

# Generative Models

Nikita Kiselev

Deep Learning, Intelligent Systems

November 17, 2025



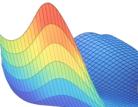
# Nikita Kiselev

## Research Interests

- 1) Generative AI
- 2) Computer Vision
- 3) Optimization

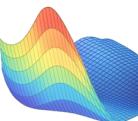
## Topics in this course

- 1) Generative Models
- 2) Acceleration



## Education

Sep 2024 – Present  
MSc in Computer Science  
Intelligent Systems, MIPT



Sep 2020 – Aug 2024  
BSc in Applied Mathematics and Physics  
Intelligent Systems, MIPT



## Experience

Jun 2024 – Present  
Research Engineer  
Kandinsky, Sber AI



Sep 2024 – Jul 2025  
Mathematician-programmer  
Research Center for Artificial Intelligence, Innopolis University



Oct 2023 – Apr 2024  
Technician  
Laboratory of Mathematical Methods of Optimization, MIPT

# Which face is real?



1



2



3

# Which face is real?



1



2

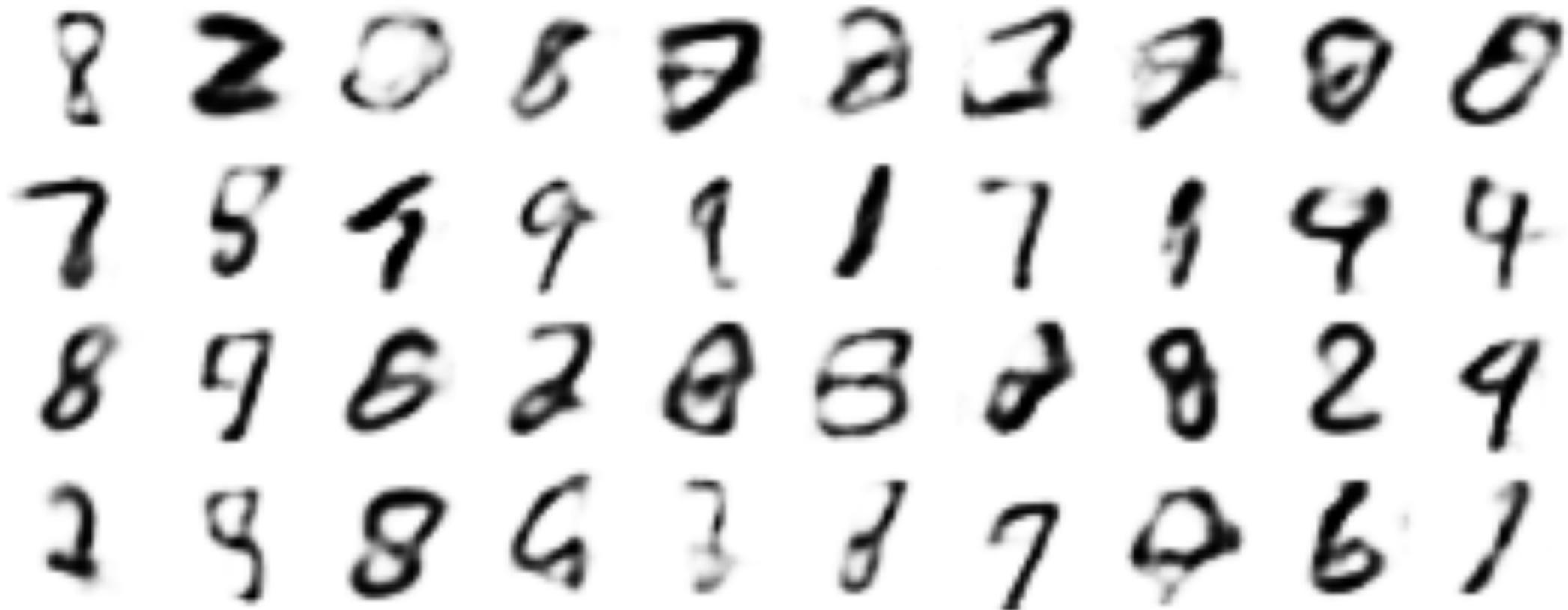


3

NONE OF THEM! ALL ARE AI-GENERATED!

# Evolution of Generative Models

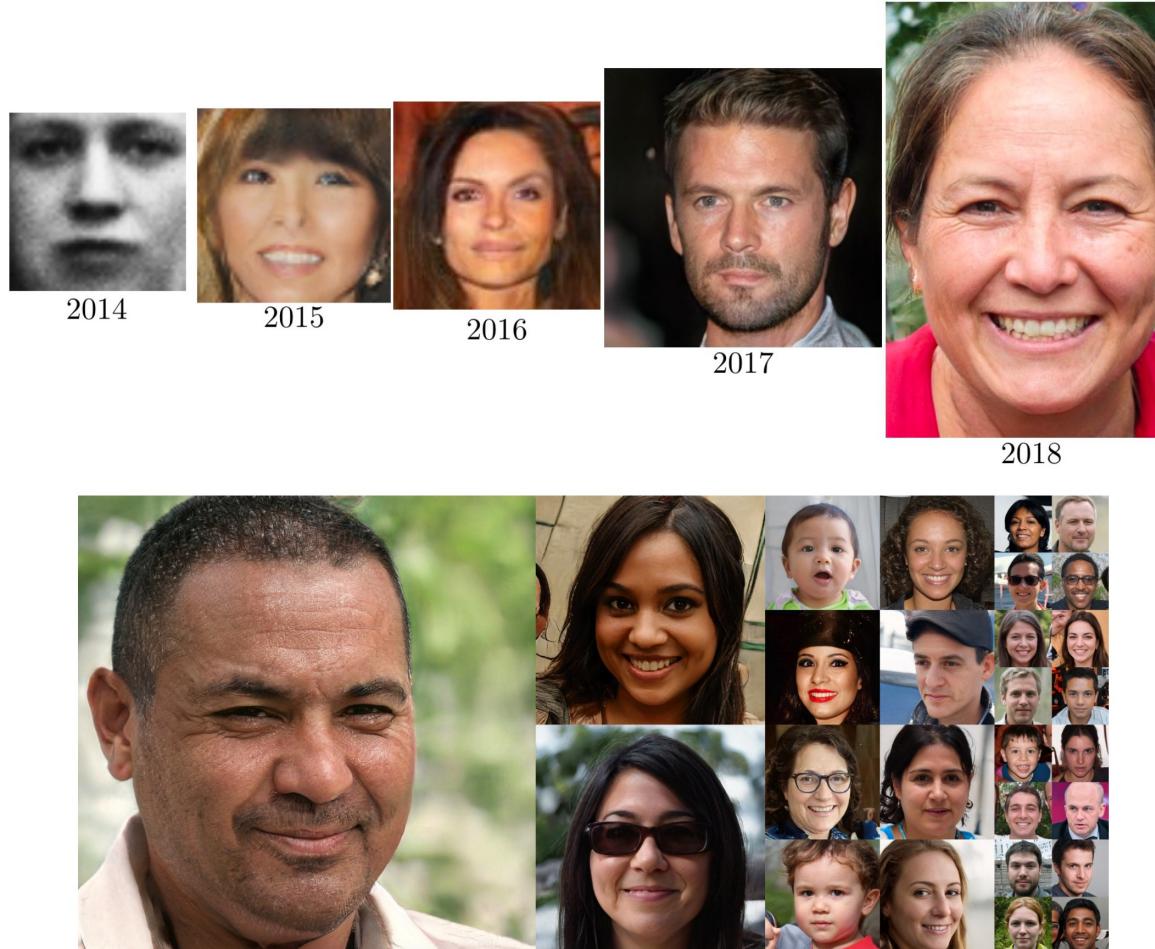
# VAE – first scalable approach



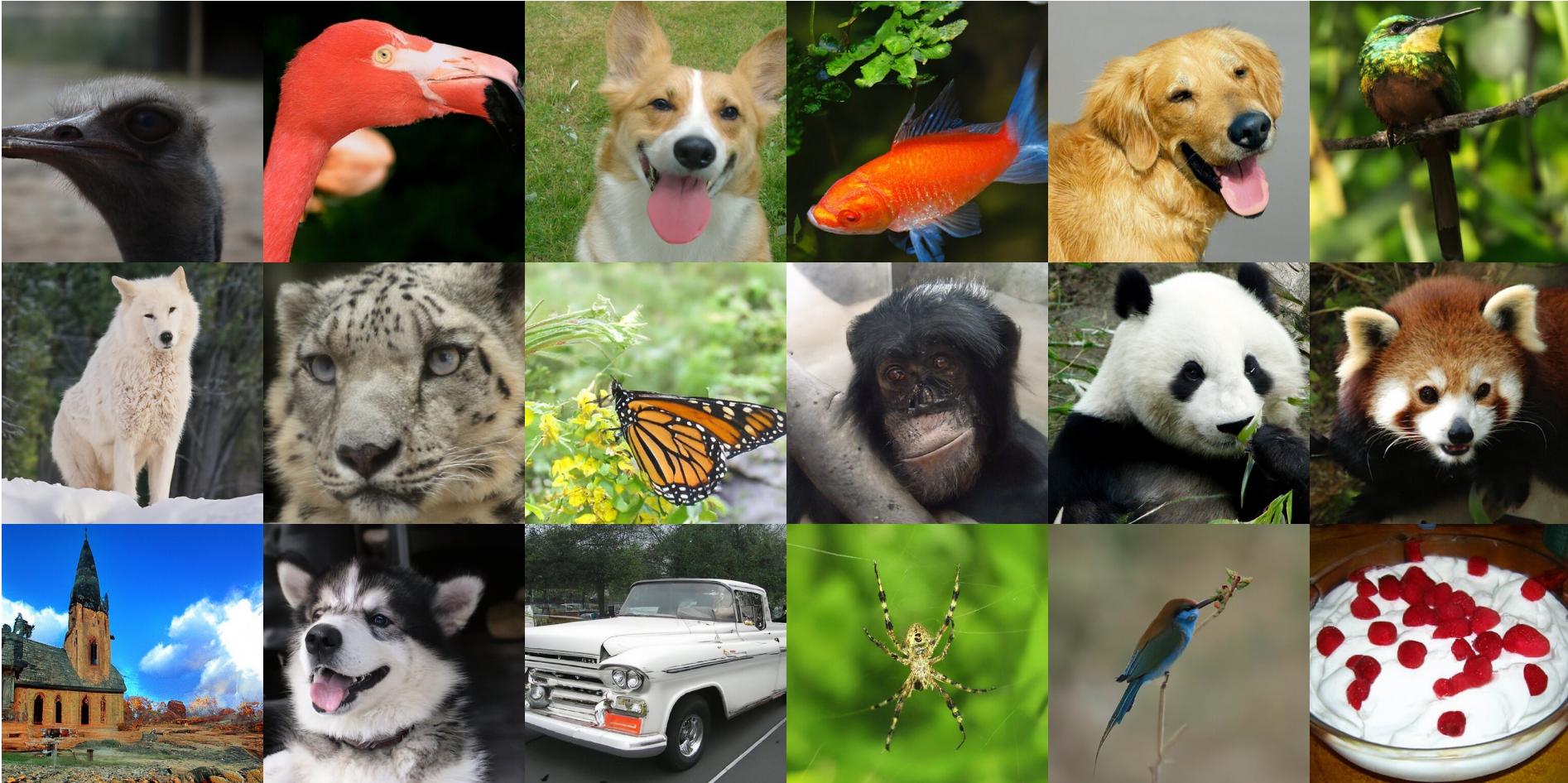
# DCGAN – first convolutional GAN



# StyleGAN – high quality faces



# Denoising Diffusion Probabilistic Models



# FLUX – cutting-edge flow matching model



# Veo3 – SOTA video + audio generation



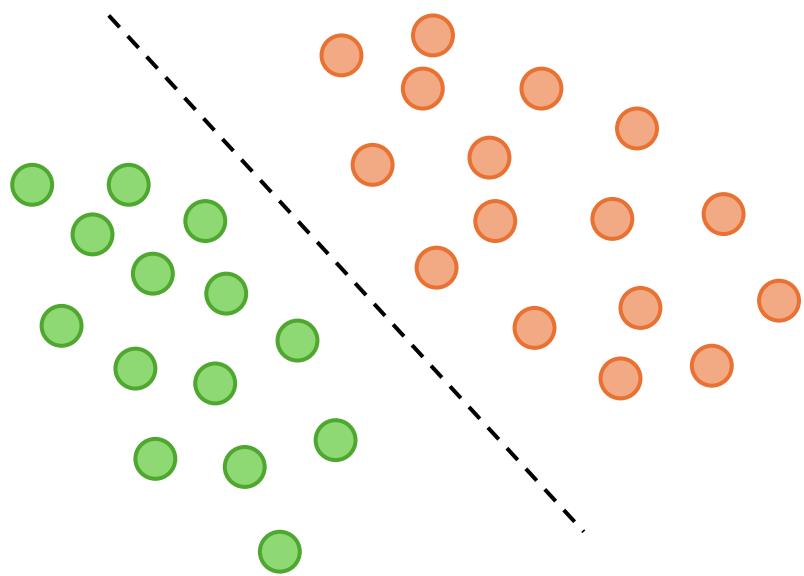
<https://deepmind.google/models/veo/>

# Contents

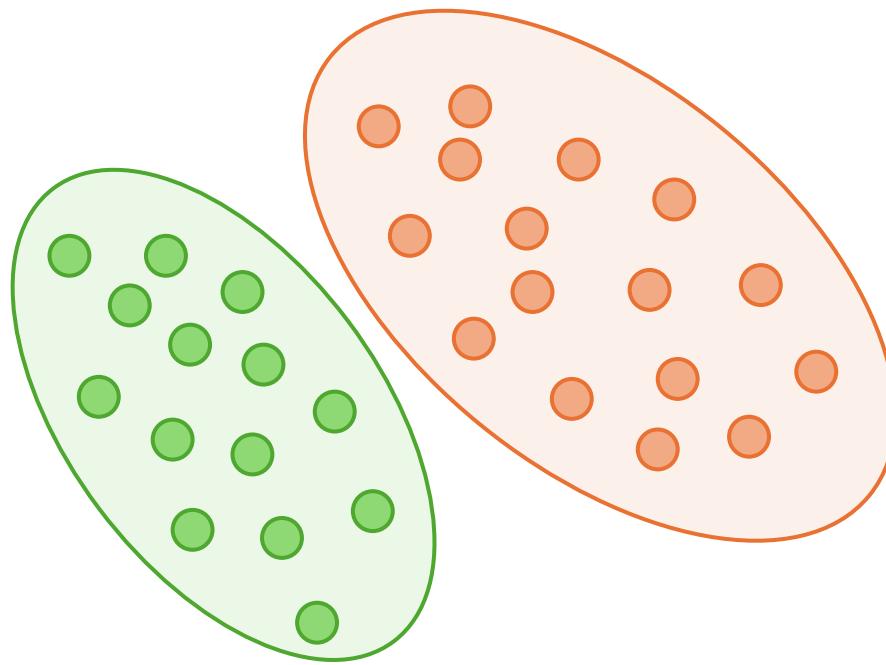
- What are Generative Models?
- Taxonomy
- Autoregressive Models
- Autoencoders
- VAEs
- GANs
- \* Extra: Modern Architectures

# What are Generative Models?

# Discriminative vs Generative Models



Discriminative  
 $p(\mathbf{y}|\mathbf{x})$



Generative  
 $p(\mathbf{x}, \mathbf{y})$

# Discriminative vs Generative Models

## Discriminative

- Learn to find the **boundary between classes**
- Model probability of a label given features
- Focus: classification, regression
- Examples: Linear Regression, Logistic Regression, SVM, ResNet, etc.

## Generative

- Learn to model the **data distribution itself**
- Model joint probability of both features and label
- Focus: generating new samples, understanding data structure
- Examples: [we will see on this lecture]

# Generative Modeling

## Given

$D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  from  $p_{\text{data}}(\mathbf{x})$

## Goals

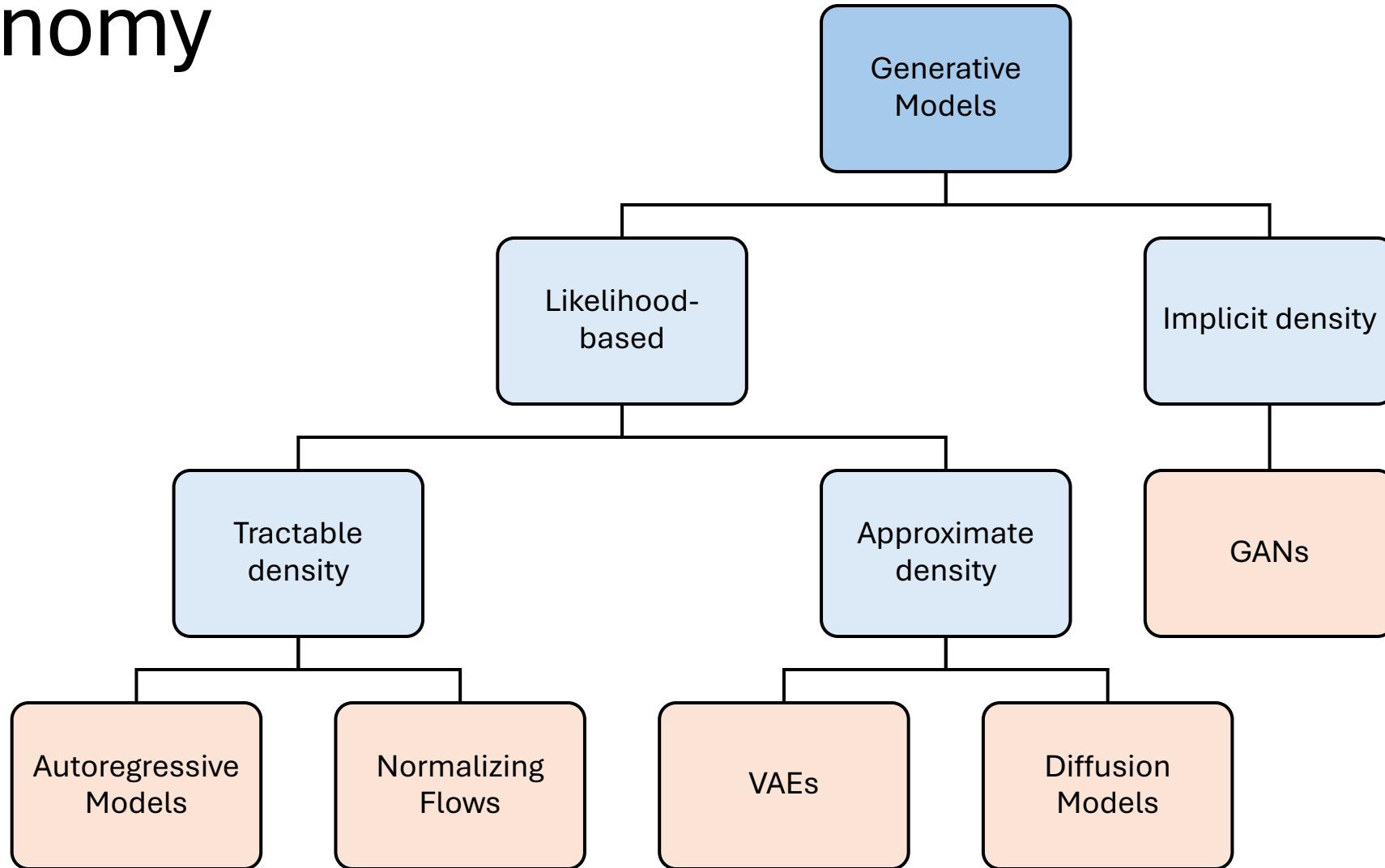
1. Produce new samples from  $p_{\text{data}}(\mathbf{x})$
2. Obtain meaningful data representations
3. Estimate density for data points



