



1

# Generative Model Zoo (part II)

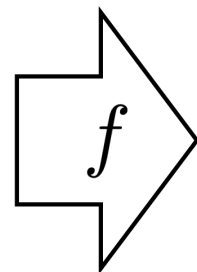
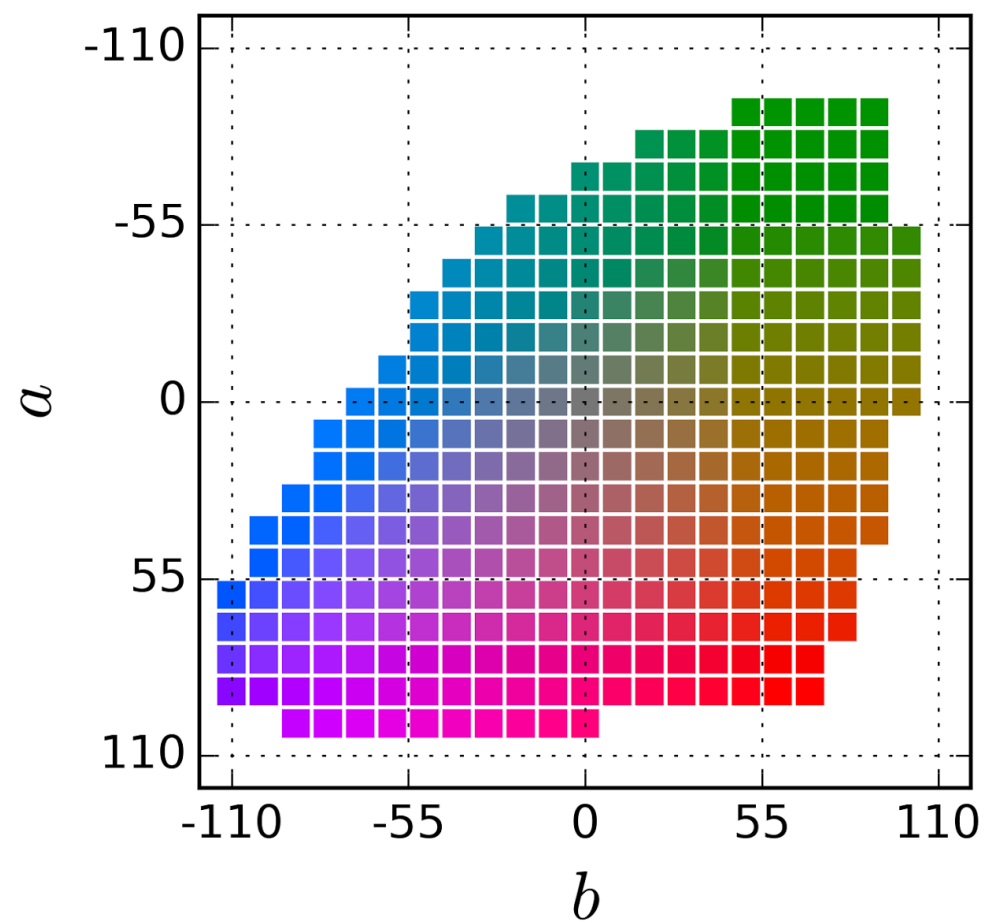
Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2023

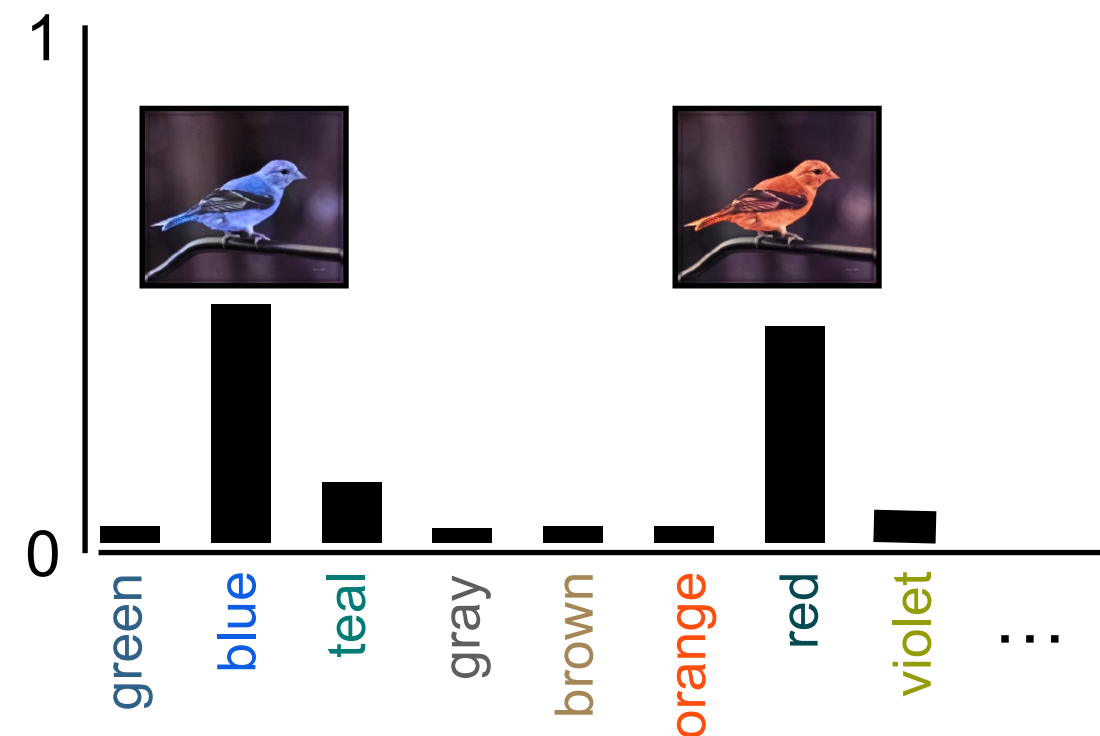
# Autoregressive Models

Recall: we can represent colors as discrete classes

$$\mathbf{y} \in \mathbb{R}^{H \times W \times K}$$

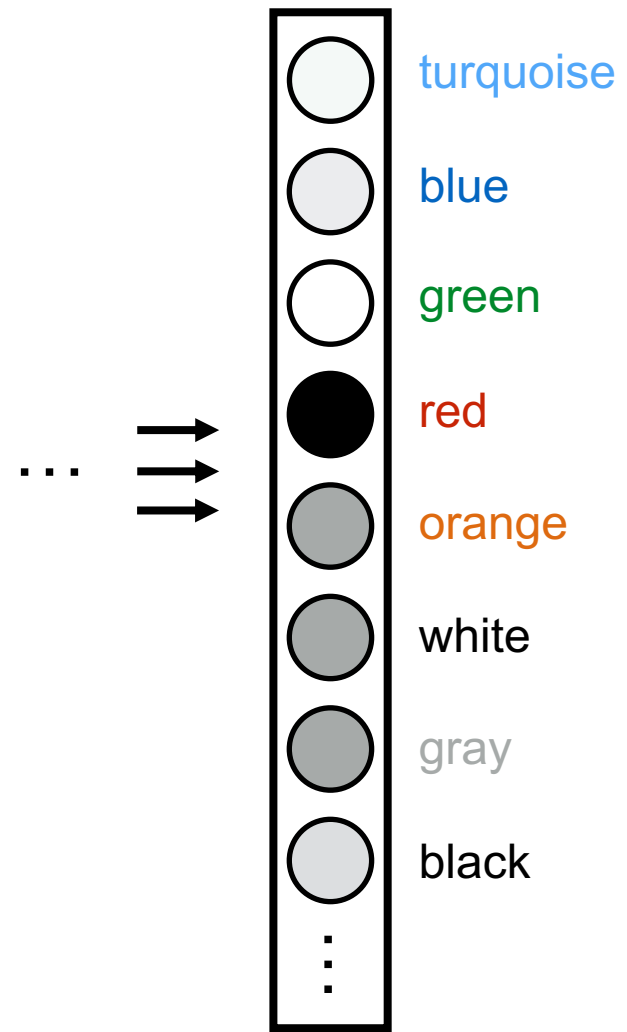


Prediction for a single pixel  $i, j$



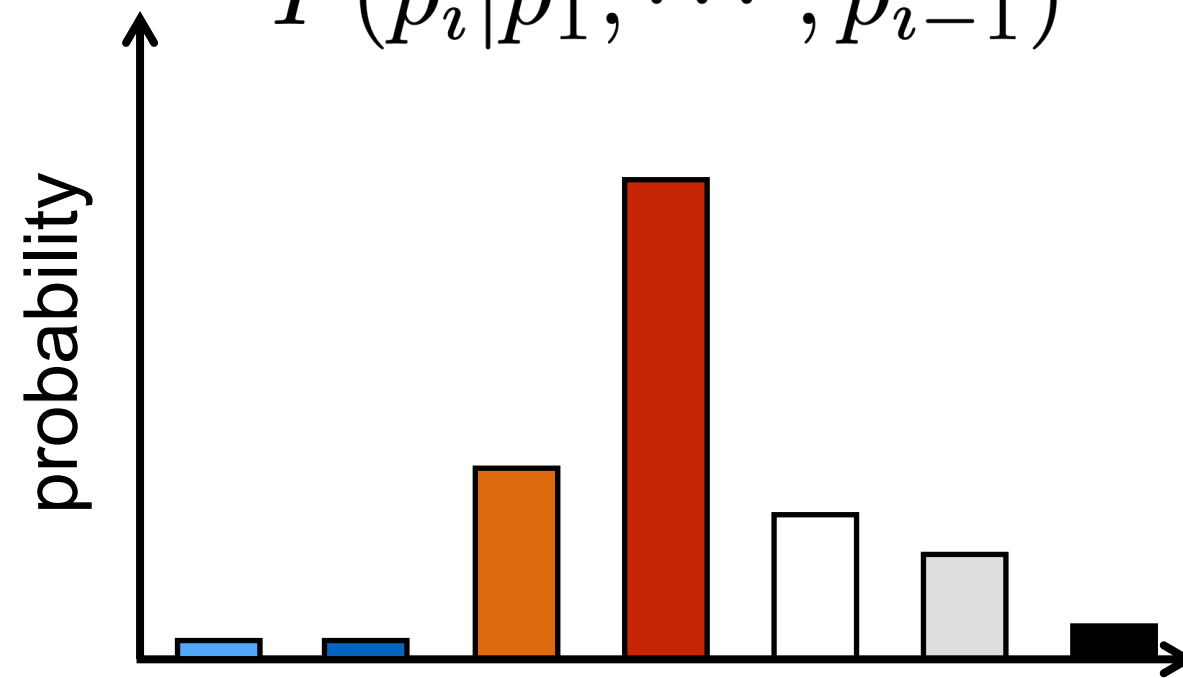
$$\mathcal{L}(\mathbf{y}, f_{\theta}(\mathbf{x})) = H(\mathbf{y}, \text{softmax}(f_{\theta}(\mathbf{x})))$$

