

# **Introduction to Diffusion Models**

**Part 1: Introduction, Formulation and Theory**

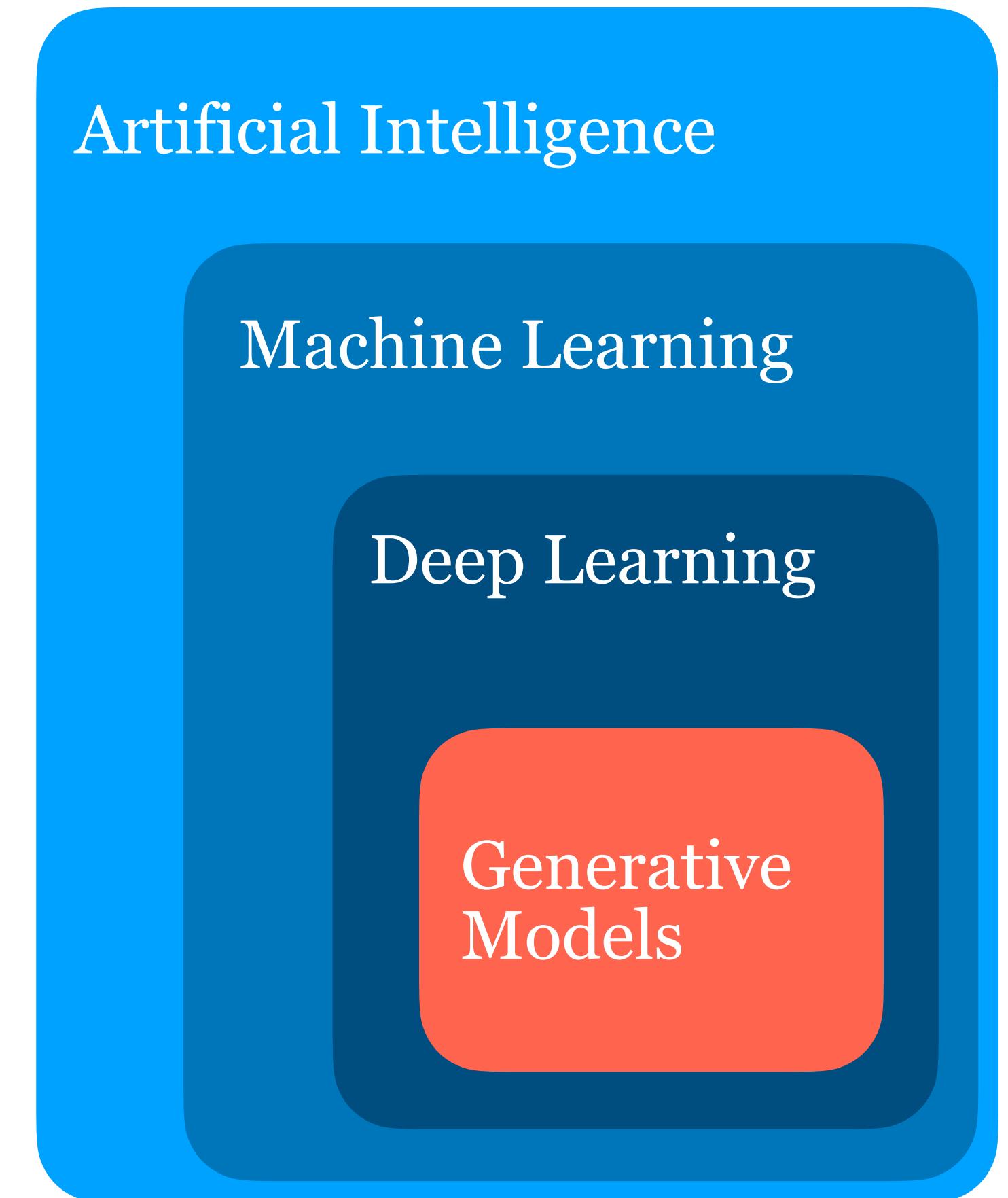
**Leonid Sigal and Siddhesh Khandelwal**

# Introduction : Generative Models

- Most AI methods are **predictive** in nature.
  - Observe and classify patterns in content or predict a new pattern
  - Ex: Label objects in an image.
- **Generative Models**, on the other hand, are capable of **creating** text, images, videos, and other types of content.
  - These models can be **conditioned** on a text prompt to get the desired output.

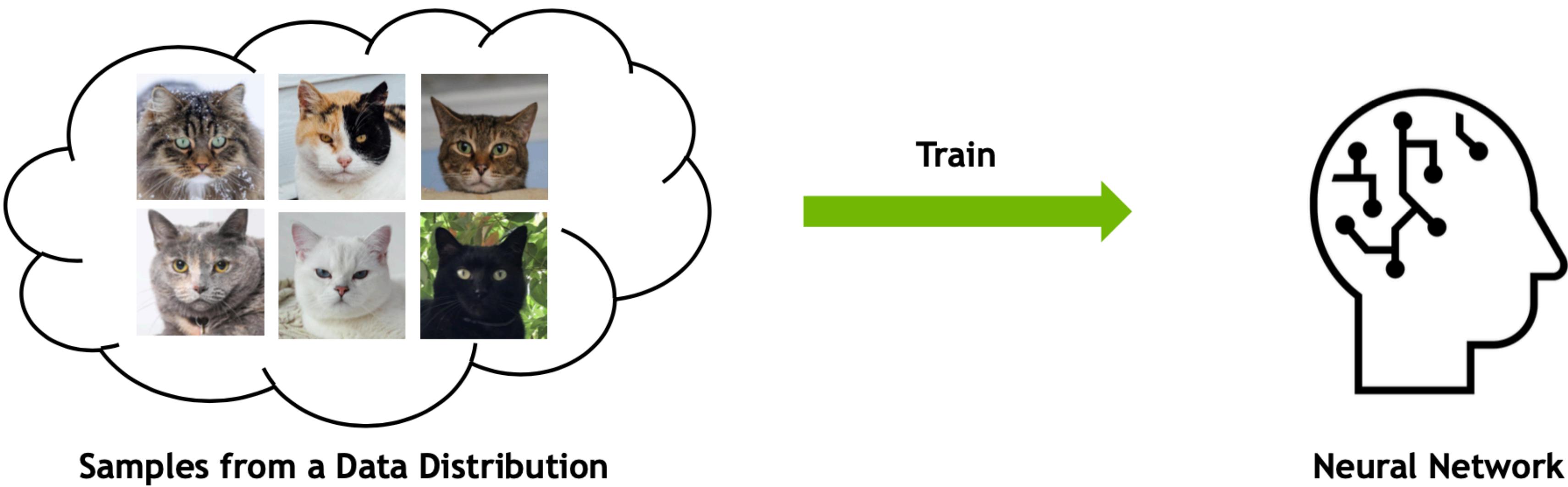
# Introduction : Generative Models

- Most AI methods are **predictive** in nature.
  - Observe and classify patterns in content or predict a new pattern
  - Ex: Label objects in an image.
- **Generative Models**, on the other hand, are capable of **creating** text, images, videos, and other types of content.
  - These models can be **conditioned** on a text prompt to get the desired output.



# **Introduction : Generative Models**

# Introduction : Generative Models



# Introduction : Generative Models



Samples from a Data Distribution

Train



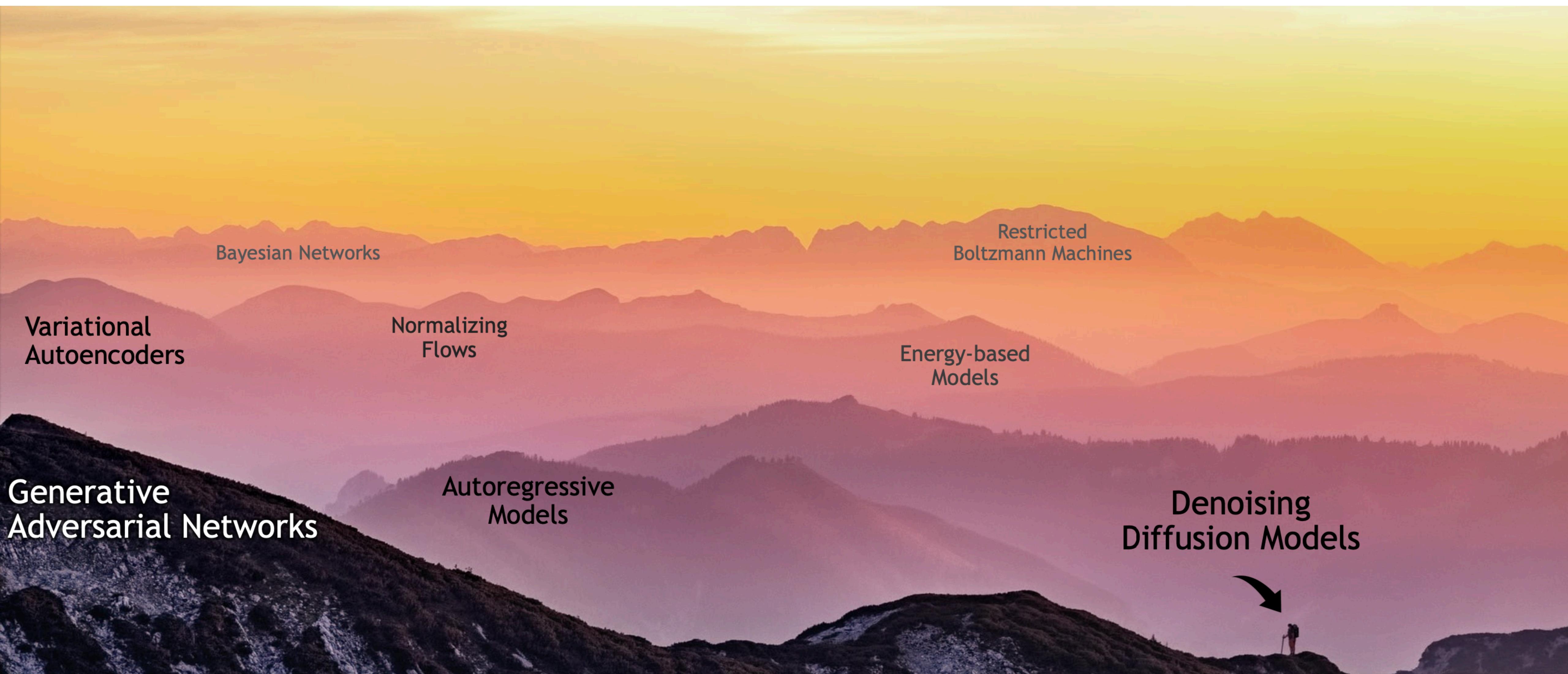
Neural Network



Sample



# Introduction : Generative Models



# **Properties of Good Generative Model**

# Properties of Good Generative Model

High Quality  
Samples

Mode Coverage  
/ Diversity

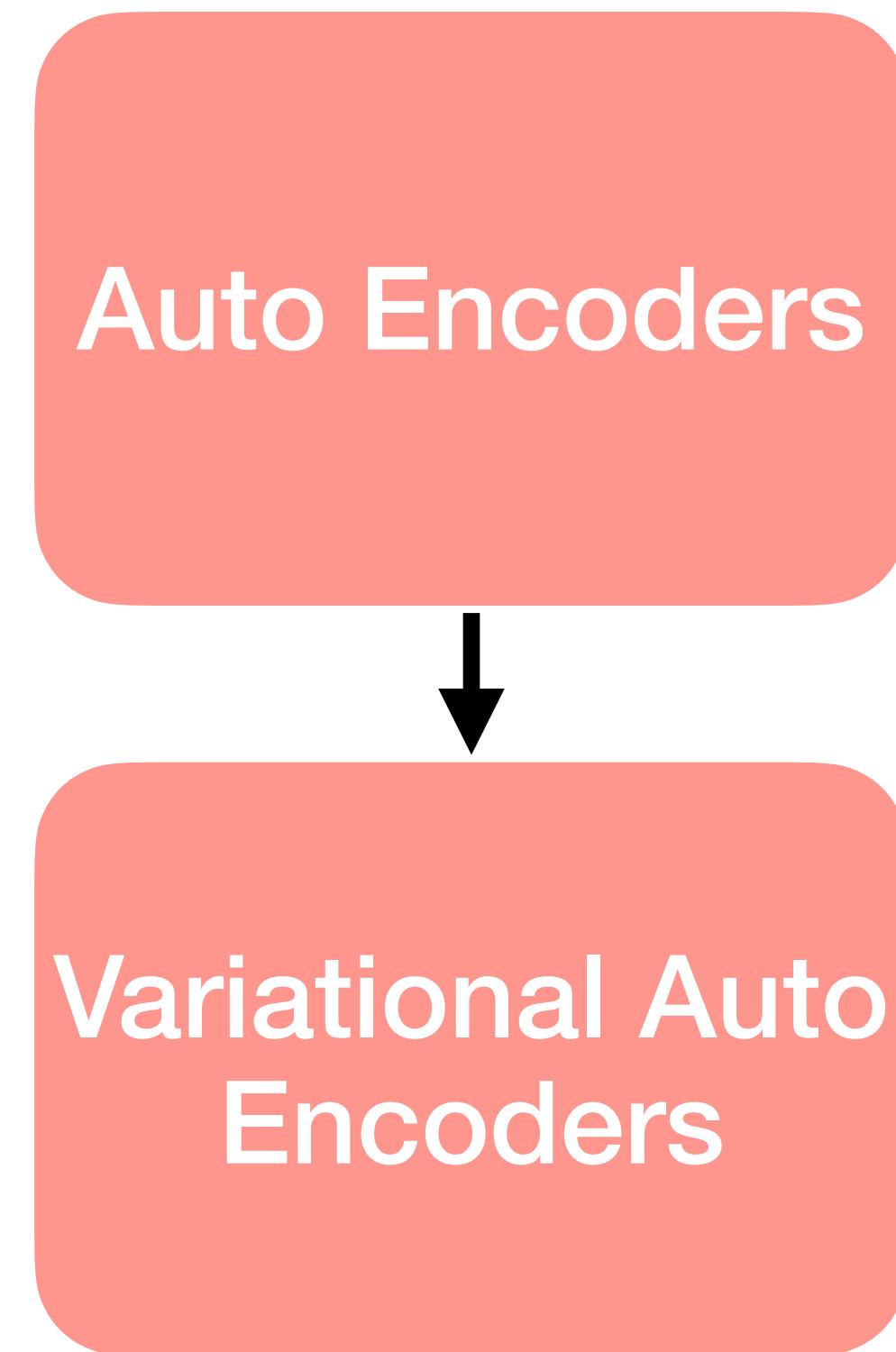
Fast Sampling

# Overview

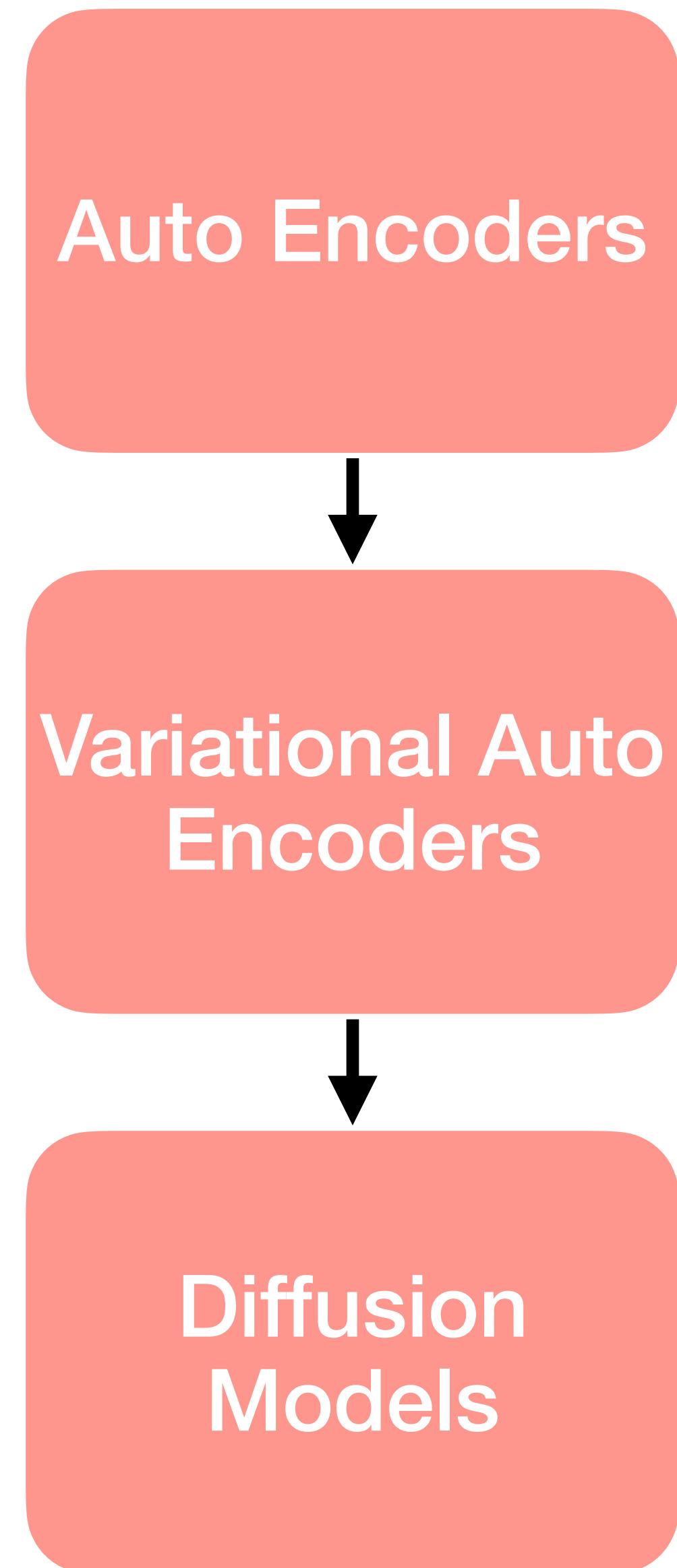
# Overview

Auto Encoders

# Overview



# Overview



# Introduction : Auto Encoders

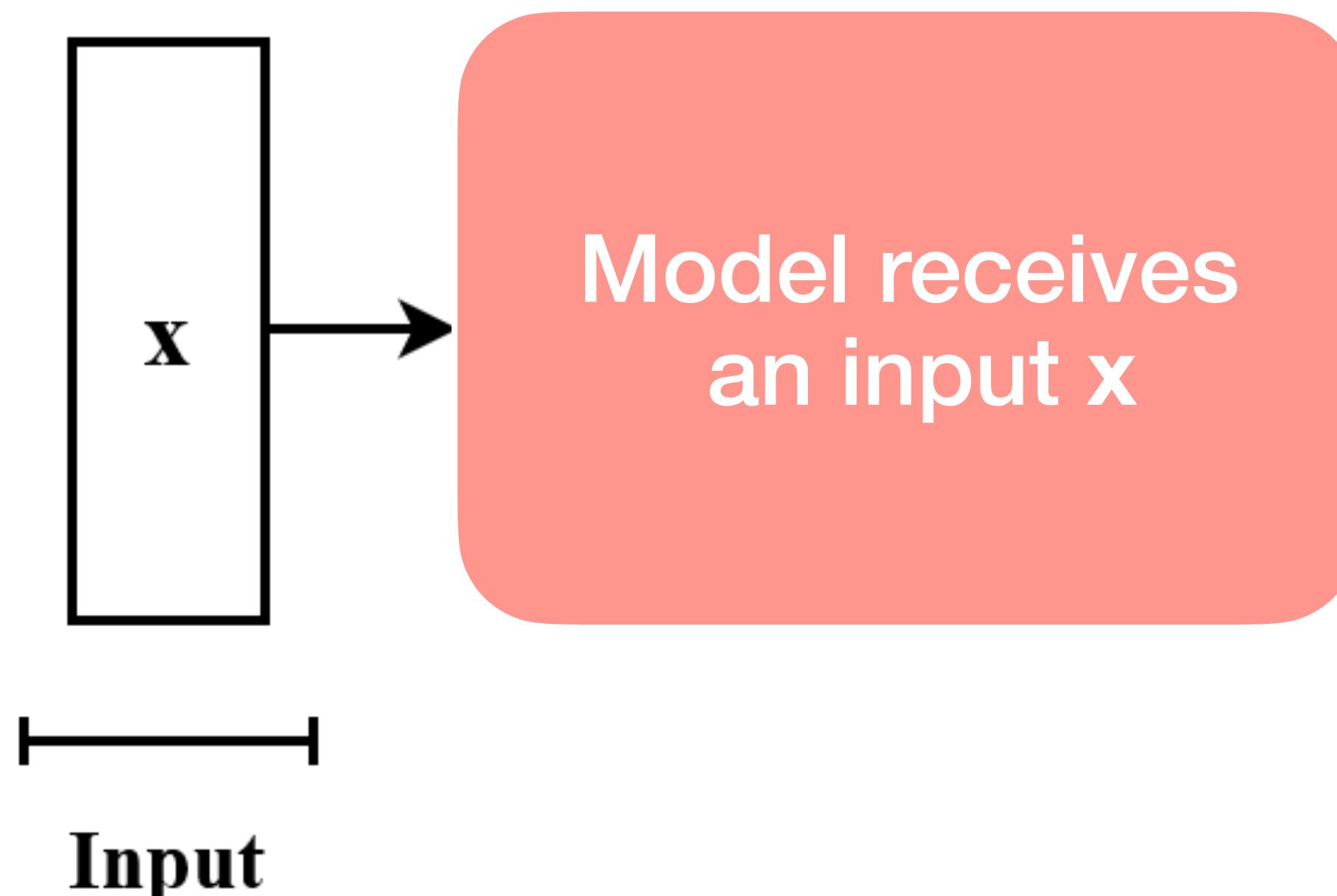
- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$

# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components  
**Encoder:** transforms the input into a latent representation

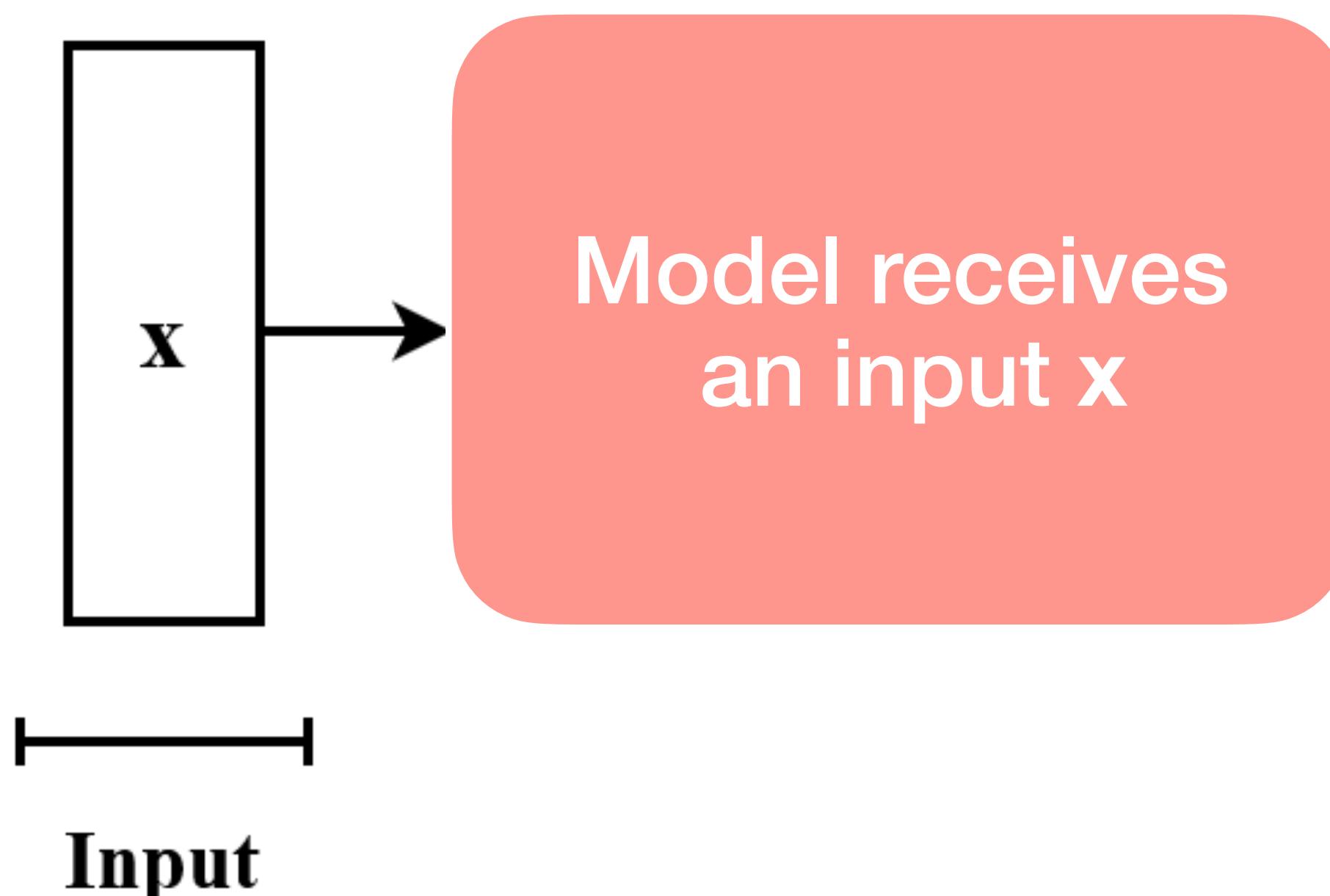
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $x$  and predict the same output  $x$ 
  - Consists of two components  
**Encoder:** transforms the input into a latent representation



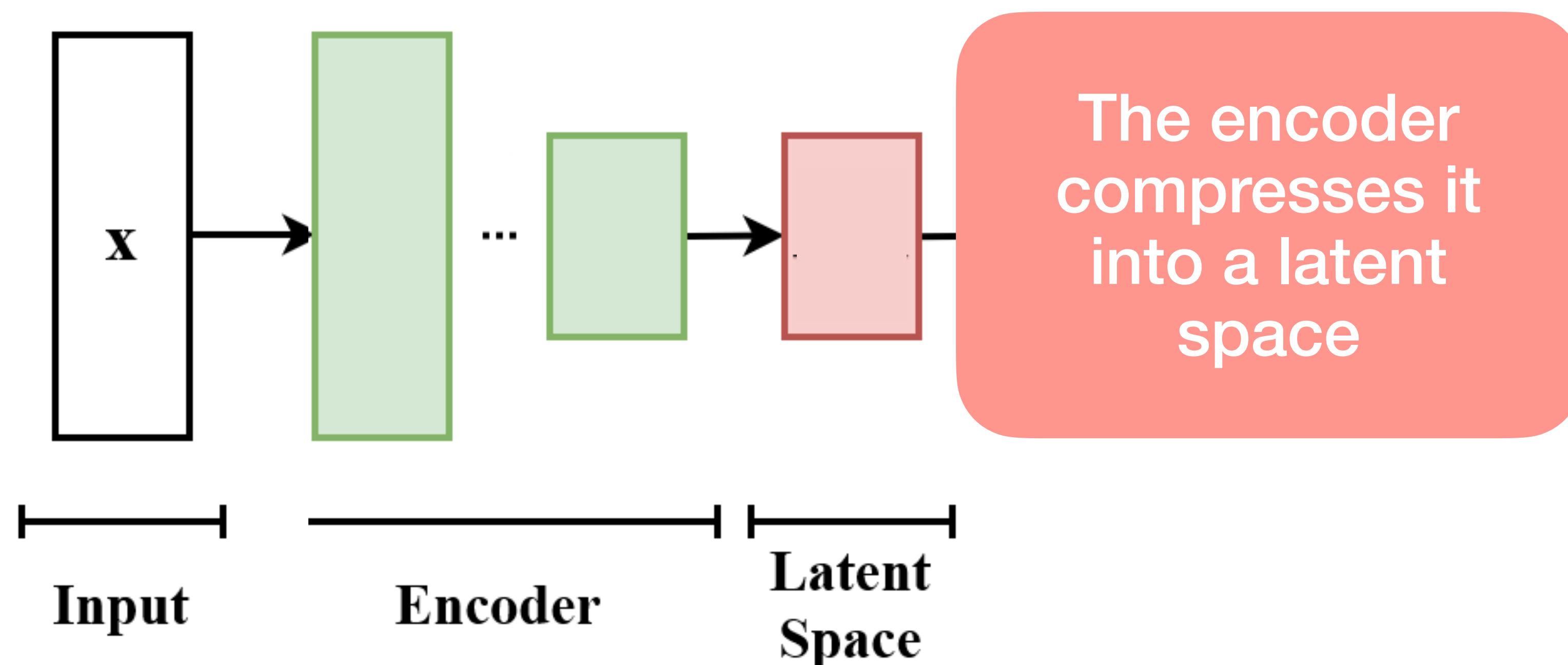
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $x$  and predict the same output  $x$ 
  - Consists of two components  
**Encoder:** transforms the input into a latent representation



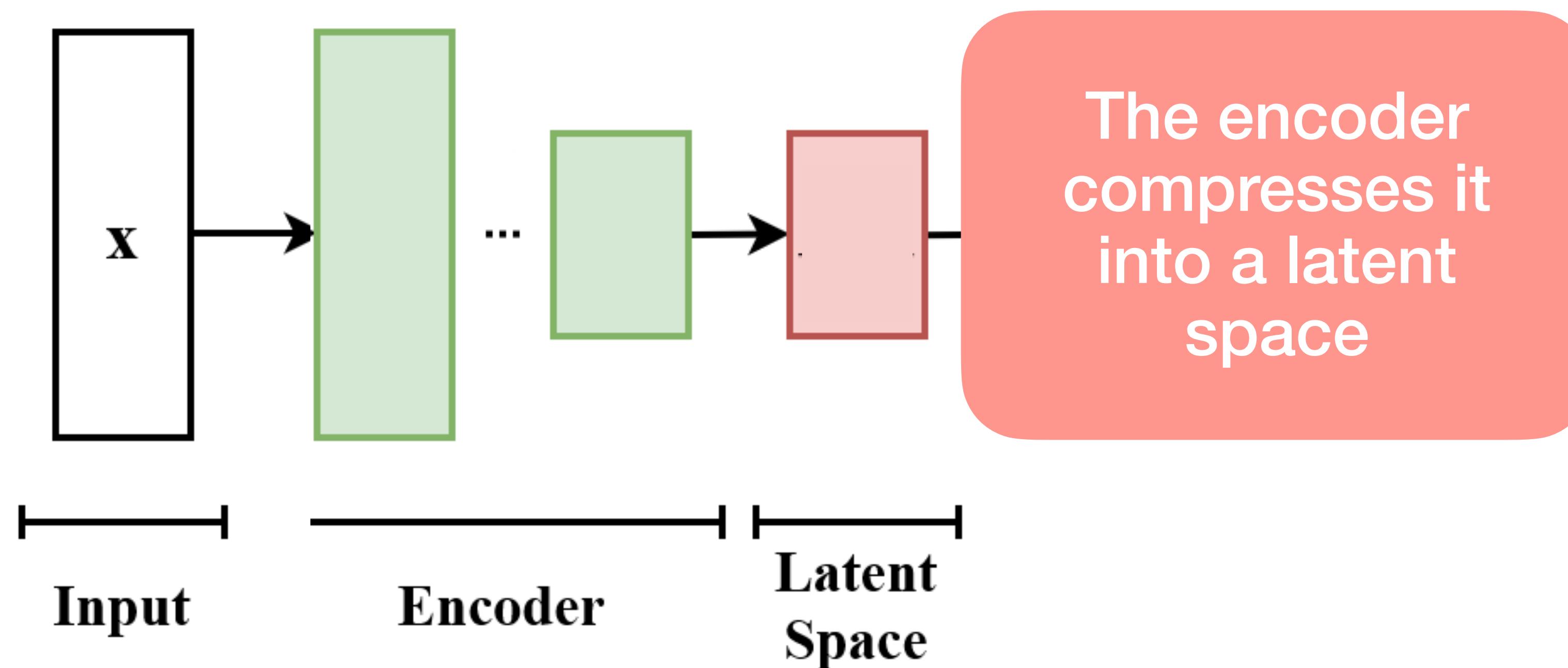
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components  
**Encoder:** transforms the input into a latent representation



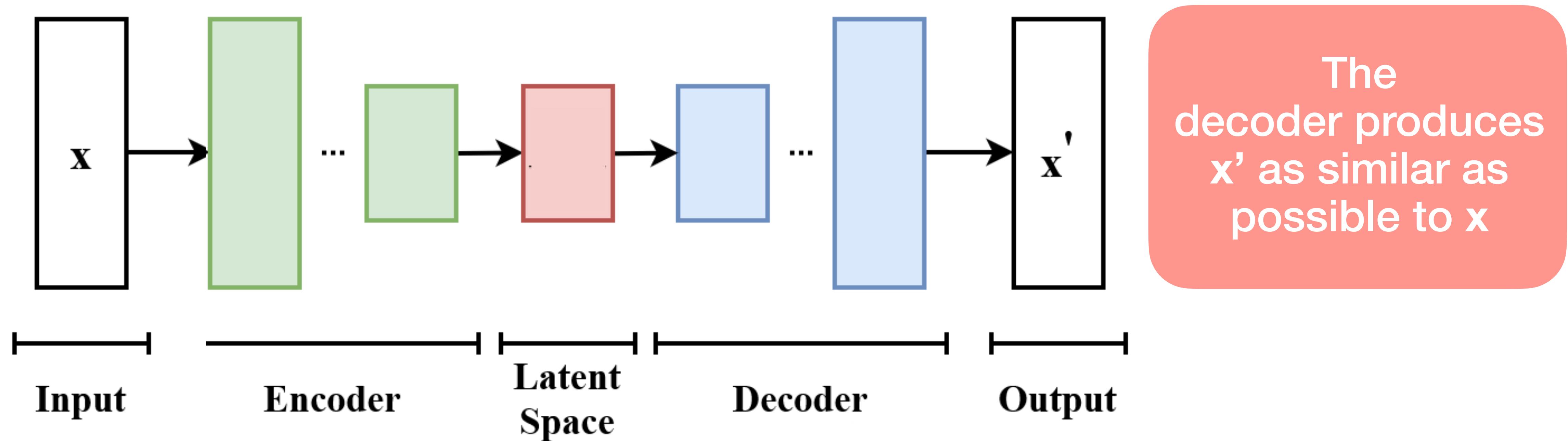
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components
    - Encoder:** transforms the input into a latent representation
    - Decoder:** recreates the input from the latent representation



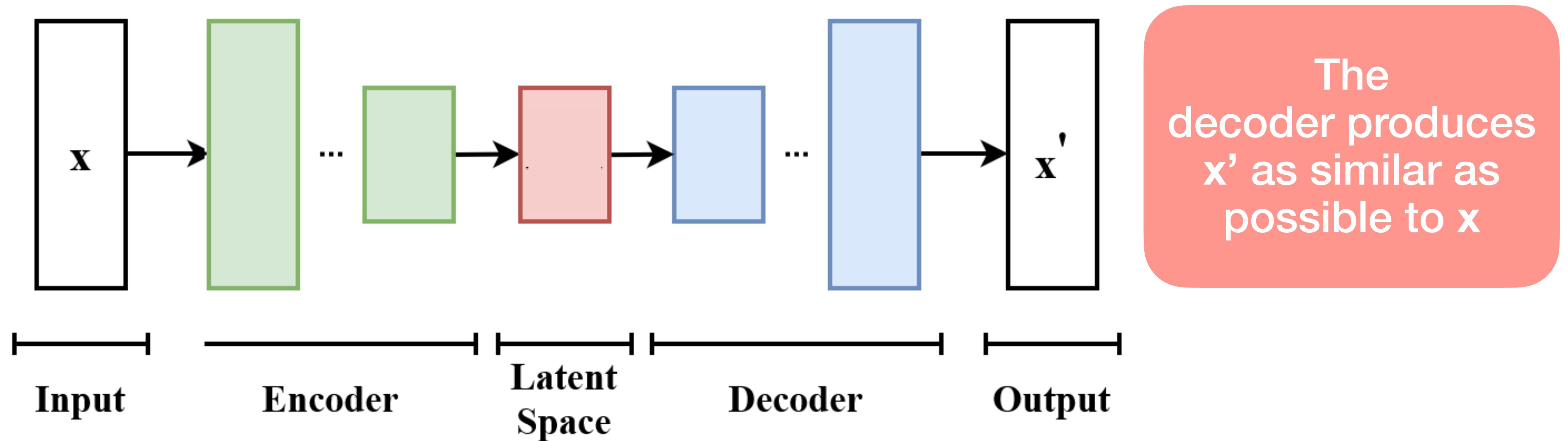
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components
    - Encoder:** transforms the input into a latent representation
    - Decoder:** recreates the input from the latent representation



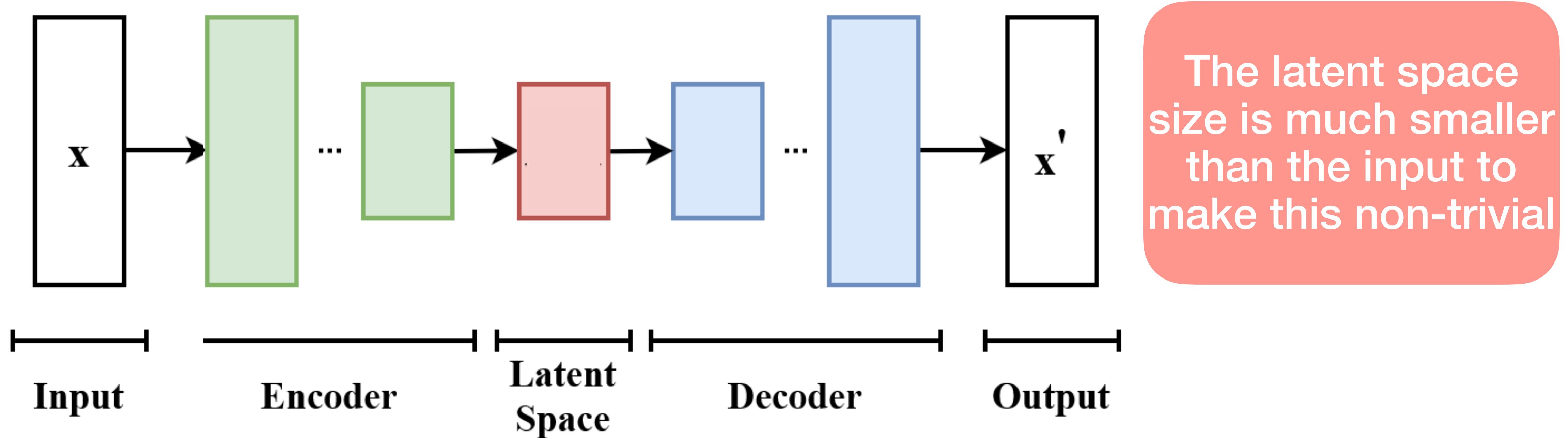
# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components
    - Encoder:** transforms the input into a latent representation
    - Decoder:** recreates the input from the latent representation



# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}'$ 
  - Consists of two components
    - Encoder:** transforms the input into a latent representation
    - Decoder:** recreates the input from the latent representation

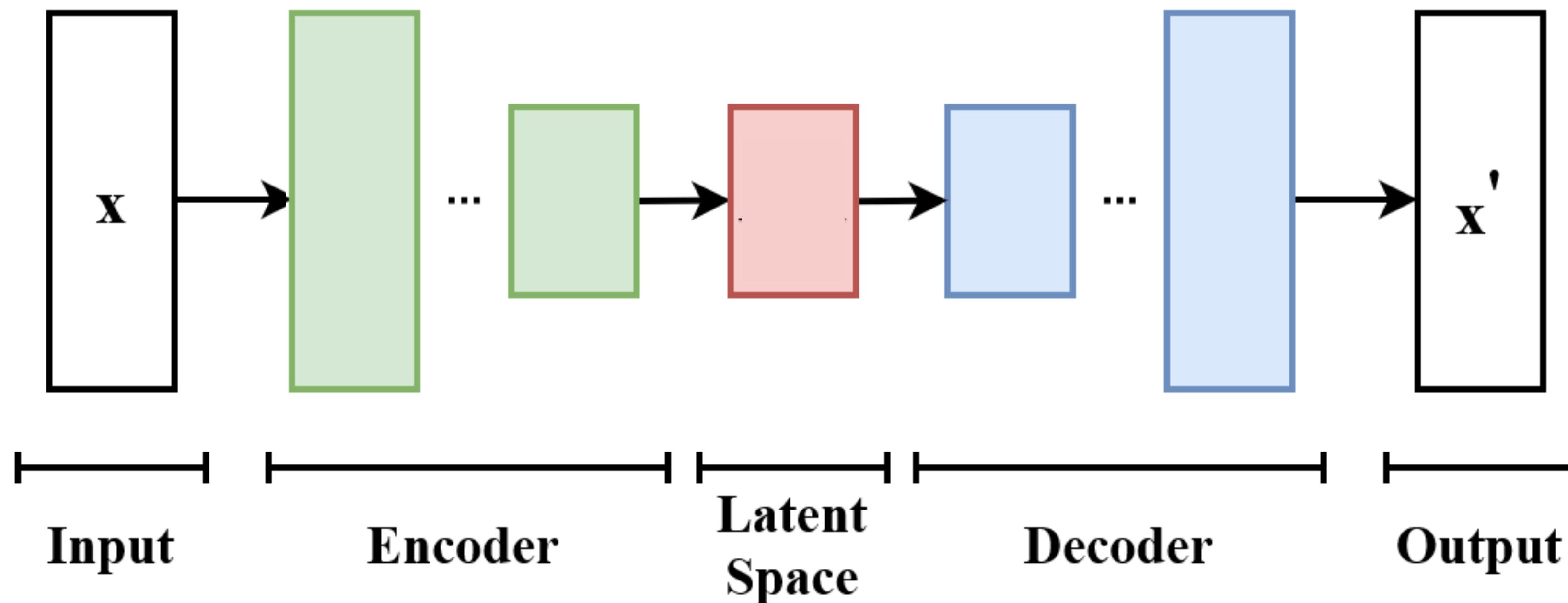


# Introduction : Auto Encoders

- An autoencoder is a neural network whose job is to take an input  $\mathbf{x}$  and predict the same output  $\mathbf{x}$ 
  - Consists of two components
    - Encoder:** transforms the input into a latent representation
    - Decoder:** recreates the input from the latent representation
- Autoencoders essentially perform **compression**
  - Map high-dimensional data to a lower-dimensional latent code.

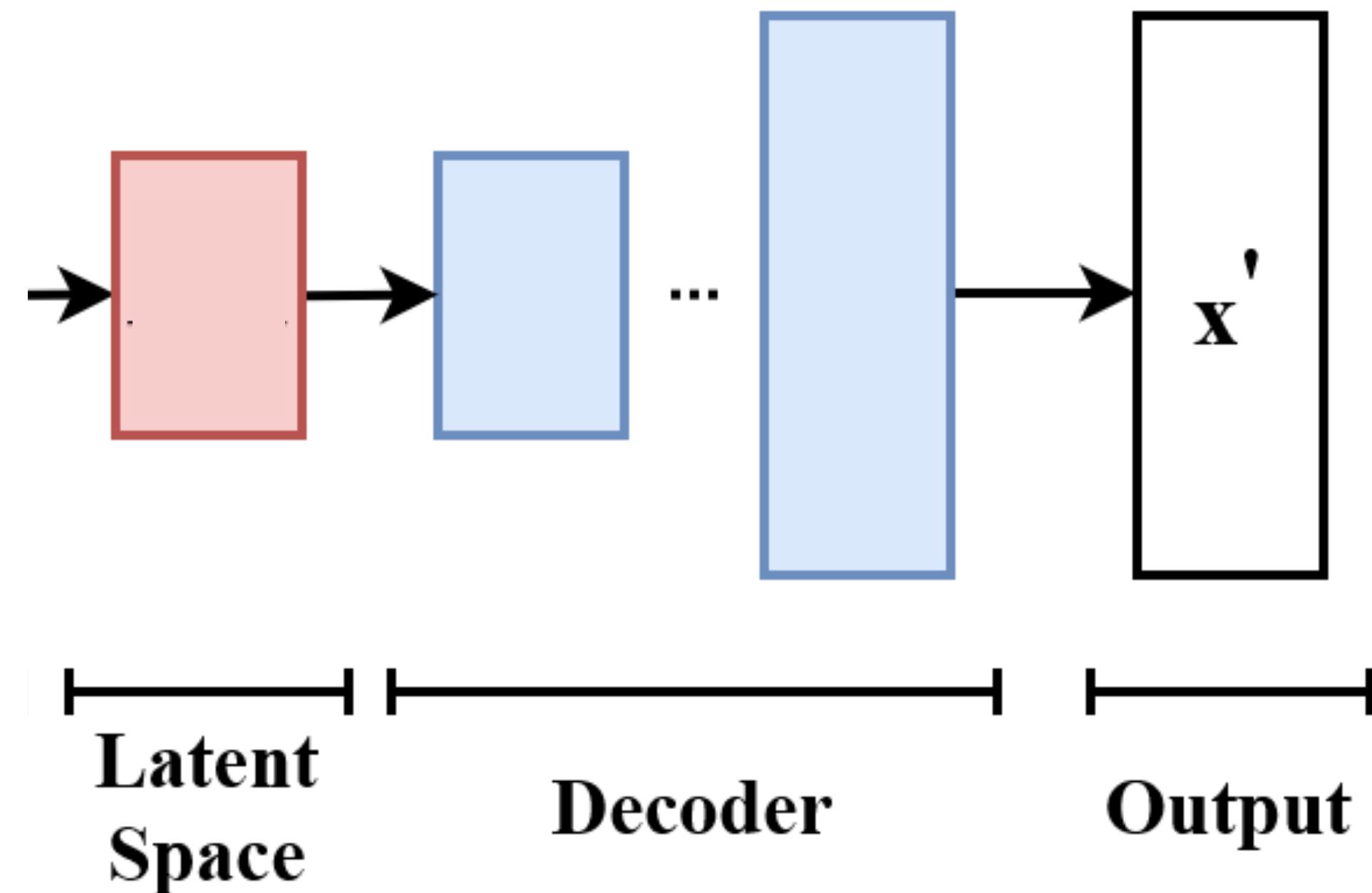
# Introduction : Auto Encoders

- Autoencoders can technically be used for data generation.
  - Pass an arbitrary latent code into the decoder.



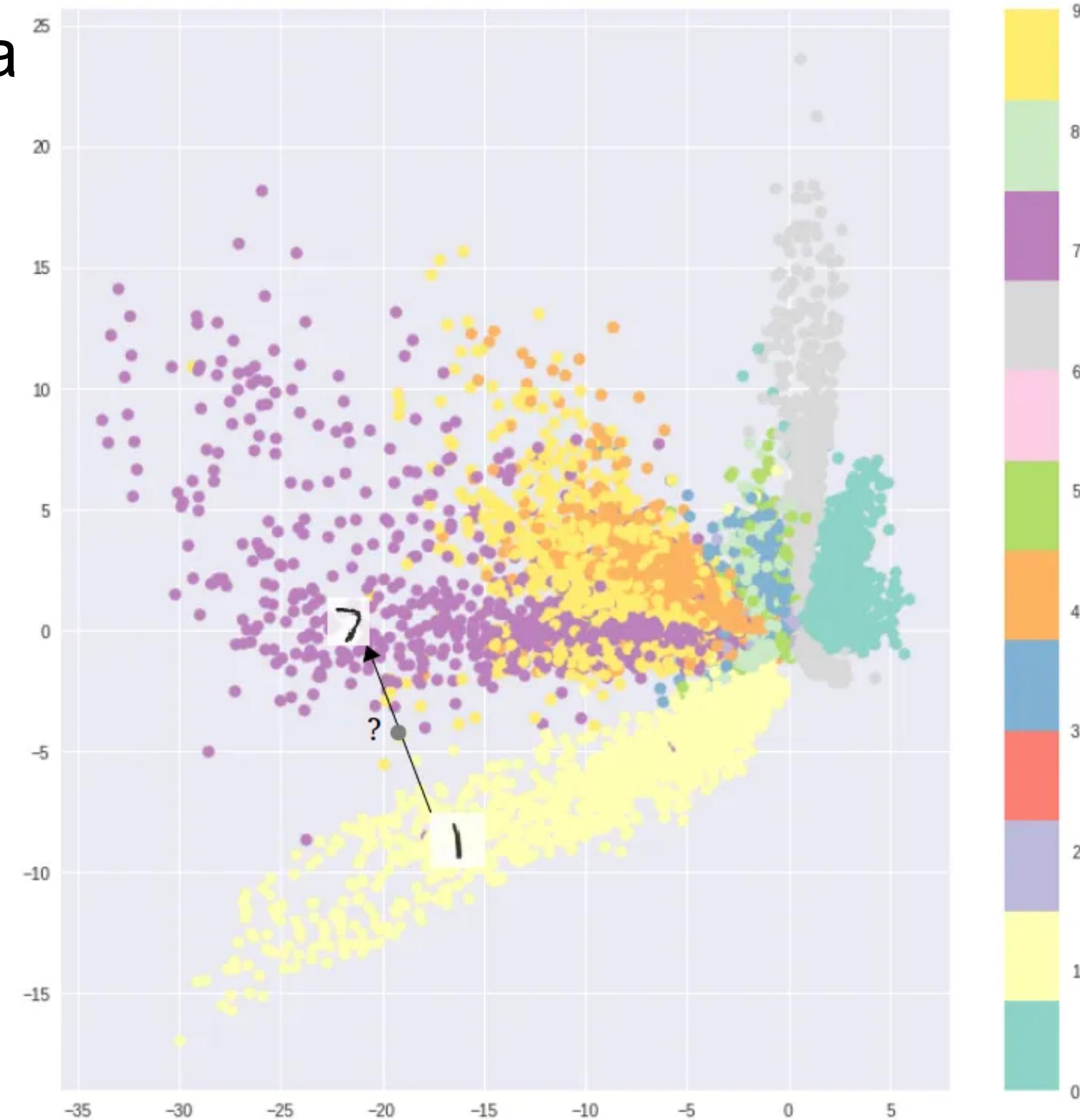
# Introduction : Auto Encoders

- Autoencoders can technically be used for data generation.
  - Pass an arbitrary latent code into the decoder.



# Introduction : Auto Encoders

- Autoencoders can technically be used for data generation.
  - Pass an arbitrary latent code into the decoder.
- However, the fundamental problem with autoencoders is that the latent space they learn may **not be continuous**, or allow easy interpolation.



# Introduction : Variational Auto Encoders

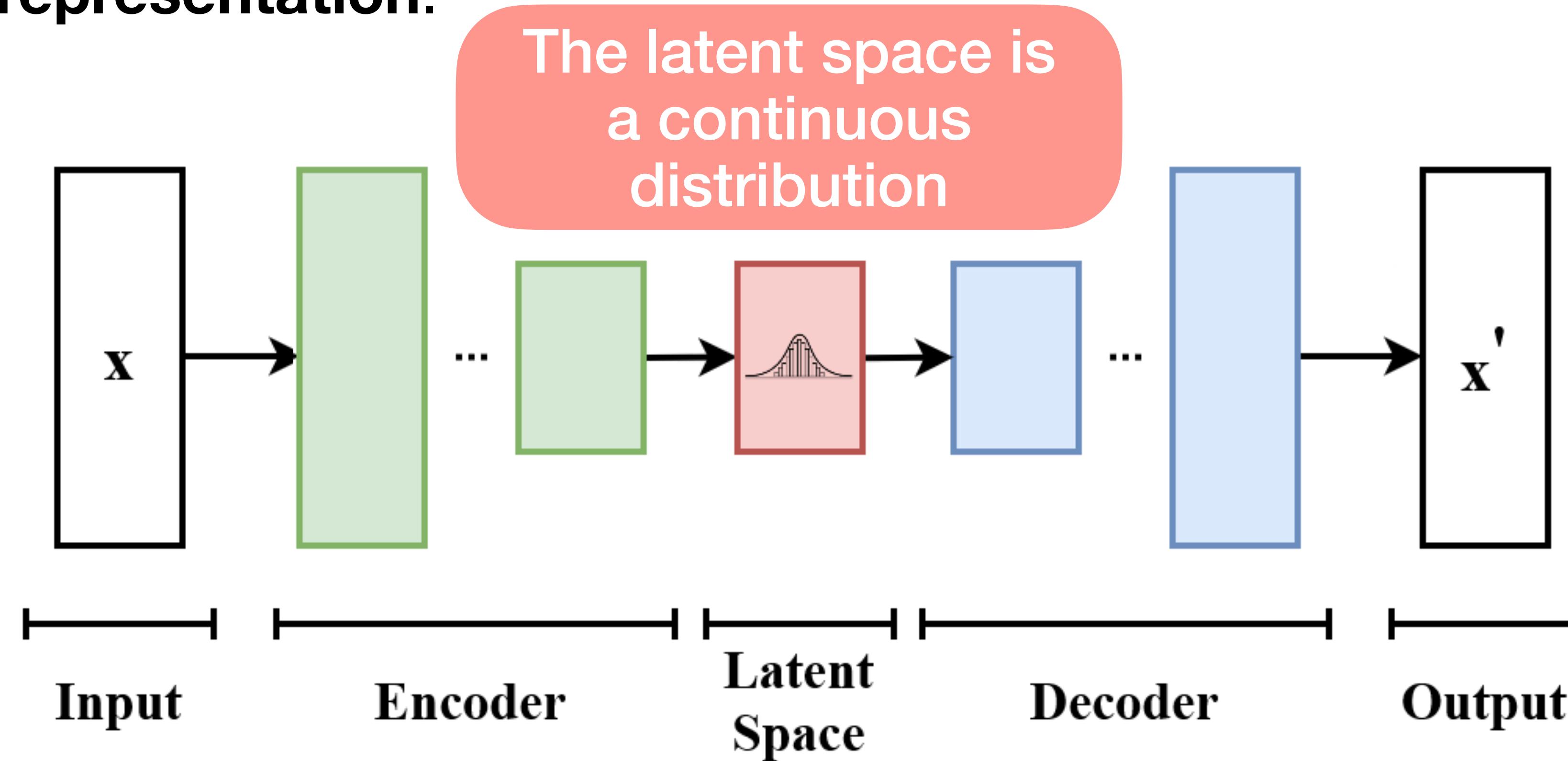
- Variational Auto Encoders (**VAEs**) are a popular class of generative models.
  - Contrary to autoencoders, they learn to map the input to a **continuous latent representation**.

# Introduction : Variational Auto Encoders

- Variational Auto Encoders (**VAEs**) are a popular class of generative models.
  - Contrary to autoencoders, they learn to map the input to a **continuous latent representation**.

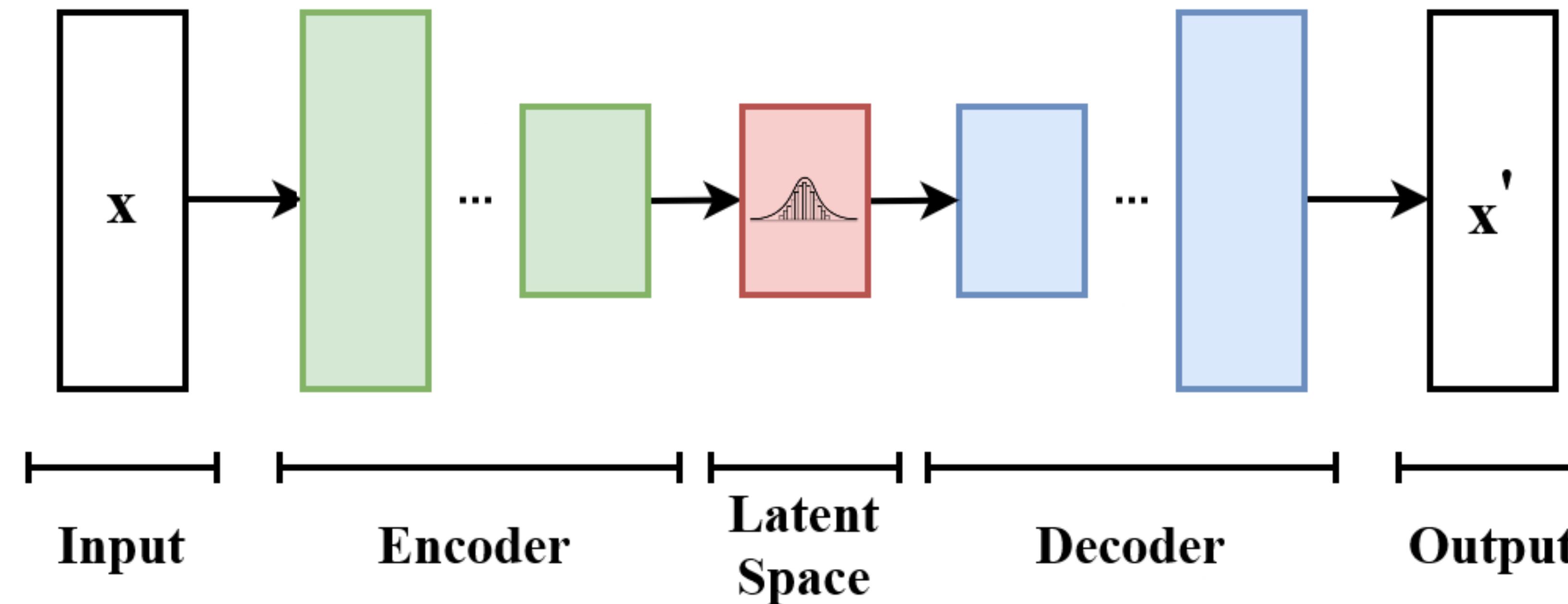
# Introduction : Variational Auto Encoders

- Variational Auto Encoders (**VAEs**) are a popular class of generative models.
  - Contrary to autoencoders, they learn to map the input to a **continuous latent representation**.



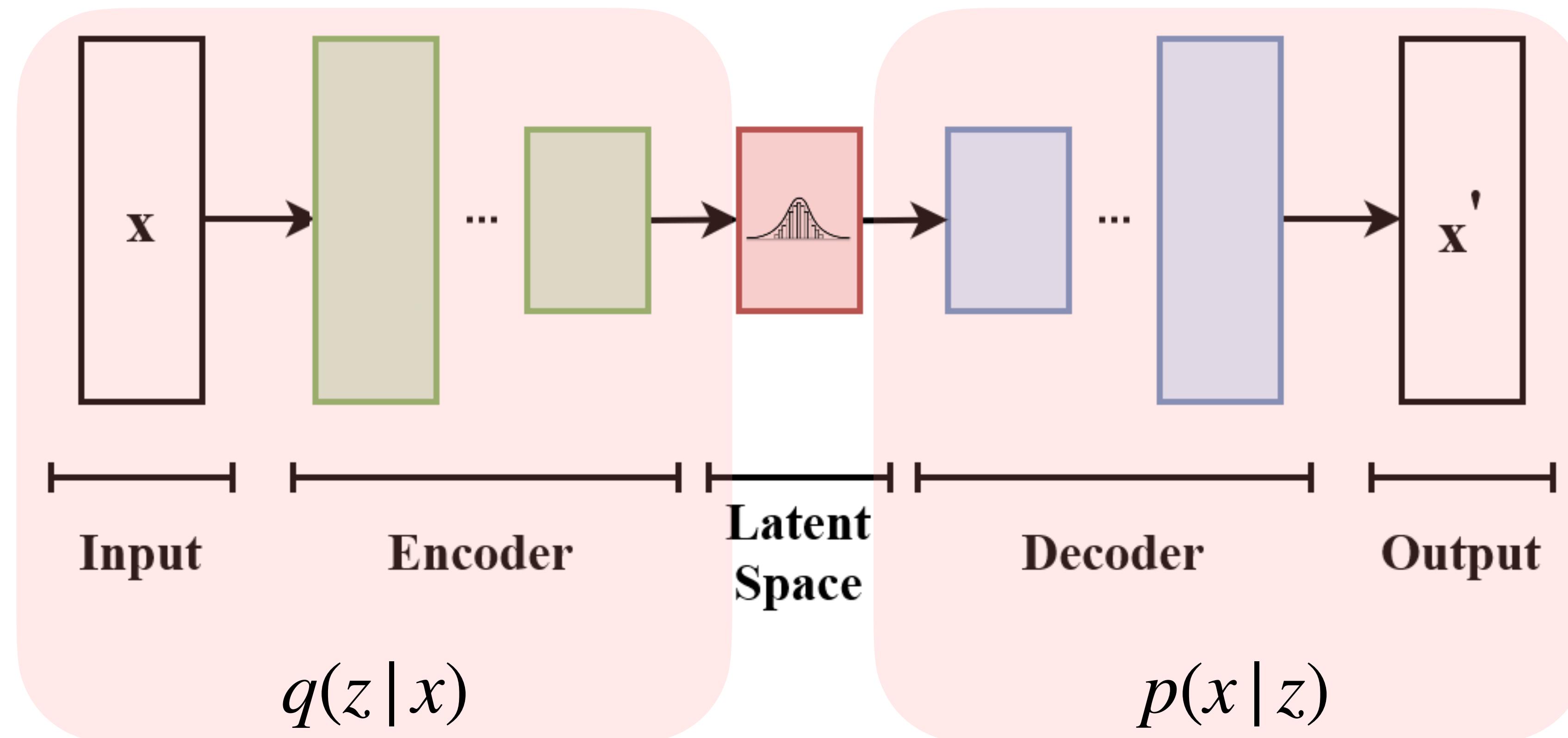
# Introduction : Variational Auto Encoders

- The encoder and decoder in a VAE are essentially modelling two different distributions.



