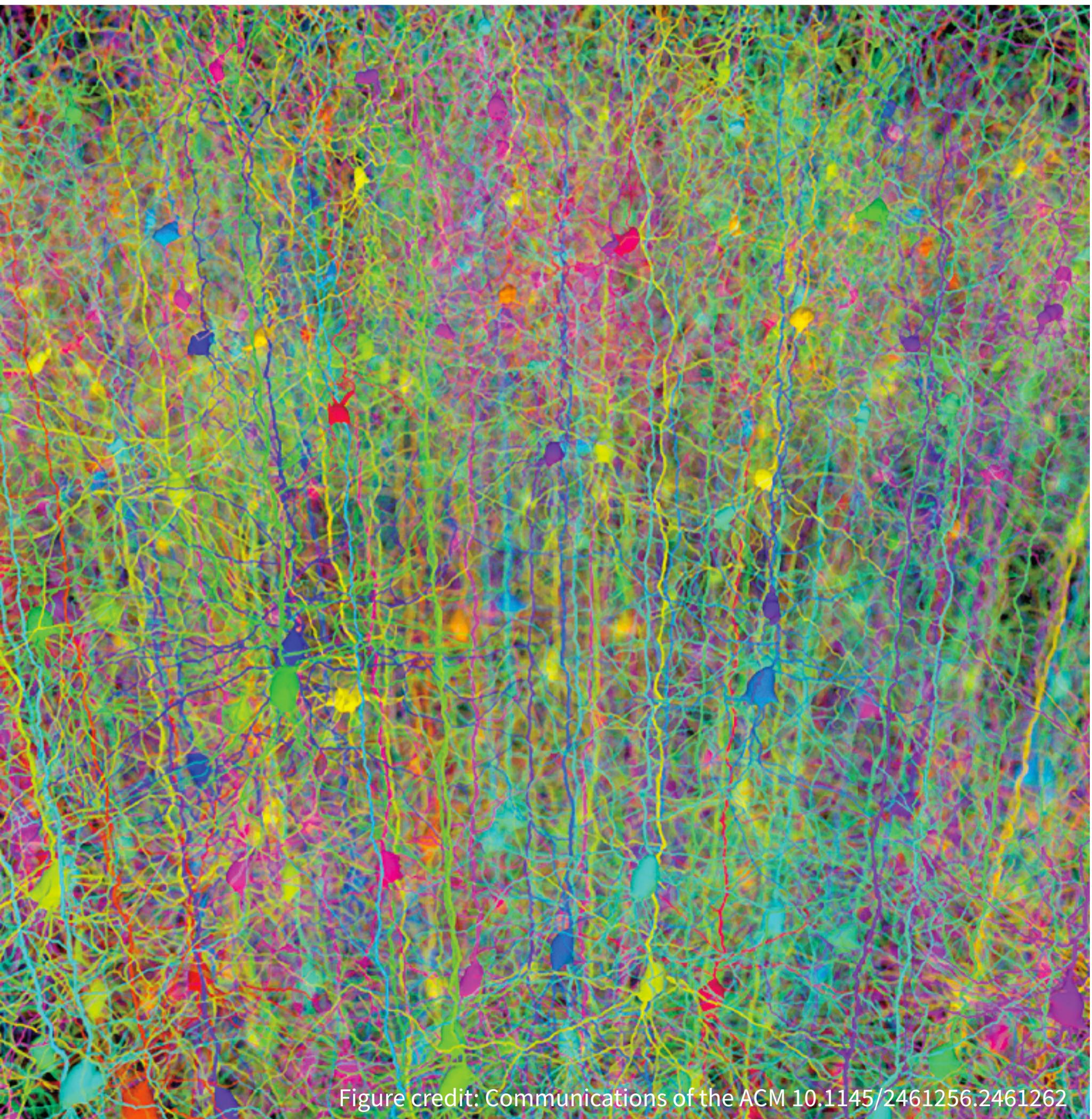


Variational Autoencoders

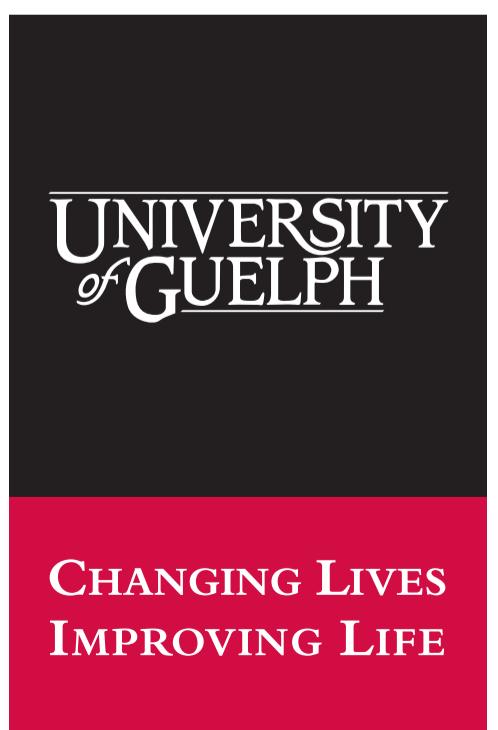


GRAHAM TAYLOR

VECTOR INSTITUTE

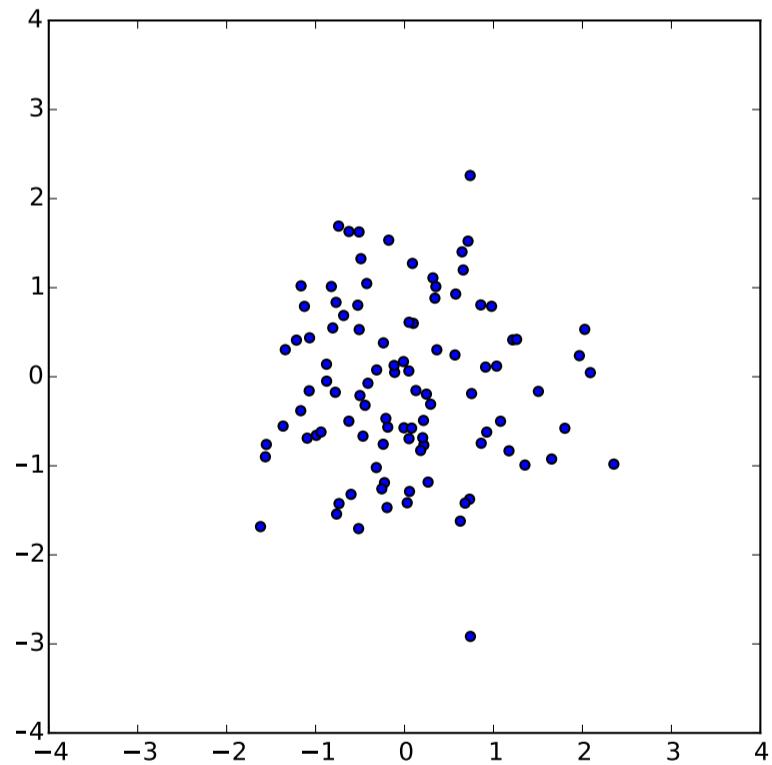
SCHOOL OF ENGINEERING
UNIVERSITY OF GUELPH