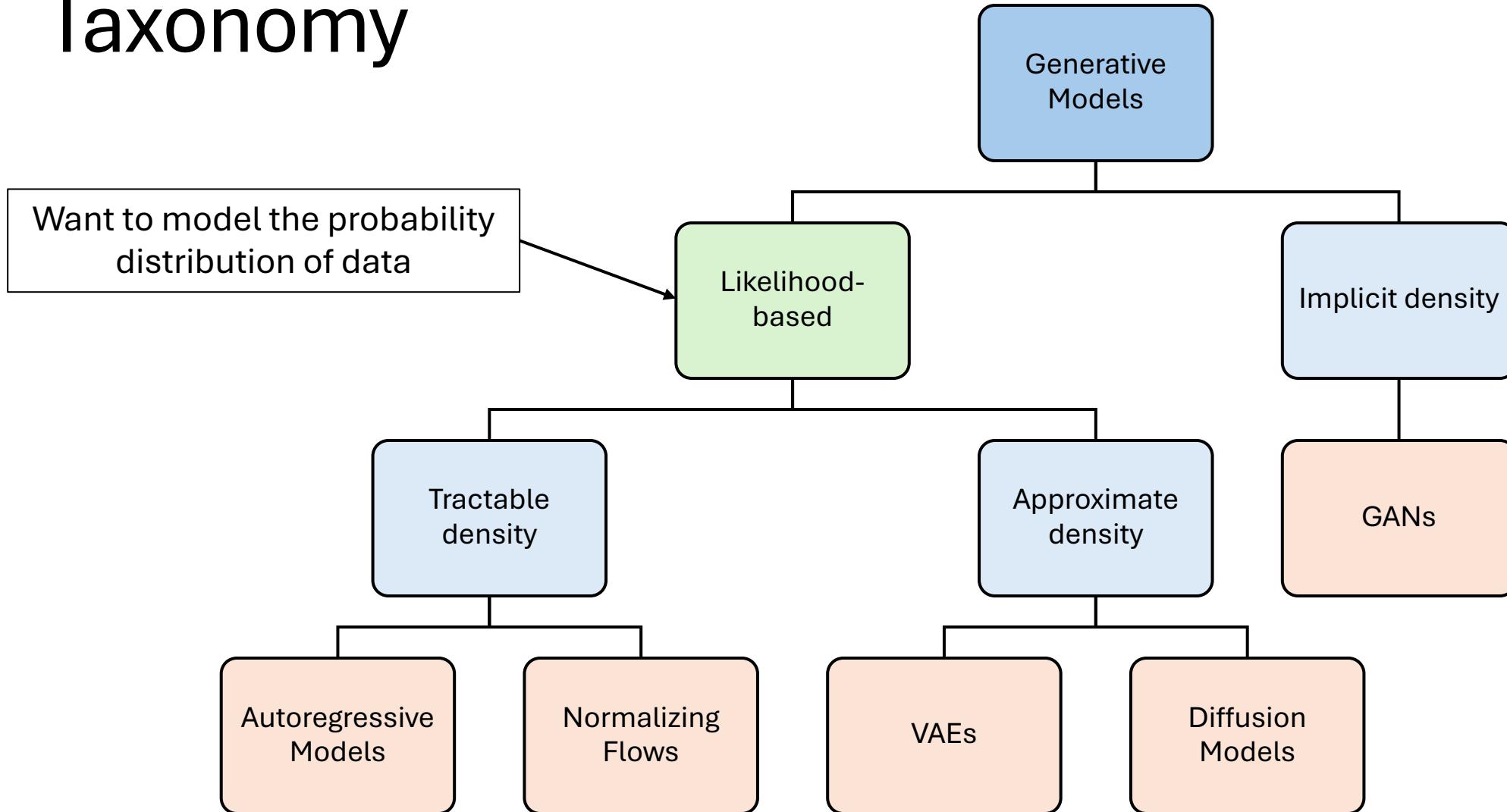
Data Distribution

# Taxonomy

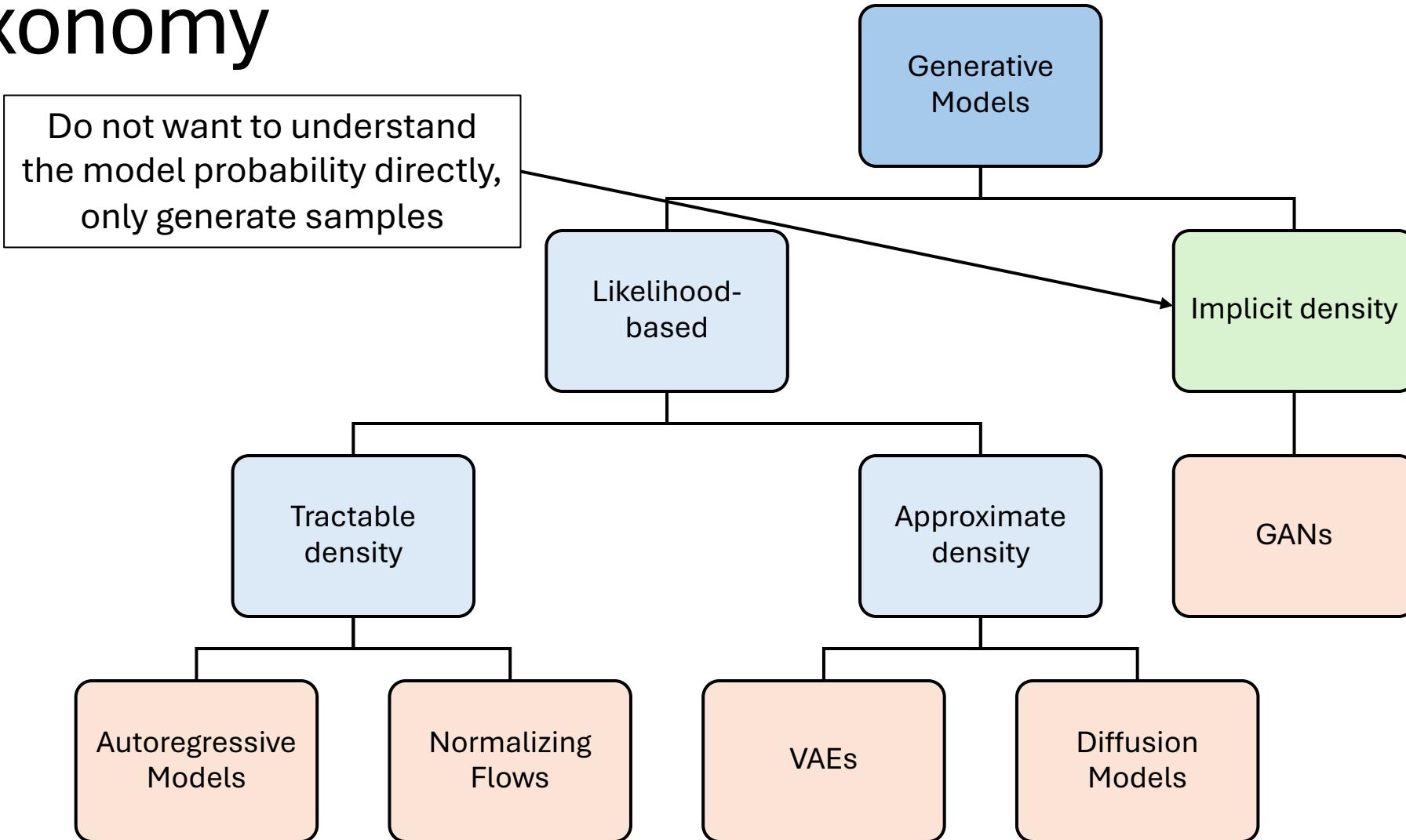
# Taxonomy



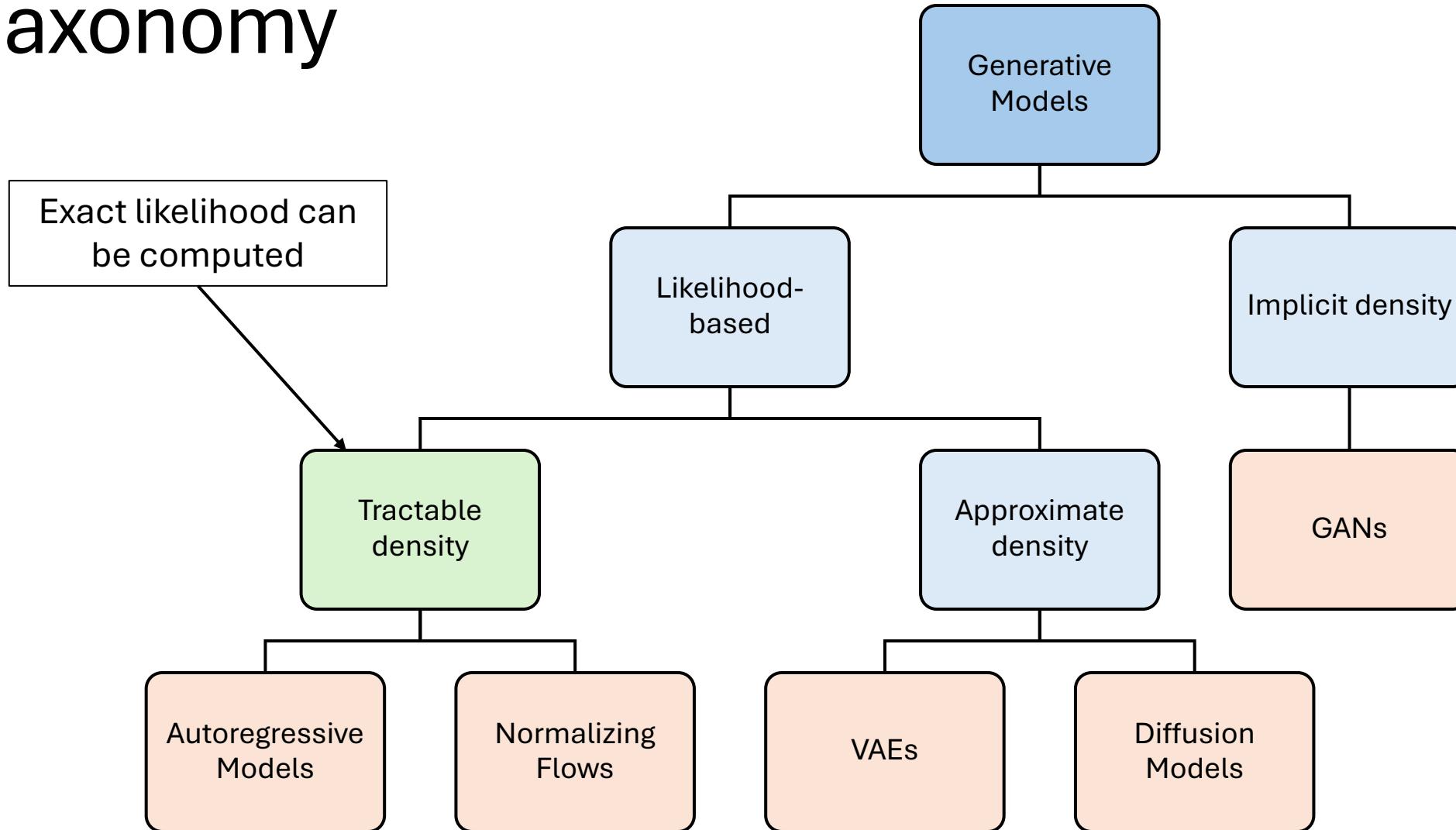
# Taxonomy



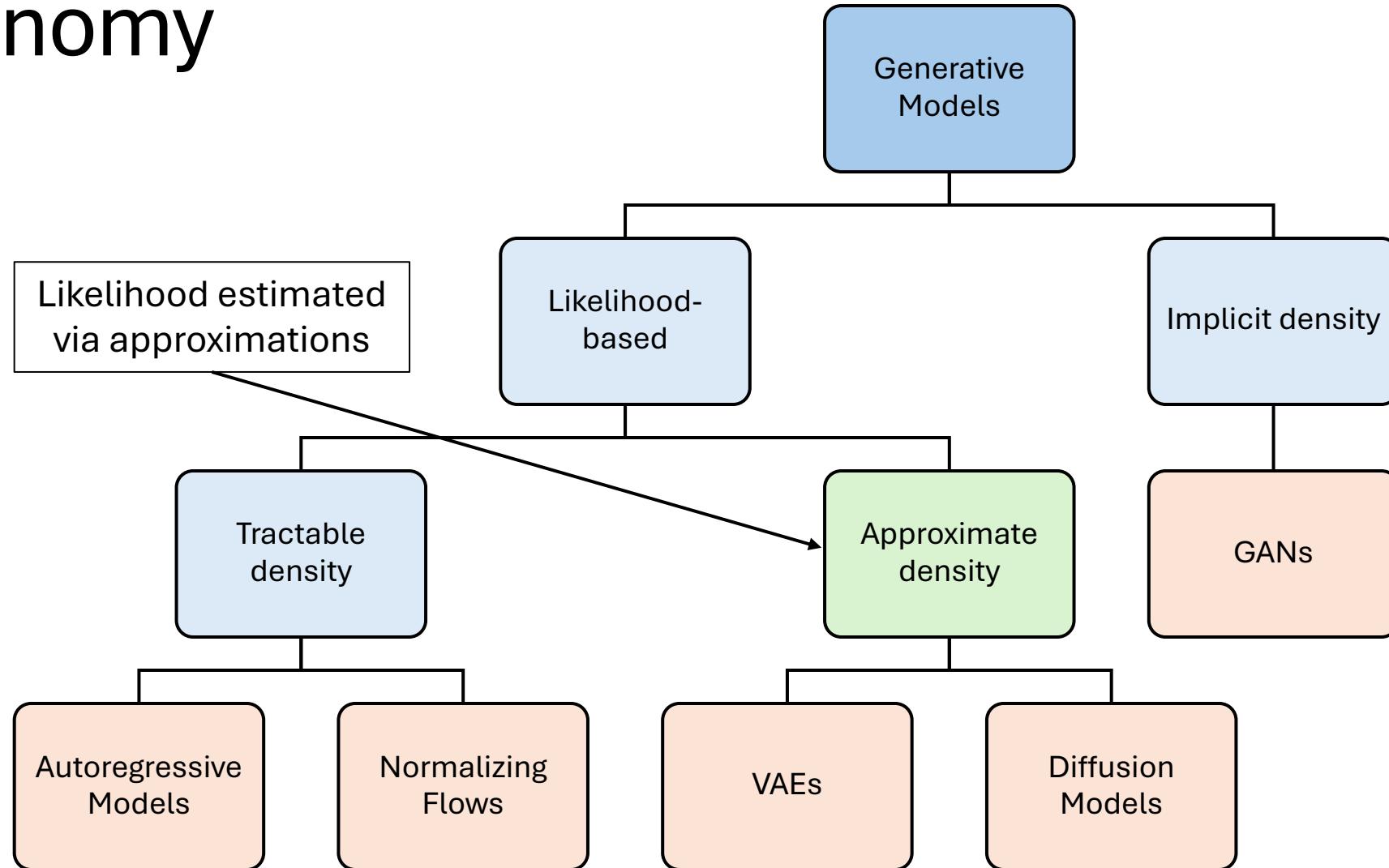
# Taxonomy



# Taxonomy

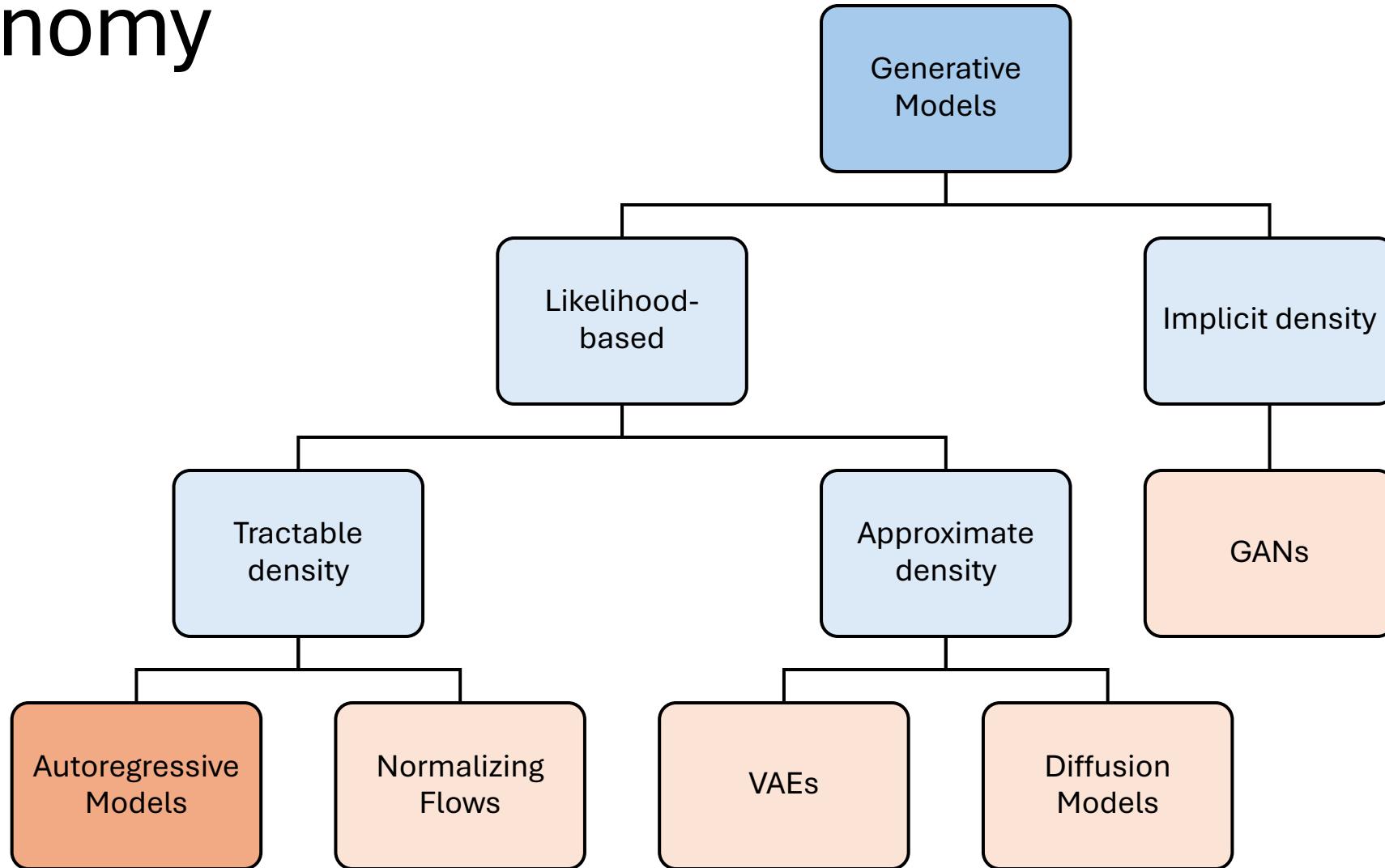


# Taxonomy



# Autoregressive Models

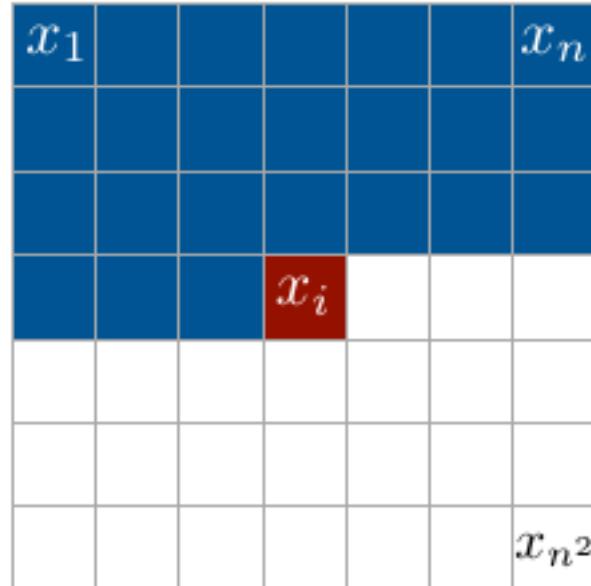
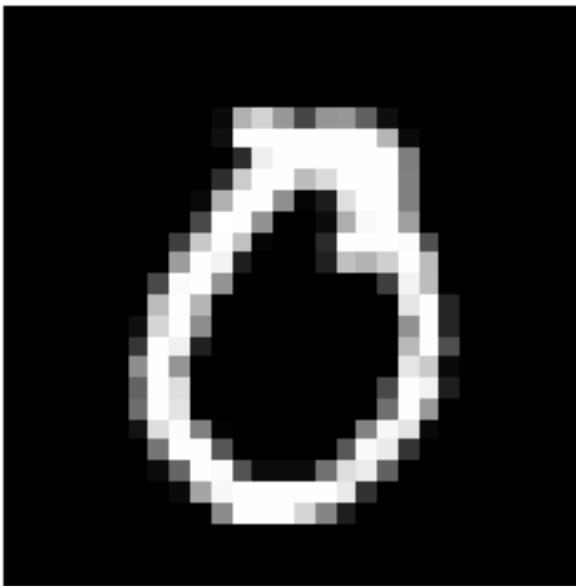
# Taxonomy



