

**Deep Learning**

**Diffusion**

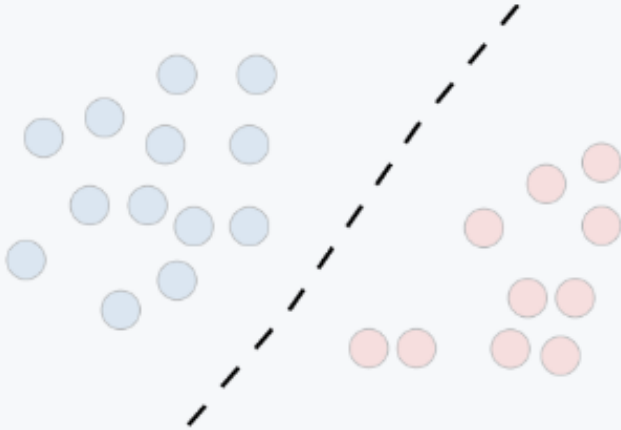
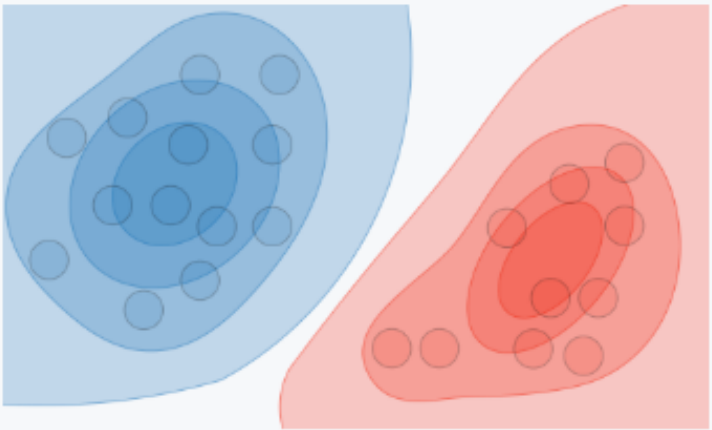
**Hao Chen**

**Fall 2024**

**Attendance: @**

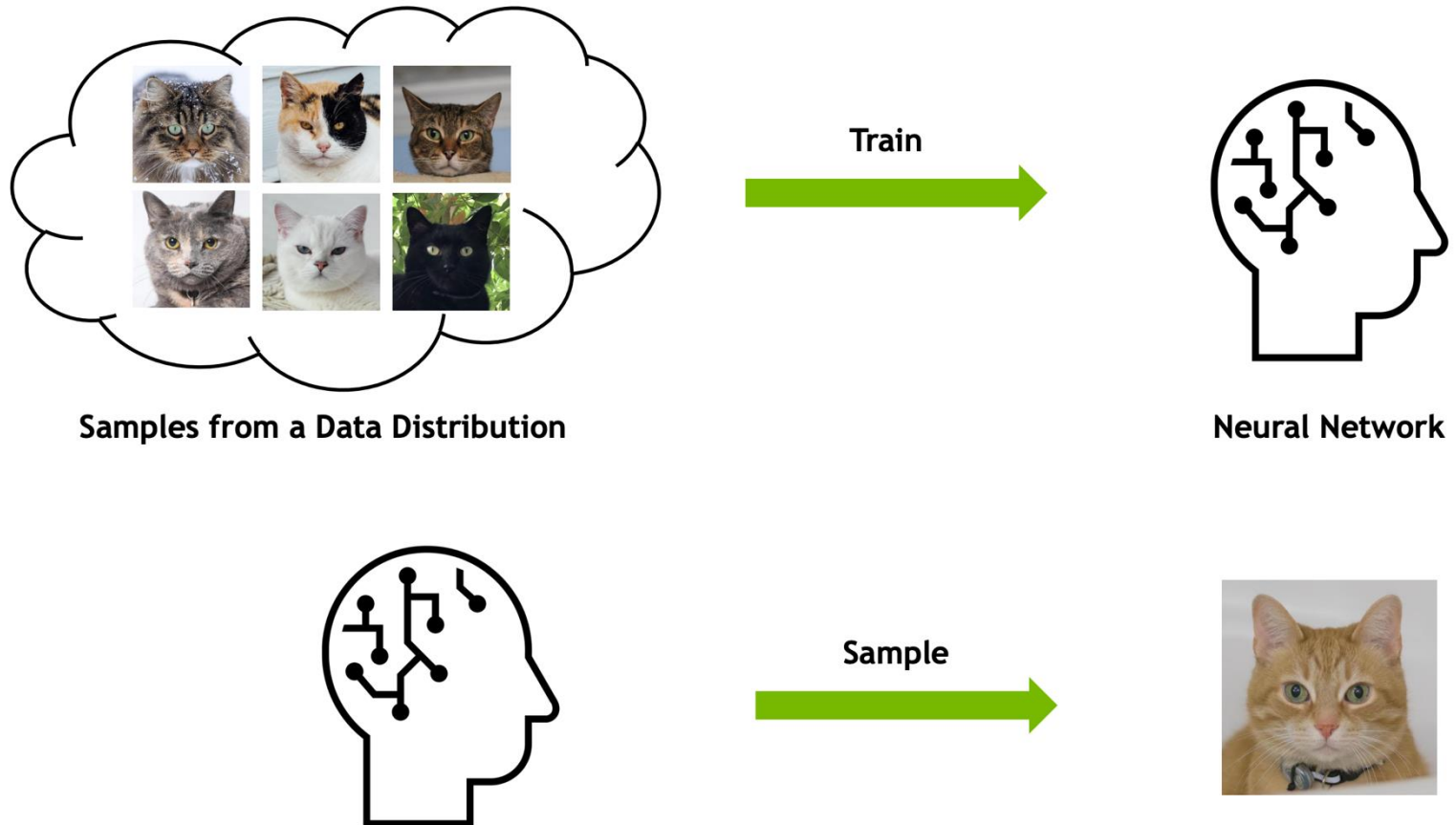
# Generative vs. Discriminative

- Generative models learn the data distribution

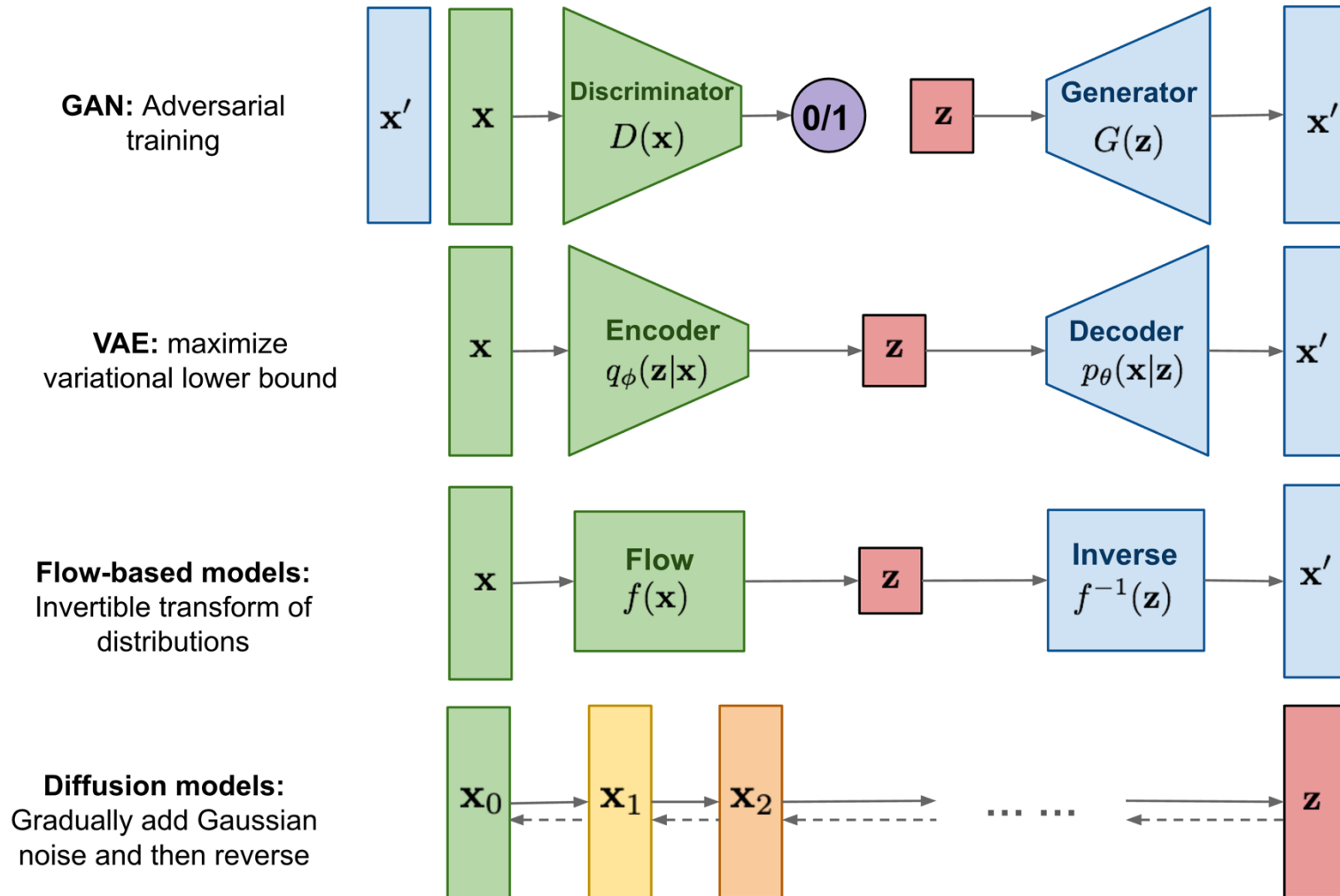
|                | Discriminative model  | Generative model   |
|----------------|---|--|
| Goal           | Directly estimate $P(y x)$  | Estimate $P(x y)$ to then deduce $P(y x)$  |
| What's learned | Decision boundary   | Probability distributions of the data  |
| Illustration   |  A scatter plot illustrating a discriminative model. It shows two classes of data points: blue circles on the left and red circles on the right. A dashed diagonal line represents the decision boundary separating the two classes. |  A scatter plot illustrating a generative model. It shows two classes of data points: blue circles on the left and red circles on the right. Each class is enclosed by a shaded region representing its probability distribution. The blue region is on the left and the red region is on the right, with some overlap between them. |

# Generative Models

- Learning to generate data

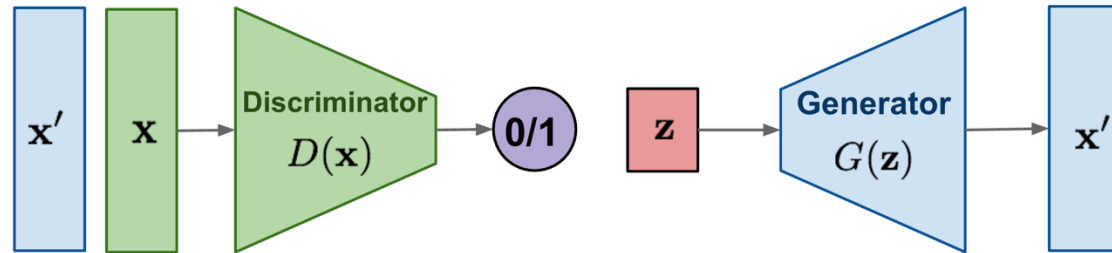


# Generative Models



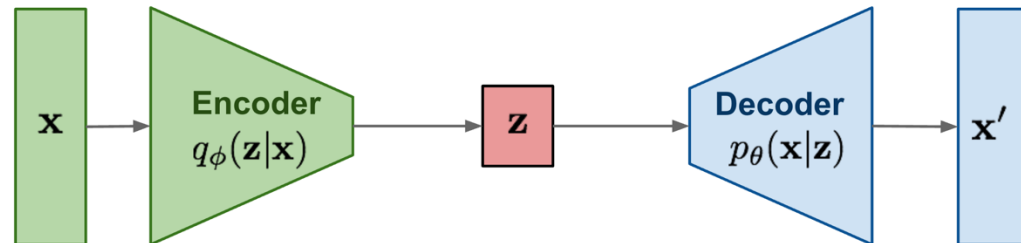
# Generative Models

**GAN:** Adversarial training

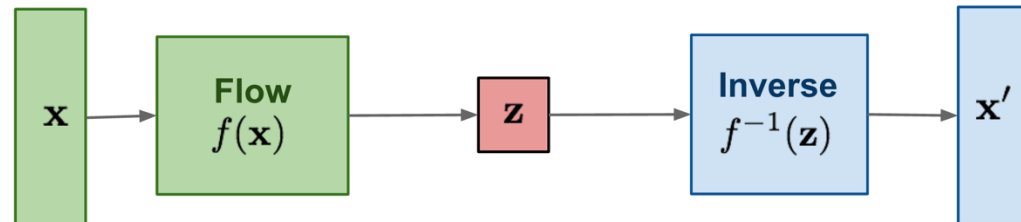


Last Lecture

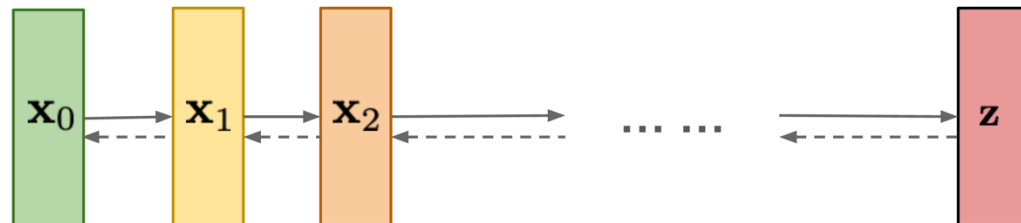
**VAE:** maximize variational lower bound



