# VAE

Aziz Temirkhanov

Laboratory for methods of big data analysis
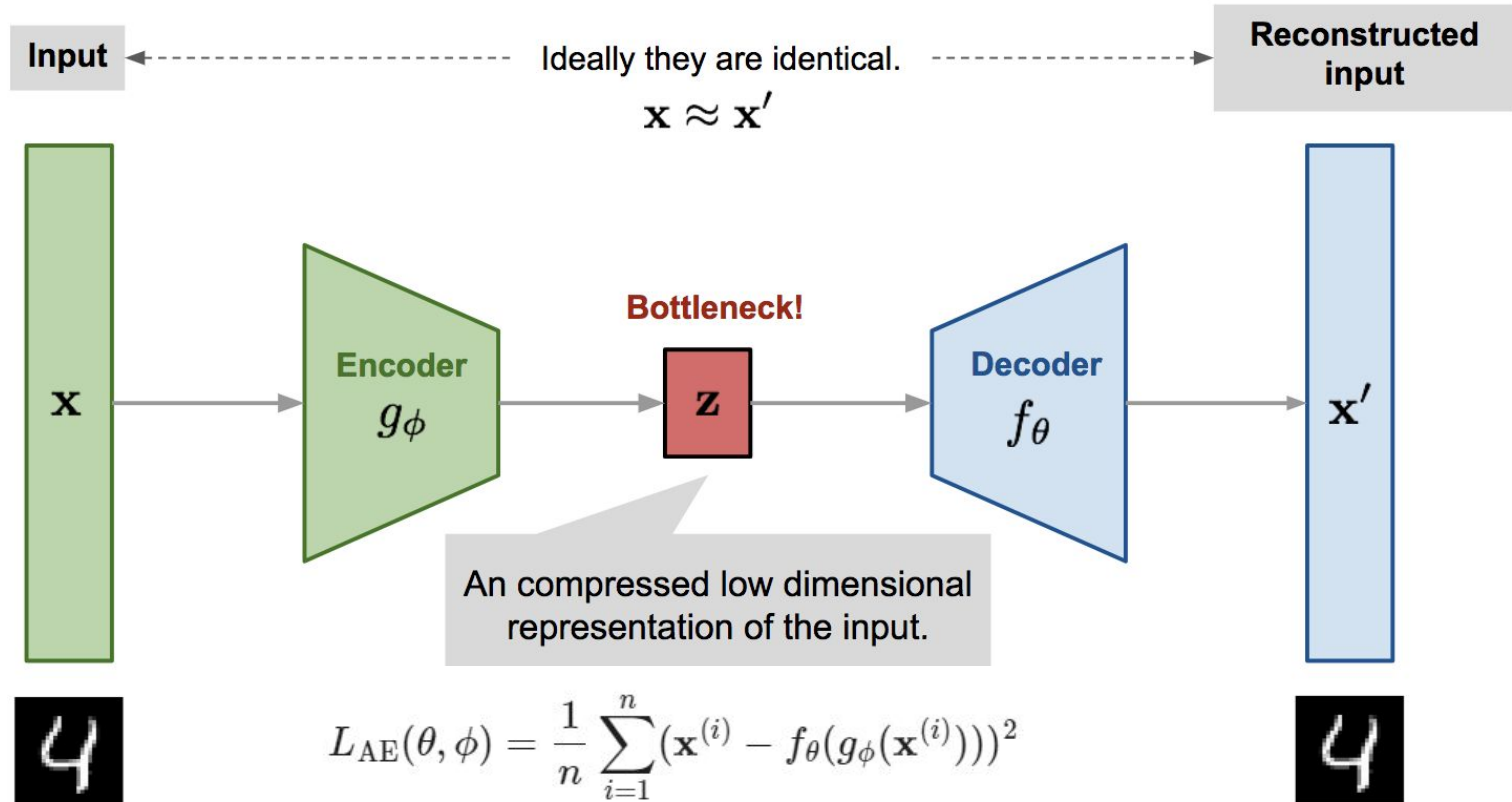
# Autoencoders

