



*deepseek*  
*API key*

*claude code*

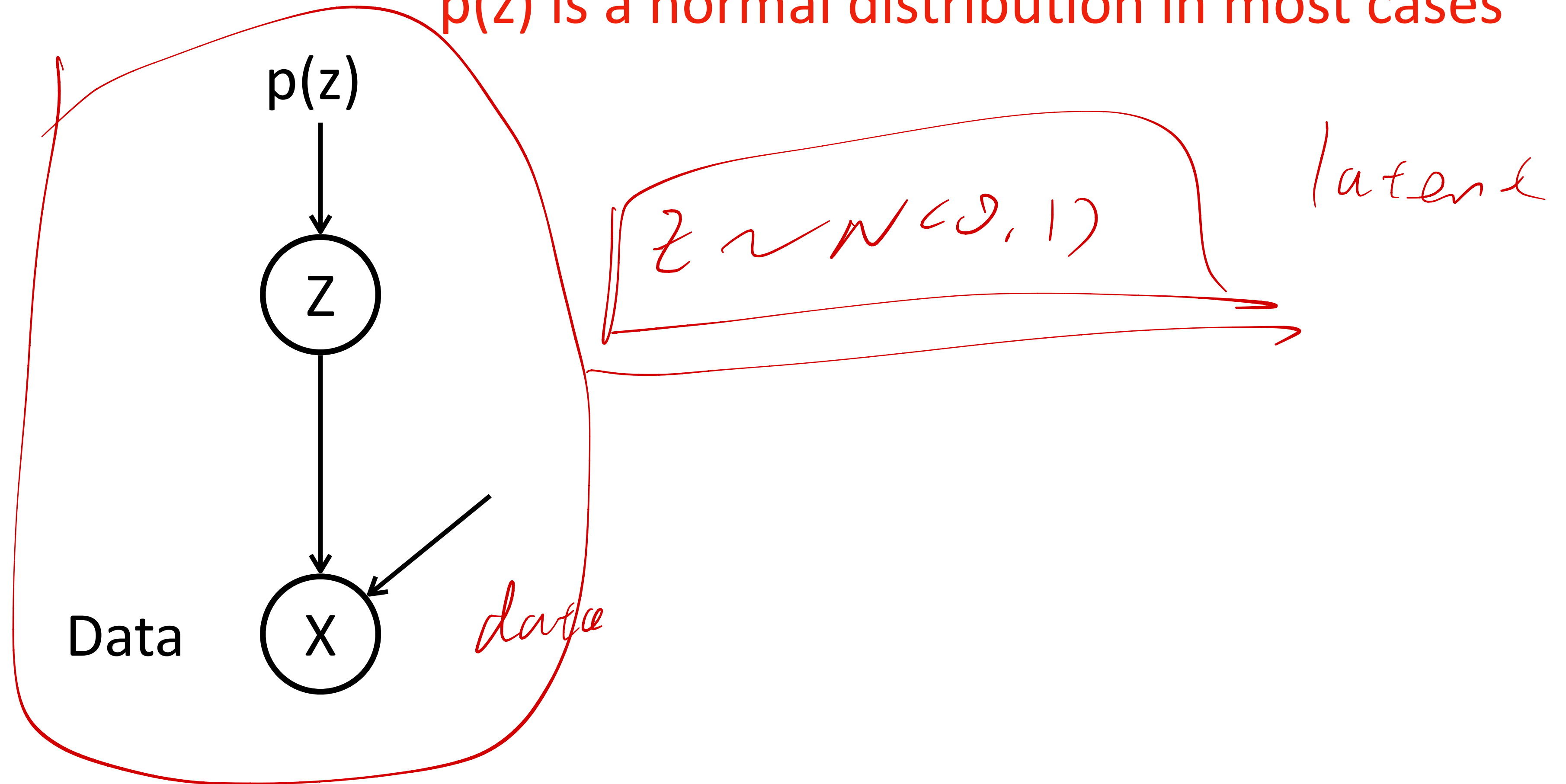
# Variational Autoencoders, Hidden Markov Models