CANADIAN INSTITUTE
FOR ADVANCED RESEARCH



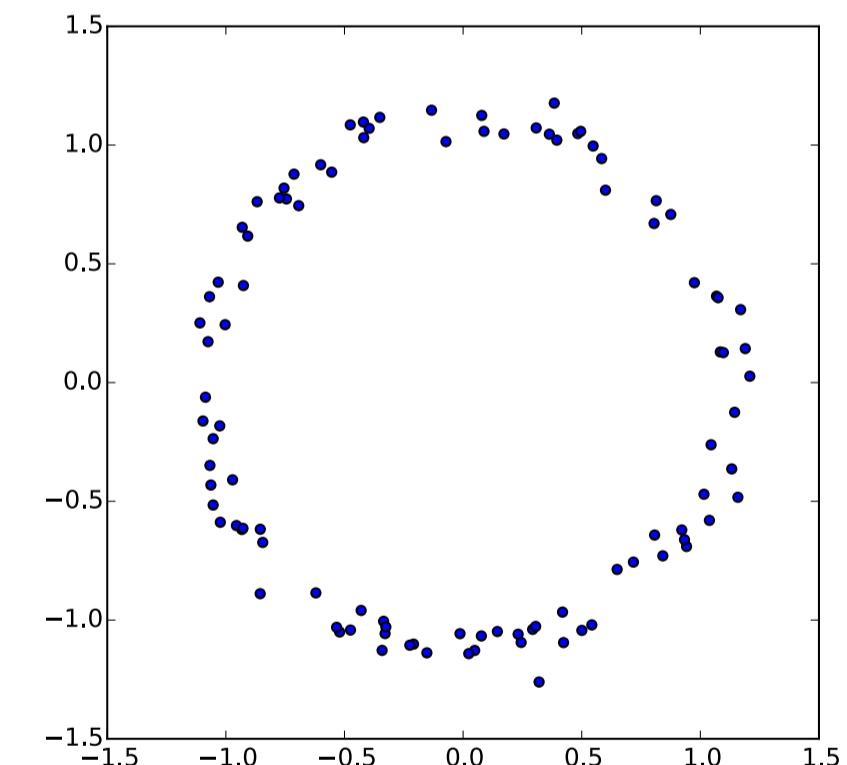
Intuition

Given a random variable with one (simple) distribution, we can create another random variable with a completely different (complex) distribution