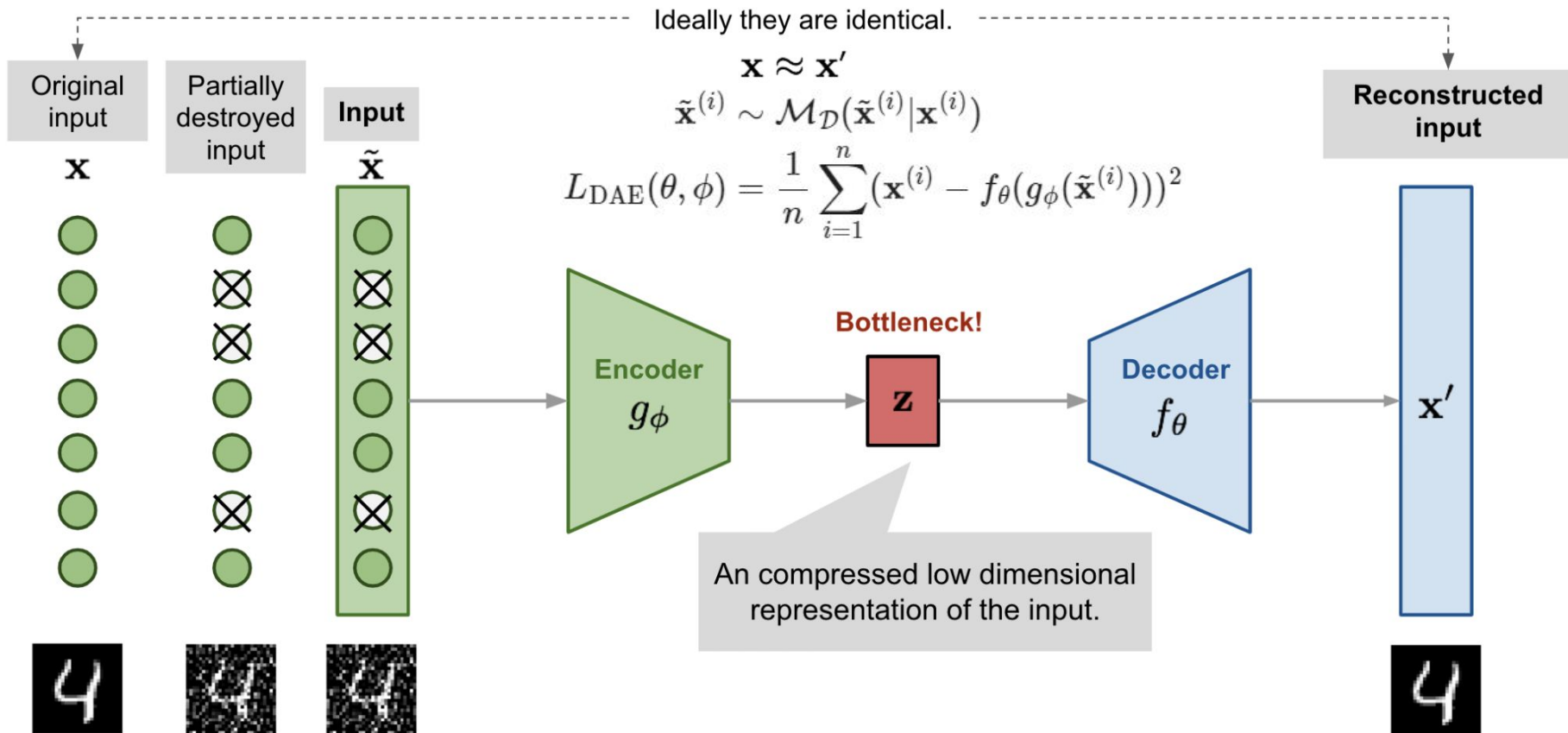# Idea



Input ← – – – – – – – – – – – – – – Ideally they are identical. – – – – – – – – – – – – → **Reconstructed input**

$$\mathbf{x} \approx \mathbf{x}'$$

**Bottleneck!**

**Encoder** $g_\phi$ — $\mathbf{z}$ — **Decoder** $f_\theta$

$\mathbf{x}$ → → $\mathbf{x}'$

An compressed low dimensional representation of the input.

$$L_{\mathrm{AE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}^{(i)} - f_\theta(g_\phi(\mathbf{x}^{(i)})))^2$$

# Denoising AE



Ideally they are identical.

