



---

# Generative Models and Neural Variational Inference

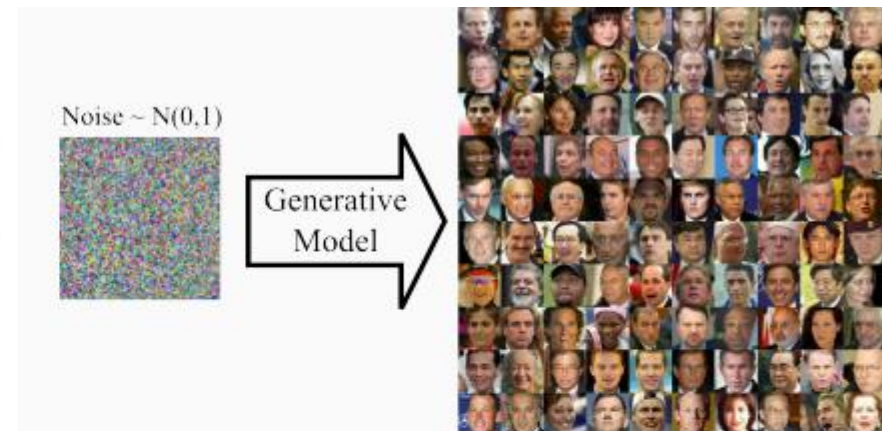
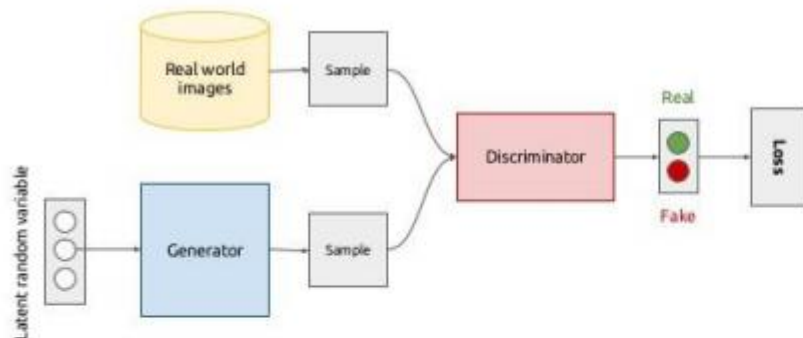
---

Alexander G. Ororbia II  
Neural Networks & Machine Learning  
CSCI-736  
2/16/2023

# Generative Modeling & Sampling

- **Task:** Given a dataset of images  $\{X_1, X_2, \dots\}$  can we learn the distribution of  $X$ ?
- Typically generative models implies modelling  $P(X)$ .
  - **Very limited, given an image the model outputs a probability**
- **More Interested** in models which we can **sample** from.
  - Can generate random examples that follow the distribution of  $P(X)$ .

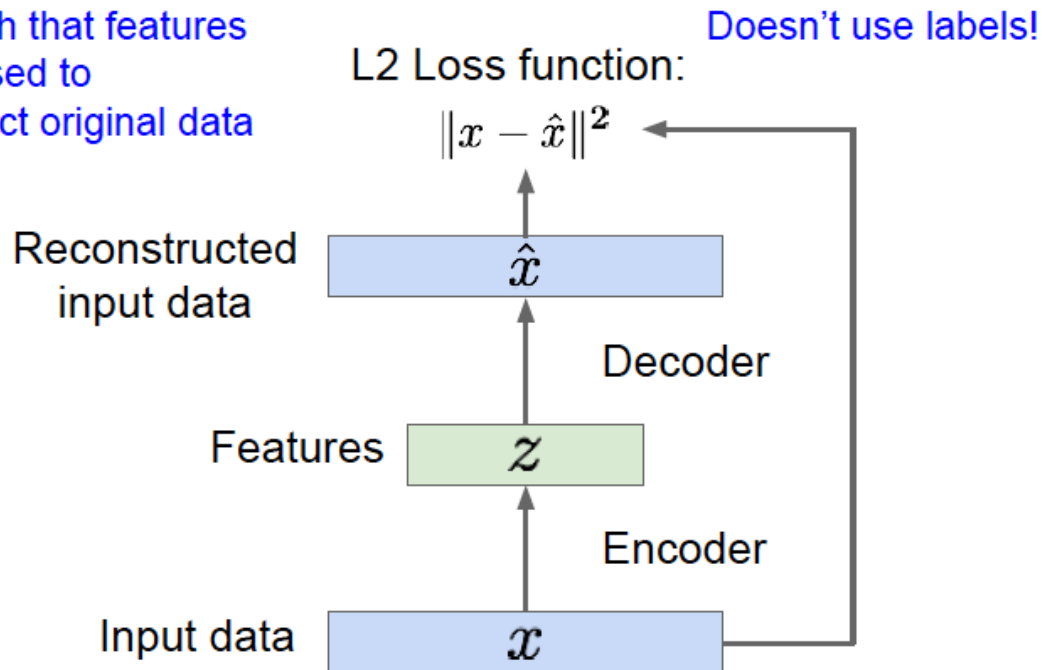
## Generative Adversarial Network:



- **Pro:** Do not have to explicitly specify a form on  $P(X|z)$ ,  $z$  is the latent space.
- **Con:** Given a desired image, difficult to map back to the latent variable.