# Network output

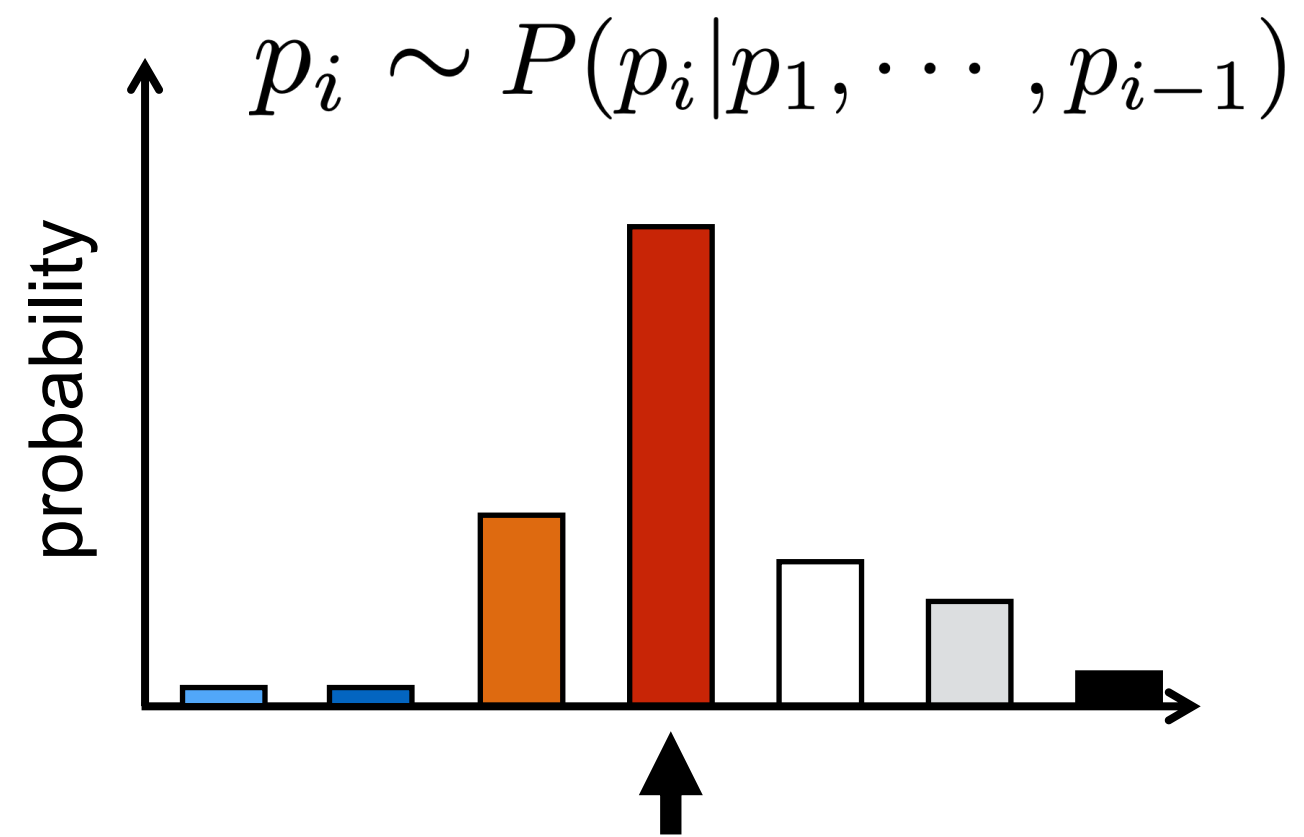
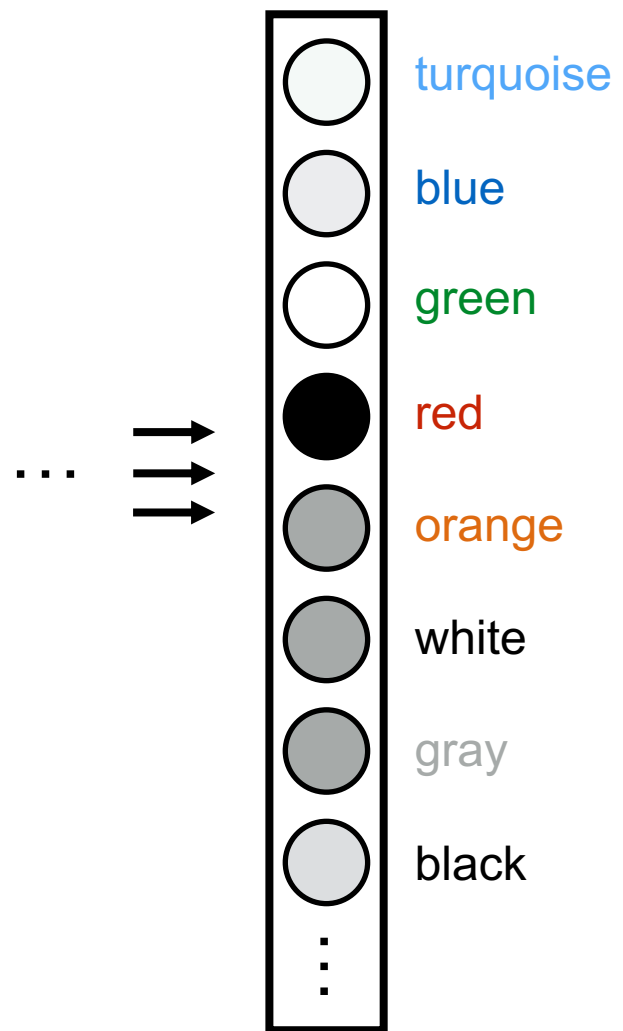


P(next pixel | previous pixels)

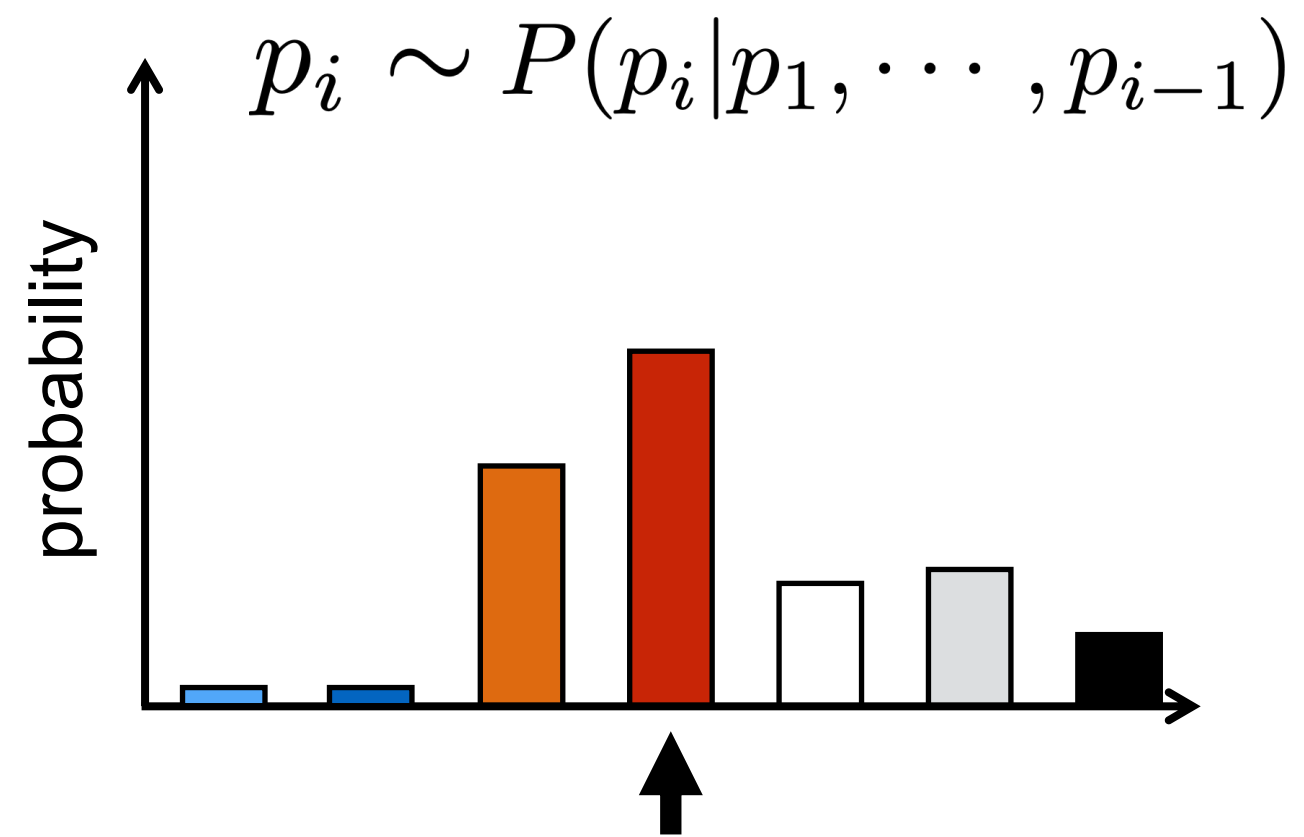
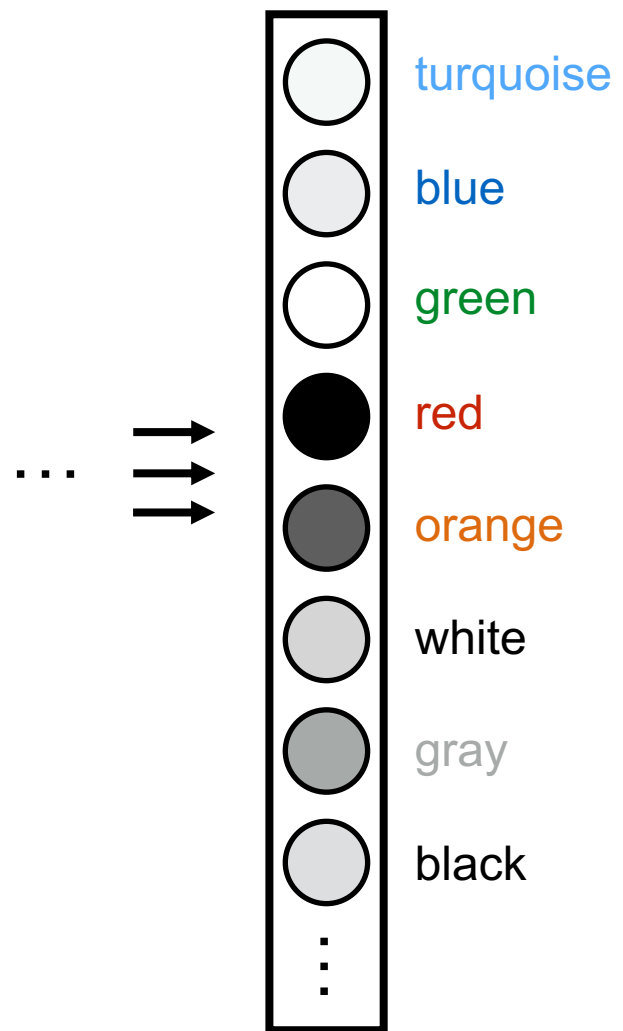
$$P(p_i | p_1, \dots, p_{i-1})$$