Original input $\mathbf{x}$

Partially destroyed input

Input $\tilde{\mathbf{x}}$

$$\mathbf{x} \approx \mathbf{x}'$$

$$\tilde{\mathbf{x}}^{(i)} \sim \mathcal{M}_{\mathcal{D}}(\tilde{\mathbf{x}}^{(i)}|\mathbf{x}^{(i)})$$

$$L_{\text{DAE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}^{(i)} - f_\theta(g_\phi(\tilde{\mathbf{x}}^{(i)})))^2$$

**Reconstructed input**

Encoder $g_\phi$

**Bottleneck!**

$\mathbf{z}$

Decoder $f_\theta$

$\mathbf{x}'$

An compressed low dimensional representation of the input.

# VAE

# VAE

▶ Instead of mapping the input into a *fixed* vector $z$, we want to map it into distribution $p_\theta$ parameterized by $\theta$

▶ In this setup, the relation between the input data $x$ and the latent encoding vector $z$ can be defined in terms of bayesian framework.

▶ Prior $p_\theta(z)$

▶ Likelihood $p_\theta(x|z)$

▶ Posterior $p_\theta(z|x)$

# VAE

▶ Assume we know $\theta^*$

▶ First, sample a $\mathbf{z^{(i)}}$ from a prior distribution $p_\theta(z)$

▶ Then a value $\mathbf{x^{(i)}}$ is generated from a conditional distribution $p_\theta(x|z=z^{(i)})$

▶ Optimality:

$$\theta^* = \arg\max_\theta \prod_{i=1}^n p_\theta(\mathbf{x}^{(i)}) \qquad \theta^* = \arg\max_\theta \sum_{i=1}^n \log p_\theta(\mathbf{x}^{(i)})$$

$$p_\theta(\mathbf{x}^{(i)}) = \int p_\theta(\mathbf{x}^{(i)}|\mathbf{z})p_\theta(\mathbf{z})d\mathbf{z}$$
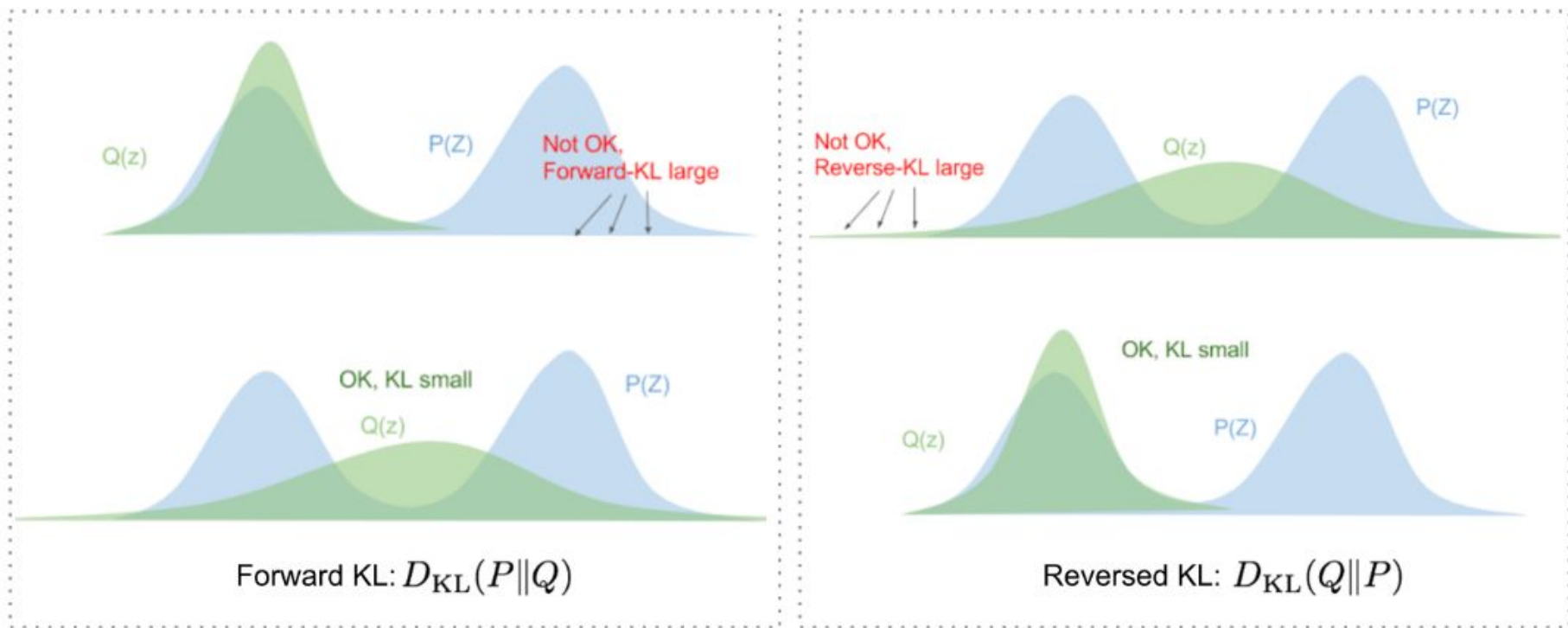
# VAE



