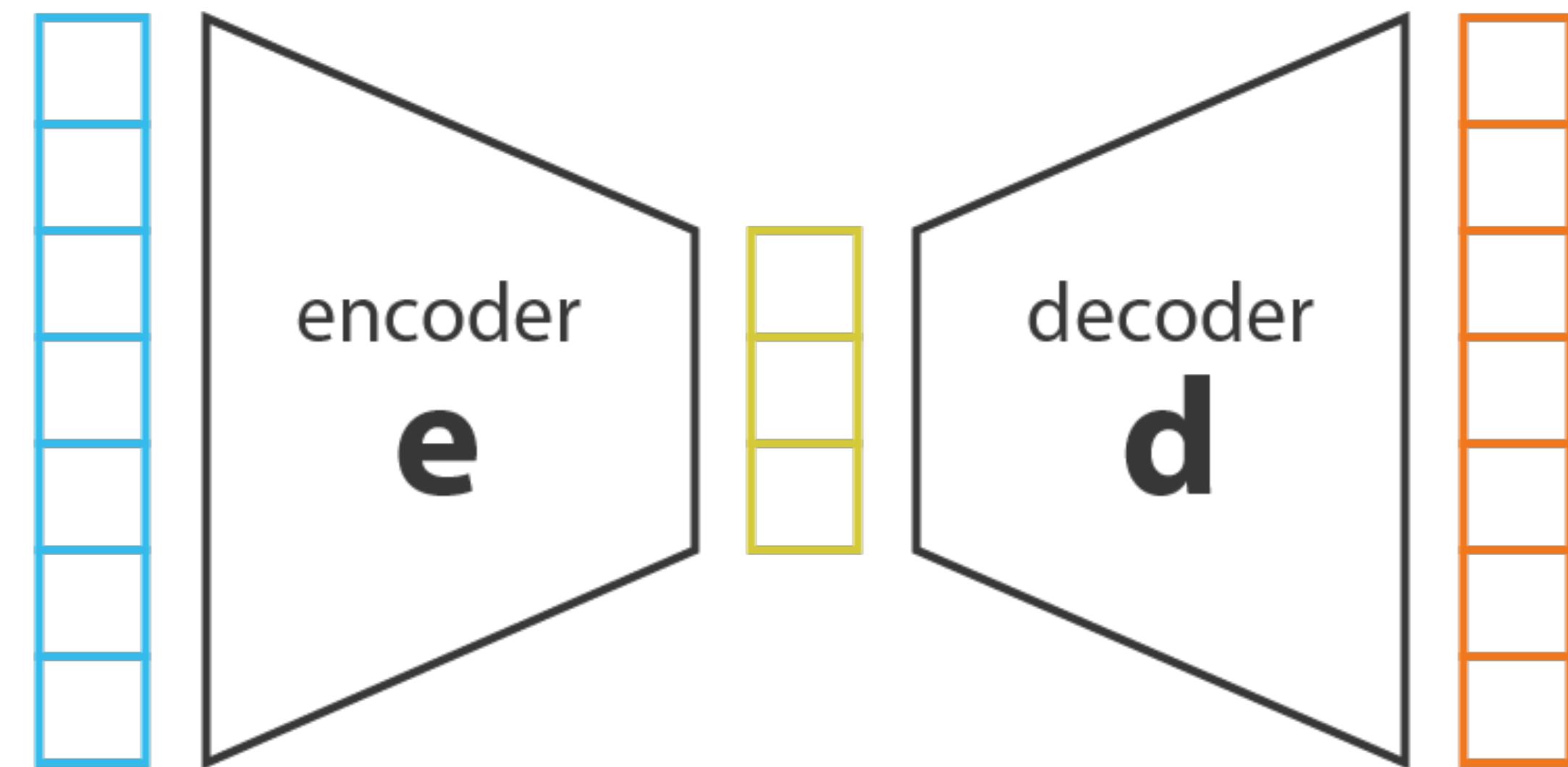


# **Variational Autoencoder**

**Dalin Guo, Kuei-Da Liao, Corey Zhou**

# Encoder and Decoder

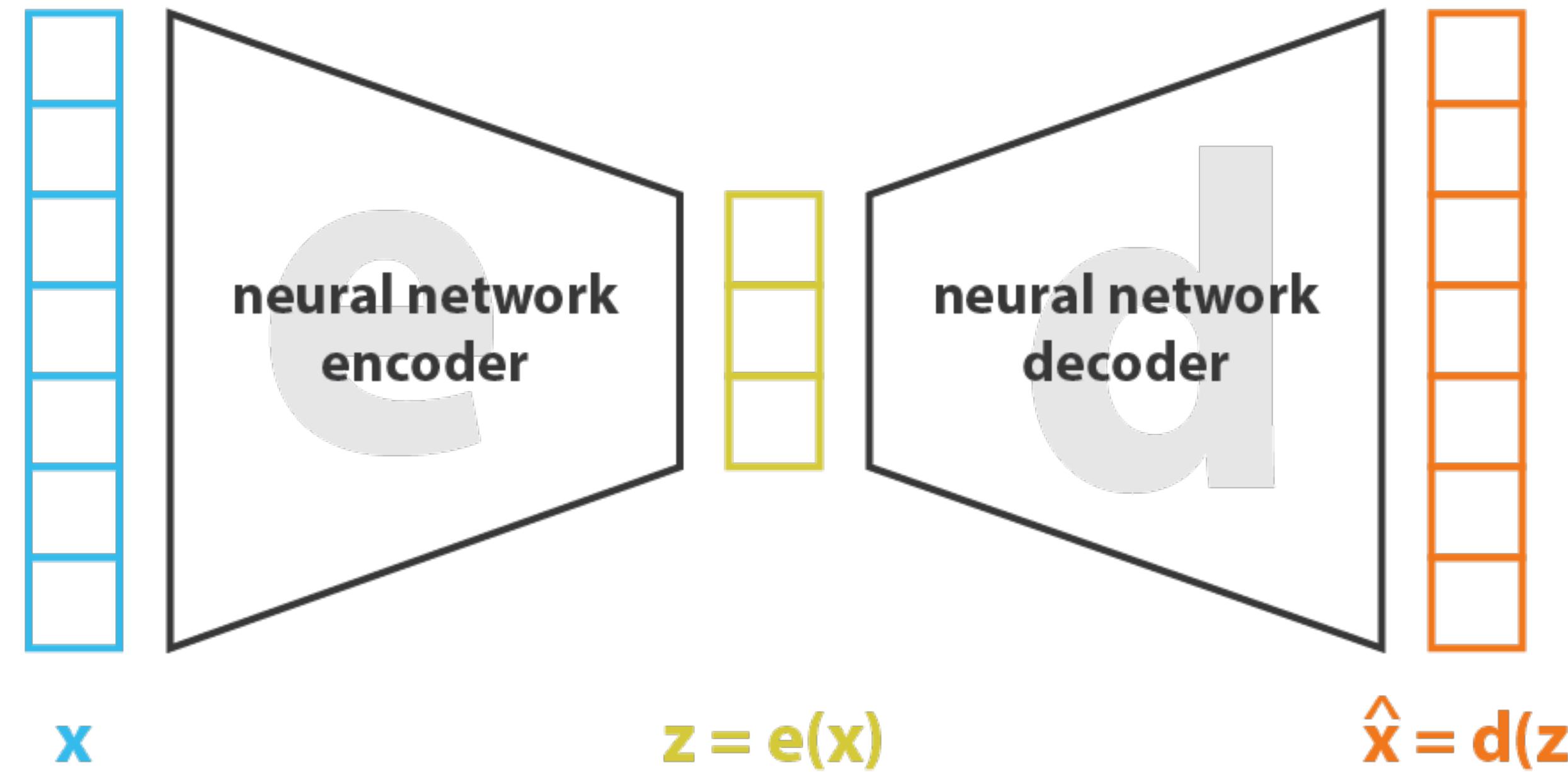


**initial data  
in space  $R^n$**

**encoded data  
in latent space  $R^m$  (with  $m < n$ )**

**encoded-decoded data  
back in the initial space  $R^n$**

