

Diffusion Model 2

CSE 849 Deep Learning
Spring 2025

Zijun Cui

Content

- Diffusion Model Basics
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- Variants of Diffusion Models
 - Denoising Diffusion Implicit Model (DDIM)
 - Latent Diffusion Models
 - Conditional Diffusion Models
- Applications of Diffusion Models

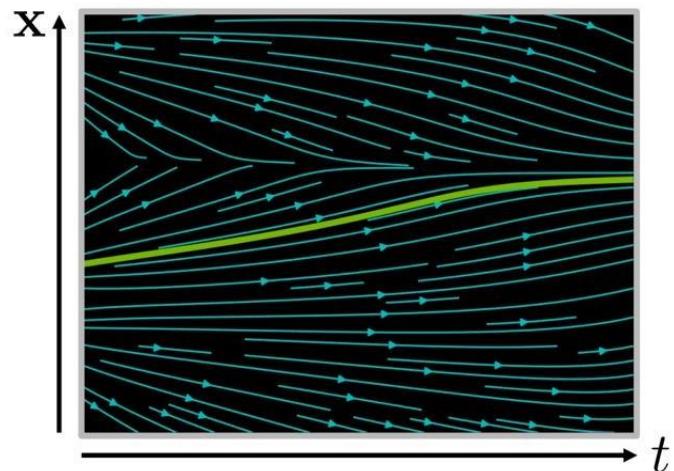
Why SDEs?

- A unified framework for interpreting diffusion models and score-based generation models
 - Variants of diffusion-based and flow-based models

Stochastic Differential Equations

Ordinary Differential Equation (ODE):

$$\frac{dx}{dt} = f(x, t) \text{ or } dx = f(x, t)dt$$



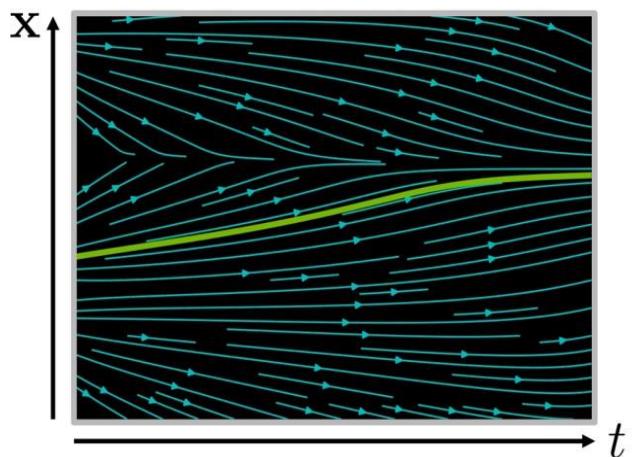
Analytical Solution: $x(t) = x(0) + \int_0^t f(x, \tau)d\tau$

Iterative Numerical $x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t$

Stochastic Differential Equations

Ordinary Differential Equation (ODE):

$$\frac{dx}{dt} = f(x, t) \text{ or } dx = f(x, t)dt$$



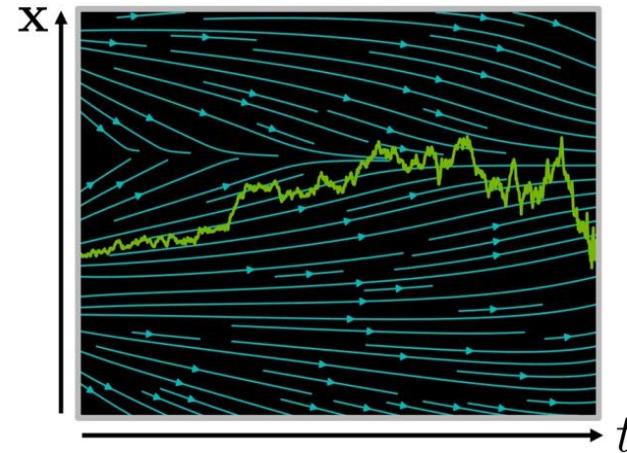
Analytical Solution: $x(t) = x(0) + \int_0^t f(x, \tau)d\tau$

Iterative Numerical $x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t$

Stochastic Differential Equation (SDE):

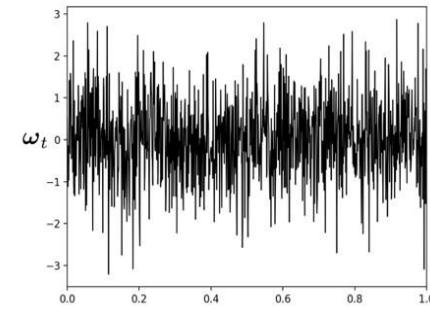
$$\frac{dx}{dt} = \underbrace{f(x, t)}_{\text{drift coefficient}} + \underbrace{\sigma(x, t)\omega_t}_{\text{diffusion coefficient}}$$

$$(dx = f(x, t)dt + \sigma(x, t)d\omega_t)$$



$$x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t + \sigma(x(t), t)\sqrt{\Delta t}\mathcal{N}(0, I)$$

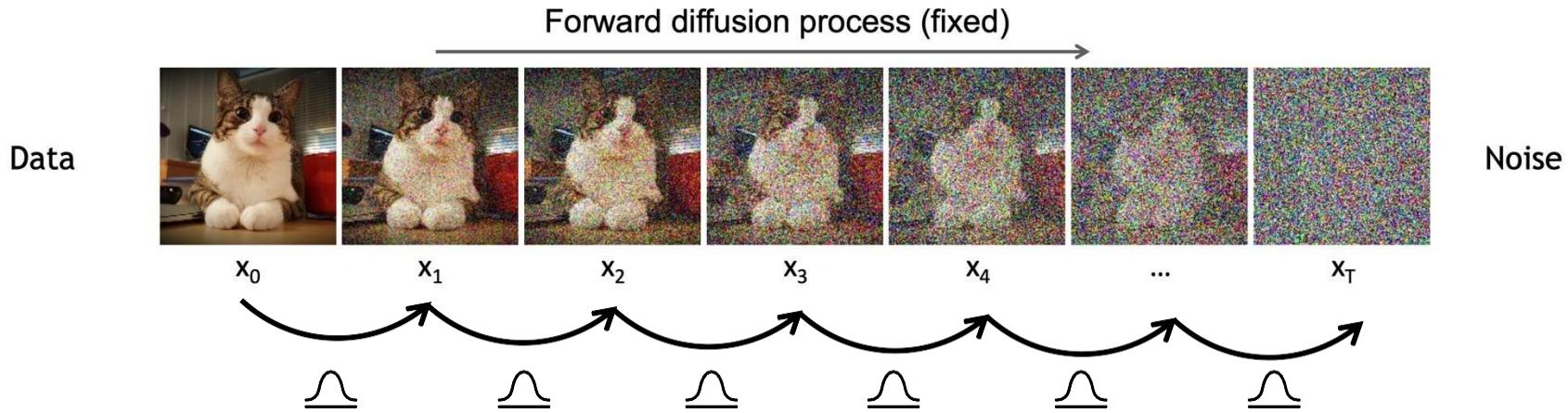
Wiener Process
(Gaussian White Noise)



$$\text{with } \omega_t \approx \frac{1}{\sqrt{\Delta t}} \cdot \mathcal{N}(0, I)$$

Continuous-time diffusion models

Stochastic differential equation framework



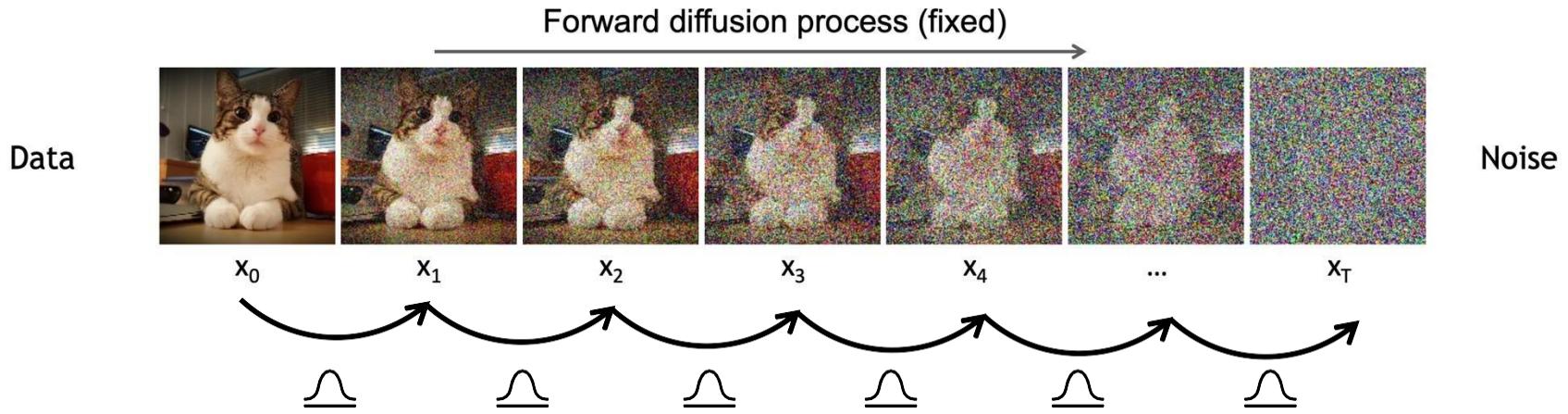
- Consider a forward process with many many small steps (still discrete)

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N} \left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I} \right)$$

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Continuous-time diffusion models

Stochastic differential equation framework



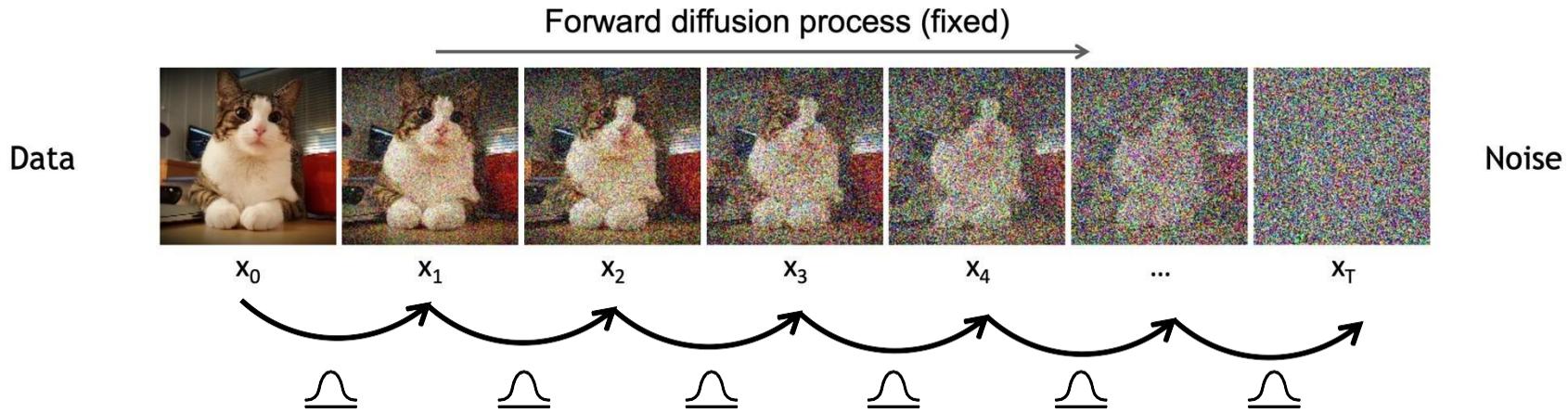
- Consider a forward process with many many small steps

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{x}_{t-1}) &= \mathcal{N} \left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I} \right) \\ \mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \sqrt{1 - \beta(t)\Delta t} \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\beta_t := \beta(t)\Delta t) \end{aligned}$$

Allows different size along t Step size

Continuous-time diffusion models

Stochastic differential equation framework



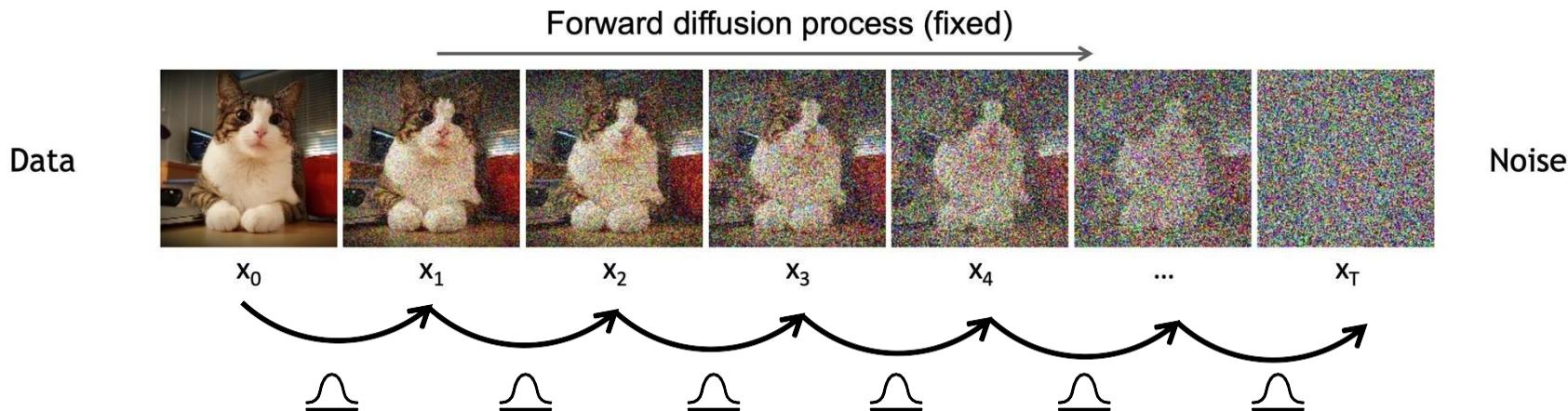
- Consider a forward process with many many small steps