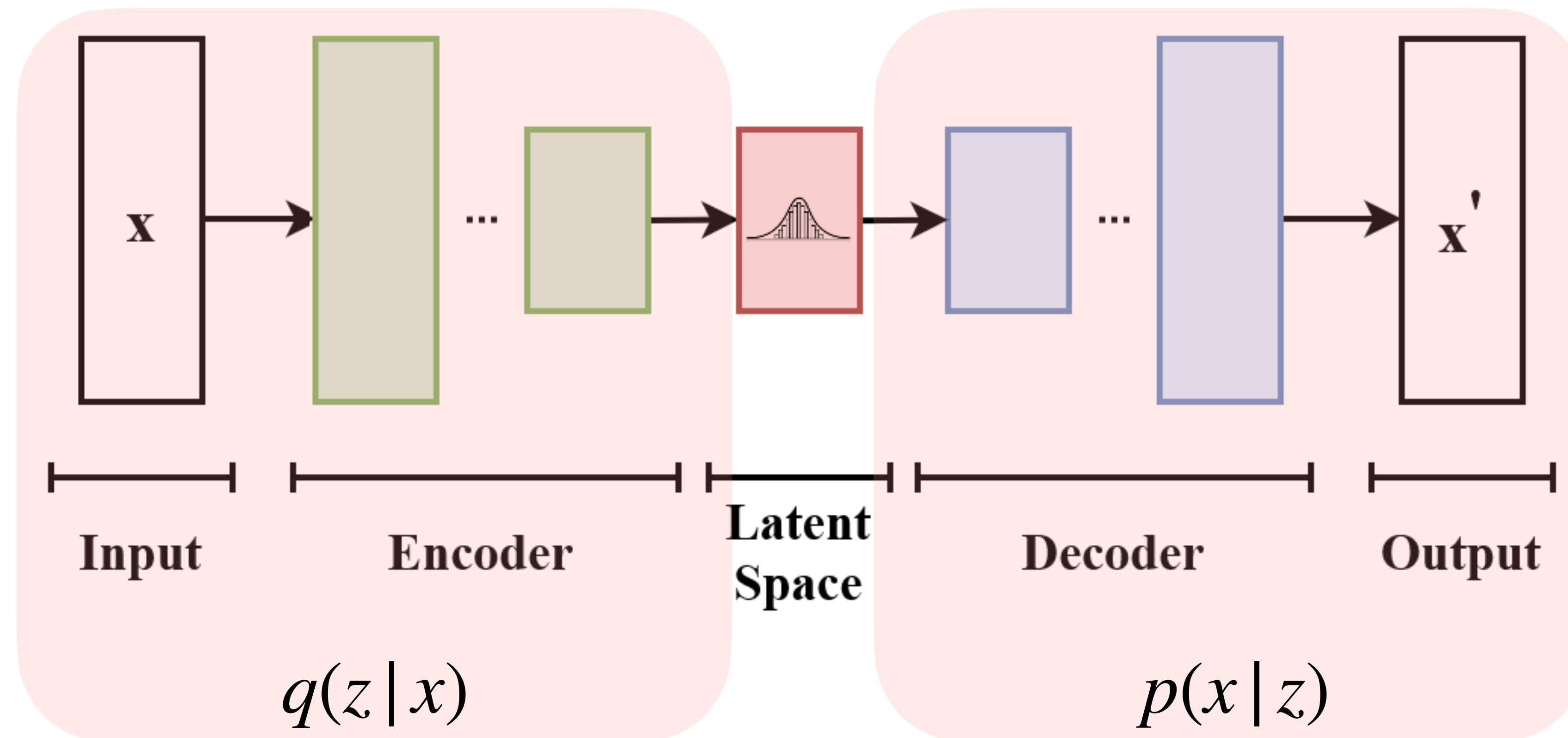
# Introduction : Variational Auto Encoders

- The encoder and decoder in a VAE are essentially modelling two different distributions.



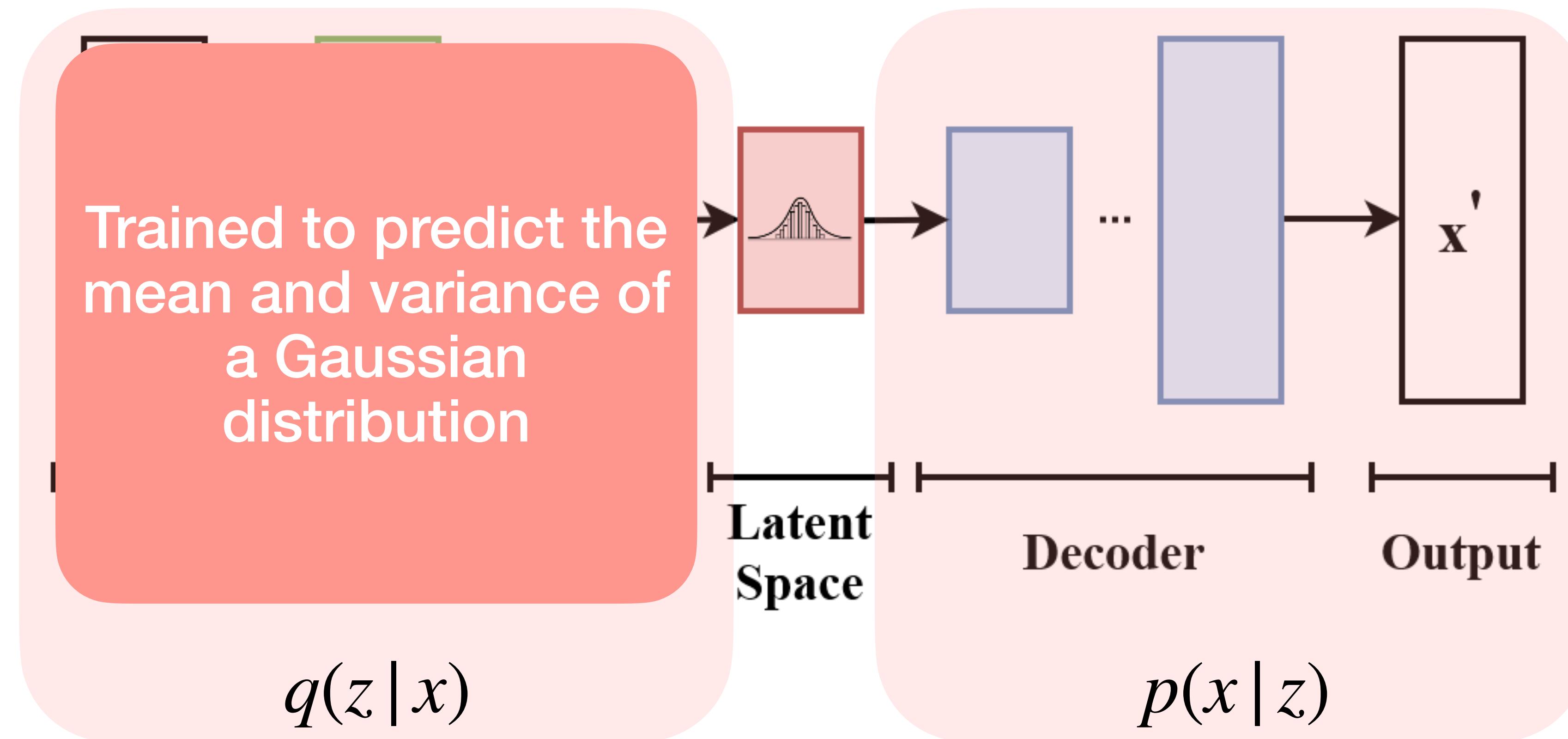
# Introduction : Variational Auto Encoders

- The encoder and decoder in a VAE are essentially modelling two different distributions.



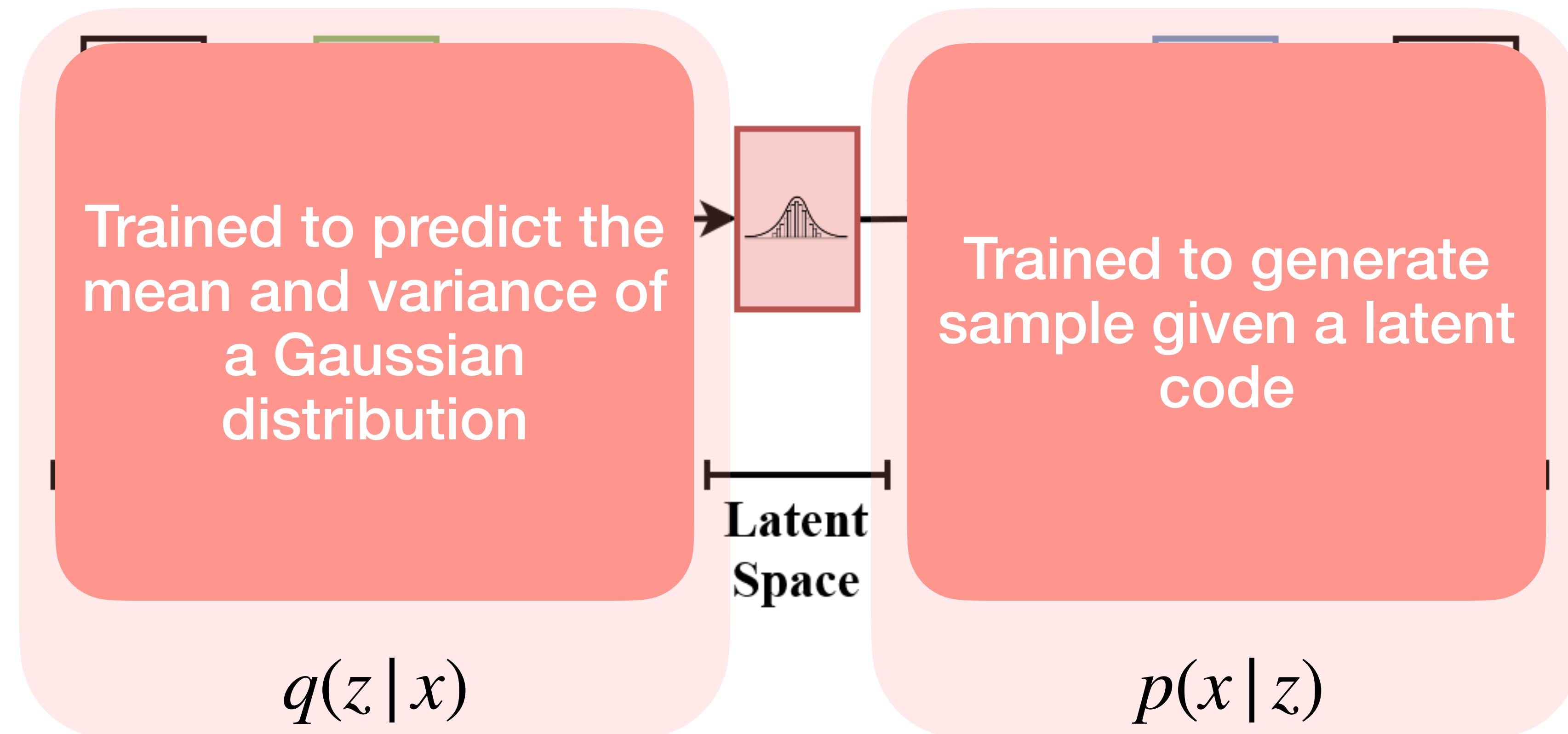
# Introduction : Variational Auto Encoders

- The encoder and decoder in a VAE are essentially modelling two different distributions.



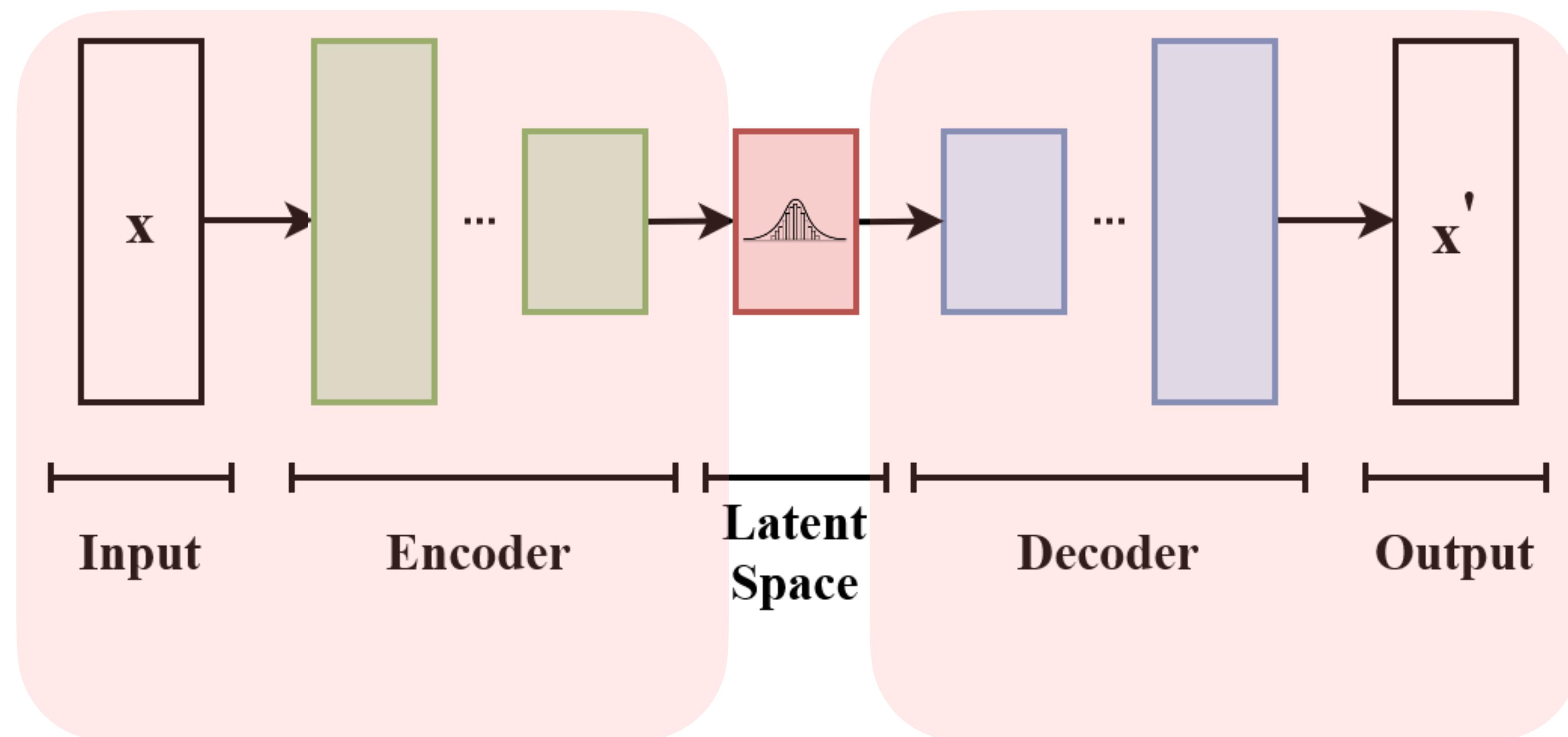
# Introduction : Variational Auto Encoders

- The encoder and decoder in a VAE are essentially modelling two different distributions.



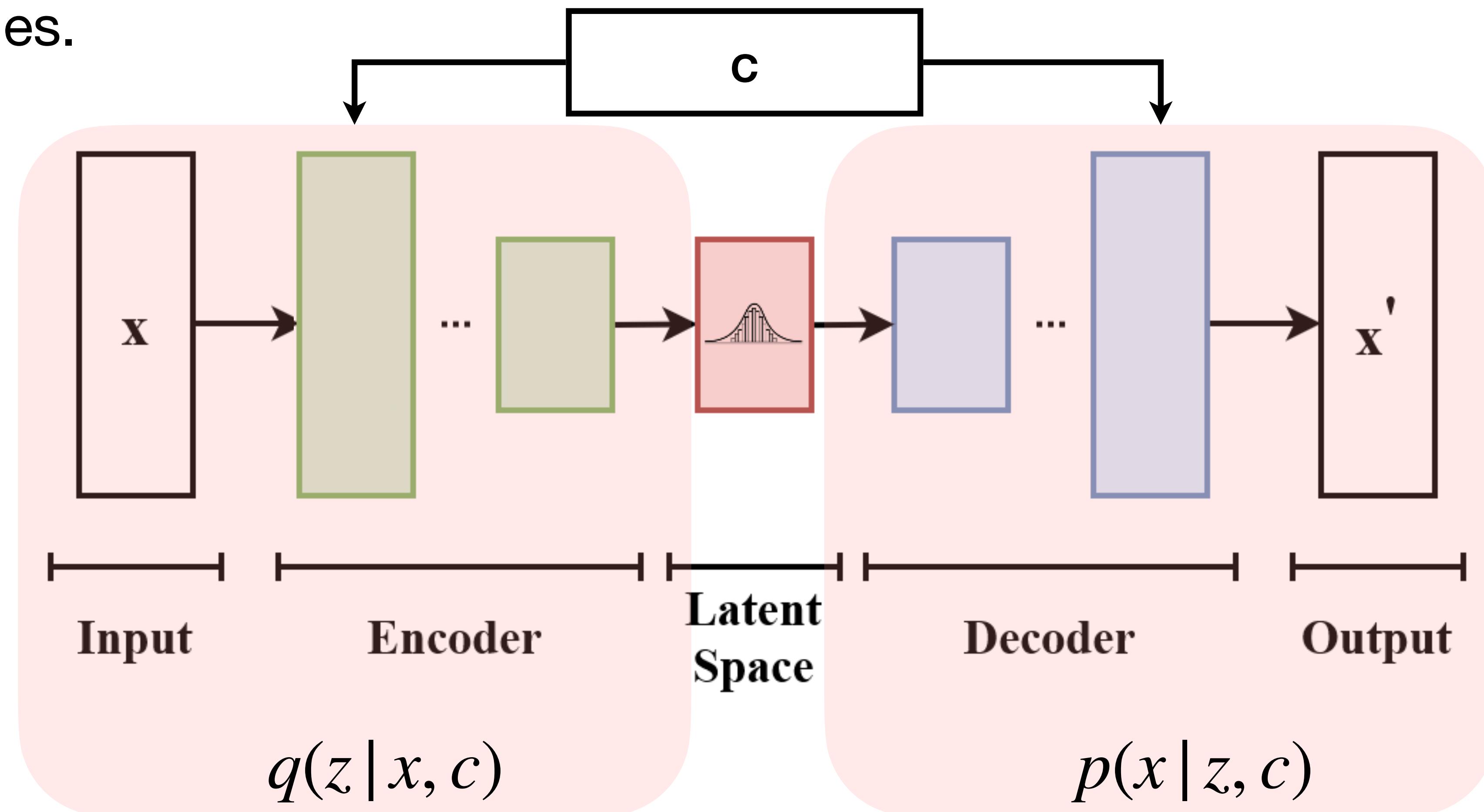
# Introduction : Variational Auto Encoders

- VAEs can also be conditioned on a label / class to generate specific output categories.



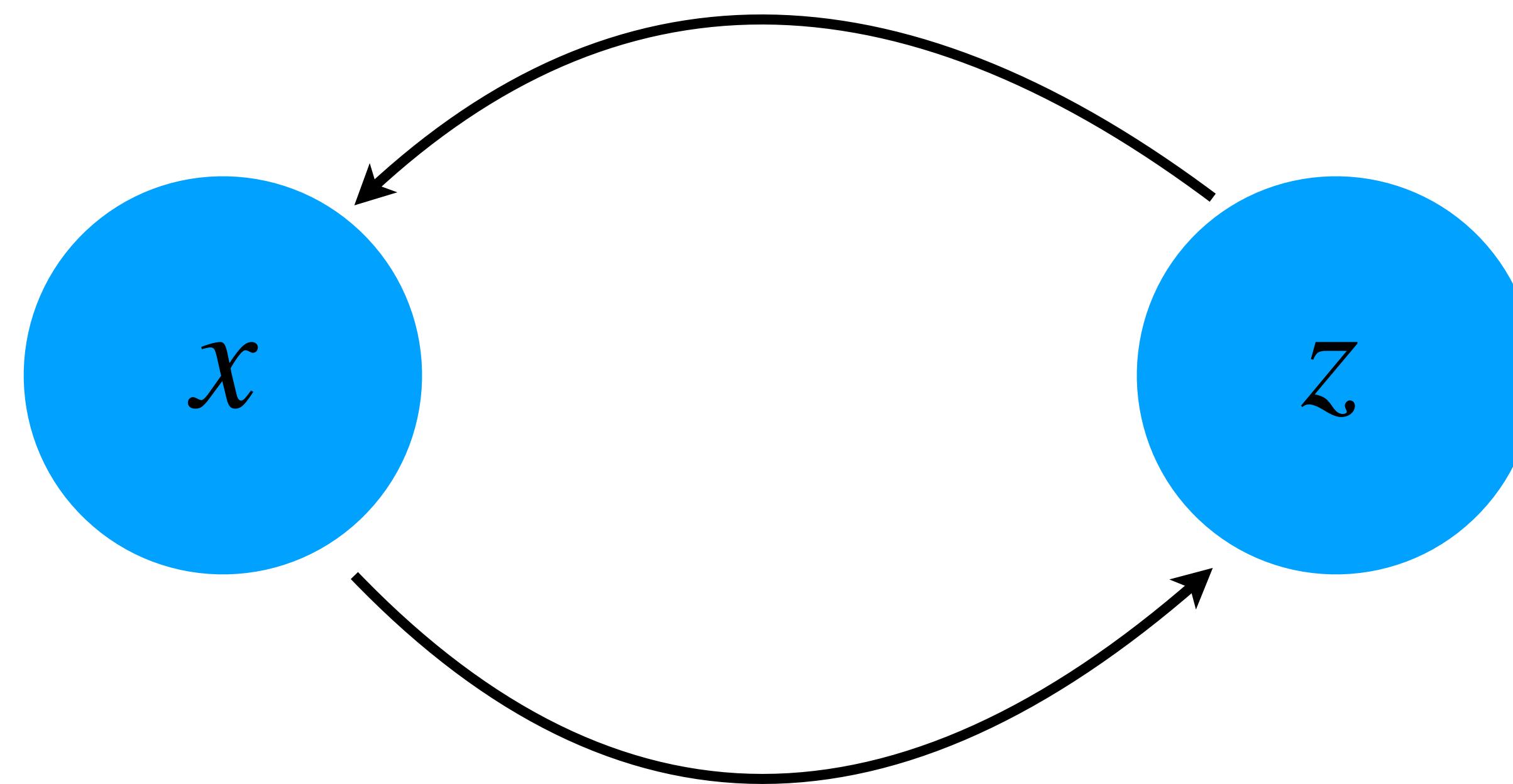
# Introduction : Variational Auto Encoders

- VAEs can also be conditioned on a label / class to generate specific output categories.



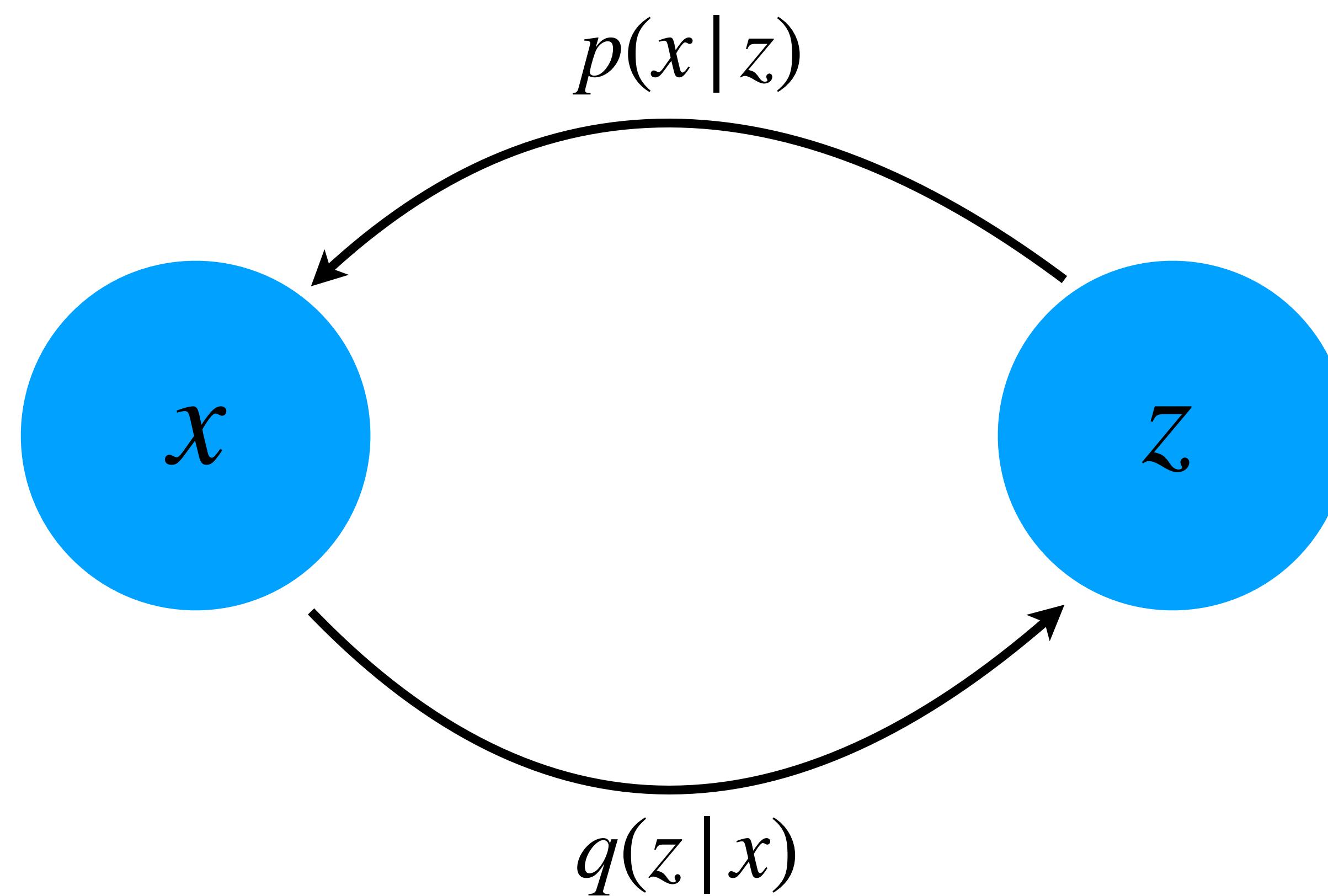
# Introduction : Variational Auto Encoders

- VAEs can be represented as a cyclic graph.



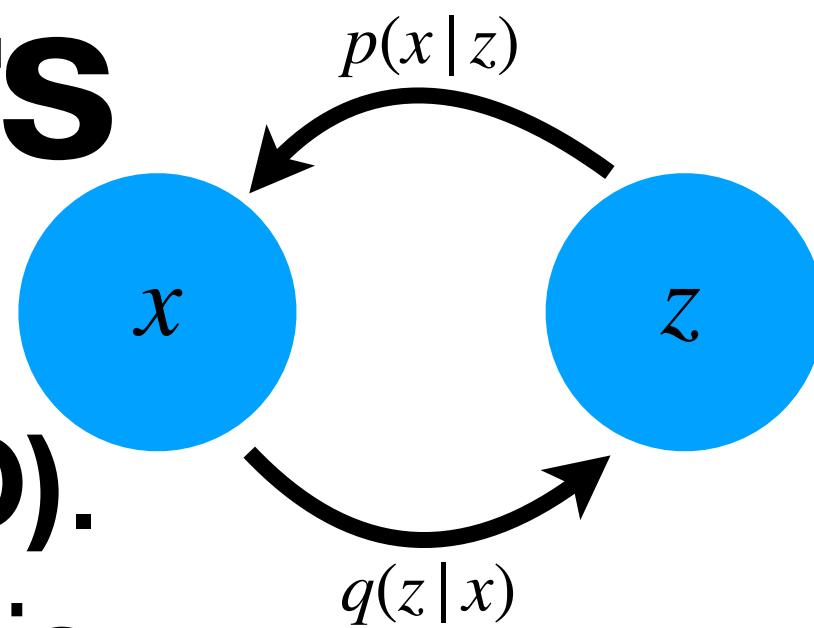
# Introduction : Variational Auto Encoders

- VAEs can be represented as a cyclic graph.



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

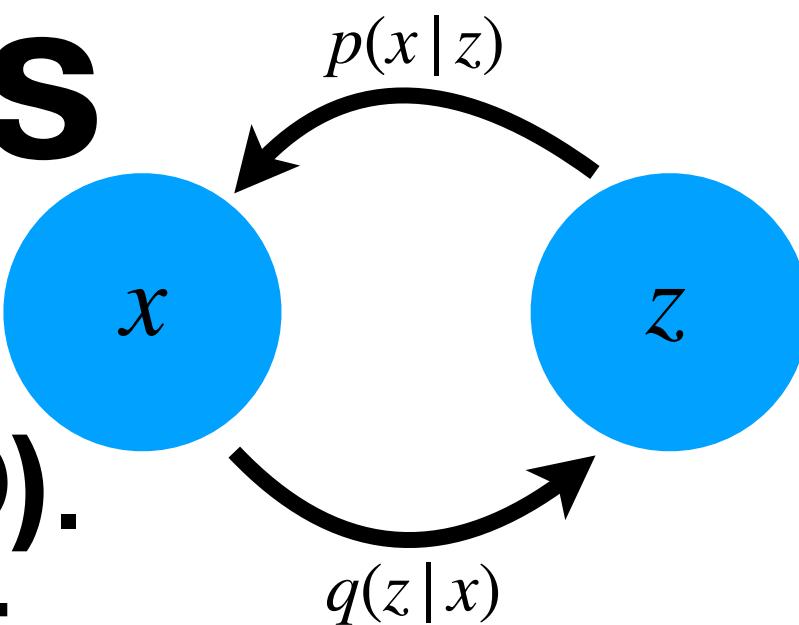


# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

Introducing additional variable



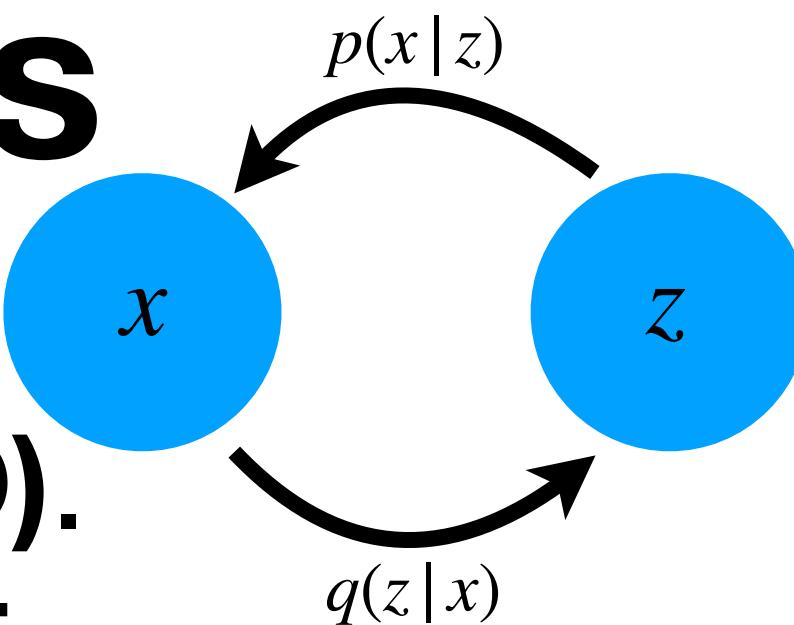
# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z} | \mathbf{x})}{q_{\phi}(\mathbf{z} | \mathbf{x})} d\mathbf{z}\end{aligned}$$

Introducing additional variable

Multiplying by 1



# Introduction : Variational Auto Encoders

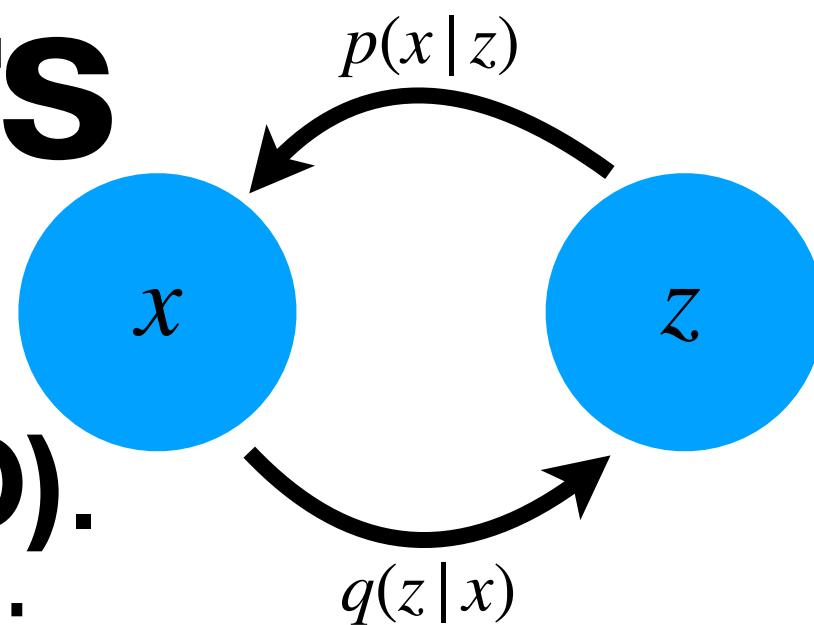
- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z} | \mathbf{x})}{q_{\phi}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \\ &= \log \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right]\end{aligned}$$

Introducing additional variable

Multiplying by 1

Definition of Expectation



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.

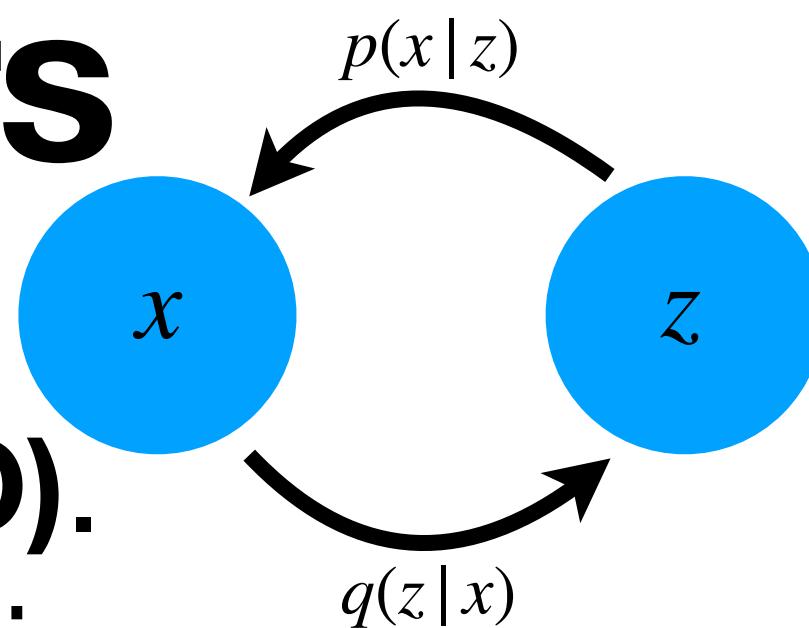
$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\ &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$

Introducing additional variable

Multiplying by 1

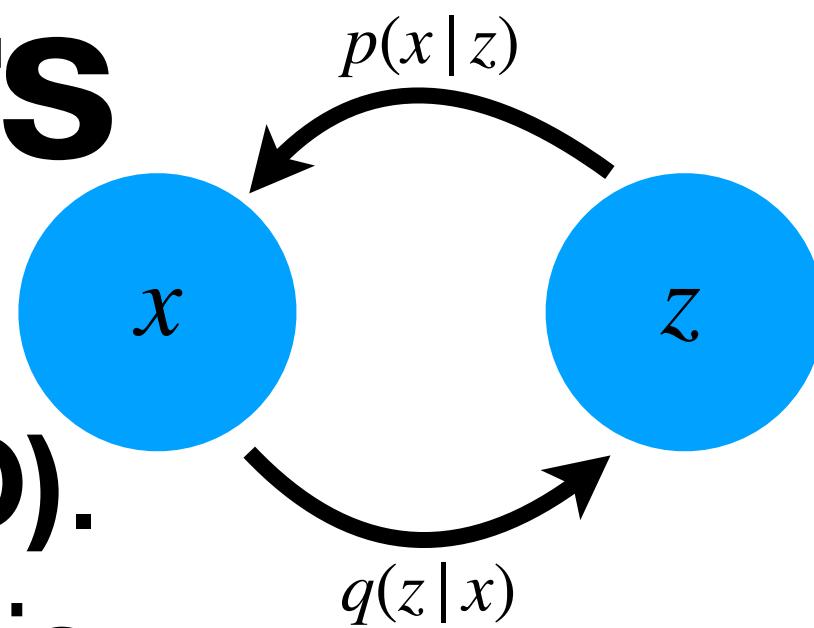
Definition of Expectation

Jensen's Inequality



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

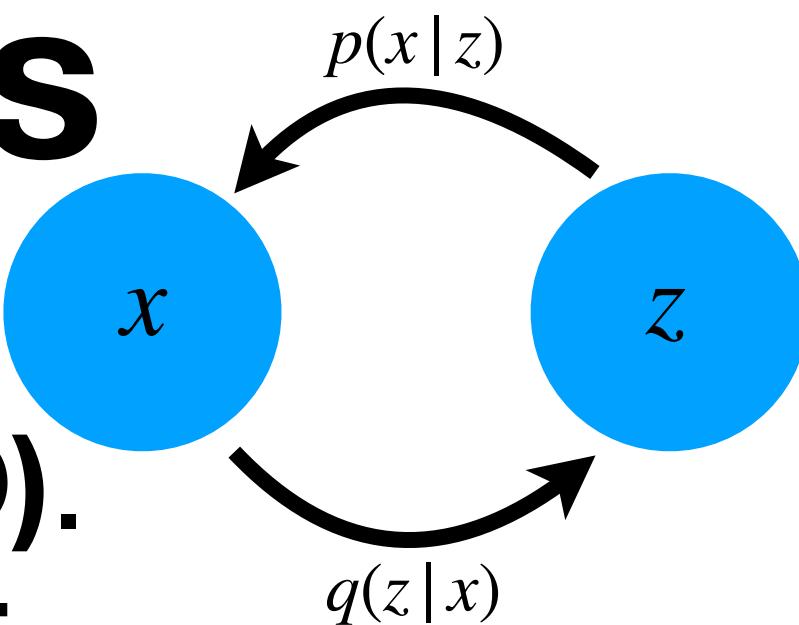


# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.

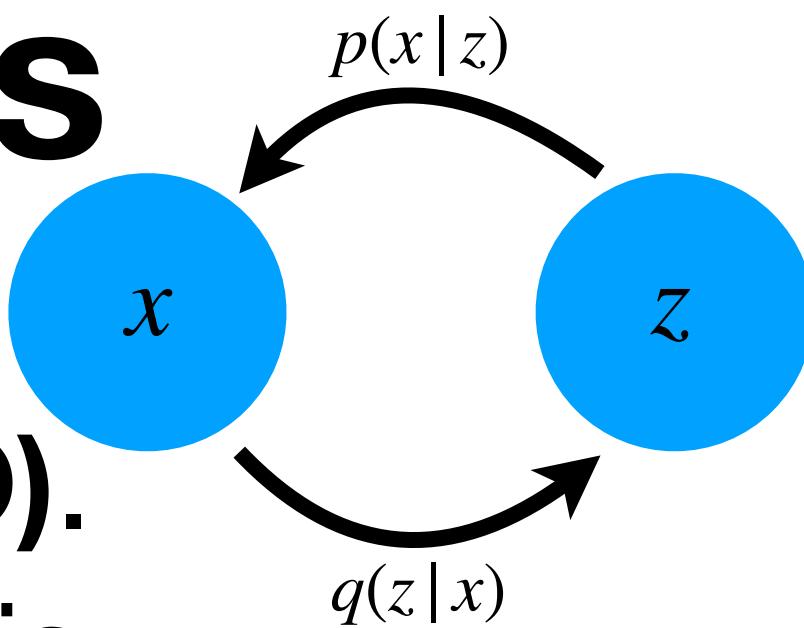
$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Chain Rule



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.



$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

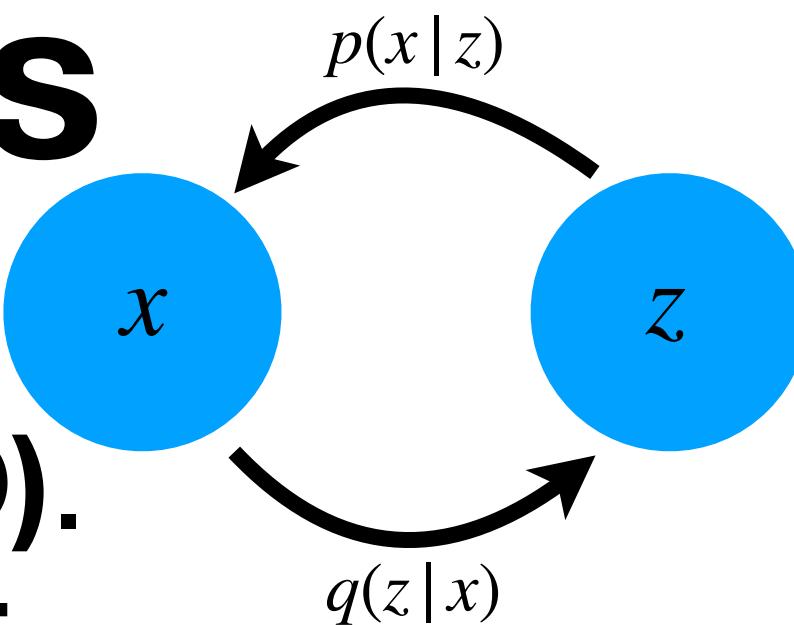
Chain Rule

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Split Expectation

# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(\mathbf{x})$  is difficult.



$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Chain Rule

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Split Expectation

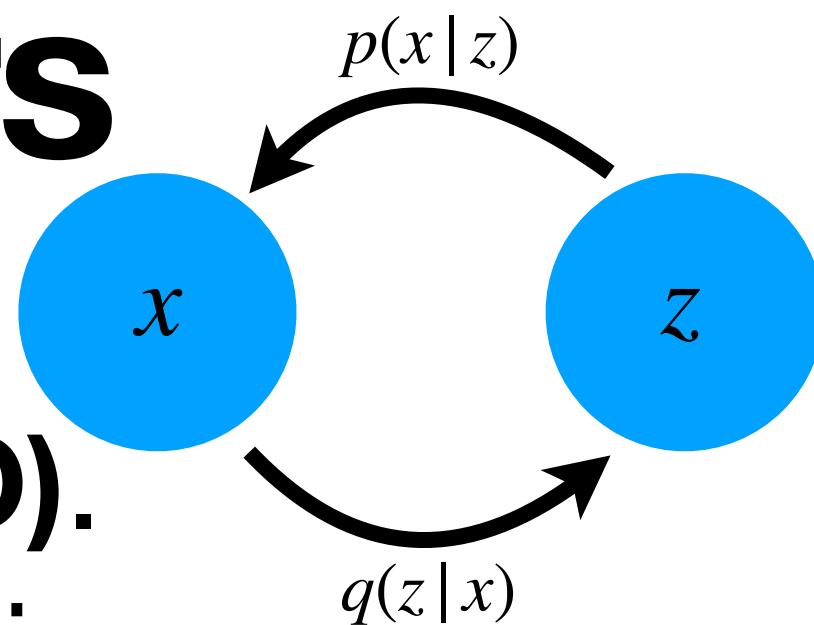
$$= \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{prior matching term}}$$

Definition of KL Divergence

# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

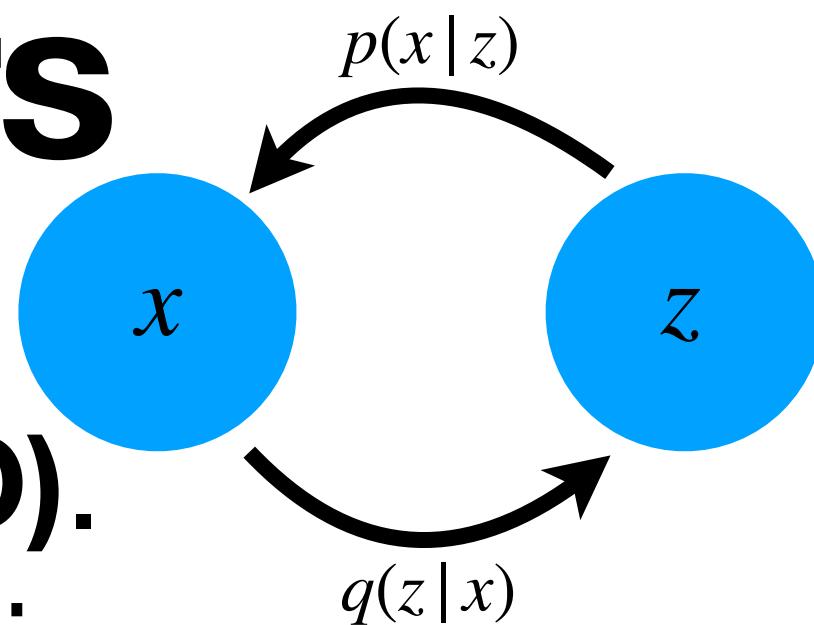
$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$

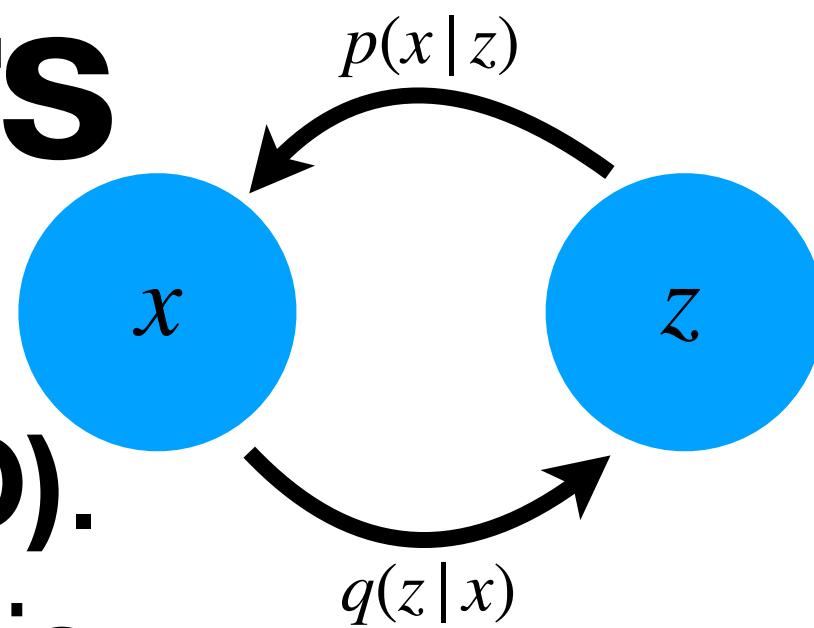


Reconstruction term that ensures the decoder regenerates the input.

# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

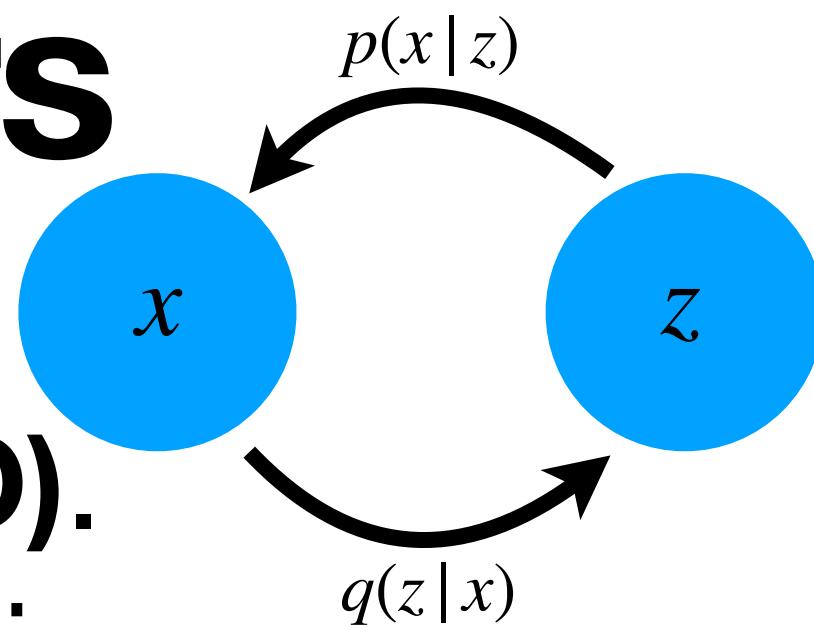
$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$



# Introduction : Variational Auto Encoders

- Training VAEs involve maximizing the **Evidence Lower Bound (ELBO)**.
  - This is because directly maximizing the likelihood of the data  $p(x)$  is difficult.

$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{\text{prior matching term}}$$



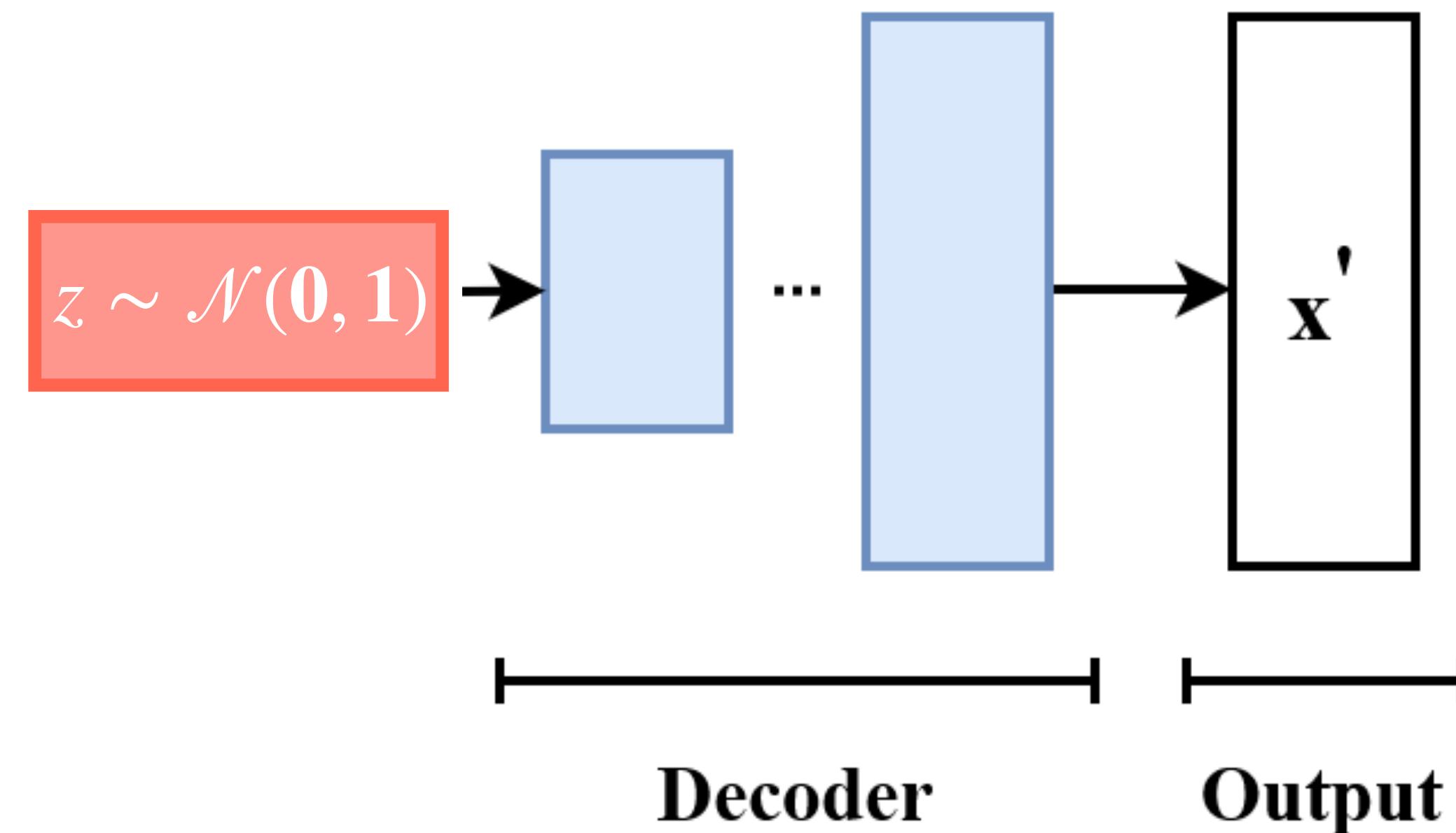
Regularization to ensure that the Gaussian parameters output by the encoder are close to a standard Gaussian distribution

# Introduction : Variational Auto Encoders

- As VAEs learn a **continuous latent representation**, it can be used for data generation.

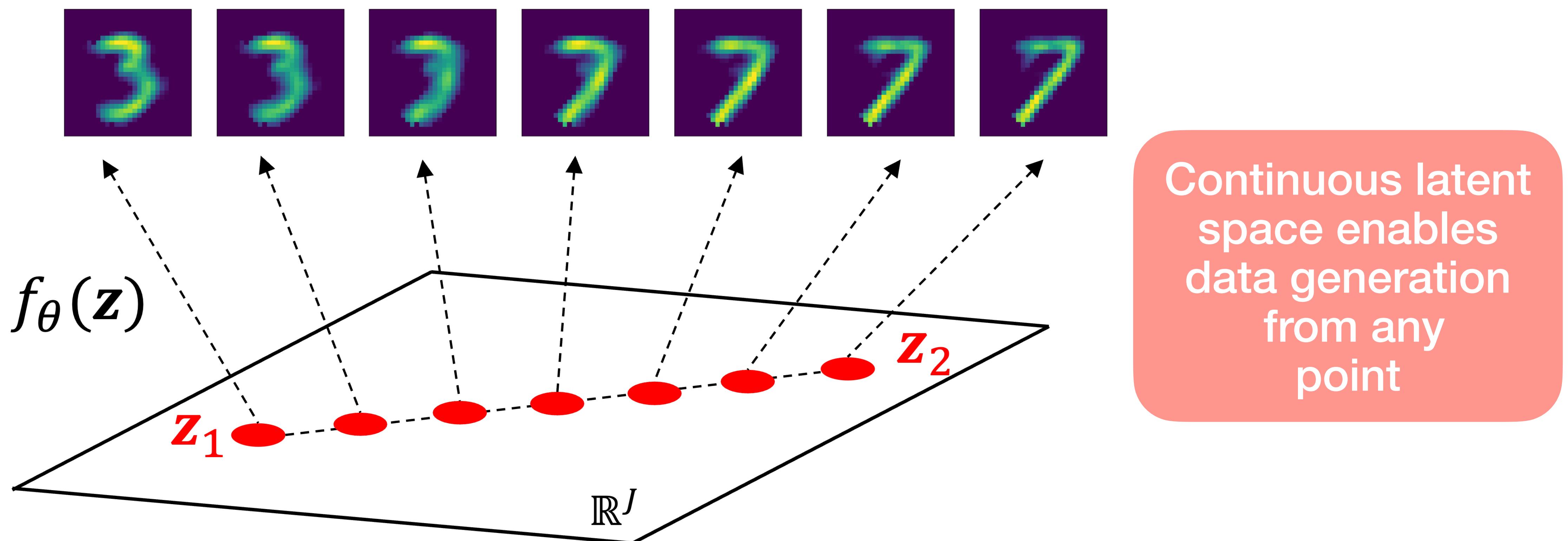
# Introduction : Variational Auto Encoders

- As VAEs learn a **continuous latent representation**, it can be used for data generation.



# Introduction : Variational Auto Encoders

- As VAEs learn a **continuous latent representation**, it can be used for data generation.

