AIR5101 / CIE6021

# Generative AI
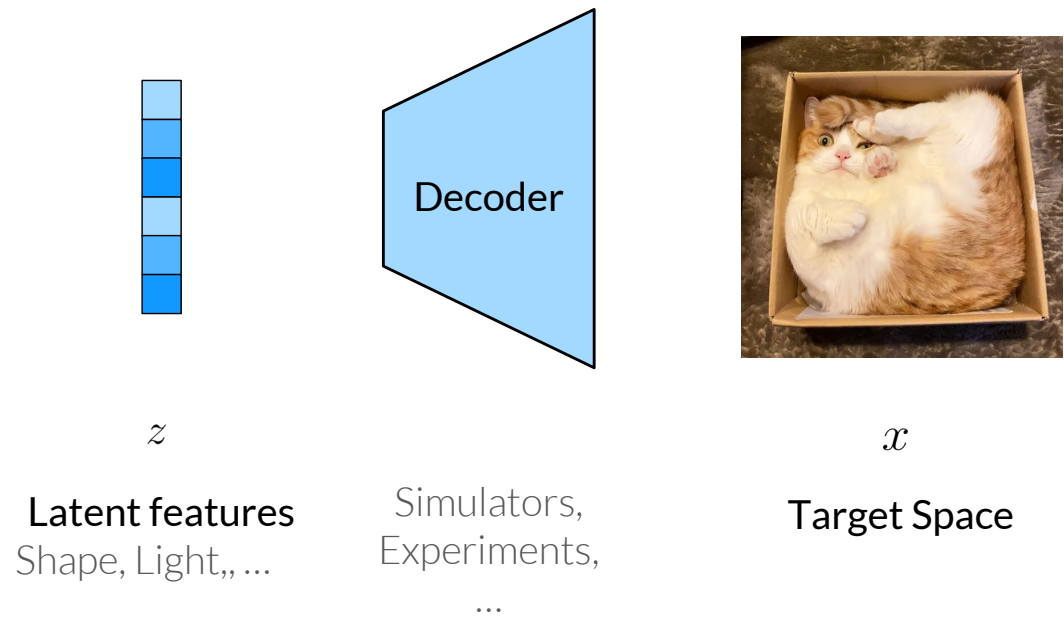
Lecture 2: Variational Autoencoder

Instructor: Zhen Liu

Spring 2025, CUHK-Shenzhen

# Some announcements

- First assignment to be released by mid of next week

- Part of the assignment: write a trainable VAE on toy datasets with PyTorch

  - Start early to see if you are comfortable in writing Python and learning PyTorch

- Find your project teammates and brainstorm ideas

- Paper list for presentations to be released by next week

# Latent Features



$z$

**Latent features**
Shape, Light,, …

Decoder

Simulators,
Experiments,
…

$x$

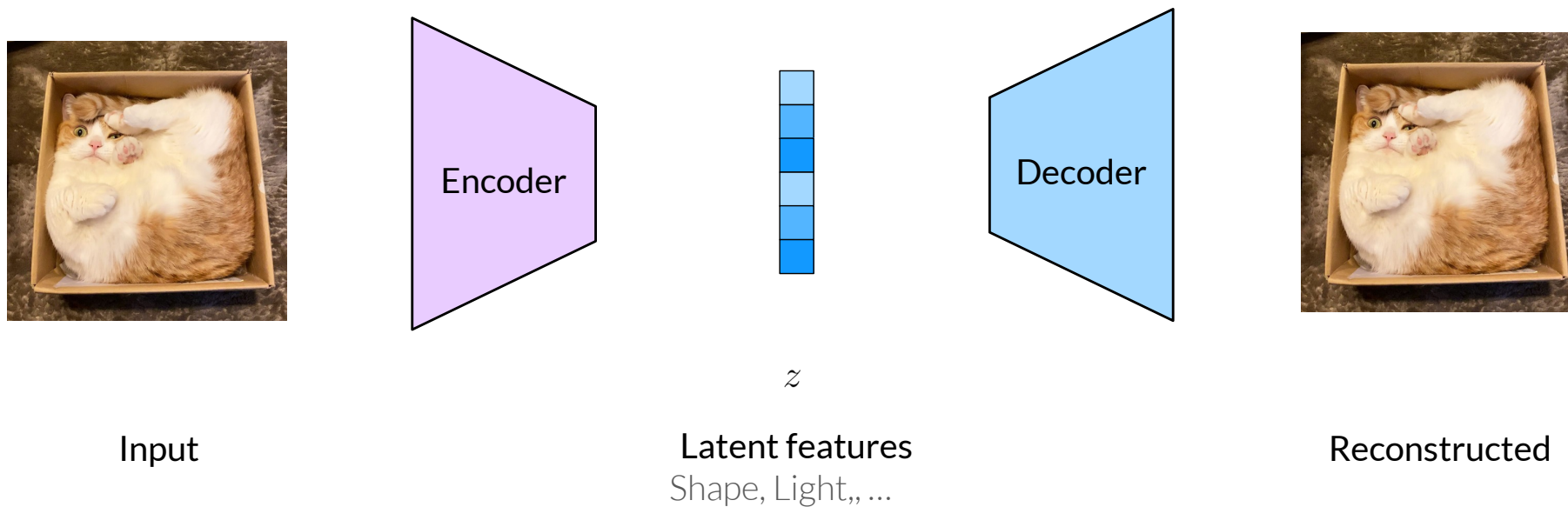Target Space

# Why Latent Spaces?

Manifold Hypothesis:

High-dimensional real-world data like images and videos
live in a low-dimensional manifold