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N} \left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I} \right)$$

$$\begin{aligned} \mathbf{x}_t &= \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \sqrt{1 - \beta(t)\Delta t} \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\beta_t := \beta(t)\Delta t) \\ &\approx \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2} \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\text{Taylor expansion } \sqrt{1 - \beta(t)\Delta t} \approx 1 - \frac{1}{2}\beta(t)\Delta t) \end{aligned}$$

Note: $\sqrt{1 - x} \approx 1 - \frac{1}{2}x$ when $x \ll 1$

Forward Diffusion Process as SDEs



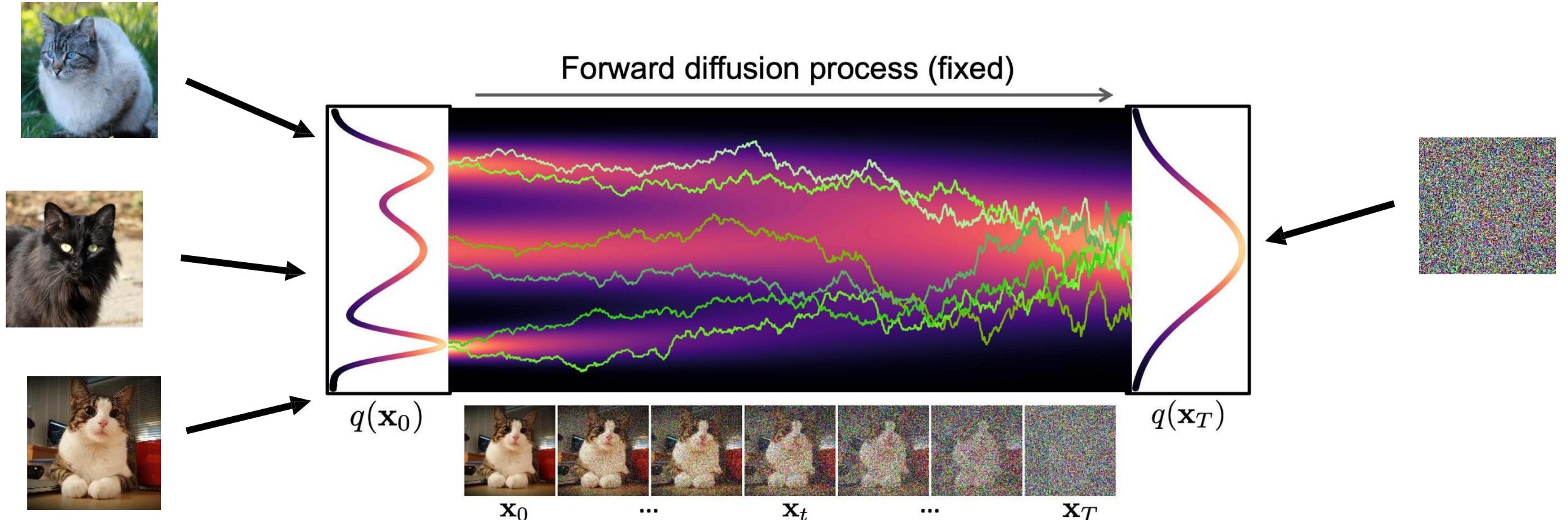
$$\mathbf{x}_t \approx \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2}\mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t dt + \sqrt{\beta(t)}d\omega_t$$

describing the diffusion in infinitesimal limit
i.e., continuous time
Stochastic Differential Equation (SDE)

An iterative update that can be viewed as SDEs

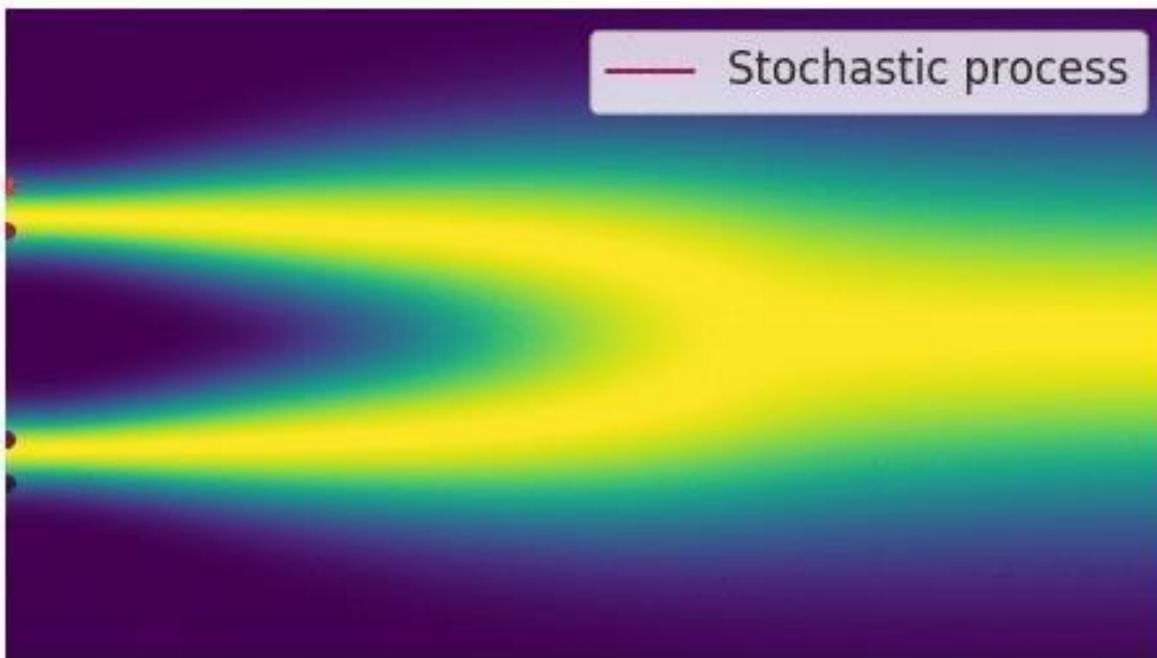
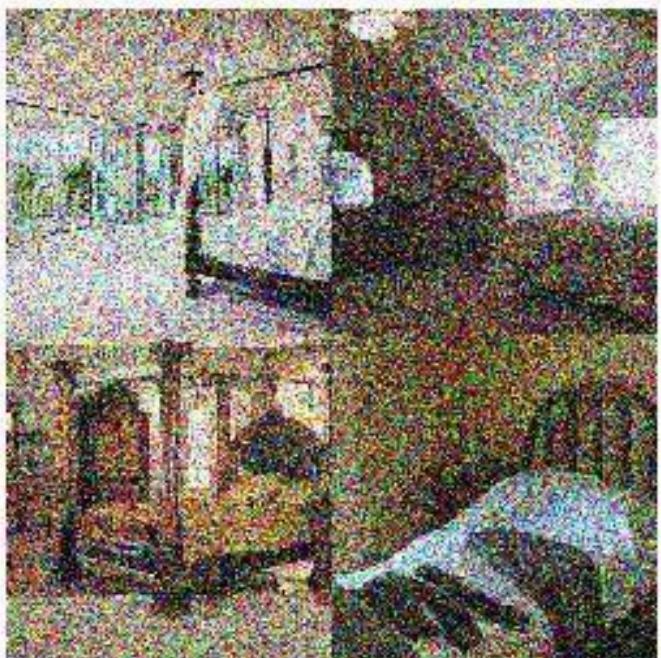
Forward Diffusion Process as SDEs



Forward Diffusion SDE:

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t dt + \sqrt{\beta(t)}d\omega_t$$

Drift Term Diffusion Term
(Pulls toward the mode) (Injects Noise)

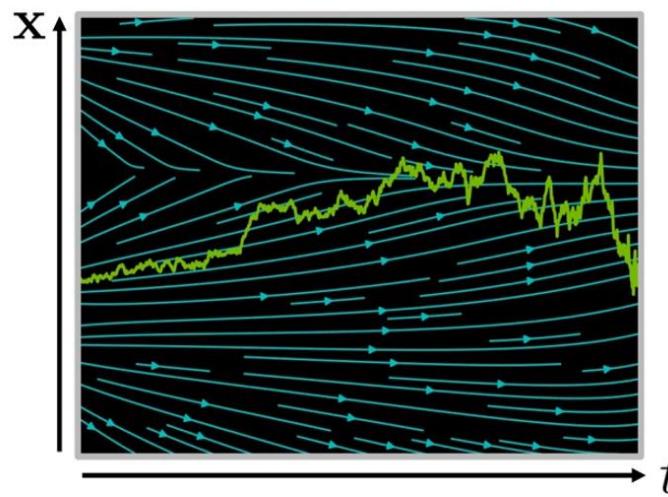


Stochastic Differential Equations

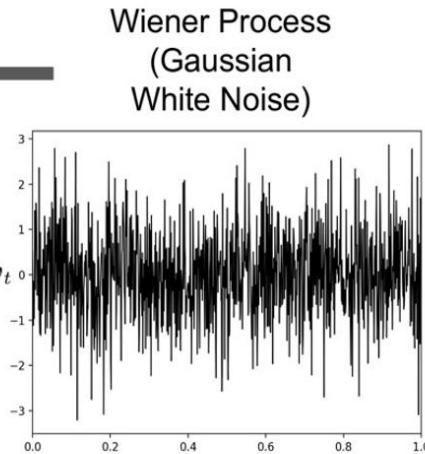
Stochastic Differential Equation (SDE):

$$\frac{dx}{dt} = \underbrace{\mathbf{f}(x, t)}_{\text{drift coefficient}} + \underbrace{\sigma(x, t)\omega_t}_{\text{diffusion coefficient}}$$

$$(dx = \mathbf{f}(x, t)dt + \sigma(x, t)d\omega_t)$$



$$x(t + \Delta t) \approx x(t) + \mathbf{f}(x(t), t)\Delta t + \sigma(x(t), t)\sqrt{\Delta t} \mathcal{N}(\mathbf{0}, \mathbf{I})$$



Given an SDE

$$dx = \mathbf{f}(x, t) dt + \sigma(x, t) d\omega_t$$

We have the **reverse-time SDE**

$$dx = [\mathbf{f}(x, t) - \sigma(x, t)\sigma^\top(x, t)\nabla_x \log p_t(x)] dt + \sigma(x, t) d\bar{\omega}_t$$

- $\nabla_x \log p_t(x)$ is the **score function**
- $d\bar{\omega}_t$ is a **reverse-time Wiener process**
- $\sigma\sigma^\top$ accounts for the diffusion covariance

Stochastic Differential Equations

Given an SDE

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) dt + \boldsymbol{\sigma}(\mathbf{x}, t) d\omega_t$$

We have the **reverse-time SDE**

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - \boldsymbol{\sigma}(\mathbf{x}, t)\boldsymbol{\sigma}^\top(\mathbf{x}, t) \boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}] dt + \boldsymbol{\sigma}(\mathbf{x}, t) d\bar{\omega}_t$$

$p_t(\mathbf{x}) = \text{density of the solution } \mathbf{x}_t \text{ at time } t$

This gradient tells us in which direction should I perturb \mathbf{x}_t to move to high-density

We call it score function, and the reverse drift can be treated as:

$$\text{drift}_{\text{reverse}} = \text{drift}_{\text{forward}} - \text{score-guided correction}$$

