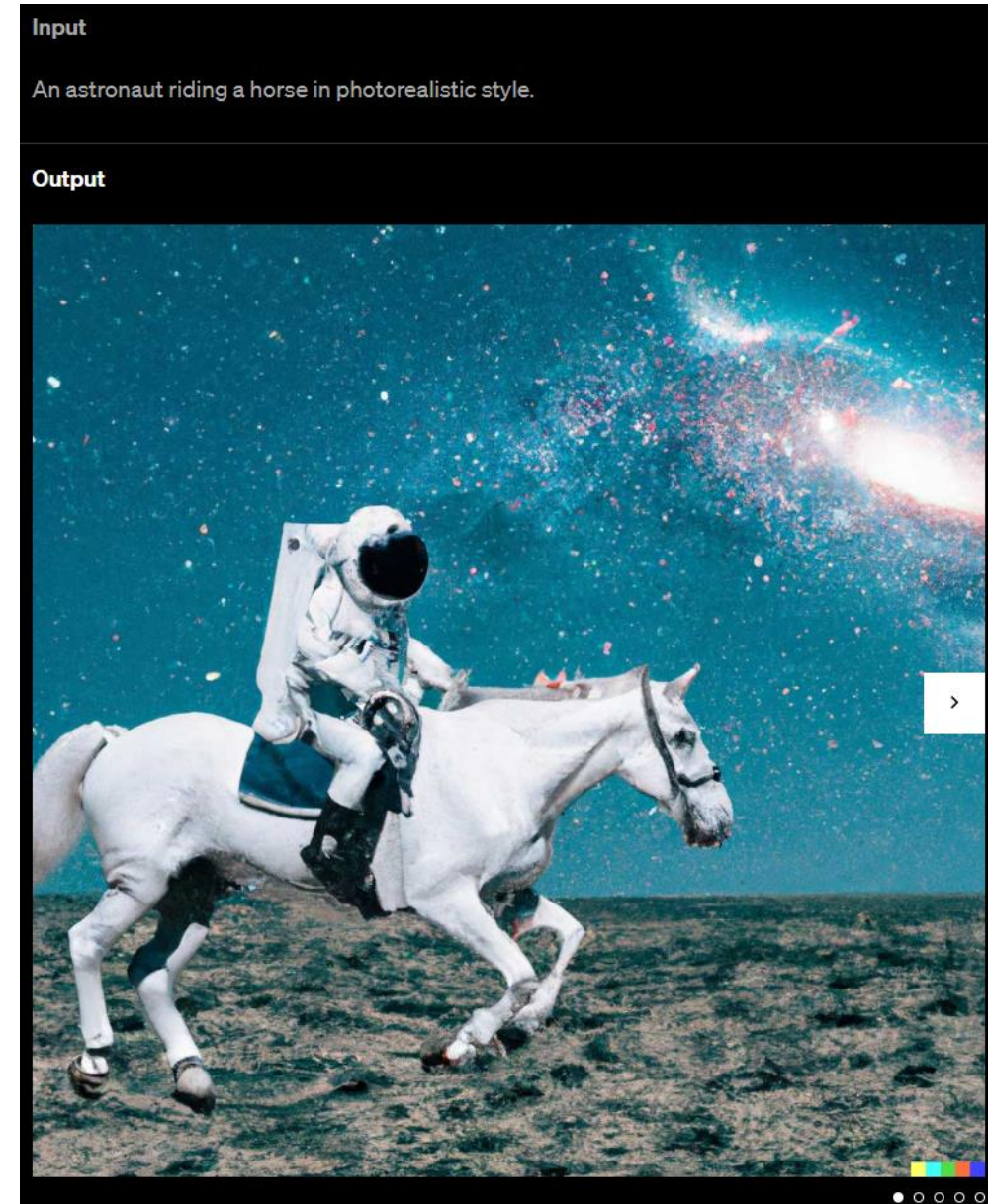# Image Synthesis

Computer Vision

idea: generate new images as variations