**Flow-based models:**  
Invertible transform of distributions

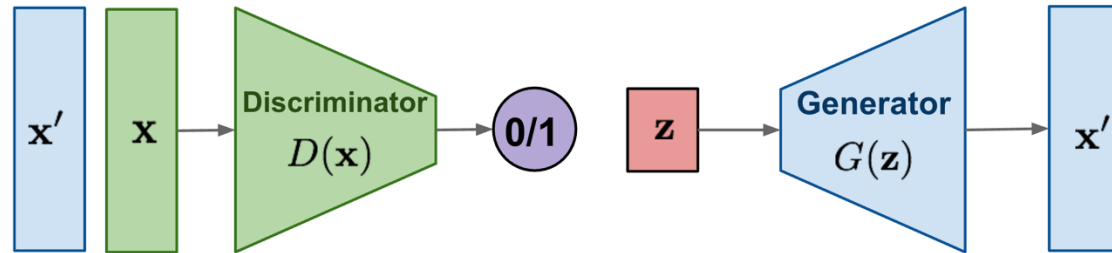


**Diffusion models:**  
Gradually add Gaussian noise and then reverse

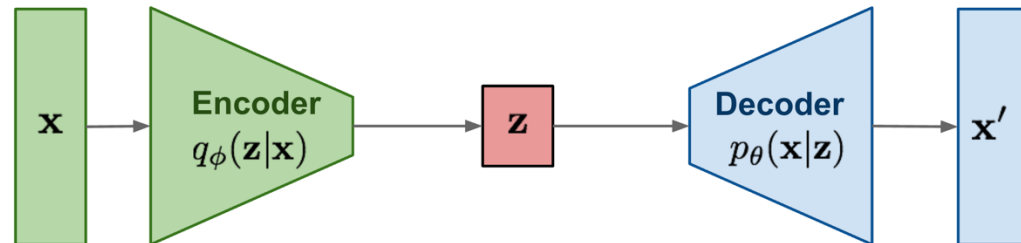


# Generative Models

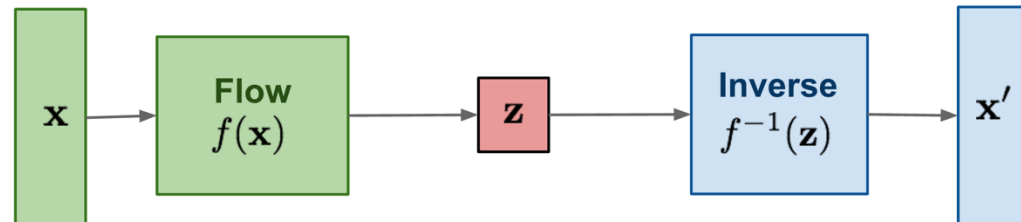
**GAN:** Adversarial training



**VAE:** maximize variational lower bound

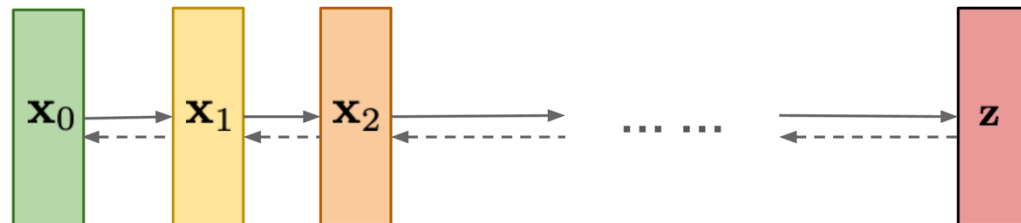


**Flow-based models:**  
Invertible transform of distributions



This Lecture

**Diffusion models:**  
Gradually add Gaussian noise and then reverse



# A Fast Evolving Field

VAEs, 2013



GANs, 2014



PixelCNN, 2016



BigGAN, 2019



Imagen, 2022



SORA 2024



# Content

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

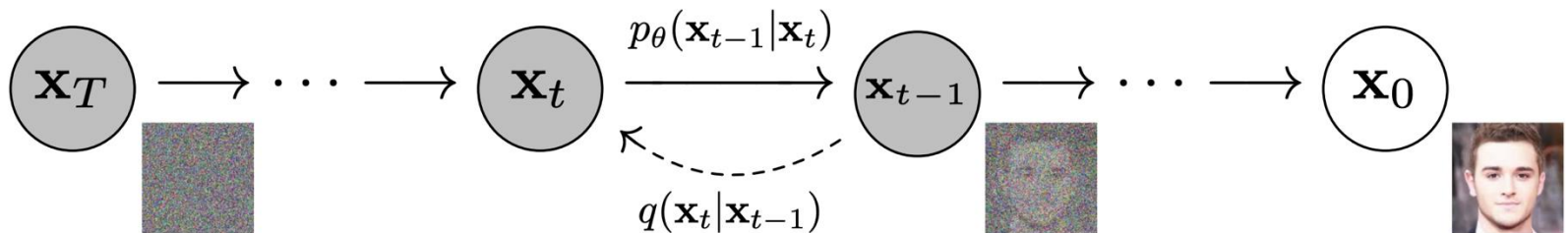


# Content

- Diffusion Model Basics
  - Diffusion Models as Stacking VAEs
  - Diffusion Models: Forward, Reverse, Training, Sampling
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- Denoising Diffusion Implicit Model (DDIM)
- Conditional Diffusion Models
- Applications of Diffusion Models

# Denoising Diffusion Models

- what we often see about diffusion models



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

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

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$$

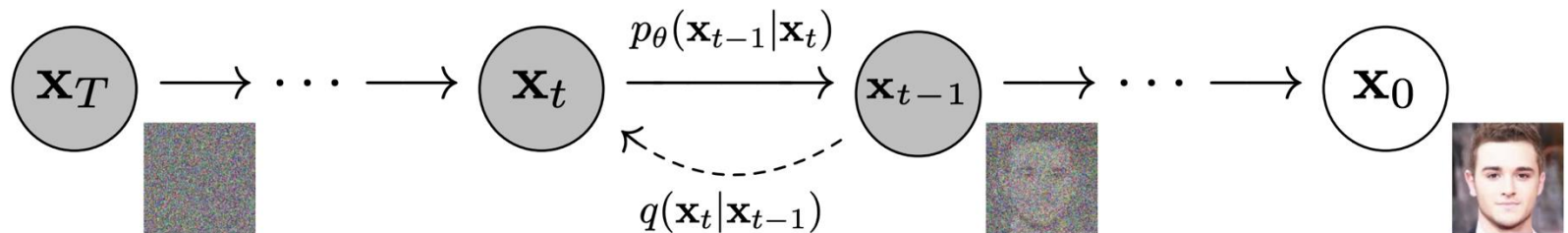
$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

Forward diffusion process

Reverse denoising process

# Denoising Diffusion Models

- what we often see about diffusion models



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

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

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$$

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

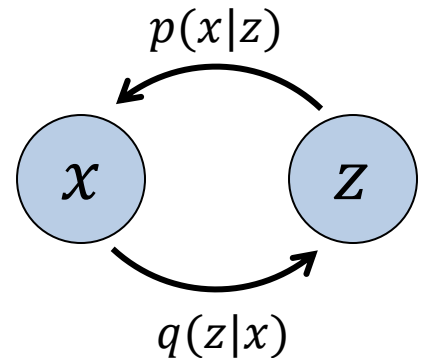
Forward diffusion process

Reverse denoising process

- this lecture: denoising diffusion is a stack of VAEs

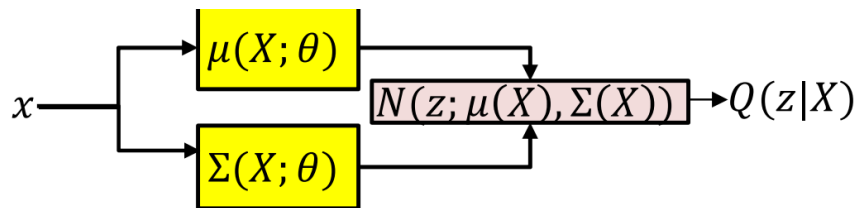
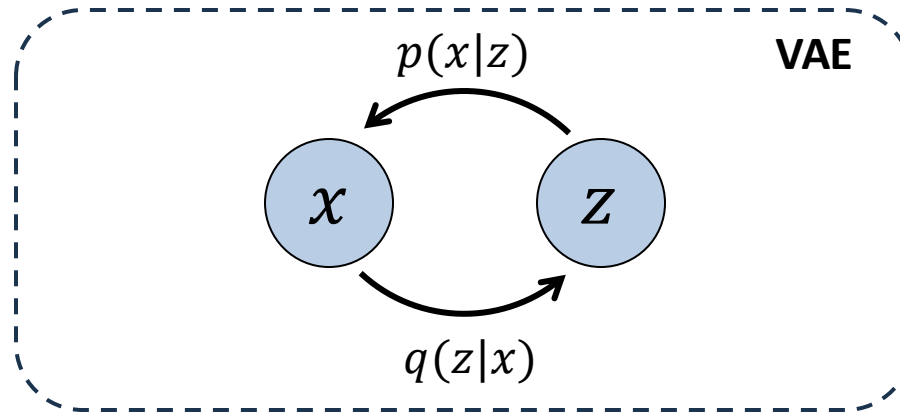
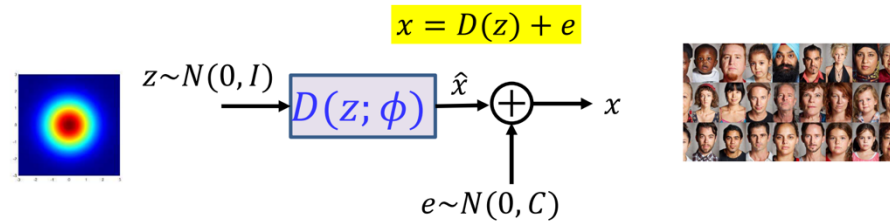
# Recap: Variational Autoencoders

- VAEs: a likelihood-based generative model
- **Encoder**: an inference model that approximates the posterior  $q(z|x)$
- **Decoder**: a generative model that transforms a Gaussian variable  $z$  to real data
- **Training**: maximize the ELBO



# Recap: Variational Autoencoders

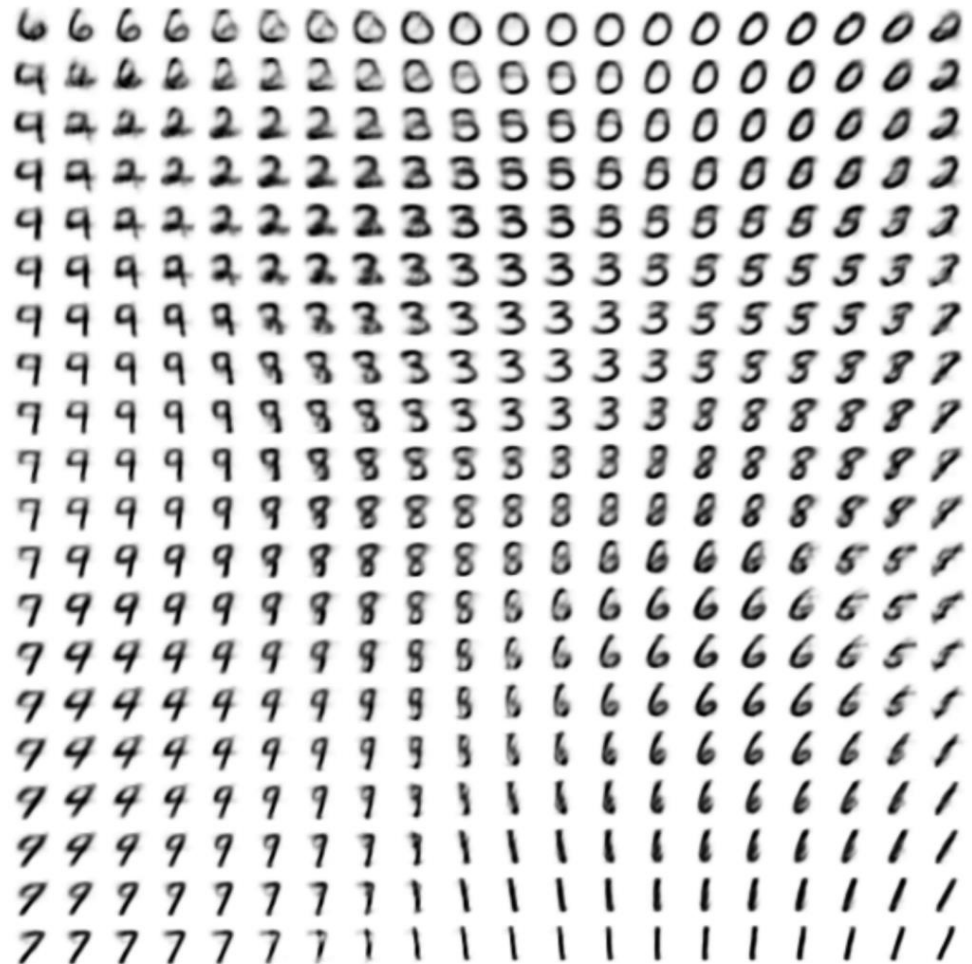
**Decoder:** transforms a Gaussian variable to real data



**Encoder:** an inference model approximates the posterior, i.e. Gaussian

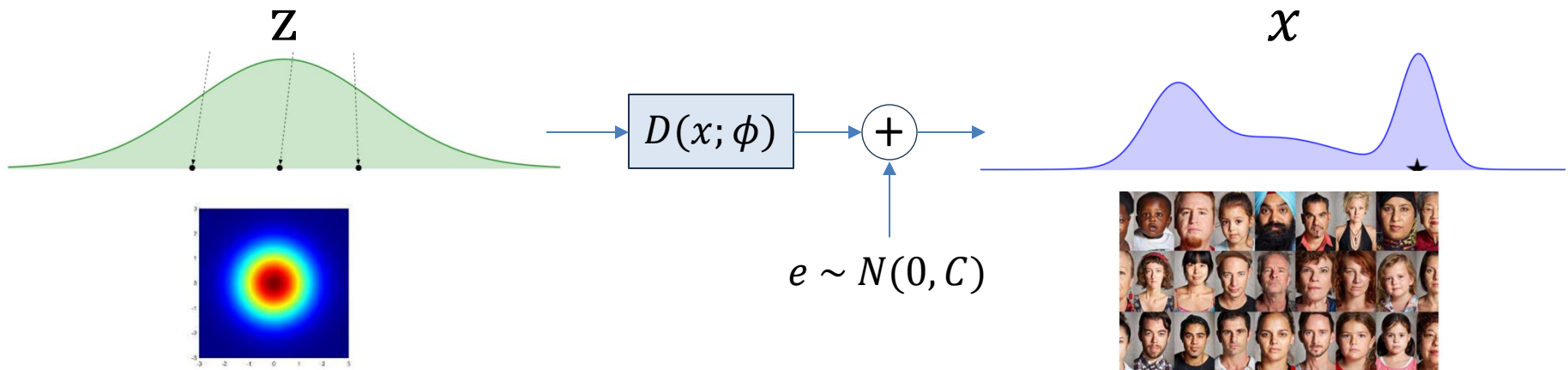
# VAEs are good, but...