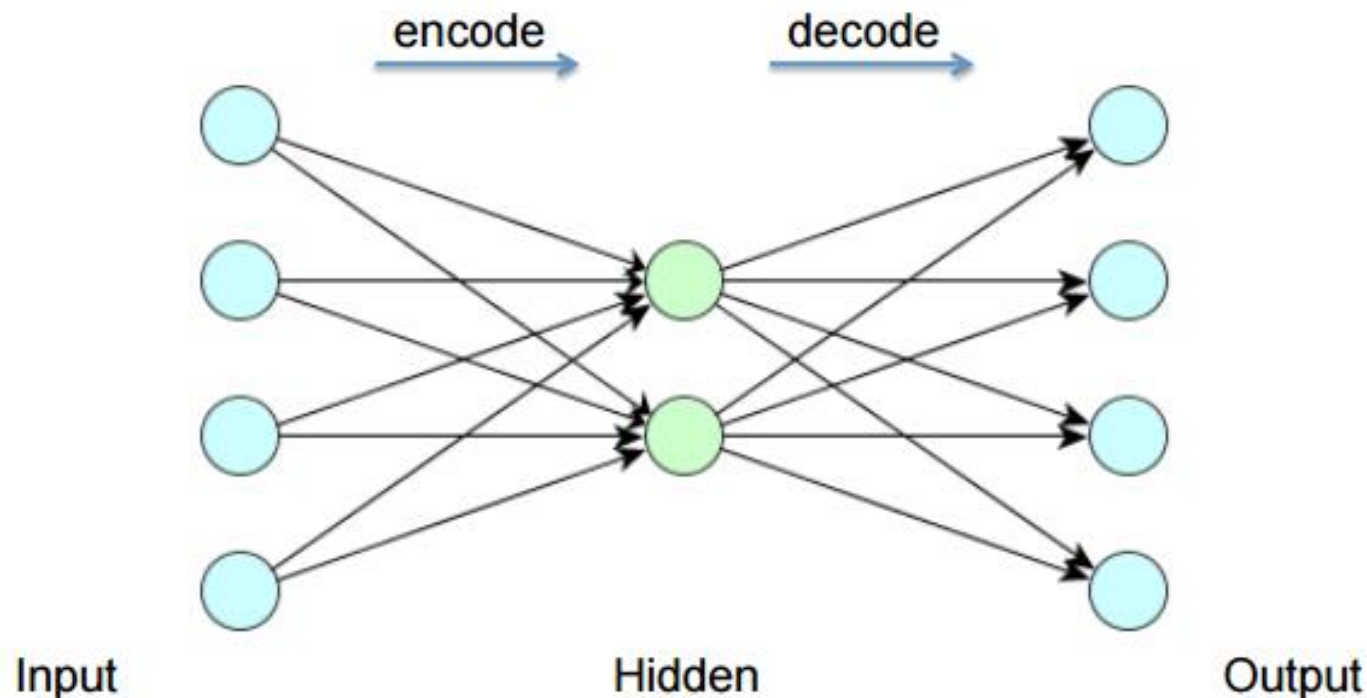
# Concept: Auto-association

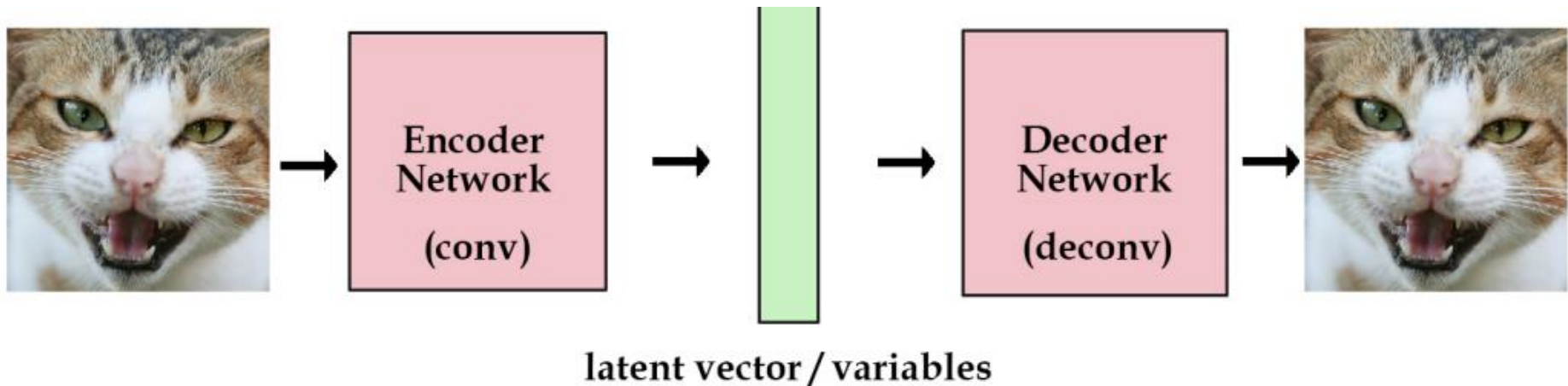
Train such that features can be used to reconstruct original data



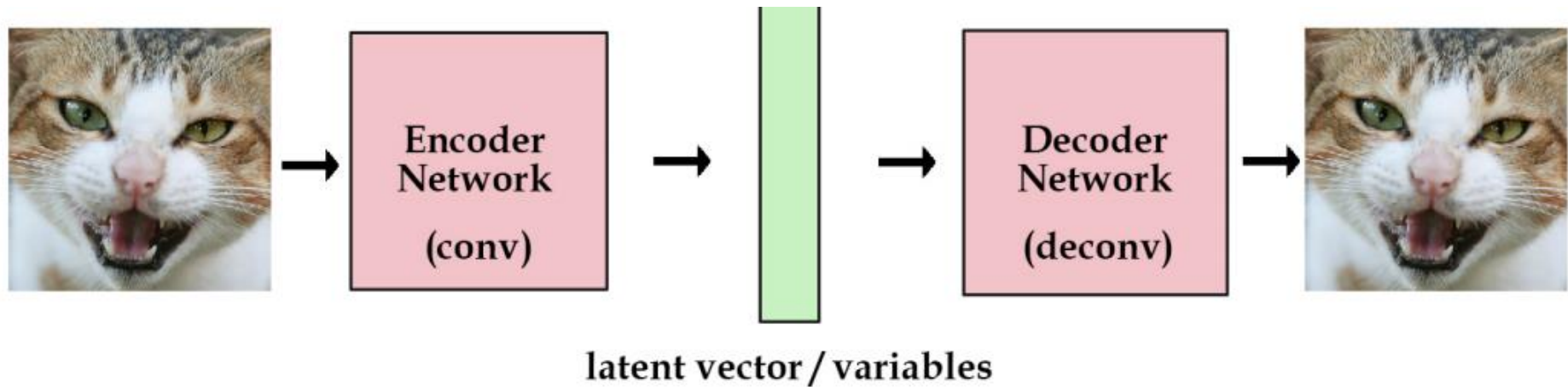
# The Encoder-Decoder Framework

- Auto-association (auto-encoding)
  - Learn a compressed representation of the input, i.e., word2vec
  - Bottleneck layer = meaningful latent space
- Can de-couple encoder & decoder
  - Each can be complex, different functions





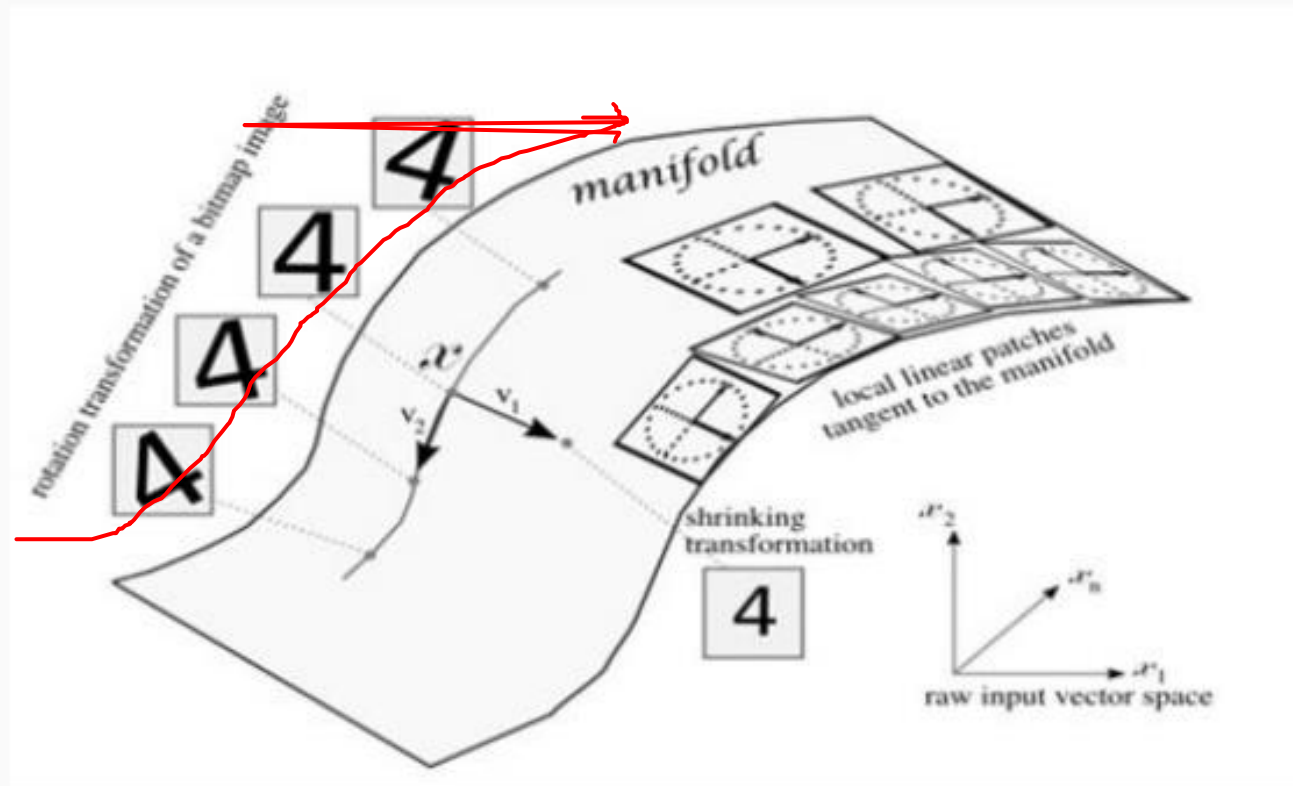
- ▶ Attempt to learn identity function
- ▶ Constrained in some way (e.g., small latent vector representation)
- ▶ Can generate new images by giving different latent vectors to trained network
- ▶ Variational: use probabilistic latent encoding



- ▶ Attempt to learn identity function
- ▶ Constrained in some way (e.g., small latent vector representation)
- ▶ Can generate new images by giving different latent vectors to trained network
- ▶ Variational: use probabilistic latent encoding

# The Manifold Hypothesis

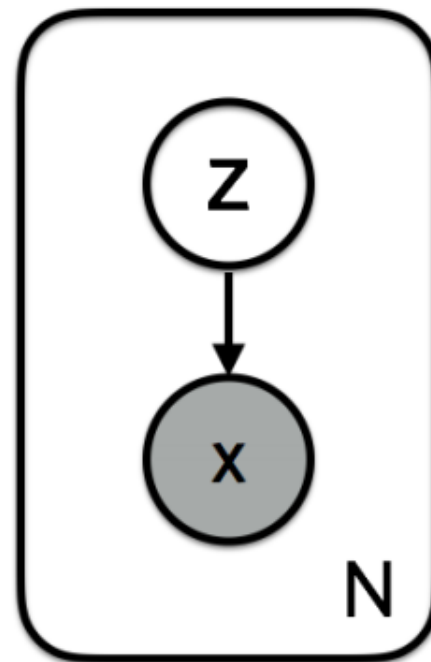
Natural data (high dimensional) actually lies in a low dimensional space.