of training data (same distribution)
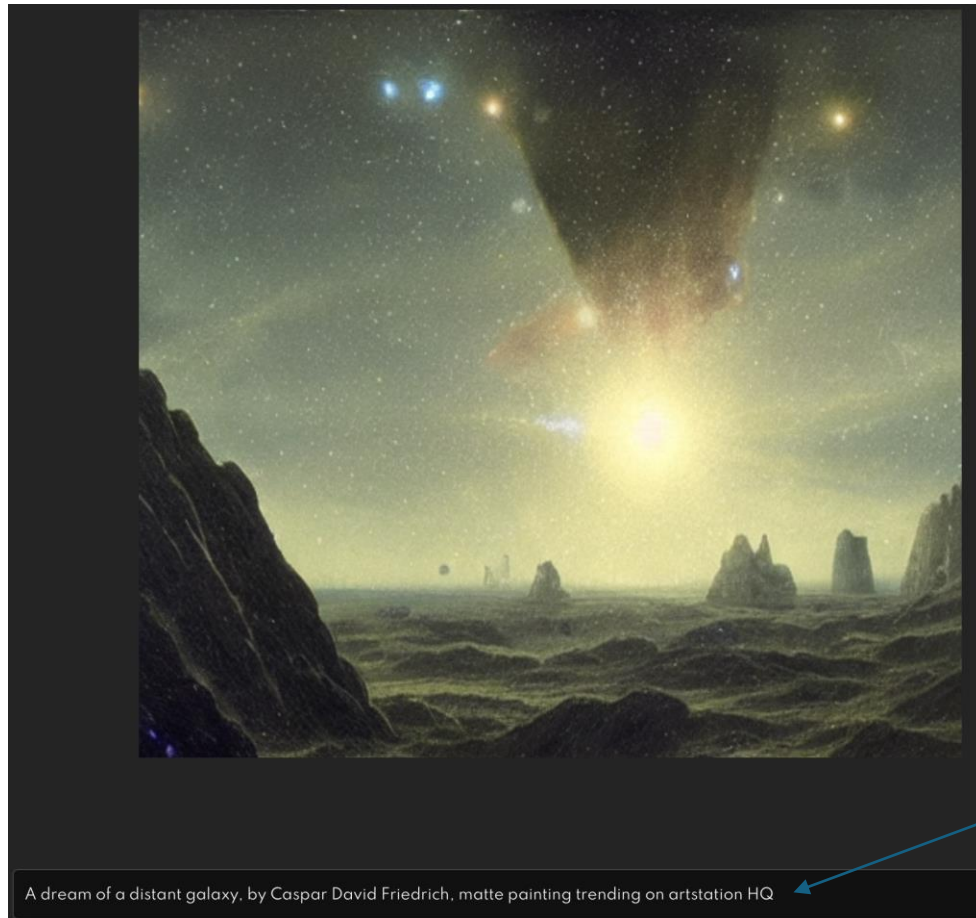
usually conditioned on text (prompt)

