# Lecture 7: Unsupervised Representation Learning and Generative Models (cont.)

## Shujian Yu
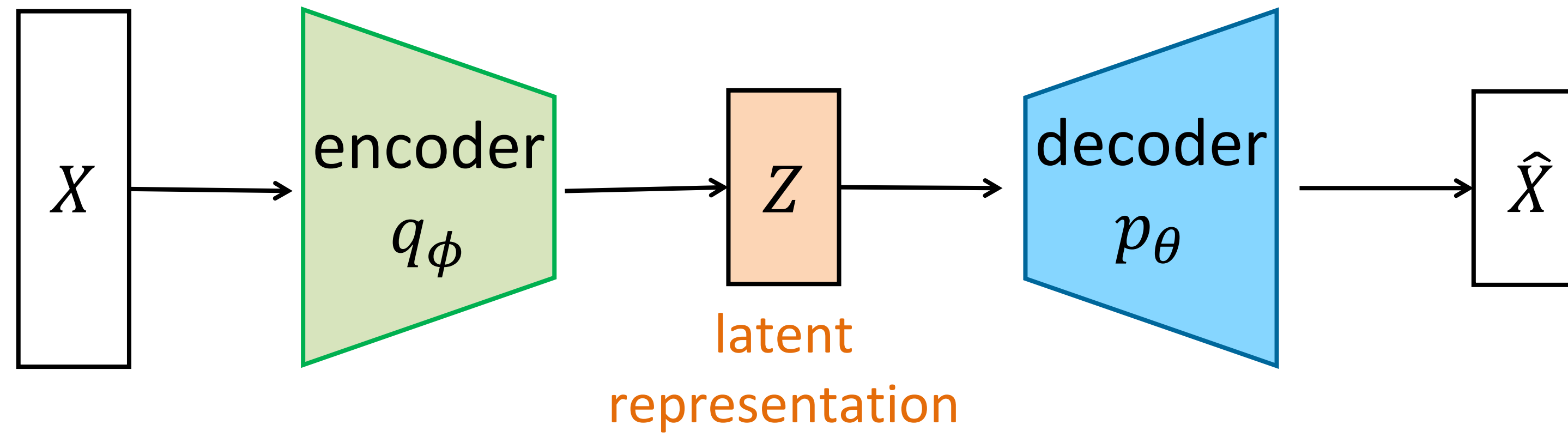Deep Learning 2023

dlvu.github.io

VRIJE UNIVERSITEIT AMSTERDAM

**part 1:** VAE implementation

**part 2:** KL divergence and maximum mean discrepancy

**part 3:** MMD-VAE and $\beta$-VAE

VU

$$X \rightarrow \text{encoder } q_\phi \rightarrow Z \rightarrow \text{decoder } p_\theta \rightarrow \hat{X}$$

latent
representation

VU

**Variational posterior**

$$\ln p_\theta(\boldsymbol{x}) = \log \int p_\theta(\boldsymbol{x}|\boldsymbol{z}) p_\lambda(\boldsymbol{z}) d\boldsymbol{z}$$

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

$$\geq \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \log \frac{p_\theta(\boldsymbol{x}|\boldsymbol{z}) p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} d\boldsymbol{z}$$

$$= \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \left[ \log p_\theta(\boldsymbol{x}|\boldsymbol{z}) + \log \frac{p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \right] d\boldsymbol{z}$$

$$= \mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})} [\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - D_{KL}\left( q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\lambda(\boldsymbol{z}) \right)$$

**Evidence Lower BOund (ELBO)**

VU

$$\ln p_\theta(\boldsymbol{x}) = \log \int p_\theta(\boldsymbol{x}|\boldsymbol{z})p_\lambda(\boldsymbol{z})d\boldsymbol{z}$$

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

$$\geq \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \log \frac{p_\theta(\boldsymbol{x}|\boldsymbol{z})p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} d\boldsymbol{z}$$

$$= \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \left[\log p_\theta(\boldsymbol{x}|\boldsymbol{z}) + \log \frac{p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\right] d\boldsymbol{z}$$

$$= \underbrace{\mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})]}_{} - \underbrace{D_{KL}\left(q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\lambda(\boldsymbol{z})\right)}_{}$$

**Reconstruction error (RE)**     **Regularization (KL)**

VU

$$\ln p_\theta(\boldsymbol{x}) = \log \int p_\theta(\boldsymbol{x}|\boldsymbol{z}) p_\lambda(\boldsymbol{z}) d\boldsymbol{z}$$

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

$$\geq \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \log \frac{p_\theta(\boldsymbol{x}|\boldsymbol{z}) p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} d\boldsymbol{z}$$

$$= \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \left[ \log p_\theta(\boldsymbol{x}|\boldsymbol{z}) + \log \frac{p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \right] d\boldsymbol{z}$$

$$= \mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})} [\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - D_{KL}\left( q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\lambda(\boldsymbol{z}) \right)$$

$$p_\theta(\boldsymbol{x}|\boldsymbol{z}) \sim \mathcal{N}(f_\theta(\boldsymbol{z}), \sigma \boldsymbol{I})$$

$$\mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})} [\log p_\theta(\boldsymbol{x}|\boldsymbol{z})]$$

$$= \mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{\|\boldsymbol{x} - f_\theta(\boldsymbol{z})\|_2^2}{2\sigma^2} \right]$$

**Reconstruction error (RE)**   **Regularization (KL)**

**VU**

$$\ln p_\theta(\boldsymbol{x}) = \log \int p_\theta(\boldsymbol{x}|\boldsymbol{z})p_\lambda(\boldsymbol{z})d\boldsymbol{z}$$

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

$$\geq \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \log \frac{p_\theta(\boldsymbol{x}|\boldsymbol{z})p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} d\boldsymbol{z}$$

$$= \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \left[ \log p_\theta(\boldsymbol{x}|\boldsymbol{z}) + \log \frac{p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \right] d\boldsymbol{z}$$

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

**decoder**

**encoder**

**marginal (prior)**

**= Variational Auto-Encoder**

VU

$$\ln p_\theta(\boldsymbol{x}) = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\ln p_\theta(\boldsymbol{x})]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{p_\theta(\boldsymbol{z}|\boldsymbol{x})p_\theta(\boldsymbol{x})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{p_\theta(\boldsymbol{x},\boldsymbol{z})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{p_\theta(\boldsymbol{x},\boldsymbol{z})q_\phi(\boldsymbol{z}|\boldsymbol{x})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})q_\phi(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{p_\theta(\boldsymbol{x},\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\right] + \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{q_\phi(\boldsymbol{z}|\boldsymbol{x})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{p_\theta(\boldsymbol{x}|\boldsymbol{z})p_\lambda(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\right] + \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{q_\phi(\boldsymbol{z}|\boldsymbol{x})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\ln p_\theta(\boldsymbol{x}|\boldsymbol{z})] - \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{q_\phi(\boldsymbol{z}|\boldsymbol{x})}{p_\lambda(\boldsymbol{z})}\right] + \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\ln \frac{q_\phi(\boldsymbol{z}|\boldsymbol{x})}{p_\theta(\boldsymbol{z}|\boldsymbol{x})}\right]$$

