# Lecture 18: Diffusion Model II: Score-based Methods

Lan Xu

SIST, ShanghaiTech

Fall, 2023

# Outline

- Diffusion: Denoising Diffusion Probabilistic Models

- Score-based Diffusion Models

- Accelerated Sampling
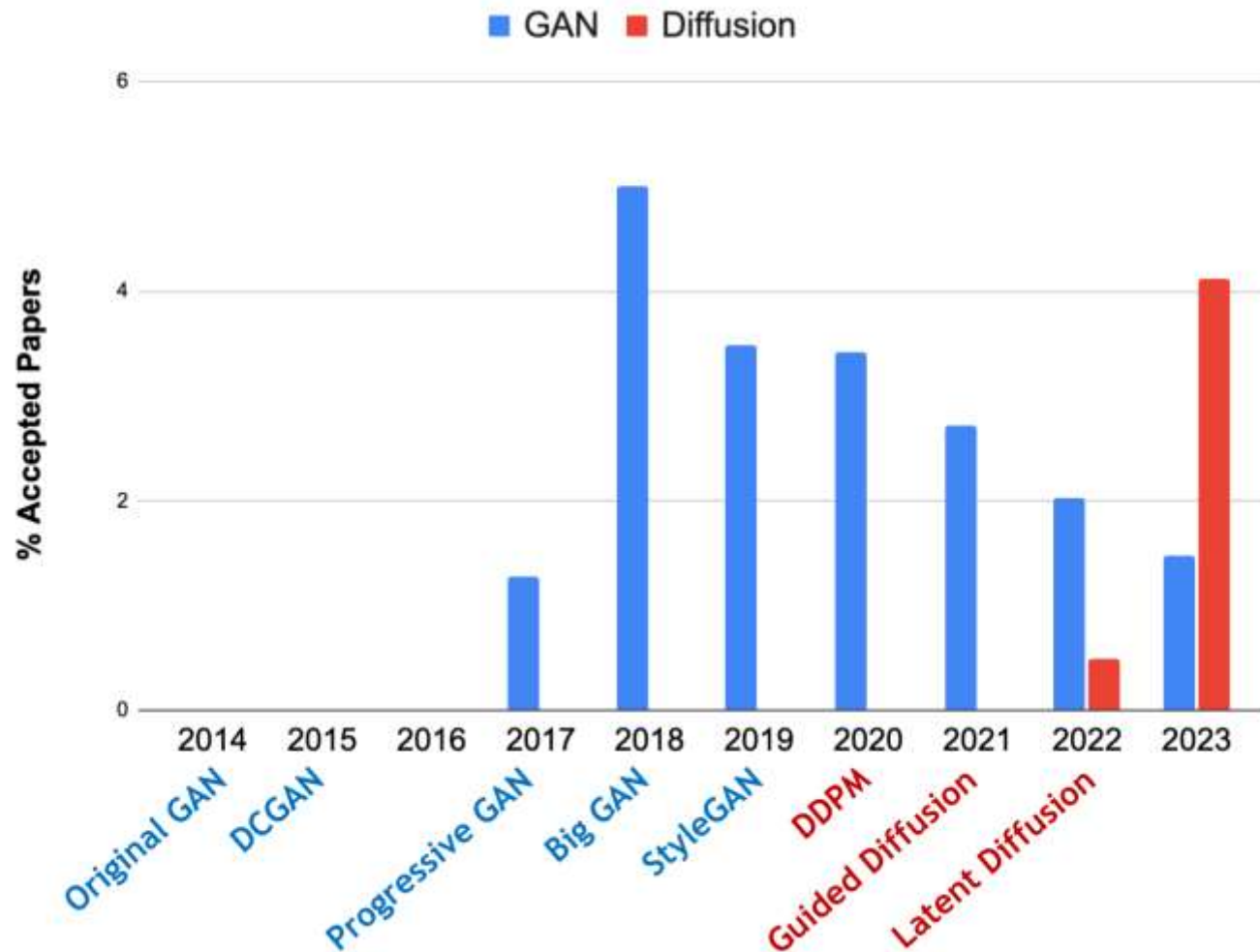
- Conditional Generation and Guidance

# The Landscape of Generative Models

# Diffusion: A Generative Learning Big Bang



Big Bang!

Big Bang!

# Diffusion: new trend in CVPR

Lan Xu – CS 280 Deep Learning

# Recall DDPM

- Denoising Diffusion Probabilistic Models
- Target: understand the training and sampling phases!

**Algorithm 1** Training

1: **repeat**
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:    $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:    Take gradient descent step on
$$\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon, t) \right\|^2$$
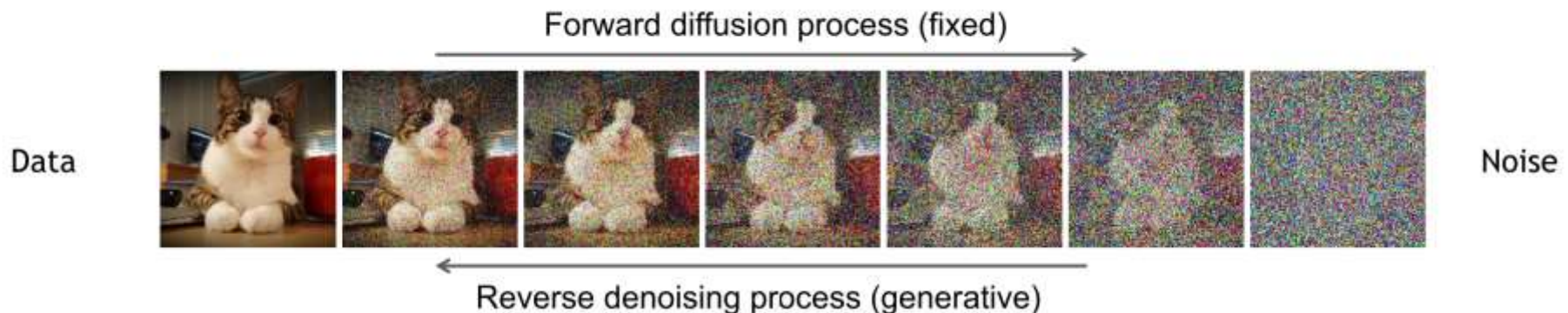6: **until** converged

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right) + \sigma_t\mathbf{z}$
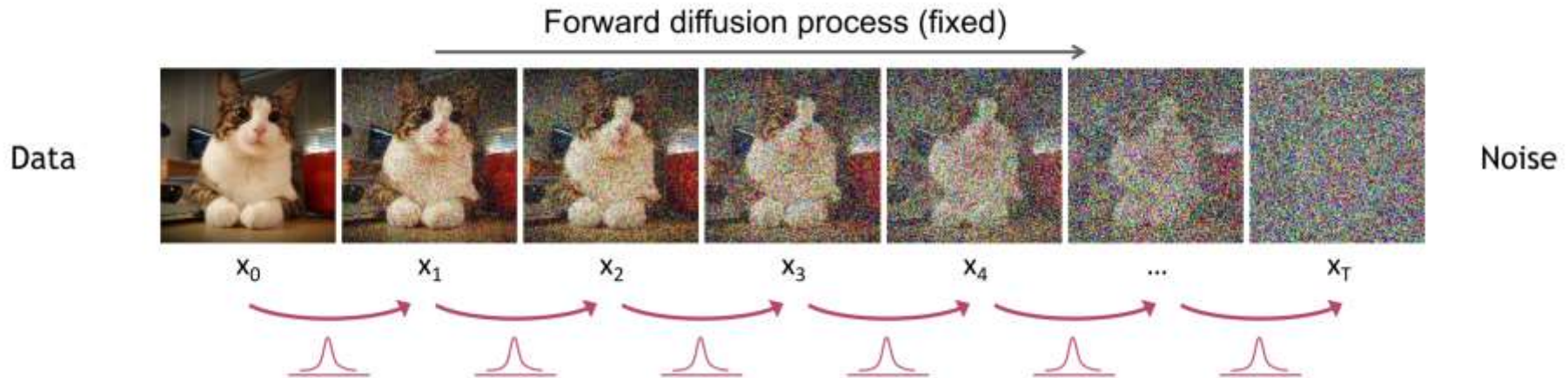5: **end for**
6: **return** $\mathbf{x}_0$

# Recall DDPM

- Learning to generate by denoising
- **Forward diffusion** process that gradually adds noise to input
- **Reverse denoising** process that learns to generate data by denoising

Forward diffusion process (fixed)

Data → ... → Noise

Reverse denoising process (generative)

# Recall DDPM

- Forward Diffusion Process in T steps

Forward diffusion process (fixed)



Data                                                                Noise
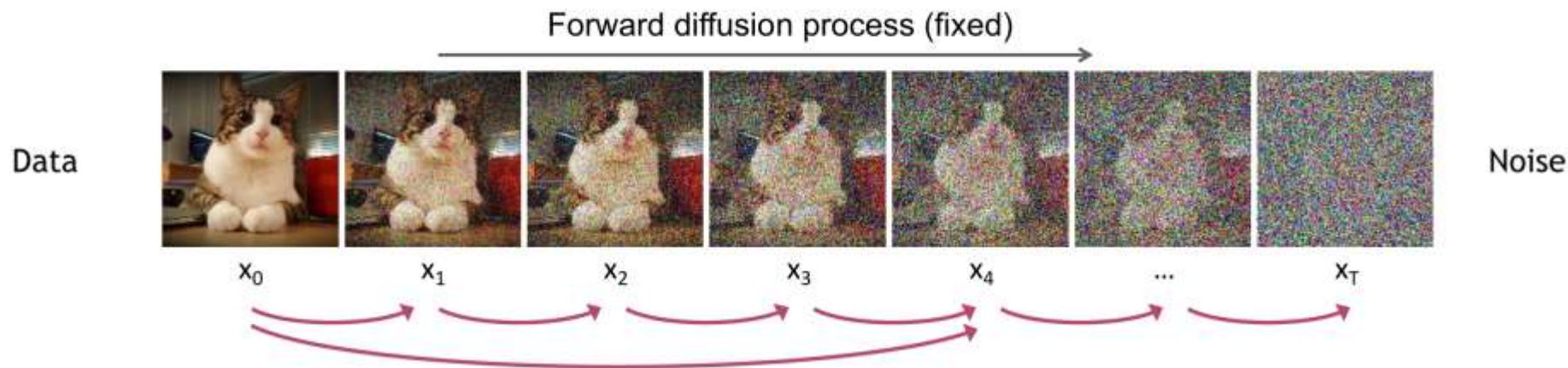
$x_0$    $x_1$    $x_2$    $x_3$    $x_4$    ...    $x_T$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$$x_t = \sqrt{1-\beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, I)$$

# Recall DDPM

■ Diffusion Kernel and Reparametrization Trick

Forward diffusion process (fixed)



Data

Noise

$x_0$  $x_1$  $x_2$  $x_3$  $x_4$  ...  $x_T$

$$\alpha_t = 1 - \beta_t \quad \bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$$

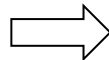$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}\left(\sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}\right)$$

$$\mathbf{x}_t = \sqrt{1-\beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$
$$= \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1-\alpha_t}\epsilon$$
$$= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1-\alpha_t\alpha_{t-1}}\epsilon$$
$$= \ldots$$
$$= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\epsilon$$

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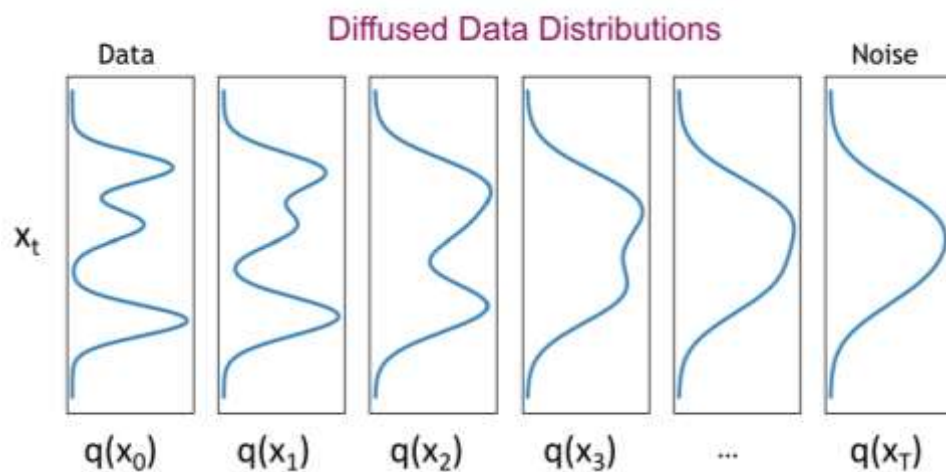
$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\,\mathbf{x}_0 + \sqrt{(1-\bar{\alpha}_t)}\,\epsilon$$

The noise schedule) is designed such that $\bar{\alpha}_T \rightarrow 0$

# Recall DDPM

- What happens to a distribution in the forward diffusion?
- The diffusion kernel $q(\mathbf{x}_t|\mathbf{x}_0)$ is **Gaussian convolution**!

$$\underbrace{q(\mathbf{x}_t)}_{\substack{\text{Diffused} \\ \text{data dist.}}} = \int \underbrace{q(\mathbf{x}_0, \mathbf{x}_t)}_{\substack{\text{Joint} \\ \text{dist.}}} d\mathbf{x}_0 = \int \underbrace{q(\mathbf{x}_0)}_{\substack{\text{Input} \\ \text{data dist.}}} \underbrace{q(\mathbf{x}_t|\mathbf{x}_0)}_{\substack{\text{Diffusion} \\ \text{kernel}}} d\mathbf{x}_0$$