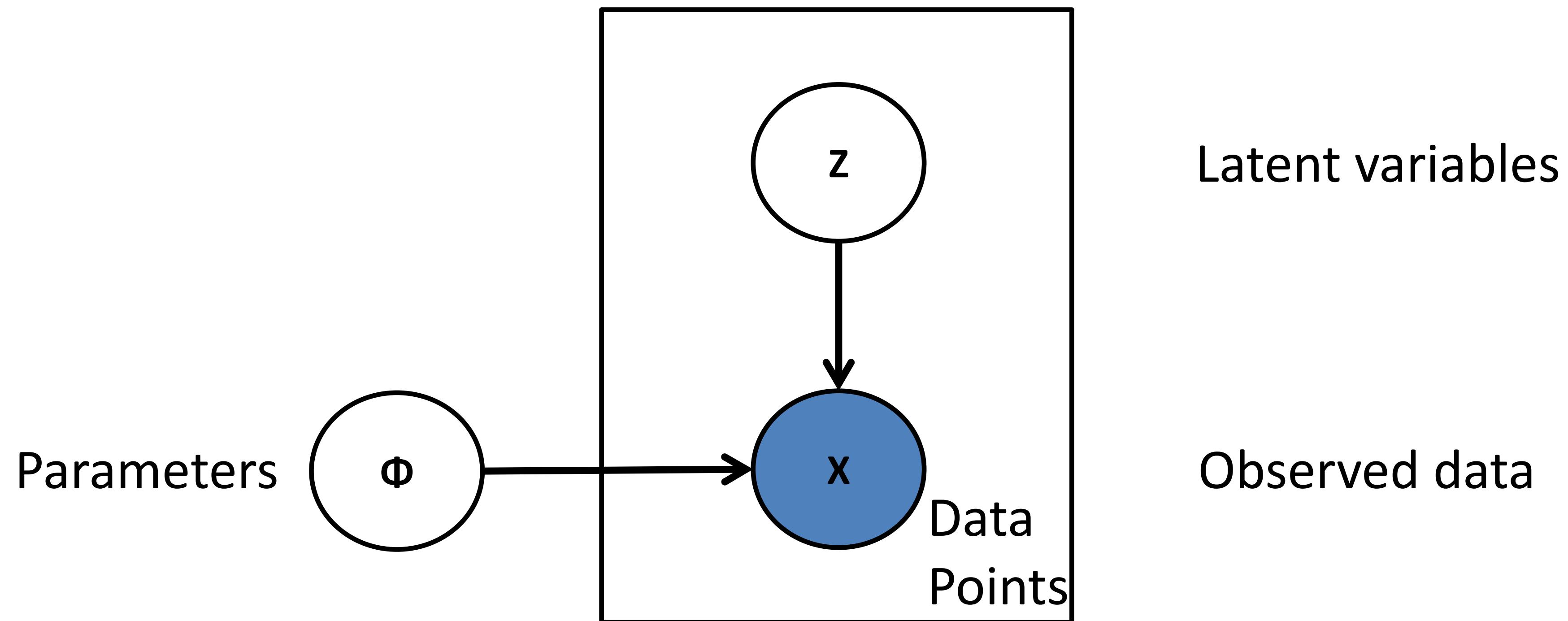
# Autoencoder



---

$$\text{loss} = \| \mathbf{x} - \hat{\mathbf{x}} \|^2 = \| \mathbf{x} - \mathbf{d}(\mathbf{z}) \|^2 = \| \mathbf{x} - \mathbf{d}(\mathbf{e}(\mathbf{x})) \|^2$$

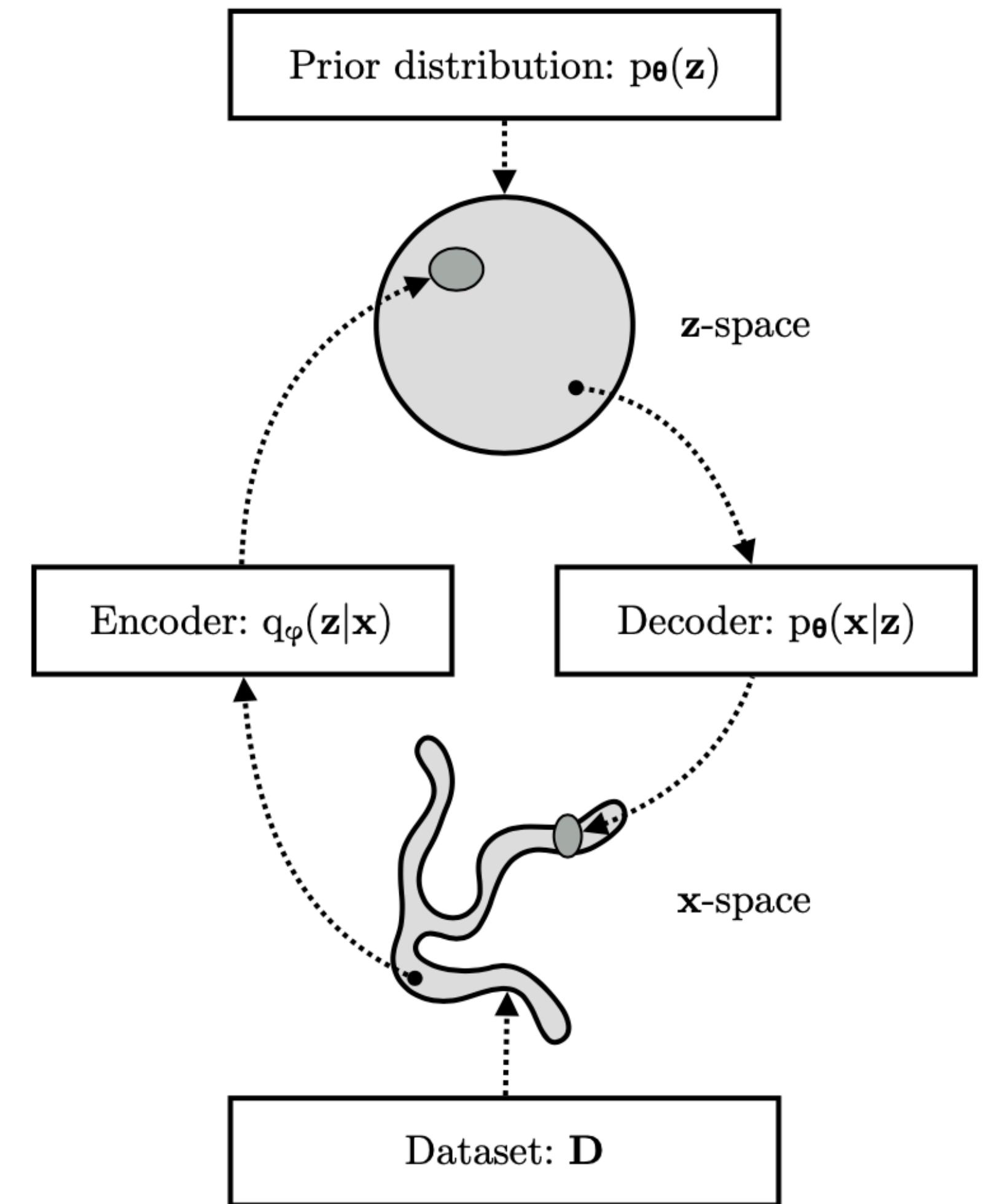
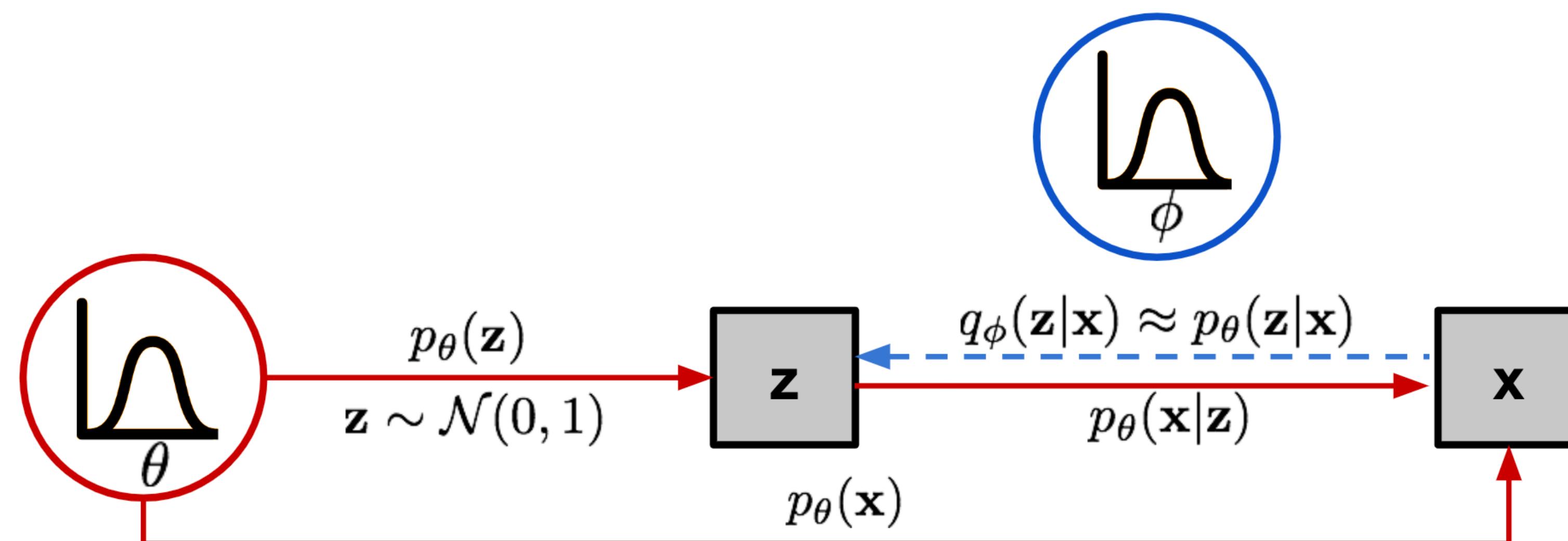
# Latent Variable Models