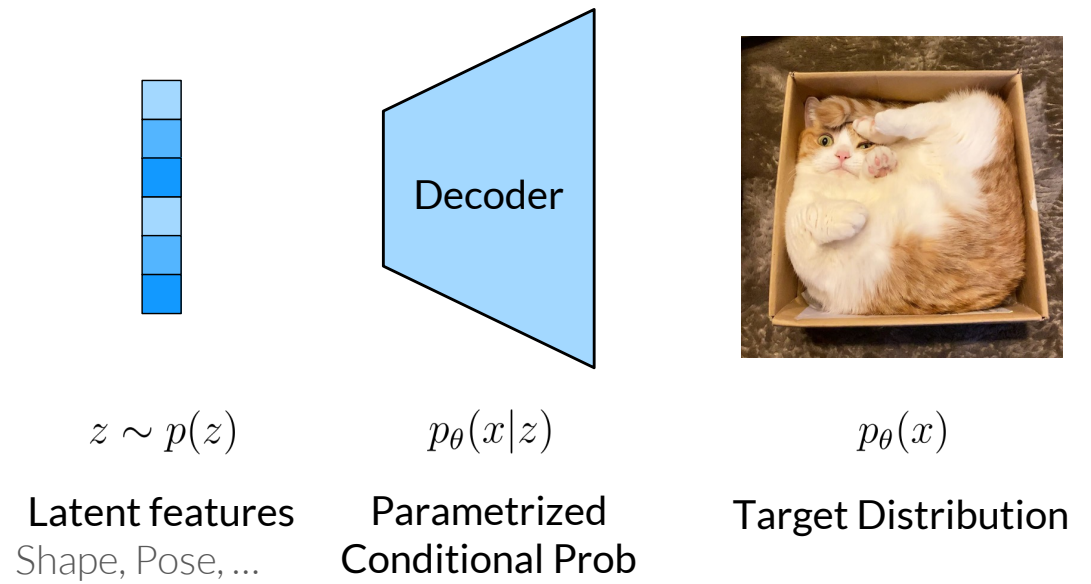Natural to learn distributions on this manifold instead.

# Autoencoder

Decoder: only see finite number of latent features



Input

Latent features
Shape, Light,, ...

$z$

Reconstructed

Issue: how to sample in the latent space?

Instead, study the problem under a probabilistic framework

# Latent Variable Model



$$z \sim p(z)$$

Latent features
Shape, Pose, ...

$$p_\theta(x|z)$$

Parametrized
Conditional Prob

$$p_\theta(x)$$

Target Distribution

Q: where does the latent feature come from?

# Recall: Maximum Likelihood Estimation

Optimize the "distance" between data and model

With "distance" as KL divergence,

$$\min_{\theta} D_{\mathrm{KL}}(p_{\mathcal{D}} \| p_{\theta})$$

$$= \min_{\theta} \mathop{\mathbb{E}}_{x \sim p_{\mathcal{D}}} \left[ \log p_{\mathcal{D}}(x) - \log p_{\theta}(x) \right] \qquad \text{(Definition of KL divergence)}$$

$$= \max_{\theta} \mathop{\mathbb{E}}_{x \sim p_{\mathcal{D}}} \log p_{\theta}(x)$$

# MLE for Latent Variable Model

Naïve attempt:

$$\max_{\theta} \; \mathbb{E}_{x \sim p_{\mathcal{D}}} \; \log p_{\theta}(x)$$

$$= \max_{\theta} \; \mathbb{E}_{x \sim p_{\mathcal{D}}} \; \log \int_{z} p_{\theta}(x|z) \; p(z) \; dz$$

$$\geq \max_{\theta} \; \mathbb{E}_{x \sim p_{\mathcal{D}}} \int_{z} \log p_{\theta}(x|z) p(z) dz \qquad \text{(Jensen's Inequality)}$$

$$= \max_{\theta} \; \mathbb{E}_{x \sim p_{\mathcal{D}}, z \sim p(z)} \; \log p_{\theta}(x|z)$$

Very bad estimate. What is the gap?

# MLE for Latent Variable Model

Decompose the log-likelihood:

$$\log p_\theta(x)$$

$$= \log \frac{p_\theta(x|z)p(z)}{p_\theta(z|x)} \qquad \text{(Bayes' Rule)}$$

$$= \log p_\theta(x|z) + \log p(z) - \log p_\theta(z|x)$$

# MLE for Latent Variable Model

Decompose the log-likelihood:

$$\log p_\theta(x)$$

$$= \log \frac{p_\theta(x|z)p(z)}{p_\theta(z|x)} \qquad \text{(Bayes' Rule)}$$

$$= \log p_\theta(x|z) + \log p(z) - \log p_\theta(z|x)$$

Therefore,

$$\underset{p(z)}{\mathbb{E}} \log p_\theta(x)$$

$$= \underset{p(z)}{\mathbb{E}} \log p_\theta(x|z) + \underset{p(z)}{\mathbb{E}} \left[\log p(z) - \log p_\theta(z|x)\right]$$

$$= \underset{p(z)}{\mathbb{E}} \log p_\theta(x|z) + D_{\mathrm{KL}}\Big(p(z)\|p_\theta(z|x)\Big) \qquad \text{Huge gap!}$$

10

# Proposal Distribution

Instead, we sample with a proposal distribution $q(z)$

$$\mathbb{E}_{q(z)} \log p_\theta(x)$$

$$= \mathbb{E}_{q(z)} \left[ \log p_\theta(x|z) + \log p(z) - \log p_\theta(z|x) \right]$$

$$= \mathbb{E}_{q(z)} \left[ \log p_\theta(x|z) \right] + \big( \log p(z) - \log q(z) \big) + \big( \log q(z) - \log p_\theta(z|x) \big) \right]$$

$$= \mathbb{E}_{q(z)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}}\Big( q(z) \| p(z) \Big) + D_{\mathrm{KL}}\Big( q(z) \| p_\theta(z|x) \Big)$$

"distance" from prior        Intractable

# Proposal Distribution

$$\mathop{\mathbb{E}}_{q(z)} \log p_\theta(x)$$

$$= \mathop{\mathbb{E}}_{q(z)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}}\Big( q(z) \| p(z) \Big) + D_{\mathrm{KL}}\Big( q(z) \| p_\theta(z|x) \Big)$$

Evidence Lower Bound (ELBO)        Non-negative (Property of KL)

$$\geq \mathop{\mathbb{E}}_{q(z)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}}\Big( q(z) \| p(z) \Big)$$

Lower bound of log-likelihood
⇒ Something we can optimize

12

# Encoder

Notice: we can pick any $q(z)$

The intractable $D_{\mathrm{KL}}\Big(q(z)\|p_\theta(z|x)\Big)$ implies we should pick