**Nice consequence:** many datasets that you think might require many variables/dimensions to describe can actually be explained with rather few of them (a subset that forms a sort of local coordinate system of the underlying manifold)

# Probabilistic Model Perspective

- ▶ Data  $x$  and latent variables  $z$
- ▶ Joint pdf of the model:  $p(x, z) = p(x|z)p(z)$
- ▶ Decomposes into likelihood:  $p(x|z)$ , and prior:  $p(z)$
- ▶ Generative process:
  - Draw latent variables  $z_i \sim p(z)$
  - Draw datapoint  $x_i \sim p(x|z)$
- ▶ Graphical model:





# Traditional Approaches to Generative Modeling

- **Explicit Modelling of  $P(X|z; \theta)$** , we will drop the  $\theta$  in the notation.
- $z \sim P(z)$ , which we can sample from, such as a Gaussian distribution.

$$P(X) = \int P(X|z; \theta)P(z)dz$$

- Maximum Likelihood --- **Find  $\theta$  to maximize  $P(X)$** , where  $X$  is the data.
- Approximate with samples of  $z$

$$P(X) \approx \frac{1}{n} \sum_{i=0}^n P(X|z_i)$$

- Approximate with samples of  $z$

$$P(X) \approx \frac{1}{n} \sum_{i=0}^n P(X|z_i)$$

- **Need a lot of samples of  $z$  and most of the  $P(X|z) \approx 0$ .**
- Not practical computationally.
- **Question: Is it possible to know which  $z$  will generate  $P(X|z) \gg 0$ ?**
  - Learn a distribution  $Q(z)$ , where  $z \sim Q(z)$  generates  $P(X|z) \gg 0$ .

# Towards Variational Inference

Assume we can learn a distribution  $Q(z)$ , where  $z \sim Q(z)$  generates  $P(X|z) \gg 0$

- We want  $P(X) = E_{z \sim P(z)} P(X|z)$ , but not practical.  $P(X) \approx \frac{1}{n} \sum_{i=0}^n P(X|z_i)$
- We **can** compute  $E_{z \sim Q(z)} P(X|z)$ , more practical.
- **Question:** How does  $E_{z \sim Q(z)} P(X|z)$  and  $P(X)$  relate?
  - We take advantage of the following relationship:

$$\log P(X) - \mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D}[Q(z) \| P(z)]$$

- Definition of KL divergence:

$$\mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)]$$

- Apply **Bayes Rule on  $P(z|X)$**  and substitute into the equation above.
  - $P(z|X) = P(X|z) P(z) / P(X)$
  - $\log(P(z|X)) = \log P(X|z) + \log P(z) - \log P(X)$
  - $P(X)$  does not depend on  $z$ , can take it outside of  $E_{z \sim Q}$

$$\mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X)$$

# Designing a Recognition Model $Q(z)$

Why is this important?

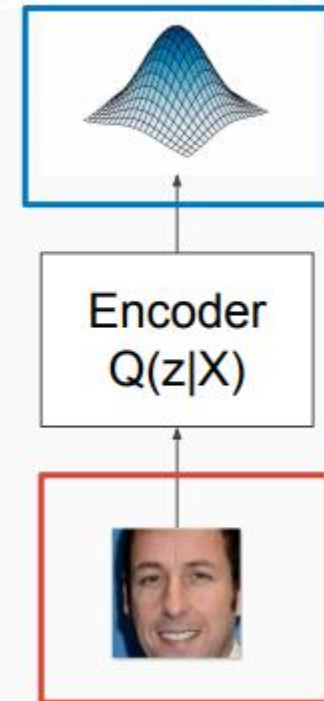
$$\log P(X) - \mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q}[\log P(X|z)] - \mathcal{D}[Q(z) \| P(z)]$$

- Recall we want to **maximize  $P(X)$**  with respect to  $\theta$ , **which we cannot do**.
- **KL divergence is always  $> 0$** .
- **$\log P(X) > \log P(X) - \mathcal{D}[Q(z) \| P(z|X)]$** .
- Maximize the **lower bound** instead.
- **Question: How do we get  $Q(z)$ ?**

$$\log P(X) - \mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q}[\log P(X|z)] - \mathcal{D}[Q(z) \| P(z)]$$

- **$Q(z)$  or  $Q(z|X)$ ?**
- Model  $Q(z|X)$  with a neural network.
- Assume  $Q(z|X)$  to be Gaussian,  $N(\mu, c \cdot I)$ 
  - Neural network **outputs** the **mean  $\mu$** , and **diagonal covariance matrix  $c \cdot I$** .
  - **Input: Image, Output: Distribution**

Let's call  $Q(z|X)$  the **Encoder**.



# Designing a Variational Encoder-Decoder

Convert the lower bound to a loss function:

$$\log P(X) - \mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D}[Q(z) \| P(z)]$$

- Model  $P(X|z)$  with a neural network, let  $f(z)$  be the network output.
- Assume  $P(X|z)$  to be i.i.d. **Gaussian**
  - $X = f(z) + \eta$ , where  $\eta \sim N(0, I)$  \*Think Linear Regression\*
  - **Simplifies to an  $l_2$  loss:**  $\|X - f(z)\|^2$

Let's call  $P(X|z)$  the Decoder.

Convert the lower bound to a loss function:

$$\log P(X) - \mathcal{D}[Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D}[Q(z) \| P(z)]$$

Assume  $P(z) \sim N(0, I)$  then  $\mathcal{D}[Q(z|X) \| P(z)]$  has a closed form solution.

Putting it all together:  $E_{z \sim Q(z|X)} \log P(X|z) \propto \|X - f(z)\|^2$

$$L = \|X - f(z)\|^2 - \lambda \cdot \mathcal{D}[Q(z) \| P(z)]$$

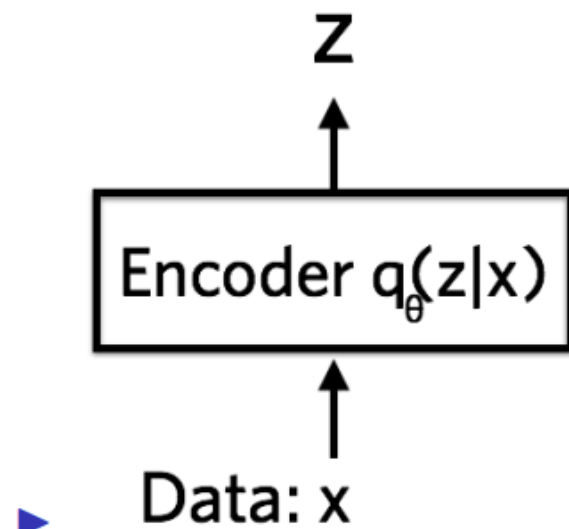
, given a  $(X, z)$  pair.

**Pixel  
difference**

**Regularization**

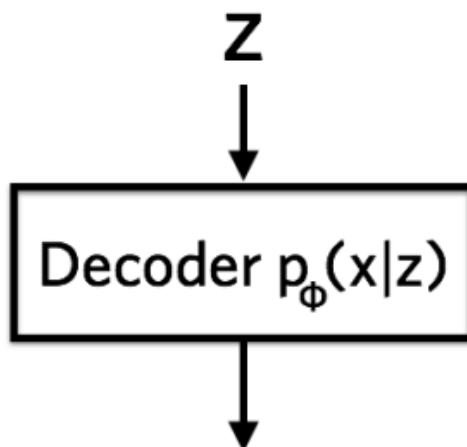
- ▶ Goal: Build a neural network that generates digits from random (Gaussian) noise
- ▶ Define two sub-networks: Encoder and Decoder
- ▶ Define a Loss Function

- ▶ A neural network  $q_{\theta}(z|x)$
- ▶ Input: datapoint  $x$  (e.g.  $28 \times 28$ -pixel digit)
- ▶ Output: encoding  $z$ , drawn from Gaussian density with parameters  $\theta$
- ▶  $|z| \ll |x|$



- ▶ Goal: Build a neural network that generates digits from random (Gaussian) noise
- ▶ Define two sub-networks: Encoder and **Decoder**
- ▶ Define a Loss Function

- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$
- ▶ Input: encoding  $z$ , output from encoder
- ▶ Output: reconstruction  $\tilde{x}$ , drawn from distribution of the data
- ▶ E.g., output parameters for  $28 \times 28$  Bernoulli variables



- ▶ Reconstruction:  $\tilde{x}$

# Neural Variational Inference (NVIL)

- Idea: Teach neural net to approximate the posterior  $p(z/x)$ 
  - $q(z/x)$  with 'variational parameters'  $\phi$
  - One-shot approximate inference
  - Also known as a recognition model
    - Construct estimator of the variational (evidence) lower bound (ELBO)
    - Can optimize jointly w.r.t.  $\phi$  jointly with  $\theta$  -> Stochastic gradient ascent

