# Strojové učení II

## Generative Models

# Supervised vs Unsupervised Learning

## Supervised

- Data: (x,y)

  x…data, y…label

- Goal: learn function to map

  $$x \rightarrow y$$

- Examples: classification, regression, object detection, semantic segmentation, ...
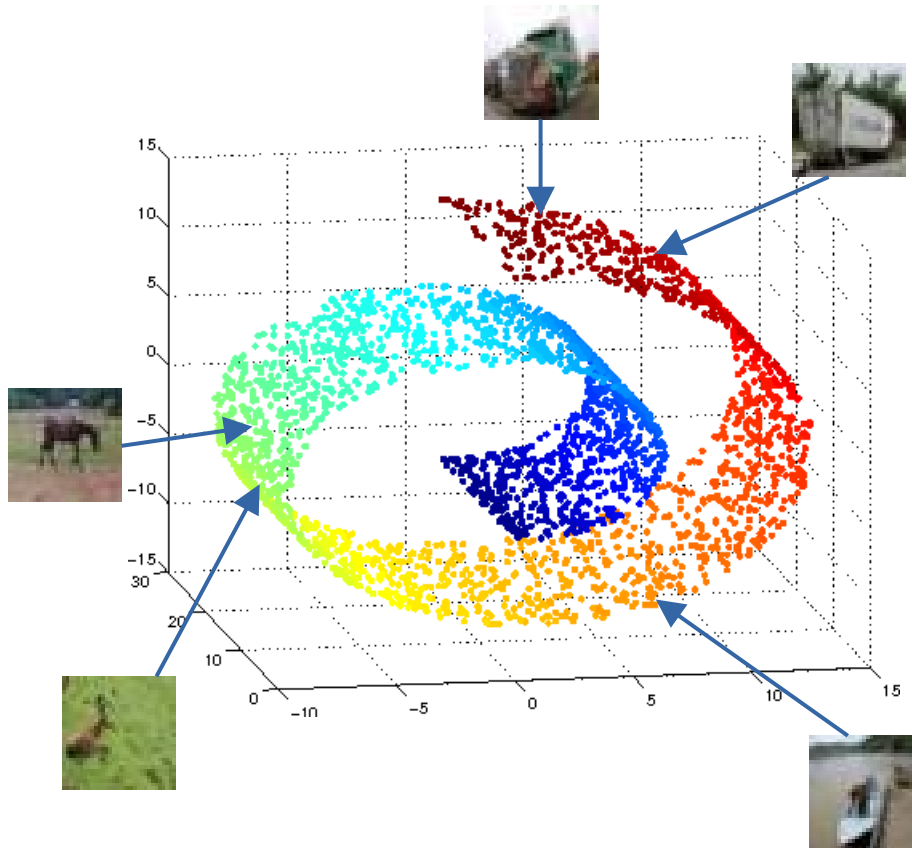
## Unsupervised

- Data: x

  x...data, no labels!

- Goal: learn hidden (underlying) structure of the data

- Examples: clustering, feature or dimensionality reduction, …

# Generative Modeling

- Natural images lie on a manifold



Input samples

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

- How to sample from this distribution?

# Generative Modeling

- Take training samples from the data distribution and learn a mapping from a simple distribution (e.g. normal) to the data distribution such that
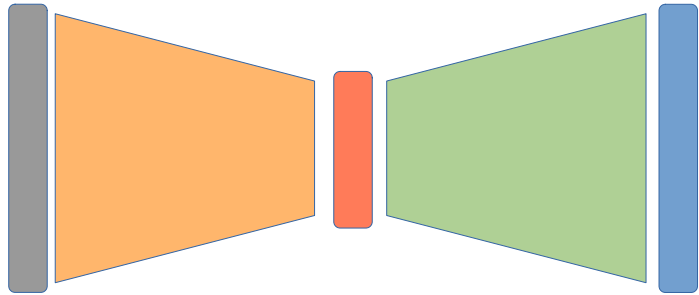


Input samples
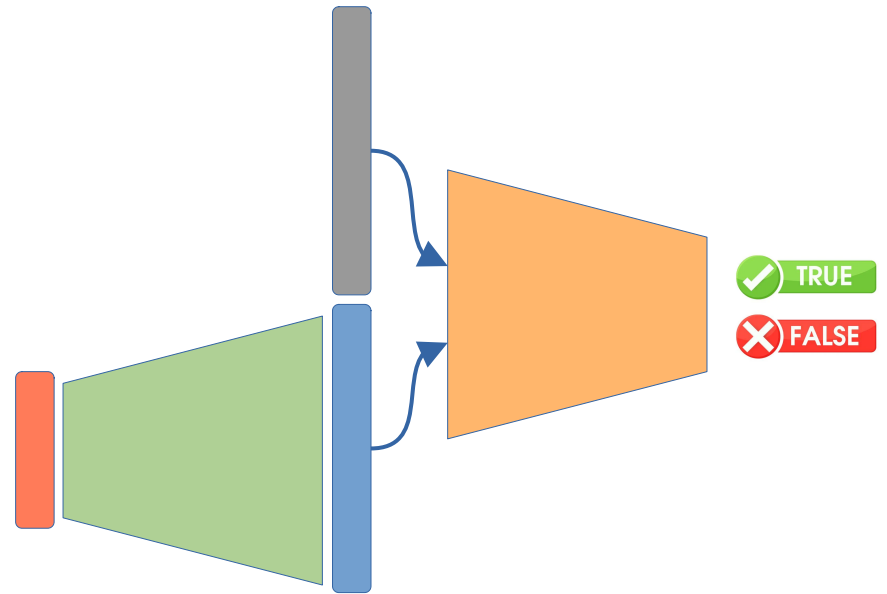$$p_{\mathrm{data}}(x)$$

$$\simeq$$
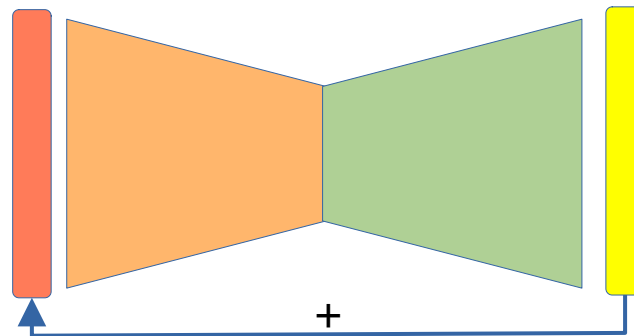
Generated samples
$$p_{\mathrm{model}}(x)$$

# Generative Models

Autoencoders, Variational
Autoencoders (VAEs)

TRUE

FALSE
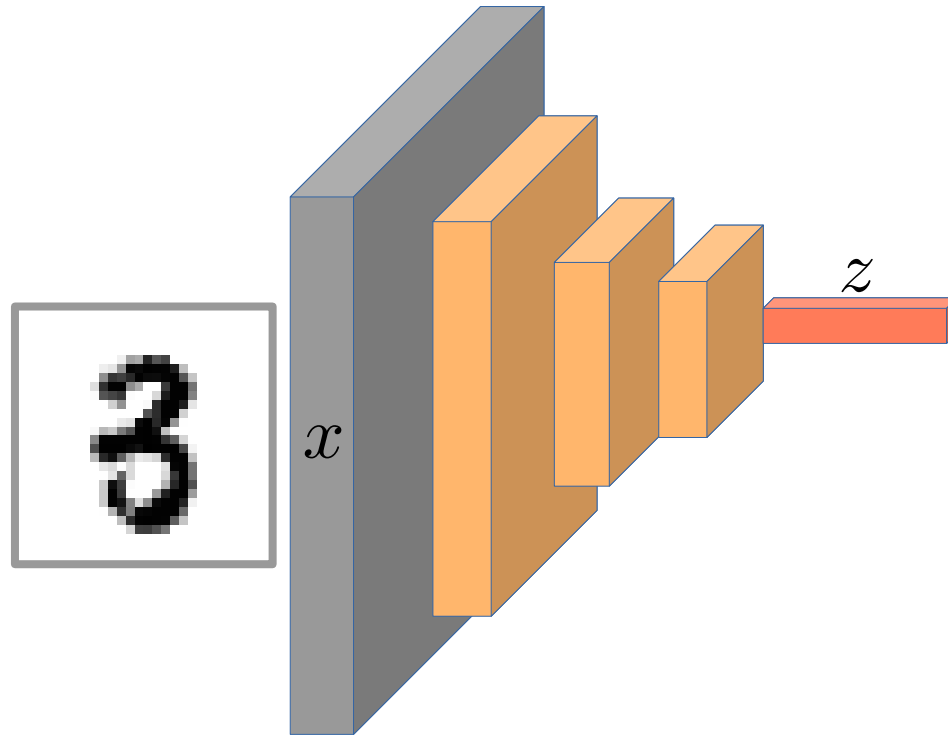
Generative Adversarial
Networks (GANs)