compared to text generation, additional mechanism needed (e.g., diffusion) due to more complex image structures
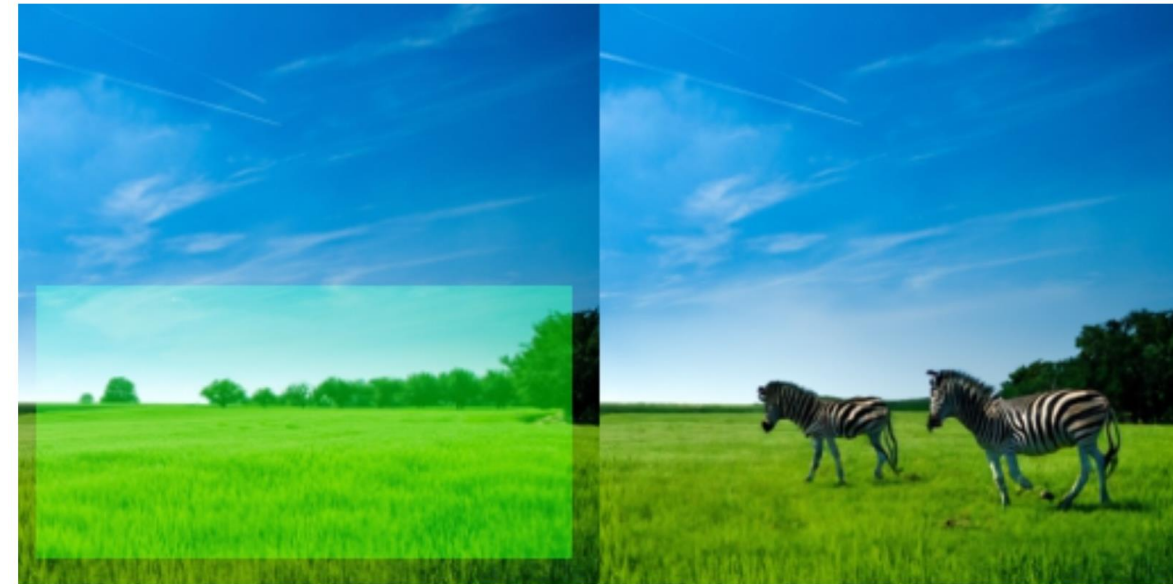


Input

An astronaut riding a horse in photorealistic style.

Output

plenty of products: DALL-E, Stable Diffusion, ImageGen, Midjourney, …

web app for Stable Diffusion: DreamStudio



A dream of a distant galaxy, by Caspar David Friedrich, matte painting trending on artstation HQ

inpainting example (GLIDE ):



prompt ⟶ "zebras roaming in the field"

source

3

# Generative vs Predictive/Discriminative Models

discriminative models:

predict conditional probability $P(Y|\boldsymbol{X})$
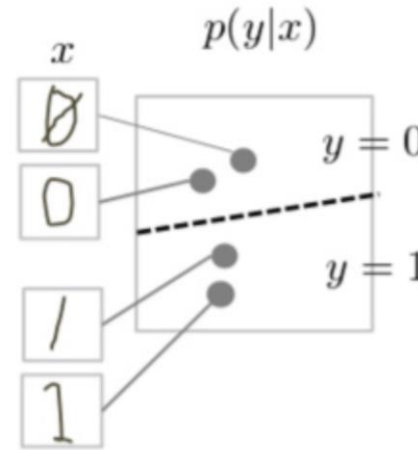
generative models:

predict joint probability $P(Y, \boldsymbol{X})$

(or just $P(\boldsymbol{X})$ → unsupervised learning)
→ allow to generate new data samples

discriminative model          generative model



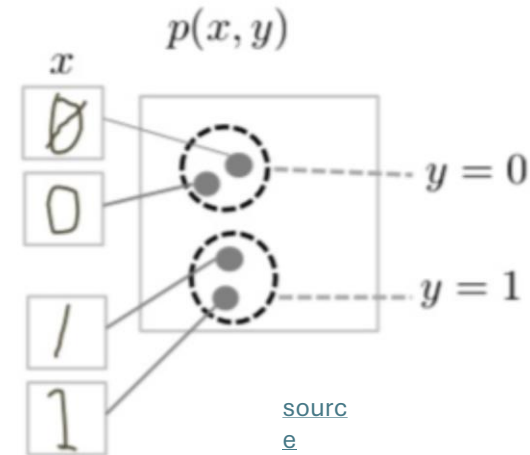task of generative models more difficult: need to model full data distribution rather than merely find patterns in inputs to distinguish outputs

Generative models can be used for predictive tasks (Bayes theorem).
But predictive models are usually better at it.