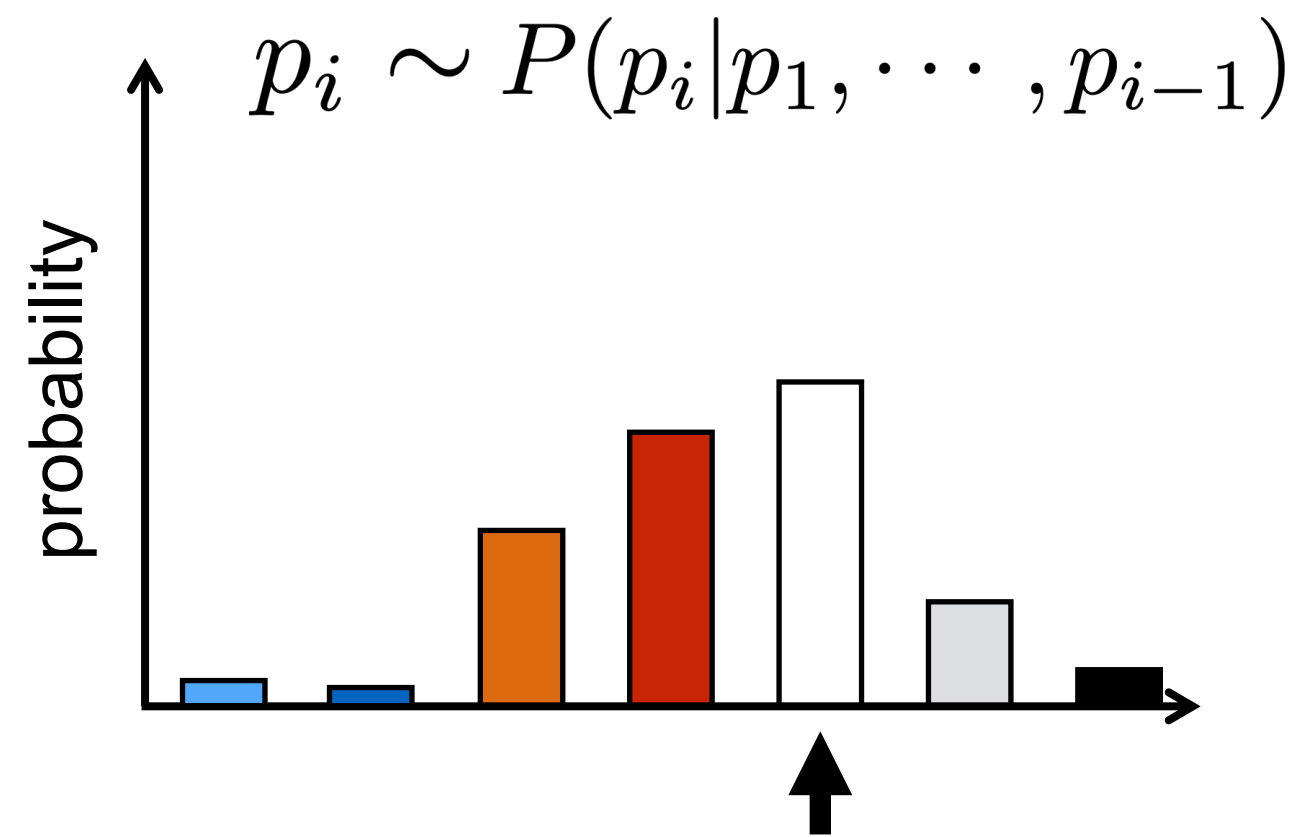
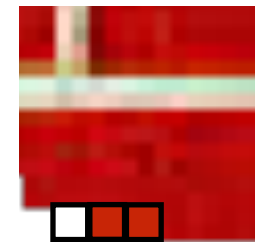
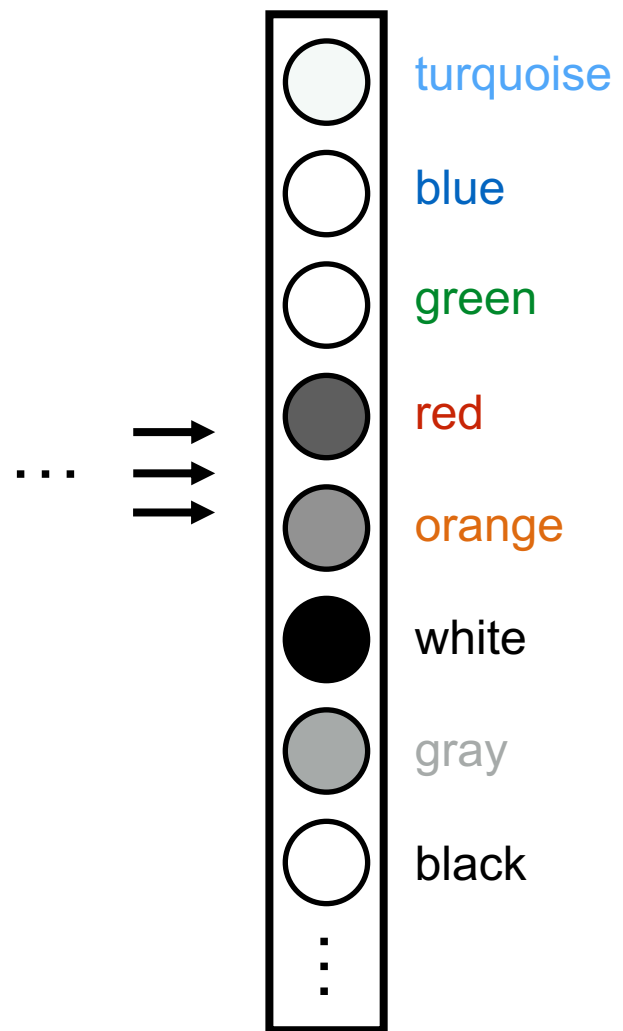
# Network output



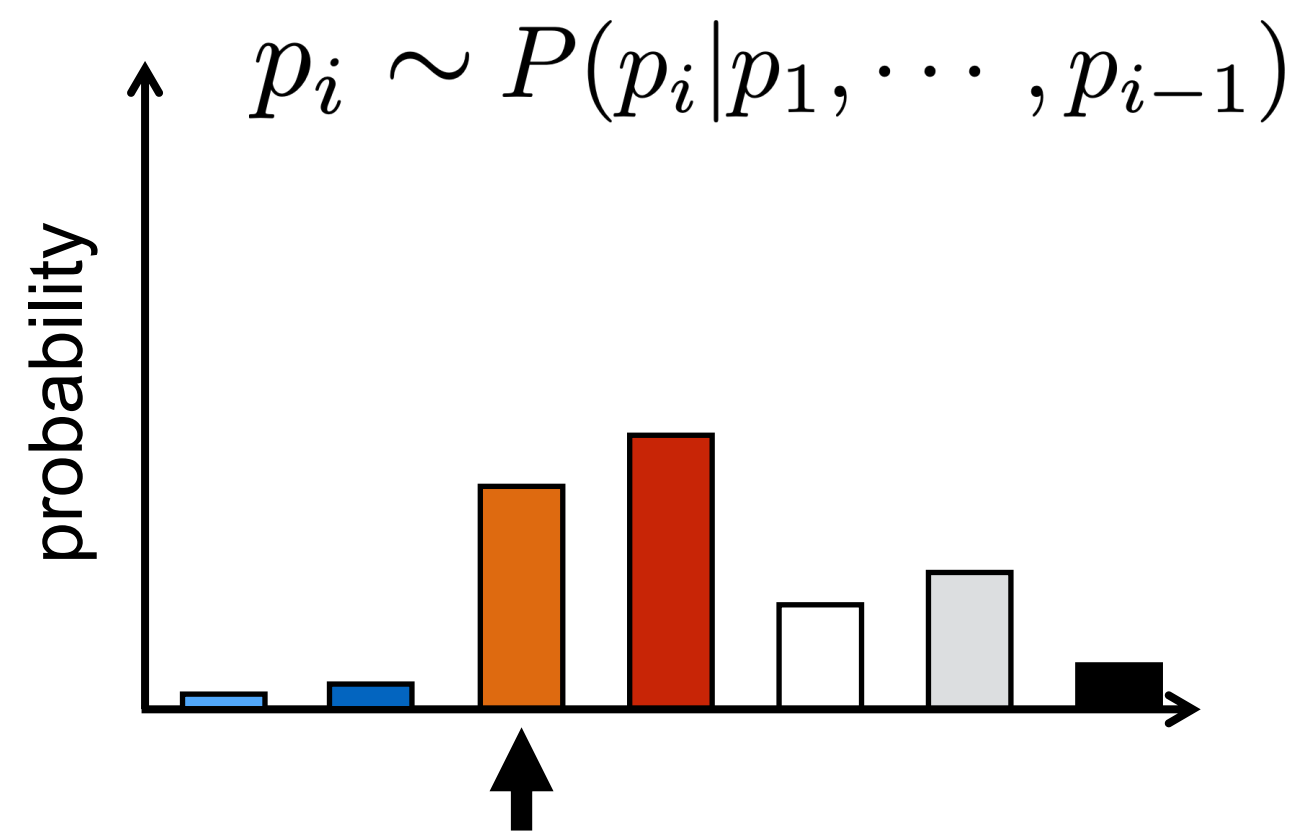
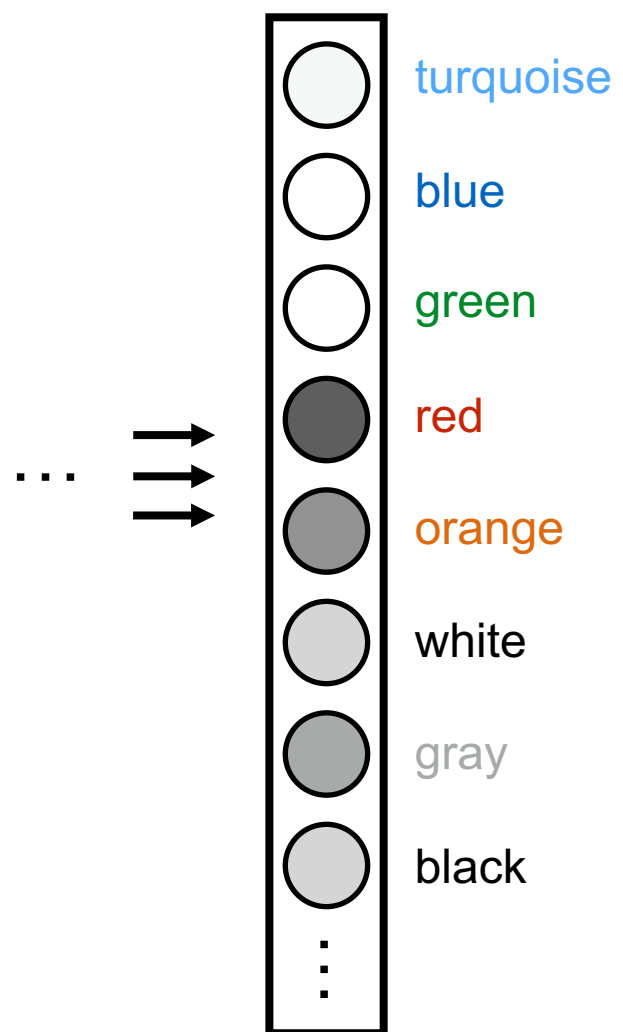
# Network output