+

Diffusion Models

# Autoencoders: background

- Unsupervised learning of a lower-dimensional feature representation from unlabeled data
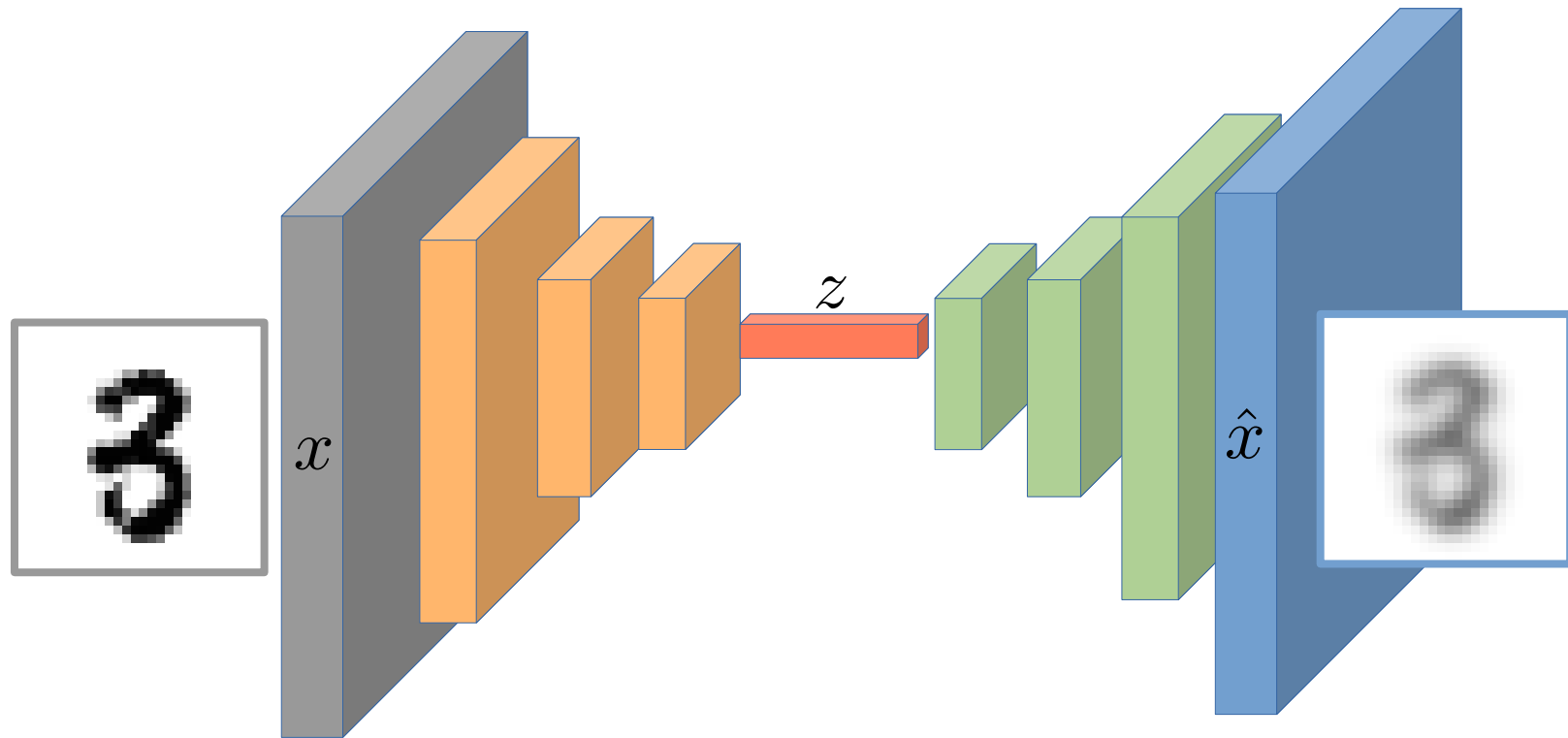


$x$

$z$

Why we want a low-dimensional latent space?

- **Encoder** learns mapping  $f(x, \phi) : x \rightarrow z$

  where $z$ is low-dim. latent space

# Autoencoders: background

- How can we learn the latent space?



- **Decoder** learns mapping $g(z, \theta) : z \to \hat{x}$

  from the latent space *z* to a reconstructed observation $\hat{x}$

# Autoencoders: background

- Train the model to use *z* to reconstruct the original *x*



$$L(x) = \|x - \hat{x}\|^2 = \|x - g(f(x))\|^2$$

- Loss is without labels

$$\min_{\phi,\psi} \sum_i L(x^{(i)})$$

# Dimensionality of latent space

- Autoencoding is a form of compression



| 2D latent space | 5D latent space | GT |

# Obstacles

- AE easily overfits and encoding in the latent space is meaningless.
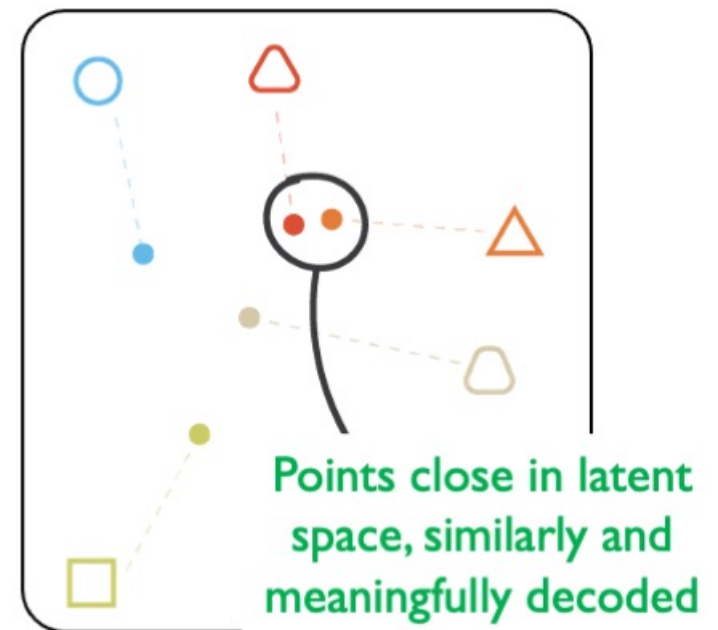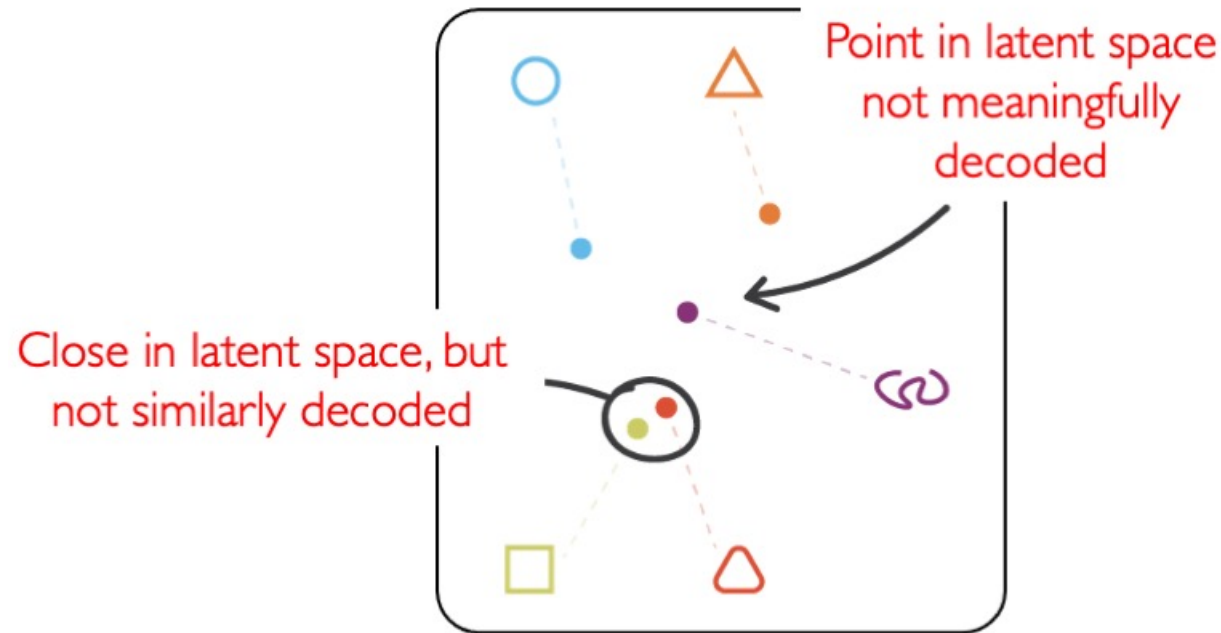
- Encoder and Decoder with sufficient capacity can learn even for the 1D latent space $z$.



Point in latent space not meaningfully decoded

Close in latent space, but not similarly decoded

Points close in latent space, similarly and meaningfully decoded

Regularization on $z$

MIT course, Introduction to deep learning
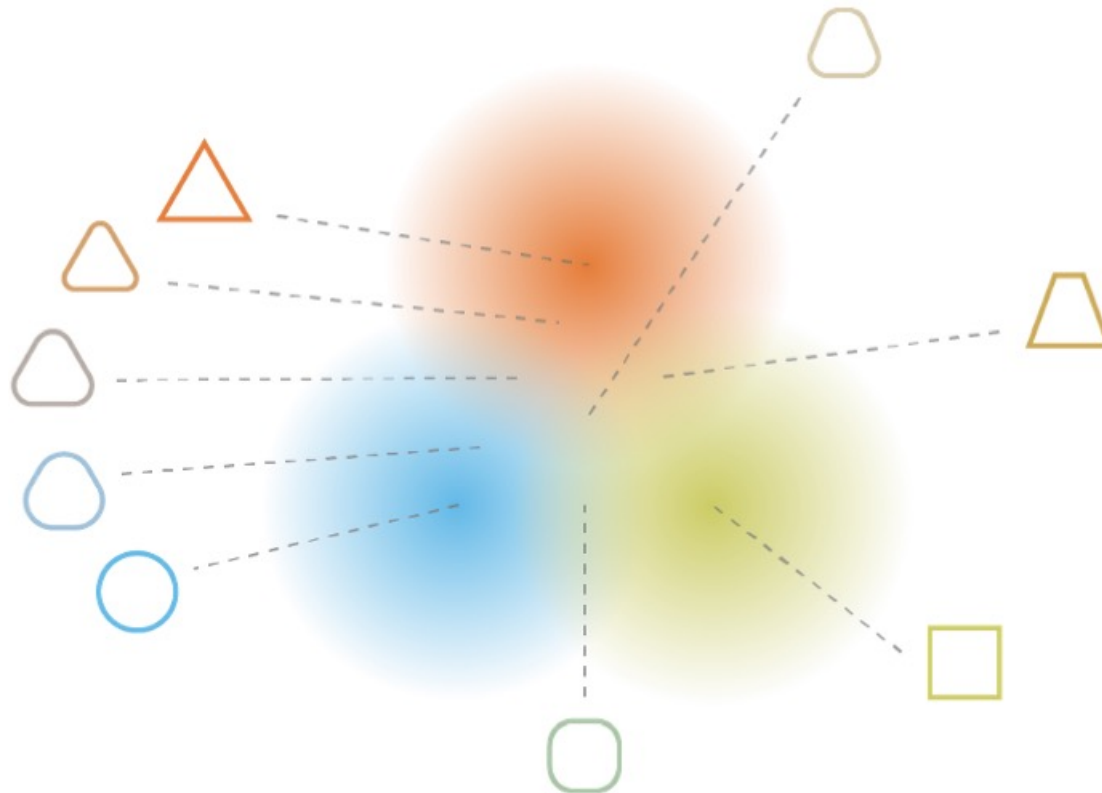
# Regularization