$D_{\text{KL}}$  KL-Divergence  $\geq 0$  depends on how good  $q(z|x)$  can approximate  $p(z|x)$

## Recall from statistics and information theory:

KL Divergence:

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

Gaussian KL Divergence:

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

$$p \sim N(\mu_1, \sigma_1)$$

$$q \sim N(\mu_2, \sigma_2)$$

# Neural Variational Inference (NVIL)

- Idea: Teach neural net to approximate the posterior  $p(z/x)$

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{z \sim q(z|x)} \log p_{\text{model}}(z, \mathbf{x}) + \mathcal{H}(q(z | \mathbf{x})) \\ &= \mathbb{E}_{z \sim q(z|x)} \log p_{\text{model}}(\mathbf{x} | z) - D_{\text{KL}}(q(z | \mathbf{x}) || p_{\text{model}}(z)) \\ &\leq \log p_{\text{model}}(\mathbf{x}).\end{aligned}$$

Where:

$p_{\text{model}}(\mathbf{x}; g(z)) = p_{\text{model}}(\mathbf{x} | z)$  and  $q(z | \mathbf{x})$  is used to obtain  $z$ .

Recall from statistics and information theory:

KL Divergence:

$$D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

Gaussian KL Divergence:

$$KL(p, q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

$$\begin{aligned}p &\sim N(\mu_1, \sigma_1) \\ q &\sim N(\mu_2, \sigma_2)\end{aligned}$$



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

Interpreting the bound:

- **Approximate posterior distribution  $q(\mathbf{z}|\mathbf{x})$ :** Best match to true posterior  $p(\mathbf{z}|\mathbf{x})$ , one of the unknown inferential quantities of interest to us.
- **Reconstruction cost:** The expected log-likelihood measures how well samples from  $q(\mathbf{z}|\mathbf{x})$  are able to explain the data  $\mathbf{x}$ .
- **Penalty:** Ensures that the explanation of the data  $q(\mathbf{z}|\mathbf{x})$  doesn't deviate too far from your beliefs  $p(\mathbf{z})$ . A mechanism for realising Ockham's razor.

# The Variational Auto-Encoder

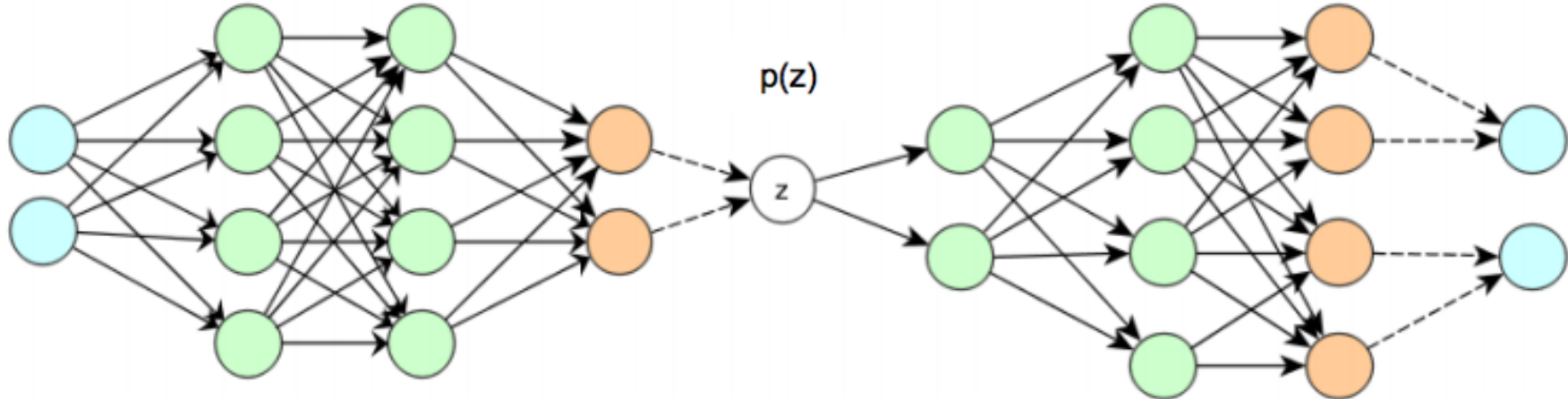
- A feed forward NN + Gaussian

$$q_{\theta}(z | x) = \mathcal{N}(z; \mu_z(x), \sigma_z(x))$$

$$q_{\theta}(z | x)$$

Just a Gaussian, with diagonal covariance.

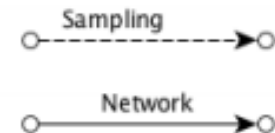
$$p_{\phi}(x | z) \quad x|z \sim \mathcal{N}(\mu_x, \sigma_x^2)$$



- For illustration  $z$  one dimensional x 2D
- Want a complex model of distribution of  $x$  given  $z$