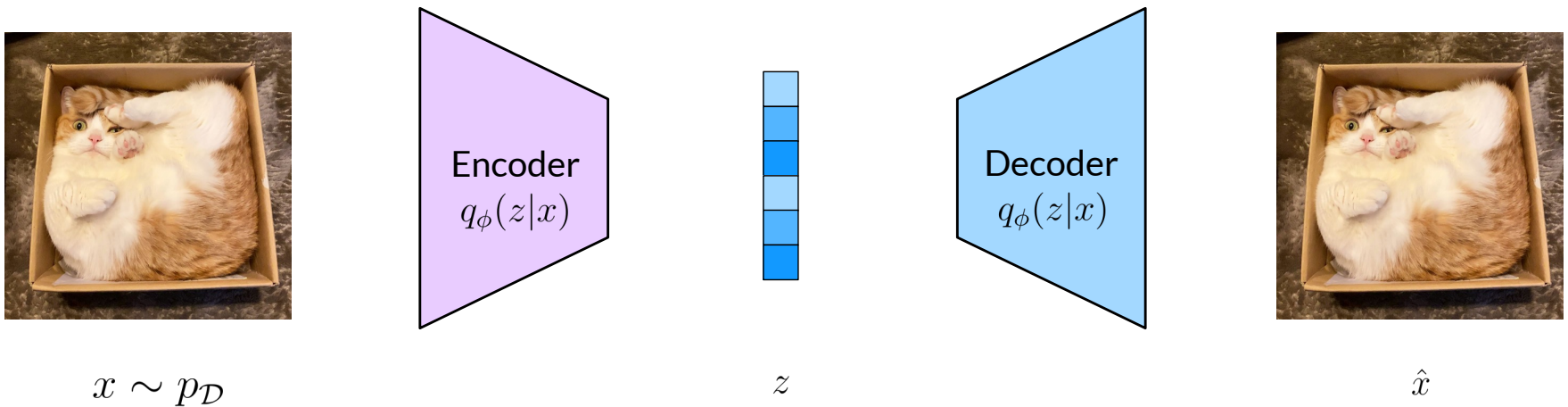$$q_\phi(z|x) \qquad \text{Encoding inputs into latents}$$

With this encoder, our log-likelihood becomes

$$\underset{q_\theta(z|x)}{\mathbb{E}}\Big[\log p_\theta(x|z)\Big] - D_{\mathrm{KL}}\Big(q_\theta(z|x)\|p(z)\Big) + D_{\mathrm{KL}}\Big(q_\theta(z|x)\|p_\theta(z|x)\Big)$$

ELBO with our encoder

# Variational Autoencoder (VAE)

$$\mathop{\mathbb{E}}_{z \sim q_\phi(z|x)} \Big[ \log p_\theta(x|z) \Big] - D_{\mathrm{KL}}\Big( q_\phi(z|x) \| p(z) \Big)$$



Encoder $q_\phi(z|x)$

Decoder $q_\phi(z|x)$

$x \sim p_\mathcal{D}$

$z$

$\hat{x}$

# Variational Autoencoder (VAE)

$$\mathbb{E}_{z \sim q_\phi(z|x)}\Big[\log p_\theta(x|z)\Big] - D_{\mathrm{KL}}\Big(q_\phi(z|x)\|p(z)\Big)$$

Encoder $q_\phi(z|x)$

Decoder $q_\phi(z|x)$

$x \sim p_{\mathcal{D}}$

$z$

$\hat{x}$

Reconstruction Loss

# Variational Autoencoder (VAE)

$$\mathop{\mathbb{E}}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}}\Big( q_\phi(z|x) \| p(z) \Big)$$



$x \sim p_{\mathcal{D}}$

Encoder $q_\phi(z|x)$

KL Regularization

$z$

Decoder $q_\phi(z|x)$

$\hat{x}$

$\mathcal{N}(0, I)$

# Choice of Decoder

$$\underset{z \sim q_\phi(z|x)}{\mathbb{E}} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}}\left( q_\phi(z|x) \| p(z) \right)$$

Typical choice: isotropic Gaussian

$$\log p_\theta(z|x) \propto \|z - g_\theta(x)\|^2 / (2\sigma_o^2)$$

Pick some fixed standard deviation ⇒ scaled L2 loss

# Choice of Encoder

$$\mathop{\mathbb{E}}_{z \sim q_\phi(z|x)} \Big[ \log p_\theta(x|z) \Big] - D_{\mathrm{KL}}\Big( q_\phi(z|x) \| p(z) \Big)$$

We should be able to easily sample $q_\phi(z|x)$

Natural choice: Gaussian

Encoder now outputs $\quad f_\phi(x) \rightarrow \mu_\phi(x), \sigma_\phi(x)$

Analytical form of KL regularization (left as exercise)

$$D_{\mathrm{KL}}\Big( \mathcal{N}(\mu_\theta(x), \sigma_\theta^2(x)) \,\Big\| \, \mathcal{N}(0, I) \Big)$$

# Reparametrization Trick

$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right]$$
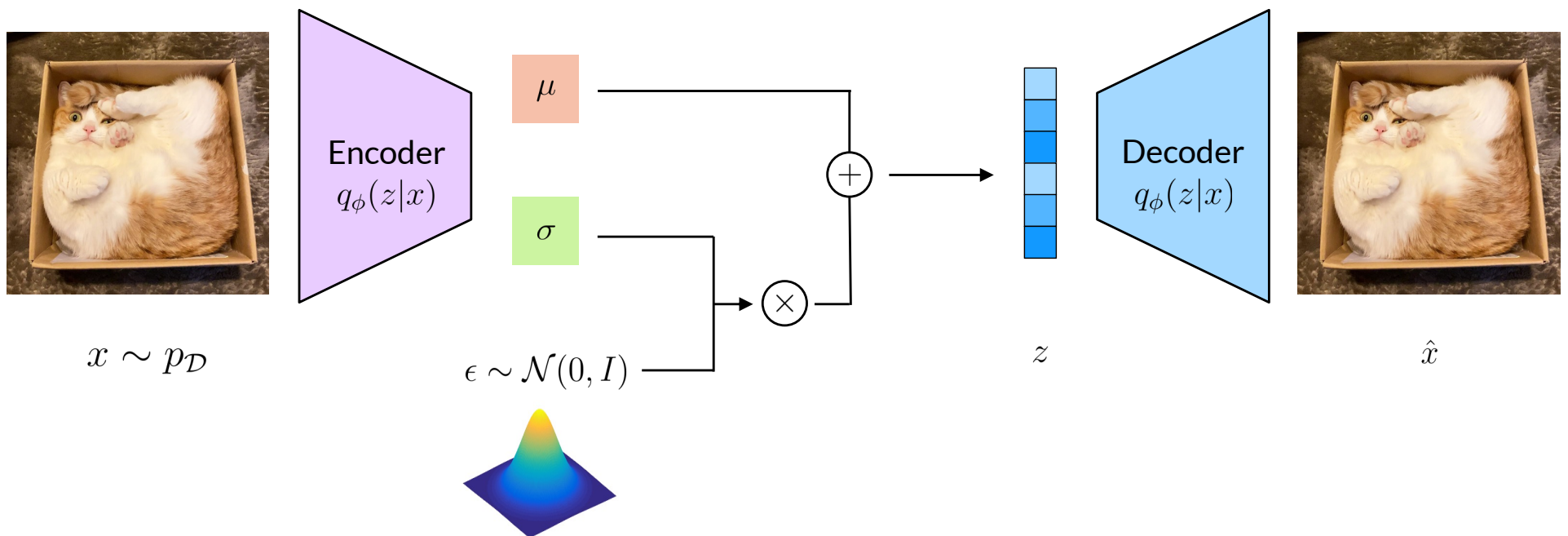
How to use backprop to optimize the encoder?