- Continuity:

  points close in latent space - > similar content after decoding

- Completeness:

  sampling from latent space - > "meaningful" content after decoding



MIT course, Introduction to deep learning

# Regularization

- Contractive autoencoder (CAE)

$$L(x) = \|x - g(f(x))\|^2 + \lambda \left\| \frac{\partial f(x)}{\partial x} \right\|_F^2$$

- Denoising autoencoder (DAE)

$$L(x) = \|x - g(f(x + \epsilon))\|^2$$

- Variational autoencoder (VAE)

# Variational Autoencoder

- Traditional autoencoder



$$f(x, \phi) = z \qquad g(z, \theta) = \hat{x}$$

$$x \qquad \phi \qquad z \qquad \theta \qquad \hat{x}$$
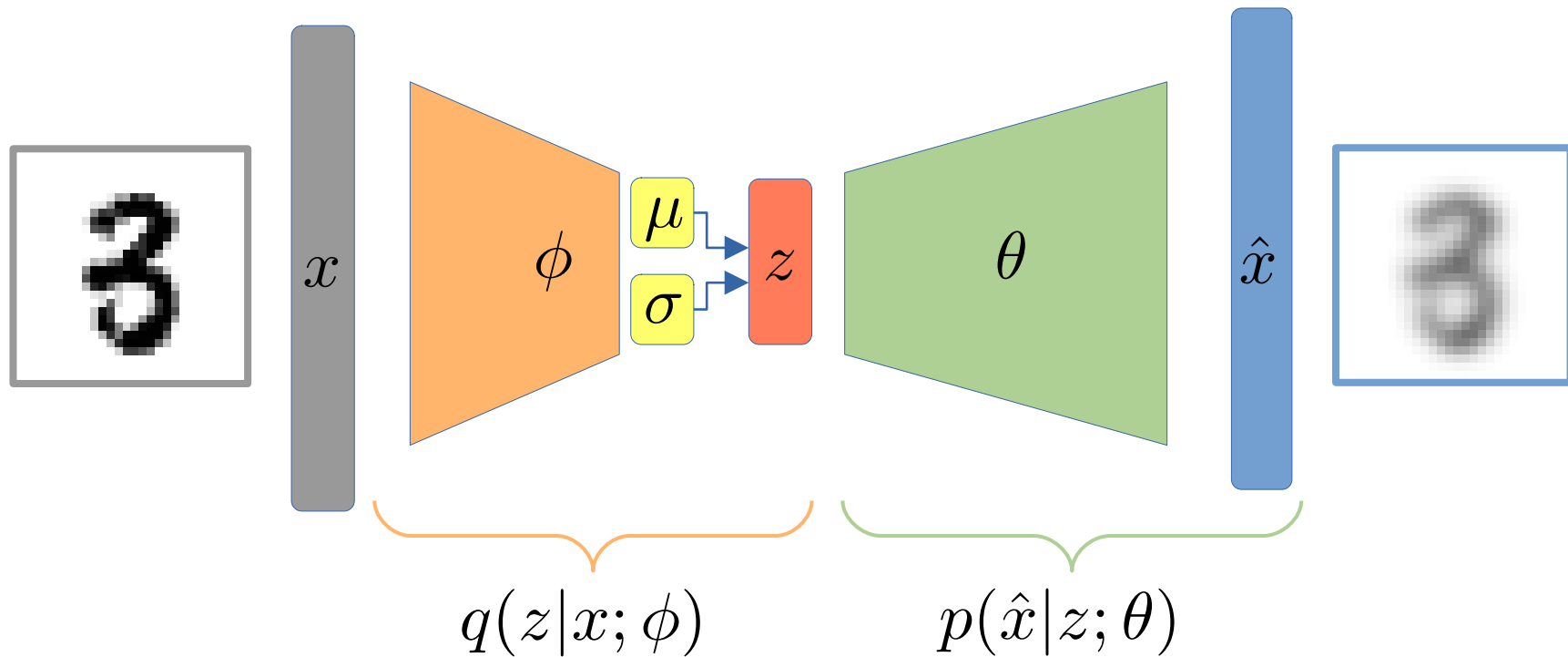
$$q(z|x; \phi) \qquad p(\hat{x}|z; \theta)$$

- Better think of probabilities –>  AE learns the means:

$$f(x) = \mathbb{E}_{q(z|x)}[z] \qquad g(z) = \mathbb{E}_{p(\hat{x}|z)}[\hat{x}]$$

# Variational Autoencoder

- VAE approximates $q(z) \approx N(z|\mu, \sigma^2)$

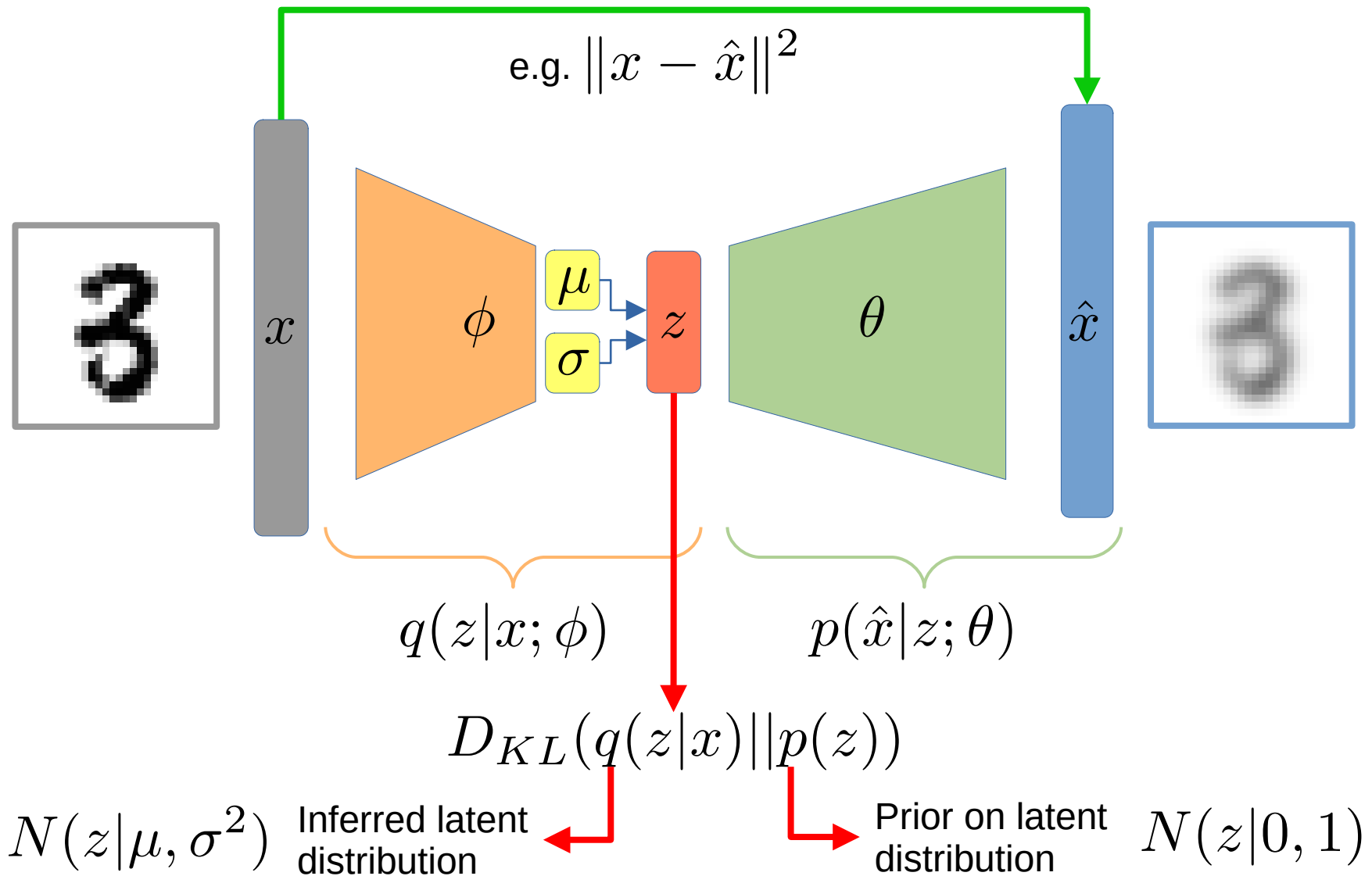  and learns both mean and standard deviation vectors



$$q(z|x; \phi) \qquad p(\hat{x}|z; \theta)$$

- Latent variable *z* is sampled from estimated $q(z|x; \phi)$

# Variational Autoencoder

- Loss: $L(x, \phi, \theta) = \textcolor{green}{\text{reconstruction loss}} + \textcolor{red}{\text{regularization term}}$