# PixelRNN & PixelCNN

**Goal:** generate binary images (e.g., MNIST)

**Idea:** represent joint distribution  $p(\mathbf{x})$  via chain rule (cond. prob)



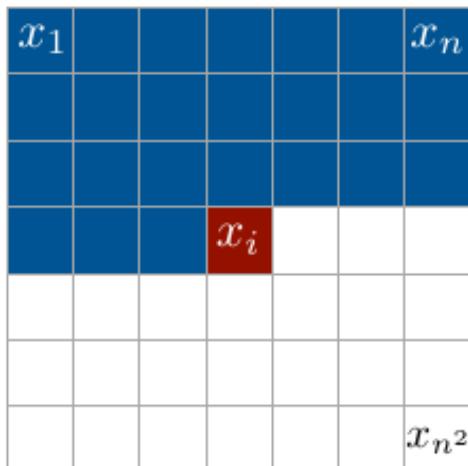
## Chain Rule

$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

# Parameterization

$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | \mathbf{x}_{<i}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

Let's parameterize history with some function (neural network  $f_\theta$ ).  
Thus we move to the conditional distribution  $p_\theta(x_i | \mathbf{x}_{<i})$ .



$$p_\theta(x_i | \mathbf{x}_{<i}) = \text{Softmax}(f_\theta(\mathbf{x}_{<i}))[x_i]$$

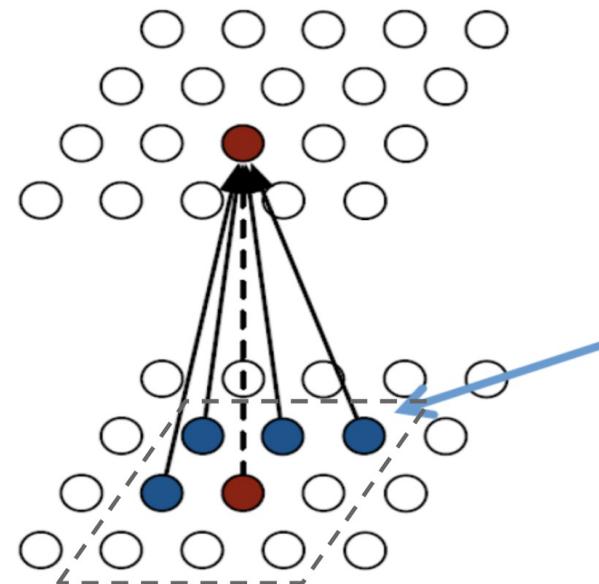
Predict the next pixel based on the previous ones

# How to choose the history function?

**Problem:** long context

**Solutions:**

- RNN
- LSTM
- Cut context → CNN



PixelCNN

Masked convolution

1	1	1
1	0	0
0	0	0

# How to train autoregressive image models?

**MLE:** maximize likelihood  $p_{\theta}(\mathbf{x})$  w.r.t. model parameters  $\theta$ .

We train by minimizing negative log-likelihood:

$$\mathcal{L}(\theta) = - \sum_{i=1}^{n^2} \log p_{\theta}(x_i | \mathbf{x}_{<i})$$

In case of binary images from MNIST it transforms into BCE loss:

$$-\log p_{\theta}(x_i | \mathbf{x}_{<i}) = - \left[ x_i \log \sigma(f_{\theta}(\mathbf{x}_{<i})) + (1 - x_i) \log (1 - \sigma(f_{\theta}(\mathbf{x}_{<i}))) \right]$$

# Pros and Cons

## Advantages

- Exact likelihood training
- Good sample diversity

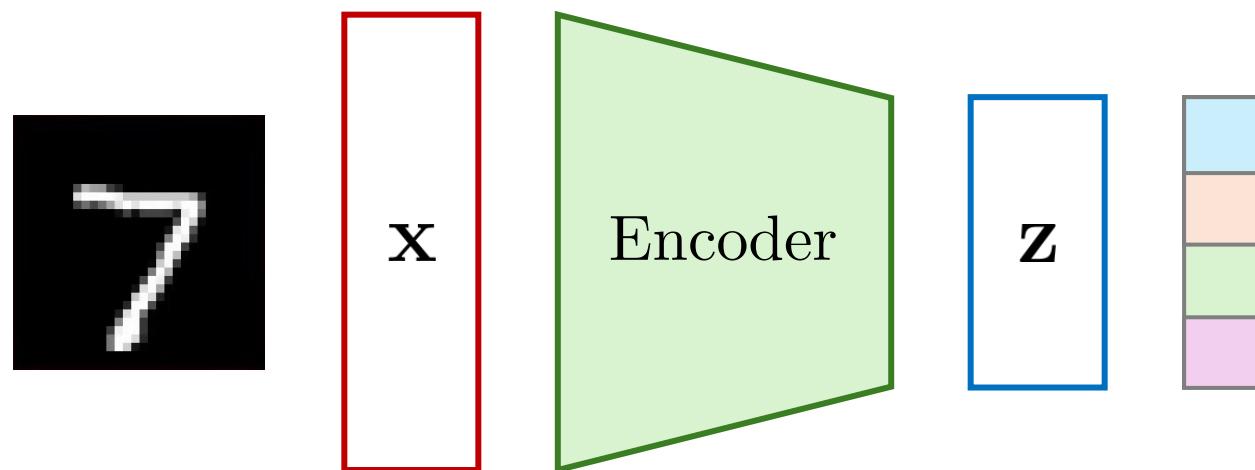
## Disadvantages

- Sequential generation – pixels generated one-by-one
- Limited global coherence – struggles to capture long-range spatial dependencies
- Scalability – computationally expensive for high-resolution data

# Autoencoders

# Autoencoders: background

Unsupervised approach for learning a **lower-dimensional** feature representation from unlabeled **training data**.

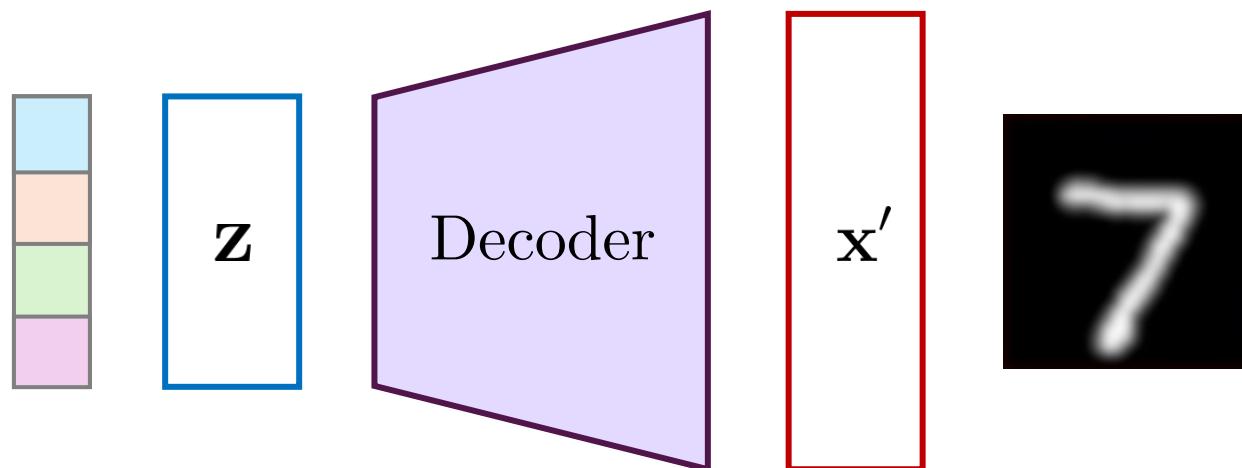


Encoder learns mapping from the data,  $x$ , to a **low-dimensional** latent space,  $z$ .

**Analogues:**  
PCA, Matrix Factorization

# How can we learn this latent space?

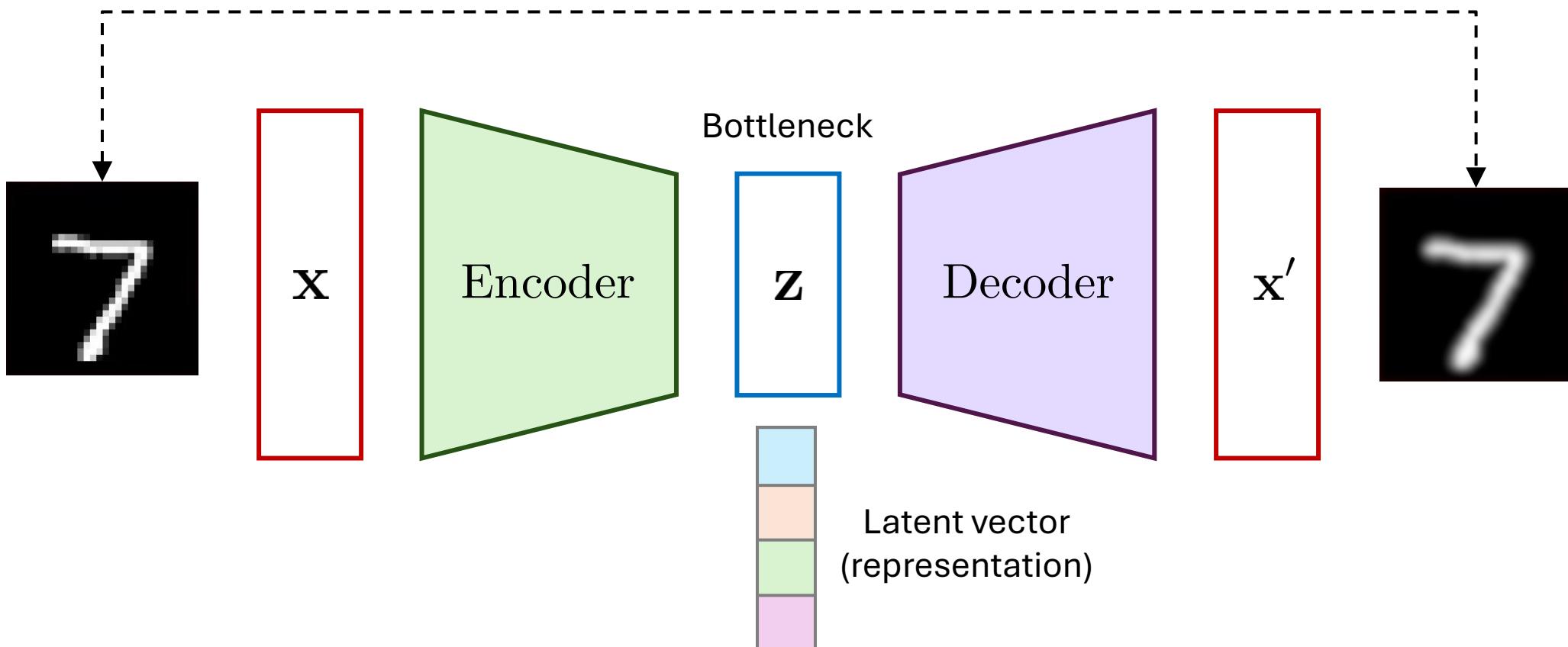
Train **another model** to use **these features** to **reconstruct the input**.



**Decoder** learns the mapping back from latent space,  $z$ , to a reconstructed observation,  $x'$ .

# Reconstruction loss

Basic MSE loss function:  $\mathcal{L} = \|\mathbf{x} - \mathbf{x}'\|_2^2$



# Representation Learning

**Bottleneck hidden layer** forces network to learn a compressed latent representation.

**Reconstruction loss**  $\mathcal{L} = \|\mathbf{x} - \mathbf{x}'\|_2^2$  forces the latent representation  $\mathbf{z}$  to capture (or encode) as much “information” about the data  $\mathbf{x}$  as possible.

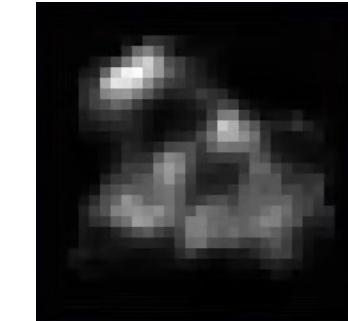
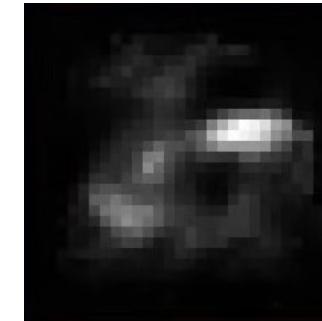
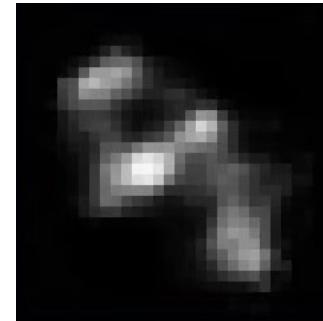
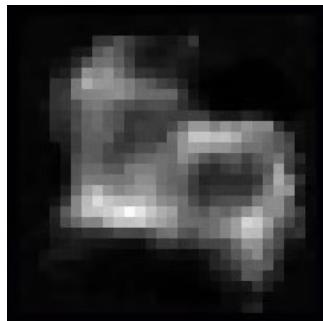
**Autoencoding** = Automatically **encoding** data  
“Auto” = **self**-encoding

# Why not simple Autoencoders?

**Goal:** generate new data using Autoencoder's latent space

**Idea:** randomly sample points from the latent space and feed them through the decoder

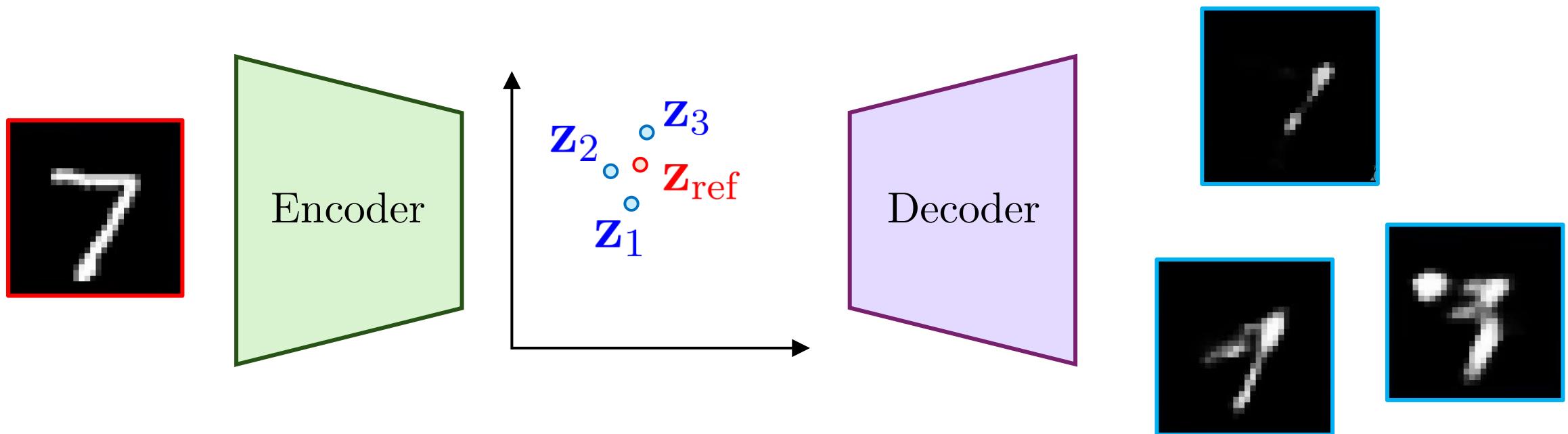
**Problem:** in most cases, this approach will not produce useful results, mainly because the **latent space is disorganized** and irregular – so large areas of it will not produce meaningful data



# Sampling latent vectors near a reference

**Goal:** obtain images that resemble the reference one

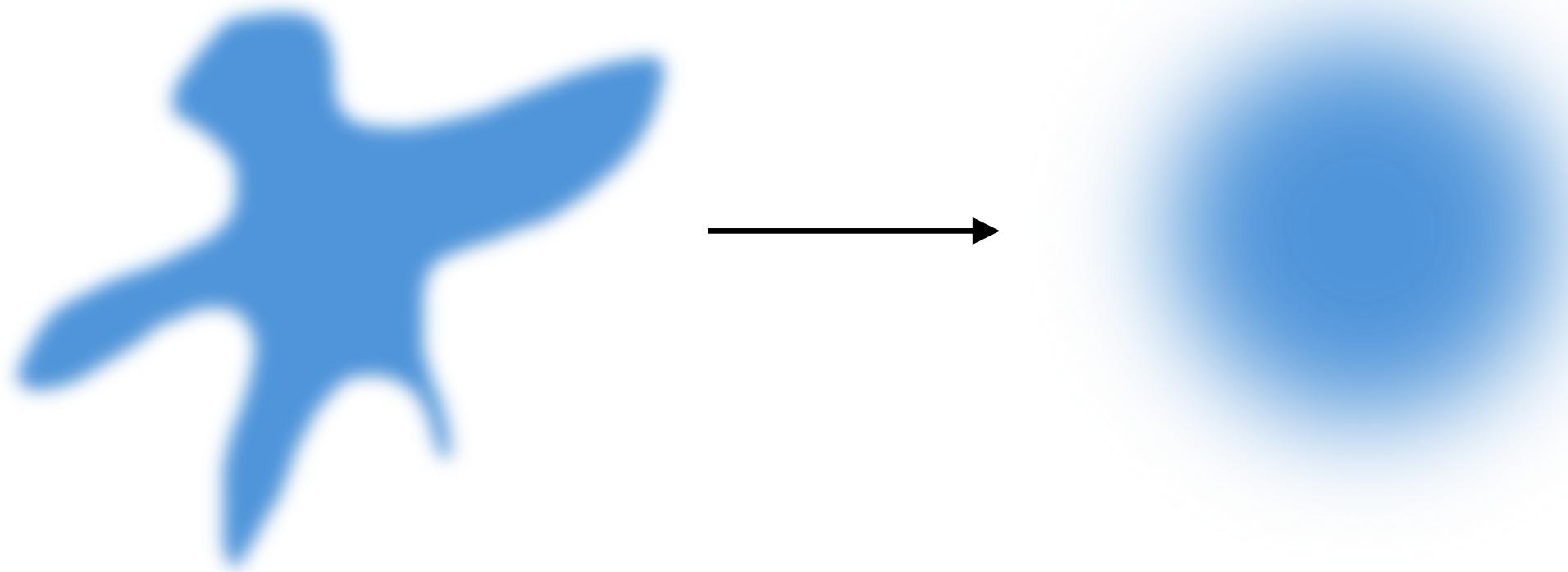
**Problem:** the latent space is so poorly structured that even nearby points often don't correspond to meaningful variations



# What is the solution?

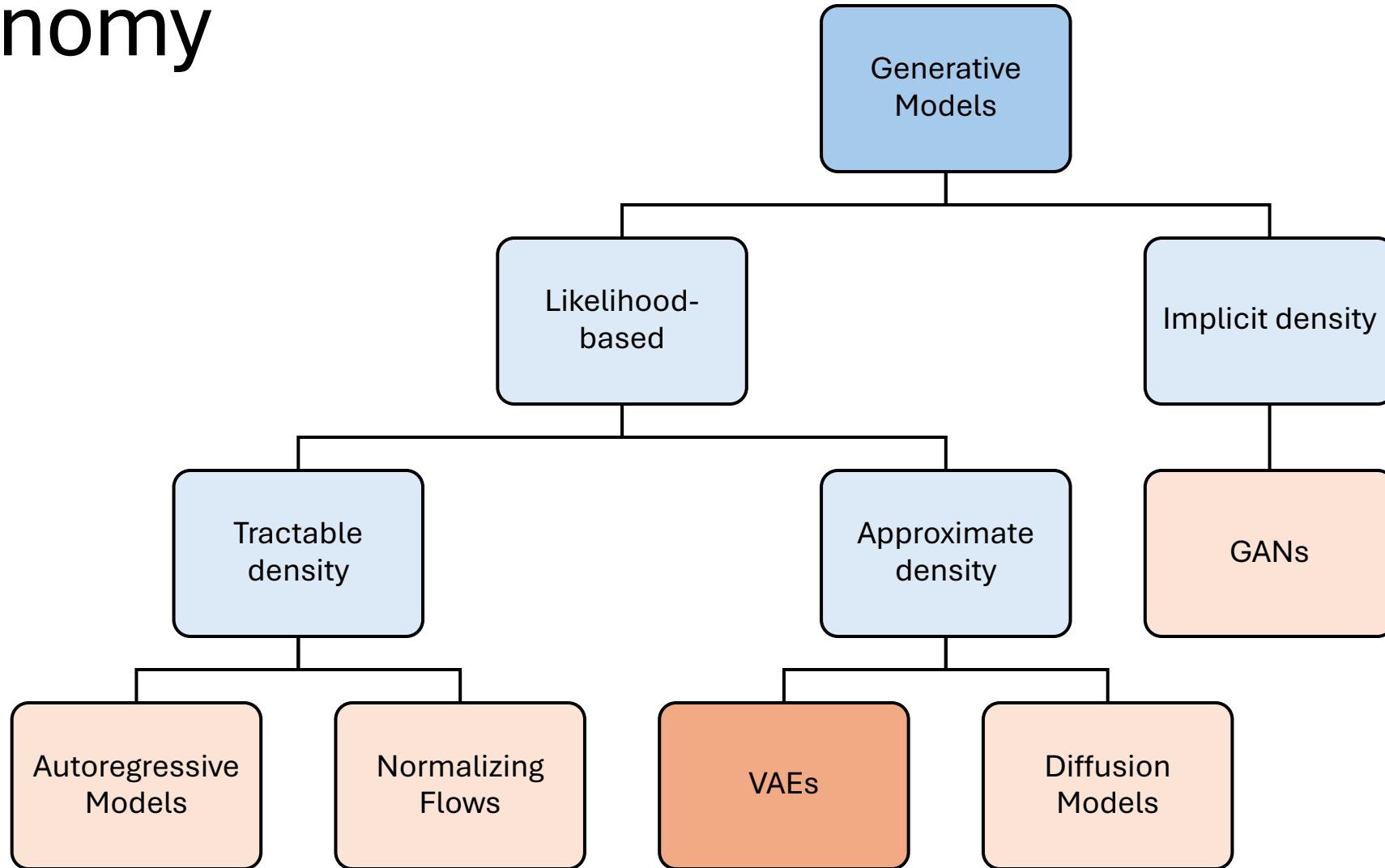
Organize the latent space!