▶ It is not easy to compute, since it is very expensive to check all the possible values of z

▶ Thus, let's introduce a new approximation function $q_{\varphi}(z|x)$

▶ Now, the structure looks a lot like AE, where conditional probability $p_{\theta}(x|z)$ defines a generative models, similar to decoder, and it is also called *probabilistic decoder*, and the approximation function $q_{\varphi}(z|x)$ is the *probabilistic encoder*.

8

# ELBO

# ELBO



Forward KL: $D_{KL}(P \| Q)$

Reversed KL: $D_{KL}(Q \| P)$

▶ The estimated posterior $q_\varphi(z|x)$ should be very close the real one $p_\theta(x|z)$

▶ We can use KL divergence $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) | p_\theta(\mathbf{z}|\mathbf{x}))$

# ELBO

$$D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x}))$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})p_\theta(\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})} d\mathbf{z} \qquad ; \text{Because } p(z|x)=p(z,x)/p(x)$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \big( \log p_\theta(\mathbf{x}) + \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})} \big) d\mathbf{z}$$

$$= \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})} d\mathbf{z} \qquad ; \text{Because } \int q(z|x)dz=1$$

$$= \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})} d\mathbf{z} \qquad ; \text{Because } p(z,x)=p(x|z)p(z)$$

$$= \log p_\theta(\mathbf{x}) + \mathbb{E}_{\mathbf{z}\sim q_\phi(\mathbf{z}|\mathbf{x})}[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} - \log p_\theta(\mathbf{x}|\mathbf{z})]$$

$$= \log p_\theta(\mathbf{x}) + D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z})) - \mathbb{E}_{\mathbf{z}\sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z})$$

# ELBO

$$D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x})) = \log p_\theta(\mathbf{x}) + D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z})) - \mathbb{E}_{\mathbf{z}\sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z})$$

$$\log p_\theta(\mathbf{x}) - D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{\mathbf{z}\sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}))$$

$$
\begin{aligned}
L_{\mathrm{VAE}}(\theta, \phi) &= -\log p_\theta(\mathbf{x}) + D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x})) \\
&= -\mathbb{E}_{\mathbf{z}\sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) + D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z})) \\
\theta^*, \phi^* &= \arg\min_{\theta,\phi} L_{\mathrm{VAE}}
\end{aligned}
$$

