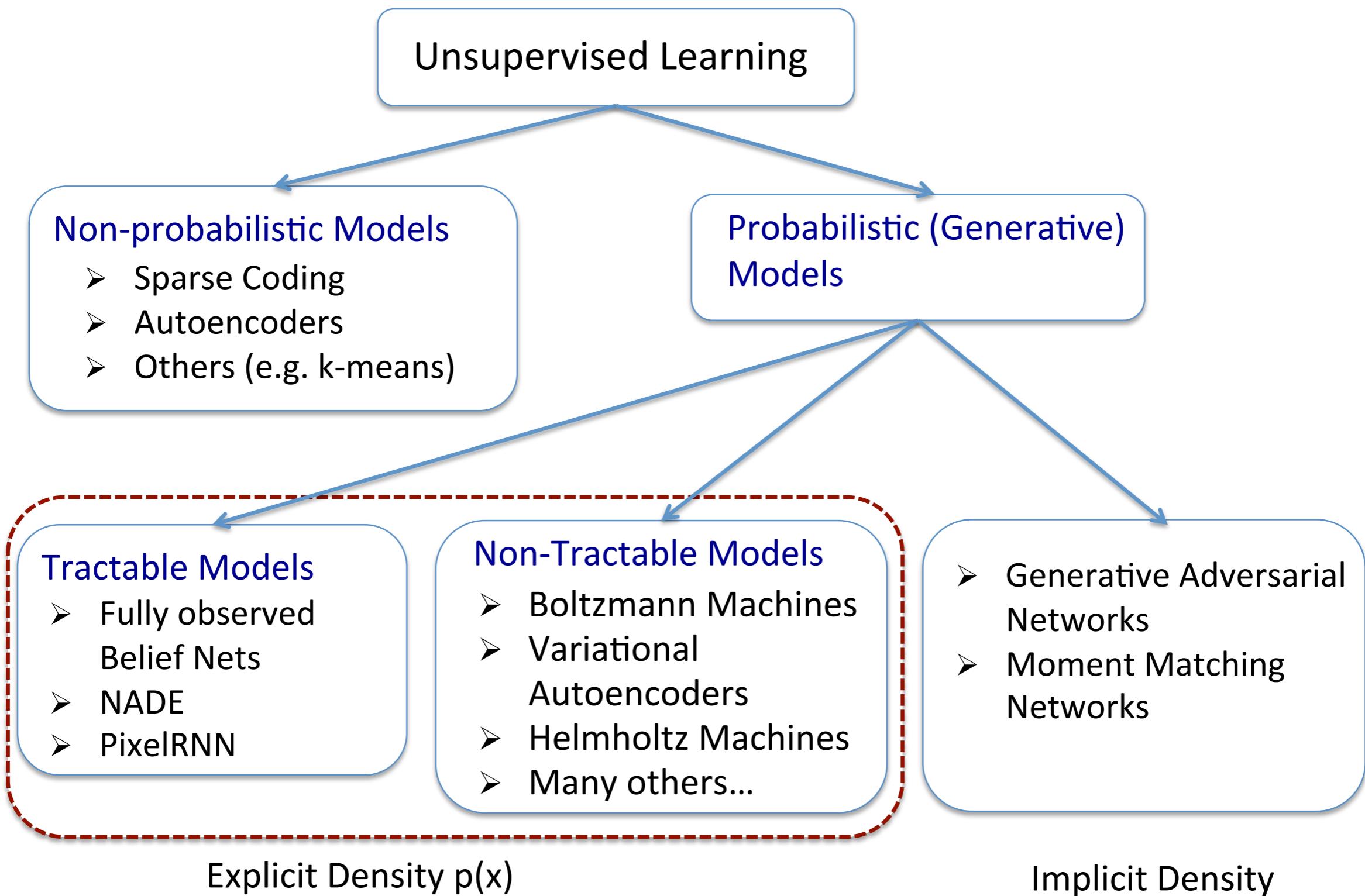


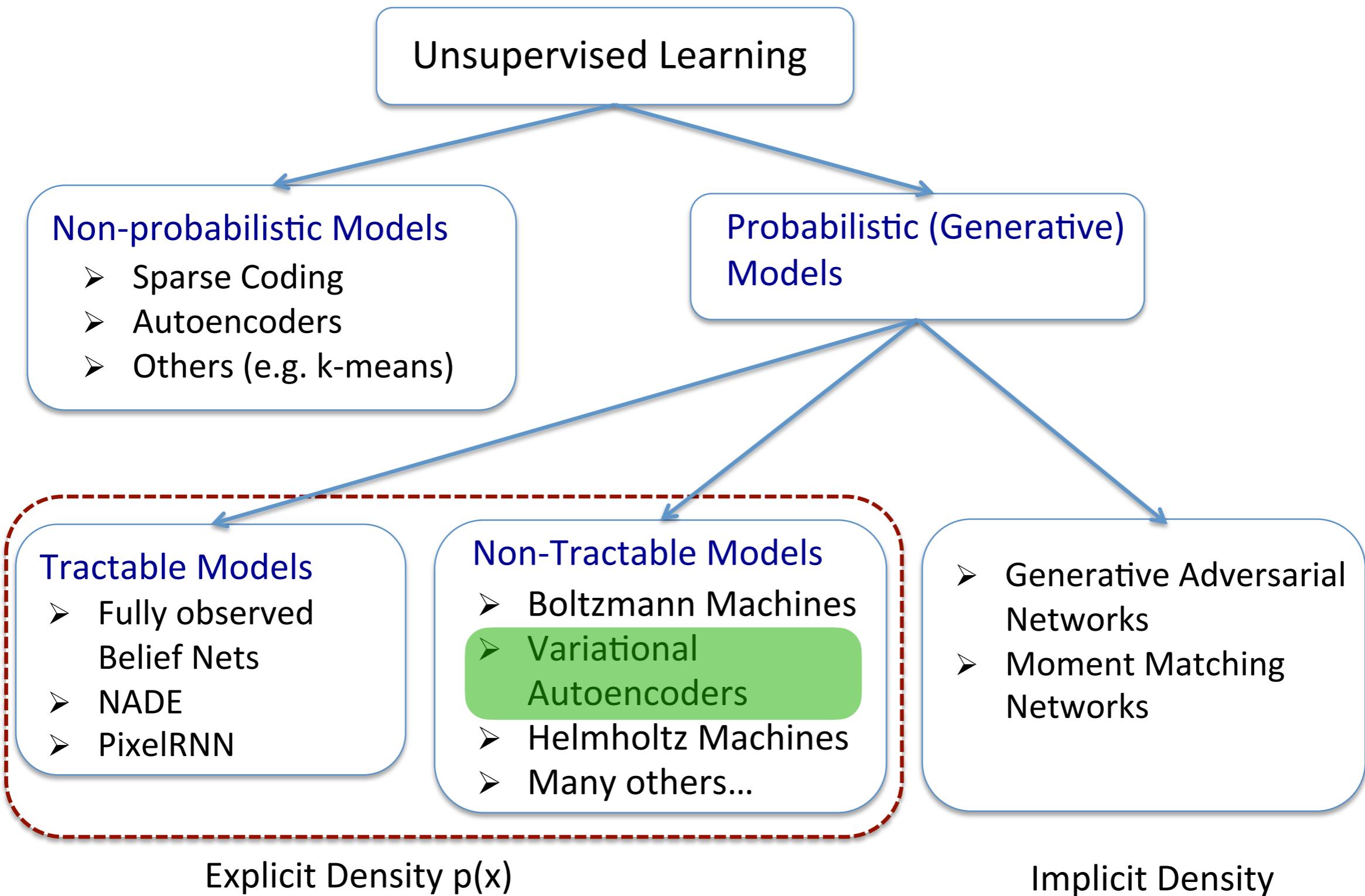
Deep Learning with Neural Networks

Deep Probabilistic Modelling: VAEs

Alejandro Lancho Serrano, alancho@ing.uc3m.es

Material original por: Pablo Martínez Olmos





Quick recap of autoencoders

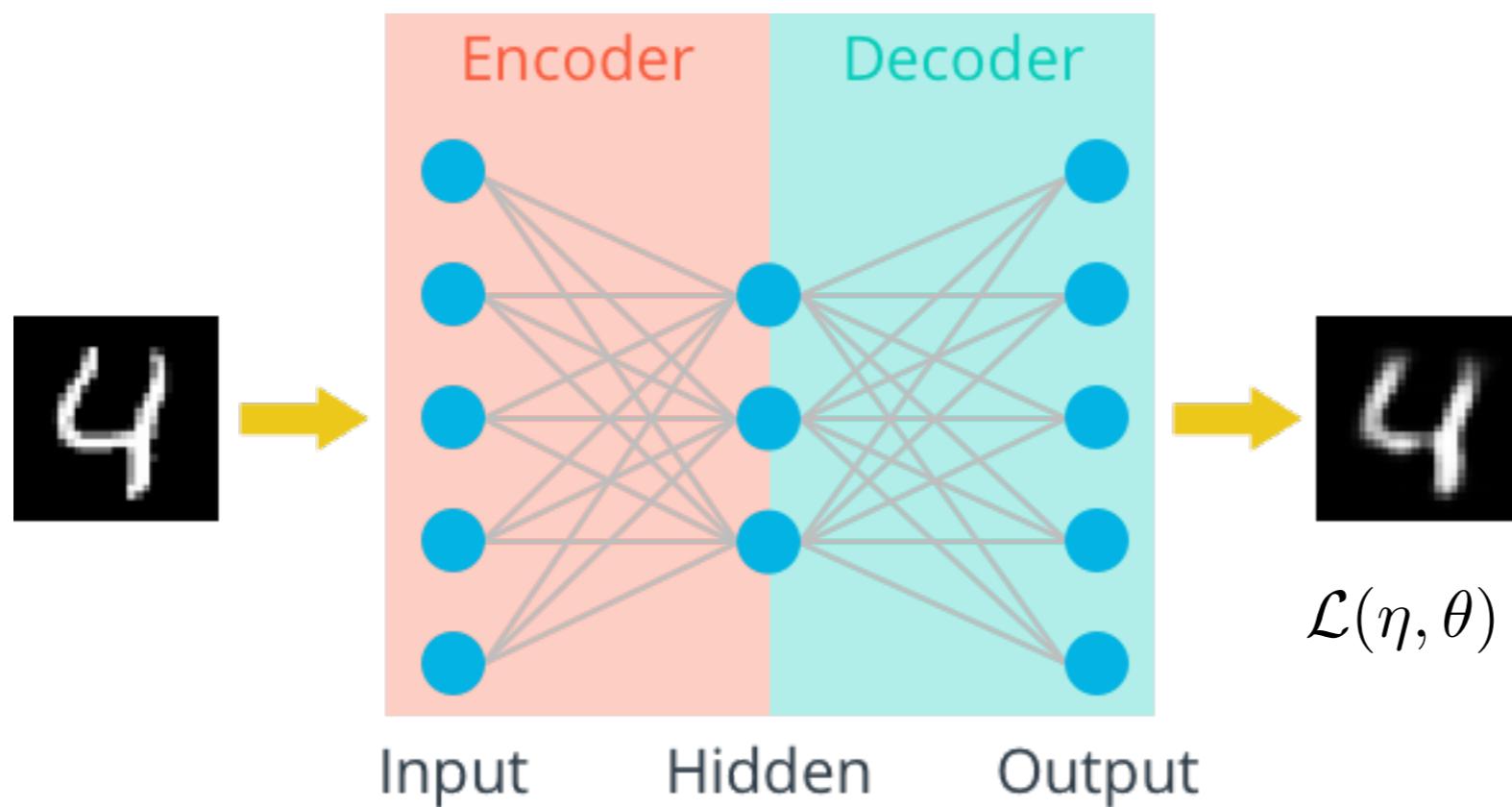
- The “encoder” learns a mapping from the data to a lower-dimensional latent space

$$\mathbf{z} = f(\mathbf{x}) = E_\eta(\mathbf{x}_n)$$

- The “decoder” learns a mapping from the latent space back to reconstruct the observation

$$\mathbf{z} = g(\mathbf{z}) = g(f(\mathbf{x})) = D_\theta(E_\eta(\mathbf{x}_n))$$