$$z \sim q_\phi(z|x)$$

$$\|$$

$$z = \mu_\theta(x) + \sigma_\theta(x)\epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

Stochasticity isolated to fixed distribution of epsilon

# Reparametrization Trick



$x \sim p_{\mathcal{D}}$

Encoder $q_\phi(z|x)$

$\mu$

$\sigma$

$\epsilon \sim \mathcal{N}(0, I)$

$z$

Decoder $q_\phi(z|x)$

$\hat{x}$

# Full Training Loss

ELBO for a single datapoint

$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}} \Big( q_\phi(z|x) \| p(z) \Big)$$

ELBO for the whole dataset

$$\mathbb{E}_{x \sim p_\mathcal{D}} \left[ \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - D_{\mathrm{KL}} \Big( q_\phi(z|x) \| p(z) \Big) \right]$$

# Generation with VAEs

Step 1: Sample from standard Gaussian $z \sim \mathcal{N}(0, I)$

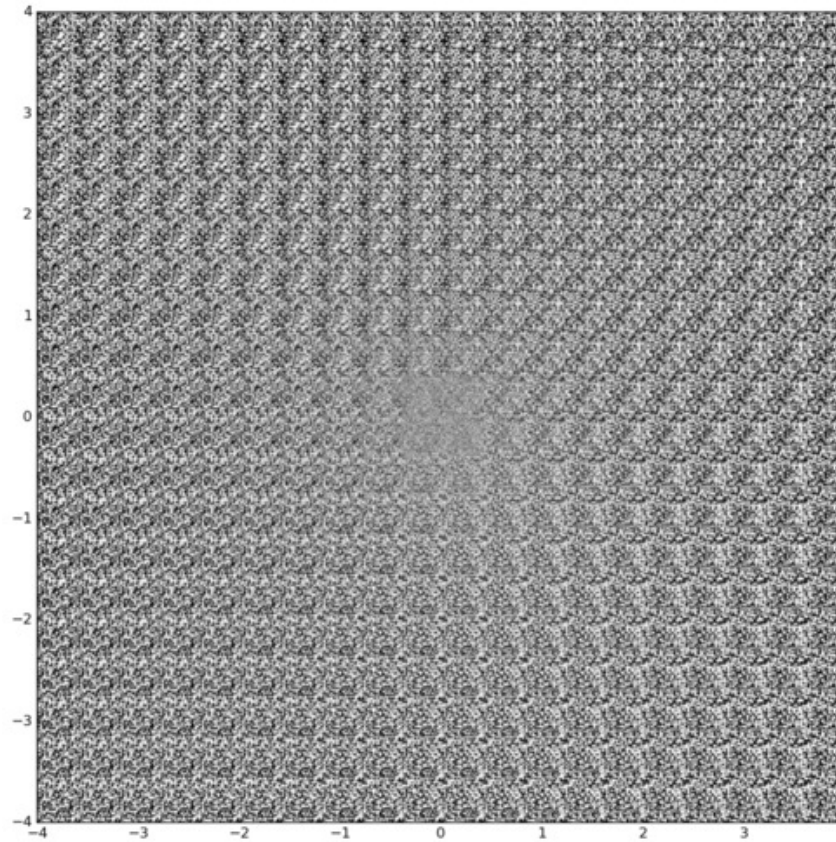Why it makes sense? $q_\phi(z|x)$ is close to $\mathcal{N}(0, I)$ due to the KL regularization