**Diffused Data Distributions**

Data            Noise

$x_t$

$q(x_0)$   $q(x_1)$   $q(x_2)$   $q(x_3)$   ...   $q(x_T)$

- We can sample $\mathbf{x}_t \sim q(\mathbf{x}_t)$ by first sampling $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, and then sampling $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$ (ancestral sampling).
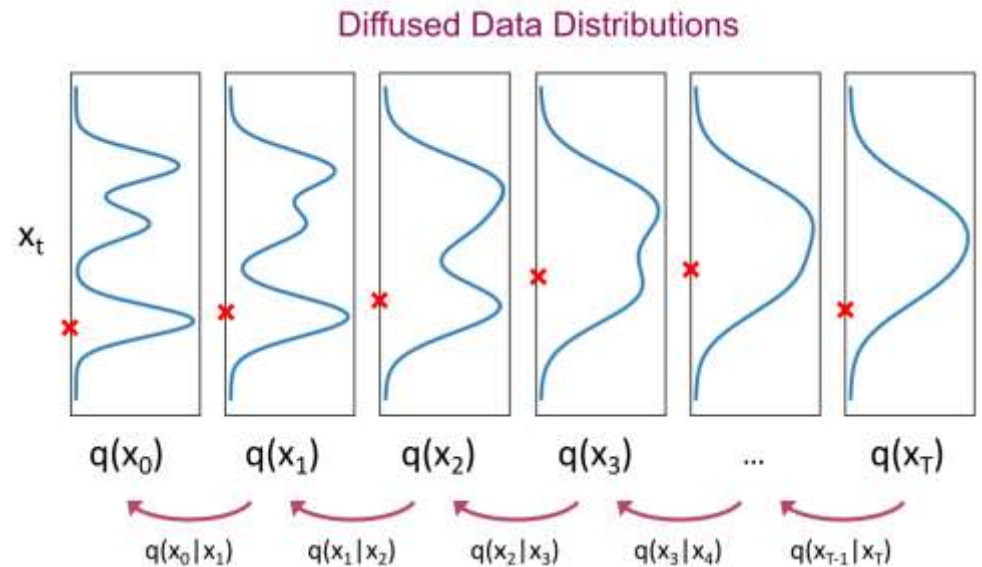
# Recall DDPM

- Reverse: Generative Learning by Denoising
- Diffusion parameters are designed: $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}))$
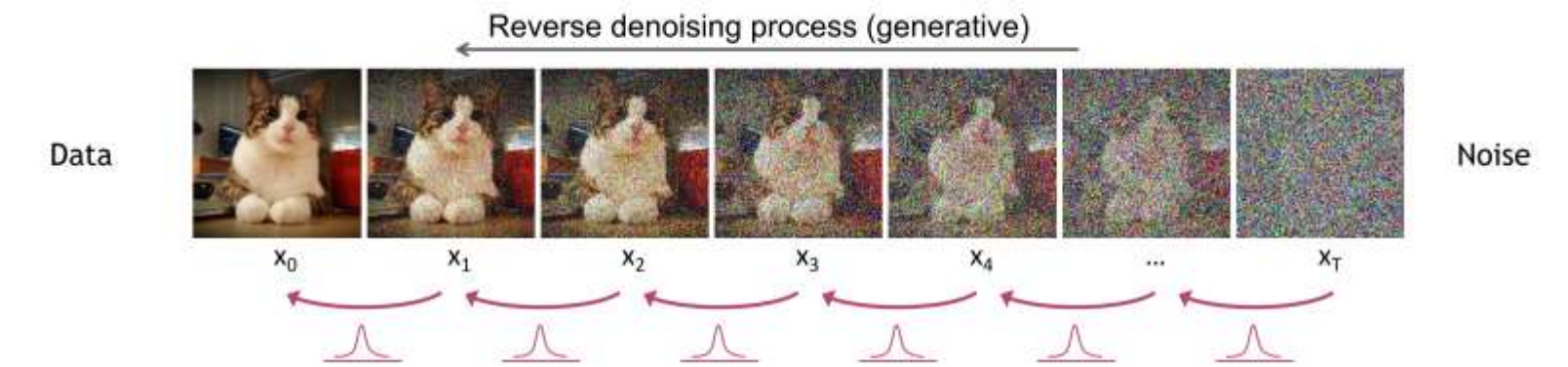
**Generation:**

Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

Iteratively sample $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1}|\mathbf{x}_t)}_{\text{True Denoising Dist.}}$

**Diffused Data Distributions**



$q(x_0)$  $q(x_1)$  $q(x_2)$  $q(x_3)$  ...  $q(x_T)$

$q(x_0|x_1)$  $q(x_1|x_2)$  $q(x_2|x_3)$  $q(x_3|x_4)$  $q(x_{T-1}|x_T)$

# Recall DDPM

- **Reverse Denoising Process**



Reverse denoising process (generative)

Data ... Noise

$$x_0 \quad x_1 \quad x_2 \quad x_3 \quad x_4 \quad \dots \quad x_T$$

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

- From $q(x_{t-1}|x_t)$ to $q(x_{t-1}|x_t, x_0)$

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_{t-1}, x_t, x_0)}{q(x_t, x_0)} = \frac{q(x_t|x_{t-1}, x_0) \cdot q(x_{t-1}, x_0)}{q(x_t, x_0)} = \frac{q(x_t|x_{t-1}, x_0) \cdot q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

$$= \frac{q(x_t|x_{t-1}) \cdot q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

# Recall DDPM

- **Learning Denoising Model**

$$\frac{q(x_t|x_{t-1}) \cdot q(x_{t-1}|x_0)}{q(x_t|x_0)}$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \boldsymbol{I})$$

where $\quad \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}\mathbf{x}_t \quad$ and $\quad \tilde{\beta}_t := \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$

- Since $\quad x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1-\bar{\alpha}_t}\bar{\epsilon}_t)$, thus: $\quad \tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\bar{\epsilon}_t\right)$

- Use NN to regress the noise:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1-\beta_t}}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right)$$

$$L = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1,T\}, \epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})}\left[\lambda_t||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\,\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\,\epsilon, t)||^2\right]$$

# Recall DDPM

- The authors of DDPM say that it's fine to drop all that baggage in the front and instead just use

$$L = \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[ \left\| \epsilon - \epsilon_\theta(\mathbf{x}_t, t) \right\|_2^2 \right]$$

$$= \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[ \left\| \epsilon - \epsilon_\theta\left( \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t \right) \right\|_2^2 \right]$$

- Note that this is not a variational lower bound on the log-likelihood anymore: in fact, you can view it as a **reweighted version of ELBO** that emphasizes reconstruction quality!

# Recall DDPM

- If we have the noise, sampling by using Gaussians:

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t \boldsymbol{I})$$

- 1) sampling $z_t$
- 2) sampling x_t-1, using the estimated noise

$$x_{t-1} = \tilde{\mu}_t + \tilde{\beta}_t \cdot z_t = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}}\overline{\epsilon}_t\right) + \frac{1-\overline{\alpha}_{t-1}}{1-\overline{\alpha}_t} \cdot \beta_t \cdot z_t$$

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}}\epsilon_\theta(x_t, t)\right) + \frac{1-\overline{\alpha}_{t-1}}{1-\overline{\alpha}_t} \cdot \beta_t \cdot z_t$$

# Recall DDPM

■ Rethinking the Training and Sampling processes…..

---

**Algorithm 1** Training

1: **repeat**
2: $\quad \mathbf{x}_0 \sim q(\mathbf{x}_0)$
3: $\quad t \sim \text{Uniform}(\{1, \ldots, T\})$
4: $\quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5: $\quad$ Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1-\bar{\alpha}_t}\boldsymbol{\epsilon}, t) \right\|^2$$
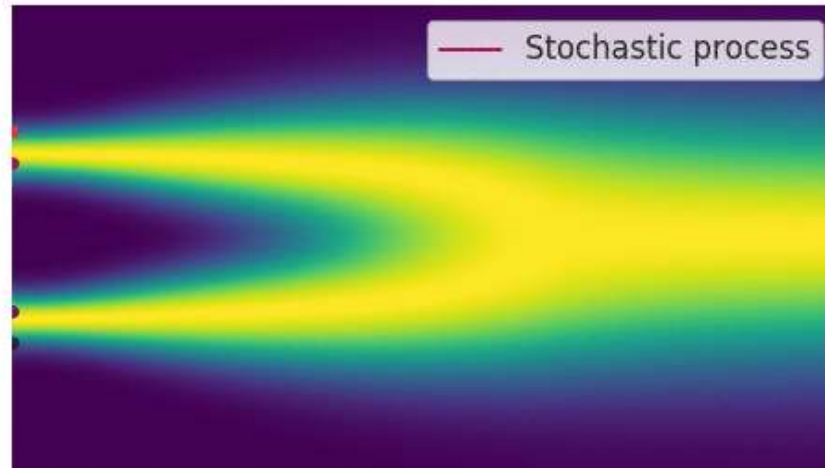6: **until** converged

---

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3: $\quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4: $\quad \mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
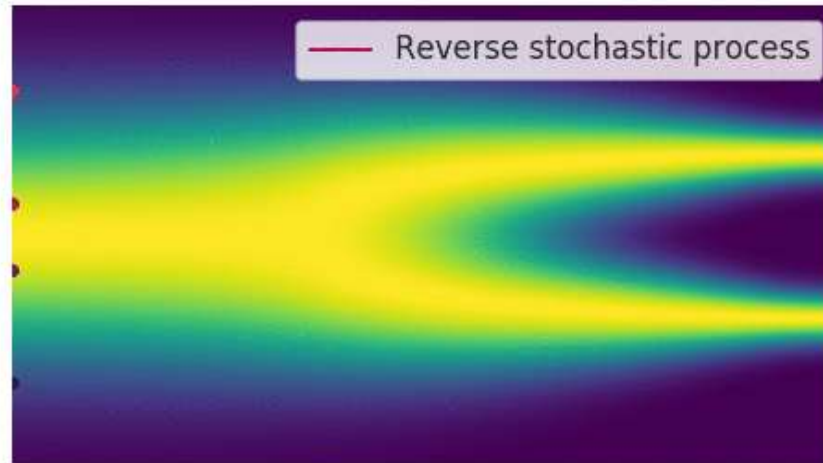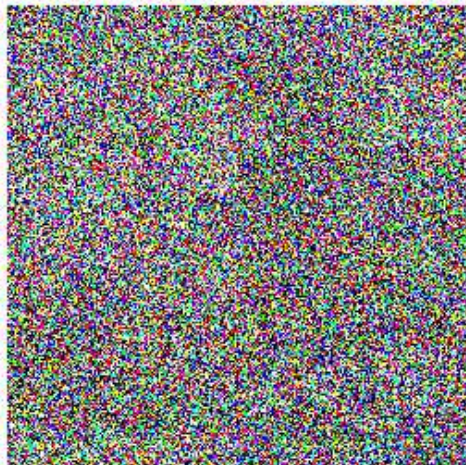5: **end for**
6: **return** $\mathbf{x}_0$

---

■ During training, add noise from 0 to t, then estimate it

■ During sampling, note that $\quad \sigma_t = \frac{1-\overline{\alpha}_{t-1}}{1-\overline{\alpha}_t} \cdot \beta_t$

■ As t increases, $\overline{\alpha}_t$ decreases, $\sqrt{1 - \overline{\alpha}_t}$ increases

■ Thus, $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ works as denoise auto-encoder for various noise levels!

# Recall DDPM

- **Forward/Reverse process for Image Generation**
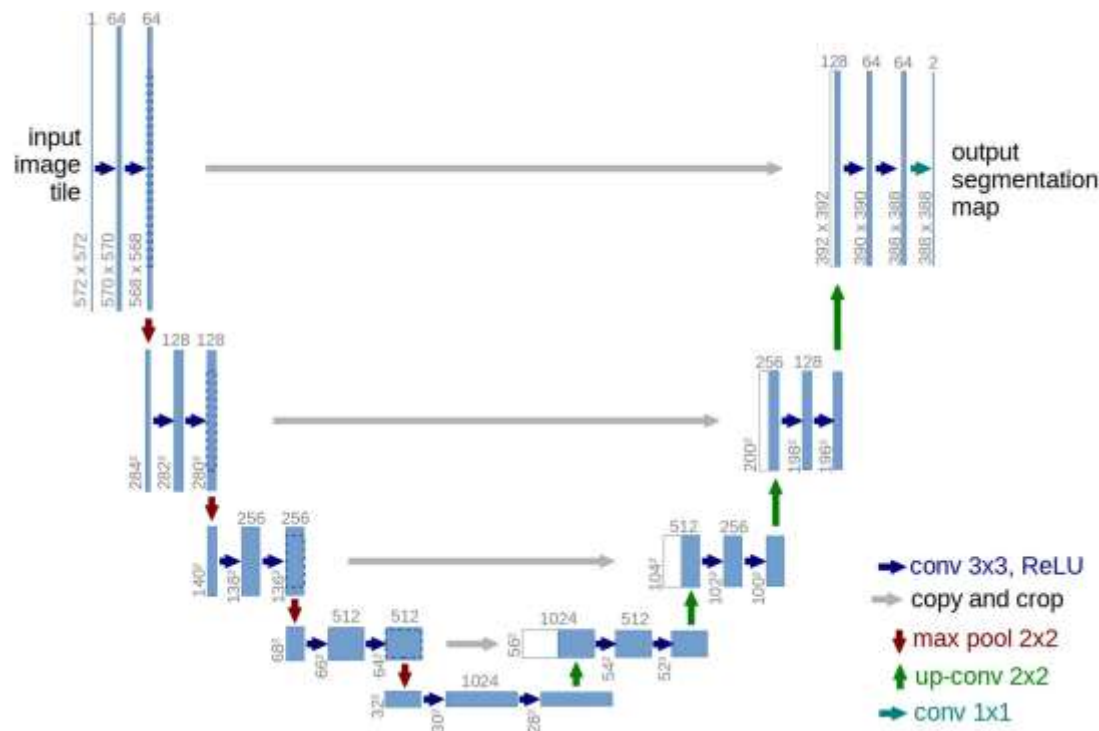