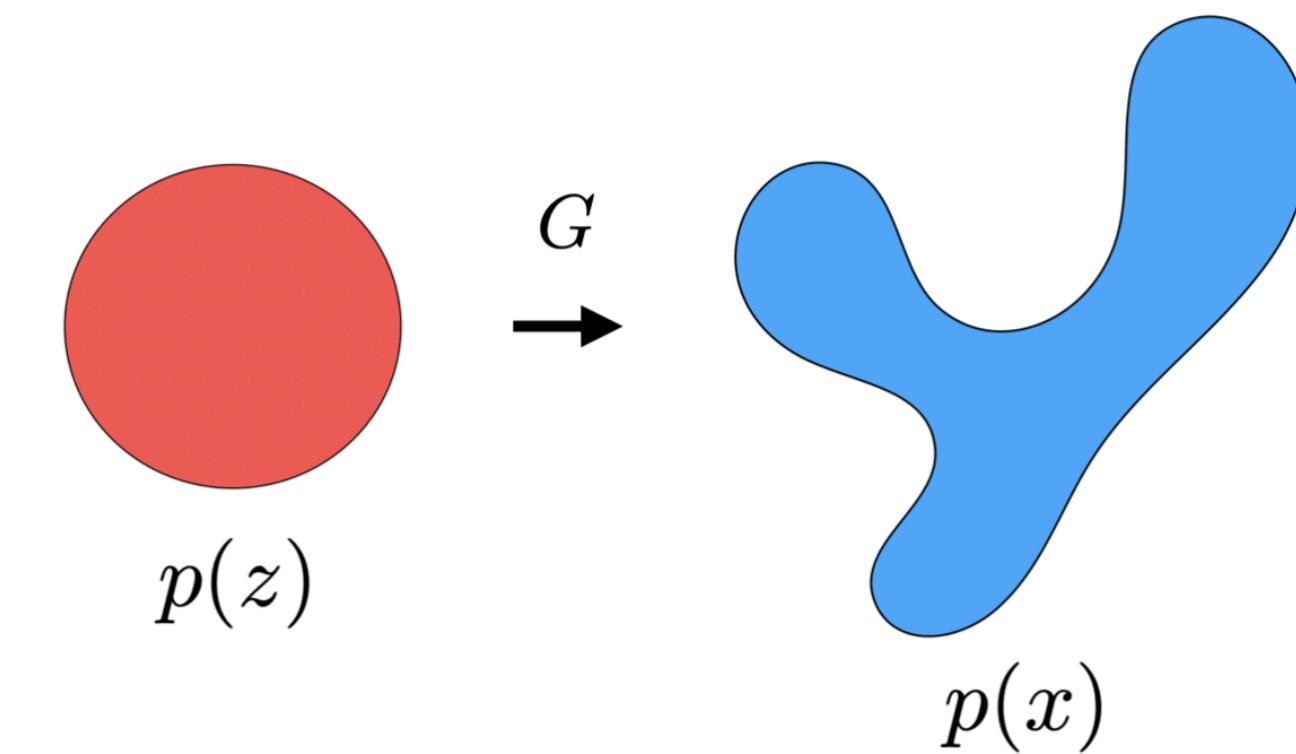


# Introduction : Variational Auto Encoders

- **VAE** effectively defines a mapping  $\mathbf{G}$  from a simple (Gaussian) latent distributions  $p(\mathbf{z})$  to an arbitrarily complex data distribution  $p(\mathbf{x})$

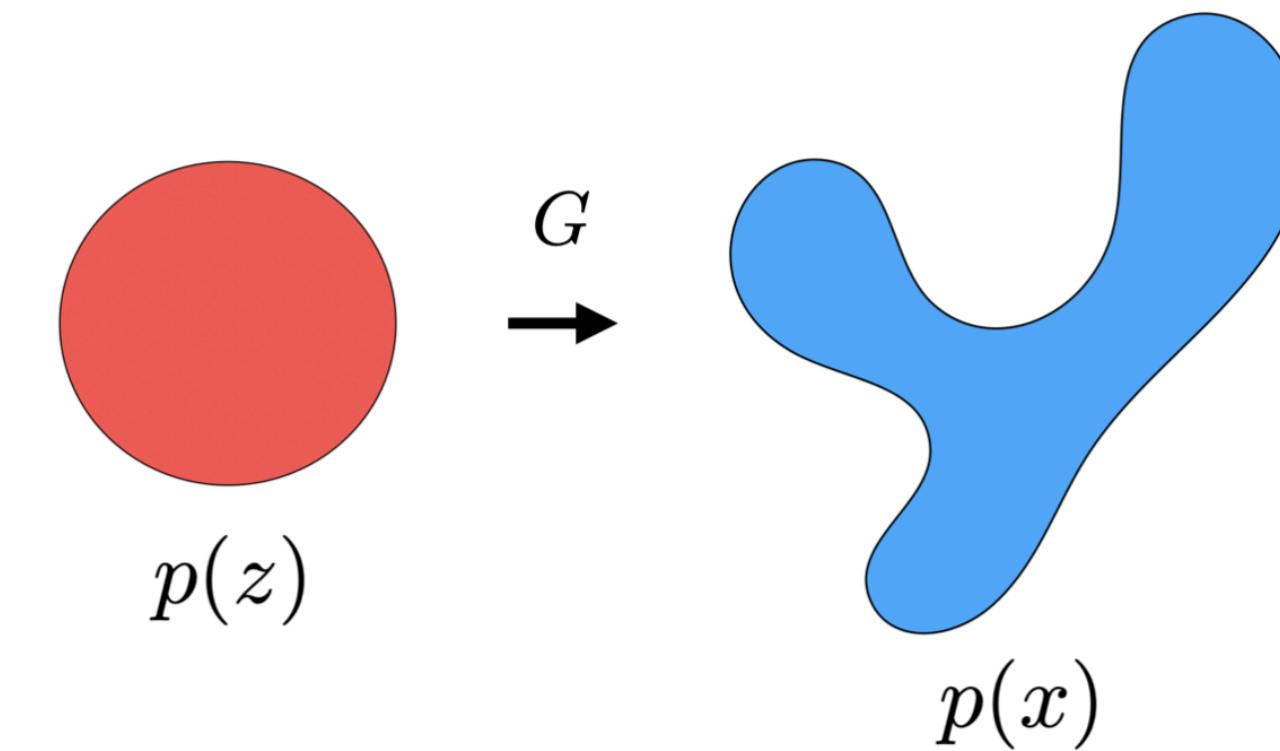
# Introduction : Variational Auto Encoders

- **VAE** effectively defines a mapping  $\mathbf{G}$  from a simple (Gaussian) latent distributions  $p(z)$  to an arbitrarily complex data distribution  $p(x)$



# Introduction : Variational Auto Encoders

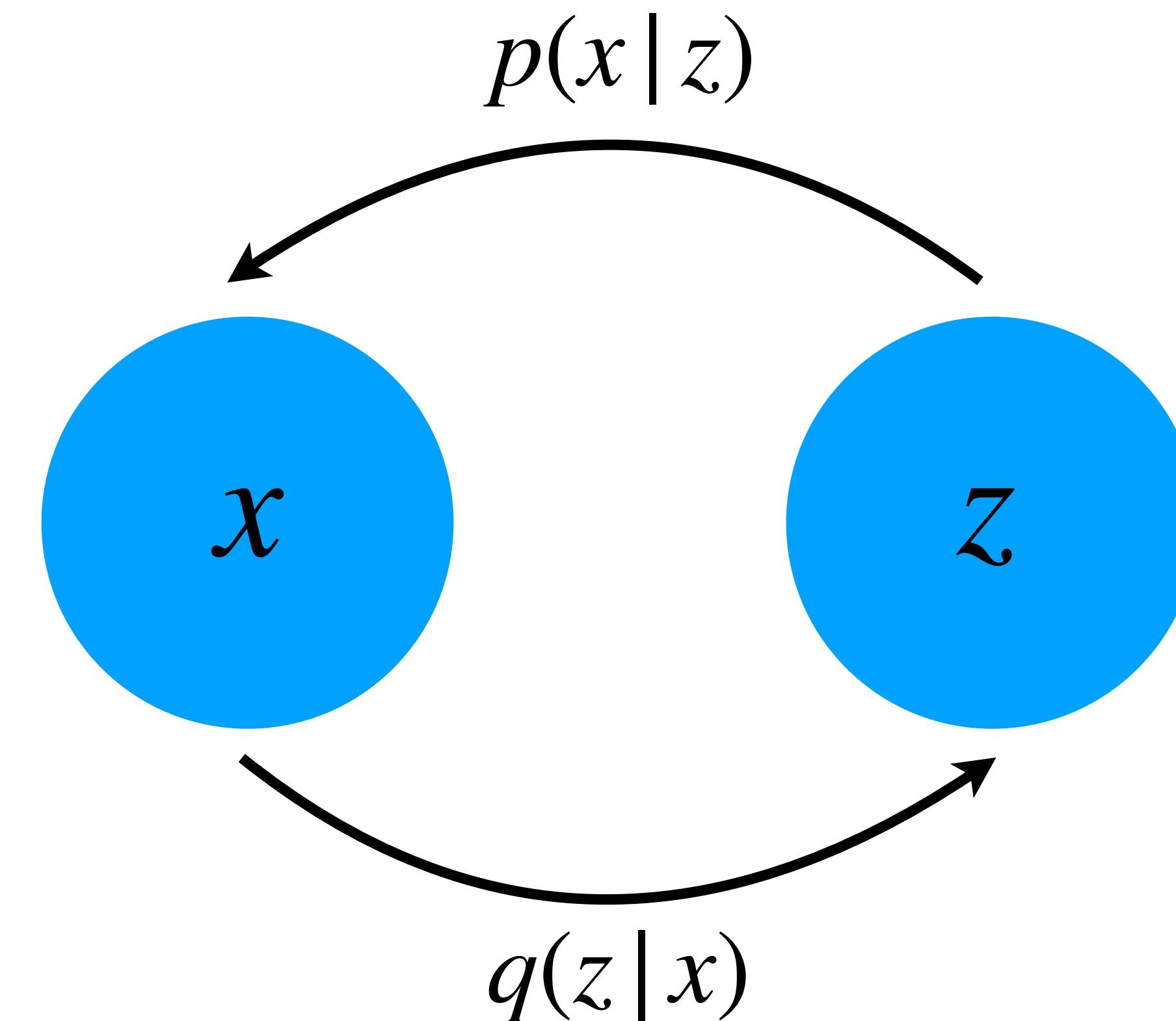
- **VAE** effectively defines a mapping  $\mathbf{G}$  from a simple (Gaussian) latent distributions  $p(z)$  to an arbitrarily complex data distribution  $p(x)$



- **Diffusion model** does something very similar, but through a series of parametrized transformations, ensuring that each transformation can be simpler and more complex data distributions can be handled more easily as a result

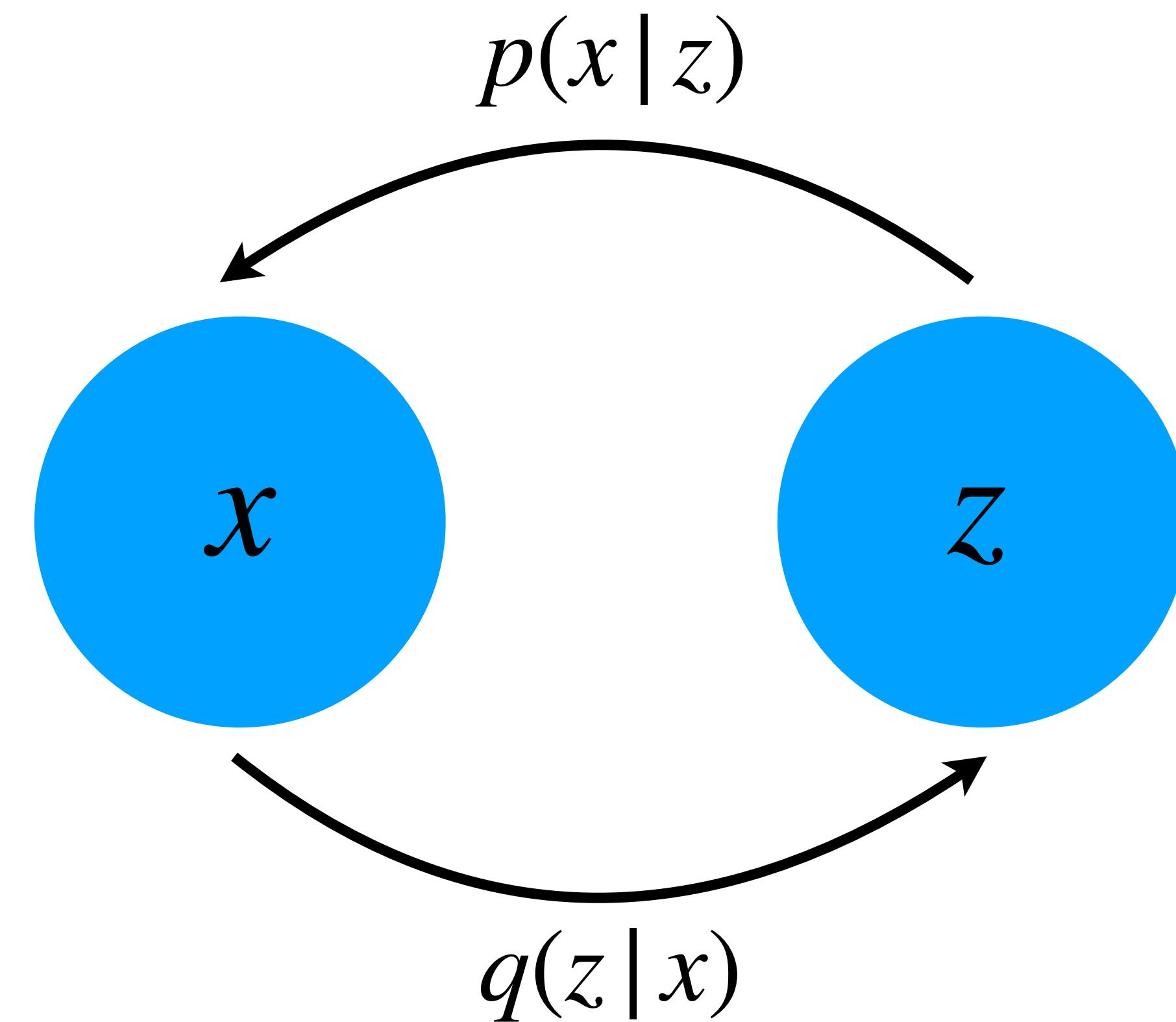
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.



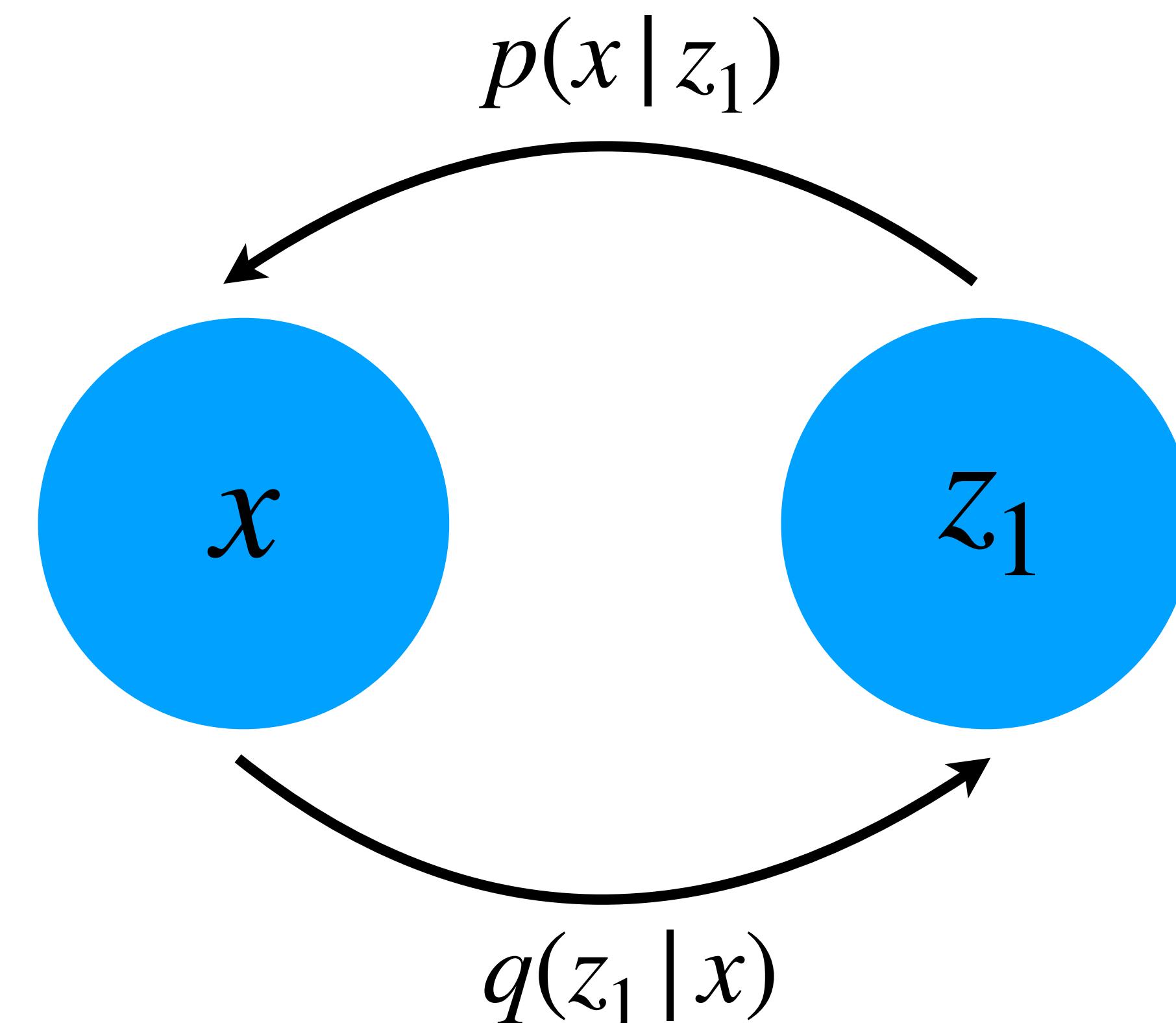
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.



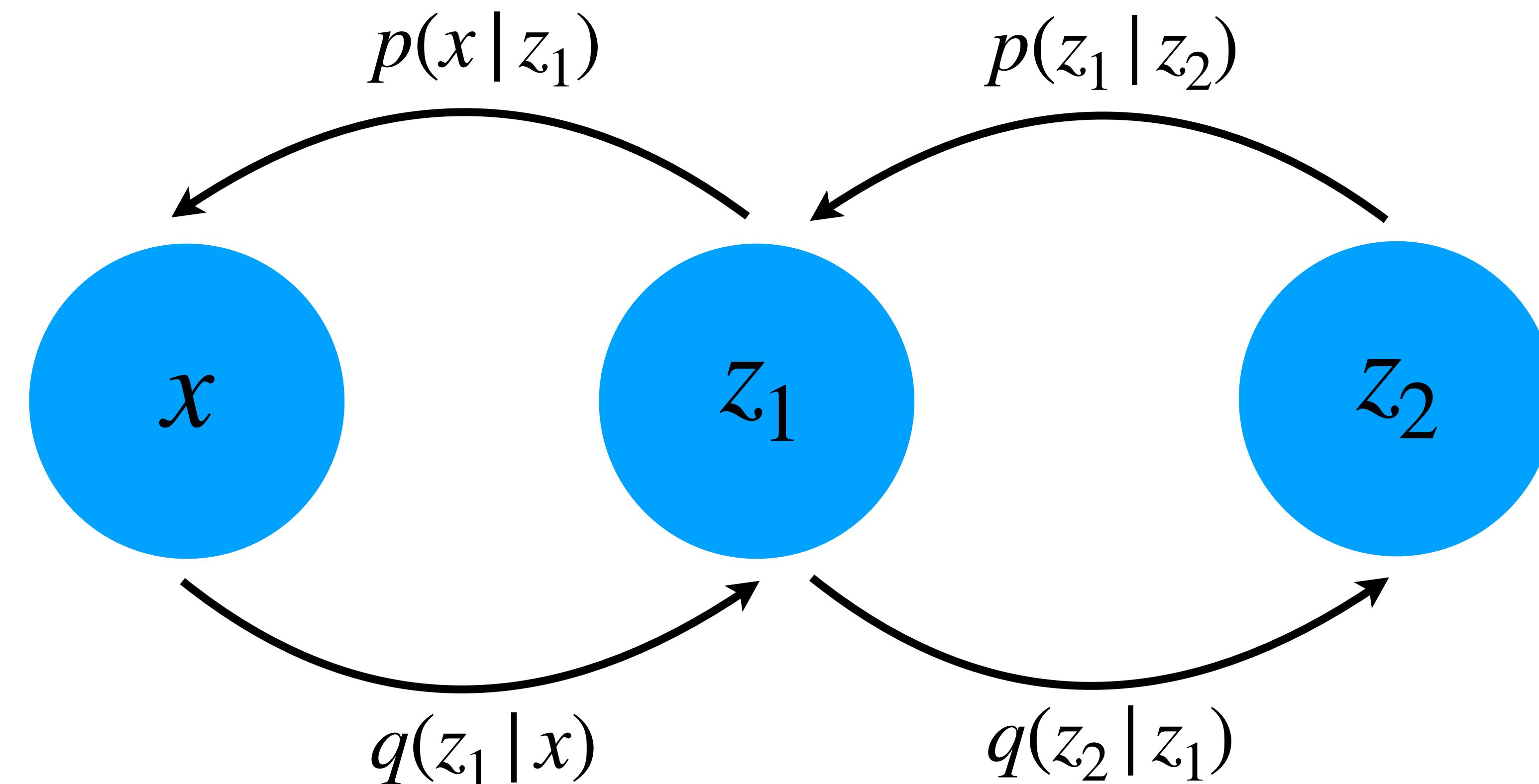
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.



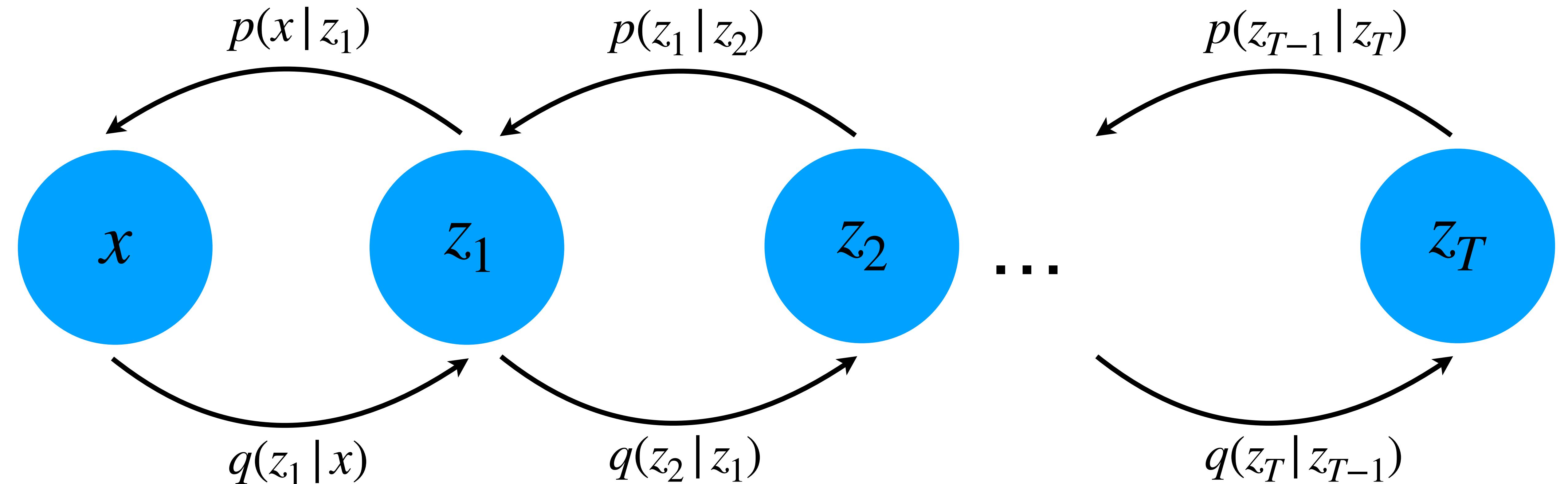
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.



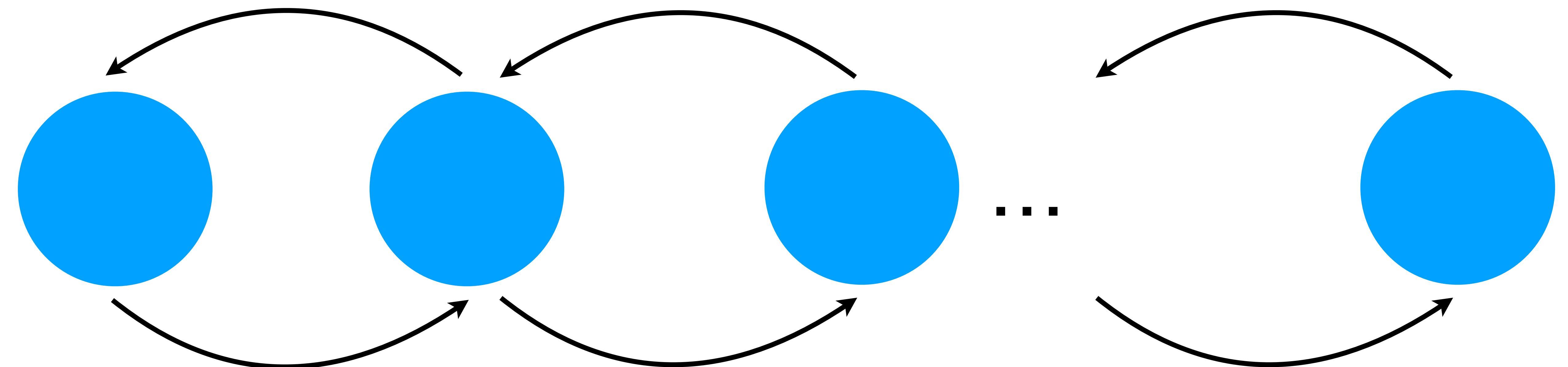
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.



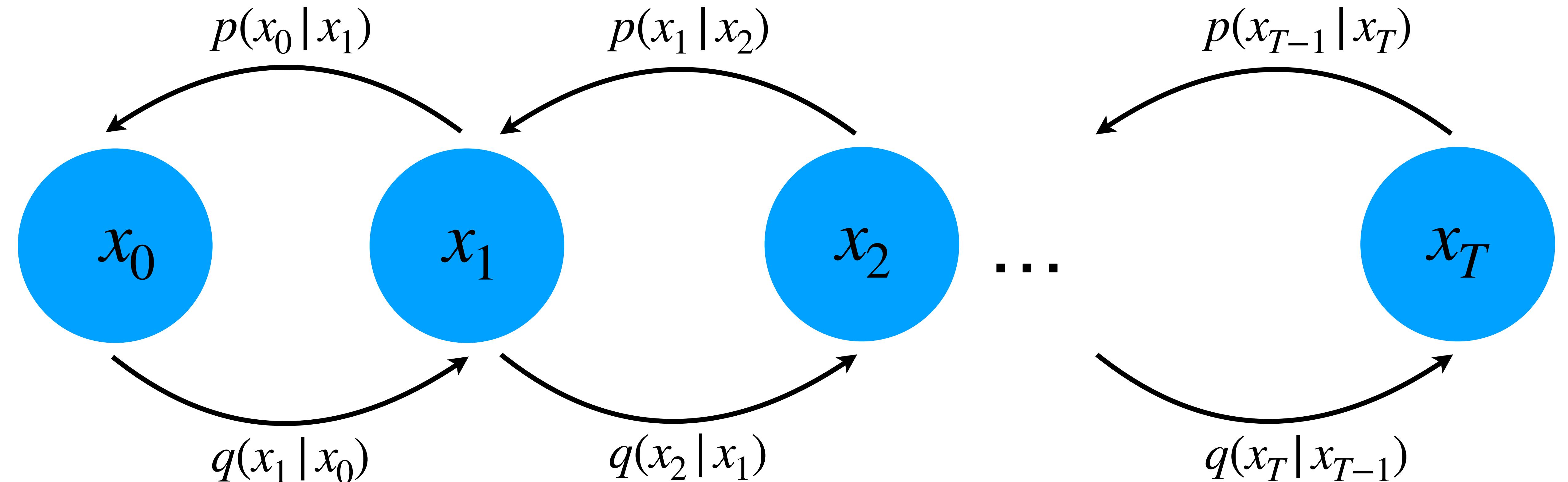
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - All latent dimensions  $\mathbf{z}$  are **the same** as the data dimension  $\mathbf{x}$



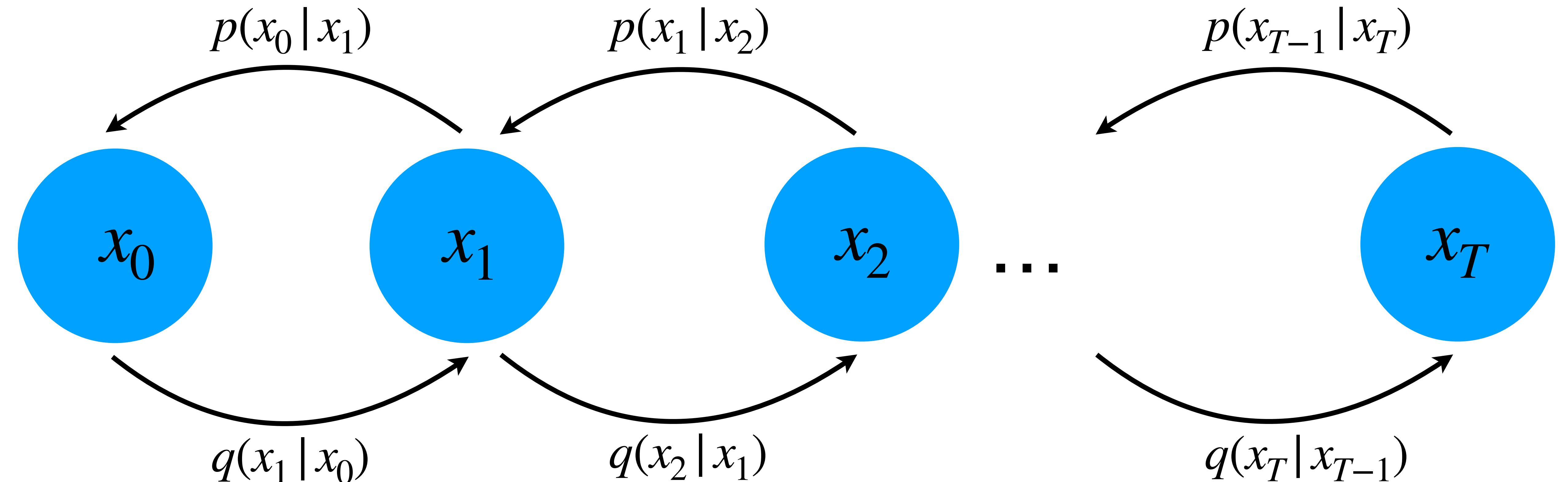
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - All latent dimensions  $\mathbf{z}$  are **the same** as the data dimension  $\mathbf{x}$



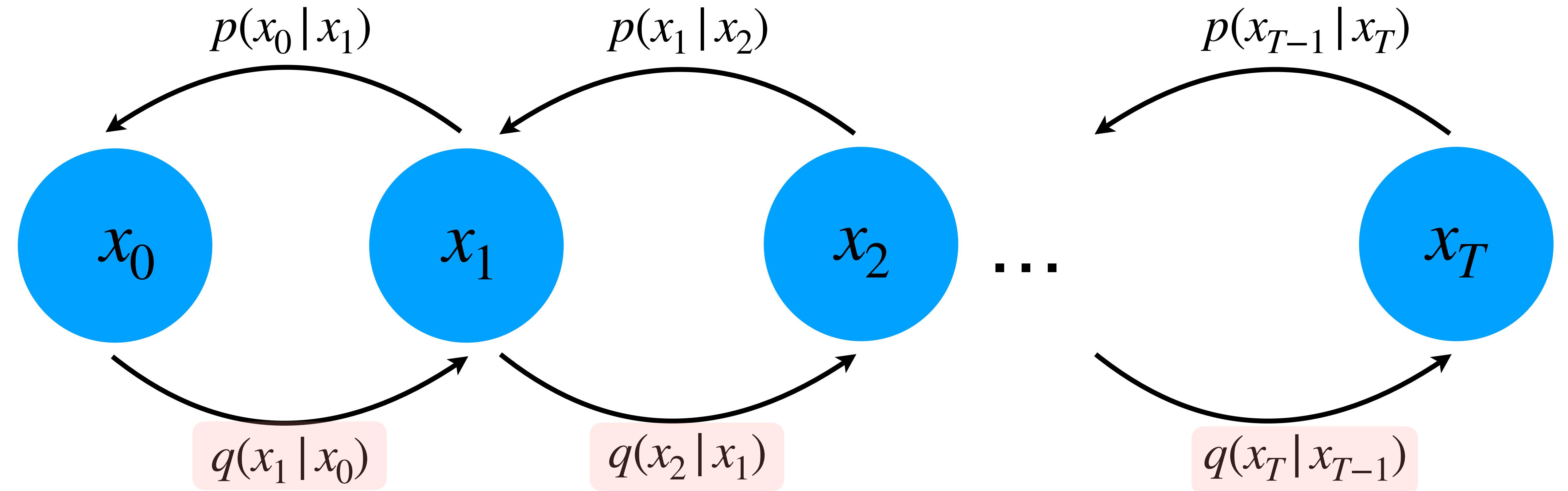
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - Encoder at each step is **not learned**; predefined as a Gaussian distribution centered around the output of the previous step.



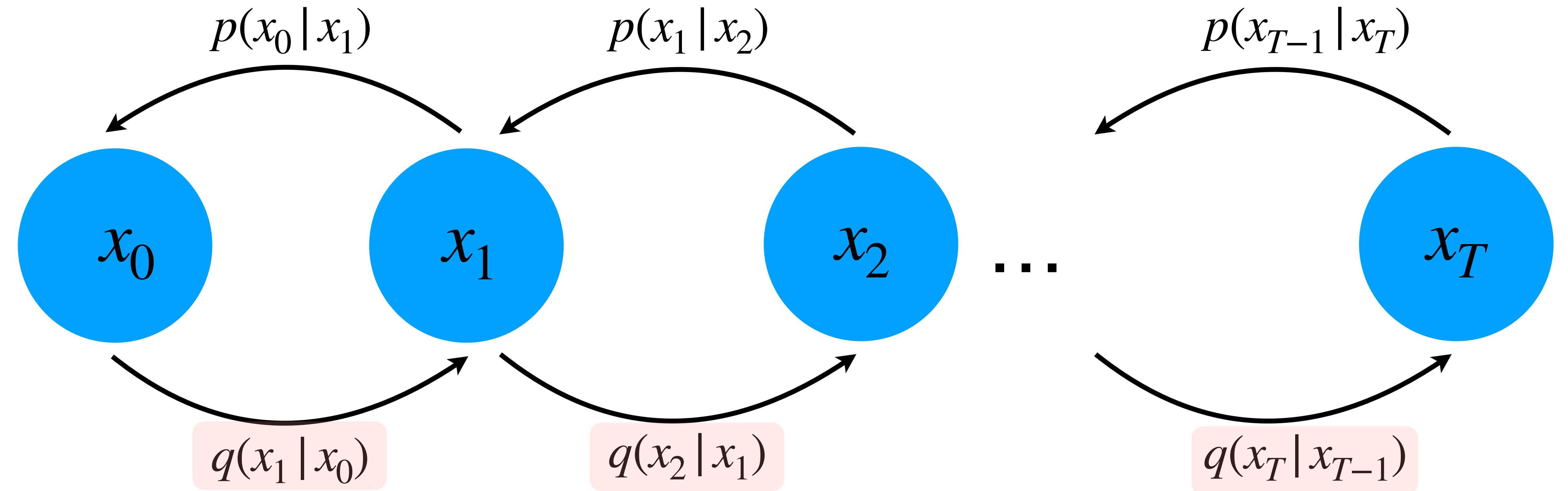
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - Encoder at each step is **not learned**; predefined as a Gaussian distribution centered around the output of the previous step.



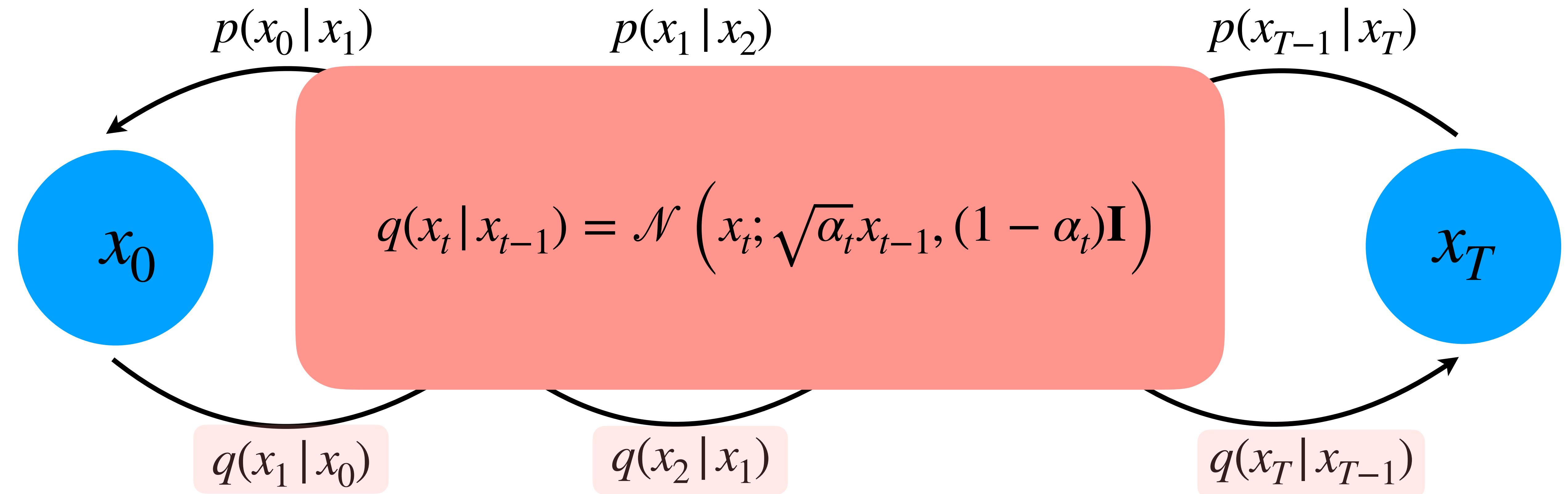
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - Encoder at each step is **not learned**; predefined as a Gaussian distribution centered around the output of the previous step.



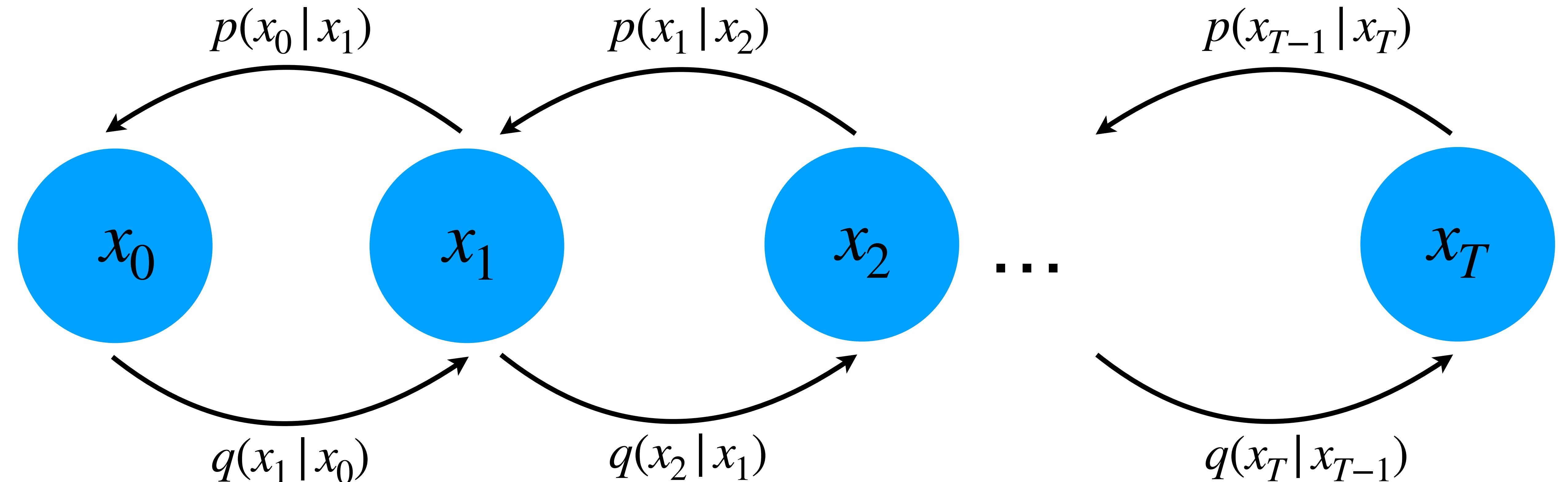
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - Encoder at each step is **not learned**; predefined as a Gaussian distribution centered around the output of the previous step.



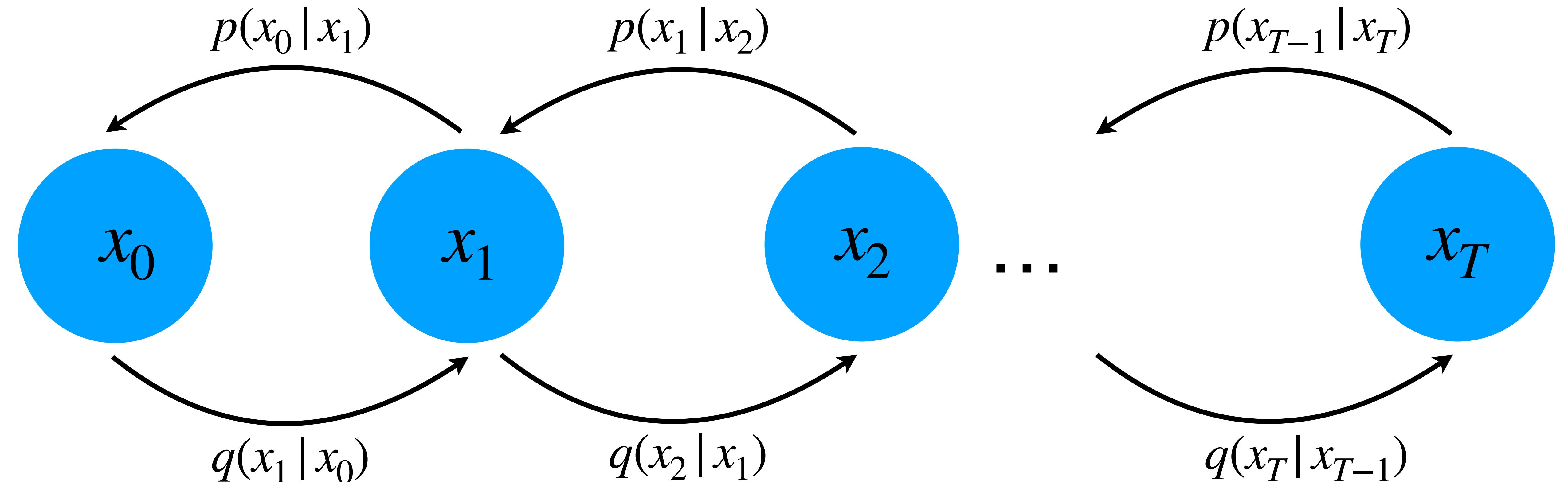
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - The Gaussian parameters of the encoder are set in a way such that the distribution at the final timestep  $\mathbf{T}$  is a standard Gaussian



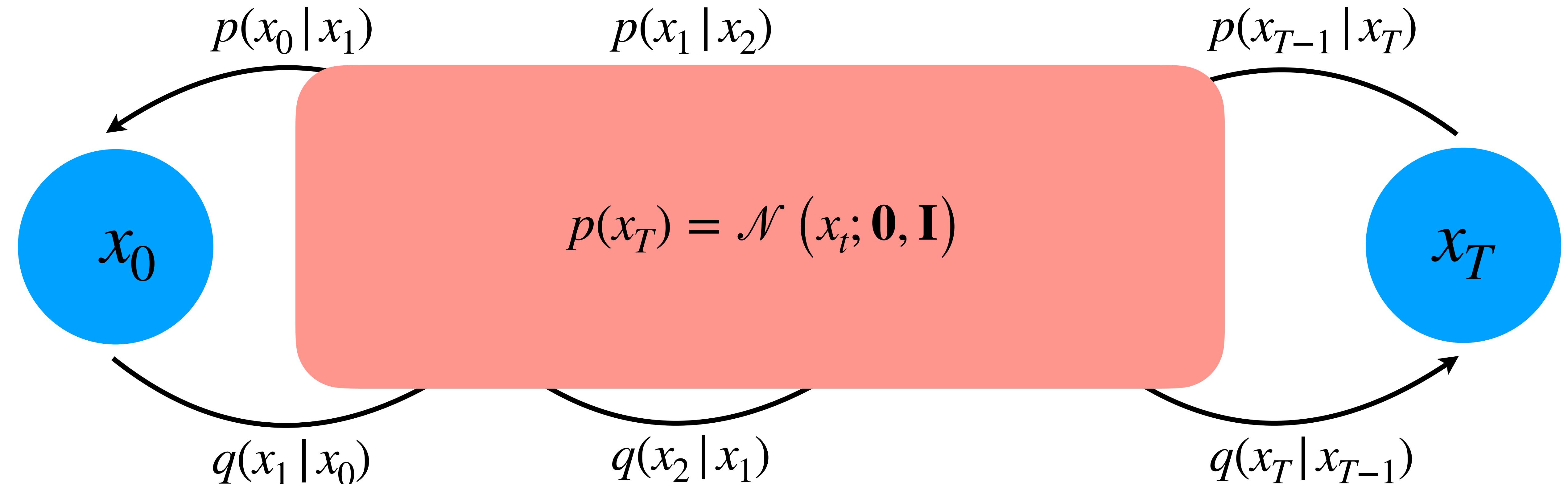
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - The Gaussian parameters of the encoder are set in a way such that the distribution at the final timestep  $\mathbf{T}$  is a standard Gaussian



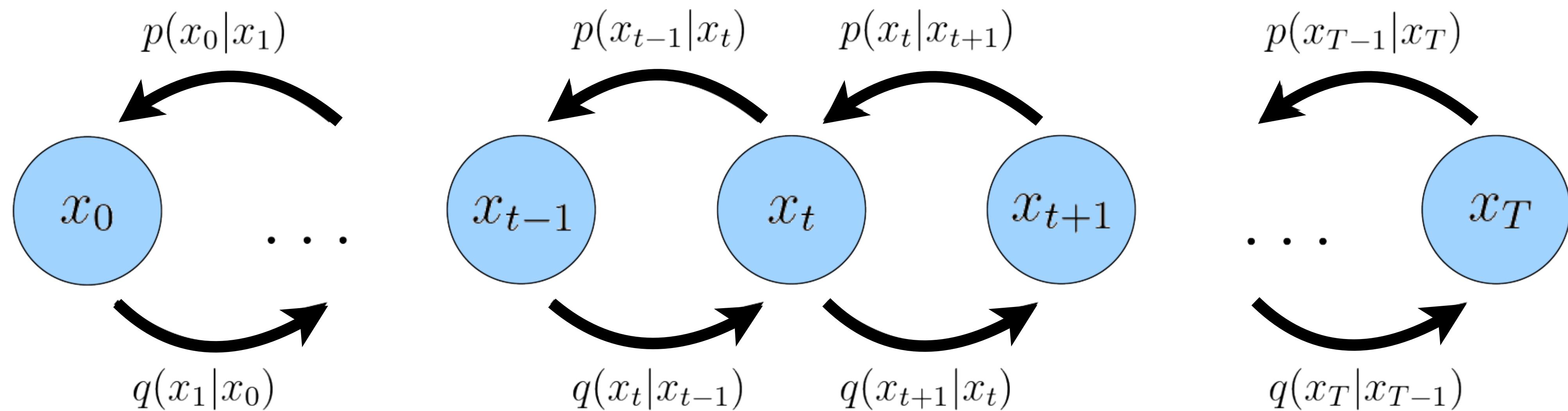
# Variational Diffusion Models

- The easiest way to think about a diffusion model is as a Markovian hierarchical variational autoencoder with some restrictions.
  - The Gaussian parameters of the encoder are set in a way such that the distribution at the final timestep  $\mathbf{T}$  is a standard Gaussian



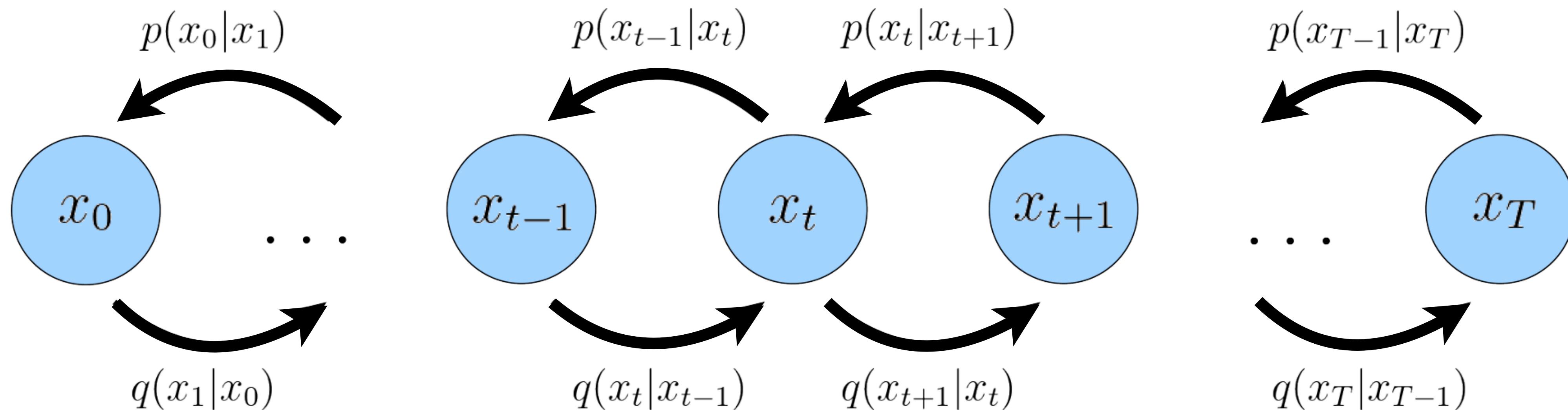
# Variational Diffusion Models

- Diffusion models can be split into two processes.



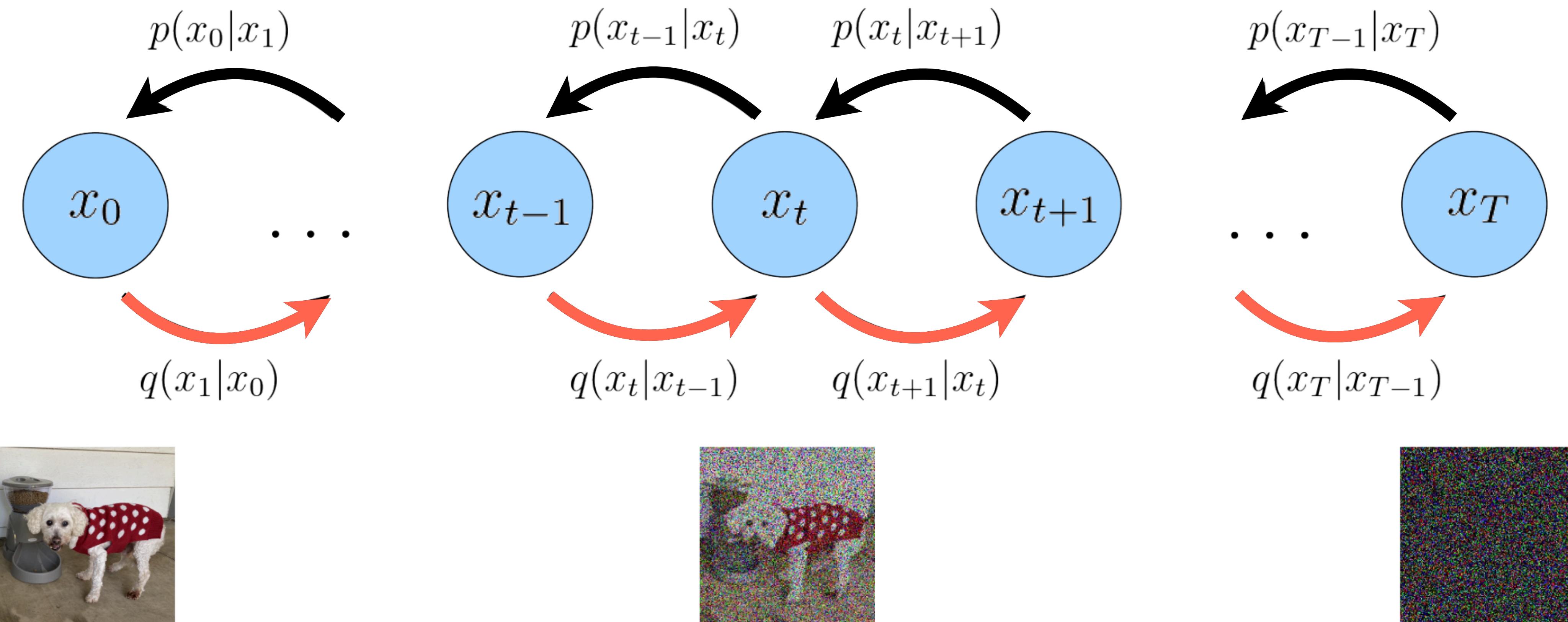
# Variational Diffusion Models

- Diffusion models can be split into two processes.
  - Forward Process:** Progressively corrupt an input by adding Gaussian noise until it becomes completely identical to pure Gaussian noise. **Not Learned.**



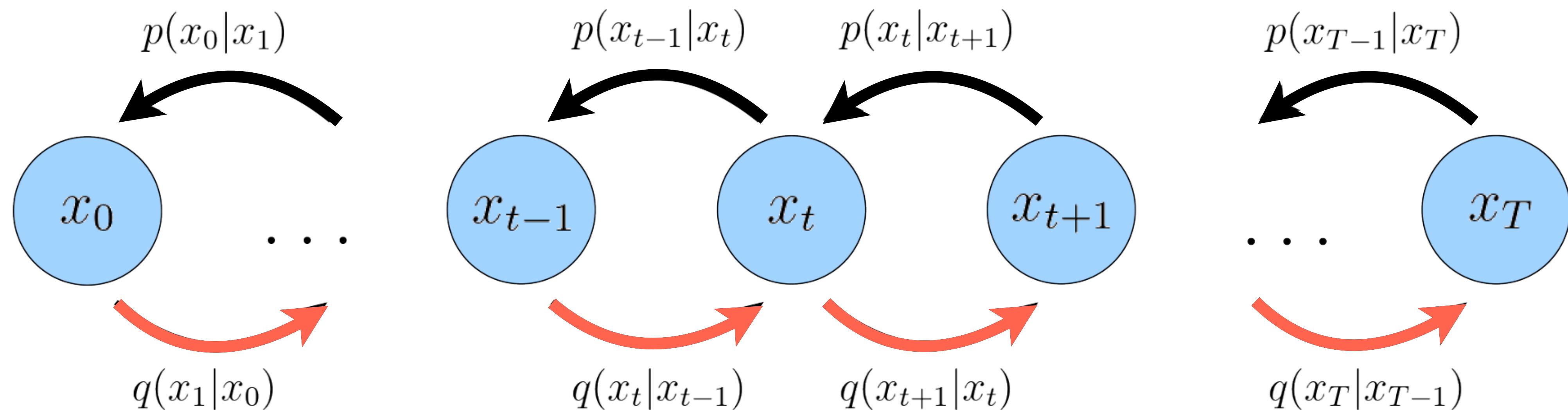
# Variational Diffusion Models

- Diffusion models can be split into two processes.
  - Forward Process:** Progressively corrupt an input by adding Gaussian noise until it becomes completely identical to pure Gaussian noise. **Not Learned.**



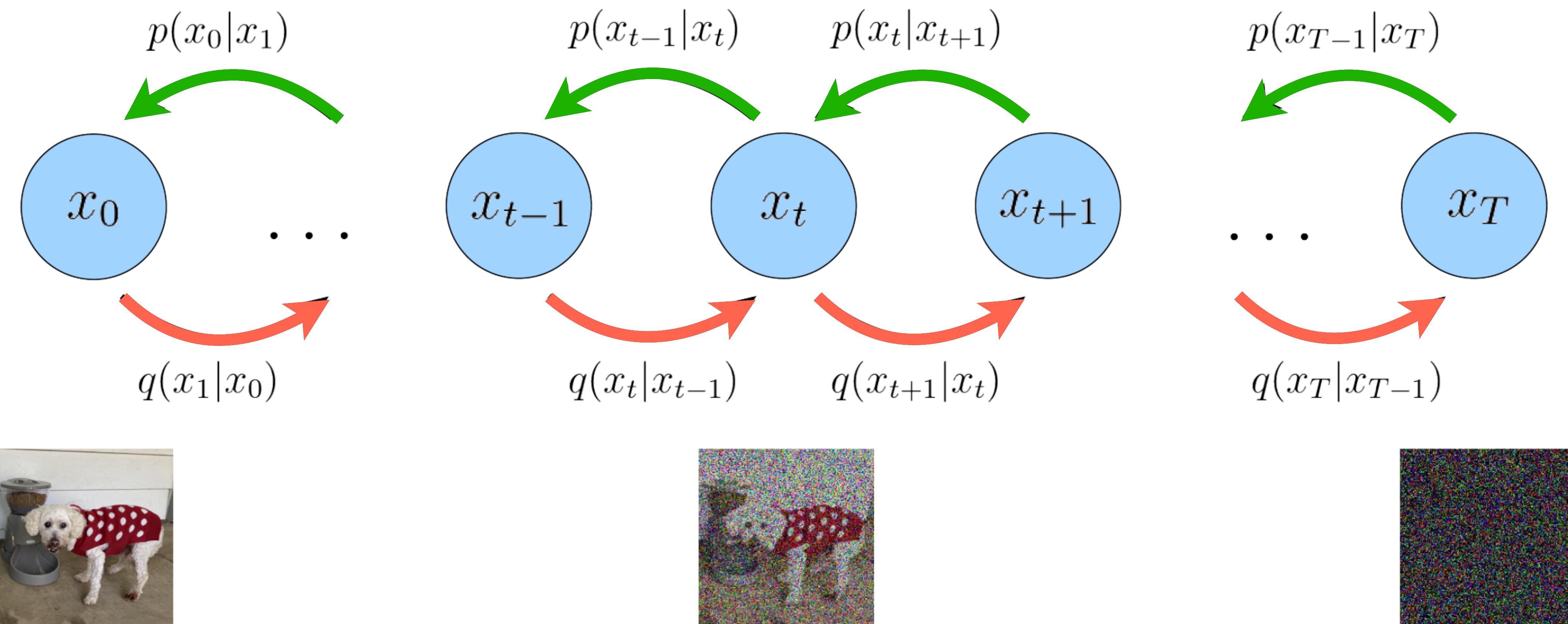
# Variational Diffusion Models

- Diffusion models can be split into two processes.
  - Reverse Process:** Noise is transformed back into a sample from the target distribution. **Learned.**

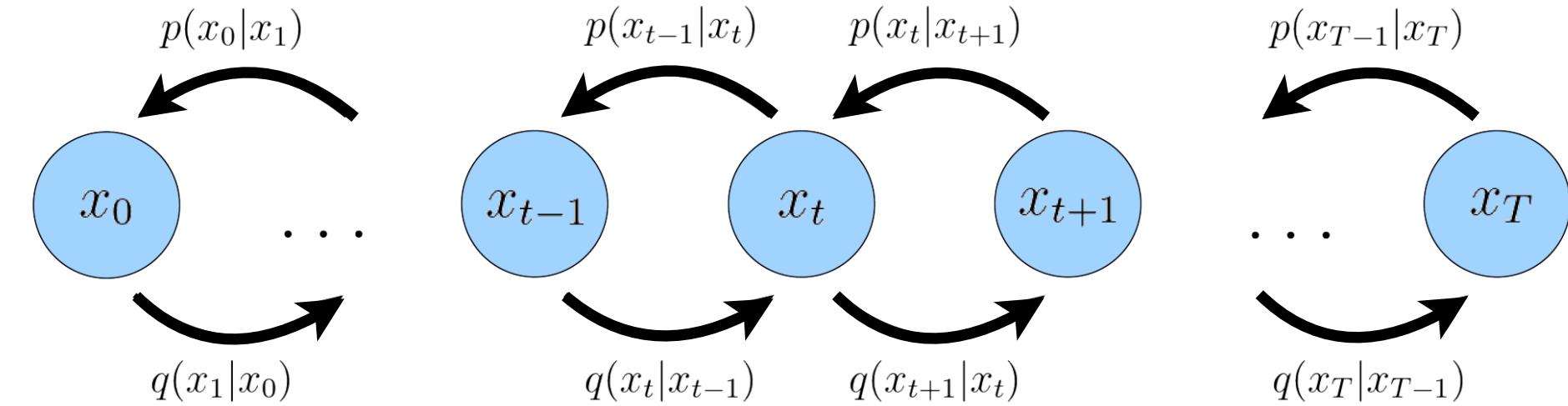


# Variational Diffusion Models

- Diffusion models can be split into two processes.
  - Reverse Process:** Noise is transformed back into a sample from the target distribution. **Learned.**



# Variational Diffusion Models: Training

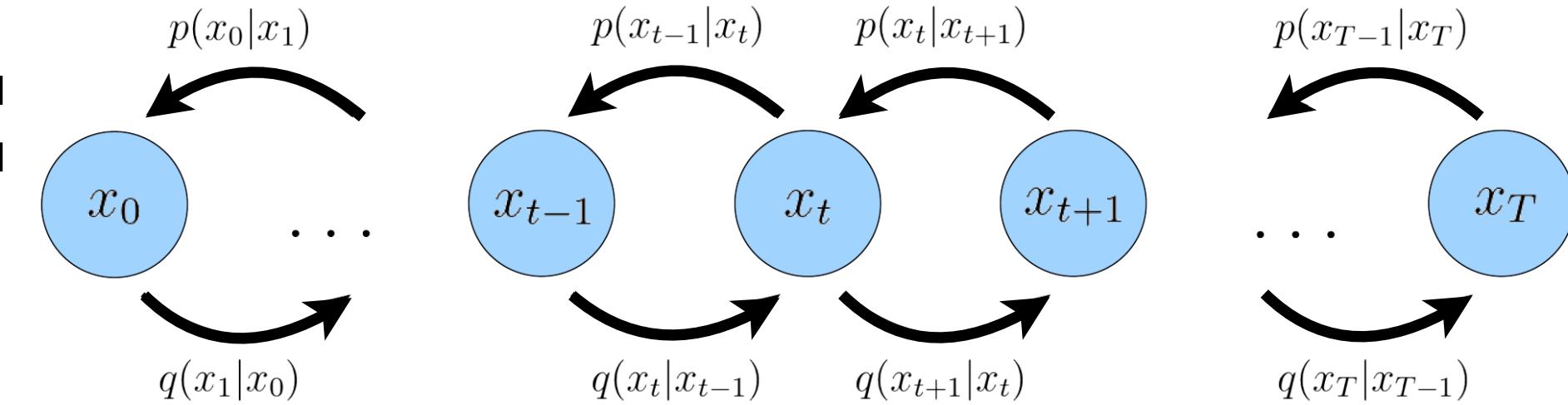


- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

# Variational Diffusion Models: Training

- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

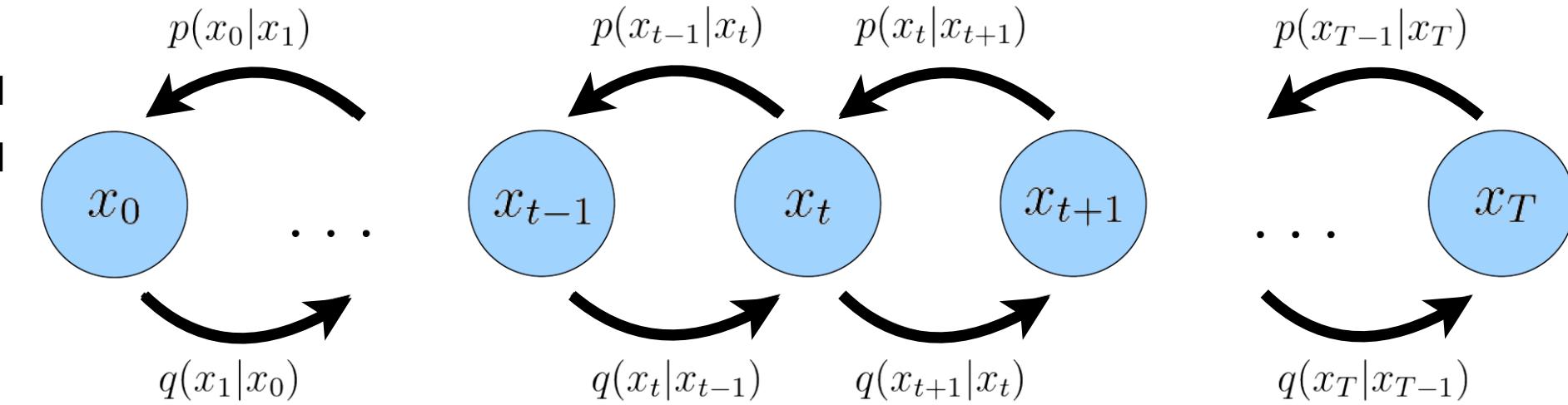
$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] - D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{reconstruction term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$



# Variational Diffusion Models: Training

- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

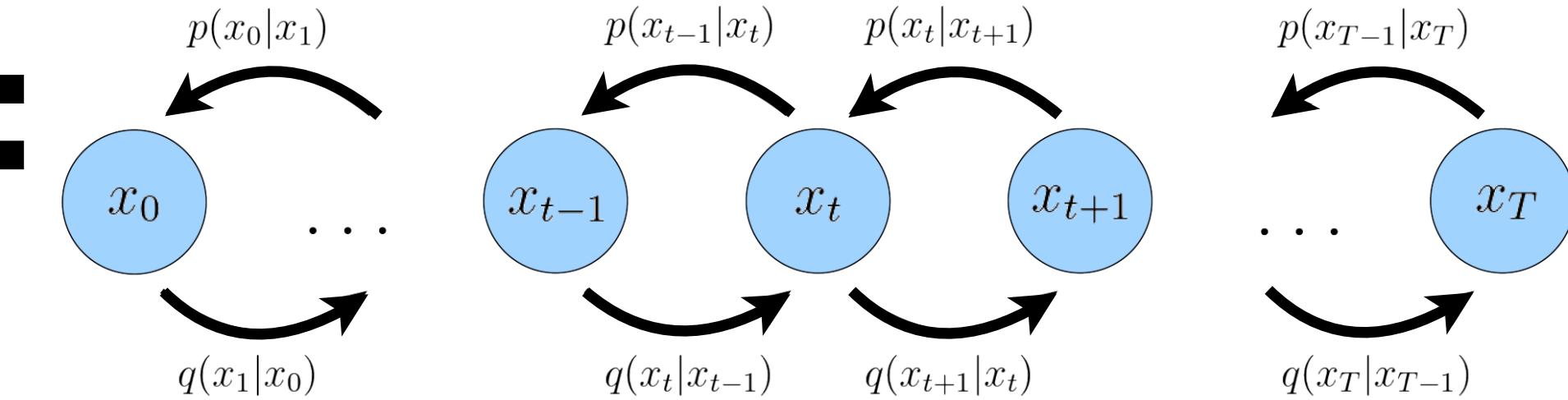
$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] - D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{reconstruction term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$



# Variational Diffusion Models: Training

- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))] }_{\text{denoising matching term}}$$

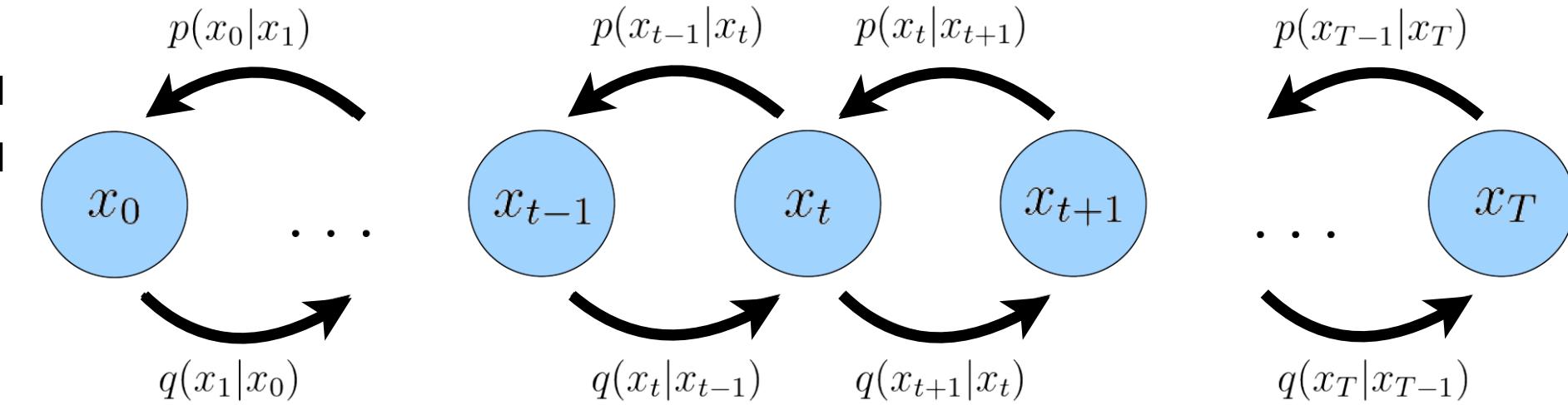


Reconstruction term that is similar to the ELBO in a VAE.