*pam*

Junxian He  
Apr 2, 2026

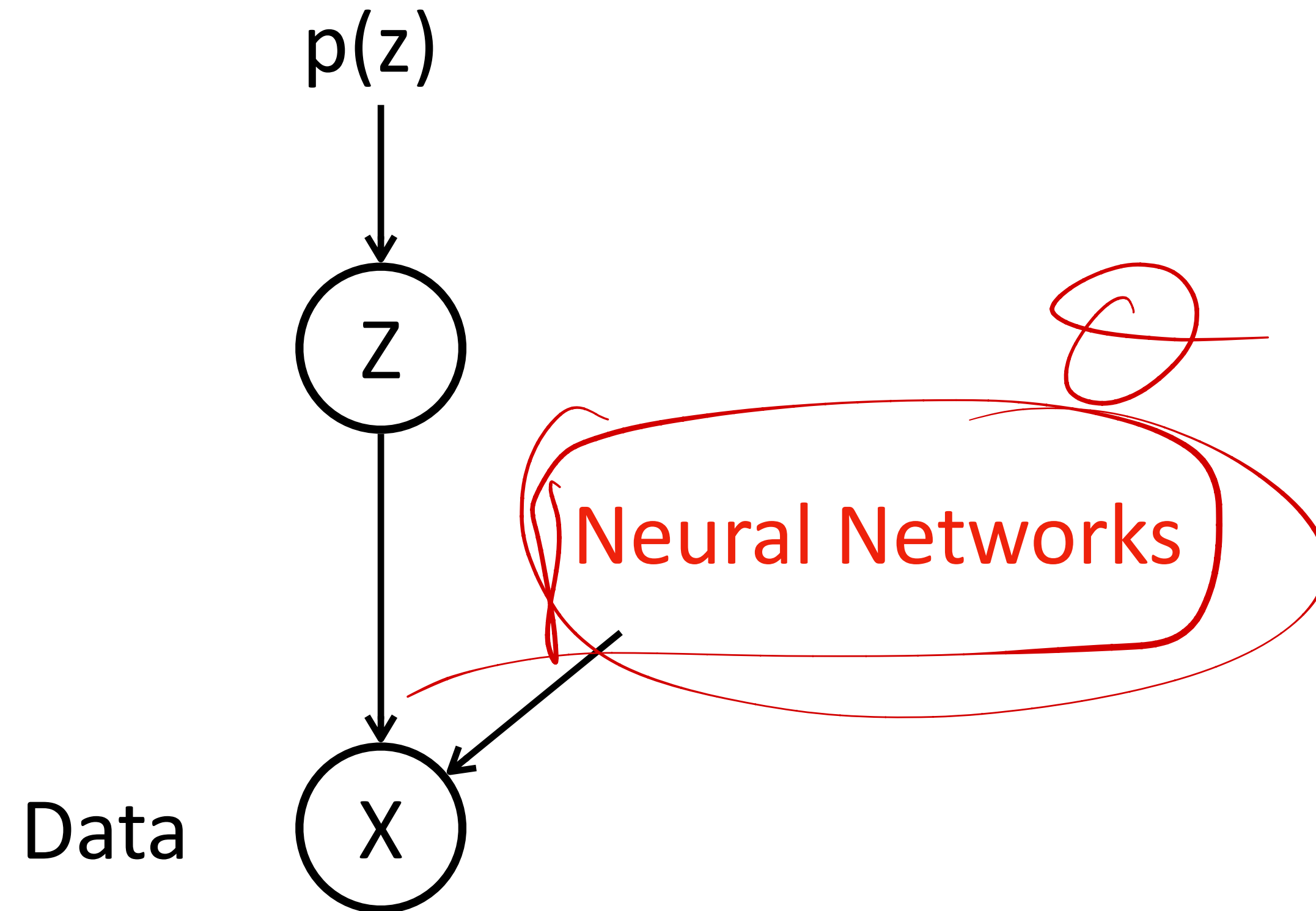
# Recap: The VAE Model

$p(z)$  is a normal distribution in most cases



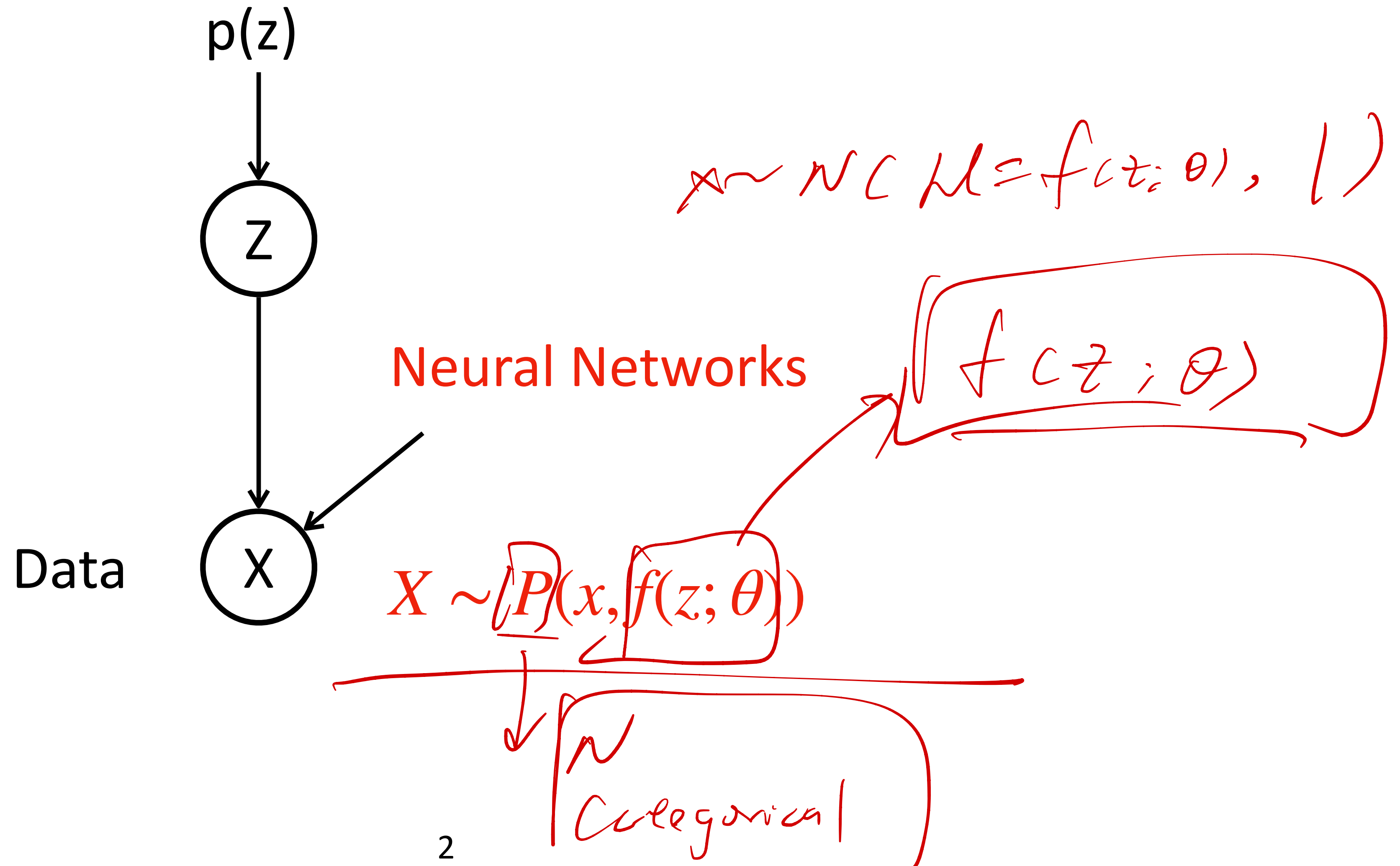
# Recap: The VAE Model

$p(z)$  is a normal distribution in most cases



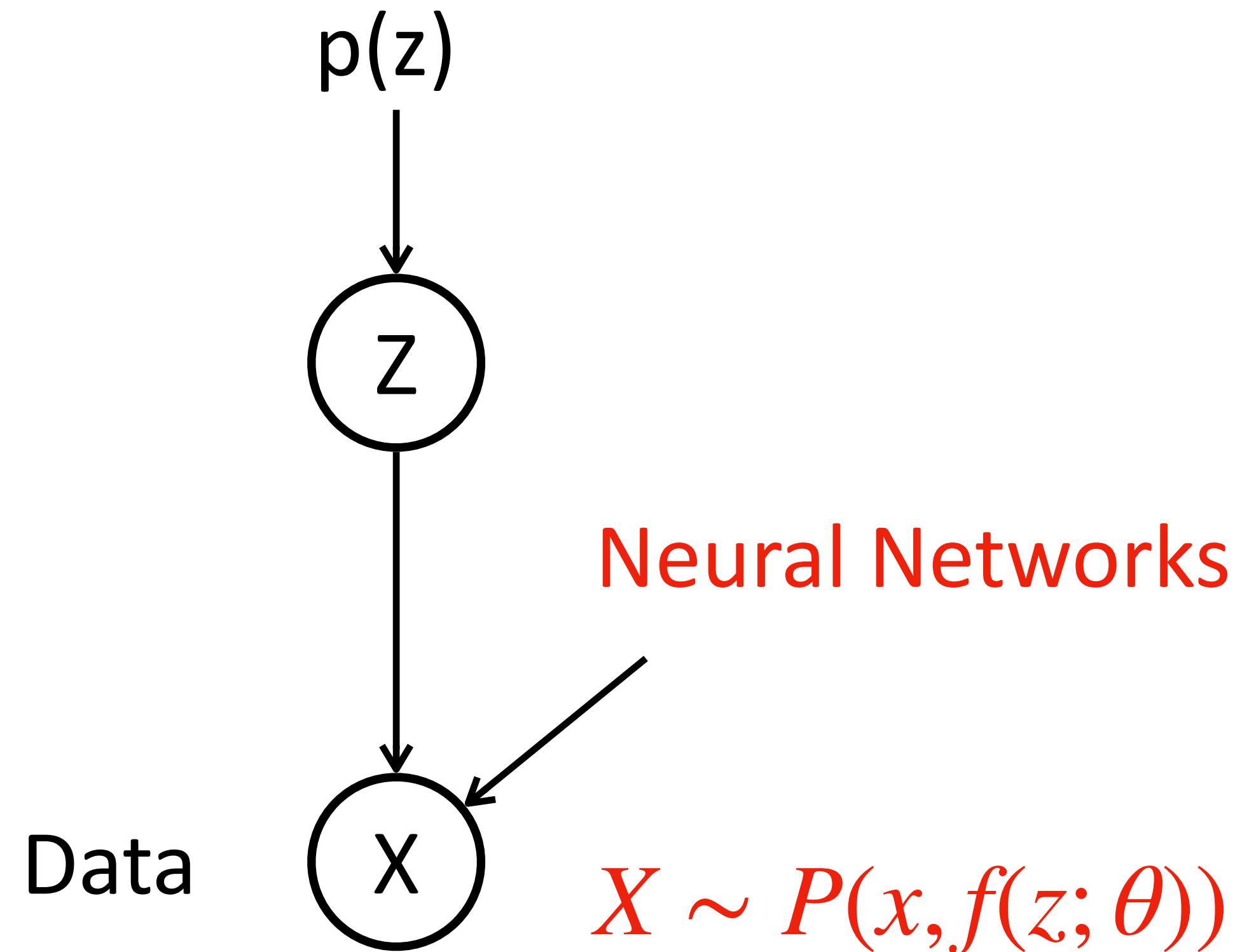
# Recap: The VAE Model

$p(z)$  is a normal distribution in most cases



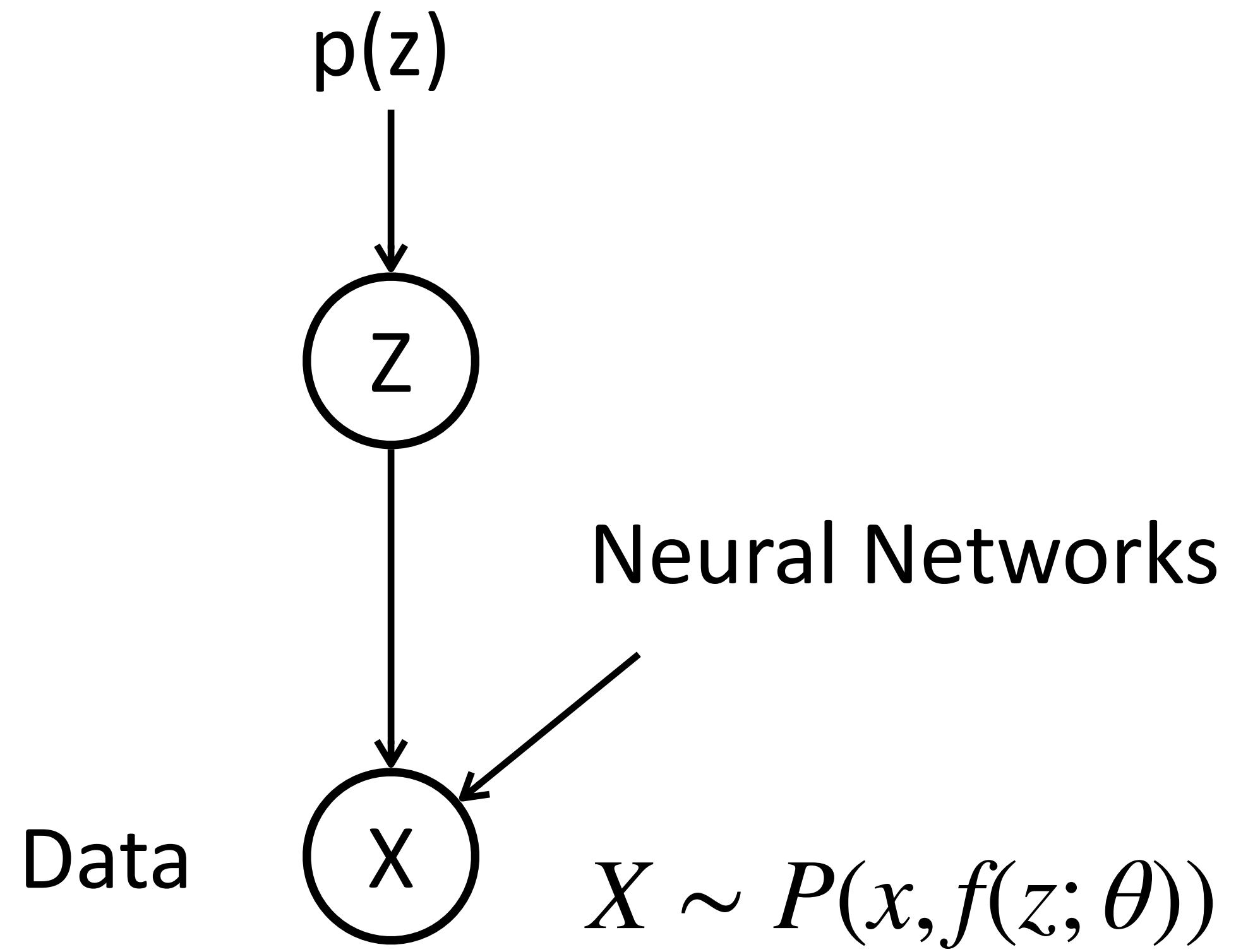
# Recap: The VAE Model

$p(z)$  is a normal distribution in most cases

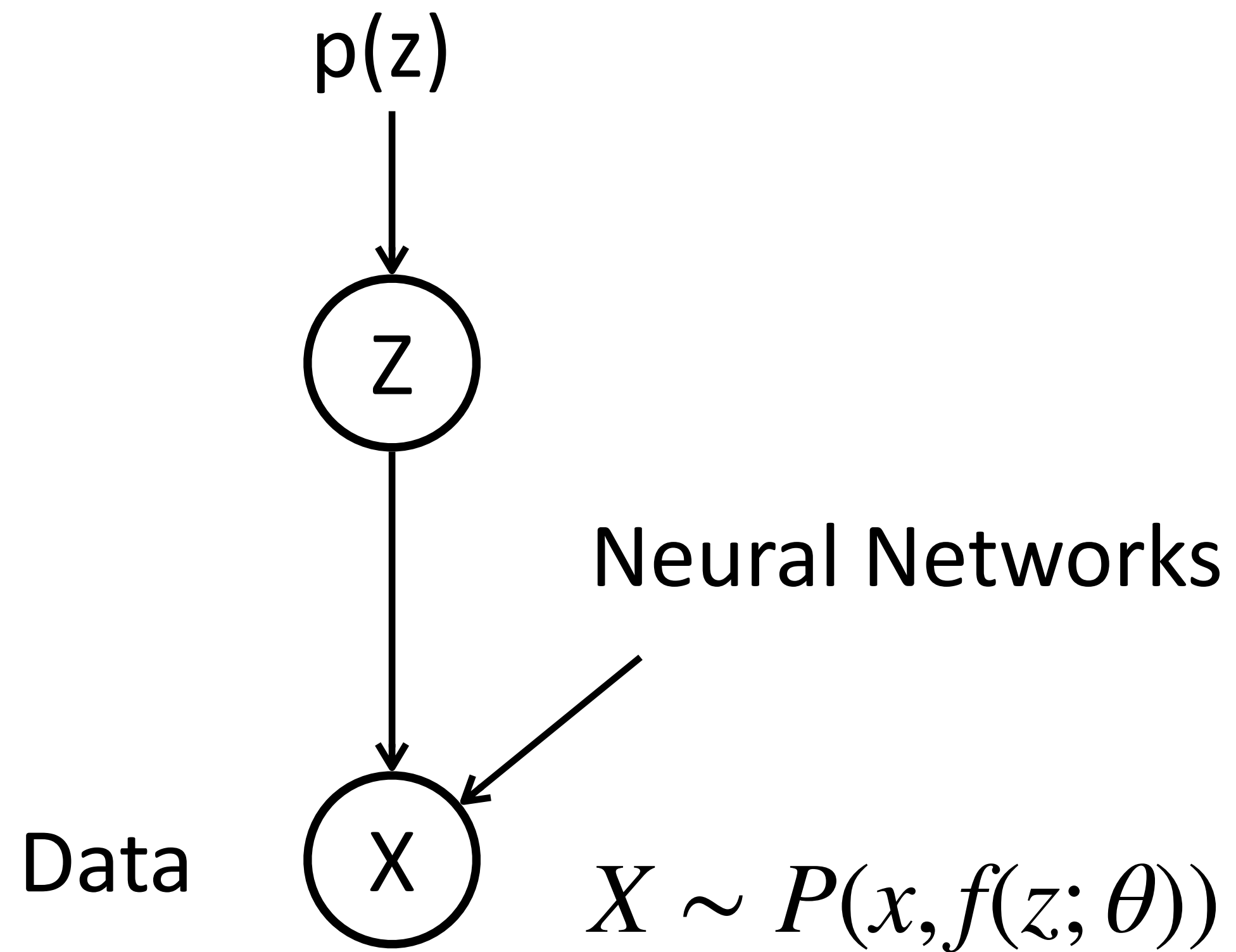


$f$  is a neural network taking  $Z$  as input

# Training

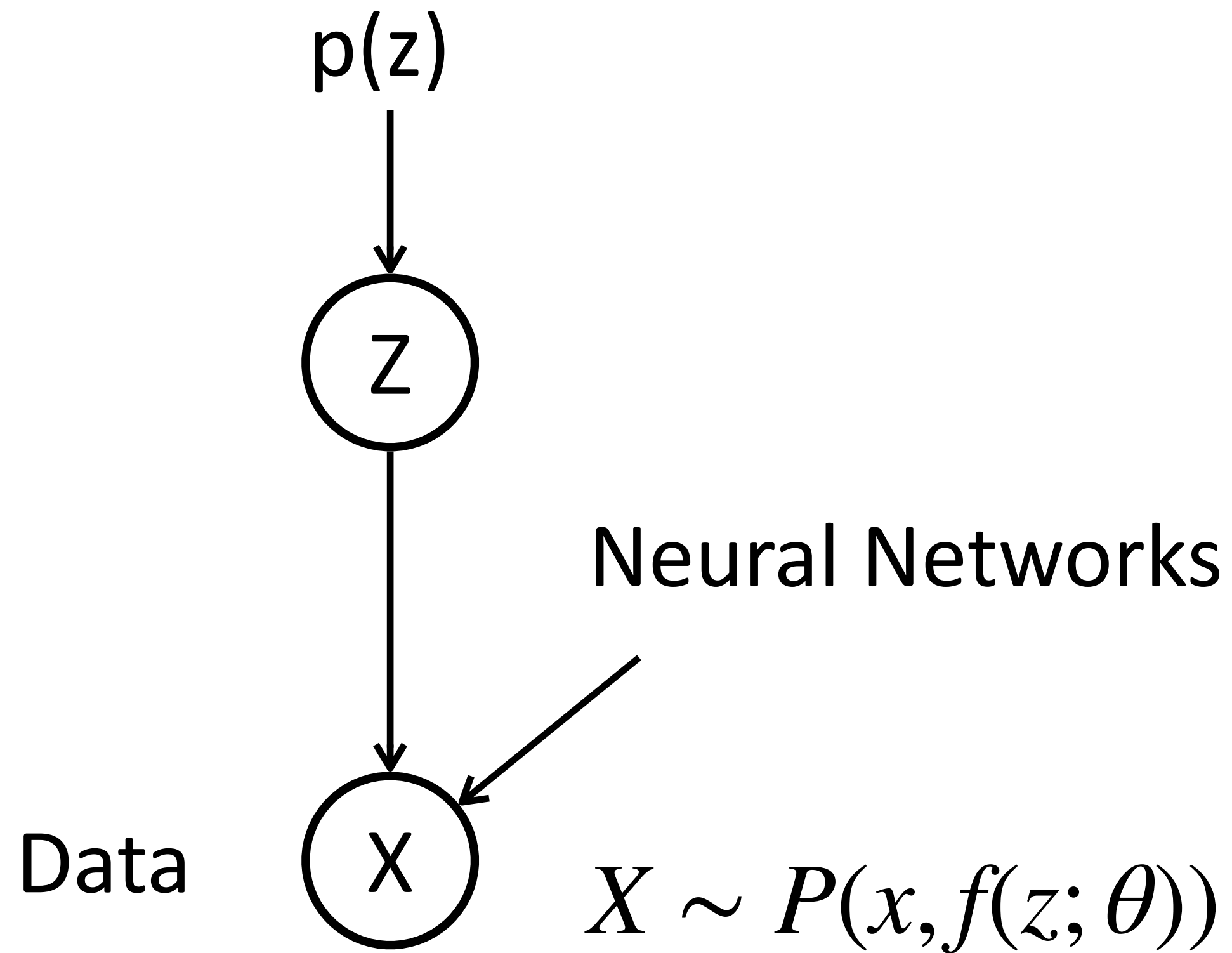


# Training



How to train the model? Can we do MLE?

# Training

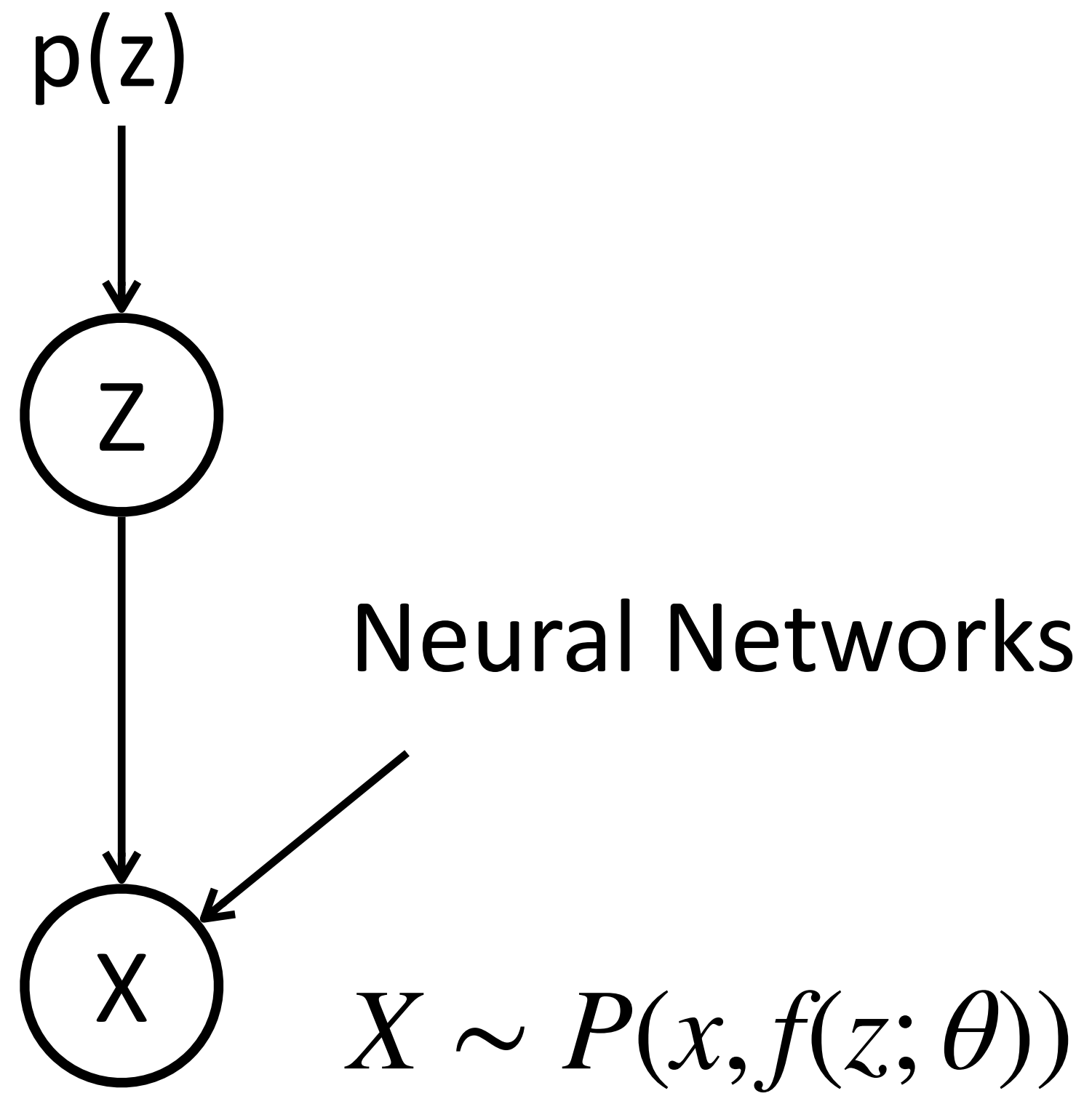


How to train the model? Can we do MLE?

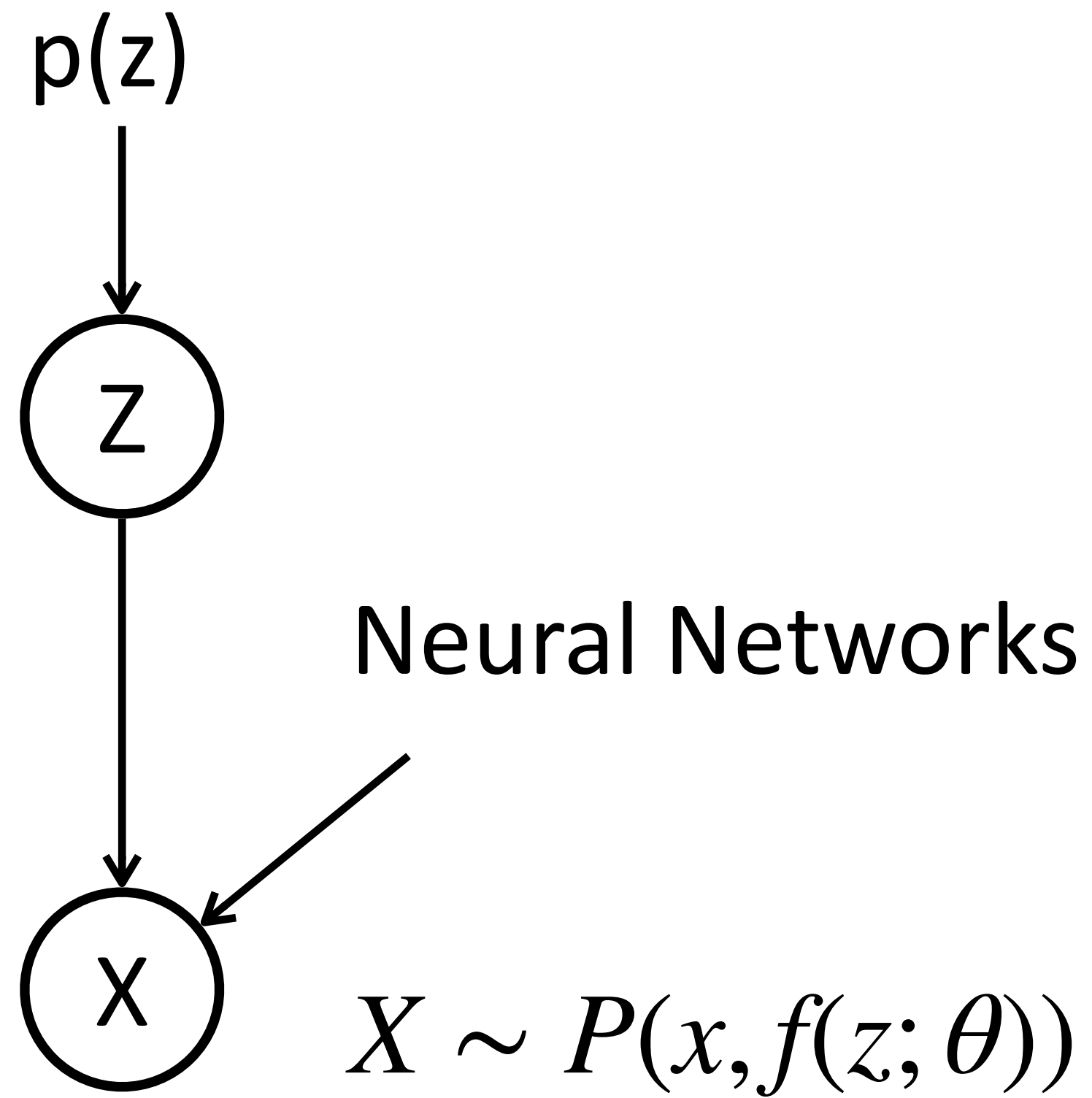
Intractable  $P(X)$ , EM algorithm?