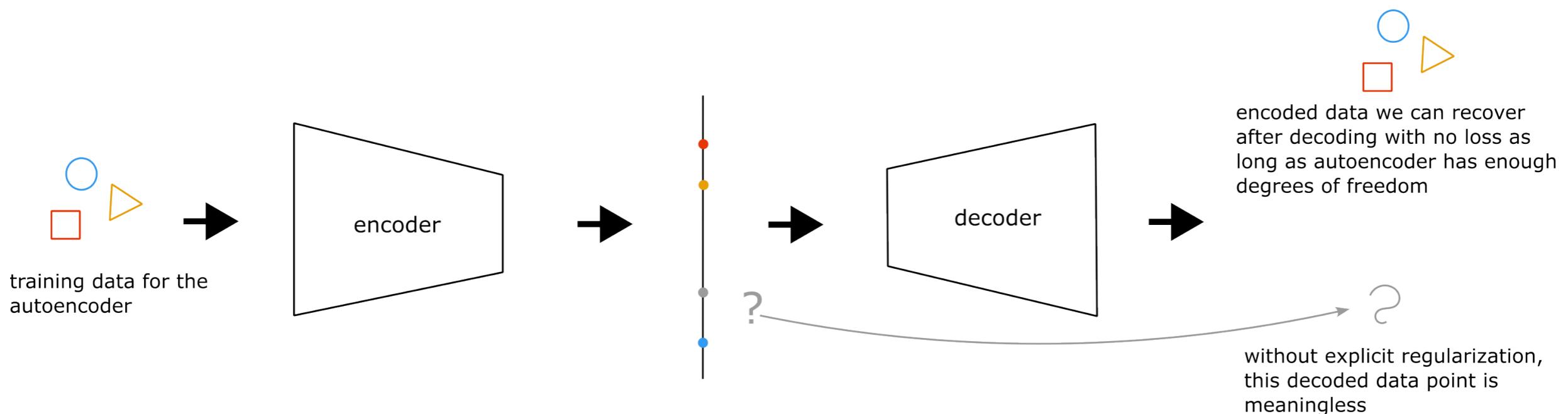
- The loss function doesn’t use any labels



$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

Quick recap of autoencoders

- The latent space is sometimes hard to regularize
- Struggle handling uncertainty
- One alternative are variational autoencoders (VAEs)
 - Regularize the latent space with variational inference principles
 - The latent space becomes rich enough for generation purposes

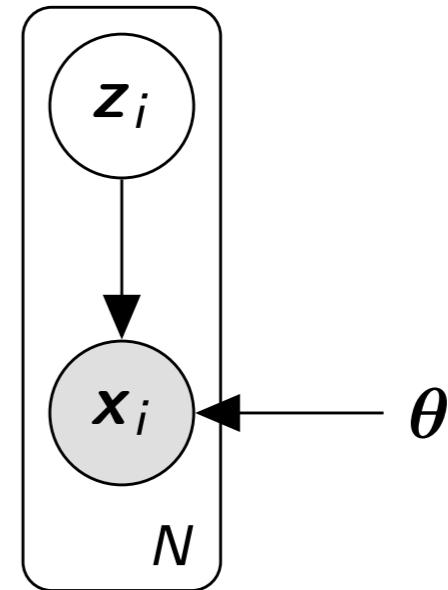


Variational Autoencoder (Kingma & Welling 2014)

Consider a set of i.i.d observations $\mathbf{x}^{(i)} \in \mathbb{R}^d$, $i = 1, \dots, N$. Let $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$.

Generative Model using a Latent Space

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|z)p(z)dz$$



- $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ ($k \leq d$) (low-dimensional embedding)
- $\mathbf{x}|\mathbf{z} \sim \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}(\mathbf{z})))$
- $\mu_{\theta}(\mathbf{z})$ and $\log(\sigma_{\theta}(\mathbf{z}))$ are the outputs of a NN $\mathbb{R}^k \rightarrow \mathbb{R}^d$ with parameter vector θ (Decoding network).

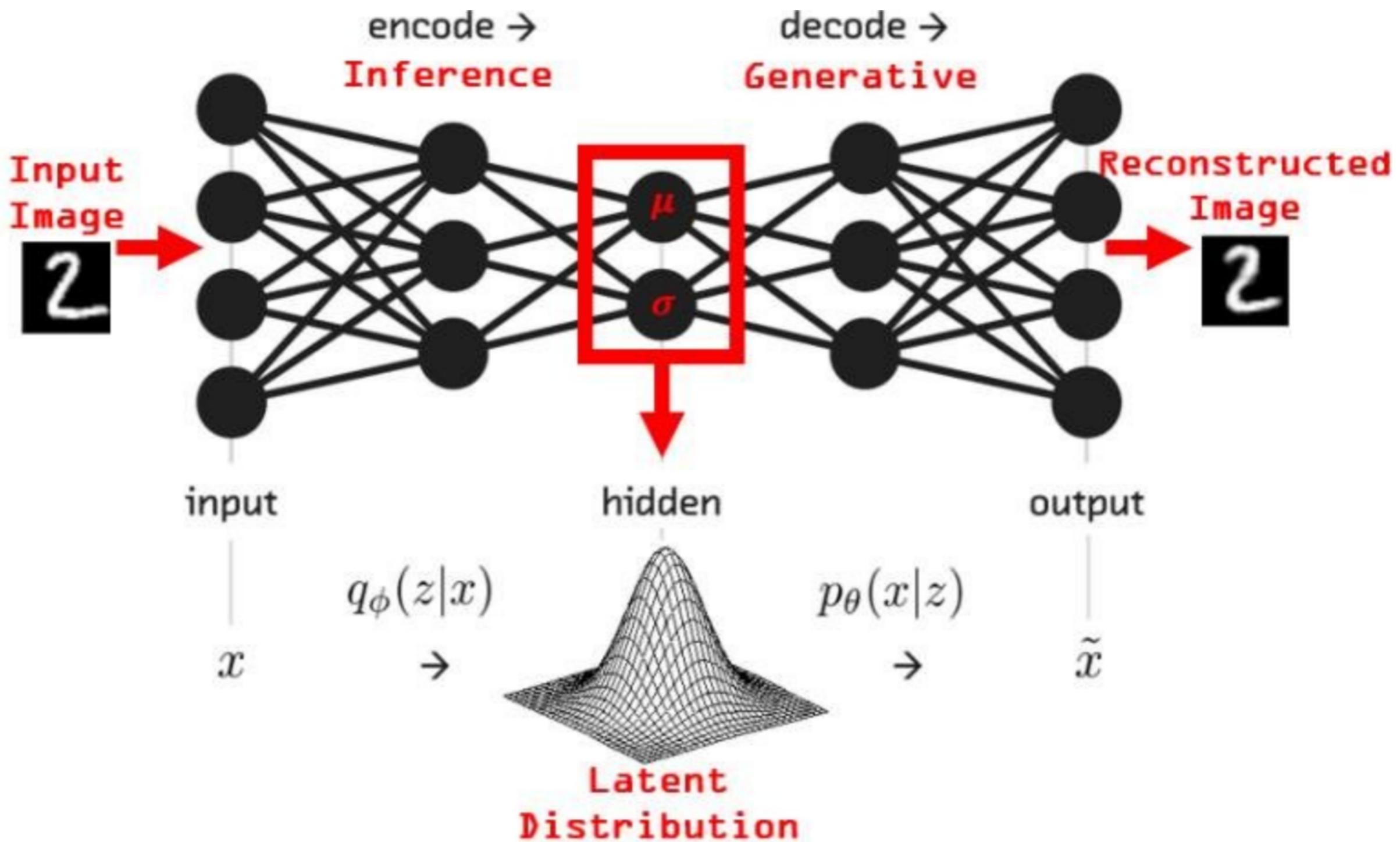
Variational Autoencoder. Approximate posterior

The Inference Network or Encoding Network

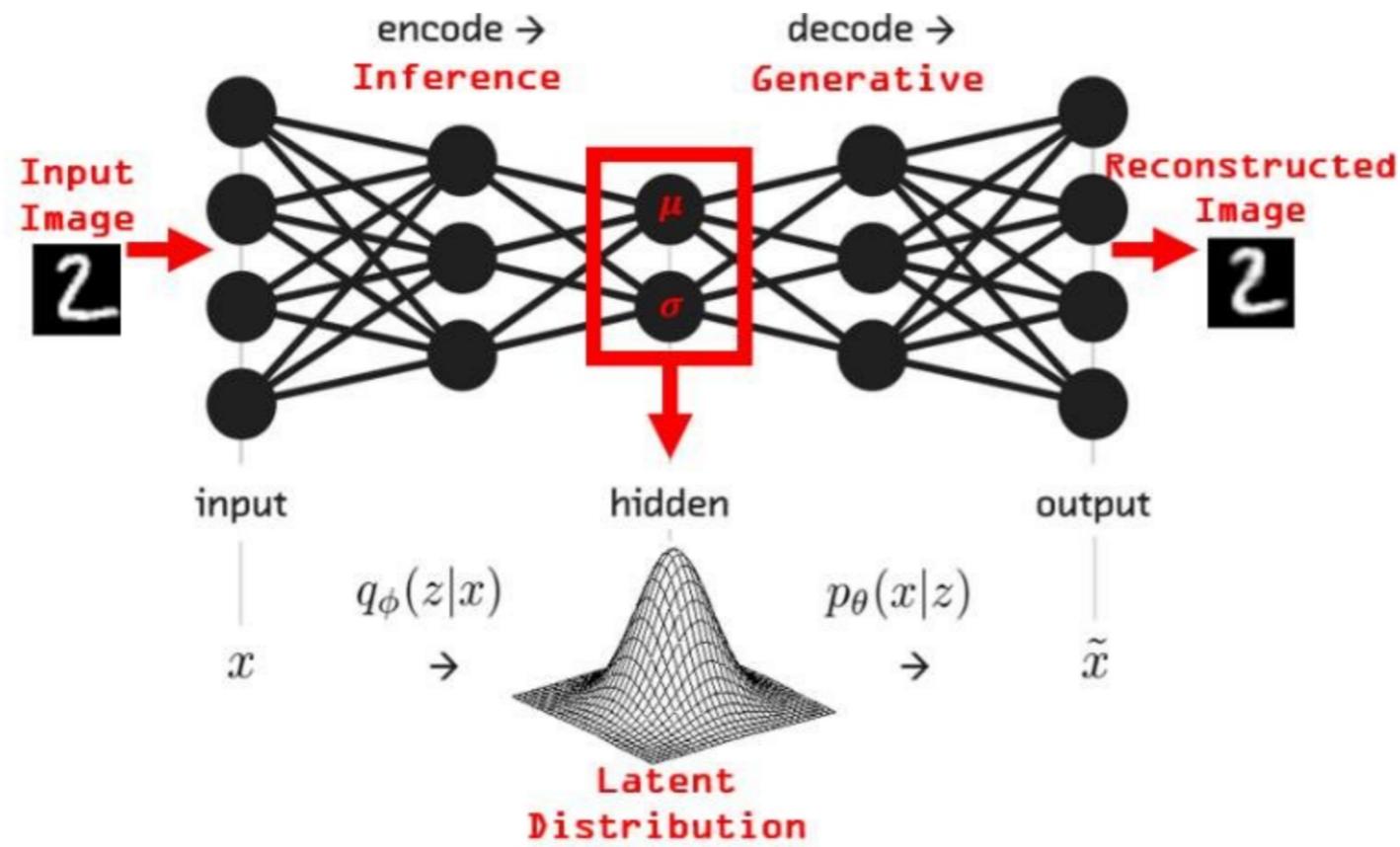
$$p_{\theta}(z|x) \approx q_{\eta,x}(z) = \mathcal{N}(\mu_{\eta}(x), \text{diag}(\sigma_{\eta}(x)))$$