But how...?



# VAEs

# Taxonomy



# Data and Latent Distributions

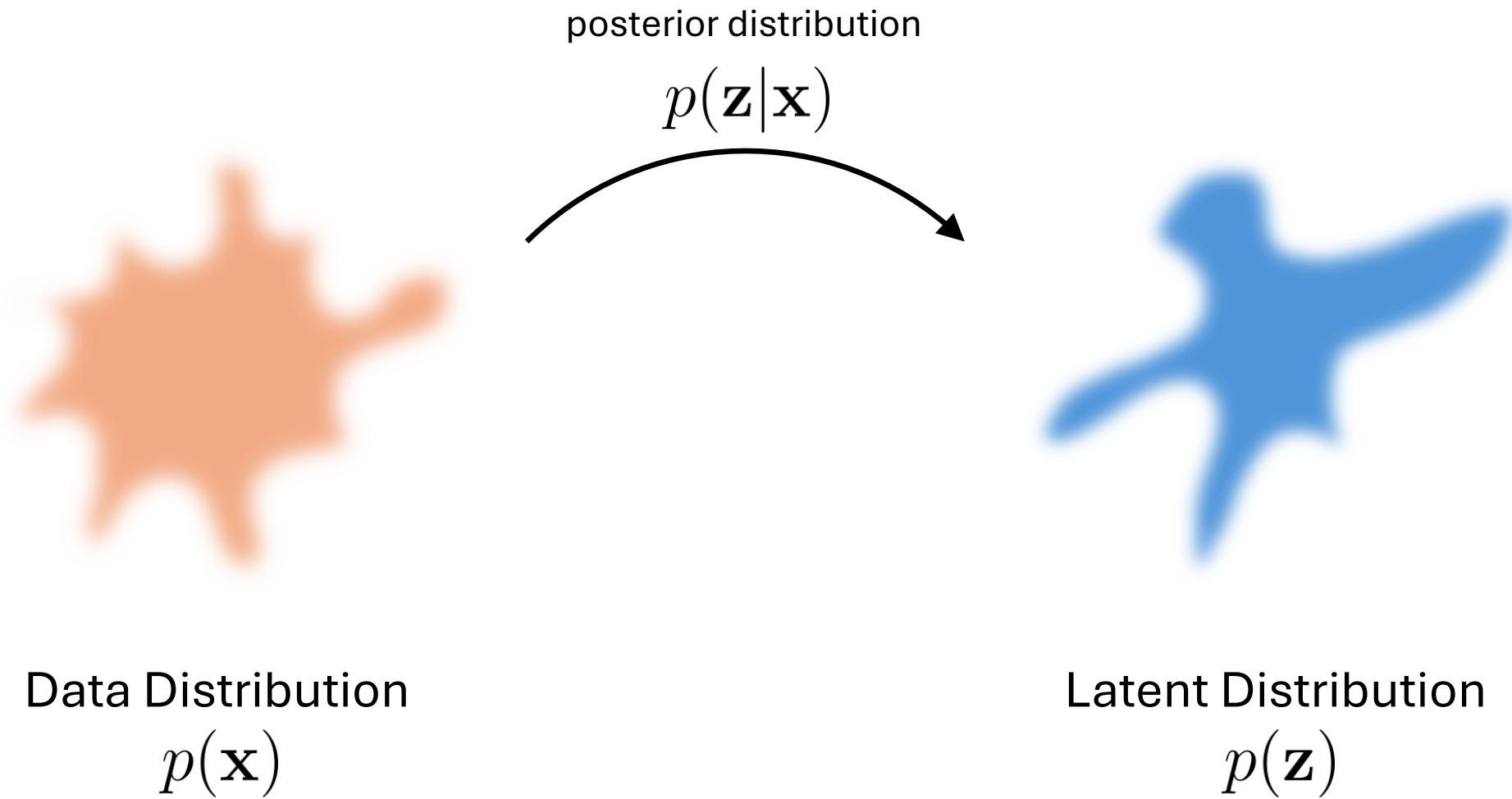


Data Distribution  
 $p(\mathbf{x})$

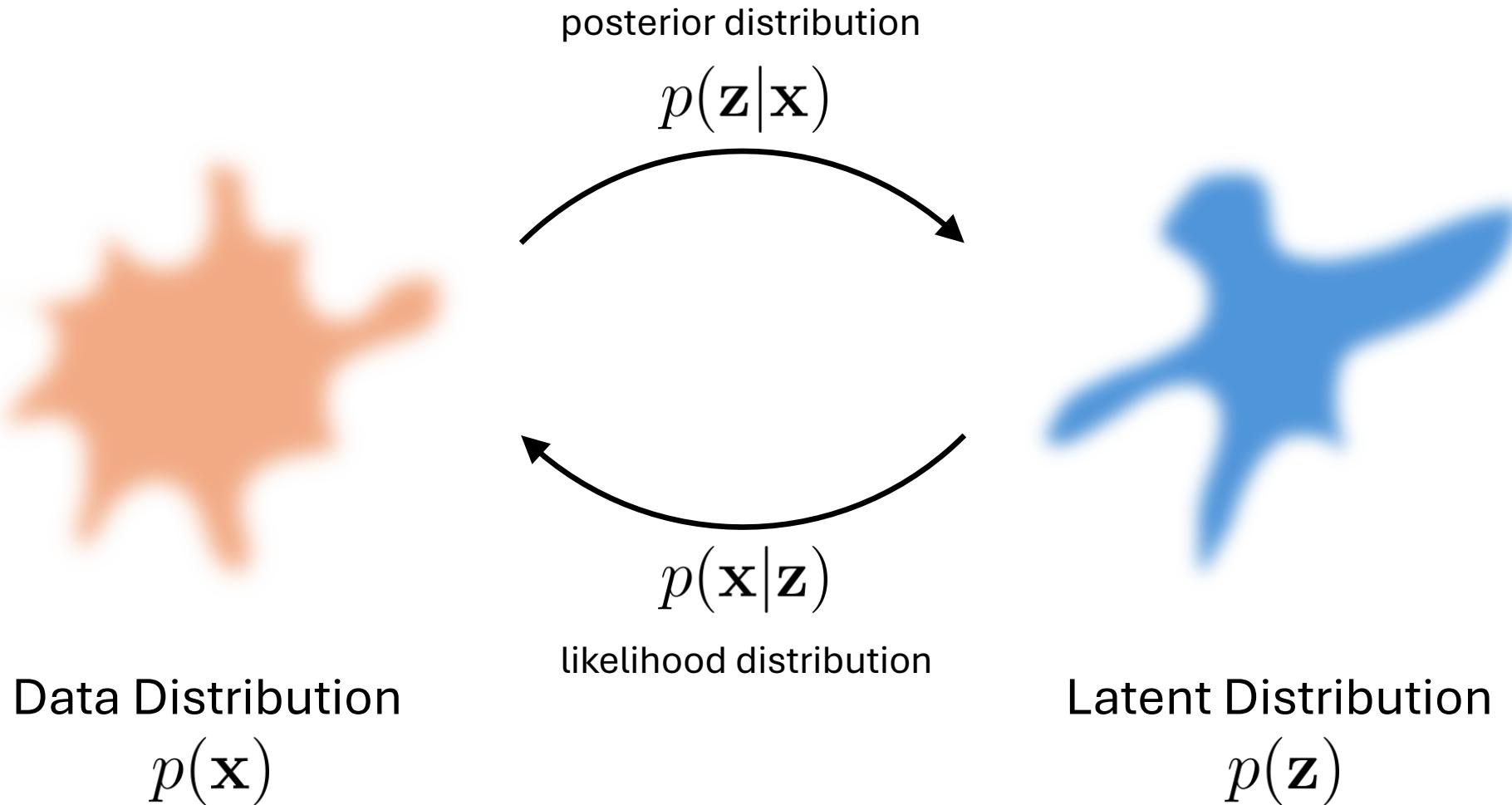


Latent Distribution  
 $p(\mathbf{z})$

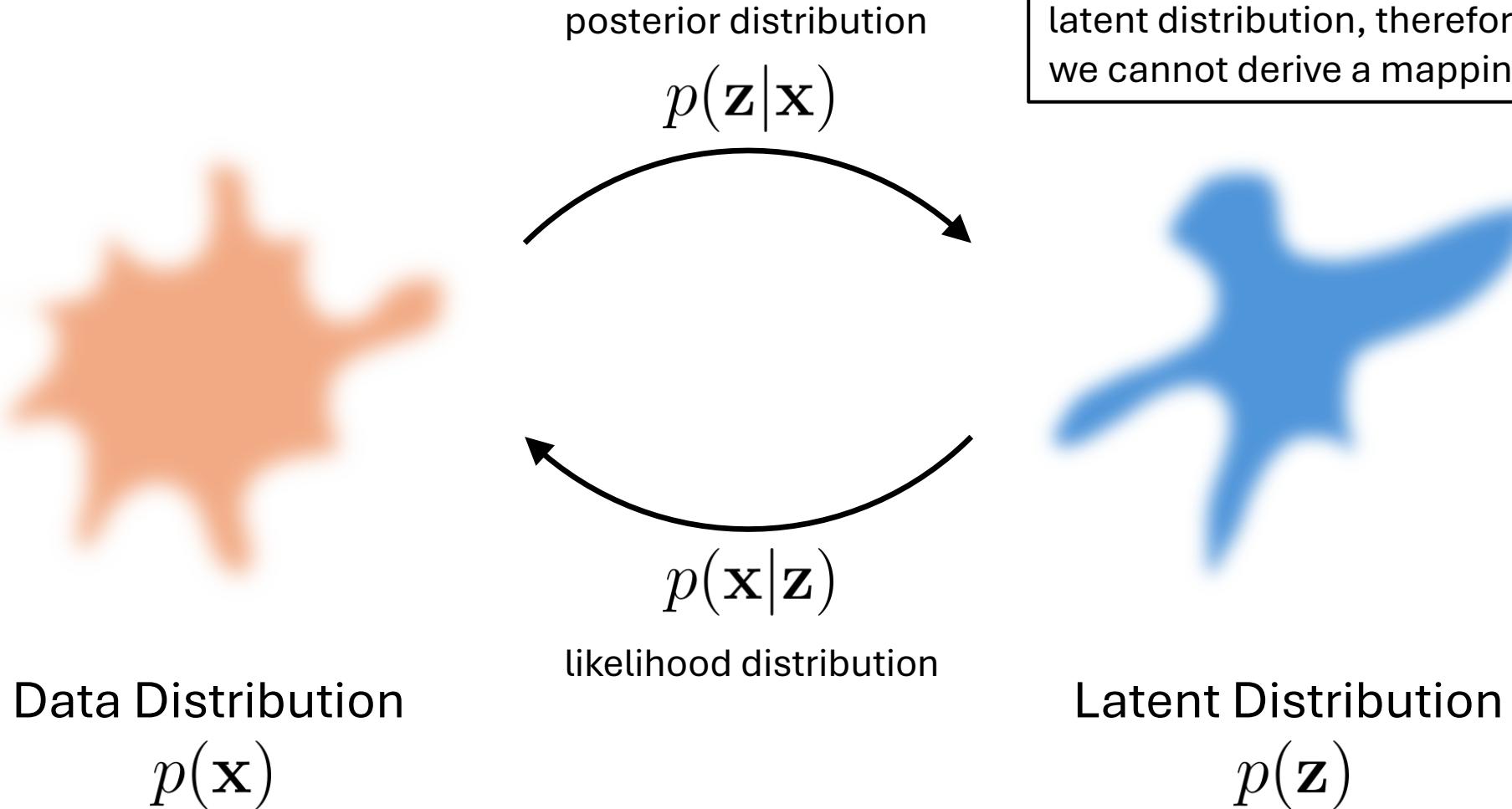
# Data and Latent Distributions



# Data and Latent Distributions



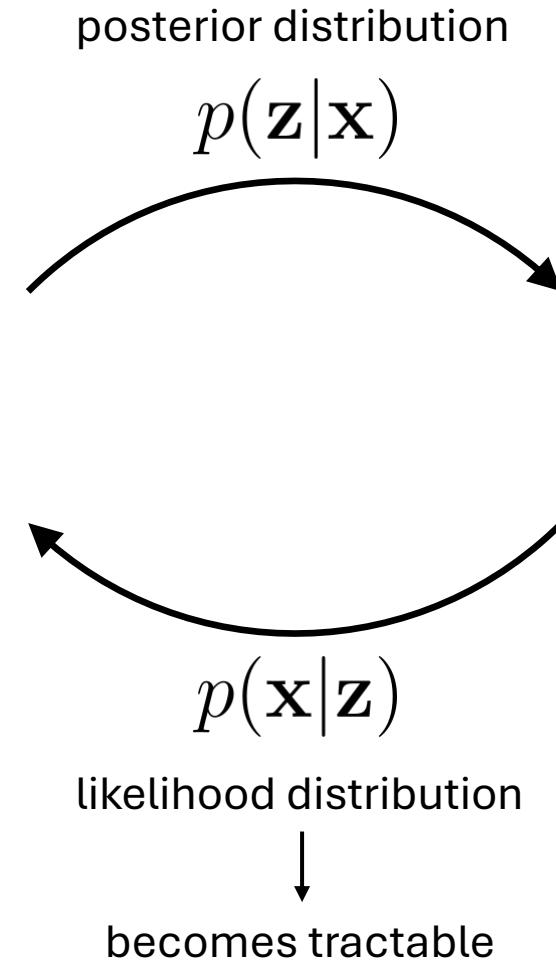
# Data and Latent Distributions



# Data and Latent Distributions

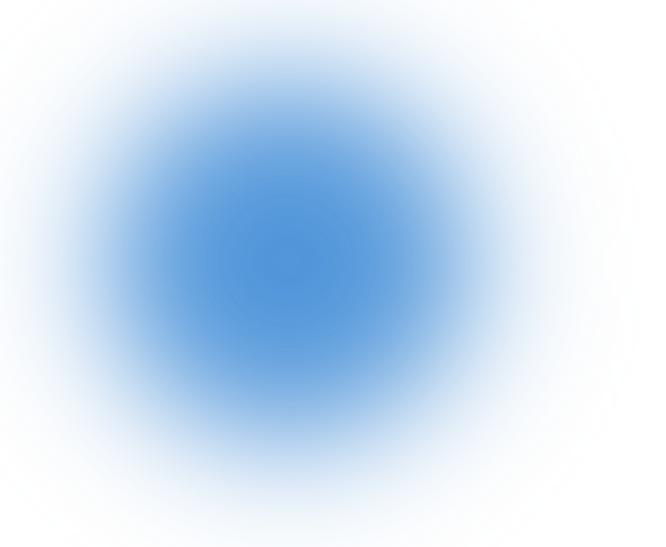


Data Distribution  
 $p(\mathbf{x})$



likelihood distribution  
↓  
becomes tractable

Idea: choose a simple latent distribution

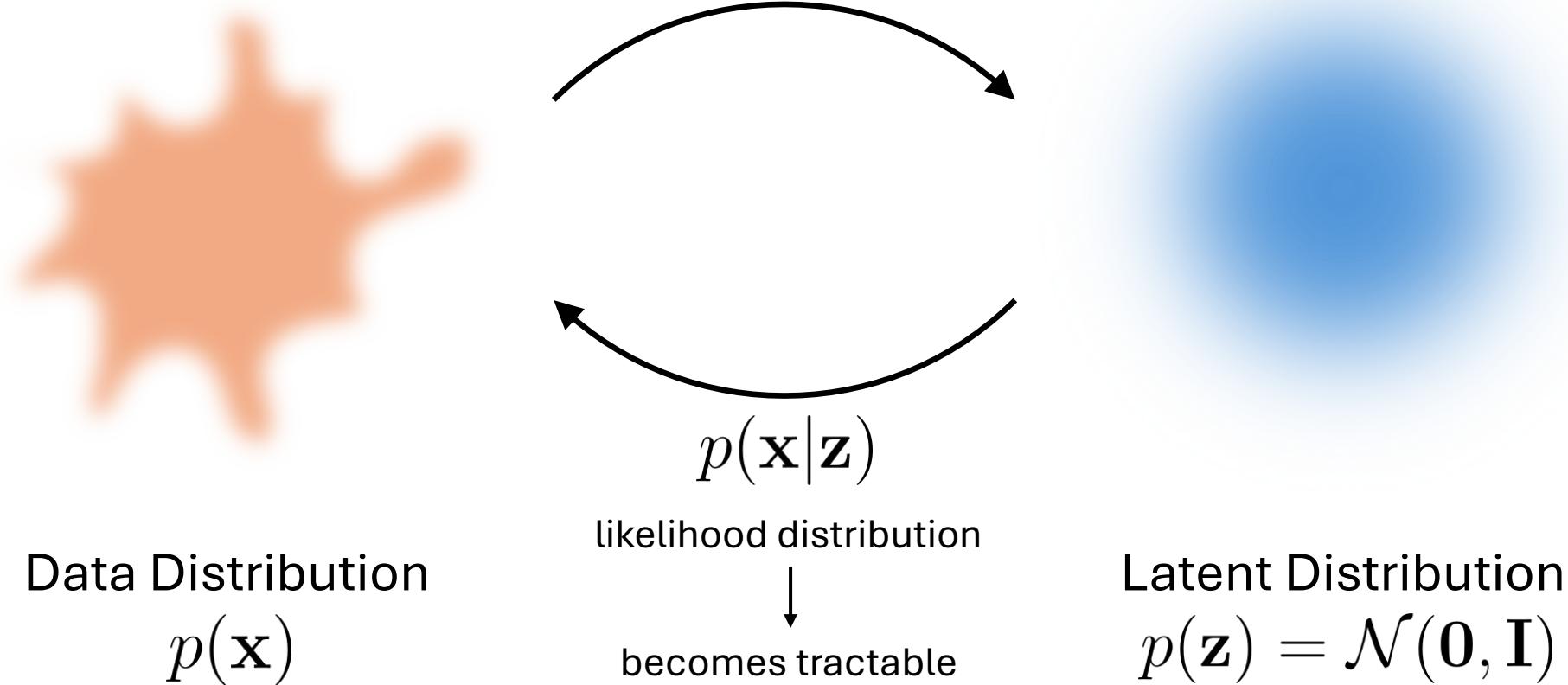


Latent Distribution  
 $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

# Variational Bayes

approximate posterior distribution with parameterized variational distribution

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$$



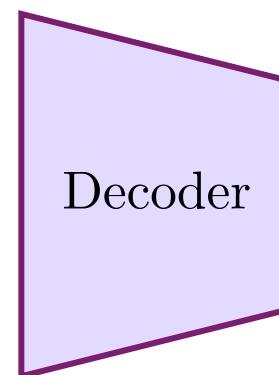
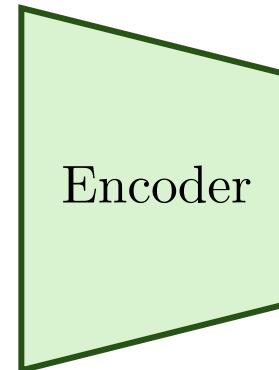
likelihood distribution  
↓  
becomes tractable