$$\int z P(x, z) dz$$

# Let's try EM



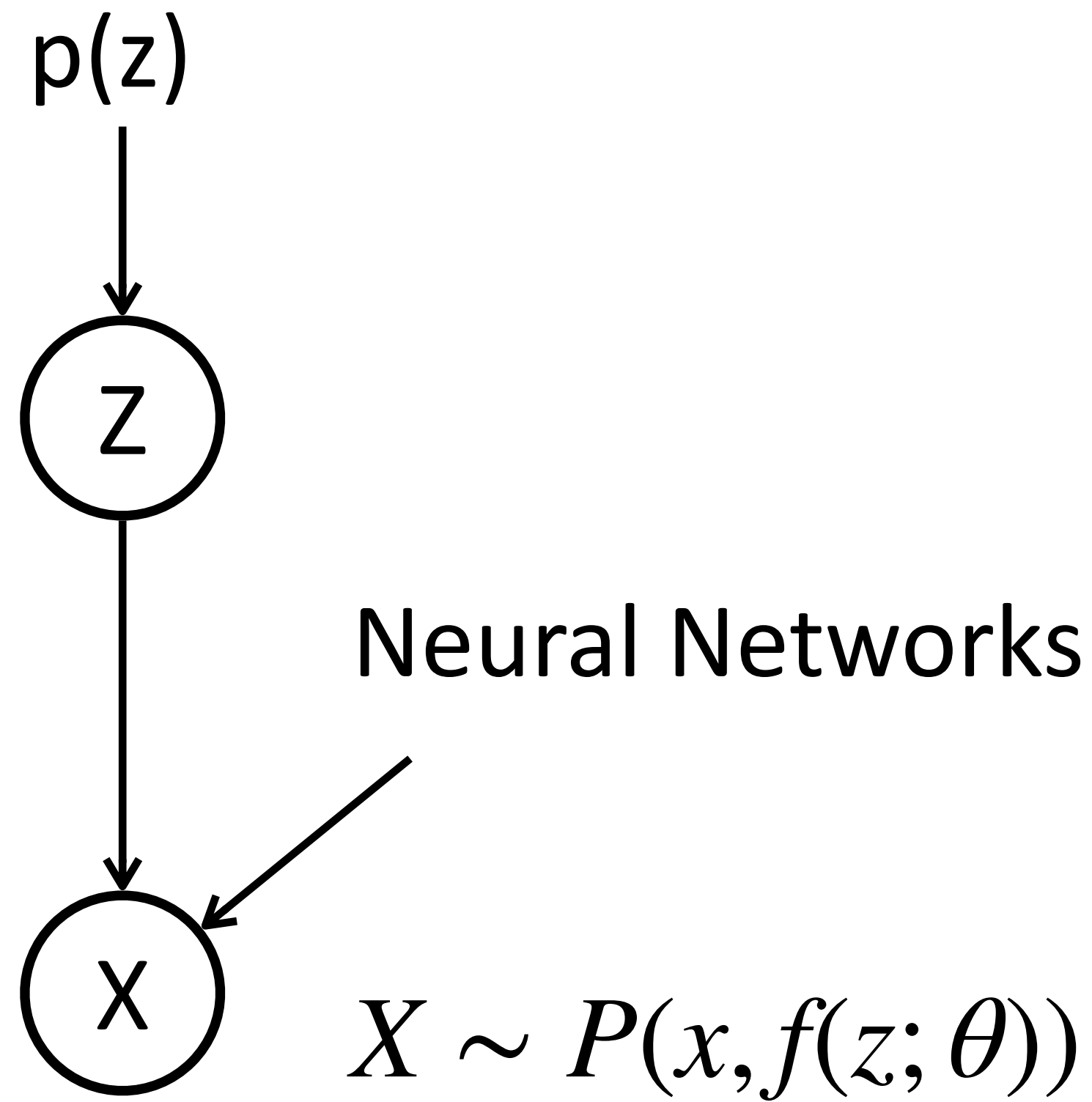
# Let's try EM



E-Step: compute  $P(z|x)$

$$Q(z) = P(z|x) \propto P(z)P(x|z)$$

# Let's try EM



E-Step: compute  $P(z|x)$

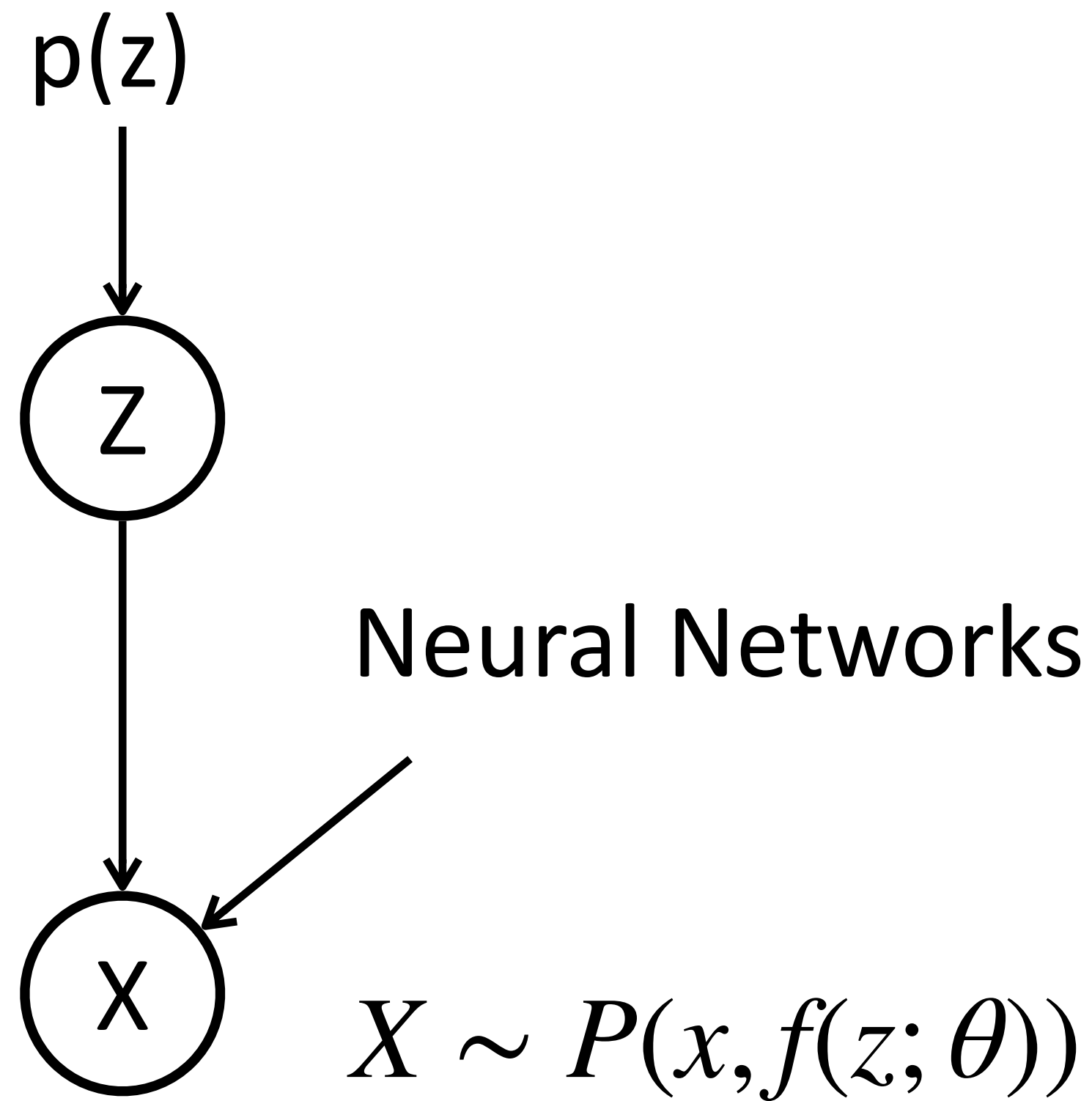
$$Q(z) = P(z|x) \propto P(z)P(x|z)$$

This is ok?

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

$P(x|z) \sim$

# Let's try EM



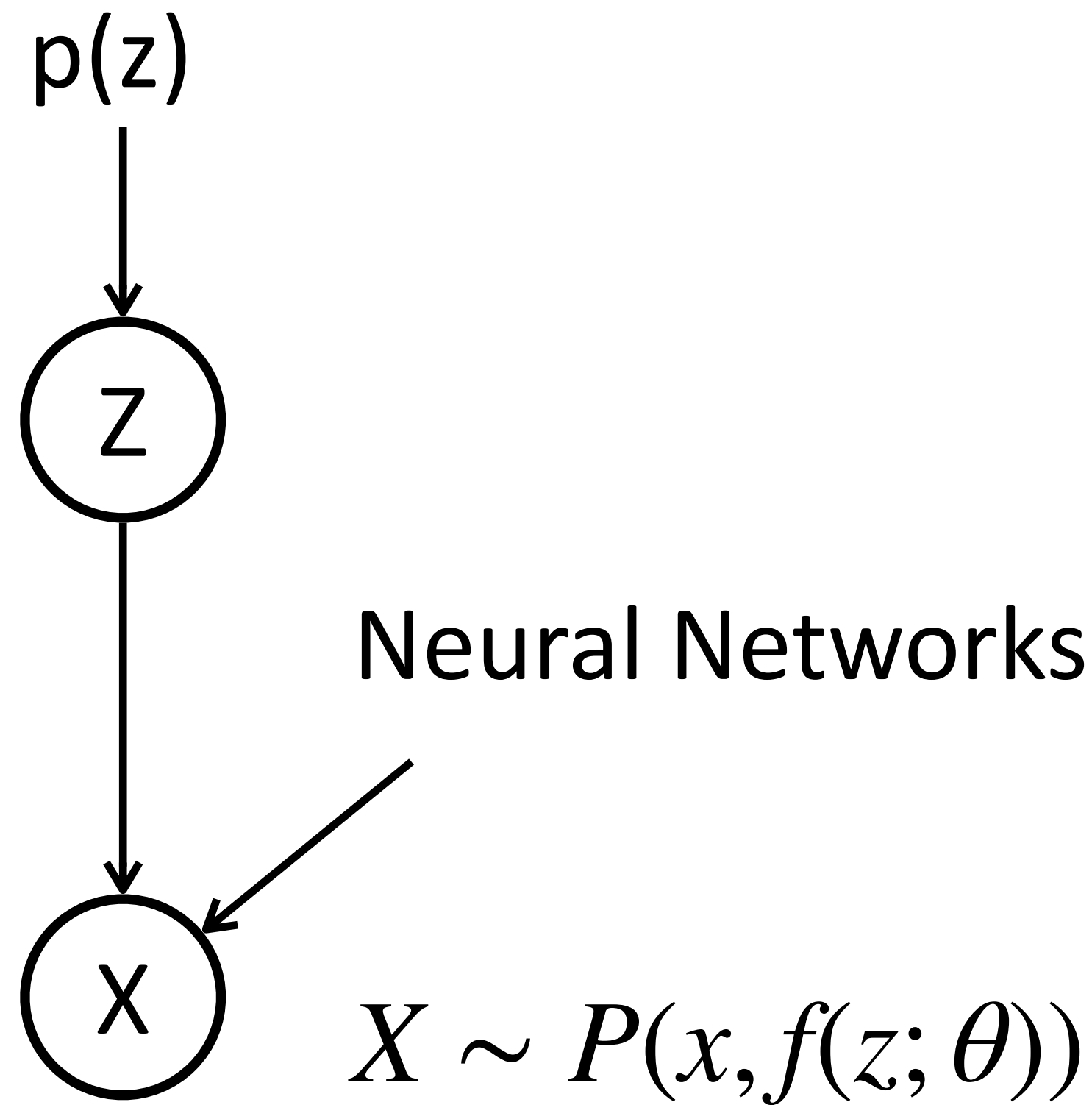
E-Step: compute  $P(z|x)$

$$Q(z) = P(z|x) \propto P(z)P(x|z) \quad \text{This is ok?}$$

M-Step: the ELBO objective

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

# Let's try EM



E-Step: compute  $P(z|x)$

$P(x|z)$

$$Q(z) = P(z|x) \propto P(z)P(x|z)$$

This is ok?

M-Step: the ELBO objective

$Q(z)$

$$\operatorname{argmax}_{\theta} \sum_z Q(z) \log p(x, z; \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim Q(z)} \log p(x, z; \theta)$$

In most cases, we cannot do the sum, and cannot easily sample from  $Q(z)$  either