where $\mu_{\eta}(x)$ and $\log(\sigma_{\eta}(x))$ are the outputs of a NN $\mathbb{R}^d \rightarrow \mathbb{R}^k$ with parameter vector η .

The variational autoencoder objective function



The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

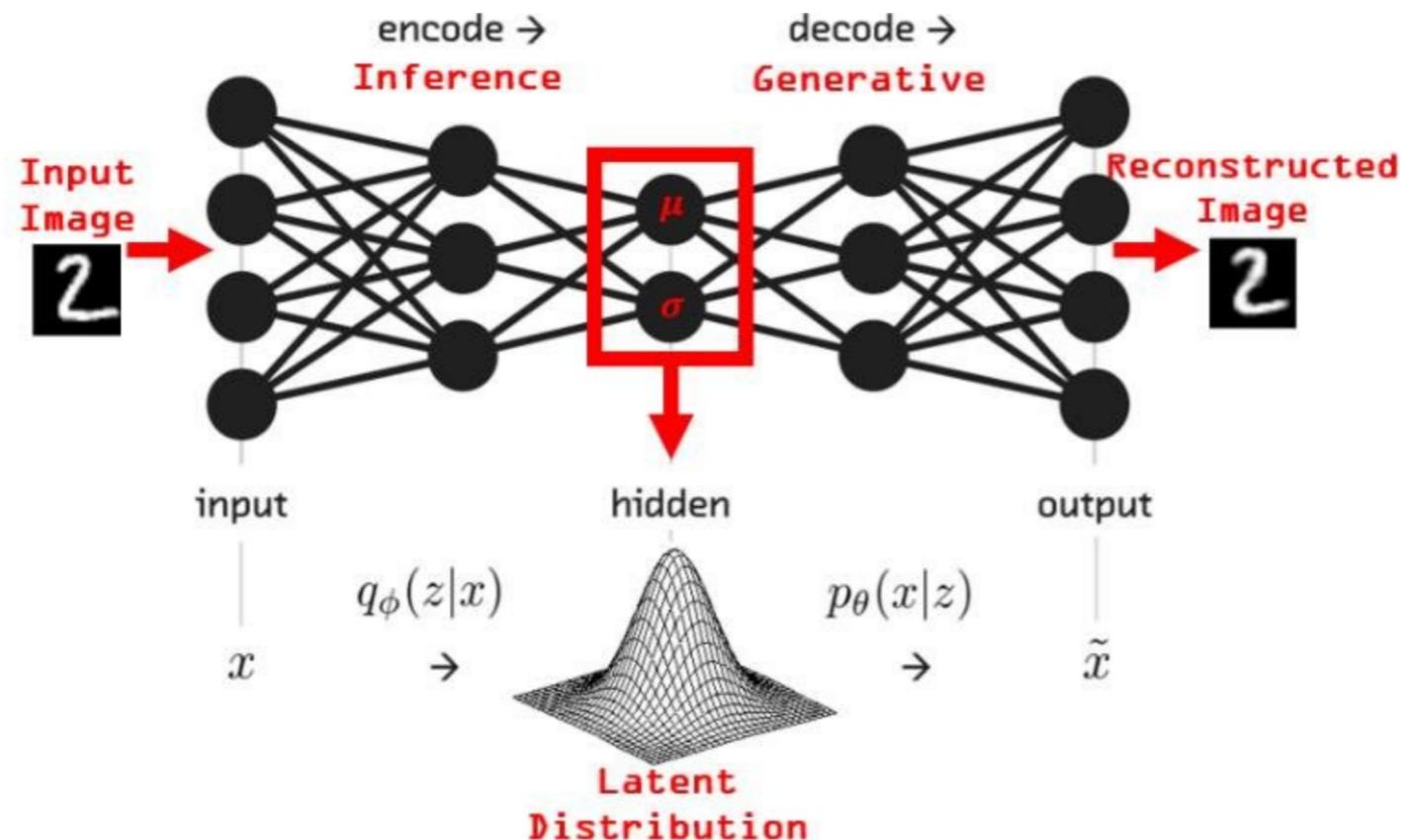
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

The variational autoencoder objective function



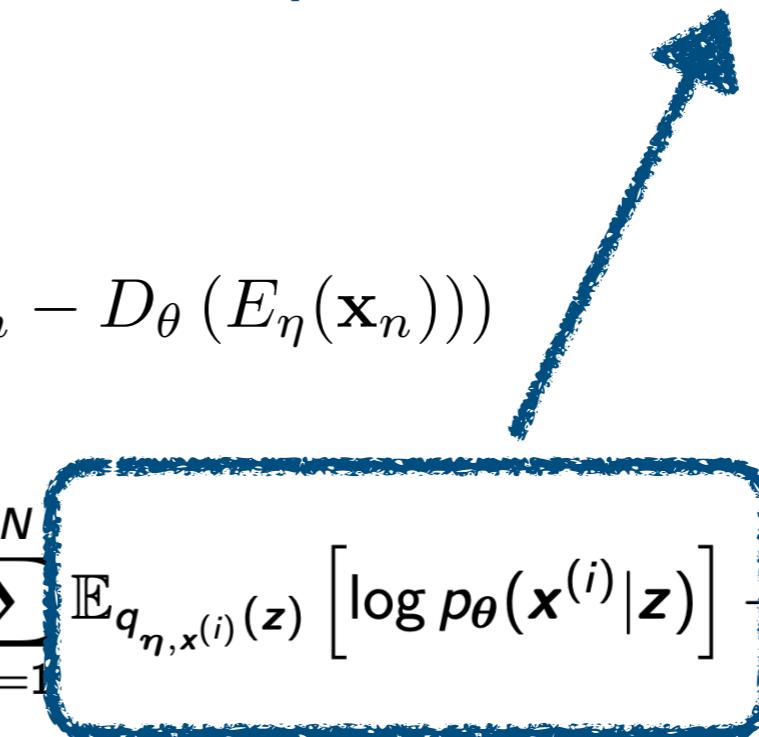
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

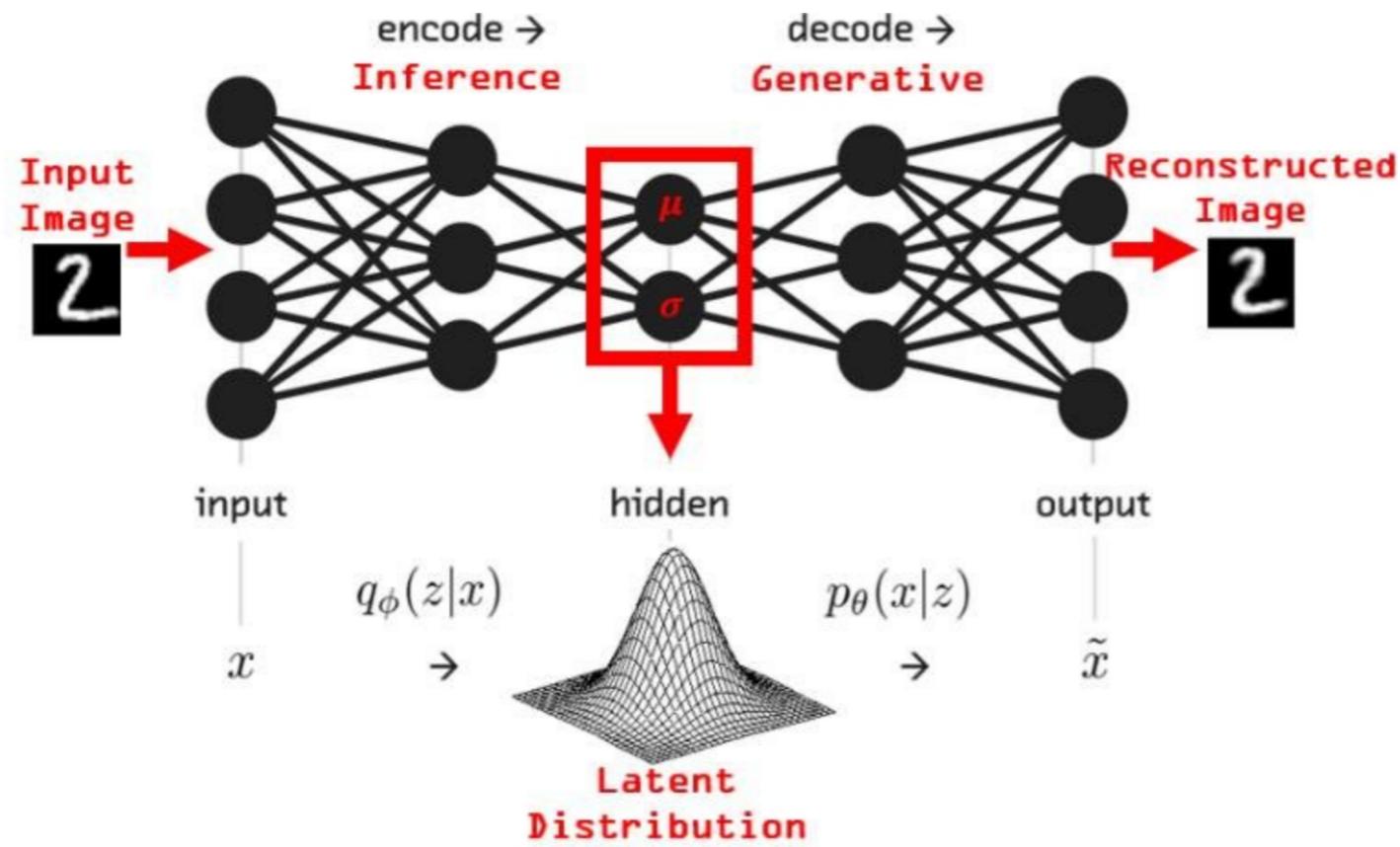
$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, x^{(i)}}(z)} [\log p_\theta(x^{(i)}|z)] - \text{KL}(q_{\eta, x^{(i)}}(z) || p(z)) \right)$$

- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

**Reconstruction term
(how well we explain training data)**



The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

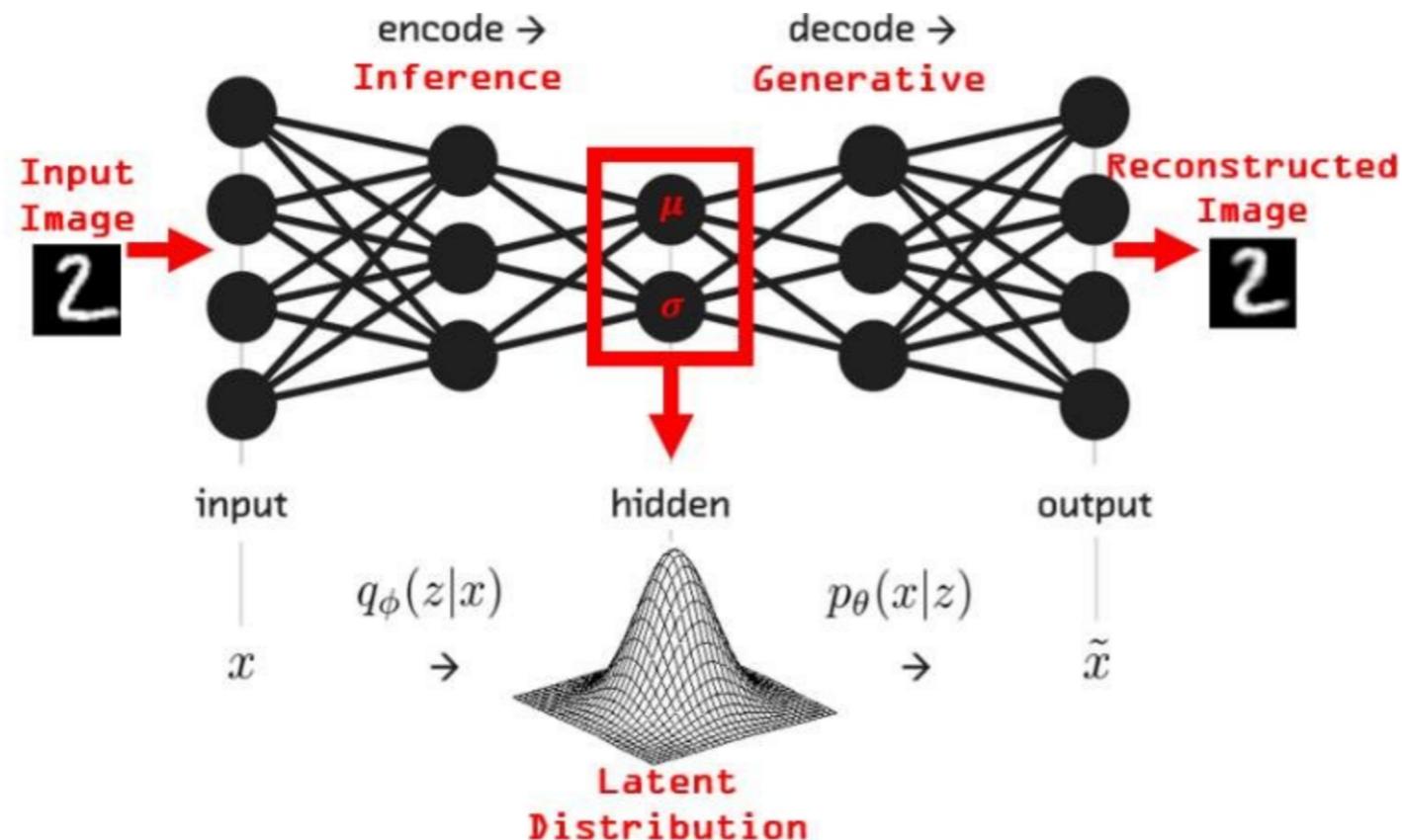
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

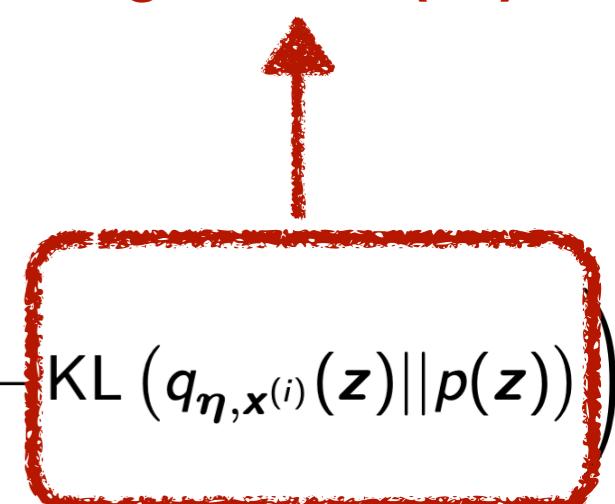
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Regularizer (>0)



Variational Autoencoder. Maximum Likelihood is intractable!

$$p_{\theta}(x) = \int p_{\theta}(x|z)p(z)dz$$

- We cannot compute such integral! ML training is infeasible
- We optimize a lower bound: the **variational** trick

Variational Autoencoder. Maximum Likelihood is intractable!

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

- We cannot compute such integral! ML training is infeasible
- We optimize a lower bound: the **variational** trick

An important concept: KL divergence

- Measures the discrepancy between two distributions
- 0 if the two distributions are identical
- Always positive

$$\text{KL}(p(\mathbf{x})||q(\mathbf{x})) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} \geq 0$$

The variational objective is a lower bound on marginal likelihood

$$p_{\theta}(x) = \int p_{\theta}(x|z)p(z)dz$$

Consider **any** approximation $q(z)$ to $p_{\theta}(z|x) \propto p_{\theta}(x|z)p(z)$

$$\log p_{\theta}(x) = \text{KL}(q(z)||p_{\theta}(z|x)) + \mathcal{L}(x, \theta)$$

The variational objective is a lower bound on marginal likelihood

$$p_{\theta}(x) = \int p_{\theta}(x|z)p(z)dz$$

Consider **any** approximation $q(z)$ to $p_{\theta}(z|x) \propto p_{\theta}(x|z)p(z)$

$$\log p_{\theta}(x) = \text{KL}(q(z)||p_{\theta}(z|x)) + \mathcal{L}(x, \theta)$$

where

$$\mathcal{L}(x, \theta) = \int q(z) \log \frac{p_{\theta}(x|z)p(z)}{q(z)} dz = \mathbb{E}_{q(z)} [\log p_{\theta}(x|z)] - \text{KL}(q(z)||p(z))$$

is the Evidence Lower Bound (ELBO).

We will optimize $\mathcal{L}(x, \theta)$ w.r.t. θ . But first we need to select a variational family for $q(z)$.

Variational Autoencoder. The variational objective

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

The Inference Network or Encoding Network

$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta,\mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

where $\mu_{\eta}(\mathbf{x})$ and $\log(\sigma_{\eta}(\mathbf{x}))$ are the outputs of a NN $\mathbb{R}^d \rightarrow \mathbb{R}^k$ with parameter vector η .

Variational Autoencoder. The variational objective

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

The Inference Network or Encoding Network

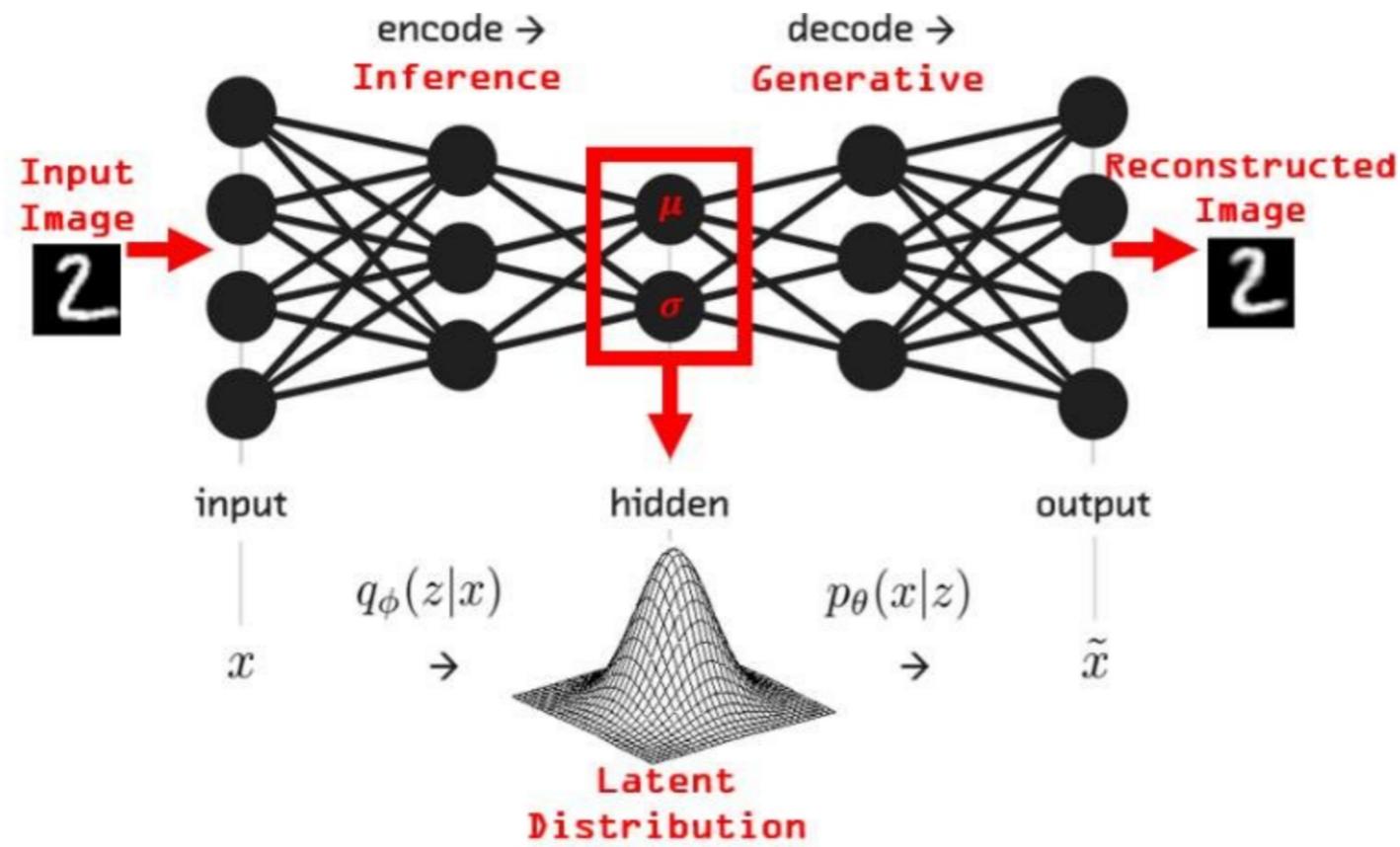
$$p_{\theta}(\mathbf{z}|\mathbf{x}) \approx q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$$

where $\mu_{\eta}(\mathbf{x})$ and $\log(\sigma_{\eta}(\mathbf{x}))$ are the outputs of a NN $\mathbb{R}^d \rightarrow \mathbb{R}^k$ with parameter vector η .