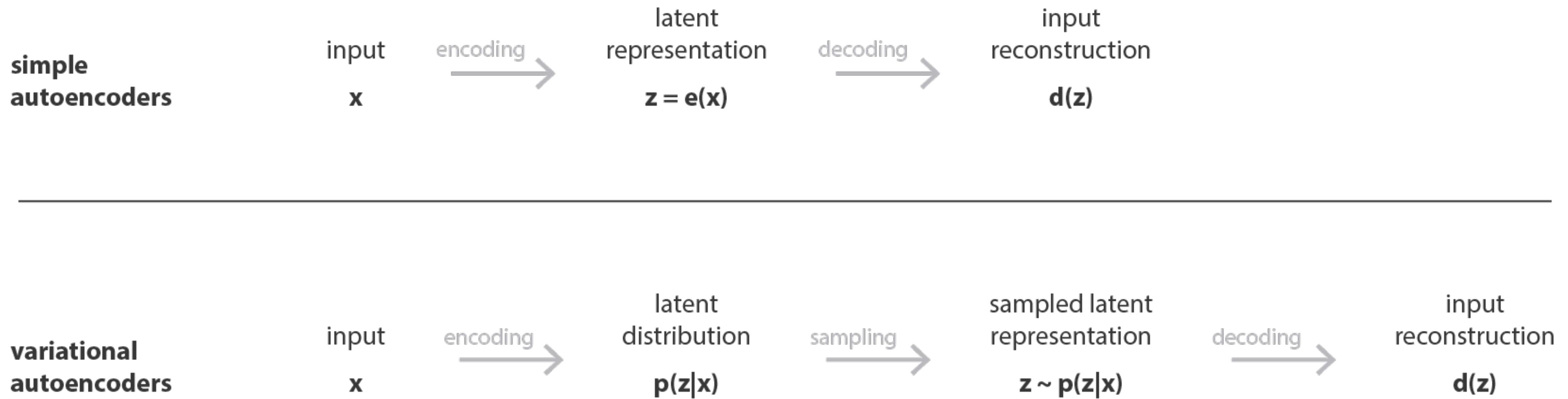
Dimensionality( $X$ )  $>>$  dimensionality( $Z$ )

$Z$  is a **bottleneck**, which finds a **compressed, low-dimensional representation** of  $X$

# Variational Autoencoder (VAE)



# AE vs. VAE

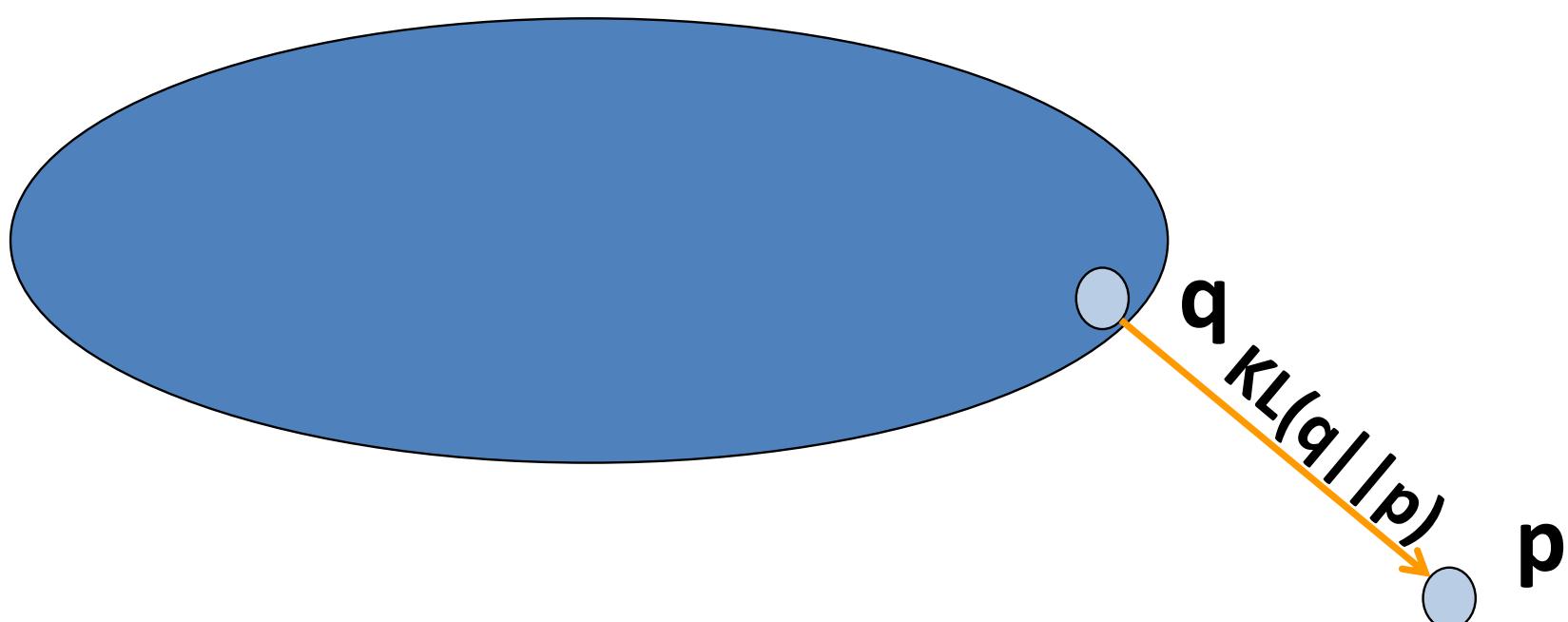
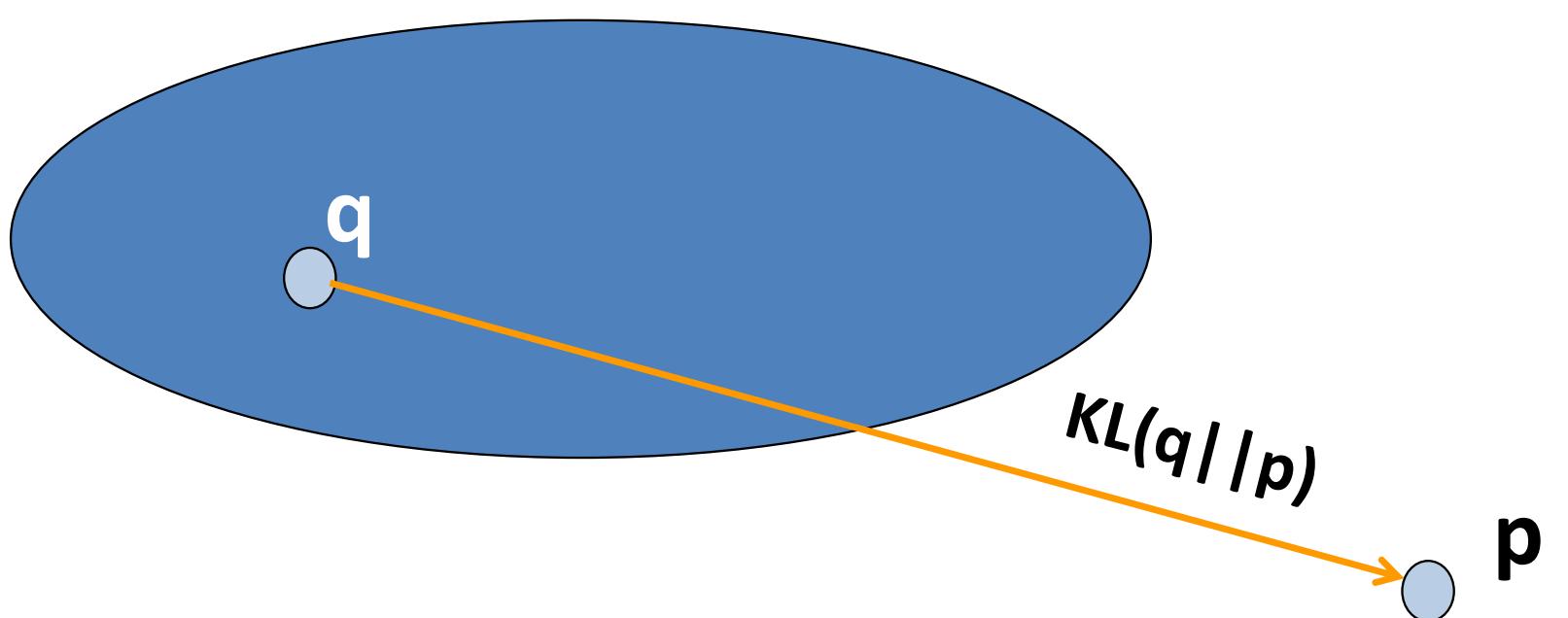


# Approximate Inference

- **Optimization approaches**
  - EM
  - Variational inference
    - **Variational Bayes**
    - Message passing
    - Laplace approximation
- **Simulation approaches (Monte Carlo methods)**
  - MCMC: Gibbs sampling, etc
  - ...