Step 2: Sample the output with $p_\theta(x|z)$

In practice, set the standard deviation of this decoder distribution to zero
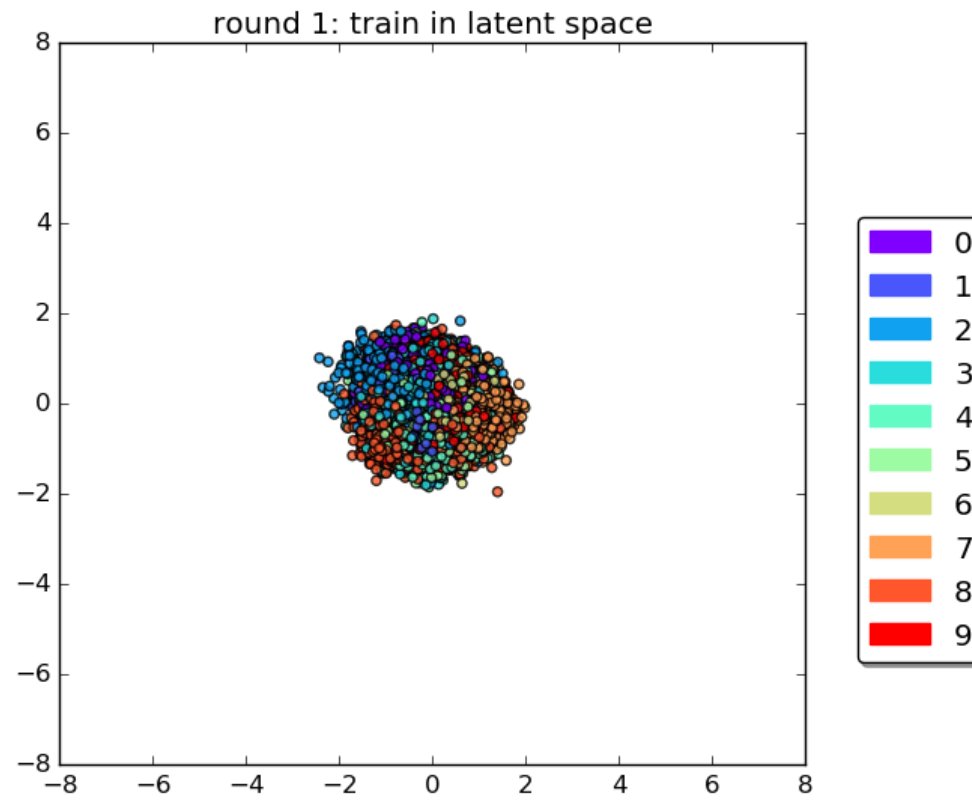
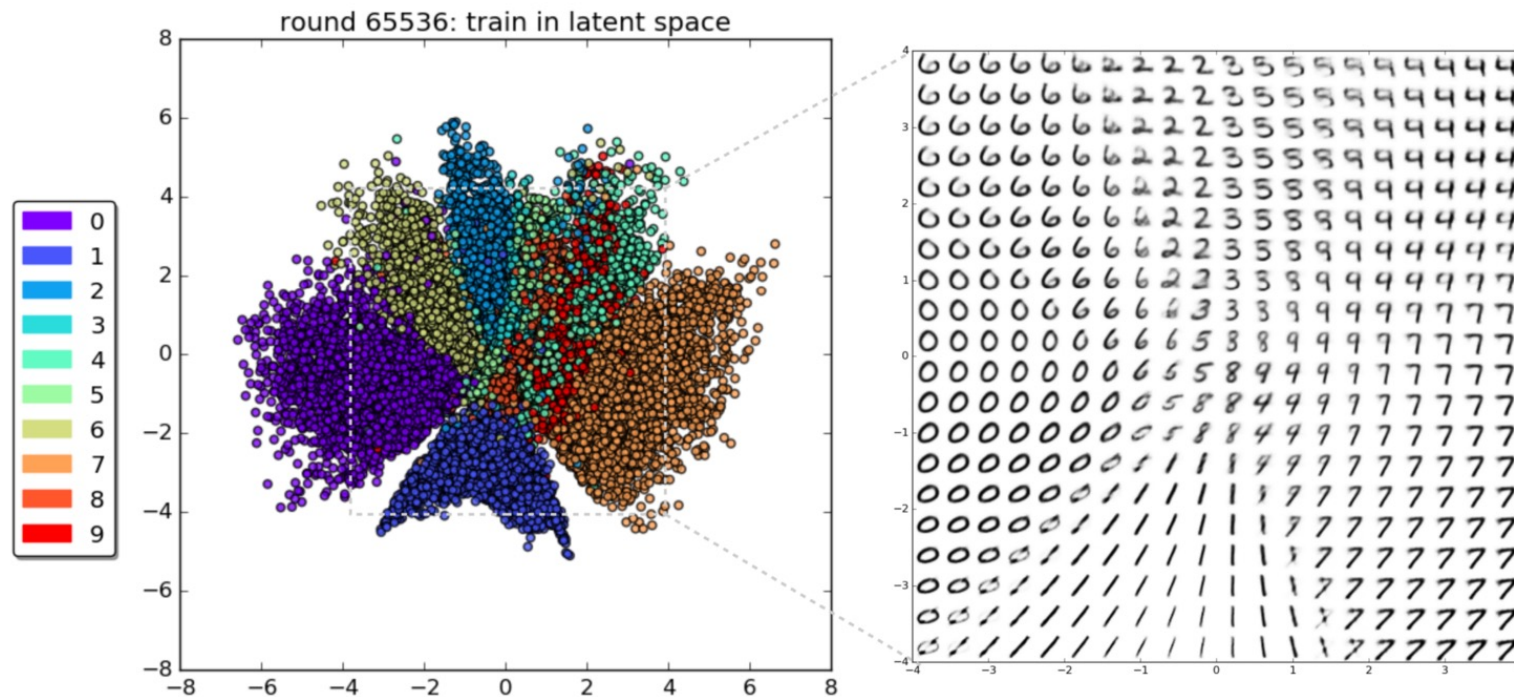# 2D Latent Space of MNIST

# 2D Latent Space of MNIST

# 2D Latent Space of MNIST



round 1: train in latent space

# 2D Latent Space of MNIST



round 65536: train in latent space

# VAE as Expectation-Maximization

Instead of gradient descent on theta and phi

Can use EM in simple cases

E step (Expectation): $\quad q^{(t)} = p_{\theta^{(t)}}(z|x)$

M step (Maximization): $\quad \theta^{(t+1)} = \arg\max_{\theta} \; \mathbb{E}_{\substack{x \sim p_{\mathcal{D}} \\ z \sim p_{\theta^{(t)}}(z|x)}} \left[ \log p_{\theta^{(t)}}(x,z) \right]$

Useful if $q^{(t)}$ has analytical forms:
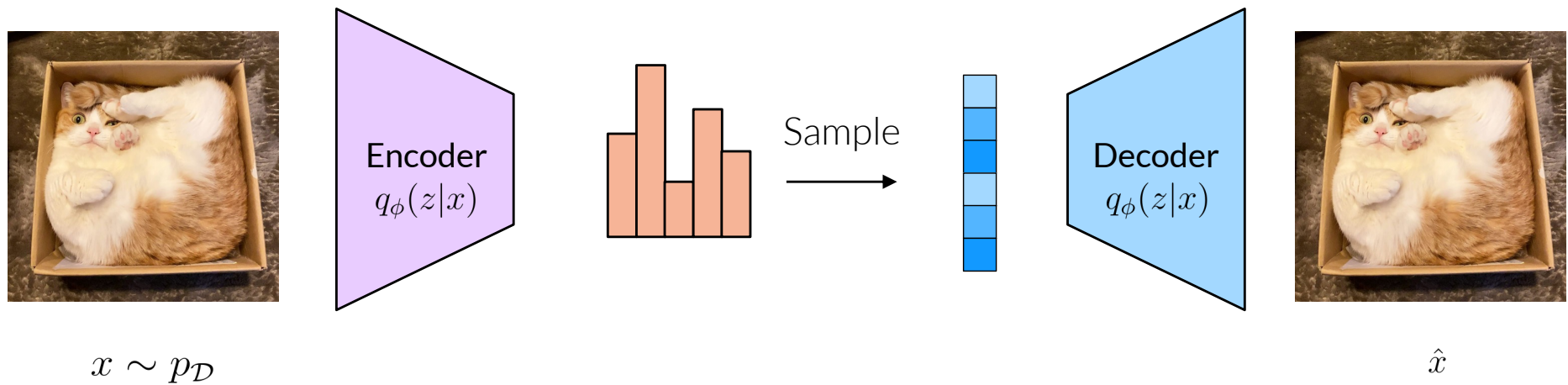mixture of Gaussians, K-means, etc.