Network output



# Network output



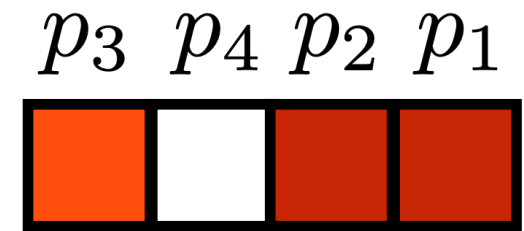


$$p_1 \sim P(p_1)$$

$$p_2 \sim P(p_2|p_1)$$

$$p_3 \sim P(p_3|p_1, p_2)$$

$$p_4 \sim P(p_4|p_1, p_2, p_3)$$



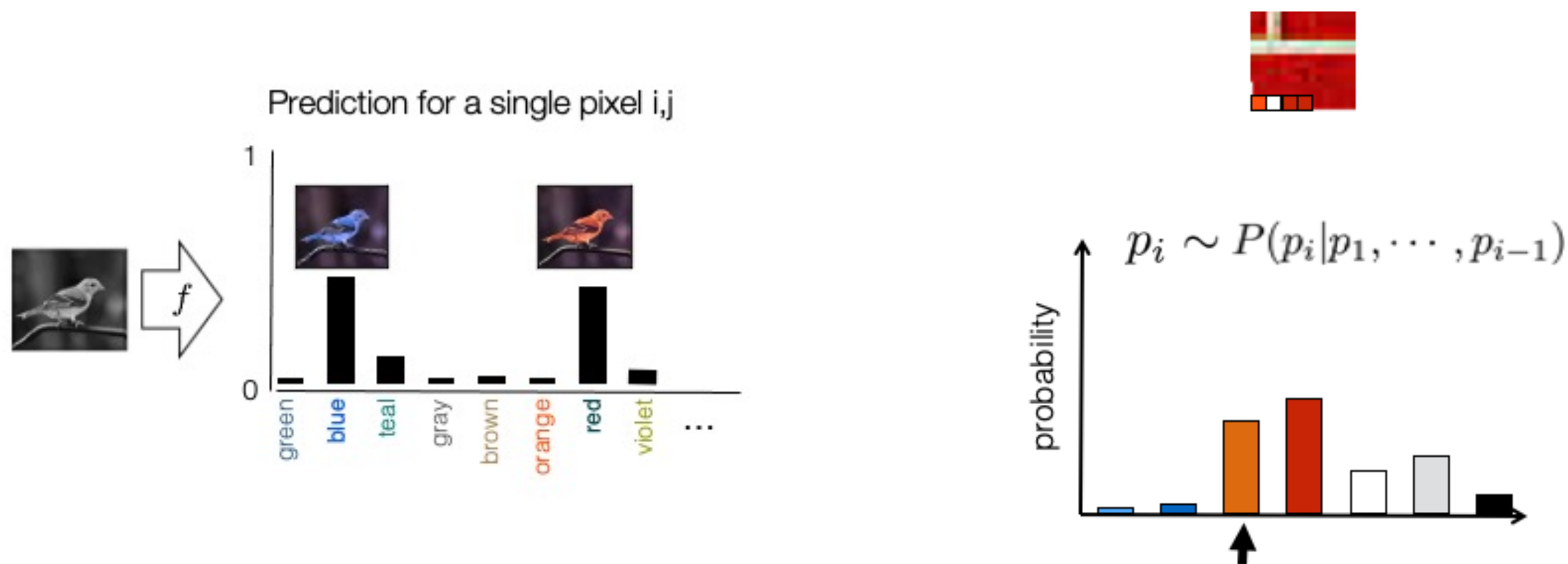
$$\{p_1, p_2, p_3, p_4\} \sim P(p_4|p_1, p_2, p_3)P(p_3|p_1, p_2)P(p_2|p_1)P(p_1)$$

$$p_i \sim P(p_i|p_1, \dots, p_{i-1})$$

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i|p_1, \dots, p_{i-1})$$

# Per-pixel classification vs. Autoregressive

- Image colorization: per-pixel classification loss
- PixelCNN, VQ-VAE2: autoregressive model
- Key idea: it can only produce discrete representation (e.g., VQ codes, color bins)



# Autoregressive model

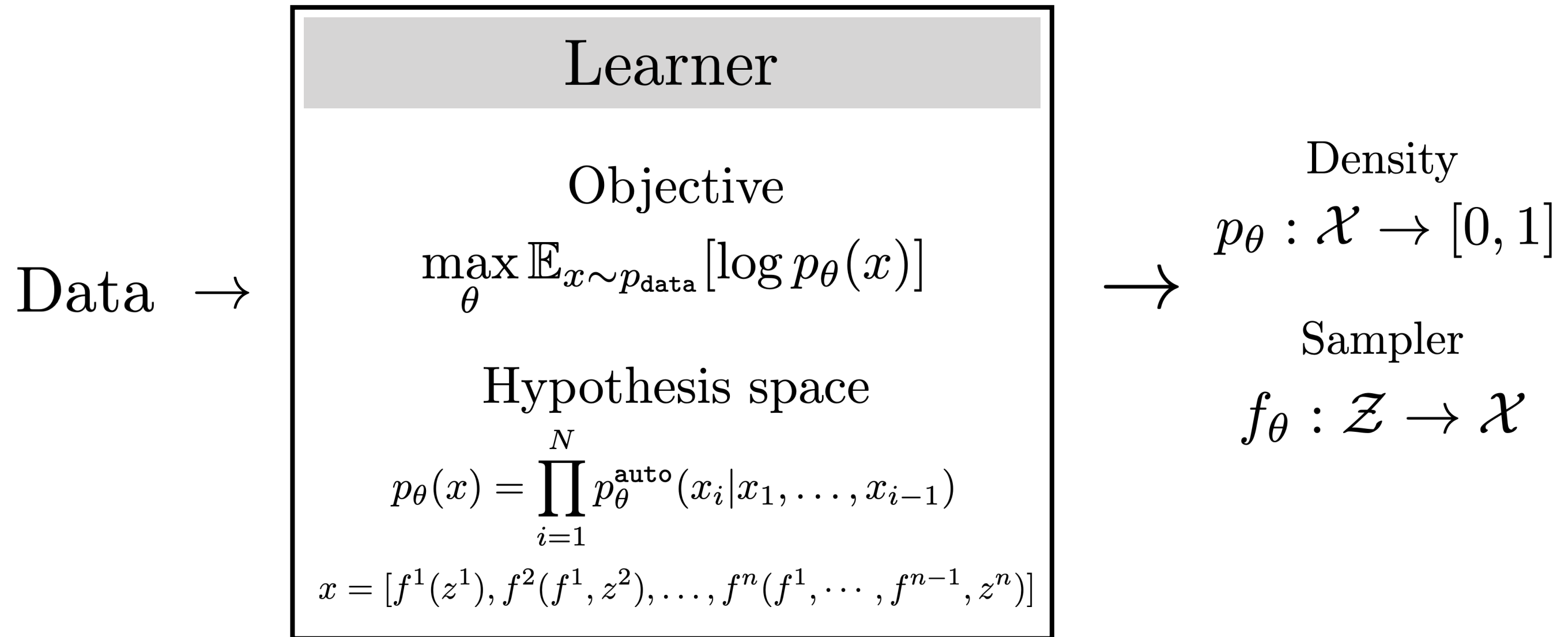
$$\mathbf{p} \sim \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1})$$

$$P(\mathbf{p}) = \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1}) \quad \leftarrow \text{General product rule}$$

The sampling procedure we defined above takes exact samples from the learned probability distribution.

Multiplying all conditionals evaluates the probability of a full joint configuration of pixels.

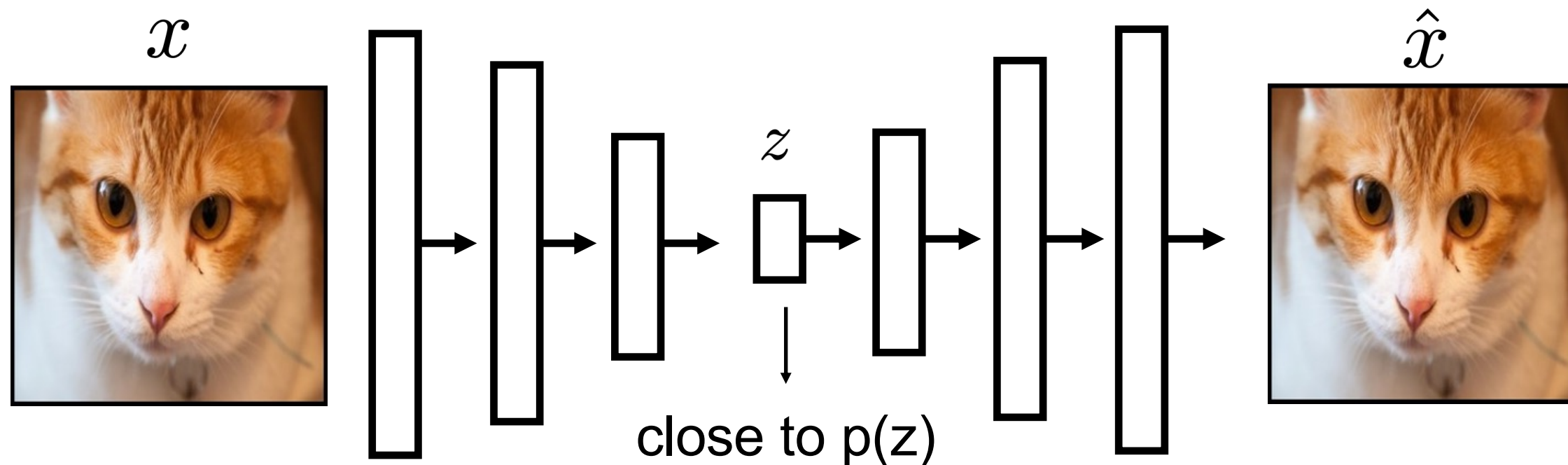
# Autoregressive model



# Variational Auto-encoder

# Variational Autoencoders (VAEs)

encoder  $q_{\psi}(z|x)$   $z = E_{\psi}^{\mu}(x) + E_{\psi}^{\sigma}(x) \cdot \epsilon_z$       generator  $p_{\theta}(x|z)$   $\hat{x} = G_{\theta}^{\mu}(z)$