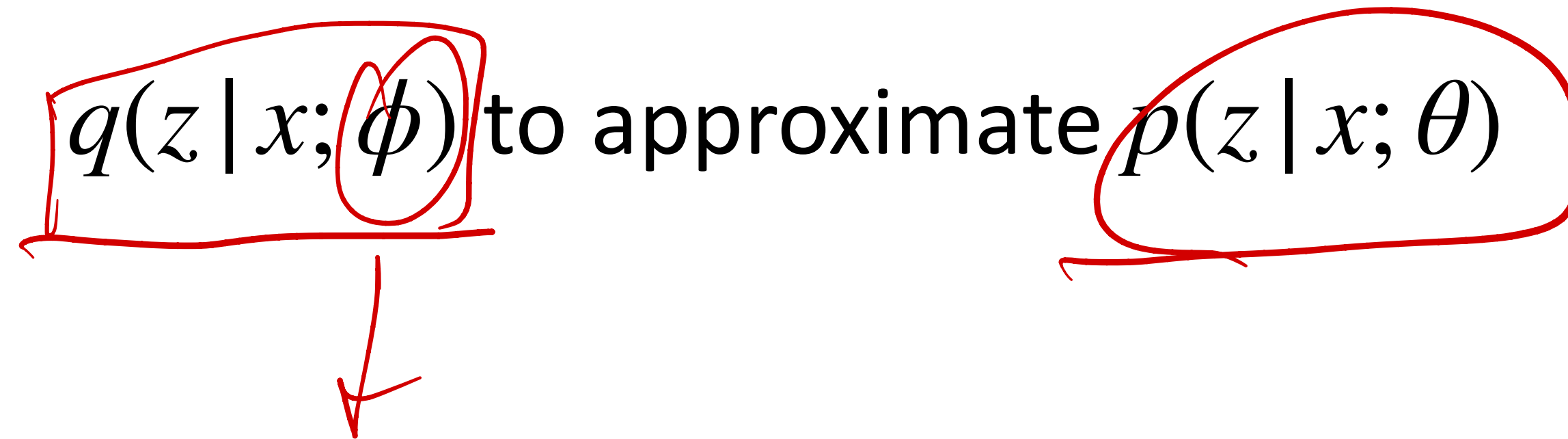
# Approximate Posterior

We need an easy-to-sample distribution to approximate  $P(z|x)$

# Approximate Posterior

We need an easy-to-sample distribution to approximate  $P(z|x)$

$q(z|x; \phi)$  to approximate  $p(z|x; \theta)$



# Approximate Posterior

We need an easy-to-sample distribution to approximate  $P(z | x)$

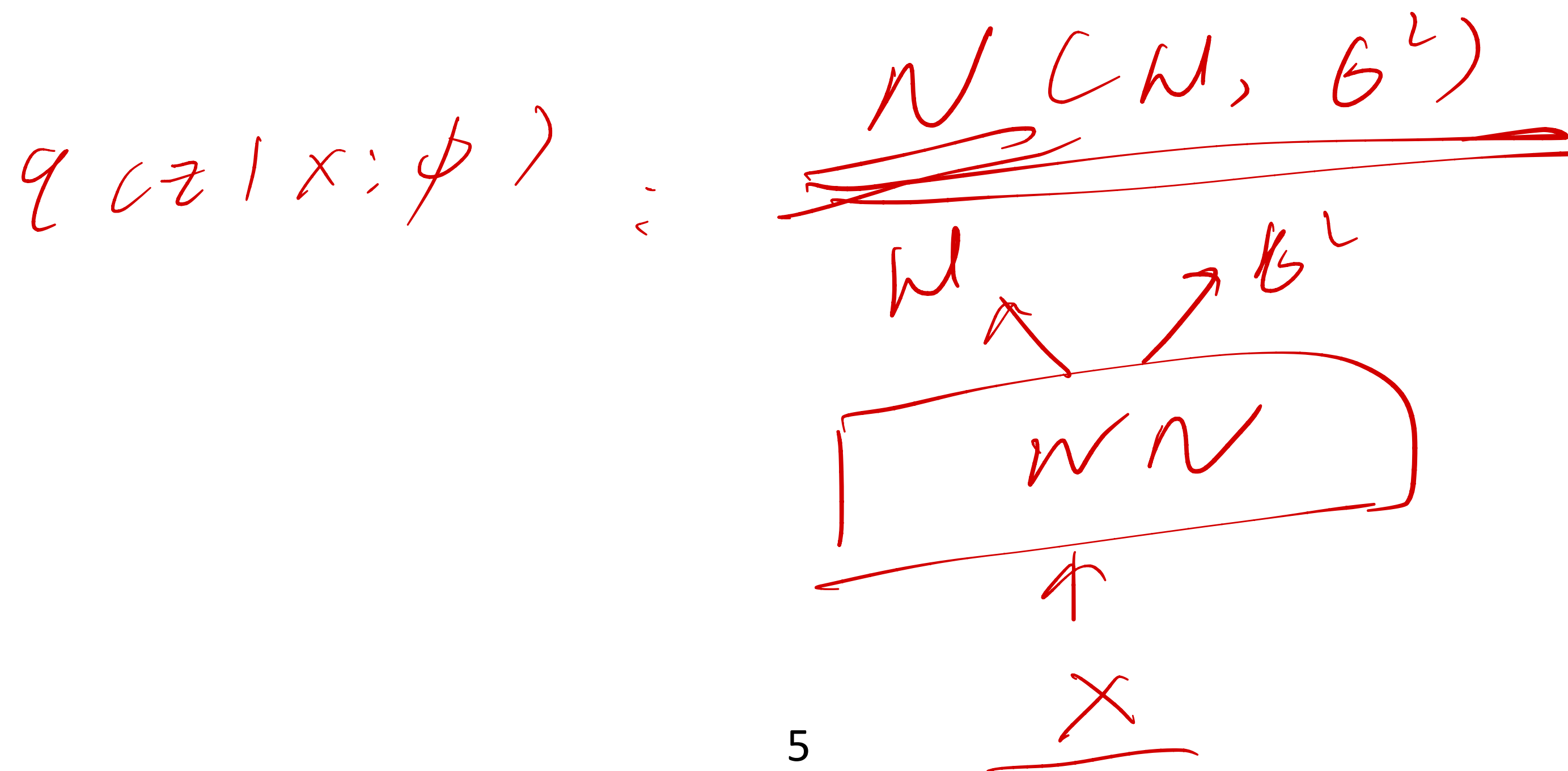
$q(z | x; \phi)$  to approximate  $p(z | x; \theta)$  Why conditioned on  $x$ ?

# Approximate Posterior

We need an easy-to-sample distribution to approximate  $P(z|x)$

$q(z|x; \phi)$  to approximate  $p(z|x; \theta)$  Why conditioned on  $x$ ?

$\phi$  is the parameter for the approximate function,  $\theta$  is the generative model parameter



# Approximate Posterior

We need an easy-to-sample distribution to approximate  $P(z|x)$

$q(z|x; \phi)$  to approximate  $p(z|x; \theta)$  Why conditioned on  $x$ ?

$\phi$  is the parameter for the approximate function,  $\theta$  is the generative model parameter

How to train  $q(z|x; \phi)$ , what would be the loss to find  $\phi$ ?

E step:  $q(z|x; \phi) \rightarrow p(z|x; \theta)$

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is  $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$ ?

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is  $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$ ?

ELBO is maximized when  $Q(z)$  is equal to  $p(z|x)$

$$Q(z) = p(z|x)$$

# Recap: ELBO

equal to

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

argmin  $\phi$   $\text{KL}(q(z|x; \phi) || p(z|x, \theta))$

What is  $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$ ?

ELBO is maximized when  $Q(z)$  is equal to  $p(z|x)$

Therefore, we can approximate the true posterior by maximizing ELBO:

$$\text{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

constant, not depend  $Q$

$$ELBO = \log P(x; \theta) - KL(Q \parallel P(z|x; \theta))$$

$$\operatorname{argmax}_Q ELBO = \operatorname{argmin}_Q KL(Q \parallel P(z|x; \theta))$$

# Recap: ELBO

$$\text{ELBO}(x; Q, \theta) = \sum_z Q(z) \log \frac{p(x, z; \theta)}{Q(z)}$$

What is  $\text{argmax}_{Q(z)} \text{ELBO}(x; Q, \theta)$ ?

ELBO is maximized when  $Q(z)$  is equal to  $p(z|x)$

Therefore, we can approximate the true posterior by maximizing ELBO:

$$\text{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

Variational Inference

$\phi$  $\theta$ 

# Training VAEs



E-Step:

$q(z|x; \phi)$  by NN

 $\mathbb{E} \text{LBO}$ 

$$\text{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

M-Step:

 $\bar{\mathbb{E}} \text{LBO}$ 

$$\text{argmax}_{\theta} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

$$\bar{\mathbb{E}} \text{LBO} = \mathbb{E}_{z \sim q(z|x; \phi)} \log \frac{p(x, z)}{q(z|x)}$$

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

Same objective, different parameters to optimize

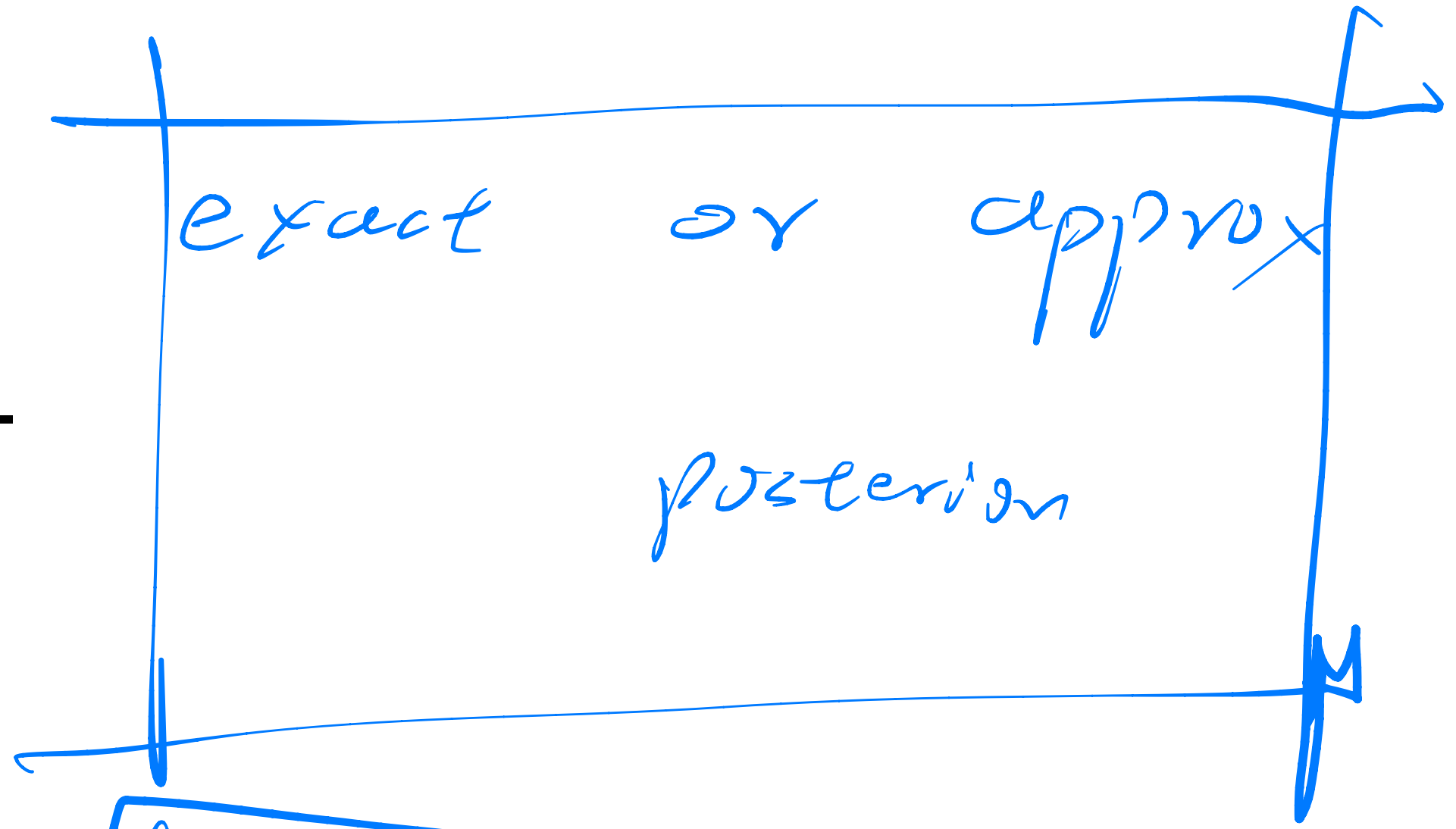
variational from "variational inference"

# Training VAEs

$p(z|x)$

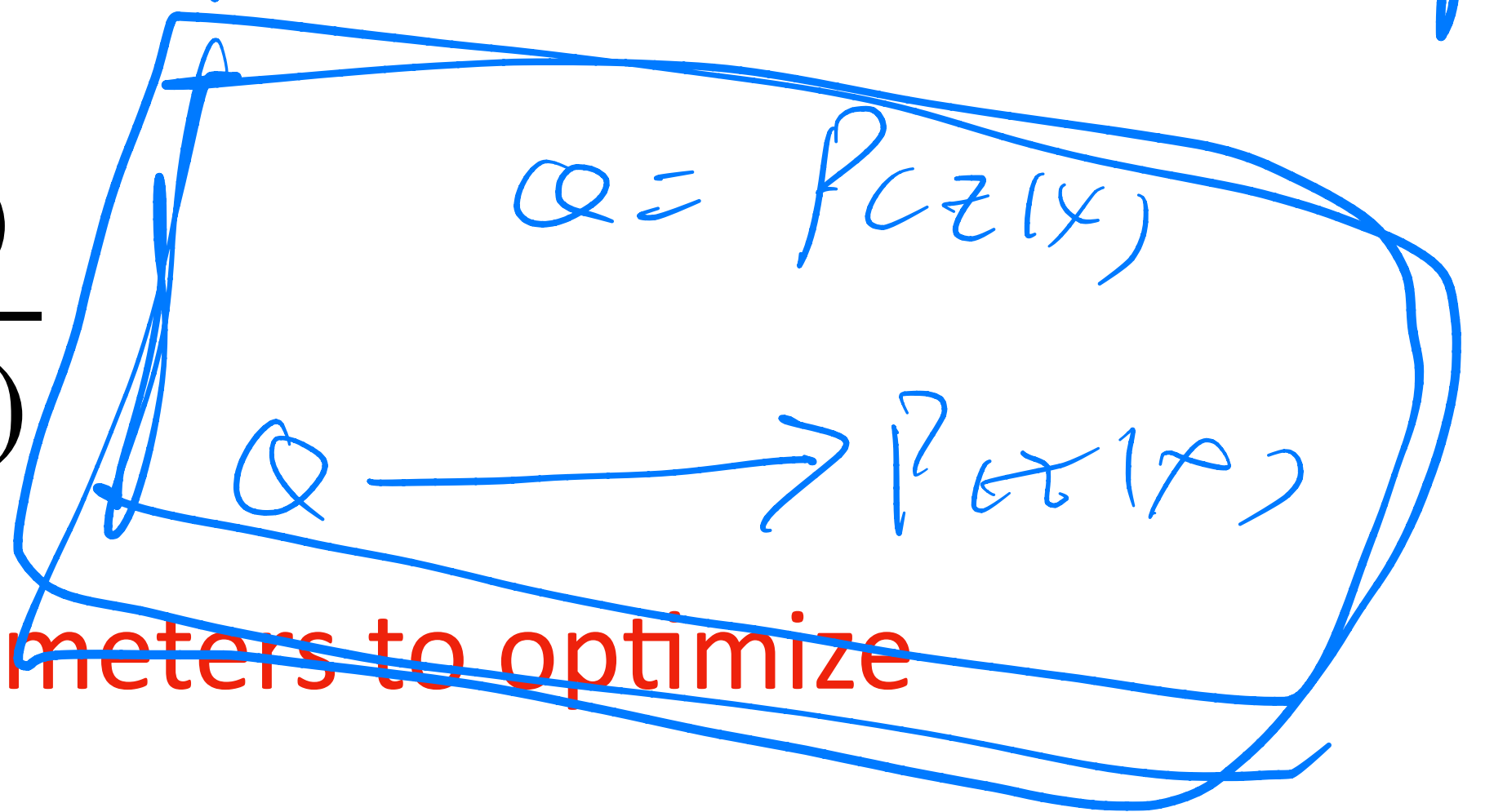
E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$



M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$



Same objective, different parameters to optimize

Because we use approximate rather than exact posterior, it is also called Variational EM

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

Same objective, different parameters to optimize

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

Same objective, different parameters to optimize

Because we use approximate rather than exact posterior, it is also called Variational EM

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

# Training VAEs

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

We use MC sampling to approximate expectation and use gradient descent to optimize  $\theta$

# Training VAEs

$$z = f(\phi)$$

$$\frac{dz}{d\phi}$$

derivative

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

Can we do gradient descent over  $\phi$ ?

M-Step:

$$\operatorname{argmax}_{\theta} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

$$z \sim \mathcal{N}(\mu, \sigma^2) = q(z|x; \phi)$$

$z$  depends on  $\phi$

We use MC sampling to approximate expectation and use gradient descent to optimize  $\theta$

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

First, we cannot do sum, but we can sample  $z_i$  from  $q(z | x; \phi)$ , which depends on  $\phi$ , how do we propagate gradients to  $\phi$ ?

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

First, we cannot do sum, but we can sample  $z_i$  from  $q(z | x; \phi)$ , which depends on  $\phi$ , how do we propagate gradients to  $\phi$ ?

Try to express  $z$  as a deterministic function  $z = g_{\phi}(\epsilon, x)$ , where  $\epsilon$  is an auxiliary random variable

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)}$$

*Handwritten notes:*  
- Above the sum:  $\frac{dz}{d\mu}$   
- To the right:  $\Sigma \sim \mathcal{N}(0, 1)$

First, we cannot do sum, but we can sample  $z_i$  from  $q(z|x; \phi)$ , which depends on  $\phi$ , how do we propagate gradients to  $\phi$ ?

Try to express  $z$  as a deterministic function  $z = g_{\phi}(\epsilon, x)$ , where  $\epsilon$  is an auxiliary random variable

$$z \sim \mathcal{N}(\mu, \sigma^2) \longrightarrow z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

*Handwritten notes:*  
- The entire equation is circled in red.  
-  $\epsilon$  is circled in red.  
-  $\mathcal{N}(0, 1)$  is circled in red.  
-  $\odot$  is circled in red.  
-  $\epsilon \sim \mathcal{N}(0, 1)$  is also written as  $\mathcal{N}(0, 1)$  in red above the equation.

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

First, we cannot do sum, but we can sample  $z_i$  from  $q(z | x; \phi)$ , which depends on  $\phi$ , how do we propagate gradients to  $\phi$ ?

Try to express  $z$  as a deterministic function  $z = g_{\phi}(\epsilon, x)$ , where  $\epsilon$  is an auxiliary random variable

$$z \sim N(\mu, \sigma^2) \longrightarrow z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(0, 1)$$

Can you verify  $z$  in this equation is Gaussian?

