbf{x}_0 ? A possible solution:

$$f_\phi(\mathbf{x}_N, t_N) = f_\phi(\mathbf{x}_{N-1}, t_{N-1}) = \dots = f_\phi(\mathbf{x}_1, t_1) = 0$$

Boundary condition:

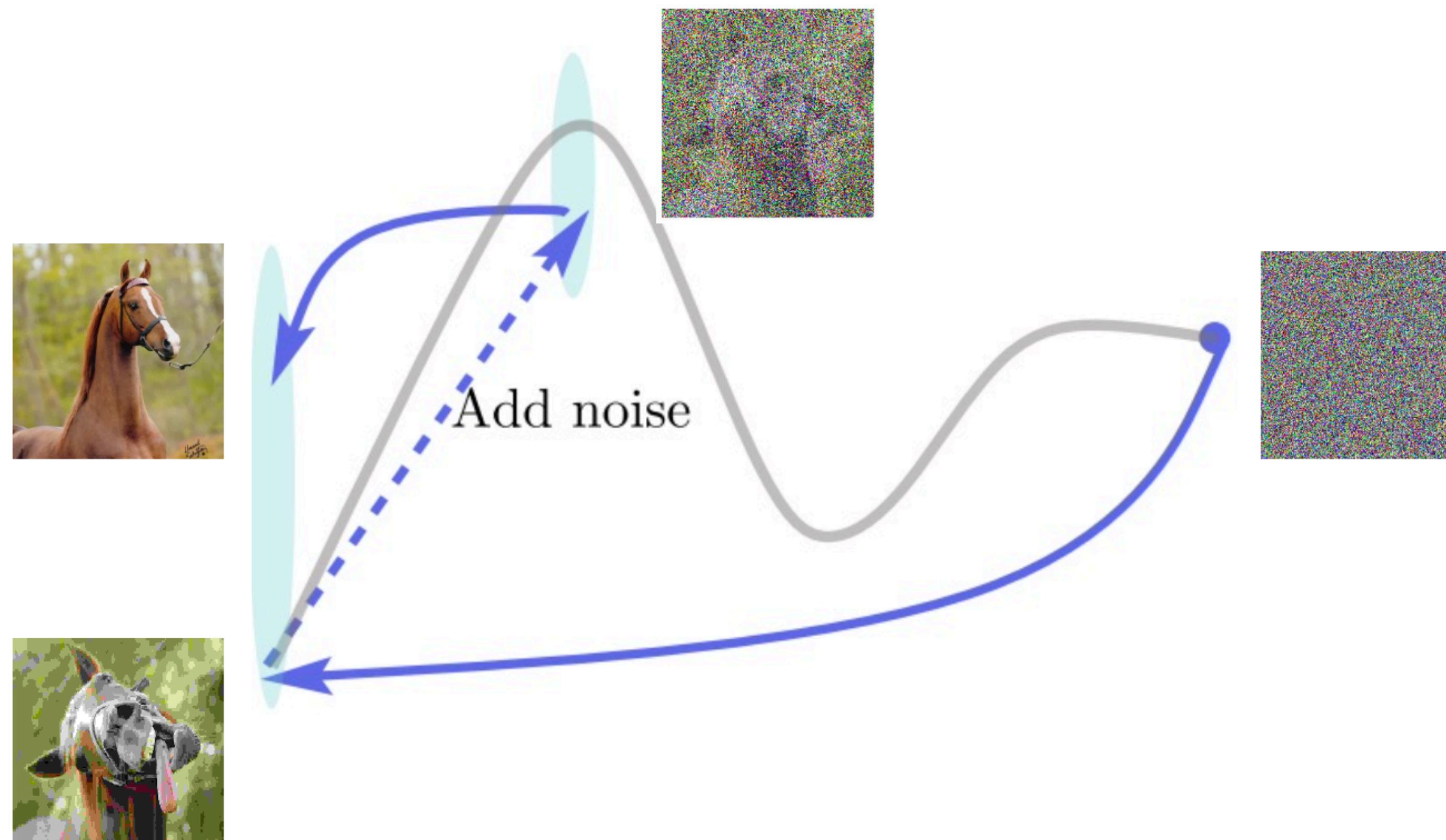
$$f_\phi(\mathbf{x}_1, t_1) = \mathbf{x}_0 - \text{sample from the dataset}$$

$$f_\phi(\mathbf{x}_N, t_N) = f_\phi(\mathbf{x}_{N-1}, t_{N-1}) = \dots = f_\phi(\mathbf{x}_1, t_1) = \mathbf{x}_0$$