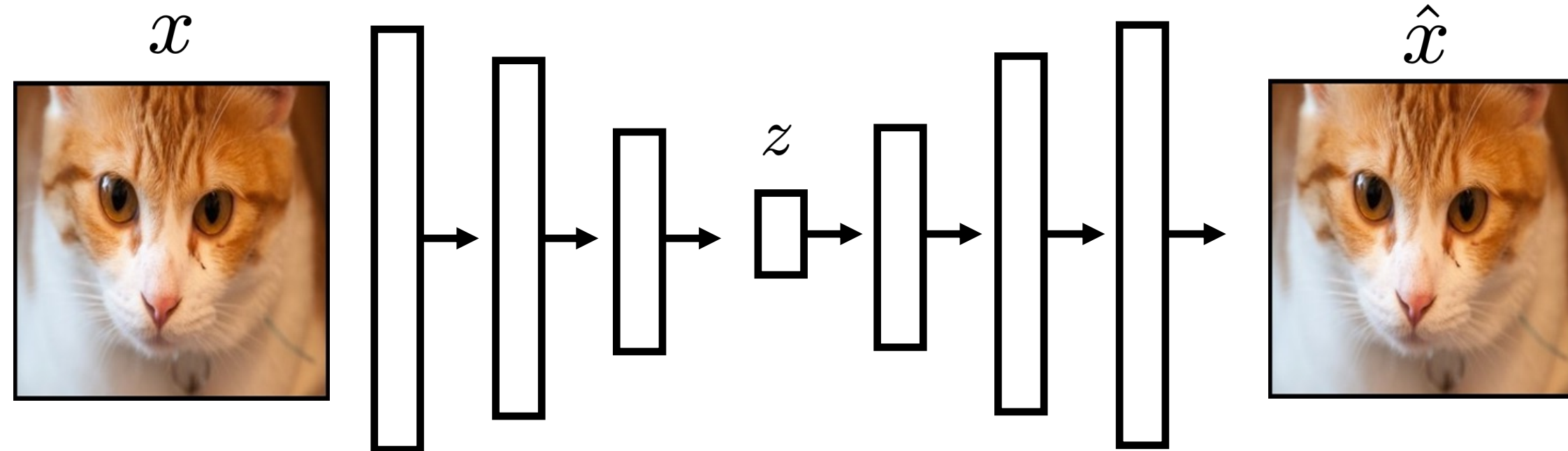
ELBO: Evidence Lower Bound

$$\begin{aligned} & \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)] && \text{Multi-variate Gaussian} \\ & \geq \max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)] - \text{KL}(q_{\psi}(z|x_i) || p(z))] \\ & \quad \uparrow && \uparrow \\ & \text{reconstruction loss} && \text{KLD loss} \\ & ||x - \hat{x}||_2 && \text{KLD}(\mathcal{N}(E_{\psi}^{\mu}(x), E_{\psi}^{\sigma}(x)) | \mathcal{N}(0, I)) \end{aligned}$$

# Autoencoders (AEs)

encoder  $z = E_{\psi}^{\mu}(x)$   
 $q_{\psi}(z|x)$

generator  $\hat{x} = G_{\theta}^{\mu}(z)$   
 $p_{\theta}(x|z)$

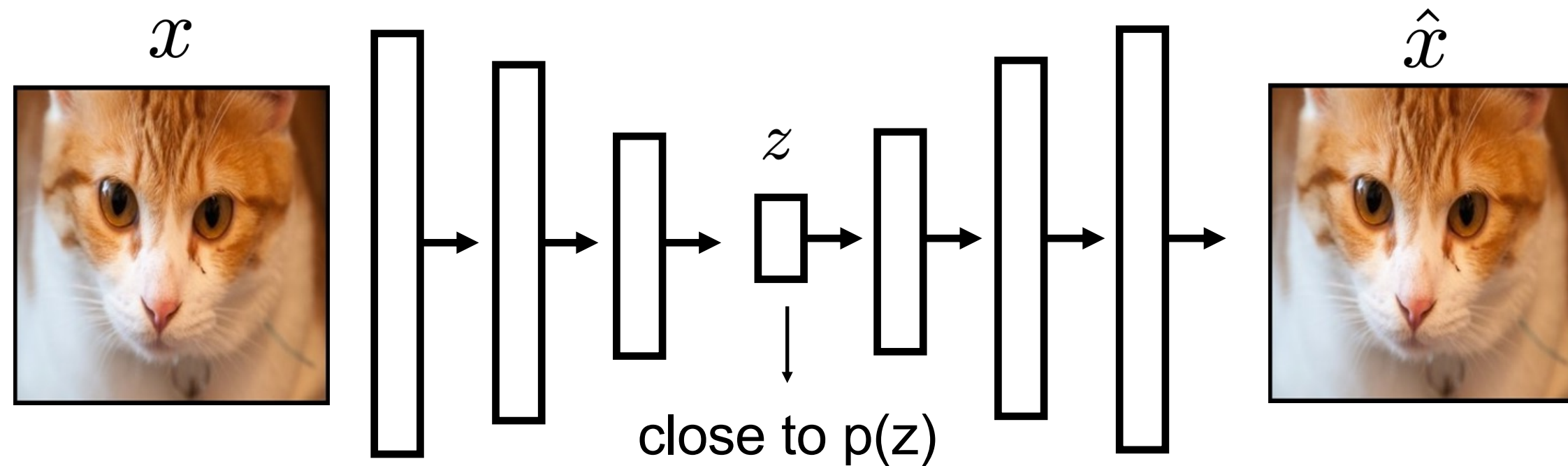


$$\max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)]]$$

↑  
reconstruction loss

$$\|x - \hat{x}\|_2$$

# Autoencoders (AEs) + Easier Sampling

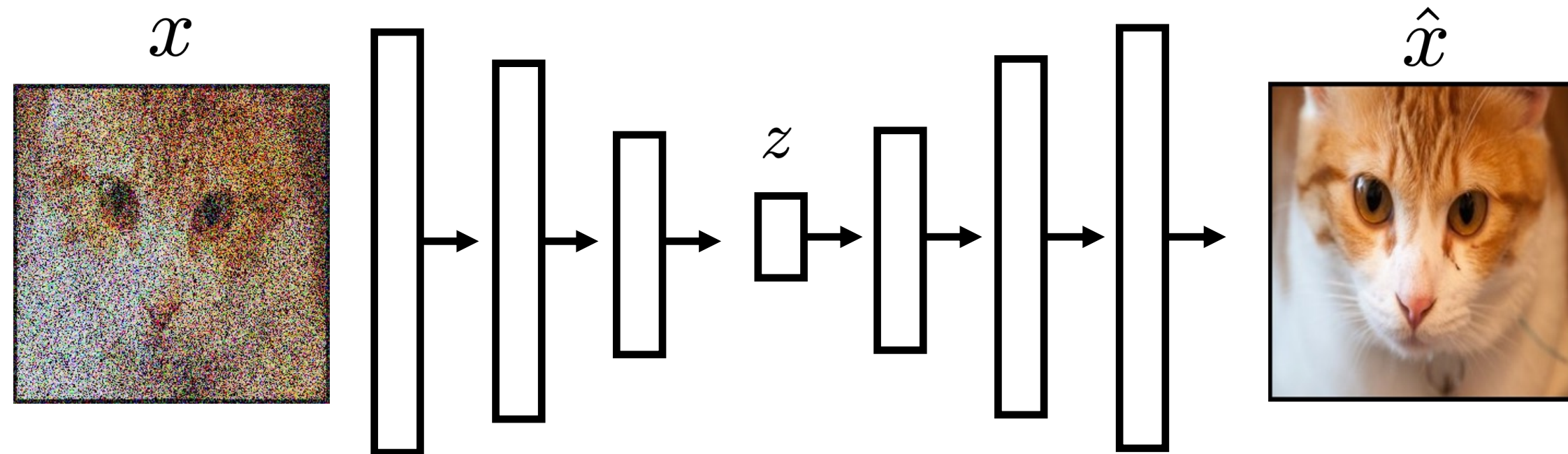


$$\max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)] - \text{KL}(q_{\psi}(z|x_i) || p(z))]$$