A quick detour: Score Function

- General form of probability density function

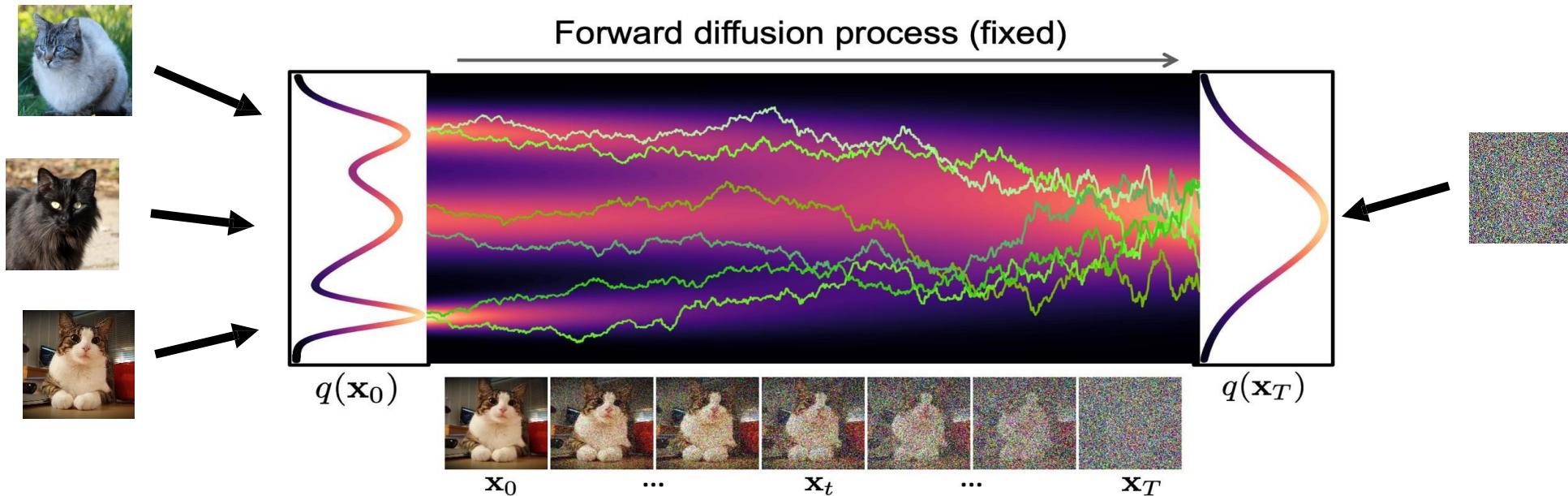
$$p_\theta(\mathbf{x}) = \frac{e^{-f_\theta(\mathbf{x})}}{Z_\theta}$$

- This is an energy-based model (EBM).
- $f_\theta(\mathbf{x})$ is the energy function. $Z_\theta = \int e^{-f_\theta(\mathbf{x})}$ is the partition function, often intractable.
- Log-likelihood is hard: $\log p_\theta(\mathbf{x}) = -f_\theta(\mathbf{x}) - \log Z_\theta$
 - Because of Z_θ
- Instead of modeling $p_\theta(\mathbf{x})$, we can model its gradient $\nabla_{\mathbf{x}} \log p(\mathbf{x})$
 - It avoids calculating Z_θ because:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) = -\nabla_{\mathbf{x}} f_\theta(\mathbf{x})$$

This gradient $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ is called score function

The Generative Reverse Stochastic Differential Equation



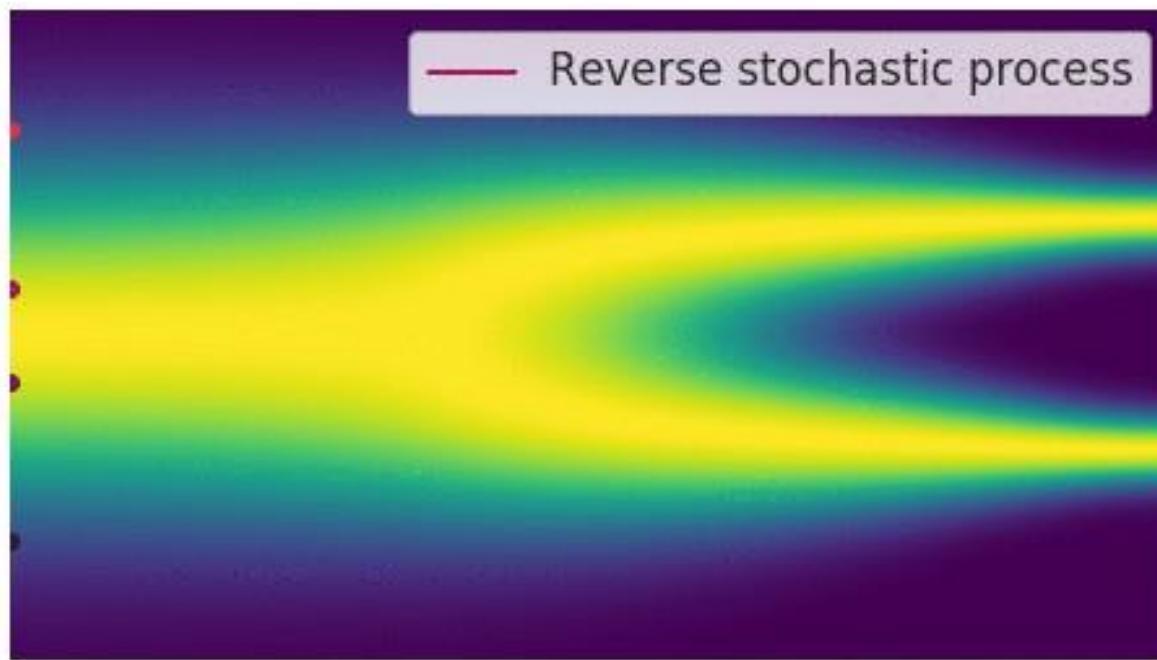
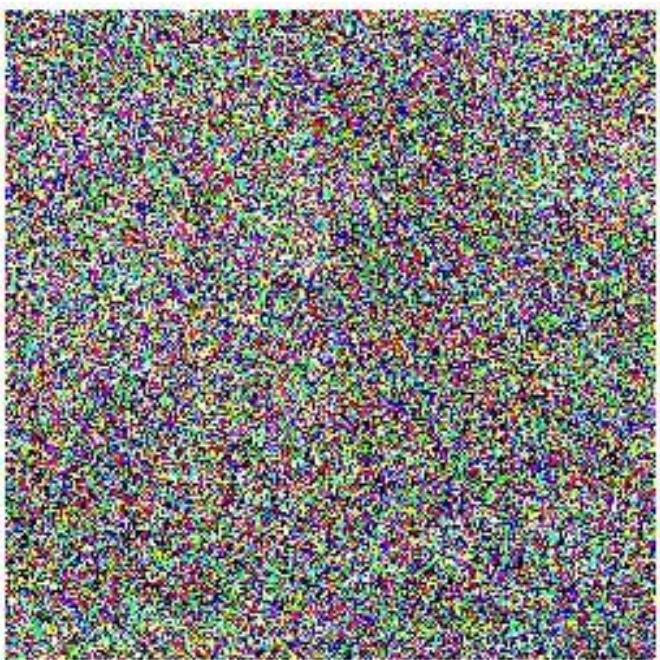
Forward Diffusion SDE:

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t dt + \sqrt{\beta(t)} d\omega_t$$

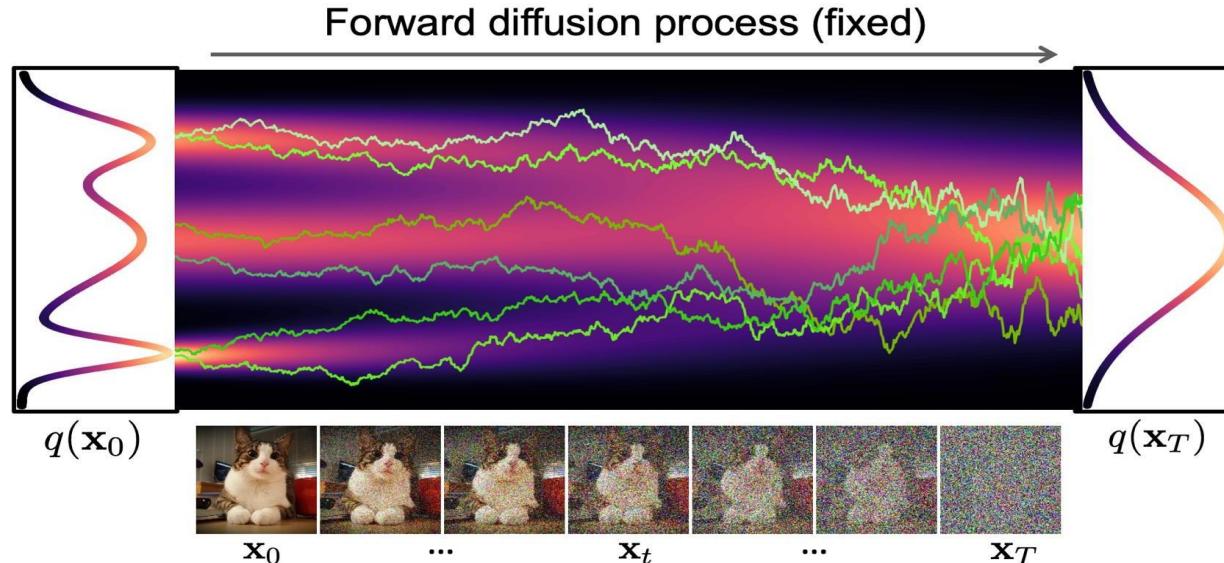
Reverse Generative Diffusion SDE:

$$d\mathbf{x}_t = \underbrace{\left[-\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \right]}_{\text{"Score Function"}} dt + \underbrace{\sqrt{\beta(t)} d\bar{\omega}_t}_{\text{diffusion term}}$$

→ Simulate reverse diffusion process: Data generation from random noise!



Generative Reverse SDEs



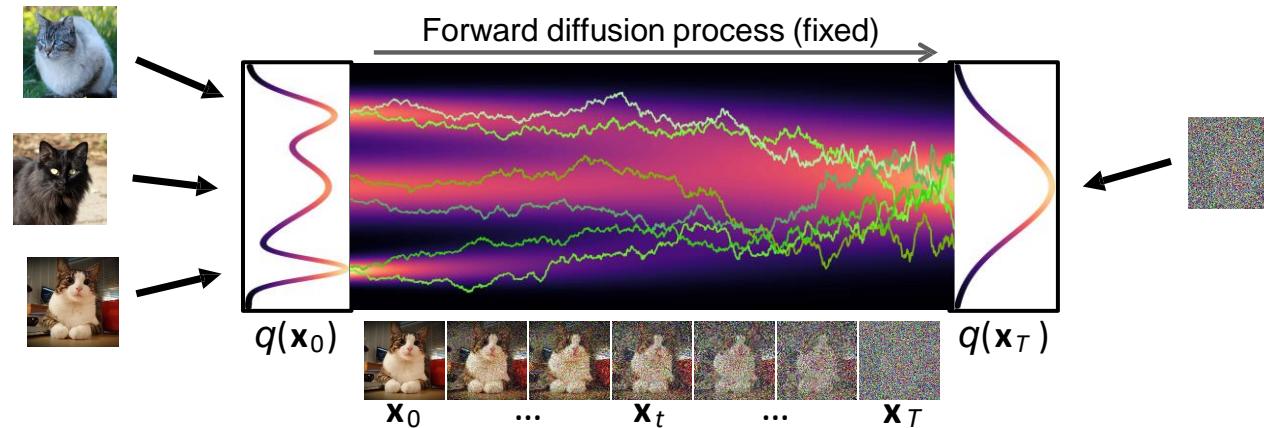
- The forward SDE has a reverse form:

$$d\mathbf{x}_t = \left[-\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t) \boxed{\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)} \right] dt + \sqrt{\beta(t)} d\omega_t$$

How to get the score function?

Score matching is the typical and foundational strategy used in diffusion models to learn the score function

Score Matching

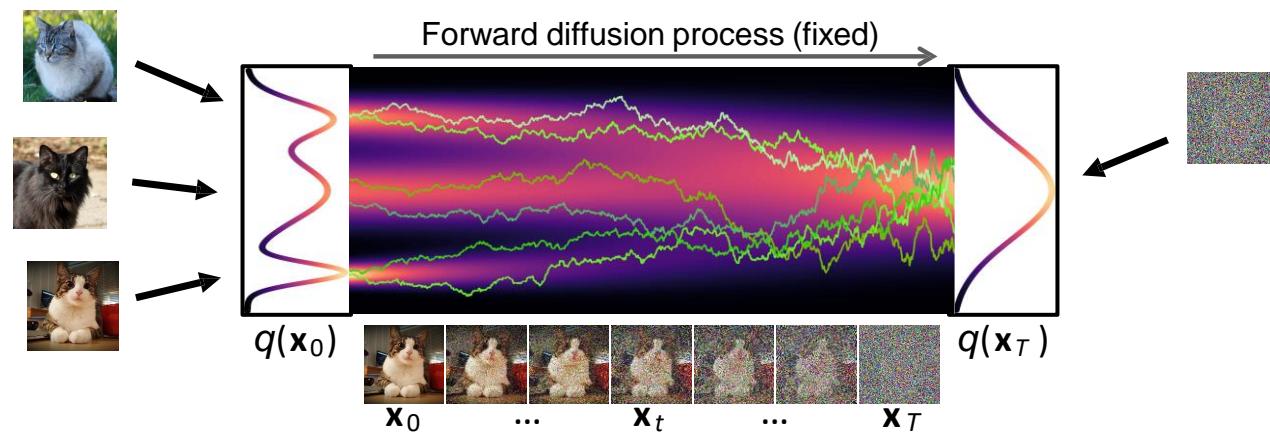


- Naïve idea, learn model for the score function by direct regression?