# Autoencoder Variational Bayes



Data Distribution  
 $p(\mathbf{x})$

$$q_{\phi}(\mathbf{z}|\mathbf{x})$$



$$p_{\theta}(\mathbf{x}|\mathbf{z})$$

Latent Distribution  
 $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

# How to train VAE?

Formally, we want to maximize the likelihood of data:  $p(\mathbf{x})$

Let's do it in a parameterized distributions family:  $p_{\theta}(\mathbf{x})$

# How to train VAE?

Formally, we want to maximize the likelihood of data:  $p(\mathbf{x})$

Let's do it in a parameterized distributions family:  $p_{\theta}(\mathbf{x})$

$$\log p_{\theta}(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (1) \text{ marginalize}$$

$$= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (2) \text{ multiply and divide}$$

$$= \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (3) \text{ expectation}$$

$$\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = \mathcal{L}_{\theta, \phi}(\mathbf{x}) \quad (4) \text{ Jensen's inequality}$$
$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$$

# ELBO optimization

We obtained lower bound for the data likelihood:

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (1) \text{ definition}$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log \left[ \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (2) \text{ Bayes' rule}$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log \frac{p_{\theta}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \quad (3) \text{ obvious}$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \text{KL} (q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z})) \quad (4) \text{ KL}$$

# VAE loss function

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$$

# VAE loss function

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Reconstruction Loss

We can model data likelihood with a Bernoulli distribution, then Reconstruction Loss becomes a simple MSE loss:

$$\mathcal{L} = \|\mathbf{x} - \mathbf{x}'\|_2^2$$

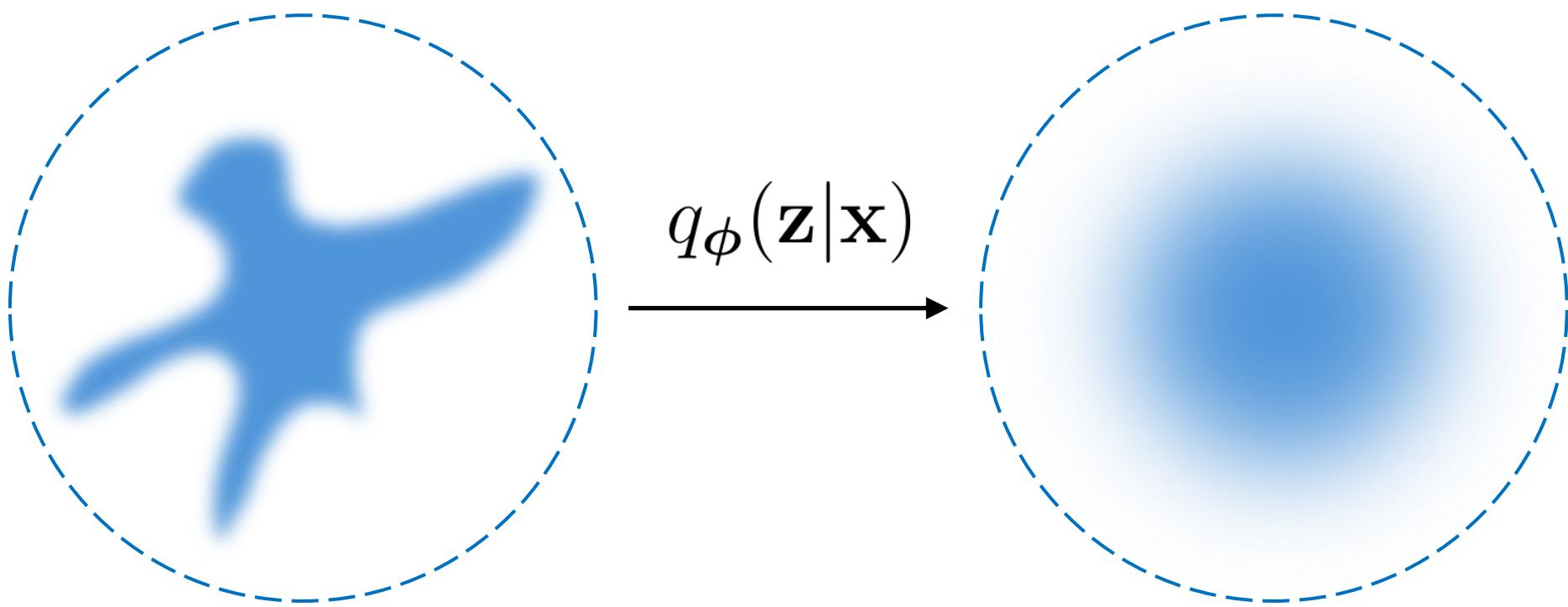


# VAE loss function

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$$

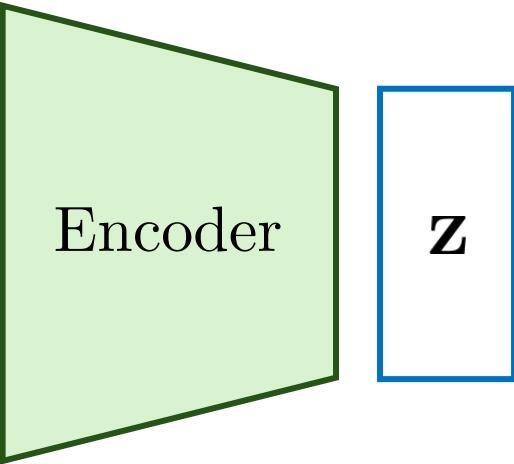
---

Regularization Loss

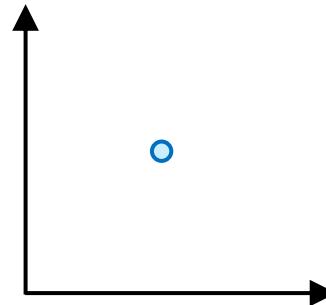


# Autoencoder vs VAE

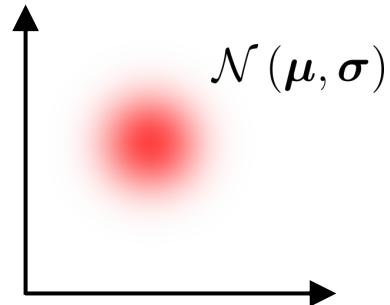
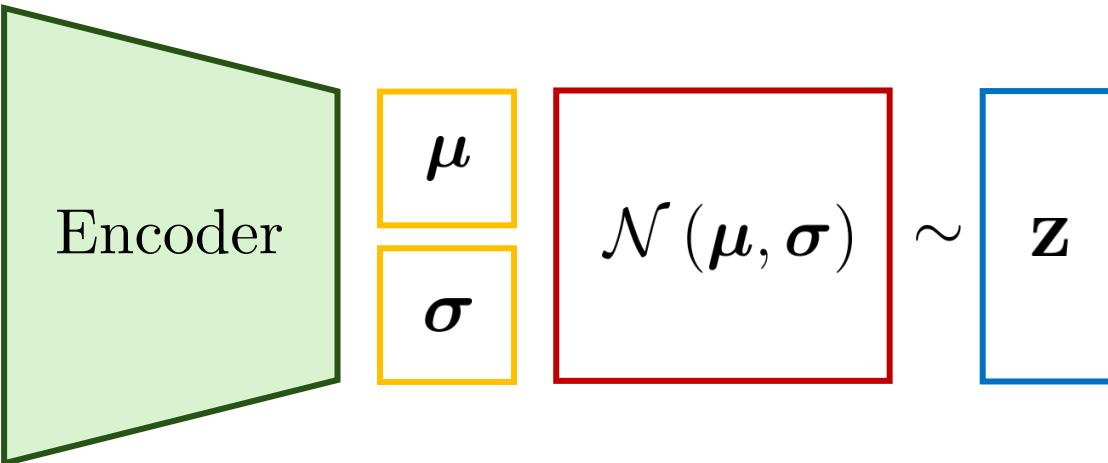
AE



Latent Space

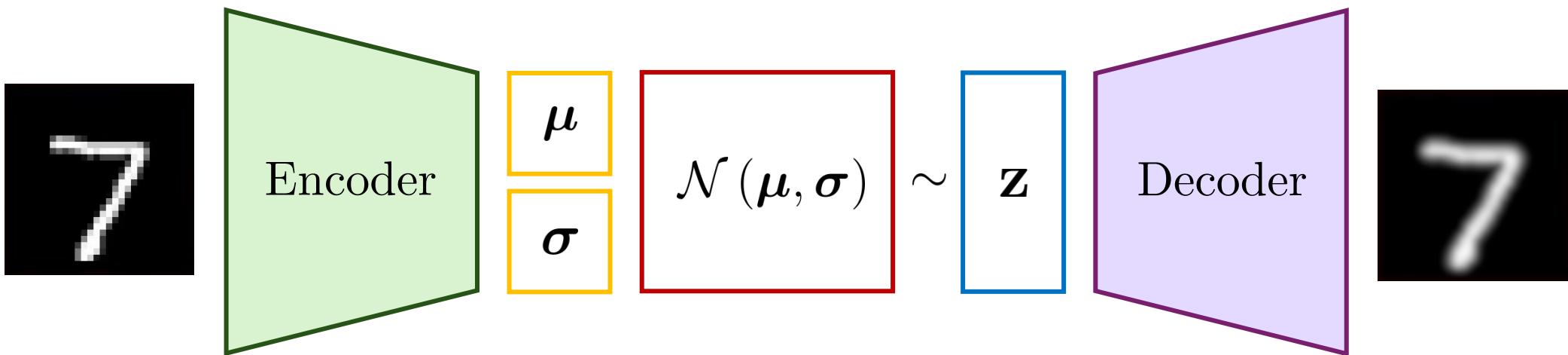


VAE



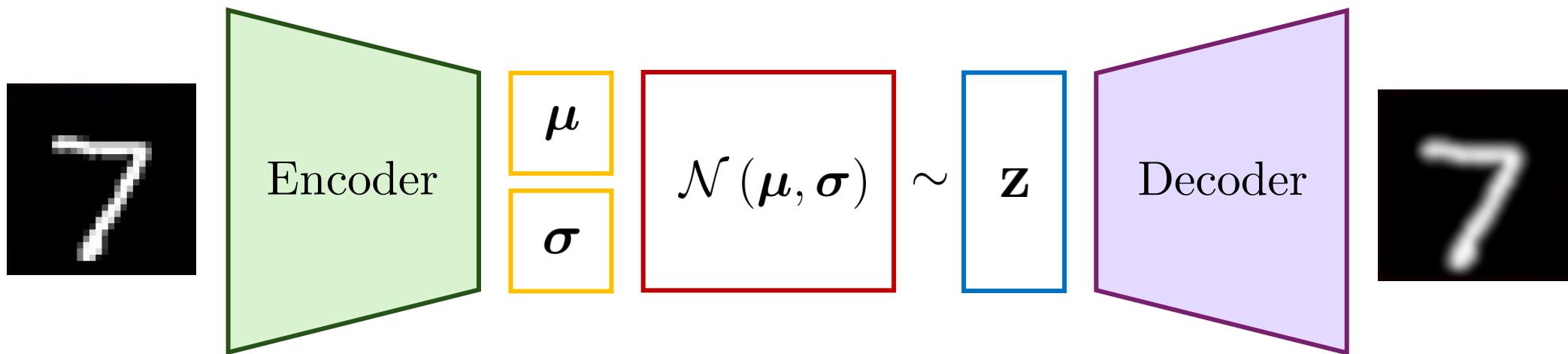
# How do we do in practice?

1. Put the image into Encoder
2. Encoder outputs mean and standard deviation  
+ Regularization Loss
3. Sample latent variable from the normal distribution
4. Put the latent vector into Decoder
5. Decoder outputs parameters of data distribution  
+ Reconstruction Loss



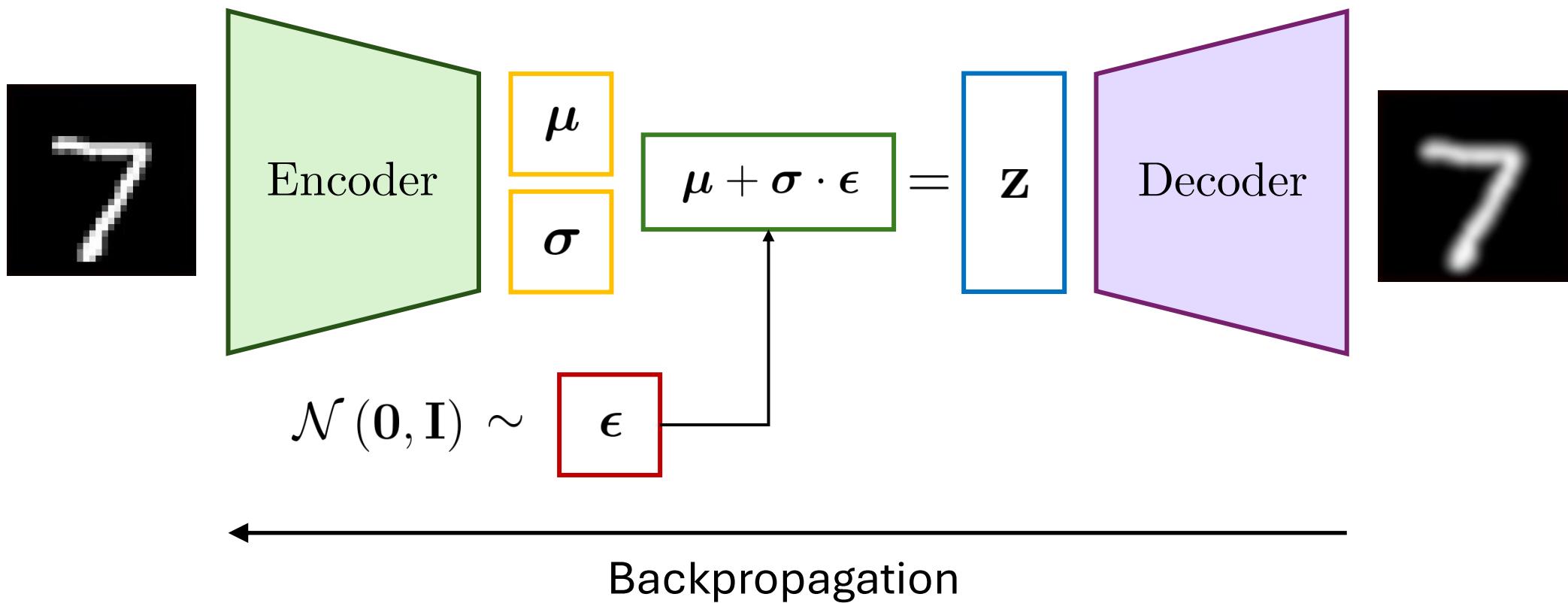
# Reparameterization Trick

How to backpropagate through the sampling process?



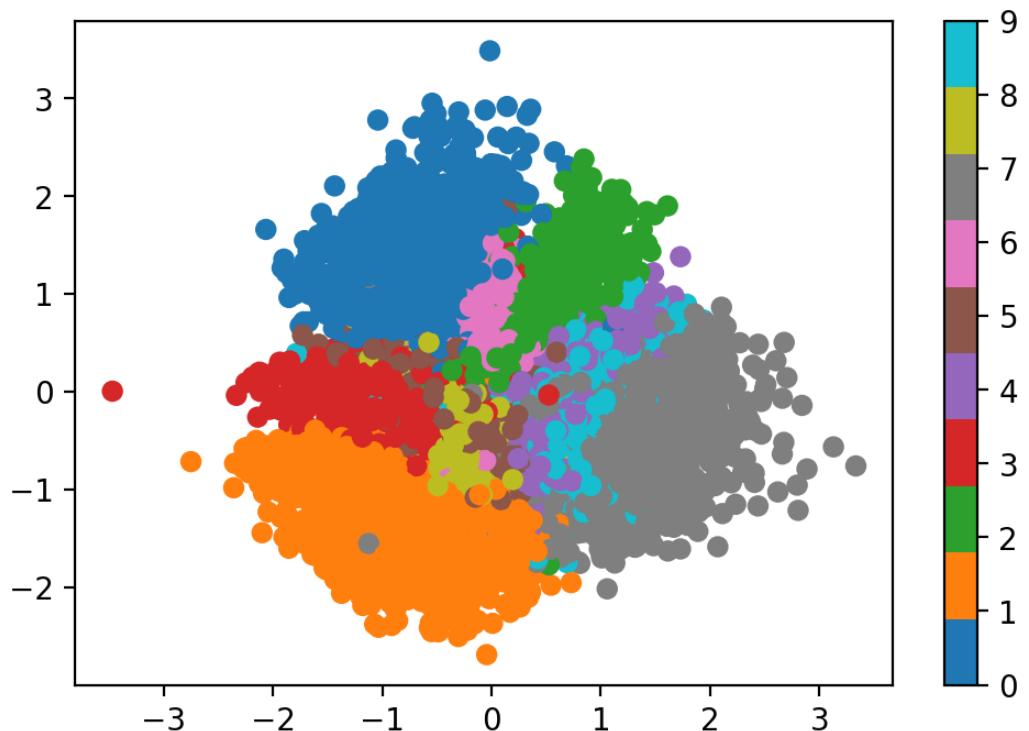
# Reparameterization Trick

Use property of normal distribution

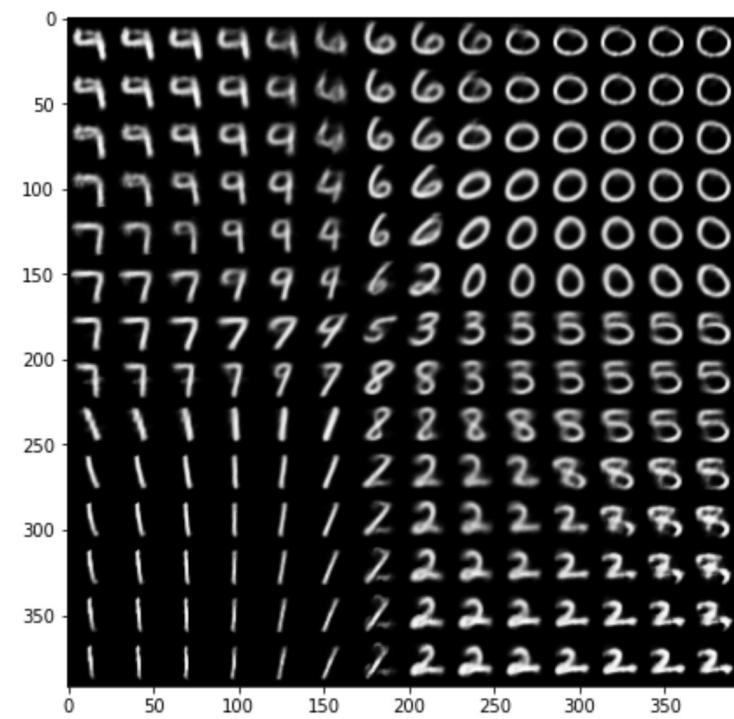


# Properties of the VAE's Latent Space

Organization

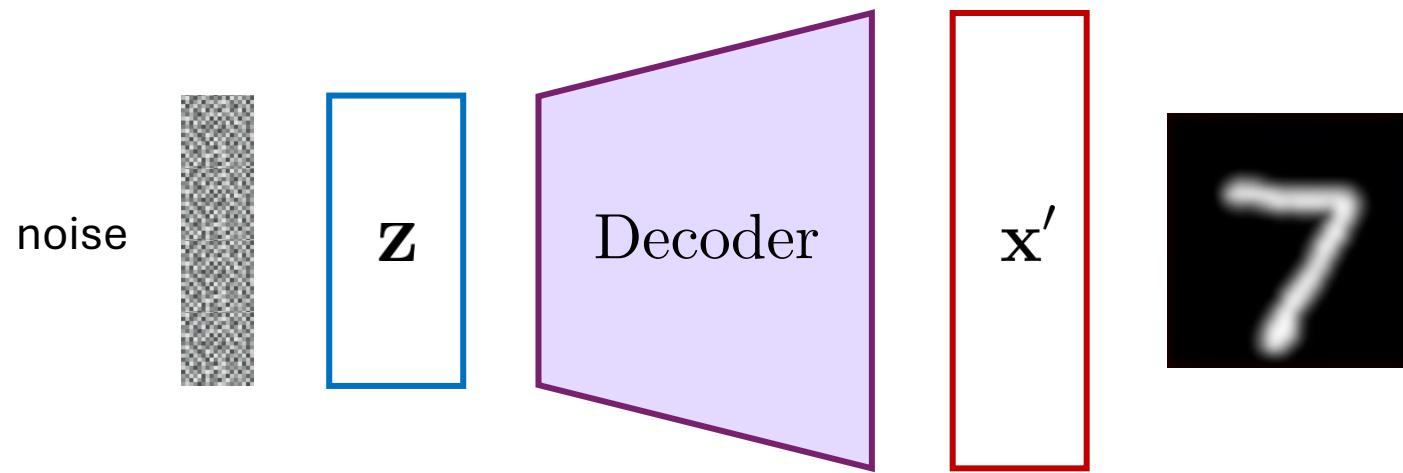


Interpolation



# VAE as a generative model

1. Sample noise latent  $\mathbf{z} \sim p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
2. Put  $\mathbf{z}$  into **Decoder** and get  $\mathbf{x}'$