- Blurry results



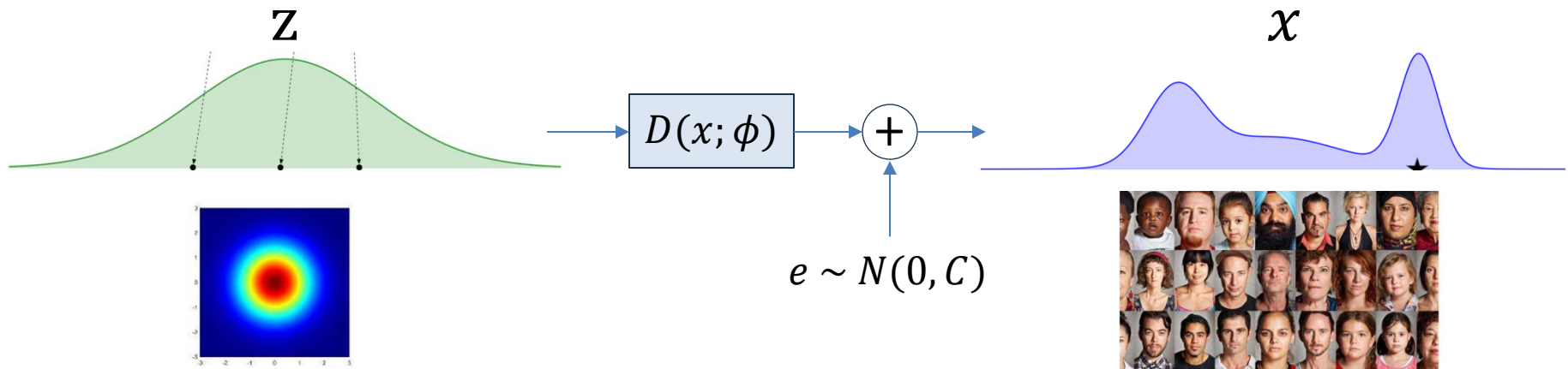
# Limitations of VAEs

- Decoder must transform a standard Gaussian all the way to the target distribution in **one-step**
  - Often too large a gap
  - Blurry results are generated



# Limitations of VAEs

- Decoder must transform a standard Gaussian all the way to the target distribution in **one-step**
  - Often too large a gap
  - Blurry results are generated



- Solution: have some intermediate latent variables to reduce the gap of each step



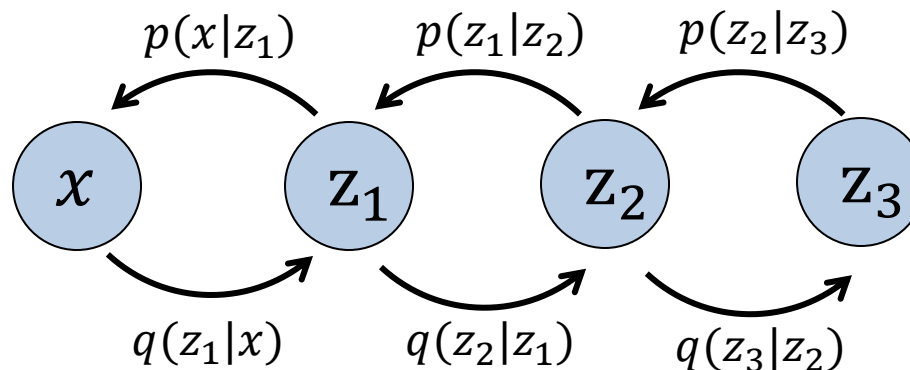
# Hierarchical VAEs

- Hierarchical VAEs – Stacking VAEs on top of each other
  - Multiple ( $T$ ) intermediate latent

- Joint distribution  $p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T) p_{\theta}(\mathbf{x} | \mathbf{z}_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t)$

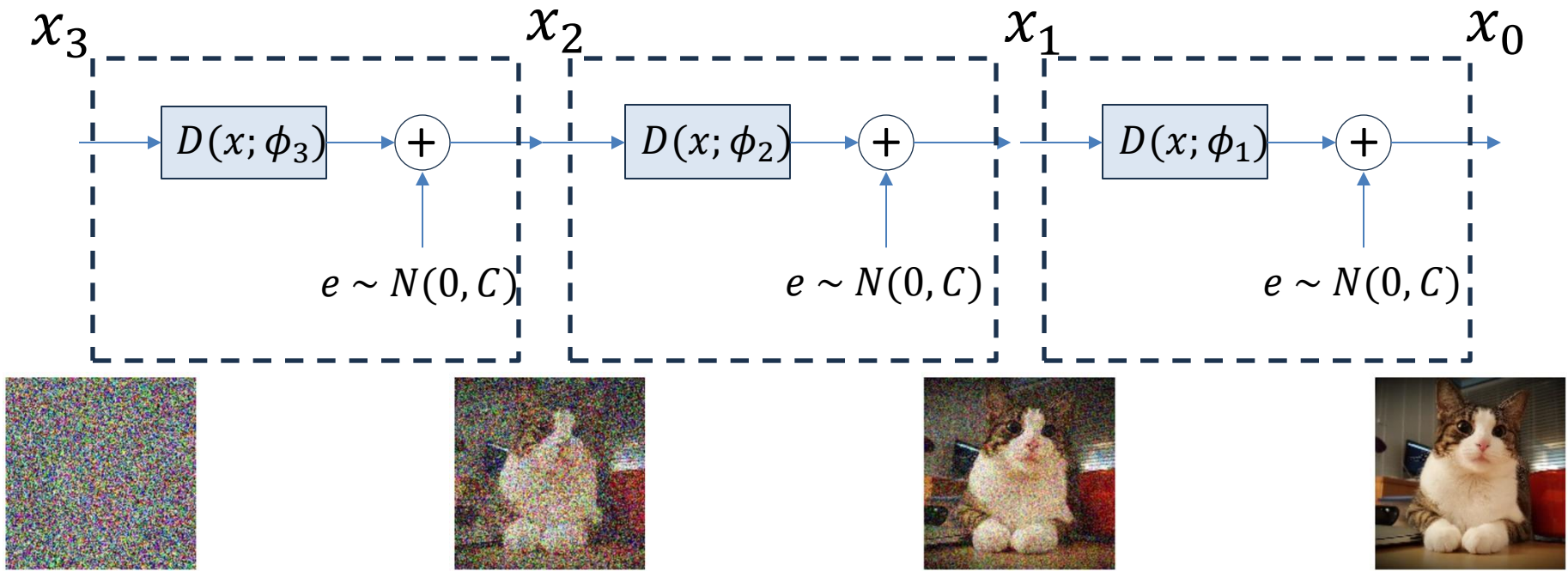
- Posterior  $q_{\phi}(\mathbf{z}_{1:T} | \mathbf{x}) = q_{\phi}(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t | \mathbf{z}_{t-1})$

- Better likelihood achieved!



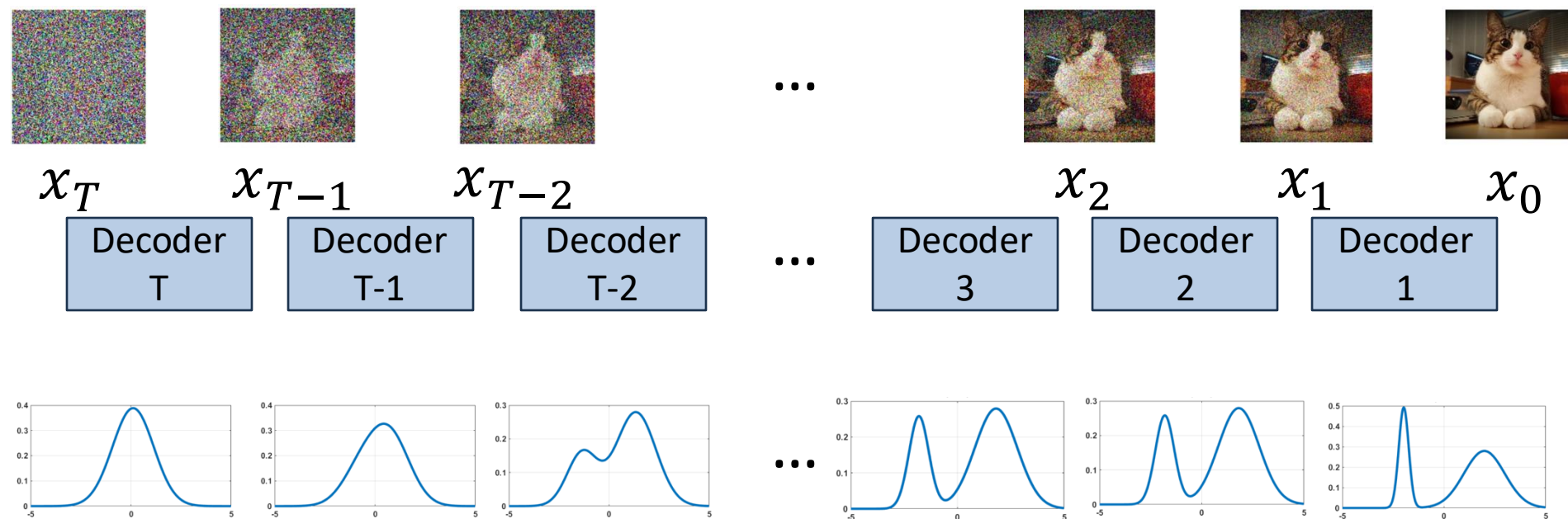
# Stacking VAEs

- Each step, the decoder removes part of the noise
- Provides a seed model closer to final distribution



# Stacking VAEs

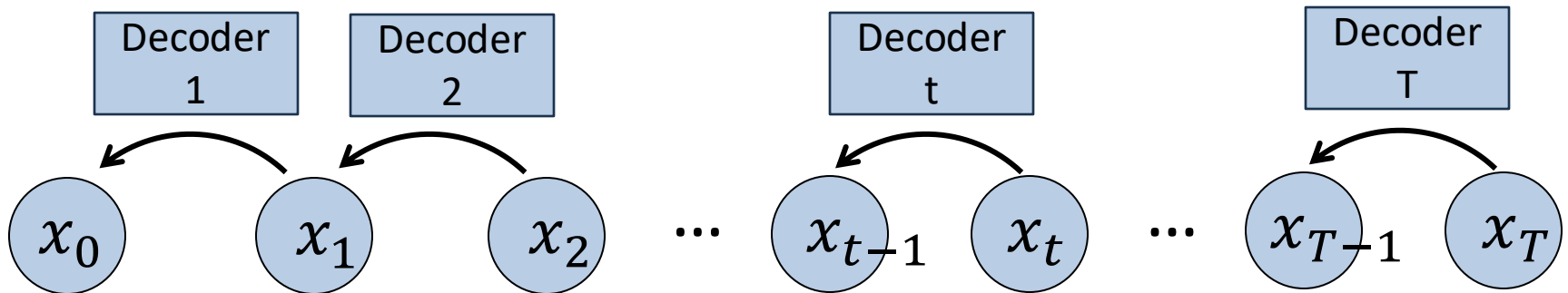
- We can have many many steps (in total  $T$ )...
- Each step incrementally recovers the final distribution



- Looks familiar?

# Diffusion Models are Stacking VAEs

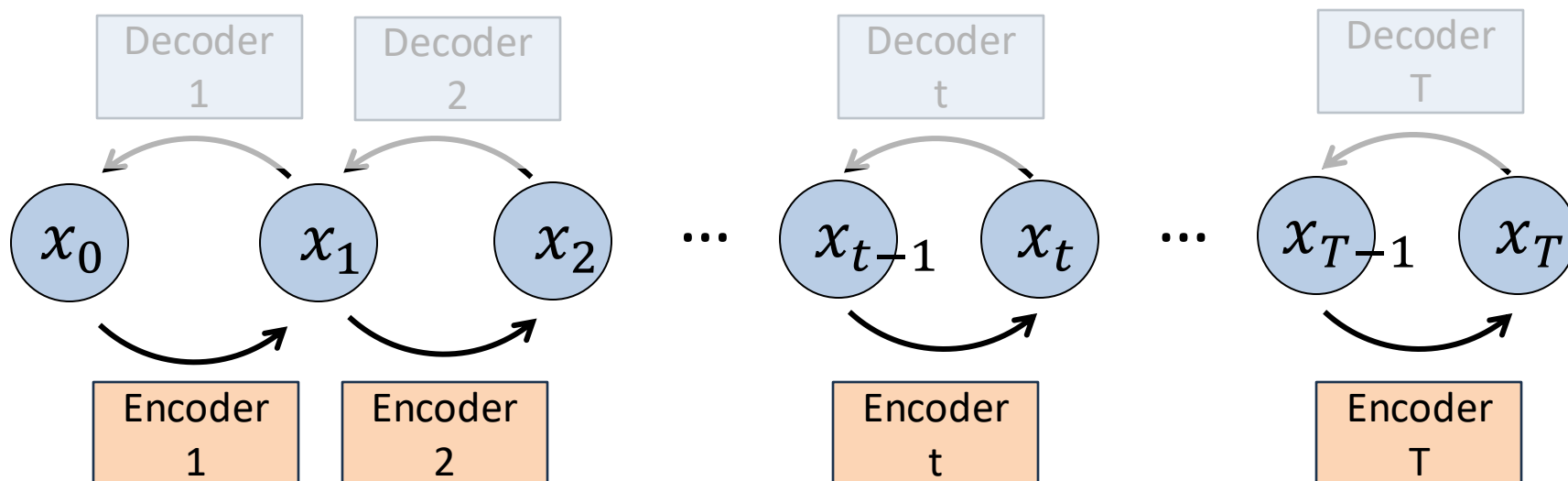
- Diffusion models are special cases of Stacking VAEs



- The reverse denoising process is the stack of decoders
- What about encoders?

# Diffusion Models are Stacking VAEs

- Diffusion models are special case of Stacking VAEs

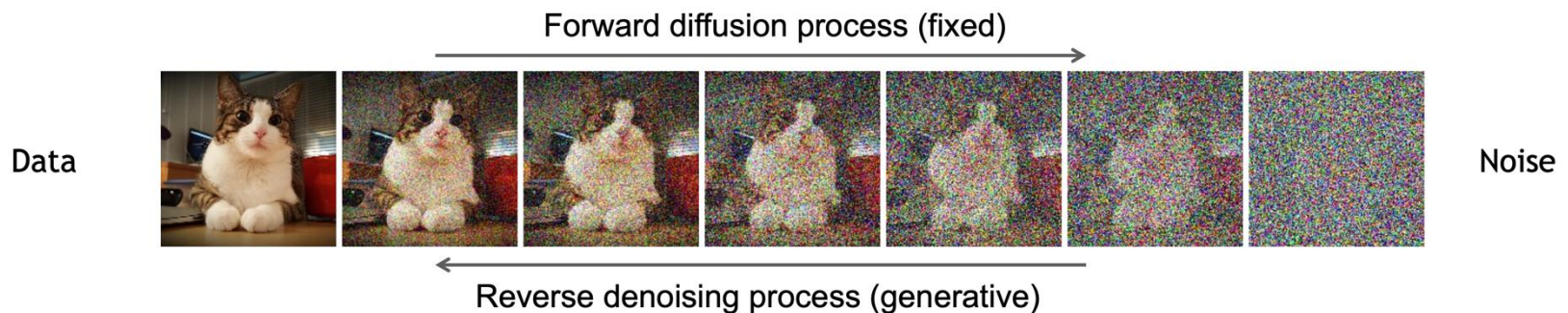


- In VAEs, encoders are learned with KL-divergence between the posterior and the prior
- Suffers from the ‘posterior-collapse’ issue
- Diffusion models use **fixed inference encoders**

# Poll

# Denoising Diffusion Models

- Diffusion models have two processes
- **Forward diffusion process** gradually adds noise to input
- **Reverse denoising process** learns to generate data by denoising

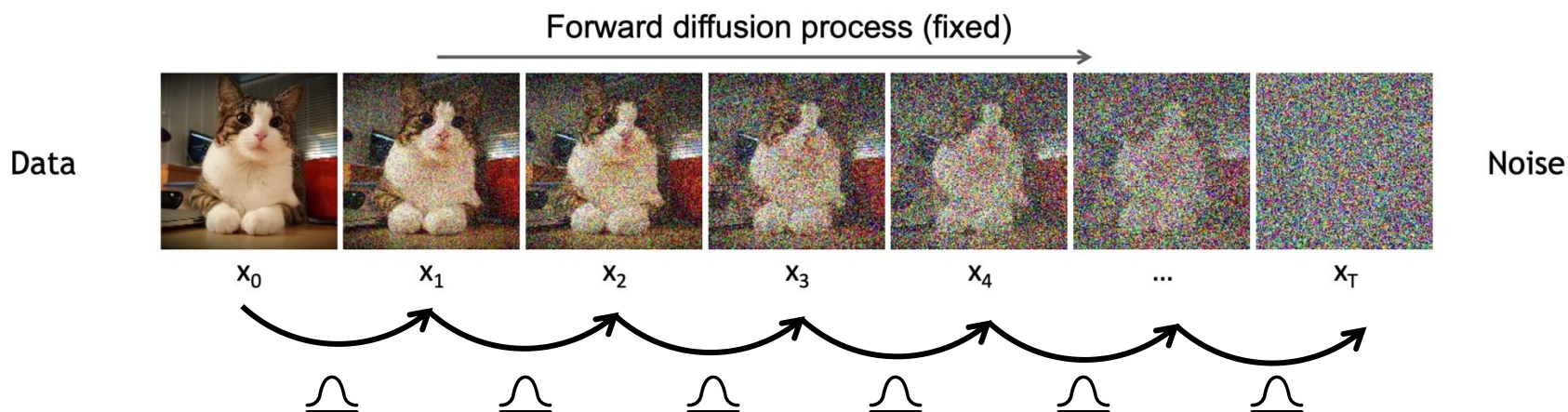


# Forward Diffusion Process

- Forward diffusion process is stacking **fixed** VAE encoders
  - gradually adding Gaussian noise according to schedule  $\beta_t$