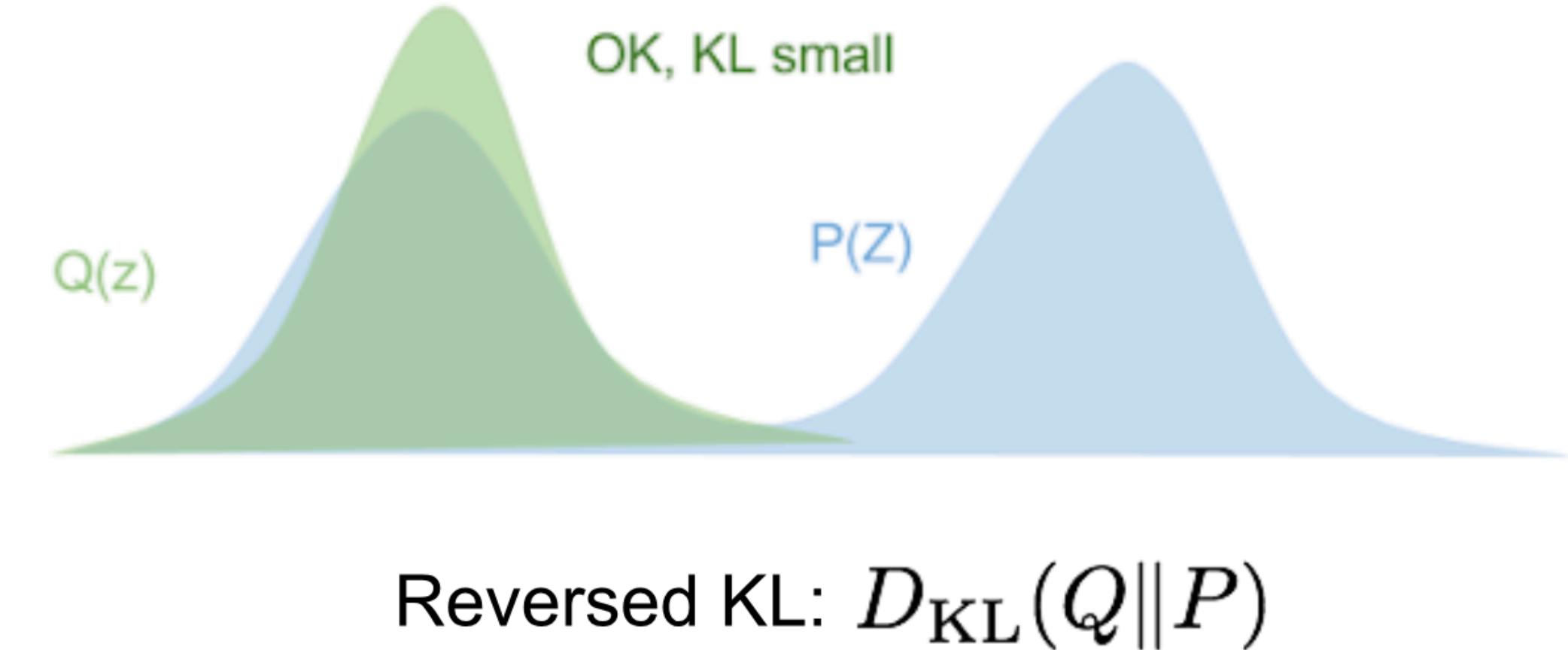
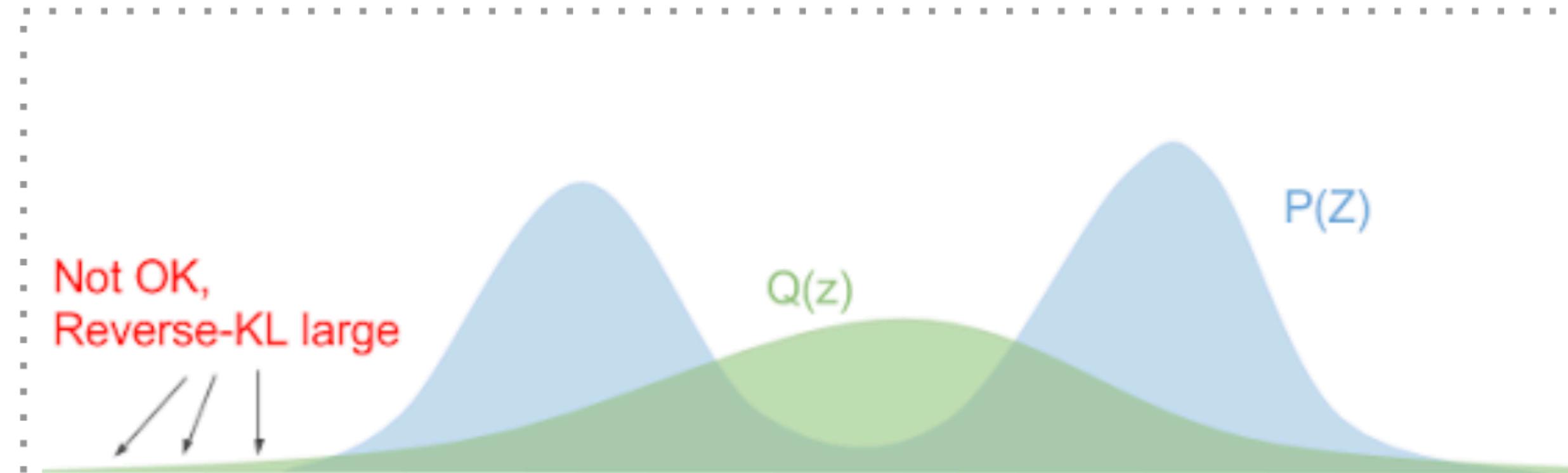
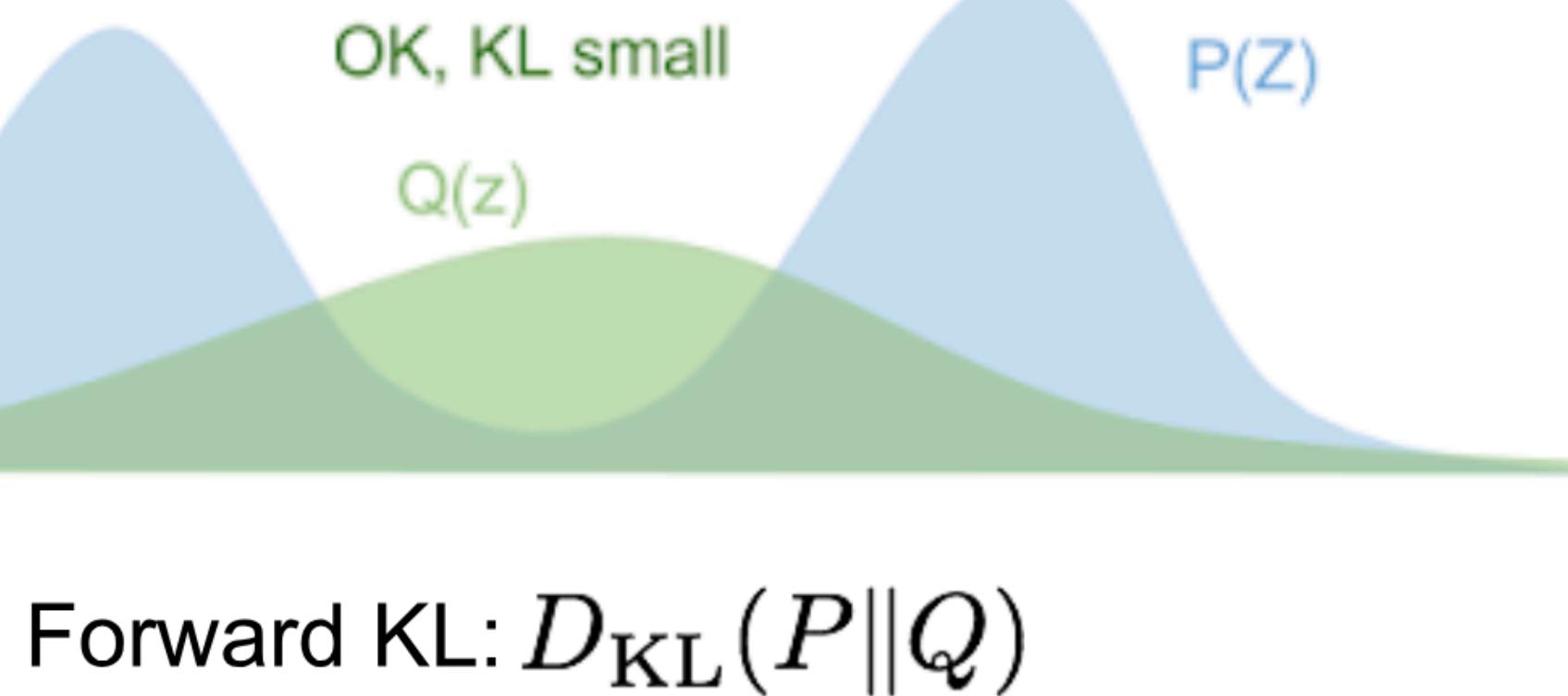
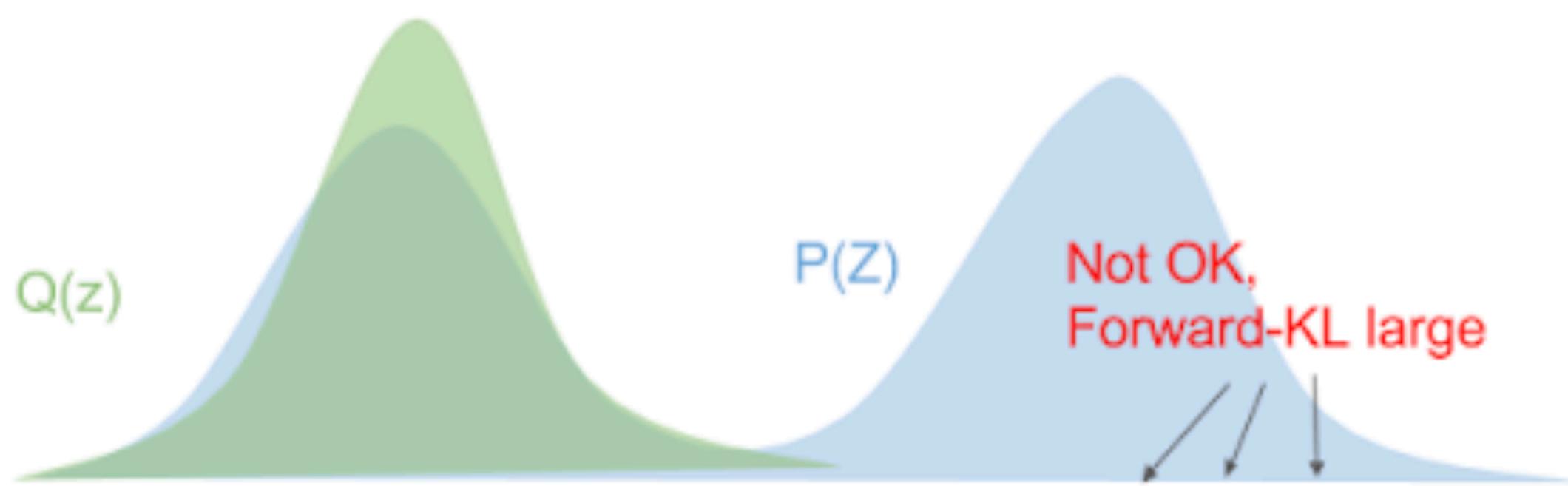
# Variational Inference