# Variational Diffusion Models: Training

- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$



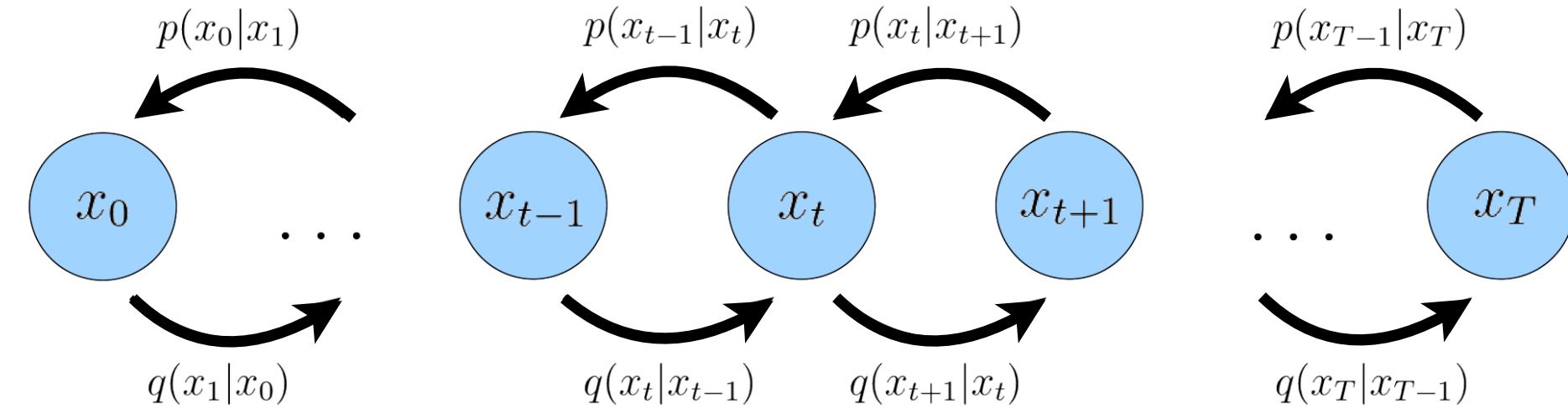
Represents how close the distribution of the final noised input  $\mathbf{x}_T$  is to a standard Gaussian prior. This term has no trainable parameters and is ignored during training.

# Variational Diffusion Models: Training

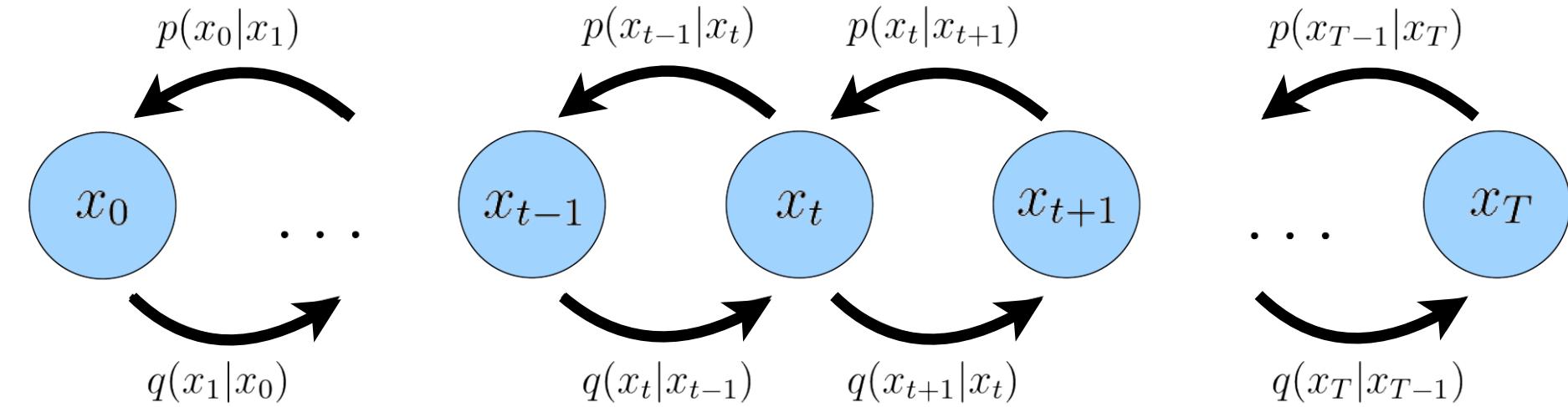
- Similar to VAEs, training diffusion models involves maximizing the **ELBO**

$$\underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)] - D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{reconstruction term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

The decoder  $p_{\theta}(x_{t-1} | x_t)$  is learned to be an approximation to the ground truth denoising transition  $q(x_{t-1} | x_t, x_0)$

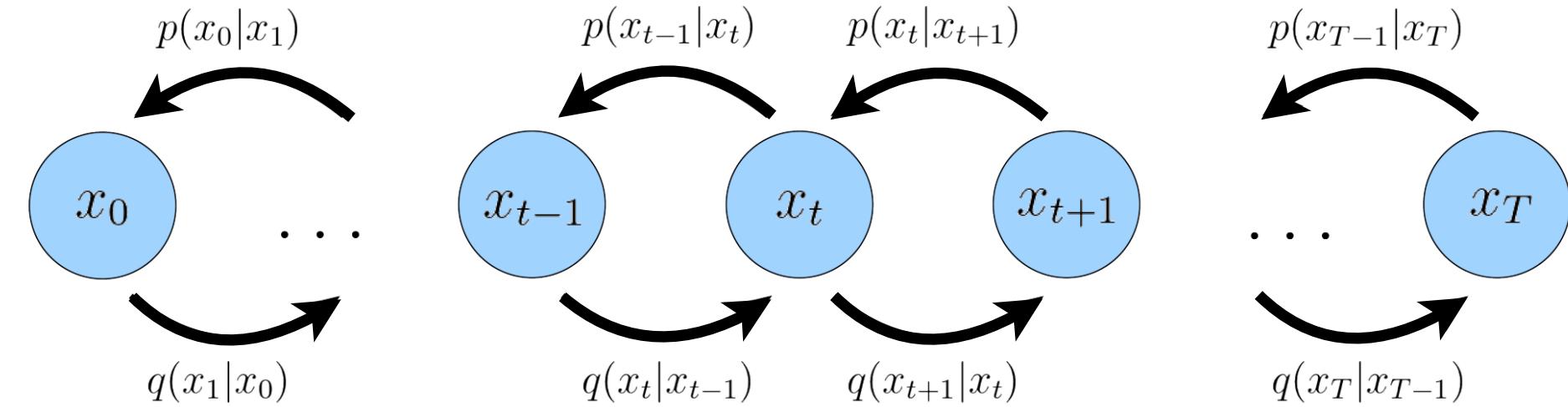


# Variational Diffusion Models: Training



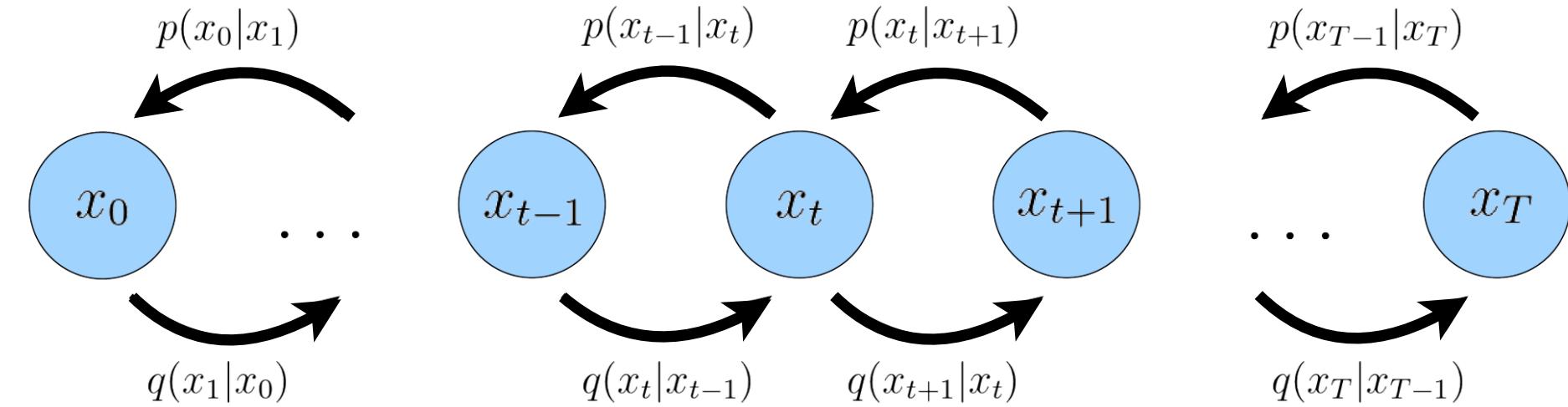
$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

# Variational Diffusion Models: Training



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

# Variational Diffusion Models: Training

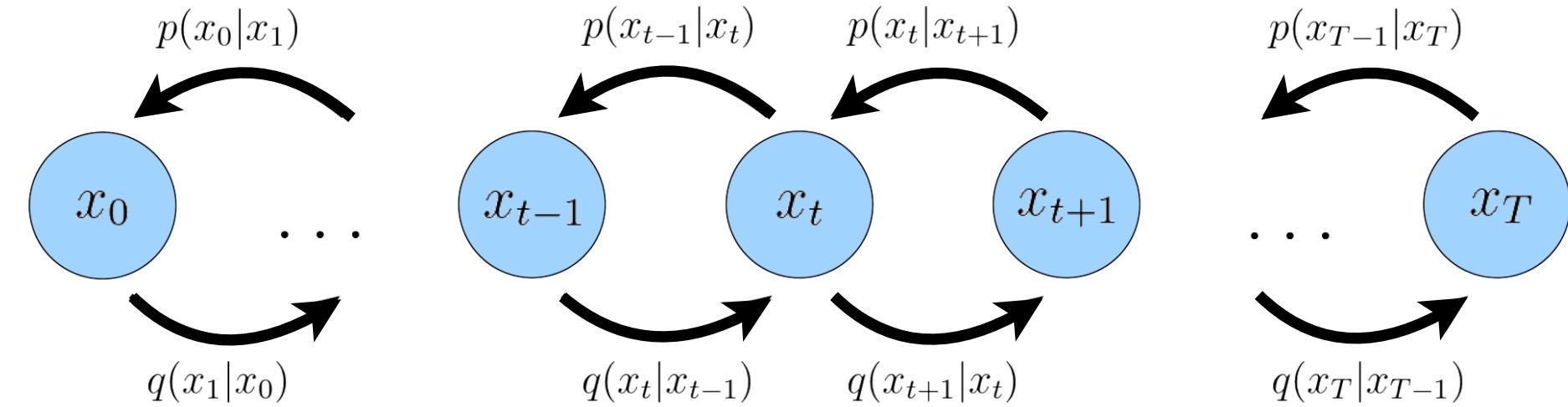


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

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

Bayes Rule

# Variational Diffusion Models: Training



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

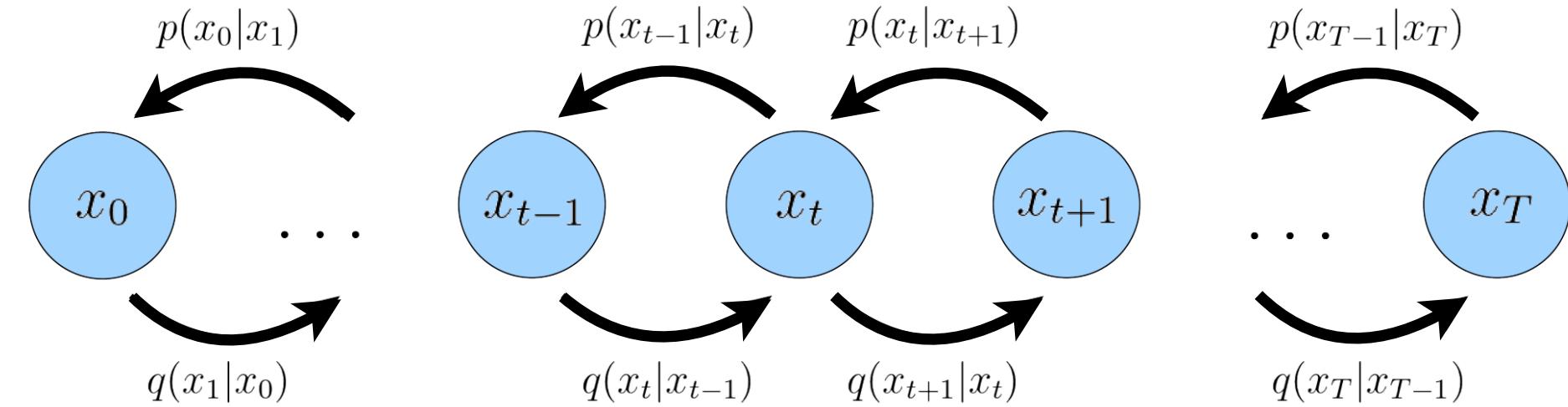
$$q(x_{t-1} | x_t, x_0) = \frac{q(x_t | x_{t-1}, x_0) q(x_{t-1} | x_0)}{q(x_t | x_0)}$$

=  $q(x_t | x_{t-1})$

Bayes Rule

Markov  
Property

# Variational Diffusion Models: Training



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

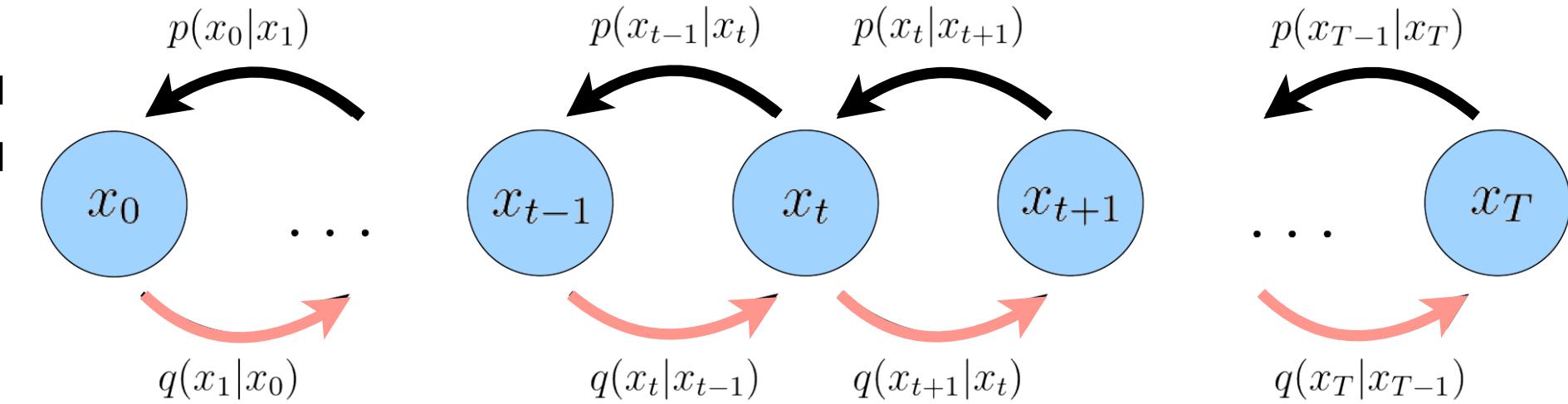
$$q(x_{t-1} | x_t, x_0) = \frac{q(x_t | x_{t-1}, x_0) q(x_{t-1} | x_0)}{q(x_t | x_0)}$$

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

Bayes Rule

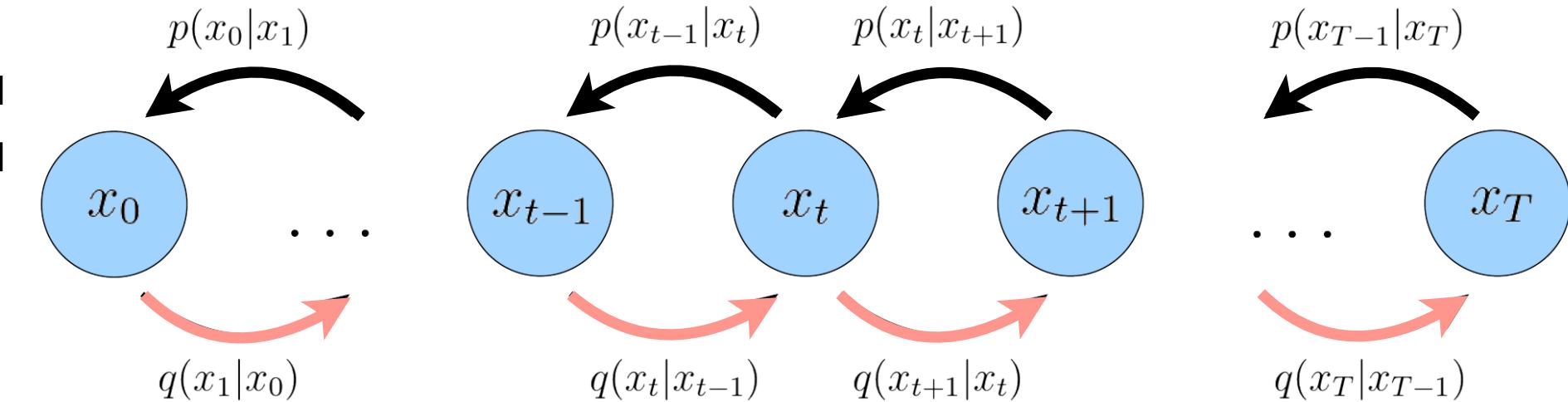
Markov  
Property

# Variational Diffusion Models: Forward Process



- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

# Variational Diffusion Models: Forward Process

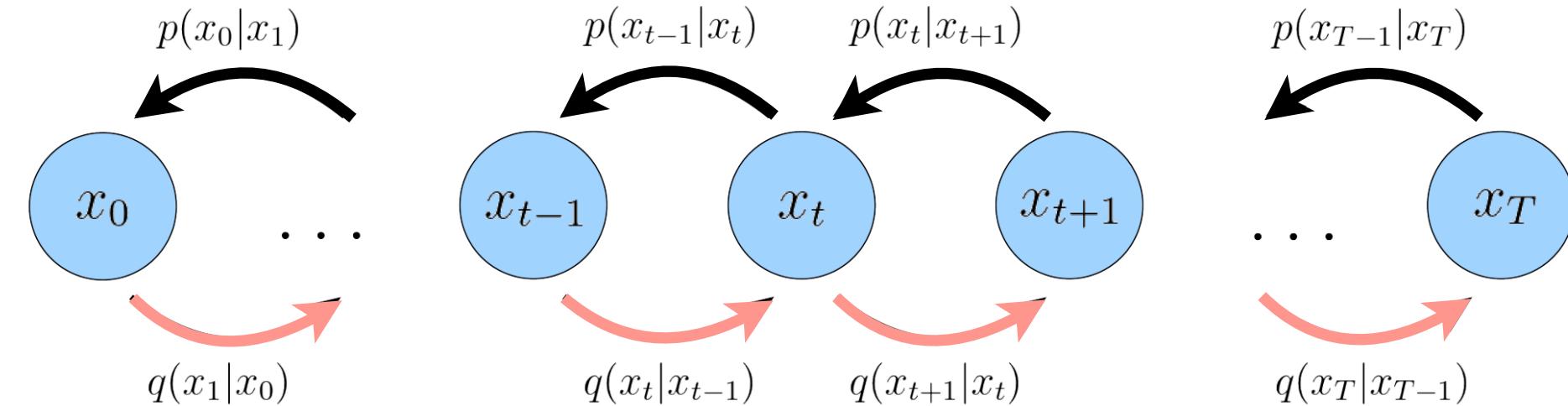


- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$

$\{\alpha_t \in (0, 1)\}_{t=1}^T$  is a predefined variance schedule

# Variational Diffusion Models: Forward Process



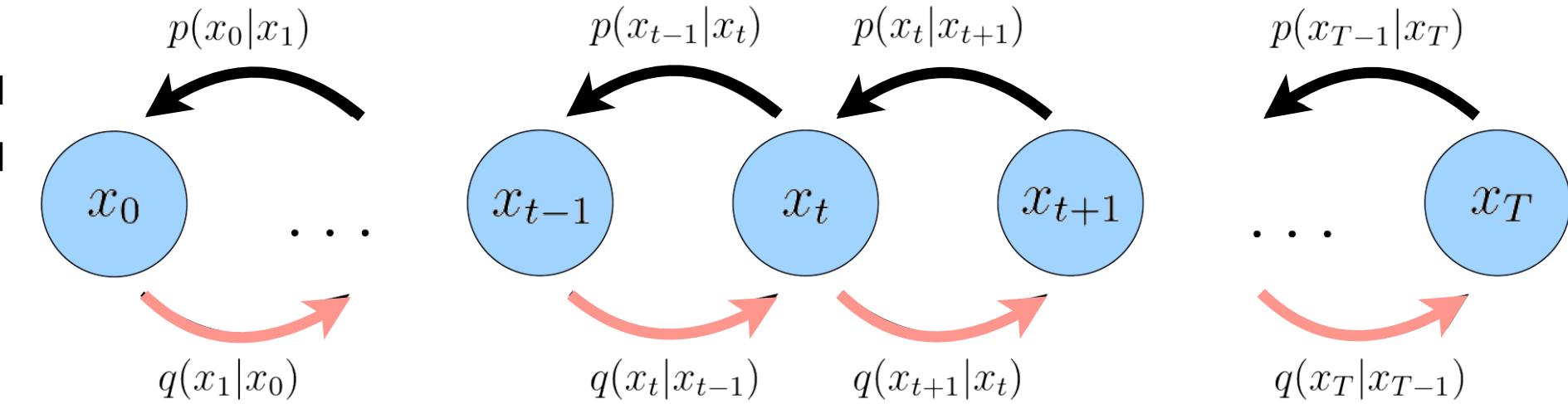
- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$

$\{\alpha_t \in (0, 1)\}_{t=1}^T$  is a predefined variance schedule

- As  $x_0$  is the initial input, obtaining the noisy input for any arbitrary timestep involves solving the following,

# Variational Diffusion Models: Forward Process



- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

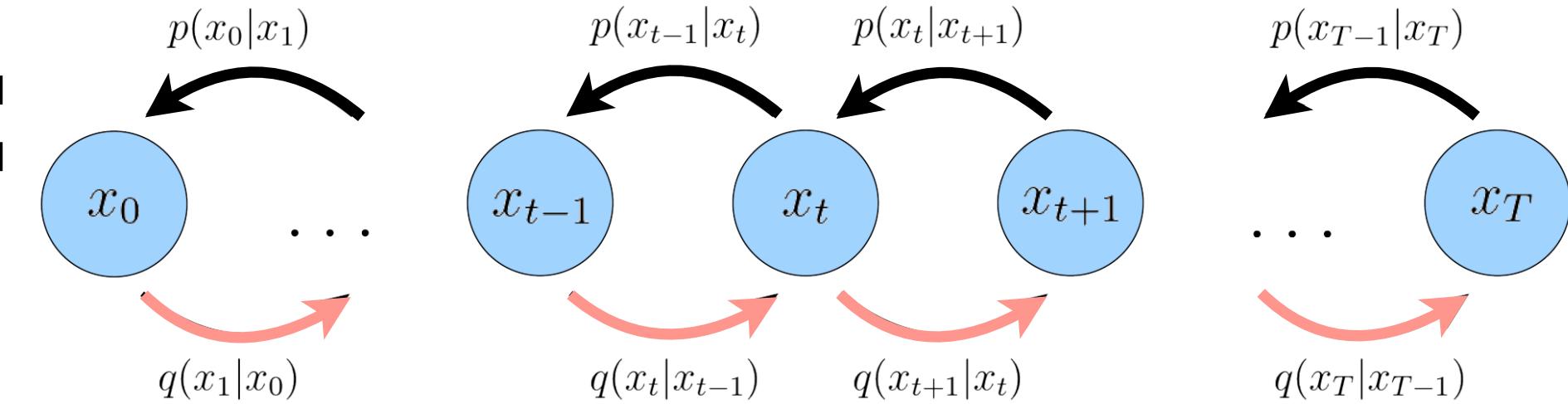
$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$

$\{\alpha_t \in (0,1)\}_{t=1}^T$  is a predefined variance schedule

- As  $x_0$  is the initial input, obtaining the noisy input for any arbitrary timestep involves solving the following,

$$q(x_t | x_0) = q(x_t | x_{t-1}) \cdot q(x_{t-1} | x_{t-2}) \dots q(x_1 | x_0)$$

# Variational Diffusion Models: Forward Process



- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

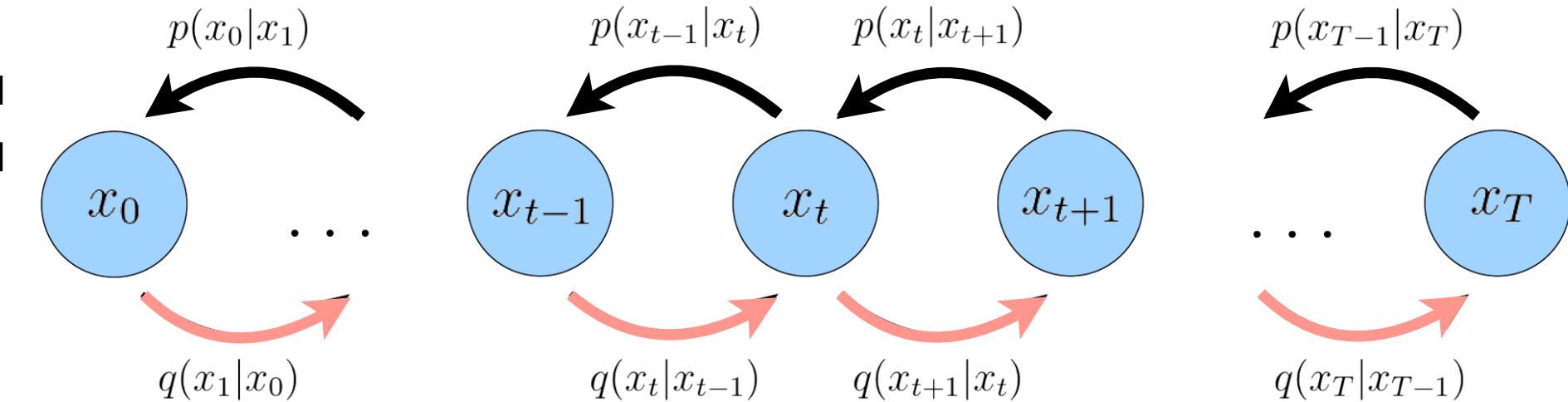
$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$

$\{\alpha_t \in (0, 1)\}_{t=1}^T$  is a predefined variance schedule

- As  $x_0$  is the initial input, obtaining the noisy input for any arbitrary timestep involves solving the following,

$$q(x_t | x_0) = q(x_t | x_{t-1}) \cdot q(x_{t-1} | x_{t-2}) \dots q(x_1 | x_0)$$

# Variational Diffusion Models: Forward Process



- **Assumption:** Encoder at each step is a Gaussian distribution centered around the output of the previous step.

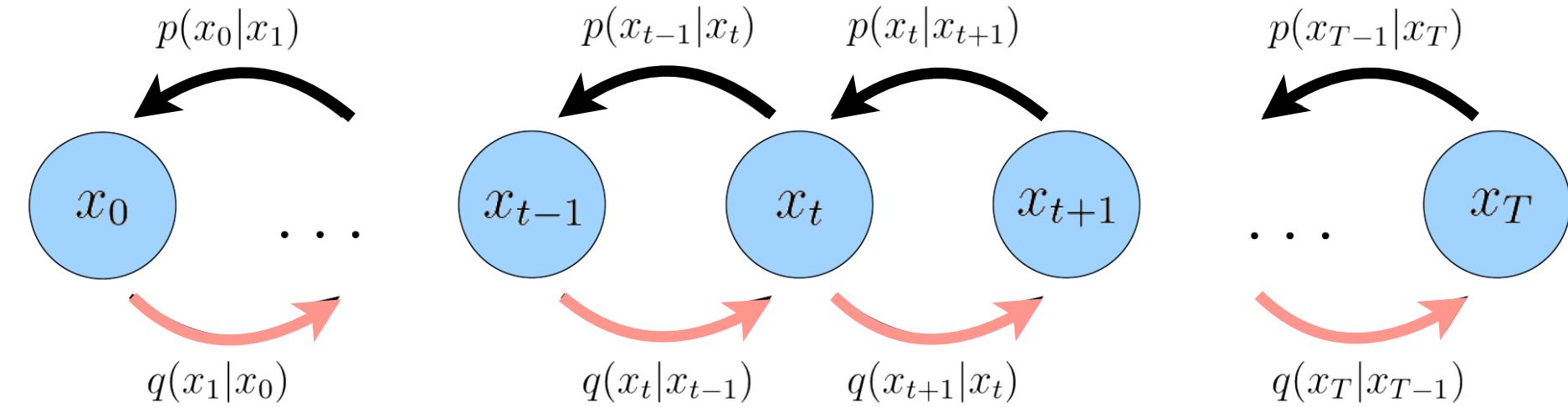
$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbf{I}\right)$$

$\{\alpha_t \in (0, 1)\}_{t=1}^T$  is a predefined variance schedule

- As  $x_0$  is the initial input, obtaining the noisy input for any arbitrary timestep involves solving the following,

$$q(x_t | x_0) = \prod_{i=1}^t q(x_i | x_{i-1})$$

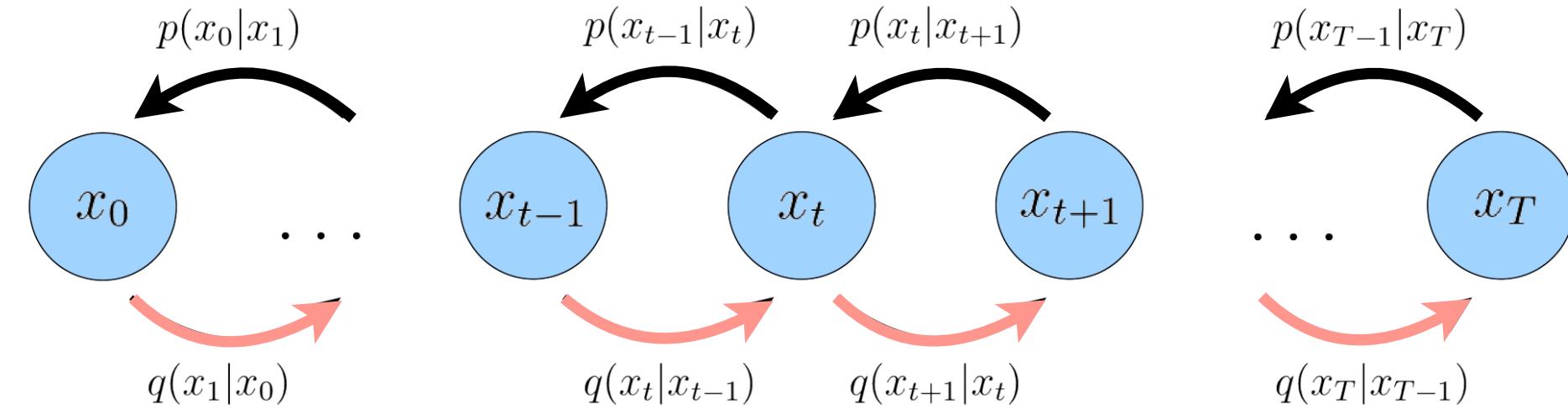
# Variational Diffusion Models: Forward Process



$$q(x_t | x_0) = \prod_{i=1}^t q(x_i | x_{i-1})$$

- Naively computing this is **slow**.
  - Given a new input, one would need to go through this chain of length  $t$ , computing intermediate noisy inputs.

# Variational Diffusion Models: Forward Process

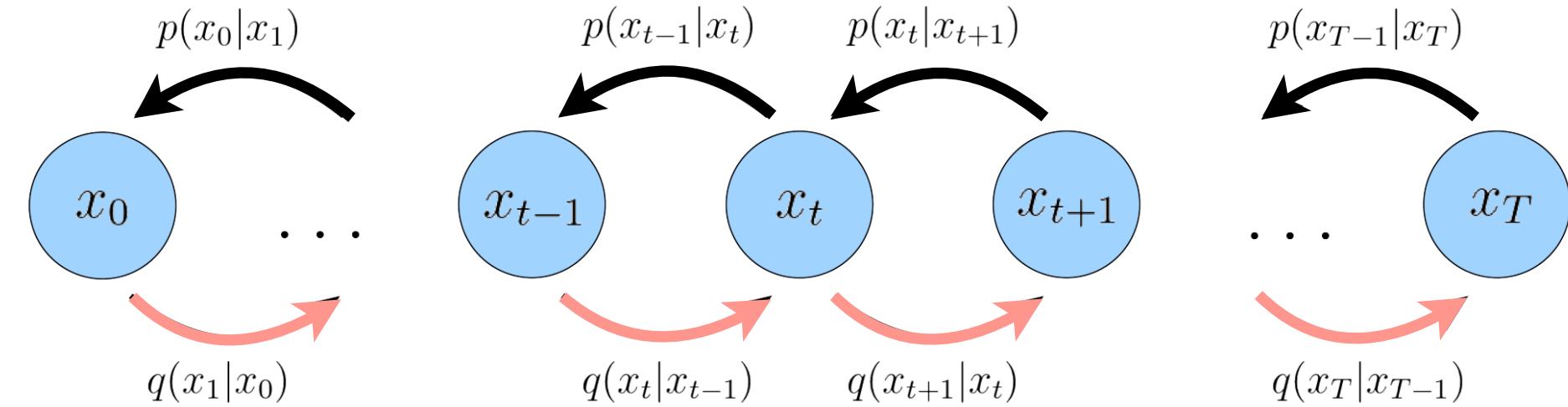


$$q(x_t | x_0) = \prod_{i=1}^t q(x_i | x_{i-1})$$

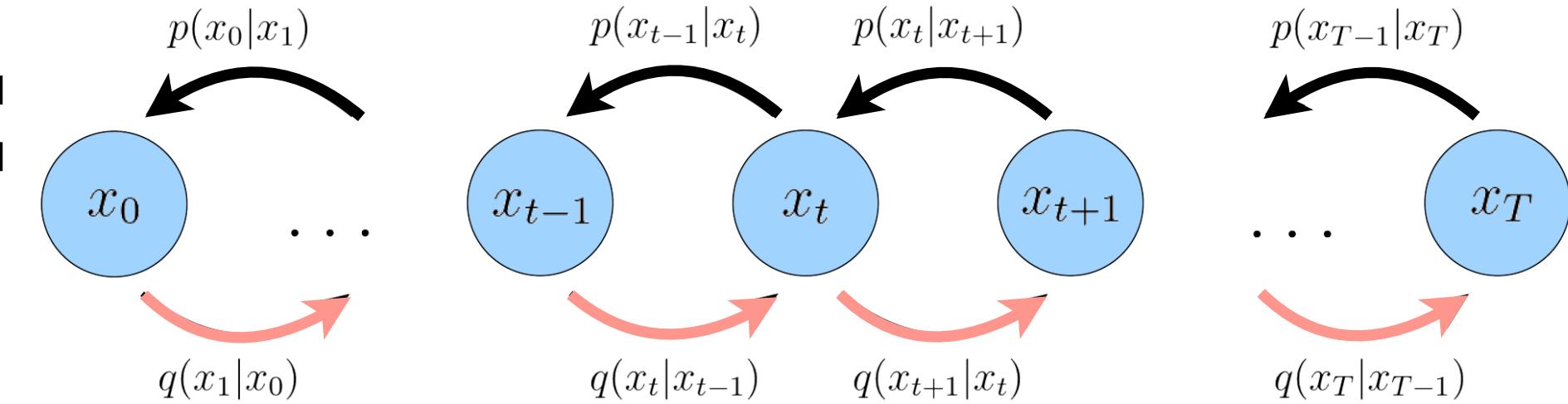
- Naively computing this is **slow**.
  - Given a new input, one would need to go through this chain of length  $t$ , computing intermediate noisy inputs.
- However, as we can leverage properties of Gaussians to make this efficient.

# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \mathcal{N} \left( x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t) \mathbf{I} \right)$$



# Variational Diffusion Models: Forward Process

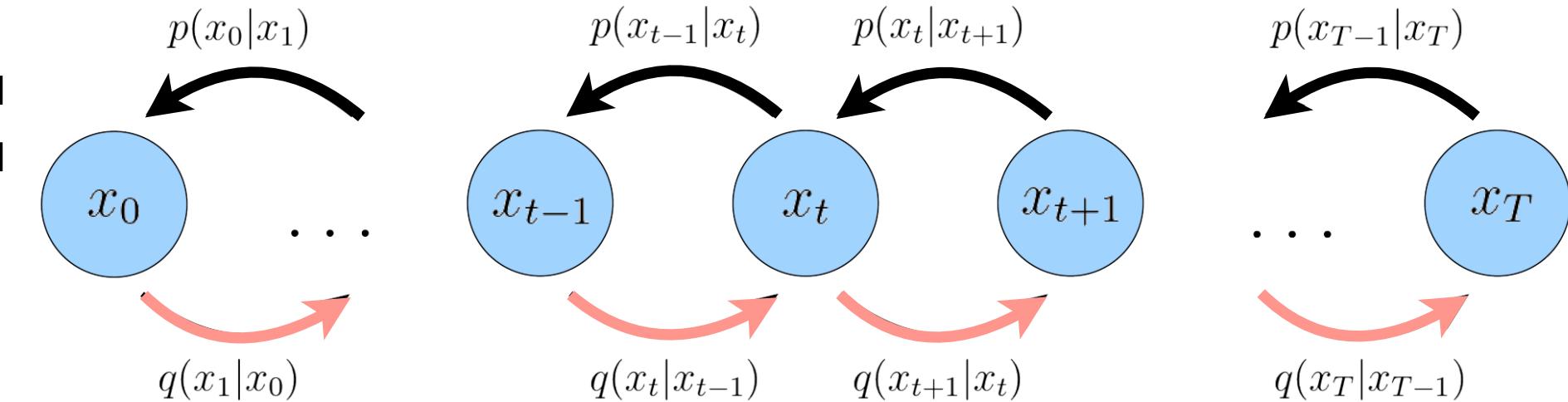


$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1} \quad \epsilon_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

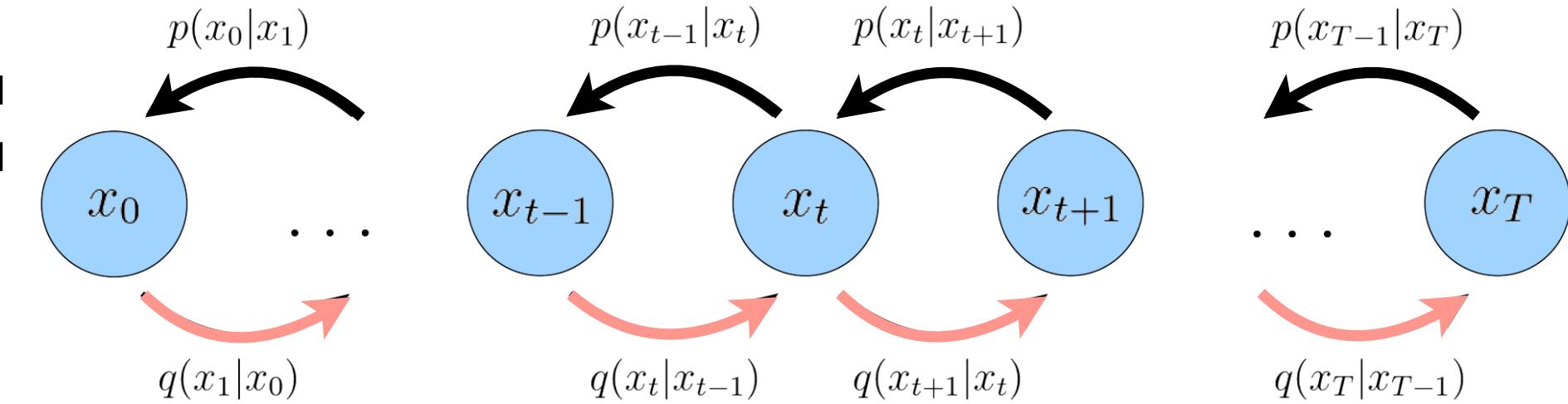
Reparametrization Trick

# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$



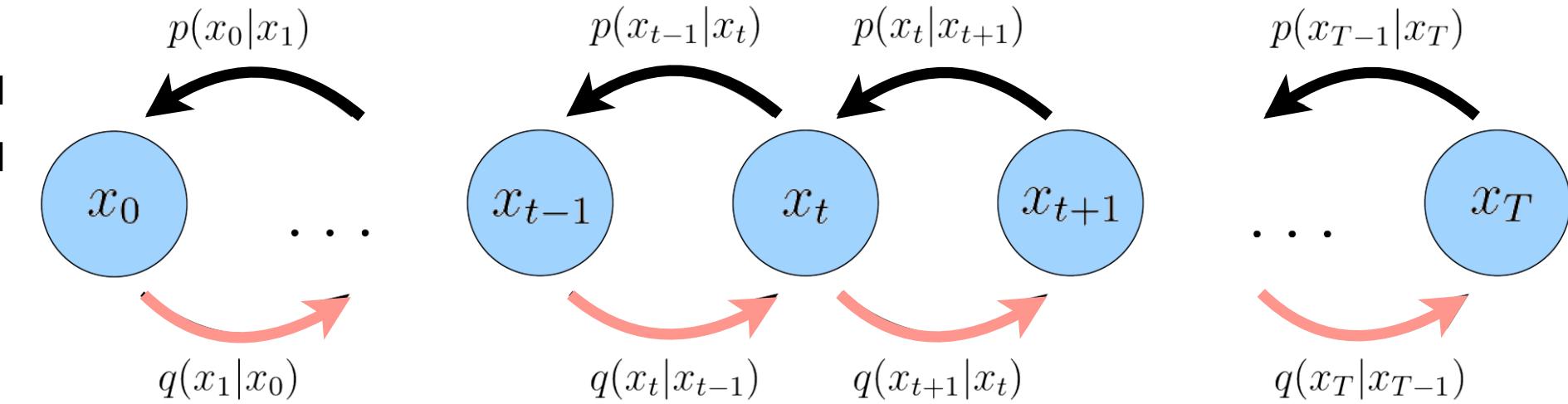
# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

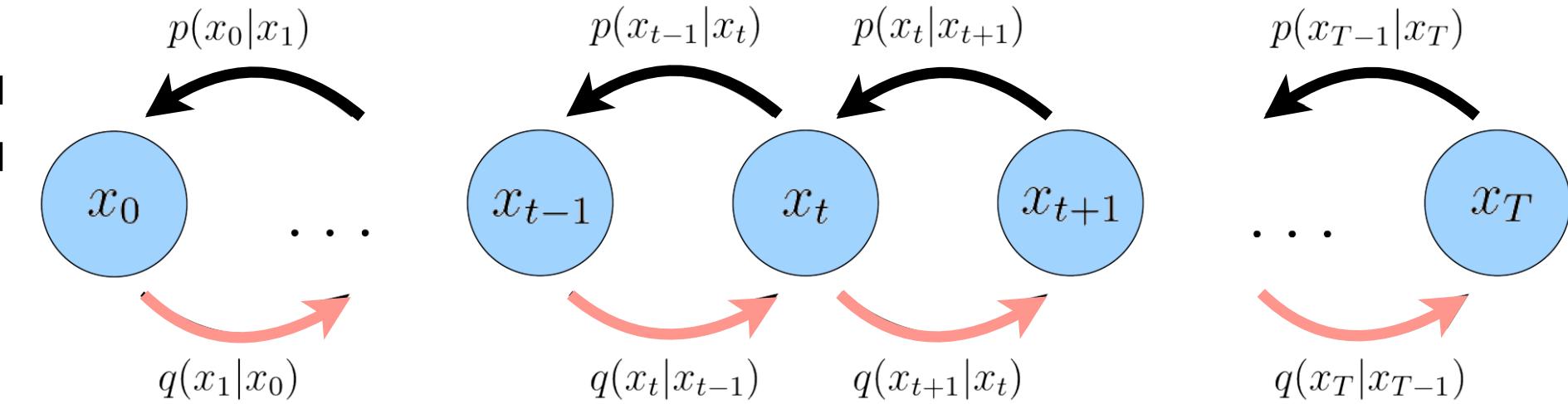
# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

# Variational Diffusion Models: Forward Process

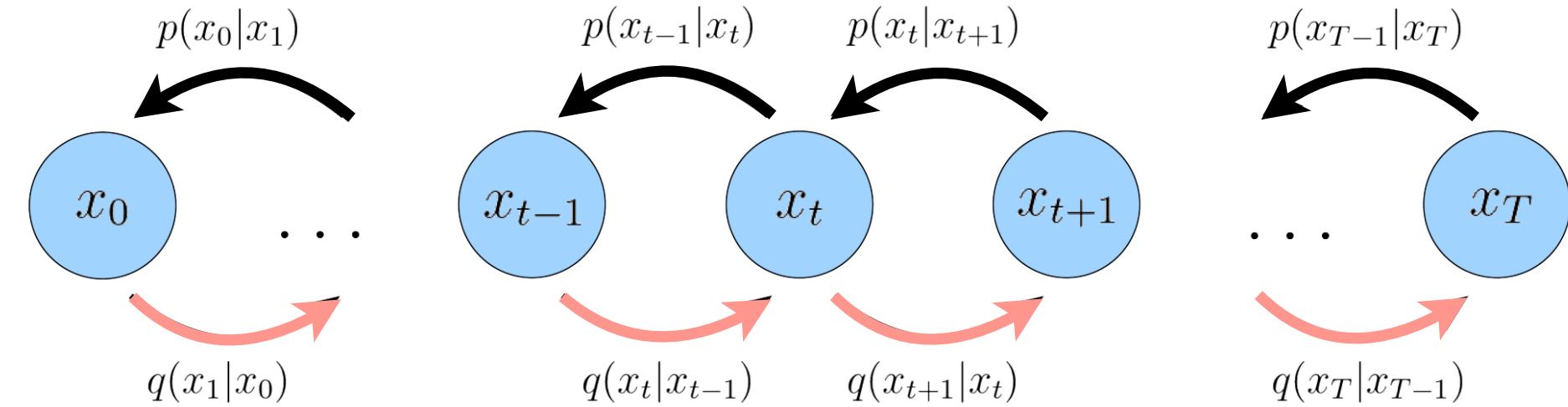


$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2} + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

# Variational Diffusion Models: Forward Process



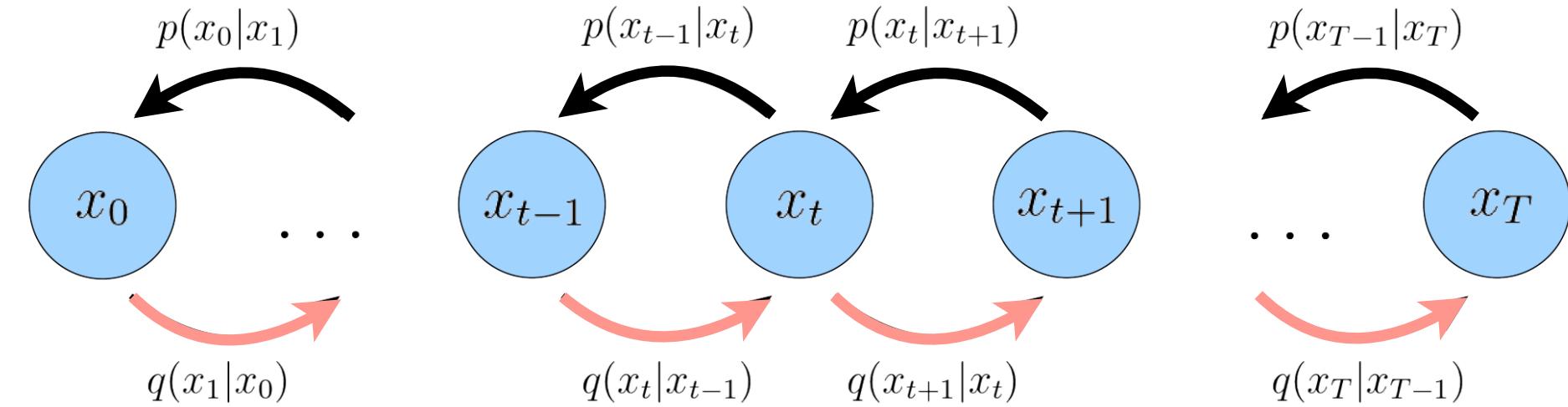
$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2} + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

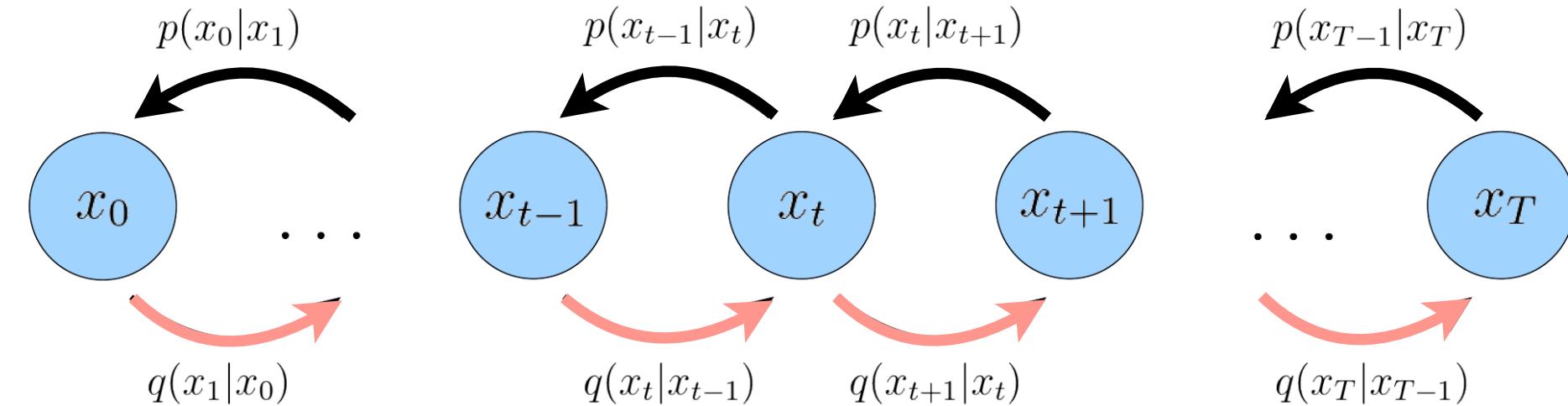
$$\sim \mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I})$$

# Variational Diffusion Models: Forward Process



$$\begin{aligned}
 x_t &\sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1} \quad \text{Recursion} \\
 &= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2} + \sqrt{1 - \alpha_t} \epsilon_{t-1} \\
 &\sim \mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I}) \quad \sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I})
 \end{aligned}$$

# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

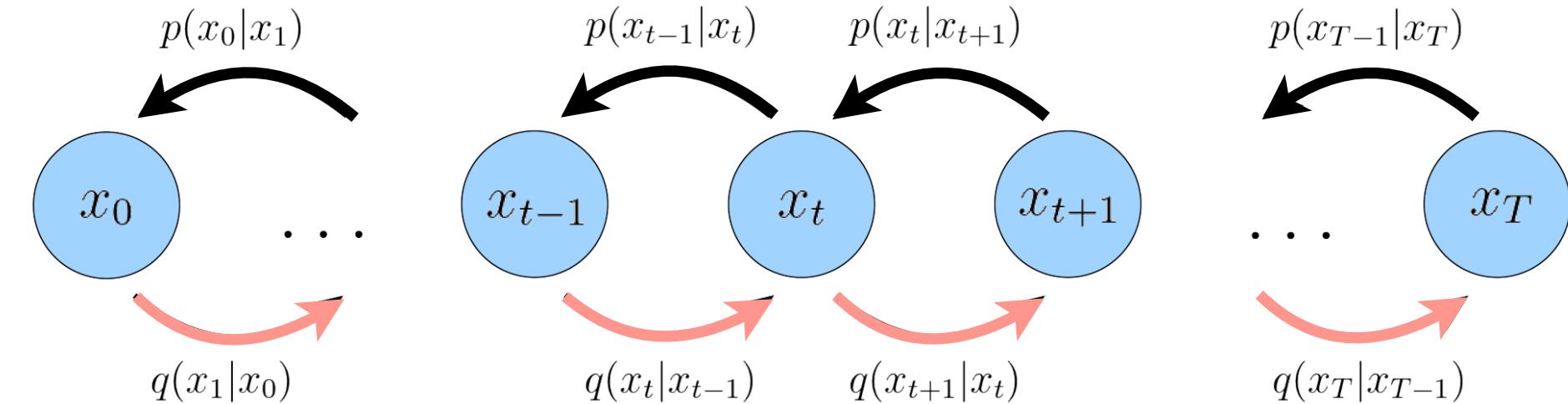
$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2} + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

$\sim \mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I})$

$\sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I})$

Sum of two independent Gaussian random variables is a Gaussian.  
 Mean = sum of the two means  
 Variance = sum of the two variances

# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2} + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

$\sim \mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1}) \mathbf{I})$

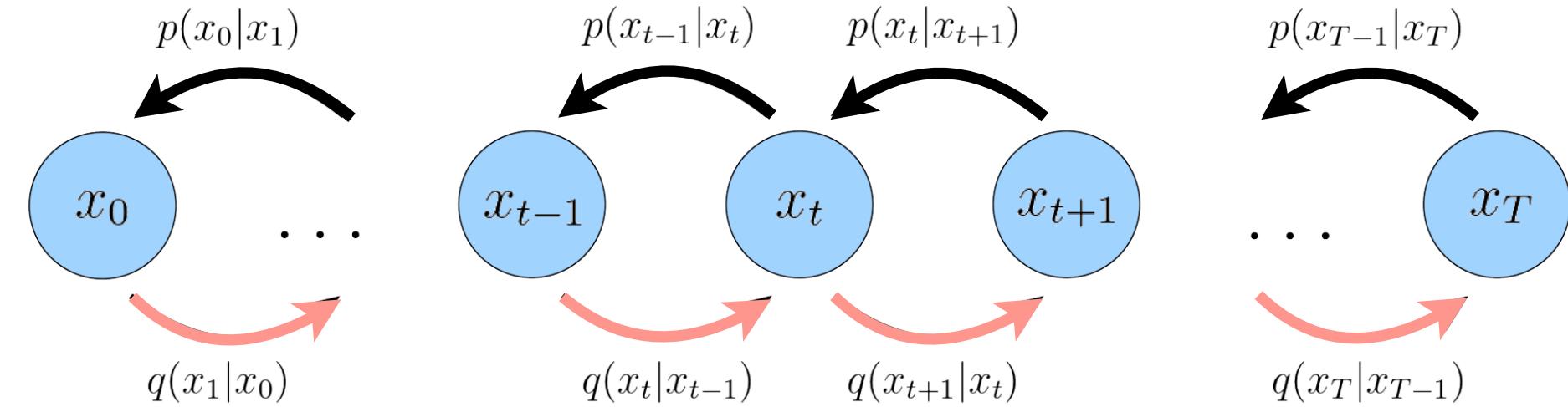
$\sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t) \mathbf{I})$

Sum of two independent Gaussian random variables is a Gaussian.

Mean = sum of the two means

Variance = sum of the two variances

# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-1}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$$

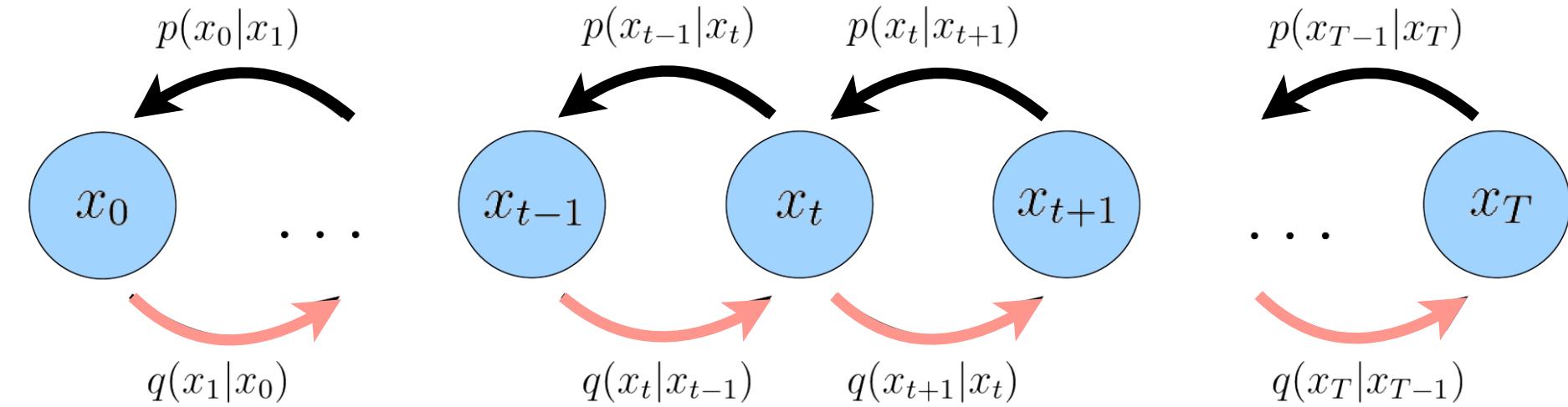
$$\sim \mathcal{N}(0, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})$$

Sum of two independent Gaussian random variables is a Gaussian.