e.g. $\|x - \hat{x}\|^2$

$q(z|x; \phi)$

$p(\hat{x}|z; \theta)$

$$D_{KL}(q(z|x)\|p(z))$$

$N(z|\mu, \sigma^2)$ Inferred latent distribution

Prior on latent distribution $N(z|0, 1)$

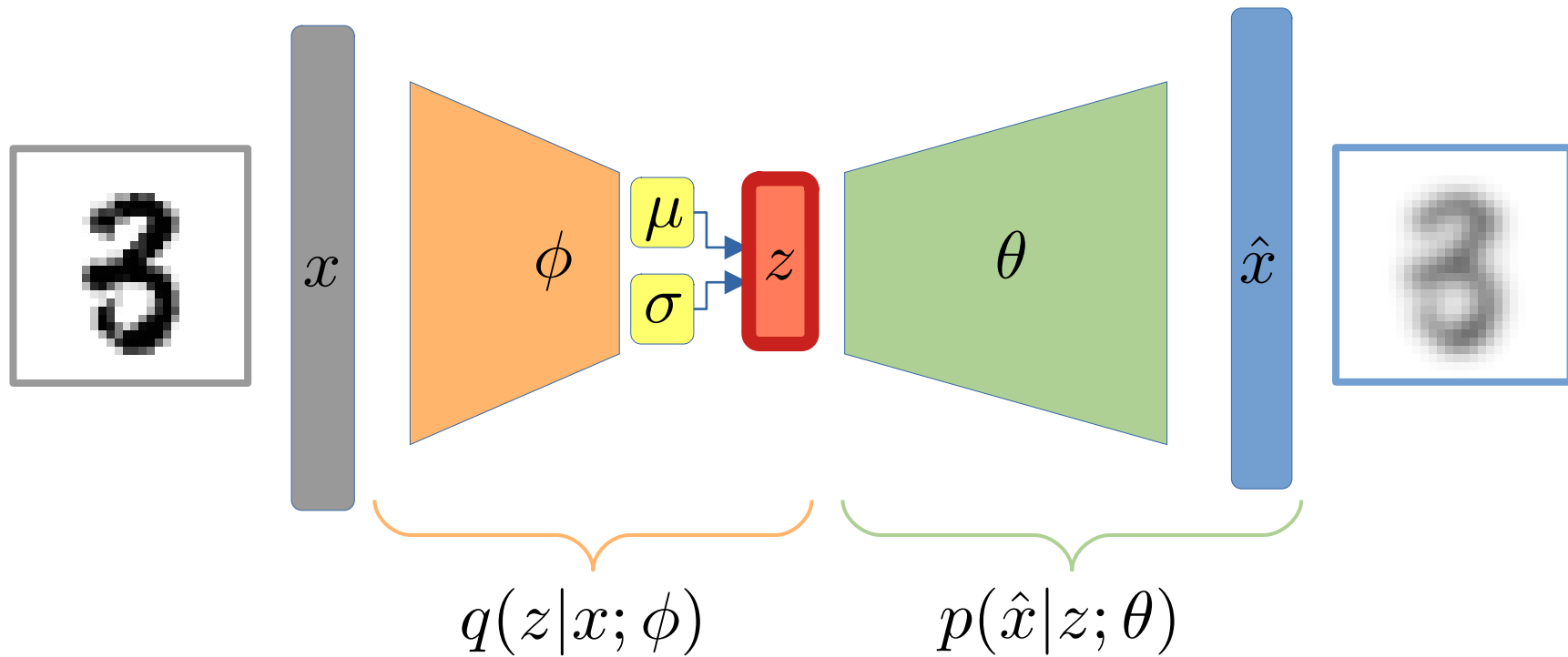# Priors on the latent distribution

$$D_{KL}(q(z|x)||p(z))$$

- q and p are normal distributions --> KL divergence has closed-form expression

$$D_{KL}(N(z|\mu,\sigma^2)||N(z|0,1)) = \frac{1}{2}(\mu^2 + \sigma^2 - 1 - \log\sigma^2)$$
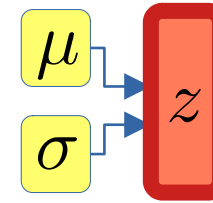
# VAE Computation Graph

- **We cannot use backpropagation if *z* is sampled.**



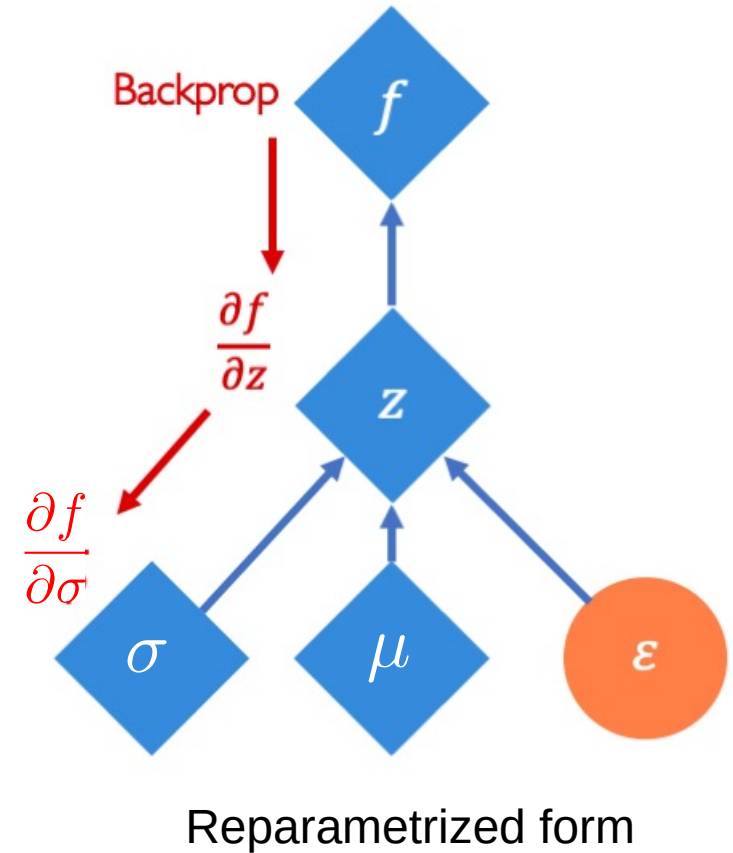$$q(z|x;\phi) \qquad p(\hat{x}|z;\theta)$$

# Reparametrization Trick
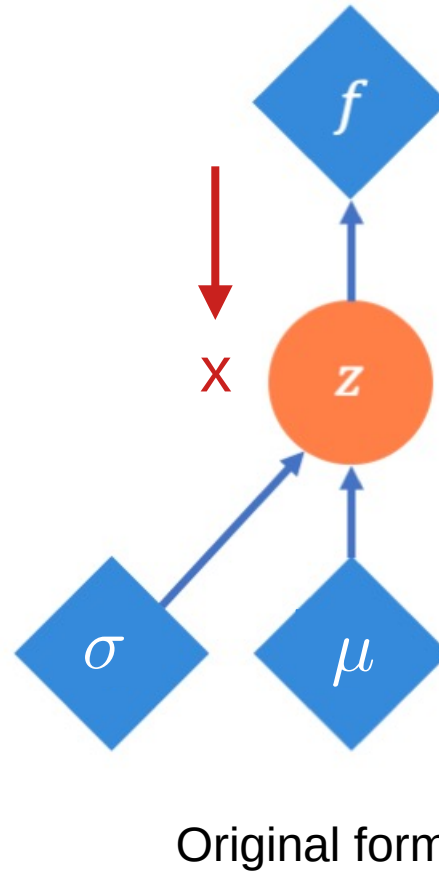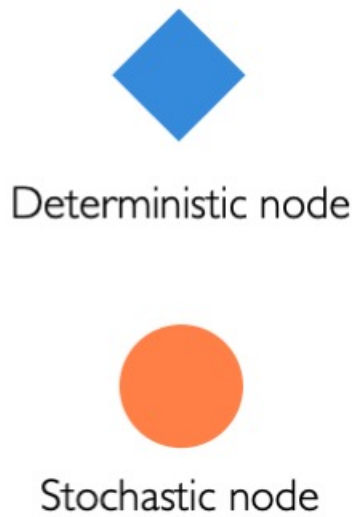
- Sample an auxiliary variable $\epsilon \sim N(0,1)$