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t)}}_{\text{diffusion time } t} \underbrace{\tilde{w}(t) \cdot ||\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)||_2^2}_{\begin{array}{l} \text{diffused data } \mathbf{x}_t \\ \text{weighting function} \\ \text{neural network} \\ \text{score of diffused data (marginal)} \end{array}}$$

➡ But $\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$ (**score of the marginal diffused density $q_t(\mathbf{x}_t)$**) is not tractable!

Denoising Score Matching



- Instead, diffuse individual data points \mathbf{x}_0 . Diffused $q_t(\mathbf{x}_t|\mathbf{x}_0)$ **is** tractable!
- **Denoising Score Matching:**

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t|\mathbf{x}_0)}}_{\text{diffusion time } t} \underbrace{\tilde{w}(t) \cdot ||\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t|\mathbf{x}_0)||_2^2}_{\text{score of diffused data sample}}$$

diffusion time t data sample \mathbf{x}_0 diffused data sample \mathbf{x}_t weighting function neural network

→ **After expectations,** $\mathbf{s}_{\theta}(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)!$

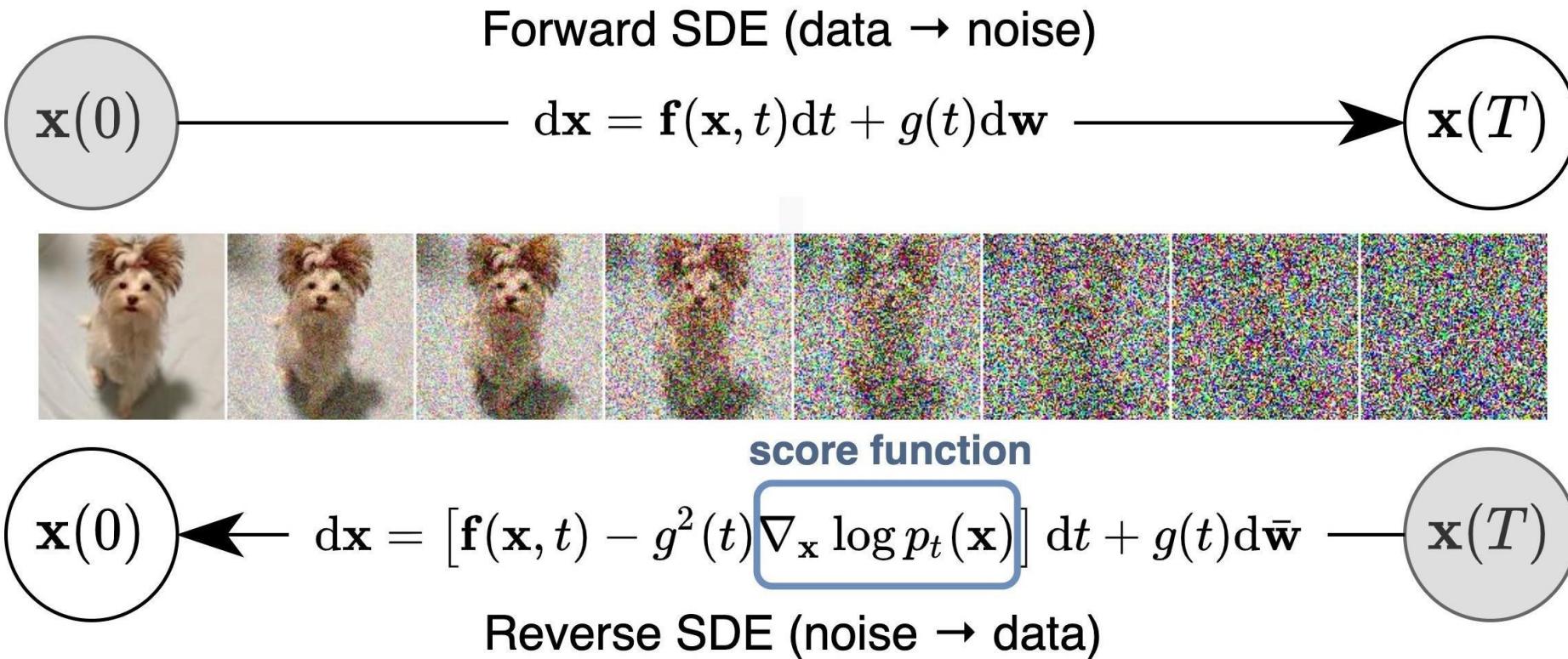
Vincent, in *Neural Computation*, 2011

Song and Ermon, *NeurIPS*, 2019

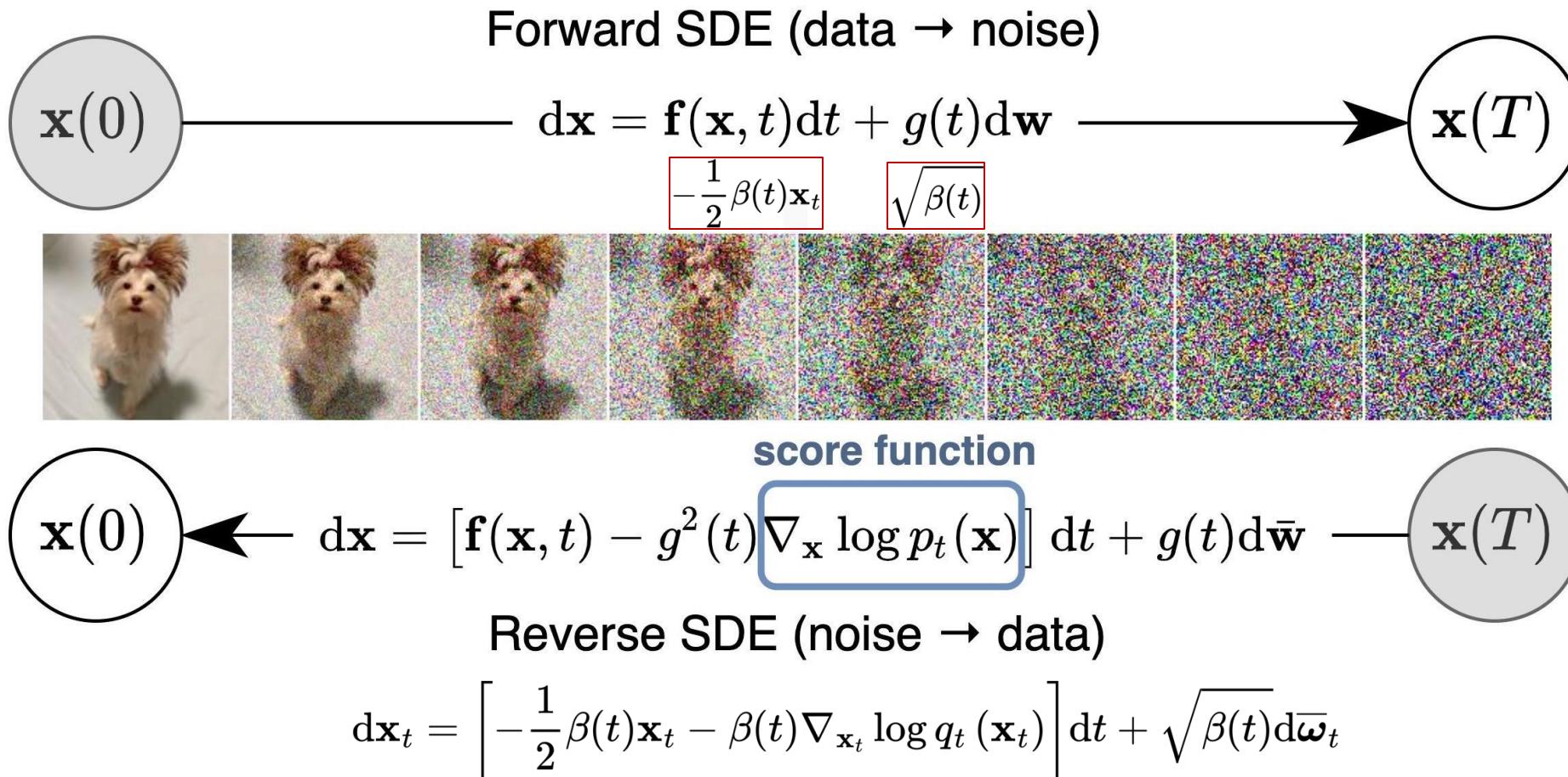
Song et al. *ICLR*, 2021

slide from <https://cvpr2022-tutorial-diffusion-models.github.io/>

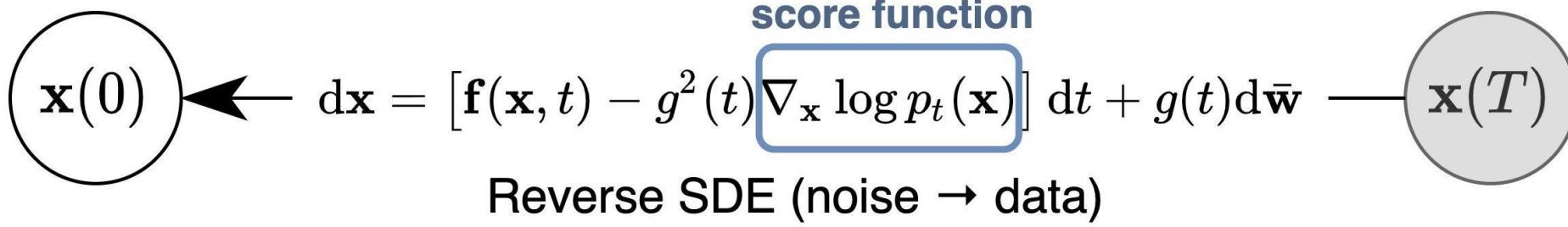
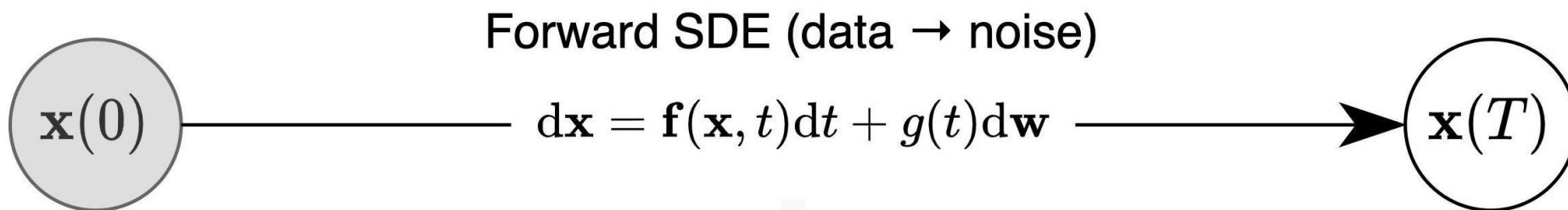
Denoising Score Matching



Denoising Score Matching



Denoising Score Matching



$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)}}_{\text{diffusion time } t} \underbrace{\tilde{w}(t) \cdot}_{\text{weighting function}} \underbrace{\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0)\|_2^2}_{\text{score of diffused data sample}}$$

Denoising Score Matching

- Denoising score matching objective

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)}}_{\substack{\text{diffusion} \\ \text{time } t}} \underbrace{\tilde{w}(t) \cdot}_{\substack{\text{weighting} \\ \text{function}}} \underbrace{\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0)\|_2^2}_{\substack{\text{neural} \\ \text{network} \\ \text{score of diffused} \\ \text{data sample}}}$$

- Re-parametrized sampling:

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Score function:

$$\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0) = -\nabla_{\mathbf{x}_t} \frac{(\mathbf{x}_t - \alpha_t \mathbf{x}_0)^2}{2\sigma_t^2} = -\frac{\mathbf{x}_t - \alpha_t \mathbf{x}_0}{\sigma_t^2} = -\frac{\alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\epsilon} - \alpha_t \mathbf{x}_0}{\sigma_t^2} = -\frac{\boldsymbol{\epsilon}}{\sigma_t}$$

- Denoising network:

$$\mathbf{s}_{\theta}(\mathbf{x}_t, t) := -\frac{\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)}{\sigma_t}$$

- Final objective:

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \hat{w}(t) \cdot \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\|_2^2 \quad \hat{w}(t) = \frac{\tilde{w}(t)}{\sigma_t}$$

What is the ELBO for continuous-time diffusion models?

- Denoising Score Matching:

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \hat{w}(t) \cdot \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t)\|_2^2$$

- [Kingma et al, 2021] showed that the variational upper bound (negative ELBO) can be reduced to a simple variant:

$$-\text{ELBO}(\theta) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} - \frac{d\lambda}{dt} \cdot \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t)\|_2^2 + \text{const}$$

where $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$

$$\lambda = \log \frac{\alpha_t^2}{\sigma_t^2} \quad (\text{signal-to-noise ratio of timestep } t)$$

- In practice, different loss weightings trade off between model with good perceptual quality vs. high log-likelihood:
 - Perceptual quality, e.g., $\hat{w}(t) = 1$
 - High likelihood: -ELBO

→ How to explain such discrepancy?
Can we still trust ELBO / maximum likelihood?