Learning the parameters  $\phi$  and  $\theta$  via backpropagation

Determining the loss function

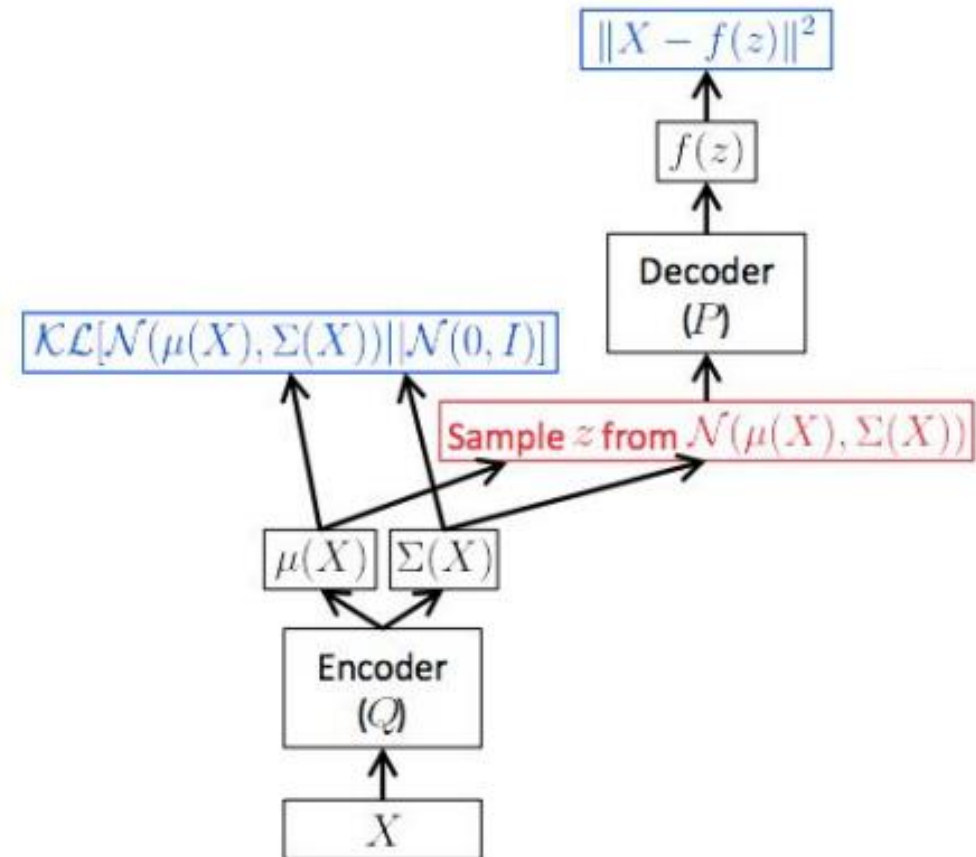


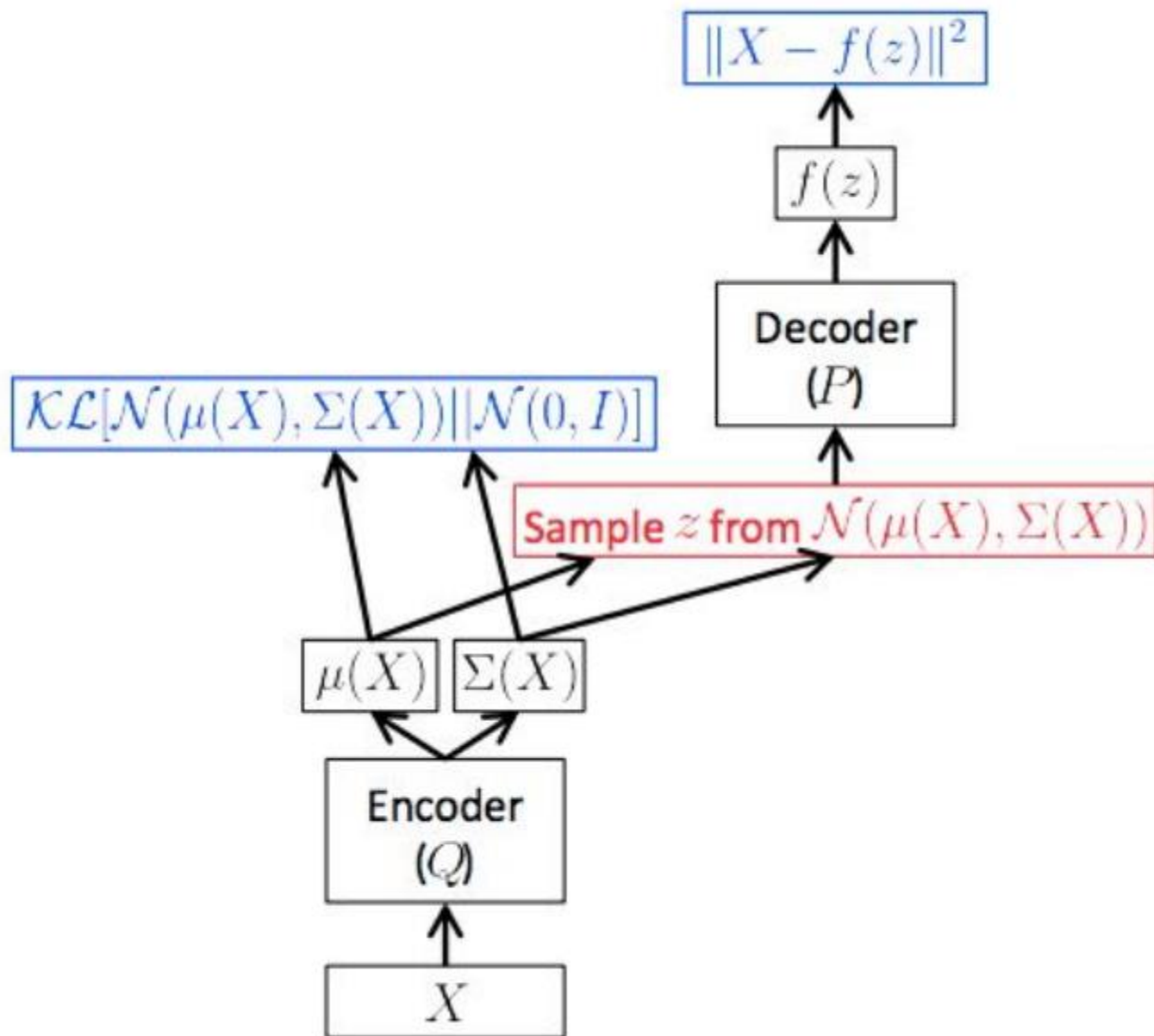
# Issue: Backprop and Sampling

Training the **Decoder** is easy, just standard backpropagation.

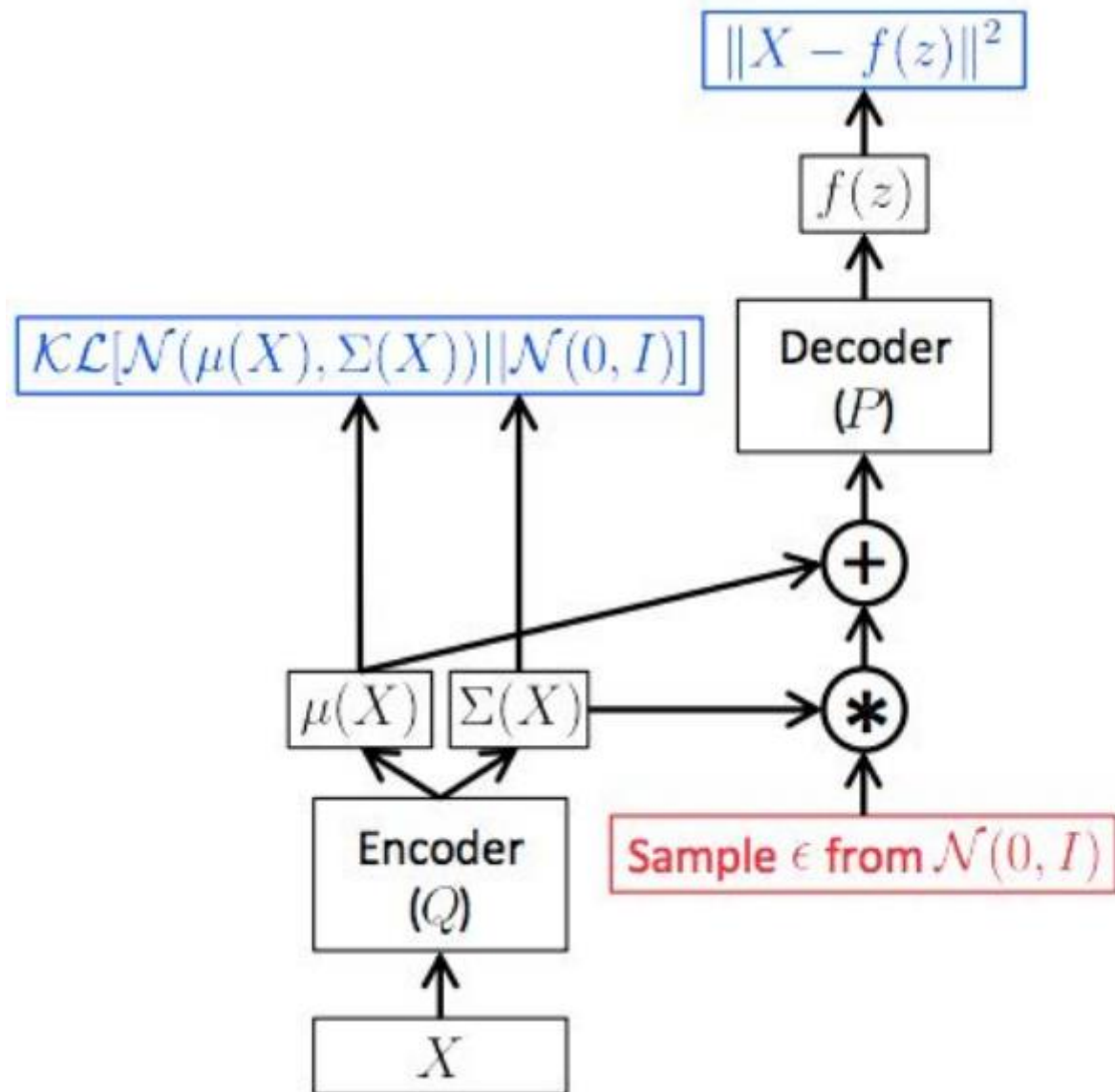
How to train the **Encoder**?

- Not obvious how to apply gradient descent through samples.





# The Reparameterization Trick



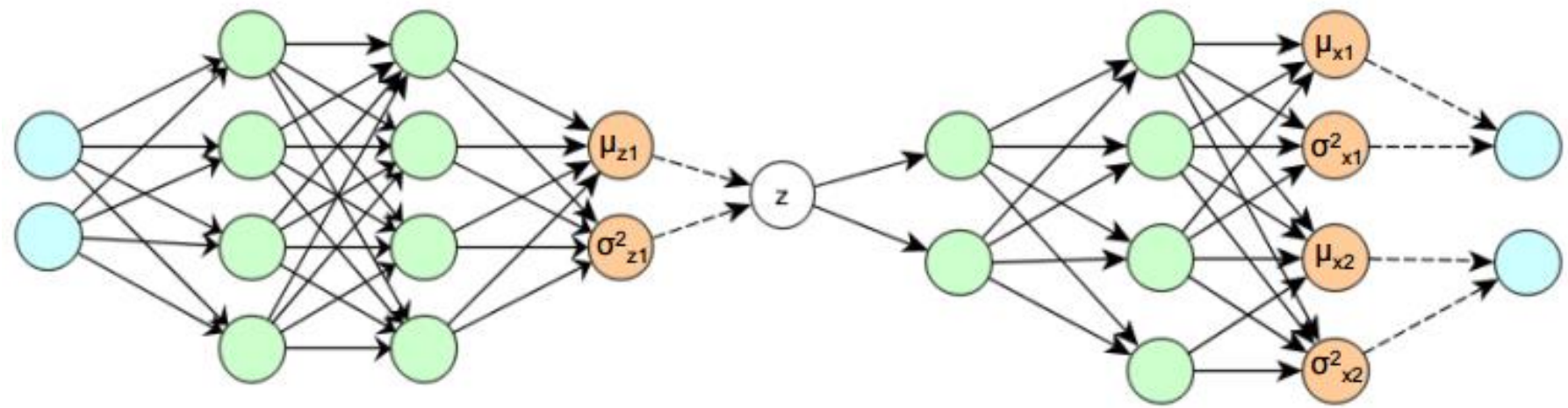
$z \sim \mathcal{N}(\mu, \sigma)$  is equivalent to  $\mu + \sigma \cdot \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 1)$