# Classic ML vs Deep Learning                    Generative vs Predictive Models

## text generation



ChatGPT

## image synthesis



Prompt: Epic anime artwork of a wizard atop a mountain at night casting a cosmic spell into the dark sky that says "Stable Diffusion 3" made out of colorful energy

Stable Diffusion 3 — Stability AI

## text-to-video



Prompt: A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She...

Sora | OpenAI

## BERT family

## tabular data



## computer vision



YOLO

# Deep Learning for Generative AI

Depending on the application, there are currently two dominant approaches for generative AI:

- text generation: decoder LLMs

- image synthesis: diffusion models

note the difference between image synthesis and multimodal understanding in LLMs
(images as additional input sequences to transformer, tokenized by splitting into patches)

# Different Model Types for Image Synthesis



**GAN:** Adversarial training

two neural networks playing a zero-sum game

**VAE:** maximize variational lower bound

learn variational distribution (not just replicating inputs)

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

more complex distributions by applying change-of-variable technique (need for specialized architecture)

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

chain of denoising autoencoders

→ generalization: flow matching

# Generative Modeling



Training data ~ $p_{data}(x)$    learning →    $p_{model}(x)$    sampling →

Objectives:
1. Learn $p_{model}(x)$ that approximates $p_{data}(x)$
2. Sampling new x from $p_{model}(x)$

Explicit density estimation: explicitly define and solve for $p_{model}(x)$
Implicit density estimation: learn model that can sample from $p_{model}(x)$ without explicitly defining it.

# Generative Adversarial Networks (GAN)

two neural networks playing a zero-sum game:

- the generator network G generating new (fake) samples

- the discriminator network D trying to distinguish between real and fake samples
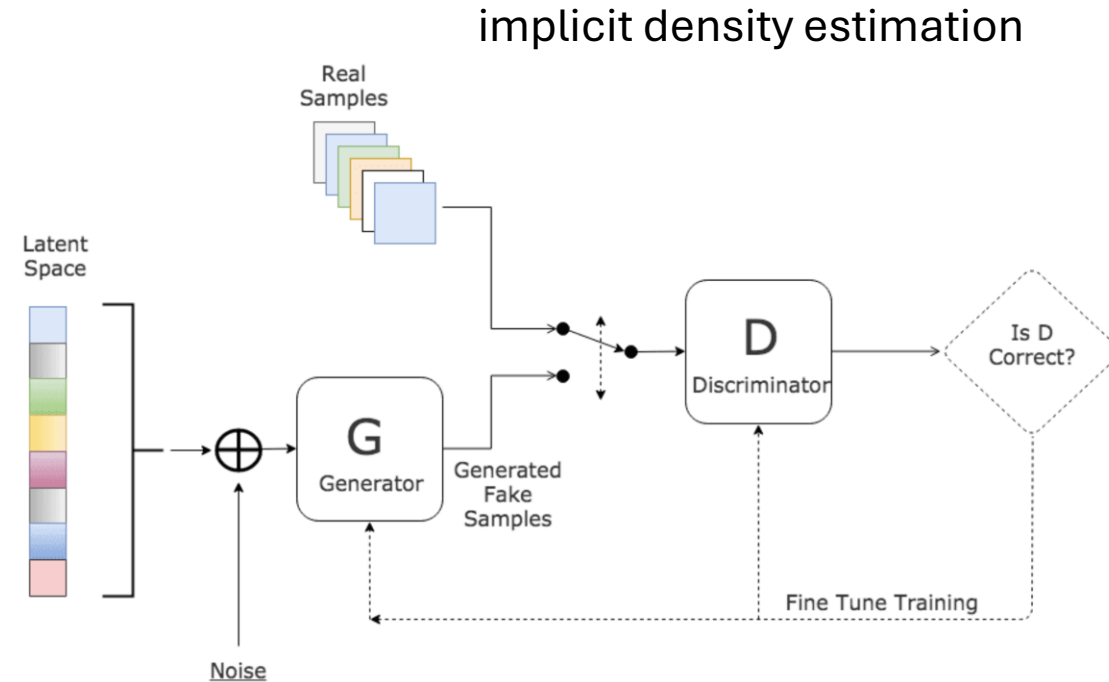
indirect training via D: G not trained directly to minimize reconstruction error of real samples, but to fool D → self-supervised approach