Weighted diffusion objectives

Understanding diffusion objectives as the ELBO with data augmentation

- Summarize different types of diffusion objectives as the following **weighted diffusion objective**:

$$\mathcal{L}_w(\theta) = \frac{1}{2} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[w(t) \cdot -\frac{d\lambda}{dt} \cdot \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2 \right]$$

- [Kingma & Gao, 2023] showed that if $w(t)$ is a monotonically increasing function, then the weighted diffusion objective is equivalent to the **negative ELBO with data augmentation** (Gaussian addition noise):

$$\mathcal{L}_w(\theta) \geq \underbrace{\mathbb{E}_{t \sim p_w(t)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t | \mathbf{x}_0)} [-\log p(\mathbf{x}_t)]}_{\text{Neg. log lik. of noise-perturbed data}} + \text{const}$$

where $p_w(t) \propto \frac{dw(t)}{dt}$

- Takeaway:** By using a certain weighting function $w(t)$ (monotonically increasing) during training, you're essentially performing maximum likelihood estimation on a noise-perturbed data distribution.
- Therefore, the ELBO objective is compatible with perceptual quality when combined with simple data augmentation (additive noise)

Weighted Diffusion Objective

- Denoising score matching objective with loss weighting in general

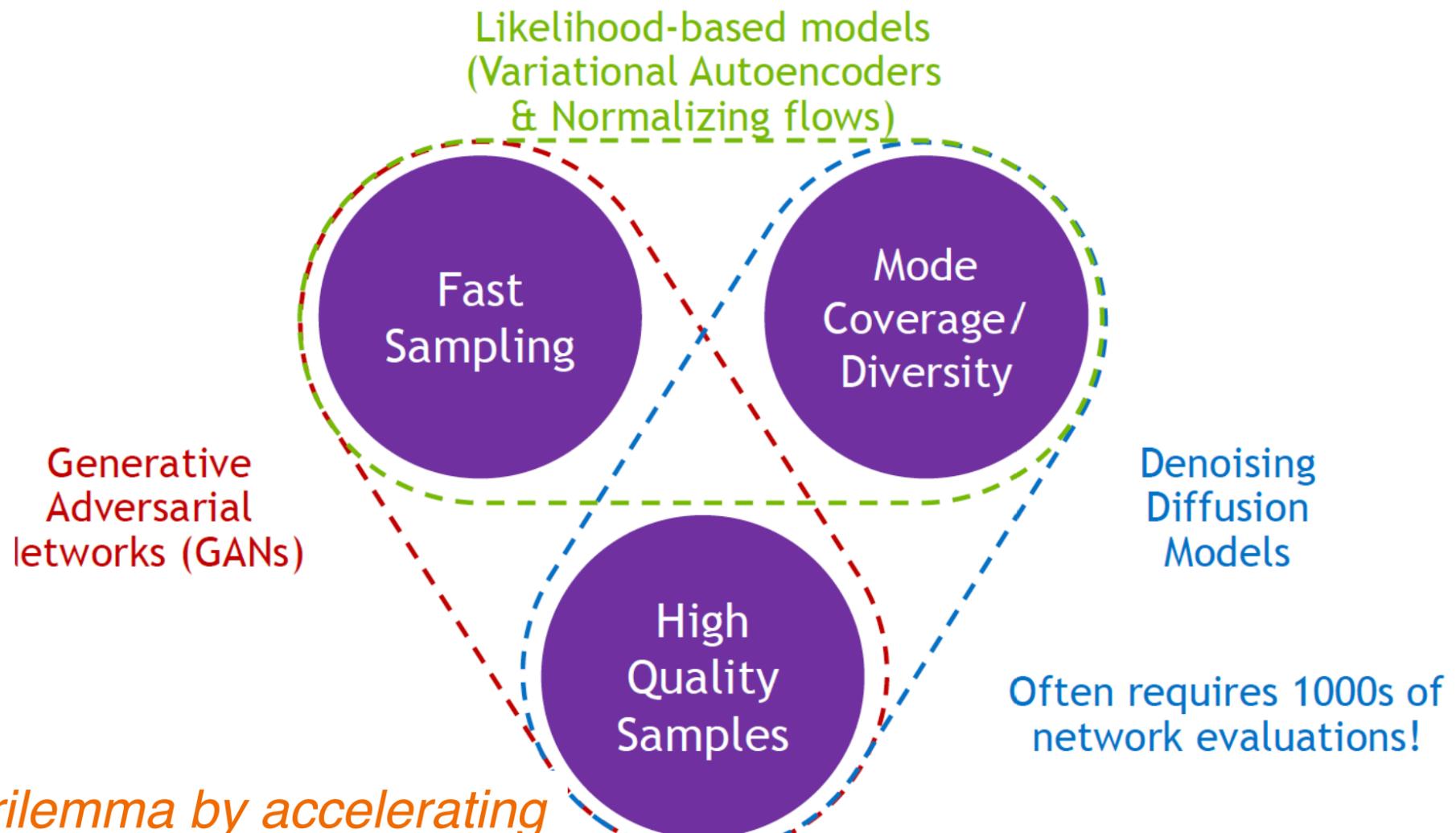
$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_0 \sim q_0(\mathbf{x}_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \frac{\lambda(t)}{\sigma_t^2} \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t)\|_2^2$$

- Loss weights trade-off between
 - good perceptual quality: $\lambda(t) = \sigma_t^2$
 - maximum likelihood: $\lambda(t) = \beta(t)$
- More complicated model parametrization and loss weighting leads to different diffusion model variants in the literature!

Content

- Diffusion Model Basics
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- **Variants of Diffusion Models**
 - Denoising Diffusion Implicit Model (DDIM); Distillation guided Diffusion
 - Latent Diffusion Models
 - Conditional Diffusion Models
- Applications of Diffusion Models

The generative learning trilemma



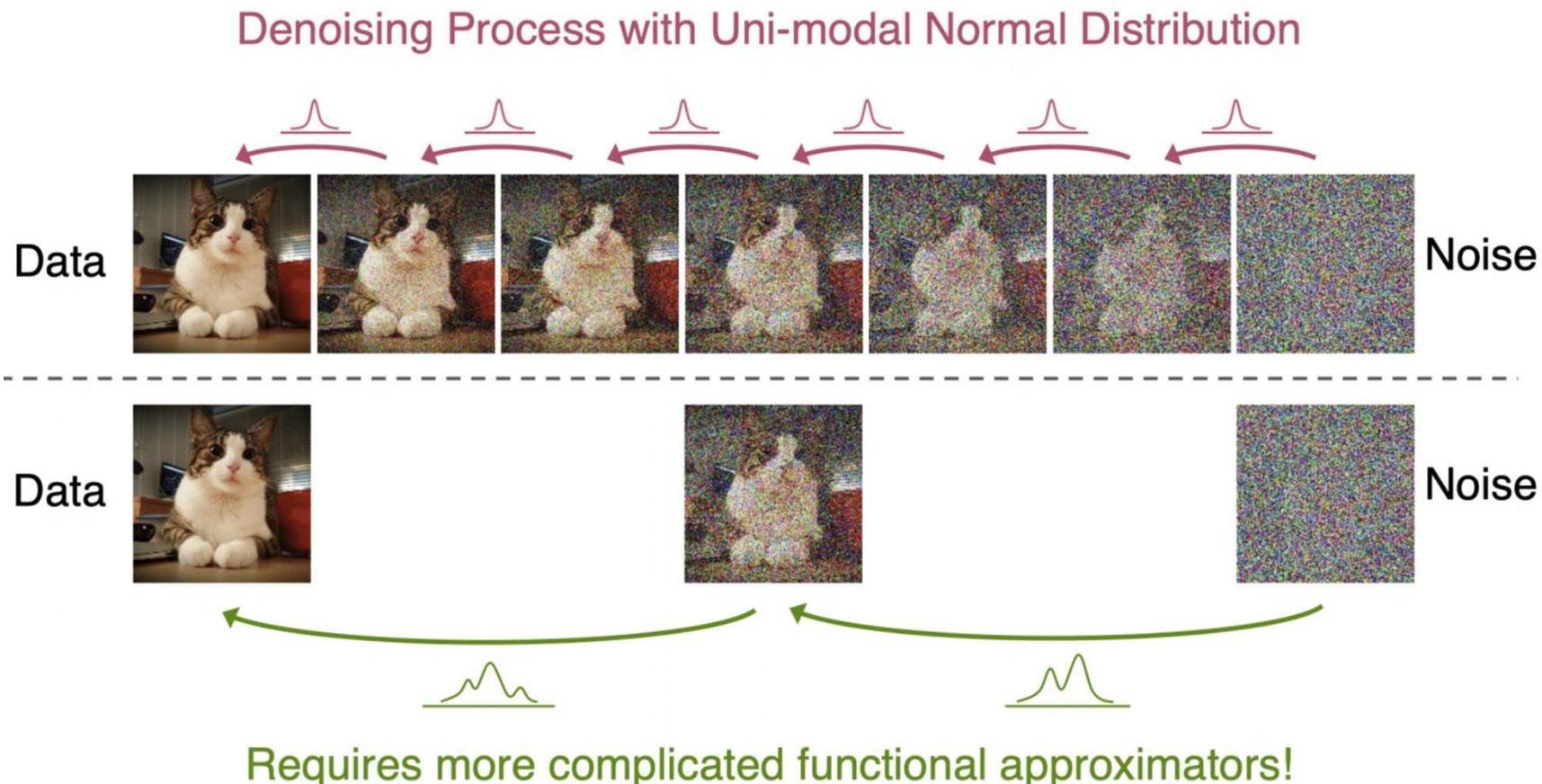
Tackle the trilemma by accelerating diffusion model sampling!

How to make sampling faster?

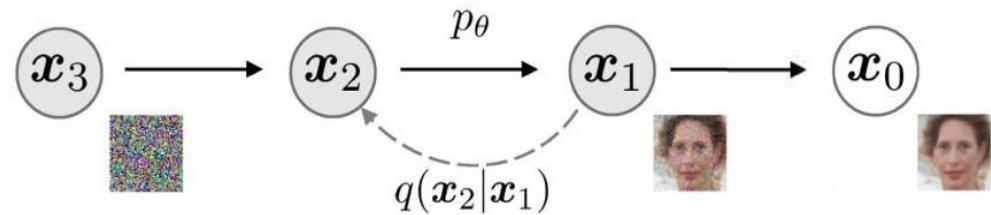
- One bottleneck of diffusion models is its slowness in sampling: need 10-1000+ steps to generate high quality samples
- Generative models need to be fast for practical use.
- Solutions:
 - Denoising Diffusion Implicit Model (DDIM) (ICLR 2021)
 - Distill diffusion models into models (using just 4-8 sampling steps)
 - *Progressive distillation for fast sampling of diffusion models, Salimans & Ho, ICLR 2022*
 - *On Distillation of Guided Diffusion Models, Meng et al., CVPR 2023*

Can we do generation with less steps?