# The Reparametrization Trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_{\theta}(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?
- ▶ We can reparameterize:  $z = \mu + \sigma \odot \epsilon$
- ▶  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\odot$  is element-wise product
- ▶ Can take derivatives of (functions of)  $z$  w.r.t.  $\mu$  and  $\sigma$
- ▶ Output of  $q_{\theta}(z|x)$  is vector of  $\mu$ 's and vector of  $\sigma$ 's

# Putting It All Together!

Prior  $p(z) \sim N(0,1)$  and  $p, q$  Gaussian, extension to  $\dim(z) > 1$  trivial



## Cost: Regularisation

$$-D_{KL}(q(z|x^{(i)})||p(z)) = \frac{1}{2} \sum_{j=1}^J \left( 1 + \log(\sigma_{z_j}^{(i)2}) - \mu_{z_j}^{(i)2} - \sigma_{z_j}^{(i)2} \right)$$

We use mini batch gradient decent to optimize the cost function over all  $x^{(i)}$  in the mini batch

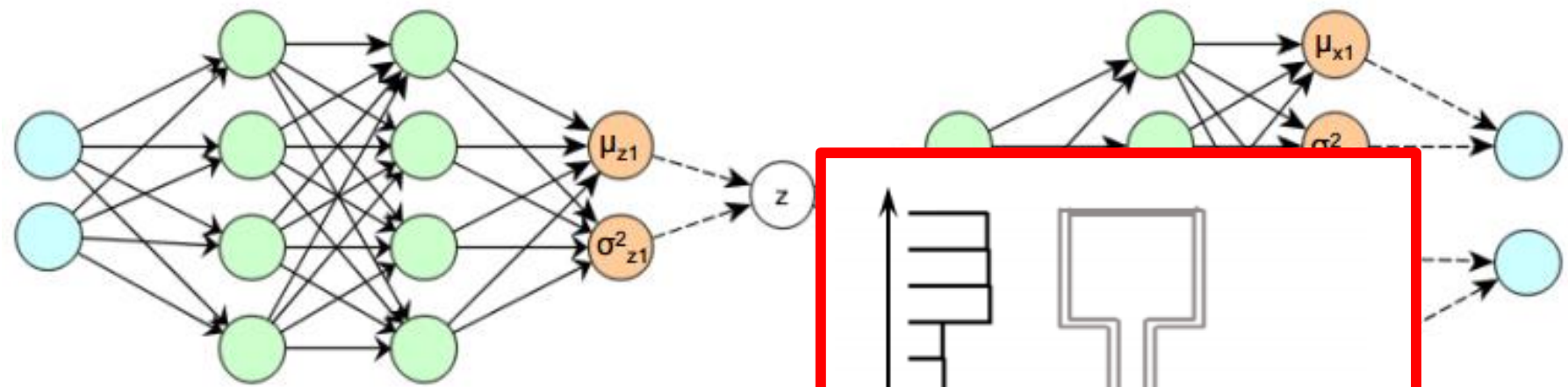
## Cost: Reproduction

$$-\log(p(x^{(i)}|z^{(i)})) = \sum_{j=1}^D \frac{1}{2} \log(\sigma_{x_j}^2) + \frac{(x_j^{(i)} - \mu_{x_j})^2}{2\sigma_{x_j}^2}$$

Least Square for constant variance

# Putting It All Together!

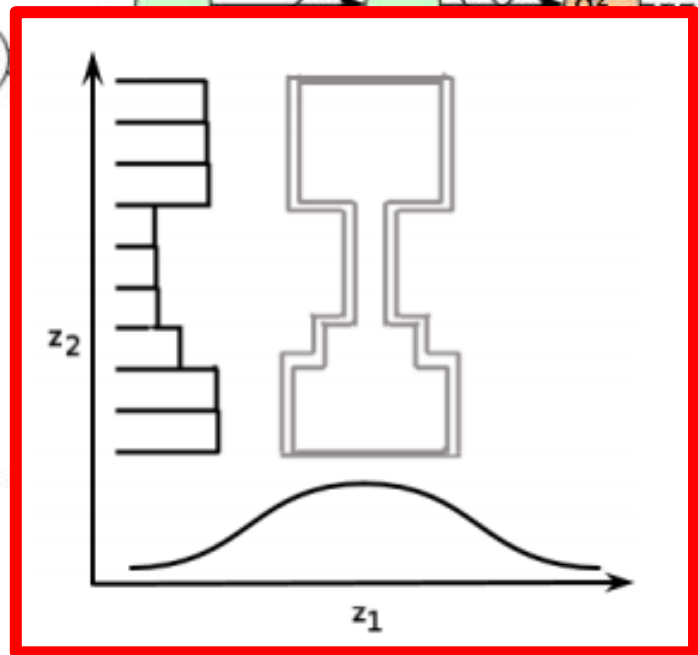
Prior  $p(z) \sim N(0,1)$  and  $p, q$  Gaussian, extension to  $\dim(z) > 1$  trivial



Cost: Regularisation

$$-D_{KL}(q(z|x^{(i)})||p(z)) = \frac{1}{2} \sum_{j=1}^J \left( 1 + \log(\sigma_{z_j}^{(i)2}) - \mu_{z_j}^{(i)2} \right)$$

**Why not capture modes (in text) with piecewise linear variables? (Serban & Ororbia, 2016)**



Least Square for constant variance

$$\sum_{j=1}^J \frac{(x_j^{(i)} - \mu_{x_j})^2}{2\sigma_{x_j}^2}$$

Gradient cost the mini



# VAE Training Algorithm

Given a dataset of examples  $\mathbf{X} = \{X_1, X_2, \dots\}$

Initialize parameters for Encoder and Decoder

**Repeat till convergence:**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  examples from  $\mathbf{X}$

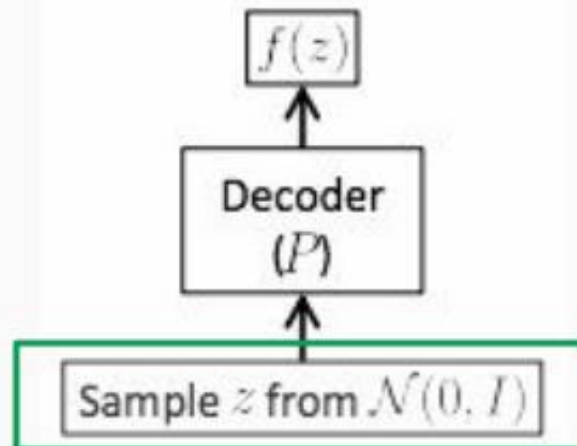
$\varepsilon \leftarrow$  Sample  $M$  noise vectors from  $N(0, I)$

Compute  $L(\mathbf{X}^M, \varepsilon, \theta)$  (i.e. run a forward pass in the neural network)

Gradient descent on  $L$  to updated Encoder and Decoder.