implicit density estimation



[source](#)

common loss for generator and discriminator:
$$L(\boldsymbol{x}_i) = E_{\boldsymbol{x} \sim p_r(\boldsymbol{x})}[\ln D(\boldsymbol{x}_i)] + E_{\boldsymbol{x} \sim p_g(\boldsymbol{x})}[\ln(1 - D(\boldsymbol{x}_i))]$$
G trying to minimize
D trying to maximize

# Conditional GANs

as discussed so far, generative methods give no control over what kind of data is generated (limited usability)

→ need for conditional approach (e.g., conditioning on describing text)

example GANs:

transform usual GAN to conditional model by feeding extra information $y$ (e.g., class labels) as additional input layer into both generator and discriminator

$$L(\boldsymbol{x}_i) = E_{\boldsymbol{x} \sim p_r(\boldsymbol{x})}[\ln D(\boldsymbol{x}_i|y_i)] + E_{\boldsymbol{x} \sim p_g(\boldsymbol{x})}[\ln(1 - D(\boldsymbol{x}_i|y_i))]$$



[source](#)

# Vector Arithmetic in GAN Latent Space



smiling woman − neutral woman + neutral man = smiling man

man with glasses − man without glasses + woman without glasses = woman with glasses

# Variational Autoencoder (VAE)

goal: generation of variations of input data rather than compressed representation

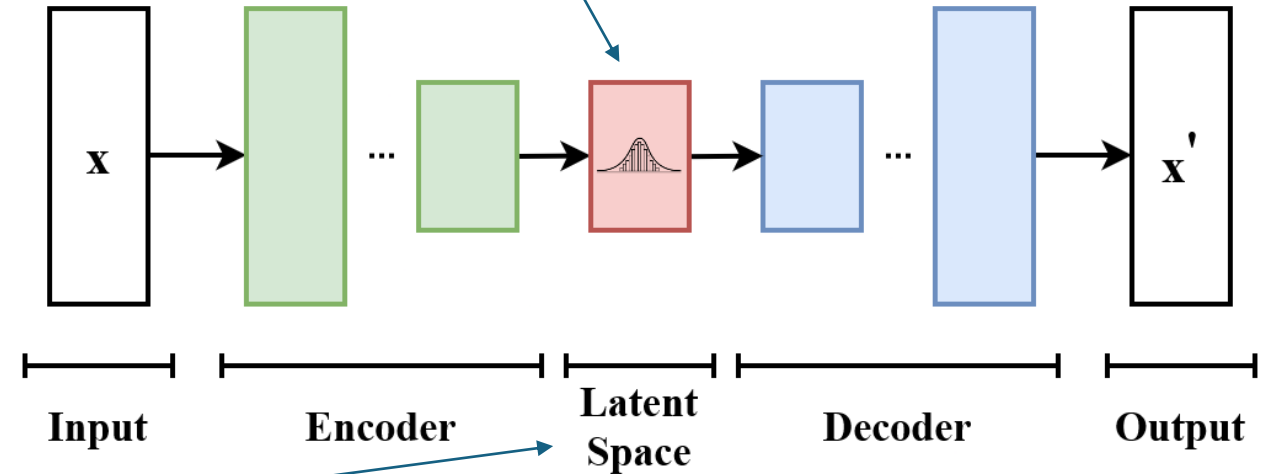→ learn variational distribution instead of identity function

to be precise: parametrized variational distribution of latent encoding variables $\boldsymbol{z}$

prior (simple distribution, in usual VAE: Gaussian): $p_{\boldsymbol{\theta}}(\boldsymbol{z})$

posterior: $p_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}) = \dfrac{p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})p_{\boldsymbol{\theta}}(\boldsymbol{z})}{\int p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})p_{\boldsymbol{\theta}}(\boldsymbol{z})d\boldsymbol{z}}$