z

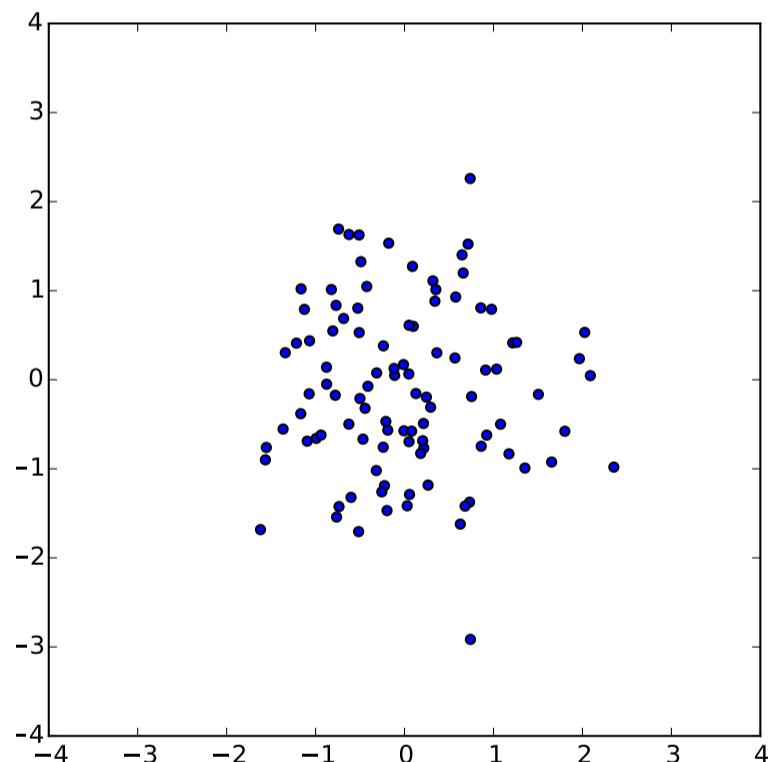
$$g(z) = z/10 + z/\|z\|$$



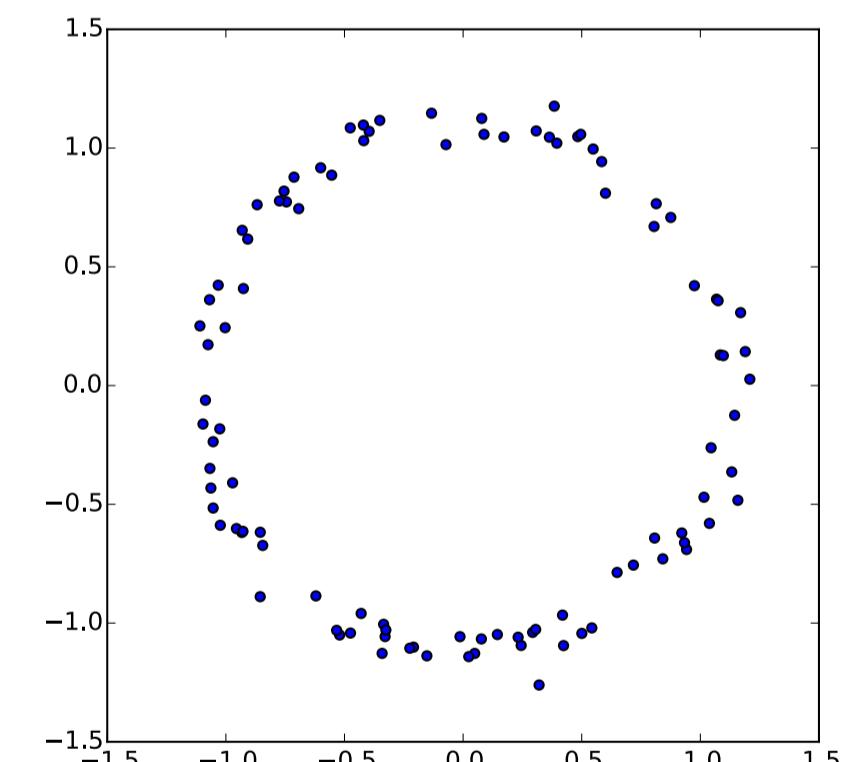
$$x = g(z)$$

Intuition

Given a random variable with one (simple) distribution, we can create another random variable with a completely different (complex) distribution



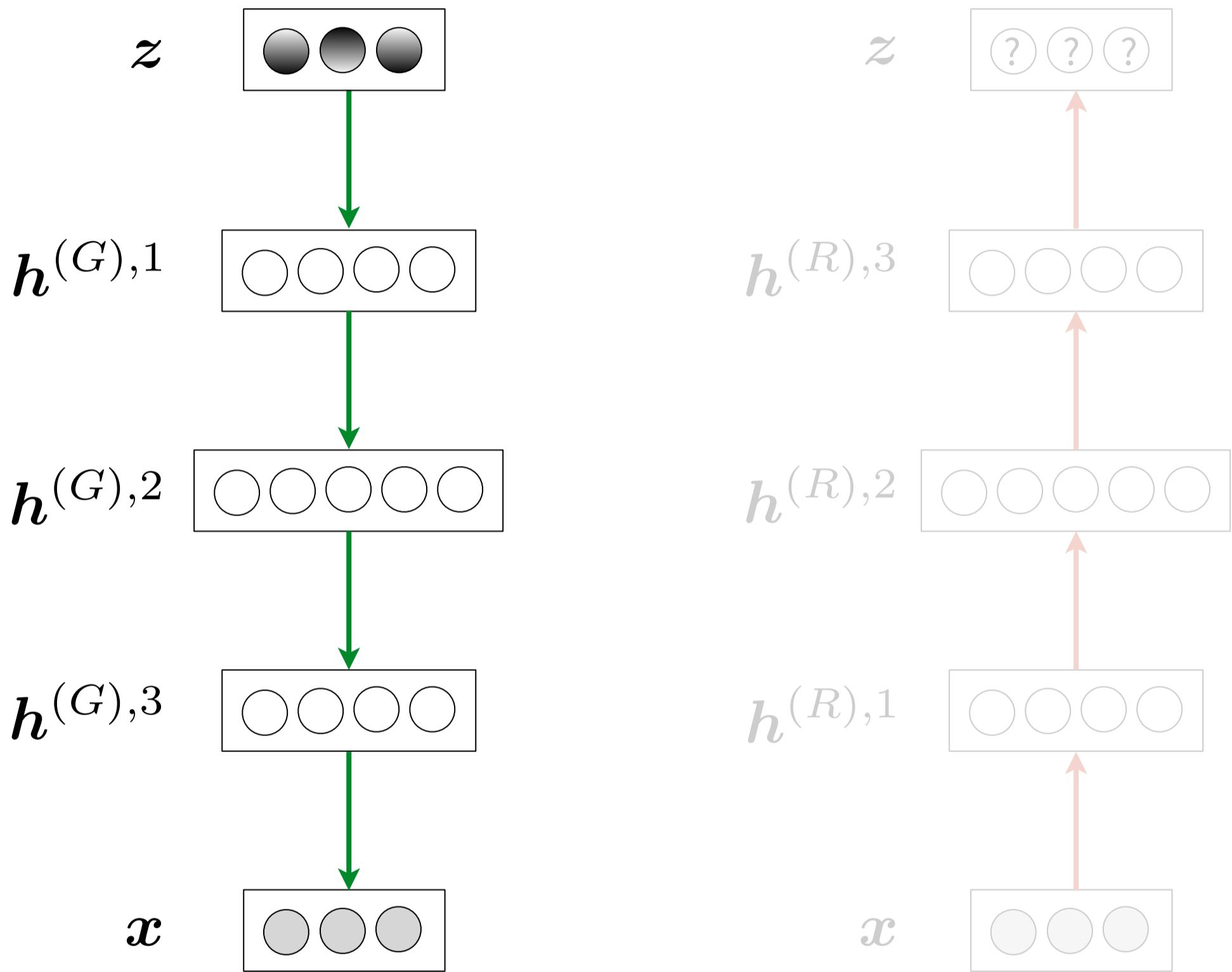
$$f(\mathbf{z}; \boldsymbol{\theta}) \longrightarrow$$