The VAE optimization problem

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

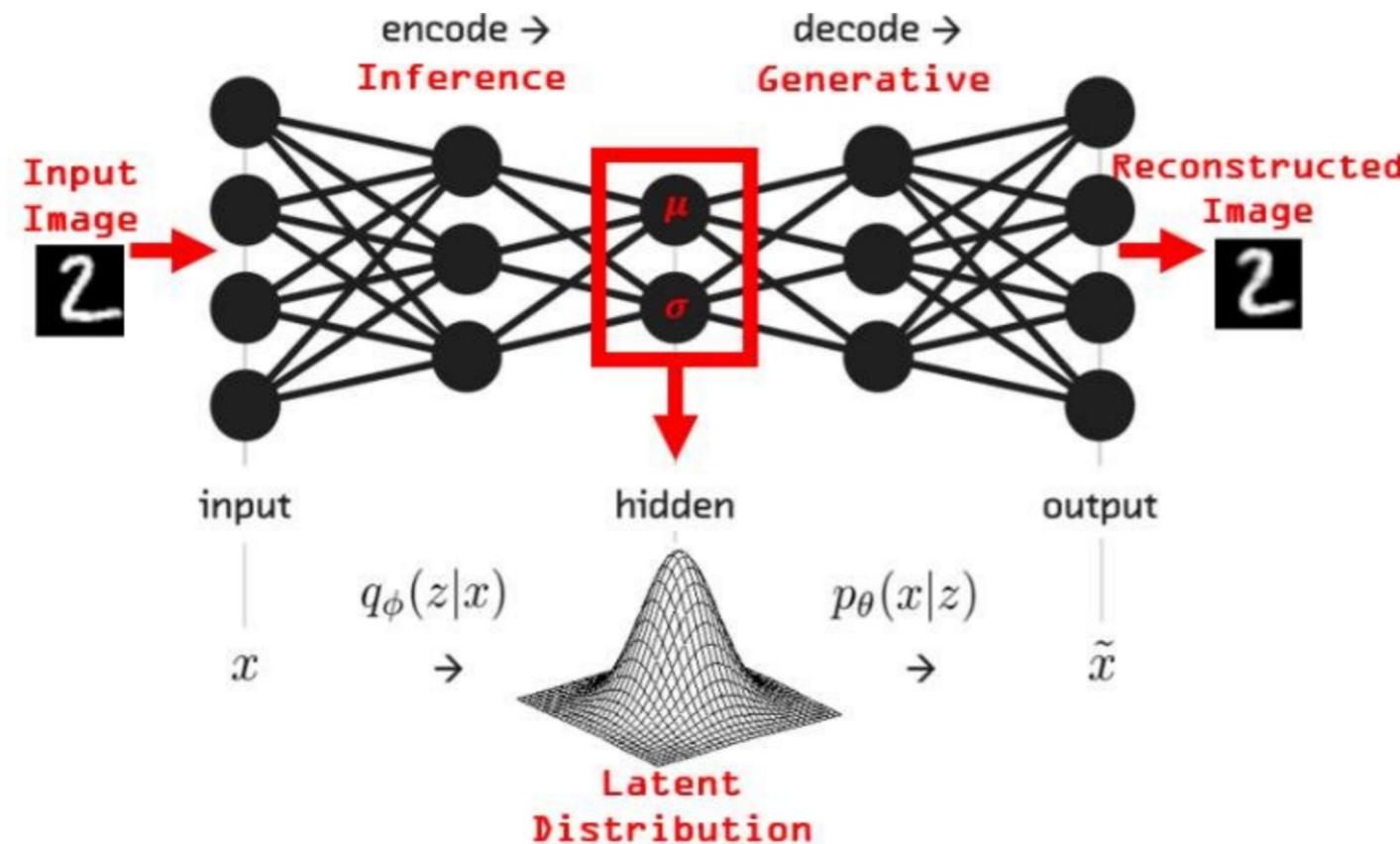
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

The variational autoencoder objective function



Autoencoder

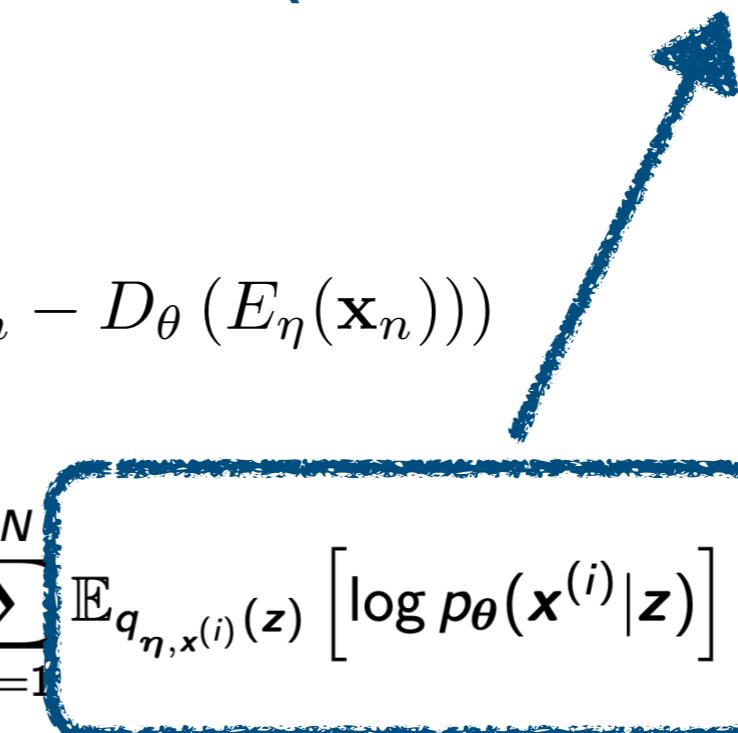
$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

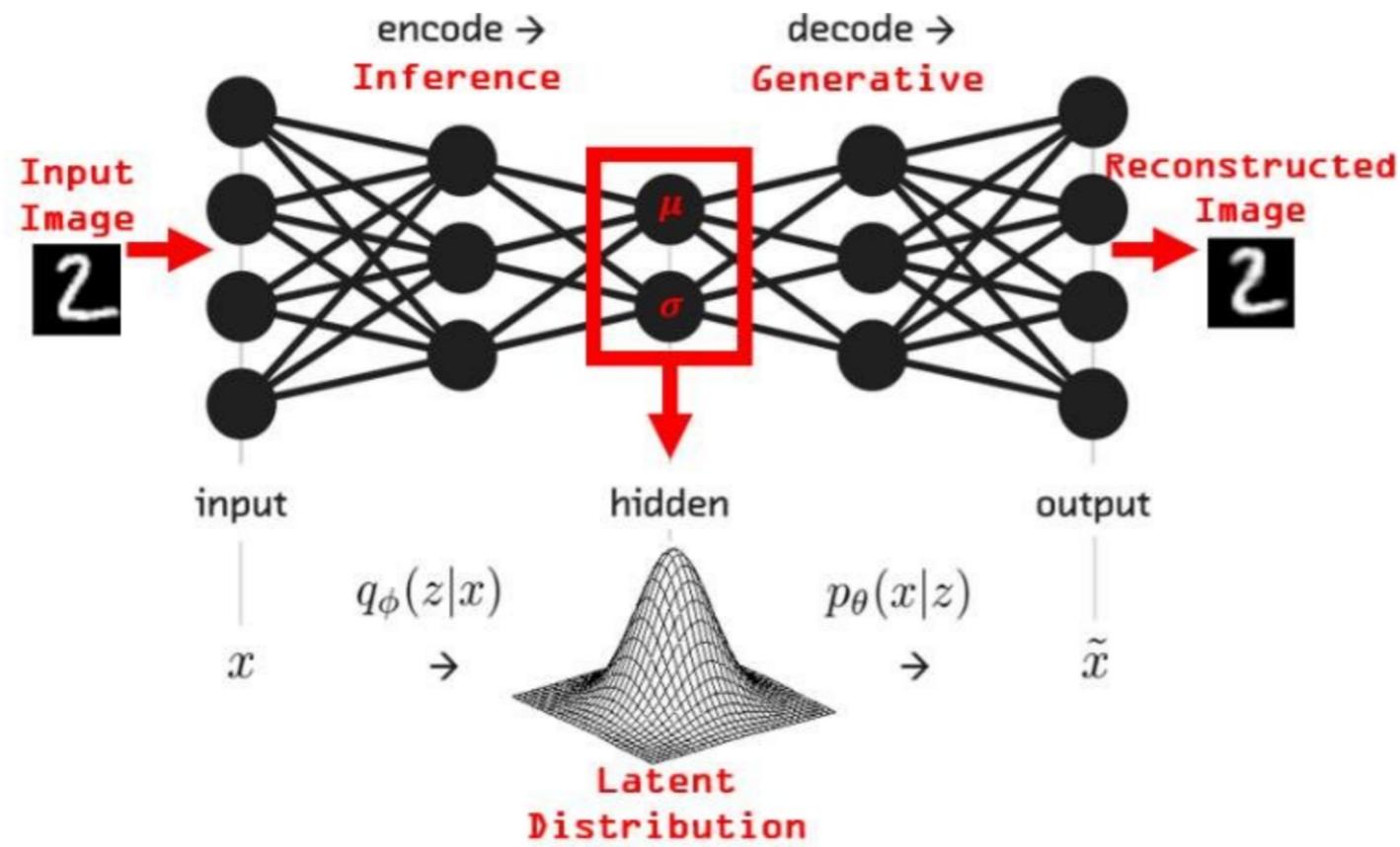
$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