# Pros and Cons

## Advantages

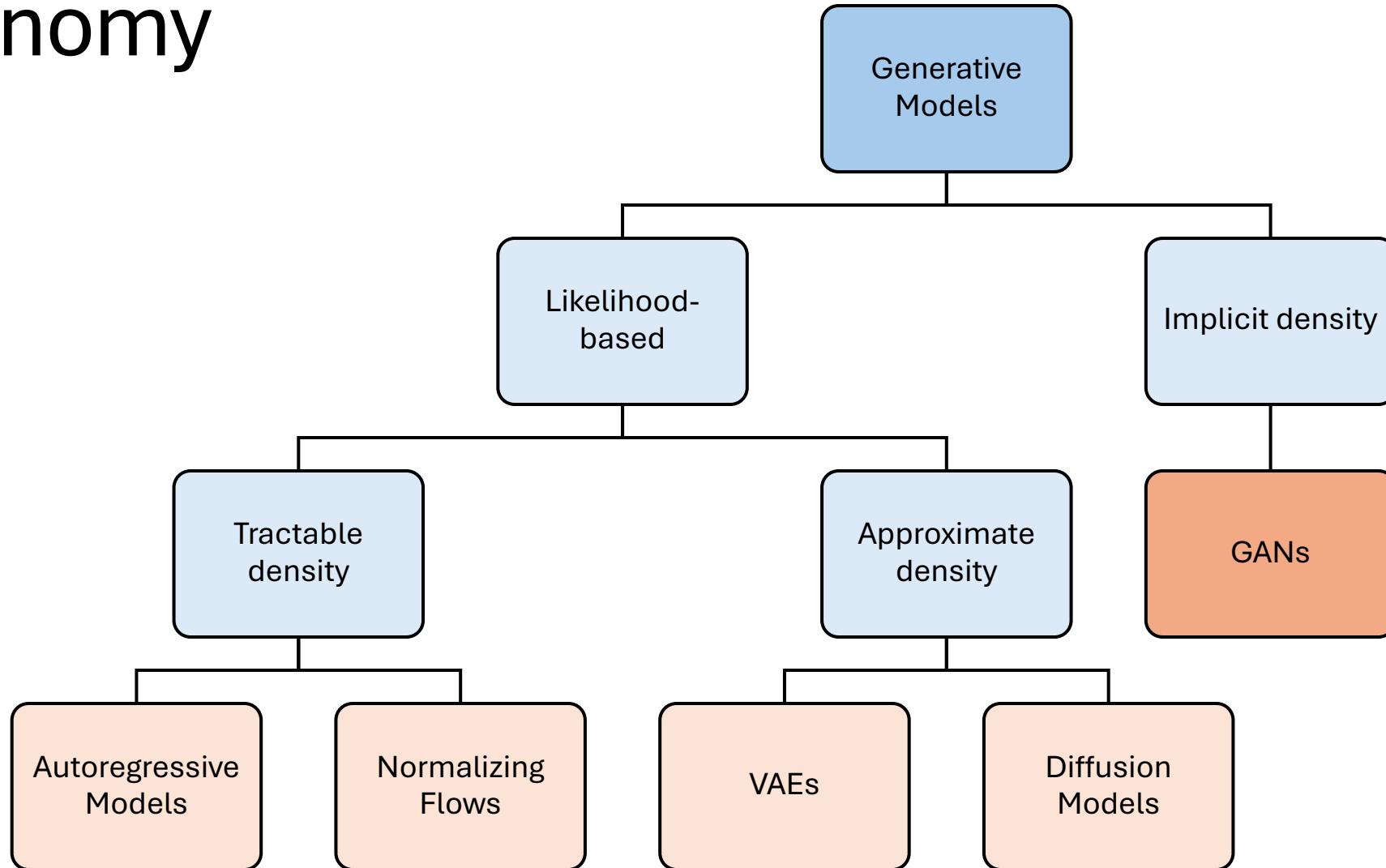
- Probabilistic, calibrated likelihoods (ELBO)
- Fast sampling (single pass), stable training
- Useful latent space for editing/interpolation

## Disadvantages

- Blur/averaging from pixel likelihoods, perceptual gap
- Posterior collapse with strong decoders (small or zero KL)
- Sensitive hyperparameters tuning

# GANs

# Taxonomy

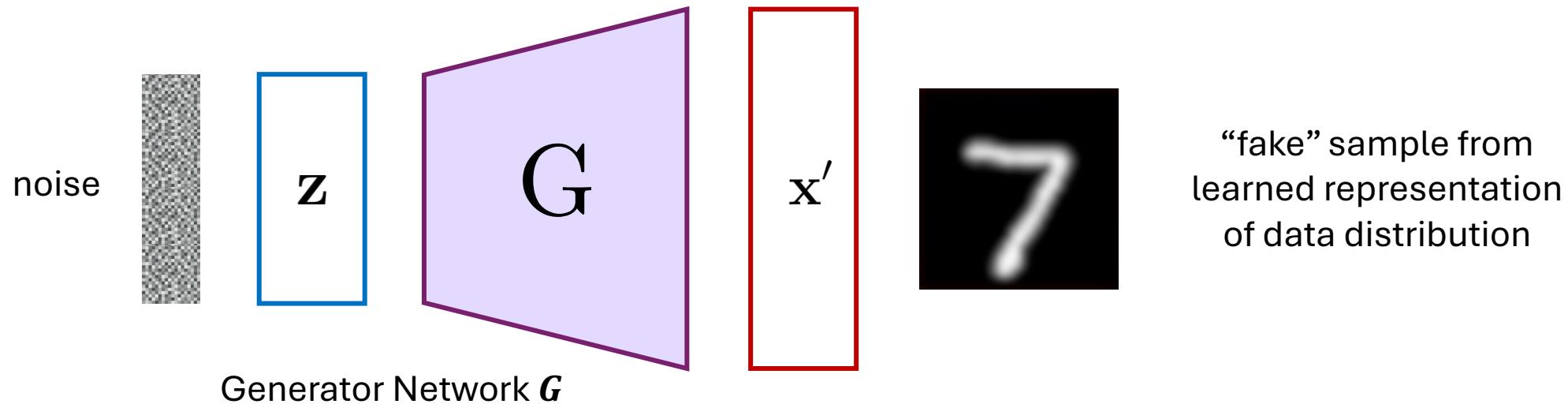


# What if we just want to sample?

**Idea:** don't explicitly model density, and instead just sample to generate new instances.

**Problem:** want to sample from complex distribution – can't do this directly!

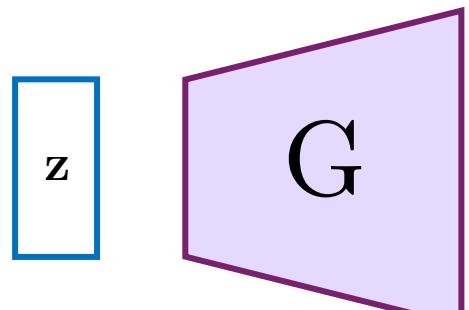
**Solution:** sample from something simple (e.g., noise), learn a transformation to the data distribution.



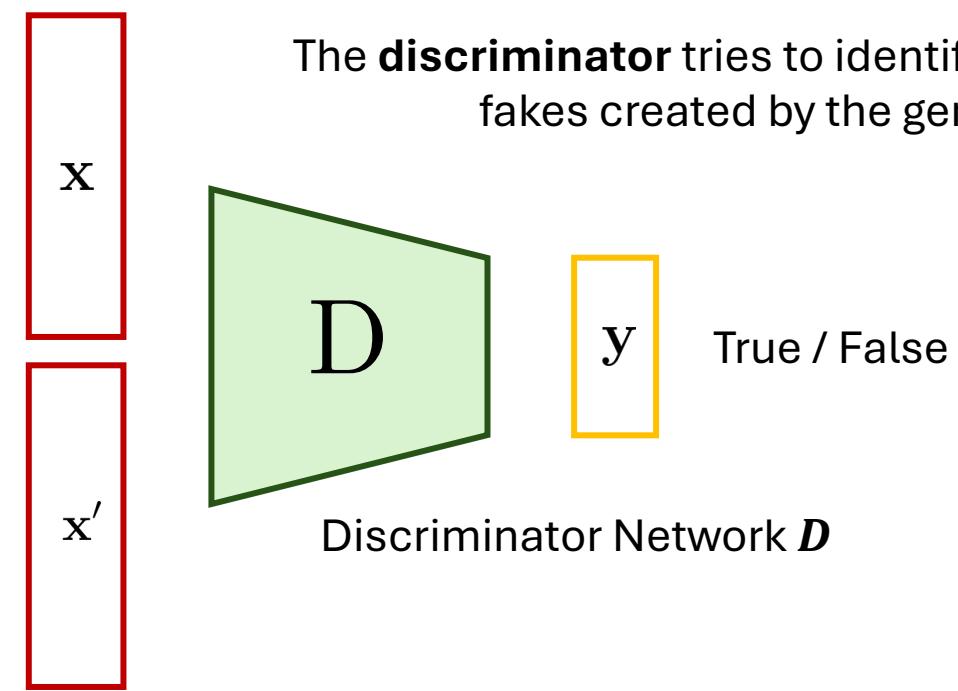
# Generative Adversarial Networks (GANs)

**Generative Adversarial Networks (GANs)** are a way to make a generative model by having two neural networks compete with each other.

The **generator** turns noise into an imitation of the data to try to trick the discriminator



Generator Network  $G$



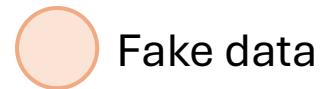
The **discriminator** tries to identify real data from fakes created by the generator

Discriminator Network  $D$

# Intuition behind GANs

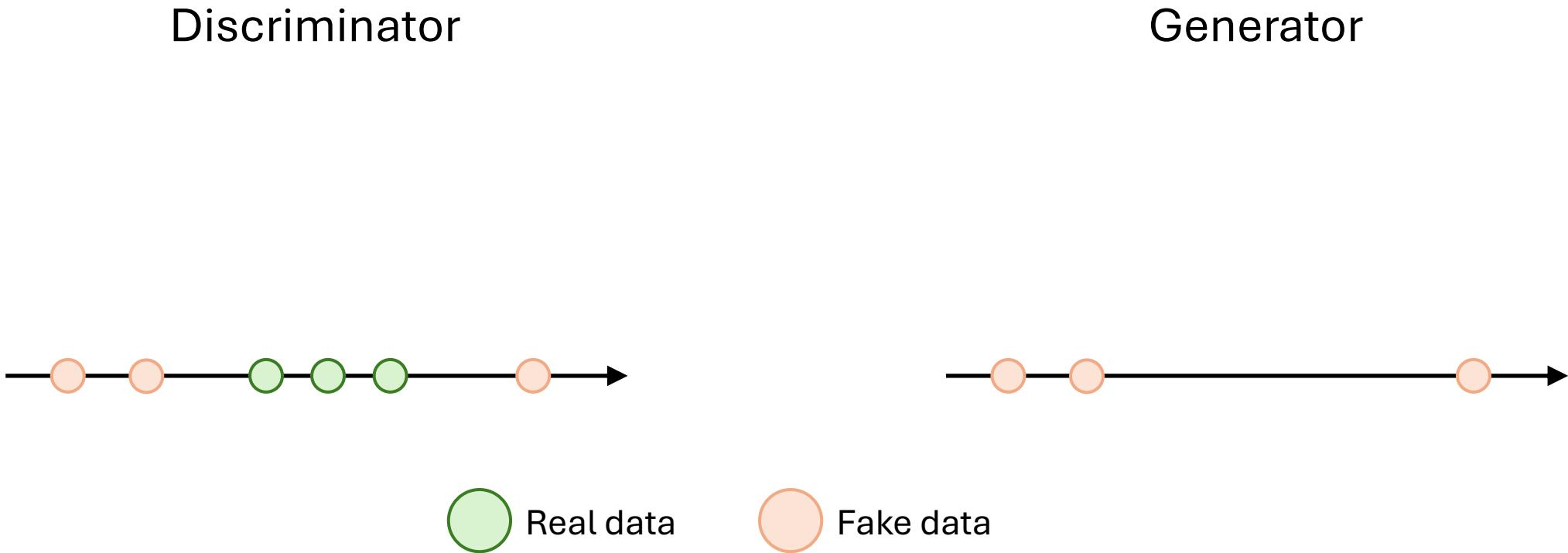
**Generator** starts from noise to try to create an imitation of the data.