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

$$x \sim \mathcal{N}\left(f(z; \boldsymbol{\theta}), \sigma^2 I\right)$$

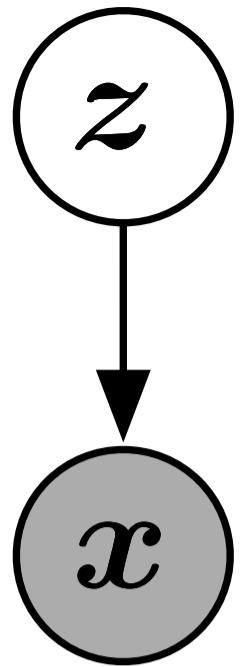
Variational Autoencoder



Generative
Network
(Decoder)

Recognition
Network
(Encoder)

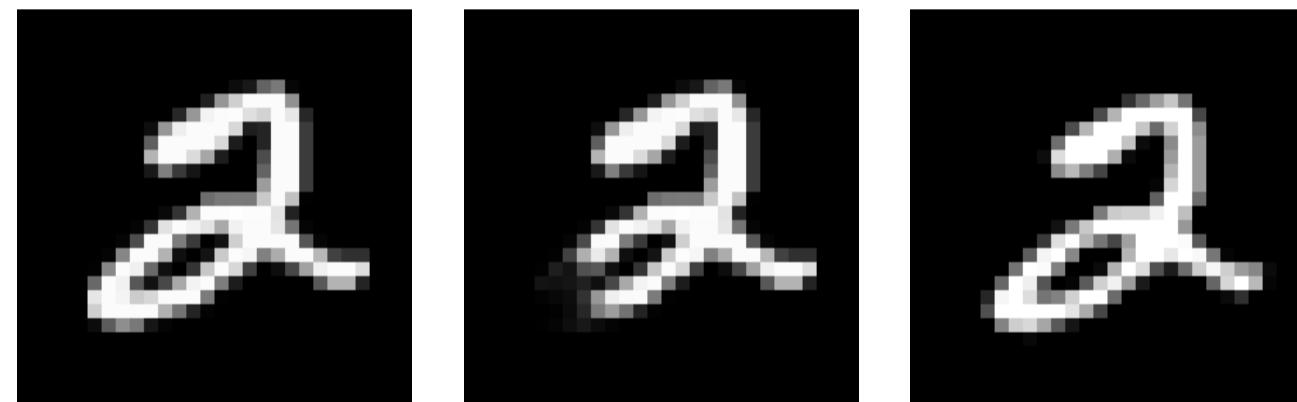
ML Objective



- Want to train to maximize $p(x) = \int p(x|z; \theta)p(z)dz$
- How to define the latent variables z ?
- How to deal with the integral over z ?

Estimating Likelihood?

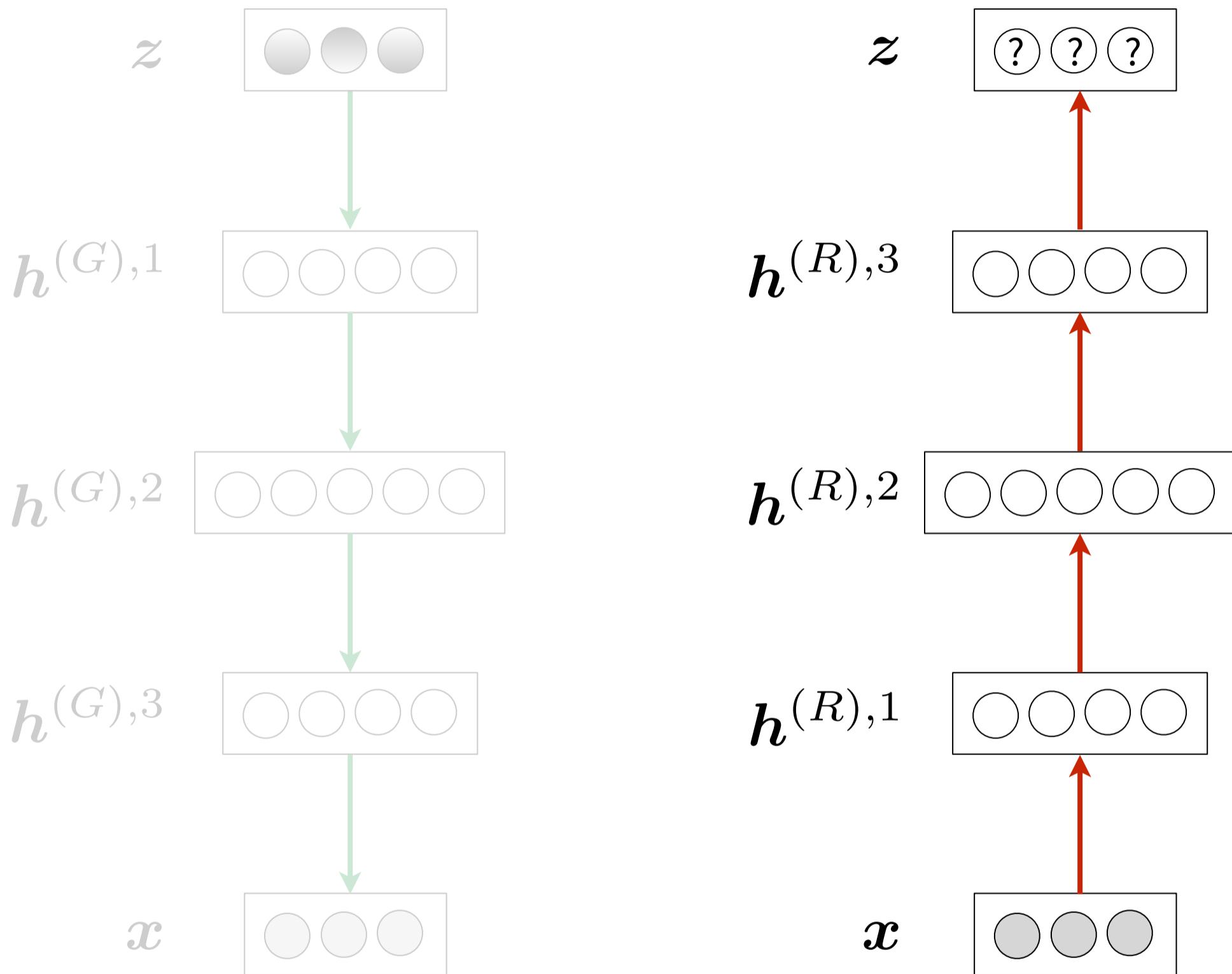
- What about estimating $p(x) \approx \frac{1}{N} \sum_n p(x|z_n)$?