- Approximate  $p(z)$  with  $q(z)$
- Make  $q(z)$  tractable to work with
- Solve an optimization problem, e.g.  $\text{KL}(q(z) \parallel p(z))$



# Variational Inference

$$\begin{aligned} KL(P||Q) &= \sum_z p(z) \log \frac{p(z)}{q(z)} \\ &= \mathbb{E}_{p(z)} \left[ \log \frac{p(z)}{q(z)} \right] \end{aligned}$$



# Monte Carlo Methods

$$E_{P(\mathbf{x})}[f(\mathbf{x})] \approx \frac{1}{S} \sum_{i=1}^S f(\mathbf{x}^{(i)}), \mathbf{x}^{(i)} \sim P(\mathbf{x})$$

- This suggests the procedure:
  - Draw S samples from P(x)
  - Compute f(x) for each of the samples
  - Approximate E[f(x)] by the sample average

# Evidence Lower BOund (ELBO)

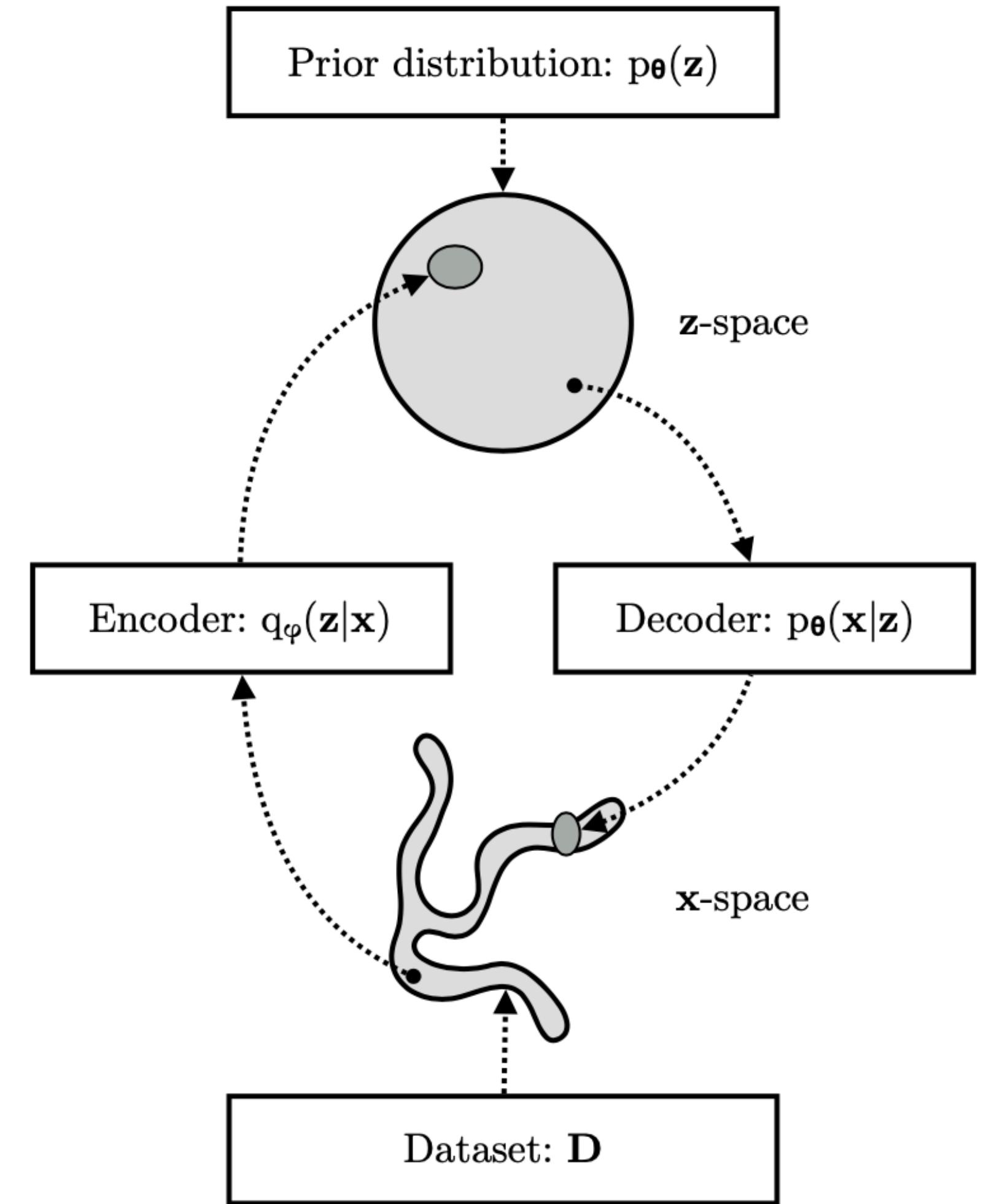
$$\log p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x})] \quad (2.5)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \right] \quad (2.6)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \right] \quad (2.7)$$

$$= \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right]}_{=\mathcal{L}_{\theta, \phi}(\mathbf{x}) \text{ (ELBO)}} + \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \right] \right]}_{=D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))} \quad (2.8)$$

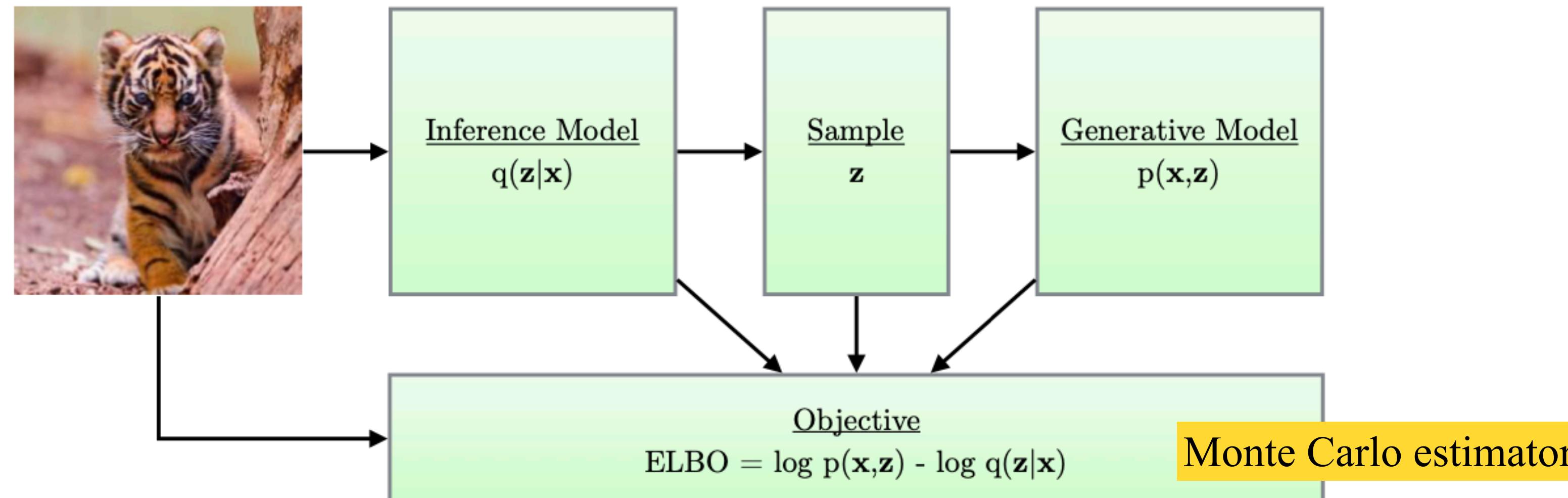
$$\begin{aligned} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \\ &\leq \log p_{\theta}(\mathbf{x}) \end{aligned}$$