Generator



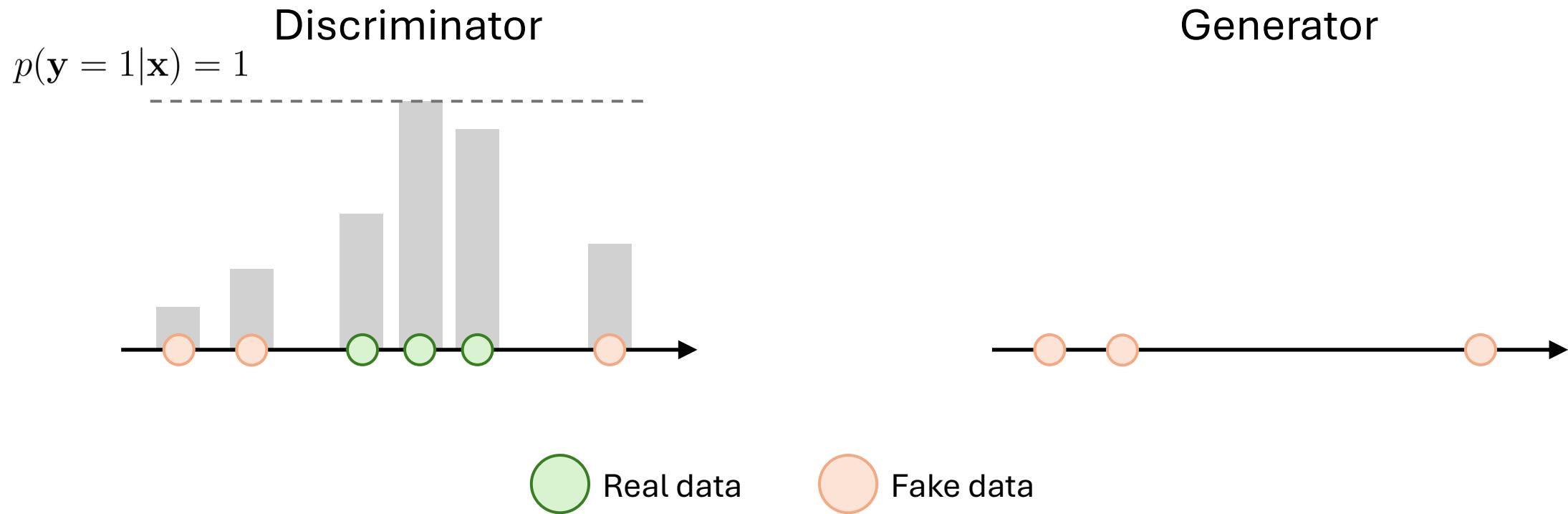
# Intuition behind GANs

**Discriminator** looks at both real data and fake data created by the generator.



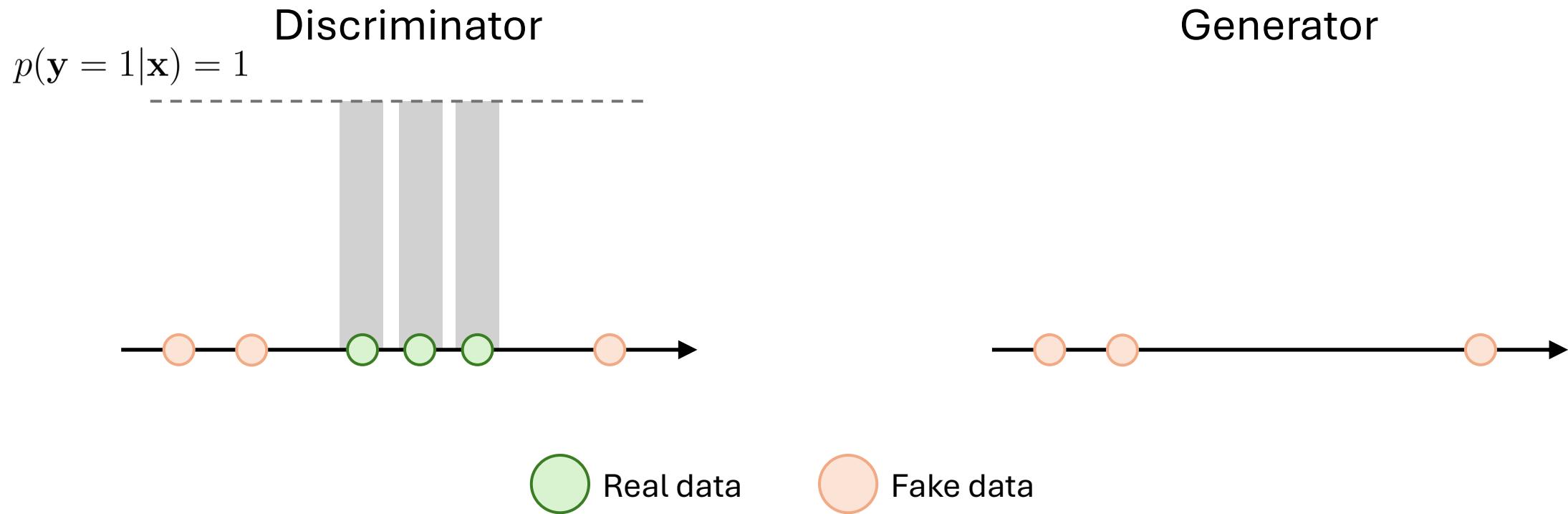
# Intuition behind GANs

**Discriminator** tries to predict what's real and what's fake.



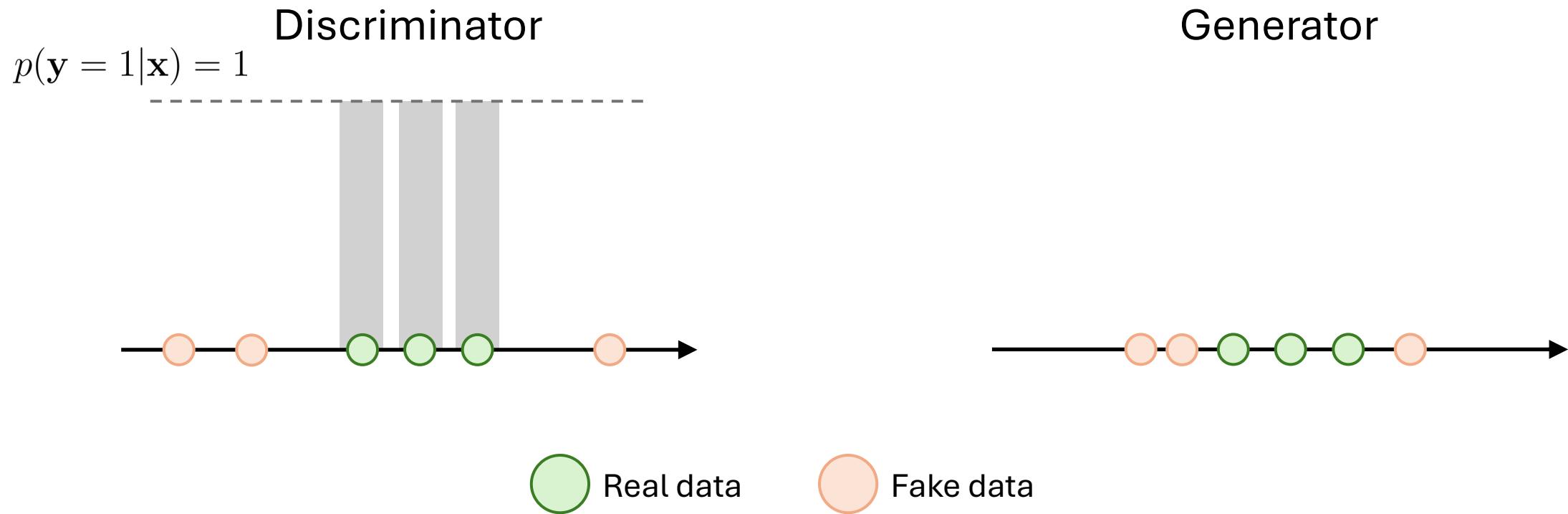
# Intuition behind GANs

**Discriminator** tries to predict what's real and what's fake.



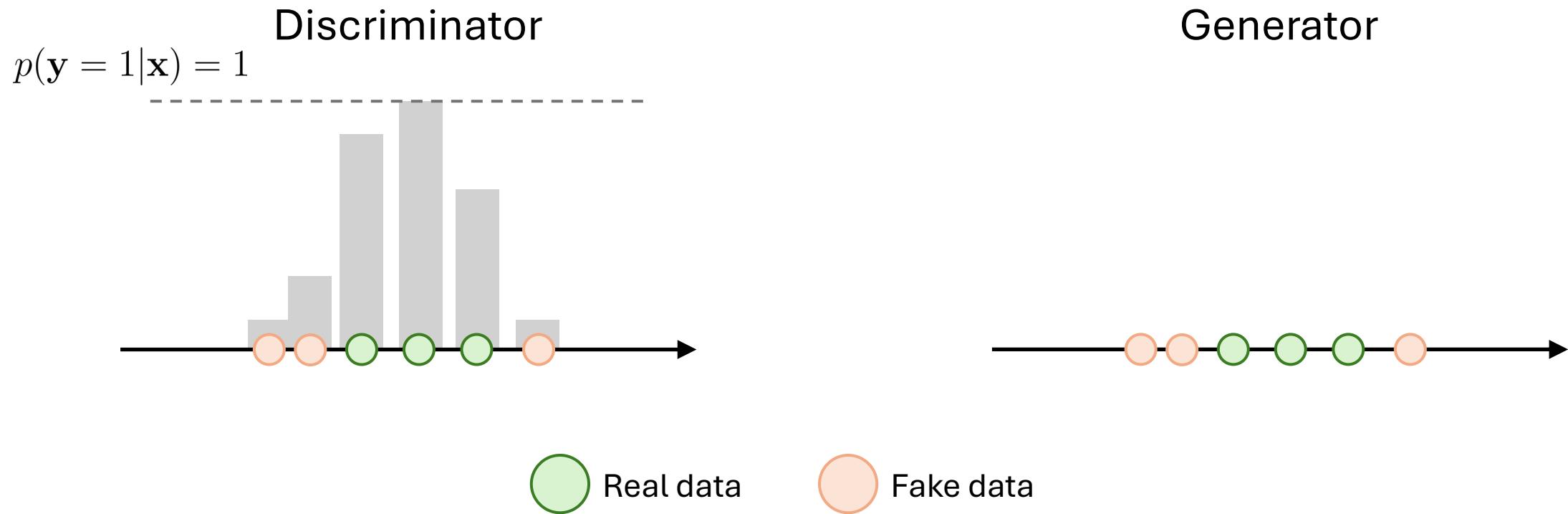
# Intuition behind GANs

**Generator** tries to improve its imitation of data.



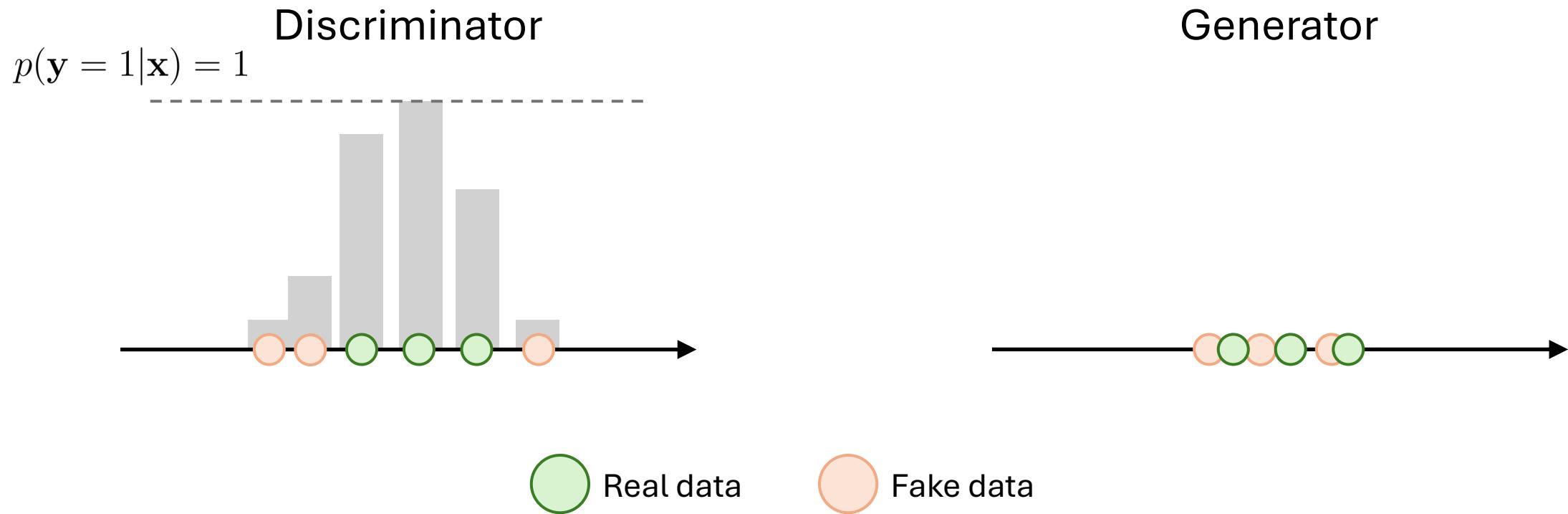
# Intuition behind GANs

**Discriminator** tries to predict what's real and what's fake.



# Intuition behind GANs

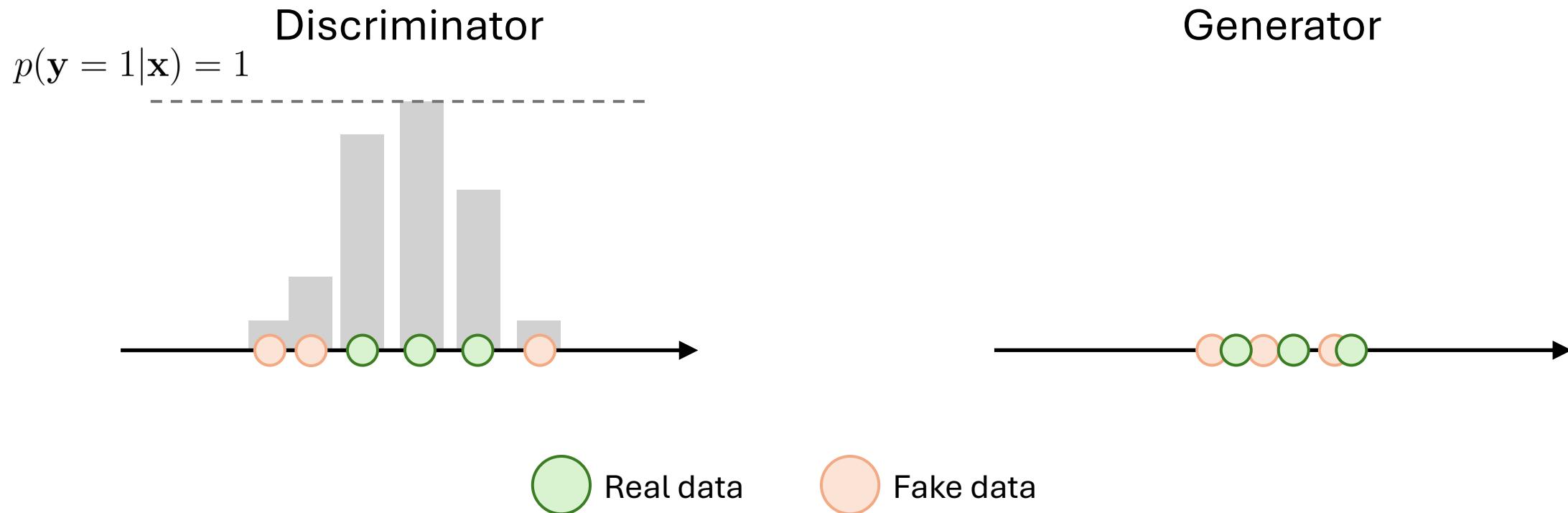
**Generator** tries to improve its imitation of data.



# Intuition behind GANs

**Discriminator** tries to identify real data from fakes created by the generator.

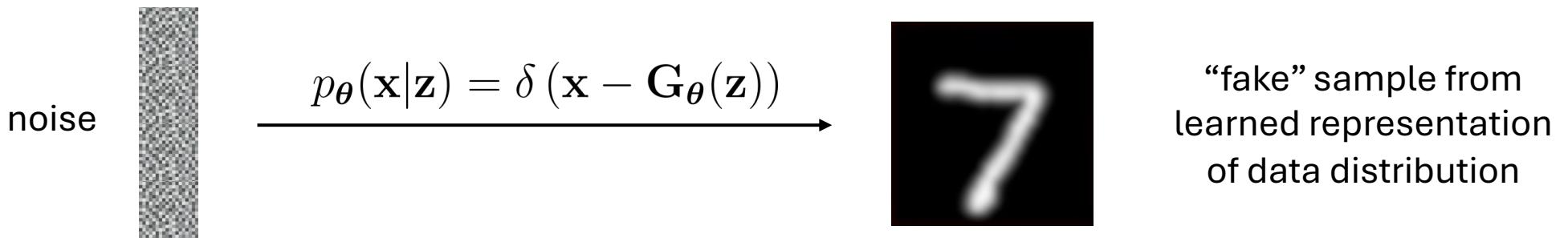
**Generator** tries to create imitations of data to trick the discriminator.



# Maths behind GANs

Assume generative model with the base latent distribution and deterministic map:

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$



# Maths behind GANs

**Discriminator:** a classifier, which distinguish real samples from data distribution  $\pi(\mathbf{x})$  and generated samples from  $p_{\theta}(\mathbf{x})$ :

$$D_{\phi}(\mathbf{x}) = p_{\phi}(y = 1 | \mathbf{x}) \in [0, 1]$$

Discriminator tries to **minimize** cross entropy.

$$\max_D [\mathbb{E}_{\pi(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z})))]$$

# Maths behind GANs

**Generator:** generative model, which makes the generated sample more realistic:

$$\mathbf{x} = G_{\theta}(\mathbf{z}), \quad \mathbf{z} \sim p(\mathbf{z})$$

Generator tries to **maximize** cross entropy.

$$\min_G \left[ \mathbb{E}_{\pi(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}))) \right]$$

# GAN objective

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

# What are optimal $G$ and $D$ ?

## Assumption

Generative distribution  $p_{\theta}(\mathbf{x})$  equals to the true distribution  $\pi(\mathbf{x})$  if we can not distinguish them using discriminative model  $D$ . It means that  $D(\mathbf{x}) = 0.5$  for each sample  $\mathbf{x}$ .

## Theorem

The minimax game  $\min_G \max_D \left[ \underbrace{\mathbb{E}_{\pi(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z})))}_{V(G,D)} \right]$

has the global optimum  $\pi(\mathbf{x}) = p_{\theta}(\mathbf{x})$ , in this case  $D^*(\mathbf{x}) = 0.5$ .

# Proof

$$\begin{aligned} V(G, D) &= \mathbb{E}_{\pi(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \log (1 - D(\mathbf{x})) \\ &= \int \underbrace{[\pi(\mathbf{x}) \log D(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x}) \log (1 - D(\mathbf{x}))]}_{y(D)} d\mathbf{x} \end{aligned}$$

$$\frac{dy(D)}{dD} = \frac{\pi(\mathbf{x})}{D(\mathbf{x})} - \frac{p_{\boldsymbol{\theta}}(\mathbf{x})}{1 - D(\mathbf{x})} = 0$$

$$D^*(\mathbf{x}) = \frac{\pi(\mathbf{x})}{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})}$$

# Proof continued (fixed $D = D^*$ )

$$\begin{aligned} V(G, D^*) &= \mathbb{E}_{\pi(\mathbf{x})} \log \left( \frac{\pi(\mathbf{x})}{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})} \right) + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \log \left( \frac{p_{\boldsymbol{\theta}}(\mathbf{x})}{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})} \right) \\ &= \text{KL} \left( \pi(\mathbf{x}) \parallel \frac{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})}{2} \right) + \text{KL} \left( p_{\boldsymbol{\theta}}(\mathbf{x}) \parallel \frac{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})}{2} \right) - 2 \log 2 \\ &= 2 \text{JSD} (\pi(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x})) - 2 \log 2 \end{aligned}$$

Jensen-Shannon divergence (symmetric KL divergence)

$$\text{JSD} (\pi(\mathbf{x}) \| p_{\boldsymbol{\theta}}(\mathbf{x})) = \frac{1}{2} \left[ \text{KL} \left( \pi(\mathbf{x}) \parallel \frac{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})}{2} \right) + \text{KL} \left( p_{\boldsymbol{\theta}}(\mathbf{x}) \parallel \frac{\pi(\mathbf{x}) + p_{\boldsymbol{\theta}}(\mathbf{x})}{2} \right) \right]$$

Could be used as a distance measure!

$$V(G^*, D^*) = -2 \log 2, \quad \pi(\mathbf{x}) = p_{\boldsymbol{\theta}}(\mathbf{x}), \quad D^*(\mathbf{x}) = 0.5$$

# GAN optimality

## Expectations

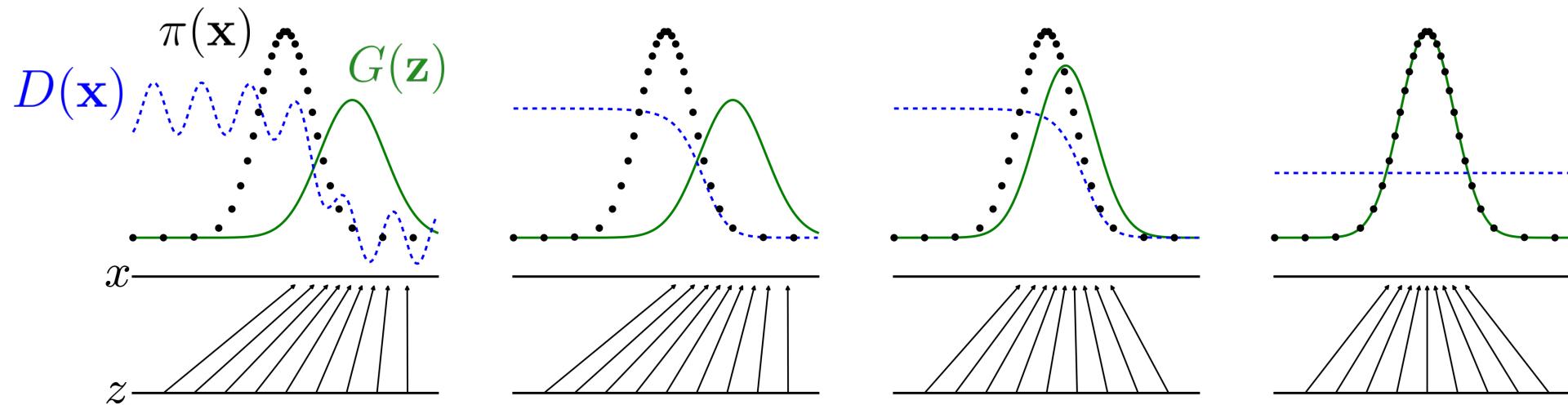
If the generator could be **any** function and the discriminator is **optimal** at every step, then the generator is **guaranteed to converge** to the data distribution.

## Reality

- Generator updates are made in parameter space, discriminator is not optimal at every step
- Generator and discriminator loss keeps oscillating during GAN training.

# GAN training

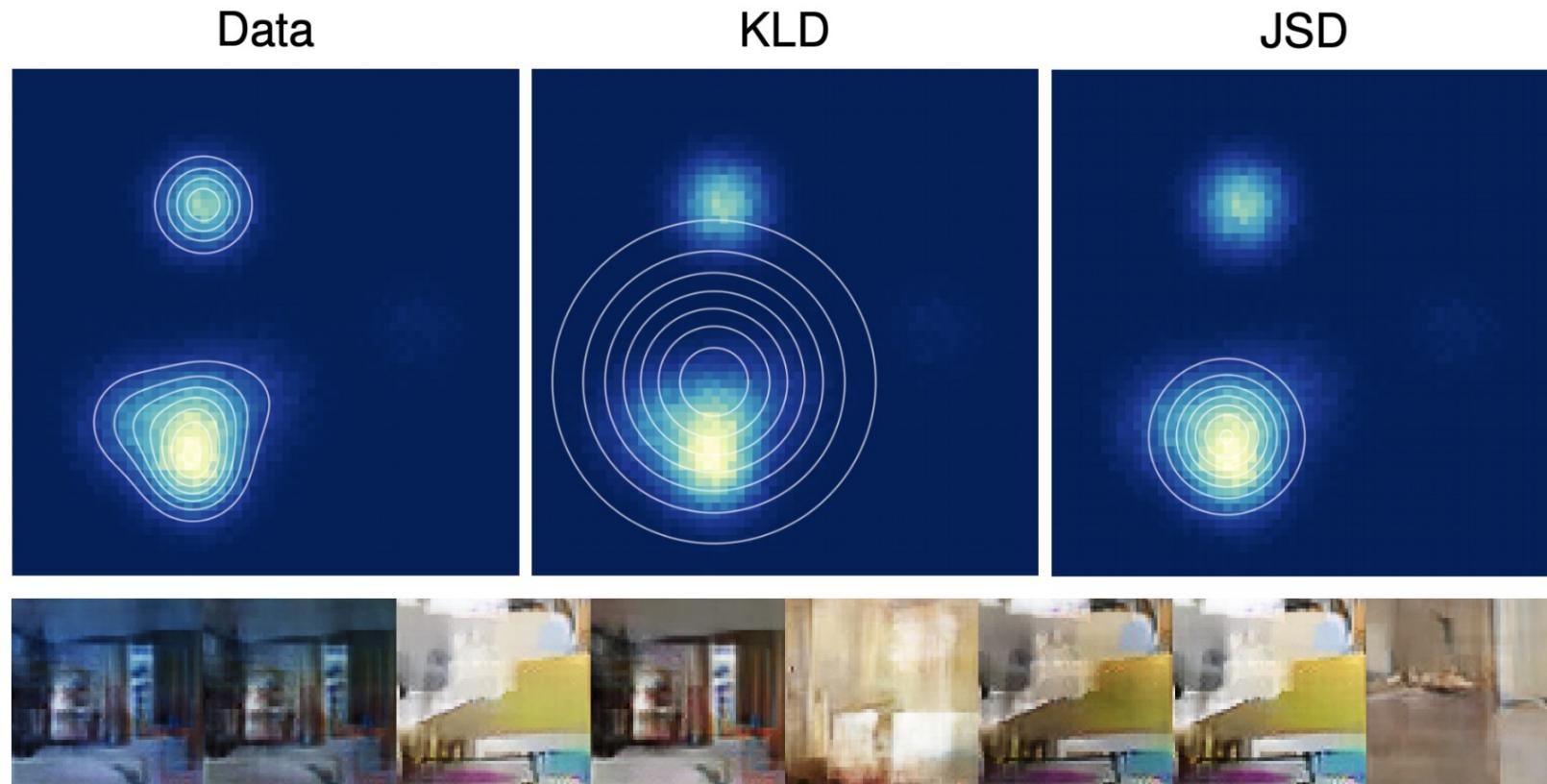
$$\min_G \max_D \left[ \mathbb{E}_{\pi(x)} \log D(x) + \mathbb{E}_{p(z)} \log(1 - D(G(z))) \right]$$



- $\mathbf{z} \sim p(\mathbf{z})$  is a latent variable.
- $p_{\theta}(\mathbf{x}|\mathbf{z}) = \delta(\mathbf{x} - G(\mathbf{z}))$  is deterministic decoder.
- We do not have encoder at all.