Forward process: converting the image distribution to pure noise

Reverse process: sampling from the image distribution, starting with pure noise

Lan Xu – CS 280 Deep Learning

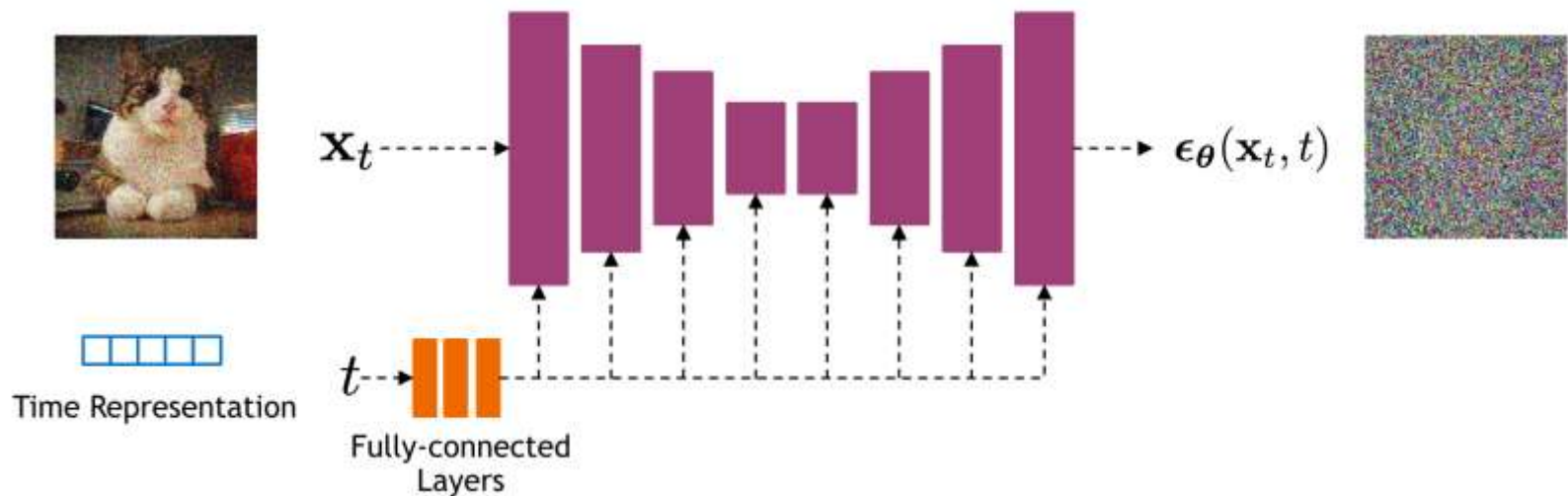# Implementation Considerations

- ## UNet + Other Stuff



Diffusion models typically use a U-Net on steroids as the noise predictive model — you take the good ol' model that you are already familiar with and add:

- Positional Embeddings
- ResNet Blocks
- ConvNext Blocks
- Attention Modules
- Group Normalization
- Swish and GeLU

It's a massive kitchen sink of modern CV tricks

# Implementation Considerations

- Diffusion models often use **U-Net architectures** with ResNet blocks and self-attention layers



- Time representation: sinusoidal positional embeddings or random Fourier features.

- Time features are fed to the residual blocks using either simple spatial addition or using adaptive group normalization layers.
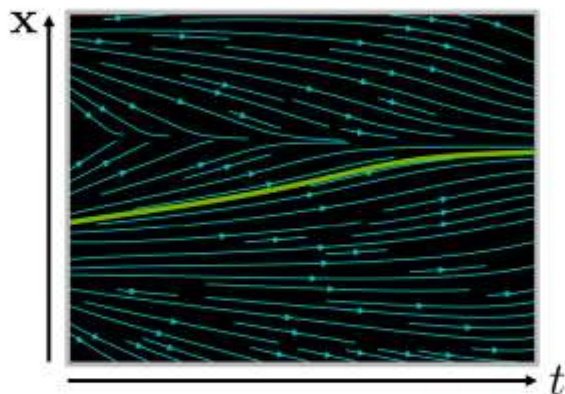
# Outline

- Diffusion: Denoising Diffusion Probabilistic Models

- **Score-based Diffusion Models**

- Accelerated Sampling

- Conditional Generation and Guidance

# Crash Course in Differential Equations

- ## ODE and SDE

**Ordinary Differential Equation (ODE):**

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



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

Iterative Numerical Solution:
$$\mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t$$
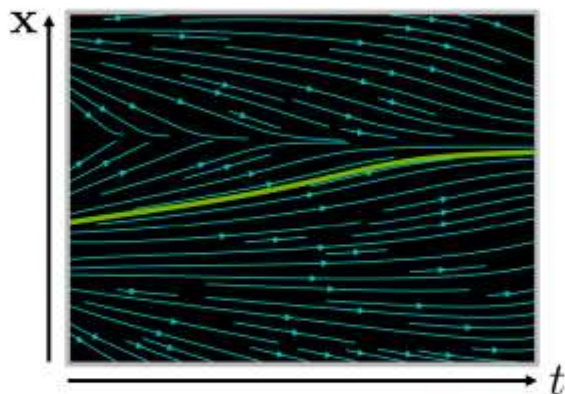
# Crash Course in Differential Equations

- ## ODE and SDE



**Ordinary Differential Equation (ODE):**

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, t) \quad \text{or} \quad d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt$$
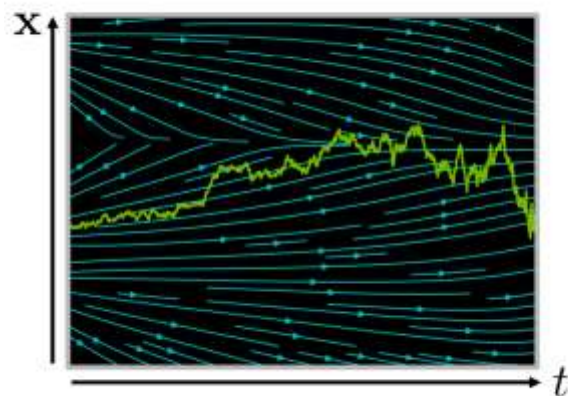
Analytical Solution:
$$\mathbf{x}(t) = \mathbf{x}(0) + \int_0^t \mathbf{f}(\mathbf{x}, \tau)d\tau$$

Iterative Numerical Solution:
$$\mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t$$
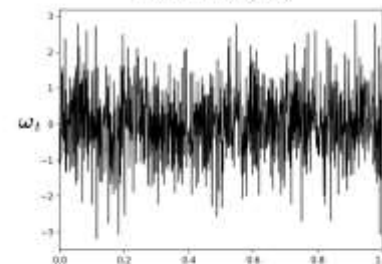
**Stochastic Differential Equation (SDE):**

$$\frac{d\mathbf{x}}{dt} = \underbrace{\mathbf{f}(\mathbf{x}, t)}_{\text{drift coefficient}} + \underbrace{\sigma(\mathbf{x}, t)}_{\text{diffusion coefficient}}\boldsymbol{\omega}_t$$

drift coefficient    diffusion coefficient

$$\left( d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\boldsymbol{\omega}_t \right)$$

Wiener Process (Gaussian White Noise)

$$\mathbf{x}(t + \Delta t) \approx \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t + \sigma(\mathbf{x}(t), t)\sqrt{\Delta t}\,\mathcal{N}(\mathbf{0}, \mathbf{I})$$
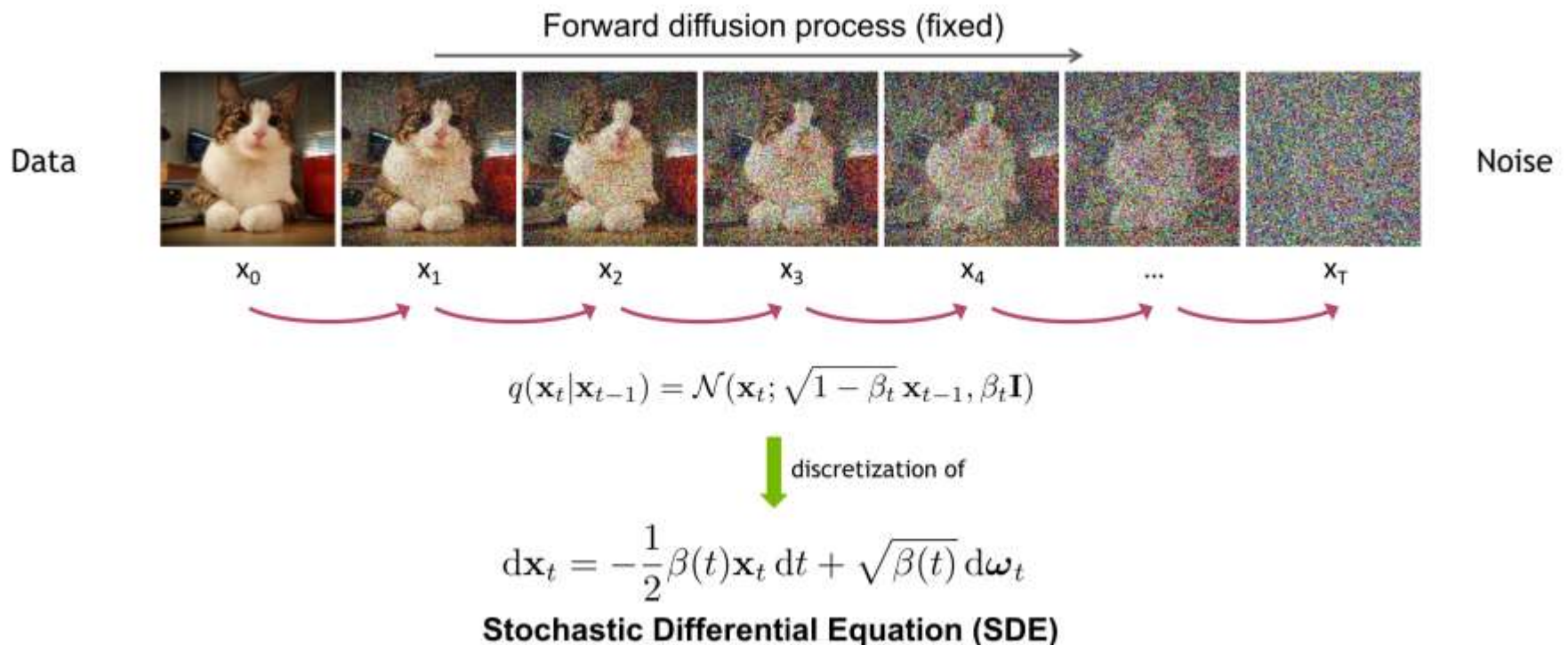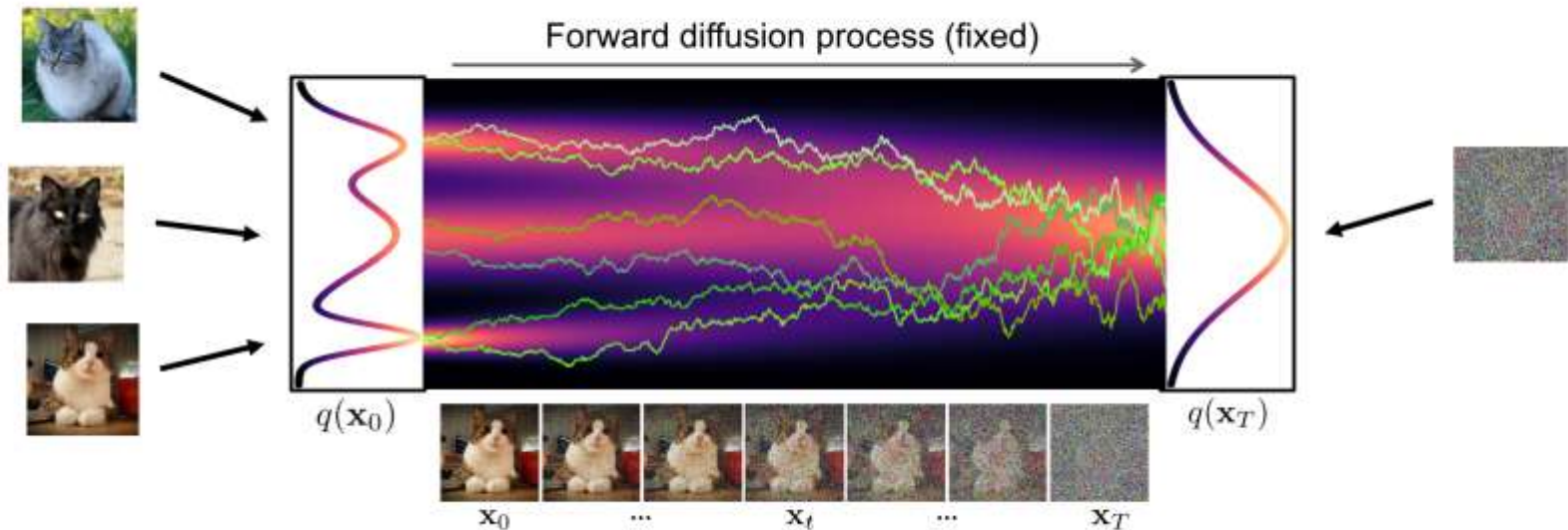
# Score-based Diffusion Models

- **Forward Diffusion Process** as Stochastic Differential Equation (continuous form)
- Consider the limit of many small steps:

Forward diffusion process (fixed)

Data

Noise

$x_0$    $x_1$    $x_2$    $x_3$    $x_4$    ...    $x_T$

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

discretization of

$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$

**Stochastic Differential Equation (SDE)**

Song et al., "Score-Based Generative Modeling through Stochastic Differential Equations", ICLR, 2021