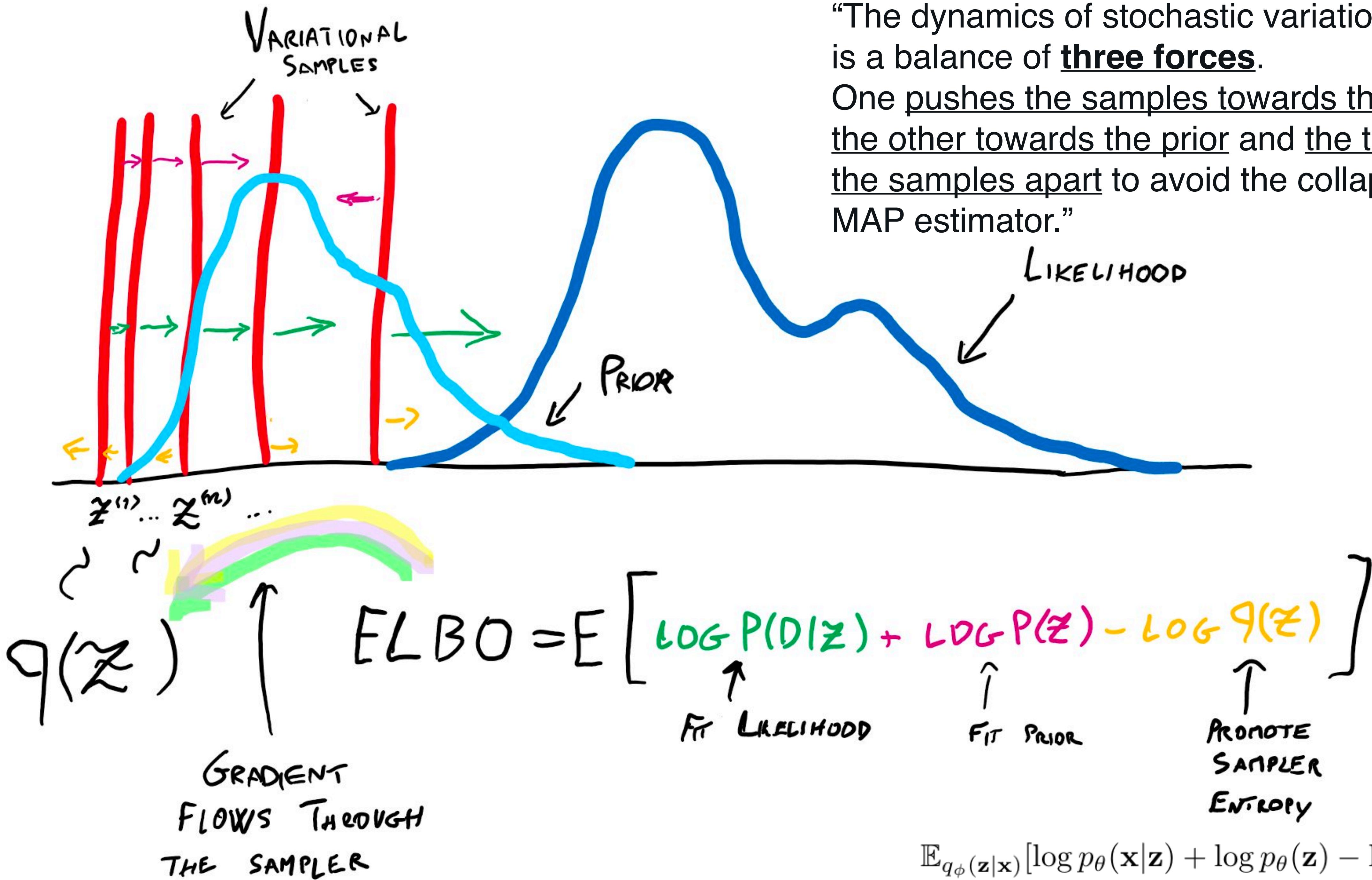


# ELBO

$$\begin{aligned}\mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) \\&= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{z} + \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{z}|\mathbf{x}) d\mathbf{z} - \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]\end{aligned}$$

Datapoint

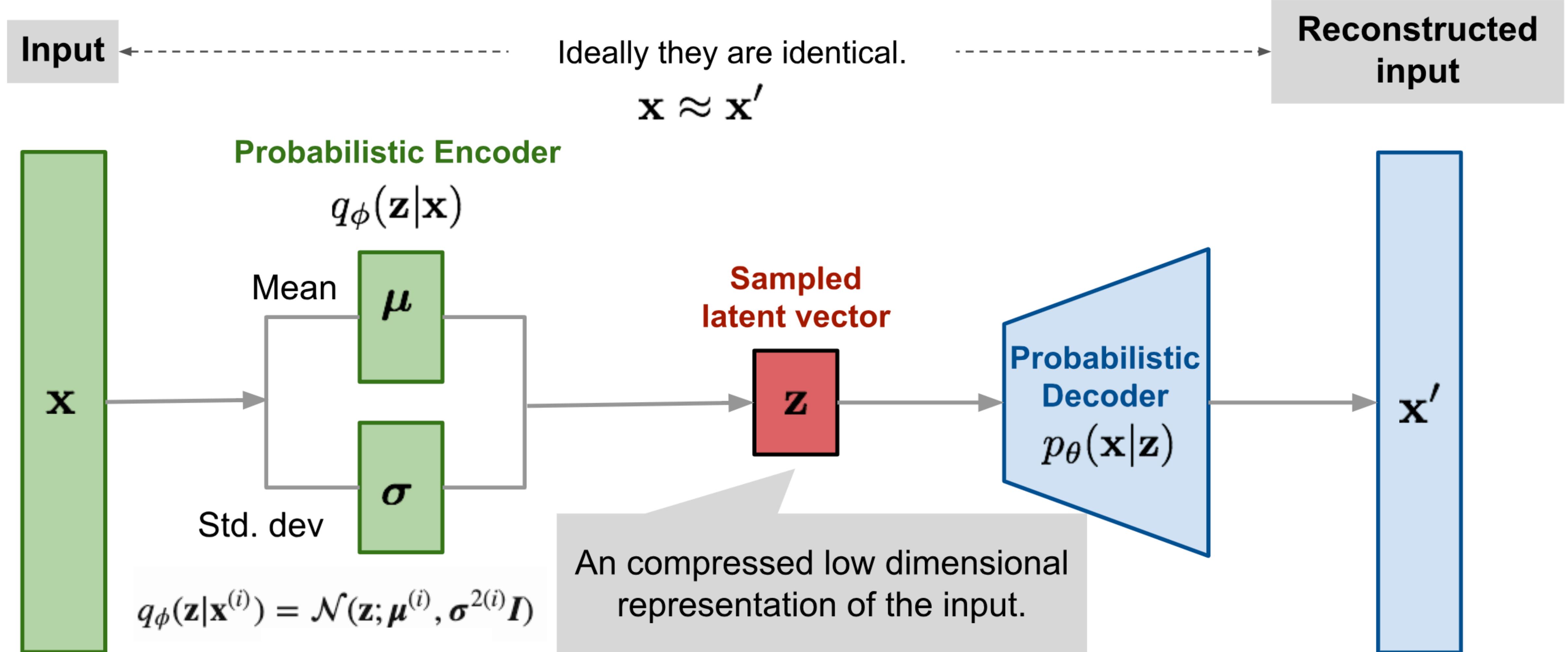




"The dynamics of stochastic variational inference is a balance of three forces. One pushes the samples towards the likelihood, the other towards the prior and the third pushes the samples apart to avoid the collapse into the MAP estimator."