# Failure of VAE

Blurry reconstructed images

One significant limitation of vanilla VAE:
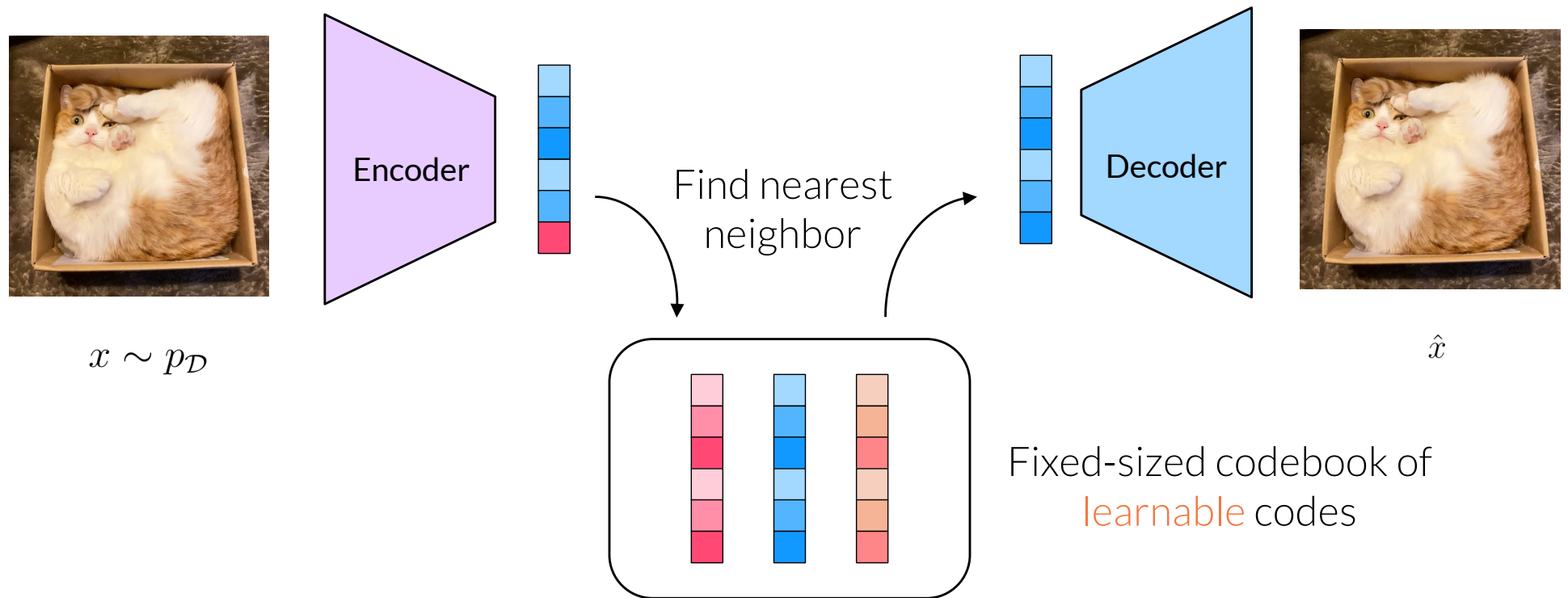- Gaussian encoder and decoder