Mean = sum of the two means

Variance = sum of the two variances

# Variational Diffusion Models: Forward Process



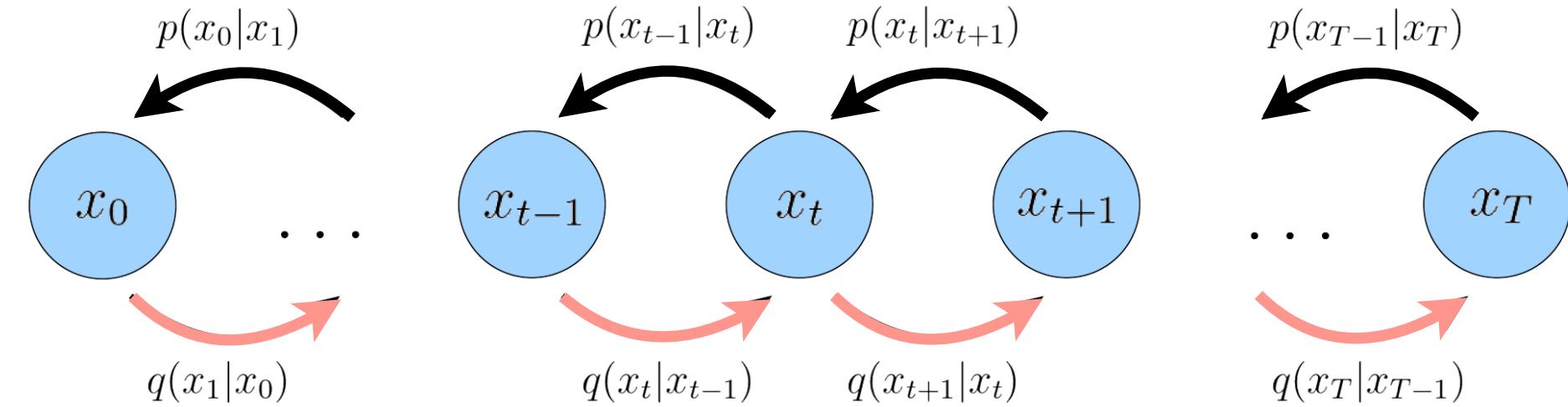
$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$$

$$\sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})$$

# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

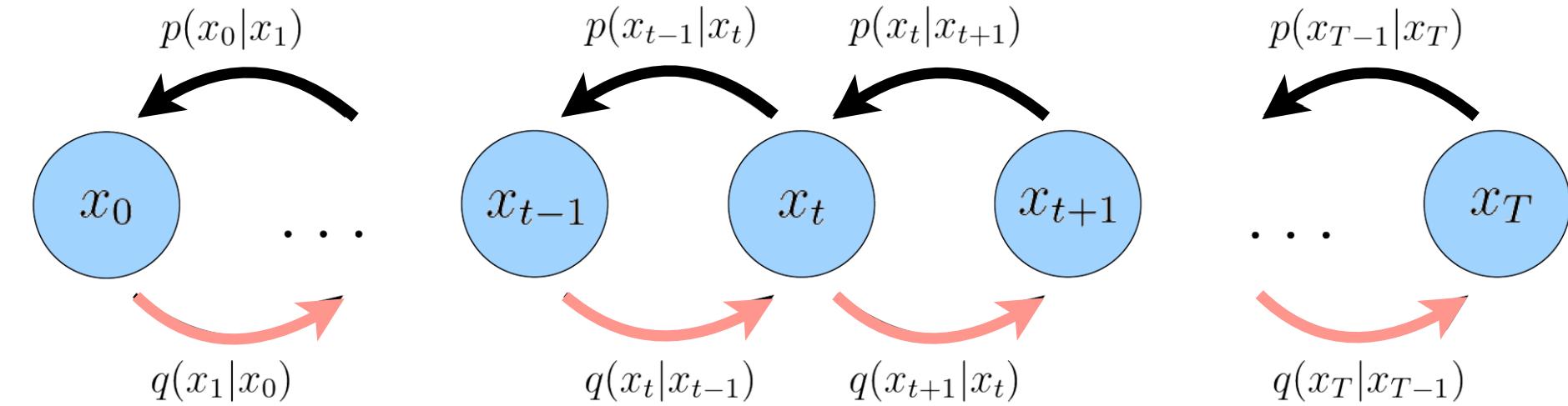
Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$$

$$\sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})$$

$$= \dots$$

# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1 - \alpha_{t-2}} \epsilon_{t-2} \right) + \sqrt{1 - \alpha_t} \epsilon_{t-1}$$

Recursion

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2}$$

$$+ \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$$

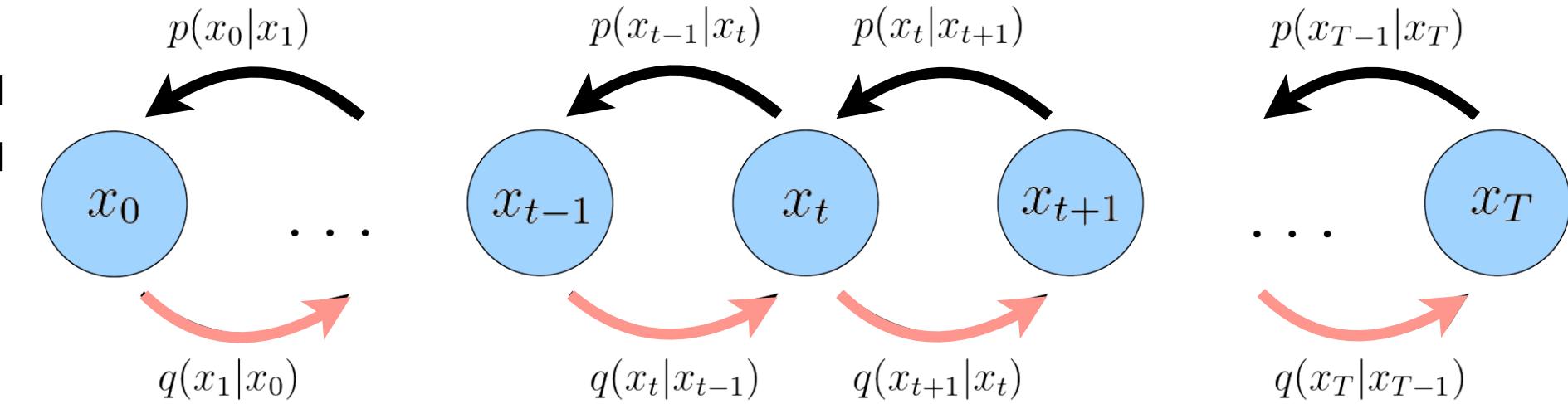
$$\sim \mathcal{N}(\mathbf{0}, (1 - \alpha_t \alpha_{t-1}) \mathbf{I})$$

$$= \dots$$

$$= \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$

# Variational Diffusion Models: Forward Process

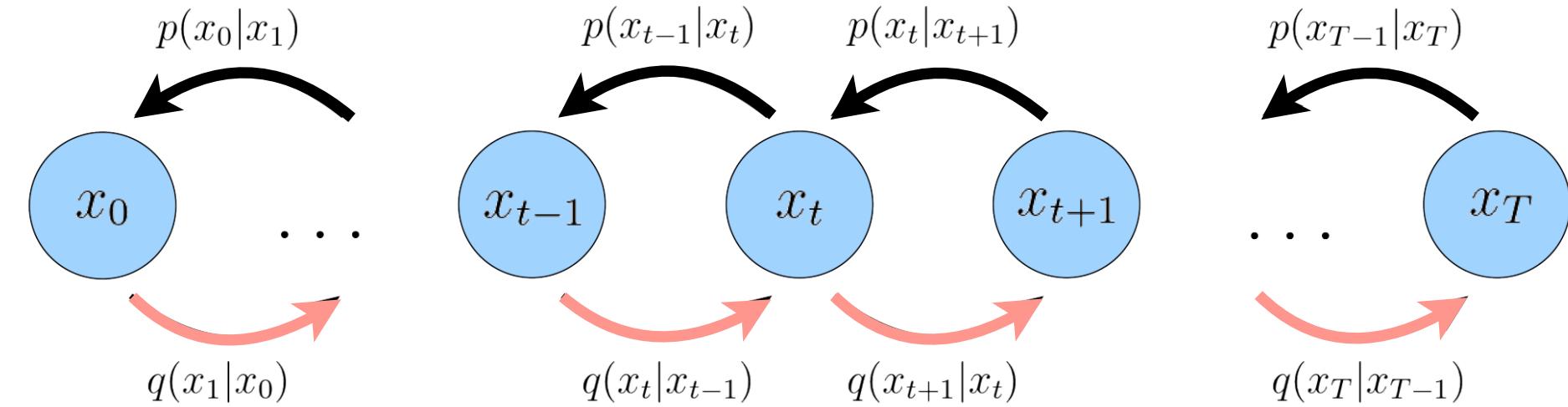
$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$



# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$

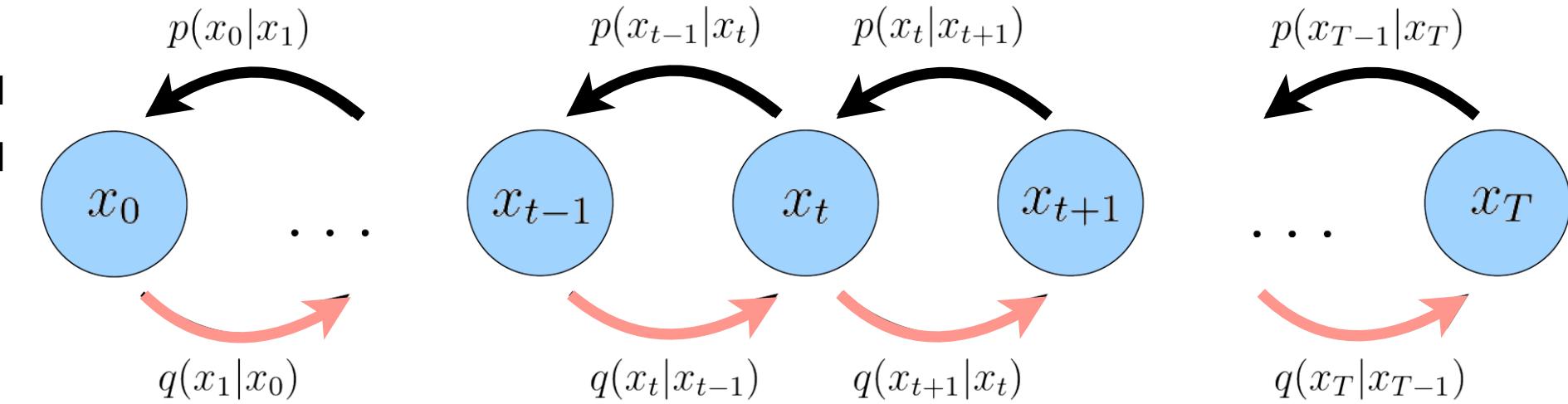
$\sqrt{\bar{\alpha}_t}$



# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$

$\sqrt{\bar{\alpha}_t}$        $\sqrt{1 - \bar{\alpha}_t}$



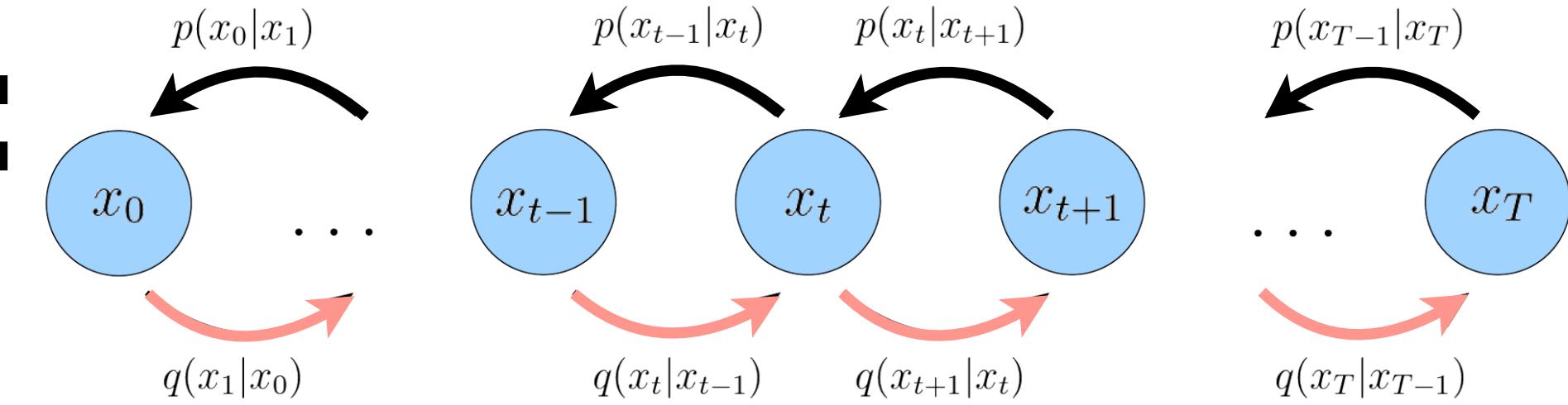
# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$

$$\sqrt{\bar{\alpha}_t}$$

$$\sqrt{1 - \bar{\alpha}_t}$$

$$= \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$



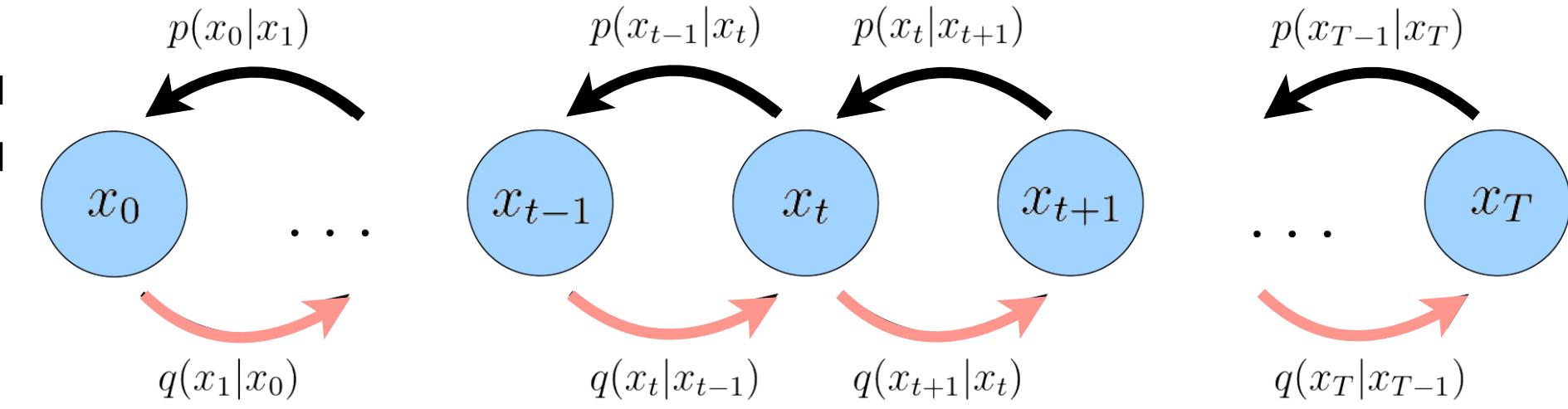
# Variational Diffusion Models: Forward Process

$$x_t \sim q(x_t | x_{t-1}) = \sqrt{\prod_{i=1}^t \alpha_i} x_0 + \sqrt{1 - \prod_{i=1}^t \alpha_i} \epsilon_0^*$$

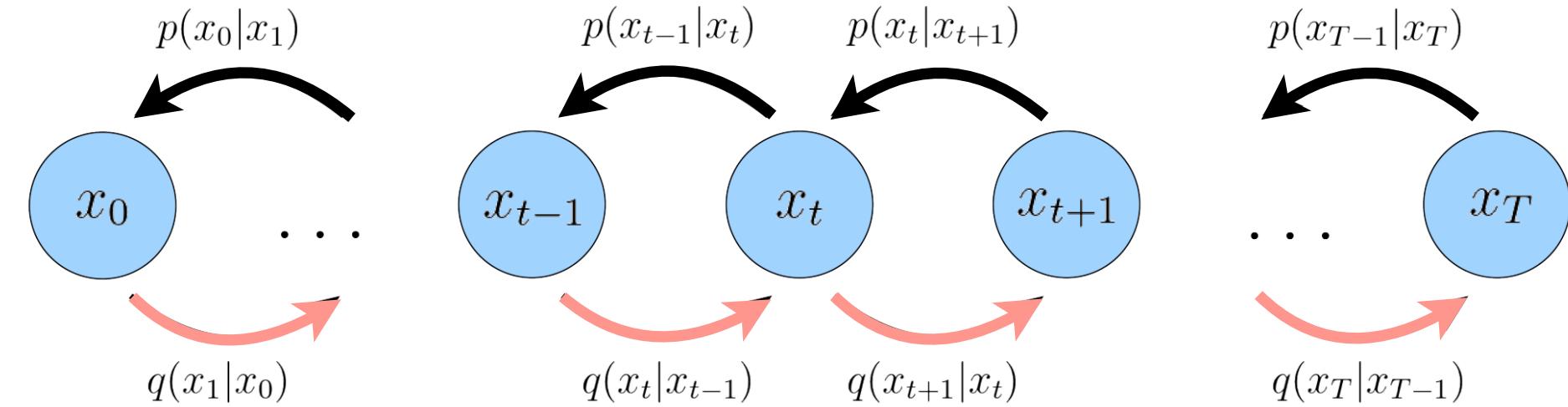
$\sqrt{\bar{\alpha}_t}$ 
 $\sqrt{1 - \bar{\alpha}_t}$

$$x_t \sim q(x_t | x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$

$\sim \mathcal{N}\left(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}\right)$



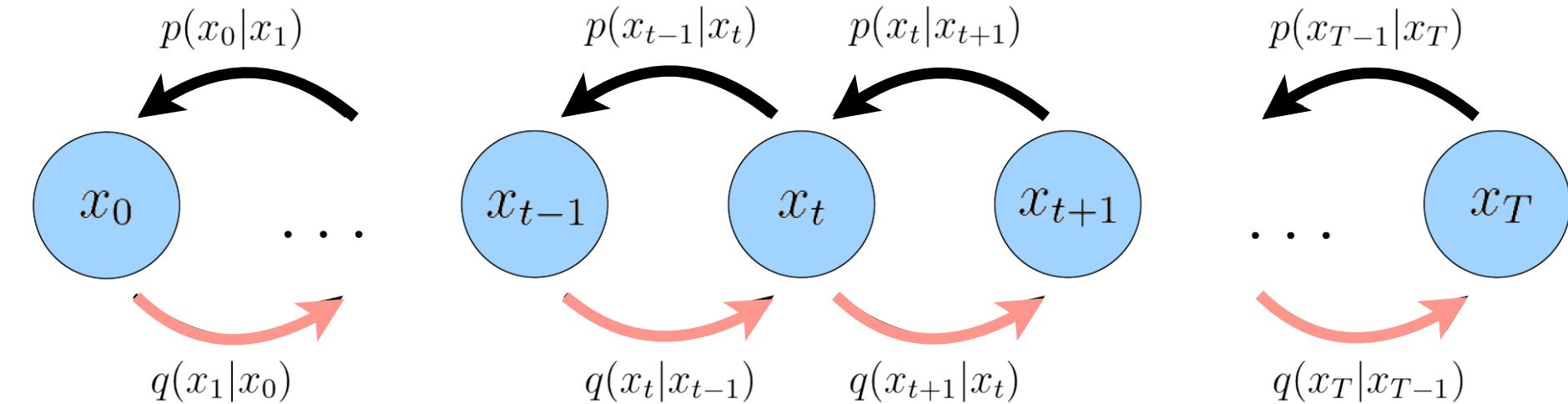
# Variational Diffusion Models: Forward Process



$$x_t \sim q(x_t | x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$

$$\sim \mathcal{N}\left(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}\right)$$

# Variational Diffusion Models: Forward Process

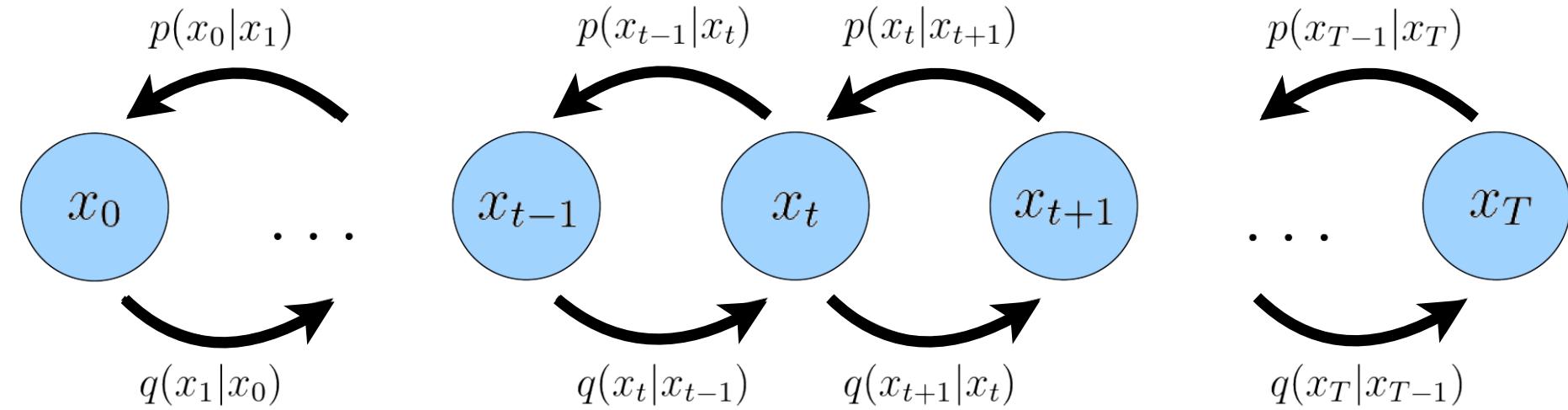


$$x_t \sim q(x_t | x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$

$$\sim \mathcal{N}\left(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}\right)$$

- This enables computing noised inputs at an arbitrary timestep directly from the denoised input.
  - Intuitively, setting  $\bar{\alpha}_1 > \bar{\alpha}_2 > \dots > \bar{\alpha}_T$  makes the sample noise with larger step size.

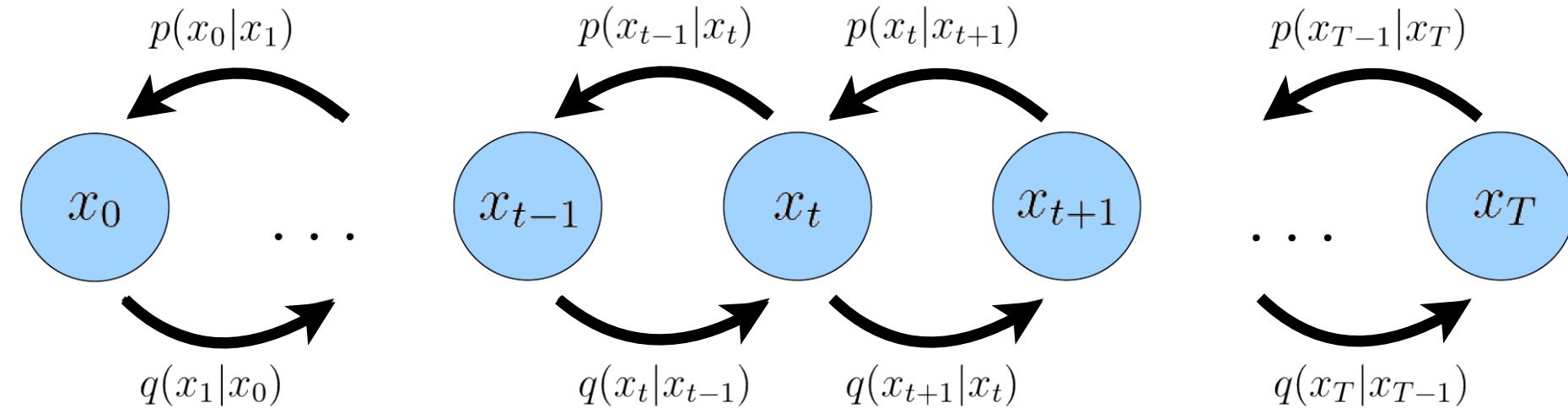
# Variational Diffusion Models: Forward Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- Recall that training diffusion models involves computing the denoising matching term.

# Variational Diffusion Models: Forward Process



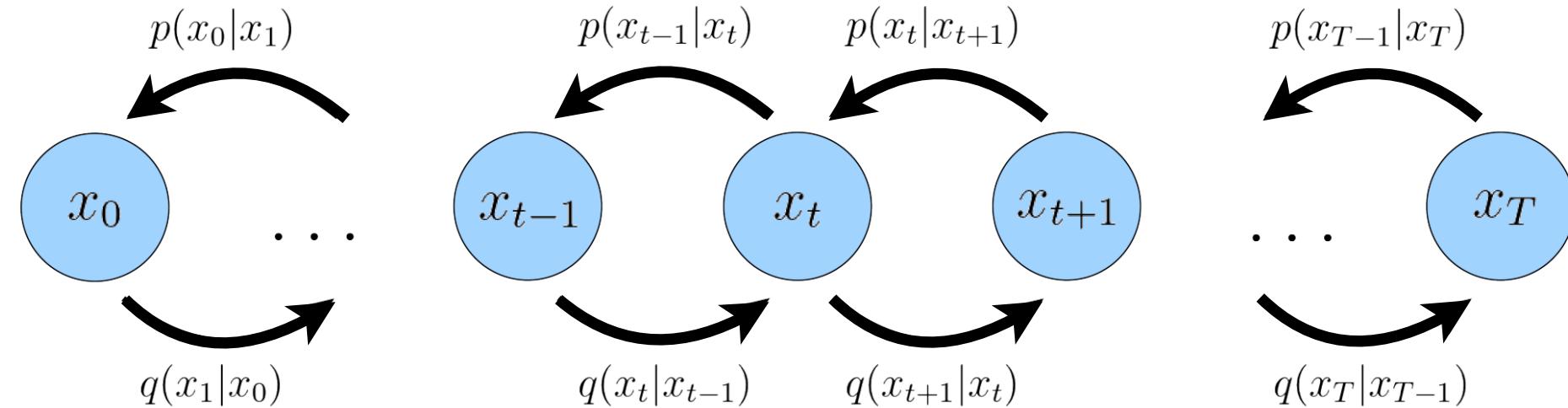
$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- Recall that training diffusion models involves computing the denoising matching term.

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

Bayes Rule

# Variational Diffusion Models: Forward Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- Recall that training diffusion models involves computing the denoising matching term.

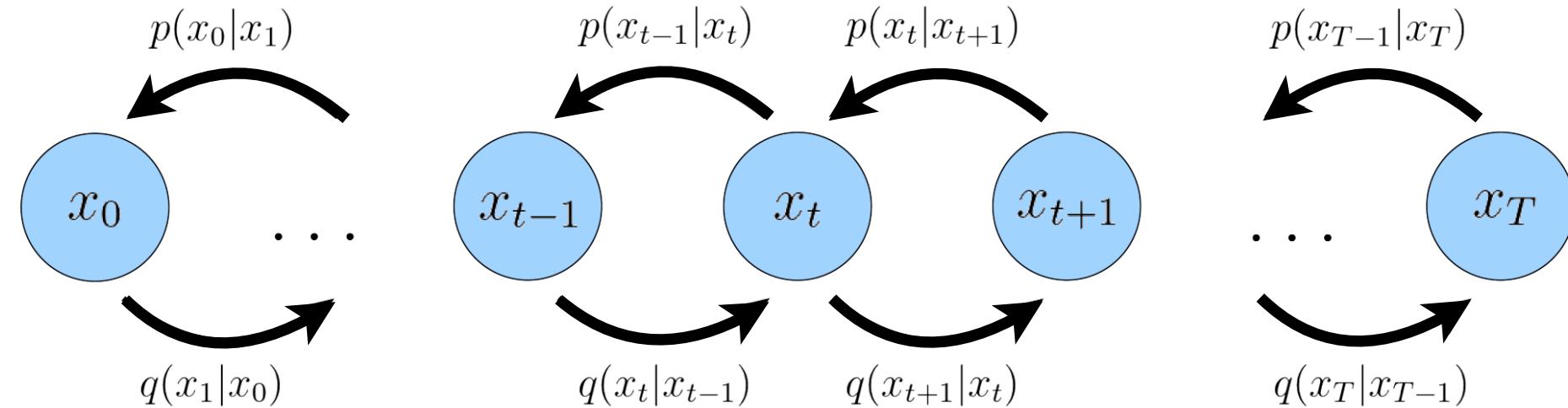
$$q(x_{t-1} | x_t, x_0) = \frac{q(x_t | x_{t-1}, x_0) q(x_{t-1} | x_0)}{q(x_t | x_0)}$$

$= q(x_t | x_{t-1})$

Bayes Rule

Markov  
Property

# Variational Diffusion Models: Forward Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- Recall that training diffusion models involves computing the denoising matching term.

$$= q(x_t | x_{t-1}) \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I})$$

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

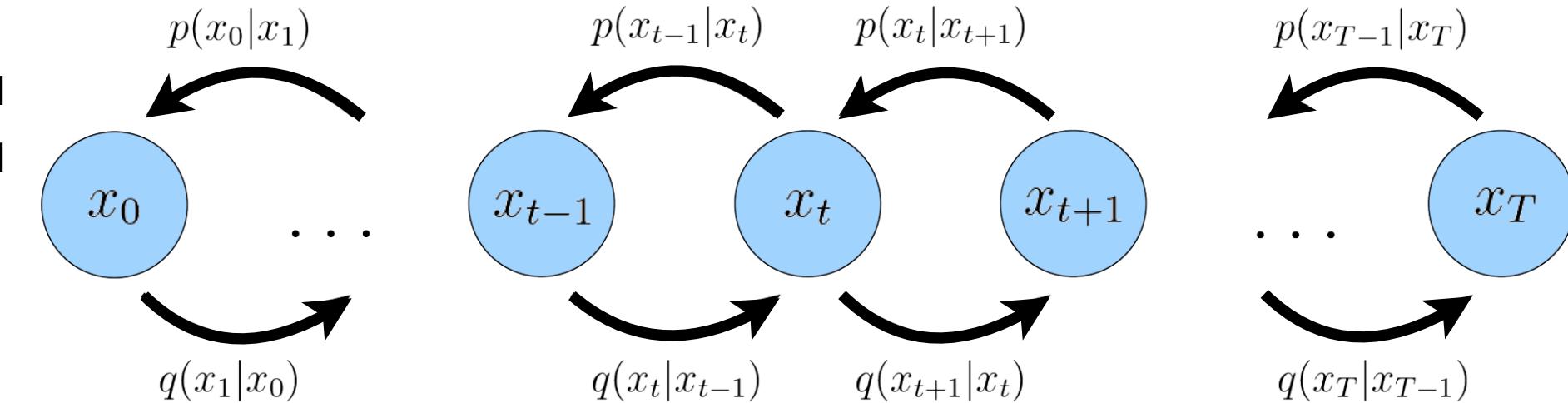
$$q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

Bayes Rule

Markov  
Property

Derivation  
Shown Earlier

# Variational Diffusion Models: Forward Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- Recall that training diffusion models involves computing the denoising matching term.

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

=  $q(x_t | x_{t-1})$ 
 $\mathcal{N}(\sqrt{\bar{\alpha}_{t-1}}x_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})$   
q(x\_t | x\_{t-1}, x\_0)
q(x\_{t-1} | x\_0)  
q(x\_t | x\_0)
 $\mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$

Bayes Rule

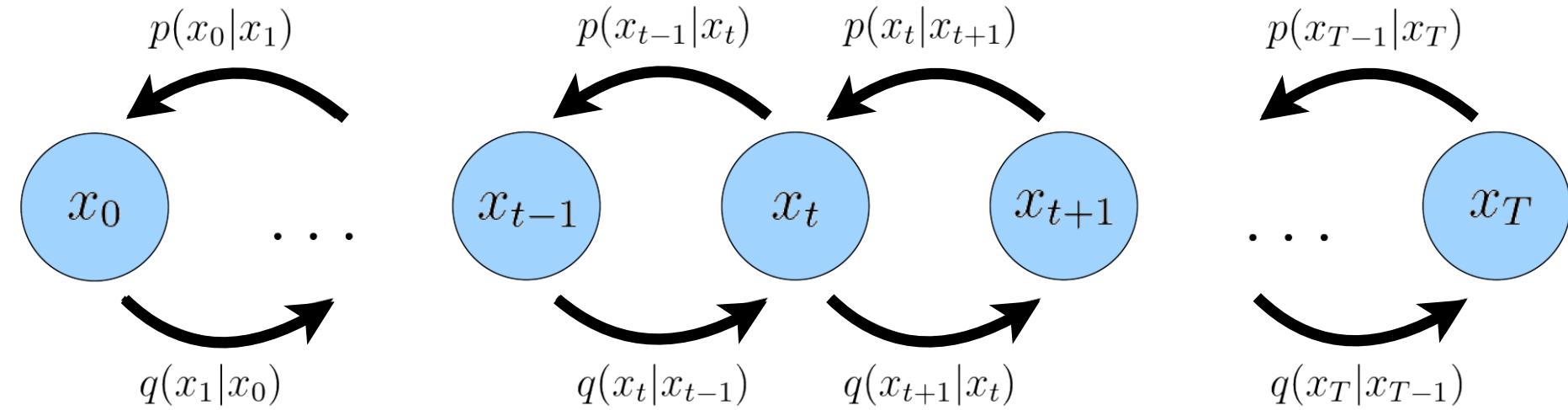
Markov  
Property

Derivation  
Shown Earlier

- Plugging in the distributions and simplifying,

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I})$$

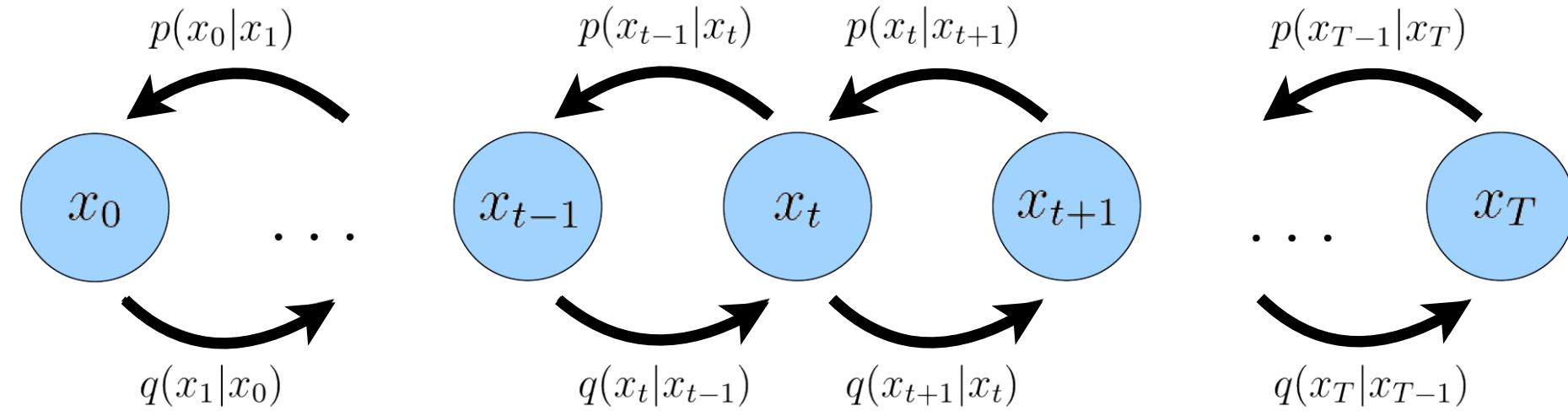
# Variational Diffusion Models: Forward Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I})$$

# Variational Diffusion Models: Forward Process

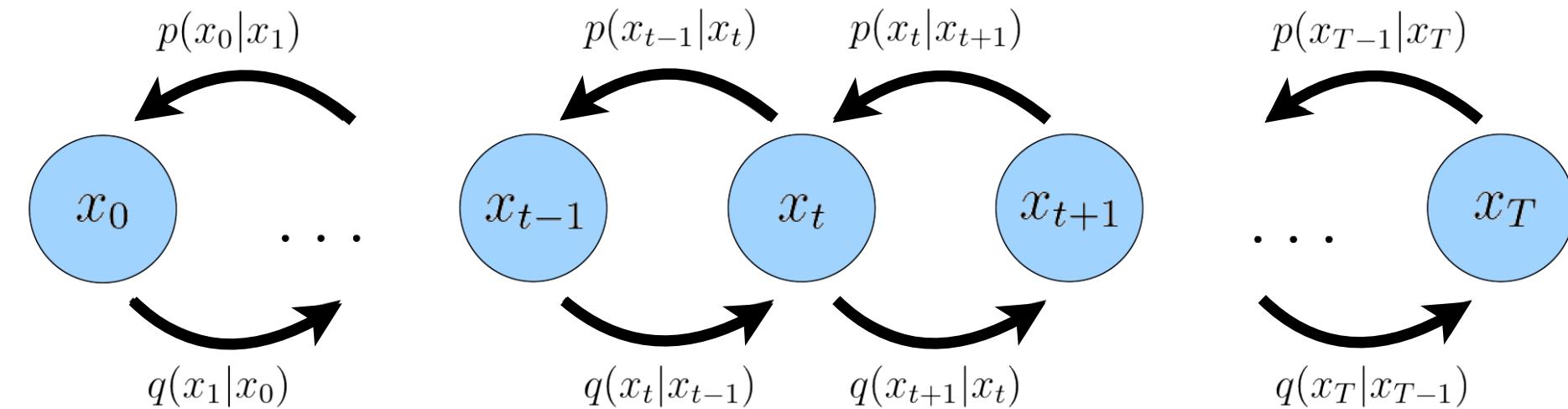


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)})$$

Dependent on noisy and ground  
truth input

# Variational Diffusion Models: Forward Process



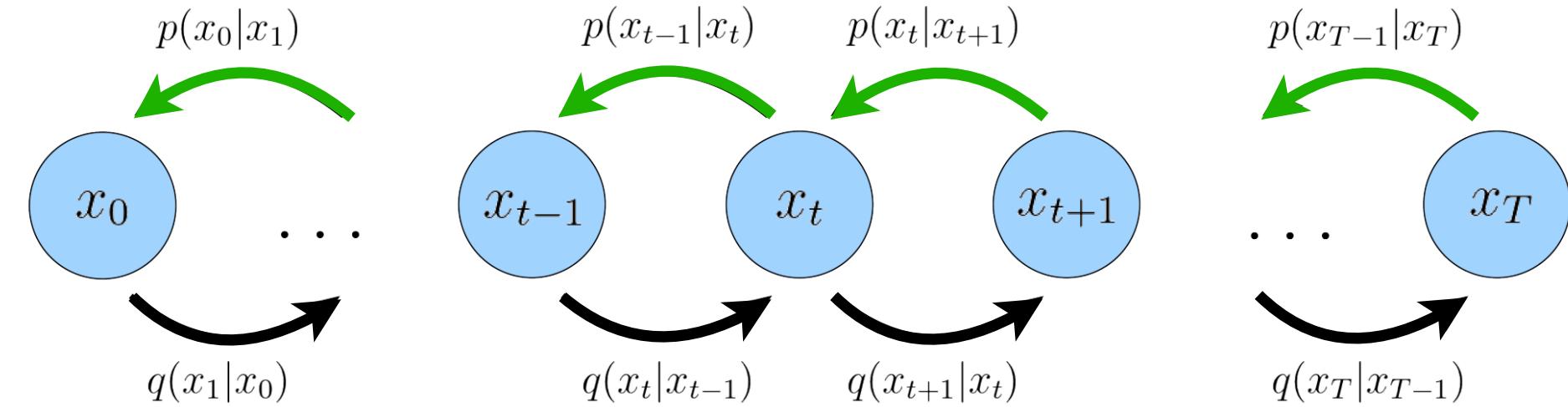
$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)})$$

Dependent on noisy and ground truth input

Only dependent on time t

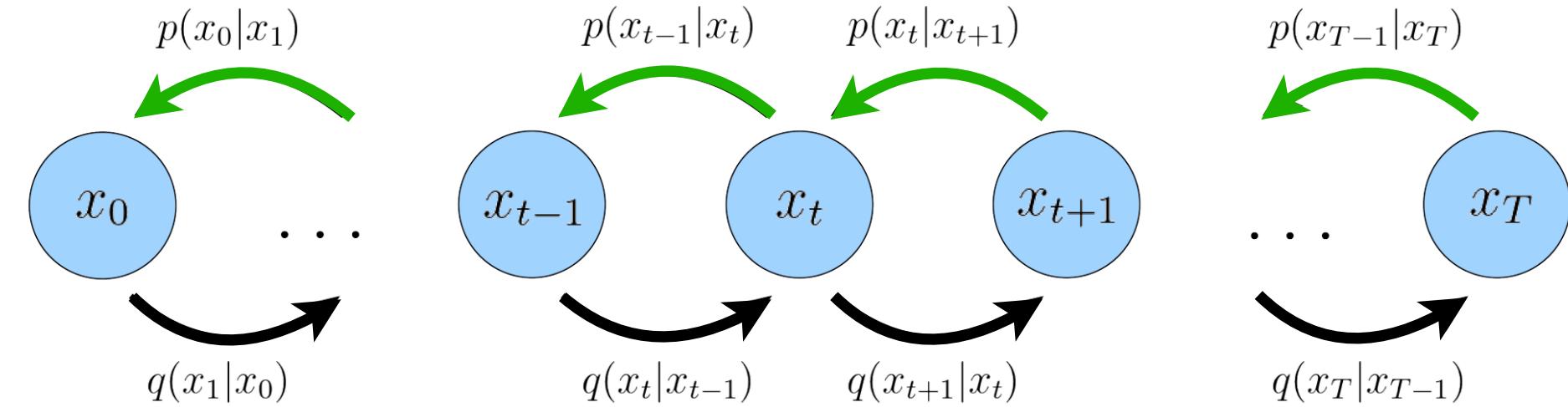
# Variational Diffusion Models: Reverse Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- The goal of the reverse process is to transform noise to a sample from the data distribution.

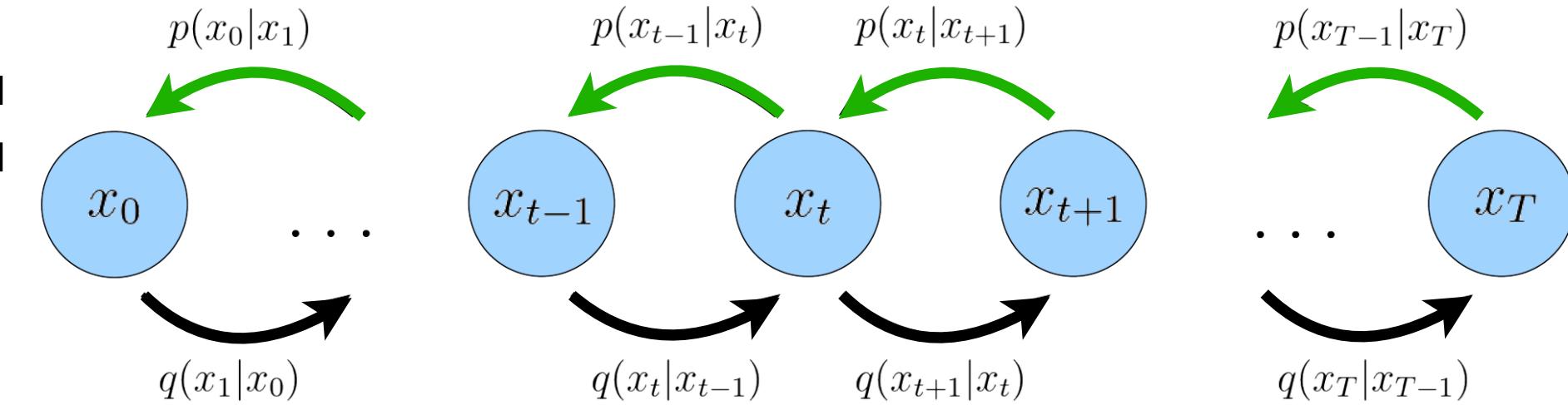
# Variational Diffusion Models: Reverse Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- The goal of the reverse process is to transform noise to a sample from the data distribution.

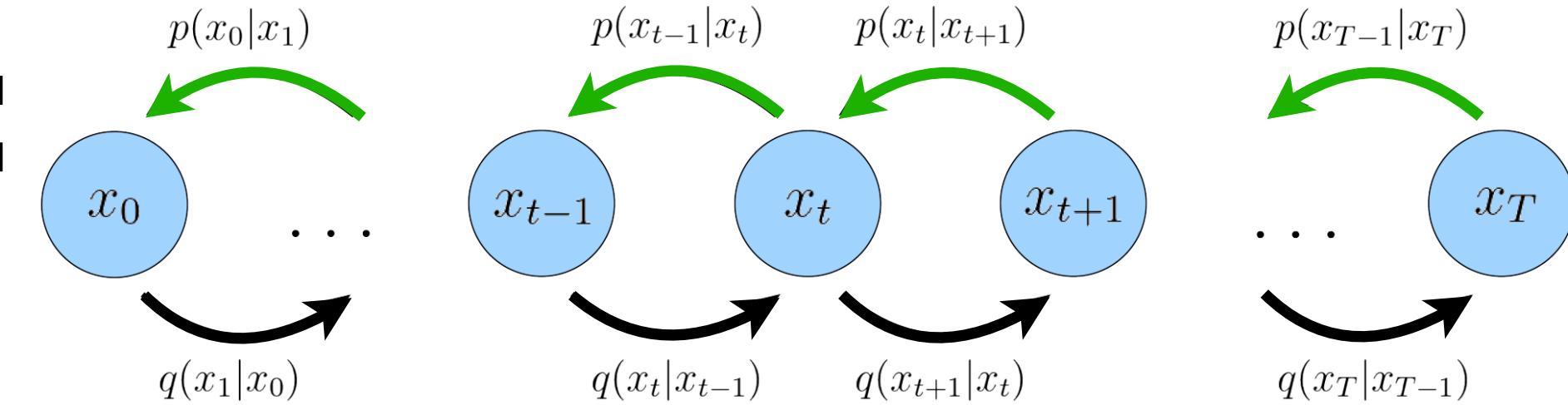
# Variational Diffusion Models: Reverse Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- The goal of the reverse process is to transform noise to a sample from the data distribution.
- Ideally if we can reverse the forward process and sample from  $q(x_{t-1} | x_t)$ , we can easily convert noise to a sample.

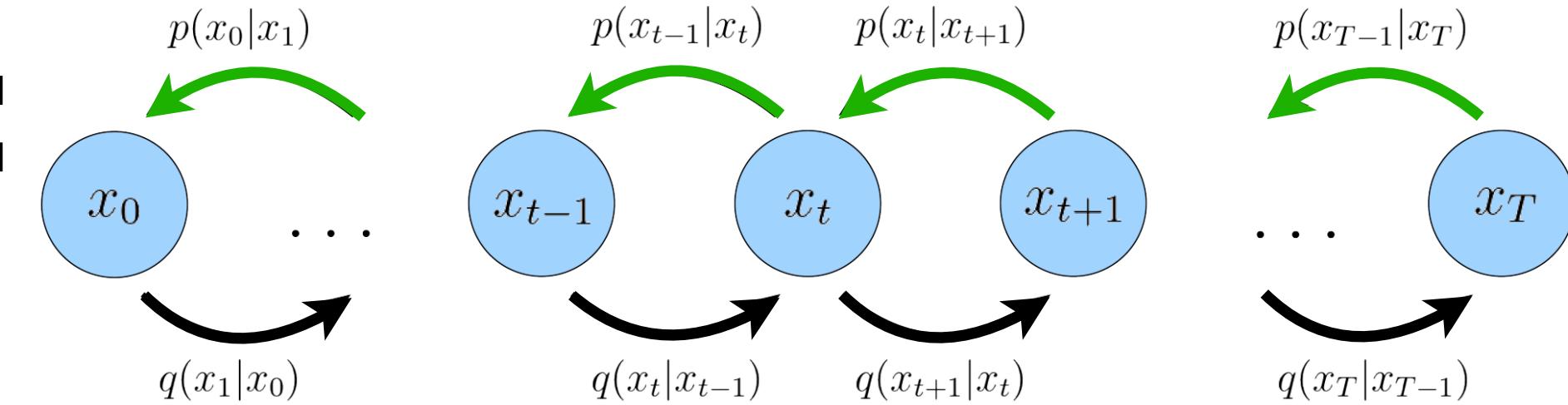
# Variational Diffusion Models: Reverse Process



$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- The goal of the reverse process is to transform noise to a sample from the data distribution.
- Ideally if we can reverse the forward process and sample from  $q(x_{t-1} | x_t)$ , we can easily convert noise to a sample.
  - $q(x_{t-1} | x_t)$  is hard to estimate as it requires computing marginals.
  - Instead learn a model  $p_{\theta}(x_{t-1} | x_t)$  to approximate the true reverse process.

# Variational Diffusion Models: Reverse Process

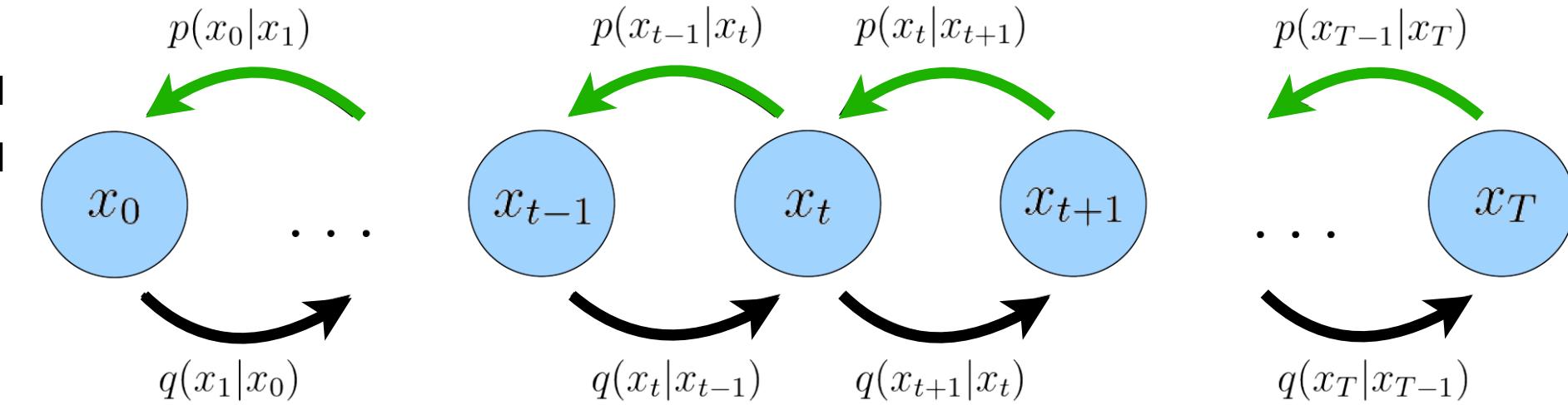


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{I}}_{\Sigma_q(t)}}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}$$

- Following the KL optimization term, in order to ensure  $p_{\theta}(x_{t-1} | x_t)$  approximates the ground-truth transition  $q(x_{t-1} | x_t, x_0)$  as closely as possible, it is modelled as a Gaussian.

# Variational Diffusion Models: Reverse Process

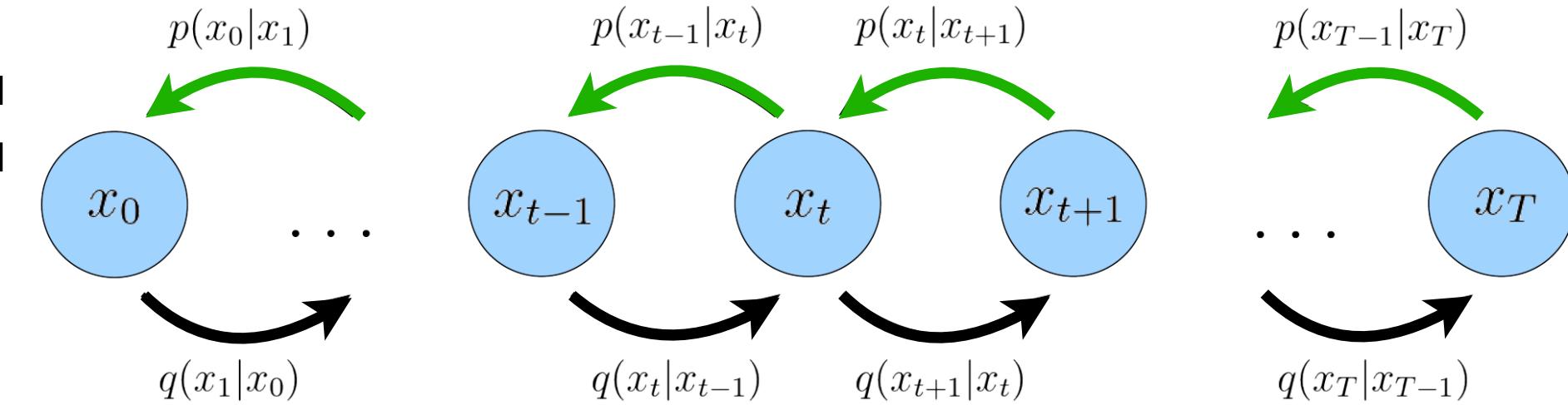


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)}}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}$$

- Following the KL optimization term, in order to ensure  $p_{\theta}(x_{t-1} | x_t)$  approximates the ground-truth transition  $q(x_{t-1} | x_t, x_0)$  as closely as possible, it is modelled as a Gaussian.
- The mean is parameterized as the function  $\mu_{\theta}(x_t, t)$

# Variational Diffusion Models: Reverse Process

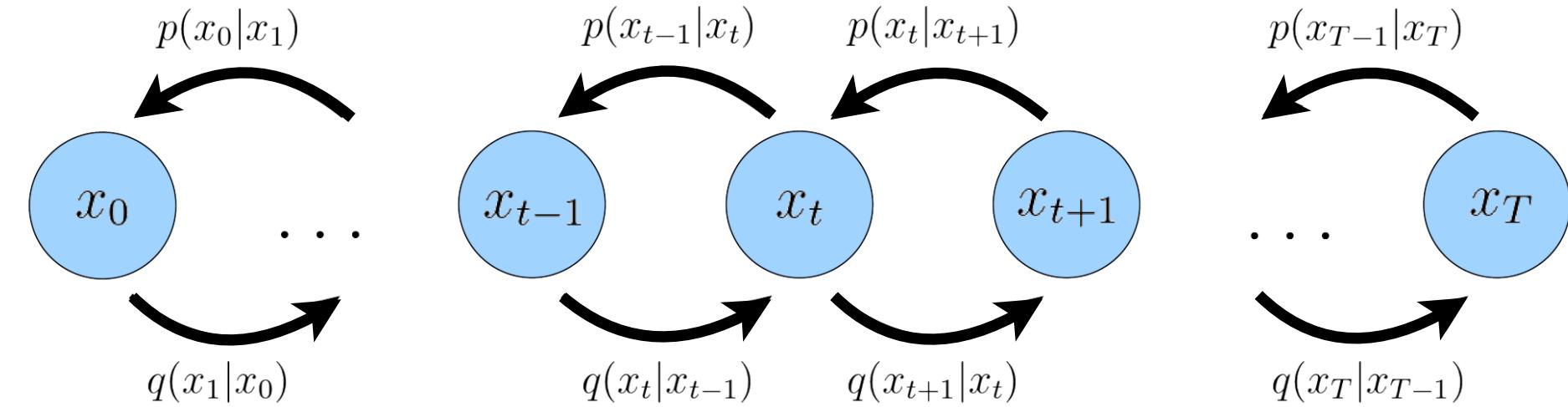


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

$$q(x_{t-1} | x_t, x_0) \propto \mathcal{N}(\mathbf{x}_{t-1}; \underbrace{\frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}, \underbrace{\frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{I}}_{\Sigma_q(t)}}}_{\mu_q(\mathbf{x}_t, \mathbf{x}_0)}$$

- Following the KL optimization term, in order to ensure  $p_{\theta}(x_{t-1} | x_t)$  approximates the ground-truth transition  $q(x_{t-1} | x_t, x_0)$  as closely as possible, it is modelled as a Gaussian.
  - The mean is parameterized as the function  $\mu_{\theta}(x_t, t)$
  - The variance is set to be the same as  $\Sigma_q(t)$

# Variational Diffusion Models: Simplifying Training

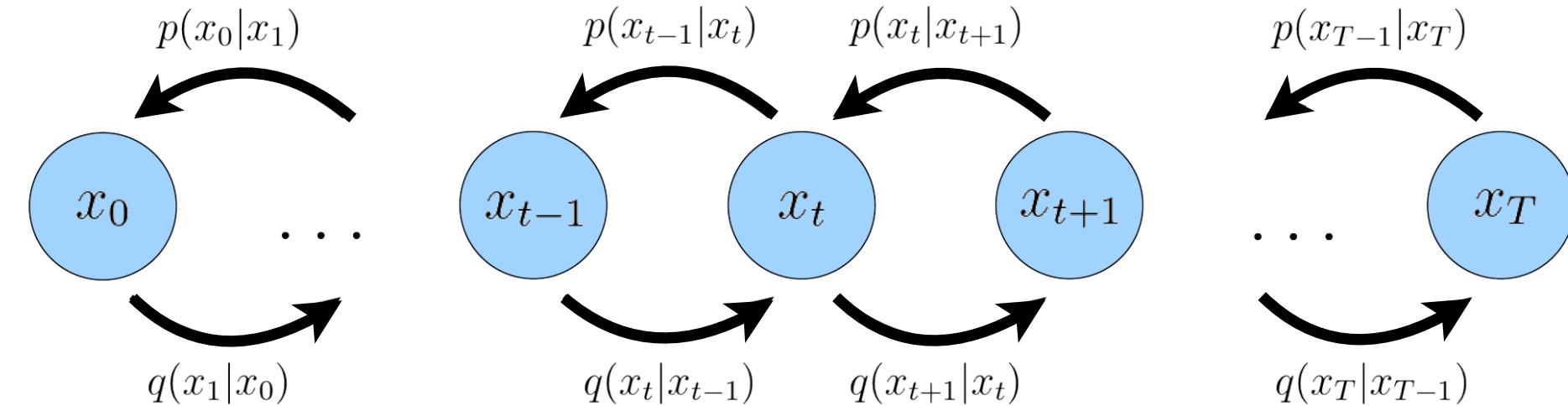


$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}$$

denoising matching term

- As it is ensured that the variances of  $p_{\theta}(x_{t-1}|x_t)$  and  $q(x_{t-1}|x_t, x_0)$  match exactly, optimizing the KL divergence can be simplified to,

# Variational Diffusion Models: Simplifying Training



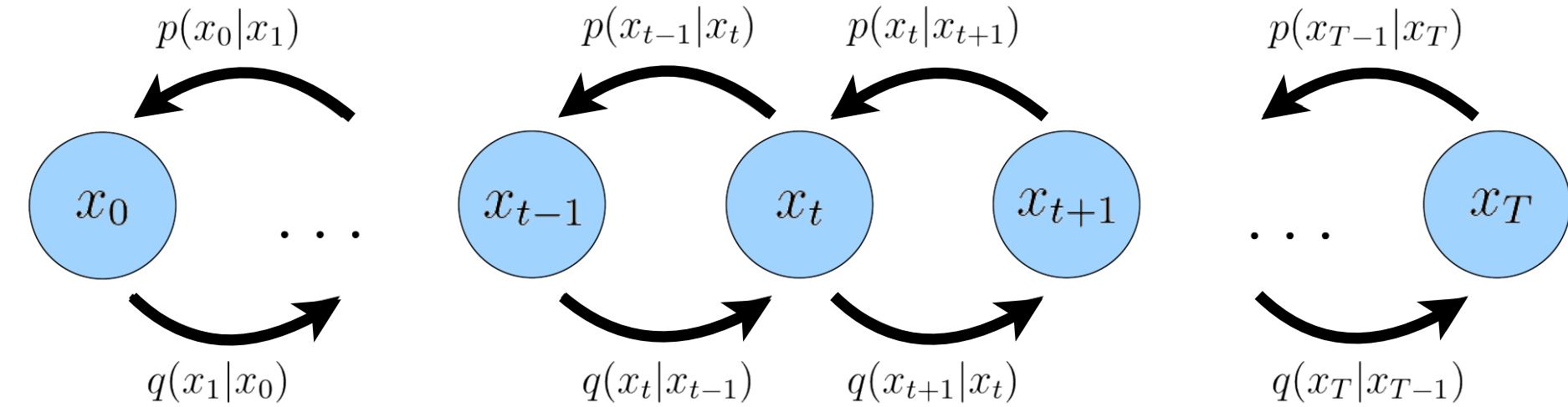
$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

- As it is ensured that the variances of  $p_{\theta}(x_{t-1}|x_t)$  and  $q(x_{t-1}|x_t, x_0)$  match exactly, optimizing the KL divergence can be simplified to,

$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} [\|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2]$$

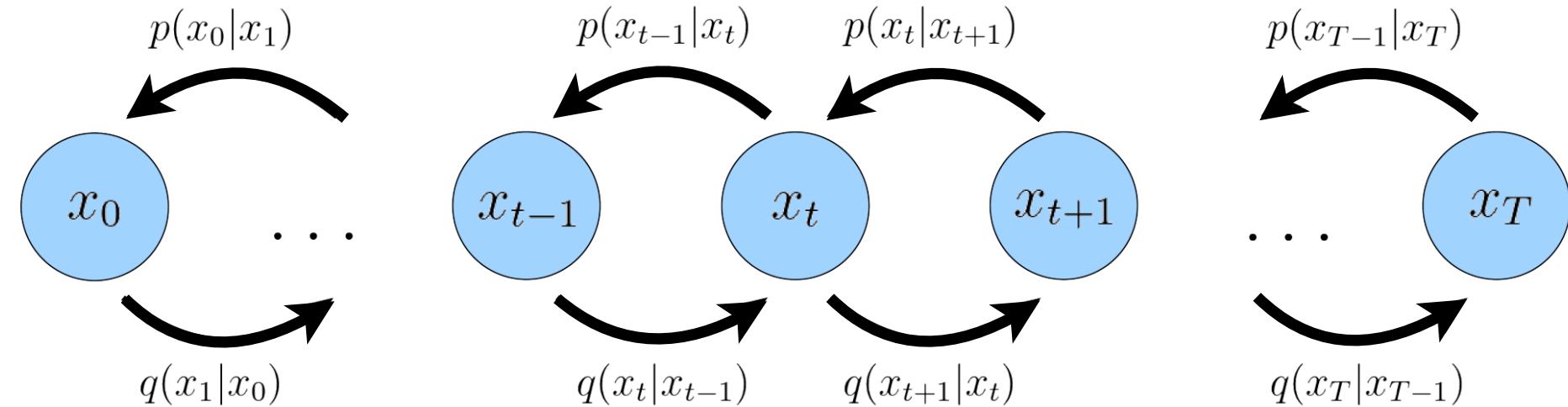
Reduces to minimizing the difference between the means of the two distributions