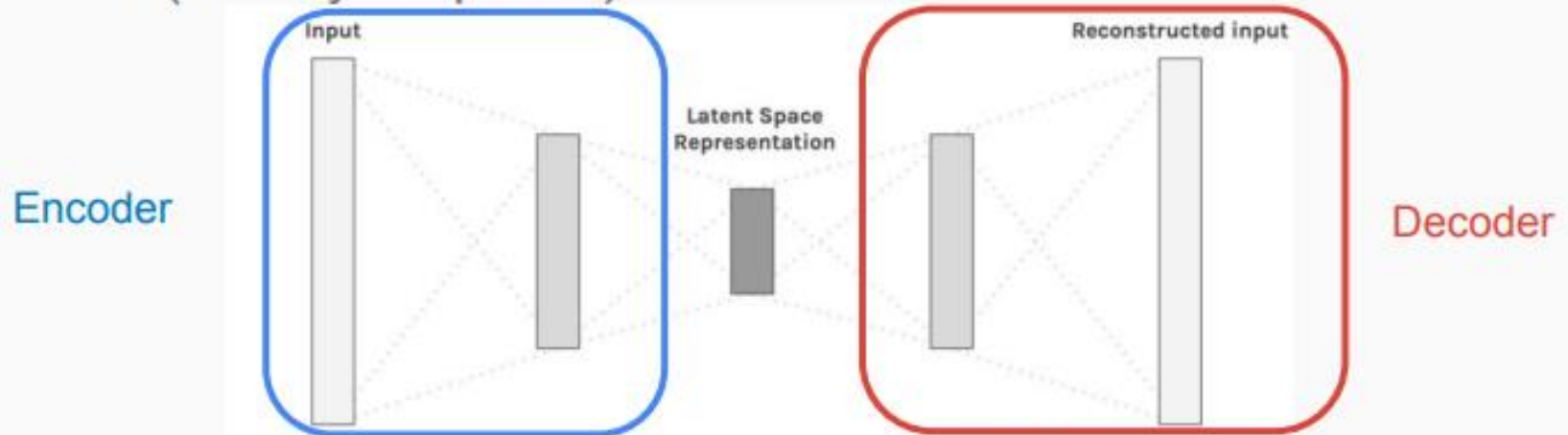
# VAE Evaluation

- At test-time, we want to evaluate the **performance of VAE to generate a new sample.**
- Remove the Encoder, as **no test-image** for generation task.
- Sample  $z \sim \mathcal{N}(0, I)$  and pass it through the Decoder.
- **No good quantitative metric, relies on visual inspection.**

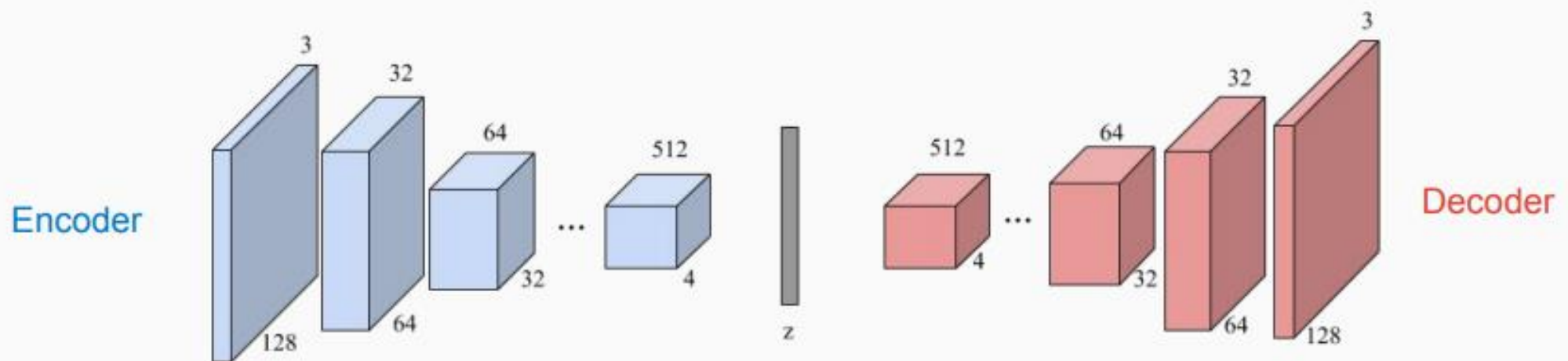


# VAE Architectures

Fully Connected (Initially Proposed)

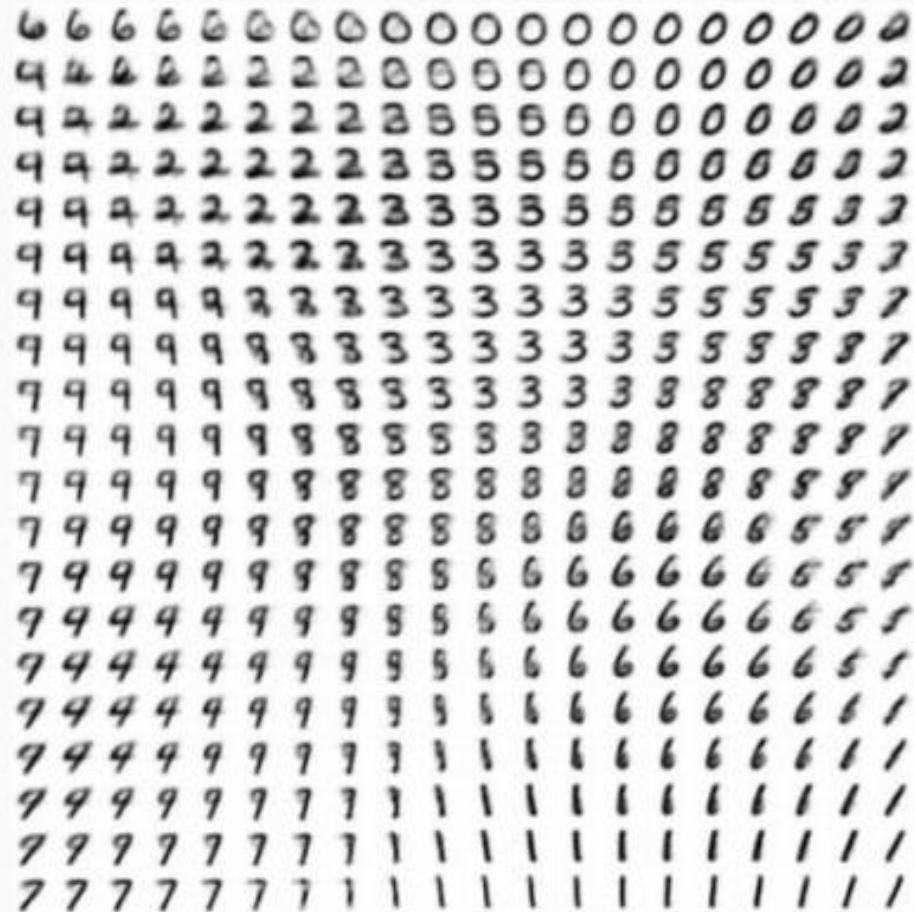


Common Architecture (convolutional) similar to DCGAN.