**Reconstruction term
(how well we explain training data)**



The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

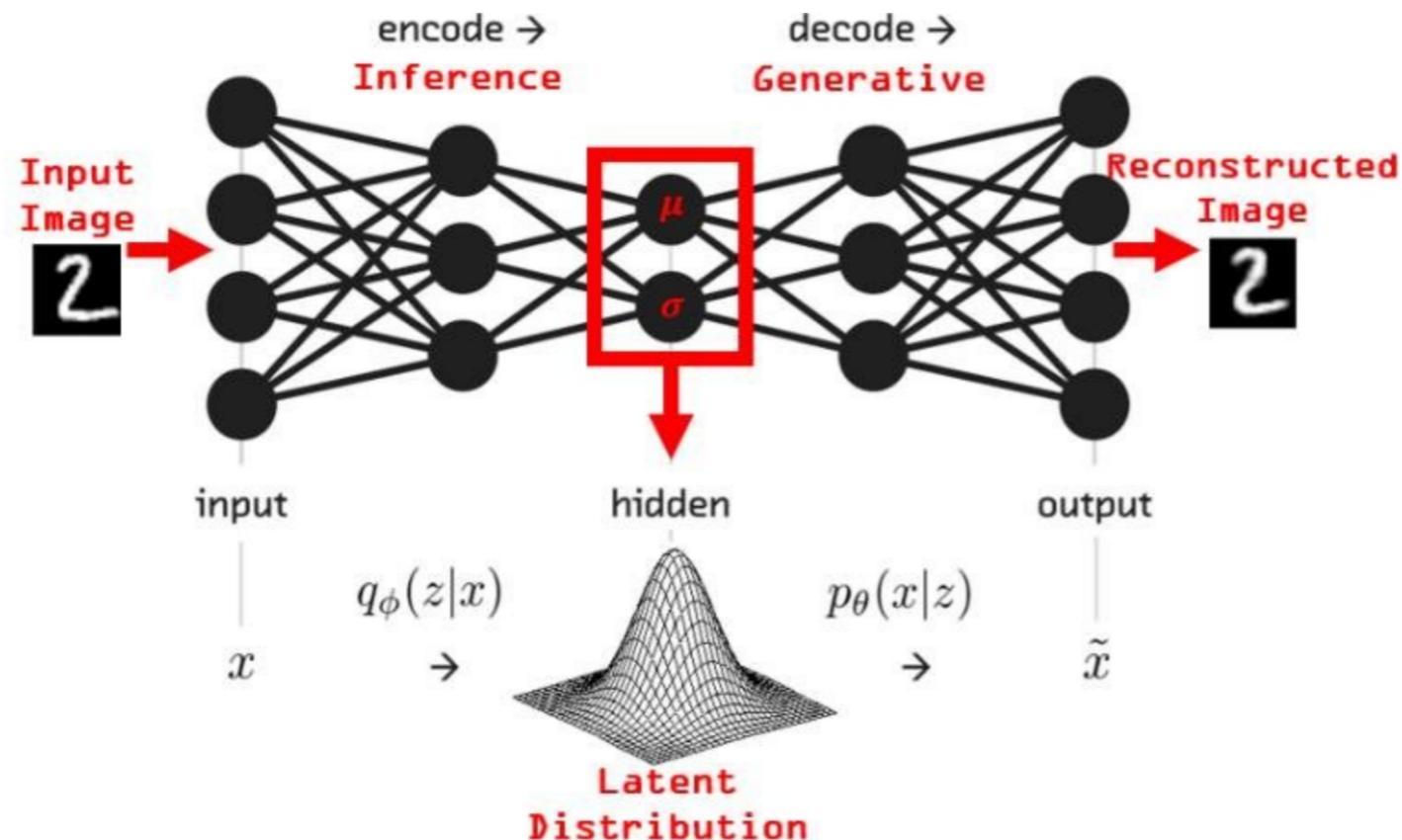
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

The variational autoencoder objective function



- Probabilistic formulation of autoencoders
- Handle uncertainty
- They perform density estimation!

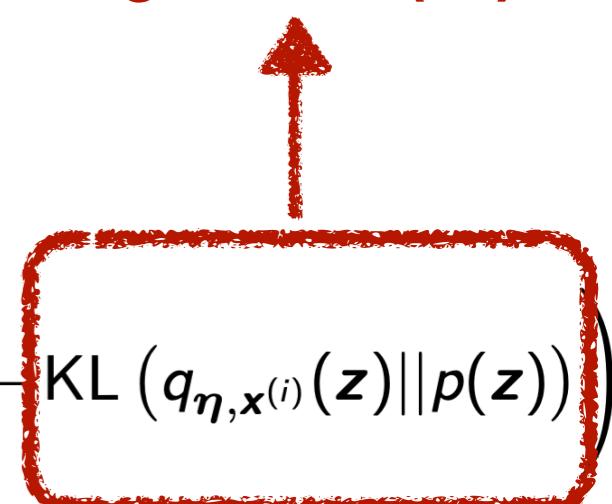
Autoencoder

$$\mathcal{L}(\eta, \theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n (\mathbf{x}_n - D_\theta (E_\eta(\mathbf{x}_n)))$$

VAE

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \frac{1}{N} \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})] - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Regularizer (>0)



Variational Autoencoder. Stochastic Optimization over the ELBO

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, x^{(i)}}(z)} \left[\log p_{\theta}(x^{(i)}|z) \right] - \text{KL} (q_{\eta, x^{(i)}}(z) || p(z)) \right)$$

Recall that $p(z) = \mathcal{N}(\mathbf{0}, I)$ and $q_{\eta, x}(z) = \mathcal{N}(\mu_{\eta}(x), \text{diag}(\sigma_{\eta}(x)))$, thus

$$\text{KL} (q_{\eta, x}(z) || p(z)) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(x) - \log \sigma_{\eta, j}(x) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

Variational Autoencoder. Stochastic Optimization over the ELBO

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Recall that $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$, thus

$$\text{KL} (q_{\eta, \mathbf{x}}(\mathbf{z}) || p(\mathbf{z})) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(\mathbf{x}) - \log \sigma_{\eta, j}(\mathbf{x}) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

We use a Monte Carlo sampling estimator

$$\mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] \approx \frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{i,s})$$

where $\mathbf{z}^{(s,i)}$, $s = 1, \dots, S$ are i.i.d. samples from $q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})$.

We typically use a **single** sample, i.e., $S = 1$ (huge estimator variance, but cheap computation).

Variational Autoencoder. Stochastic Optimization over the ELBO

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Recall that $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$, thus

$$\text{KL} (q_{\eta, \mathbf{x}}(\mathbf{z}) || p(\mathbf{z})) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(\mathbf{x}) - \log \sigma_{\eta, j}(\mathbf{x}) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

Variational Autoencoder. Stochastic Optimization over the ELBO

$$\max_{\theta, \eta} \mathcal{L}(\theta, \eta) = \max_{\theta, \eta} \left(\sum_{i=1}^N \mathbb{E}_{q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z})} \left[\log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \text{KL} (q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right)$$

Recall that $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $q_{\eta, \mathbf{x}}(\mathbf{z}) = \mathcal{N}(\mu_{\eta}(\mathbf{x}), \text{diag}(\sigma_{\eta}(\mathbf{x})))$, thus

$$\text{KL} (q_{\eta, \mathbf{x}}(\mathbf{z}) || p(\mathbf{z})) = \frac{1}{2} \left(-k + \sum_{j=1}^k \sigma_{\eta, j}(\mathbf{x}) - \log \sigma_{\eta, j}(\mathbf{x}) + \mu_{\eta, j}^2 \right)$$

→ Differentiable w.r.t. η

Reparameterization Trick

Express each sample $\mathbf{z}^{(s, i)}$ as a deterministic function of $\mathbf{x}^{(i)}$ and some noise vector $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ that is independent of η . For a Gaussian distribution we have

$$\mathbf{z}^{(s, i)} = f_{\eta}(\mathbf{x}^{(i)}, \epsilon) = \mu_{\eta}(\mathbf{x}^{(i)}) + \sqrt{\sigma_{\eta}(\mathbf{x}^{(i)})} \cdot \epsilon$$

Variational Autoencoder. Algorithm

Algorithm 1 The Variational Autoencoder (S=1)

- 1: $\theta \leftarrow \theta_0, \eta \leftarrow \eta_0$
- 2: $\ell \leftarrow 0$
- 3: **while** not converged **do**
- 4: Sample minibatch $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ from \mathcal{D}
- 5: Sample $\{\epsilon^1, \dots, \epsilon^M\}$ from $p(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: Compute noisy gradients:

$$\mathbf{g}_\theta, \mathbf{g}_\eta \leftarrow \frac{1}{M} \sum_{i=1}^N \nabla_{\theta, \eta} \left[\log p_\theta(\mathbf{x}^{(i)} | f_\eta(\epsilon^i, \mathbf{x}^{(i)})) - \text{KL}(q_{\eta, \mathbf{x}^{(i)}}(\mathbf{z}) || p(\mathbf{z})) \right]$$

- 7: Perform SGD-updates:

$$\begin{aligned}\theta_{\ell+1} &\leftarrow \theta_{\ell+1} + h_\ell \mathbf{g}_\theta \\ \eta_{\ell+1} &\leftarrow \eta_{\ell+1} + h_\ell \mathbf{g}_\eta\end{aligned}$$

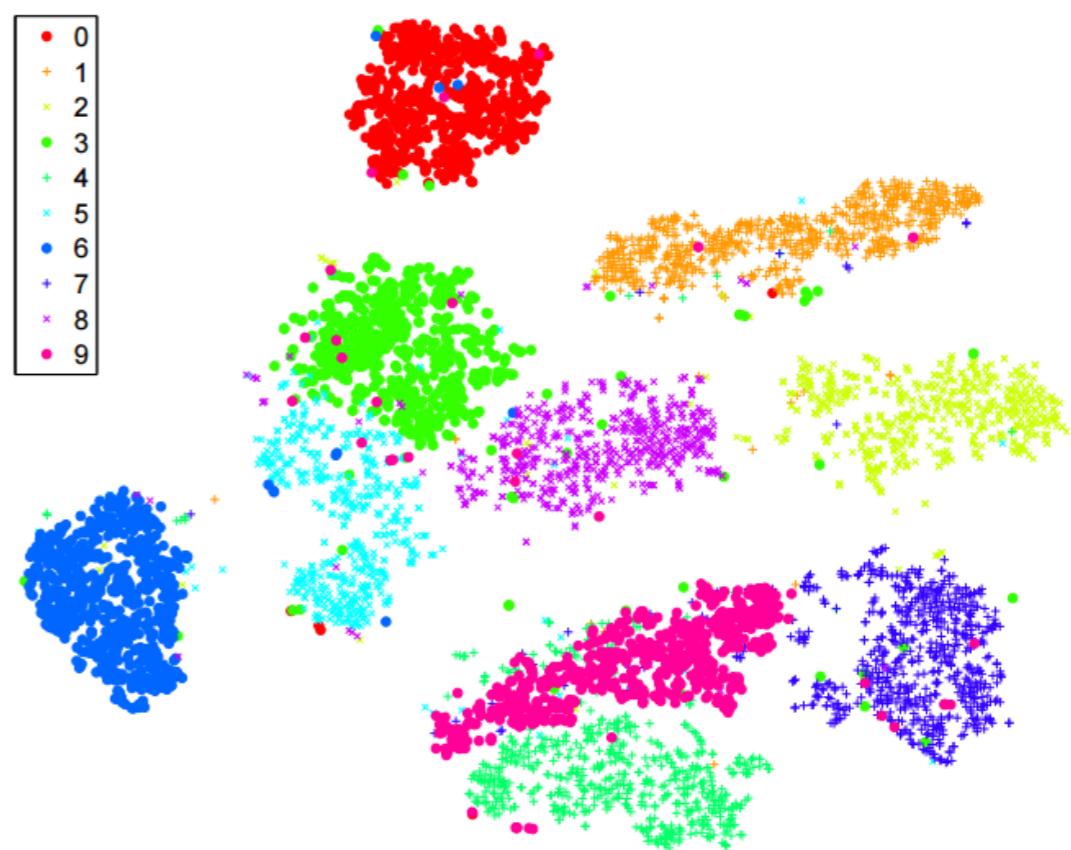
- 8: $\ell \leftarrow \ell + 1$
 - 9: **end while**
-