Approximate Posterior

- In practice, for most z , $p(x|z)$ will be **nearly zero** and contribute almost nothing to the estimate for $p(x)$
- The key idea behind the VAE is to attempt to sample values of z that are likely to have produced x and compute $p(x)$ just from those z
- This means we **need a new function** $q(z|x)$ which can take a value of x and give a distribution over z values that are likely to produce x

Variational Autoencoder



Generative
Network
(Decoder)

Recognition
Network
(Encoder)

The Math (1)

The Math (1)

- Start with the Kullback-Leibler (KL) divergence between $p(z|x)$ and $q(z)$ for some arbitrary $q(\cdot)$

$$\mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(z|x)]$$

The Math (1)

- Start with the Kullback-Leibler (KL) divergence between $p(z|x)$ and $q(z)$ for some arbitrary $q(\cdot)$

$$\mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(z|x)]$$

- We can get both $p(x)$ and $p(x|z)$ into this equation by applying Bayes' rule to $p(z|x)$

$$\mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(x|z) - \log p(z)] + \log p(x)$$



Comes out of the expectation,
does not depend on z

The Math (2)

The Math (2)

- Starting point:

$$\mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(x|z) - \log p(z)] + \log p(x)$$

The Math (2)

- Starting point:

$$\mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(x|z) - \log p(z)] + \log p(x)$$

- Negating both sides, rearranging, and spotting another KL-divergence yields

$$\log p(x) - \mathcal{D} [q(z) || p(z|x)] = \mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D} [q(z) || p(z)]$$

The Math (2)

- Starting point:

$$\mathcal{D} [q(z) \parallel p(z|x)] = \mathbb{E}_{z \sim q} [\log q(z) - \log p(x|z) - \log p(z)] + \log p(x)$$

- Negating both sides, rearranging, and spotting another KL-divergence yields

$$\log p(x) - \mathcal{D} [q(z) \parallel p(z|x)] = \mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D} [q(z) \parallel p(z)]$$

- Note that x is fixed, and $q(\cdot)$ can be any distribution. It makes sense to construct a $q(\cdot)$ that depends on x

$$\log p(x) - \mathcal{D} [\underline{q(z|x)} \parallel p(z|x)] = \mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D} [\underline{q(z|x)} \parallel p(z)]$$

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Why is this a sensible objective?

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Why is this a sensible objective?

- The LHS has the **thing we want to maximize** plus an **error term**
 - The error term will be small for high-capacity $q(\cdot)$

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Why is this a sensible objective?

- The LHS has the **thing we want to maximize** plus an **error term**
 - The error term will be small for high-capacity $q(\cdot)$
- The RHS is something we can optimize using SGD

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Reconstruction term Regularizer term
Decode \mathbf{z} into \mathbf{x} Encode \mathbf{x} into \mathbf{z}

Why is this a sensible objective?

- The LHS has the **thing we want to maximize** plus an **error term**
 - The error term will be small for high-capacity $q(\cdot)$
- The RHS is something we can optimize using SGD
- Note that the RHS looks something like an autoencoder

The VAE Objective

$$\log p(\mathbf{x}) - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x})] = \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Reconstruction term Regularizer term
Decode \mathbf{z} into \mathbf{x} Encode \mathbf{x} into \mathbf{z}

Why is this a sensible objective?

- The LHS has the **thing we want to maximize** plus an **error term**
 - The error term will be small for high-capacity $q(\cdot)$
- The RHS is something we can optimize using SGD
- Note that the RHS looks something like an autoencoder

Optimizing VAE Objective

Optimizing VAE Objective

- How can we perform SGD on the RHS of our objective?

$$\mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}[q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})]$$

Optimizing VAE Objective

- How can we perform SGD on the RHS of our objective?

$$\mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D}[q(z|x) || p(z)]$$

- We have to get a bit more formal for $q(z|x)$

Optimizing VAE Objective

- How can we perform SGD on the RHS of our objective?

$$\mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D}[q(z|x) || p(z)]$$

- We have to get a bit more formal for $q(z|x)$
- The usual choice is to let: $q(z|x) = \mathcal{N}(z|\mu(x;\theta), \Sigma(x;\theta))$

so

$$\begin{aligned}\mathcal{D}[\mathcal{N}(\mu_0, \Sigma_0) || \mathcal{N}(\mu_1, \Sigma_1)] &= \\ \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)\end{aligned}$$

Optimizing VAE Objective

- How can we perform SGD on the RHS of our objective?

$$\mathbb{E}_{z \sim q} [\log p(x|z)] - \mathcal{D}[q(z|x) || p(z)]$$

- We have to get a bit more formal for $q(z|x)$
- The usual choice is to let: $q(z|x) = \mathcal{N}(z|\mu(x;\theta), \Sigma(x;\theta))$

so

$$\begin{aligned}\mathcal{D}[\mathcal{N}(\mu_0, \Sigma_0) || \mathcal{N}(\mu_1, \Sigma_1)] &= \\ \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)\end{aligned}$$