Denoising Diffusion Implicit Model

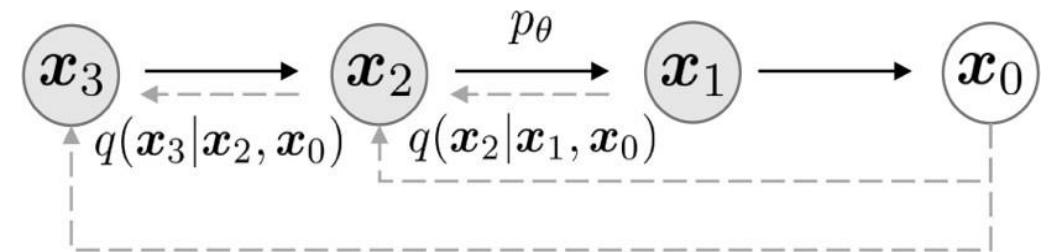


DDPM



$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N} \left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I} \right)$$

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N} \left(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I} \right)$$



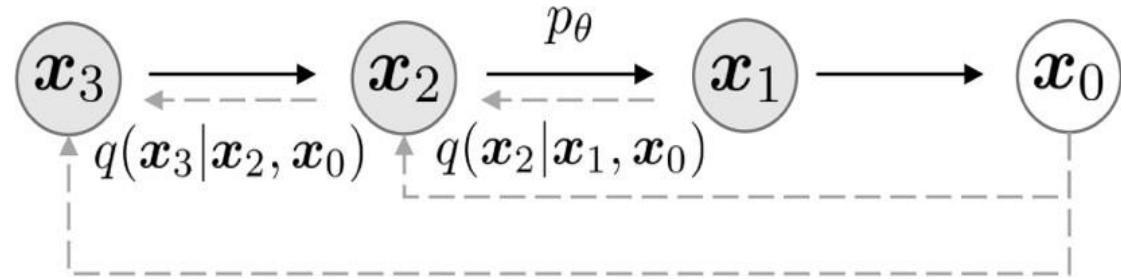
$$\begin{aligned} q_\sigma(\mathbf{x}_{1:T} \mid \mathbf{x}_0) &:= q_\sigma(\mathbf{x}_T \mid \mathbf{x}_0) \prod_{t=2}^T q_\sigma(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \\ q_\sigma(\mathbf{x}_T \mid \mathbf{x}_0) &= \mathcal{N} \left(\sqrt{\alpha_T} \mathbf{x}_0, (1 - \alpha_T) \mathbf{I} \right) \end{aligned}$$

- A Non-Markovian Forward Process

$$q_\sigma(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q_\sigma(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) q_\sigma(\mathbf{x}_t \mid \mathbf{x}_0)}{q_\sigma(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}$$

$$q_\sigma(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N} \left(\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I} \right)$$

DDIM



- Reverse process

$$p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \begin{cases} \mathcal{N}\left(f_{\theta}^{(1)}(\mathbf{x}_1), \sigma_1^2 \mathbf{I}\right) & \text{if } t=1 \\ q_{\sigma}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_t, f_{\theta}^{(t)}(\mathbf{x}_t)\right) & \text{otherwise,} \end{cases}$$

$f_{\theta}^{(t)}(\mathbf{x}_t) := \left(\mathbf{x}_t - \sqrt{1 - \alpha_t} \cdot \epsilon_{\theta}^{(t)}(\mathbf{x}_t)\right) / \sqrt{\alpha_t}$

Trying to estimate \mathbf{x}_0

DDPM vs DDIM

Algorithm DDPM Sampling

```
 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
     $\mu \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t))$ 
     $\mathbf{x}_{t-1} \leftarrow \mu + \sigma_t \epsilon$  Stochastic
end for
return  $\mathbf{x}_0$ 
```

Algorithm DDIM Sampling

```
 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\bar{\epsilon} \leftarrow \epsilon_\theta(\mathbf{x}_t, t)$ 
     $\bar{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\bar{\epsilon}}{\sqrt{\bar{\alpha}_t}}$  Estimate  $\mathbf{x}_0$ 
     $\mathbf{x}_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}}\bar{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\bar{\epsilon}$ 
end for
return  $\mathbf{x}_0$ 
```

DDIM with Fewer Steps Sampling

DDIM

Algorithm Original DDIM Sampling

```
 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\bar{\epsilon} \leftarrow \epsilon_{\theta}(\mathbf{x}_t, t)$ 
     $\bar{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \bar{\epsilon}}{\sqrt{\bar{\alpha}_t}}$ 
     $\mathbf{x}_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \bar{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \bar{\epsilon}$ 
end for
return  $\mathbf{x}_0$ 
```

Increasing
Sub-sequence

$[1, \dots, T] \Rightarrow [\tau_0 = 0, \dots, \tau_S = T]$
E.g., $\tau = [0, 10, 20, 30, \dots, 1000]$

Algorithm Fewer-Steps DDIM Sampling

```
 $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
for all  $s$  from  $S$  to 1 do
     $t \leftarrow \tau_s$ 
     $t' \leftarrow \tau_{s-1}$ 
     $\bar{\epsilon} \leftarrow \epsilon_{\theta}(\mathbf{x}_t, t)$ 
     $\bar{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \bar{\epsilon}}{\sqrt{\bar{\alpha}_t}}$ 
     $\mathbf{x}_{t'} \leftarrow \sqrt{\bar{\alpha}_{t'}} \bar{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_{t'}} \bar{\epsilon}$ 
end for
return  $\mathbf{x}_0$ 
```

DDIM Results

Table 1: CIFAR10 and CelebA image generation measured in FID. $\eta = 1.0$ and $\hat{\sigma}$ are cases of DDPM (although Ho et al. (2020) only considered $T = 1000$ steps, and $S < T$ can be seen as simulating DDPMs trained with S steps), and $\eta = 0.0$ indicates DDIM.

S	CIFAR10 (32×32)					CelebA (64×64)				
	10	20	50	100	1000	10	20	50	100	1000
η	0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53
	0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79
	0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09
	1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93
$\hat{\sigma}$		367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20
										3.26

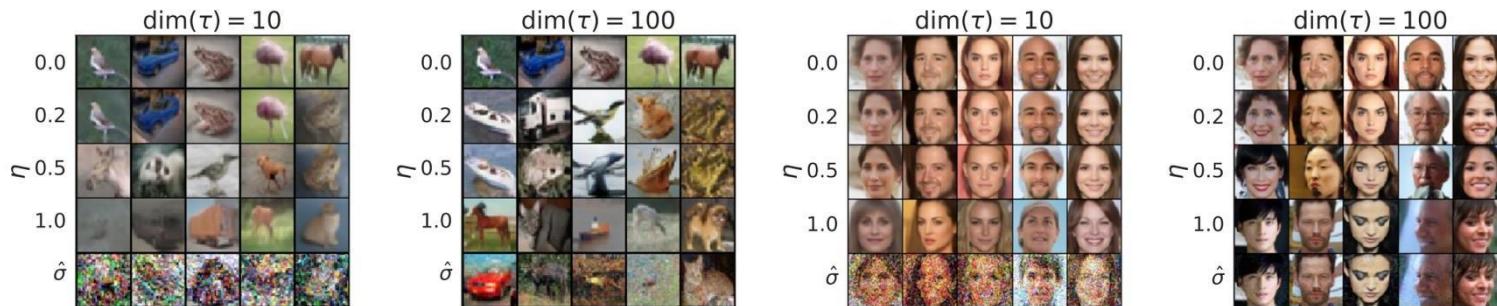
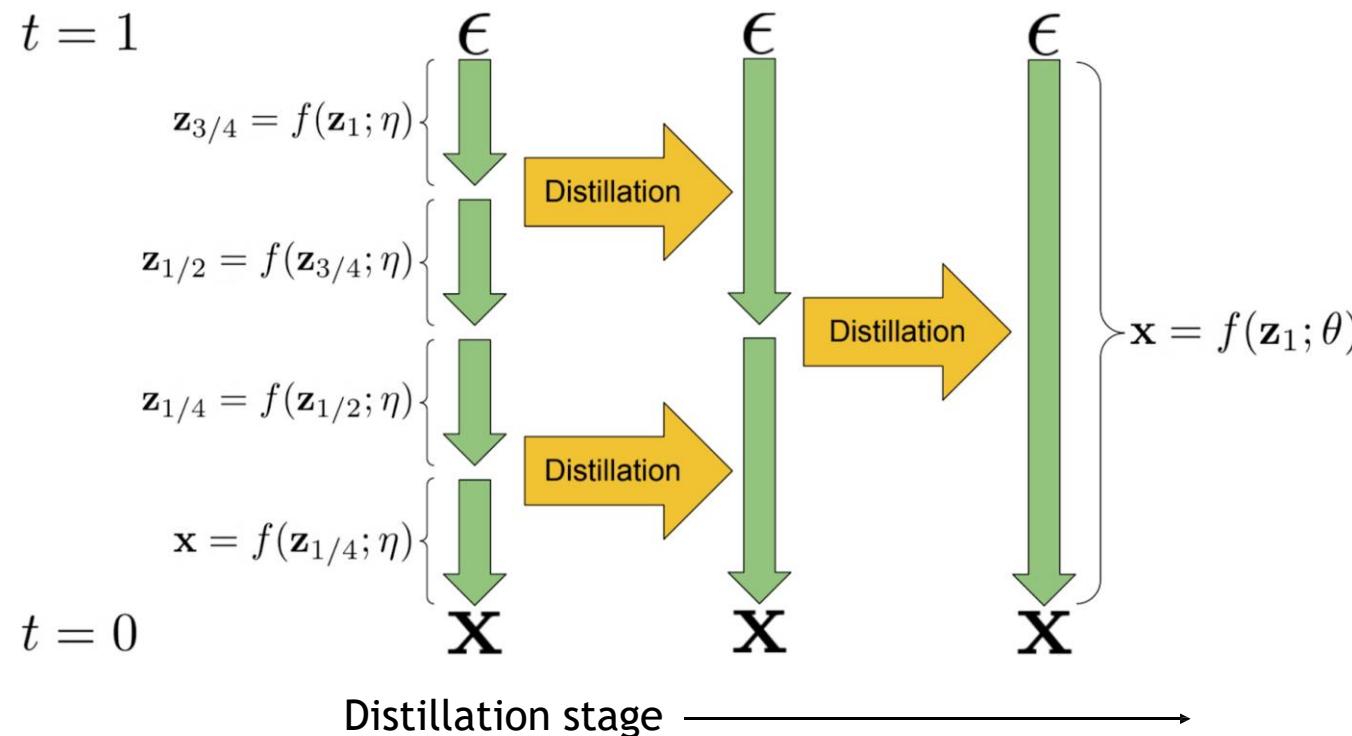


Figure 3: CIFAR10 and CelebA samples with $\text{dim}(\tau) = 10$ and $\text{dim}(\tau) = 100$.

Progressive distillation

How to make sampling faster?

- Distill a deterministic ODE sampler (i.e. DDIM sampler) to the same model architecture.
- At each stage, a “student” model is learned to distill two adjacent sampling steps of the “teacher” model to one sampling step.
- At next stage, the “student” model from previous stage will serve as the new “teacher” model.

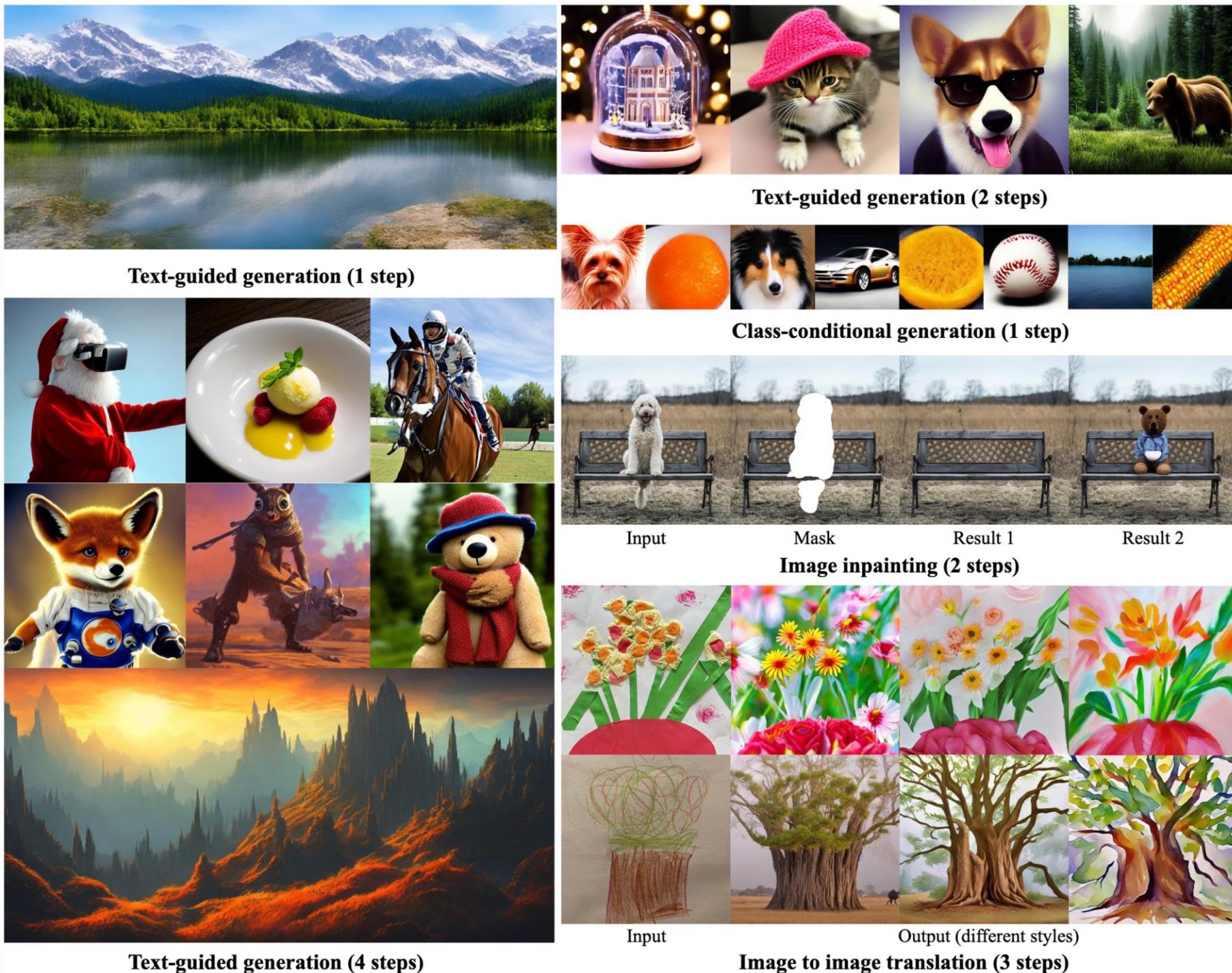


On Distillation of Guided Diffusion Models

Meng et al., CVPR 2023 award nominated

Now also works with

- CF-Guidance
- Stochastic sampling
- Text-to-image/video
- Image-to-image
- Inpainting
- Latent Diffusion