# Proposal: Discretized Latent Space



$$x \sim p_{\mathcal{D}}$$

Encoder $q_\phi(z|x)$

Sample

Decoder $q_\phi(z|x)$

$$\hat{x}$$

## How to discretize the latent space?

# Vector Quantization



$x \sim p_{\mathcal{D}}$

Encoder

Find nearest neighbor

Decoder

$\hat{x}$

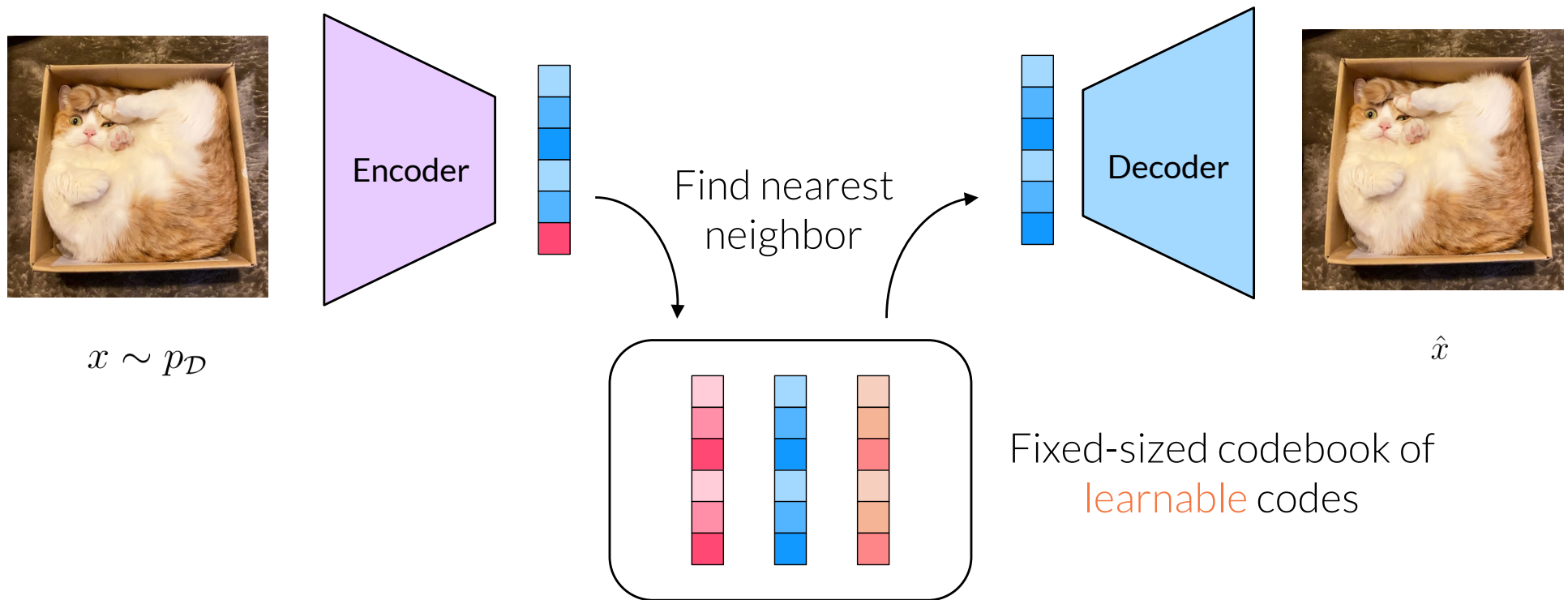Fixed-sized codebook of learnable codes

# Training VQ-VAE

$$\mathcal{L} = \underbrace{\|x - \hat{x}\|^2}_{\text{Reconstruction Loss}} + \beta \underbrace{\|z_e(x) - \text{sg}[z_q(x)]\|^2}_{\text{Commitment Loss}} + \underbrace{\|\text{sg}[z_e(x)] - z_q(x)\|^2}_{\text{Codebook Loss}},$$

sg = stop-gradient

- Commitment loss: encoded feature move towards quantized feature

- Codebook loss: quantized feature towards encoded feature

Exercise: prove that the KL term in ELBO is a constant (with a uniform prior)

# Repametrization Trick for VQ-VAE?



Encoder

Find nearest
neighbor

Decoder

$x \sim p_{\mathcal{D}}$

$\hat{x}$
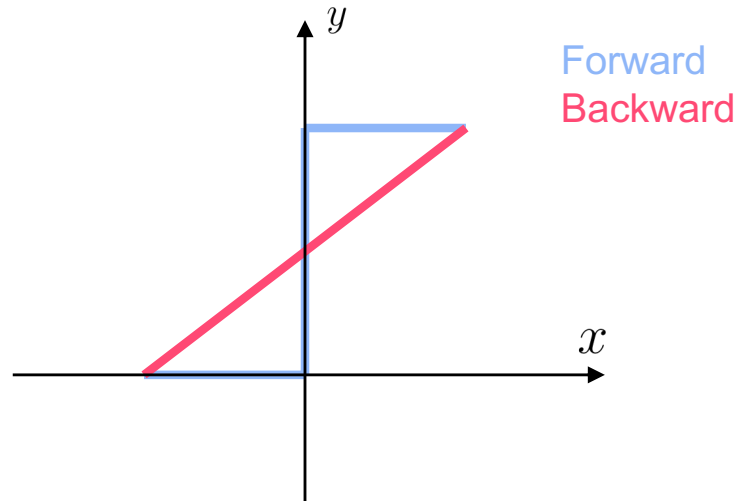
Fixed-sized codebook of
learnable codes

## Non-differentiable operation!

# Straight-through Estimator

Forward: thresholding function

Backward: pretend that it is an identity function (or some other proxy functions)
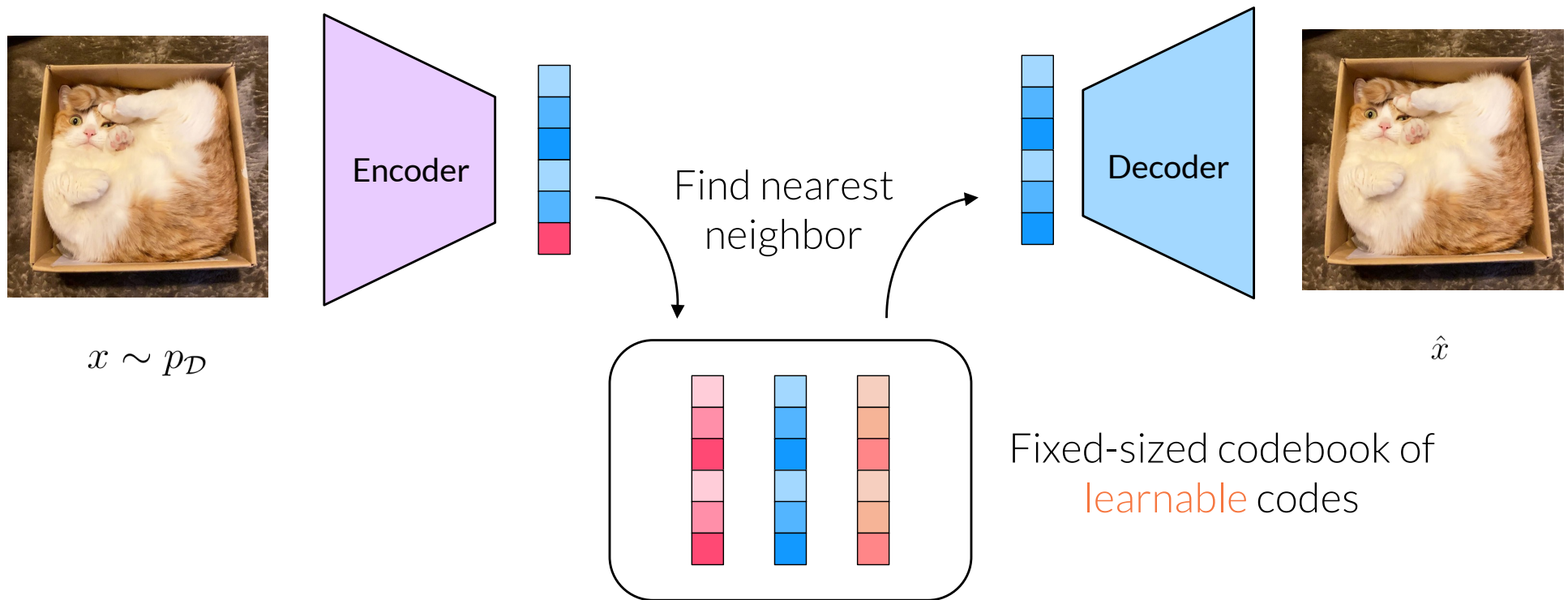
# Straight-through Estimator (STE)

Sampling from categorical distribution = "thresholding" on softmax values

Implementation to use softmax gradients for argmax:

$$\text{stop-grad}\Big(\text{argmax}_k(x_{1:K}) - \text{softmax}_{k,\tau}(x_{1:K})\Big) + \text{softmax}_{k,\tau}(x_{1:K})$$

Output of forward      Gradient for backward

Can we do better?