$$= \underbrace{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\ln p_\theta(\boldsymbol{x}|\boldsymbol{z})] - D_{KL}\big(q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\lambda(\boldsymbol{z})\big)}_{\textit{evidence lower bound } (\textbf{ELBO})} + D_{KL}\big(q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\theta(\boldsymbol{z}|\boldsymbol{x})\big)$$
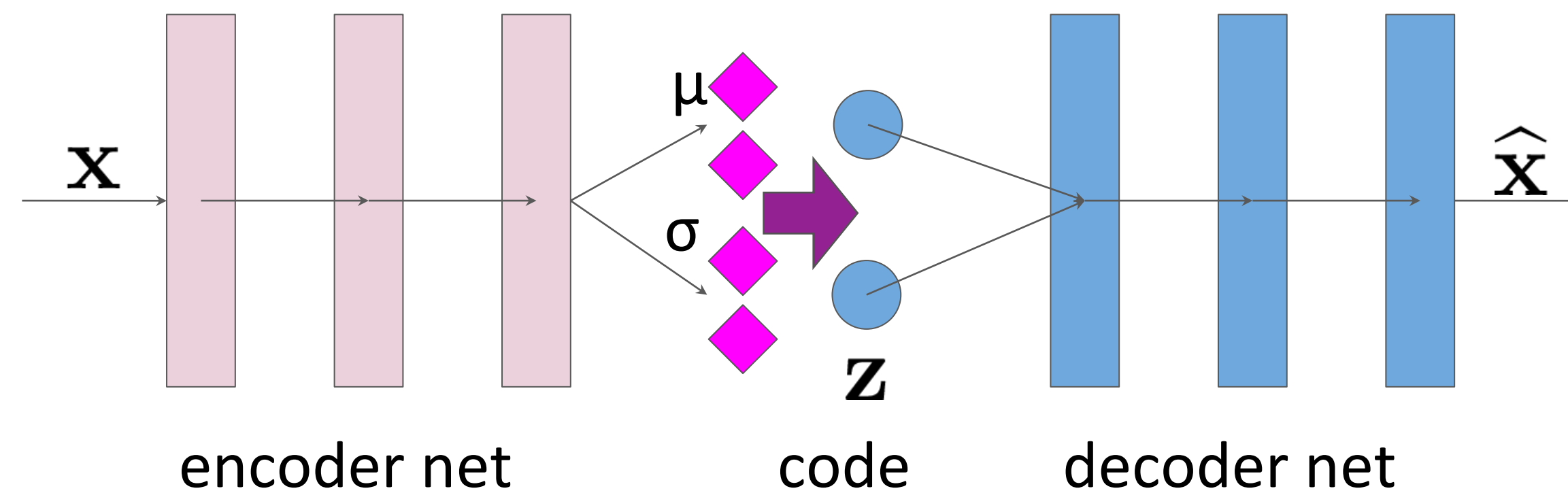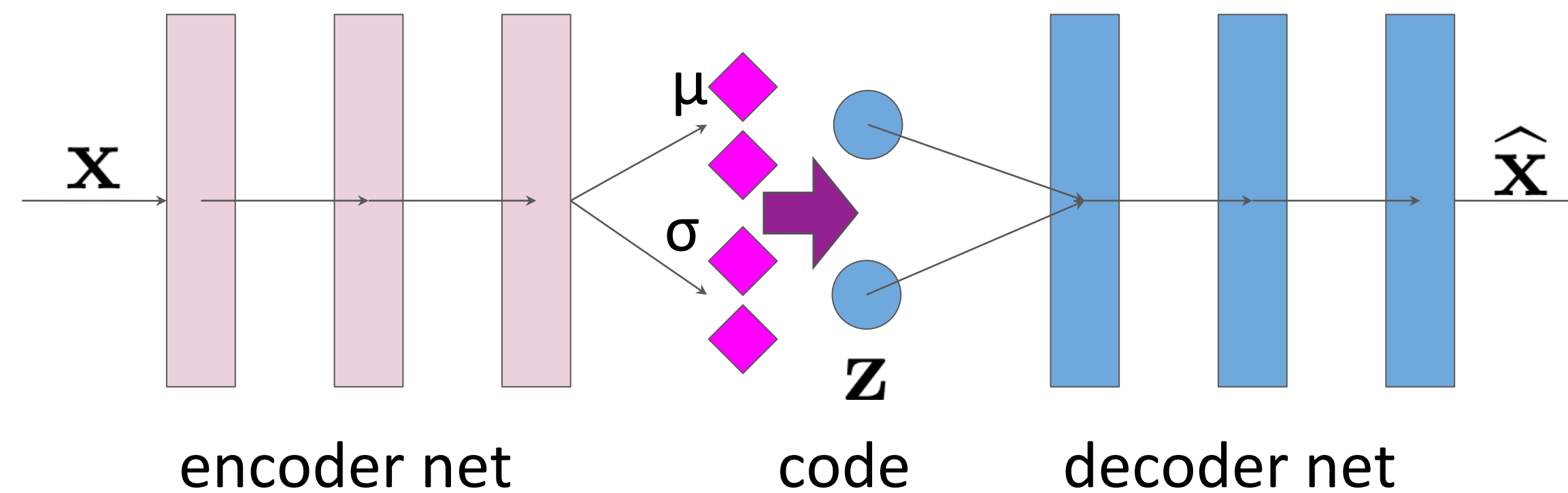
$\geq 0$

**VU** 🦅

# PART ONE: VAE IMPLEMENTATION

Variational posterior (**encoder**) and the likelihood function (**decoder**) are parameterized by neural networks.
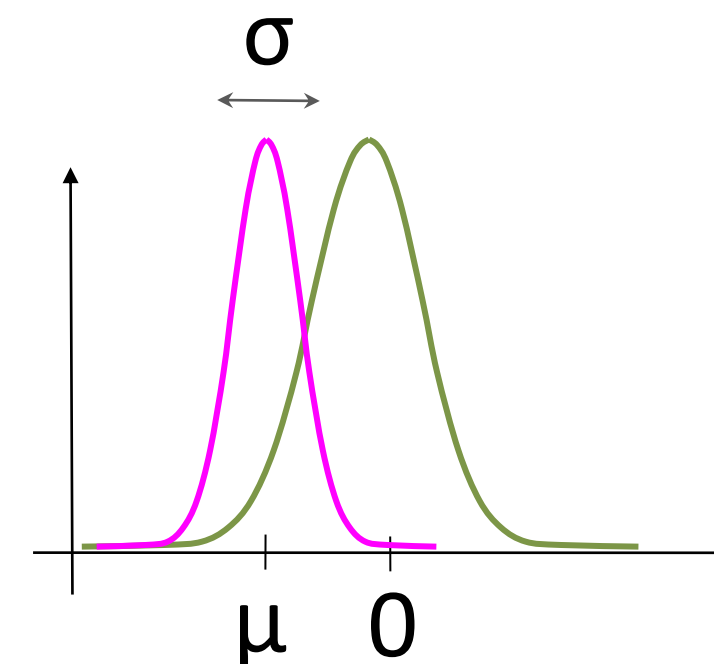


encoder net     code     decoder net

Variational posterior (**encoder**) and the likelihood function (**decoder**) are parameterized by neural networks.



encoder net      code      decoder net

***Reparameterization trick***:

$$z \sim \mathcal{N}(\mu, \sigma)$$

$$\downarrow$$

$$z = \mu + \sigma \cdot \varepsilon, \varepsilon \sim \mathcal{N}(0,1)$$

VAE copies input to output through a **bottleneck**.

VAE learns a **code** of the data.

VAE copies input to output through a **bottleneck**.

VAE learns a **code** of the data.

VAE has a **marginal** on the latent code.

VAE can **generate** new data.

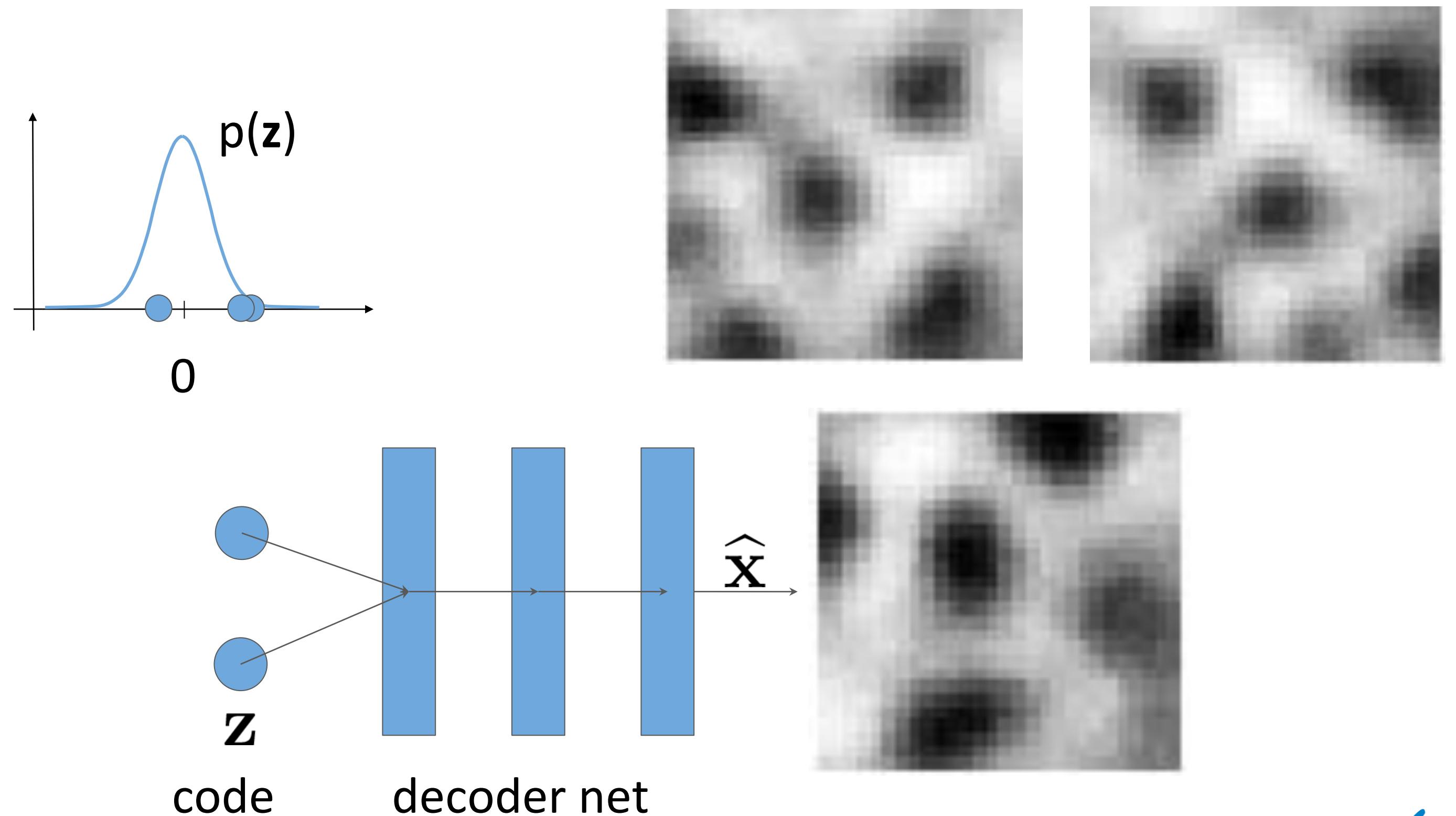VAE has a **marginal** on the latent code.

VAE can **generate** new data.

VAE has a **marginal** on the latent code.