Since KL divergence is always non-negative and thus $-L_{\mathrm{VAE}}$ is the lower bound of $\log p_\theta(\mathbf{x})$
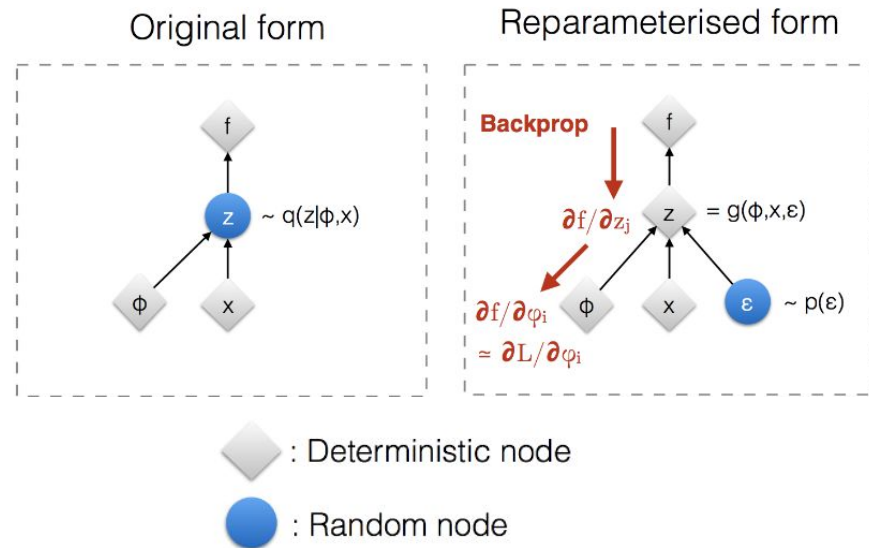
$$-L_{\mathrm{VAE}} = \log p_\theta(\mathbf{x}) - D_{\mathrm{KL}}(q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z}|\mathbf{x})) \leq \log p_\theta(\mathbf{x})$$
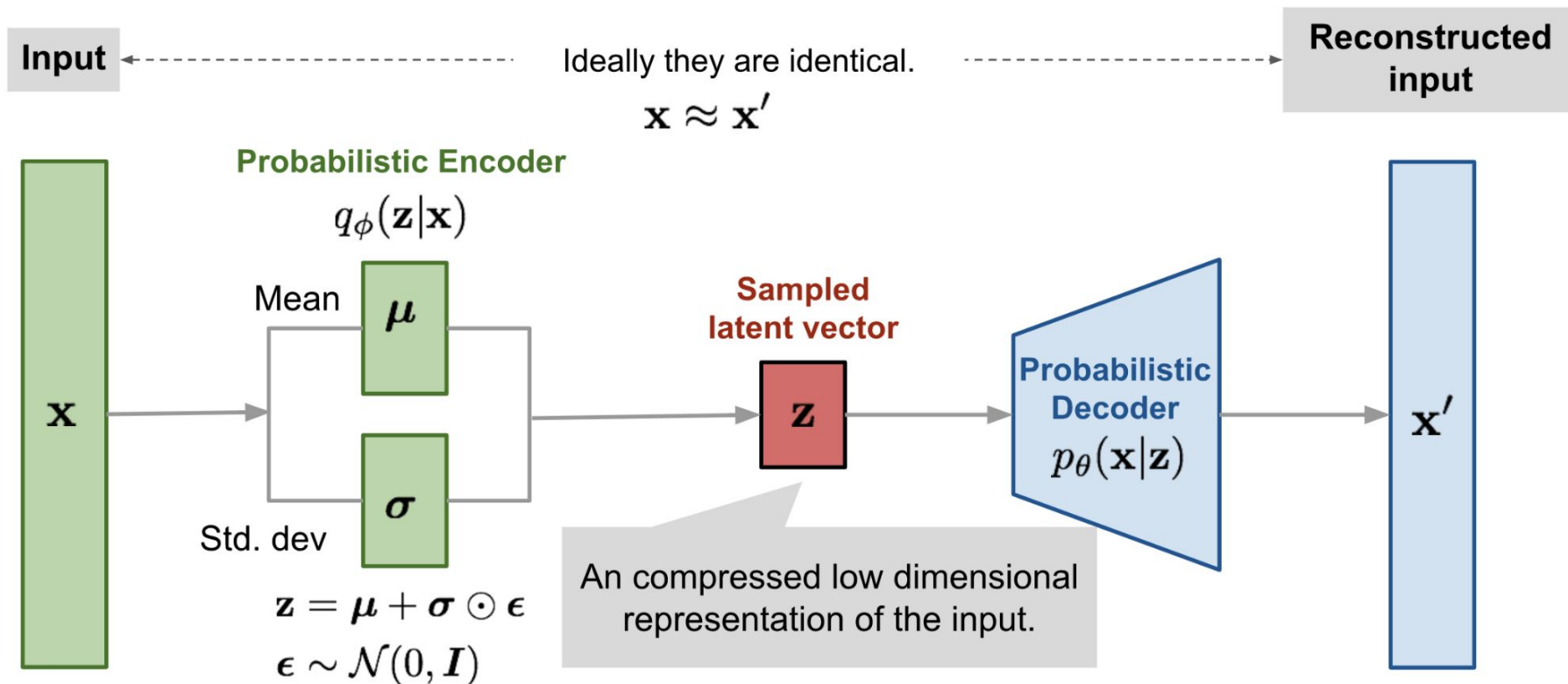
# Reparametrization trick