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

For every gradient step (assuming batch size=1):

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

For every gradient step (assuming batch size=1):

1. Randomly sample  $\epsilon^{(i)} \sim N(0,1)$

# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

For every gradient step (assuming batch size=1):

1. Randomly sample  $\epsilon^{(i)} \sim N(0,1)$
2. Obtain z sample as  $z^{(i)} = \mu + \sigma \odot \epsilon^{(i)}$

# Reparameterization Trick

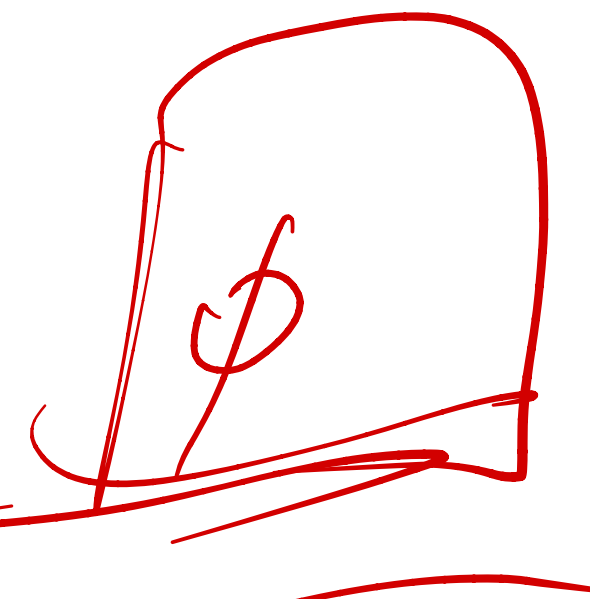
E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

For every gradient step (assuming batch size=1):

1. Randomly sample  $\epsilon^{(i)} \sim N(0,1)$
2. Obtain z sample as  $z^{(i)} = \mu + \sigma \odot \epsilon^{(i)}$

3. Perform gradient descent w.r.t.  $\log \frac{p(x, z^{(i)}; \theta)}{q(z^{(i)} | x; \phi)}$



# Reparameterization Trick

E-Step:

$$\operatorname{argmax}_{\phi} \sum_z q(z | x; \phi) \log \frac{p(x, z; \theta)}{q(z | x; \phi)}$$

For every gradient step (assuming batch size=1):

1. Randomly sample  $\epsilon^{(i)} \sim N(0,1)$
2. Obtain z sample as  $z^{(i)} = \mu + \sigma \odot \epsilon^{(i)}$
3. Perform gradient descent w.r.t.  $\log \frac{p(x, z^{(i)}; \theta)}{q(z^{(i)} | x; \phi)}$

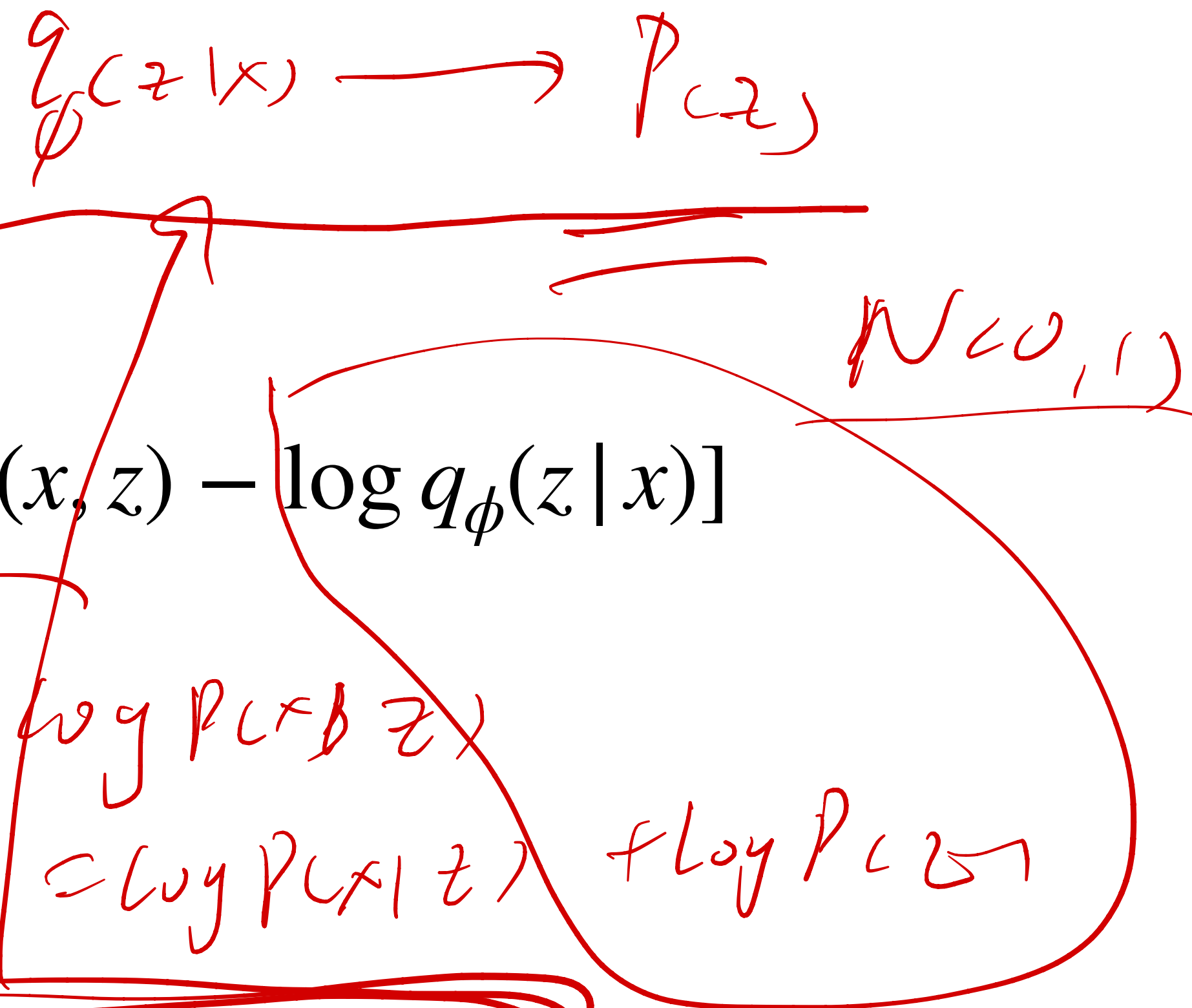
We can now propagate  
gradients from z to  $\phi$

# ELBO

$$\sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)} = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

# ELBO

$$\sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)} = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$



ELBO is implemented with the following form:

$$\max \underbrace{\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(\mathbf{x}|z)]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) || p(z))}_{\text{KL Regularizer}}$$



Handwritten notes:

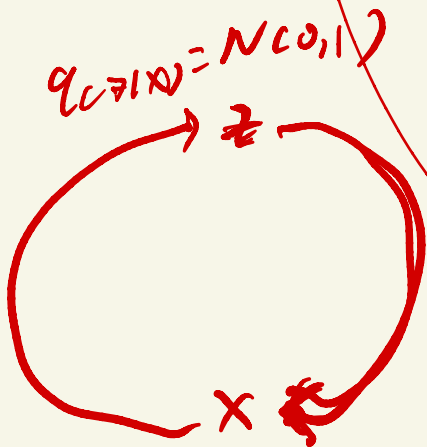
- $\mathbb{E}_{z \sim q_\phi(z|x)} \log \frac{p(x|z)}{p(z)}$

$$E_{q(z|x)} \log p(x|z)$$

Reconstruction

$$K L C q(z|x) \| p(z)$$

regularizer



$$q(z|x) = p(z)$$

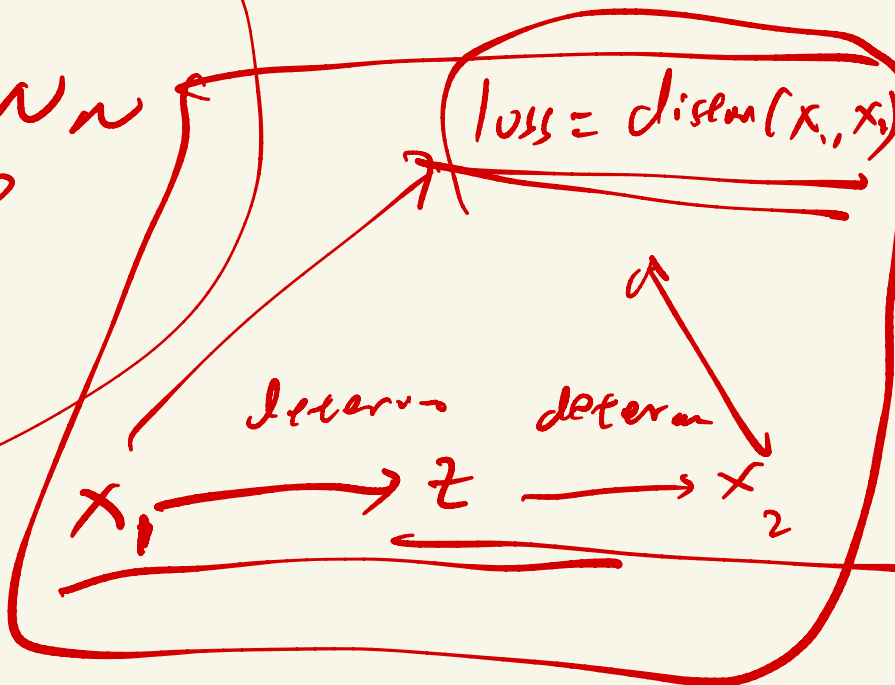
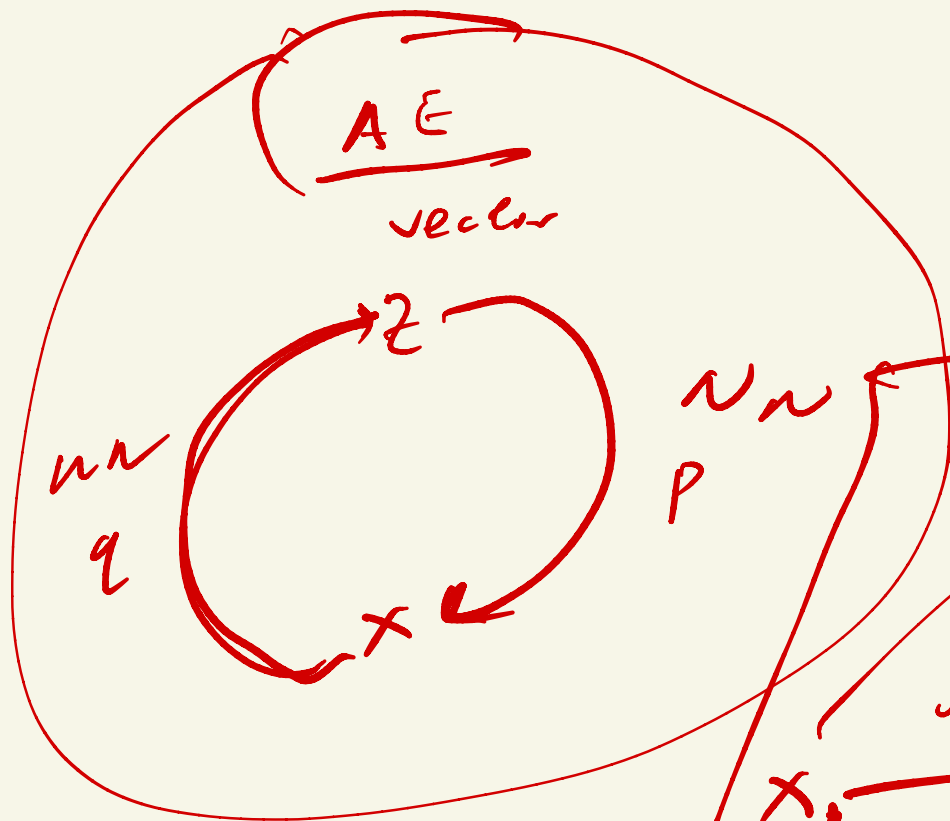
# ELBO

$$\sum_z q(z|x; \phi) \log \frac{p(x, z; \theta)}{q(z|x; \phi)} = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)]$$

ELBO is implemented with the following form:

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Autoencoder



# ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

# ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Autoencoder Loss

# ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Autoencoder Loss

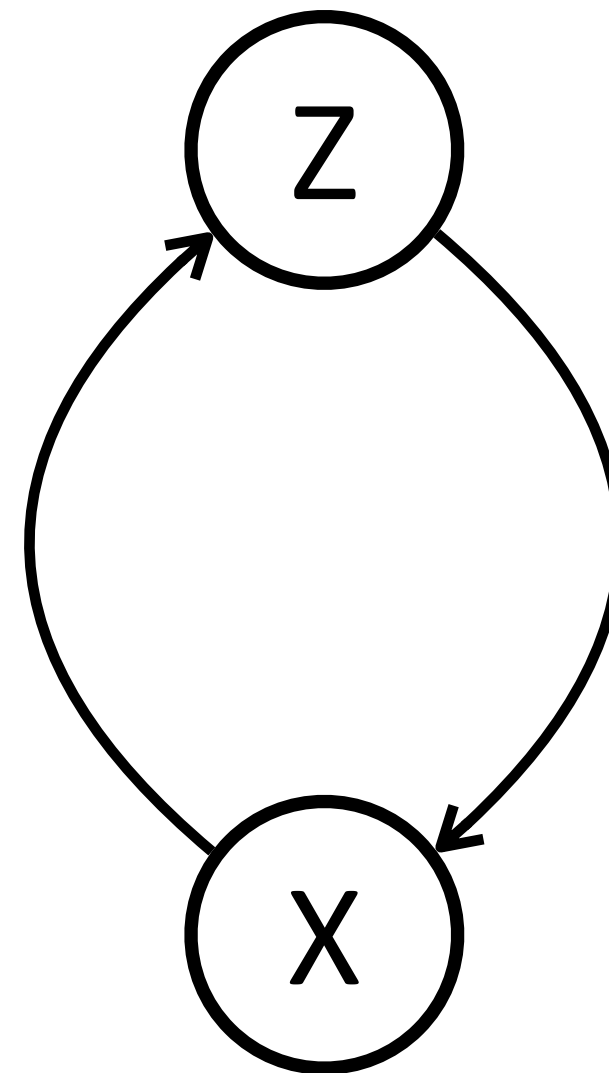
$q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z})$  are both Gaussian,  
there is a closed-form for this

# ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Autoencoder Loss

$q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z})$  are both Gaussian,  
there is a closed-form for this

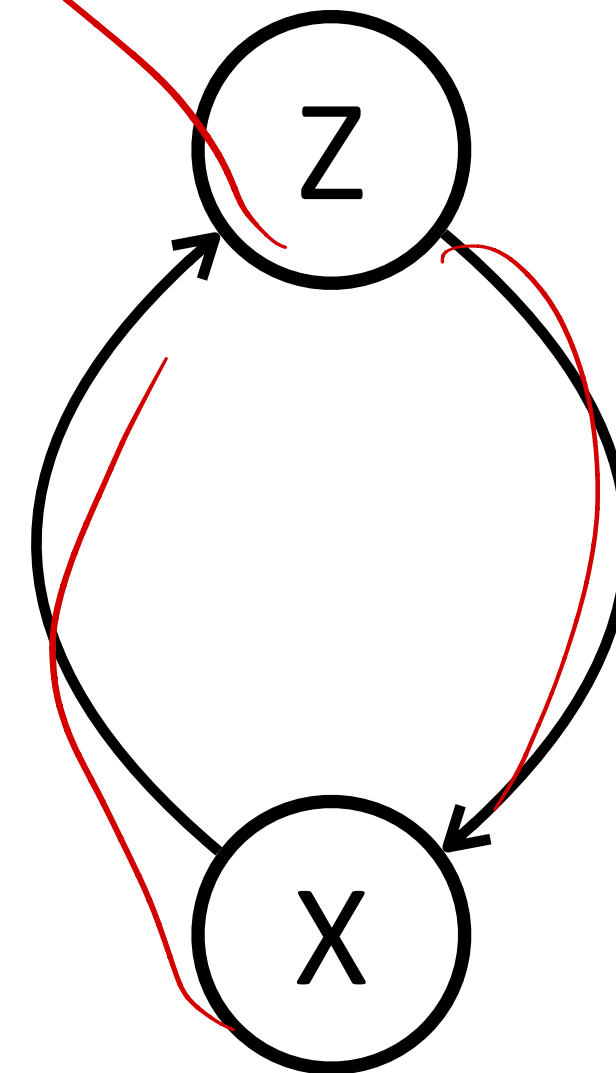


# ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Autoencoder Loss

$q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z})$  are both Gaussian, there is a closed-form for this



$$\text{ELBO} = \log p(\mathbf{x}|\theta) - \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

This is why it is called variational "autoencoder"

# ELBO

$$D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

# ELBO

$$D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

$$\int q_{\theta}(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} = \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) d\mathbf{z}$$
$$= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\mu_j^2 + \sigma_j^2)$$

J is the dimensionality of z

# ELBO

$$D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})).$$

$$\begin{aligned}\int q_{\theta}(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\mu_j^2 + \sigma_j^2)\end{aligned}$$

J is the dimensionality of z

$$\begin{aligned}\int q_{\theta}(\mathbf{z}) \log q_{\theta}(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2)\end{aligned}$$

# ELBO

$$D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

$$\begin{aligned}\int q_{\theta}(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\mu_j^2 + \sigma_j^2)\end{aligned}$$

J is the dimensionality of z

$$\begin{aligned}\int q_{\theta}(\mathbf{z}) \log q_{\theta}(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2)\end{aligned}$$

$$-D_{\text{KL}}(q_{\phi}(\mathbf{z})||p_{\theta}(\mathbf{z})) = \int q_{\theta}(\mathbf{z}) (\log p_{\theta}(\mathbf{z}) - \log q_{\theta}(\mathbf{z})) d\mathbf{z}$$

$$= \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2)$$

# Training VAEs

*until convergence*

$q$   $p(z|x)$

*parameter trick*

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

M-Step:

$$\operatorname{argmax}_{\theta} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

# Training VAEs

loss  $\mathbb{E} \log p(\mathbf{z})$

E-Step:

$\text{argmax}_{\phi}$

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

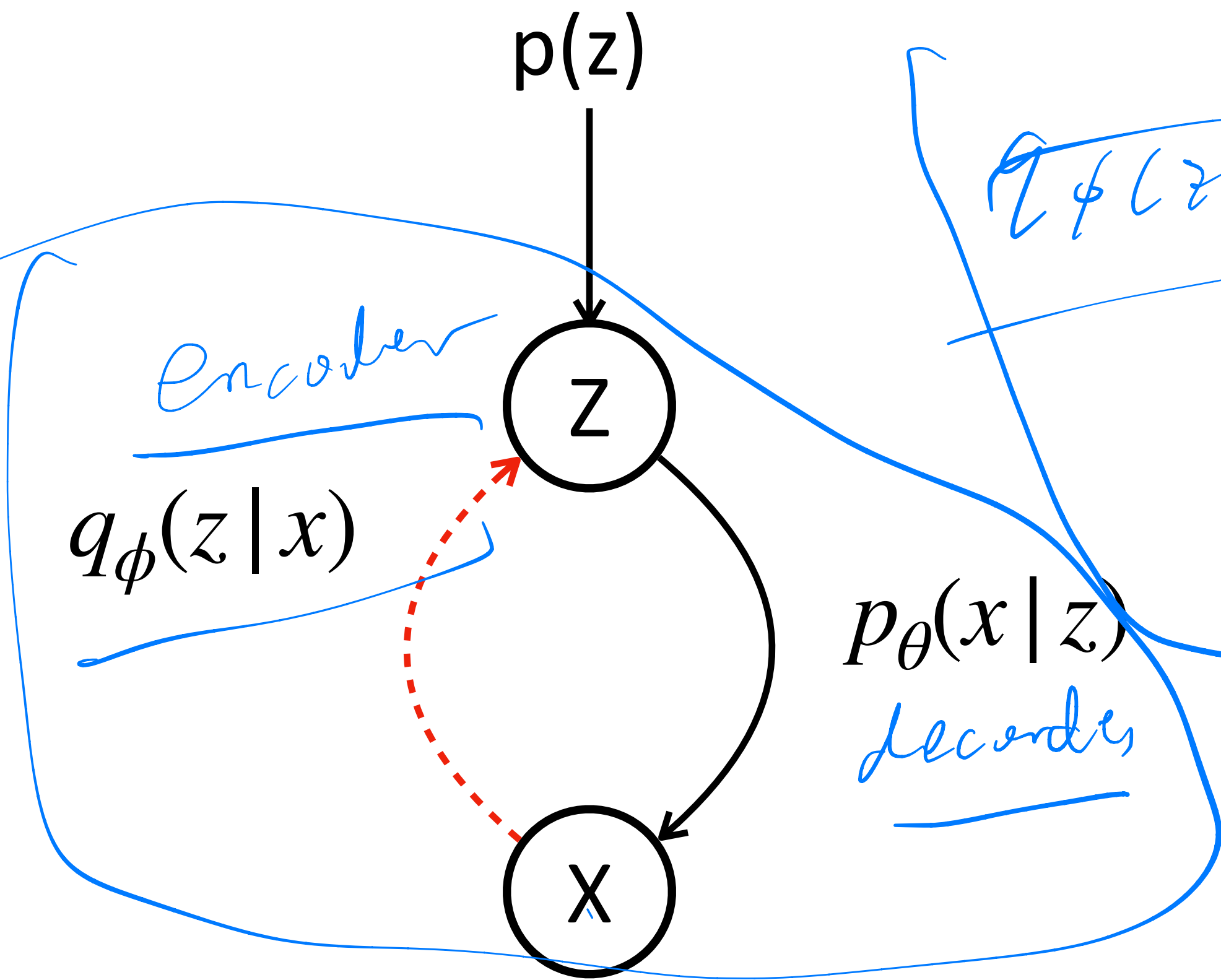
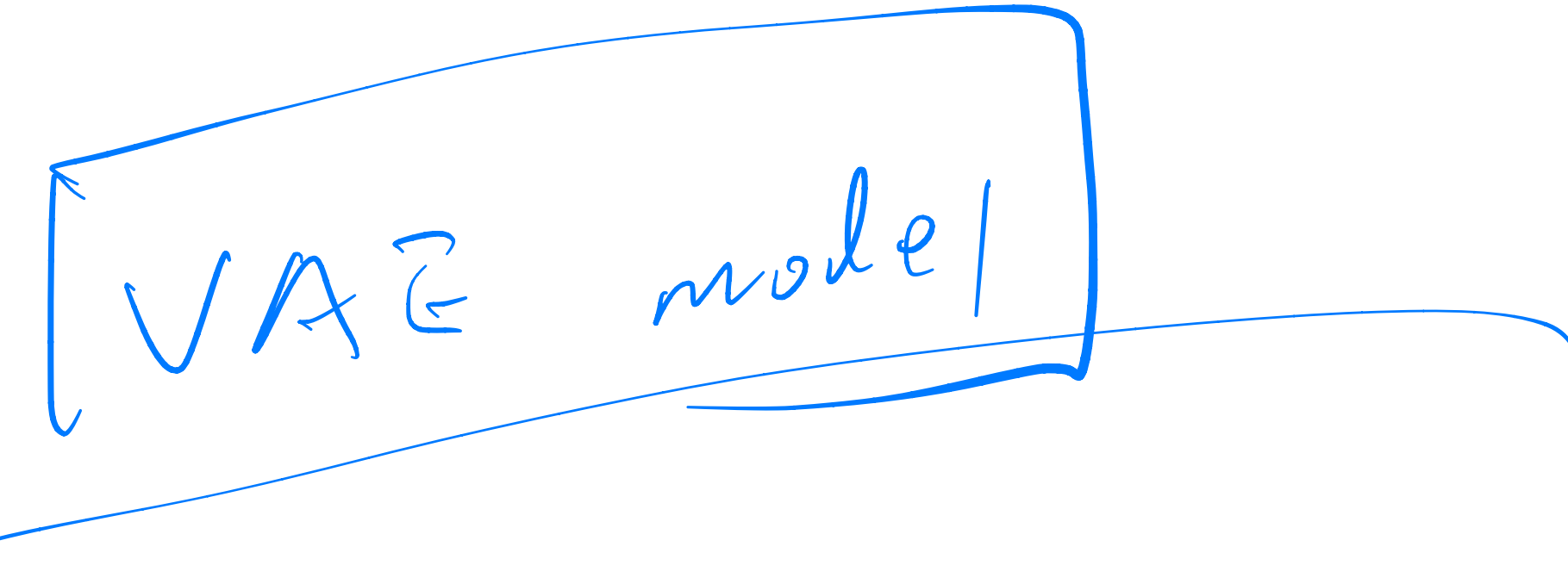
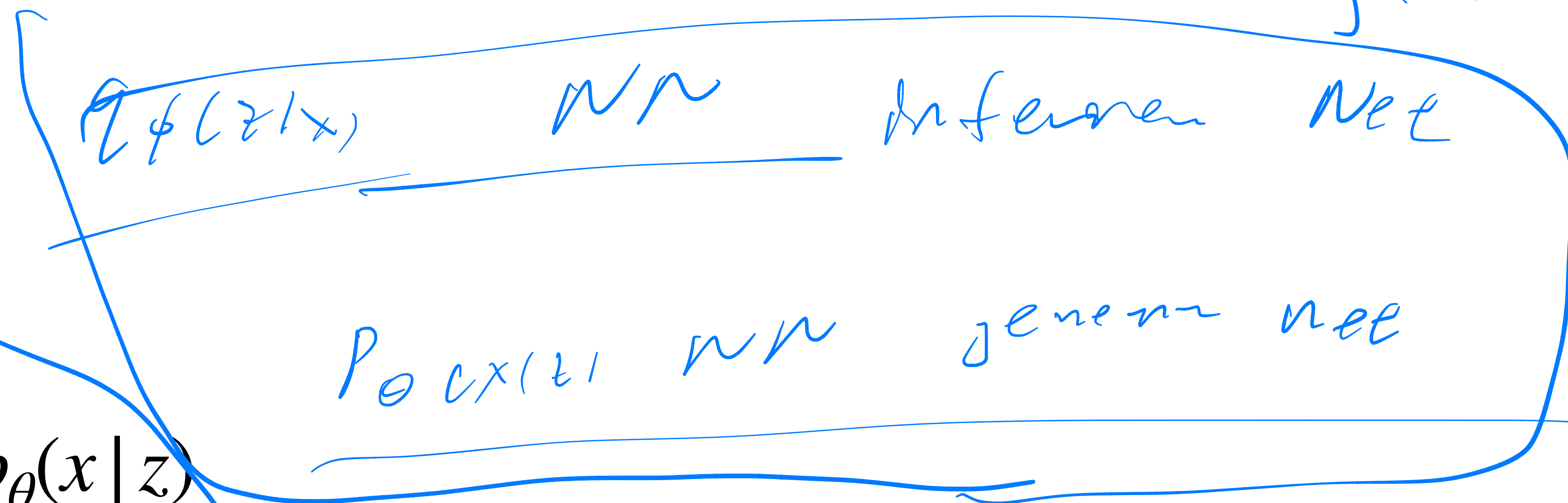
M-Step:

$\text{argmax}_{\theta}$

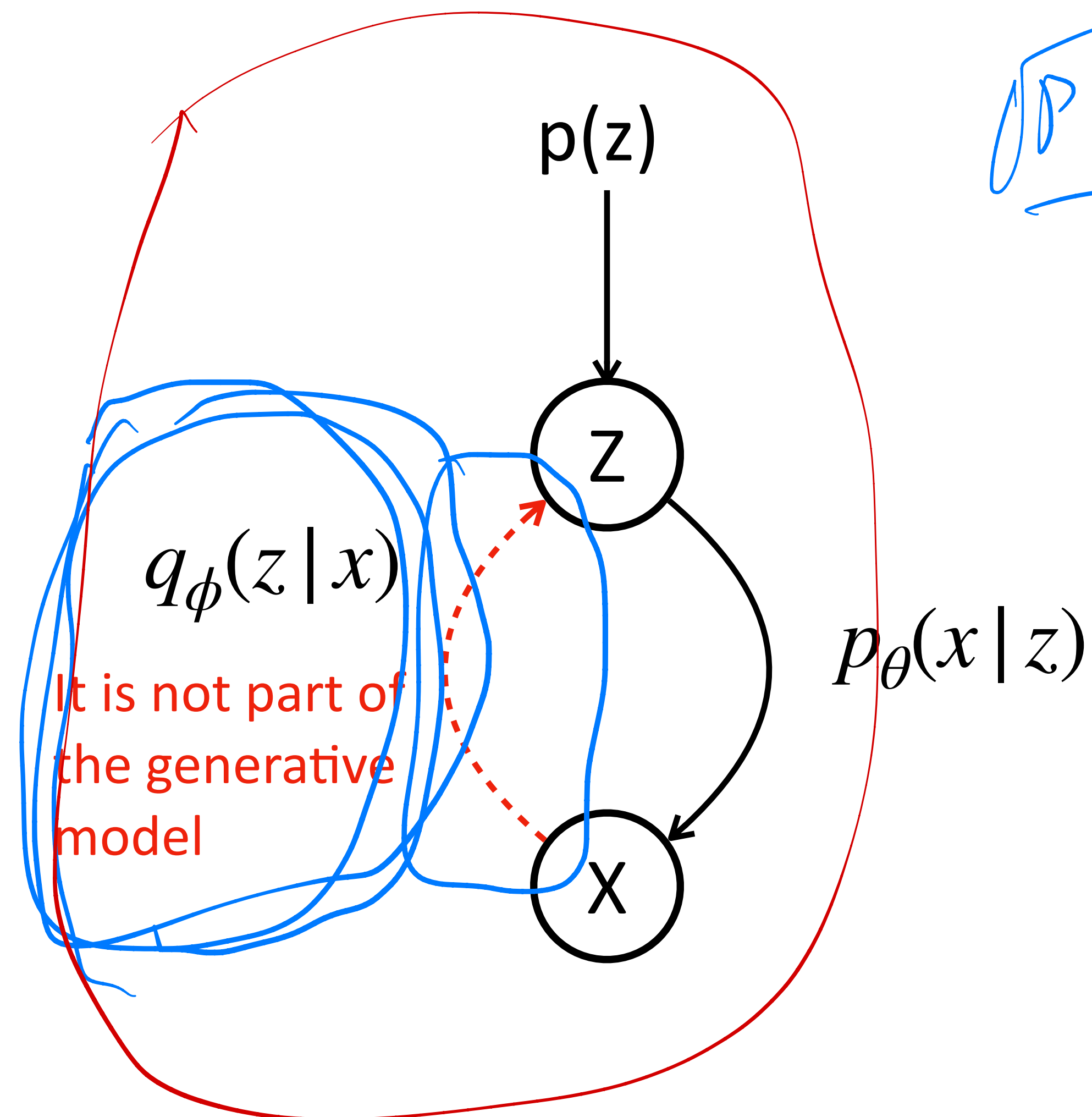
$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

Intuitively we hope to approximate  $p(\mathbf{z}|\mathbf{x})$  with  $q(\mathbf{z}|\mathbf{x})$  accurately in the E-step, to approximate the true EM algorithm