↑ reconstruction loss                      ↑ KLD loss



# Denoising Autoencoders (AEs)

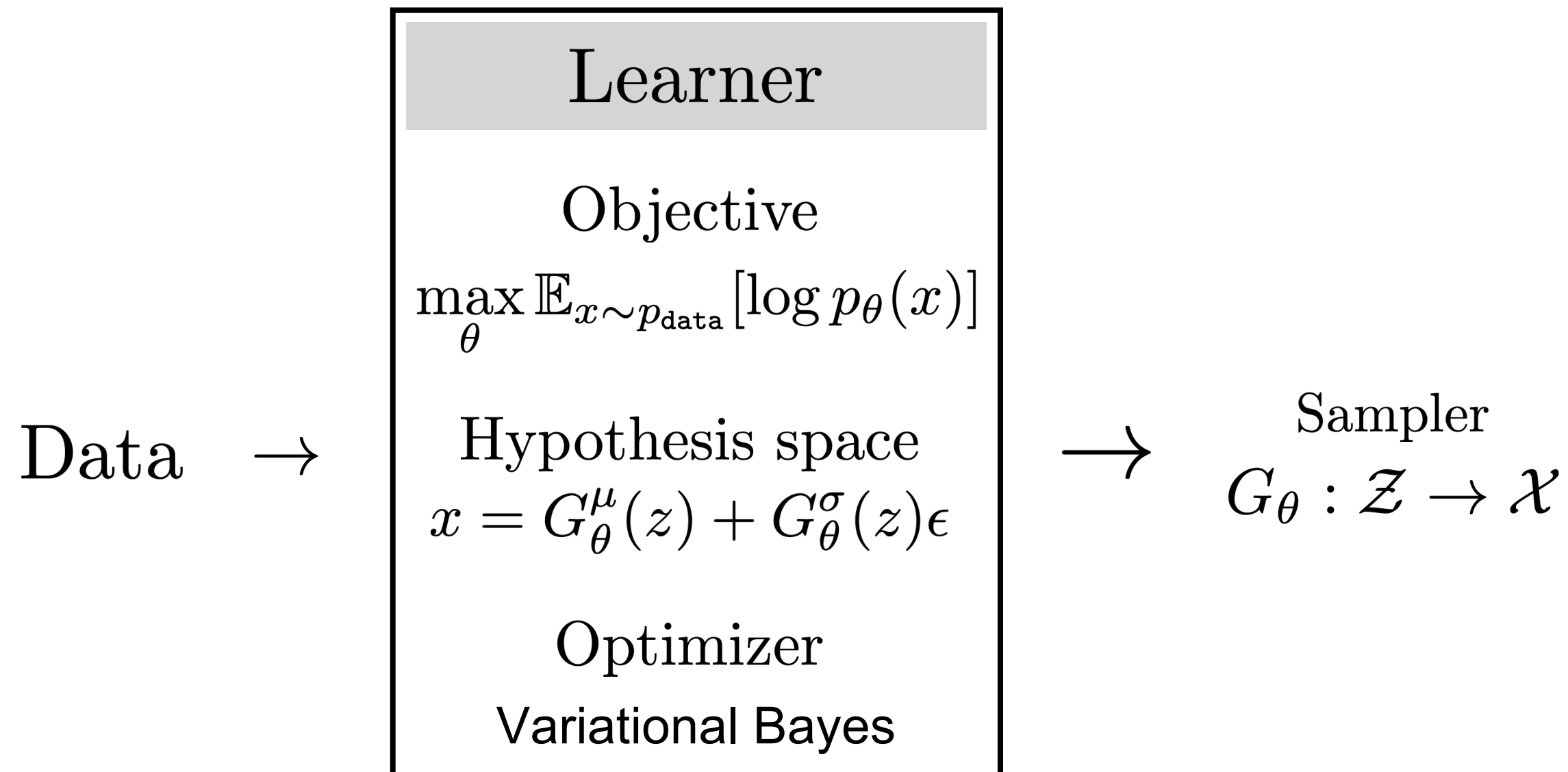


$$\max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)]]$$

↑  
reconstruction loss  
corrupt input

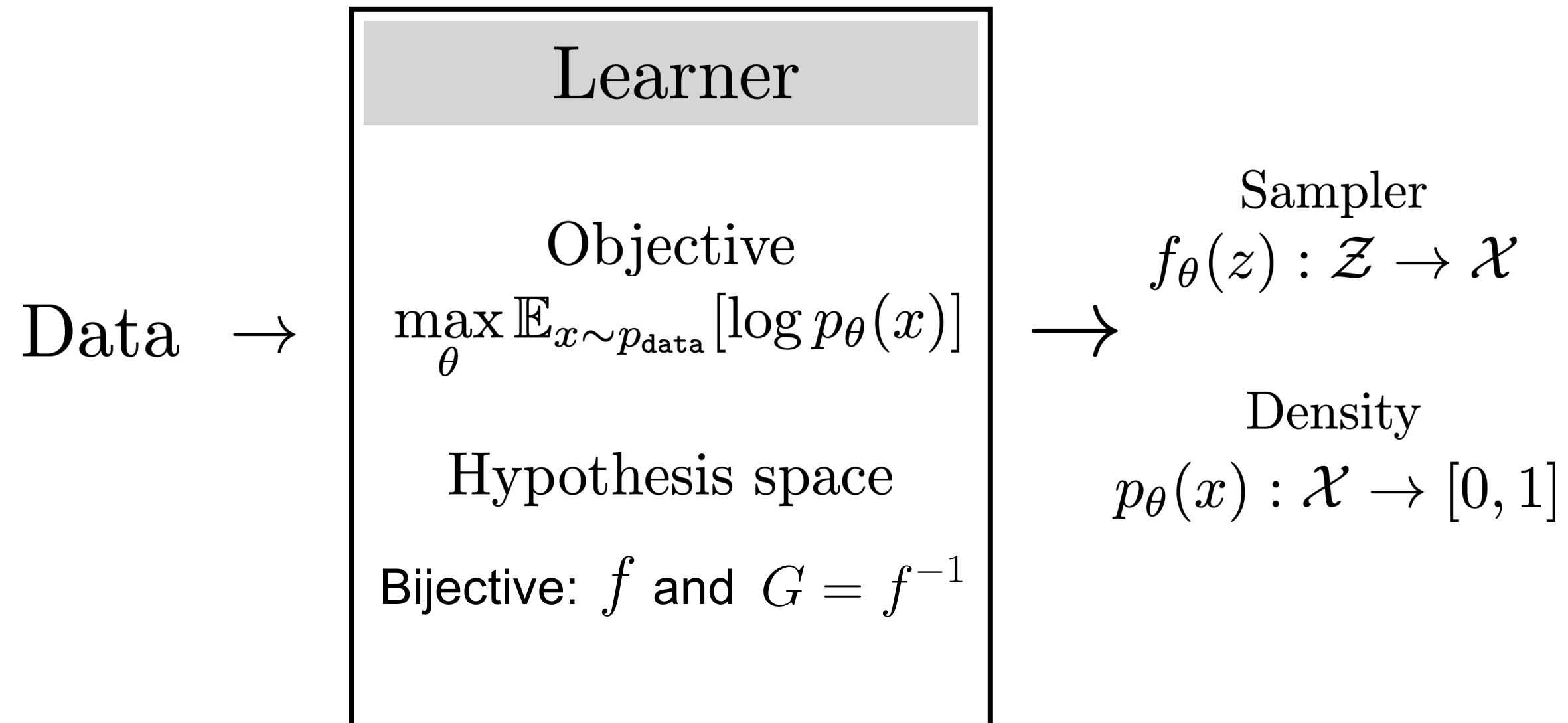
Denoising vs.  
Compression

# Variational Autoencoder (VAE)



# Flow-based Models

# Flow-based models



- x and z have the same number of dimensions (memory; training speed)
- limited choices of f and G
- + Fast sample; accurate density estimate

# Flow-based models

- Density estimate

$$x \sim p_{data}(x)$$

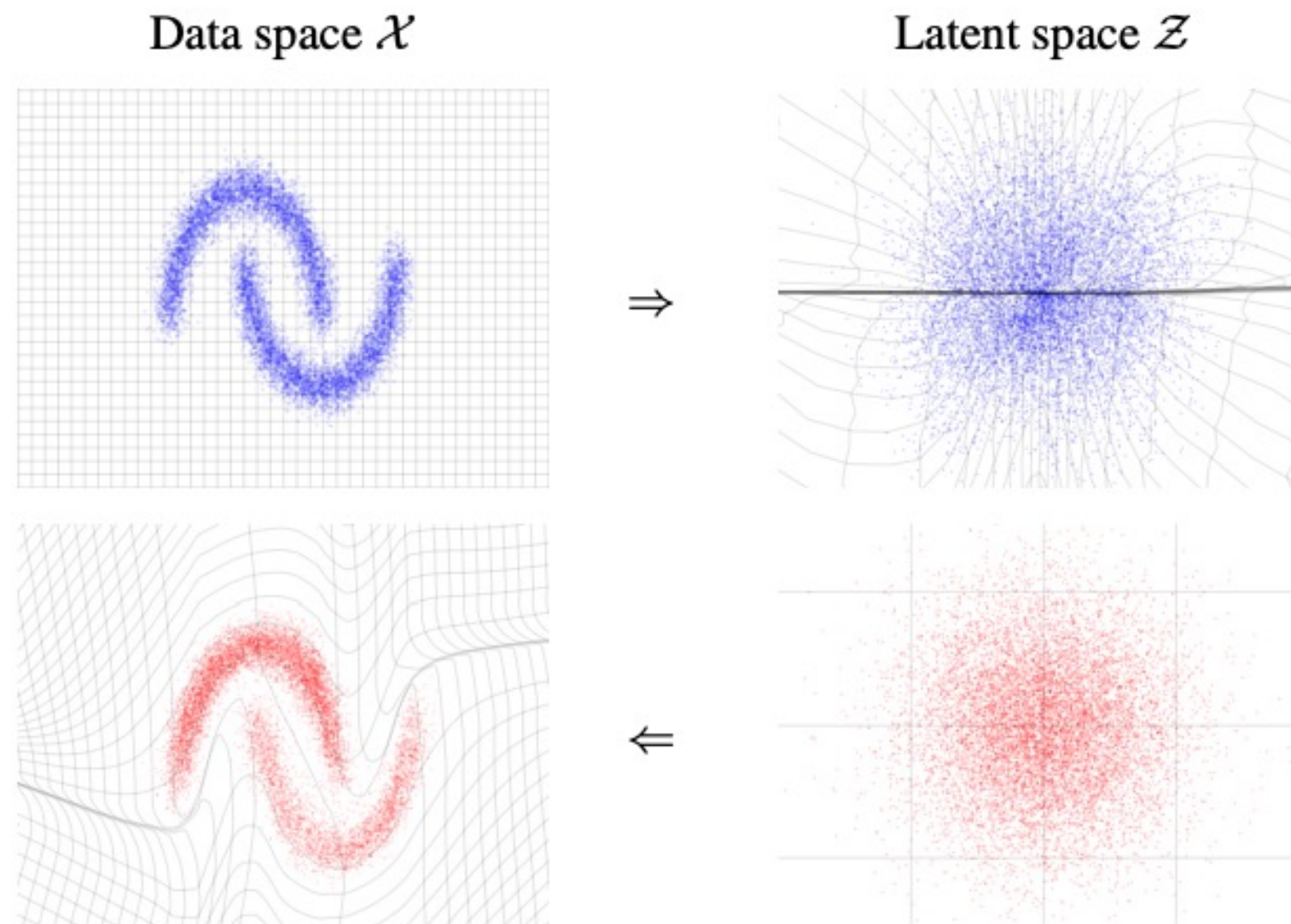
$$z = f(x)$$

- Sampling

$$z \sim p(z)$$

$$x = G(z)$$

Generator  $G = f^{-1}$



# Flow-based models

## Training objective

- Density estimate

$$x \sim p_{data}(x)$$

$$z = f(x)$$

Change of variable formula

$$p_{data}(x) = p_z(f(x)) \left| \det\left(\frac{\partial f(x)}{\partial x^T}\right) \right|$$

$$\log p_{data}(x) = \log(p_z(f(x))) + \log\left(\left| \det\left(\frac{\partial f(x)}{\partial x^T}\right) \right|\right)$$

- Sampling

$$z \sim p(z)$$

$$x = G(z)$$

Generator  $G = f^{-1}$

Easy to compute  
as  $z$  follows Gaussian distribution

hard to compute  
Jacobian determinant  
for most layers

design layers whose Jacobian  
is a triangular matrix

# Design invertible layers

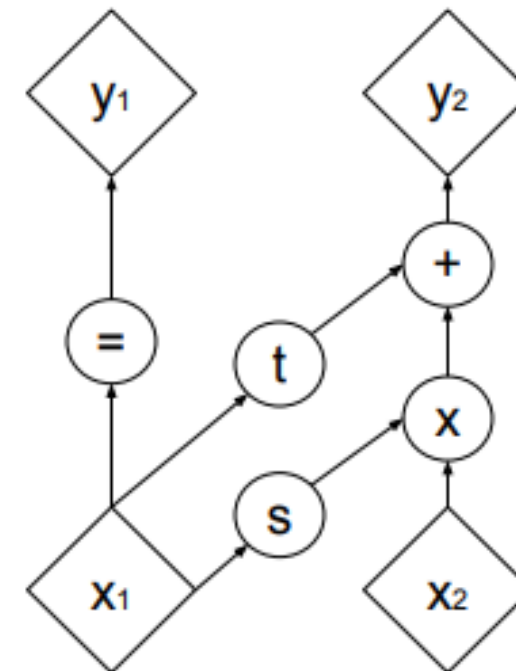
$$y_{1:d} = x_{1:d}$$

$$y_{d+1:D} = x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d}),$$

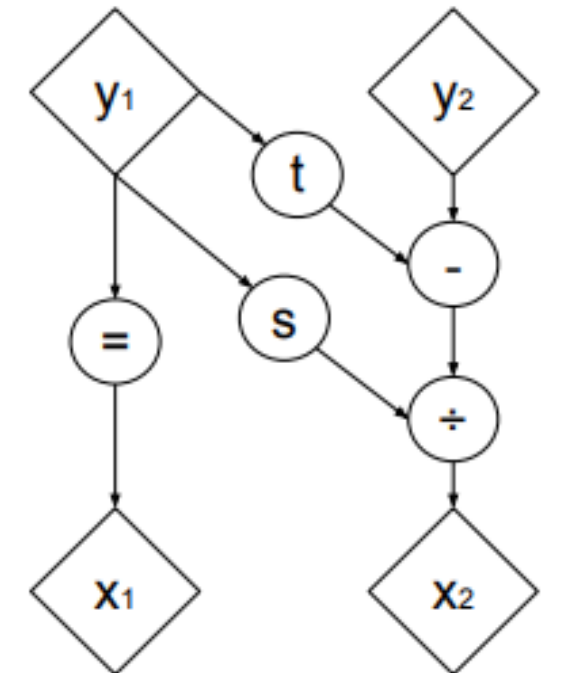
Jacobian matrix

$$\frac{\partial y}{\partial x^T} = \begin{bmatrix} \mathbb{I}_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}^T} & \text{diag}(\exp[s(x_{1:d})]) \end{bmatrix},$$

s, t: affine transformation  
with parameters



(a) Forward propagation



(b) Inverse propagation

# Flow-based Models



Real NVP [Dinh et al., ICLR 2017]