▶ The expectation term in the loss function invokes generating samples from $z \sim q_\varphi(z|x)$

▶ Sampling is a stochastic process, and cannot be backpropagated

▶ Thus, let's introduce the reparametrization trick:



Original form    Reparameterised form

$\square$ : Deterministic node

$\bullet$ : Random node

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \boldsymbol{I})$$
$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I}) \qquad ; \text{Reparameterization trick.}$$
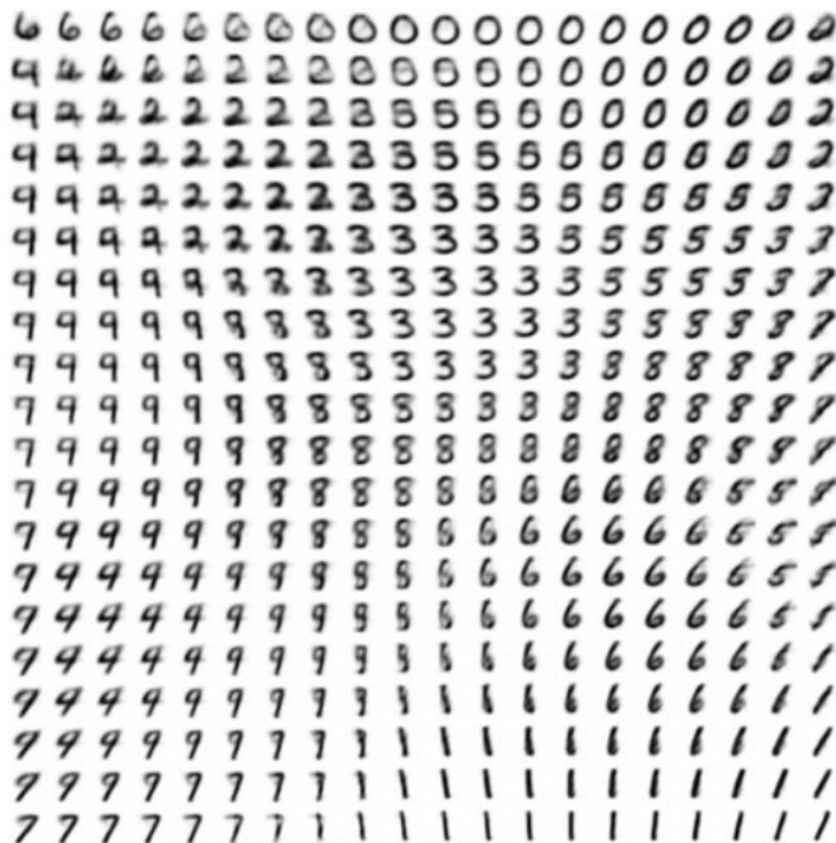
# VAE

**Input** <--------------------------------- Ideally they are identical. -------------------------> **Reconstructed input**

$$\mathbf{x} \approx \mathbf{x}'$$

**Probabilistic Encoder**

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

Mean $\boldsymbol{\mu}$

Std. dev $\boldsymbol{\sigma}$

**Sampled latent vector**

$\mathbf{z}$

**Probabilistic Decoder**

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

$\mathbf{x}$

$\mathbf{x}'$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

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

An compressed low dimensional representation of the input.