- Calculate: $z = \mu + \sigma \epsilon$

- Then $z \sim N(\mu, \sigma)$

# Reparametrization Trick



Deterministic node

Stochastic node

Original form

Reparametrized form

Backprop

$\frac{\partial f}{\partial z}$

$\frac{\partial f}{\partial \sigma}$

# VAE Latent Perturbation



Smile

Head Pose
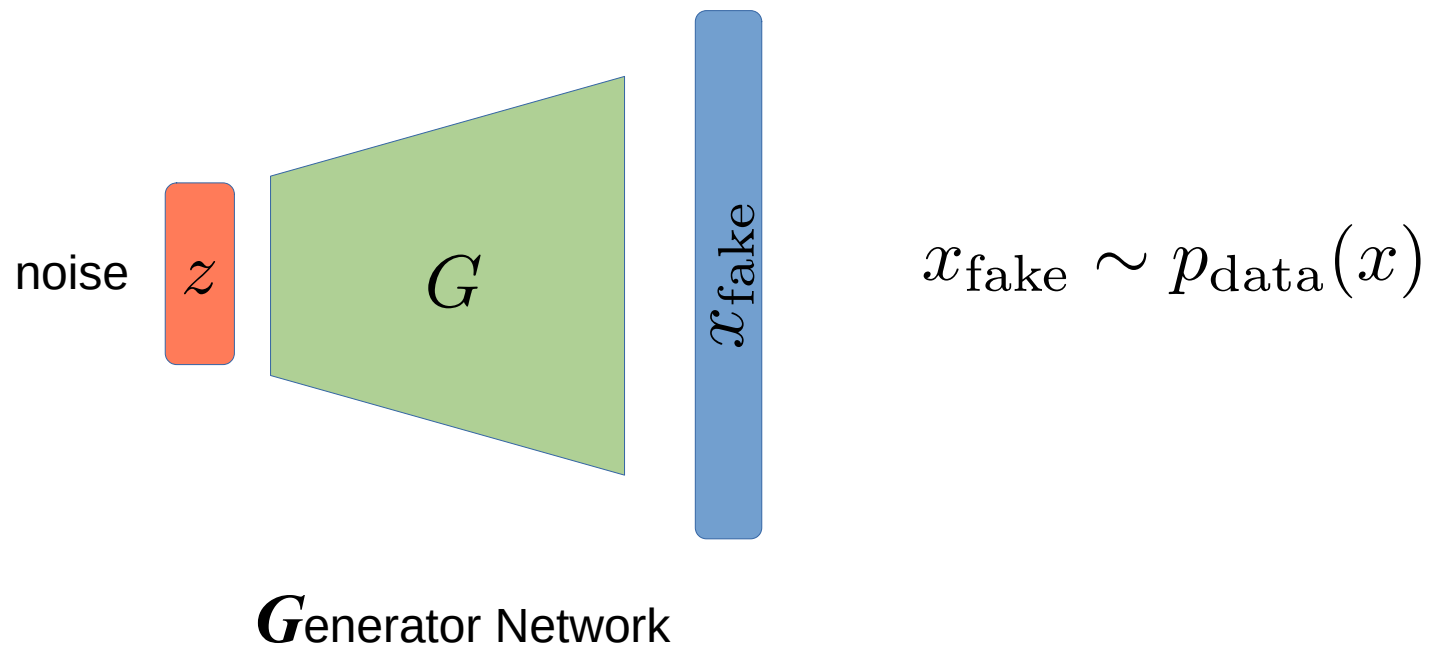
Kingma et al., VAE, 2014

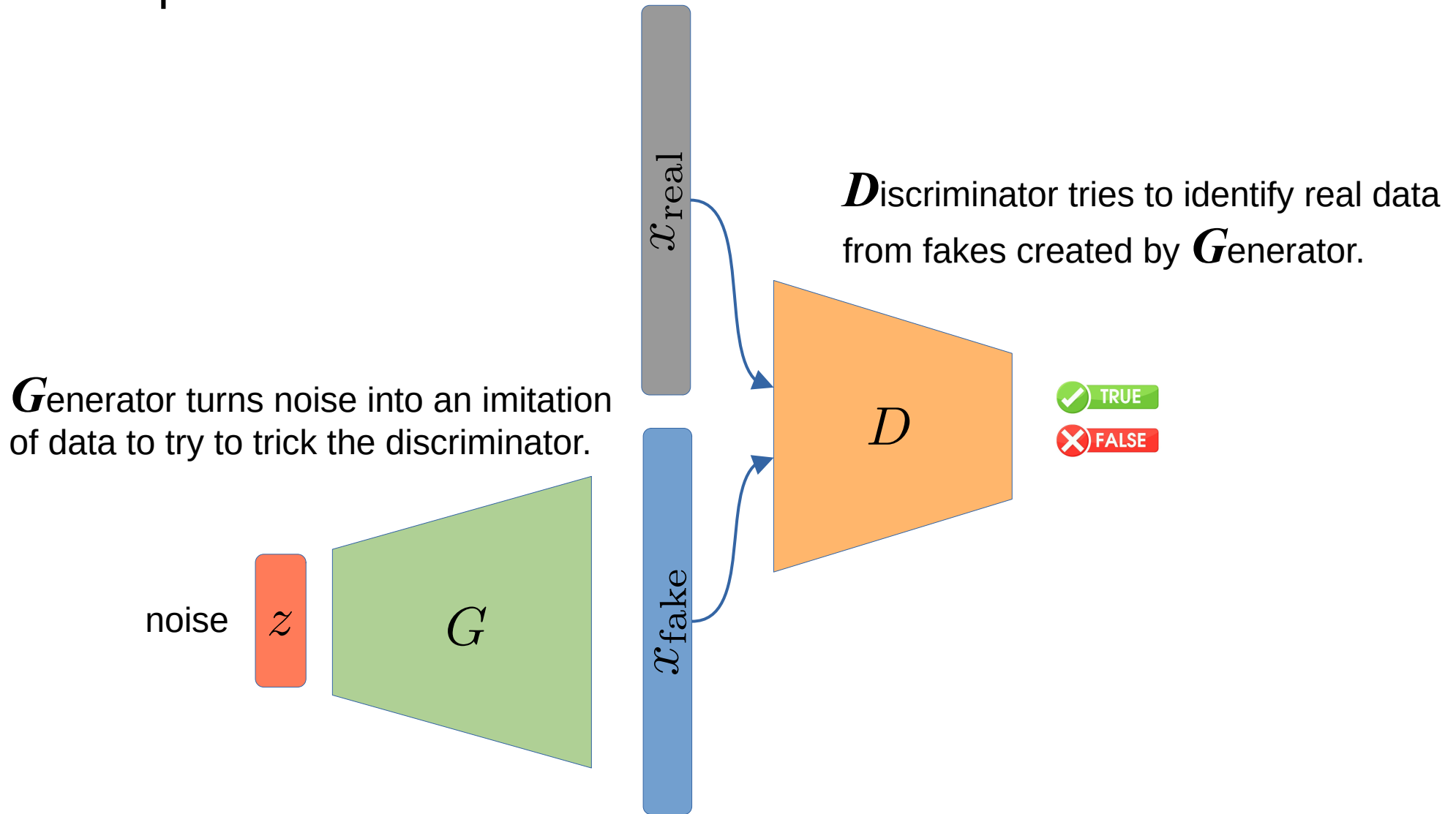# Generative Adversarial Networks - GANs

- What if we just want to sample?

- Idea: Don't explicitly model density, just sample

- Solution: Sample from something simple (white noise) and learn a transformation to the data distribution.

noise

$z$

$G$

$x_{\text{fake}}$

$x_{\text{fake}} \sim p_{\text{data}}(x)$

**G**enerator Network

# Generative Adversarial Networks - GANs

- GAN: create a generative model by having two networks compete with each other.

$x_{\text{real}}$

$\boldsymbol{D}$iscriminator tries to identify real data from fakes created by $\boldsymbol{G}$enerator.

$\boldsymbol{G}$enerator turns noise into an imitation of data to try to trick the discriminator.

noise   $z$   $G$

$x_{\text{fake}}$

$D$

✓ TRUE
✗ FALSE

# Intuition behind GANs



Generator

Real data    Fake data

# Training GANs

- Adversarial objectives for *D* and *G*

$$\min_{G} \max_{D} \mathbb{E}_{z,x}[\log(1 - D(G(z))) + \log D(x)]$$

$\boldsymbol{D}$ tries to identify the synthesized instances.

$\boldsymbol{G}$ tries to synthesize fake instances that fool $\boldsymbol{D}$.

$x_{\text{real}}$

$x_{\text{fake}}$

$z$

noise

$G$

$D$

✔ TRUE
✖ FALSE

# Training GANs

- Loss:

$$\max_D \mathbb{E}_{z,x} \left[ \log(1 - D(G(z))) + \log D(x) \right]$$

fake          real



$D$ tries to identify the synthesized instances.

$D$(fake) = 0
$D$(real) = 1

$x_{\text{real}}$

$x_{\text{fake}}$

$D$

TRUE

FALSE

- Loss:

$$\min_{G} \mathbb{E}_{z,x}[\log(1 - D(G(z))) + \log D(x)]$$
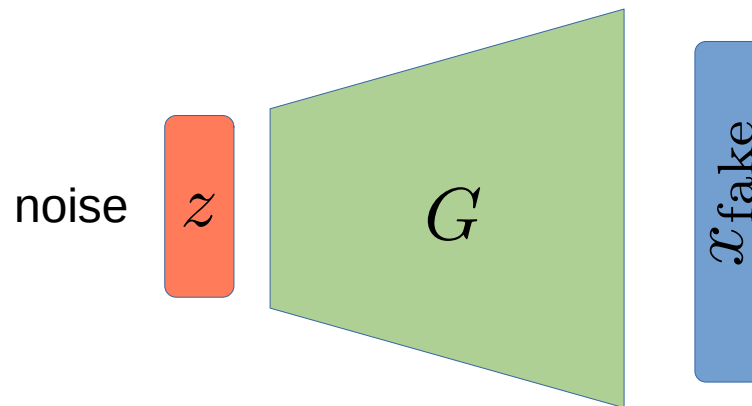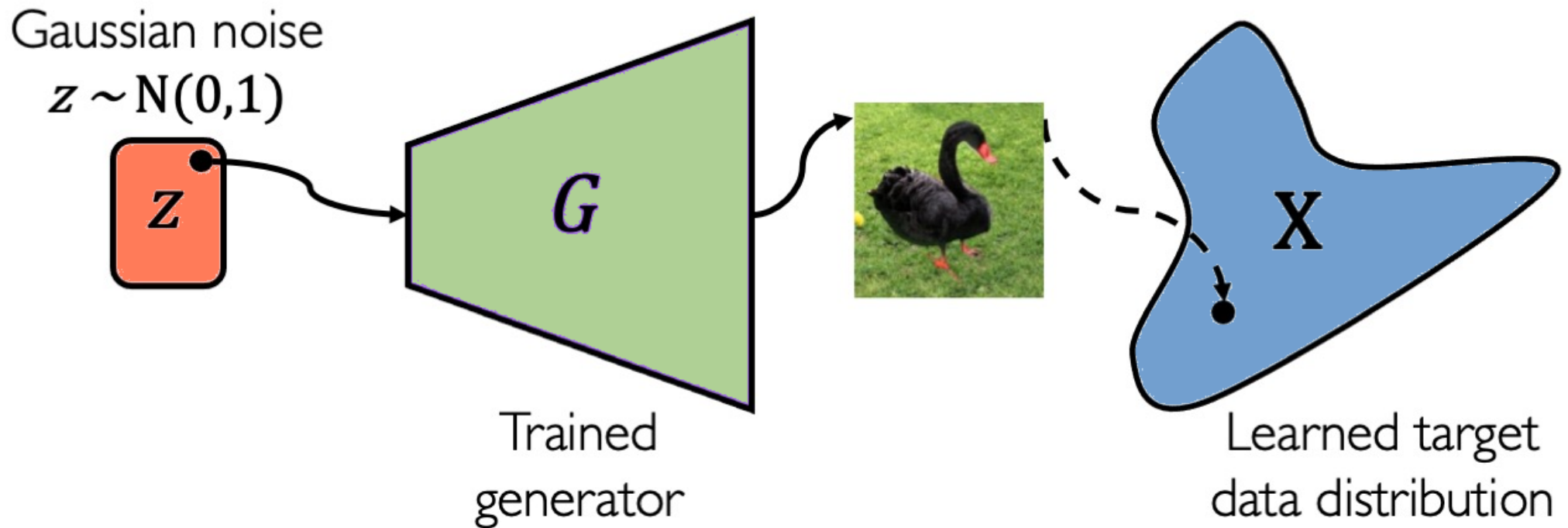
$G$ tries to synthesize fake

instances that fool $D$.

noise $z$ | $G$ | $x_{\text{fake}}$

# GANs are distribution transformers



Gaussian noise
$z \sim N(0,1)$

$z$

$G$

Trained generator

$X$

Learned target data distribution

# GANs are distribution transformers



Gaussian noise
$z \sim N(0,1)$

$z$

$G$

Trained generator

$X$

Learned target data distribution

# GANs are distribution transformers



Gaussian noise
$z \sim N(0,1)$

$z$

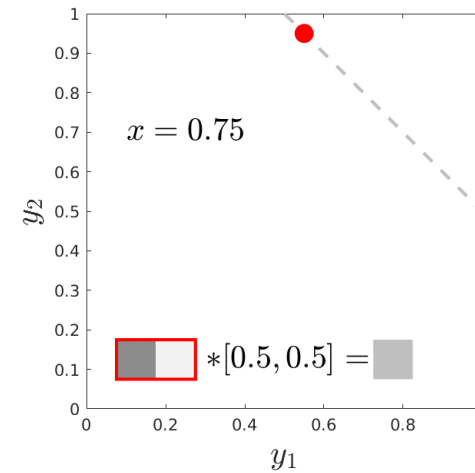$G$

Trained generator

$?$

$X$

Learned target data distribution

# GANs vs AEs



$$\frac{1}{2}(y_1 + y_2) = x$$



MSE $- (h * y - x)^2$

MAE $- |h * y - x|$

MAP $-$ Adversarial loss

Sonderby et al., ICLR, 2017

# Diffusion Models

- Diffusion



- Forward process

$$q(x_t|x_{t-1}) \equiv N(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I)$$



$$\mathbf{x}_0 \qquad \mathbf{x}_1 \qquad \mathbf{x}_{t-1} \qquad \mathbf{x}_t \qquad \mathbf{x}_T$$