# Variational Diffusion Models: Simplifying Training



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

# Variational Diffusion Models: Simplifying Training

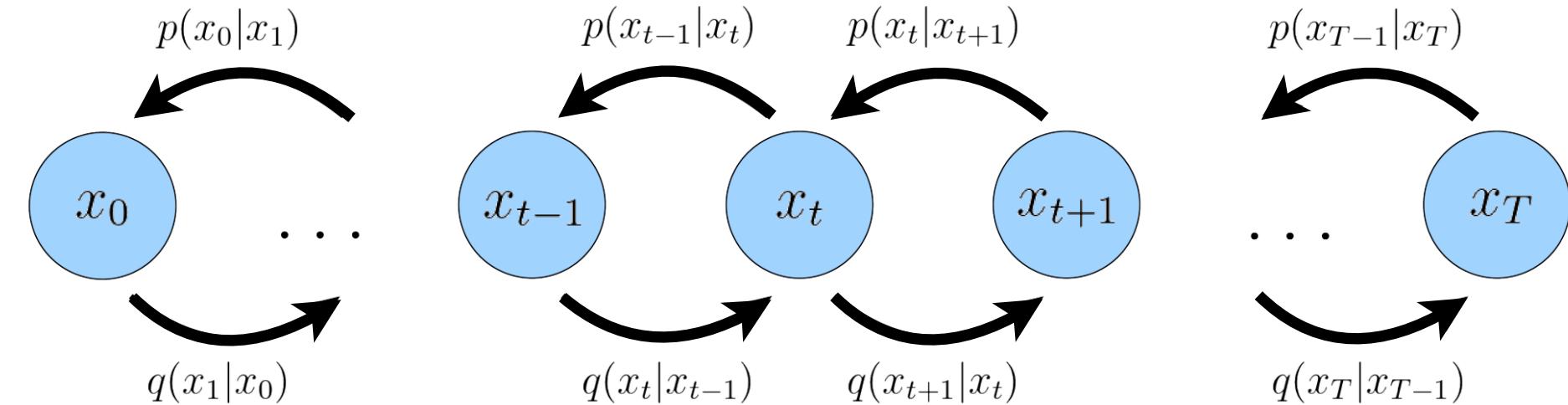


$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

From derivation shown earlier

# Variational Diffusion Models: Simplifying Training



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

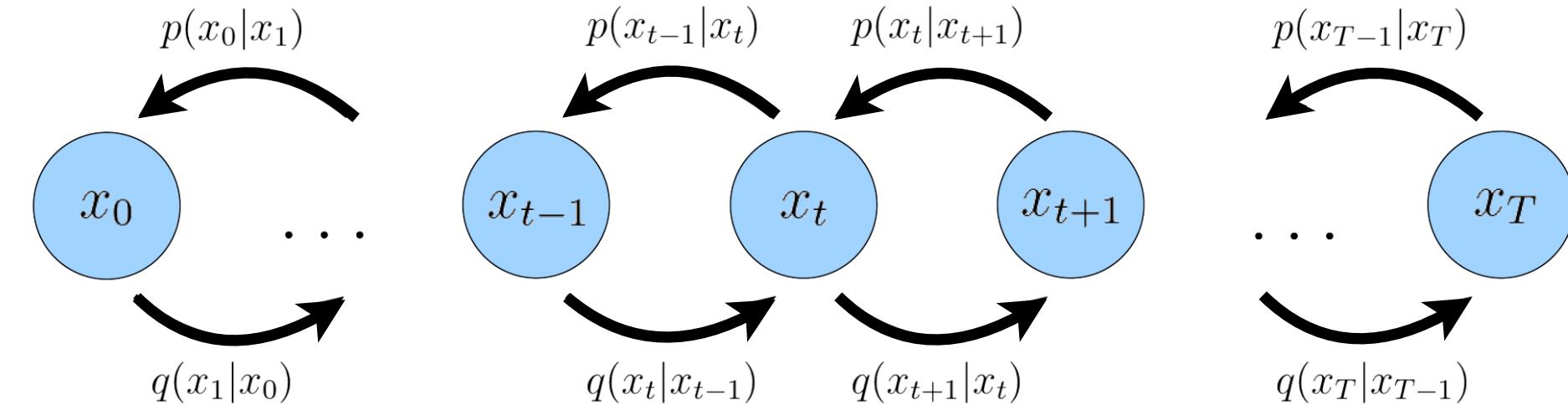
$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

From derivation shown earlier

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Defined to match the above mean closely

# Variational Diffusion Models: Simplifying Training



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\mathbf{x}_0}{1 - \bar{\alpha}_t}$$

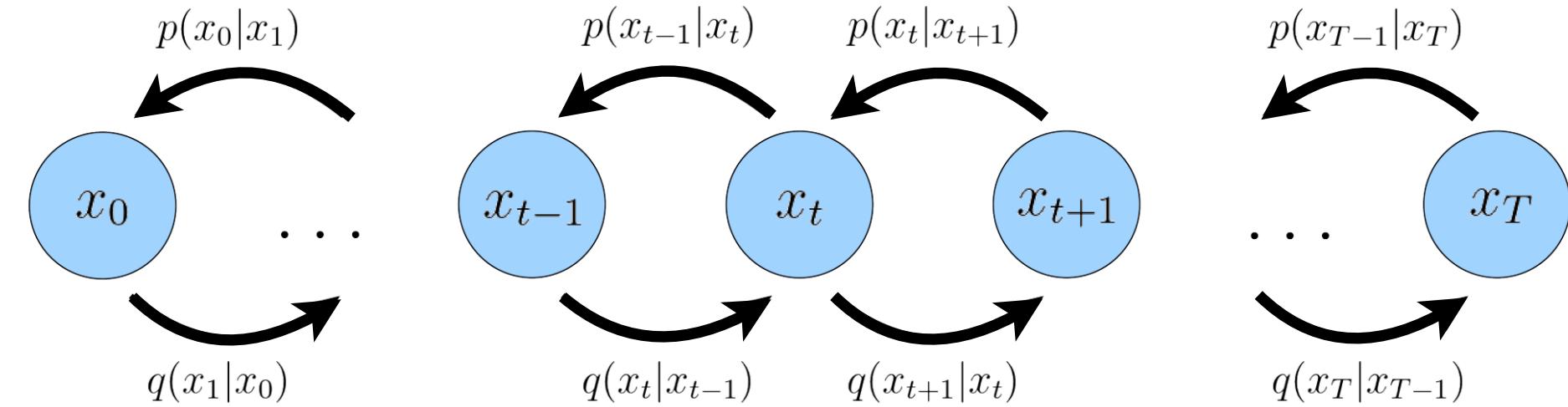
From derivation shown earlier

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{x}_{\theta}(\mathbf{x}_t, t)}{1 - \bar{\alpha}_t}$$

Defined to match the above mean closely

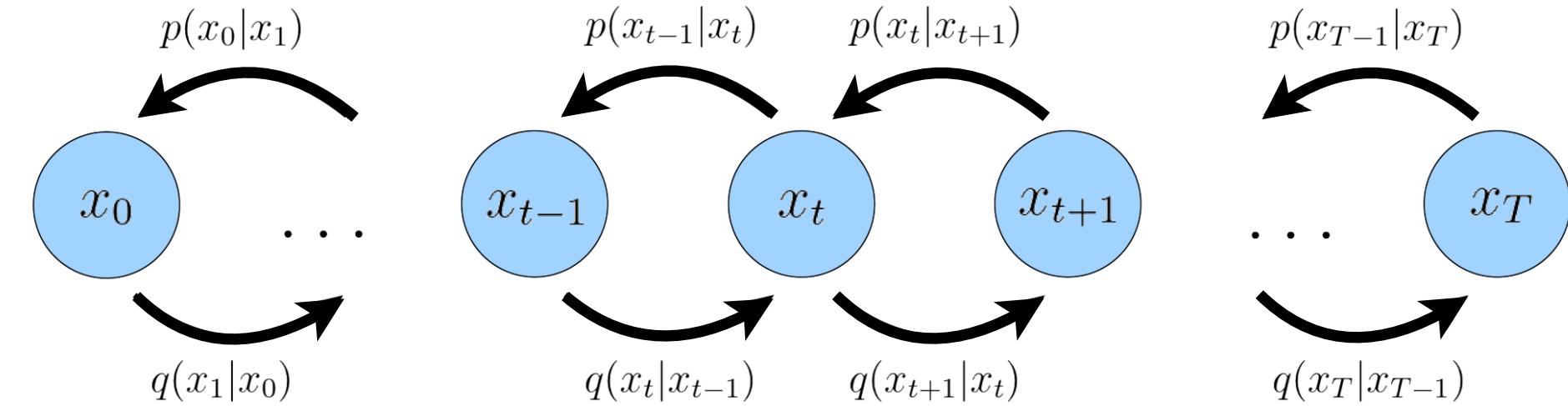
- $\hat{x}_{\theta}(\mathbf{x}_t, t)$  is parametrized by a neural network.
  - Predicts the original input  $\mathbf{x}_0$  from a noised input  $\mathbf{x}_t$  at time  $t$

# Variational Diffusion Models: Simplifying Training



$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right] \end{aligned}$$

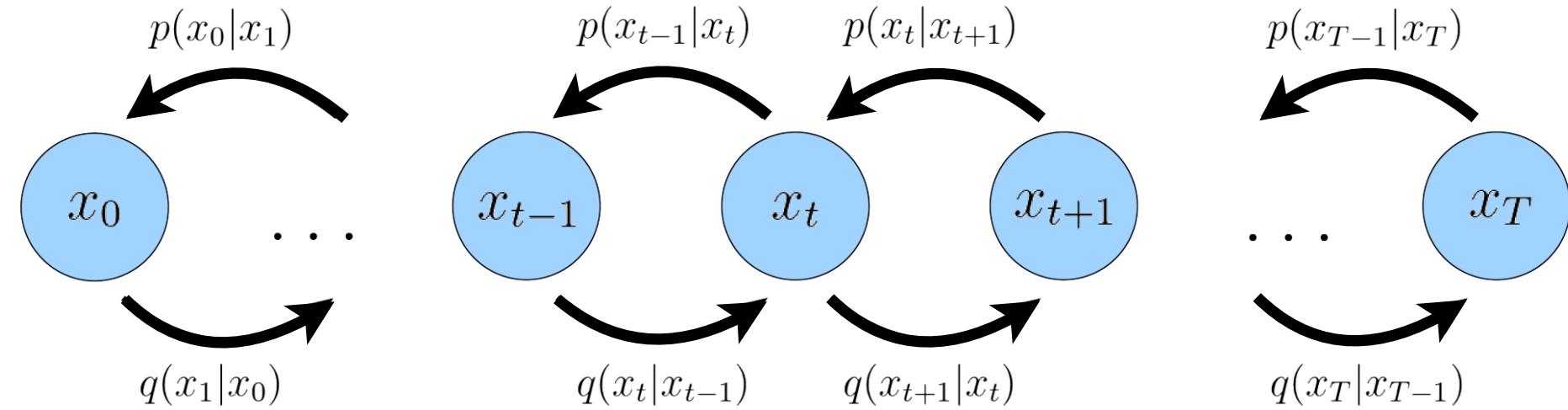
# Variational Diffusion Models: Simplifying Training



$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)^2} \left[ \|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|_2^2 \right] \end{aligned}$$

Optimizing a diffusion model simplifies to learning a neural network to predict the original ground truth input from an arbitrary noised version of it

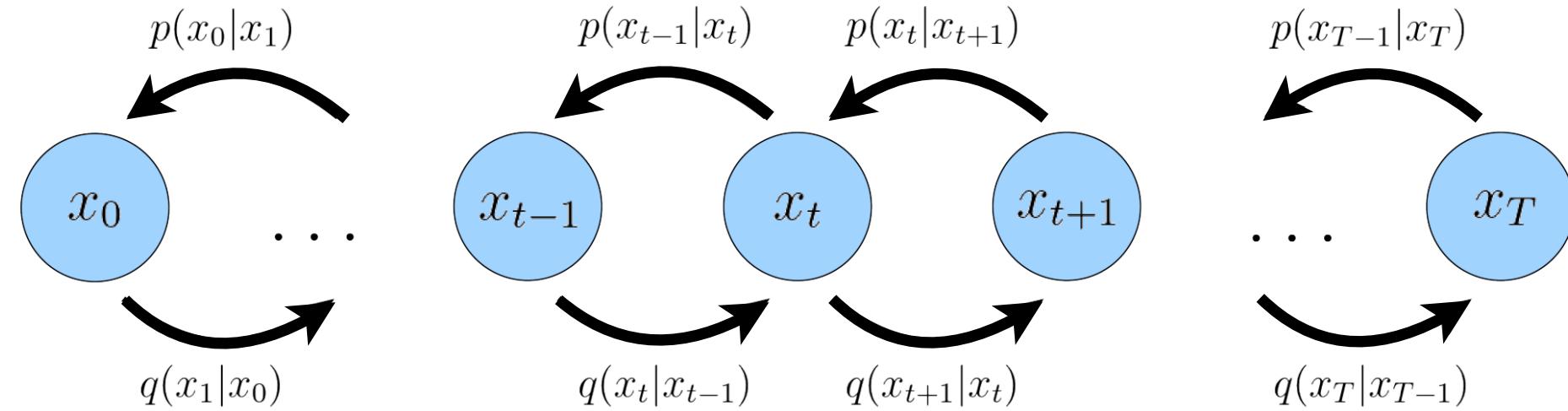
# Variational Diffusion Models: Equivalent View



$$x_t \sim q(x_t | x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$

$$\sim \mathcal{N}\left(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)\mathbf{I}\right)$$

# Variational Diffusion Models: Equivalent View



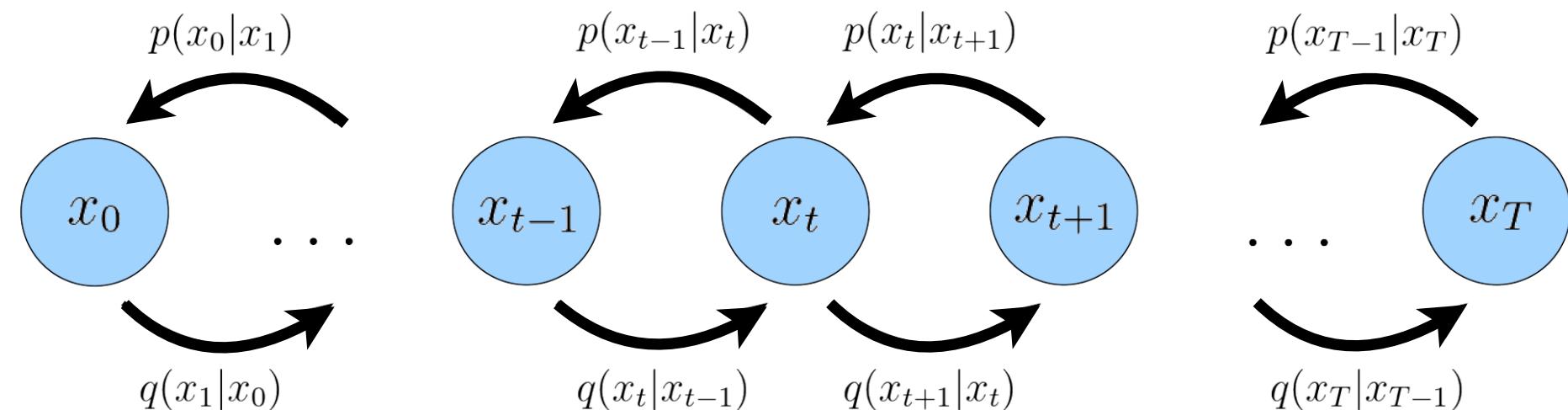
$$x_t \sim q(x_t | x_0) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*$$

$$\sim \mathcal{N}\left(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}\right)$$

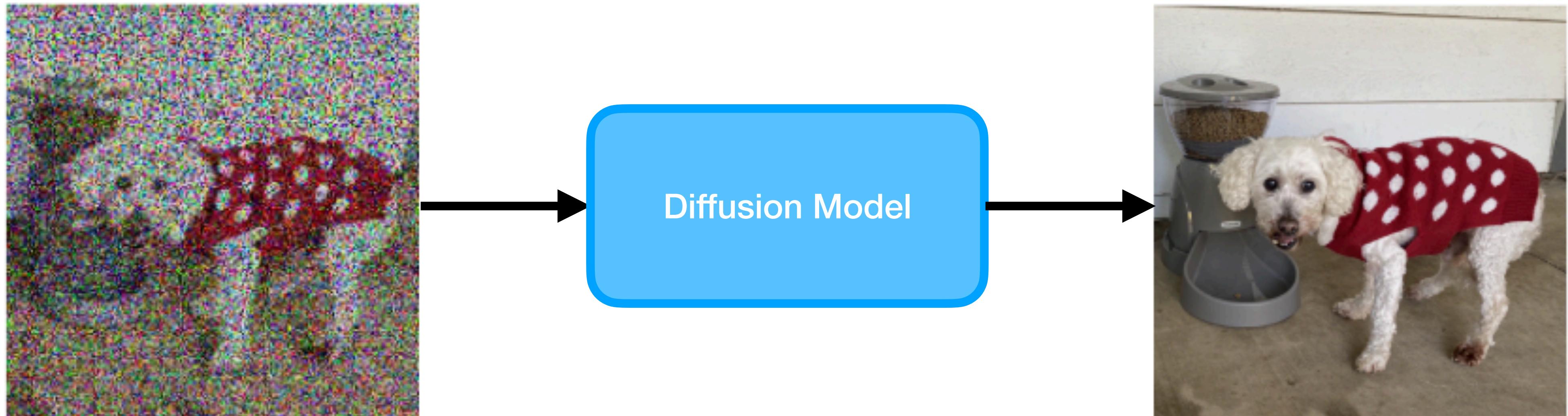
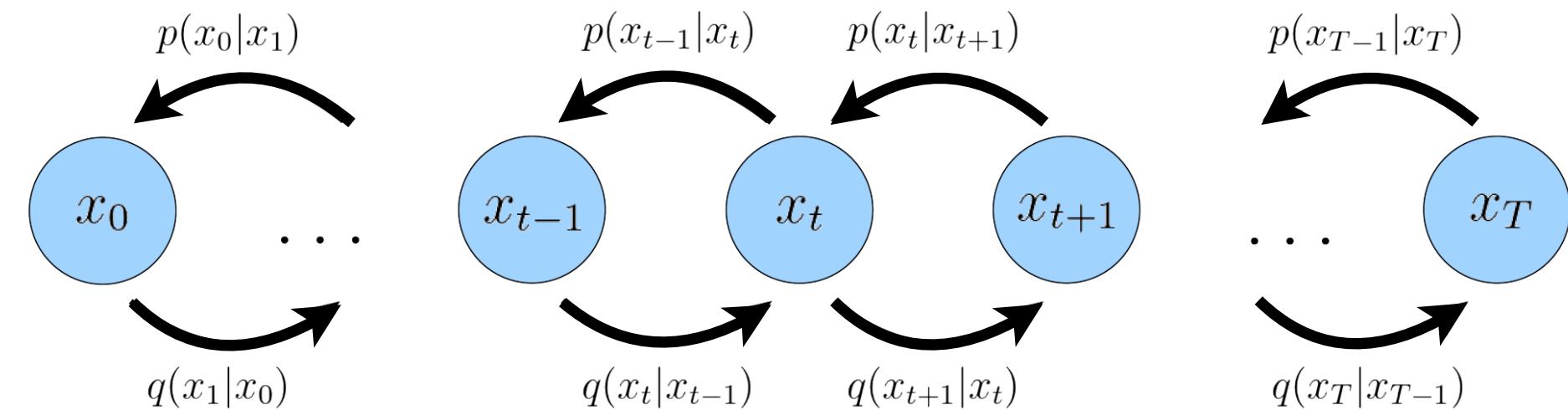
$$x_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_0^*}{\sqrt{\bar{\alpha}_t}}$$

Knowing the noise added and the noisy input is sufficient to give the denoised input

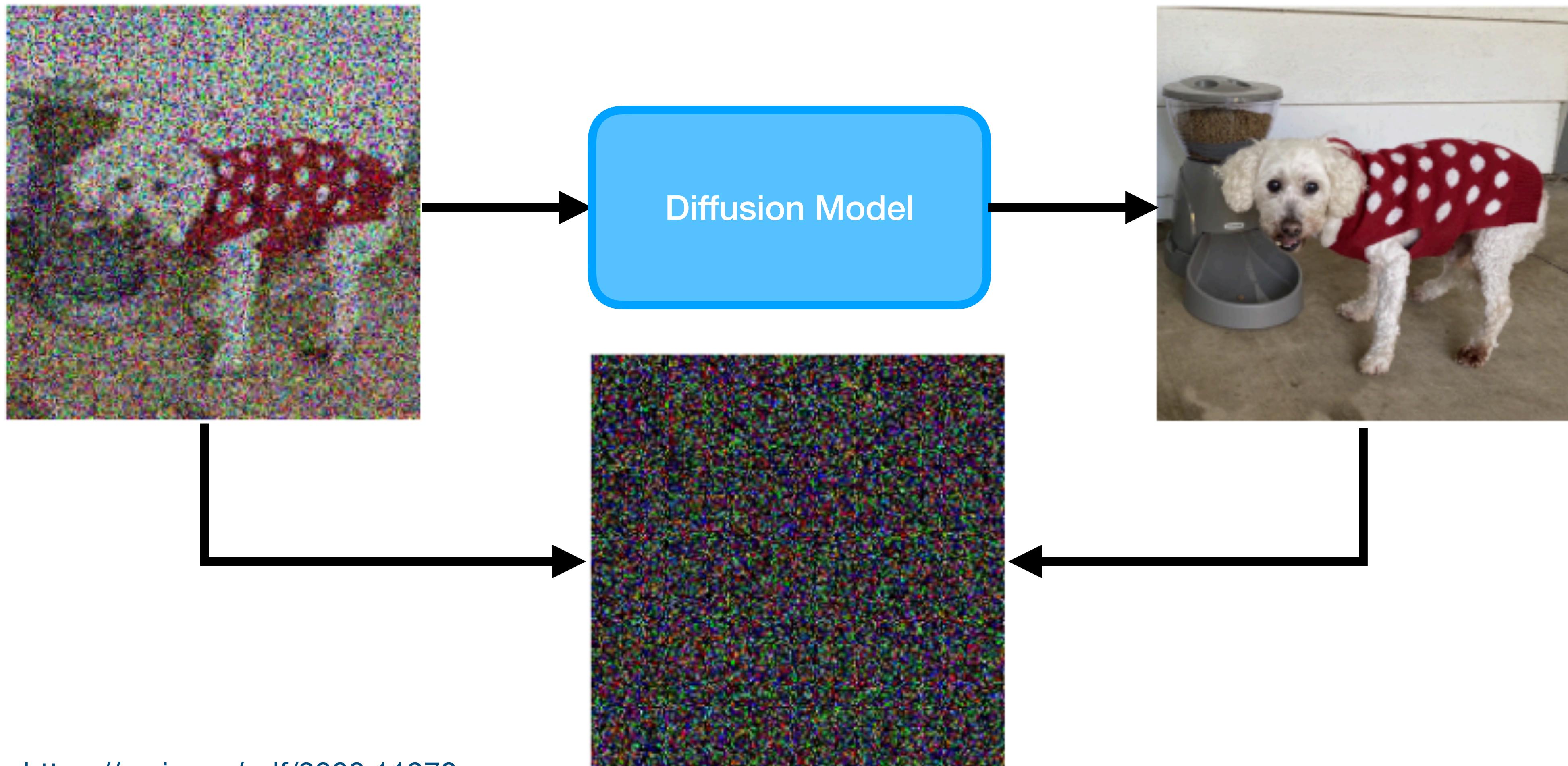
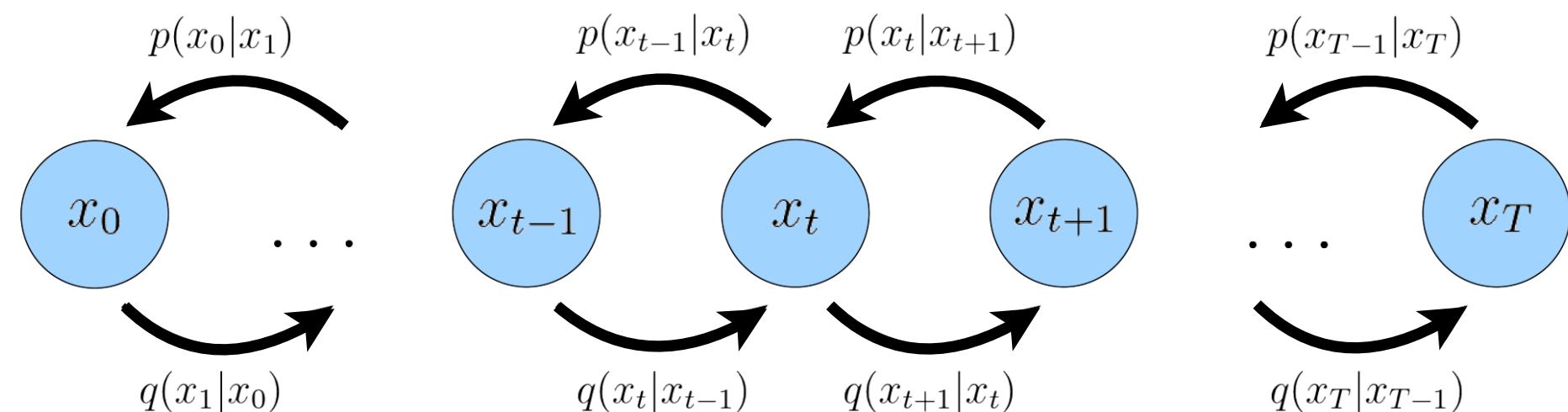
# Variational Diffusion Models: Equivalent View



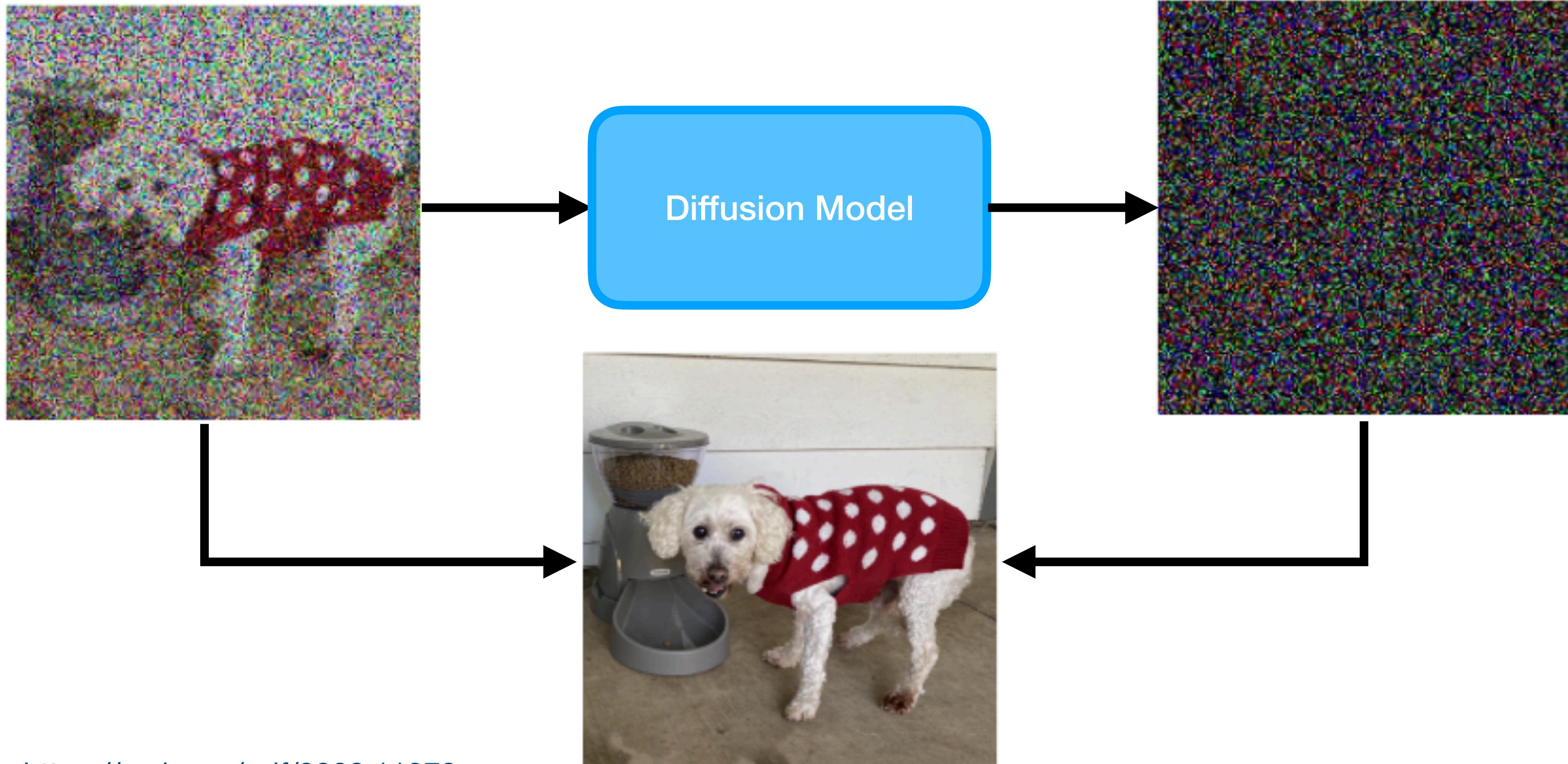
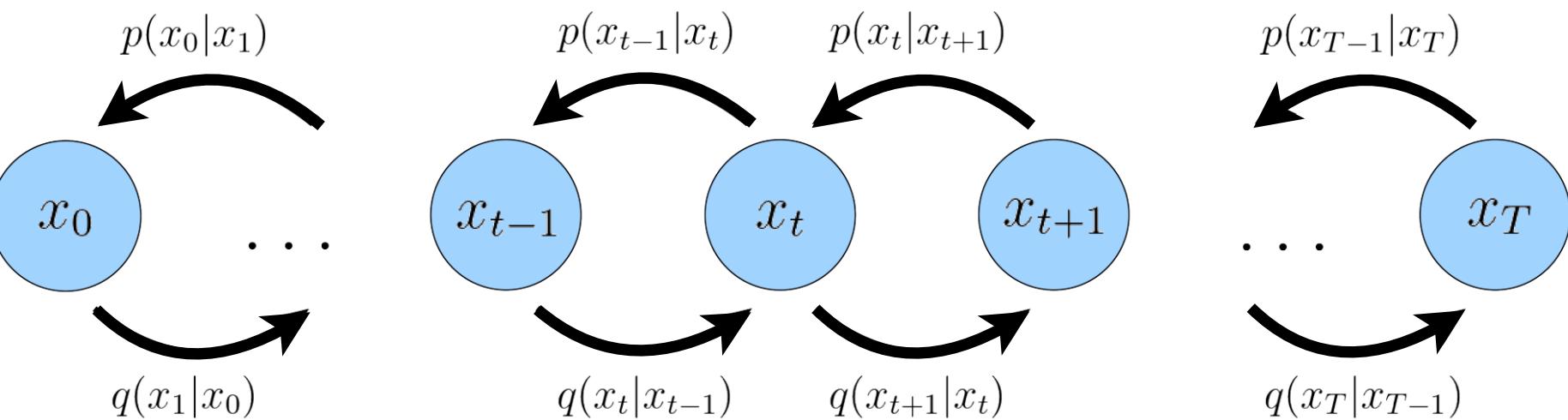
# Variational Diffusion Models: Equivalent View



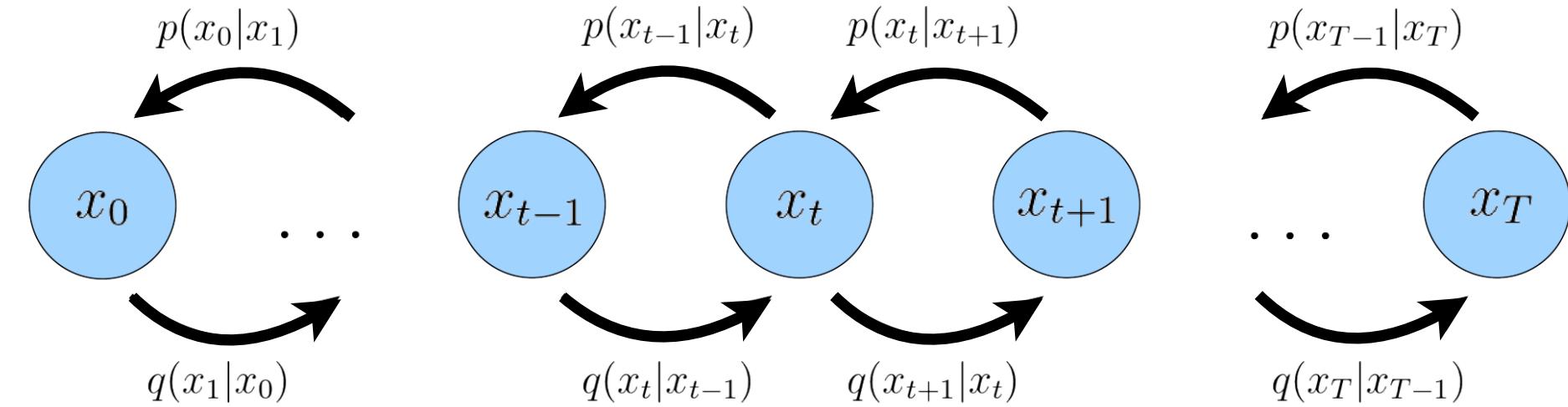
# Variational Diffusion Models: Equivalent View



# Variational Diffusion Models: Equivalent View

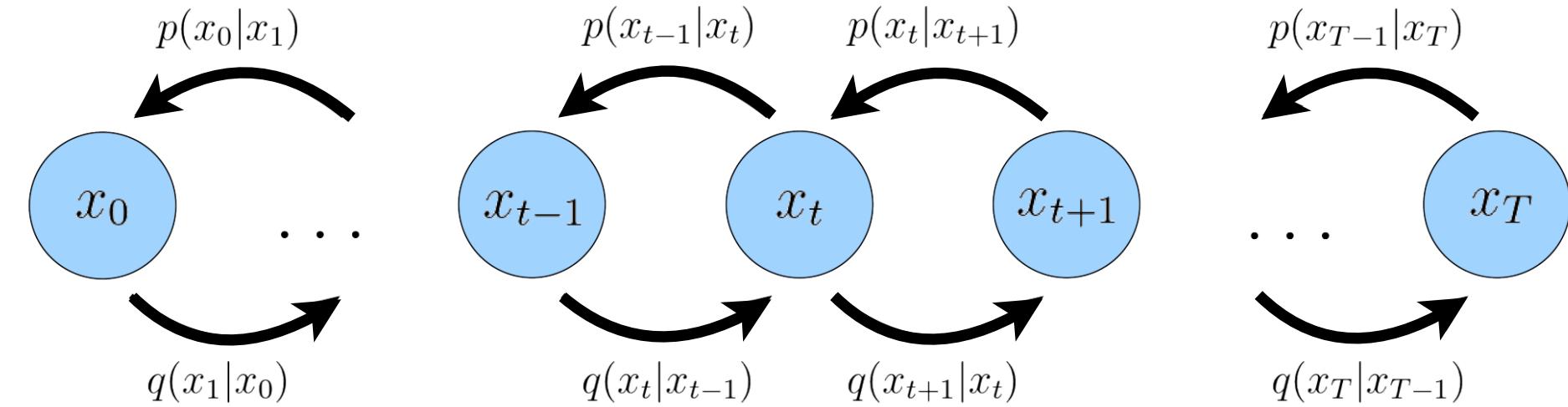


# Variational Diffusion Models: Equivalent View



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

# Variational Diffusion Models: Equivalent View

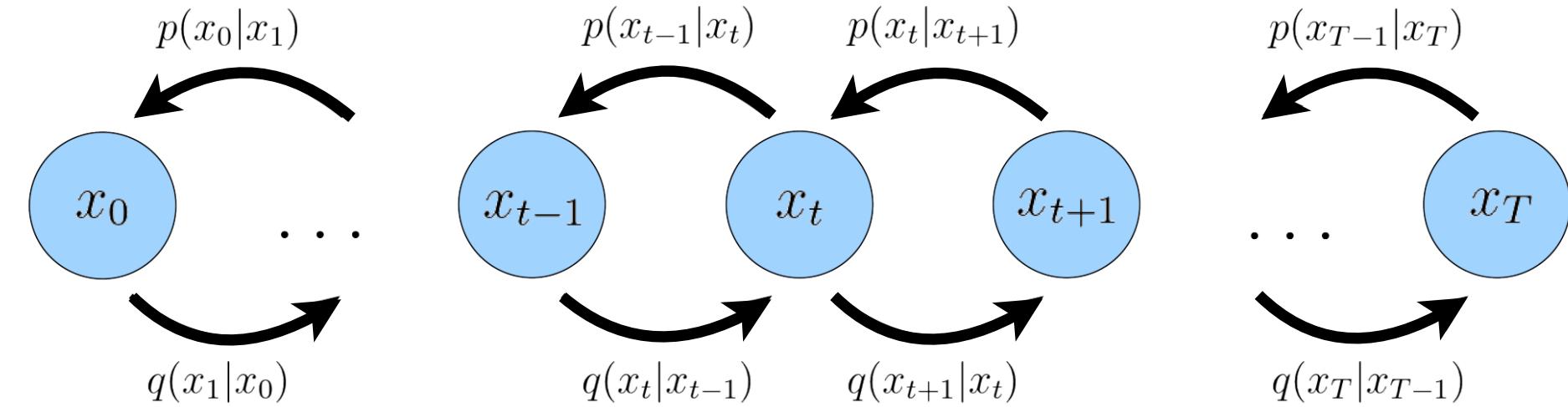


$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0$$

Alternate Formulation using equivalence

# Variational Diffusion Models: Equivalent View



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

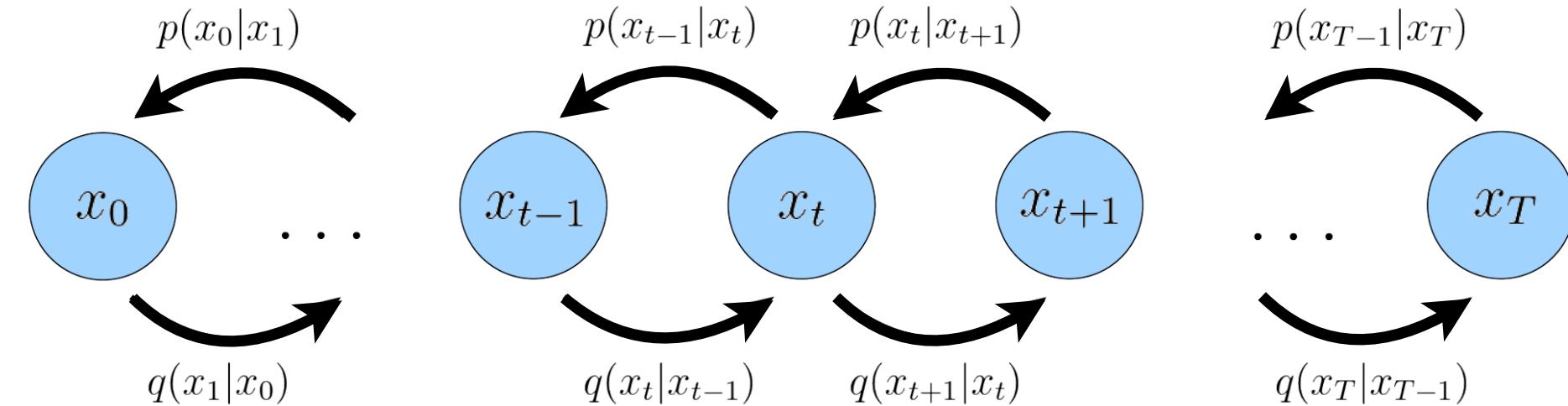
$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0$$

Alternate Formulation using equivalence

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)$$

Defined to match the above mean closely

# Variational Diffusion Models: Equivalent View



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \boldsymbol{\epsilon}_0$$

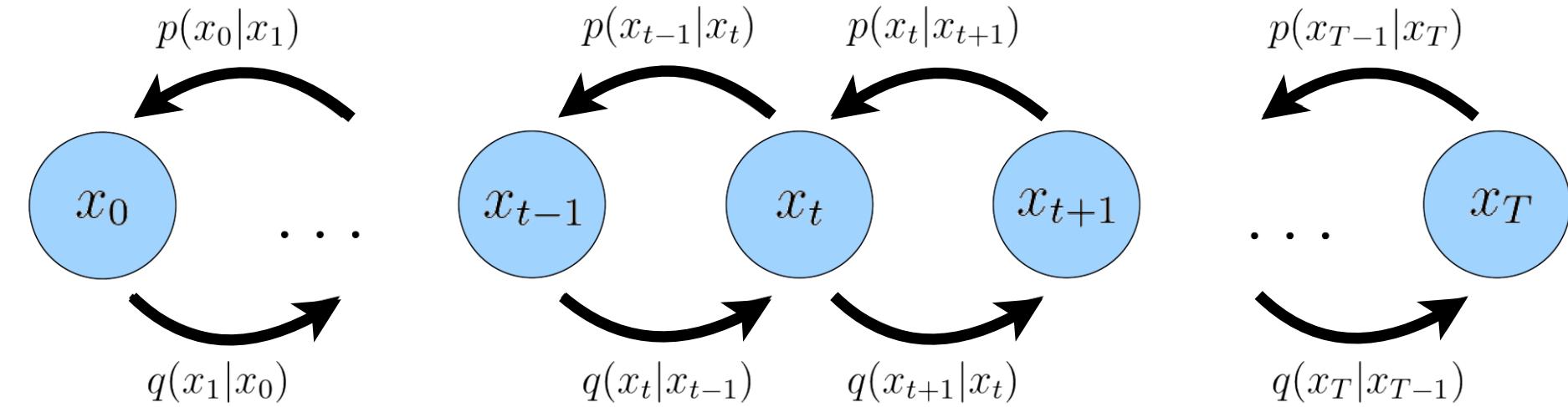
Alternate Formulation using equivalence

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)$$

Defined to match the above mean closely

- $\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)$  is parametrized by a neural network.
- Predicts the added noise  $\boldsymbol{\epsilon}_0^*$  from a noised input  $\mathbf{x}_t$  at time t

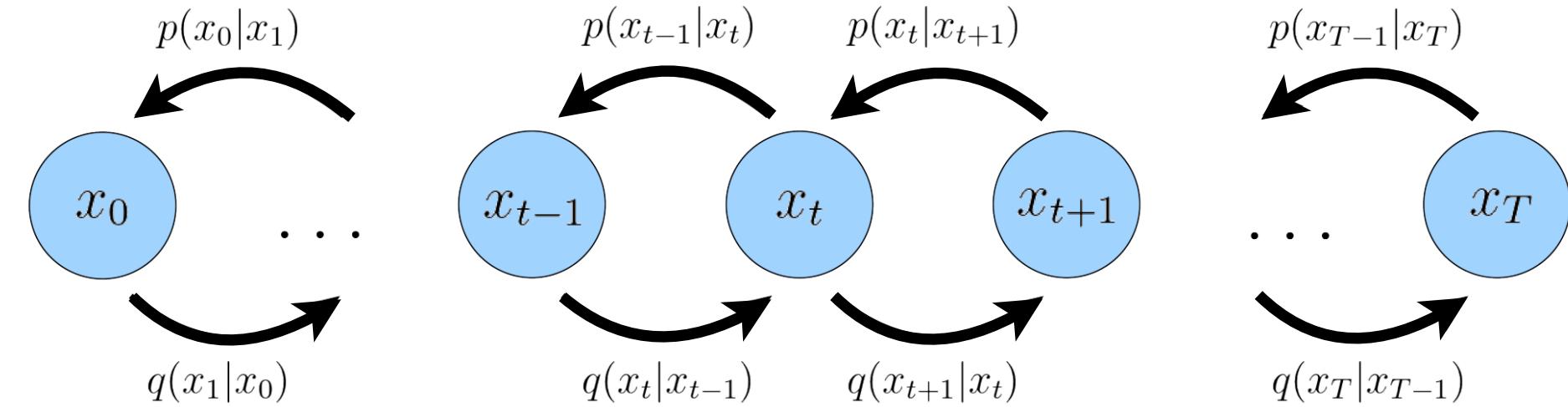
# Variational Diffusion Models: Equivalent View



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

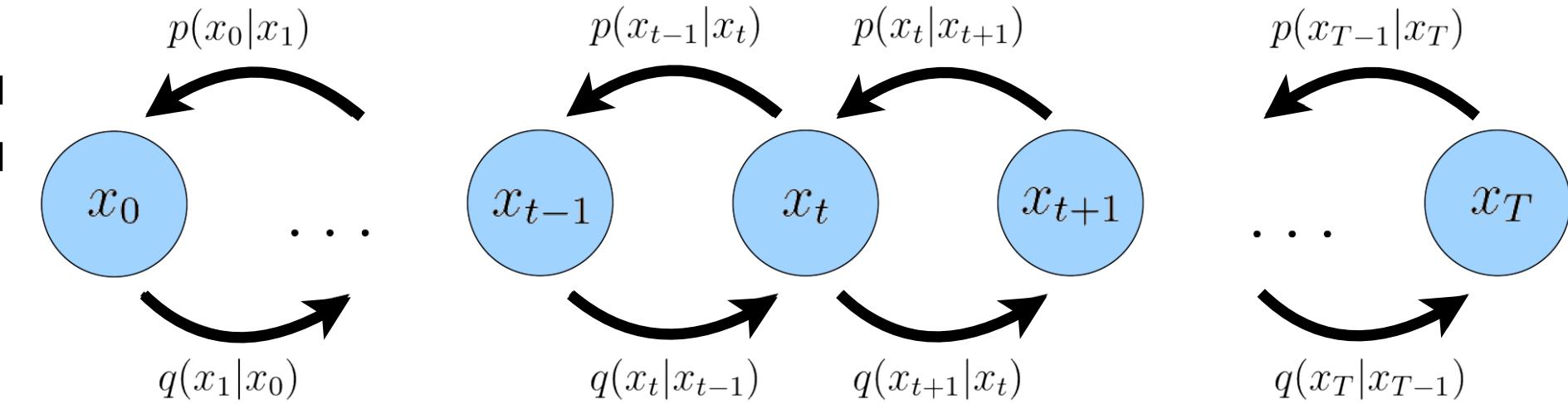
# Variational Diffusion Models: Equivalent View



$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \left[ \|\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \end{aligned}$$

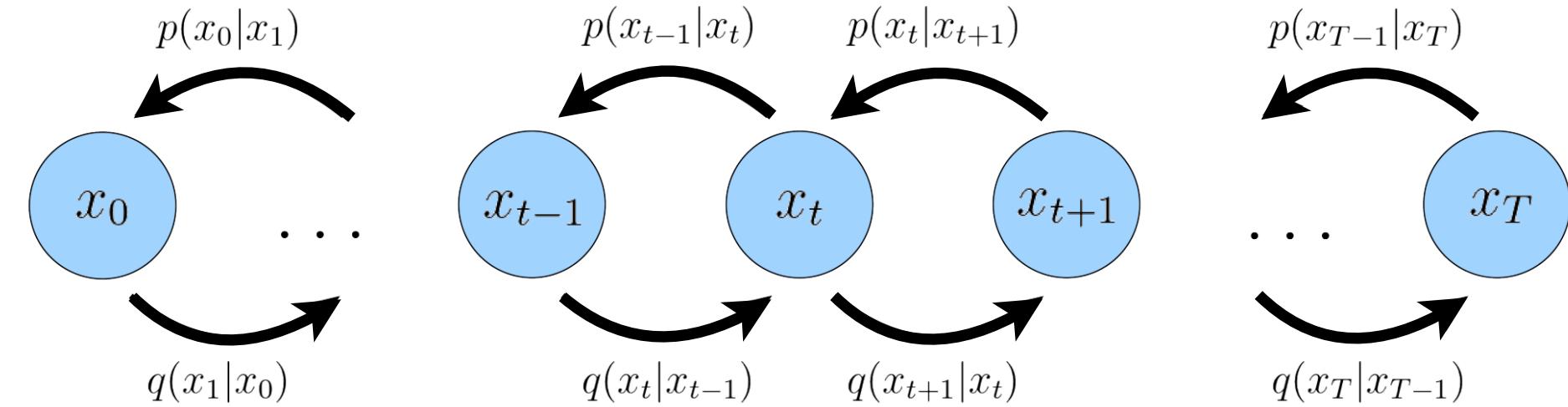
Equivalently, optimizing a diffusion model can also be simplified to learning a neural network to predict the source noise added to  $\mathbf{x}_0$  to get  $\mathbf{x}_t$

# Variational Diffusion Models: Third Interpretation



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

# Variational Diffusion Models: Third Interpretation

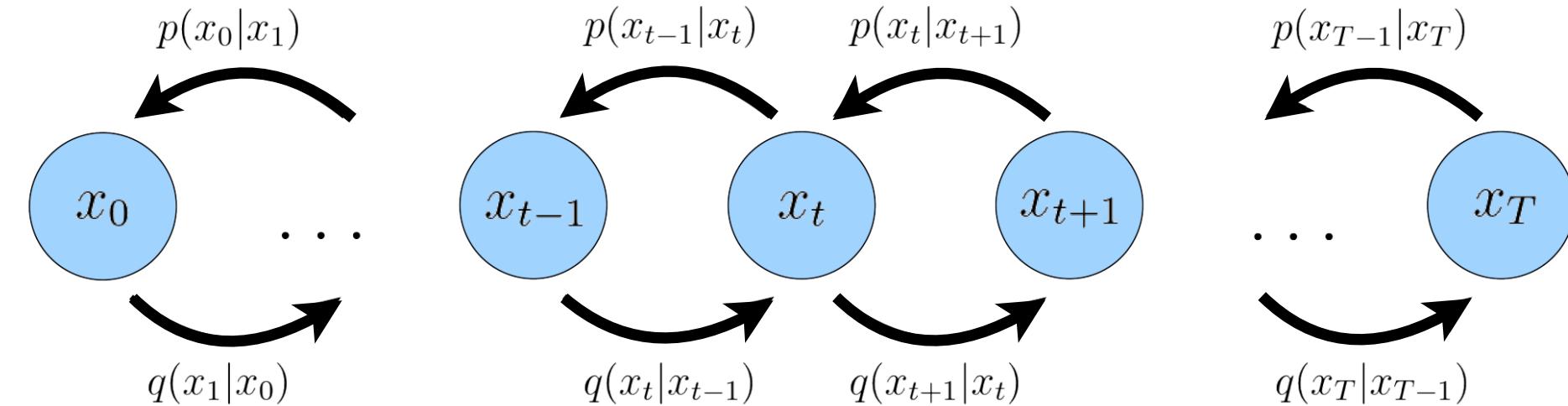


$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t)$$

Alternate Formulation using score function

# Variational Diffusion Models: Third Interpretation



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

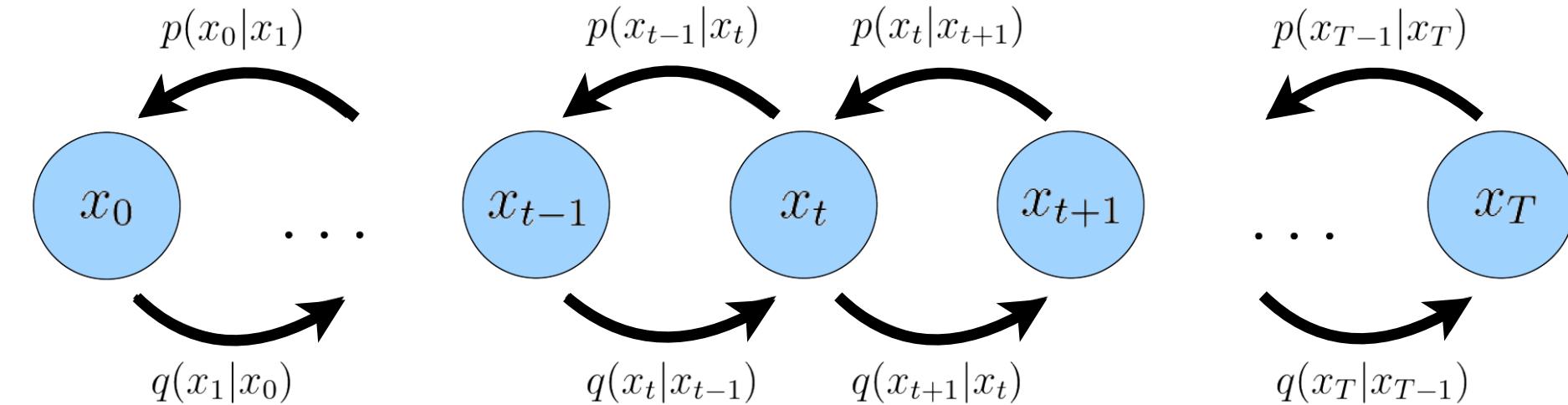
$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t)$$

Alternate Formulation using score function

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \mathbf{s}_{\theta}(\mathbf{x}_t, t)$$

Defined to match the above mean closely

# Variational Diffusion Models: Third Interpretation



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla \log p(\mathbf{x}_t)$$

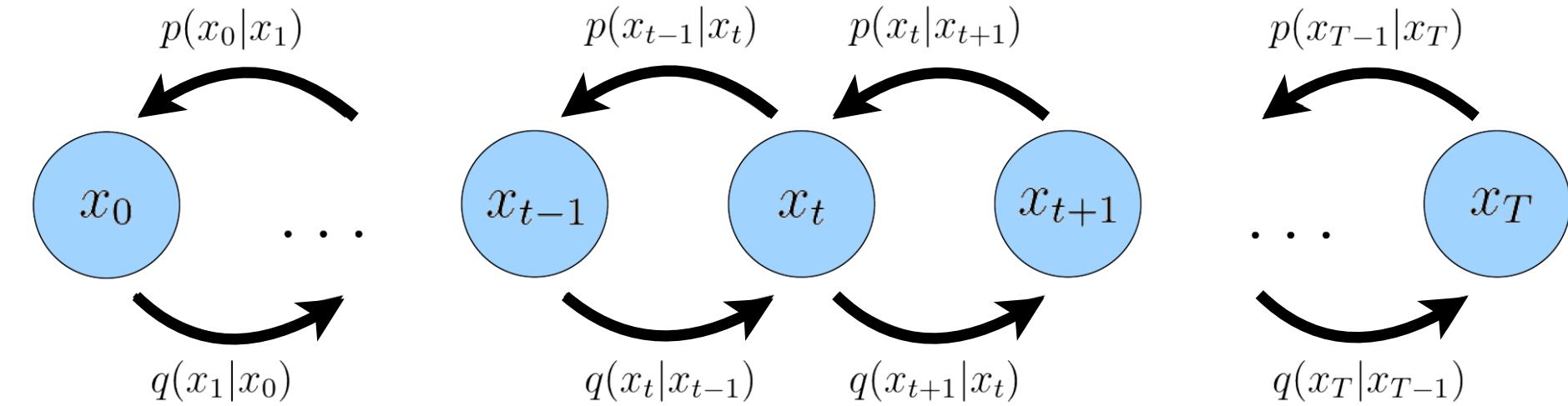
Alternate Formulation using score function

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \mathbf{s}_{\theta}(\mathbf{x}_t, t)$$

Defined to match the above mean closely

- $s_{\theta}(x_t, t)$  is parametrized by a neural network.
- Predicts the score function  $\nabla \log p(x_t)$  from a noised input  $\mathbf{x}_t$  at time  $t$

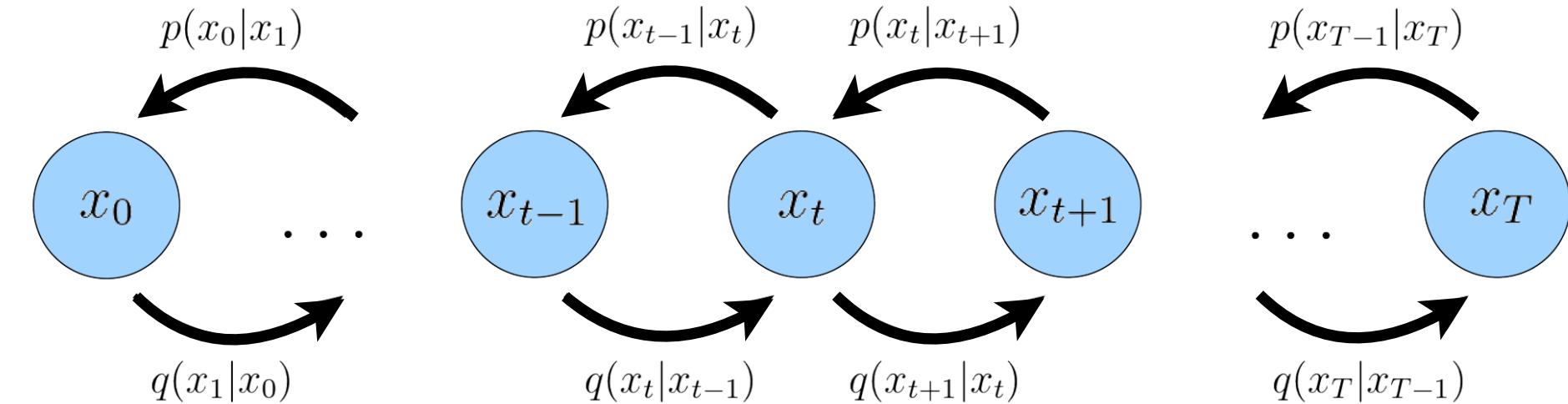
# Variational Diffusion Models: Third Interpretation



$$\arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[ \|\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q\|_2^2 \right]$$

# Variational Diffusion Models: Third Interpretation

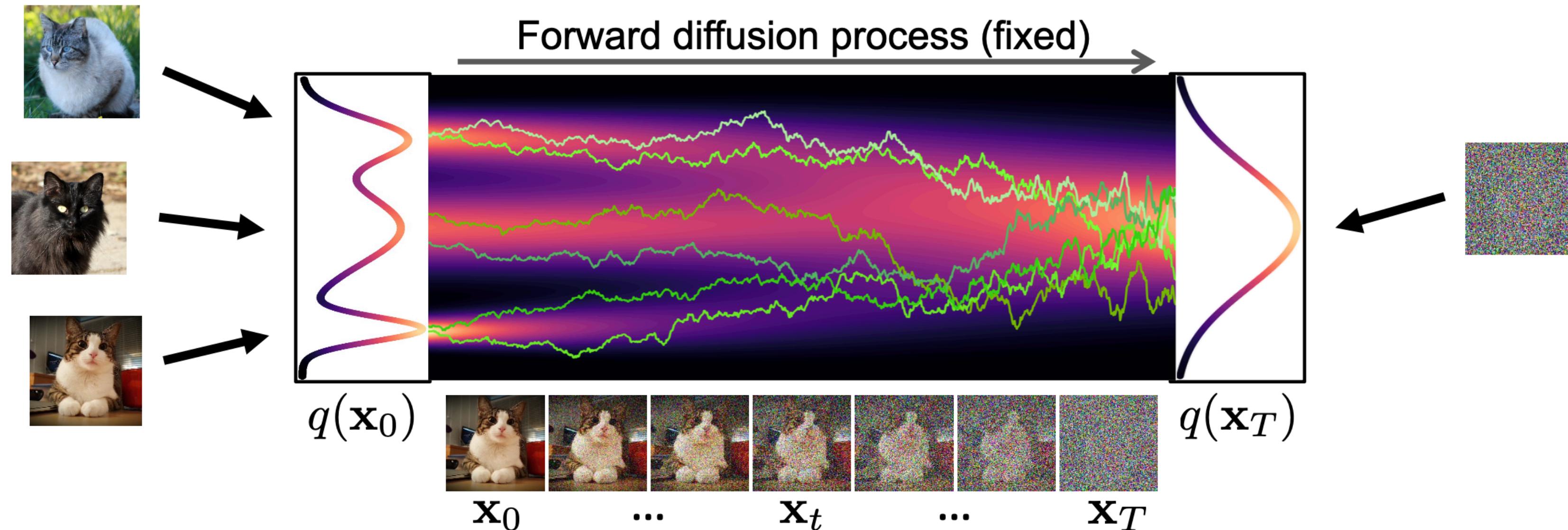


$$\begin{aligned} & \arg \min_{\theta} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)) \\ &= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \left[ \|s_{\theta}(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|_2^2 \right] \end{aligned}$$

Equivalently, optimizing a diffusion model can also be simplified to learning a neural network to predict the score function, which is the gradient of  $\mathbf{x}_t$  in data space, for a given  $t$ .