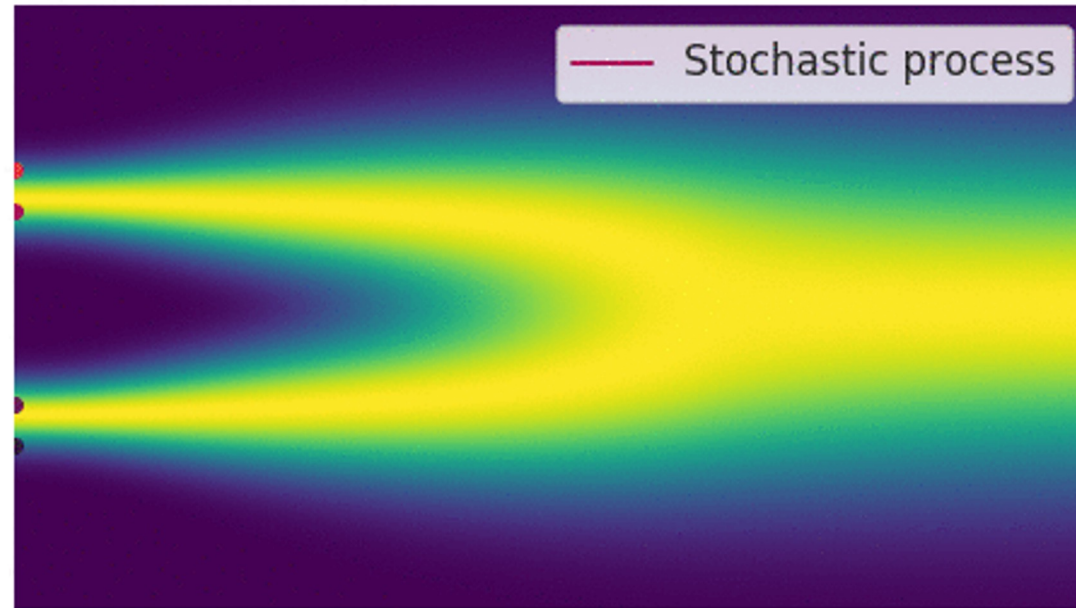


Glow [Kingma and Dhariwal, NeurIPS 2018 ]

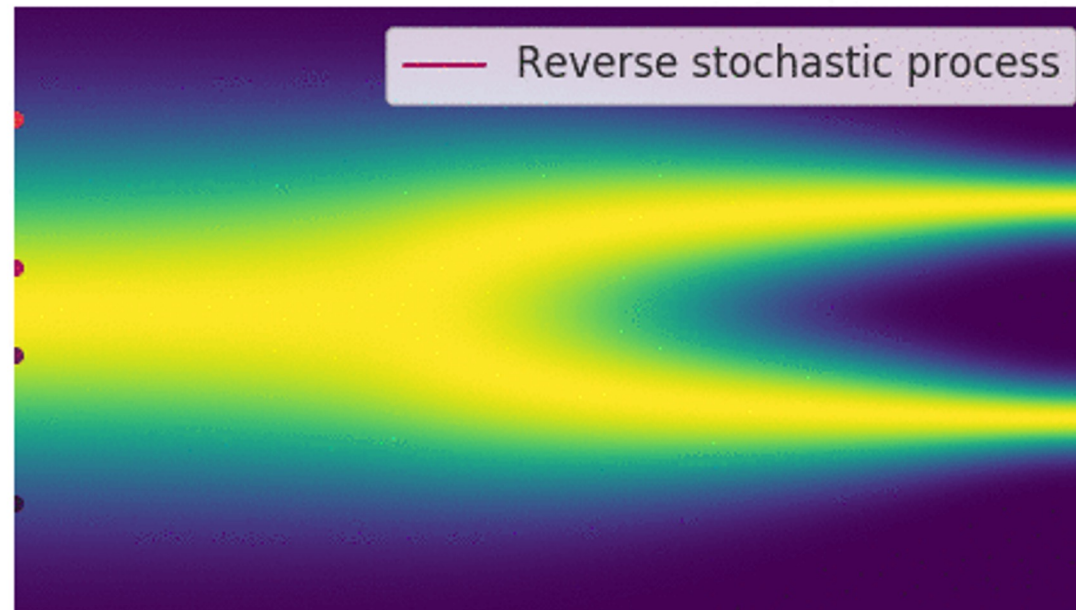


# Diffusion Models

# Diffusion Models



- “destroy” the data by gradually adding small amounts of gaussian noise



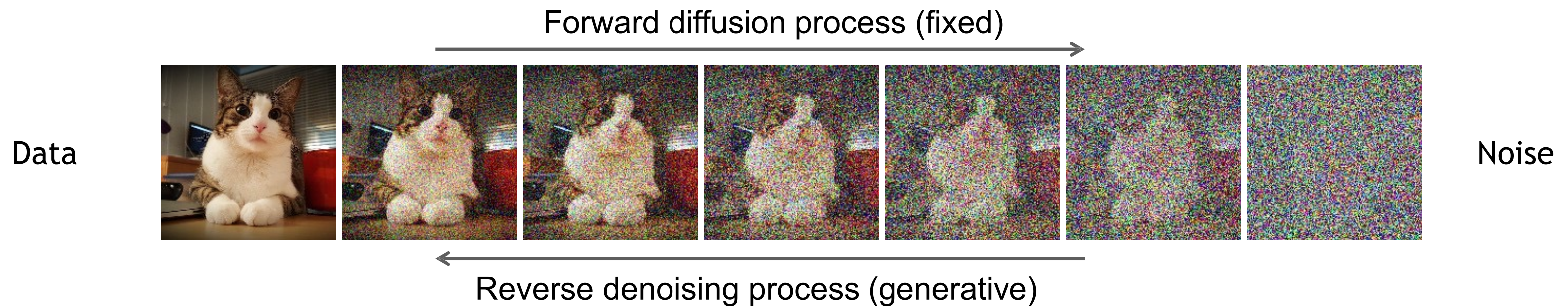
- “create” data by gradually denoising a noisy code from a stationary distribution

# Denoising Diffusion Models

Learning to generate by denoising

Denoising diffusion models consist of two processes:

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



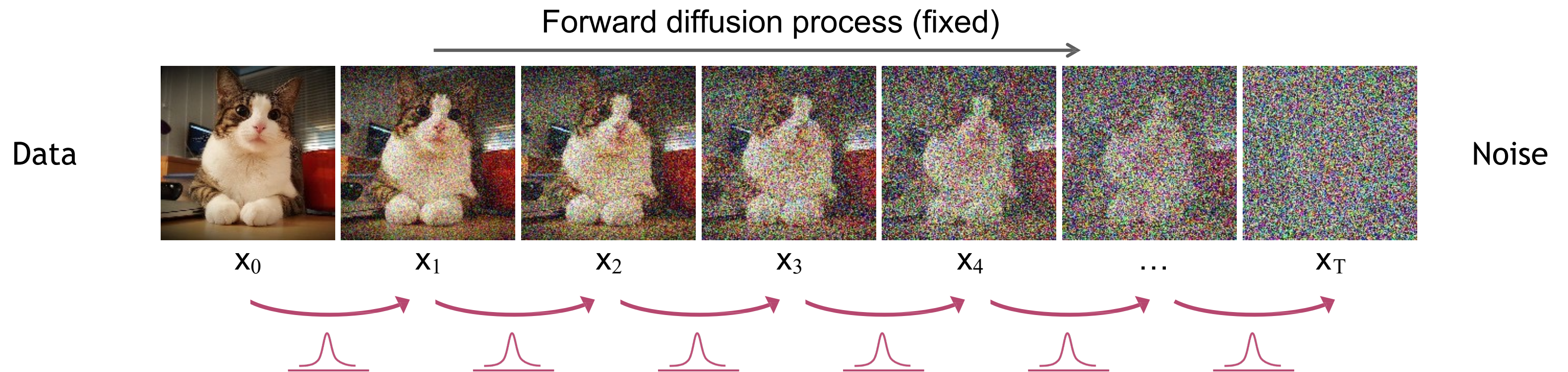
[Sohl-Dickstein et al., Deep Unsupervised Learning using Nonequilibrium Thermodynamics, ICML 2015](#)

[Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020](#)

[Song et al., Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021](#)

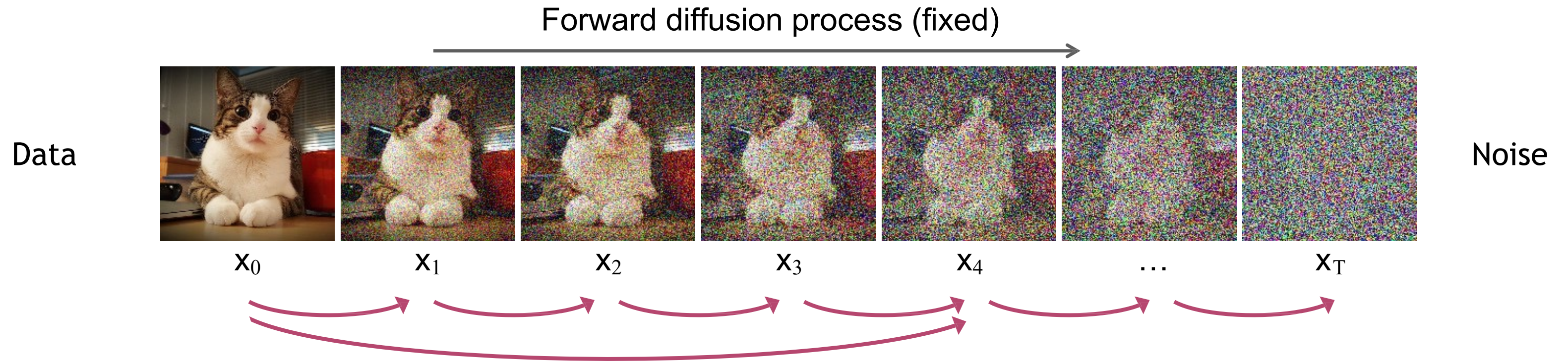
# Forward Diffusion Process

The formal definition of the forward process in T steps:



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

# Diffusion Kernel



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

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

$\beta_t$  values schedule (i.e., the noise schedule) is designed such that  $\bar{\alpha}_T \rightarrow 0$  and  $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

## Mathematic details: Merge Multiple Gaussian

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

## Mathematic details: Merge Multiple Gaussian

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

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1}$$

## Mathematic details: Merge Multiple Gaussian

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

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2}\end{aligned}$$



# Mathematic details: Merge Multiple Gaussian

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

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}\end{aligned}$$

# Mathematic details: Merge Multiple Gaussian

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

$$\begin{aligned}\mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}\end{aligned}$$

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

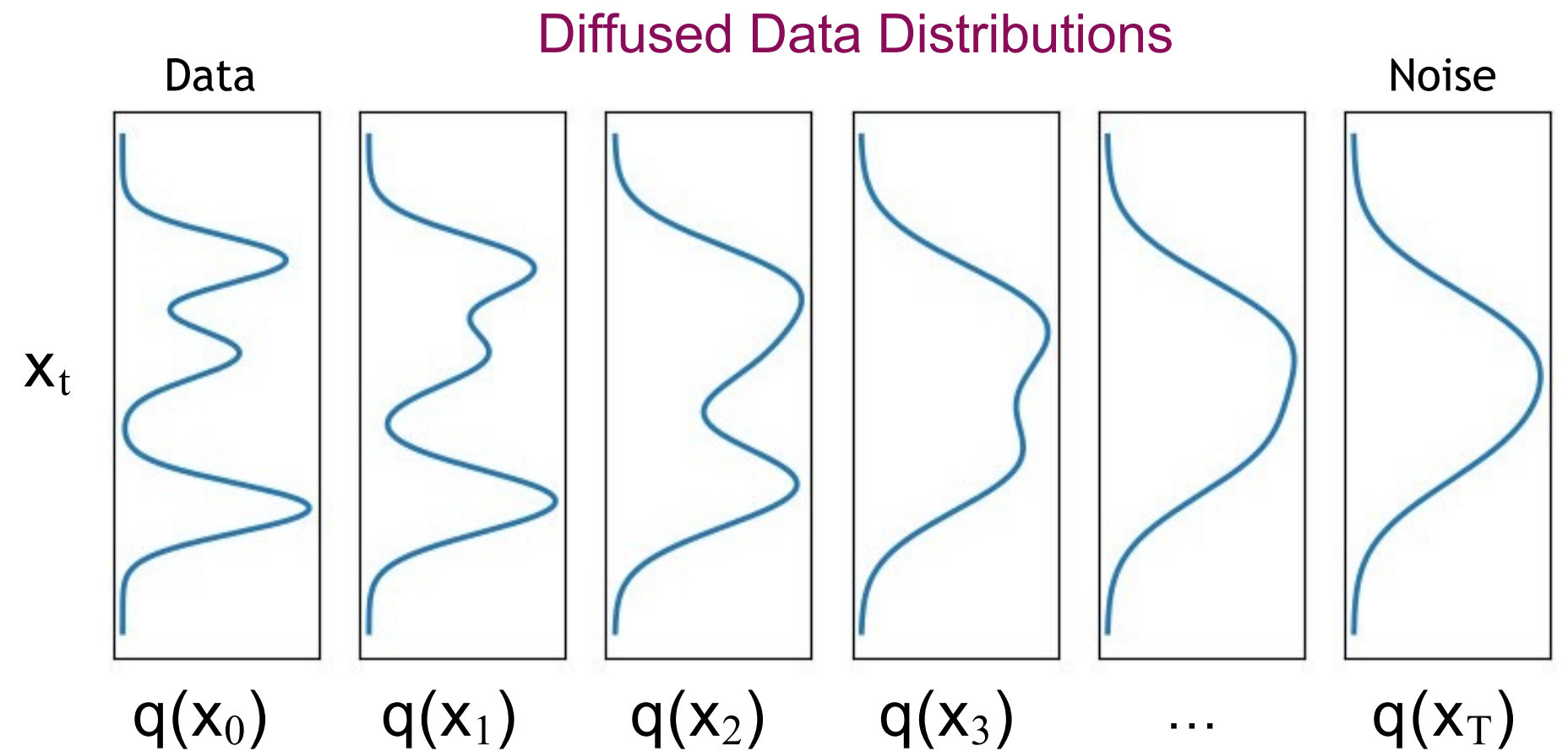
Direct sampling from  $0 \rightarrow t$

# What happens to a distribution in the forward diffusion?

So far, we discussed the diffusion kernel  $q(\mathbf{x}_t|\mathbf{x}_0)$  but what about  $q(\mathbf{x}_t)$ ?

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

The diffusion kernel is Gaussian convolution.



We can sample  $\mathbf{x}_t \sim q(\mathbf{x}_t)$  by first sampling  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$  and then sampling  $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$

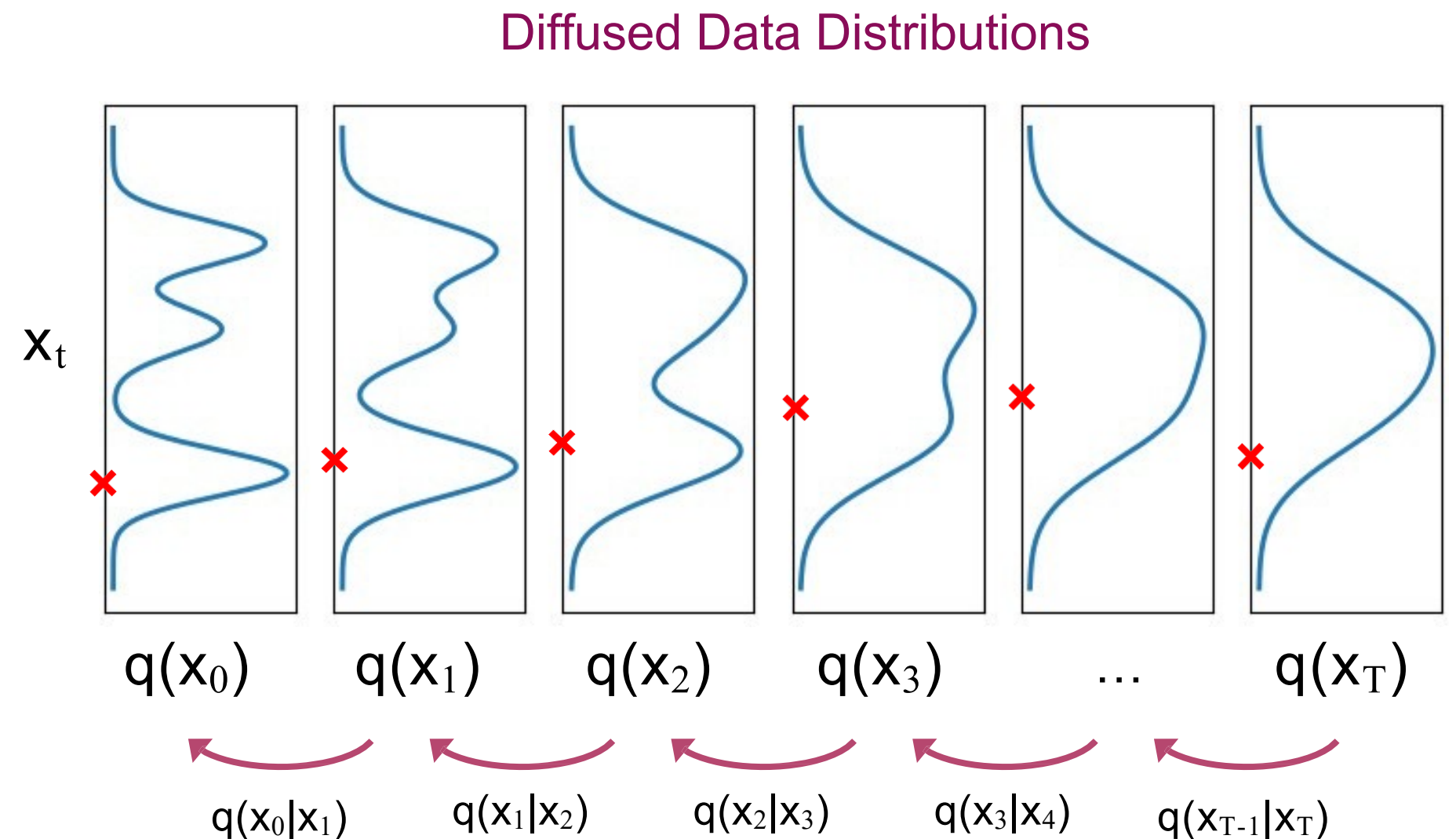
# Generative Learning by Denoising

Recall, that the diffusion parameters are designed such that  $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

**Generation:**

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

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

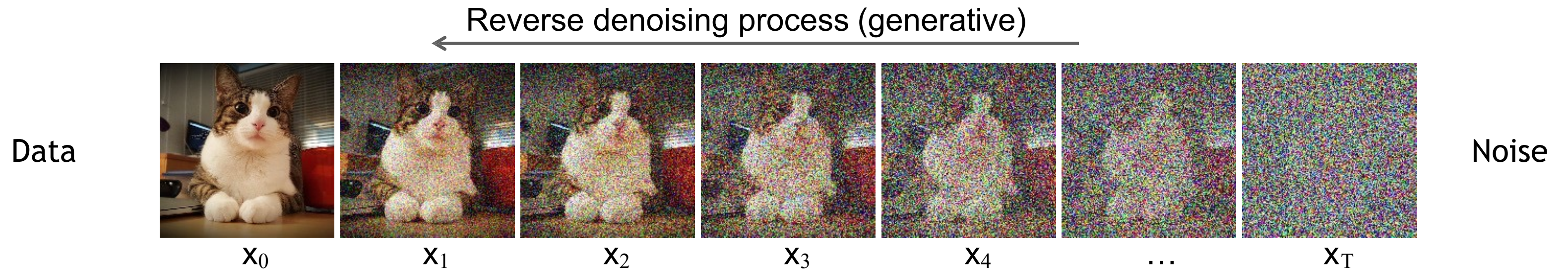


In general,  $q(\mathbf{x}_{t-1}|\mathbf{x}_t) \propto q(\mathbf{x}_{t-1})q(\mathbf{x}_t|\mathbf{x}_{t-1})$  is intractable.

Can we approximate  $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ ? Yes, we can use a **Normal distribution** if  $\beta_t$  is small in each forward diffusion step.

# Reverse Denoising Process

Formal definition of forward and reverse processes in T steps:



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

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

Trainable network  
(U-net, Denoising Autoencoder)