Diffusion Model in Pixel Space

- Pros
 - Intuitive Understanding: Diffusion in pixel space directly affects image pixels, making the changes visually easy to understand.
- Cons
 - Computational Cost

The larger the number of pixels, the greater the computation.
 - Memory Usage

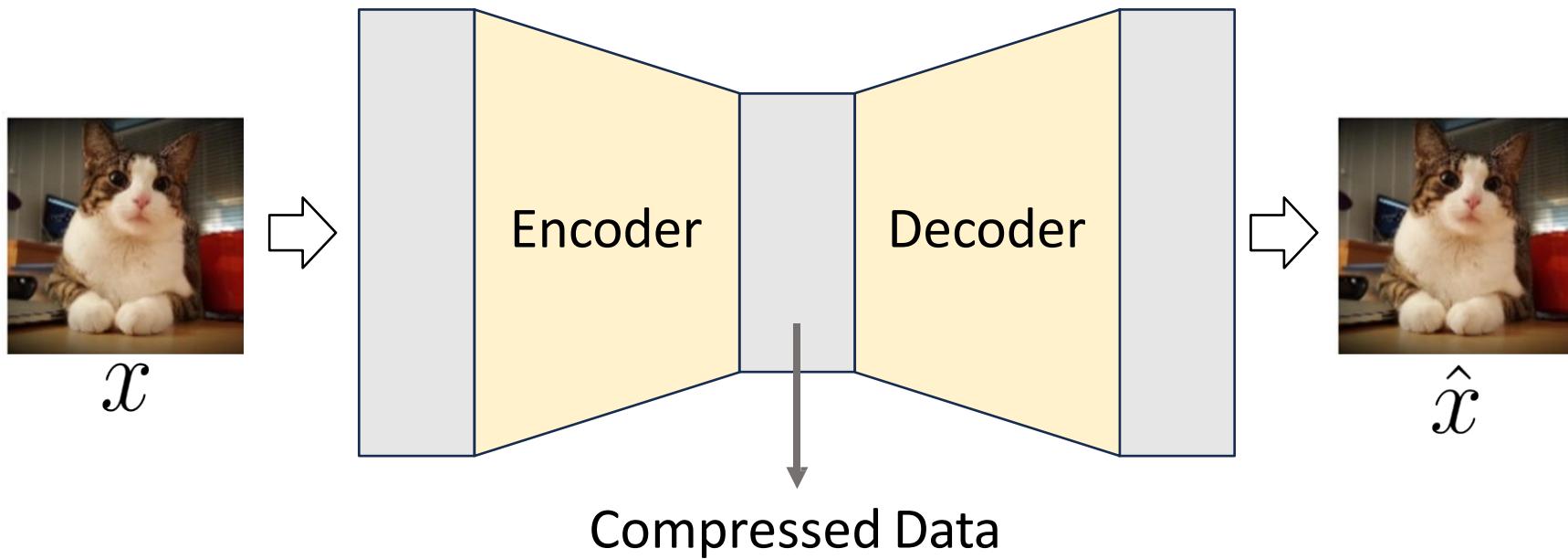
Handling high-resolution images requires substantial memory.

Latent Diffusion Model

- Latent spaces typically have lower dimensions than pixel spaces, resulting in lower computational costs.
 - Pixel Space >> Latent Space
- Runs the diffusion process in the latent space instead of pixel space

Latent Diffusion Model

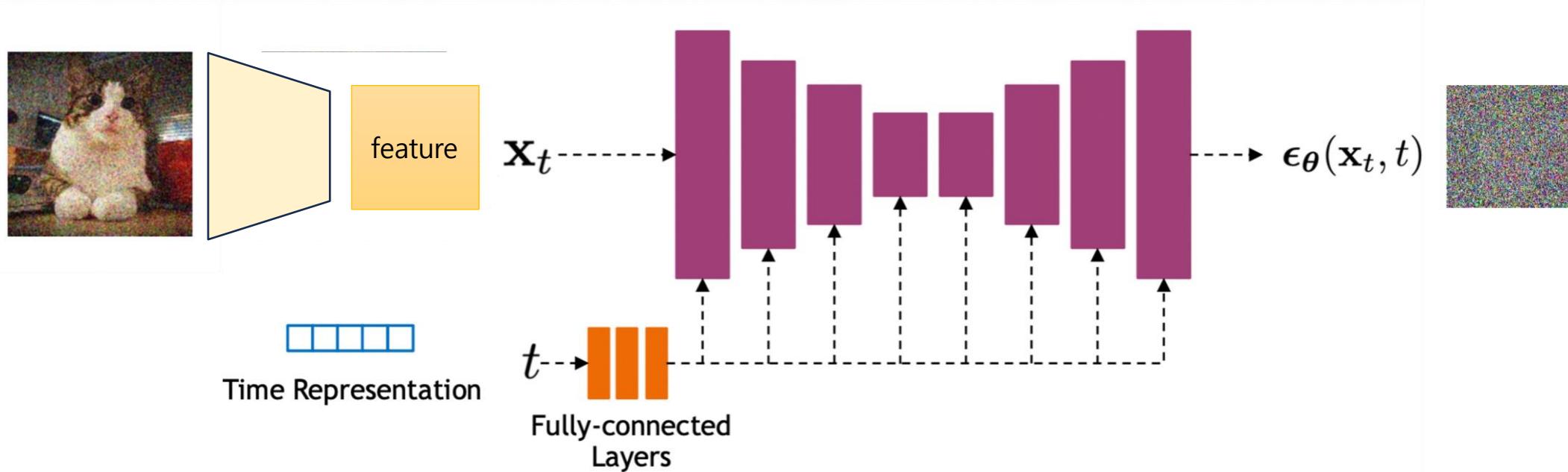
- Autoencoders can be particularly valuable as they enable a compressed yet remaining semantic and conceptual meaning of an image.



$$loss = \frac{1}{n} \sum_{i=0}^n (x_i - \hat{x}_i)^2$$

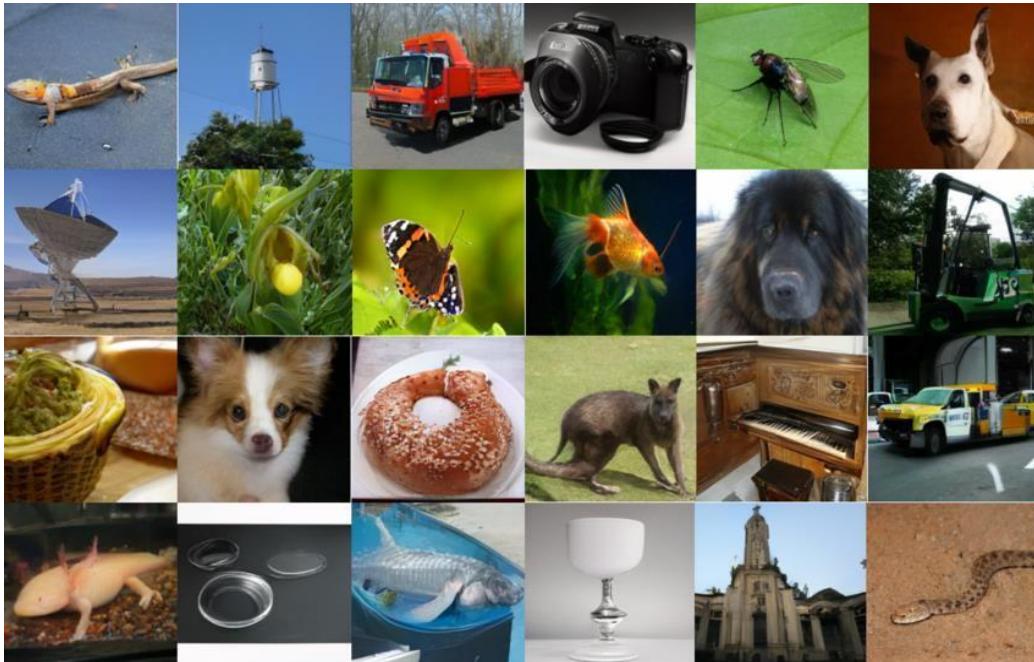
Latent Diffusion Model

- 2 Stage Training : Auto-Encoder + Latent Diffusion



Conditional Diffusion Models

- Un-conditional
 - Conditional



$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$$

- Conditional



$$p(\mathbf{x}_{0:T} \mid y) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, y)$$

More controllable!

Conditional Score Matching

- Score matching with conditional information

$$\begin{aligned}\nabla \log p(\mathbf{x}_t | y) &= \nabla \log \left(\frac{p(\mathbf{x}_t)p(y | \mathbf{x}_t)}{p(y)} \right) \\ &= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y | \mathbf{x}_t) - \nabla \log p(y) \\ &= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y | \mathbf{x}_t)}_{\text{adversarial gradient}}\end{aligned}$$

Classifier Guidance

- Use a discriminative classifier for $\nabla \log p(y | \mathbf{x}_t)$

$$\nabla \log p(\mathbf{x}_t | y) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(y | \mathbf{x}_t)$$

- γ controls the strength of the condition
- Limitations:
 - Need a separate classifier
 - Conditioning depends on the performance of classifier

Classifier-Free Guidance

- Score matching with conditional information

$$\nabla \log p(\mathbf{x}_t | y) = \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(y | \mathbf{x}_t)$$

- By Bayes' rule we can define an implicit classifier

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

- Classifier-free guidance

$$\begin{aligned}\nabla \log p(\mathbf{x}_t | y) &= \nabla \log p(\mathbf{x}_t) + \gamma (\nabla \log p(\mathbf{x}_t | y) - \nabla \log p(\mathbf{x}_t)) \\ &= \nabla \log p(\mathbf{x}_t) + \gamma \nabla \log p(\mathbf{x}_t | y) - \gamma \nabla \log p(\mathbf{x}_t) \\ &= \underbrace{\gamma \nabla \log p(\mathbf{x}_t | y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}}\end{aligned}$$

Training of Classifier-Free Guidance

- For conditional embeddings
 - Randomly drop p original conditionals with an additional unconditional class

$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

Embedding of label y

↓

Replace it with unconditional case
with probability p , i.e.,

$$\epsilon_{\theta}(z_t, t, \emptyset)$$

Content

- Diffusion Model Basics
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- Variants of Diffusion Models
 - Denoising Diffusion Implicit Model (DDIM)
 - Latent Diffusion Models
 - Conditional Diffusion Models
- Applications of Diffusion Models

Case study: Imagen

Imagen: text-to-image diffusion models

By Google (imagen.research.google)

Input: text; Output: 1kx1k images

- An unprecedented degree of photorealism
 - SOTA automatic scores & human ratings
- A deep level of language understanding
- Extremely simple
 - no latent space, no quantization



A brain riding a rocketship heading towards the moon.

Imagen

By Google (imagen.research.google)



Imagen

A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

Imagen

By Google (imagen.research.google)



A dragon fruit wearing karate belt in the snow.



“toilet paper with real cactus spikes”

What goes to Imagen?

Data

- Image-text pairs
- LAION-400M
- Internal (~500M images)

Model

- Diffusion models
- Cascading super-res
- Frozen text encoders

Sampler

- Classifier-free guidance
- Maximizing text-alignment

Scaling up

What goes to Imagen?