# Alternative View of Diffusion Processes

- Both forward and reverse diffusion processes can be written as stochastic differential equations (SDEs)
- Reverse SDE governs warping from known simple noise distribution to the data distribution (can also be formulated using ODE)



# **Introduction to Diffusion Models**

**Part 2: Practice and Advance Topics**

**Leonid Sigal and Siddhesh Khandelwal**

# Practical Summary: Training

---

## 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  
      
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
  - 6: **until** converged
-

# Practical Summary: Sampling

---

## Algorithm 2 Sampling

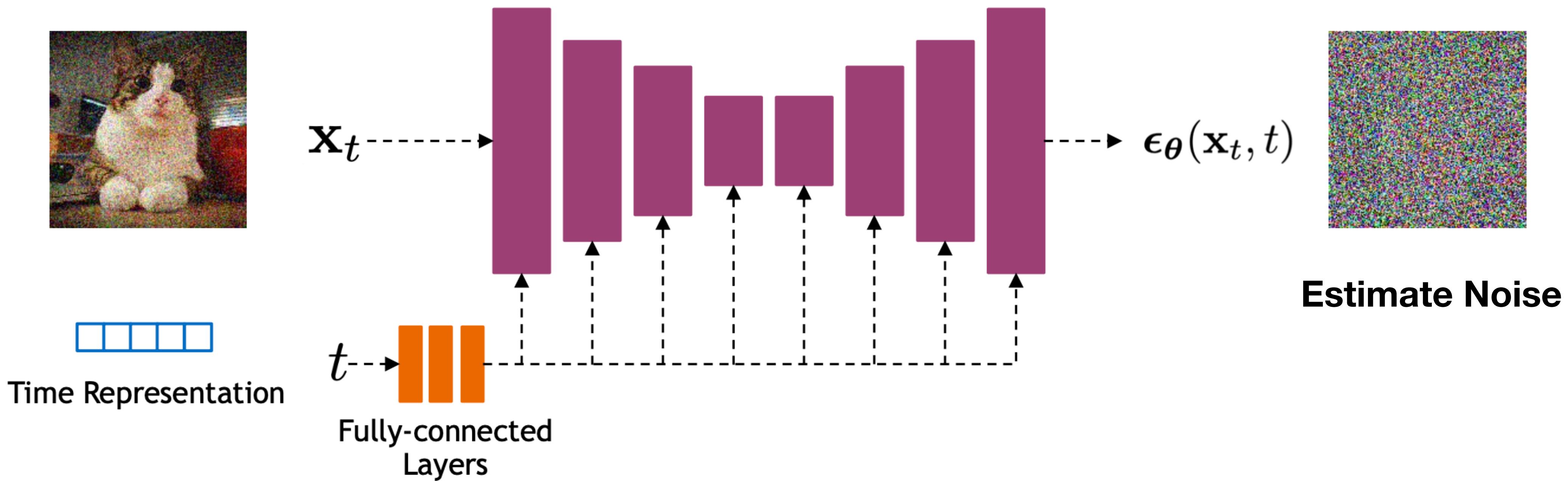
---

- 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**  $\mathbf{x}_0$

---

# Real Example: Image Generation

U-Net Neural Architecture



# Real Example: Training for Image Generation

For each training step:

1. Randomly select a time step & encode it



---

## 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
```

---

From Steins (medium.com)

- **Step 1:** Choose and encode time step (usually uniformly from 1 to  $T=1000$ )

# Real Example: Training for Image Generation

For each training step:

1. Randomly select a time step & encode it



	$\sin(t/1)$	$\cos(t/1)$	$\sin(t/10)$	$\cos(t/10)$	...
$t=14$	0.24	0.97	0.02	1.00	...

From Steins (medium.com)

- **Step 1:** Choose and encode time step (usually uniformly from 1 to  $T=1000$ )

---

## 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
```

---

# Real Example: Training for Image Generation

For each training step:

1. Randomly select a time step & encode it



2. Add noise to image

A diagram illustrating the addition of noise. It shows three images side-by-side: a "Noisy image" (a blurry photo of a landscape), an "Original image" (a clear photo of the same landscape), and a "Gaussian noise" image (a dark gray square). Below these, an equation shows the noisy image as the sum of the original image and Gaussian noise:  $\text{Noisy image} = \text{Original image} + \text{Gaussian noise}$ . A formula at the bottom shows the generation of the noisy image:  $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$ .

Adjust the amount of noise according to the time step  $t$

$$\varepsilon \sim \mathcal{N}(0, 1)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

From Steins (medium.com)

- **Step 2:** Add noise corresponding to sampled step  $t$ , this creates training data

---

## Algorithm 1 Training

---

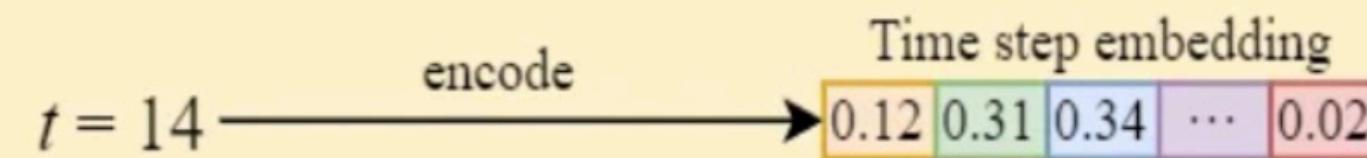
```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
      $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t}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

---

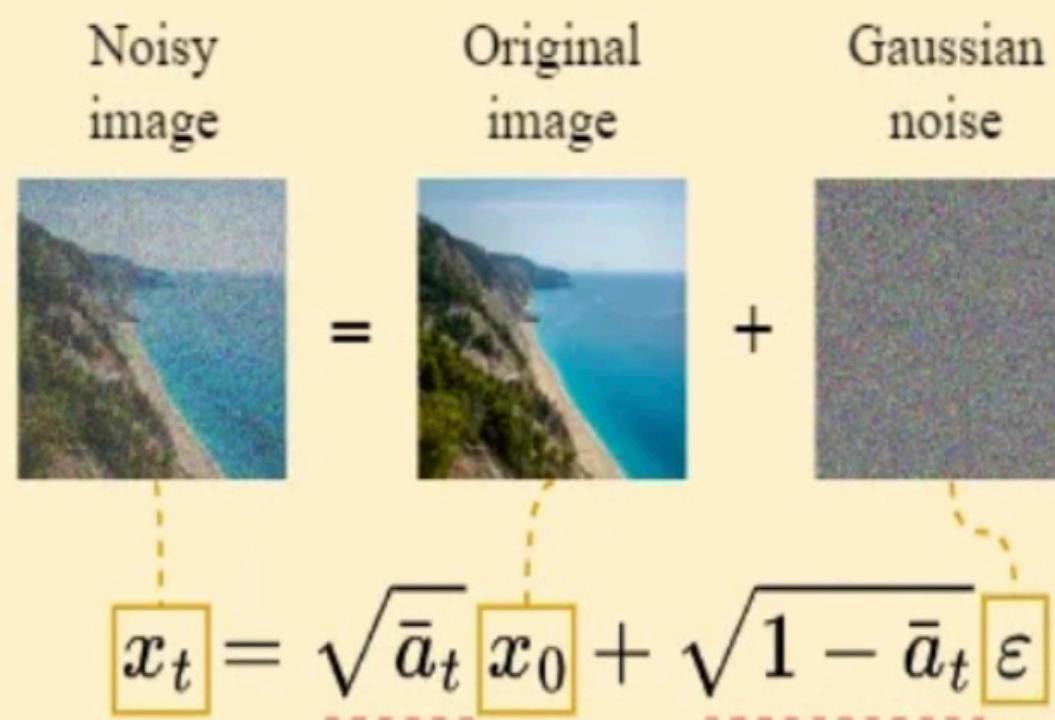
# Real Example: Training for Image Generation

For each training step:

1. Randomly select a time step & encode it



2. Add noise to image



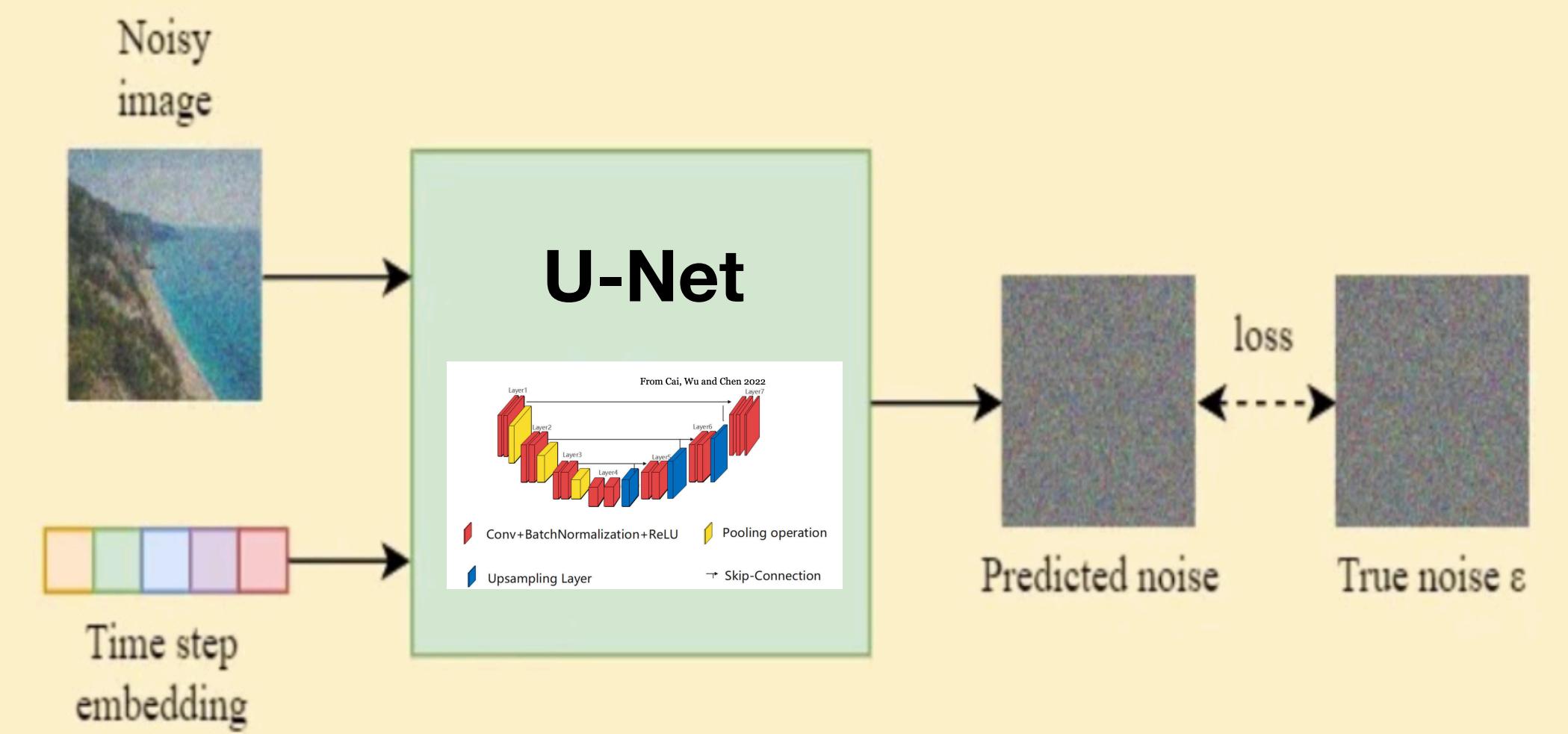
Adjust the amount of noise according to the time step  $t$

$$\varepsilon \sim \mathcal{N}(0, 1)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

3. Train the UNet



From Steins (medium.com)

\*step encoding simply added to the noisy sample representation

- **Step 3:** Train network to recover noise from corrupted sample

---

## Algorithm 1 Training

---

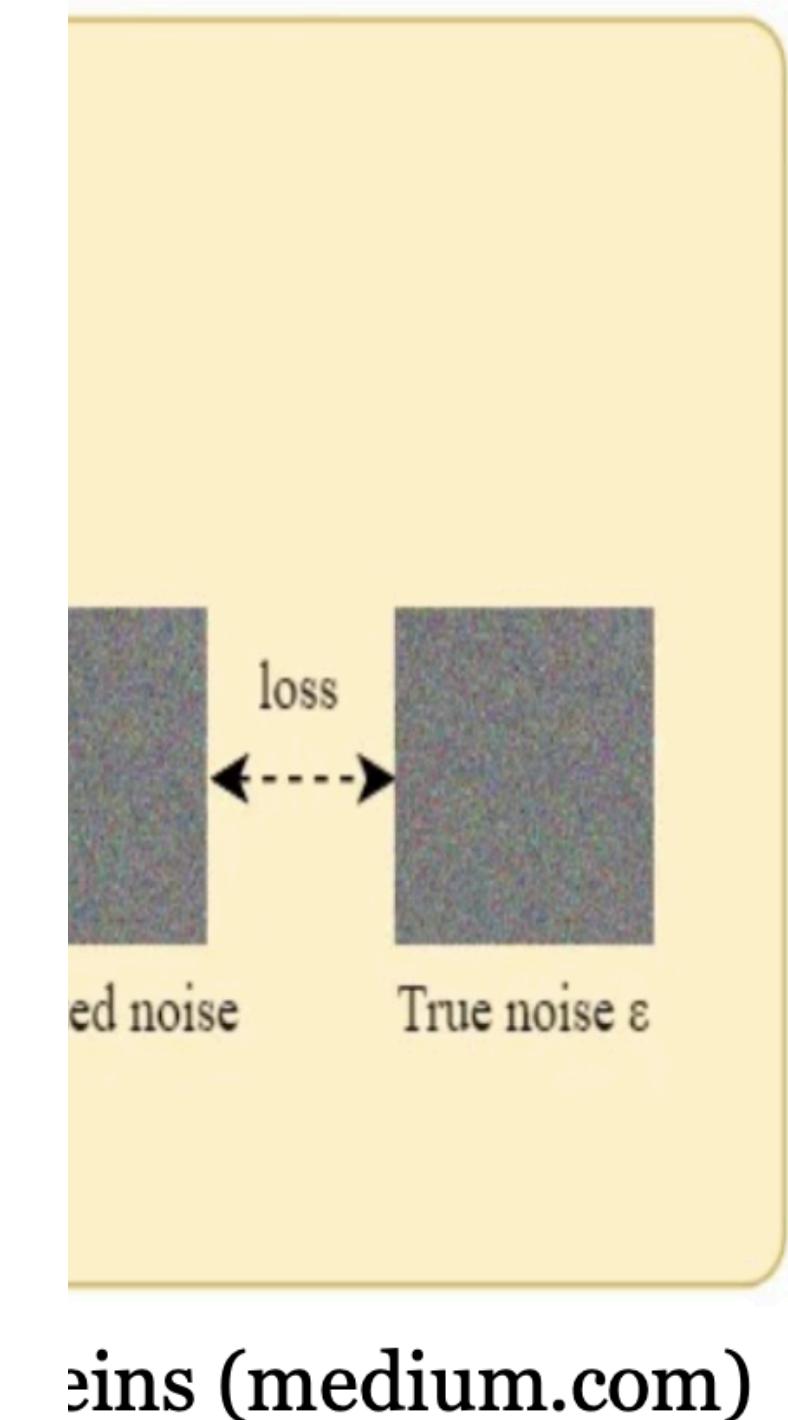
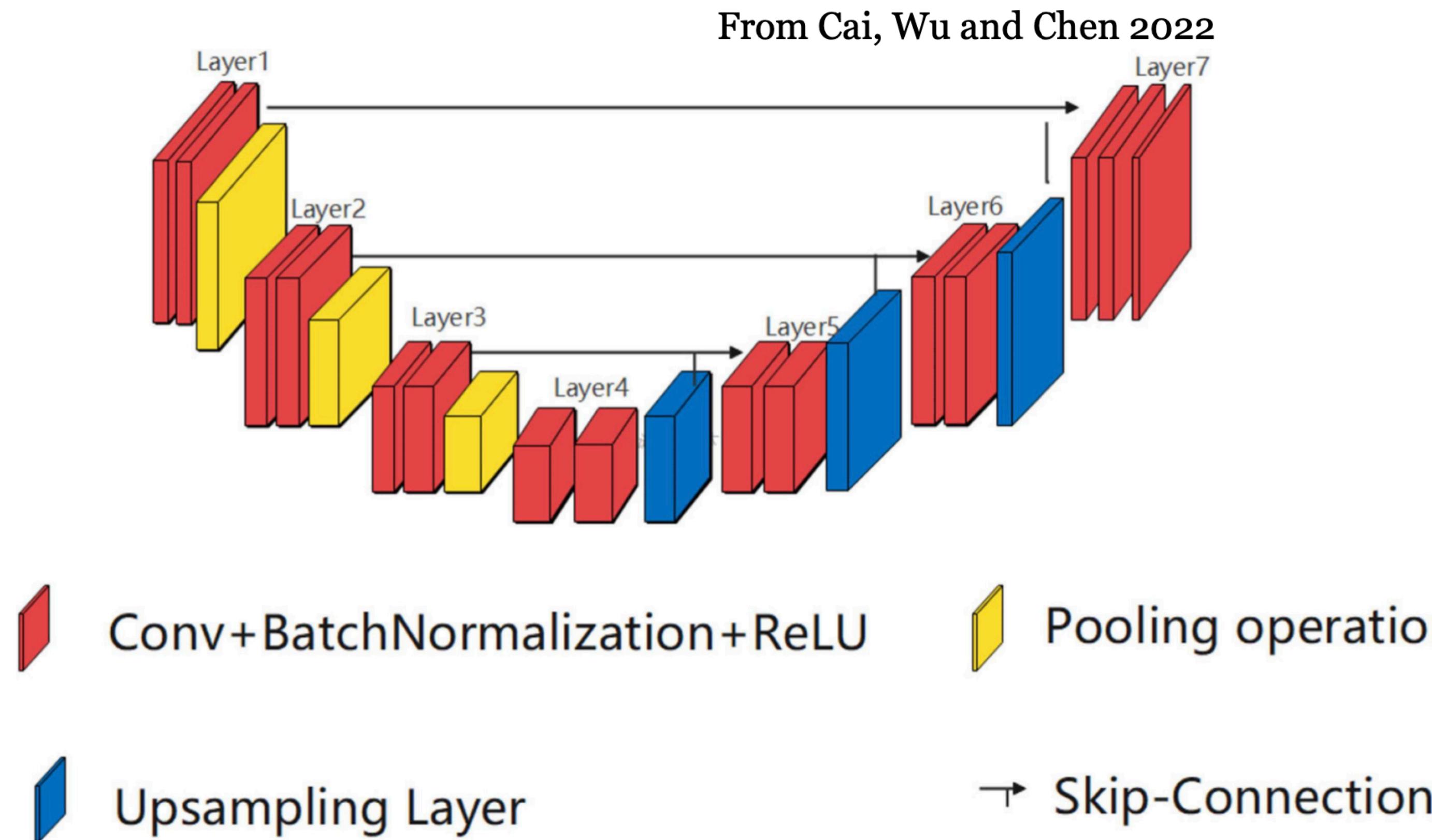
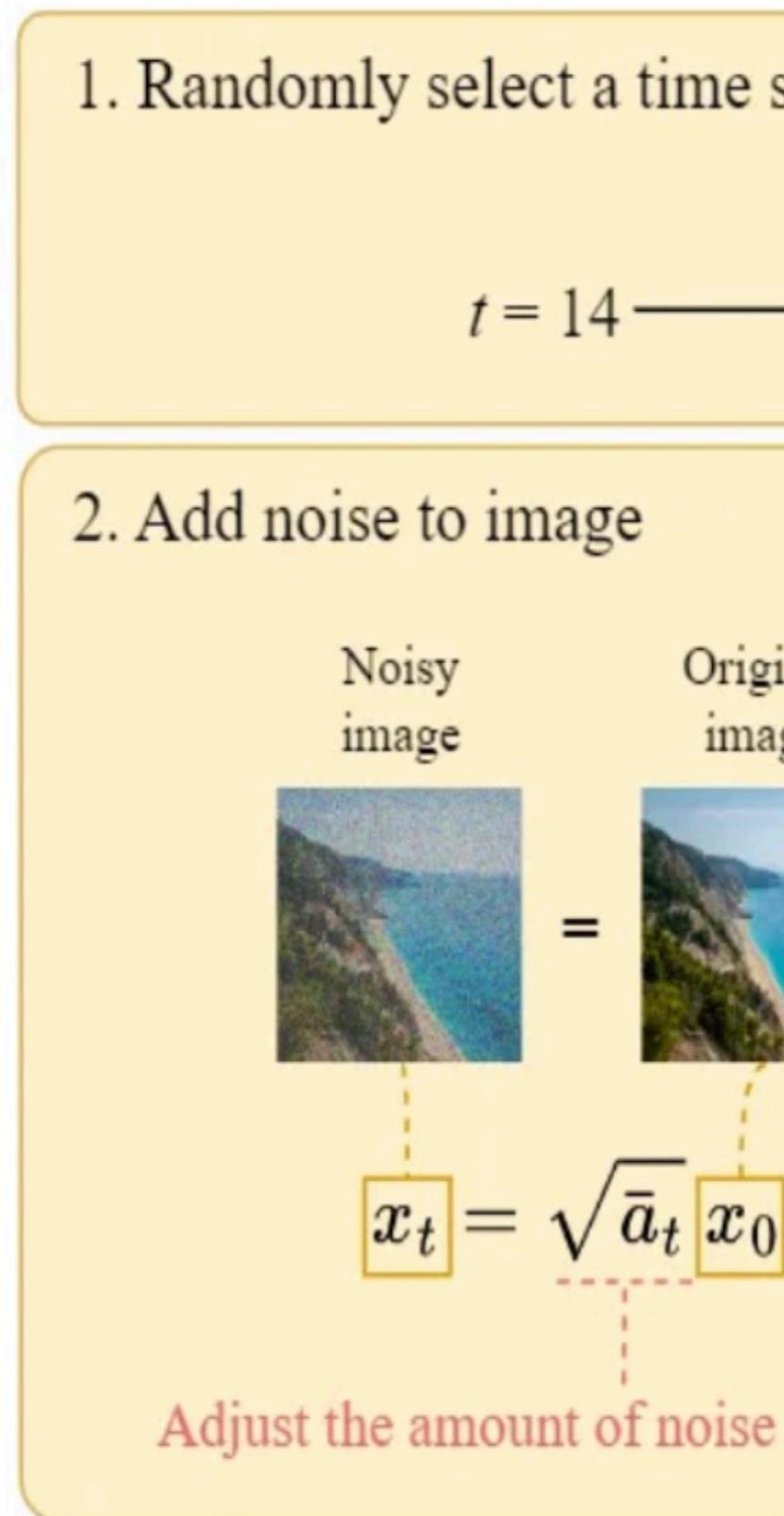
```

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
  
```

---

# Real Example: Training for Image Generation

For each training step:




---

## Algorithm 1 Training

---

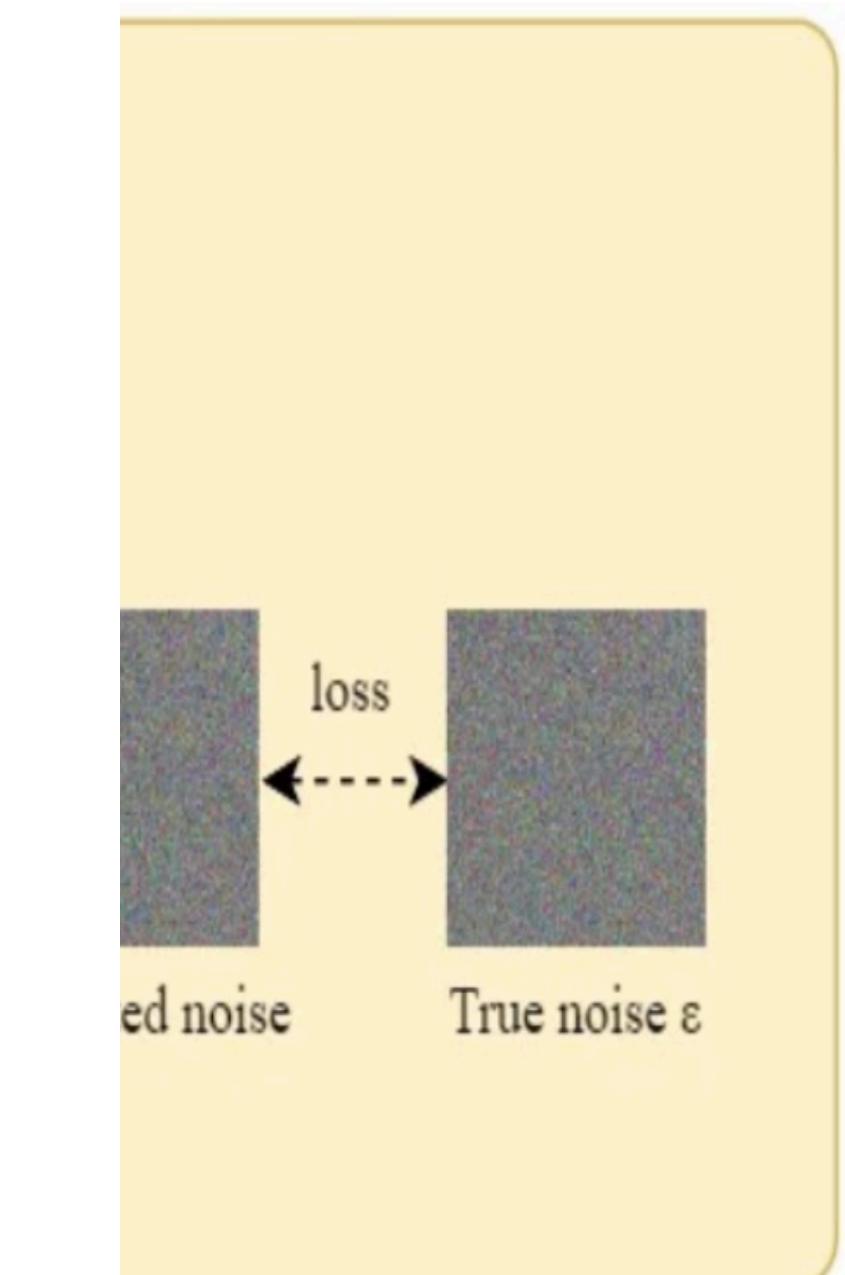
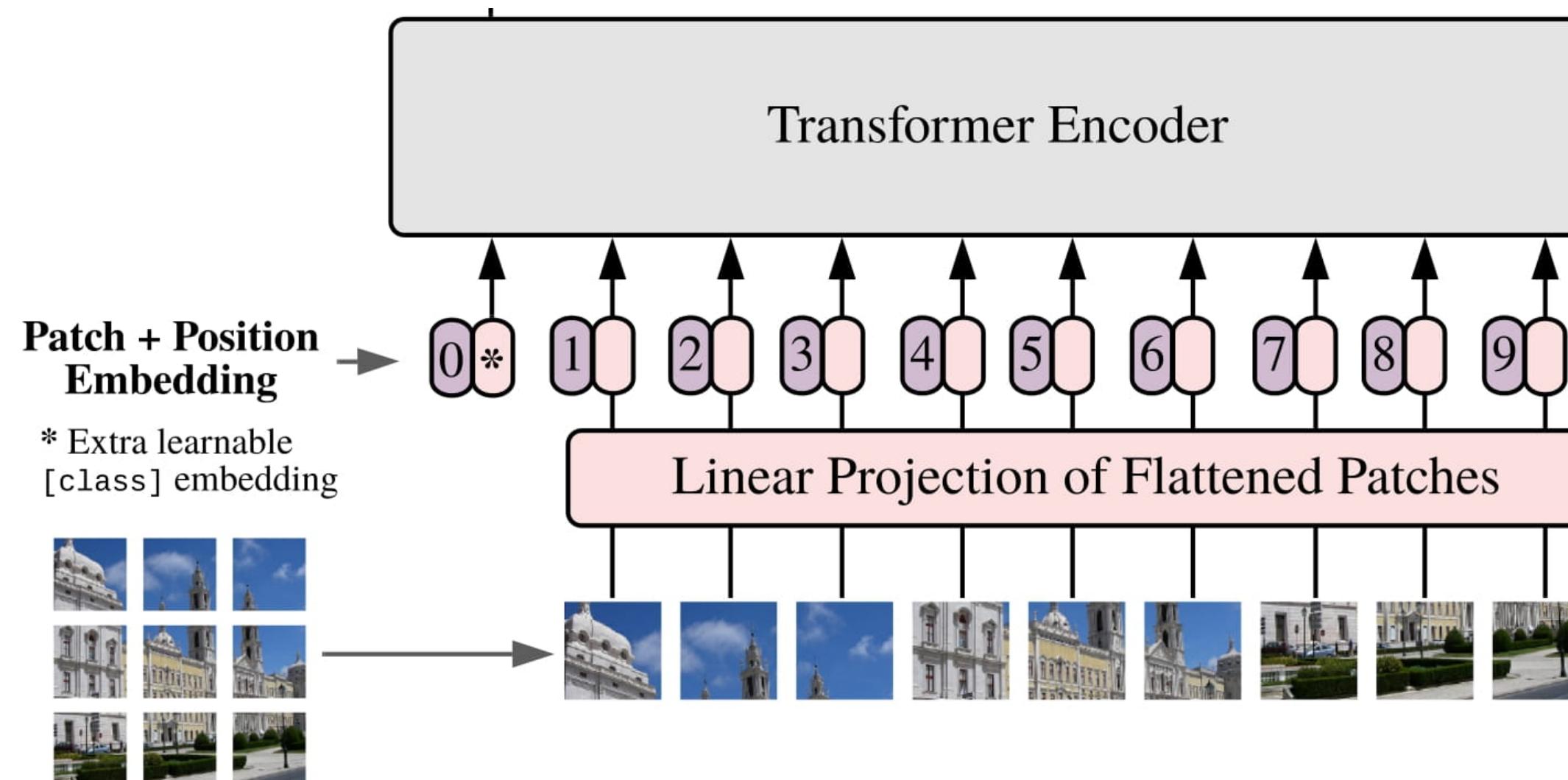
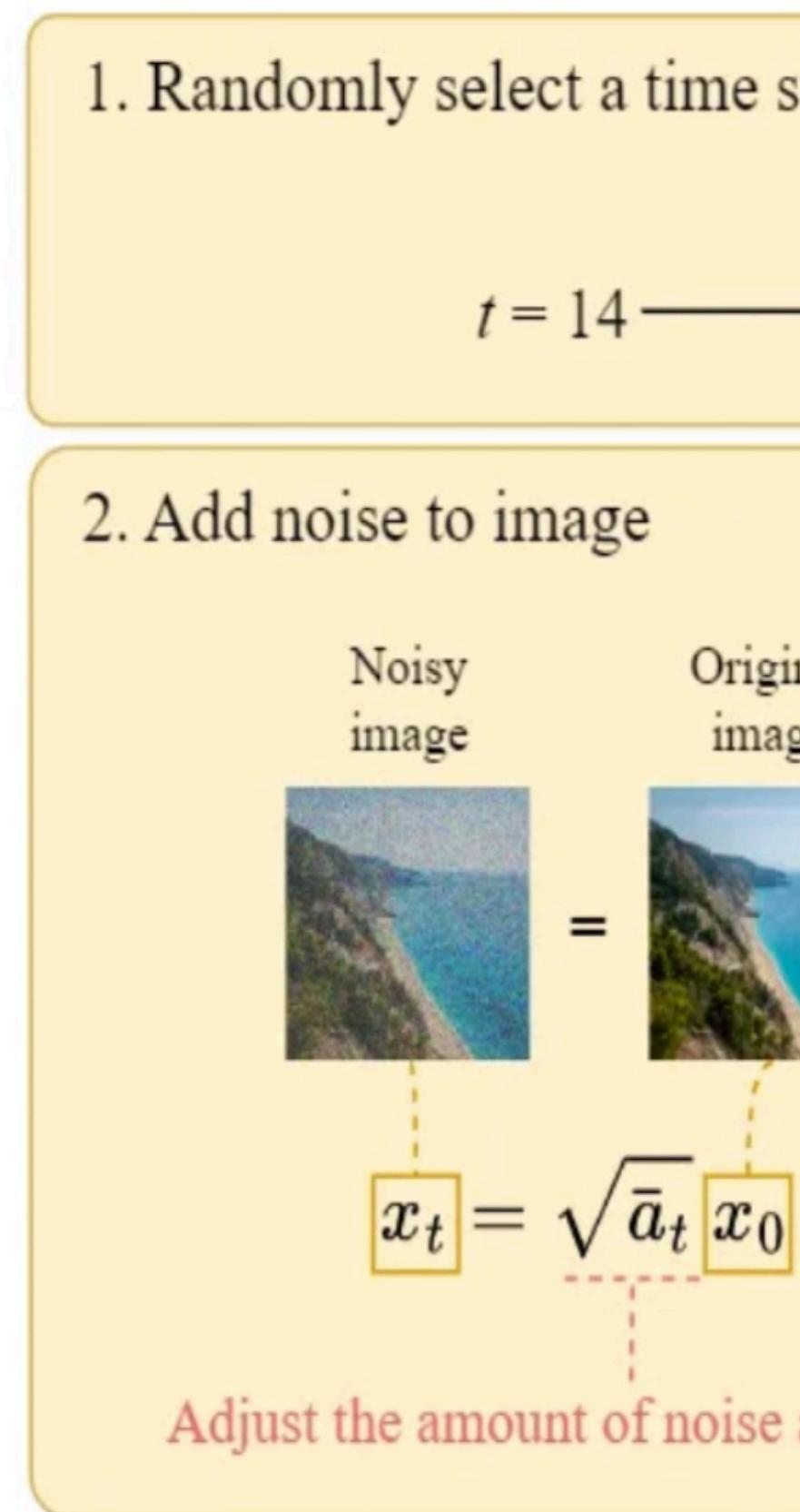
```

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
       $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{a}_t} \mathbf{x}_0 + \sqrt{1 - \bar{a}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged
  
```

---

# Real Example: Training for Image Generation

For each training step:



Source: [medium.com](https://medium.com)

---

## Algorithm 1 Training

---

```

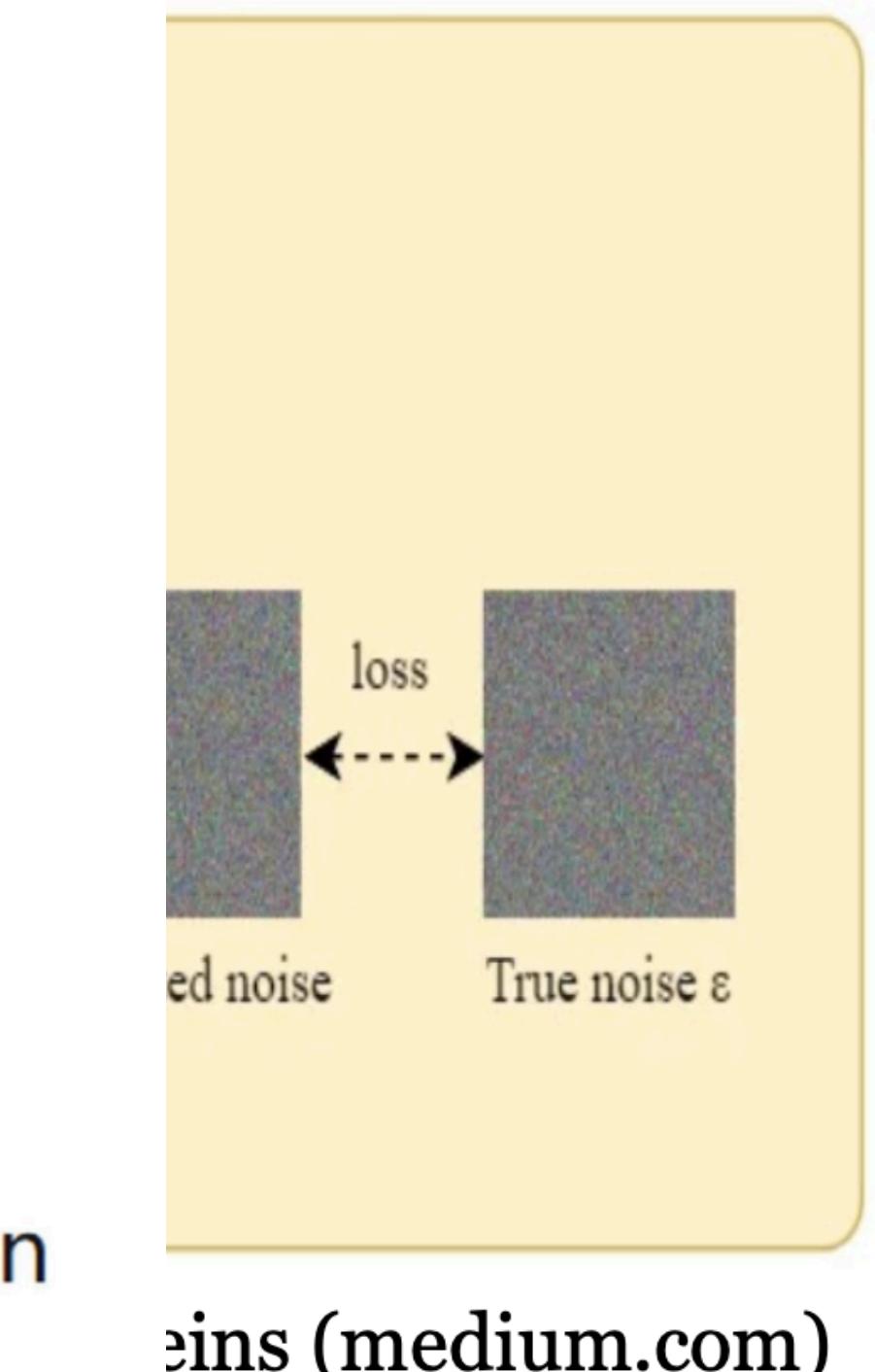
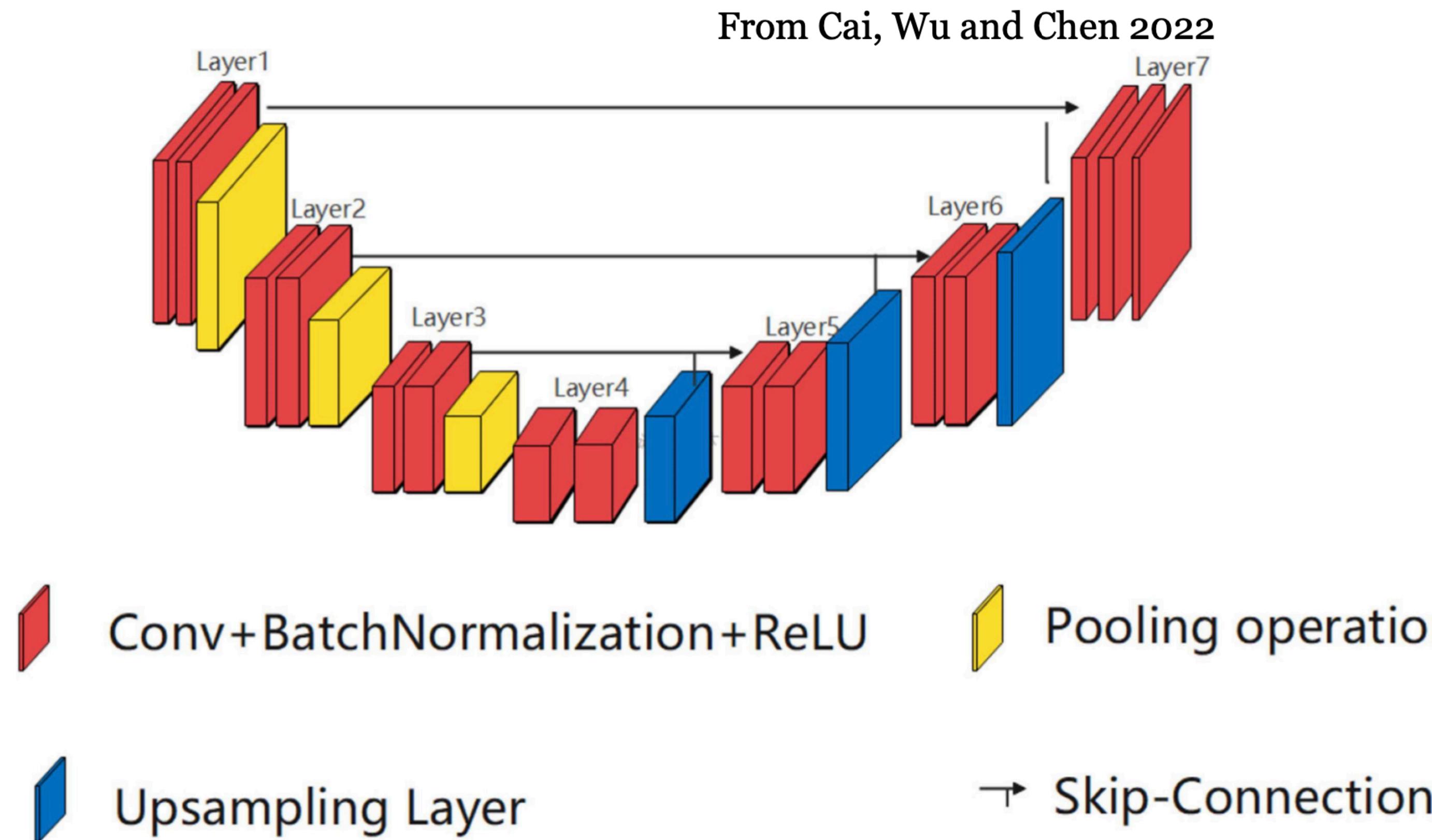
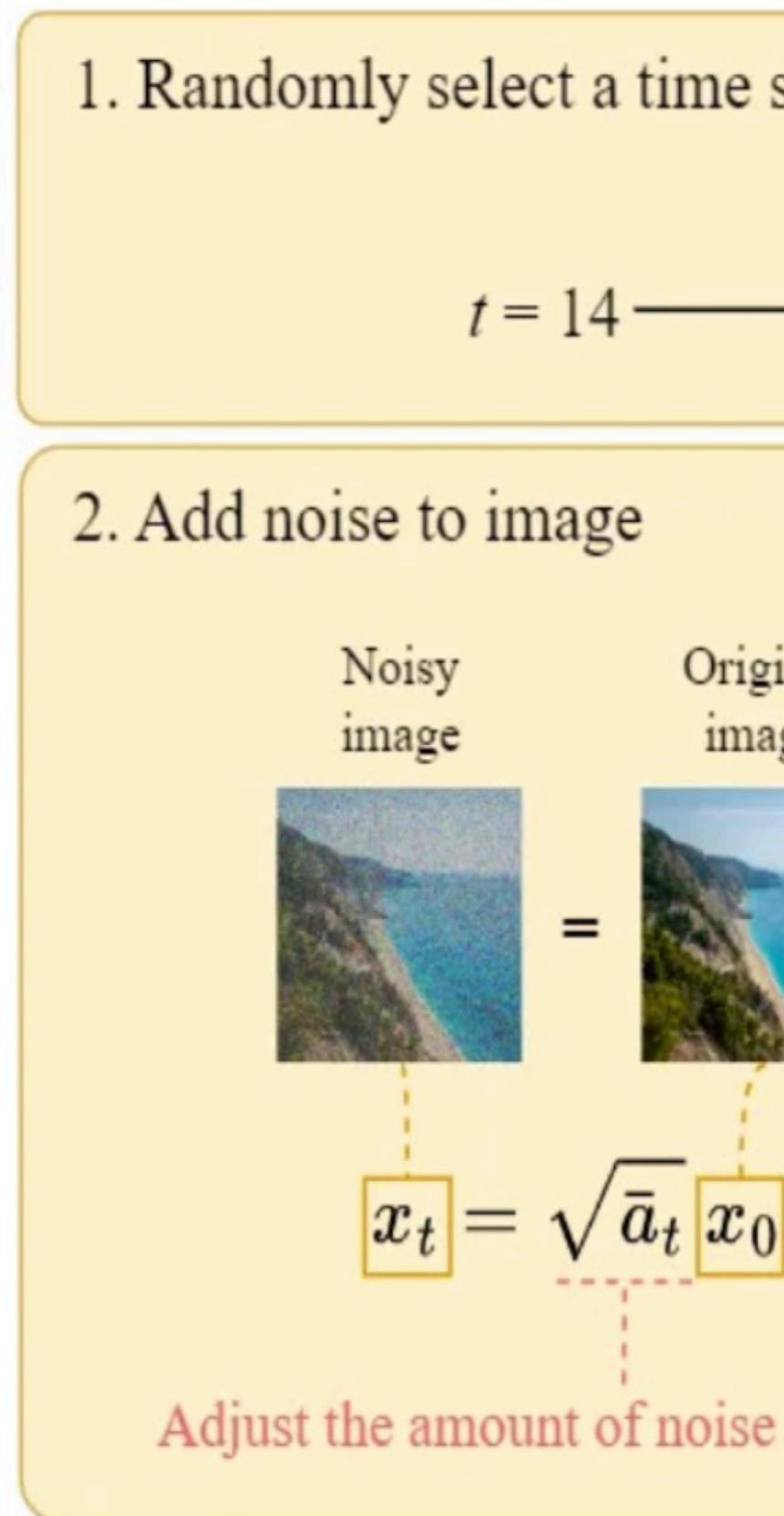
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
       $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{a}_t}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

```

---

# Real Example: Training for Image Generation

For each training step:




---

## 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
       $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{a}_t} \mathbf{x}_0 + \sqrt{1 - \bar{a}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged

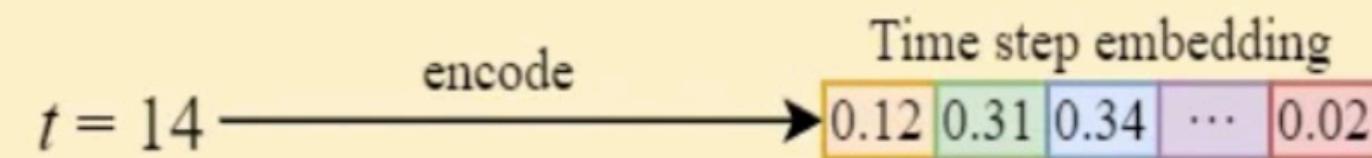
```

---

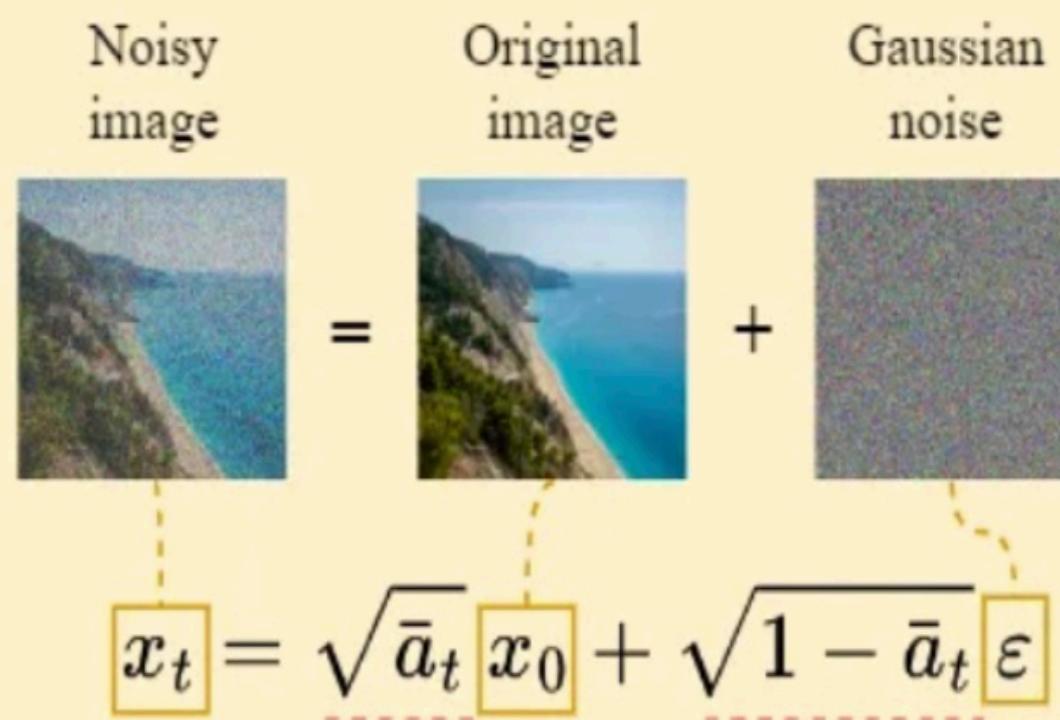
# Real Example: Training for Image Generation

For each training step:

1. Randomly select a time step & encode it



2. Add noise to image



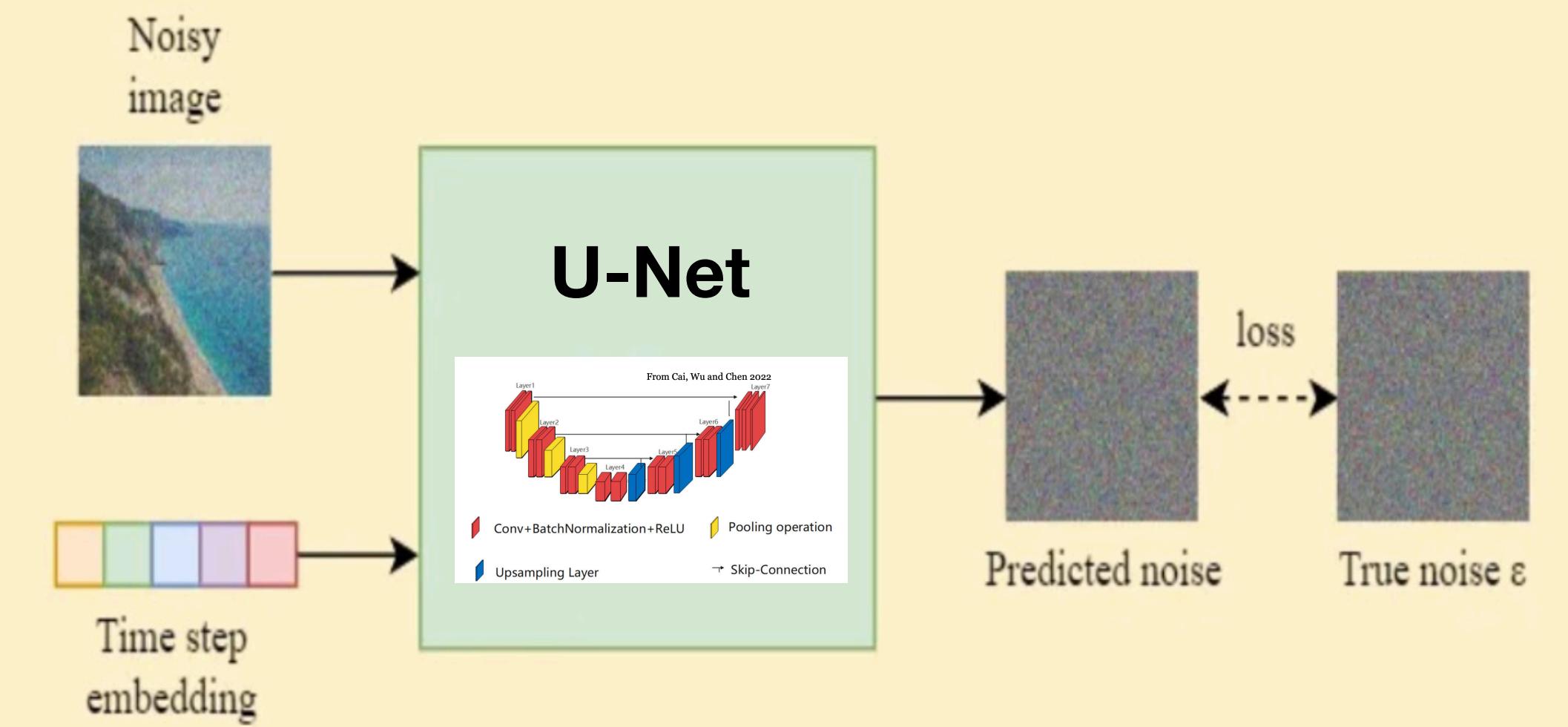
Adjust the amount of noise according to the time step  $t$

$$\varepsilon \sim \mathcal{N}(0, 1)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

3. Train the UNet



From Steins (medium.com)

- **Important:** We are training one network (shared params) across all steps

---

## 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
  
```

---

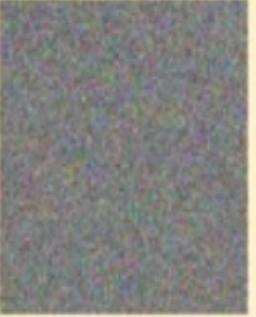
# Real Example: Inference for Image Generation

1. Sample a Gaussian noise

$$x_T \sim N(0, I)$$

E.g.  $T = 1000$

$$x_{1000} \sim N(0, I)$$



---

## 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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

---

- Step 1: Start from sampling pure noise

From Steins (medium.com)

# Real Example: Inference for Image Generation

---

## 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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

---

1. Sample a Gaussian noise

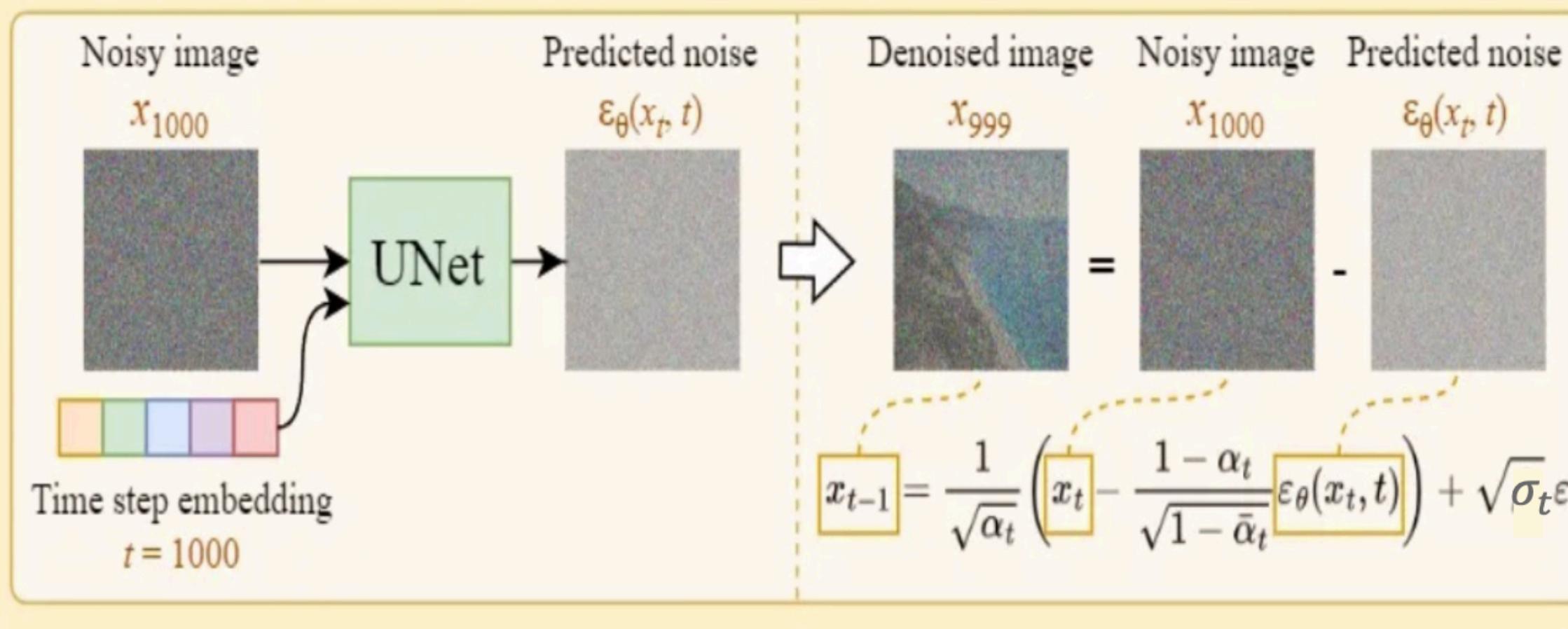
$$\mathbf{x}_T \sim N(\mathbf{0}, \mathbf{I})$$

E.g.  $T = 1000$

$$\mathbf{x}_{1000} \sim N(\mathbf{0}, \mathbf{I})$$



2. Iteratively denoise the image



From Steins (medium.com)

- **Step 2:** Predict the noise using a trained (U-Net) model for  $t-1$

# Real Example: Inference for Image Generation

---

## 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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```

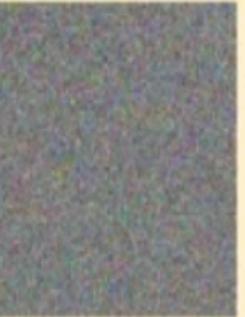
---

1. Sample a Gaussian noise

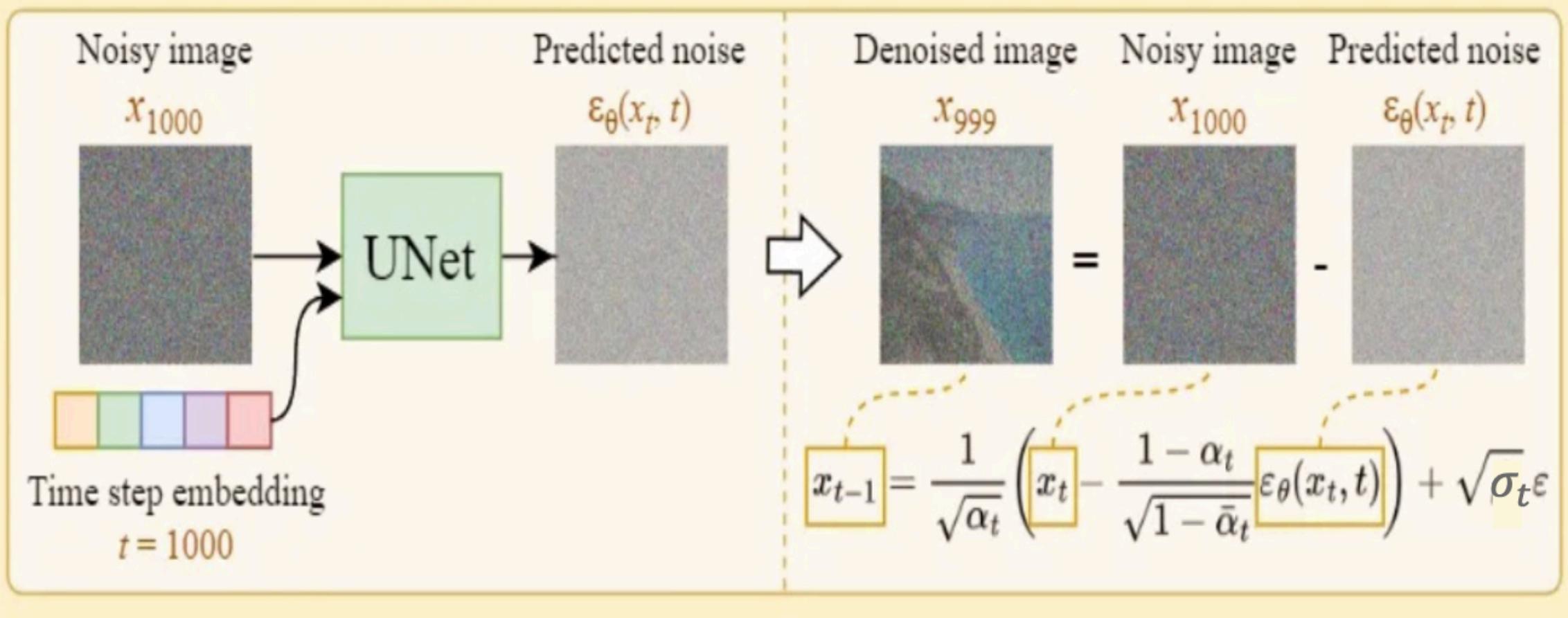
$$\mathbf{x}_T \sim N(\mathbf{0}, \mathbf{I})$$

E.g.  $T = 1000$

$$\mathbf{x}_{1000} \sim N(\mathbf{0}, \mathbf{I})$$



2. Iteratively denoise the image



From Steins (medium.com)

- **Step 2:** Subtract predicted noise to obtain cleaner version of the sample

# Real Example: Inference for Image Generation

---

**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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
  