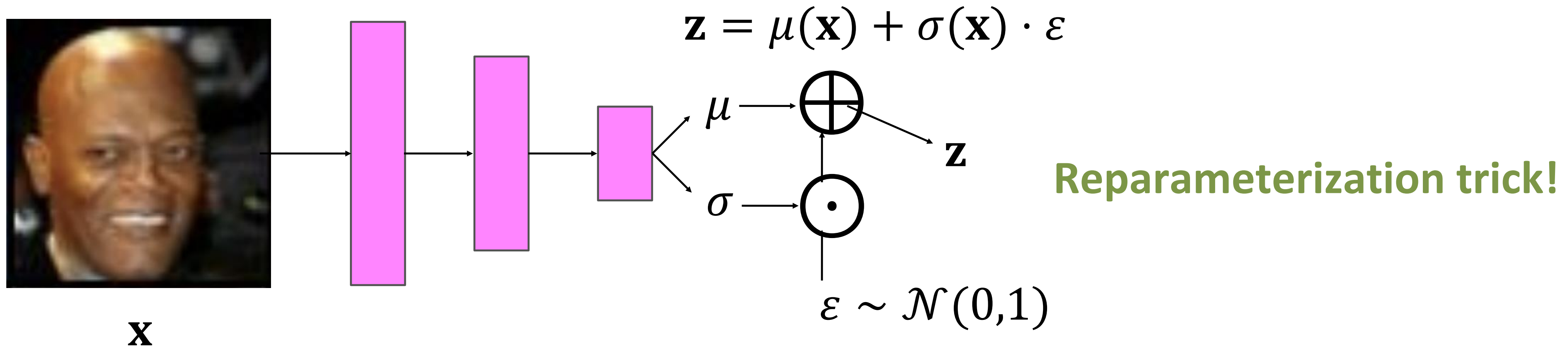VAE can **generate** new data.



p(**z**)

0

**z**

code          decoder net

$\widehat{\mathbf{x}}$

VU

$\mathbf{x}$

**x**

Example architecture for the encoder:

**x → Linear(D, 300) → ReLU → Linear(300, 2$M$) → split to 2 vectors**

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0,1)$$

**Reparameterization trick!**

Example architecture for the encoder:

**x → Linear(D, 300) → ReLU → Linear(300, 2*M*) → split to 2 vectors**

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0,1)$$

**x**

Example architecture for the encoder:

**x → Linear(D, 300) → ReLU → Linear(300, 2$M$) → split to 2 vectors**

**No non-linearity here!
We model means and log-std
for Gaussian.**

VU

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0,1)$$

$\mathbf{x}$

$\hat{\mathbf{x}} = \theta(\mathbf{z})$

Example architecture for the encoder:

**x → Linear(D, 300) → ReLU → Linear(300, 2*M*) → split to 2 vectors**

Example architecture for the decoder:

**z → Linear(*M*, 300) → ReLU → Linear(300, *D*) → means**

**No non-linearity here!**
**We model means only.**

VU

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

$\mu$

$\sigma$

$\mathbf{z}$

$\varepsilon \sim \mathcal{N}(0,1)$
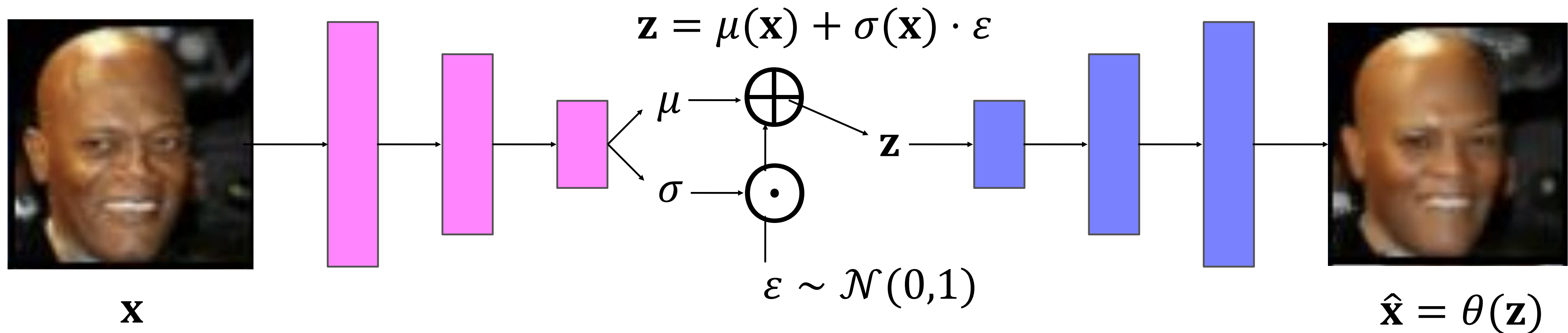
$\mathbf{x}$

$\hat{\mathbf{x}} = \theta(\mathbf{z})$

We approximate expected values using a single sample:

$$ELBO = \ln\mathcal{N}(\mathbf{x}|\theta(\mathbf{z}), 1) - [\ln\mathcal{N}(\mathbf{z}|\mu(\mathbf{x}), \sigma^2(\mathbf{x})) - \ln\mathcal{N}(\mathbf{z}|0,1)]$$

$$p_\theta(\mathbf{x}|\mathbf{z}) \qquad q_\phi(\mathbf{z}|\mathbf{x}) \qquad p_\lambda(\mathbf{z})$$

VU

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

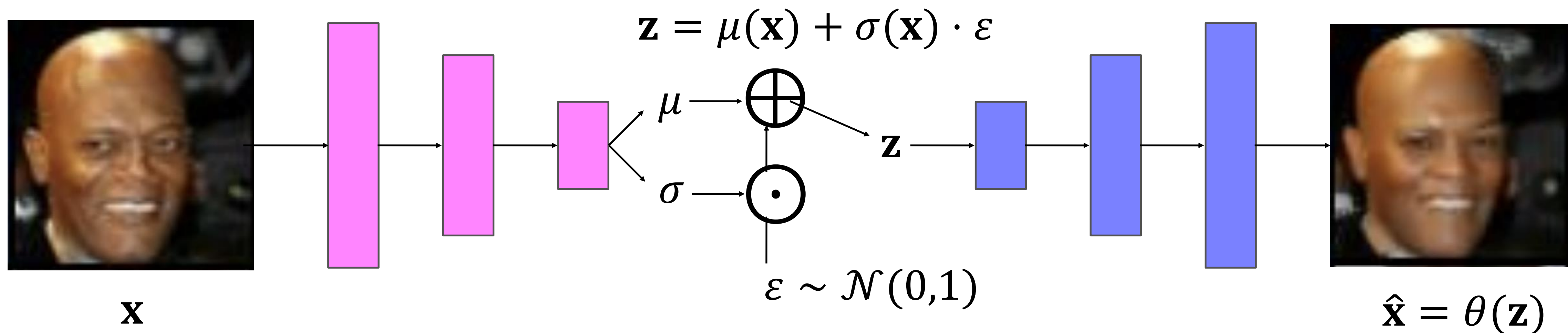$$\varepsilon \sim \mathcal{N}(0,1)$$

$\mathbf{x}$

$\hat{\mathbf{x}} = \theta(\mathbf{z})$

We approximate expected values using a single sample:

$$ELBO = \ln\mathcal{N}(\mathbf{x}|\theta(\mathbf{z}), 1) - [\ln\mathcal{N}(\mathbf{z}|\mu(\mathbf{x}), \sigma^2(\mathbf{x})) - \ln\mathcal{N}(\mathbf{z}|0,1)]$$

RE

KL

VU

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \cdot \varepsilon$$

$$\varepsilon \sim \mathcal{N}(0,1)$$

$\mathbf{x}$

$\hat{\mathbf{x}} = \theta(\mathbf{z})$

We approximate expected values using a single sample:

**We assume a Gaussian variational posterior.**

$$ELBO = \ln\mathcal{N}(\mathbf{x}|\theta(\mathbf{z}), 1) - [\ln\mathcal{N}(\mathbf{z}|\mu(\mathbf{x}), \sigma^2(\mathbf{x})) - \ln\mathcal{N}(\mathbf{z}|0,1)]$$