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

$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I) \qquad \bar{\alpha}_t = \prod_{i=1}^{t}(1-\beta_i)$$

# Diffusion Model

- Forward process

$$\equiv N(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

$$q(x_1 | x_0)$$

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



$$\mathbf{x}_0 \qquad \mathbf{x}_1 \qquad \mathbf{x}_{t-1} \qquad \mathbf{x}_t \qquad \mathbf{x}_T$$

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

$$p(x_0 | x_1) \qquad p(x_{t-1} | x_t)$$

$$= N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

- Reverse process

  – Training (ML): $\min \mathbb{E}_{q(x_0)} [- \log p(x_0)]$

$$p(x_0) = \int p(x_0, \ldots, x_T) dx_1 \ldots dx_T = \int p(x_T) \prod_{i=1}^{T} p(x_{i-1} | x_i) dx_1 \ldots dx_T$$
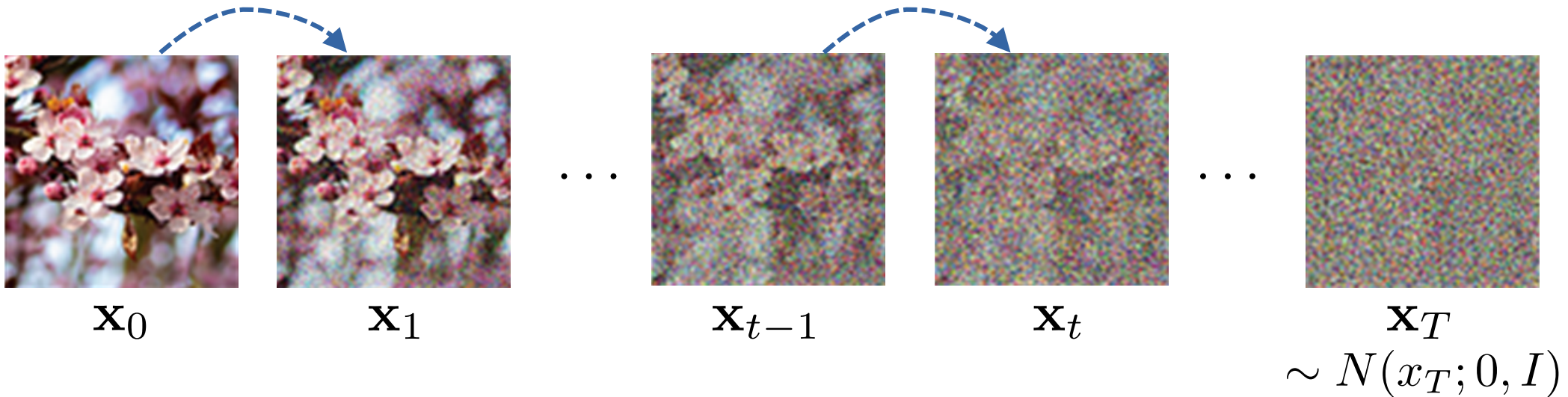
- Forward

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$
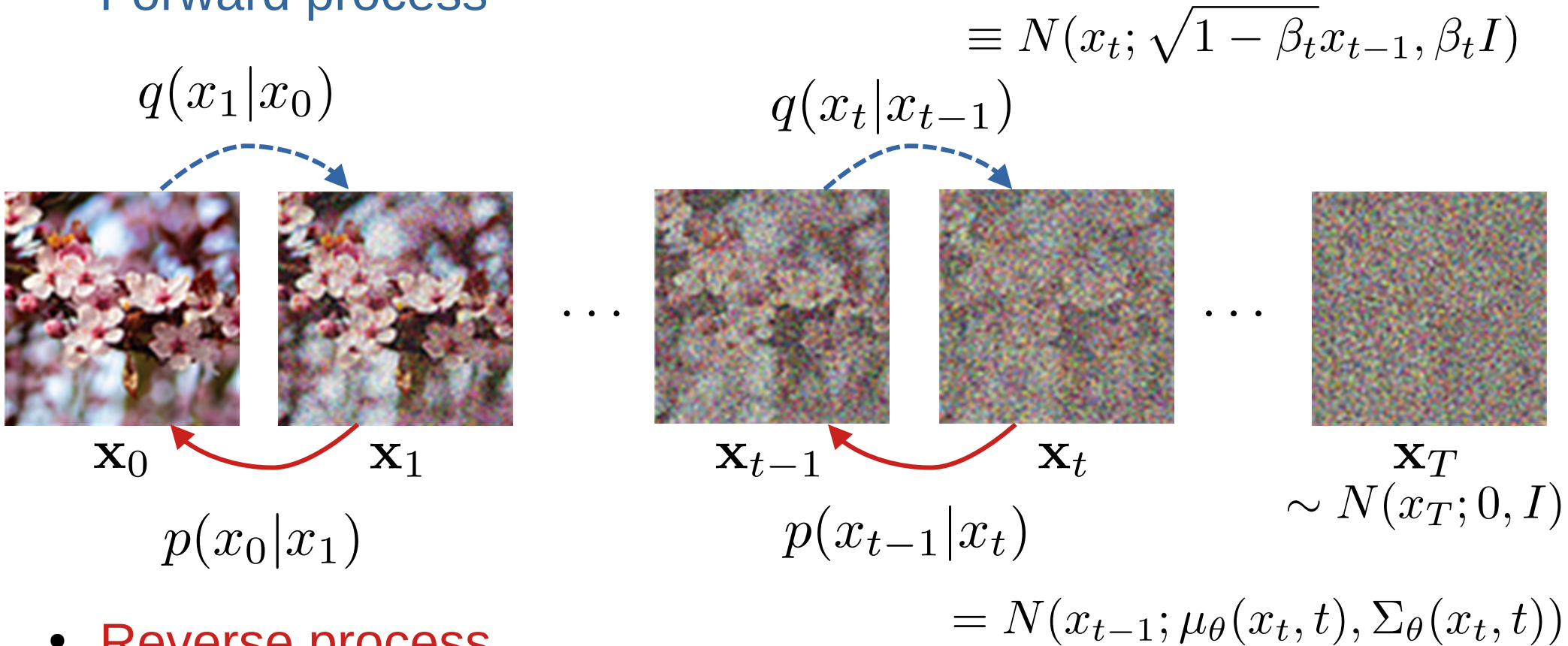
$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I) \qquad \bar{\alpha}_t = \prod_{i=1}^{t}(1 - \beta_i)$$

reparametrization: $\quad x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon \qquad \epsilon \sim N(0, I)$

- Reverse

$$p(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

- Training
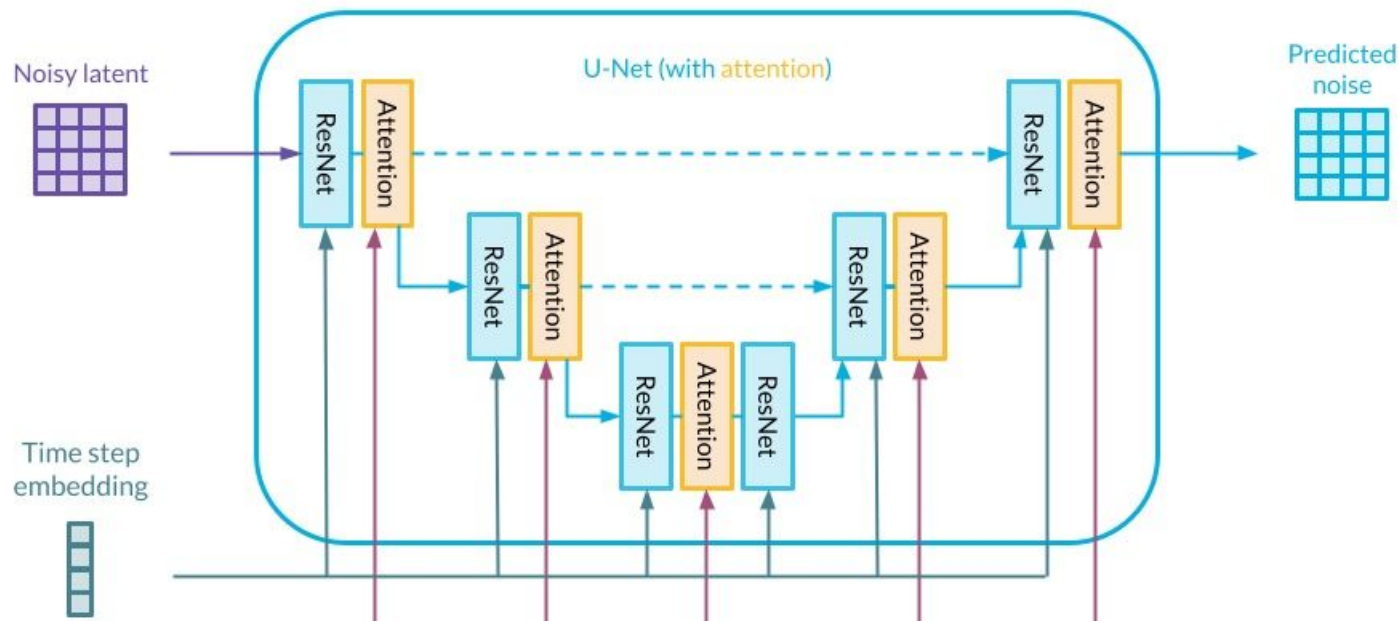
$$\min_\theta \mathbb{E}_{q(x_0)}[-\log p(x_0)] = \dots$$

$$\min_\theta \mathbb{E}_{x_0, \epsilon, t}\left[\frac{1}{2\|\Sigma_\theta\|^2}\left\|\frac{1}{\sqrt{1 - \beta_t}}\left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon\right) - \mu_\theta(x_t, t)\right\|^2\right]$$

or predict noise $\quad \min_\theta \mathbb{E}_{x_0, \epsilon, t}\left[\frac{1}{2\|\Sigma_\theta\|^2}\|\epsilon - \epsilon_\theta(x_t, t)\|^2\right]$

- Prediction of $\mu_\theta(x_t, t)$ or $\epsilon_\theta(x_t, t)$ is done by

  U-Net($\theta$) with residual and attention blocks

  and $t$ implemented as sinusoid positional encoding.