Variational Autoencoder. Clustering in the latent space



Generated Samples

Projection over latent space



VQ-VAE

- Uses vector quantization (VQ) to obtain a **discrete latent representation**
- How it differs from traditional VAEs:
 - ▶ Encoder network output is discrete
 - ▶ Prior is learnt rather than static
- VQ ideas used to learn the discrete latent space
 - ▶ Help improve posterior collapse

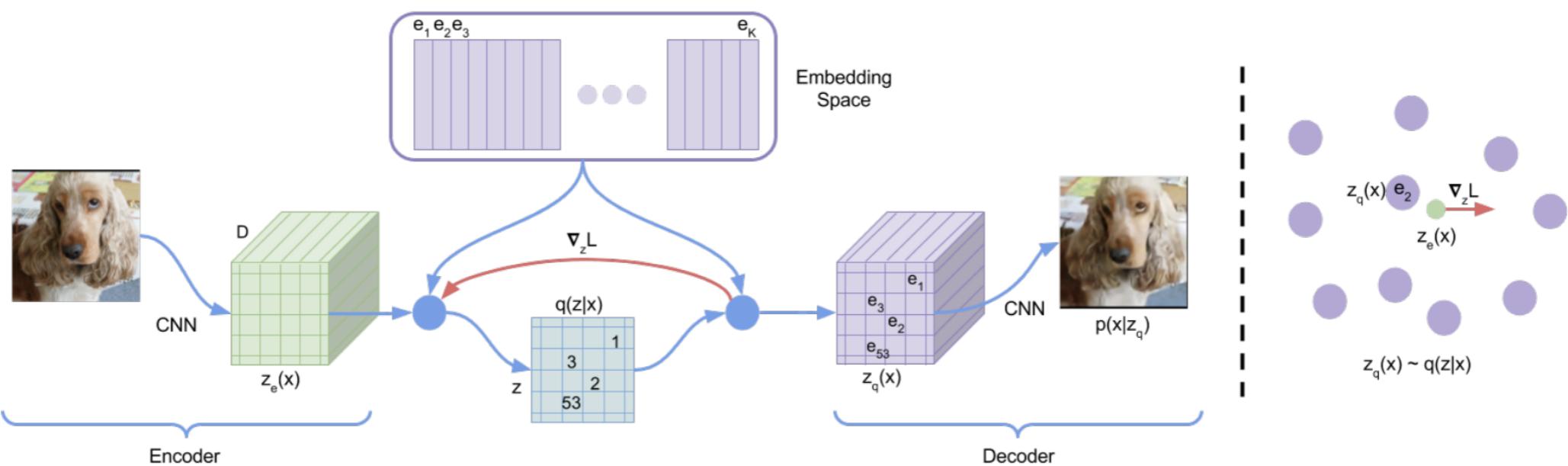


Figure 1: Left: A figure describing the VQ-VAE. Right: Visualisation of the embedding space. The output of the encoder $z(x)$ is mapped to the nearest point e_2 . The gradient $\nabla_z L$ (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.

VQ-VAE

- Uses vector quantization (VQ) to obtain a **discrete latent representation**
- How it differs from traditional VAEs:
 - ▶ Encoder network output is discrete
 - ▶ Prior is learnt rather than static
- VQ ideas used to learn the discrete latent space
 - ▶ Help improve posterior collapse

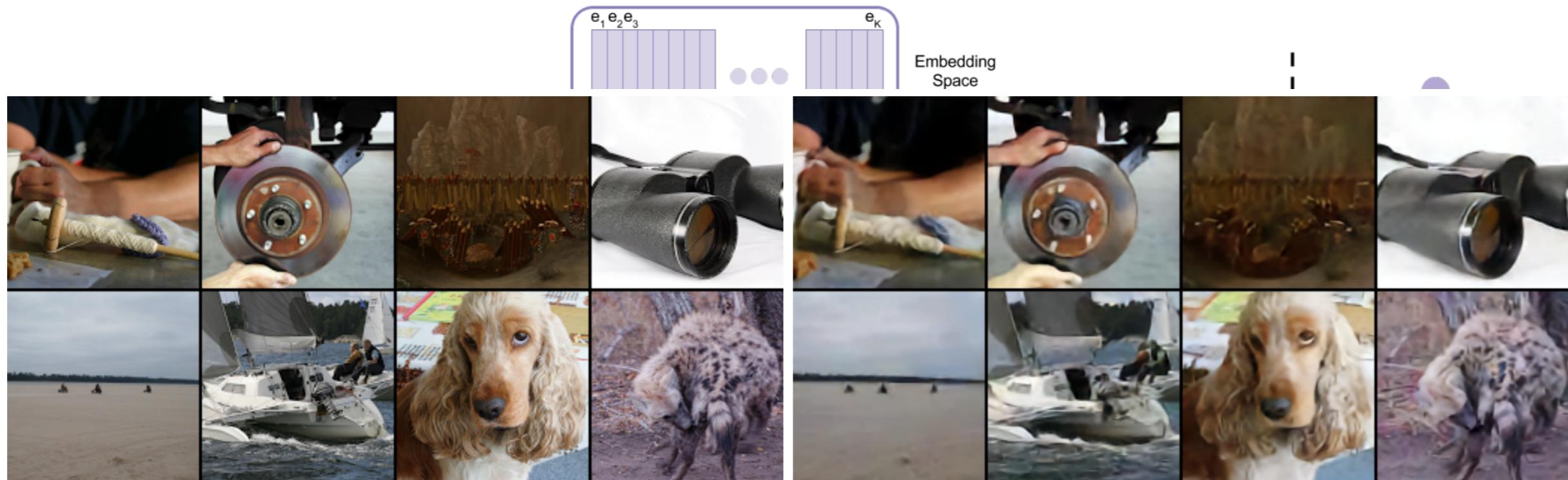
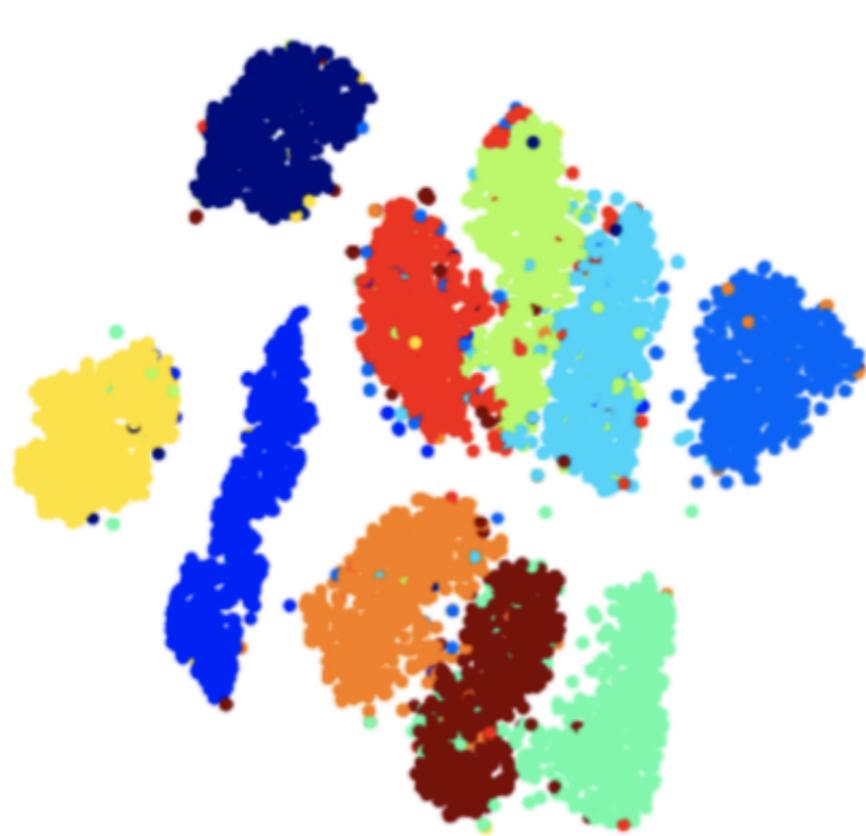


Figure 2: Left: ImageNet 128x128x3 images, right: reconstructions from a VQ-VAE with a 32x32x1 latent space, with K=512.

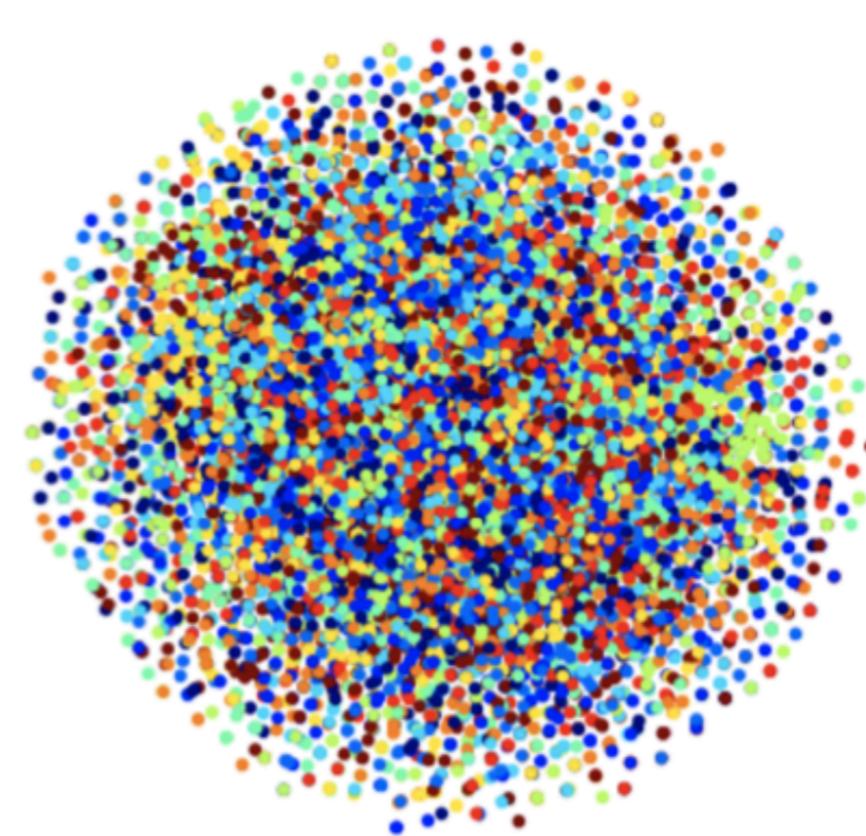
[Aaron van den Oord, Oriol Vinyals, Koray Kavukcuoglu, "Neural Discrete Representation Learning"](#) NeurIPS 2017