$p_{\boldsymbol{\theta}}(\boldsymbol{x})$: mixture of Gaussians

from which to sample



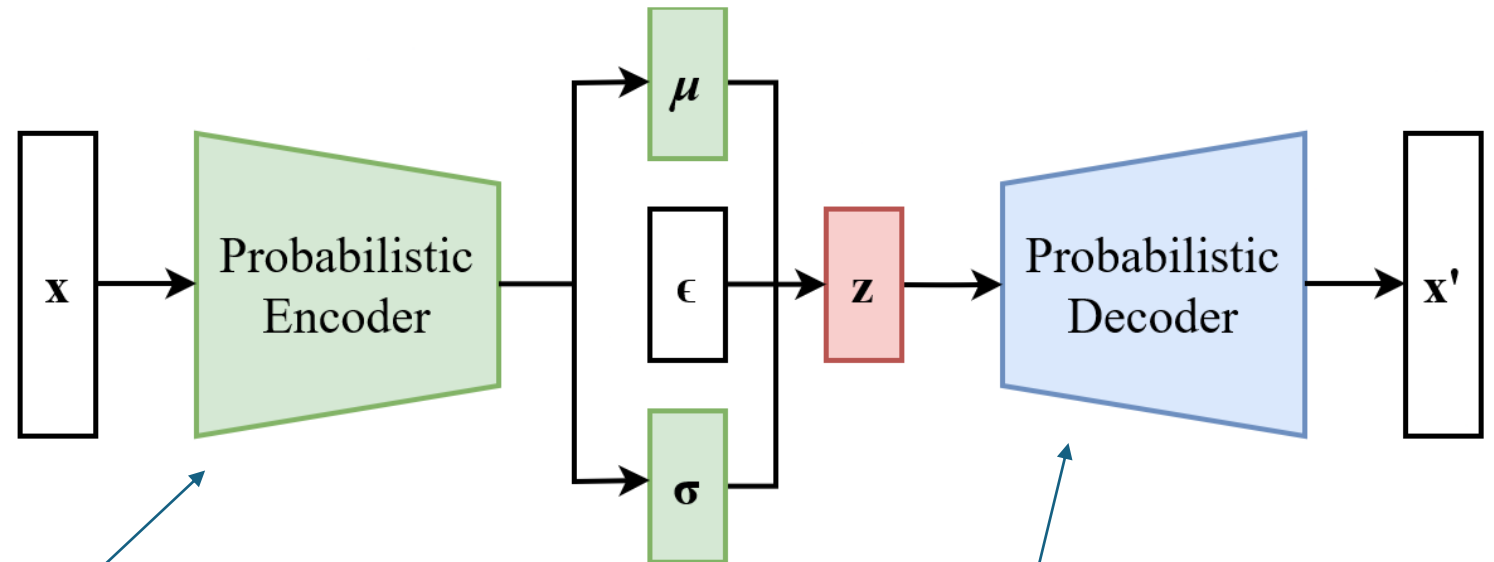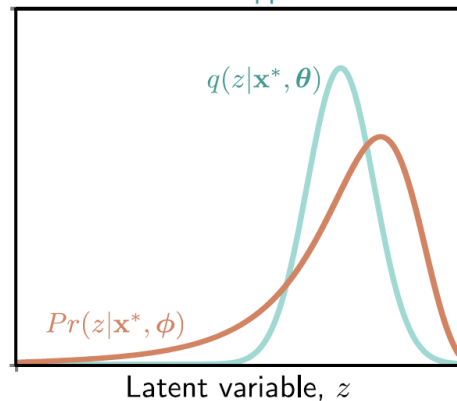| Input | Encoder | Latent Space | Decoder | Output |