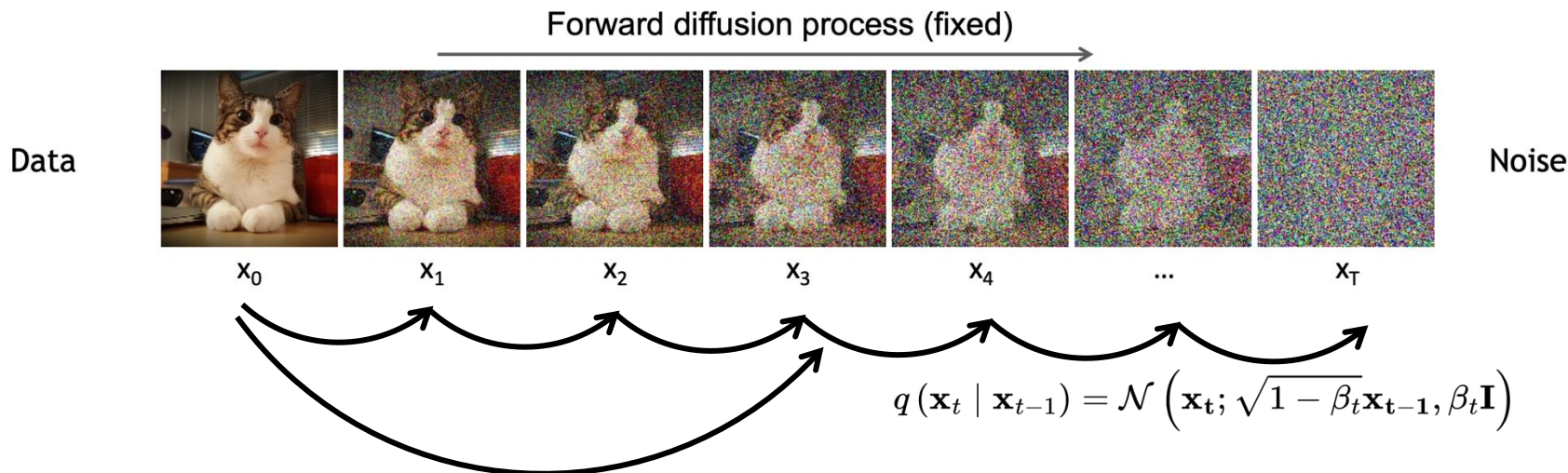
$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$





# Forward Diffusion Process



- The forward process allows sampling of  $\mathbf{x}_t$  at arbitrary timestep  $t$  in closed form:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad \bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$$

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- The noise schedule ( $\beta_t$  values) is designed such that

$$q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$

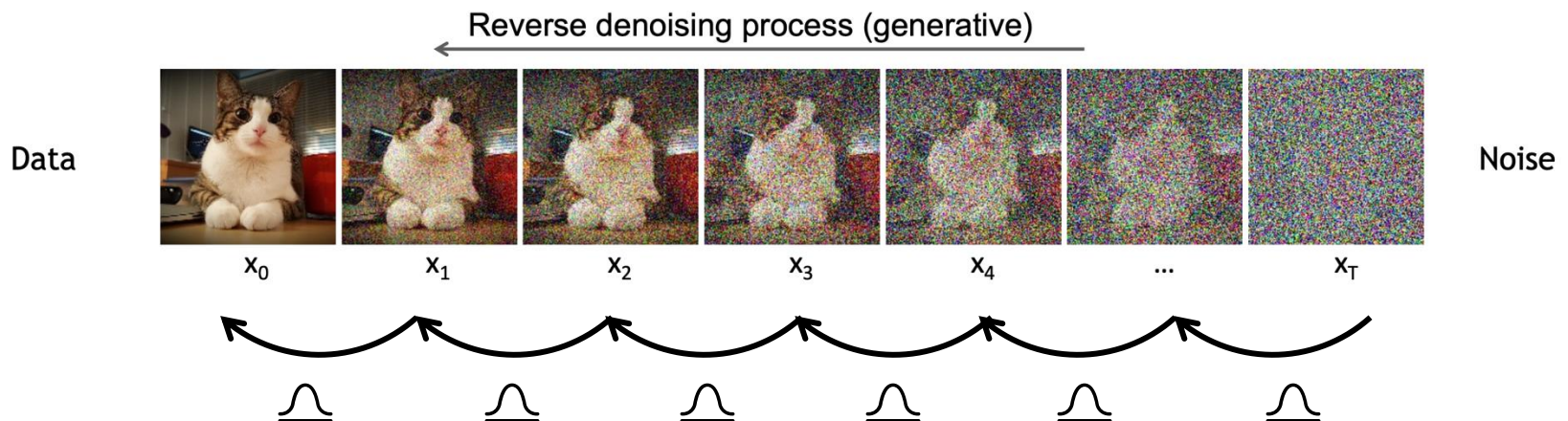
# Reverse Denoising Process

- Generation process
  - Sample  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$
  - Iteratively sample  $\mathbf{x}_{t-1} \sim q(\mathbf{x}_{t-1} | \mathbf{x}_t)$
- $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  not directly tractable
- But can be estimated with a Gaussian distribution if  $\beta_t$  is small at each step
  - The purpose of our stack of VAE decoders!

# Reverse Denoising Process

- Reverse diffusion process is stacking **learnable** VAE decoders
  - Predicting the mean and std of added Gaussian Noise

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \quad p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$
$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

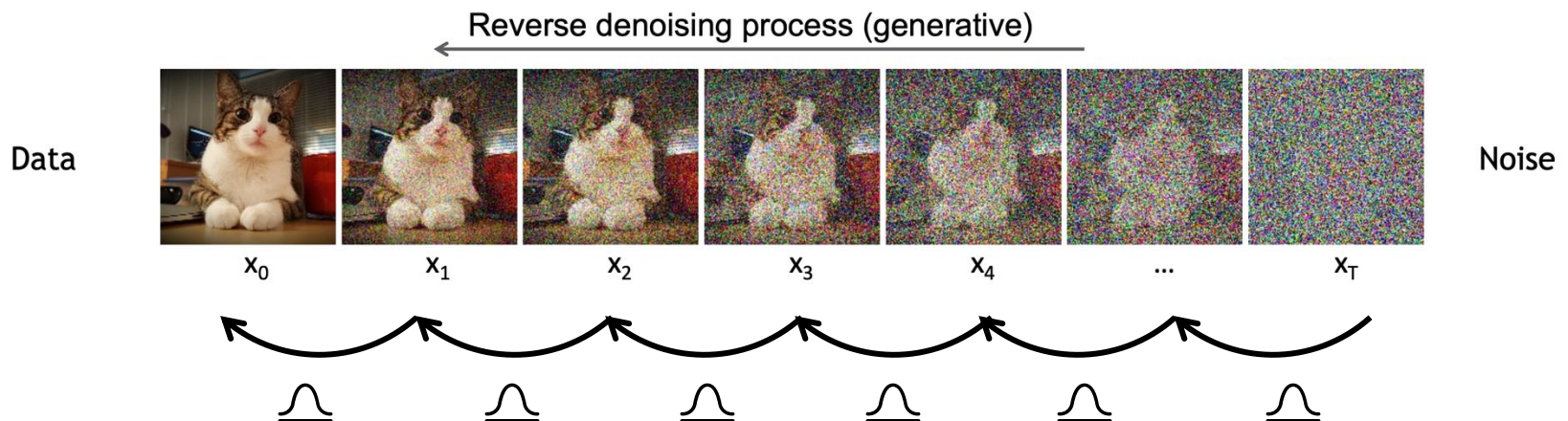


# Reverse Denoising Process

- Reverse diffusion process is stacking **learnable** VAE decoders
  - Predicting the mean and std of added Gaussian Noise

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

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$



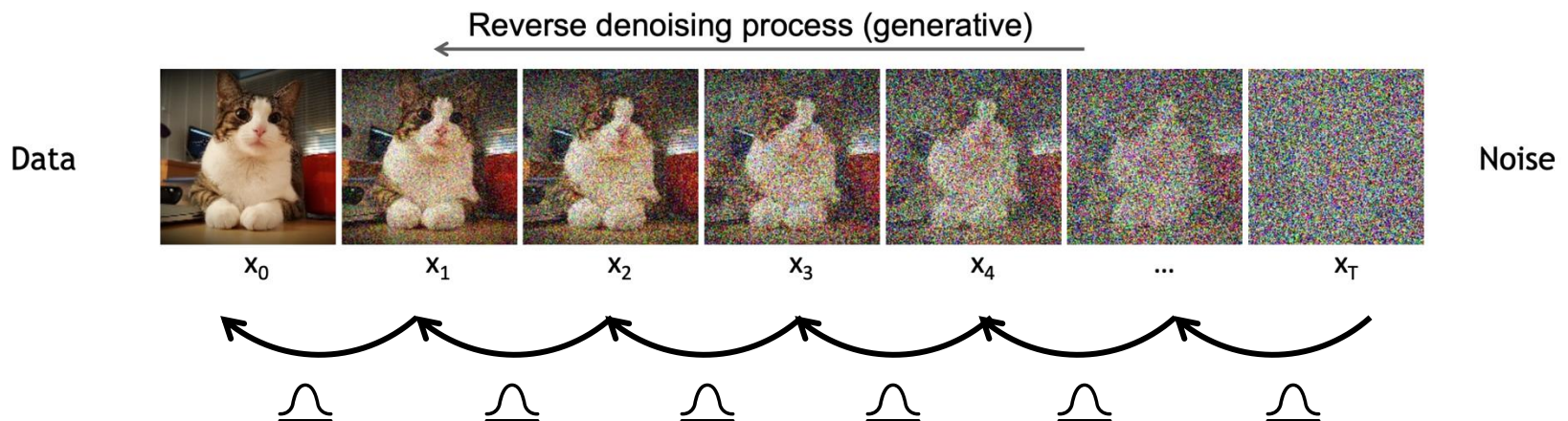
# Reverse Denoising Process

- Reverse diffusion process is stacking **learnable** VAE decoders
  - Predicting the mean and std of added Gaussian Noise

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

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}}_{\text{Trainable Network, Shared Across All Timesteps}})$$

Trainable Network, Shared Across All Timesteps



# Learning the Denoising Model

- Denoising models are trained with variational upper bound (negative ELBO), as VAEs

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] =: L$$

- which derives to:

$$L = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]$$

- tractable posterior distribution (closed-form)

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

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



# Learning the Denoising Model

- Denoising models are trained with variational upper bound (negative ELBO), as VAEs

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] =: L$$

- which derives to:

$$L = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T))}_{\substack{L_T \\ \text{constant}}} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} \underbrace{- \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{\substack{L_0 \\ \text{Scaling}}} \right]$$

- tractable posterior distribution (closed-form)

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

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

# Learning the Denoising Model

- Denoising models are trained with variational upper bound (negative ELBO), as VAEs

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] =: L$$

- which derives to:

$$L = \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p(\mathbf{x}_T))}_{L_T} + \sum_{t>1} \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]$$

- tractable posterior distribution (closed-form)

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

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



# Parameterizing the Denoising Model

- KL divergence has a simple form between Gaussians

$$L_{t-1} = D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) = \mathbb{E}_q \left[ \frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_{\theta}(\mathbf{x}_t, t)\|^2 \right] + C$$

- Recall that:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$

$$\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{1 - \beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

- Trainable network predicts the noise mean

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

- Final Objective

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \frac{\beta_t^2}{2\sigma_t^2 (1 - \beta_t) (1 - \bar{\alpha}_t)} \left\| \epsilon - \epsilon_{\theta}(\underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{\mathbf{x}_t}, t) \right\|^2 \right] + C$$

# Simplified Training Objective

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

- $\lambda_t$  ensures the weighting for correct maximum likelihood estimation
- In DDPM, this is further simplified to:

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

# Summary: Training and Sampling

---

## Algorithm 1 Training

---

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

---

---

## Algorithm 2 Sampling

---

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

---

# Summary: Noise Schedule

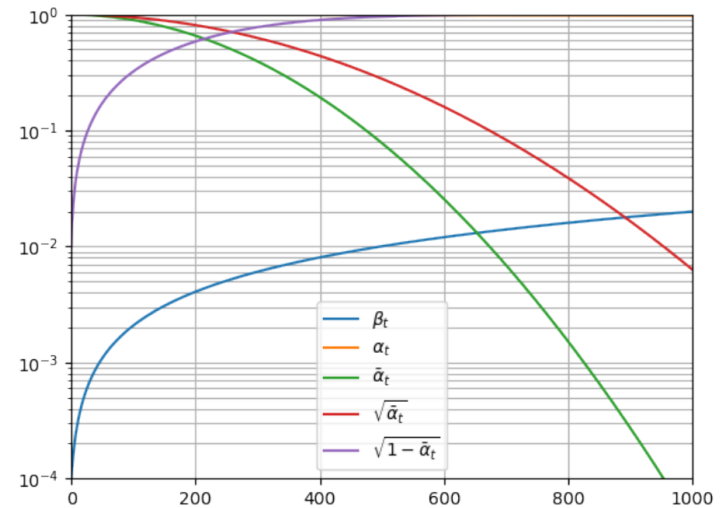
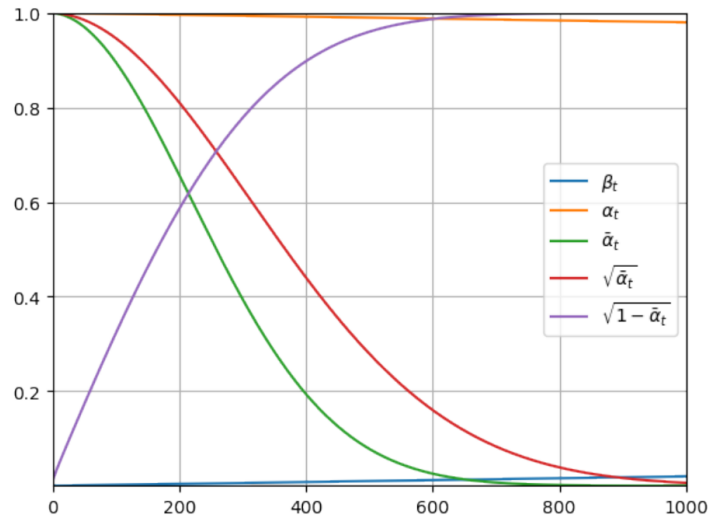
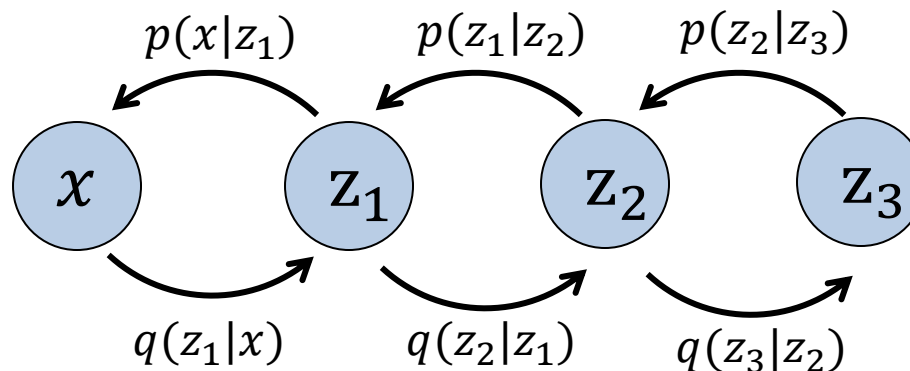


Figure 2: Parameter values for  $\beta = [10^{-4}, 0.02]$  over 1000 time steps  $t$  using a linear schedule. The information in the two figures are the same, but the right-hand side uses log-scale on the  $y$ -axis to show the speed of which  $\tilde{\alpha}_t$  goes towards zero.