# Score-based Diffusion Models

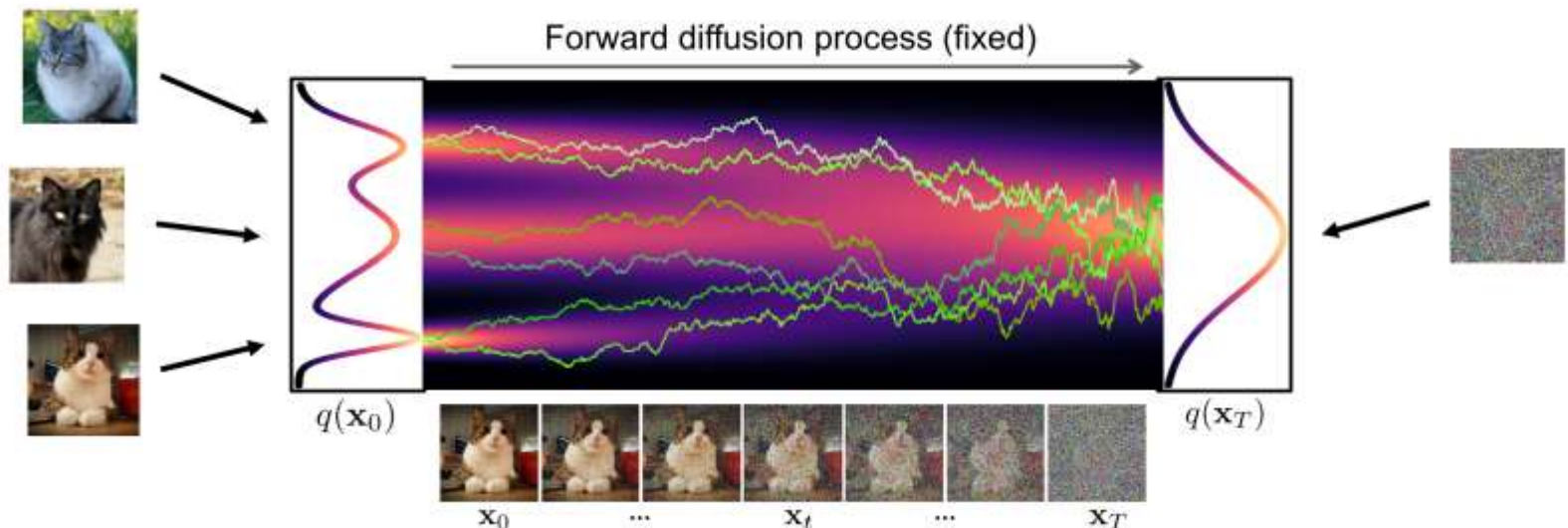- **Forward Diffusion Process** as Stochastic Differential Equation (continuous form)



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$ ... $q(\mathbf{x}_T)$

$\mathbf{x}_0$ ... $\mathbf{x}_t$ ... $\mathbf{x}_T$

**Forward Diffusion SDE:**

$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$

drift term
(pulls towards mode)

diffusion term
(injects noise)

# Score-based Diffusion Models

- The Generative Reverse Stochastic Differential Equation



**Forward Diffusion SDE:**
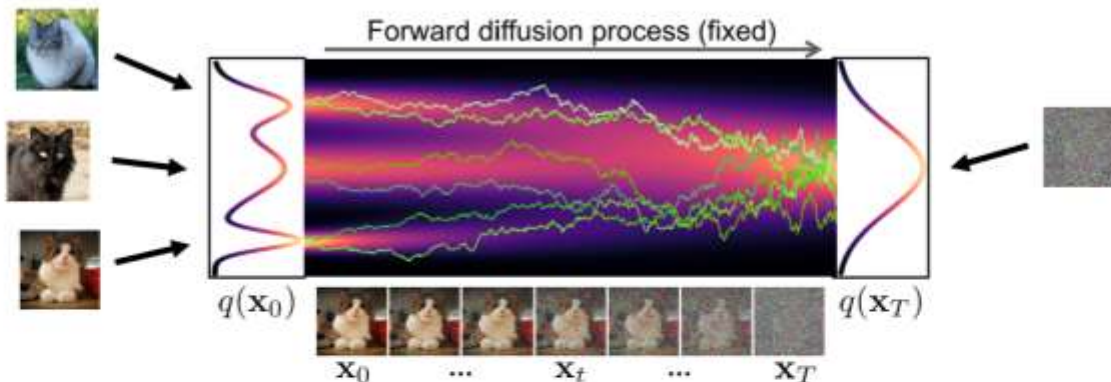$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$

- But what about the **reverse direction**, necessary for generation?

# Score-based Diffusion Models

- The Generative Reverse Stochastic Differential Equation



**Forward Diffusion SDE:**
$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,dt + \sqrt{\beta(t)}\,d\boldsymbol{\omega}_t$$

**Reverse Generative Diffusion SDE:**
$$d\mathbf{x}_t = \left[-\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t)\right]dt + \sqrt{\beta(t)}\,d\bar{\boldsymbol{\omega}}_t$$

drift term — "Score Function"    diffusion term

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

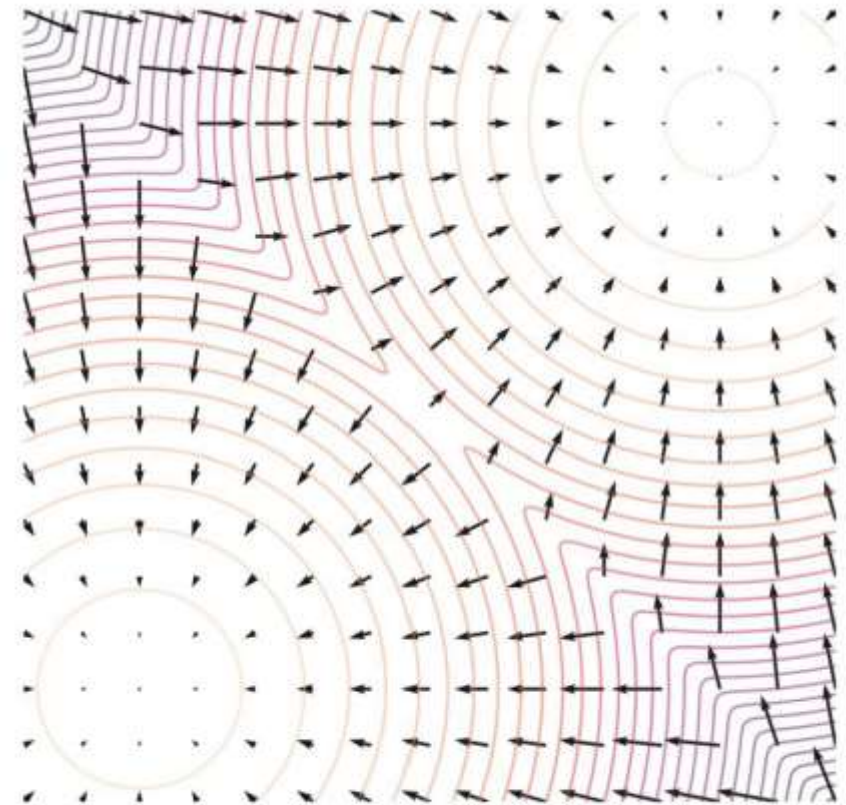# Score function example

- Probability density function (pdf)
$$p(\mathrm{X})$$

- Score function
$$\nabla_{\mathbf{x}} \log p(\mathrm{X})$$

- e.g. Gaussian distribution

$$p_\theta(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
$$\longrightarrow s_\theta(x) = -\frac{x-\mu}{\sigma^2}$$



Density function: Contours
Score function: Vector field

# Score-based Diffusion Models

- **The Generative Reverse Stochastic Differential Equation**



**But how to get the score function** $\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$**?**

Forward Diffusion SDE:
$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\, dt + \sqrt{\beta(t)}\, d\omega_t$$
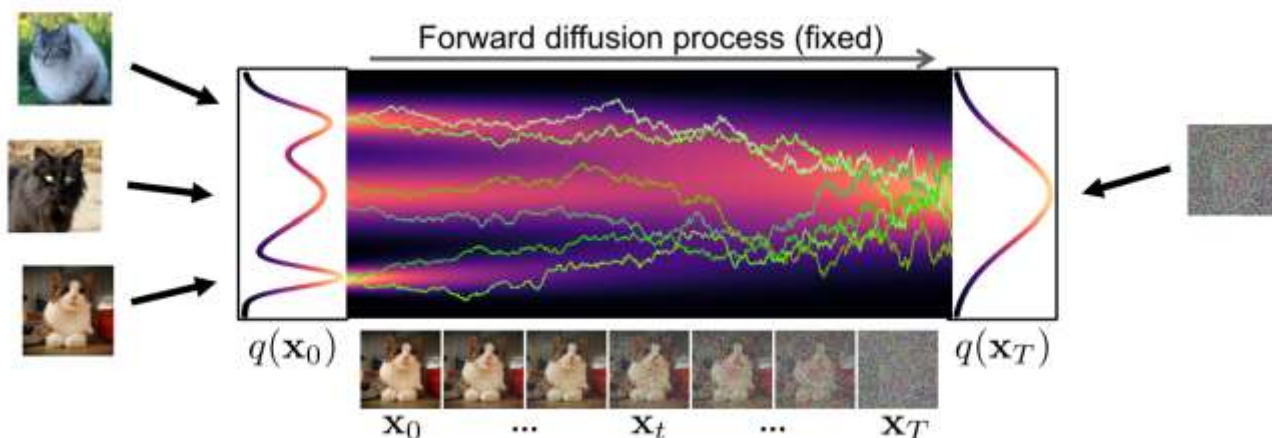
drift term

diffusion term

Reverse Generative Diffusion SDE:
$$d\mathbf{x}_t = \left[ -\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t)\nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t)\right] dt + \sqrt{\beta(t)}\, d\bar{\omega}_t$$

"Score Function"

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

# Score Matching

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



$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{\mathbf{x}_t \sim q_t(\mathbf{x}_t)} ||\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)||_2^2$$

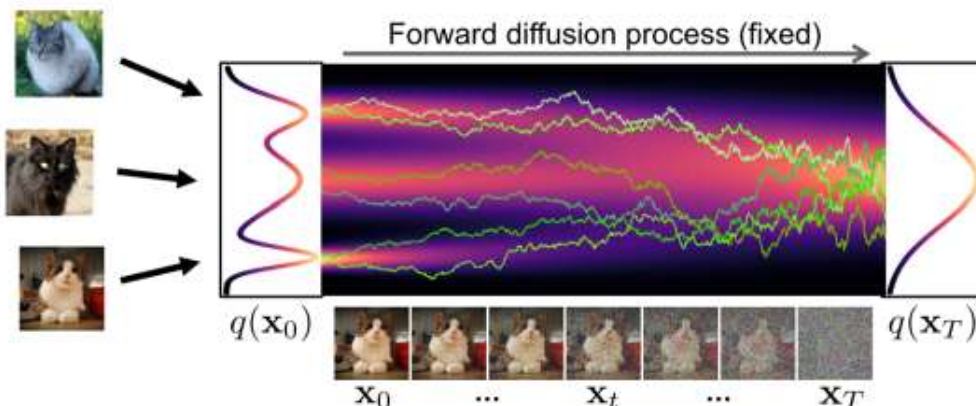| diffusion time $t$ | diffused data $\mathbf{x}_t$ | neural network | score of diffused data (marginal) |

➡ **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 Trackable!



"Variance Preserving" SDE:
$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,dt + \sqrt{\beta(t)}\,d\boldsymbol{\omega}_t$$
$$q_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma_t\mathbf{x}_0, \sigma_t^2\mathbf{I})$$

- Denoising Score Matching:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim q_0(\mathbf{x}_0)}\mathbb{E}_{\mathbf{x}_t\sim q_t(\mathbf{x}_t|\mathbf{x}_0)}||\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t|\mathbf{x}_0)||_2^2$$

| diffusion time $t$ | data sample $\mathbf{x}_0$ | diffused data sample $\mathbf{x}_t$ | neural network | score of diffused data sample |
|---|---|---|---|---|

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

# Denoising Score Matching

- Matching the score of a noise-perturbed distribution

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

$$\mathbf{X}$$

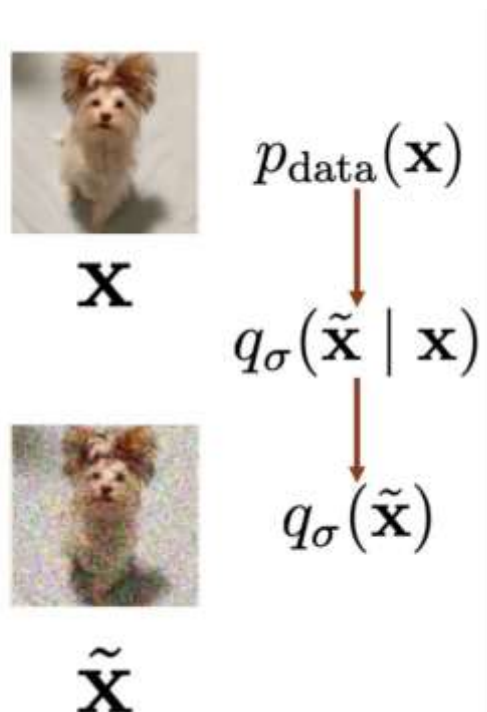$$q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})$$

$$q_\sigma(\tilde{\mathbf{x}})$$

$$\tilde{\mathbf{X}}$$

$$\frac{1}{2} E_{\tilde{\mathbf{x}} \sim q_\sigma}[\|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}) - \boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2]$$

$$= \frac{1}{2} \int q_\sigma(\tilde{\mathbf{x}}) \|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}) - \boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2 \, d\tilde{\mathbf{x}}$$

$$= \frac{1}{2} \int q_\sigma(\tilde{\mathbf{x}}) \|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})\|_2^2 \, d\tilde{\mathbf{x}} + \frac{1}{2} \int q_\sigma(\tilde{\mathbf{x}}) \|\boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2 \, d\tilde{\mathbf{x}}$$

$$- \int q_\sigma(\tilde{\mathbf{x}}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} \boldsymbol{s}_\theta(\tilde{\mathbf{x}}) \, d\tilde{\mathbf{x}}$$

$$= \text{const.} + \frac{1}{2} E_{\tilde{\mathbf{x}} \sim q_\sigma}[\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2] - \int q_\sigma(\tilde{\mathbf{x}}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} \boldsymbol{s}_\theta(\tilde{\mathbf{x}}) \, d\tilde{\mathbf{x}}$$

# Denoising Score Matching

■ Matching the score of a noise-perturbed distribution



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

$$\mathbf{x}$$

$$q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})$$

$$q_\sigma(\tilde{\mathbf{x}})$$

$$\tilde{\mathbf{x}}$$

$$-\int q_\sigma(\tilde{\mathbf{x}}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\int q_\sigma(\tilde{\mathbf{x}}) \frac{1}{q_\sigma(\tilde{\mathbf{x}})} \nabla_{\tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\int \nabla_{\tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\int \nabla_{\tilde{\mathbf{x}}} \left( \int p_{\text{data}}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x}) \, \mathrm{d}\mathbf{x} \right)^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\int \left( \int p_{\text{data}}(\mathbf{x}) \nabla_{\tilde{\mathbf{x}}} q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x}) \, \mathrm{d}\mathbf{x} \right)^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\int \left( \int p_{\text{data}}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x}) \, \mathrm{d}\mathbf{x} \right)^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -\iint p_{\text{data}}(\mathbf{x}) q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\mathbf{x} \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= -E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})} \left[ \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})^\mathsf{T} s_\theta(\tilde{\mathbf{x}}) \right]$$