@LucaAmb

# (vanilla) VAE



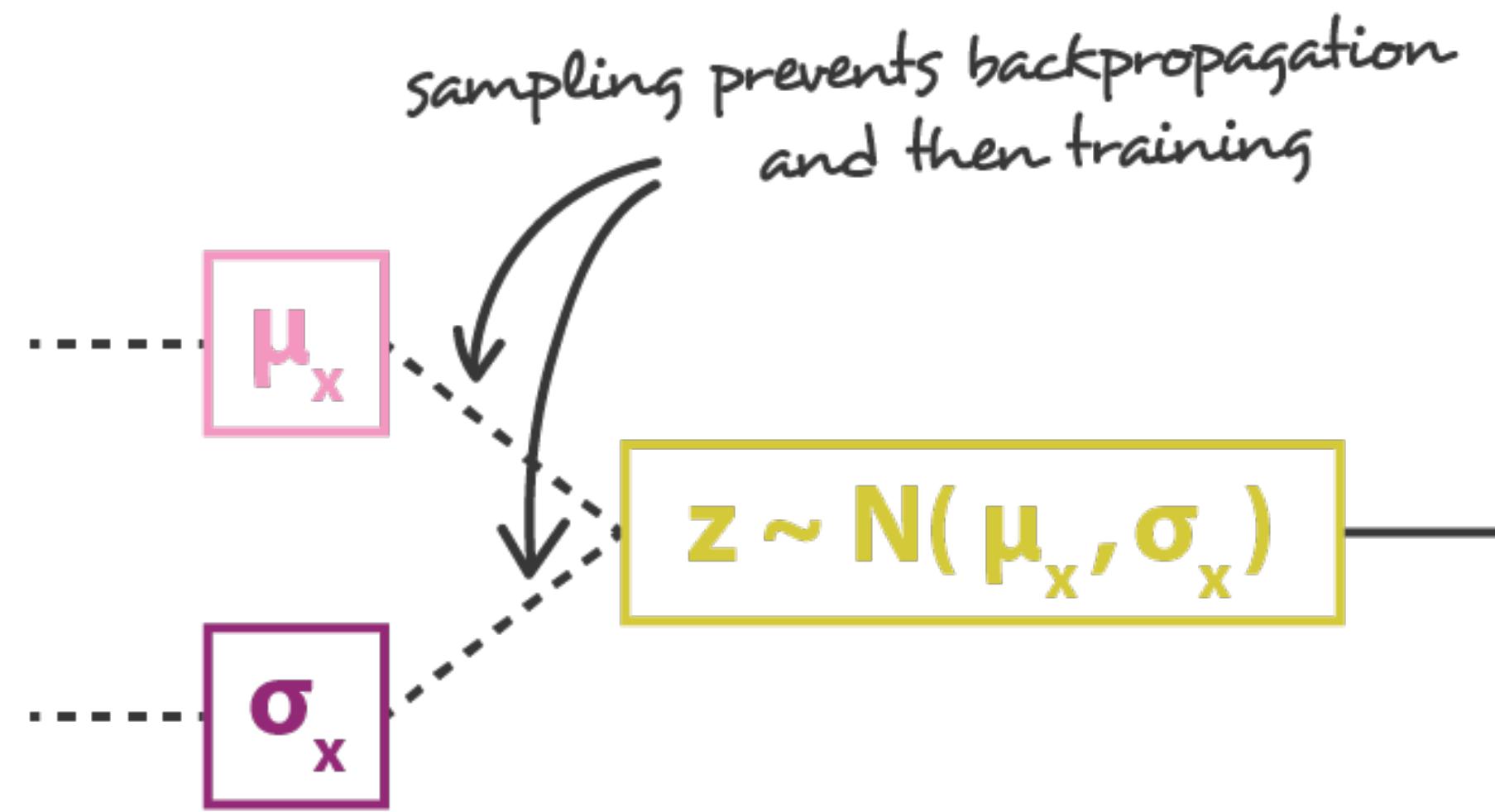
# Reparameterization trick



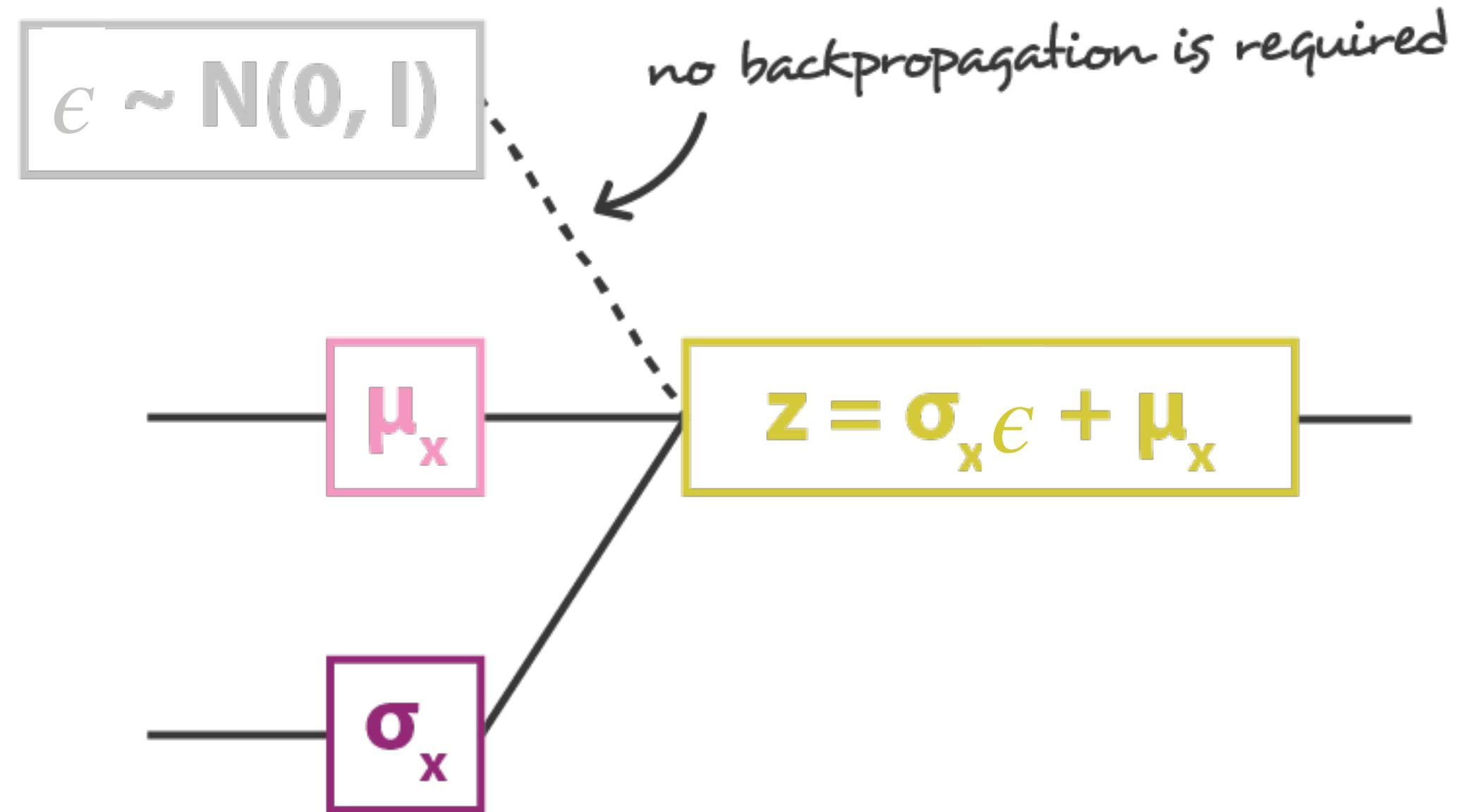
no problem for backpropagation



backpropagation is not possible due to sampling

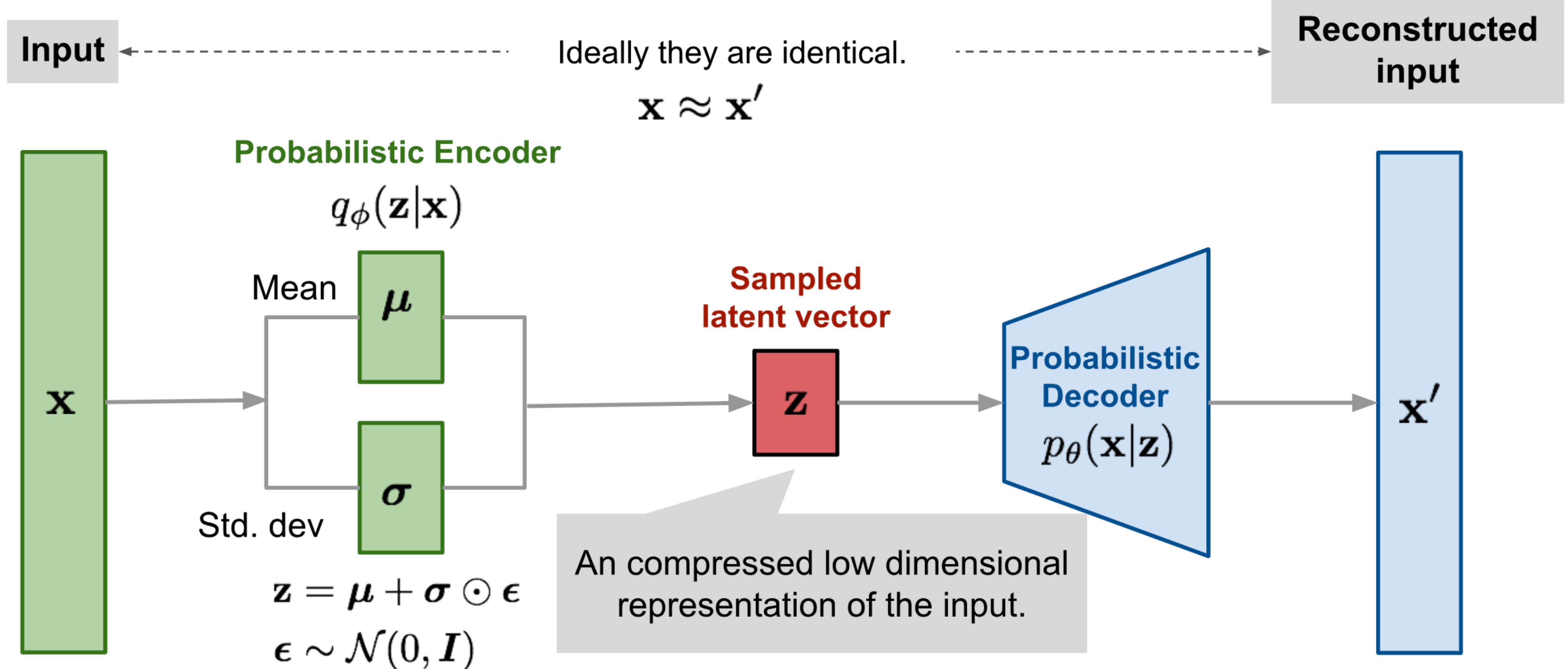


sampling without reparametrisation trick



sampling with reparametrisation trick

# (vanilla) VAE



# ELBO

$$\begin{aligned}\mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) \\&= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{z} + \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{z}|\mathbf{x}) d\mathbf{z} - \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\&= \mathbb{E}_{p(\epsilon)} [\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})], \text{ where } \mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}\end{aligned}$$

# ELBO

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})], \text{ where } \mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$$

Assuming a Gaussian likelihood for the decoder, i.e.  $\mathbf{x}|\mathbf{z}=z \sim \mathcal{N}(\boldsymbol{\mu}_{\theta,z}, \boldsymbol{\Sigma}_{\theta,z})$ , we have

$$\begin{aligned}\log p_{\theta}(\mathbf{x}|\mathbf{z}) &= \log \left[ (2\pi)^{-\frac{J}{2}} |\boldsymbol{\Sigma}_{\theta,z}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\theta,z})^T \boldsymbol{\Sigma}_{\theta,z}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{\theta,z})\right) \right] \\ &= -\frac{1}{2} \left[ (\mathbf{x} - \boldsymbol{\mu}_{\theta,z})^T \boldsymbol{\Sigma}_{\theta,z}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{\theta,z}) + J \log(2\pi) + \log |\boldsymbol{\Sigma}_{\theta,z}| \right]\end{aligned}$$

In the isotropic Gaussian case, this reduces to

$$\begin{aligned}\log p_{\theta}(\mathbf{x}|\mathbf{z}) &= \sum_{i=1}^J \log \left[ (2\pi)^{-\frac{1}{2}} \sigma_i^{-1} e^{-\frac{1}{2\sigma_i^2} (x_i - \mu_{\theta,z,i})^2} \right] \\ &= -\frac{1}{2} \sum_{i=1}^J \left[ \frac{(x_i - \mu_{\theta,z,i})^2}{\sigma_i^2} + \log(2\pi) + \log \sigma_i^2 \right],\end{aligned}$$

where  $\mu_{\theta,z,i}$  is the  $i^{\text{th}}$  element of  $\boldsymbol{\mu}_{\theta,z}$ , and  $\sigma_i^2$  is the  $i^{\text{th}}$  diagonal entry of  $\boldsymbol{\Sigma}_{\theta,z}$ .

# ELBO

$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$ , where  $\mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$

$$\begin{aligned}\log p_{\theta}(\mathbf{z}) &= \sum_{i=1}^J \log \mathcal{N}(z_i; 0, 1) \\ &= \sum_{i=1}^J \log \left[ (2\pi)^{-\frac{1}{2}} e^{-\frac{1}{2} z_i^2} \right] \\ &= -\sum_{i=1}^J \frac{1}{2} (z_i^2 + \log(2\pi))\end{aligned}$$

# ELBO

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})], \text{ where } \mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$$

Reparameterization

$\boldsymbol{\epsilon}$ ->  $\mathbf{z}$ :

Change of variables formula (vector case):

$$p(\boldsymbol{\epsilon}) = q_{\phi}(\mathbf{g}(\boldsymbol{\epsilon})) \left| \det \left[ \frac{\partial \mathbf{g}(\boldsymbol{\epsilon})}{\partial \boldsymbol{\epsilon}} \right] \right|$$

$$\log p(\boldsymbol{\epsilon}) = \log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log \left| \det \left[ \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right] \right|$$