# Connection with Hierarchical VAEs

- Diffusion models are special case of Hierarchical VAEs
  - Fixed inference models in forward process
  - Latent variables have same dimension as data
  - ELBO is decomposed to each timestep: faster to train
  - Model is trained with some weighting of ELBO



# Poll

# Content

- Diffusion Model Basics
  - Diffusion Models as Stacking VAEs
  - Diffusion Models: Forward, Reverse, Training, Sampling
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- Classifier-Free Guidance for Conditional Models
- Applications of Diffusion Models

# Why SDEs?

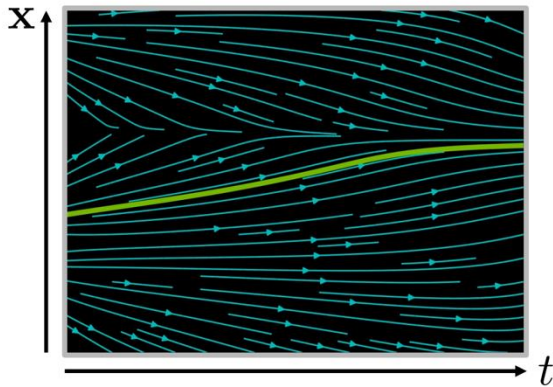
- A unified framework for interpreting diffusion models and score-based generation models
  - Variants of diffusion-based and flow-based models



# Stochastic Differential Equations

Ordinary Differential Equation (ODE):

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



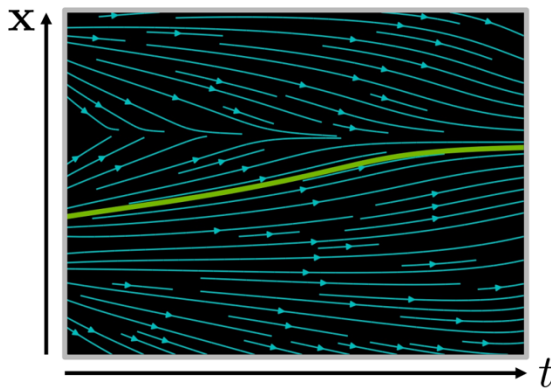
Analytical Solution: 
$$x(t) = x(0) + \int_0^t f(x, \tau) d\tau$$

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

# Stochastic Differential Equations

Ordinary Differential Equation (ODE):

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



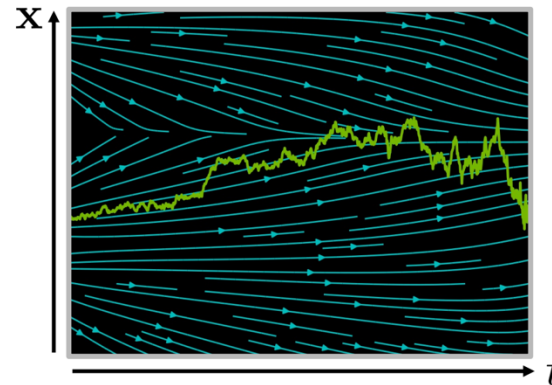
Analytical Solution:  $x(t) = x(0) + \int_0^t f(x, \tau) d\tau$

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

Stochastic Differential Equation (SDE):

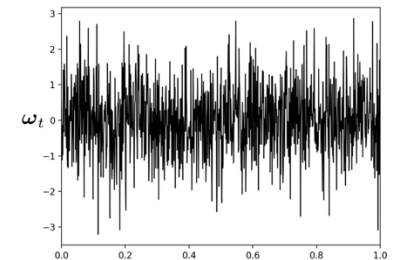
$$\frac{dx}{dt} = \underbrace{f(x, t)}_{\text{drift coefficient}} + \underbrace{\sigma(x, t)\omega_t}_{\text{diffusion coefficient}}$$

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



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

Wiener Process  
(Gaussian  
White Noise)



# Score Matching

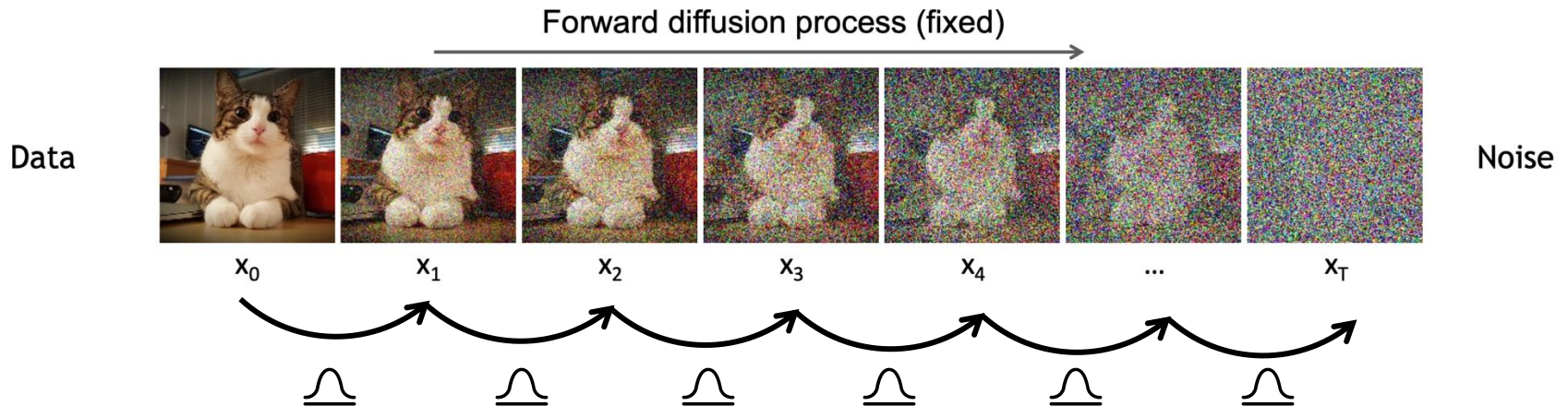
- General form of probability density function

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

- Maximizing the log-likelihood requires us to know  $Z_{\theta}$ 
  - Often intractable
- Instead, we can model the score function

$$\nabla_{\mathbf{x}} \log p(\mathbf{x})$$

# Forward Diffusion Process as SDEs

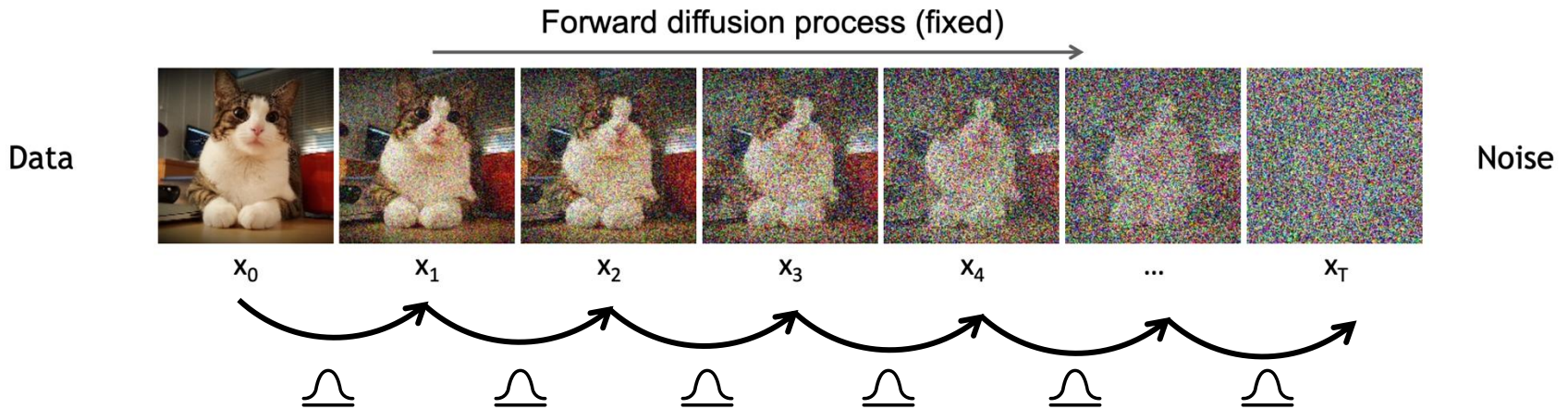


- Consider a forward process with many many small steps (continuous time)

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

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

# Forward Diffusion Process as SDEs



- Consider a forward process with many many small steps

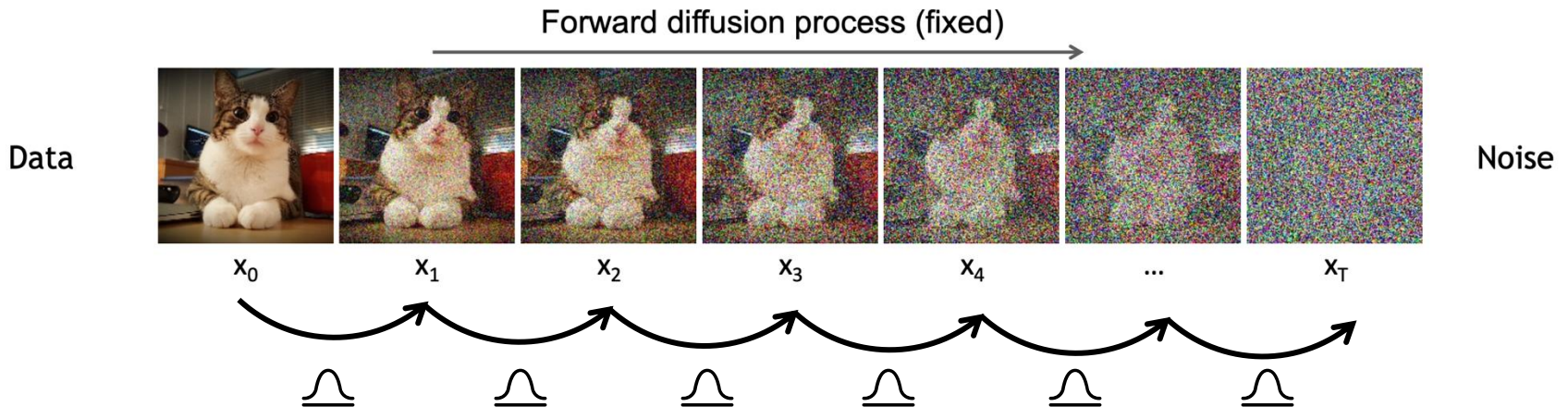
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

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

Allows different size along  $t$

Step size

# Forward Diffusion Process as SDEs



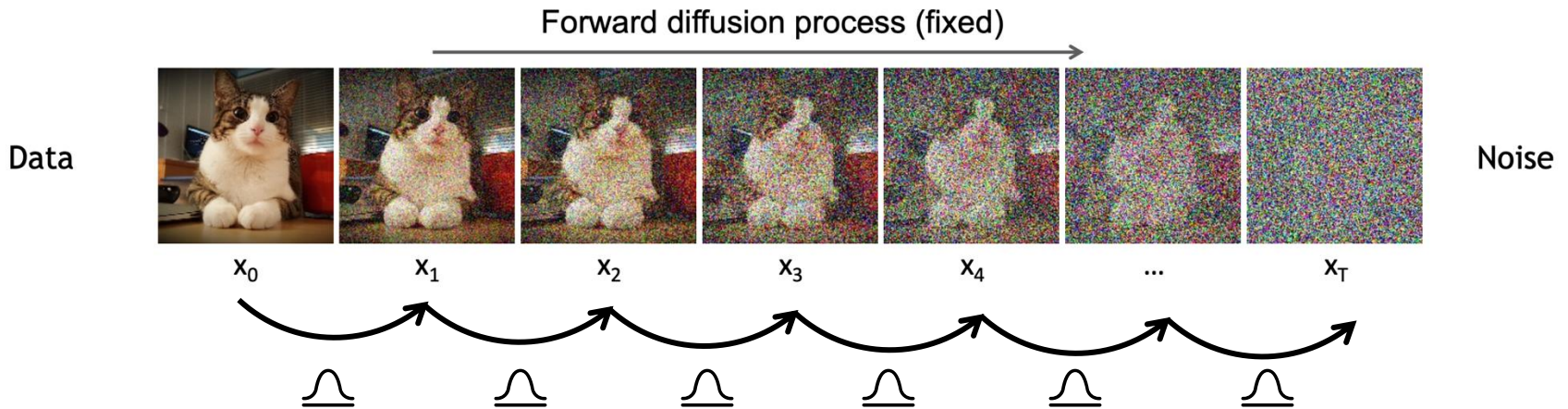
- Consider a forward process with many many small steps