# Denoising Score Matching

- Matching the score of a noise-perturbed distribution

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

$\mathbf{x}$

$q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})$

$q_\sigma(\tilde{\mathbf{x}})$
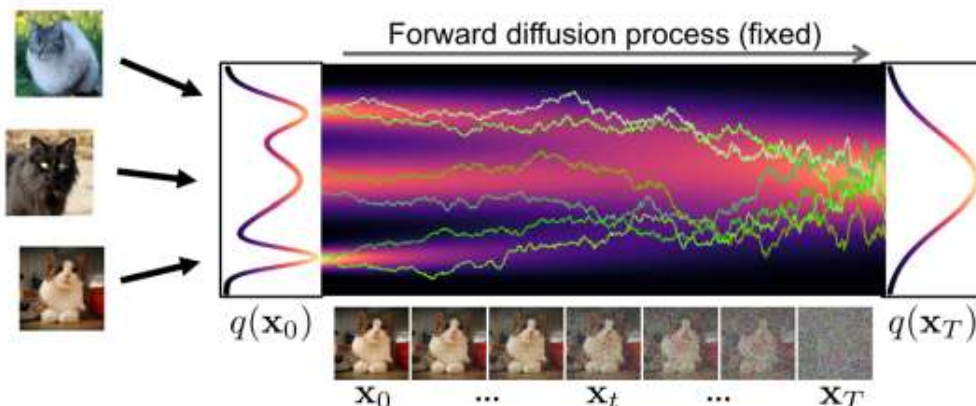
$\tilde{\mathbf{x}}$

$$\frac{1}{2} E_{\tilde{\mathbf{x}} \sim q_\sigma} [\|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}) - \boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2]$$

$$= \text{const.} + \frac{1}{2} E_{\tilde{\mathbf{x}} \sim q_\sigma} [\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2] - \int q_\sigma(\tilde{\mathbf{x}}) \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}})^\mathsf{T} \boldsymbol{s}_\theta(\tilde{\mathbf{x}}) \, \mathrm{d}\tilde{\mathbf{x}}$$

$$= \text{const.} + \frac{1}{2} E_{\tilde{\mathbf{x}} \sim q_\sigma} [\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}})\|_2^2] - E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})^\mathsf{T} \boldsymbol{s}_\theta(\tilde{\mathbf{x}})]$$

$$= \text{const.} + \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})\|_2^2]$$

$$- \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})\|_2^2]$$

$$= \text{const.} + \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})\|_2^2] + \text{const.}$$

$$= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x}), \tilde{\mathbf{x}} \sim q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} [\|\boldsymbol{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} \mid \mathbf{x})\|_2^2] + \text{const.}$$

# Denoising Score Matching

- Instead, diffuse individual data points $\mathbf{x}_0$

  Diffused $q_t(\mathbf{x}_t|\mathbf{x}_0)$ is Trackable!



"Variance Preserving" SDE:
$$\mathrm{d}\mathbf{x}_t = -\frac{1}{2}\beta(t)\mathbf{x}_t\,\mathrm{d}t + \sqrt{\beta(t)}\,\mathrm{d}\boldsymbol{\omega}_t$$
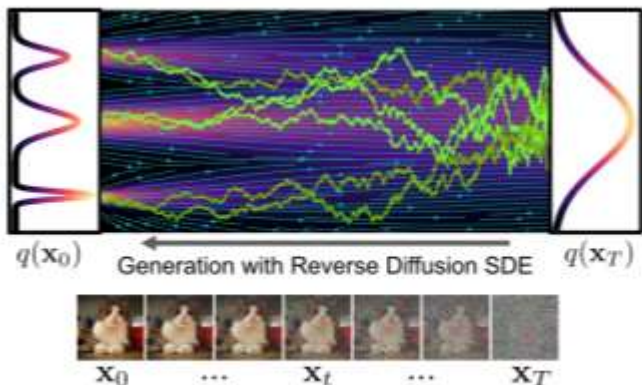$$q_t(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t;\gamma_t\mathbf{x}_0,\sigma_t^2\mathbf{I})$$

Forward diffusion process (fixed)

$q(\mathbf{x}_0)$ $\quad$ $q(\mathbf{x}_T)$

$\mathbf{x}_0$ $\quad$ ... $\quad$ $\mathbf{x}_t$ $\quad$ ... $\quad$ $\mathbf{x}_T$

- Denoising Score Matching:

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim q_0(\mathbf{x}_0)}\mathbb{E}_{\mathbf{x}_t\sim q_t(\mathbf{x}_t|\mathbf{x}_0)}||\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}_t,t) - \nabla_{\mathbf{x}_t}\log q_t(\mathbf{x}_t|\mathbf{x}_0)||_2^2$$

| diffusion time $t$ | data sample $\mathbf{x}_0$ | diffused data sample $\mathbf{x}_t$ | neural network | score of diffused data sample |

$$\Rightarrow \quad \min_{\boldsymbol{\theta}} \mathbb{E}_{t\sim\mathcal{U}(0,T)}\mathbb{E}_{\mathbf{x}_0\sim q_0(\mathbf{x}_0)}\mathbb{E}_{\boldsymbol{\epsilon}\sim\mathcal{N}(\mathbf{0},\mathbf{I})}\frac{1}{\sigma_t^2}||\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t,t)||_2^2$$

# Synthesis with SDE vs. ODE
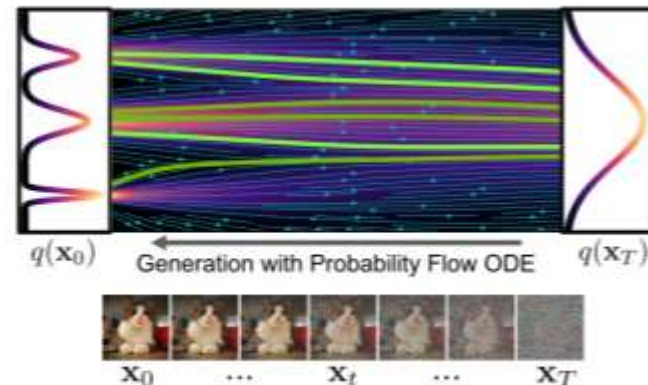
■ How to solve the generative SDE or ODE in practice?



$q(\mathbf{x}_0)$   Generation with Reverse Diffusion SDE   $q(\mathbf{x}_T)$

$\mathbf{x}_0$   ...   $\mathbf{x}_t$   ...   $\mathbf{x}_T$

**Generative Diffusion SDE:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\left[\mathbf{x}_t + 2\mathbf{s}_\theta(\mathbf{x}_t, t)\right] dt + \sqrt{\beta(t)}\, d\bar{\boldsymbol{\omega}}_t$$

➡ *Euler-Maruyama:*

$$\mathbf{x}_{t-1} = \mathbf{x}_t + \frac{1}{2}\beta(t)\left[\mathbf{x}_t + 2\mathbf{s}_\theta(\mathbf{x}_t, t)\right]\Delta t + \sqrt{\beta(t)\Delta t}\,\mathcal{N}(\mathbf{0}, \mathbf{I})$$

➡ *Ancestral Sampling* (Part 1) is also a generative SDE sampler!



$q(\mathbf{x}_0)$   Generation with Probability Flow ODE   $q(\mathbf{x}_T)$

$\mathbf{x}_0$   ...   $\mathbf{x}_t$   ...   $\mathbf{x}_T$

**Probability Flow ODE:**

$$d\mathbf{x}_t = -\frac{1}{2}\beta(t)\left[\mathbf{x}_t + \mathbf{s}_\theta(\mathbf{x}_t, t)\right] dt$$

➡ *Euler's Method:*

$$\mathbf{x}_{t-1} = \mathbf{x}_t + \frac{1}{2}\beta(t)\left[\mathbf{x}_t + \mathbf{s}_\theta(\mathbf{x}_t, t)\right]\Delta t$$

➡ In practice: Higher-Order ODE solvers (*Runge-Kutta, linear multistep methods, exponential integrators, ...*)
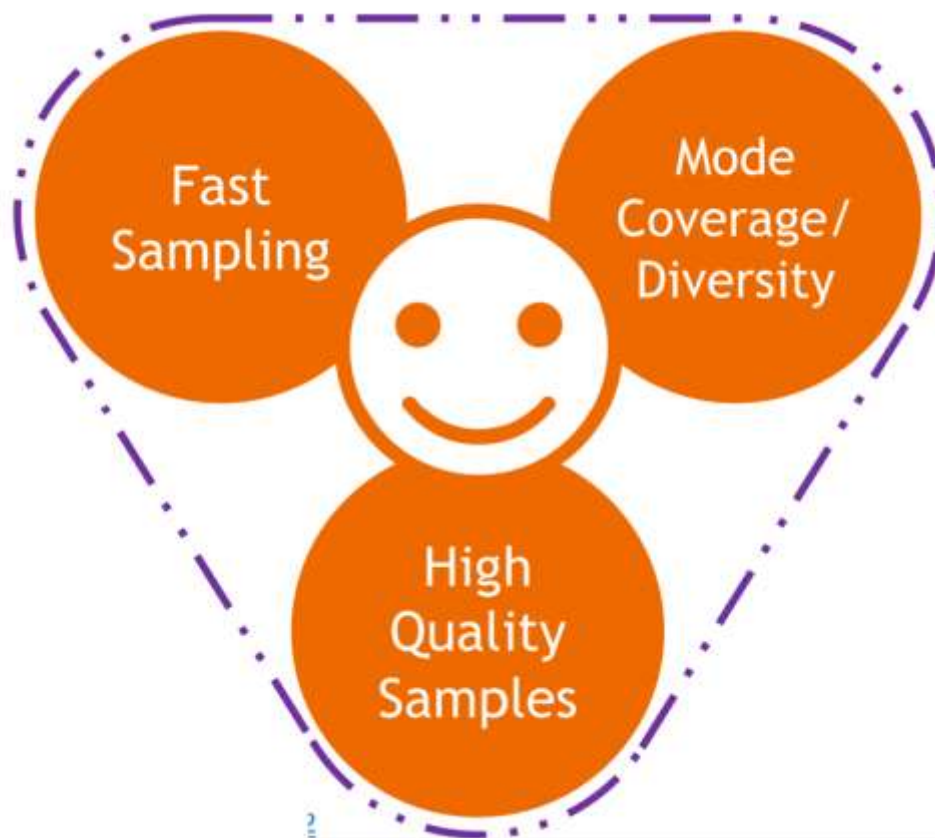
# Synthesis with SDE vs. ODE

- How to solve the generative SDE or ODE in practice?

- Runge-Kutta adaptive step-size ODE solver
- Higher-Order adaptive step-size SDE solver
- Reparametrized, smoother ODE
- Higher-Order ODE solver with linear multistepping
- Exponential ODE Integrators
- Higher-order ODE solver with Heun's Method
- …….