Posterior Collapse



Ideal

Model fits well - good predictive likelihood and good generation

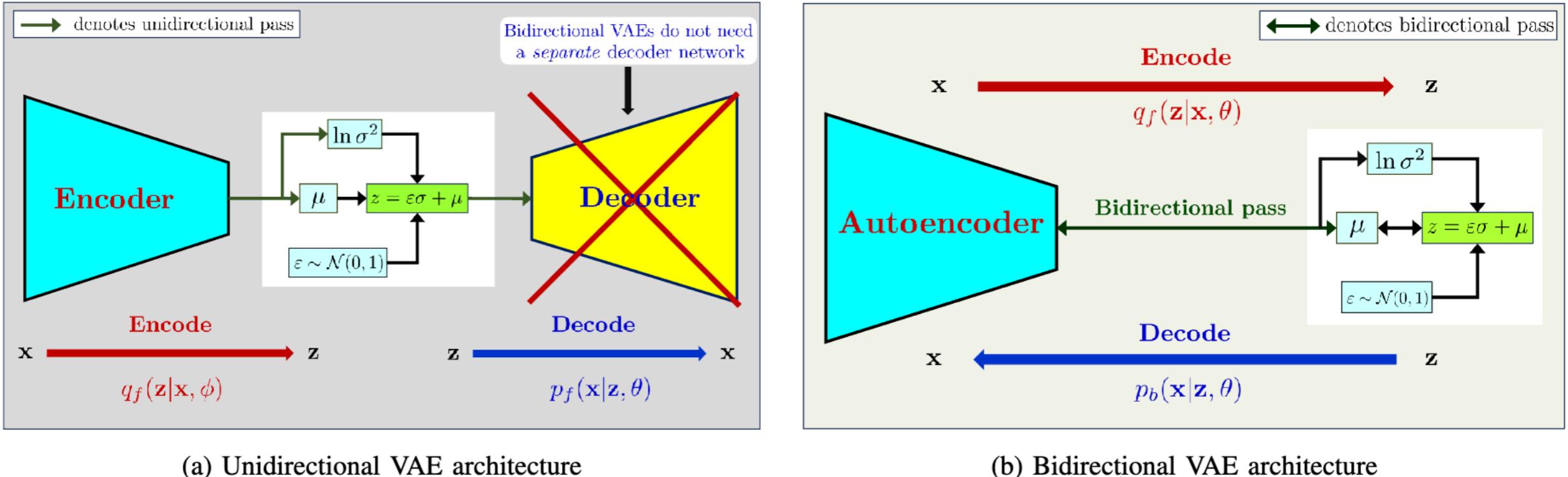


Reality

Jha et al. (CVPR, 2018)

Posterior equal to prior - non informative, useless as representation

VAEs in complex datasets

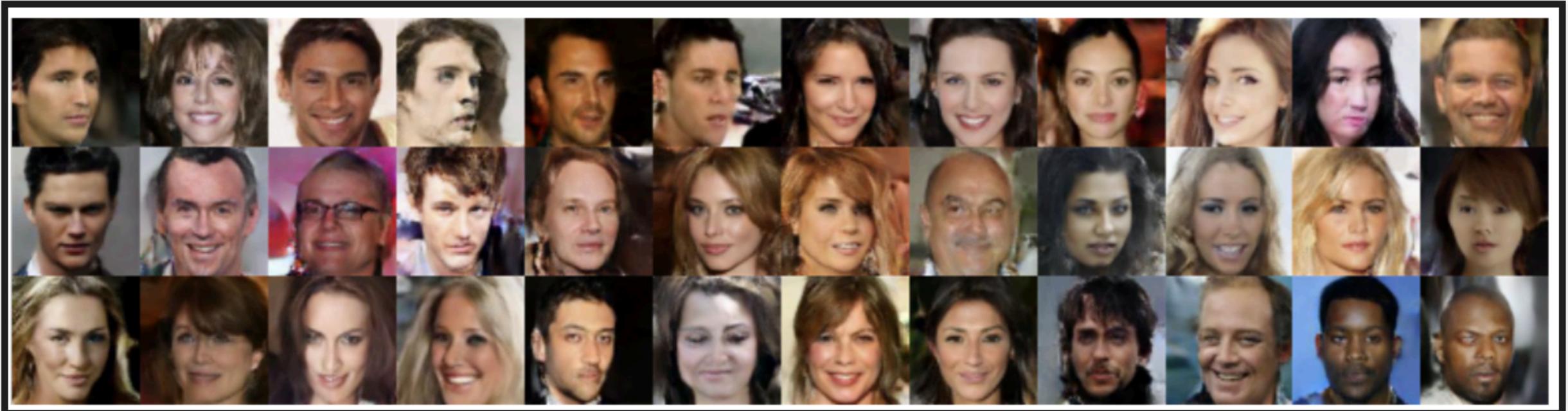


(a) Unidirectional VAE architecture

(b) Bidirectional VAE architecture

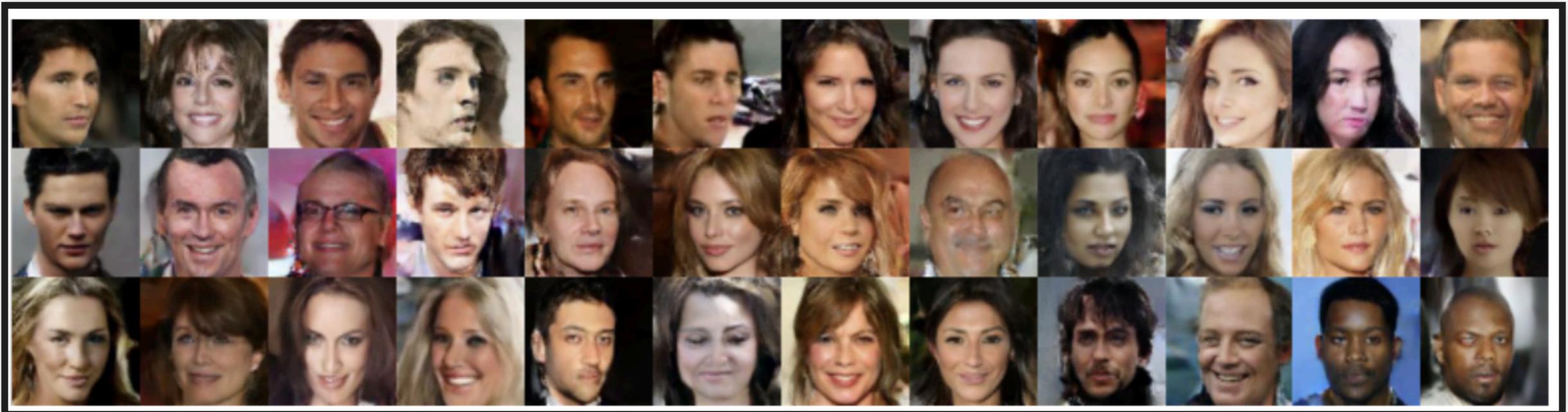
Fig. 1: Bidirectional vs. unidirectional variational autoencoders: Unidirectional VAEs use the forward passes of two separate networks for encoding and decoding. Bidirectional VAEs run their encoding on the forward pass and decoding on the backward pass with the same synaptic webs–weight matrices in both directions. This cuts the number of tunable parameters roughly in half. (a) The decoder network with parameter θ approximates $p(x|z, \theta)$ and the encoder network with parameter ϕ approximates $q(z|x, \theta)$. (b) Bidirectional VAEs use the forward pass of a network with parameter θ to approximate $q(z|x, \theta)$ and the backward pass of the network to approximate $p(x|z, \theta)$.

VAEs in complex datasets

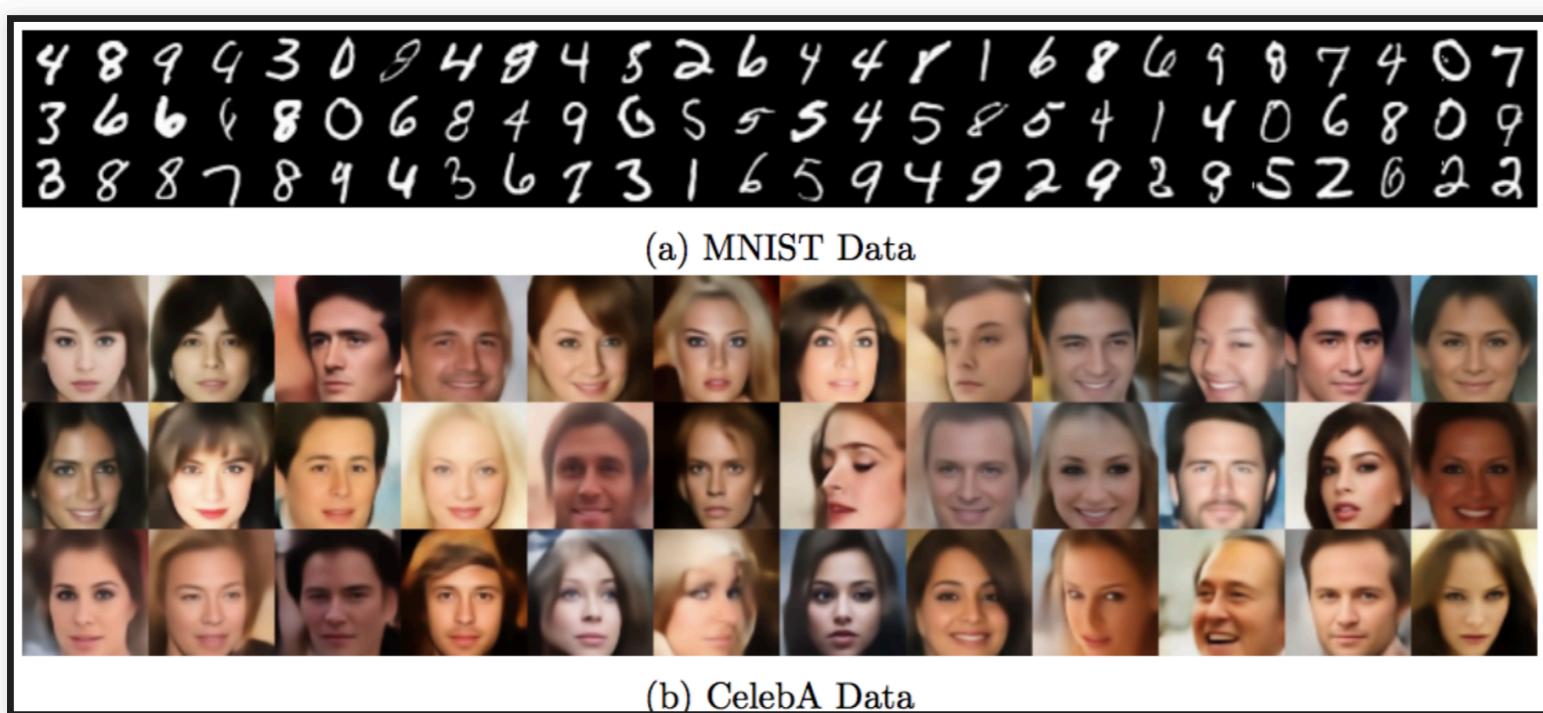


Bidirectional Inference VAE (Maaloe et al. 2019)

VAEs in complex datasets



Bidirectional Inference VAE (Maaloe et al. 2019)



Two-stage VAE (Dai & Wipf 2019)