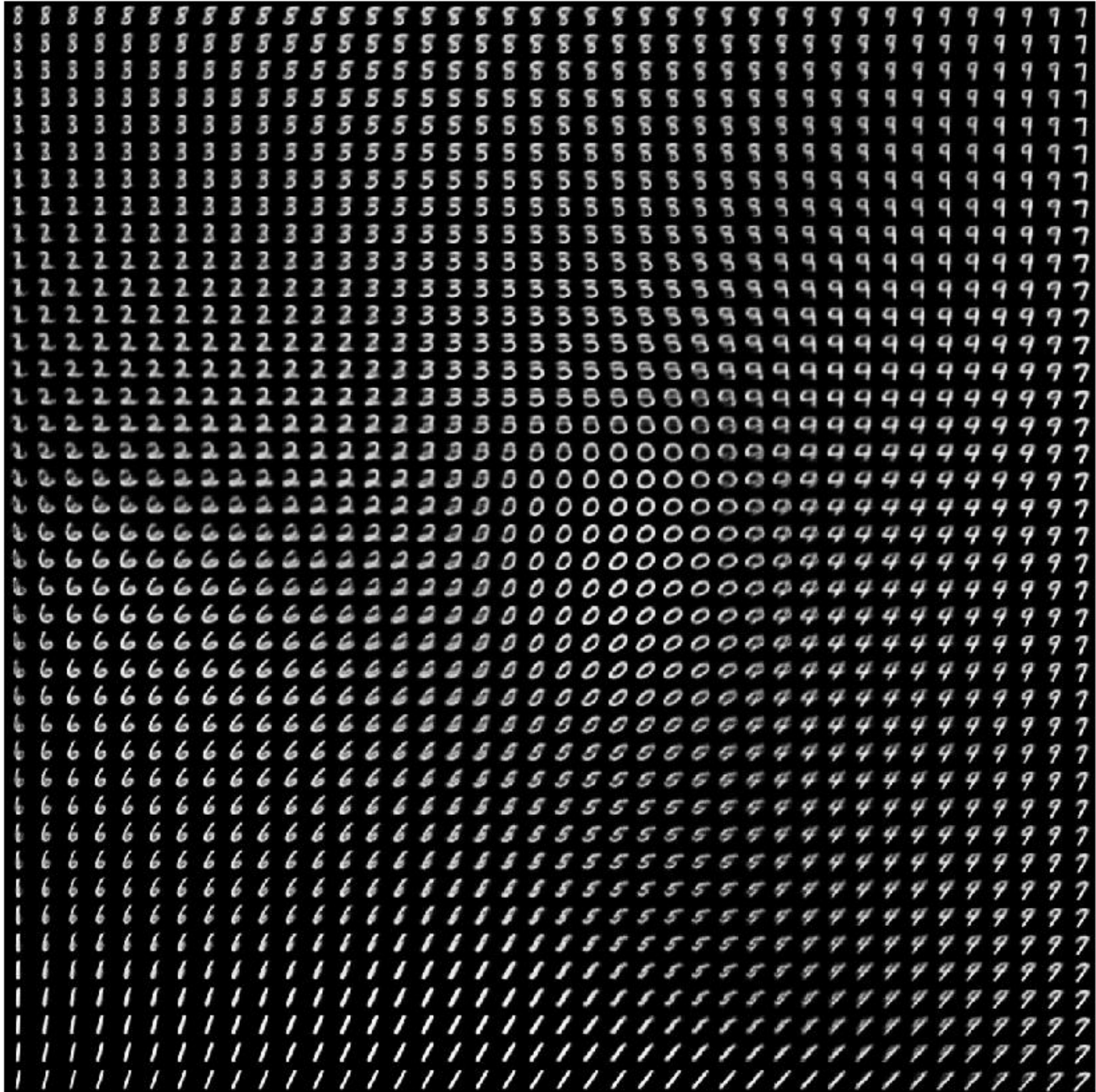
# Fantasies/Dreams of the VAE

- VAEs can disentangle potential factors of variation



[http://www.dpkingma.com/sgvb\\_mnist\\_demo/demo.html](http://www.dpkingma.com/sgvb_mnist_demo/demo.html)

Deep digit  
dreams...



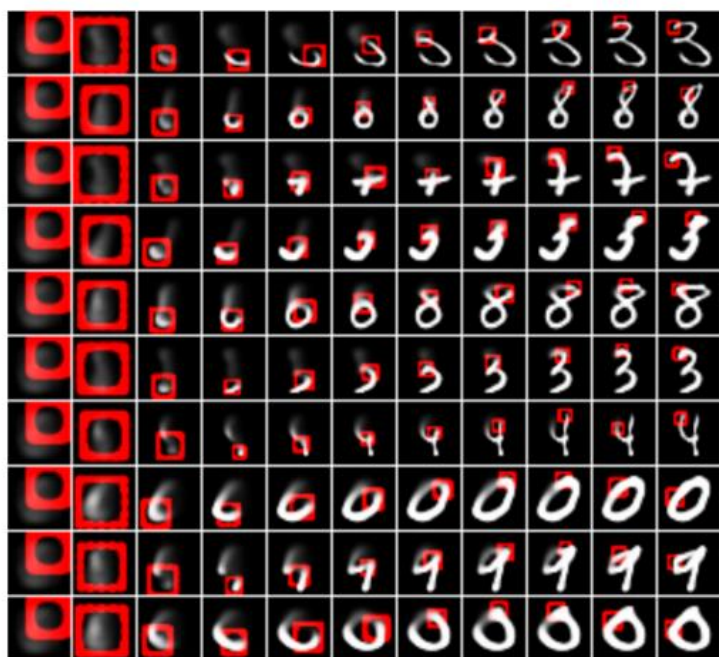
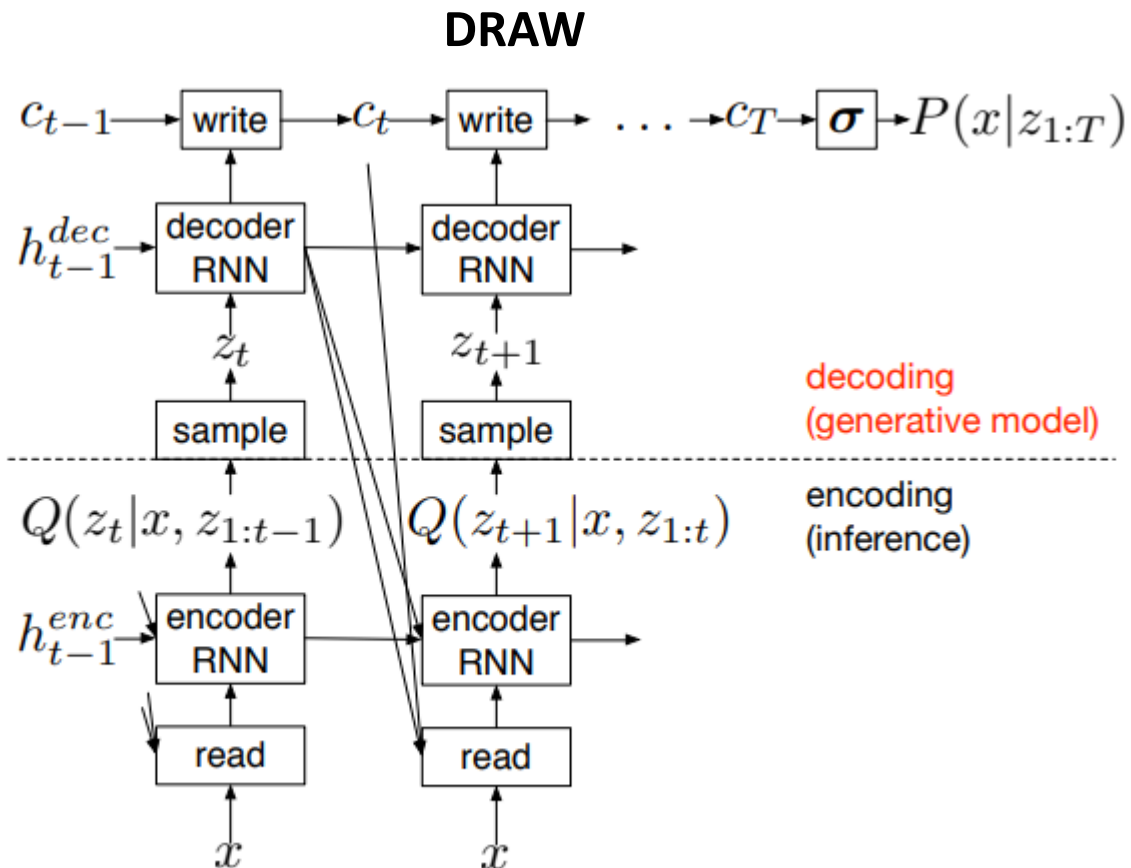
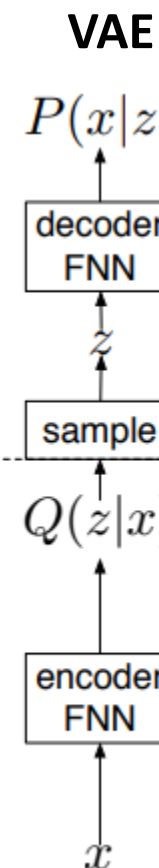
# Issues/Limitations with VAEs

- Images are blurry (compared to models based on GANs, for example)
  - Result of likelihood objective? (places probability mass on training images and nearby points, which include blurry images)
    - Has tendency to ignore input features that occupy few pixels or that cause only small change in brightness of pixels (that they occupy)
  - Uses only small subset of latent variables? (struggles to find enough transformation directions to match factorized prior over  $z$ )

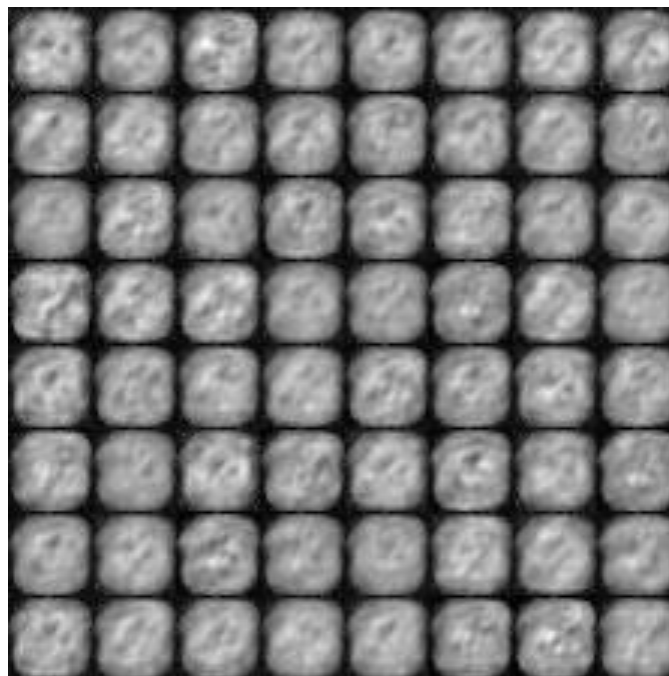
# DRAW

## The Deep Recurrent Attentive Writer

Combining latents, variational inference, and RNNs



Time →



# QUESTIONS?