# Outline

- Diffusion: Denoising Diffusion Probabilistic Models

- Score-based Diffusion Models

- **Accelerated Sampling**

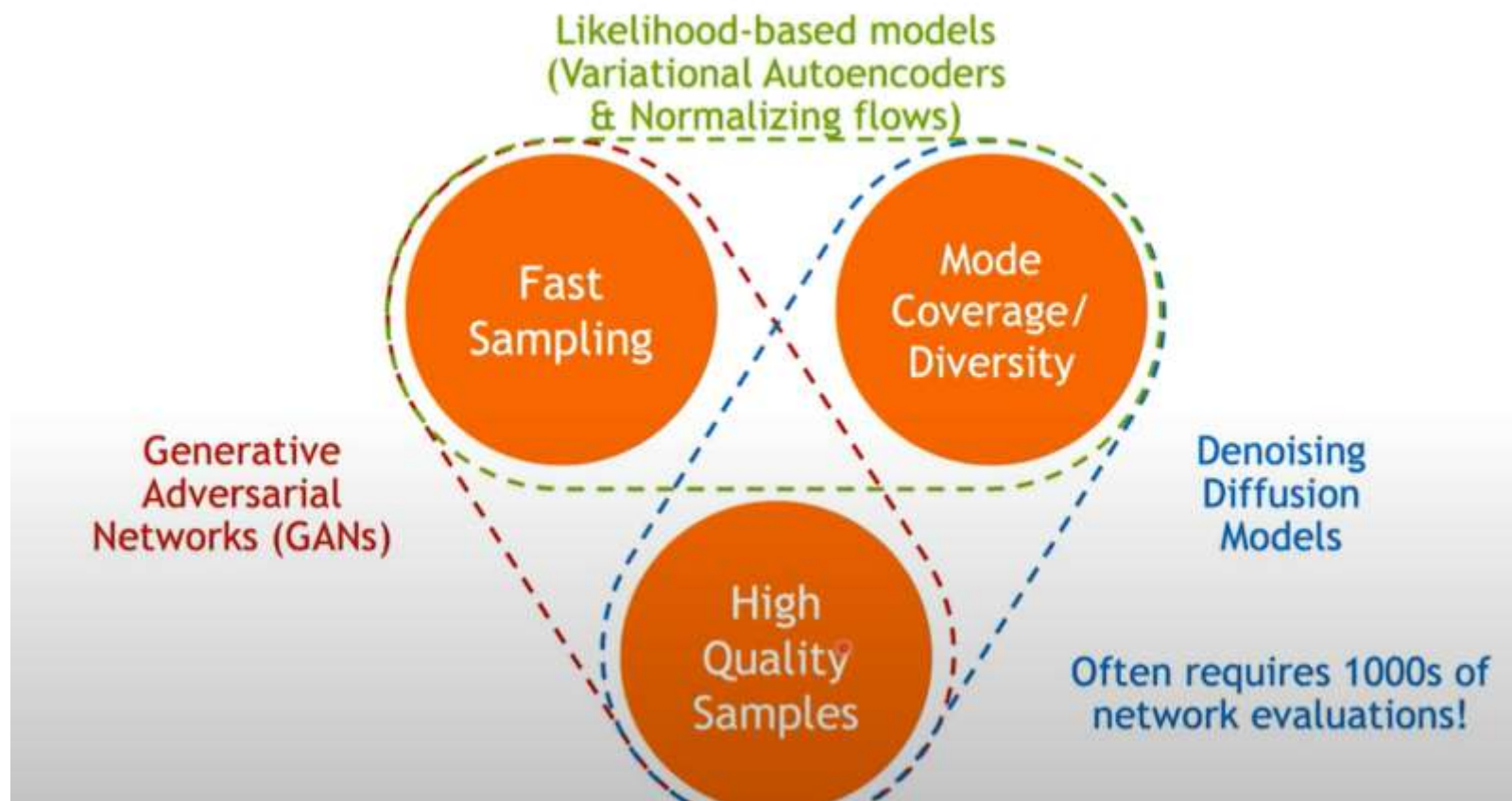- Conditional Generation and Guidance

# What makes a good generative model?

- The generative learning trilemma
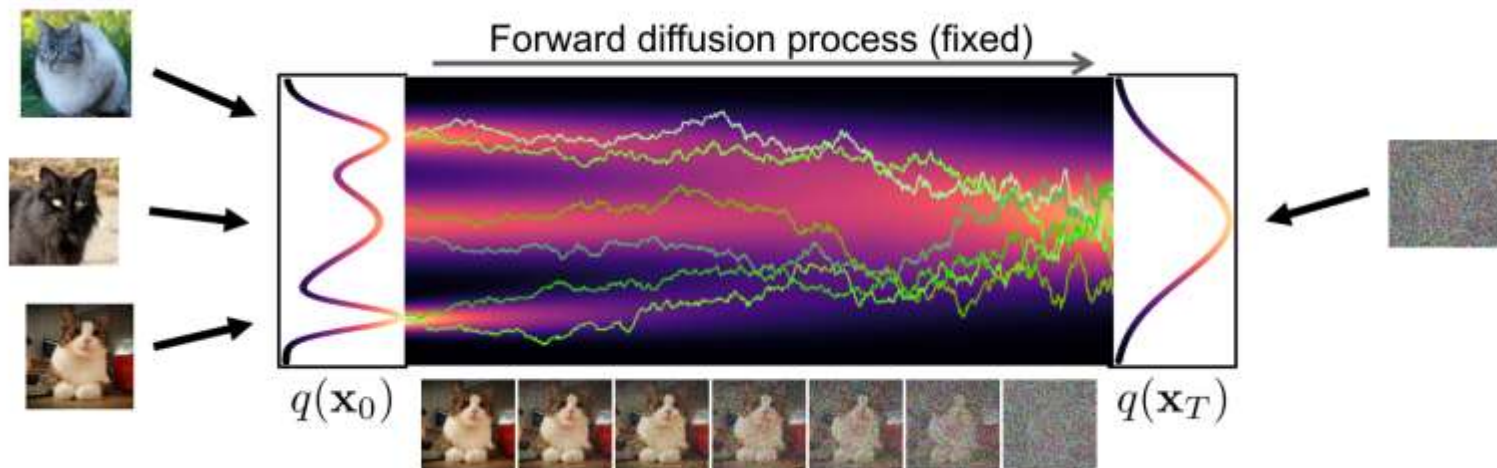- Tackle the trilemma by accelerating diffusion models

# What makes a good generative model?

- The generative learning trilemma
- Tackle the trilemma by accelerating diffusion models
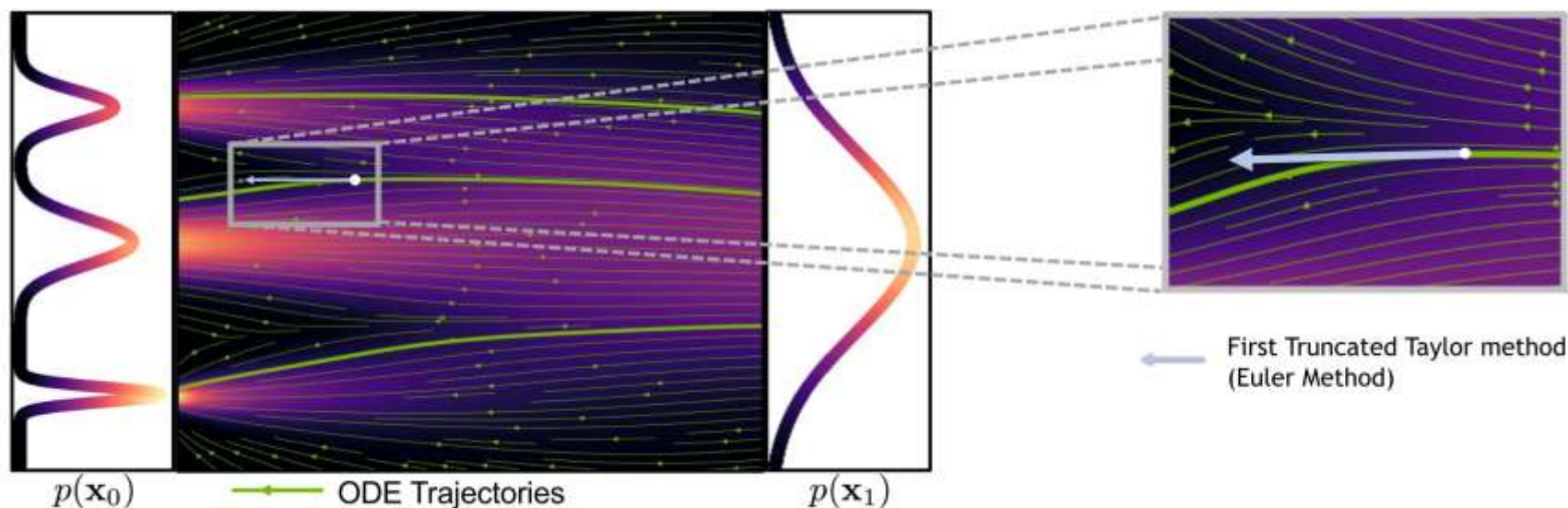
# Accelerated Sampling

- Advanced ODE/SDE Solvers
- Distillation Techniques
- Low-dim Diffusion Processes

- … …



Forward diffusion process (fixed)

$q(\mathbf{x}_0)$ $q(\mathbf{x}_T)$

# Generative ODEs

- Solve ODEs with as little function evaluations as possible!

$$dx = \epsilon_\theta(x, t)dt$$



First Truncated Taylor method (Euler Method)

$p(\mathbf{x}_0)$    ODE Trajectories    $p(\mathbf{x}_1)$

Song et al., "Denoising Diffusion Implicit Models (DDIM)", ICLR 2021
https://arxiv.org/abs/2010.02502

# Generative ODEs

- Solve ODEs with as little function evaluations as possible!
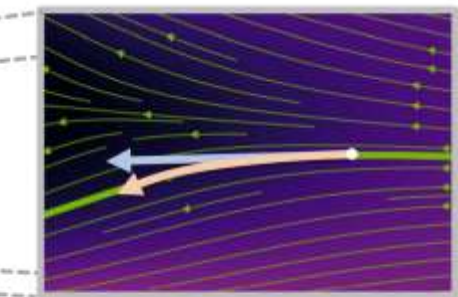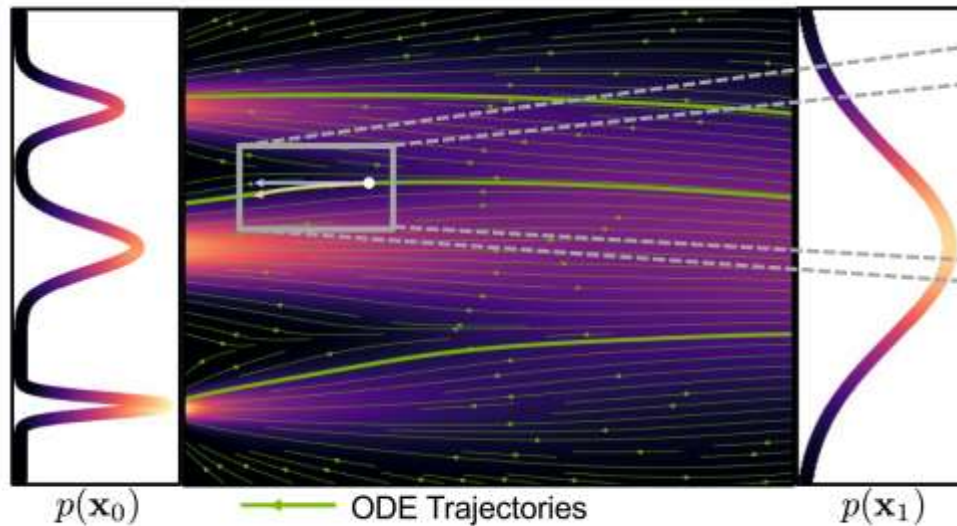
$$dx = \epsilon_\theta(x, t)dt$$



First Truncated Taylor method (Euler Method)

Higher order methods (RK4, Multistepping, Heun)

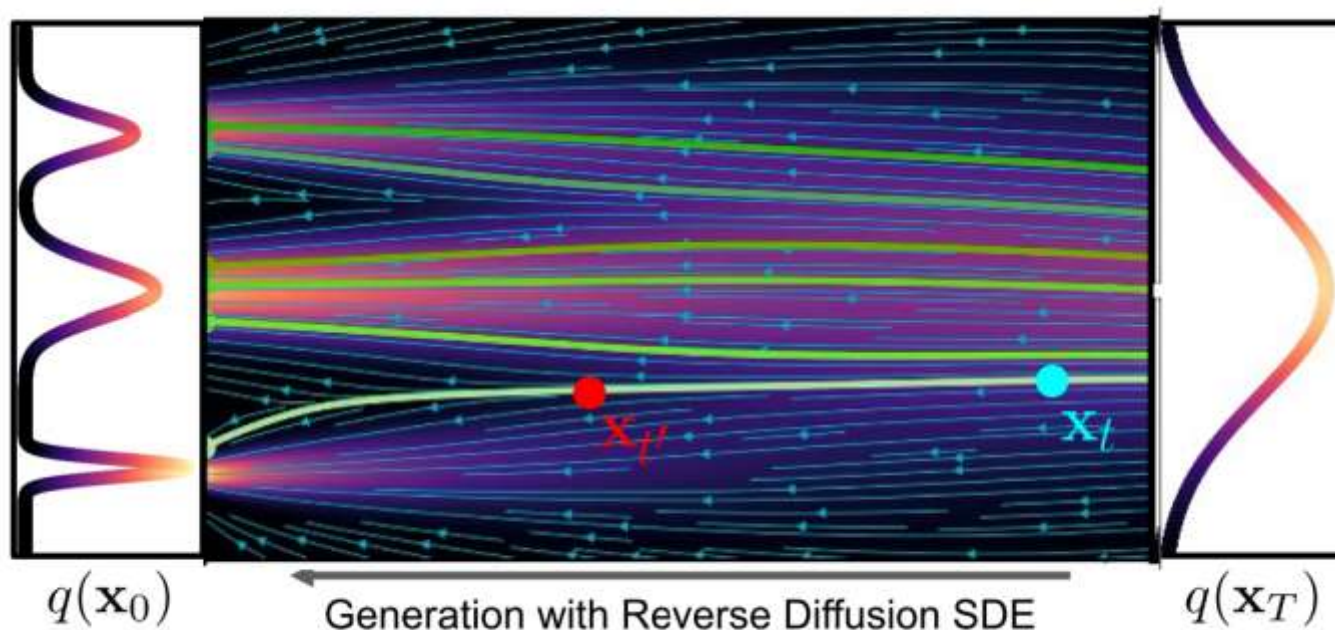Approximate the higher order terms numerically

# Generative ODEs

- Solve ODEs with as little function evaluations as possible!

$$dx = \epsilon_\theta(x, t)dt$$

  □ Runge-Kutta adaptive step-size ODE solver:

  https://arxiv.org/abs/2011.13456

  □ Higher-Order adaptive step-size SDE solver:

  https://arxiv.org/abs/2105.14080

  □ Reparametrized, smoother ODE:  gDDIM

  https://arxiv.org/abs/2206.05564

  □ Higher-Order ODE solver with linear multistepping:

  https://arxiv.org/abs/2202.09778

  □ Exponential ODE Integrators: DPM, DPM-Solver++

  https://arxiv.org/abs/2206.00927

  □ Higher-Order ODE solver with Heun's Method:

  https://arxiv.org/abs/2206.00364

  □ Many more:  https://arxiv.org/abs/2305.19947

# Distillation Techniques

- ODE Distillation



$q(\mathbf{x}_0)$     Generation with Reverse Diffusion SDE     $q(\mathbf{x}_T)$
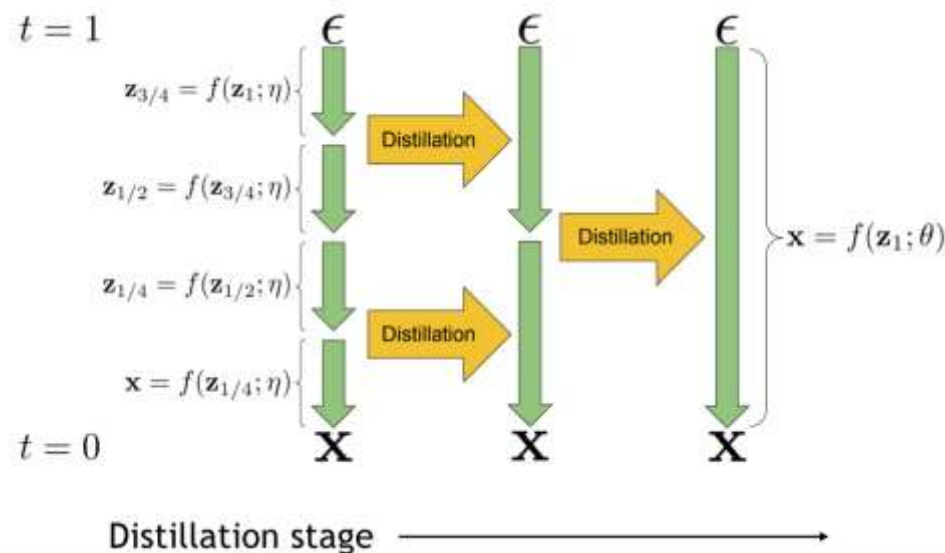
Can we train a neural network to directly predict $\mathbf{x}_{t'}$ given $\mathbf{x}_t$ ?