Variational Bayesian Method

12

# Gaussian Approximation

learn mean and variance of multivariate Gaussian with diagonal covariance structure



good approximation:

poor approximation:

Posterior and approximation

Posterior and approximation

$q(z|\mathbf{x}^*, \boldsymbol{\theta})$

$Pr(z|\mathbf{x}^*, \boldsymbol{\phi})$

$q(z|\mathbf{x}^*, \boldsymbol{\theta})$

$Pr(z|\mathbf{x}^*, \boldsymbol{\phi})$

Latent variable, $z$

Latent variable, $z$

$Pr(x, z)$

$z = 1$
$z = 2$
$z = 3$

Marginalize over latent variable, $z$

$Pr(x)$

$x$

source

13

# Diffusion

training: distort training data by successively adding random noise, then learn to reverse this process (denoising)

generation: sample random noise and run through the learned denoising process
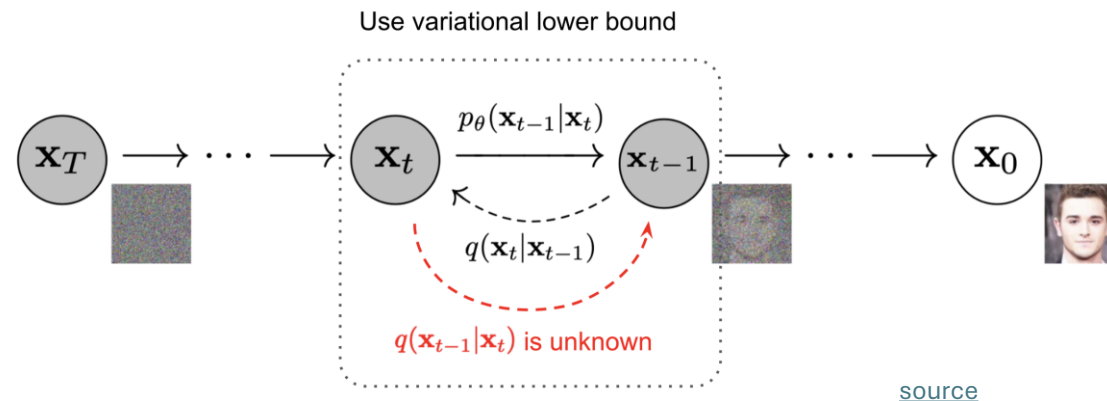


Use variational lower bound

$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$

$q(\mathbf{x}_t|\mathbf{x}_{t-1})$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is unknown

source

advantages: easy to train, produce high-quality/realistic samples

can be interpreted as special case of hierarchical VAE (one latent variable generates another) with fixed encoder and latent space of same size as the data