$$\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \log p(\boldsymbol{\epsilon}) - \log \left| \det \left[ \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right] \right|$$

$$\log \left| \det \left[ \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right] \right| = \sum_{i=1}^J \log \sigma_i = \frac{1}{2} \sum_{i=1}^J \log \sigma_i^2$$

# ELBO

$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$ , where  $\mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$

$$\begin{aligned}\mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{z})] &= \mathbb{E}_{p(\epsilon)}\left[-\frac{1}{2}\sum_{i=1}^J(z_i^2 + \log(2\pi))\right] \\ &= -\frac{1}{2}\sum_{i=1}^J\mathbb{E}_{p(\epsilon)}[(\mu_i + \sigma_i \cdot \epsilon_i)^2] - \frac{1}{2}\sum_{i=1}^J\log(2\pi) \\ &= -\frac{1}{2}\sum_{i=1}^J(\mathbb{E}_{p(\epsilon)}[\mu_i^2] + \mathbb{E}_{p(\epsilon)}[\sigma_i^2\epsilon_i^2] + \mathbb{E}_{p(\epsilon)}[2\mu_i\sigma_i\epsilon_i]) - \frac{J}{2}\log(2\pi) \\ &= -\frac{1}{2}\sum_{i=1}^J(\mu_i^2 + \sigma_i^2\mathbb{E}_{p(\epsilon)}[\epsilon_i^2] + 2\mu_i\sigma_i\mathbb{E}_{p(\epsilon)}[\epsilon_i]) - \frac{J}{2}\log(2\pi) \\ &= -\frac{1}{2}\sum_{i=1}^J(\mu_i^2 + \sigma_i^2 + 2\mu_i\sigma_i \cdot 0) - \frac{J}{2}\log(2\pi) \\ &= -\frac{J}{2}\log(2\pi) - \frac{1}{2}\sum_{i=1}^J(\mu_i^2 + \sigma_i^2)\end{aligned}$$

# ELBO

$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$ , where  $\mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$

$$\begin{aligned}\mathbb{E}_{p(\epsilon)}[\log q_{\phi}(\mathbf{z}|\mathbf{x})] &= \mathbb{E}_{p(\epsilon)}\left[-\frac{1}{2}\sum_{i=1}^J(\epsilon_i^2 + \log(2\pi) + \log\sigma_i^2)\right] \\ &= -\frac{1}{2}\sum_{i=1}^J\mathbb{E}_{p(\epsilon)}[\epsilon_i^2] - \frac{1}{2}\sum_{i=1}^J\log(2\pi) - \frac{1}{2}\sum_{i=1}^J\log\sigma_i^2 \\ &= -\frac{J}{2}\log(2\pi) - \frac{1}{2}\sum_{i=1}^J(1 + \log\sigma_i^2).\end{aligned}$$

# ELBO

$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$ , where  $\mathbf{z} = \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} + \boldsymbol{\mu}$

$$\mathcal{L}_{\theta,\phi}(x_j) \simeq \frac{1}{2} \sum_{i=1}^J (1 + \log((\sigma_i^{(j)})^2) - (\mu_i^{(j)})^2 - (\sigma_i^{(j)})^2) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(x_j|z_j)]$$

# ELBO: full covariance

- Reparameterization

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

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}$$

$$\begin{aligned}\boldsymbol{\Sigma} &= \mathbb{E} [(\mathbf{z} - \mathbb{E} [\mathbf{z}])(\mathbf{z} - \mathbb{E} [\mathbf{z}])^T] \\ &= \mathbb{E} [\mathbf{L}\boldsymbol{\epsilon}(\mathbf{L}\boldsymbol{\epsilon})^T] = \mathbf{L}\mathbb{E} [\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T]\mathbf{L}^T \\ &= \mathbf{L}\mathbf{L}^T\end{aligned}$$

$$\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \mathbf{L}$$

$$\log |\det\left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}}\right)| = \sum_i \log |L_{ii}|$$

$$\begin{aligned}(\boldsymbol{\mu}, \log \boldsymbol{\sigma}, \mathbf{L}') &\leftarrow \text{EncoderNeuralNet}_{\phi}(\mathbf{x}) \\ \mathbf{L} &\leftarrow \mathbf{L}_{mask} \odot \mathbf{L}' + \text{diag}(\boldsymbol{\sigma})\end{aligned}$$

$$\log \left| \det \left( \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| = \sum_i \log \sigma_i$$

Note that  $\sigma_i$  here is not standard deviation