# Progressive Distillation

- Distill a deterministic ODE 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.



Salimans & Ho, "Progressive distillation for fast sampling of diffusion models", ICLR 2022. https://arxiv.org/abs/2202.00512
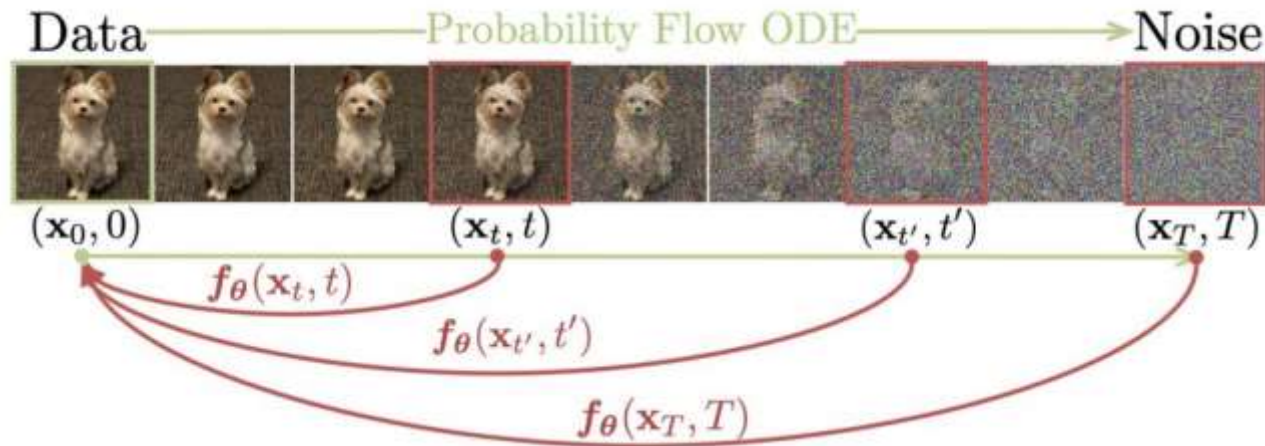
# Progressive Distillation in Latent Space



(a) 2 denoising steps    (b) 4 denoising steps    (c) 8 denoising steps

Meng et al., "On Distillation of Guided Diffusion Models", CVPR 2023 (Award Candidate). https://arxiv.org/abs/2210.03142

# Consistency Distillation



Points on the same trajectory should generated the same $\mathbf{x}_0$

Assume $f_\theta(\mathbf{x}_t, t)$ is the current estimation of $\mathbf{x}_0$

Basic idea:

- Find $\mathbf{x}_t$ and $\mathbf{x}_{t'}$ on a trajectory by solving generative ODE in $[t, t']$

- Minimize:

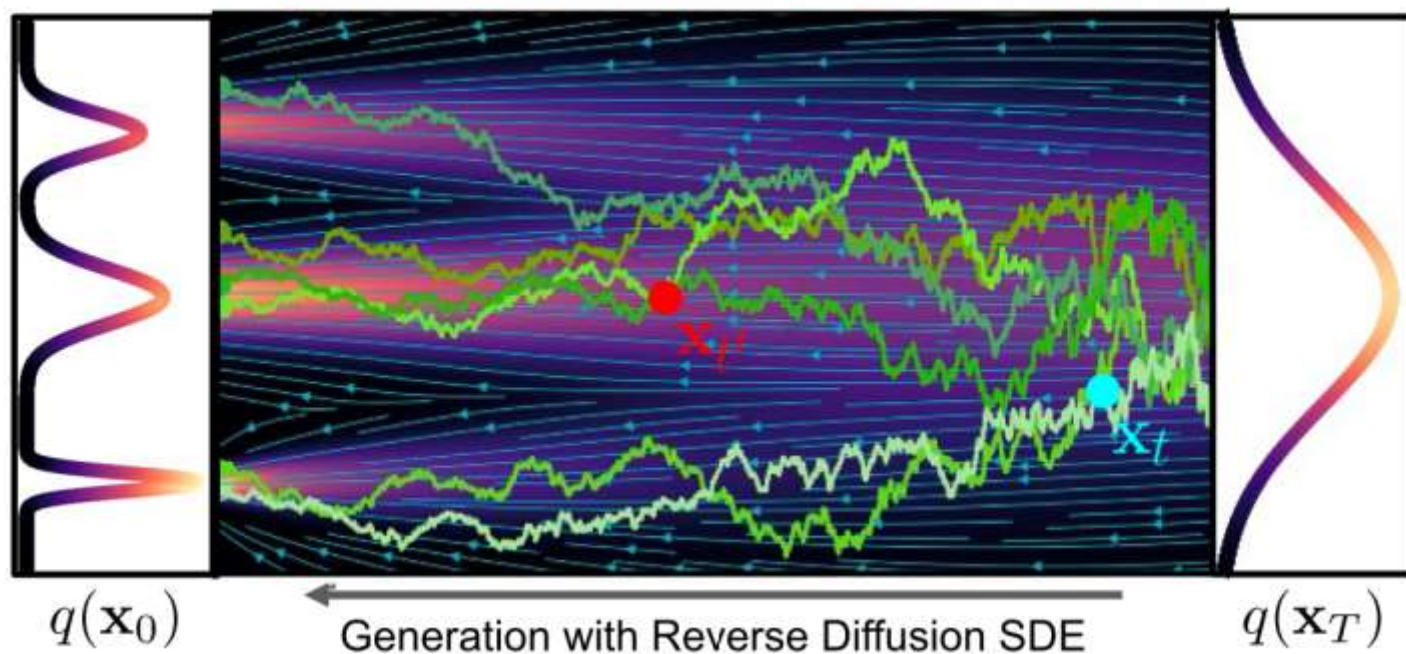$$\min_\theta \ ||f_{\text{EMA}}(\mathbf{x}_t, t) - f_\theta(\mathbf{x}_t', t')||_2^2$$

Song et al., Consistency Models, ICML 2023,
https://arxiv.org/abs/2303.01469

# Distillation Techniques

- ## SDE Distillation



$q(\mathbf{x}_0)$ ← Generation with Reverse Diffusion SDE → $q(\mathbf{x}_T)$
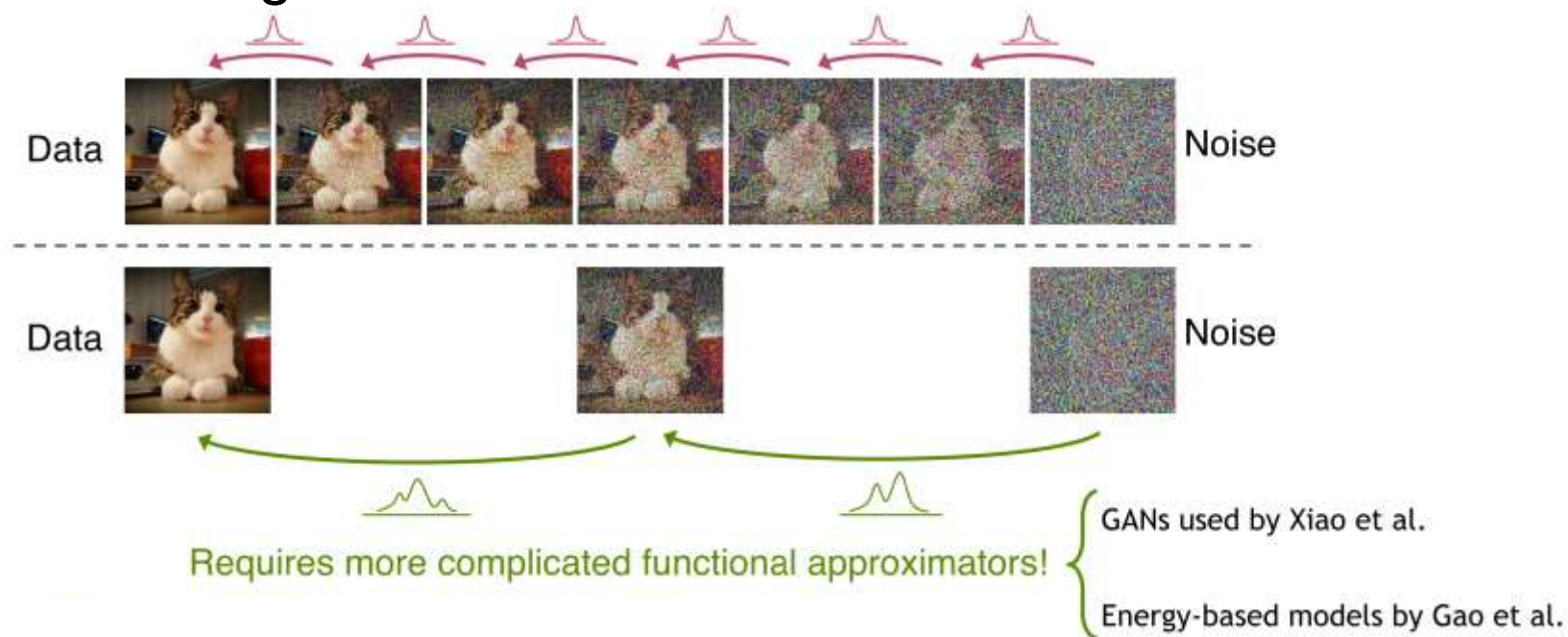
Can we train a neural network to directly predict $\mathbf{x}_{t'}$ given $\mathbf{x}_t$ ?

# Approximation of Reverse Process

- Normal assumption in denoising distribution holds only for small step
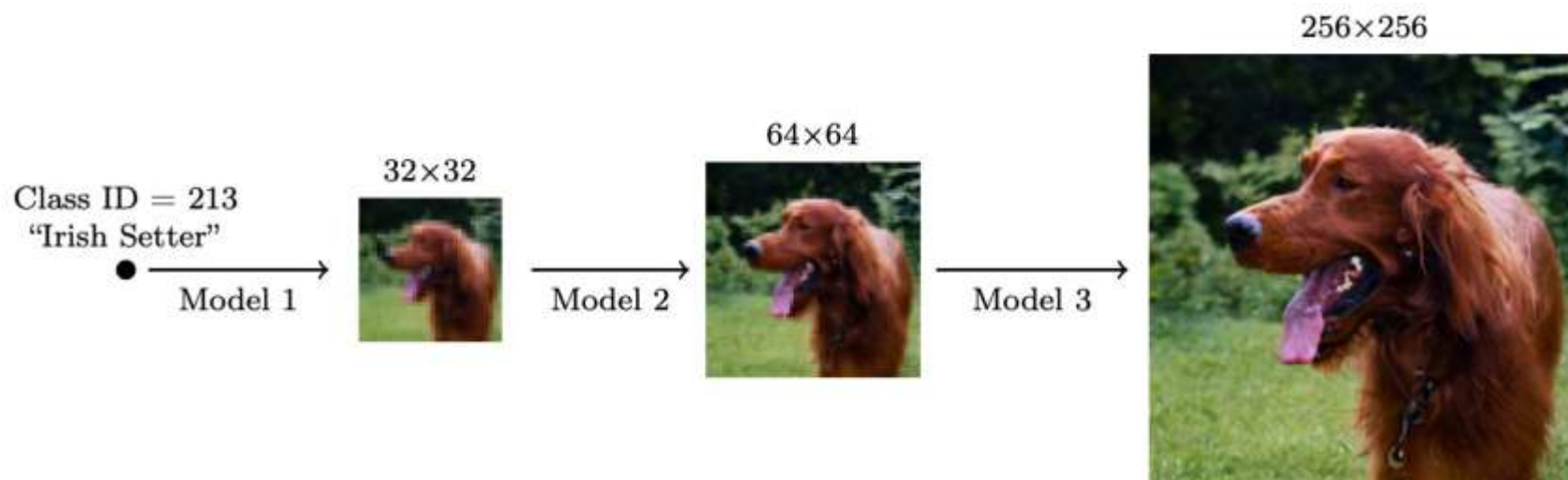- Denoising Process with Uni-modal Normal Distribution



Xiao et al., "Tackling the Generative Learning Trilemma with Denoising Diffusion GANs", ICLR 2022.
Gao et al., "Learning energy-based models by diffusion recovery likelihood", ICLR 2021.
https://arxiv.org/abs/2112.07804

# Low-dim. Diffusion Processes

- Cascaded Generation
- Cascaded Diffusion Models outperform Big-GAN in FID and IS and VQ-VAE2 in Classification Accuracy Score.