- This simplifies to:

$$\mathcal{D}[\mathcal{N}(\mu(x), \Sigma(x)) || \mathcal{N}(\mathbf{0}, I)] = \frac{1}{2} \left(\text{tr}(\Sigma(x)) + \mu(x)^T \mu(x) - k + \log \det(\Sigma(x)) \right)$$

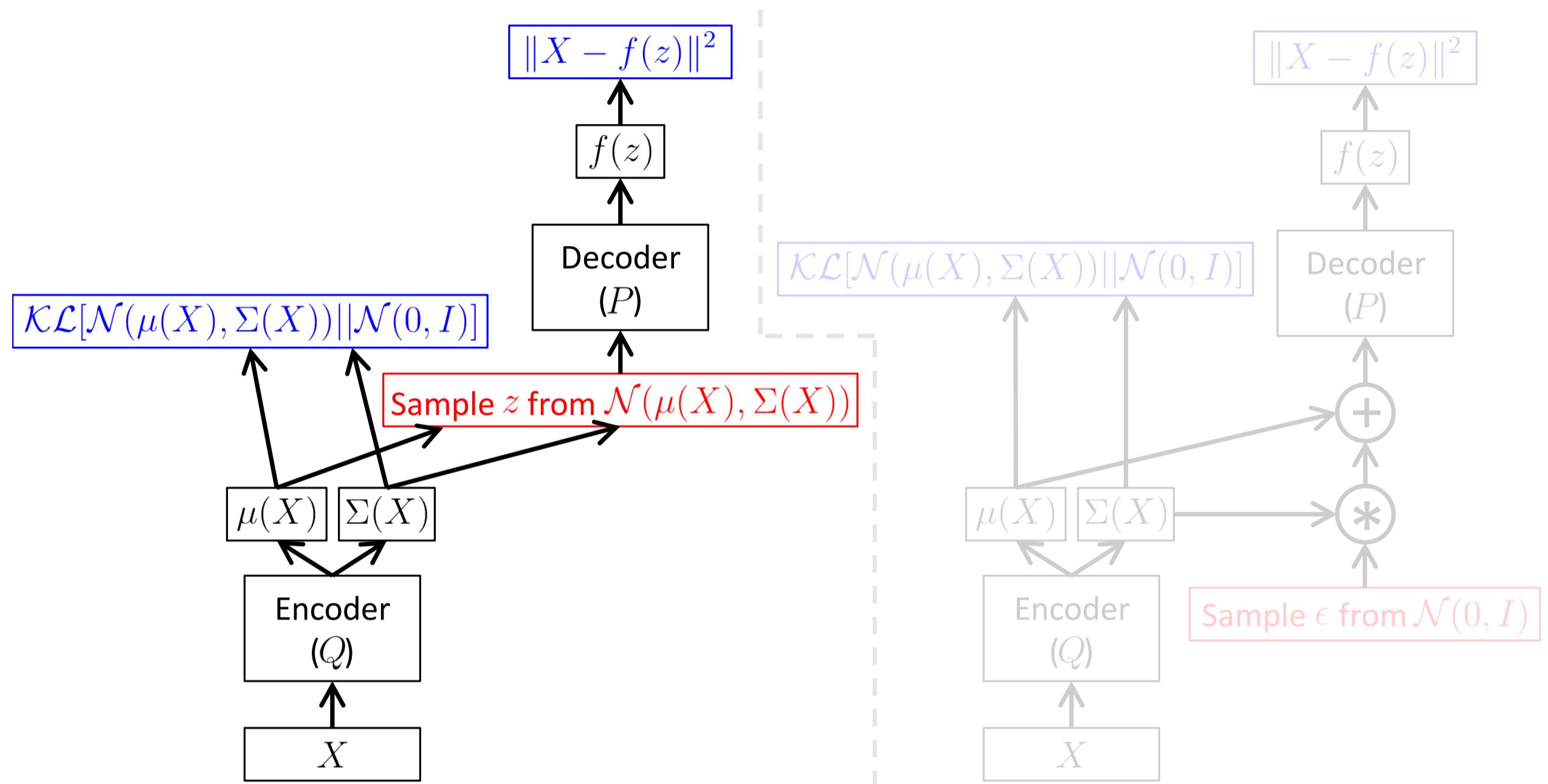
Stochastic Gradient Descent

- Think about what we really want to optimize:

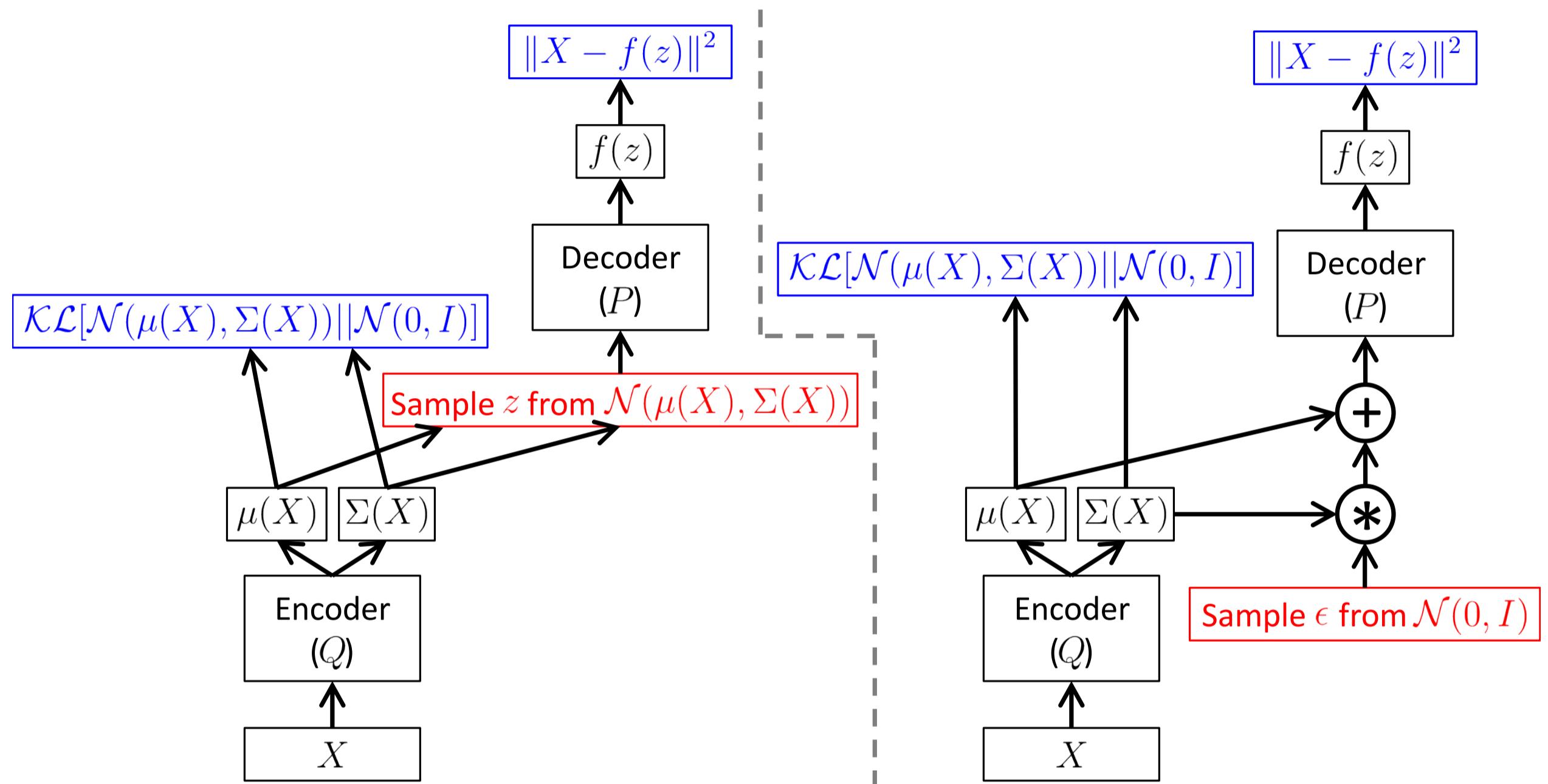
$$\begin{aligned}\mathbb{E}_{\mathbf{x} \sim D} [\log p(\mathbf{x}) - \mathcal{D}[q(z|\mathbf{x}) || p(z|\mathbf{x})]] &= \\ \mathbb{E}_{\mathbf{x} \sim D} [\mathbb{E}_{z \sim q} [\log p(\mathbf{x}|z)] - \mathcal{D}[q(z|\mathbf{x}) || p(z)]]\end{aligned}$$

- If we take the gradient of this equation, the **gradient can be moved inside the expectations**
- Therefore, as usual in SGD, we can sample a single value of \mathbf{x} and a single value of $z \sim q(z|\mathbf{x})$ and compute the gradient of $\log p(\mathbf{x}|z) - \mathcal{D}[q(z|\mathbf{x}) || p(z)]$