RE

KL

**We assume a standard Gaussian prior.**

VU

```python
import torch.nn as nn

class VAE(nn.Module):
    def __init__(self, D, M):
        super(LinearVAE, self).__init__()
        self.D = D
        self.M = M

        self.enc1 = nn.Linear(in_features=self.D, out_features=300)
        self.enc2 = nn.Linear(in_features=300, out_features=self.M*2)

        self.dec1 = nn.Linear(in_features=self.M, out_features=300)
        self.dec2 = nn.Linear(in_features=300, out_features=self.D)

    def reparameterize(self, mu, log_std):
        std = torch.exp(log_std)
        eps = torch.randn_like(std)
        Z = mu + (eps * std)
        return z
```
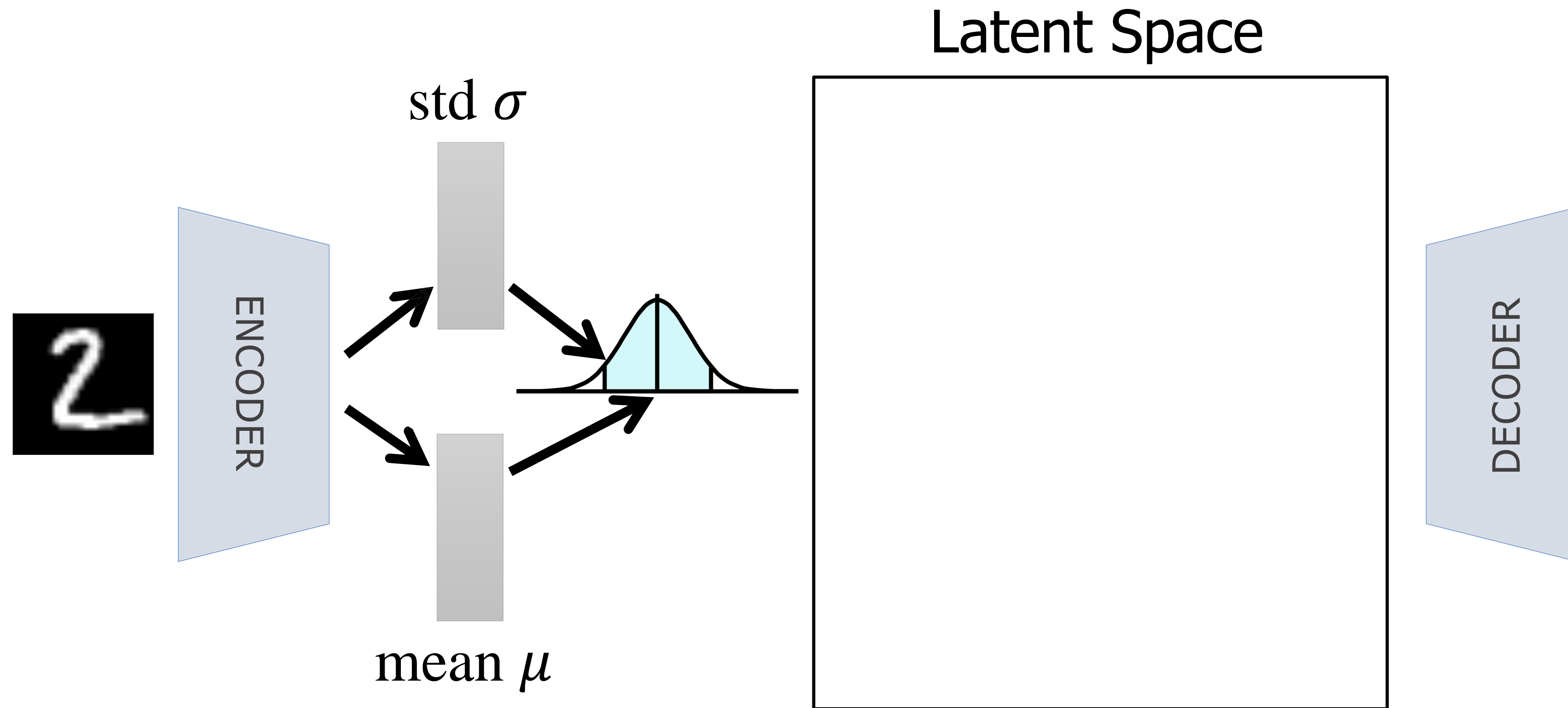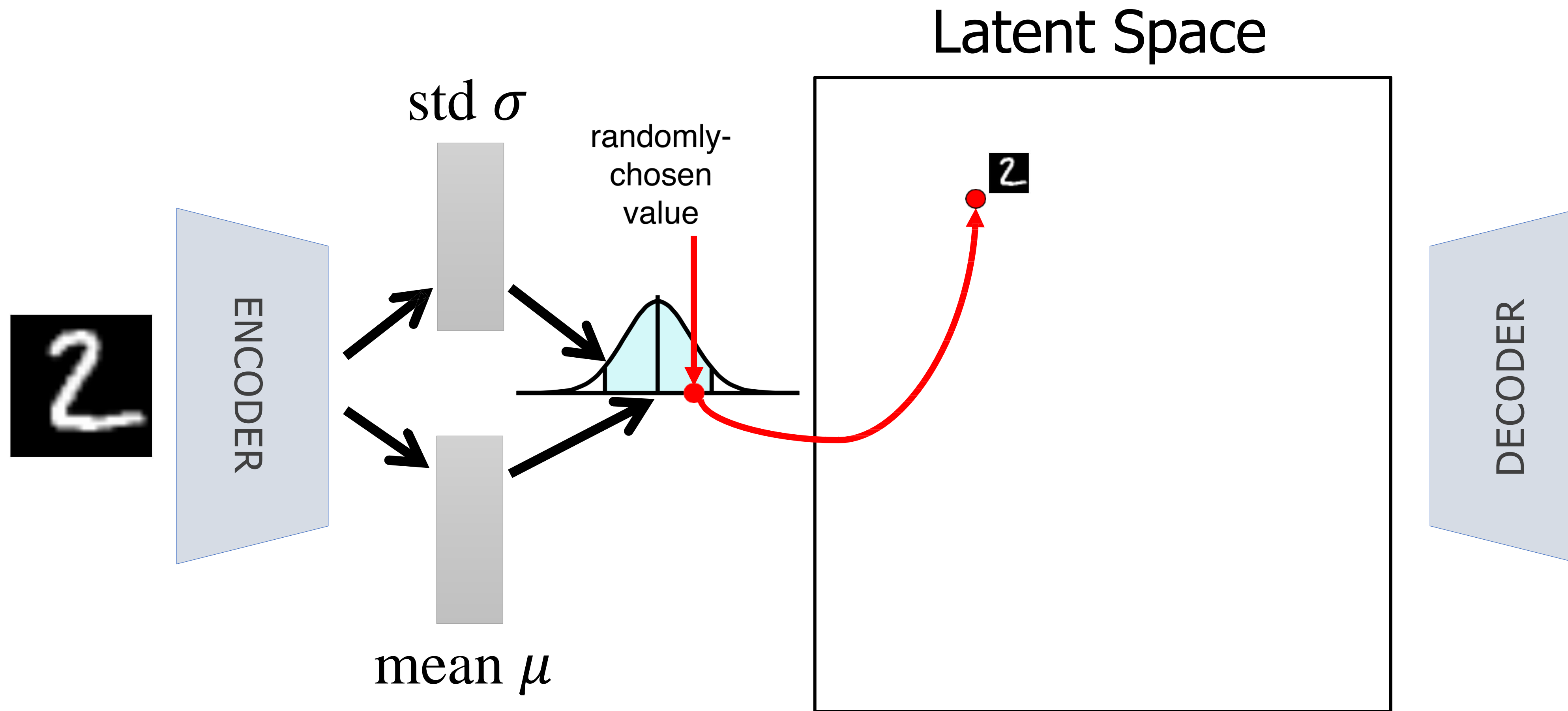
VU

```python
def forward(self, x):
    # encoder
    x = nn.functional.relu(self.enc1(x))
    x = self.enc2(x).view(-1, 2, self.M)

    # get mean and log-std
    mu = x[:, 0, :]
    log_var = x[:, 1, :]

    # reparameterization
    z = self.reparameterize(mu, log_std)

    # decoder
    x_hat = nn.functional.relu(self.dec1(z))
    x_hat = self.dec2(x)
    return x_hat, mu, log_std
```

```python
def elbo(self, x, x_hat, z, mu, log_std):
    # reconstruction error
    RE = nn.loss.mse(x, x_hat)

    # kl-regularization
    # We assume here that log_normal is implemented
    KL = log_normal(z, mu, log_std) - log_normal(z, 0, 1)

    # REMEMBER! We maximize ELBO, but optimizers minimize.
    # Therefore, we need to take the negative sign!
    return -(RE - KL)
```