**Algorithm 1** Training

1: **repeat**
2:   $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:   $t \sim \text{Uniform}(\{1, \ldots, T\})$
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

**Algorithm 2** Sampling

1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
2: **for** $t = T, \ldots, 1$ **do**
3:   $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
4:   $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
5: **end for**
6: **return** $\mathbf{x}_0$
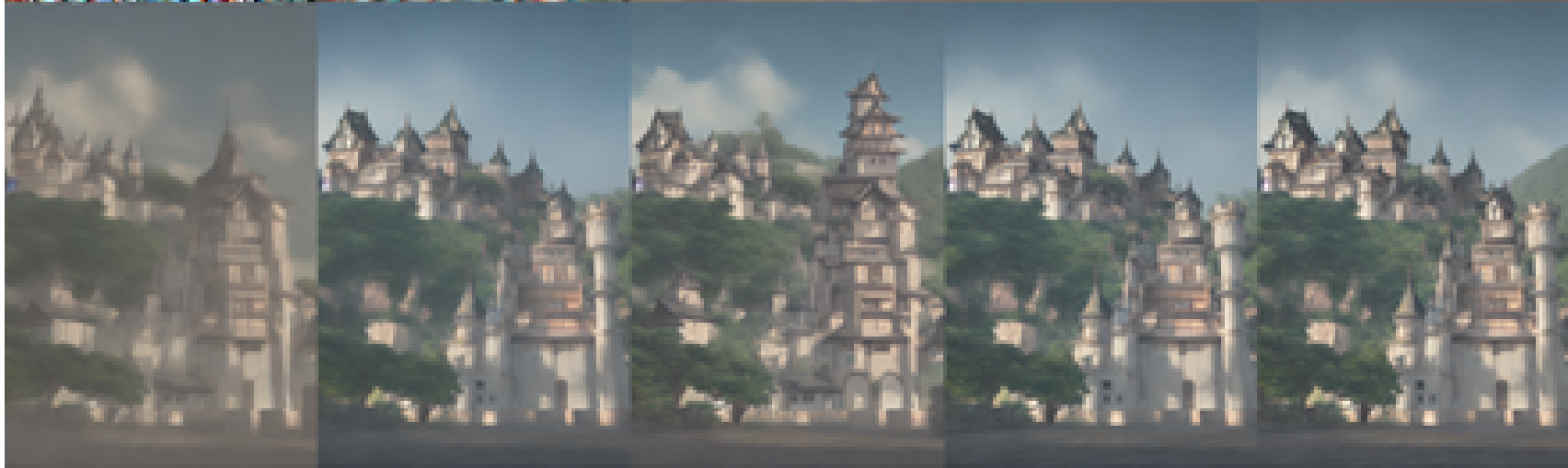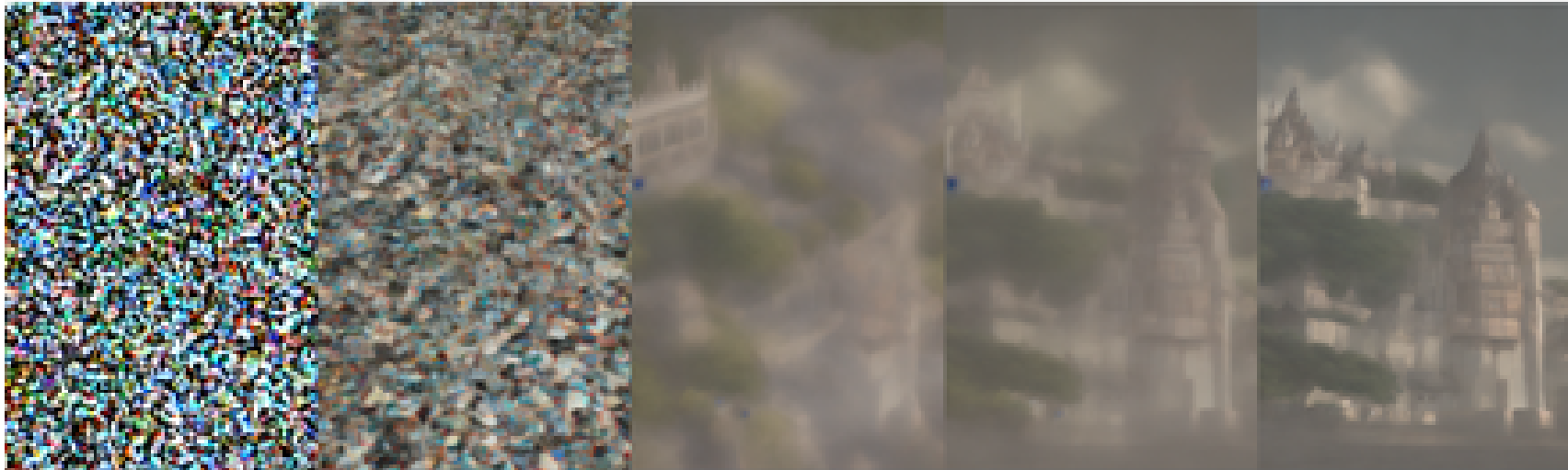
Steps: 1      Steps: 2      Steps: 3      Steps: 5      Steps: 8
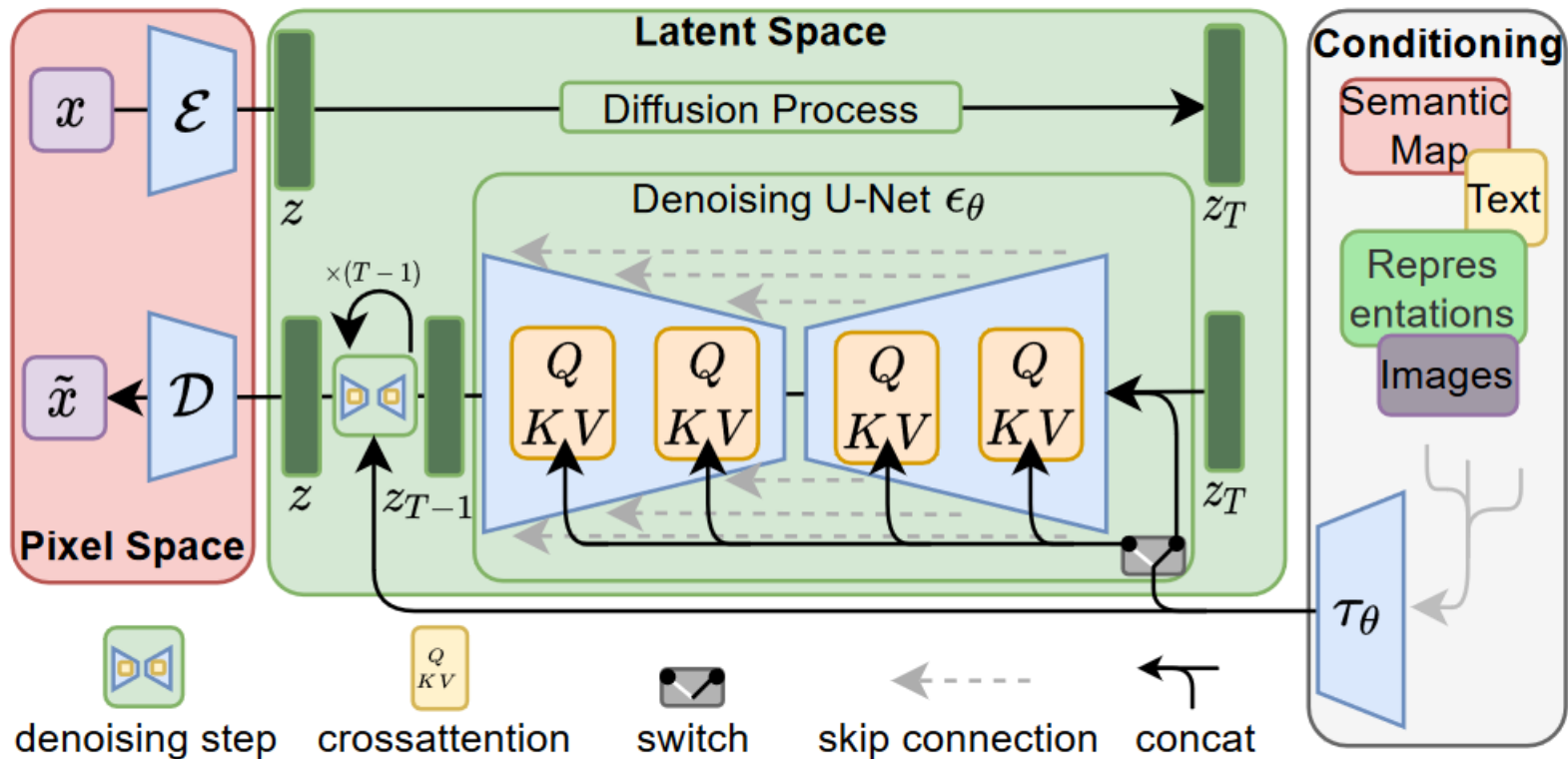
Steps: 10      Steps: 15      Steps: 20      Steps: 30      Steps: 40