3. Sampling from CD

$f_\phi(x_n, t_n)$ is not a score function anymore
→ we cannot use solvers to generate data

Stochastic multistep sampling:

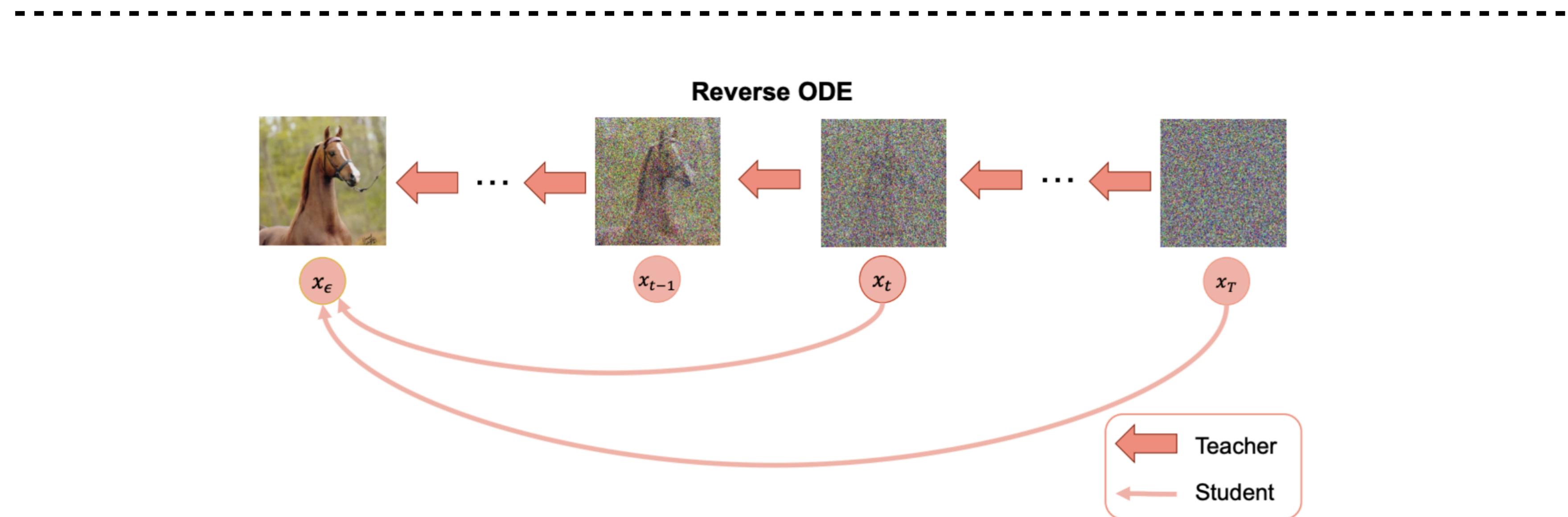
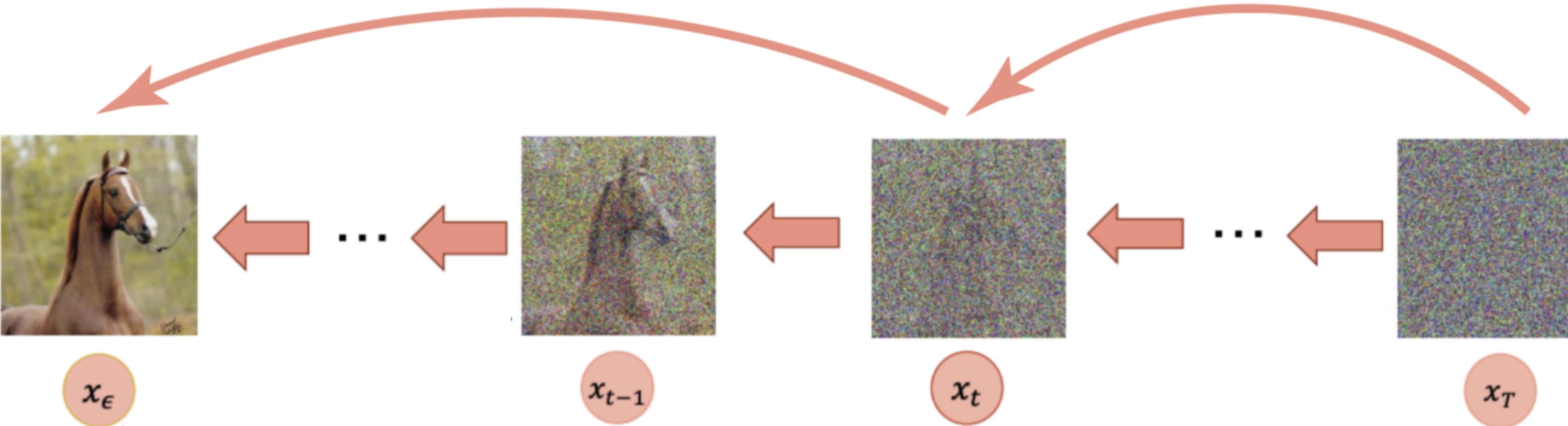


We need about 4-5 steps.