VU

Latent Space

std $\sigma$

ENCODER

mean $\mu$

DECODER

Encode the first sample (a "2") and find $\mu_1, \sigma_1$

VU

## Latent Space

std $\sigma$

randomly-chosen value

ENCODER

DECODER

mean $\mu$

Sample $\boldsymbol{z}_1 \sim N(\mu_1, \sigma_1)$

VU

std $\sigma$

randomly-chosen value

Latent Space

ENCODER

mean $\mu$

DECODER

Denote to $\widehat{\boldsymbol{x}}_1$

VU

Latent Space

std $\sigma$

randomly-chosen value

mean $\mu$

ENCODER
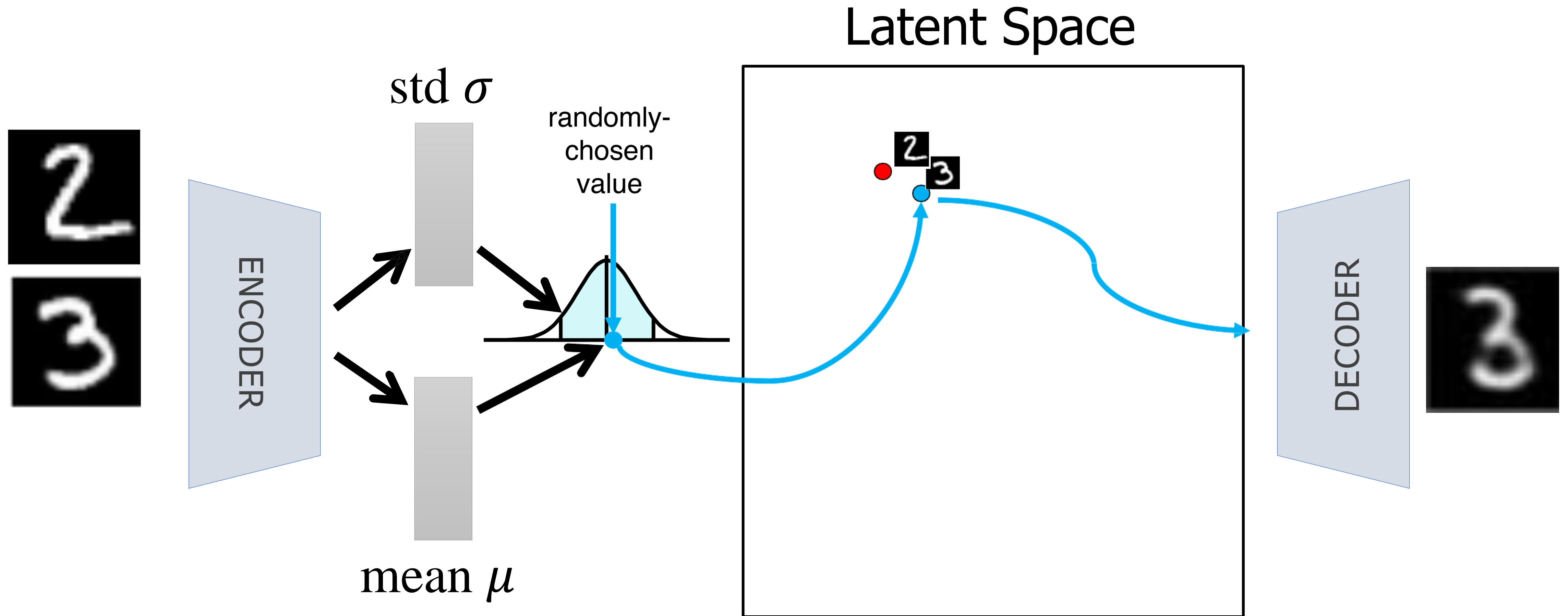
DECODER

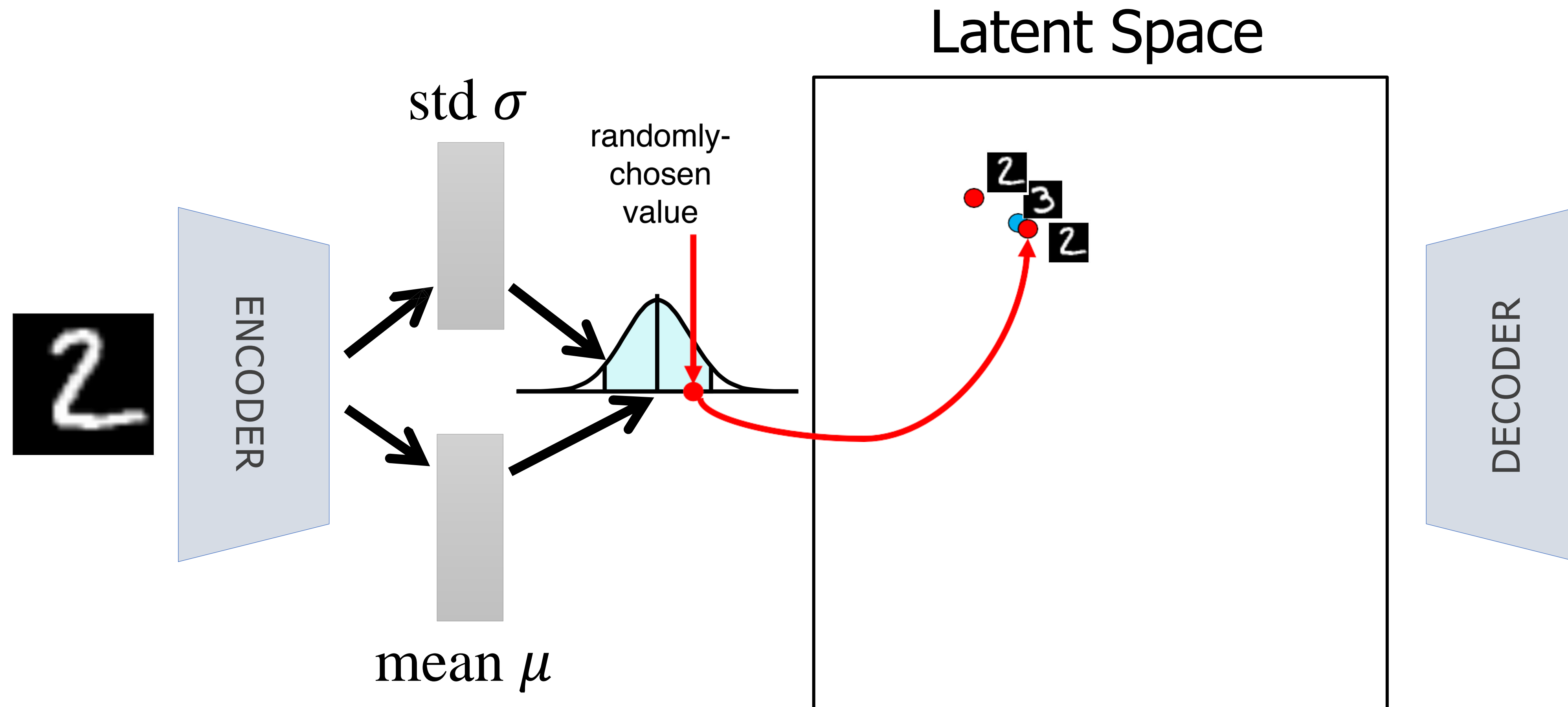Encode the first sample (a "3") and find $\mu_2, \sigma_2$, and sample $z_2 \sim N(\mu_2, \sigma_2)$

VU

Latent Space

std $\sigma$

randomly-chosen value

ENCODER

mean $\mu$

DECODER

Decode to $\widehat{x}_2$

VU

Latent Space

std $\sigma$

randomly-
chosen
value
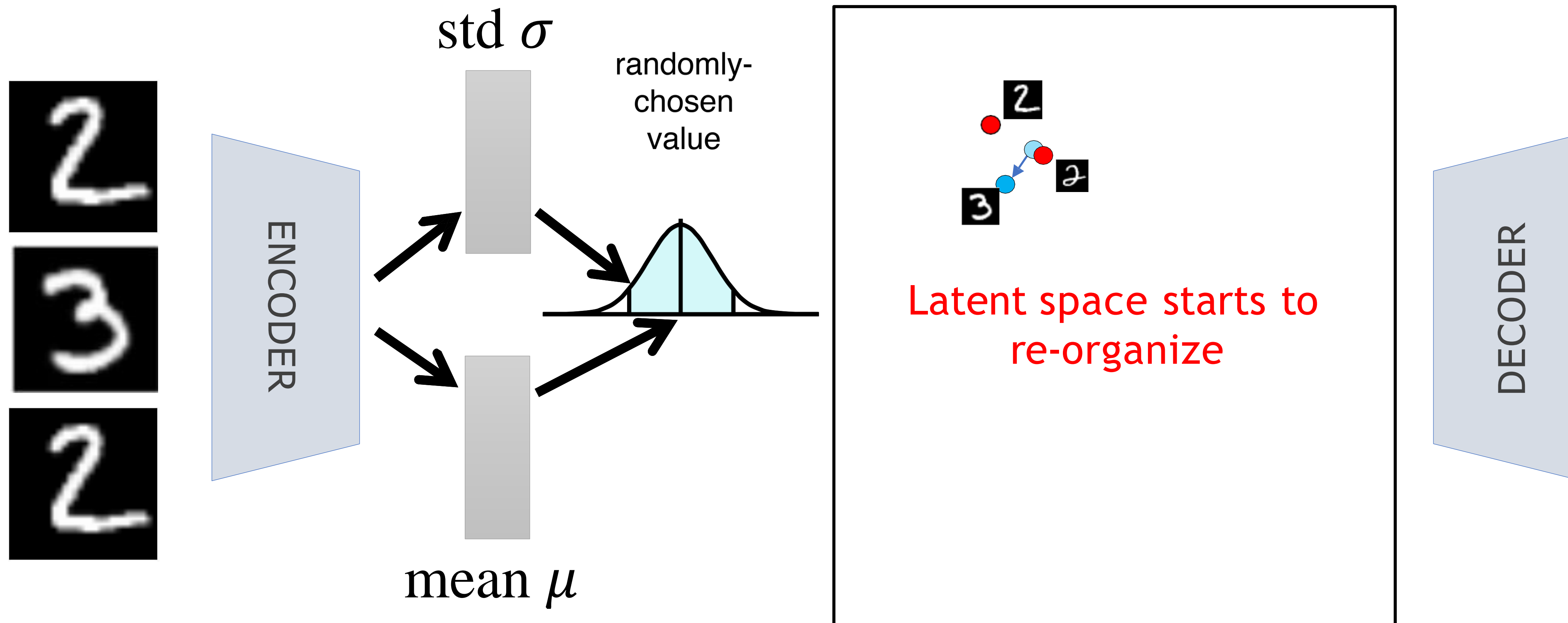
ENCODER

mean $\mu$

DECODER

Train with the first sample (a "2") again and find $\mu_1, \sigma_1$. However, $\boldsymbol{z}_1 \sim N(\mu_1, \sigma_1)$ will not be the same. It can happen to be close to the "3" in latent space.

VU

## Latent Space

std $\sigma$

randomly-
chosen
value

ENCODER

DECODER

mean $\mu$

Decode to $\hat{x}_1$. Since the decoder only knows how to map from latent space to $\hat{x}$ space, it will return a "3".