# Stable Diffusion Model

- Latent space

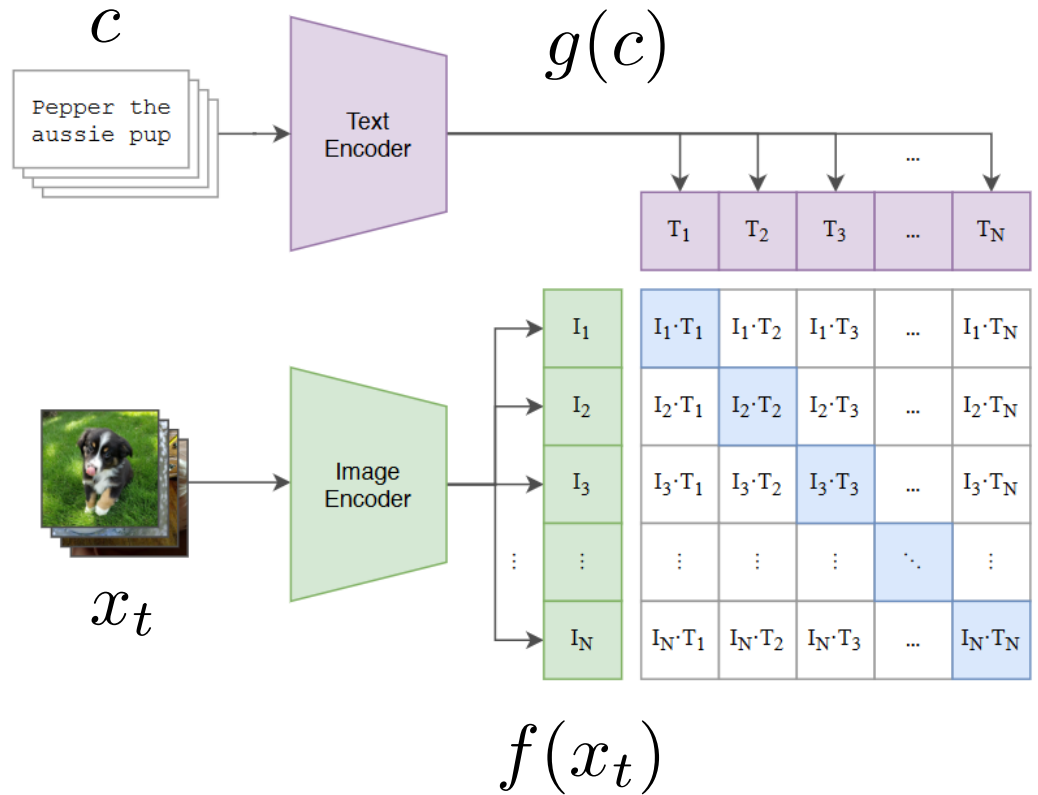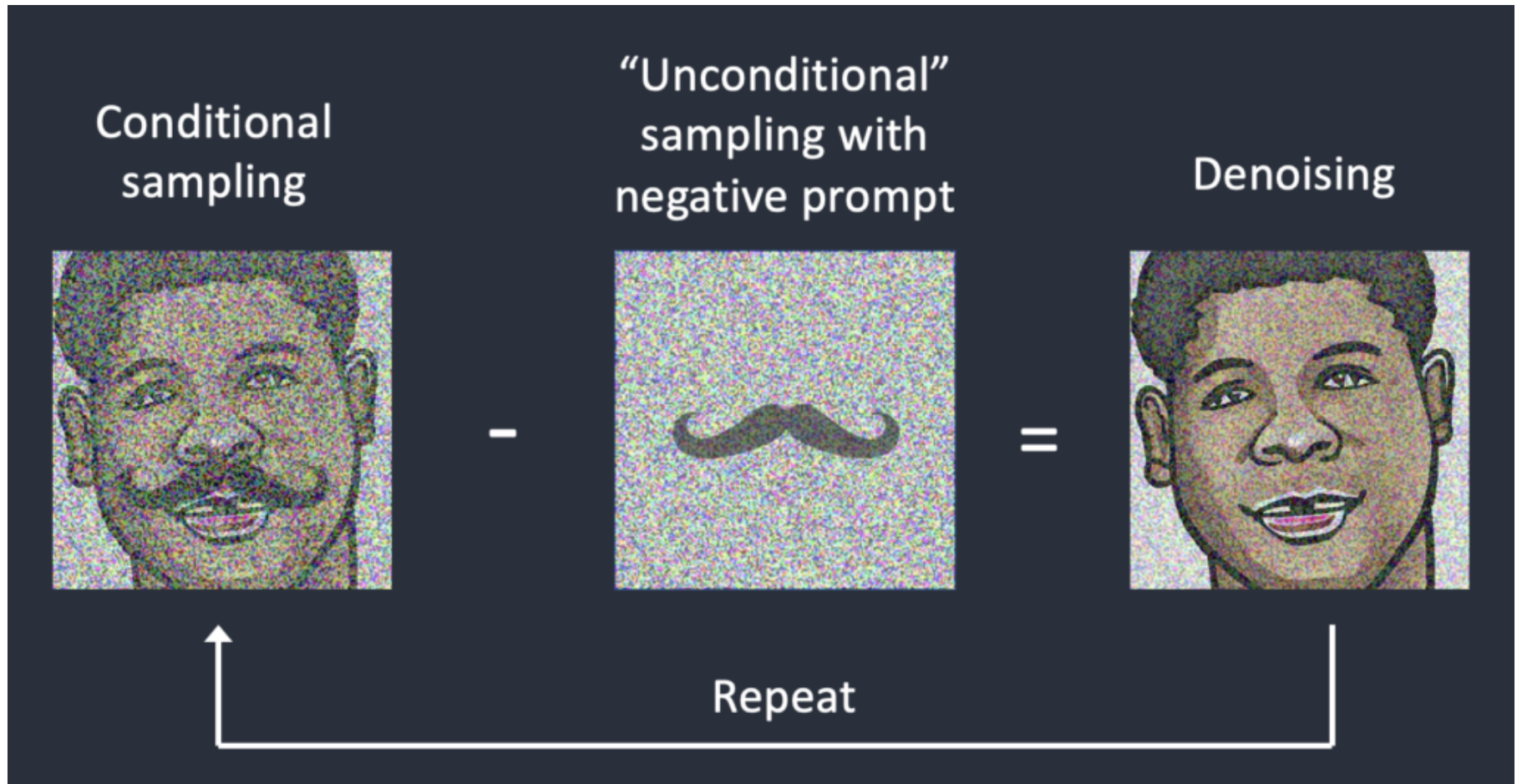- Conditional diffusion (cross-attention)



Rombach et al. LDM, 2022

# DALL-E

gradient scale

$$\hat{\mu}_\theta(x_t, t) = \mu_\theta(x_t, t) + s\Sigma_\theta(x_t, t)\nabla_{x_t}\underbrace{\langle f(x_t), g(c)\rangle}_{\approx \log p(c|x_t)}$$



$c$

Pepper the aussie pup

Text Encoder

$g(c)$

| $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |

$x_t$

Image Encoder

$f(x_t)$

| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ |
|---|---|---|---|---|---|
| $I_1$ | $I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$ |
| $I_2$ | $I_2 \cdot T_1$ | $I_2 \cdot T_2$ | $I_2 \cdot T_3$ | ... | $I_2 \cdot T_N$ |
| $I_3$ | $I_3 \cdot T_1$ | $I_3 \cdot T_2$ | $I_3 \cdot T_3$ | ... | $I_3 \cdot T_N$ |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| $I_N$ | $I_N \cdot T_1$ | $I_N \cdot T_2$ | $I_N \cdot T_3$ | ... | $I_N \cdot T_N$ |

Nichol et al. GLIDE, 2022
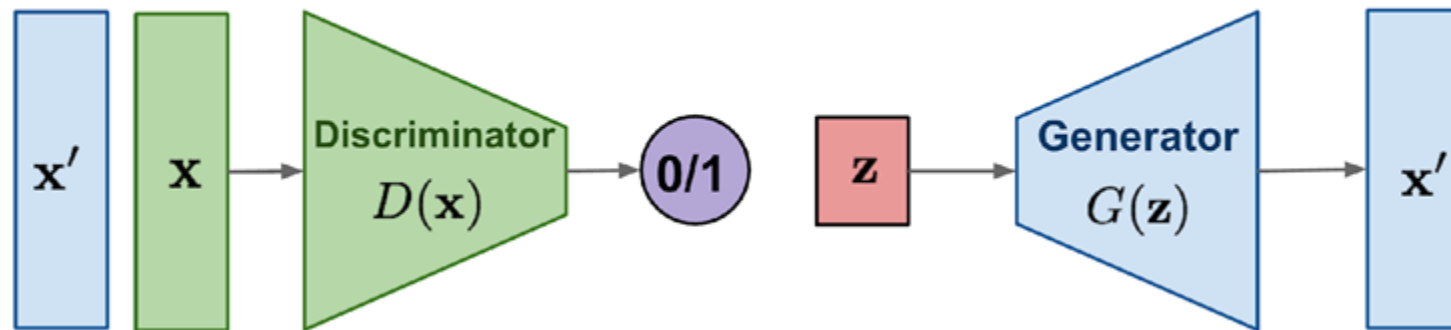
# Negative Prompt

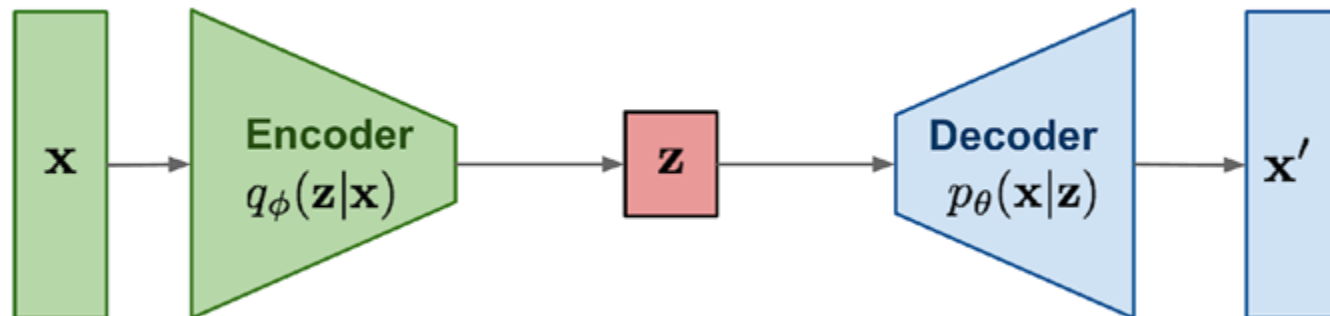detailed meadow with colorful flowers



detailed meadow with colorful flowers
–no blue

**GAN:** Adversarial training

**VAE:** maximize variational lower bound

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