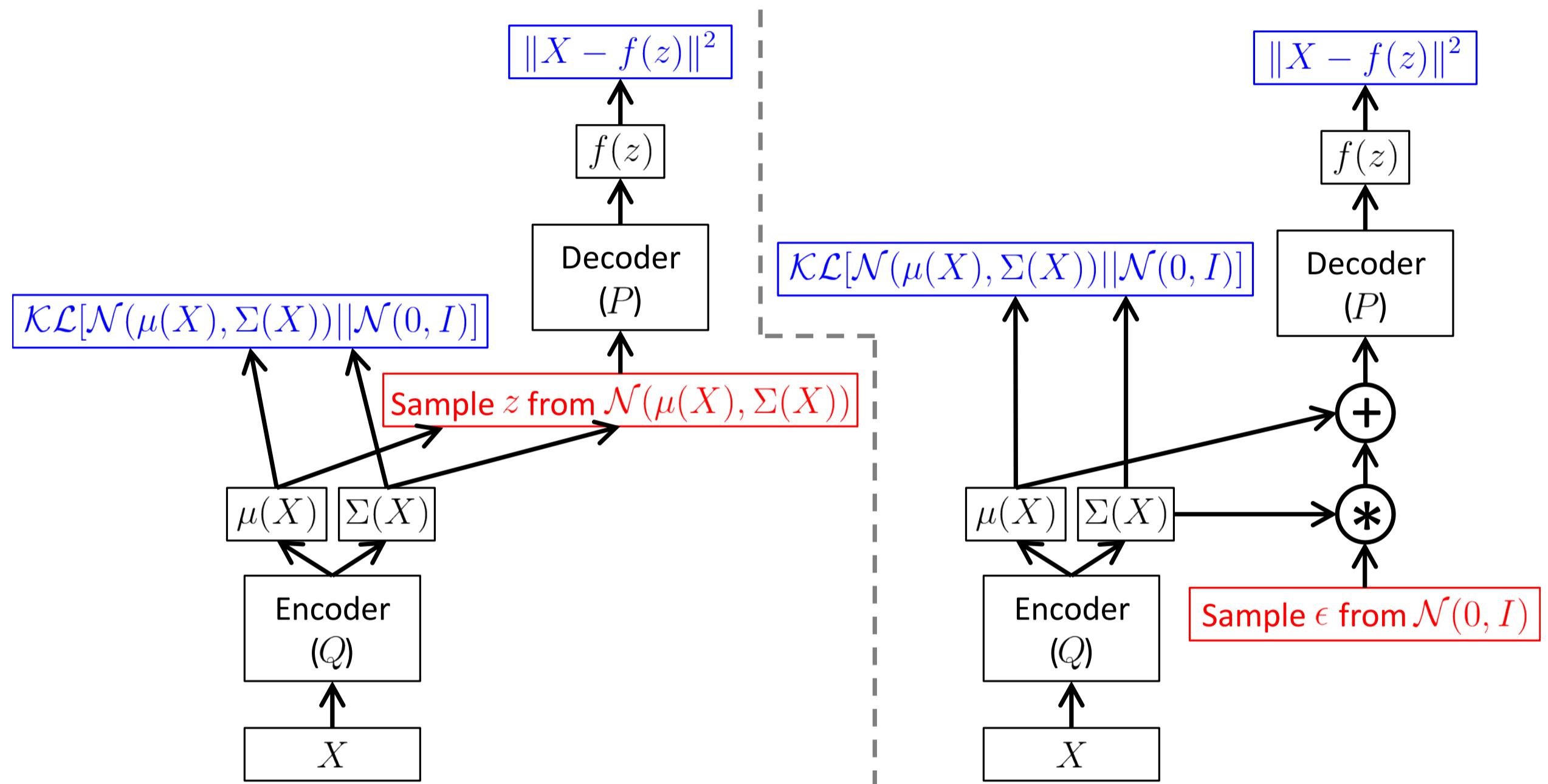
The Reparameterization Trick



The Reparameterization Trick

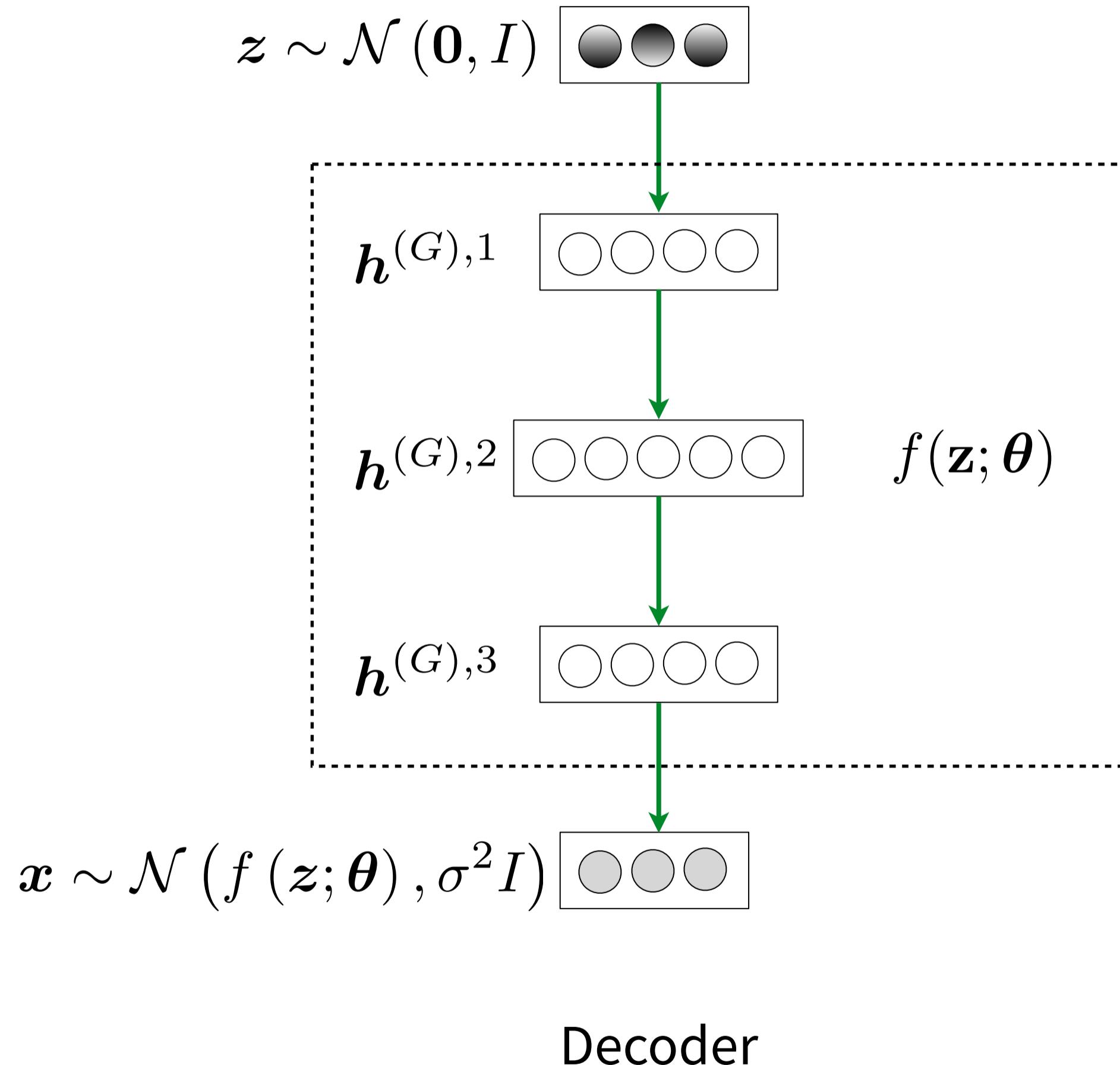


The Reparameterization Trick



$$\mathbb{E}_{\mathbf{x} \sim D} \left[\mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, I)} \left[\log p \left(\mathbf{x} | \mathbf{z} = \boldsymbol{\mu}(\mathbf{x}) + \boldsymbol{\Sigma}^{1/2}(\mathbf{x}) \cdot \boldsymbol{\epsilon} \right) \right] - \mathcal{D} [q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})] \right]$$

Test Time / Sampling



References and Demos

- Carl Doersch's tutorial (these slides mainly follow his presentation):
<https://arxiv.org/abs/1606.05908>
- Practical tips on implementation (and way more links):
<http://bjlkeng.github.io/posts/a-variational-autoencoder-on-the-svhn-dataset/>



Generation from VAE
trained on SVHN
Courtesy of Brian Keng