VU

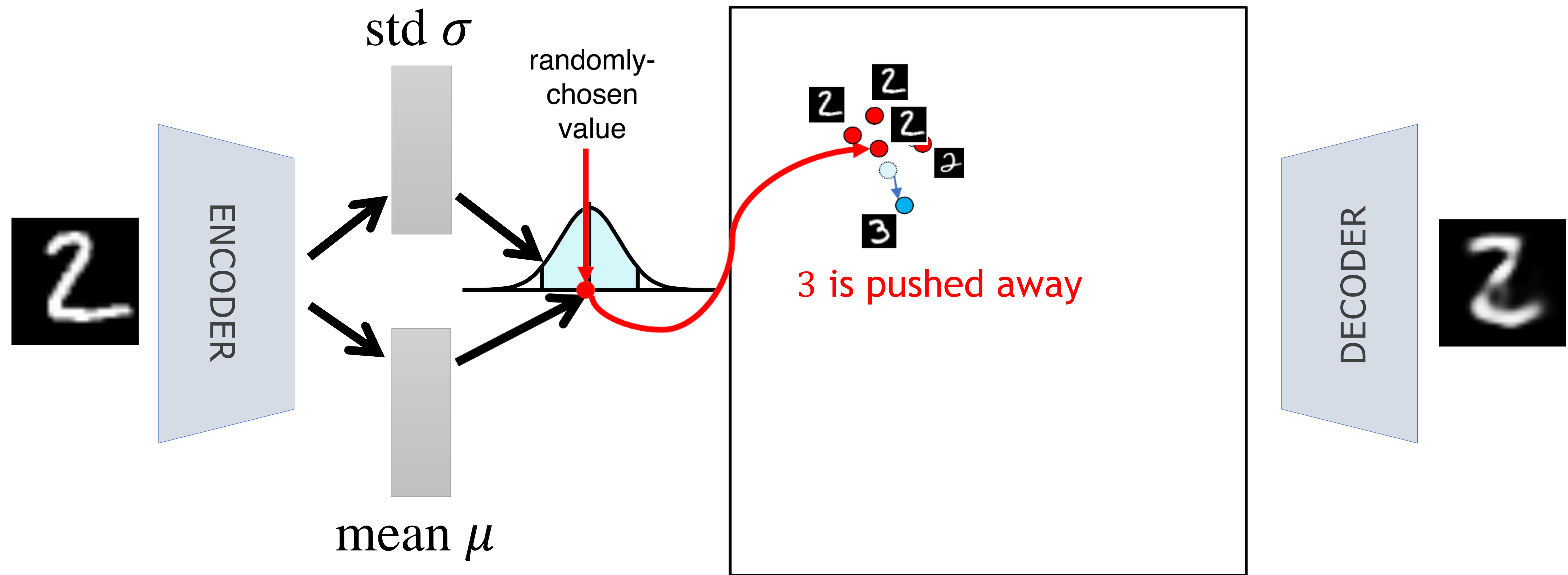Train with 1$^{st}$ sample again

std $\sigma$
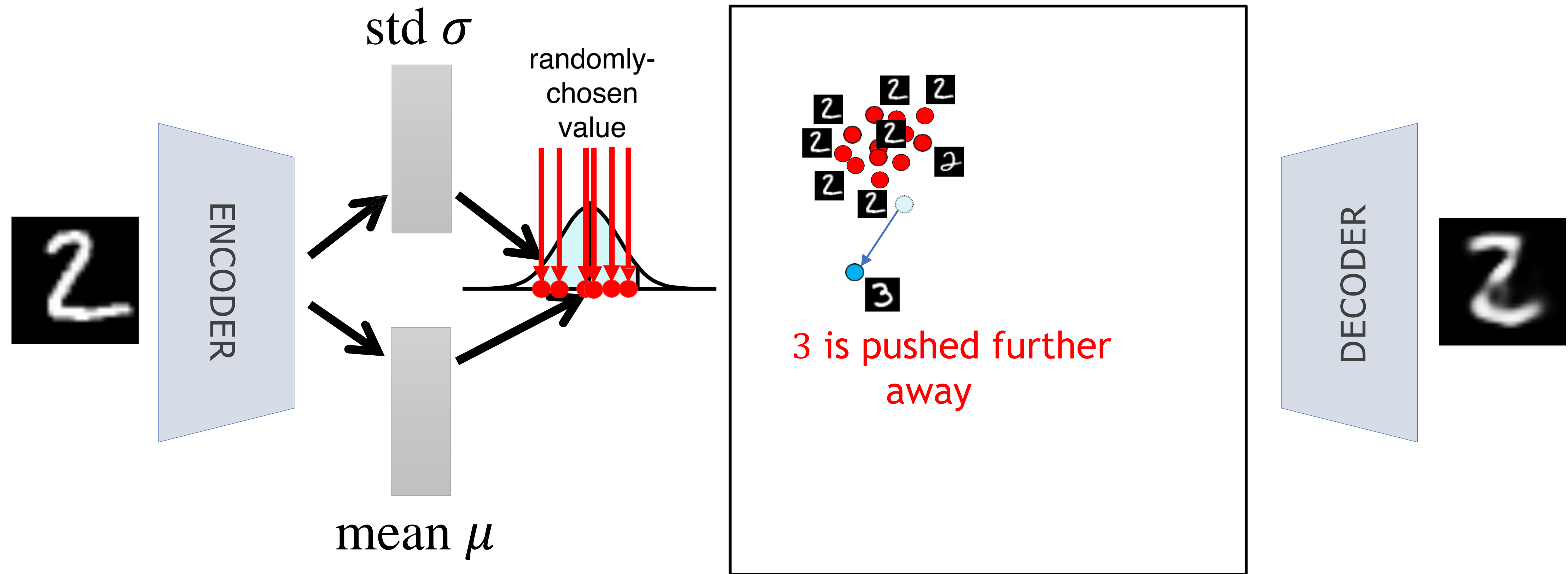
randomly-chosen value

ENCODER

mean $\mu$

## Latent Space

Latent space starts to re-organize

DECODER

VU

And again …

Latent Space

std $\sigma$

randomly-chosen value

mean $\mu$

ENCODER

DECODER

3 is pushed away

VU

Many times

Latent Space

std $\sigma$

randomly-chosen value

mean $\mu$

ENCODER

DECODER

3 is pushed further away

Now, let's test again

std $\sigma$

randomly-chosen value

mean $\mu$

## Latent Space

ENCODER

DECODER

Try on 3's again

std $\sigma$

randomly-chosen value

Latent Space

mean $\mu$

ENCODER

DECODER

VU

Many times …

std $\sigma$

randomly-chosen value

Latent Space

ENCODER
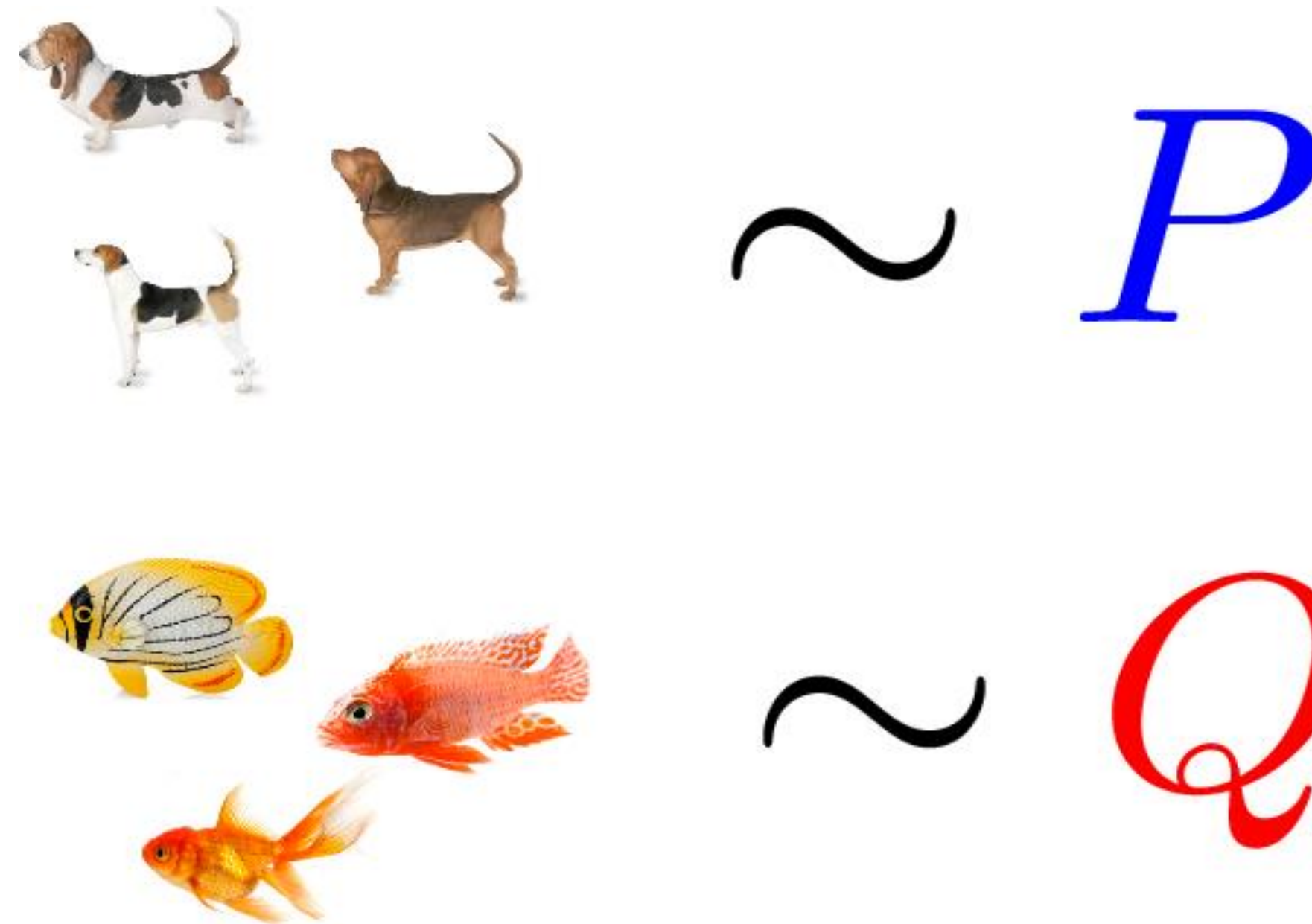
mean $\mu$

DECODER

VU

# PART TWO: KL DIVERGENCE AND MAXIMUM MEAN DISCREPANCY

How to measure the distance/divergence?

VU

**Given**: samples from unknown distributions $P$ and $Q$

**Goal**: do $P$ and $Q$ differ?

$$\sim P$$

$$\sim Q$$

VU

**Given**: samples from unknown distributions *P* and *Q*

**Goal**: do *P* and *Q* differ?



real MNIST samples



generated samples from VAE

Are there significant difference in VAE and MNIST?