# Repametrization Trick for VQ-VAE?



Encoder

Find nearest neighbor

Decoder

$x \sim p_{\mathcal{D}}$

$\hat{x}$

Fixed-sized codebook of learnable codes

## How to construct the "epsilon"?

# Gumbel-Max Trick

Suppose that we have a categorical distribution $p_\theta(i)$ to pick indices

Define

$$i^* = \arg\max_i \left( \log p_\theta(i) - \log \left( -\log \epsilon_i \right) \right), \quad \epsilon_i \sim \mathrm{Uniform}(0,1)$$

Gumbel distribution

Exercise: show that

$$i^* \sim p_\theta(i)$$
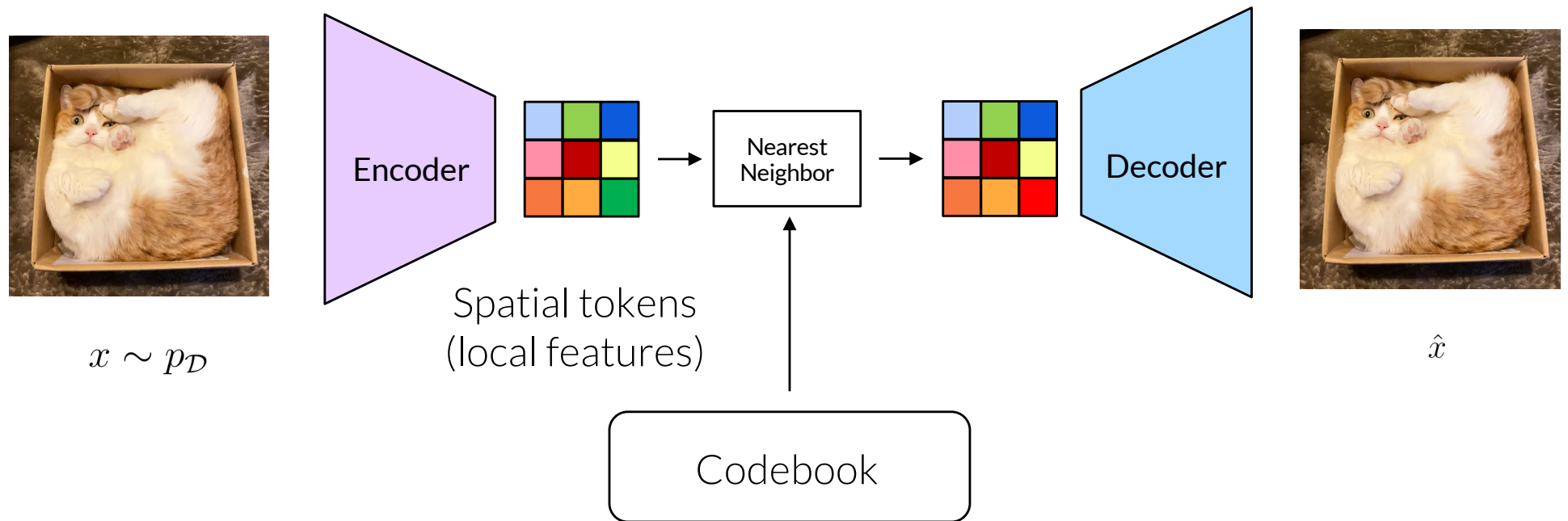
# Gumbel-Softmax

Argmax is not differentiable either

⇒ Replace argmax with softmax

$$i^* = \operatorname*{softmax}_{i,\tau} \Big( \log p_\theta(i) - \log \big( -\log \epsilon_i \big) \Big), \quad \epsilon_i \sim \mathrm{Uniform}(0,1)$$

where

$$\mathrm{softmax}_{i,\tau}(x_i) = \frac{\exp(x_i/\tau)}{\sum_i \exp(x_i/\tau)}$$

# VQ-VAE as Tokenizers



$x \sim p_{\mathcal{D}}$

Encoder

Spatial tokens
(local features)

Nearest
Neighbor
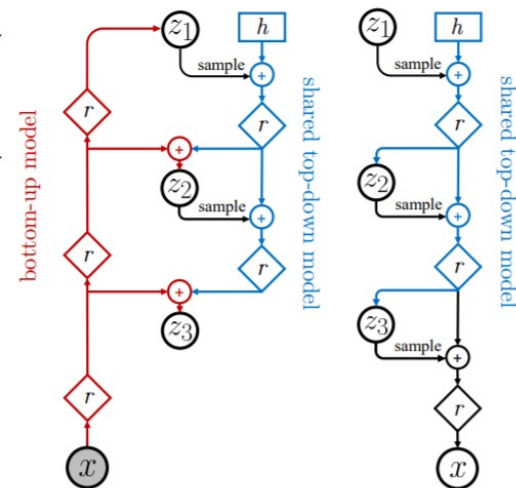
Codebook

Decoder

$\hat{x}$

$\Rightarrow$ **Autoregressive modeling on tokens** *(e.g., VQ-VAE-2 and some recent work)*

# Alternative Solution: Hierarchical VAEs

VAEs on latents

Harder (than vanilla VAEs) to train

Will see its connections to other models



(a) Bidirectional Encoder (b) Generative Model

Figure 2: The neural networks implementing an encoder $q(z|x)$ and generative model $p(x, z)$ for a 3-group hierarchical VAE. ⟨r⟩ denotes residual neural networks, ⊕ denotes feature combination (e.g., concatenation), and [h] is a trainable parameter.

NVAE: A Deep Hierarchical Variational Autoencoder. Vahdat et al. NeurIPS 2020.