# Results

- The latent space is well organized

- The results does not suffer from mode collapse

- The sampling speed is decent

- Images are blurred

# Results



(a) 2-D latent space     (b) 5-D latent space     (c) 10-D latent space     (d) 20-D latent space

# Results

Due to minimization of KL
$$KL(q(z; \phi)||p(z|x; \theta)) \to 0.$$

# VQ-VAE

# Latemd Spaces

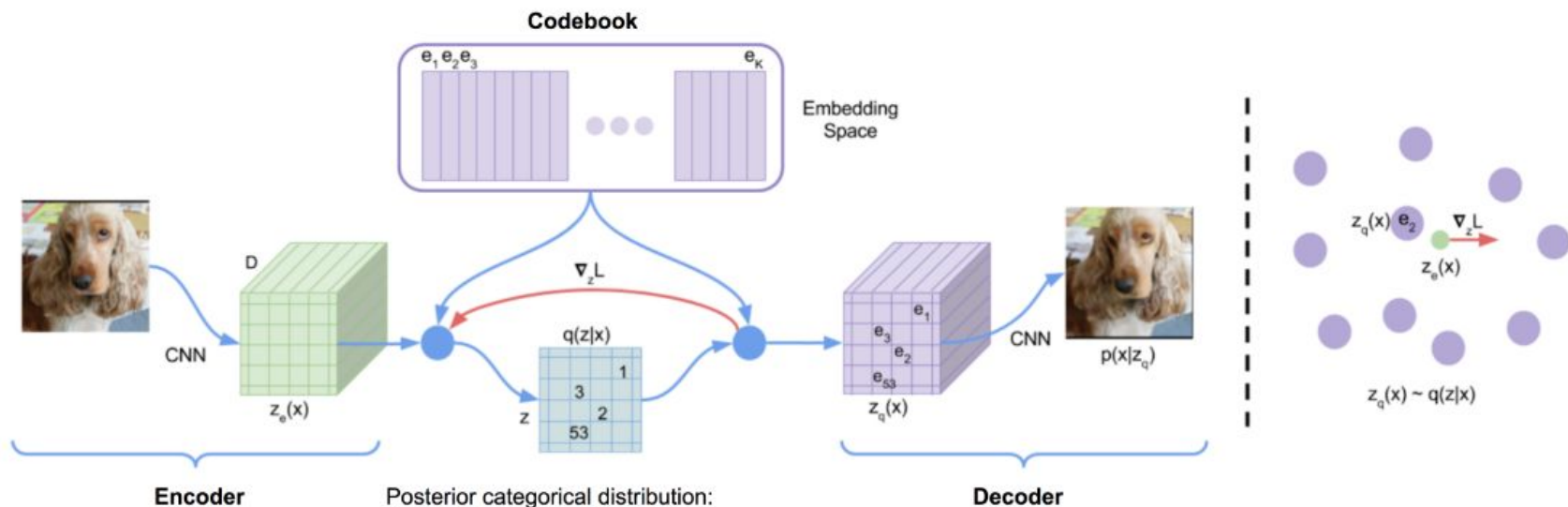VAE gives a good representations

Can we study it?

**Encoder**



image to
discrete codes

| 56 | 73 | 67 | 23 | 81 | 19 | ... |

**Decoder**

| 56 | 73 | 67 | 23 | 81 | 19 | ... |

discrete codes
to image



An Oversimplified Example of a Cat/Dog Image Latent Space



time of day of the image

daytime

night time

More cat-like

More dog-like

how dog-like versus cat-like an image is

# VQ-VAE



Posterior categorical distribution:

$$q(\mathbf{z} = \mathbf{e}_k | \mathbf{x}) = \begin{cases} 1 & \text{if } k = \arg\min_i \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_i\|_2 \\ 0 & \text{otherwise.} \end{cases}$$

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