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

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



# Forward Diffusion Process as SDEs



- An iterative update that can be viewed as SDEs

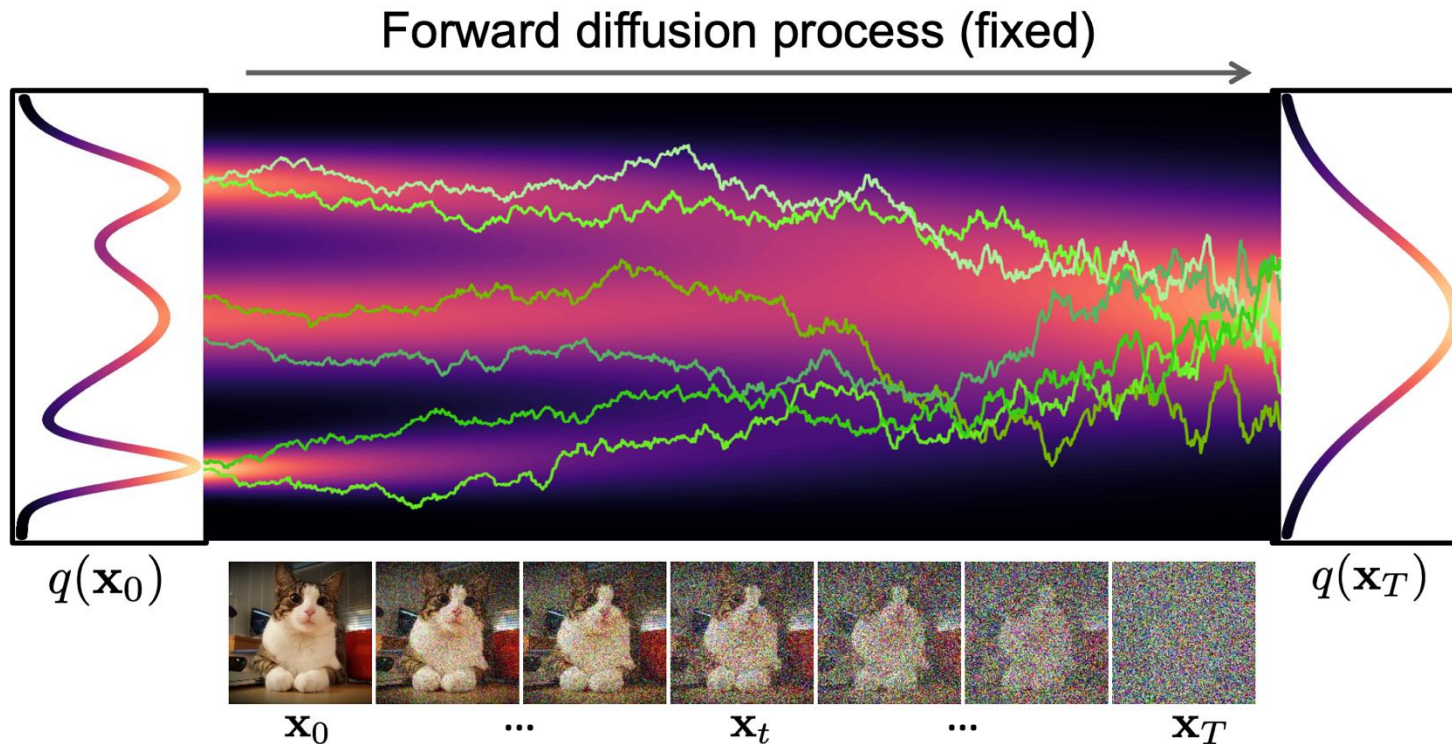
$$\mathbf{x}_t \approx \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2}\mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I})$$



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

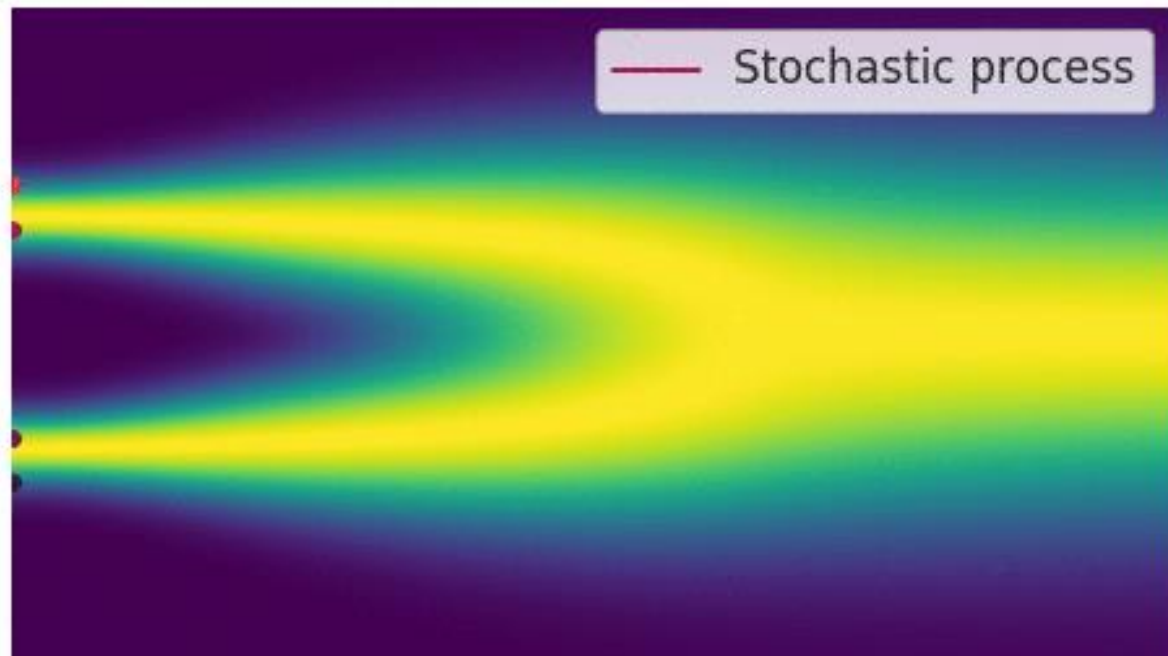
Stochastic Differential Equation (SDE)

# Forward Diffusion Process as SDEs

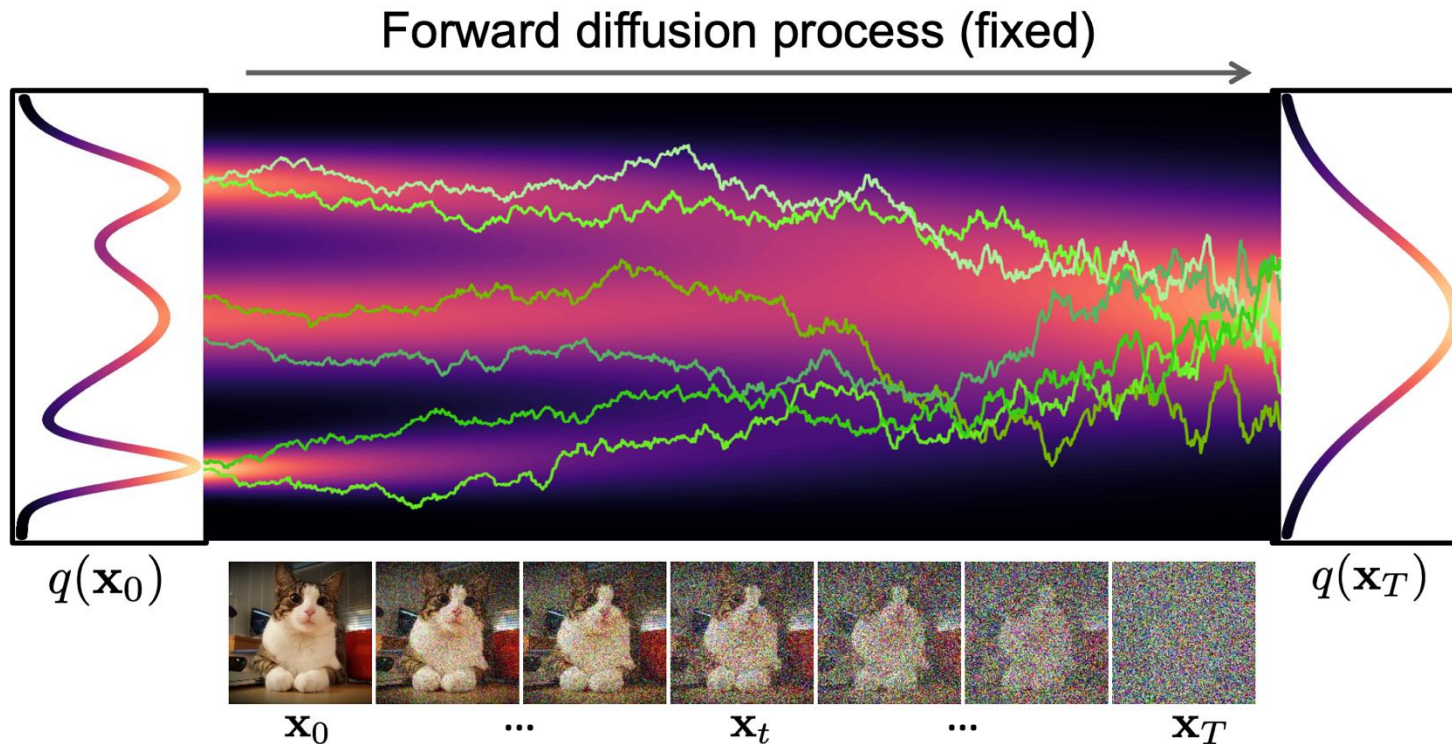


$$d\mathbf{x}_t = \underbrace{-\frac{1}{2}\beta(t)\mathbf{x}_t dt}_{\text{Drift Term (Pulls toward the mode)}} + \underbrace{\sqrt{\beta(t)}d\boldsymbol{\omega}_t}_{\text{Diffusion Term (Injects Noise)}}$$



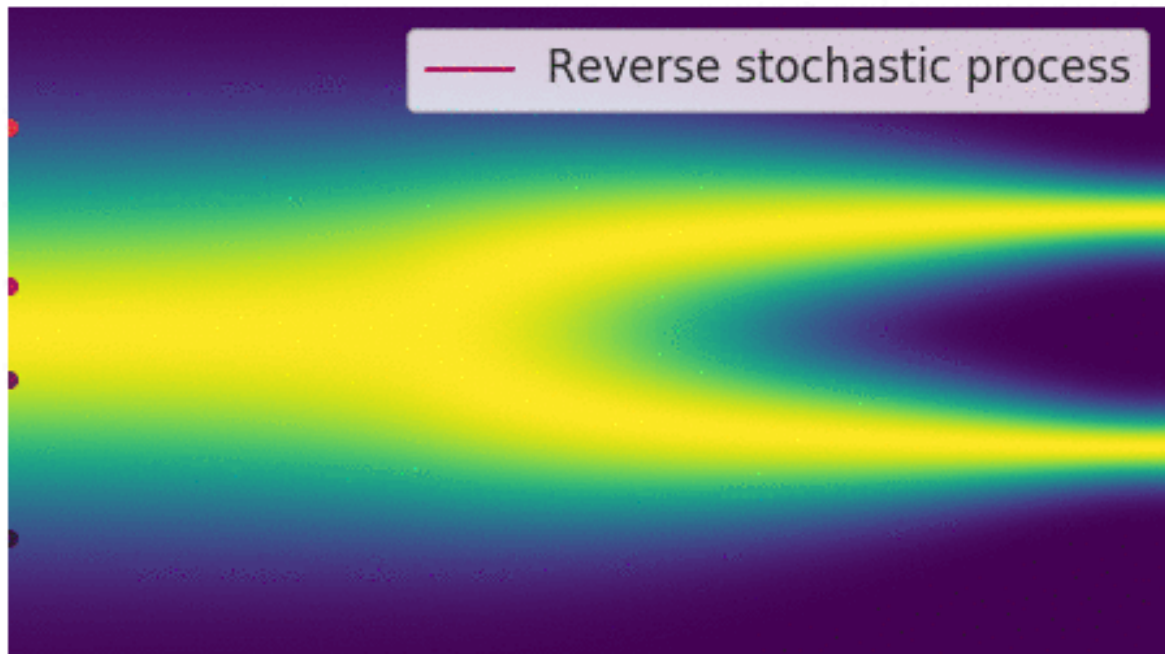


# Generative Reverse SDEs



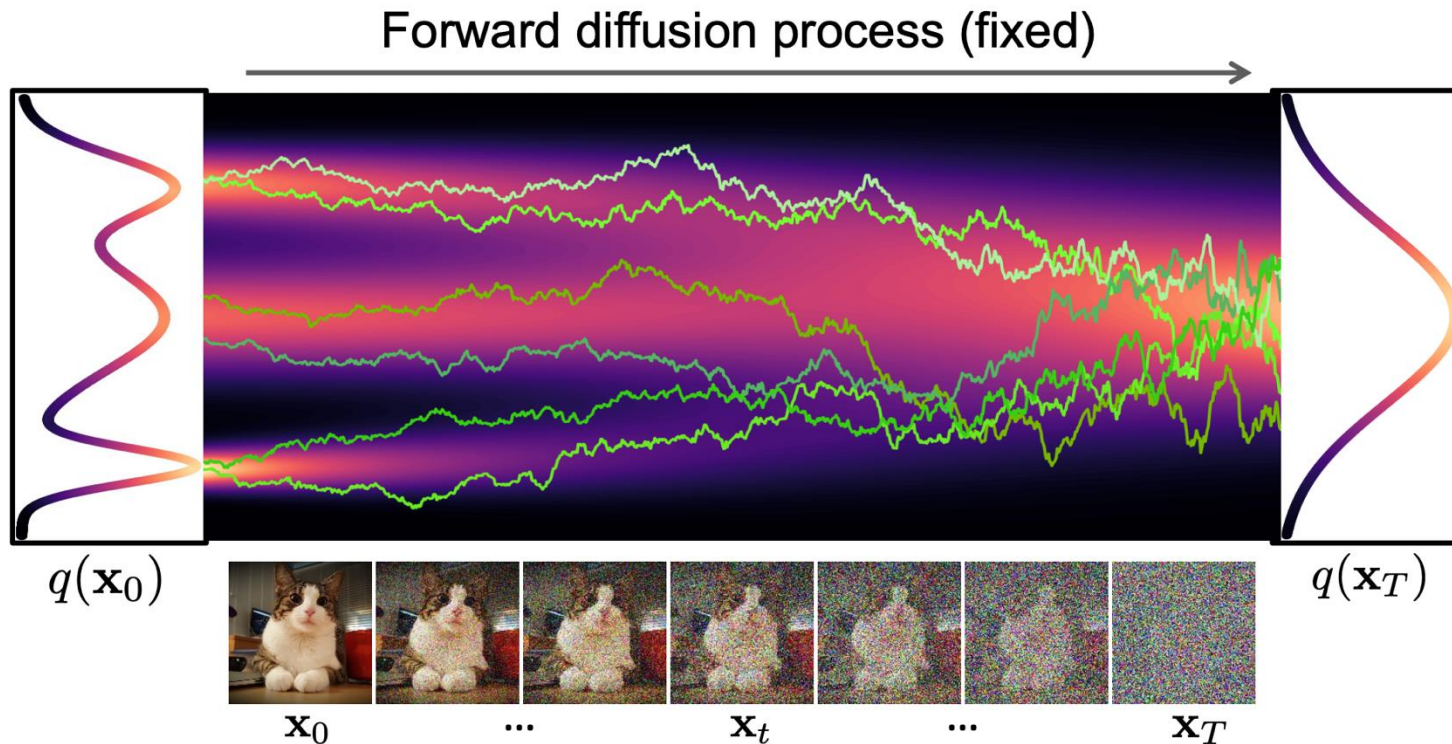
- The forward SDE has a reverse form:

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





# Generative Reverse SDEs



- The forward SDE has a reverse form:

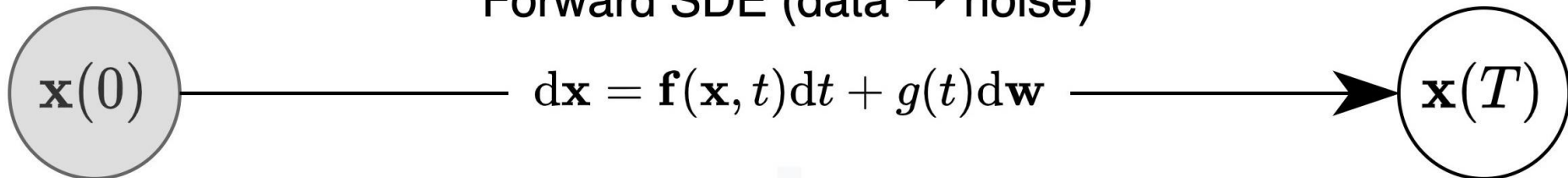
$$d\mathbf{x}_t = \left[ -\frac{1}{2}\beta(t)\mathbf{x}_t - \beta(t) \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \right] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t$$

Score function

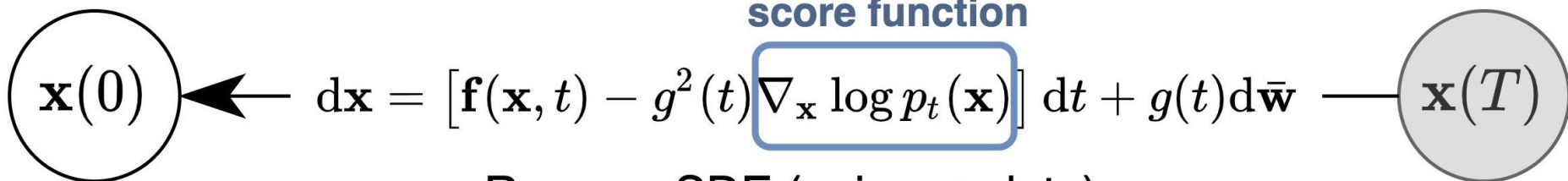
How to get it?