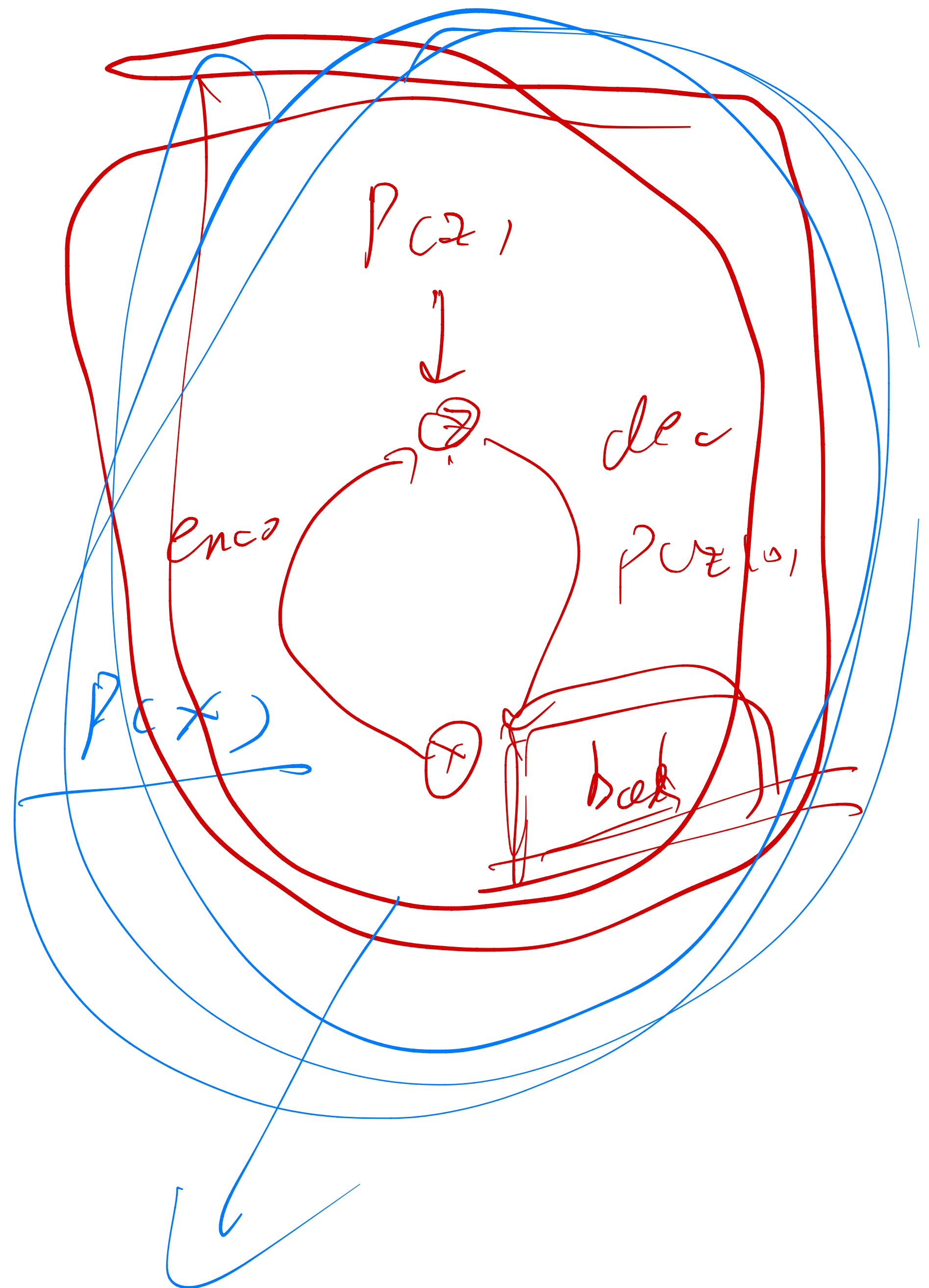
# Review VAE



# Review VAE

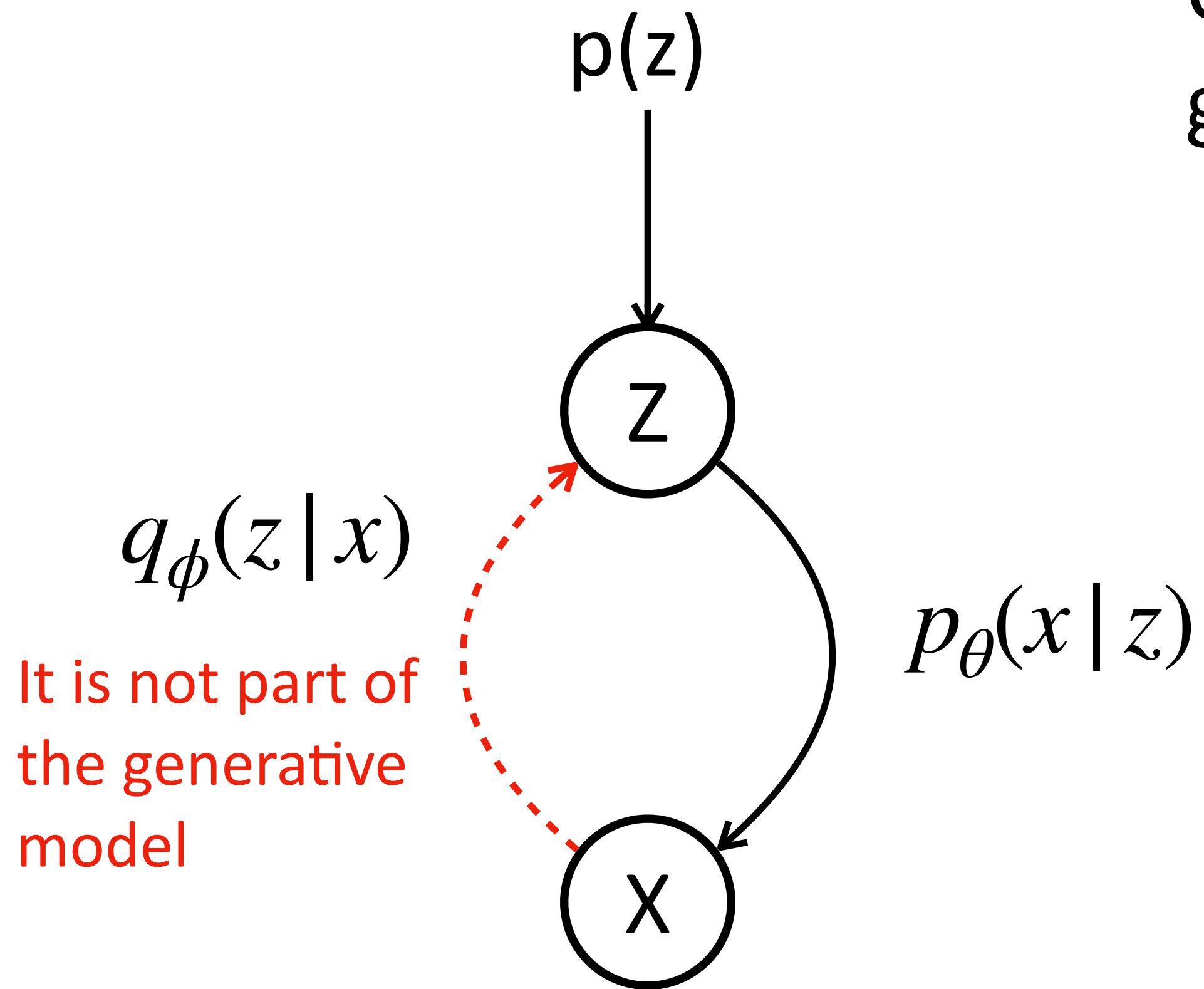


PGM



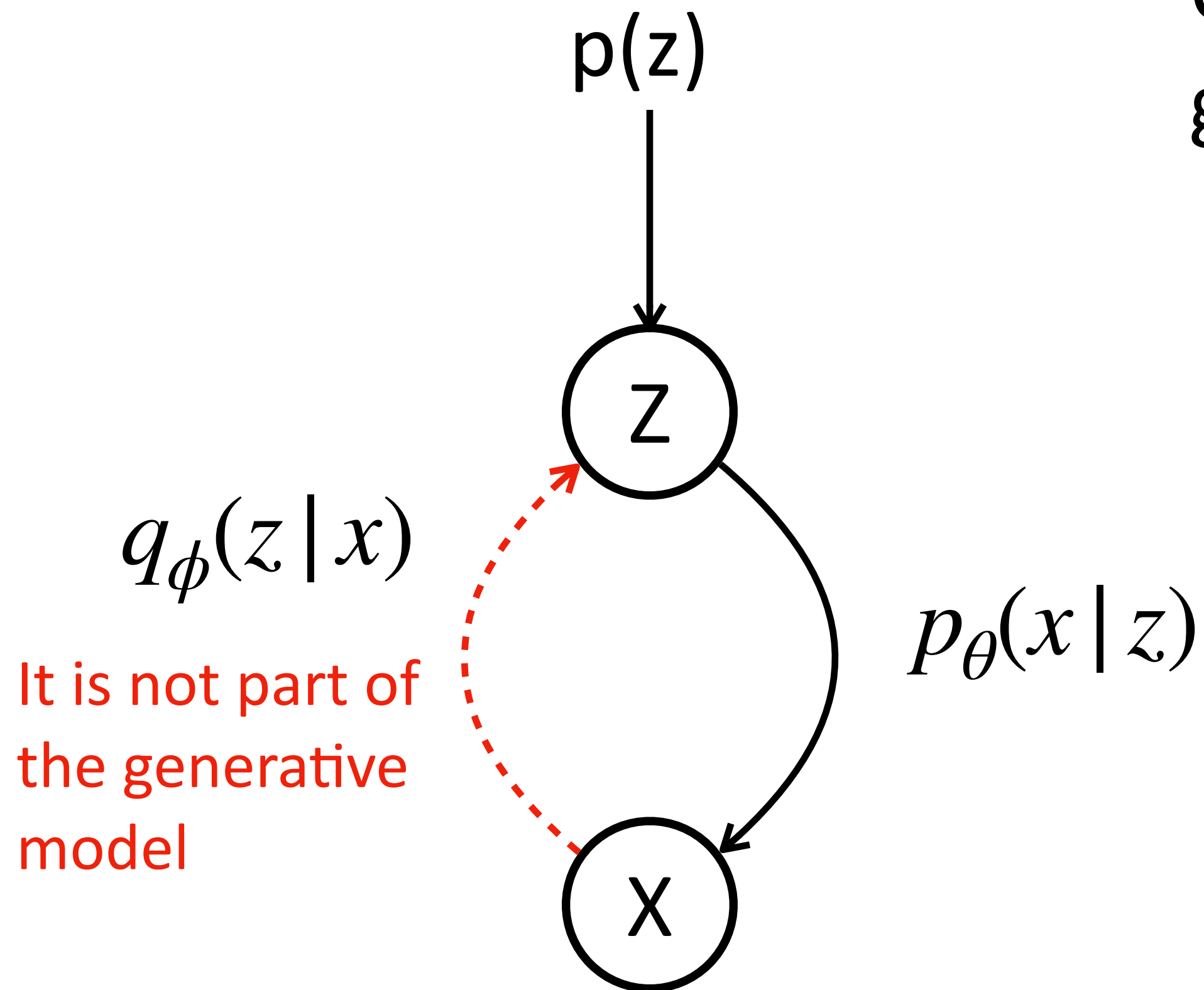
# Review VAE

Only the right (black) part defines the generative model, and the distribution



# Review VAE

Only the right (black) part defines the generative model, and the distribution



$p_{\theta}(x | z)$ : generative network/decoder

$q_{\phi}(z | x)$ : inference network/encoder

# Review VAE

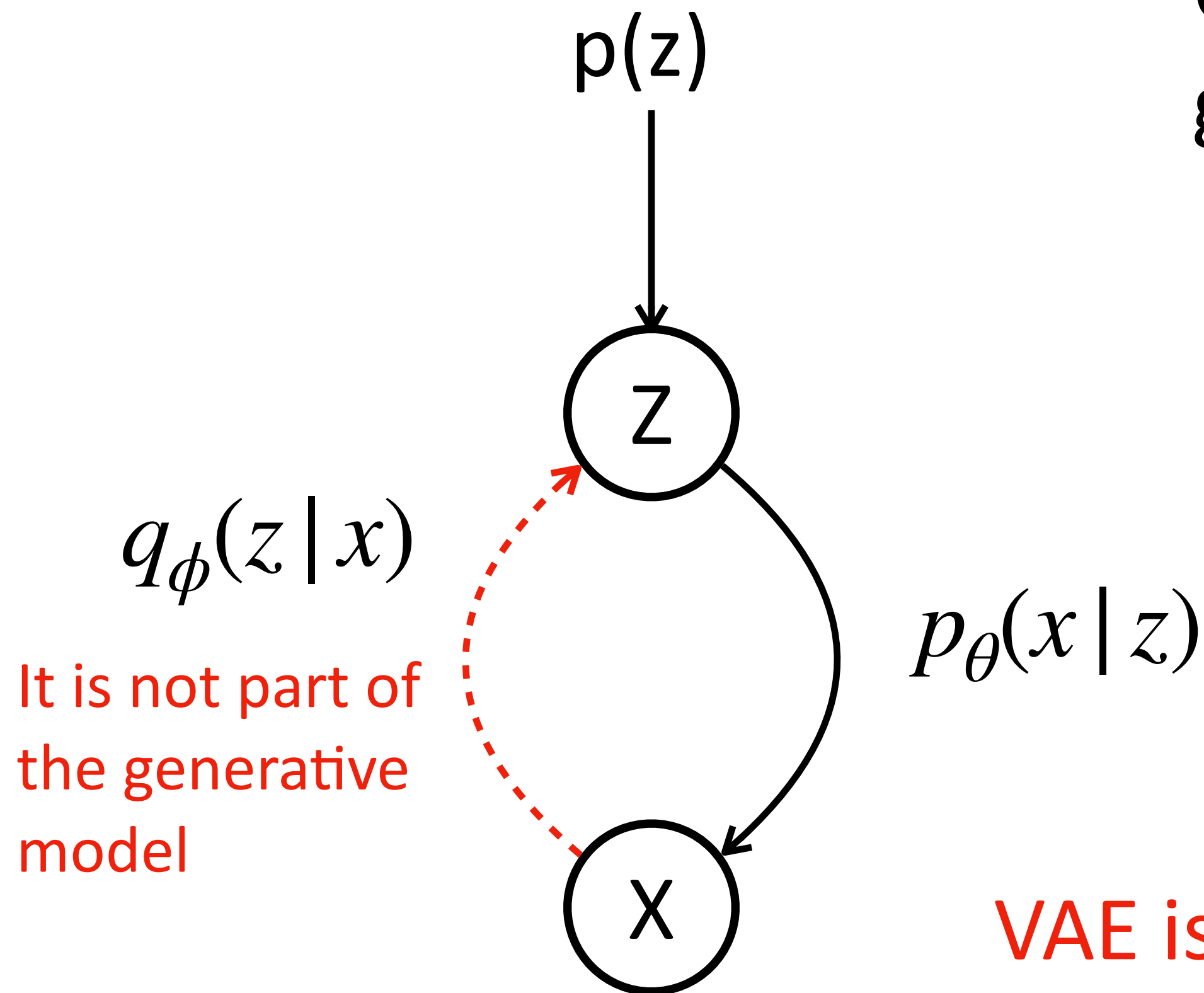
Only the right (black) part defines the generative model, and the distribution

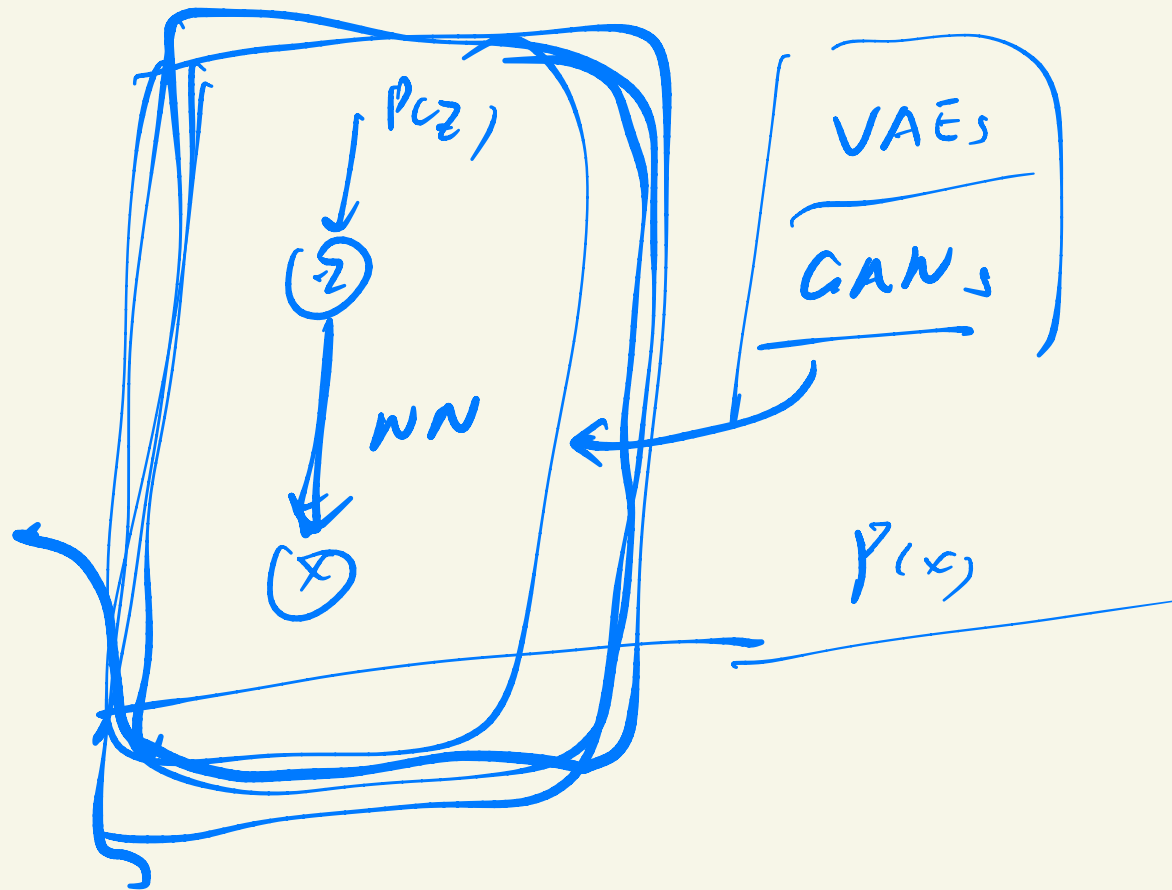
$p_{\theta}(x | z)$ : generative network/decoder

$q_{\phi}(z | x)$ : inference network/encoder

VAE := model + training method

VAE is a name to represent both the model  $p(x)$  and the inference network that is used to train the model, but do not confuse them together





# Training VAEs

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$

---

$$\mathbb{E} \text{ELBO} \leq \log p_{\text{UX}}(\theta)$$

# Training VAEs

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$

---

# Training VAEs

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\Sigma \quad M$

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$

---

End-to-end, because the objectives are the same (ELBO)

# Training VAEs

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$


---

End-to-end, because the objectives are the same (ELBO)

VAE training is optimizing ELBO with gradient descent

true goal

# Recap: EM is Hill Climbing




$\log p(x; \theta)$



ELBO

$\log p(x; \theta)$

EM



Larger




$ELBO \leq \log p(x)$

# Recap: EM is Hill Climbing



$\log p(x; \theta)$

Only related to  $\theta$ , no  $z$



Larger



ELBO



# Recap: EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

# Recap: EM is Hill Climbing



$\log p(x; \theta)$



ELBO



E-step:  $Q(z) = p(z | x; \theta)$ , making ELBO tight



Larger

$ELBO = \log p(x)$

# Recap: EM is Hill Climbing



$\log p(x; \theta)$



ELBO



E-step:  $Q(z) = p(z | x; \theta)$ , making ELBO tight

“dog” doesn't change, because  $\theta$  does not change



Larger

# Recap: EM is Hill Climbing



$\log p(x; \theta)$



ELBO



Larger

E-step:  $Q(z) = p(z | x; \theta)$ , making ELBO tight  
“dog” doesn’t change, because  $\theta$  does not change