```

---

1. Sample a Gaussian noise

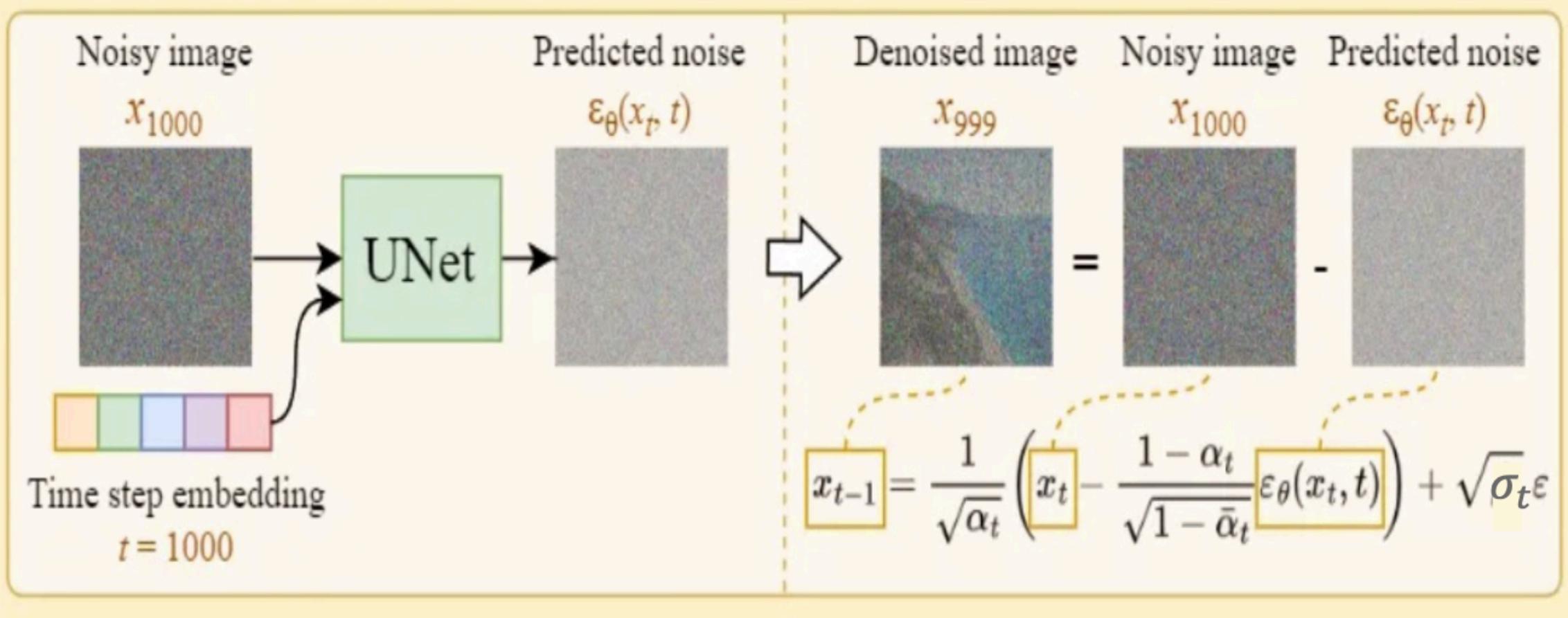
$$\mathbf{x}_T \sim N(\mathbf{0}, \mathbf{I})$$

E.g.  $T = 1000$

$$\mathbf{x}_{1000} \sim N(\mathbf{0}, \mathbf{I})$$



2. Iteratively denoise the image



- Repeat many ( $T$ ) times

From Steins (medium.com)

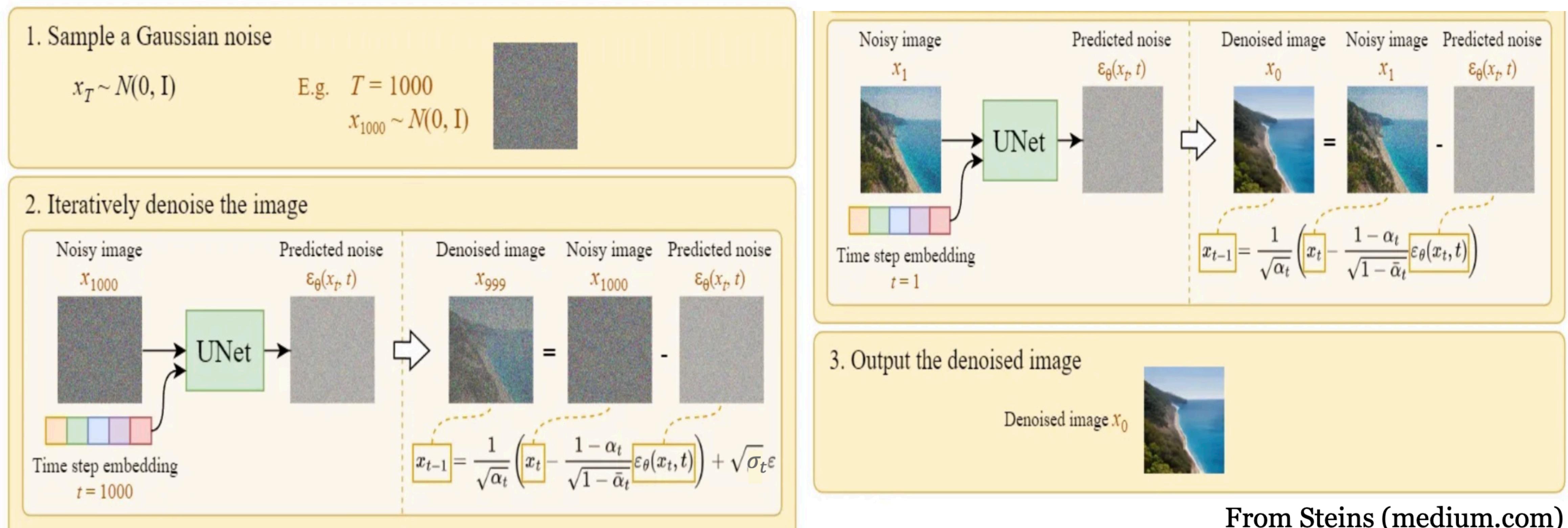
# Real Example: Inference for Image Generation

## 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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

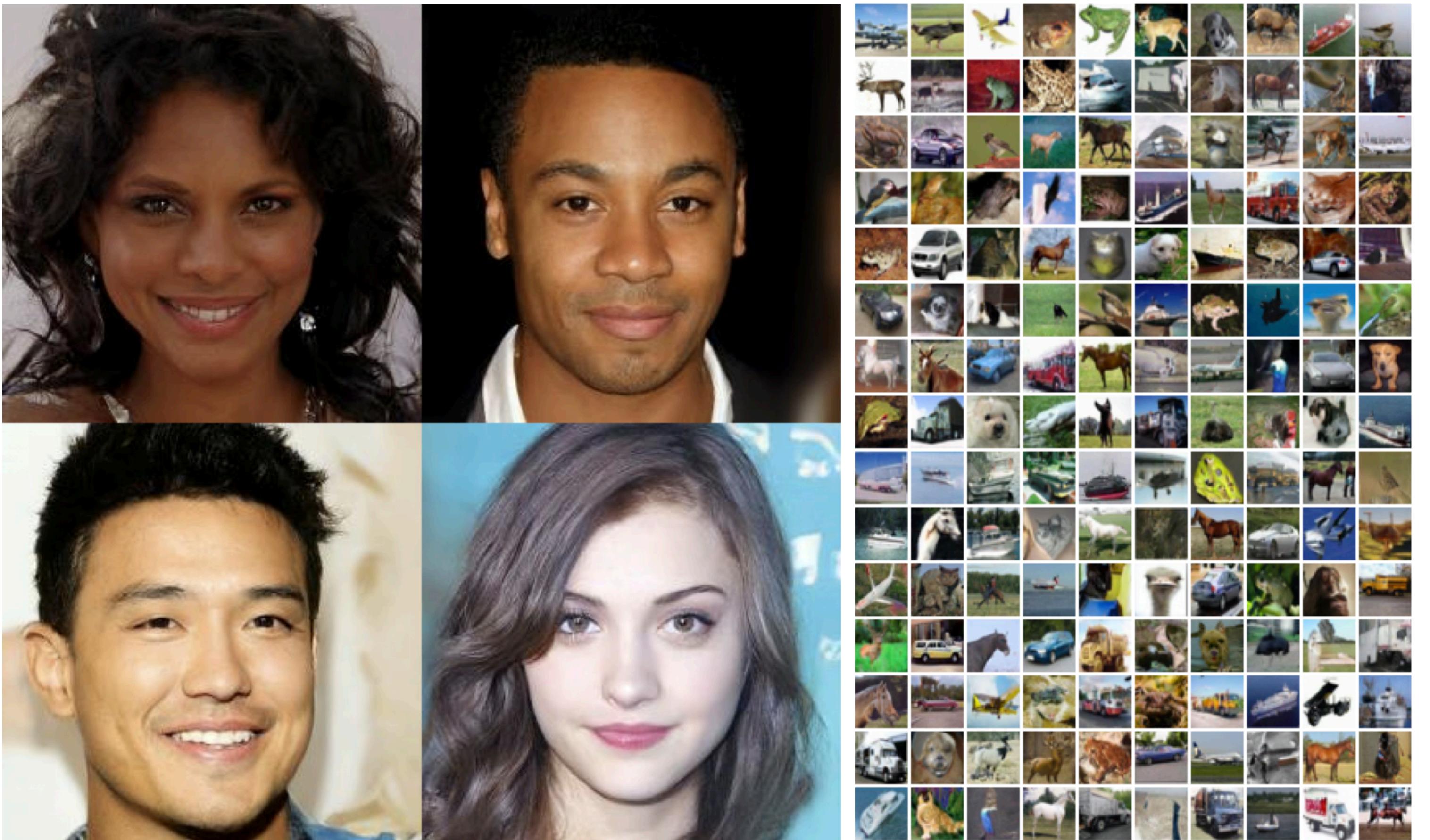
```



From Steins (medium.com)

- **Step 3:** The final output is a clean image corresponding to  $t=1$

# Real Example: Image Generation



# Properties of Good Generative Model

High Quality  
Samples

Mode Coverage  
/ Diversity

Fast Sampling

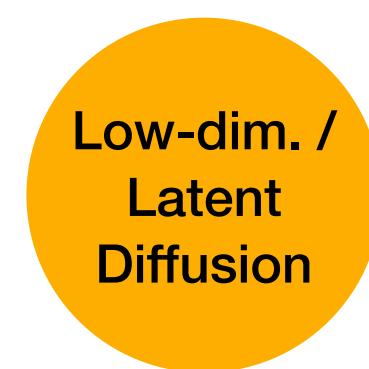
# Accelerating Sampling

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



# Accelerated Sampling: Denoising Diffusion Implicit Models (DDIM)

- **Theory:** Assume non-Markovian process (conditioning on original sample) and re-derive forward and backwards equations
- **Practice:**
  - Init: Let  $t = T = 1000$ , sample noise
    - **Step 1** : always predict  $t=1$
    - **Step 2** : add noise to get to an intermediate step (like in training)
    - Continue this process
  - **Nice property:** does not require any changes to training

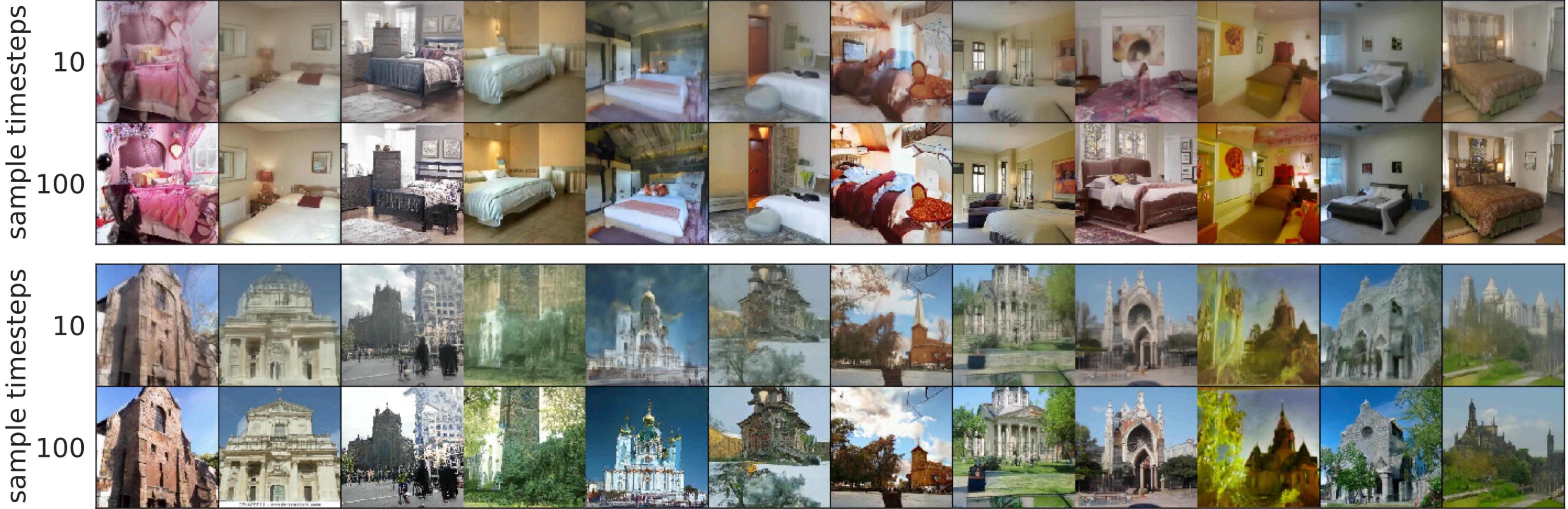
# Accelerated Sampling: Denoising Diffusion Implicit Models (DDIM)

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



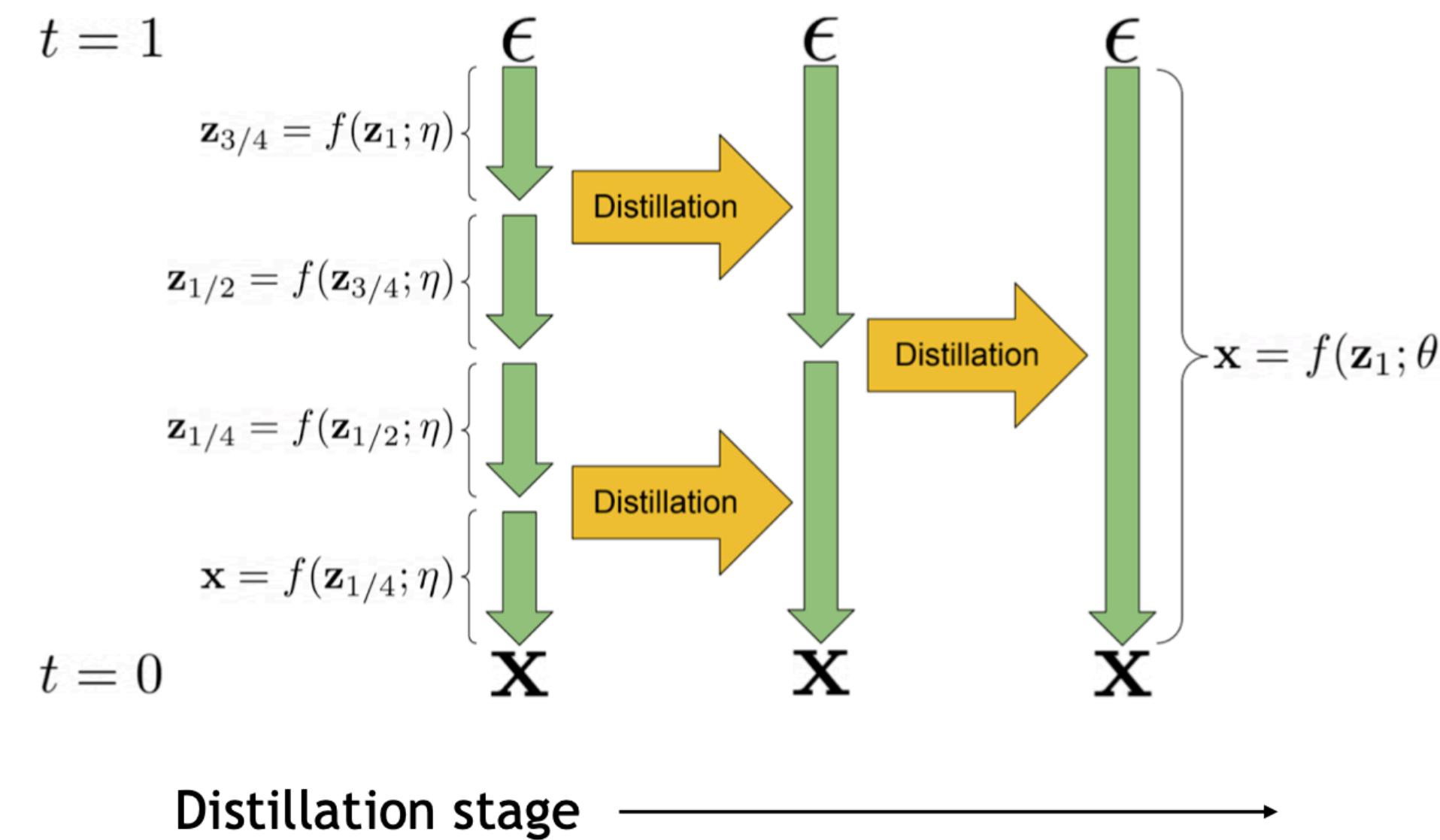
# Accelerated Sampling: Progressive Distillation

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



- Distill a sampler to the same model architecture but with fewer steps
  - e.g., a “student” model is learned to distill two adjacent sampling steps of the “teacher” model

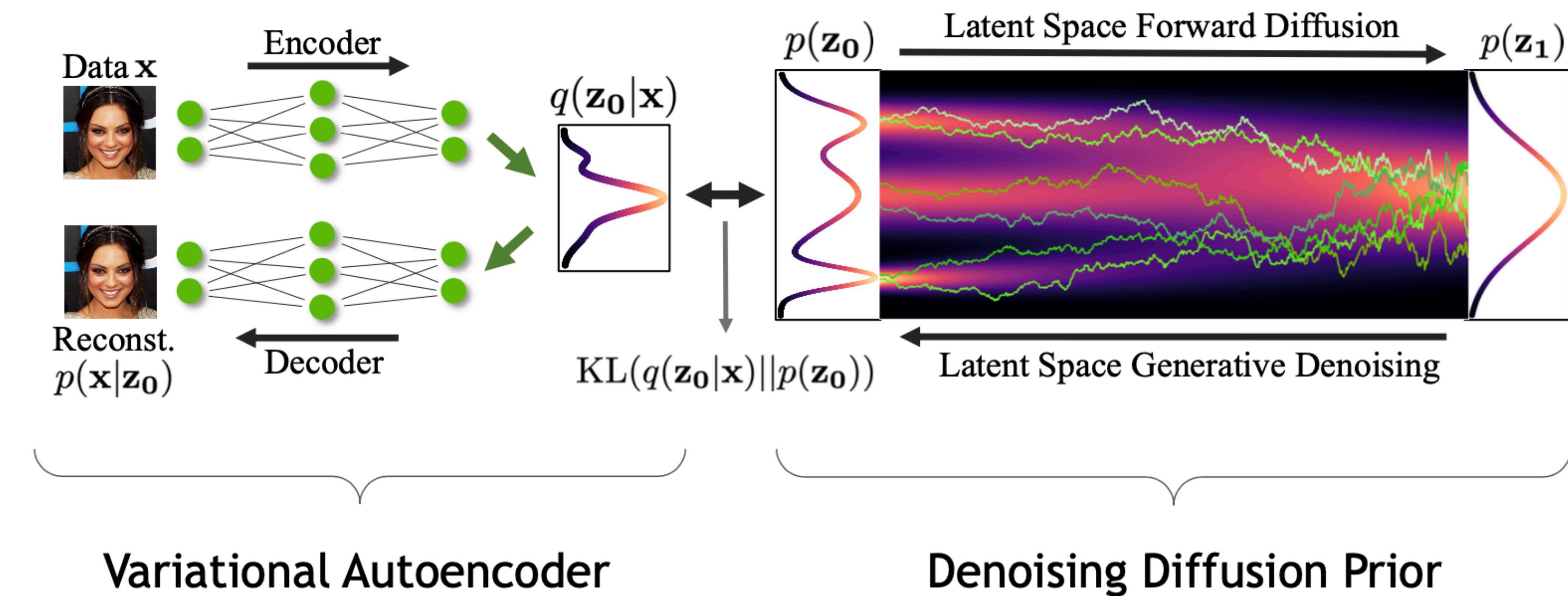
# Accelerated Sampling: Latent Diffusion Model

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



- **Main idea:** Encode data into an embedding space and apply latent diffusion in this space

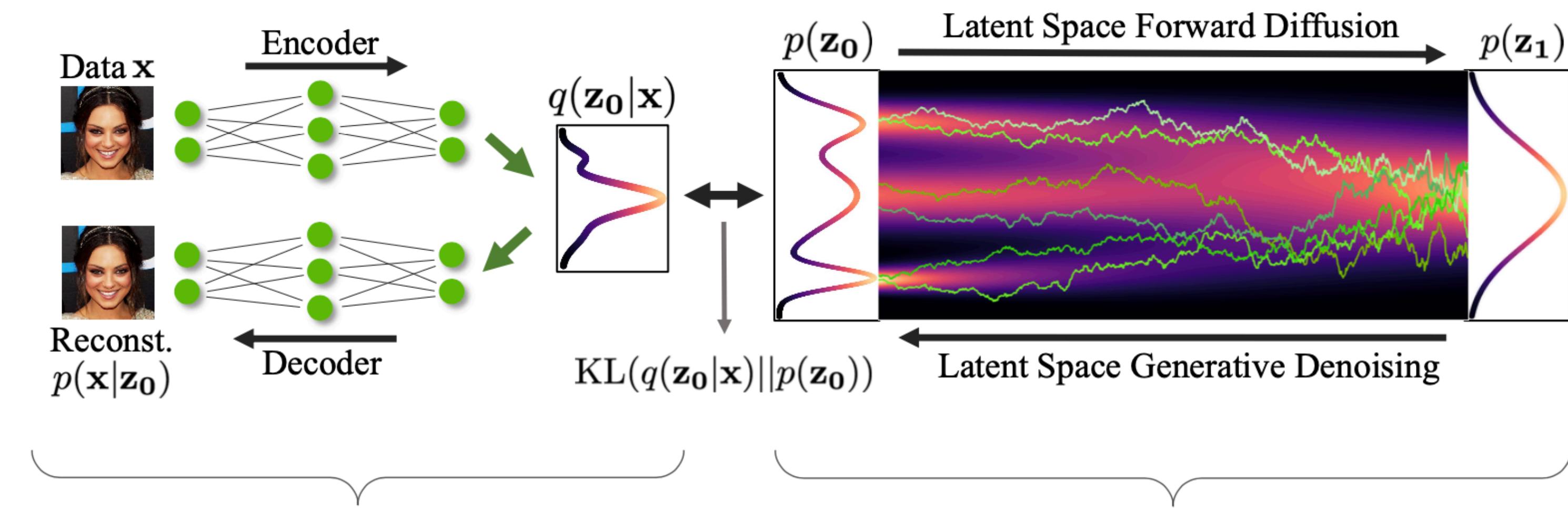
# Accelerated Sampling: Latent Diffusion Model

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



- **Benefits:**

Variational Autoencoder

Denoising Diffusion Prior

- Distribution in latent space normal  $\rightarrow$  simpler denoising
- Distribution in latent space low-dimensional  $\rightarrow$  faster sampling
- Easier to apply to different data types (graphs, tables, text, 3d)

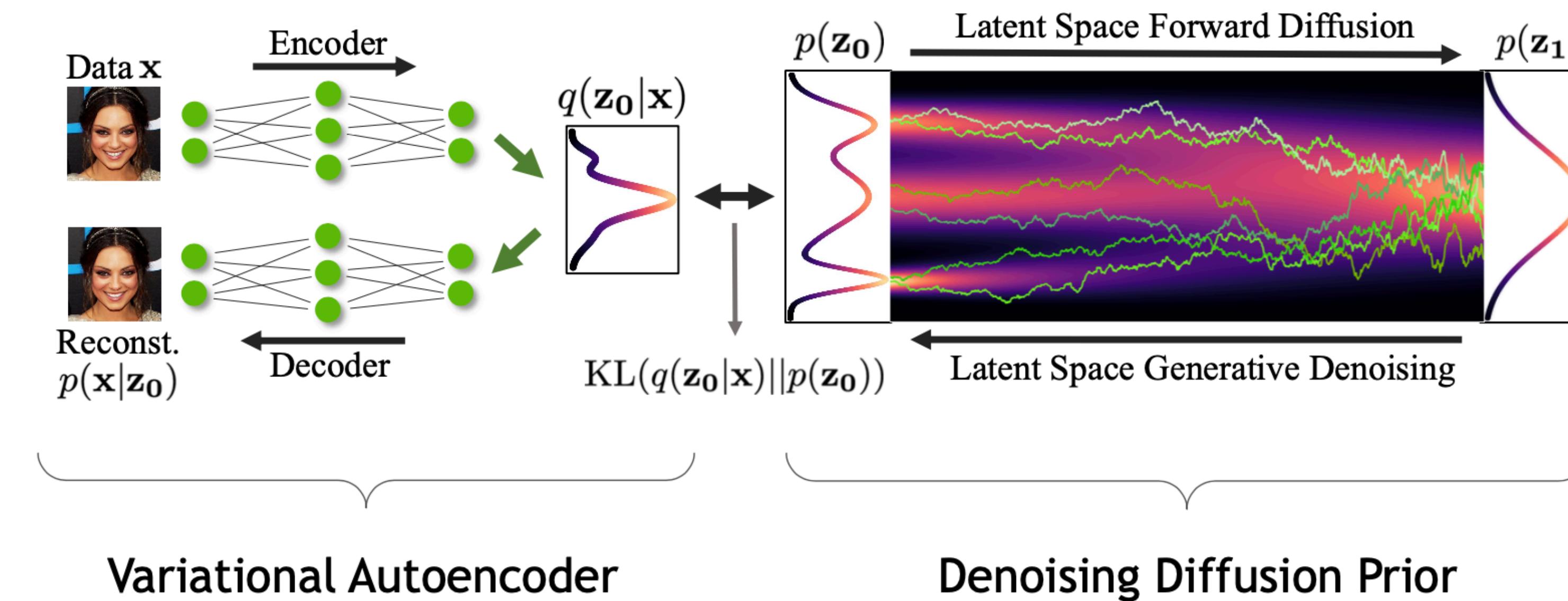
# Accelerated Sampling: Latent Diffusion Model

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

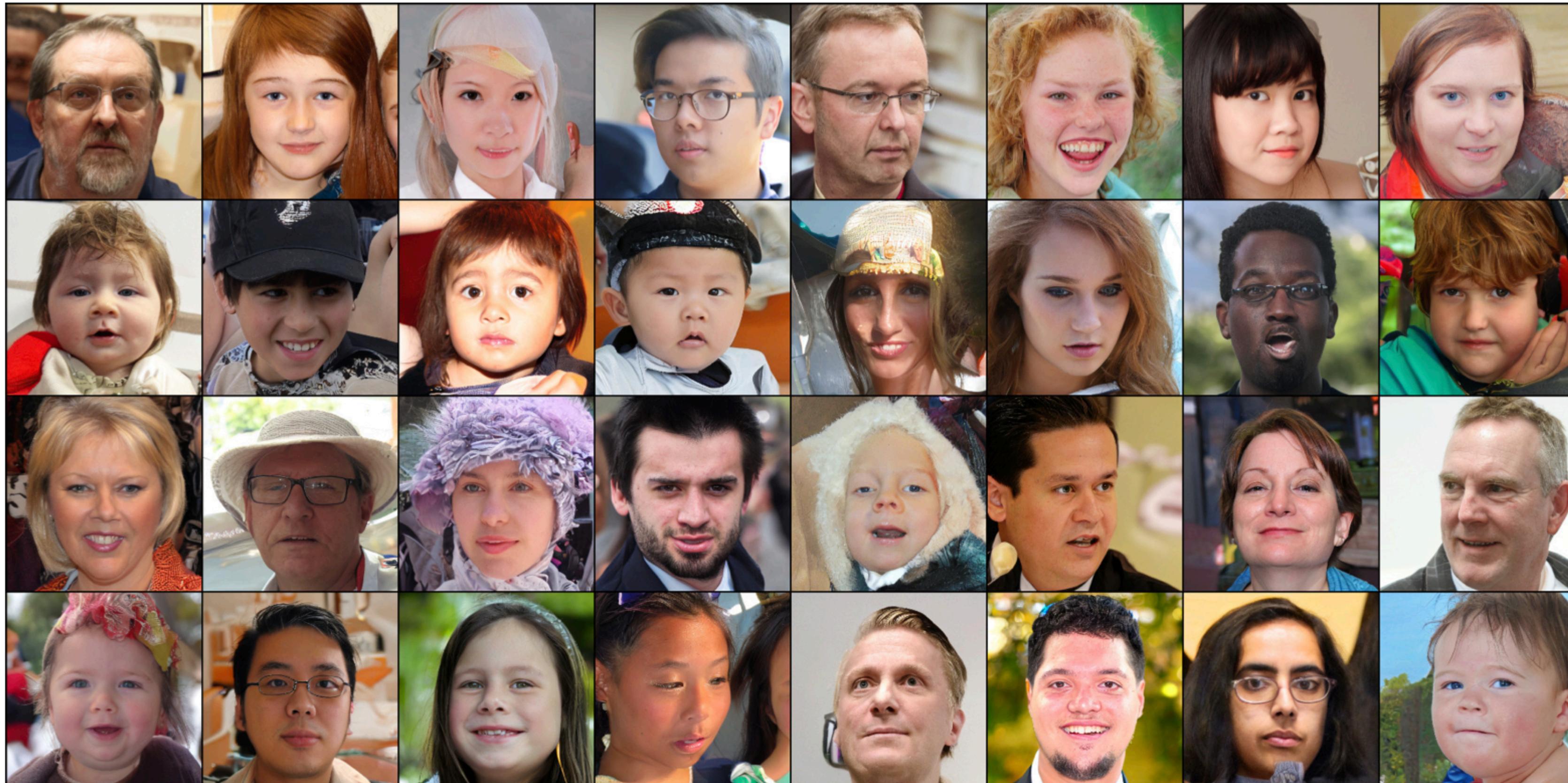
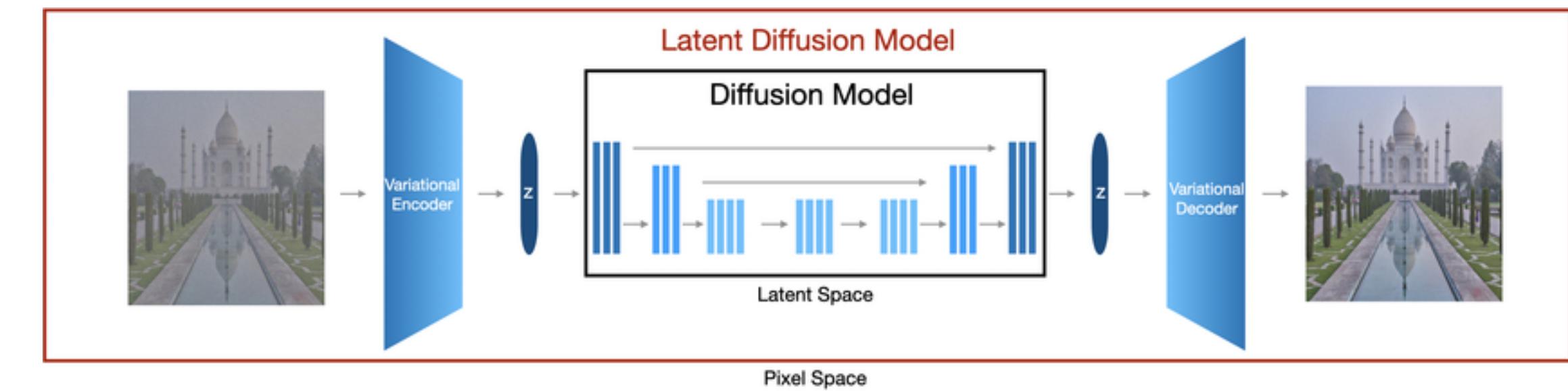
Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



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

# Accelerated Sampling: Latent Diffusion Model



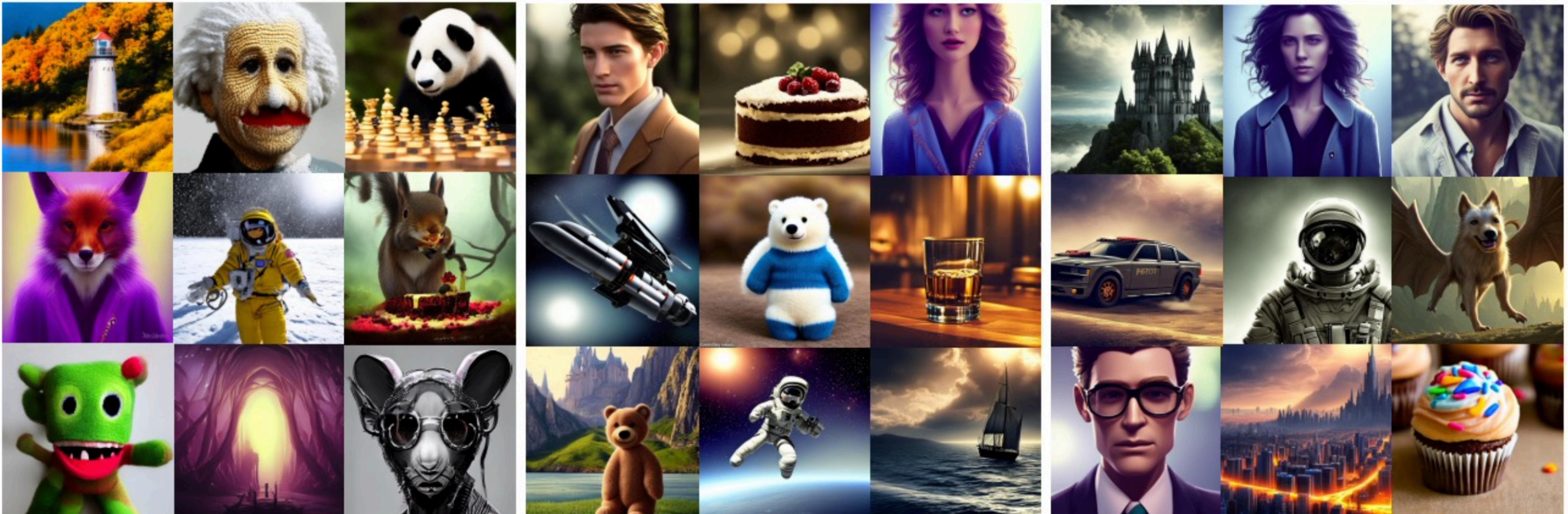
# Accelerated Sampling: Latent Diffusion Model

Advanced  
ODE / SDE  
Solvers

Distillation  
Techniques

Low-dim. /  
Latent  
Diffusion

Advanced  
Diffusion  
Processes



(a) 2 denoising steps

(b) 4 denoising steps

(c) 8 denoising steps

# Accelerating Sampling : Summary

Advanced  
ODE / SDE  
Solvers

e.g., **DDIM**

Simple  
Requires no training  
Could be inaccurate

Distillation  
Techniques

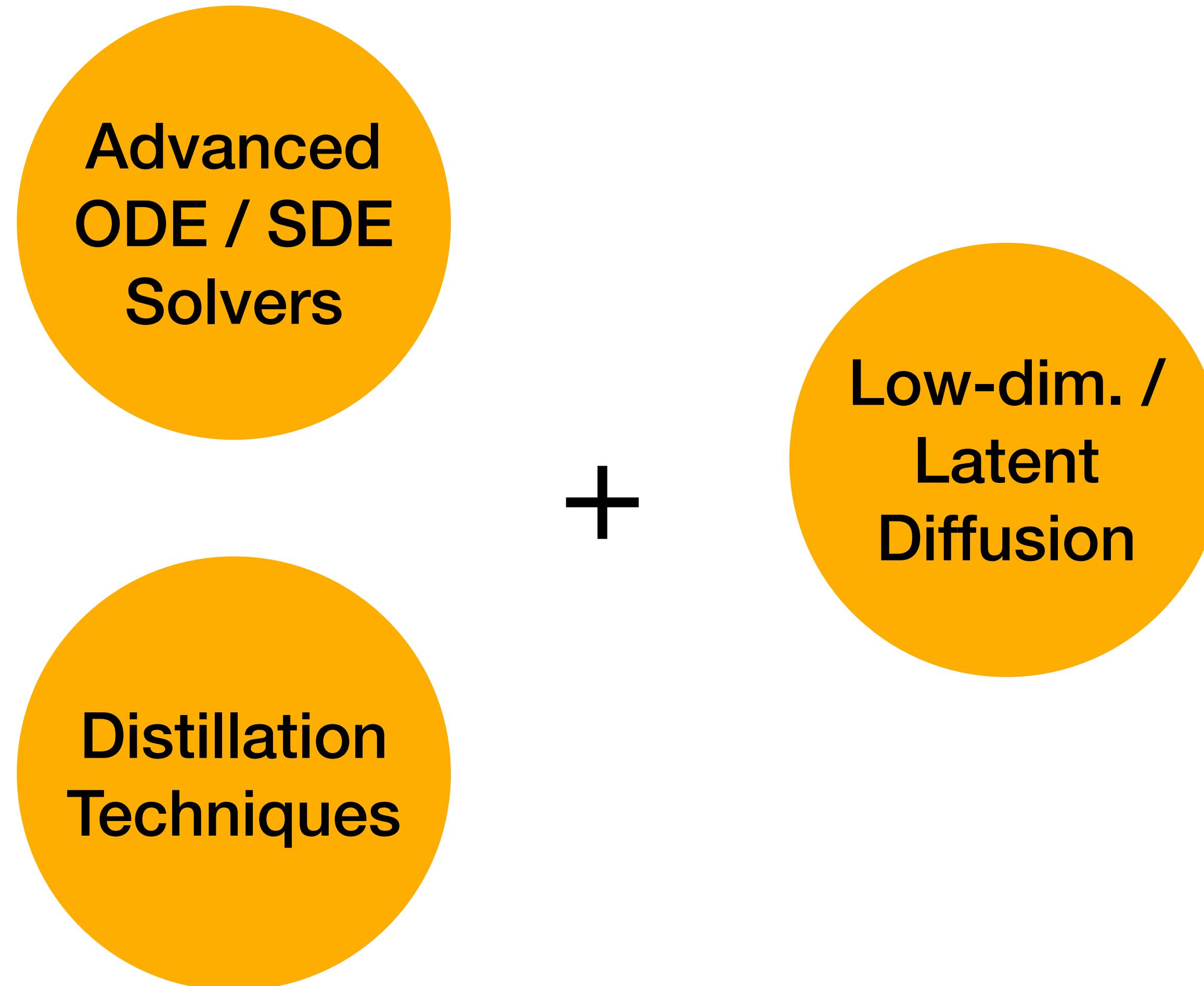
More complex to setup  
Requires training  
Usually more accurate

Low-dim. /  
Latent  
Diffusion

Moderate  
Requires training of encoder  
Usually works well

Advanced  
Diffusion  
Processes

# Accelerating Sampling : Summary



# Conditional Diffusion Models

- Generating raw samples of data is often not very useful in practice, much more practically useful is to generate data conditioned on some input (e.g., text to image generation)

DALL·E 2

*“a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese”*

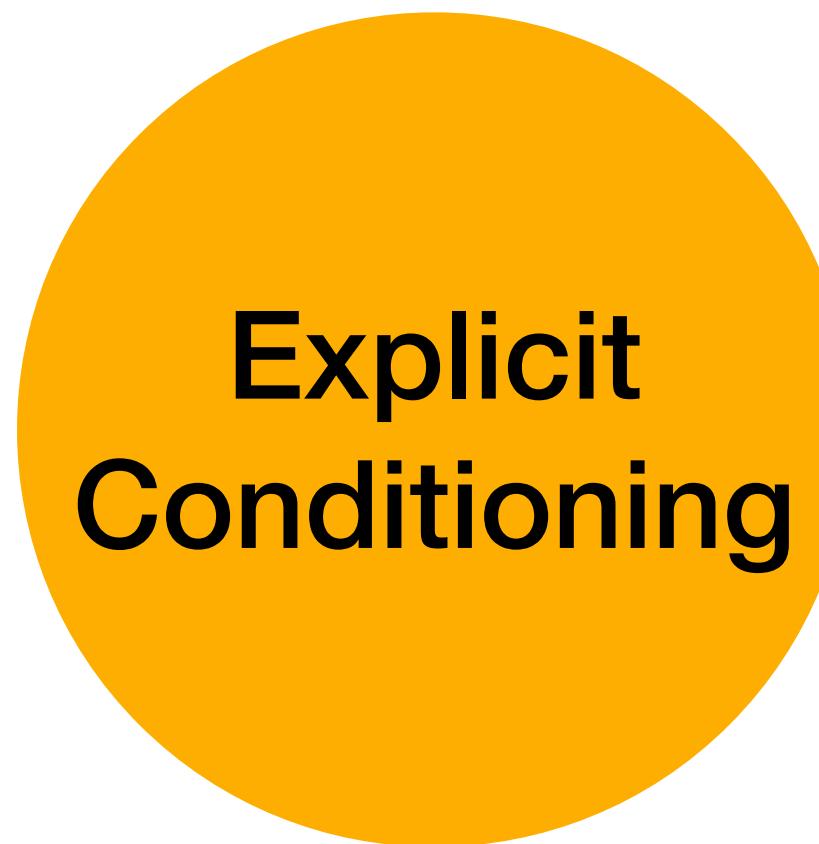


IMAGEN

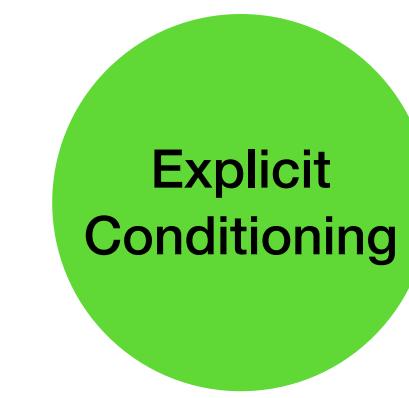
*“A photo of a raccoon wearing an astronaut helmet, looking out of the window at night.”*



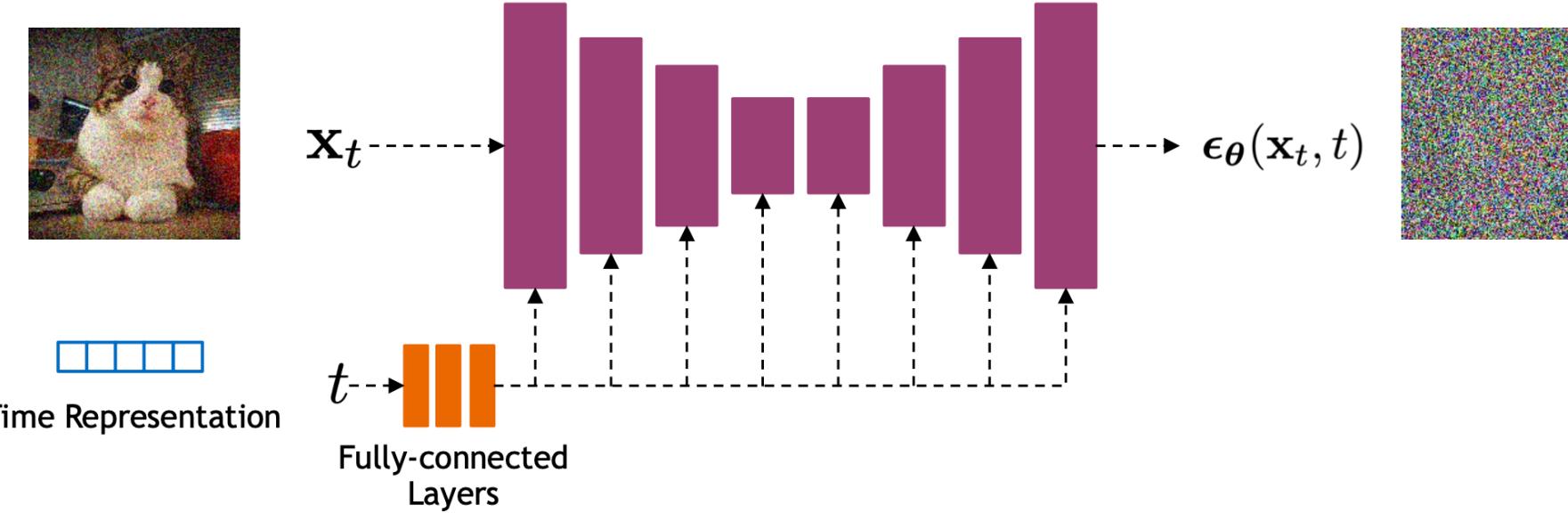
# Conditional Diffusion Models:



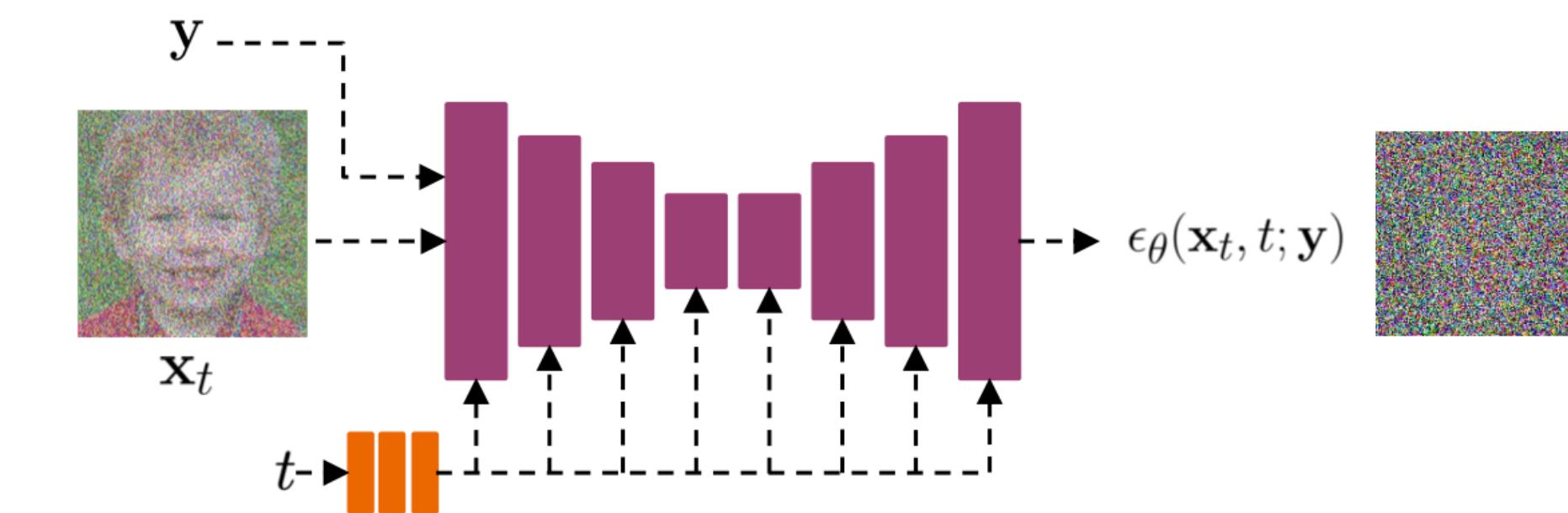
# Conditional Diffusion Models: Explicit Conditioning



- Instead of training diffusion model to model  $p(x)$ , train it to model conditional distribution  $p(x|y)$ , where  $y$  is the conditioning variable (e.g., text) and  $x$  is the output (e.g., image) corresponding to the condition.
- Effectively requires no changes to the diffusion process (just data and fusion)

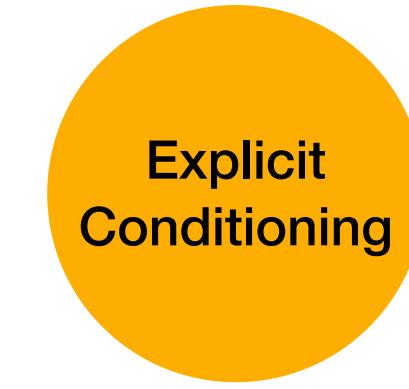


**Unconditional Diffusion Model**



**Conditional Diffusion Model**

# Conditional Diffusion Models: Classifier Guidance

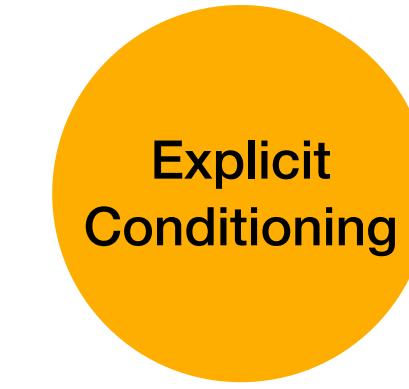


- When sampling from a diffusion model, we need an estimate of gradient  $p(x)$  with respect to the sample
- Using Bayes' rule we can express the gradient for conditional diffusion as

$$p(x_t|y) = \frac{p(y|x_t)p(x_t)}{p(y)}$$

$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

# Conditional Diffusion Models: Classifier Guidance



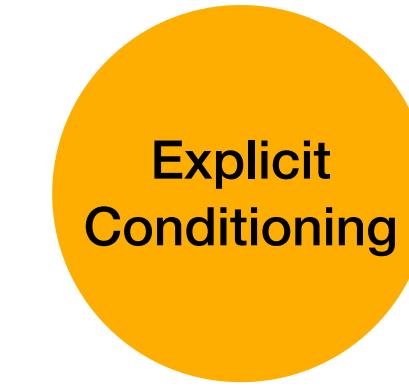
- When sampling from a diffusion model, we need an estimate of gradient  $p(x)$  with respect to the sample
- Using Bayes' rule we can express the gradient for conditional diffusion as

$$p(x_t|y) = \frac{p(y|x_t)p(x_t)}{p(y)}$$

$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

Gradient of pre-trained classifier

# Conditional Diffusion Models: Classifier Guidance



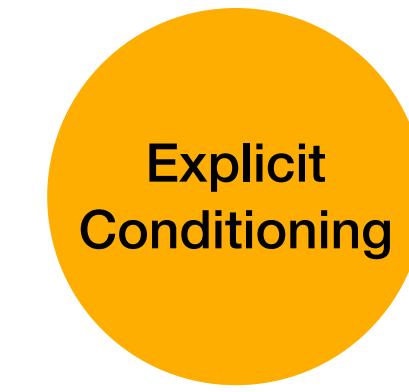
- When sampling from a diffusion model, we need an estimate of gradient  $p(x)$  with respect to the sample
- Using Bayes' rule we can express the gradient for conditional diffusion as

$$p(x_t|y) = \frac{p(y|x_t)p(x_t)}{p(y)}$$

$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

Standard diffusion gradient

# Conditional Diffusion Models: Classifier Guidance



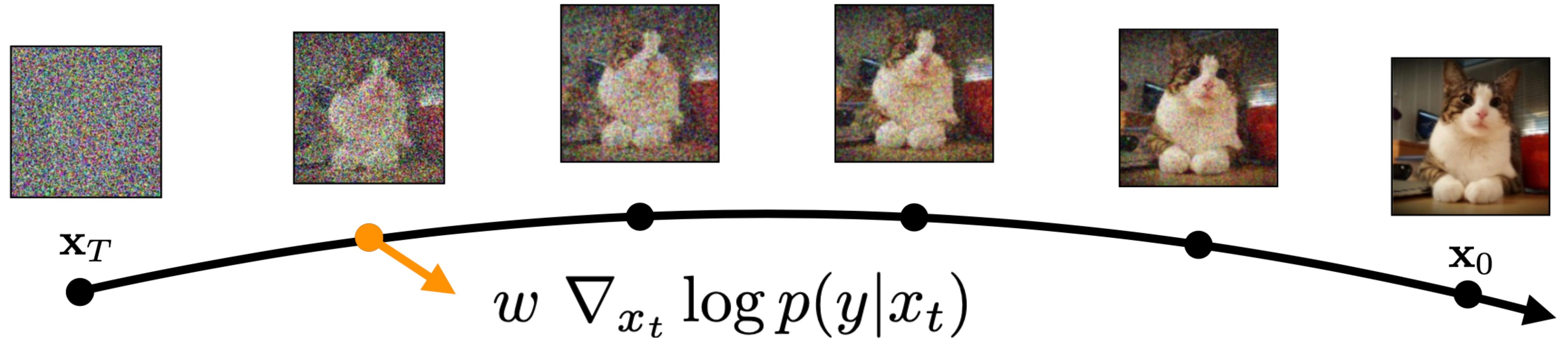
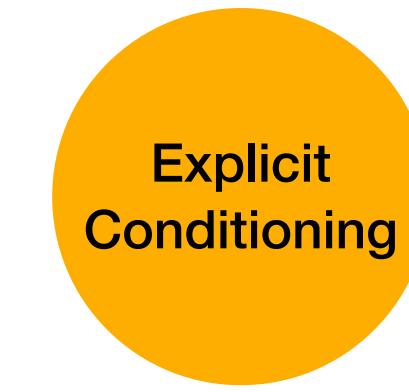
- When sampling from a diffusion model, we need an estimate of gradient  $p(x)$  with respect to the sample
- Using Bayes' rule we can express the gradient for conditional diffusion as

$$p(x_t|y) = \frac{p(y|x_t)p(x_t)}{p(y)}$$

$$\nabla_{x_t} \log p(x_t|y) = w \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

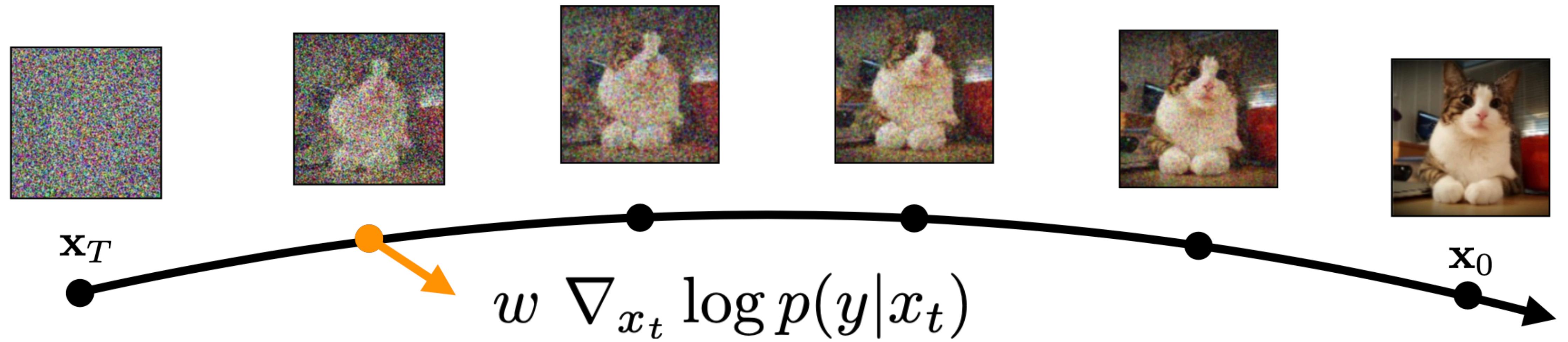
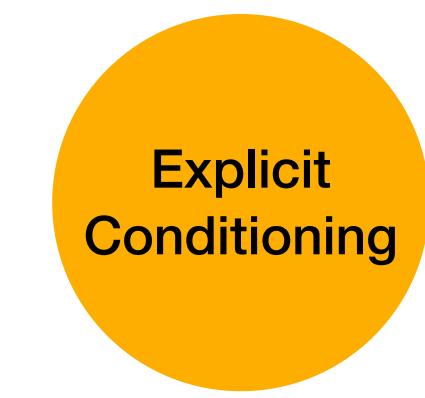
Strength of guidance

# Conditional Diffusion Models: Classifier Guidance

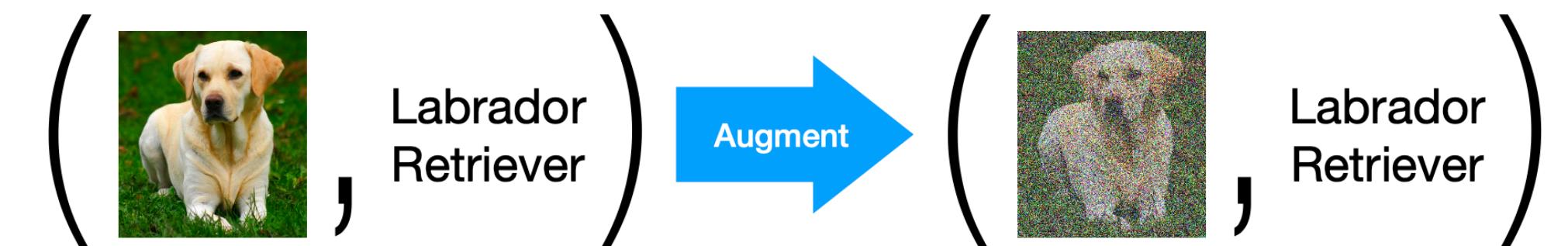


- **Intuition:** Perturb the inference trajectory (series of de-noising steps) at each step using feedback from a pre-trained classifier.

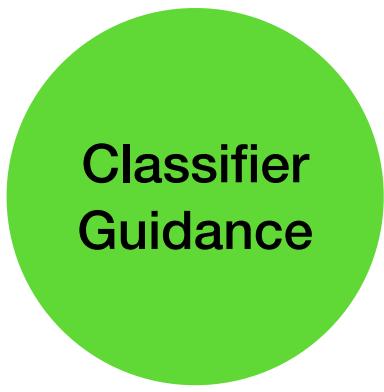
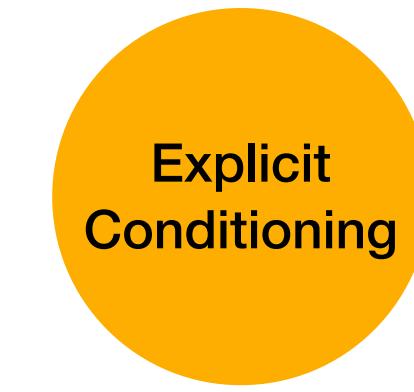
# Conditional Diffusion Models: Classifier Guidance



- **Issue:** Need to fine-tune classifier with noisy samples



# Conditional Diffusion Models: Classifier Guidance

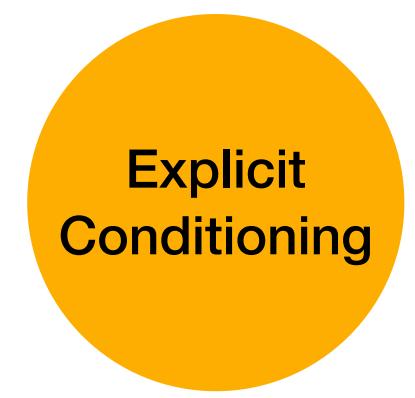


Guidance Weight 1.0



Guidance Weight 10.0

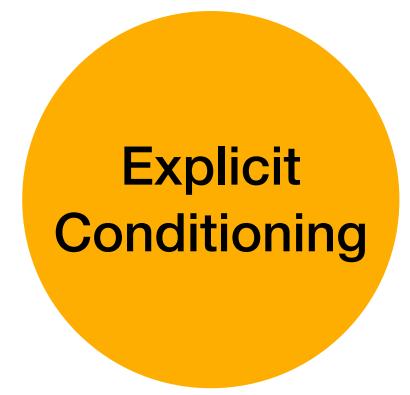
# Conditional Diffusion Models: Classifier-free Guidance



$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

Bayes' Rule

# Conditional Diffusion Models: Classifier-free Guidance



$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

Bayes' Rule

Bayes' Rule (again)

$$\nabla_{x_t} \log p(y|x_t) = \nabla_{x_t} \log p(x_t|y) - \nabla_{x_t} \log p(x_t)$$

# Conditional Diffusion Models: Classifier-free Guidance



$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

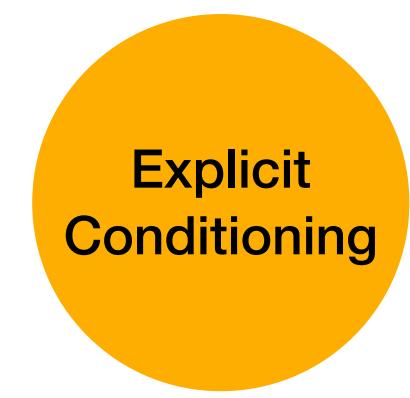
Bayes' Rule

Bayes' Rule (again)

$$\nabla_{x_t} \log p(y|x_t) = \nabla_{x_t} \log p(x_t|y) - \nabla_{x_t} \log p(x_t)$$

Standard diffusion gradient

# Conditional Diffusion Models: Classifier-free Guidance



$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

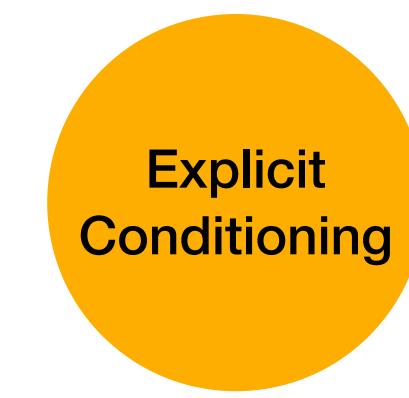
Bayes' Rule

Bayes' Rule (again)

$$\nabla_{x_t} \log p(y|x_t) = \nabla_{x_t} \log p(x_t|y) - \nabla_{x_t} \log p(x_t)$$

Conditional diffusion gradient

# Conditional Diffusion Models: Classifier-free Guidance



$$\nabla_{x_t} \log p(x_t|y) = \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t)$$

Bayes' Rule

Bayes' Rule (again)

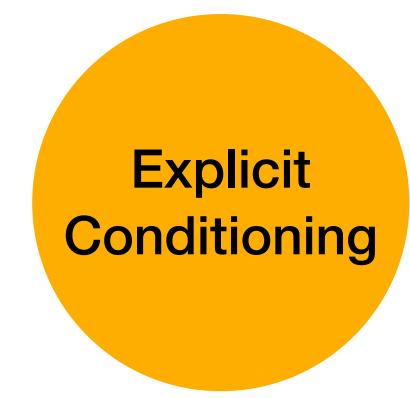
$$\nabla_{x_t} \log p(y|x_t) = \nabla_{x_t} \log p(x_t|y) - \nabla_{x_t} \log p(x_t)$$

---

$$= (1 + w) \nabla_{x_t} \log p(x_t|y) + w \nabla_{x_t} \log p(x_t)$$

- **Intuition:** instead of training a separate classifier, obtain an implicit classifier by jointly training conditional and unconditional diffusion at the same time
- Amounts to dropping the condition (e.g., zeroing it out) some fraction of the time

# Conditional Diffusion Models: Classifier-free Guidance



Non-guidance



$\omega = 1$



$\omega = 3$