Ho et al., "Cascaded Diffusion Models for High Fidelity Image Generation", 2021. https://arxiv.org/abs/2106.15282
Ramesh et al., "Hierarchical Text-Conditional Image Generation with CLIP Latents", 2022,
https://arxiv.org/abs/2204.06125
Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", 2022,
https://arxiv.org/abs/2205.11487

# Latent Diffusion Models

- VAE + score-based prior



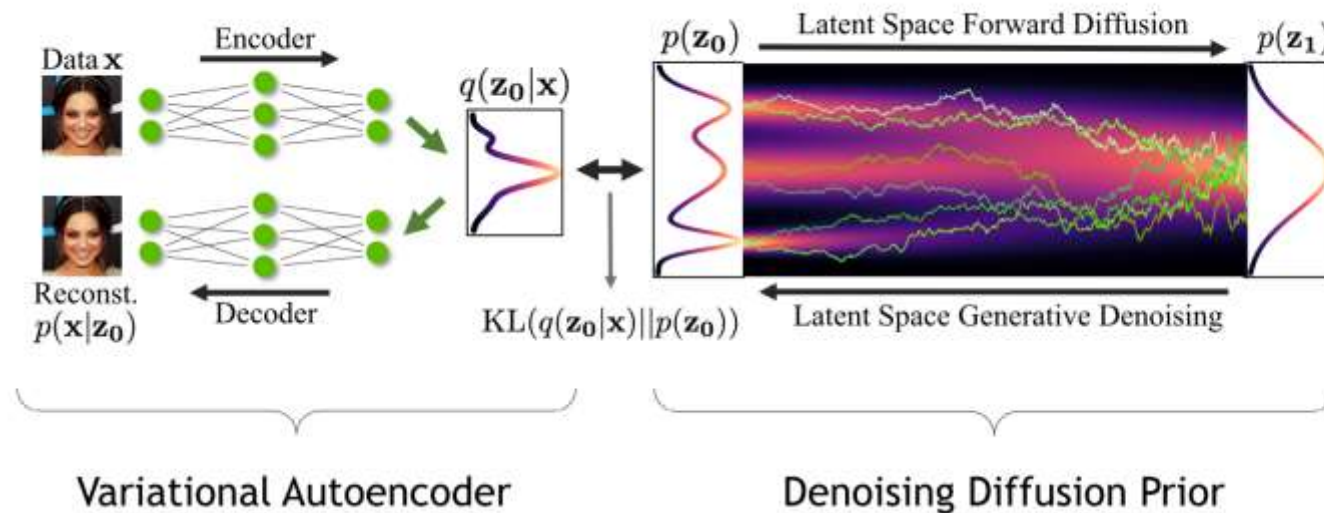**Main Idea**

Encoder maps the input data to an embedding space

Denoising diffusion models are applied in the latent space

Vahdat et al., "Score-based generative modeling in latent space",
NeurIPS 2021. https://arxiv.org/abs/2106.05931

# Latent Diffusion Models

- ## VAE + score-based prior



| | |
|---|---|
| Data $\mathbf{x}$ — Encoder → | $q(\mathbf{z_0}|\mathbf{x})$ |
| Reconst. $p(\mathbf{x}|\mathbf{z_0})$ ← Decoder | $KL(q(\mathbf{z_0}|\mathbf{x})||p(\mathbf{z_0}))$ |

Latent Space Forward Diffusion — $p(\mathbf{z_0})$ ... $p(\mathbf{z_1})$

Latent Space Generative Denoising

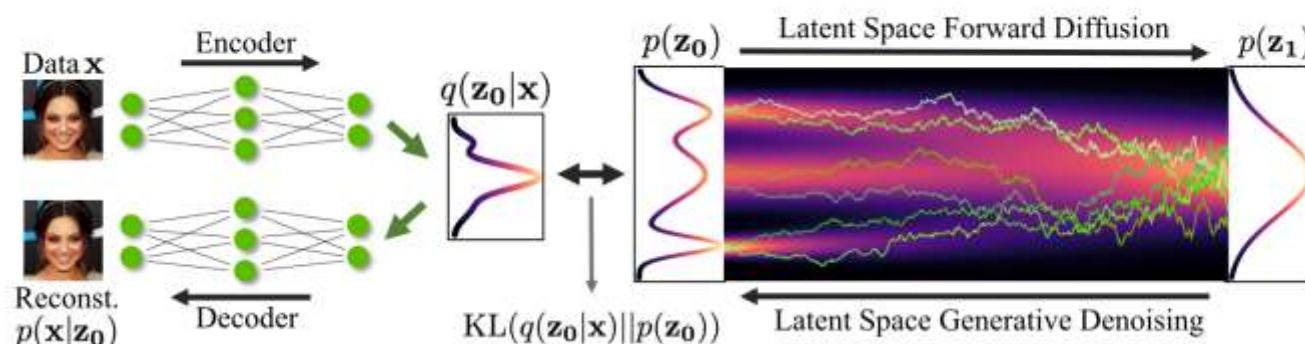**Variational Autoencoder**          **Denoising Diffusion Prior**

### Advantages:

(1) The distribution of latent embeddings close to Normal distribution → *Simpler denoising, Faster synthesis!*

(2) Latent space → *More expressivity and flexibility in design!*

(3) Tailored Autoencoders → *More expressivity, Application to any data type (graphs, text, 3D data, etc.) !*

Vahdat et al., "Score-based generative modeling in latent space",
NeurIPS 2021. https://arxiv.org/abs/2106.05931

# Latent Diffusion Models

- End-to-End Training objective



$$\mathcal{L}(\mathbf{x}, \boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\psi}) = \mathbb{E}_{q_\phi(\mathbf{z_0}|\mathbf{x})} \left[ -\log p_\psi(\mathbf{x}|\mathbf{z_0}) \right] + \mathrm{KL}(q_\phi(\mathbf{z_0}|\mathbf{x})||p_\theta(\mathbf{z_0}))$$

$$= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z_0}|\mathbf{x})} \left[ -\log p_\psi(\mathbf{x}|\mathbf{z_0}) \right]}_{\text{reconstruction term}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z_0}|\mathbf{x})} \left[ \log q_\phi(\mathbf{z_0}|\mathbf{x}) \right]}_{\text{negative encoder entropy}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z_0}|\mathbf{x})} \left[ -\log p_\theta(\mathbf{z_0}) \right]}_{\text{cross entropy}}$$
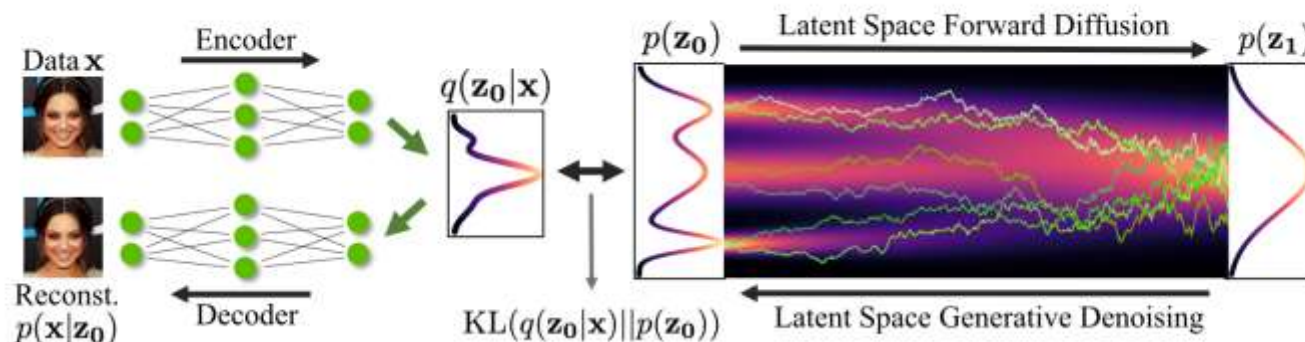
$$CE(q(\mathbf{z_0}|\mathbf{x})||p(\mathbf{z_0})) = \mathbb{E}_{t\sim\mathcal{U}[0,1]} \left[ \underbrace{\frac{g(t)^2}{2}}_{\text{time sampling}} \underbrace{\mathbb{E}_{q(\mathbf{z}_t,\mathbf{z_0}|\mathbf{x})}}_{\text{Forward diffusion}} \left[ ||\underbrace{\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t|\mathbf{z_0})}_{\text{Diffusion kernel}} - \underbrace{\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)}_{\substack{\text{Trainable} \\ \text{score function}}}||_2^2 \right] \right] + \underbrace{\frac{D}{2} \log \left( 2\pi e \sigma_0^2 \right)}_{\text{Constant}}$$

Vahdat et al., "Score-based generative modeling in latent space",
NeurIPS 2021. https://arxiv.org/abs/2106.05931

# Latent Diffusion Models

- ## Stable Diffusion: the efficiency and expressivity of LDM+ open-source access fueled a large body of work in the community



- Two stage training: train autoencoder first, then train the diffusion prior

- Focus on compression without of any loss in reconstruction quality

- Demonstrated the expressivity of latent diffusion models on many conditional problems

Rombach et al., "High-Resolution Image Synthesis with Latent Diffusion Models", CVPR 2022. https://arxiv.org/abs/2112.10752
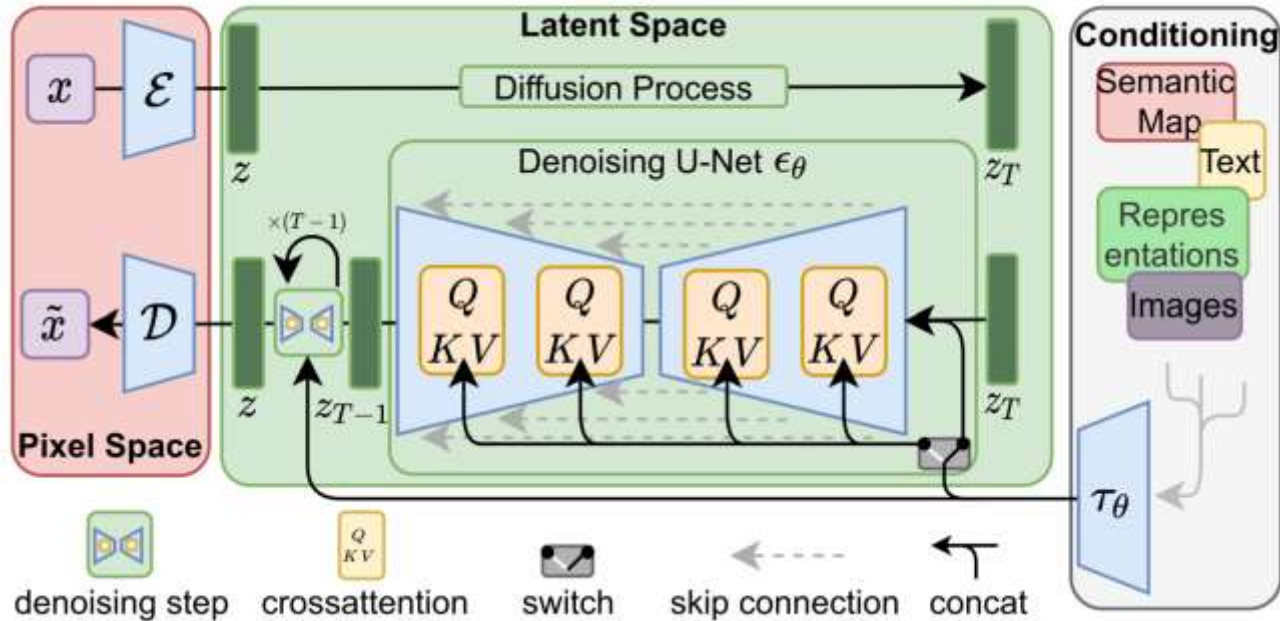
# Latent Diffusion Models

- Training models in the pixel space is excessively computationally expensive (can easily multiple days on a V100 GPU)

  - □ Even image synthesis is very slow compared to GANs
  - □ Images are high dimensional → more things to model

- Researchers observed that most "bits" of an image contribute to its perceptual characteristics since aggressively compressing it usually maintains its semantic and conceptual composition

  - □ In layman's terms, there are more bits for describing pixel-level details while less bits for describing "the meaning" within an image
  - □ Generative models should learn the latter

- Can we separate these two components?

# Huge success of text-2-img!

- Stable Diffusion model (CVPR2022)



**High-Resolution Image Synthesis with Latent Diffusion Models**

Robin Rombach[1] *      Andreas Blattmann[1] *      Dominik Lorenz[1]      Patrick Esser[R]      Björn Ommer[1]

[1]Ludwig Maximilian University of Munich & IWR, Heidelberg University, Germany      [R]Runway ML

https://github.com/CompVis/latent-diffusion

# Latent Diffusion Models

Latent Diffusion Models can be divided **into two stages**:

1）Training perceptual compression models that strip away irrelevant high-level details and learn a latent space that is semantically equivalent to the high level image pixel-space
   a. The loss is a combination of a reconstruction loss, an adversarial loss (remember GANs?) that promotes high quality decoder reconstruction, and regularization terms

$$L_{\text{Autoencoder}} = \min_{\mathcal{E},\mathcal{D}} \max_{\psi} \Big( L_{rec}(x, \mathcal{D}(\mathcal{E}(x))) - L_{adv}(\mathcal{D}(\mathcal{E}(x))) + \log D_{\psi}(x) + L_{reg}(x; \mathcal{E}, \mathcal{D}) \Big)$$

2）Performing a diffusion process *in this latent space*. There are several benefits to this:
   a. The diffusion process is only focusing on the relevant semantic bits of the data
   b. Performing diffusion in a low dimensional space is significantly more efficient

# Huge success of text-2-img!

- **Stable Diffusion model (CVPR2022)**
- **Long story between** Stability AI, Runway ML and LAION-5B
- AI paradigm: data + algorithm + computing resource



Computer Vision & Learning Group
Ludwig Maximilian University of Munich
(LMU)

~4000 A100 from
Stability AI

Huge text-image
dataset from
LAION

# Huge success of text-2-img!

- The multi-modality framework is important
- The trend continues: big data, big modal ……
- Enjoy better text-encoder and suitable generator



**Google: Parti Model**, "Scaling

Autoregressive Models for Content-

Rich Text-to-Image Generation"



**Stablity AI: DeepFloyd IF Model**

T5-XXL as Text-encoder; pixel-level

Diffusion

# Huge success of text-2-img!

- Enjoy cross-modality abilities
- Enjoy downstream conditioning abilities



Input images

Subject-Driven Generation using **DreamBooth**

in the Acropolis  in a doghouse  in a bucket  getting a haircut

Conditioning using **ControlNet**

"Swap sunflowers with roses"  "Add fireworks to the sky"  "Replace the fruits with cake"

Editing Instructions using **InstructPix2Pix** (based on GPT-3)

# Huge success of text-2-img-3D!

- Use NeRF as inherent representation to bridge 2D-DM with 3D scene
- More explicit disentanglement towards geometry, color, lighting …



3D Editing Instructions
using **InstructNeRF2NeRF**



NVDIA **Magic3D**    18 Nov 2022



OpenAI **Point-E**   21 Dec 2022

# Summary

- Preliminary theory of diffusion (don't worry if this is confusing!)

- Some tricks that modern diffusion models employ for image generation:

  - □ A U-Net architecture equipped with all kinds of modifications

  - □ Other architecture improvements

  - □ Several implementation tricks (different noise schedules, covariance parametrizations)

- Latent diffusion models for improving diffusion quality and efficiency