→ more sophisticated latent space than just Gaussian mixture in VAE

# Noise Prediction

**Algorithm 1** Training

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

**Algorithm 2** Sampling

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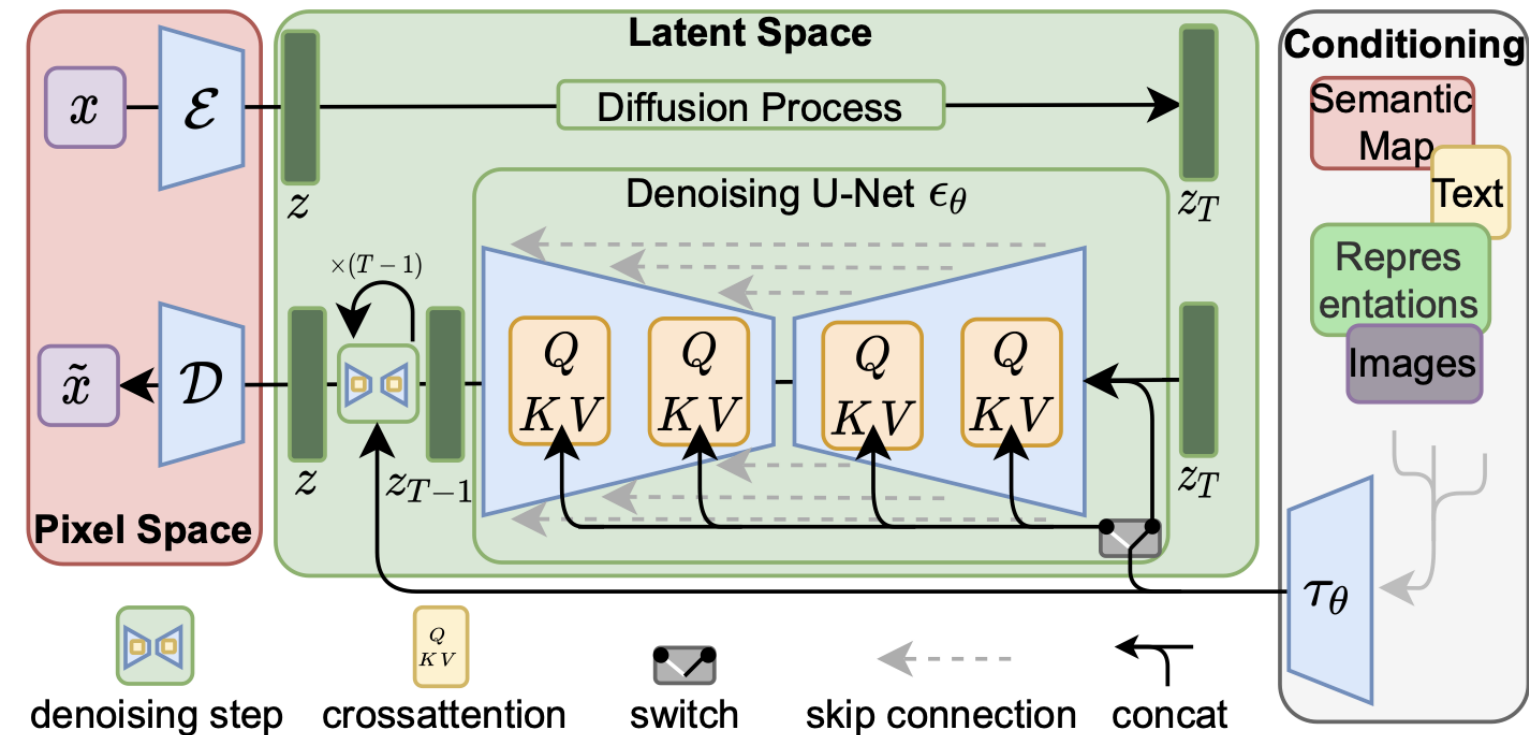
source

# Diffusion as Chain of Denoising Autoencoders

differences of diffusion models to typical denoising autoencoders:

- no bottleneck (care about output here, not internal representation): latent space with high dimensionality (same as original data)

- handle many different noise levels with single set of shared parameters

important application: AlphaFold 3 uses diffusion-based architecture for protein structure prediction

# Latent Diffusion

add noise to latent representation rather than raw data
→ significant speedup



- convolution and transposed convolution layers
- skip connections between layers operating at the same scale
- use of attention mechanism for flexible conditioning

alternative to convolutional U-nets: vision transformers (e.g., DiT)

source

# Guided Diffusion

condition diffusion process on class information (label or just text)

$$\epsilon_\theta(x_t|\emptyset) + s \cdot \left(\epsilon_\theta(x_t|y) - \epsilon_\theta(x_t|\emptyset)\right)$$

$s$: hyperparameter to control tradeoff between diversity (unconditioned) and fidelity (guidance)

similar idea as softmax temperature in auto-regressive LLMs

"Pembroke Welsh corgi"

source

# Outlook: Text2Anything

next step: text-to-video (Make-A-Video, Lumiere, Sora, …)
→ rudimentary physics understanding


at some point maybe also generation of proteins, materials, …