**Given**: A model *P* and samples from *Q*

**Goal**: is one Gaussian or two Gaussians more fit for *Q*?

**Given**: samples from a (joint) distribution $P_{X,Y}$

**Goal**: are $X$ and $Y$ independent? If not, the strength of dependence?



independent     $p(x, y) = p(x)p(y)$

dependent     $p(x, y) \neq p(x)p(y)$

$$I = D_{KL}(p(x, y); p(x)p(y))$$

$$P - Q$$

$$\frac{P}{Q}$$

VU

Integral probability metrics

$f$-divergence

$$D_{\mathcal{H}}(P; Q)$$
$$= \sup_{g \in \mathcal{H}} \left| \mathbb{E}_{X \sim P} g(X) - \mathbb{E}_{Y \sim Q} g(Y) \right|$$

**maximum mean discrepancy (MMD)**
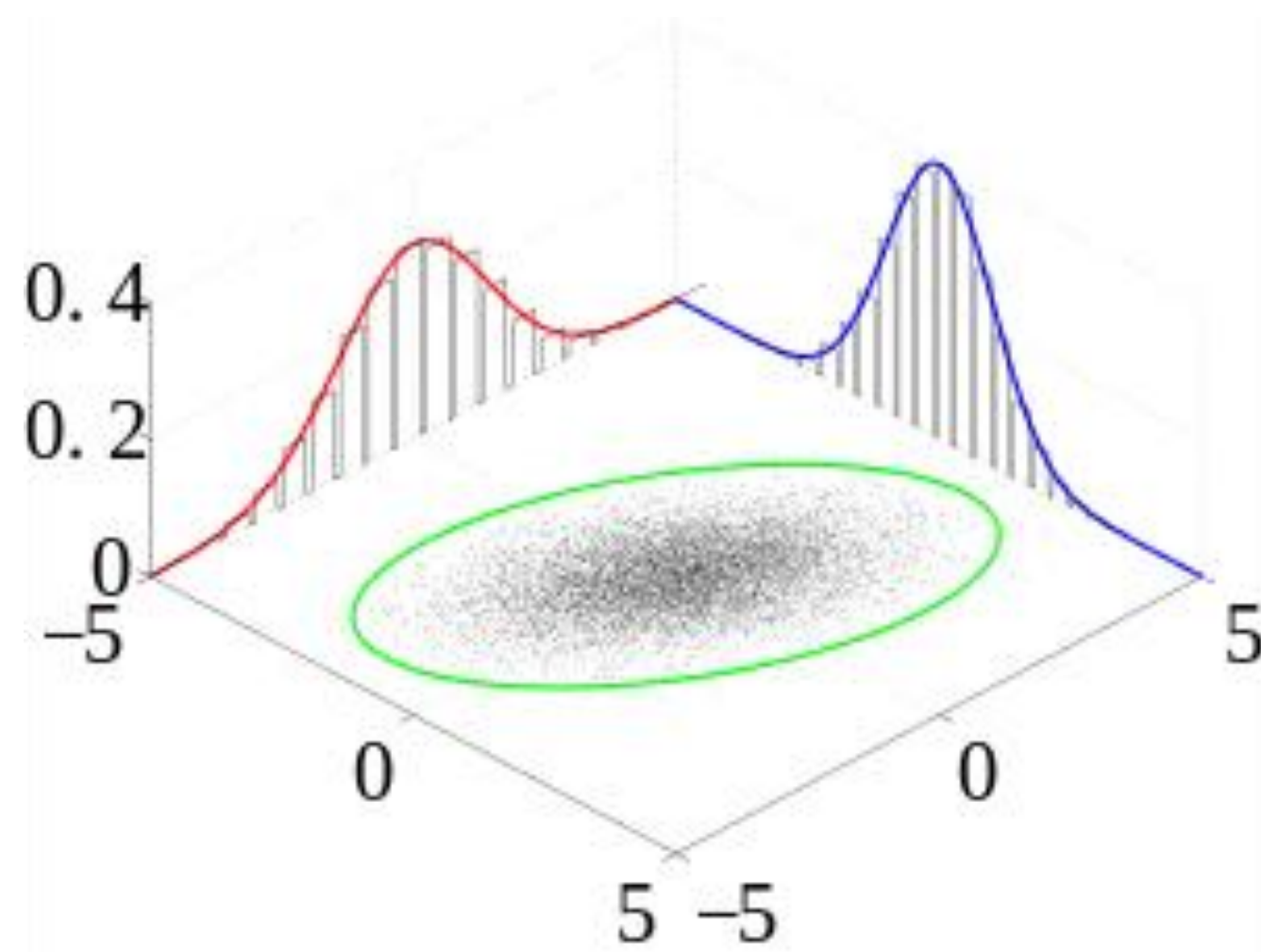
$$D_{KL}(P; Q)$$
$$= \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

**Kullback–Leibler (KL) divergence**

Gretton, Arthur, et al. "A kernel two-sample test." *The Journal of Machine Learning Research* 13.1 (2012): 723-773. https://www.jmlr.org/papers/volume13/gretton12a/gretton12a.pdf

VU

# Two Gaussians with different means

Two Gaussians with same mean but different variances

Idea: look at difference in means of features of the random variables

In Gaussian case: second order features of form $\varphi(x) = x^2$

Gaussian and Laplacian distributions

Same mean *and* same variance

Difference in means using <span style="color:red">higher order features</span>



Gaussian and Laplacian densities

For a feature map $\varphi\colon \mathcal{X} \to \mathcal{H}$, representing distances between distributions as distances between mean embeddings of features



$$\mathrm{MMD}^2(P; Q) = \left\| \mathbb{E}_{X\sim P}\,\varphi(X) - \mathbb{E}_{Y\sim Q}\,\varphi(Y) \right\|_{\mathcal{H}}^2$$

Muandet, Krikamol, et al. "Kernel mean embedding of distributions: A review and beyond." *Foundations and Trends® in Machine Learning* 10.1-2 (2017): 1-141. https://www.nowpublishers.com/article/Details/MAL-060

For a feature map $\varphi: \mathcal{X} \to \mathcal{H}$, representing distances between distributions as distances between mean embeddings of features



$p(\mathbf{x})$

$\mathbb{P}$

$\mathbb{Q}$

$\mathbf{x}$

Reproducing Kernel Hilbert Space

RKHS embedding of $\mathbb{Q}$

RKHS embedding of $\mathbb{P}$

$$\text{MMD}^2(P; Q) = \left\| \mathbb{E}_{X \sim P} \varphi(X) - \mathbb{E}_{Y \sim Q} \varphi(Y) \right\|_{\mathcal{H}}^2$$

$$= \left\langle \mathbb{E}_{X \sim P} \varphi(X), \mathbb{E}_{X' \sim P} \varphi(X') \right\rangle_{\mathcal{H}} + \left\langle \mathbb{E}_{Y \sim Q} \varphi(Y), \mathbb{E}_{Y' \sim Q} \varphi(Y') \right\rangle_{\mathcal{H}}$$

$$-2 \left\langle \mathbb{E}_{X \sim P} \varphi(X), \mathbb{E}_{Y \sim Q} \varphi(Y) \right\rangle_{\mathcal{H}} \qquad (x - y)^2 = x^T x + y^T y - 2x^T y$$

$$= \mathbb{E}_{X,X' \sim P} \kappa(X, X') + \mathbb{E}_{Y,Y' \sim Q} \kappa(Y, Y') - 2\mathbb{E}_{X \sim P, Y \sim Q} \kappa(X, Y)$$

The kernel trick: $\kappa(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$

VU

For a feature map $\varphi\colon \mathcal{X} \to \mathcal{H}$, representing distances between distributions as distances between mean embeddings of features

Can be inner product in <span style="color:red">infinite</span> dimensional space

Assume $x \in R^1$ and $\gamma > 0$.

$$e^{-\gamma\|x_i - x_j\|^2} = e^{-\gamma(x_i - x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2}$$

$$= e^{-\gamma x_i^2 - \gamma x_j^2}\left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \cdots\right)$$

$$= e^{-\gamma x_i^2 - \gamma x_j^2}\left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}}x_i \cdot \sqrt{\frac{2\gamma}{1!}}x_j + \sqrt{\frac{(2\gamma)^2}{2!}}x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}}x_j^2\right.$$

$$\left. + \sqrt{\frac{(2\gamma)^3}{3!}}x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}}x_j^3 + \cdots\right) = \phi(x_i)^T \phi(x_j),$$

where

$$\phi(x) = e^{-\gamma x^2}\left[1, \sqrt{\frac{2\gamma}{1!}}x, \sqrt{\frac{(2\gamma)^2}{2!}}x^2, \sqrt{\frac{(2\gamma)^3}{3!}}x^3, \cdots\right]^T$$

VU

For a feature map $\varphi\colon \mathcal{X} \to \mathcal{H}$, representing distances between distributions as distances between mean embeddings of features

$$\text{MMD}^2(P; Q) = \mathbb{E}_{X,X'\sim P}\kappa(X, X') + \mathbb{E}_{Y,Y'\sim Q}\kappa(Y, Y') - 2\mathbb{E}_{X\sim P, Y\sim Q}\kappa(X, Y)$$

$$\widehat{\text{MMD}}^2(P; Q) = \underbrace{\frac{1}{N^2}\sum_{i=1}^{N}\sum_{j=1}^{N}G_\sigma(x_i - x_j)}_{\textbf{within distribution similarity}} + \underbrace{\frac{1}{M^2}\sum_{i=1}^{M}\sum_{j=1}^{M}G_\sigma(y_i - y_j)}_{\textbf{within distribution similarity}} - \underbrace{\frac{2}{NM}\sum_{i=1}^{N}\sum_{j=1}^{M}G_\sigma(x_i - y_j)}_{\textbf{cross-distribution similarity}}$$

Gretton, Arthur, et al. "A kernel two-sample test." *The Journal of Machine Learning Research* 13.1 (2012): 723-773.