# Denoising Score Matching

Forward SDE (data  $\rightarrow$  noise)



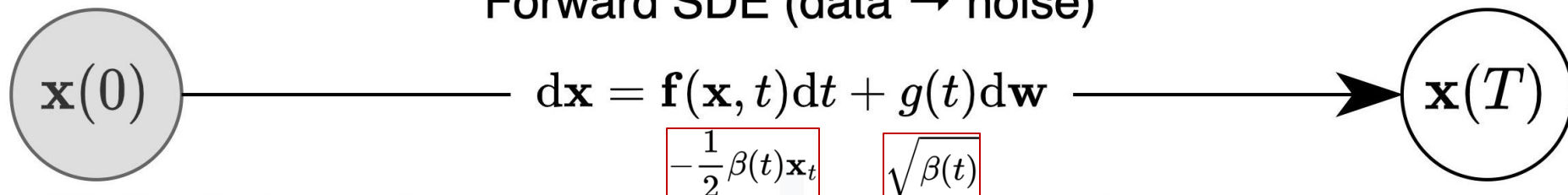
score function



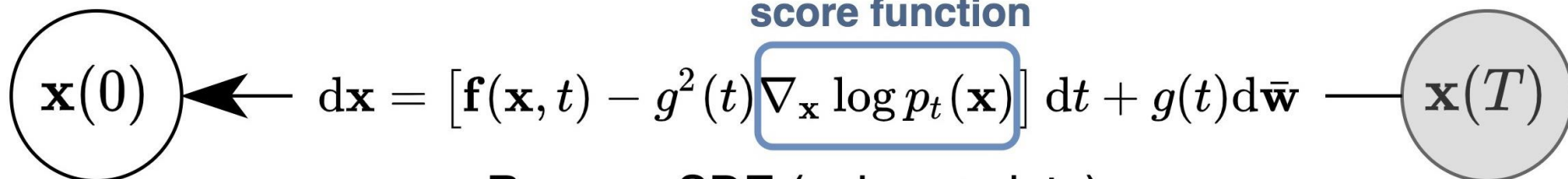
Reverse SDE (noise  $\rightarrow$  data)

# Denoising Score Matching

Forward SDE (data  $\rightarrow$  noise)



score function



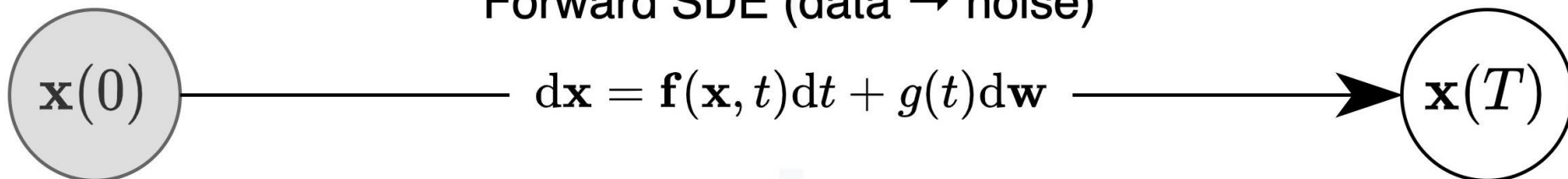
Reverse SDE (noise  $\rightarrow$  data)

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

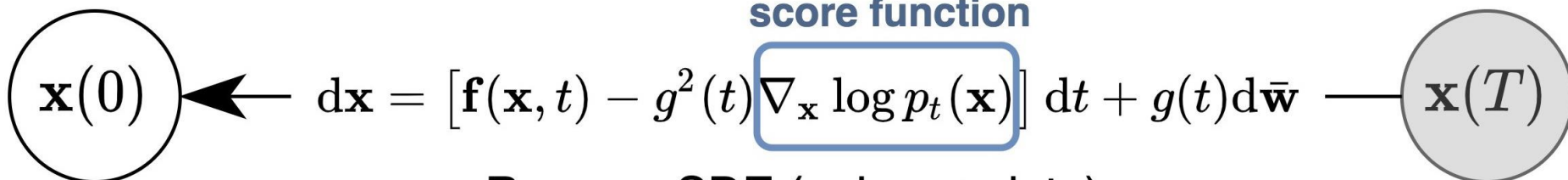


# Denoising Score Matching

Forward SDE (data  $\rightarrow$  noise)



score function



Reverse SDE (noise  $\rightarrow$  data)

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

Looks similar?

# Denoising Score Matching

- Denoising score matching objective

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

- 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:

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



# Weighted Diffusion Objective

- Denoising score matching objective with loss weighting

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

- 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!

# Poll

# Content

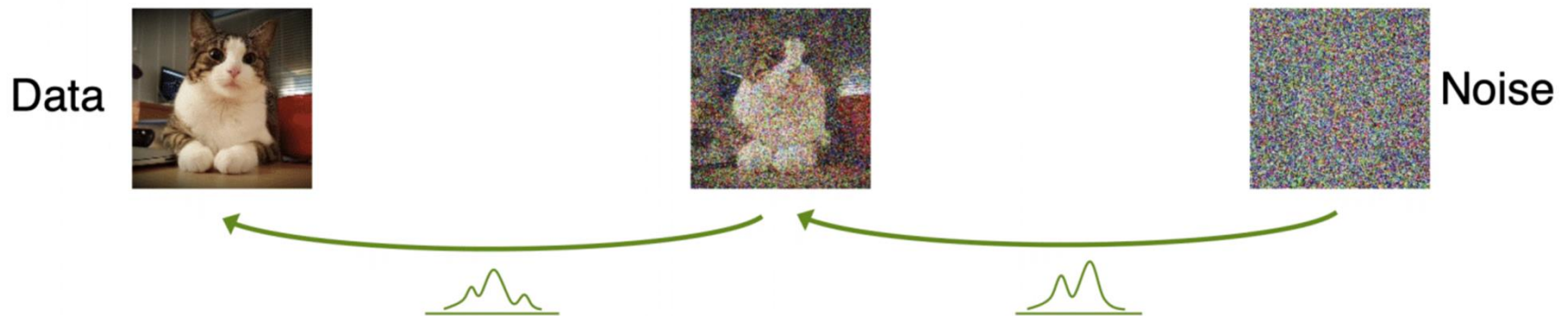
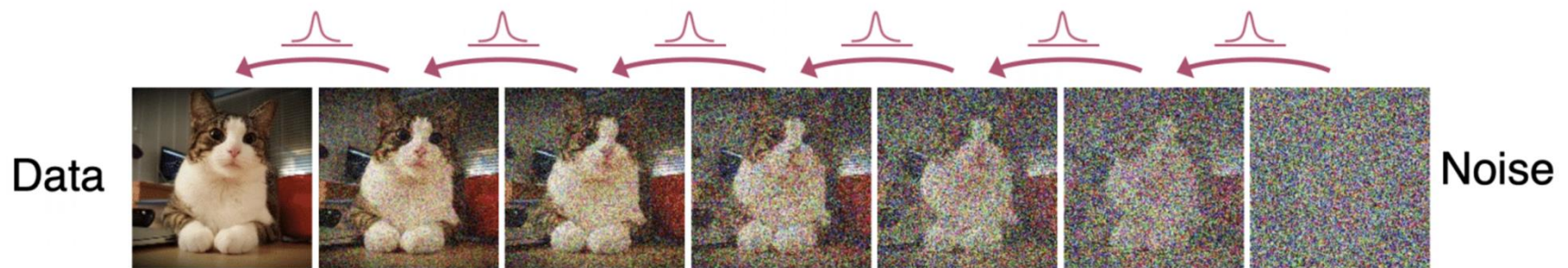
- Diffusion Model Basics
- Diffusion Models from Stochastic Differential Equations and Score Matching Perspective
- **Denoising Diffusion Implicit Model (DDIM)**
- Conditional Diffusion Models
- Applications of Diffusion Models

# Many Steps in Diffusion

- Slow in generation
- In Training, we randomly sample one time step
- But in inference, we must transit from  $T$  to  $0$ 
  - 1000 steps
  - extremely slow for raw images/signals

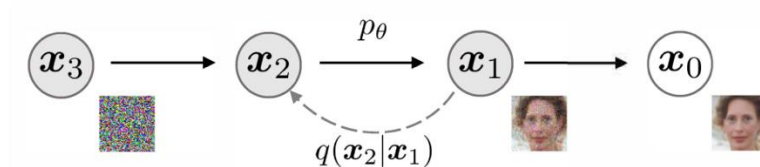
# Can we do generation with less steps?

## Denoising Process with Uni-modal Normal Distribution



Requires more complicated functional approximators!

# DDPM

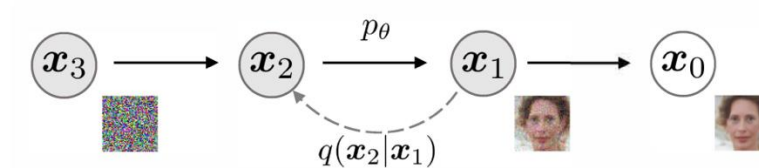


$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

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

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

# DDPM



Only depends on previous step

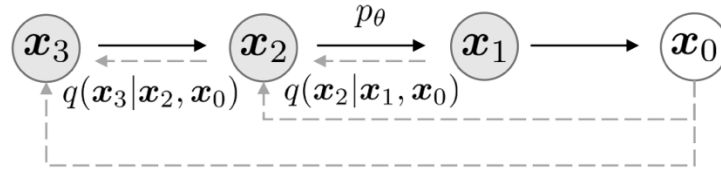
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

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

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

Only used during training

# DDIM



$$q_\sigma(\mathbf{x}_{1:T} | \mathbf{x}_0) := q_\sigma(\mathbf{x}_T | \mathbf{x}_0) \prod_{t=2}^T q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

$$q_\sigma(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_T} \mathbf{x}_0, (1 - \alpha_T) \mathbf{I})$$

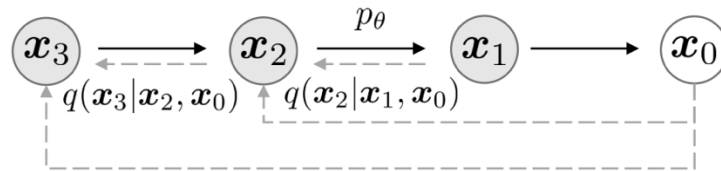
$$q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right)$$

- A Non-Markovian Forward Process

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



# DDIM



- Backward process

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

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

# DDPM vs DDIM

---

## Algorithm DDPM Sampling

---

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

**for all**  $t$  from  $T$  to 1 **do**

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

$$\mu \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right)$$

$$\mathbf{x}_{t-1} \leftarrow \mu + \boxed{\sigma_t \epsilon} \text{ Stochastic}$$

**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**

$$\boxed{\bar{\epsilon}} \leftarrow \epsilon_{\theta}(\boxed{\mathbf{x}_t}, t)$$

$$\bar{\mathbf{x}}_0 \leftarrow \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \bar{\epsilon}}{\sqrt{\bar{\alpha}_t}} \quad \text{Estimate } \mathbf{x}_0$$

$$\boxed{\mathbf{x}_{t-1}} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \boxed{\bar{\mathbf{x}}_0} + \sqrt{1 - \bar{\alpha}_{t-1}} \boxed{\bar{\epsilon}}$$

**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

$[1, \dots, T] \implies [\tau_0 = 0, \dots, \tau_S = T]$

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 ( $32 \times 32$ ) |             |             |             |             | CelebA ( $64 \times 64$ ) |              |             |             |             |
|----------------|----------------------------|-------------|-------------|-------------|-------------|---------------------------|--------------|-------------|-------------|-------------|
|                | 10                         | 20          | 50          | 100         | 1000        | 10                        | 20           | 50          | 100         | 1000        |
| $\eta = 0.0$   | <b>13.36</b>               | <b>6.84</b> | <b>4.67</b> | <b>4.16</b> | 4.04        | <b>17.33</b>              | <b>13.73</b> | <b>9.17</b> | <b>6.53</b> | 3.51        |
| $\eta = 0.2$   | 14.04                      | 7.11        | 4.77        | 4.25        | 4.09        | 17.66                     | 14.11        | 9.51        | 6.79        | 3.64        |
| $\eta = 0.5$   | 16.66                      | 8.35        | 5.25        | 4.46        | 4.29        | 19.86                     | 16.06        | 11.01       | 8.09        | 4.28        |
| $\eta = 1.0$   | 41.07                      | 18.36       | 8.01        | 5.78        | 4.73        | 33.12                     | 26.03        | 18.48       | 13.93       | 5.98        |
| $\hat{\sigma}$ | 367.43                     | 133.37      | 32.72       | 9.99        | <b>3.17</b> | 299.71                    | 183.83       | 71.71       | 45.20       | <b>3.26</b> |

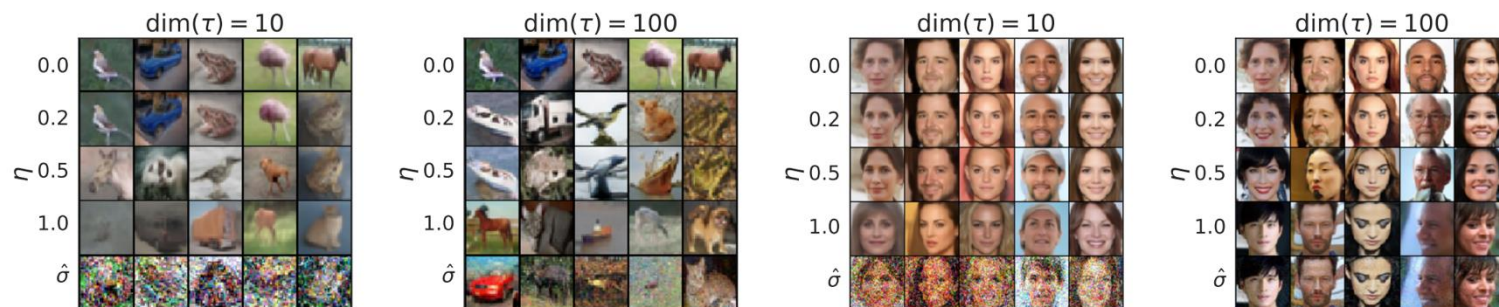


Figure 3: CIFAR10 and CelebA samples with  $\dim(\tau) = 10$  and  $\dim(\tau) = 100$ .