# Mode collapse

The phenomena where the generator of a GAN collapses to one or few distribution modes.



# How do we do in practice?

---

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

---

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by descending its stochastic gradient:

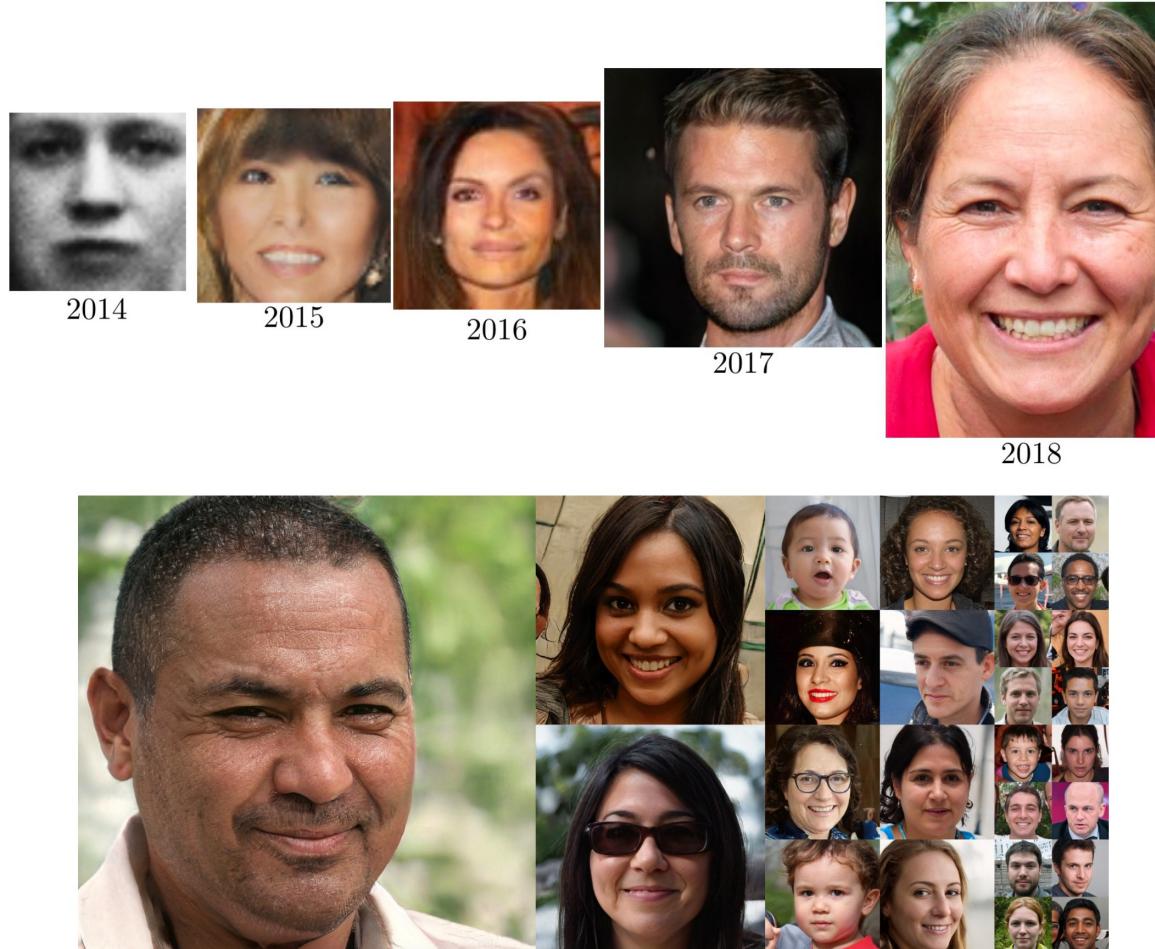
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

---

# StyleGAN – high quality faces



# Pros and Cons

## **Advantages**

- Capable of producing highly realistic synthetic data
- Flexible framework adaptable to images, text, audio and more
- Encouraged rapid progress in generative modeling research

## **Disadvantages**

- Training is unstable and prone to mode collapse
- Difficult to evaluate model performance quantitatively
- Require large datasets and computational resources

# Recap

- What are Generative Models?
- Taxonomy
- Autoregressive Models
- Autoencoders
- VAEs
- GANs
- \* Extra: Modern Architectures

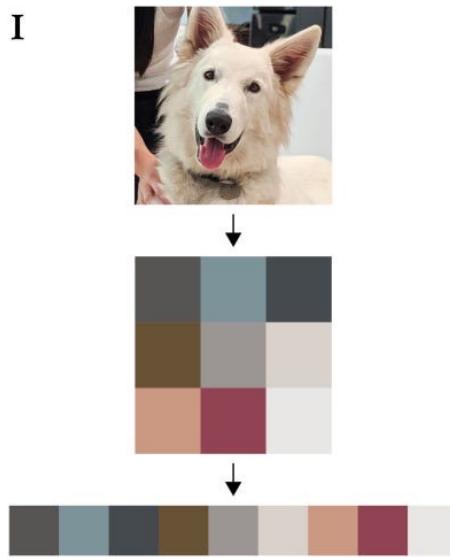


**TIME FOR A BREAK**

# Extra: Modern Architectures

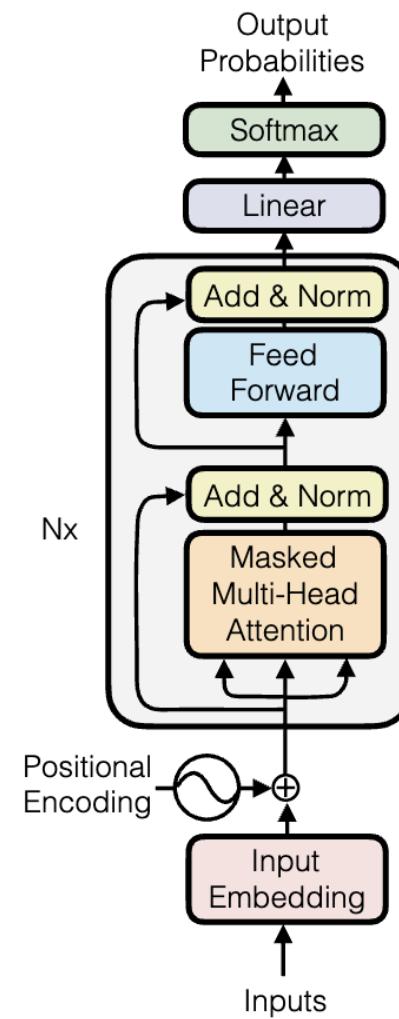
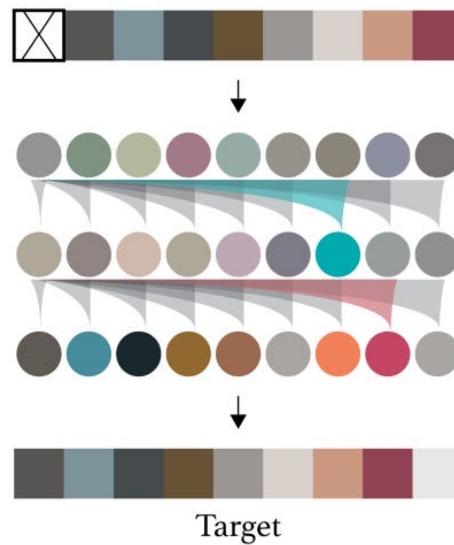
# Autoregression

# Image GPT

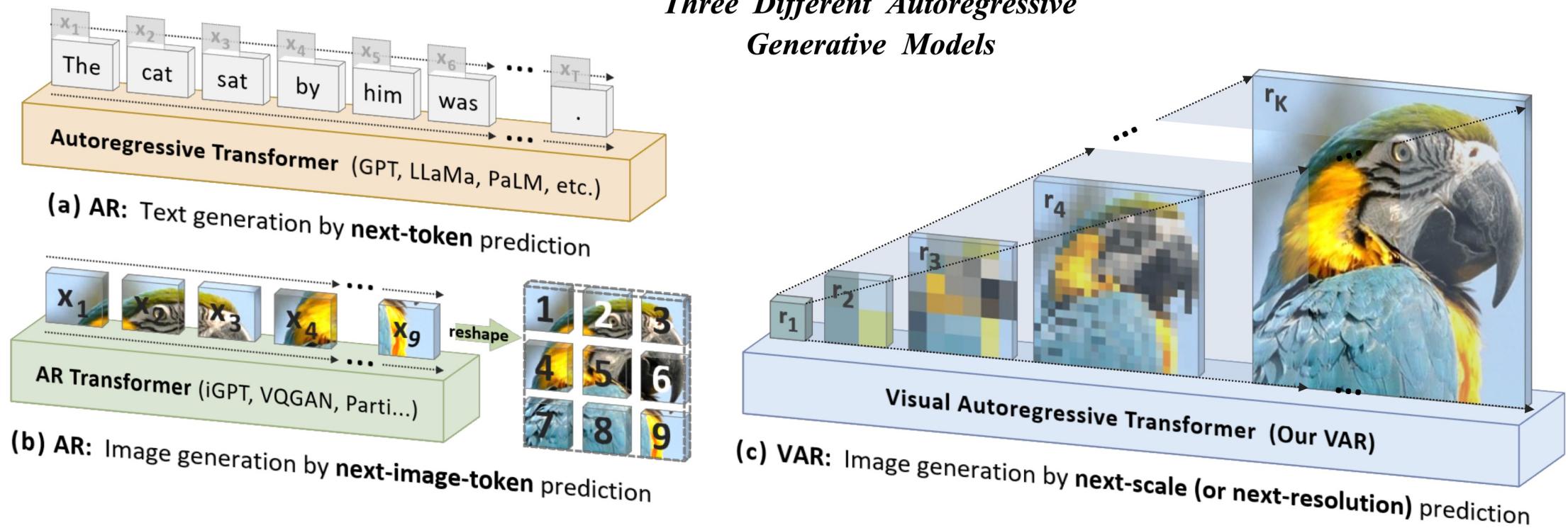


2

(a) Autoregressive

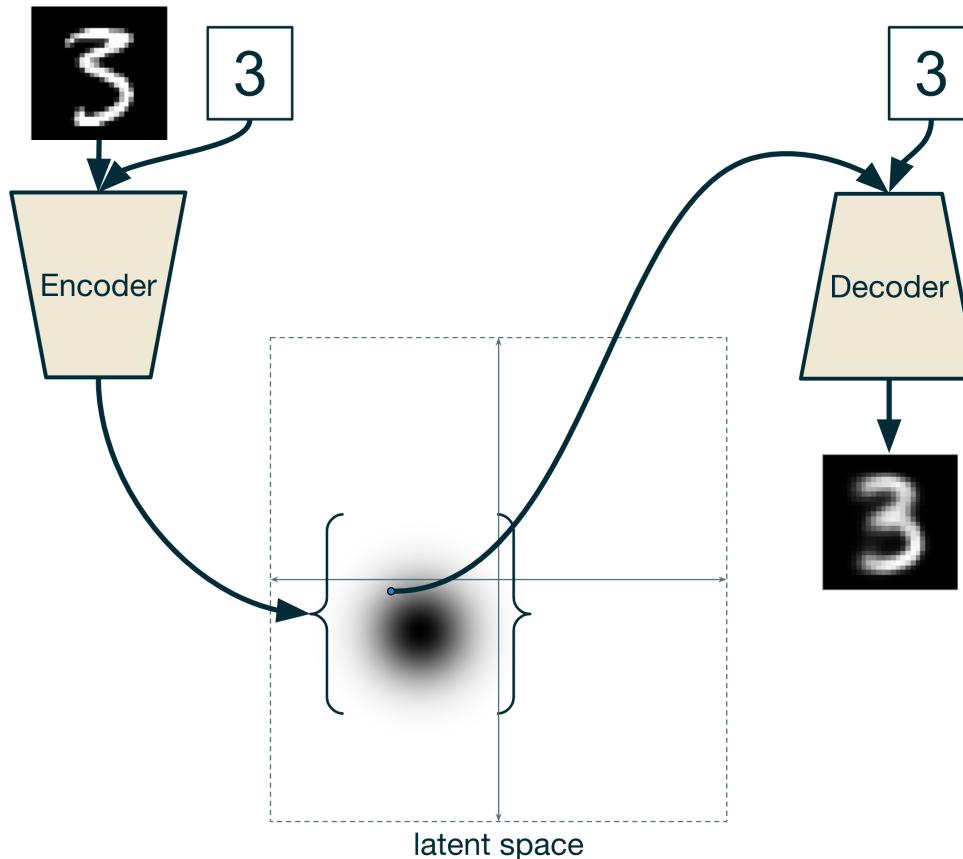


# VAR: Visual Autoregressive Modeling



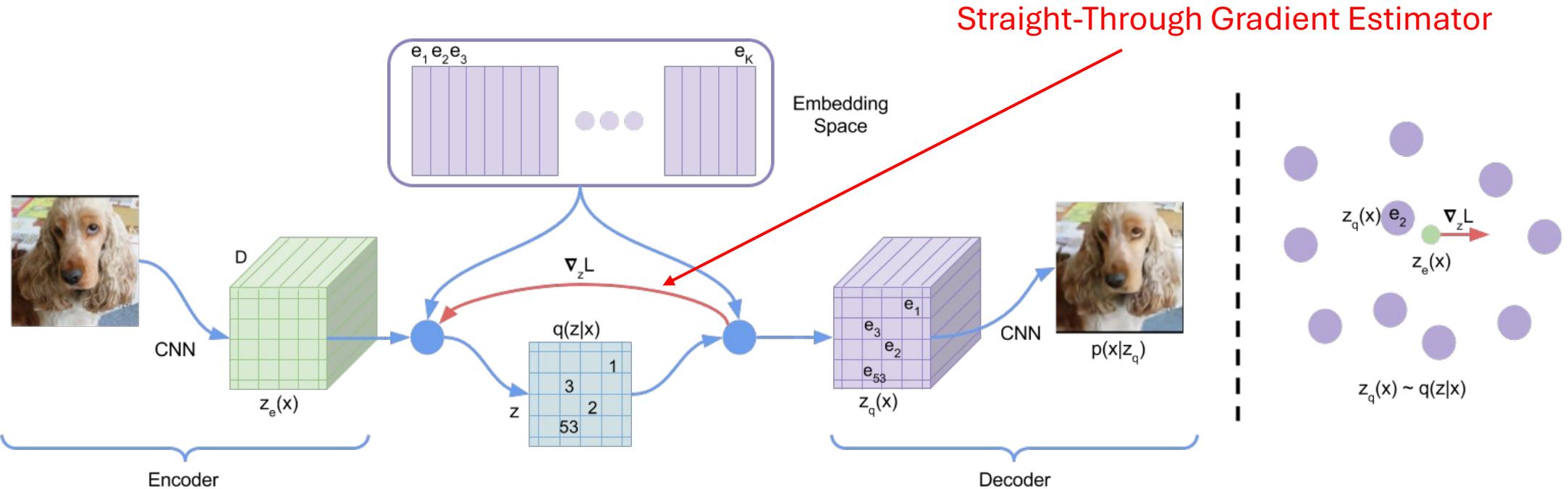
VAE

# cVAE: Conditional VAE



$$\mathcal{L}_{\text{cVAE}}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} \log p_{\theta}(\mathbf{x}|\mathbf{z}, \mathbf{y}) - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})||p(\mathbf{z}))$$

# VQ-VAE: Vector Quantized VAE



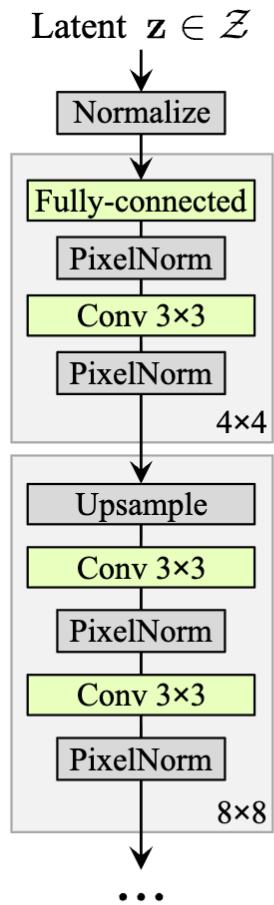
$$\mathcal{L}_{\text{VQ-VAE}}(\mathbf{x}) = \underbrace{\|\mathbf{x} - \mathcal{D}(\mathbf{e}_k)\|_2^2}_{\text{reconstruction loss}} + \underbrace{\|\text{sg}[\mathcal{E}(\mathbf{x})] - \mathbf{e}_k\|_2^2}_{\text{VQ loss}} + \underbrace{\beta \cdot \|\mathcal{E}(\mathbf{x}) - \text{sg}[\mathbf{e}_k]\|_2^2}_{\text{commitment loss}}$$

# GAN

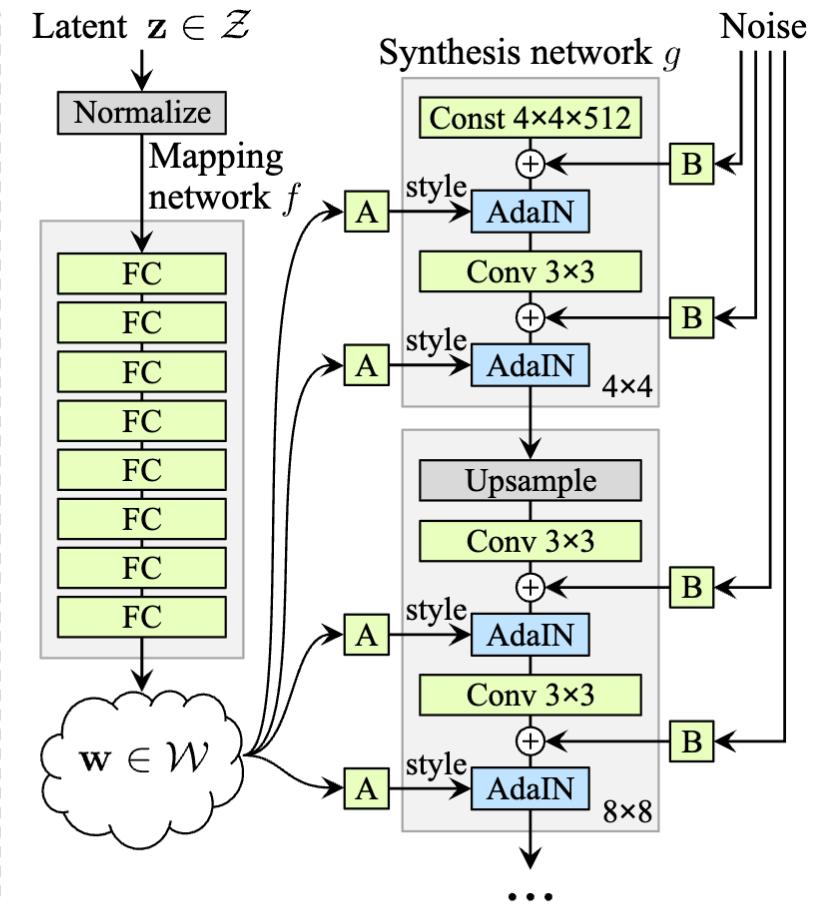
# StyleGAN

1. Style-based generator: maps latent  $\mathbf{z}$  to intermediate latent  $\mathbf{w}$  that controls image attributes
  2. Stochastic variation: adds per-pixel noise to each layer creating realistic micro-details (pores, hair, textures)
  3. Disentangled latent control: smooth and semantically meaningful manipulations in  $\mathbf{w}$ -space (e.g., pose, lighting, identity)

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i}$$



(a) Traditional



### (b) Style-based generator