VU

**PART THREE: MMD-VAE**

Uninformative latent code

- $D_{KL}\left(q_\phi(\boldsymbol{z}|\boldsymbol{x}); p_\lambda(\boldsymbol{z})\right)$ might be too restrictive

- If the decoder is sufficiently flexible, failed to learn meaningful representation

A poor prior distribution $p_\lambda(\boldsymbol{z})$

Which divergence to choose?

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

MMD-VAE objective:

$$\mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})] - \text{MMD}\left(q_\phi(\boldsymbol{z}); p_\lambda(\boldsymbol{z})\right)$$

Zhao, Shengjia, Jiaming Song, and Stefano Ermon. "Infovae: Information maximizing variational autoencoders." *arXiv preprint arXiv:1706.02262* (2017). https://arxiv.org/abs/1706.02262

VU

# Which divergence to choose?



VAE

MMD-VAE

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

$$D_{KL}\left(q_\phi(\mathbf{z}|\mathbf{x}); p_\lambda(\mathbf{z})\right)$$

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

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{\color{blue}q_\phi(\mathbf{z})} \frac{\color{blue}q_\phi(\mathbf{z})}{\color{red}\prod_j q_\phi(\mathbf{z}_j)} \frac{\color{red}\prod_j q_\phi(\mathbf{z}_j)}{p_\lambda(\mathbf{z})} d\mathbf{z}$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{\color{blue}q_\phi(\mathbf{z})} \frac{\color{blue}q_\phi(\mathbf{z})}{\color{red}\prod_j q_\phi(\mathbf{z}_j)} \frac{\color{red}\prod_j q_\phi(\mathbf{z}_j)}{\color{purple}\prod_j p_\lambda(\mathbf{z}_j)} d\mathbf{z}$$

$$= D_{KL}\left(q_\phi(\mathbf{z}|\mathbf{x}); {\color{blue}q_\phi(\mathbf{z})}\right) + D_{KL}\left({\color{blue}q_\phi(\mathbf{z})}; {\color{red}\prod_j q_\phi(\mathbf{z}_j)}\right) + D_{KL}\left({\color{red}\prod_j q_\phi(\mathbf{z}_j)}; {\color{purple}\prod_j p_\lambda(\mathbf{z}_j)}\right)$$

$$= I(\mathbf{z}; \mathbf{x}) + \underbrace{D_{KL}\left({\color{blue}q_\phi(\mathbf{z})}; {\color{red}\prod_j q_\phi(\mathbf{z}_j)}\right)}_{\textbf{Total correlation}} + \underbrace{\sum_{j=1}^d D_{KL}\left({\color{red}q_\phi(\mathbf{z}_j)}; {\color{purple}p_\lambda(\mathbf{z}_j)}\right)}_{\substack{\textbf{dimension-wise} \\ \textbf{KL divergence}}}$$

VU

$$D_{KL}\left(q_\phi(\mathbf{z}|\mathbf{x}); p_\lambda(\mathbf{z})\right)$$

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

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

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

$$= D_{KL}\left(q_\phi(\mathbf{z}|\mathbf{x}); q_\phi(\mathbf{z})\right) + D_{KL}\left(q_\phi(\mathbf{z}); \prod_j q_\phi(\mathbf{z}_j)\right) + D_{KL}\left(\prod_j q_\phi(\mathbf{z}_j); \prod_j p_\lambda(\mathbf{z}_j)\right)$$

$$= I(\mathbf{z}; \mathbf{x}) + D_{KL}\left(q_\phi(\mathbf{z}); \prod_j q_\phi(\mathbf{z}_j)\right) + \sum_{j=1}^{d} D_{KL}\left(q_\phi(\mathbf{z}_j); p_\lambda(\mathbf{z}_j)\right)$$

**independence between each dimension of latent codes**

**Total correlation (TC)**

**dimension-wise KL divergence**

VU

A $\beta$-VAE optimizes the following function:

$$\mathbb{E}_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})] - \beta D_{KL}\left(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}); p_{\lambda}(\boldsymbol{z})\right)$$

$$\mathbb{E}_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})] - \beta \left\{ \underline{I(\boldsymbol{z}; \boldsymbol{x})} + \underline{TC(\boldsymbol{z})} + \underline{\sum_{j=1}^{d} D_{KL}\left(q_{\phi}(\boldsymbol{z}_j); p_{\lambda}(\boldsymbol{z}_j)\right)}\right\}$$

<span style="color:blue">Minimality</span>

<span style="color:green">closeness to prior distribution</span>

<span style="color:orange">Disentanglement</span>

Assuming a factorized prior for $\boldsymbol{z}$, a $\beta$-VAE optimizes both for the information bottleneck (IB) Lagrangian and for disentanglement.

Higgins, Irina, et al. "beta-vae: Learning basic visual concepts with a constrained variational framework." *International conference on learning representations*. 2016. https://openreview.net/forum?id=Sy2fzU9gl
Burgess, Christopher P., et al. "Understanding disentangling in $\beta$-VAE." *arXiv preprint arXiv:1804.03599* (2018). https://arxiv.org/abs/1804.03599

VU

Start with very high $\beta$ and slowly decrease during training.
Beginning: Very strict bottleneck, only encode most important factor
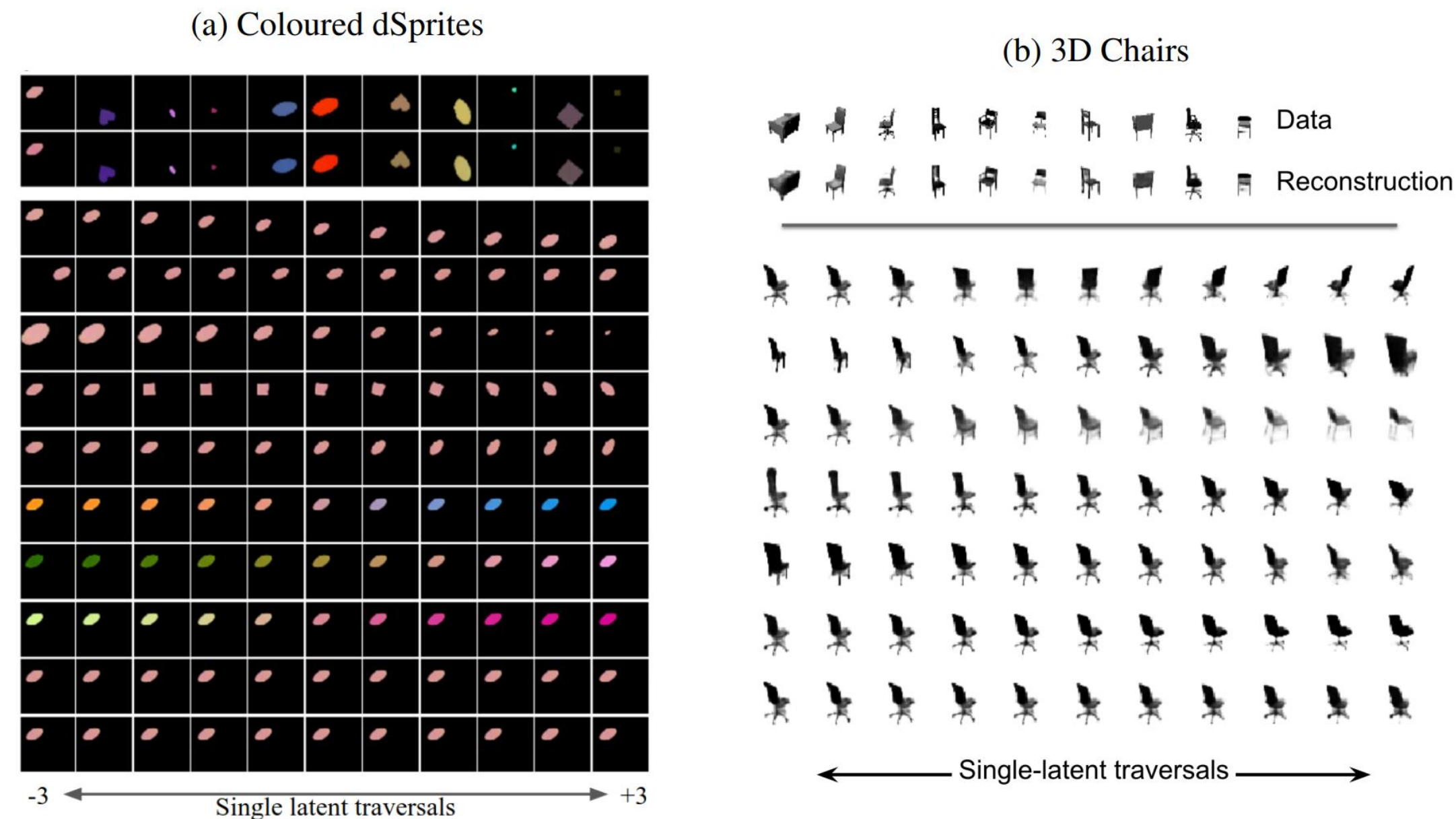End: Very large bottleneck, encode all remaining factors



Figure 4: **Disentangling and reconstructions from $\beta$-VAE with controlled capacity increase. (a)** Latent traversal plots for a $\beta$-VAE trained with controlled capacity increaseon the coloured dSprites dataset. The top two rows show data samples and corresponding reconstructions. Subsequent rows show single latent traversals, ordered by their average KL divergence with the prior (high to low). To generate the traversals, we initialise the latent representation by inferring it from a seed image (left data sample), then traverse a single latent dimension (in $[-3, 3]$), whilst holding the remaining latent dimensions fixed, and plot the resulting reconstruction. The corresponding reconstructions are the rows of this figure. The disentangling is evident: different latent dimensions independently code for position, size, shape, rotation, and colour. **(b)** Latent traversal plots, as in (a), but trained on the Chairs dataset [3].