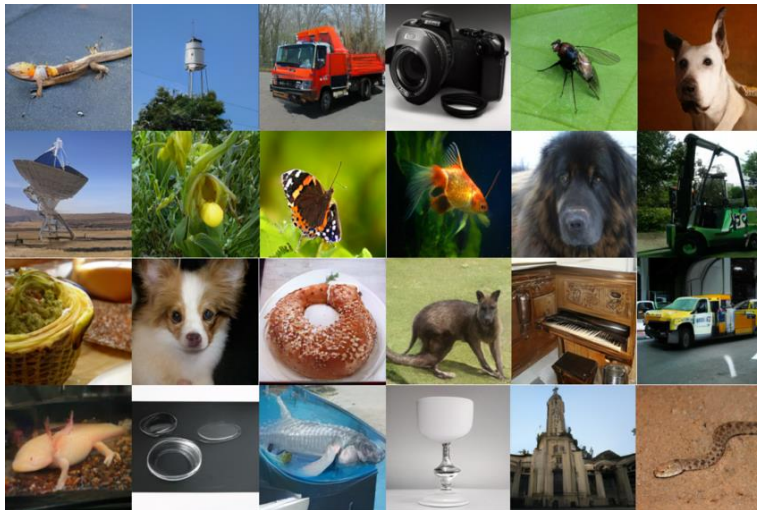
# Poll

# Content

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

# Conditional Diffusion Models

- Un-conditional



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

- Conditional



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

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

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

- $\gamma$  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

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

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

- 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  $\mathbf{p}$  original conditionals with an additional unconditional class

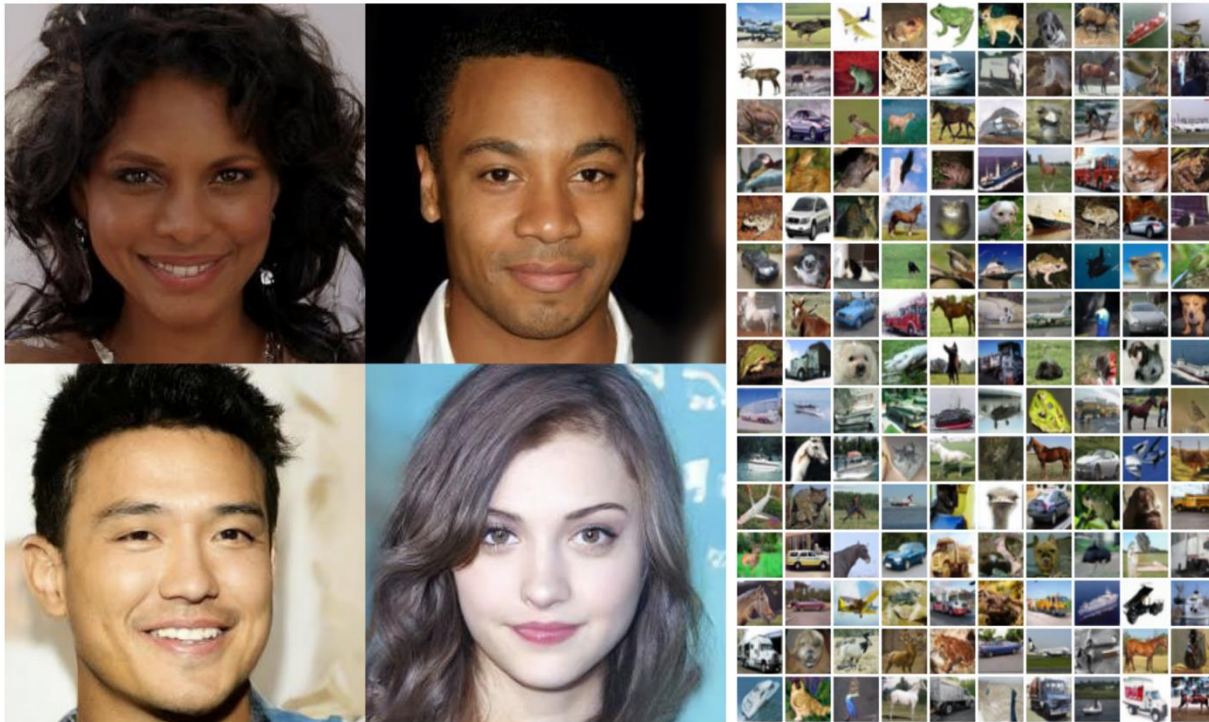
$$\mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right]$$

# Content

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

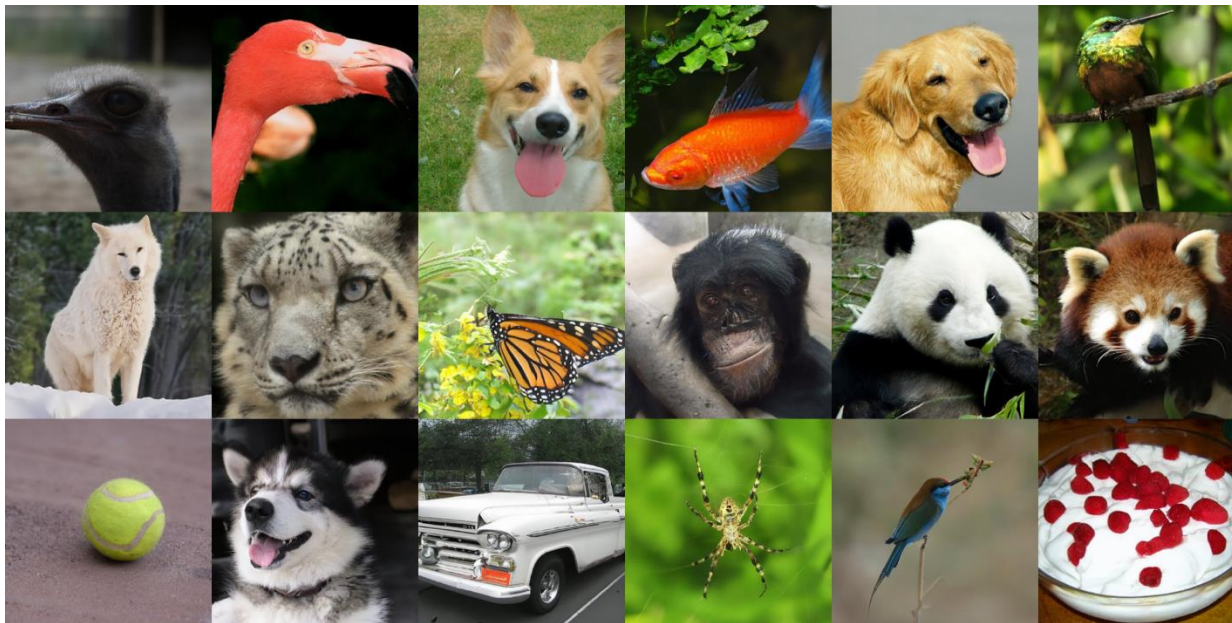
# DDPM

- Training diffusion models on raw images with a U-Net model



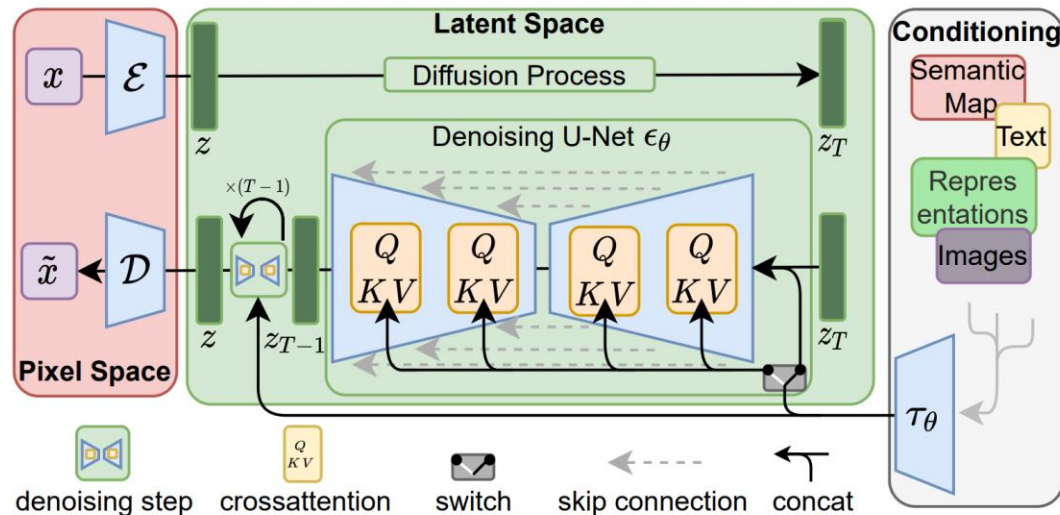
# Diffusion Models Beat GANs

- Larger denoising model with sophisticated design
  - Adaptive group normalization
  - Attention layers in U-Net



# Latent Diffusion Models (LDMs)

- Learn diffusion on VAE's latent
  - Yet another VAE! Except pre-trained.





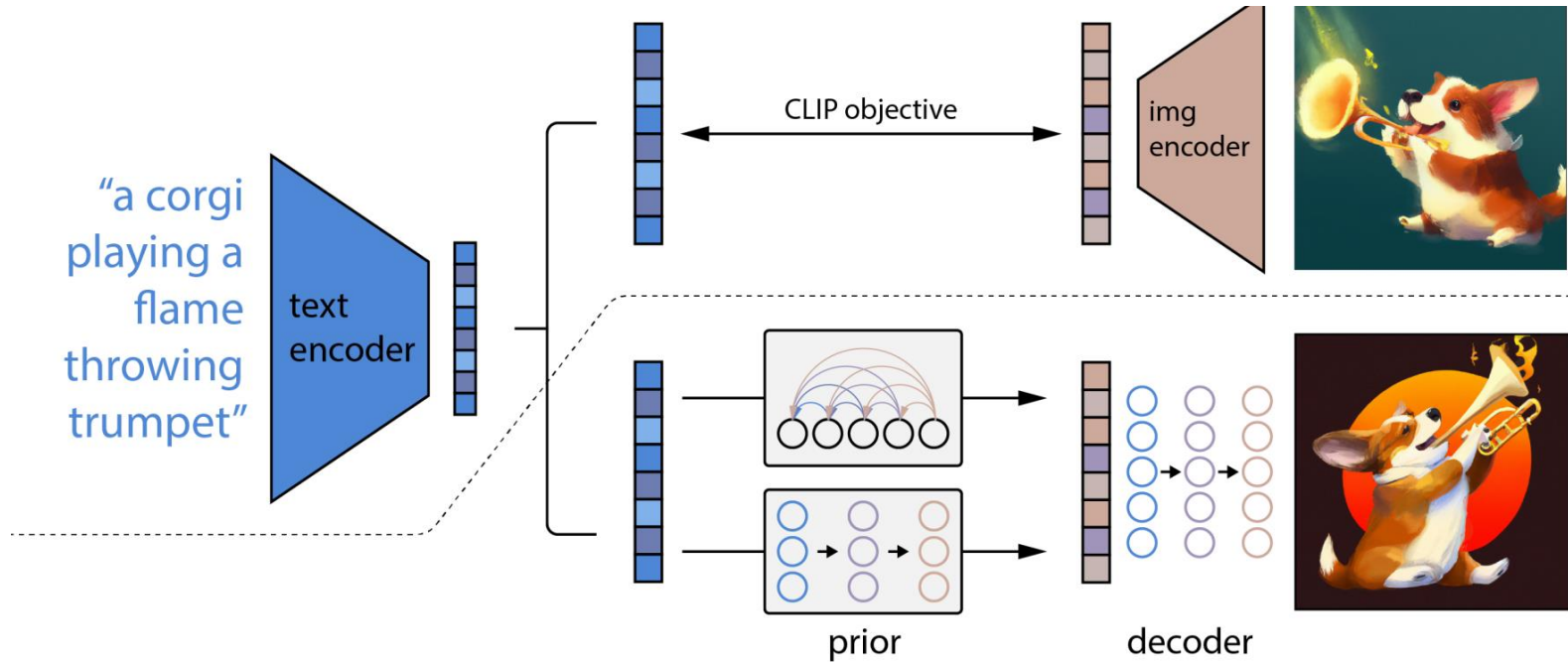
# Stable Diffusion

- Large-scale text-conditional LDMs
  - With VAEs trained also on larger datasets



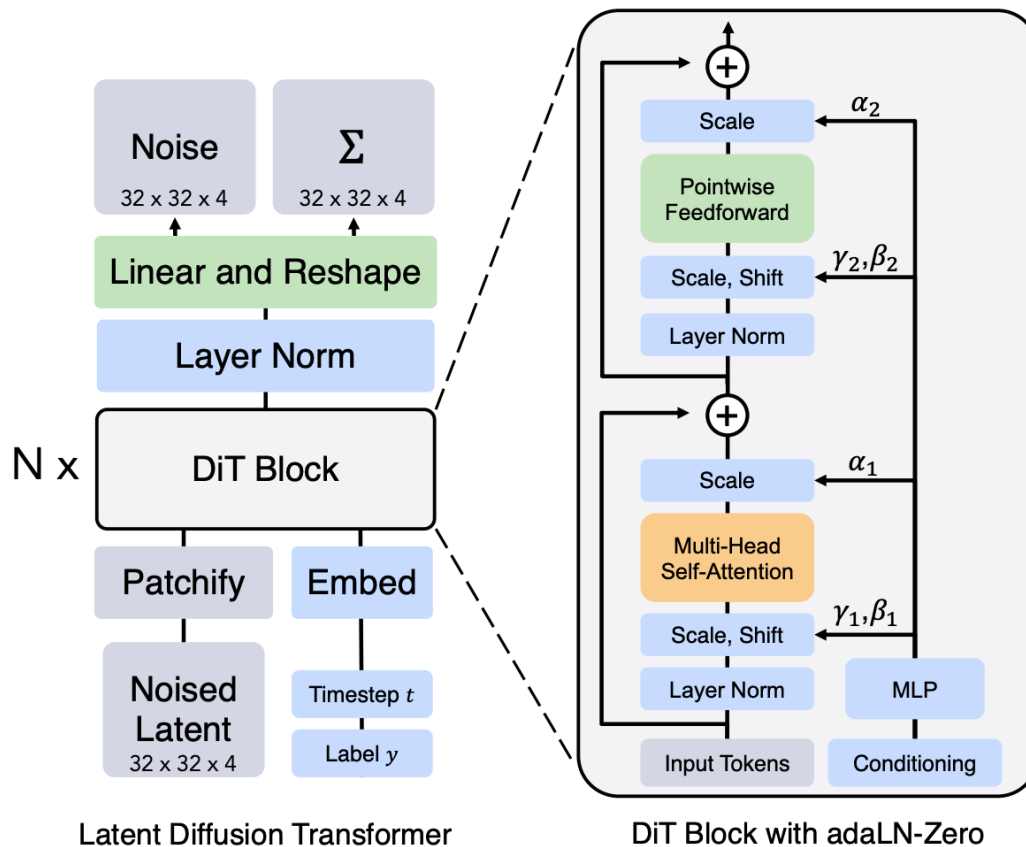


# DALLE



# DiT

- A transformer architecture for diffusion models



# MAR

- An autoregressive model with diffusion loss