Data

- Image-text pairs
- LAION-400M
- Internal (~500M images)

Model

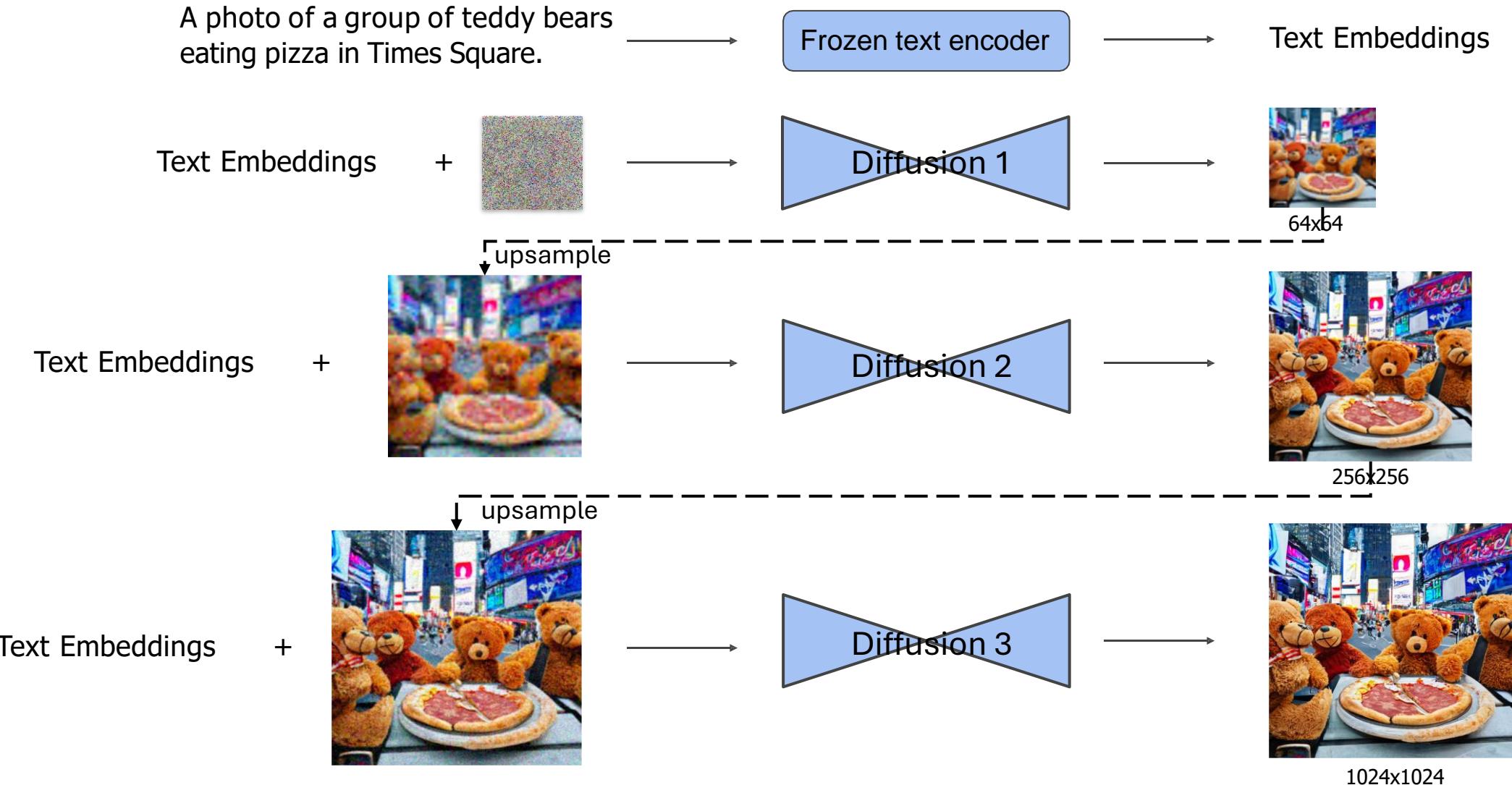
- Diffusion models
- Cascading super-res
- Frozen text encoders

Sampler

- Classifier-free guidance
- Maximizing text-alignment

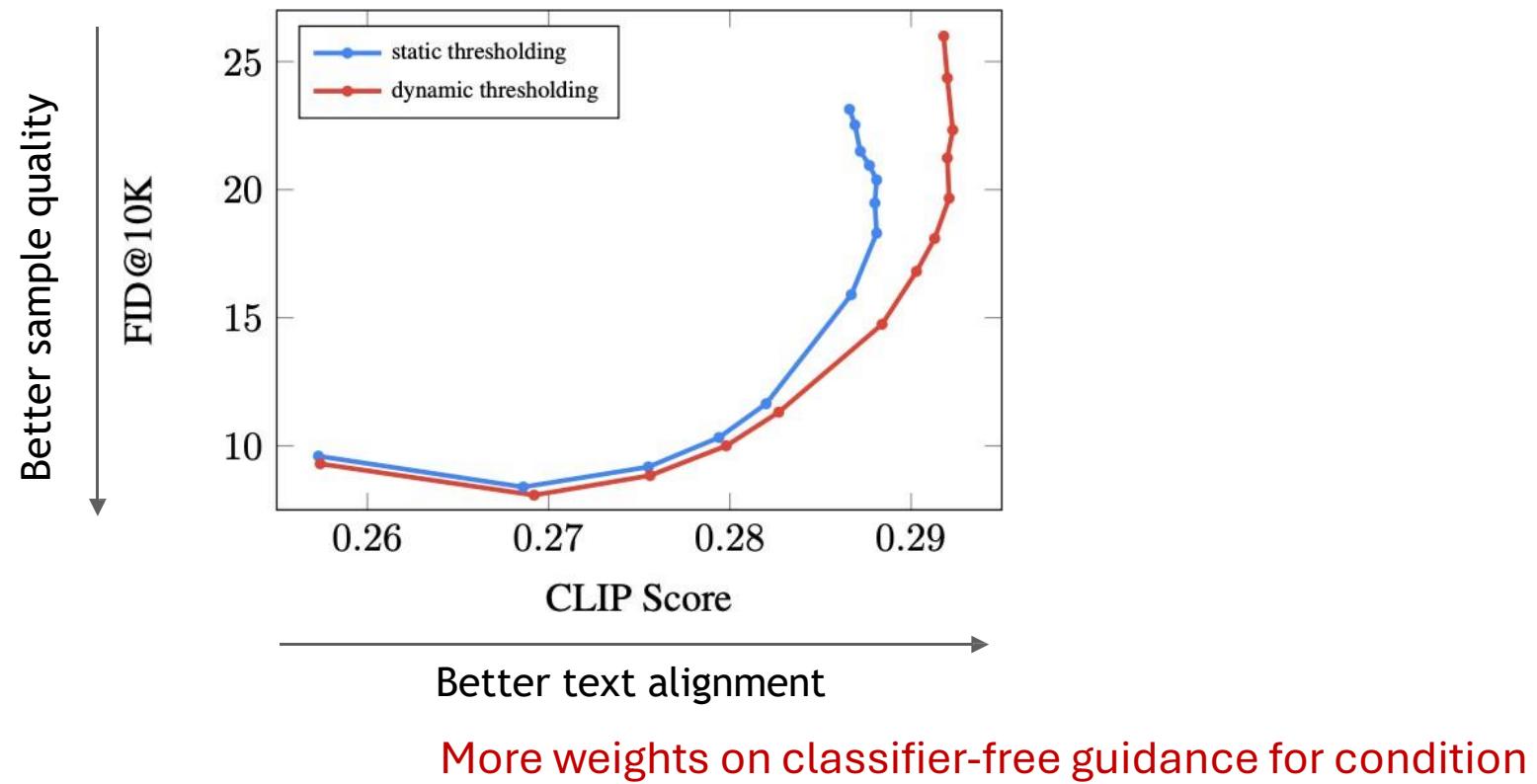
Scaling up

Imagen: Cascaded generation pipeline

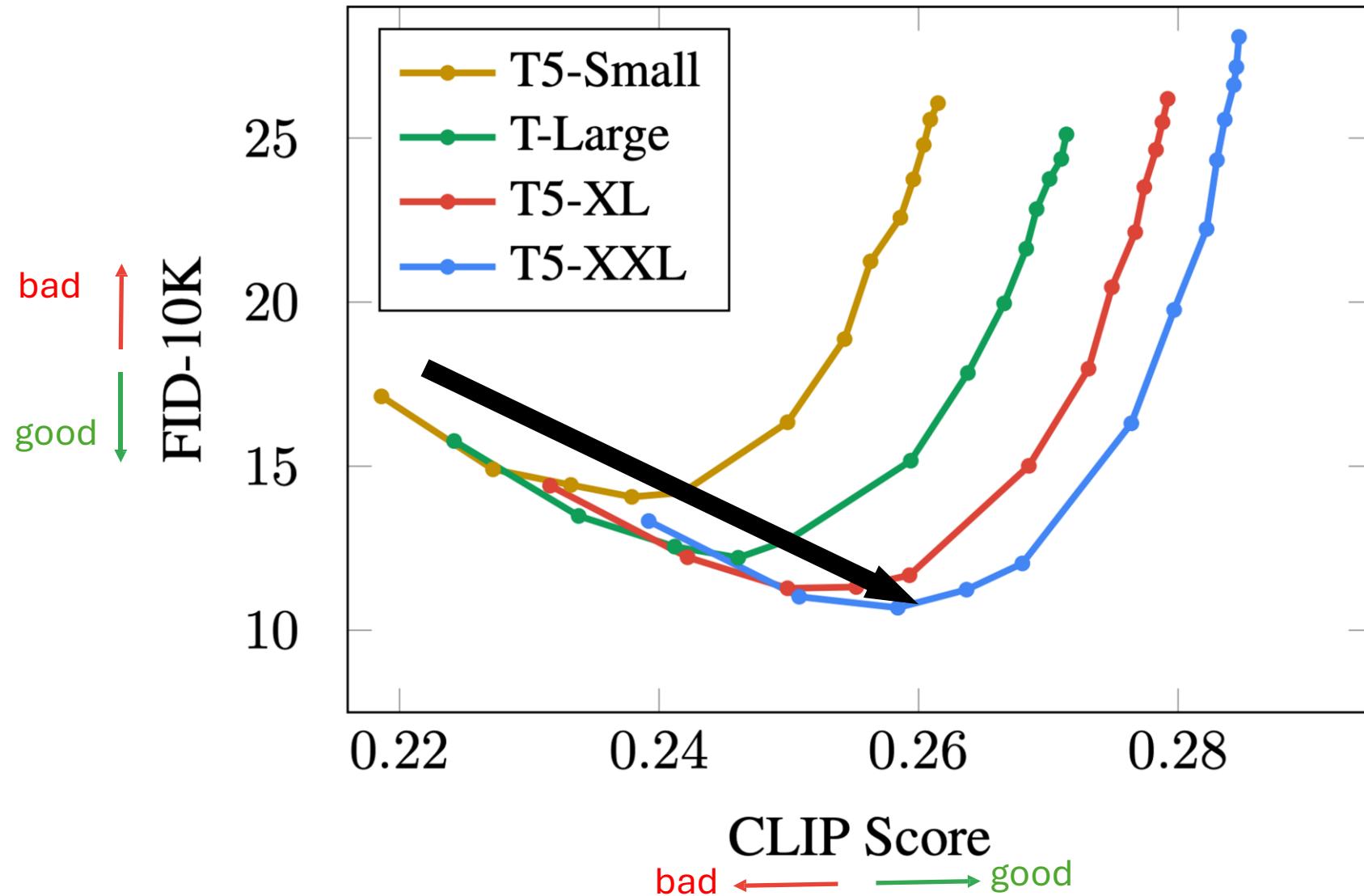


Classifier-free guidance in Imagen

- Large classifier-free guidance weights → better text alignment, worse image fidelity



Larger Text Encoders → Better Alignment, Better Fidelity

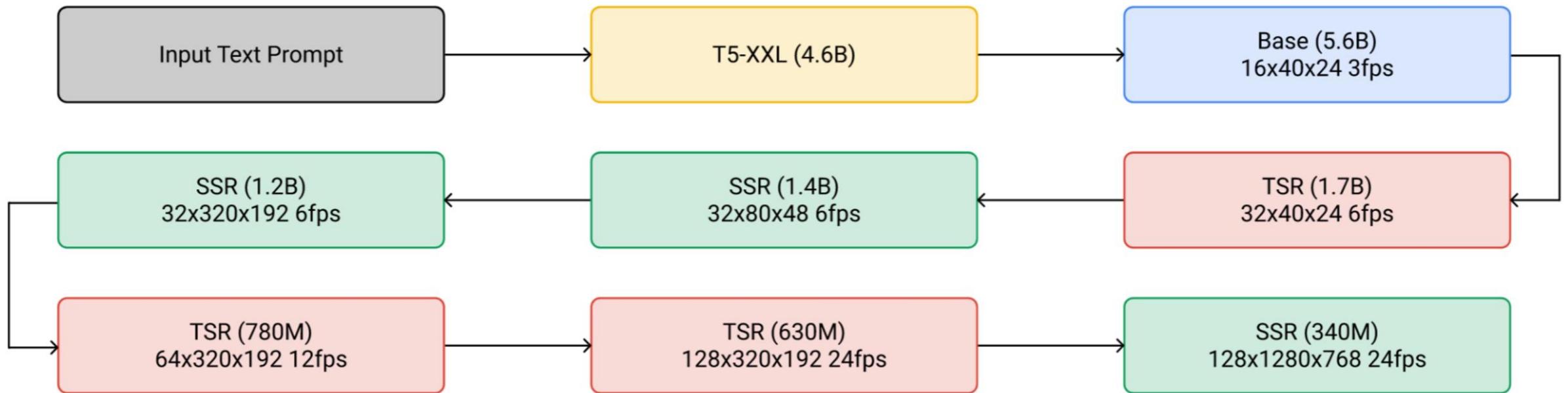


Can we generalize this to video?

Imagen Video



Imagen Video generalizes **Imagen** to the video domain using a cascade of super-resolution diffusion models in space and time



TSR: Temporal Super Resolution

SSR: Spatial Super Resolution

Output shape: frames × width × height

fps: frames per second

Video diffusion models

Ho & Salimans, et al. <https://video-diffusion.github.io/>

- Image → Video: Just add another dimension to the data tensor
- Image architecture: 2D UNet
- Video architecture: 3D UNet, space-time separable
 - repeat the 2D UNet over frames
 - additional layers to mix over time using attention or convolution