# Recap: EM is Hill Climbing



$$\log p(x; \theta)$$

$$\geq ELBO$$



Larger



ELBO



M-step:  $\max_{\theta} ELBO$

$\theta$



ELBO becomes larger, and it is not tight anymore because posterior changes

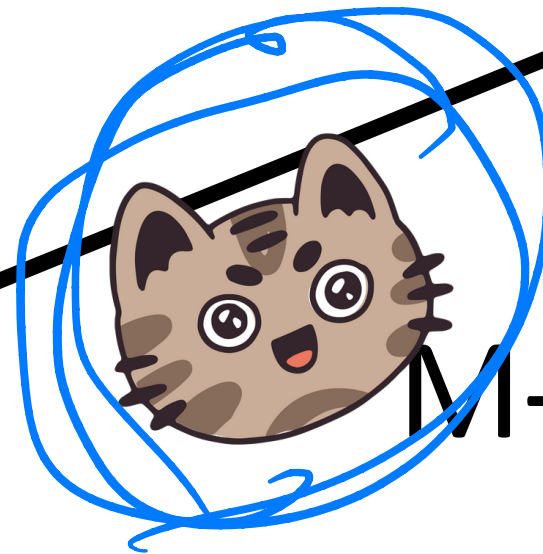
# Recap: EM is Hill Climbing



$\log p(x; \theta)$



ELBO



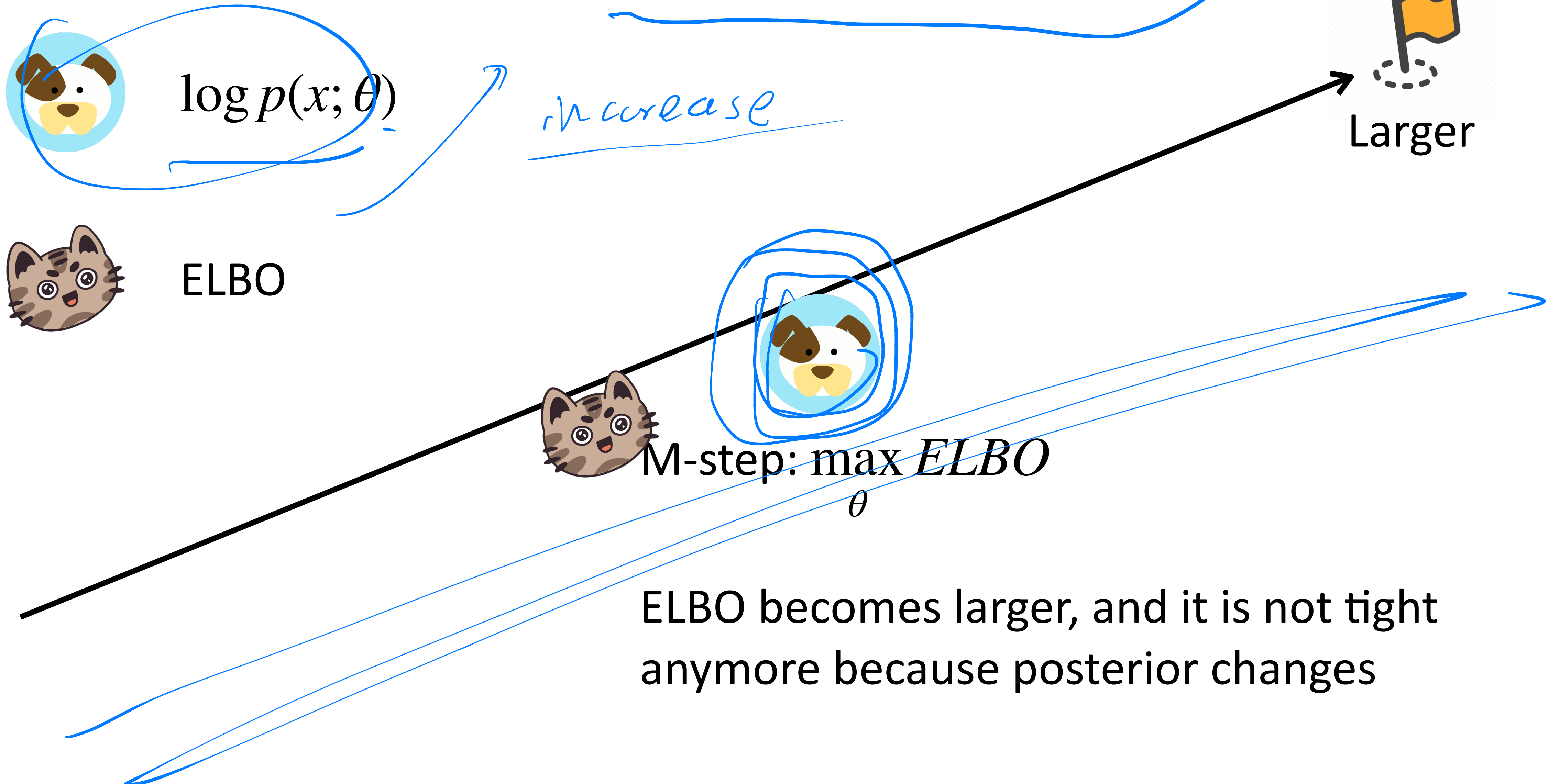
M-step:  $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes



Larger

# Recap: EM is Hill Climbing



$\log p(x; \theta)$

*increase*

Larger

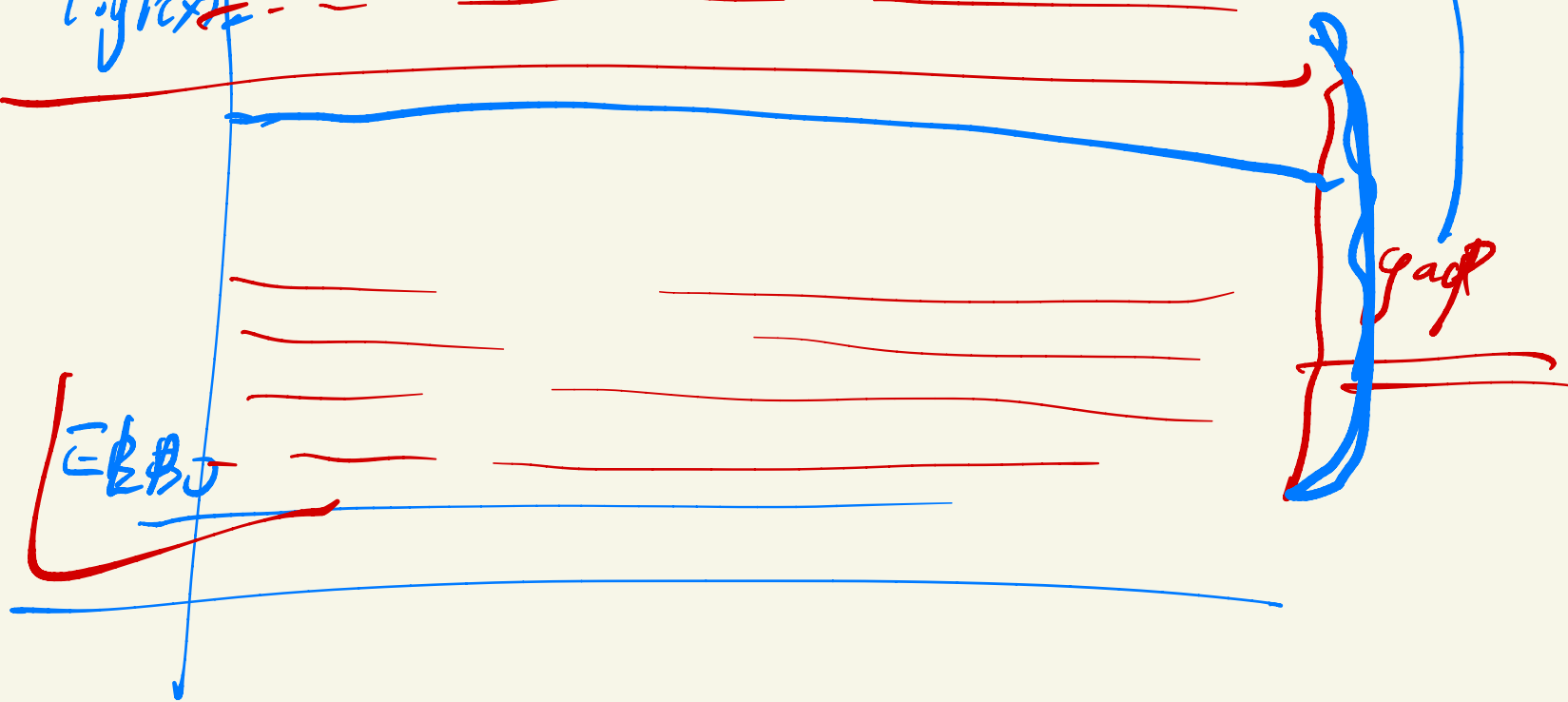
ELBO

M-step:  $\max_{\theta} ELBO$

ELBO becomes larger, and it is not tight anymore because posterior changes

$$ELBO = \log P(x) + \underbrace{KL(P(z|x) || P(z))}_{\text{gap}} - \underbrace{KL(Q || P(z|x))}_{\text{gap}}$$

$\log P(x)$



# Is VAE training still Hill Climbing?

$q_\phi$        $p_\theta$

$$\text{gap} = \text{KLC}[q_\phi(z|x), p_\theta(z|x)]$$

# Is VAE training still Hill Climbing?

It is not, because  $q(z|x)$  may not be accurate to approximate  $p(z|x)$

# Is VAE training still Hill Climbing?

It is not, because  $q(z|x)$  may not be accurate to approximate  $p(z|x)$

In VAE training, there is no guarantee that  $\log p(x)$  is monotonically increasing

# Is VAE training still Hill Climbing?

It is not, because  $q(z|x)$  may not be accurate to approximate  $p(z|x)$

In VAE training, there is no guarantee that  $\log p(x)$  is monotonically increasing

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

According to EM,  $\phi$  should be optimized to convergence to have a good approximation for  $p(z|x)$  before conducting the M-step, but VAE does not

# Is VAE training still Hill Climbing?

It is not, because  $q(z|x)$  may not be accurate to approximate  $p(z|x)$

In VAE training, there is no guarantee that  $\log p(x)$  is monotonically increasing

It just works in many cases

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

According to EM,  $\phi$  should be optimized to convergence to have a good approximation for  $p(z|x)$  before conducting the M-step, but VAE does not

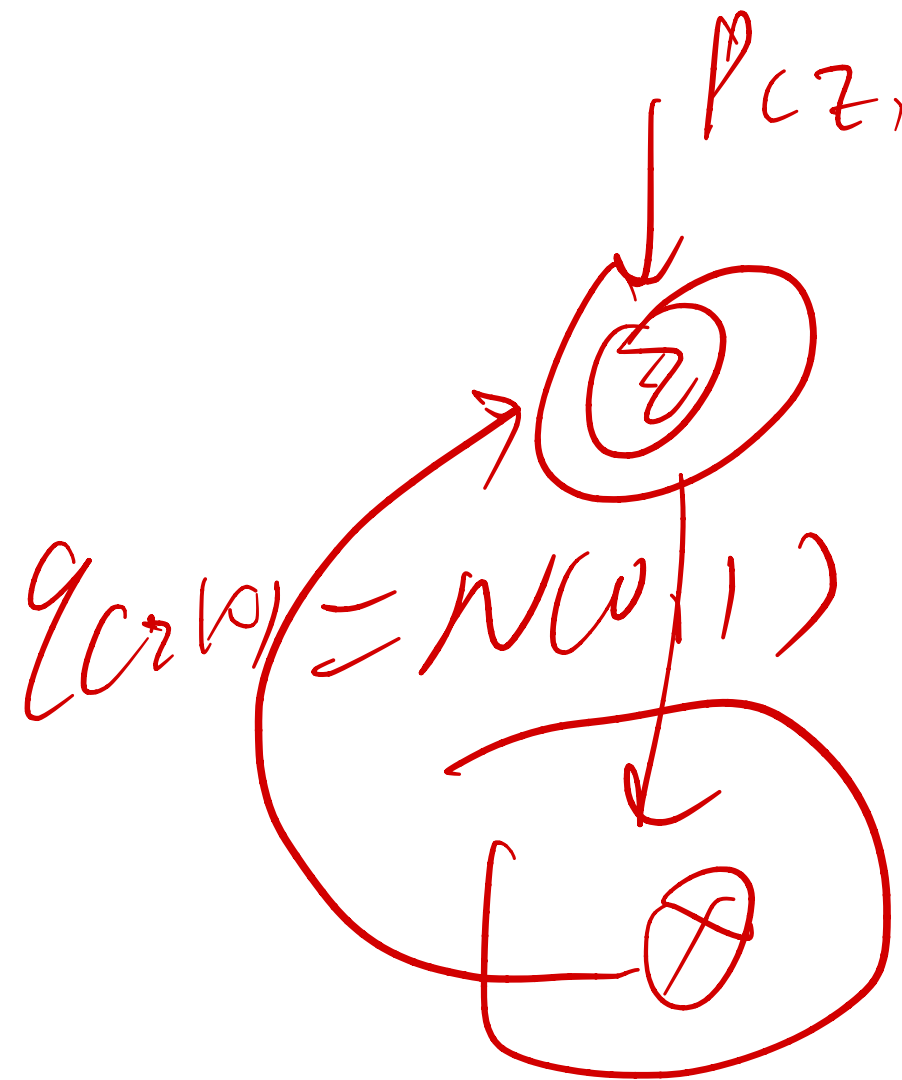
# The Posterior Collapse Issue

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

# The Posterior Collapse Issue

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

In practice, it is often found that after training,  $q_{\phi}(z|x) = p(z)$  and  $z$  and  $x$  becomes independent (especially in applications of NLP)



# The Posterior Collapse Issue

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

In practice, it is often found that after training,  $q_{\phi}(z|x) = p(z)$  and  $z$  and  $x$  becomes independent (especially in applications of NLP)

*Z does not affect x, the model degenerates to a generative model without latent variables*

# The Posterior Collapse Issue

$p(x)$

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

In practice, it is often found that after training,  $q_{\phi}(z|x) = p(z)$  and  $z$  and  $x$  becomes independent (especially in applications of NLP)

$z$  does not affect  $x$ , the model degenerates to a generative model without latent variables

Researchers commonly blame that the KL regularizer is too strong for this and use a weight  $0 < \lambda < 1$  to control it:

# The Posterior Collapse Issue

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

In practice, it is often found that after training,  $q_{\phi}(z|x) = p(z)$  and  $z$  and  $x$  becomes independent (especially in applications of NLP)

*Z does not affect x, the model degenerates to a generative model without latent variables*

Researchers commonly blame that the KL regularizer is too strong for this and use a weight  $0 < \lambda < 1$  to control it:

Reconstruction Loss -  $\lambda$  \* KL regularizer

# The Posterior Collapse Issue

ELBO

$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

In practice, it is often found that after training,  $q_{\phi}(z|x) = p(z)$  and  $z$  and  $x$  becomes independent (especially in applications of NLP)

$z$  does not affect  $x$ , the model degenerates to a generative model without latent variables

Researchers commonly blame that the KL regularizer is too strong for this and use a weight  $0 < \lambda < 1$  to control it:

Reconstruction Loss -  $\lambda$  \* KL regularizer

could be  $> \log p(x)$   
 $> \text{ELBO}$

This is not a lower-bound of  $\log p(x)$  anymore and it breaks MLE, but what is wrong with MLE?

# Is VAE training still Hill Climbing?

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

According to EM,  $\phi$  should be optimized to convergence to have a good approximation for  $p(\mathbf{z}|\mathbf{x})$  before conducting the M-step, but VAE does not

Variational EM

$q(\mathbf{z}|\mathbf{x})$   $p(\mathbf{z}|\mathbf{x})$

why lower bound

# Is VAE training still Hill Climbing?

E-Step:

$$\operatorname{argmax}_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{KL Regularizer}}$$

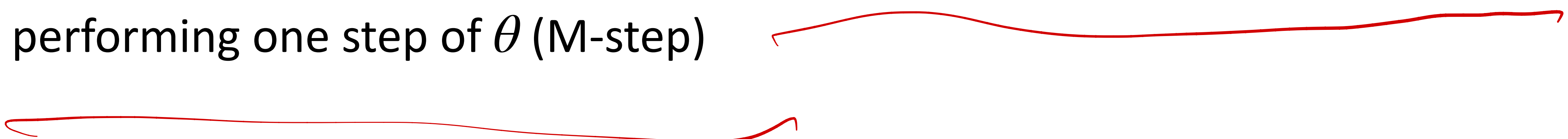
According to EM,  $\phi$  should be optimized to convergence to have a good approximation for  $p(\mathbf{z}|\mathbf{x})$  before conducting the M-step, but VAE does not

Can we make it closer to EM to have good guarantees?

# VAE training that is Closer to EM

# VAE training that is Closer to EM

At every iteration, perform multiple gradient updates of  $\phi$  (E-step) before performing one step of  $\theta$  (M-step)



# VAE training that is Closer to EM

At every iteration, perform multiple gradient updates of  $\phi$  (E-step) before performing one step of  $\theta$  (M-step)

Published as a conference paper at ICLR 2019

---

## LAGGING INFERENCE NETWORKS AND POSTERIOR COLLAPSE IN VARIATIONAL AUTOENCODERS

**Junxian He, Daniel Spokoyny, Graham Neubig**  
Carnegie Mellon University  
{junxianh, dspokoyn, gneubig}@cs.cmu.edu

**Taylor Berg-Kirkpatrick**  
University of California San Diego  
tberg@eng.ucsd.edu