Still works bad for one step
(but ODE-free approaches work better)

3. Multiboundary CD

Let's split the solution interval

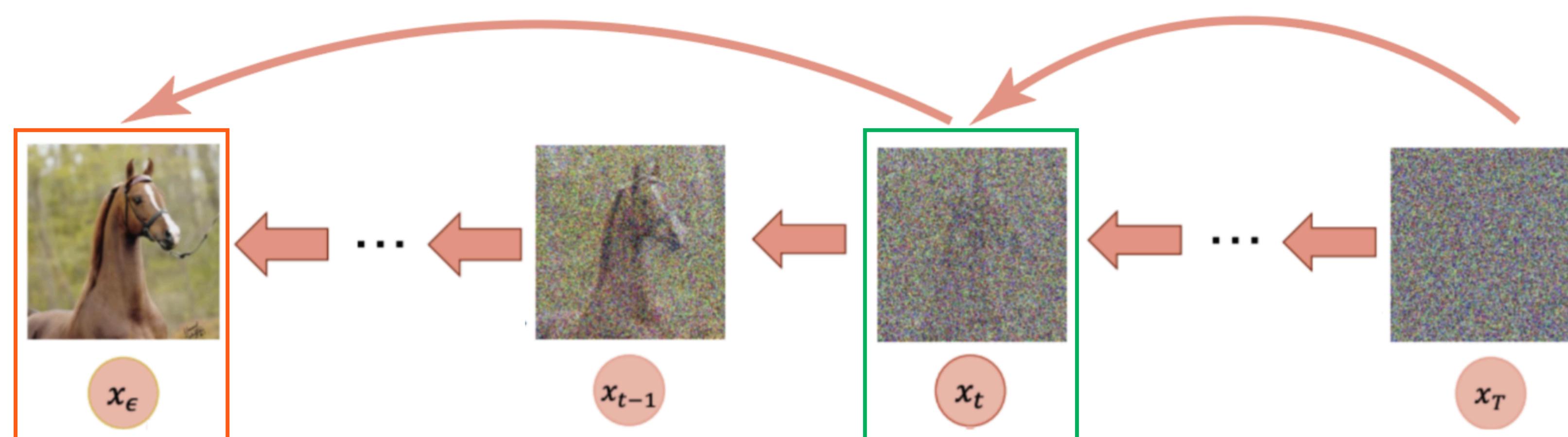


3. Multiboundary CD

For example, we can consider 2-boundary consistency distillation

$$f_{\phi}(x_1, t_1) = \dots = f_{\phi}(x_n, t_n) = x_0 - \text{first boundary}$$

$$f_{\phi}(x_{n+1}, t_{n+1}) = \dots = f_{\phi}(x_N, t_N) = x_{t_n} - \text{second boundary}$$



$$1. x_N \sim \mathcal{N}(x_N | 0, \sigma^2), \quad x_n = f_{\phi}(x_N, t_N); \quad 2. x_0 = f_{\phi}(x_n, t_n)$$

Summary

01

Why DMs are slow and how
to accelerate them?

DM trajectories are curved
Diffusion distillation
methods are needed

02

Knowledge Distillation

Requires a lot of
synthetic samples and
inherits discretization
errors

03

Consistency Distillation

Does not require
synthetic data
Does not inherit
discretization errors
But the quality could
be better even for 5
steps

04

Multi-boundary
Consistency Distillation

Looks like a good
solution, but still
struggles with one step

Ours, 4 steps



Teacher, 50 steps



Ours, 4 steps



Teacher, 50 steps



A donkey in a clown costume giving a lecture at the front of a lecture hall. The blackboard has mathematical equations on it.

an airplane taxiing on a runway with the sun behind it



A castle made of tortilla chips, in a river made of salsa. There are tiny burritos walking around the castle



A raccoon wearing formal clothes, wearing a tophat and holding a cane. The raccoon is holding a garbage bag. Oil painting in the style of cubism.