# Variational autoencoders



$q(\boldsymbol{z}|\boldsymbol{x})$

$\mu_{\boldsymbol{z}}$

$\mu_{\boldsymbol{z}}$

$\sigma_{\boldsymbol{z}}$

$\sigma_{\boldsymbol{z}}$

$\boldsymbol{x}$

$\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{z}; \mu_{\boldsymbol{z}}, \sigma_{\boldsymbol{z}})$

$p(\boldsymbol{x}|\boldsymbol{z})$

$\boldsymbol{z}$

$\boldsymbol{z}$

$\boldsymbol{x}$
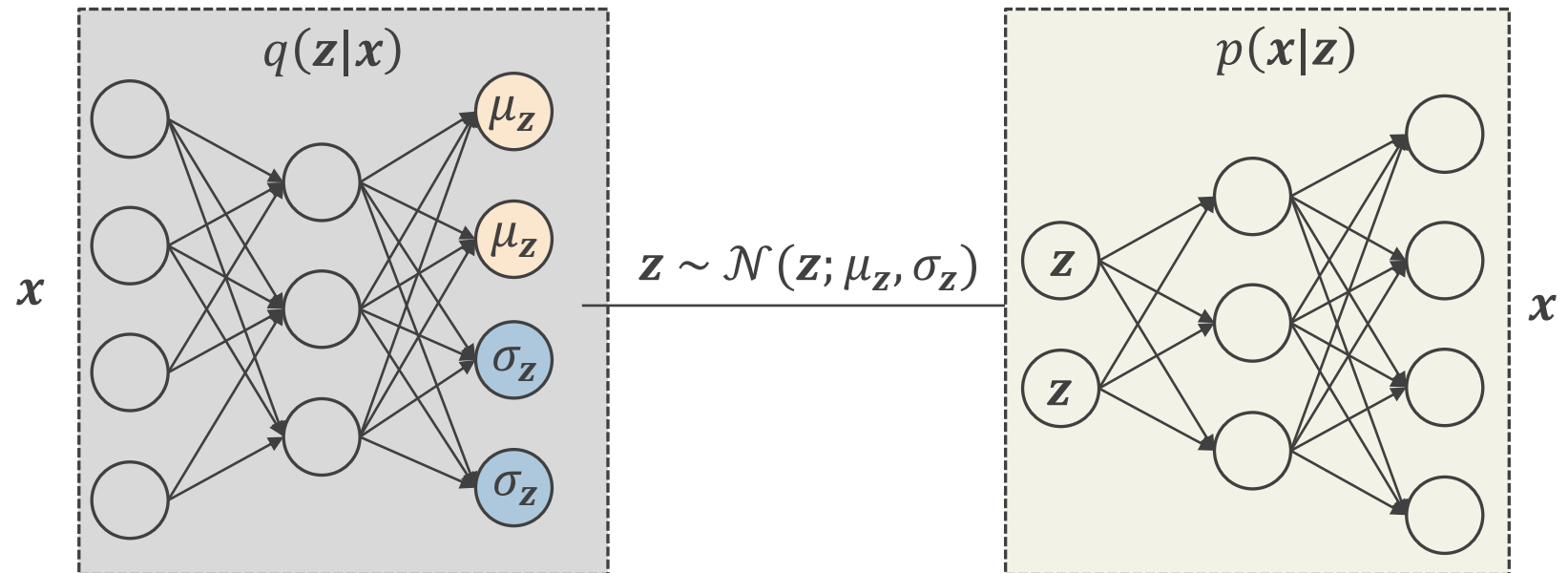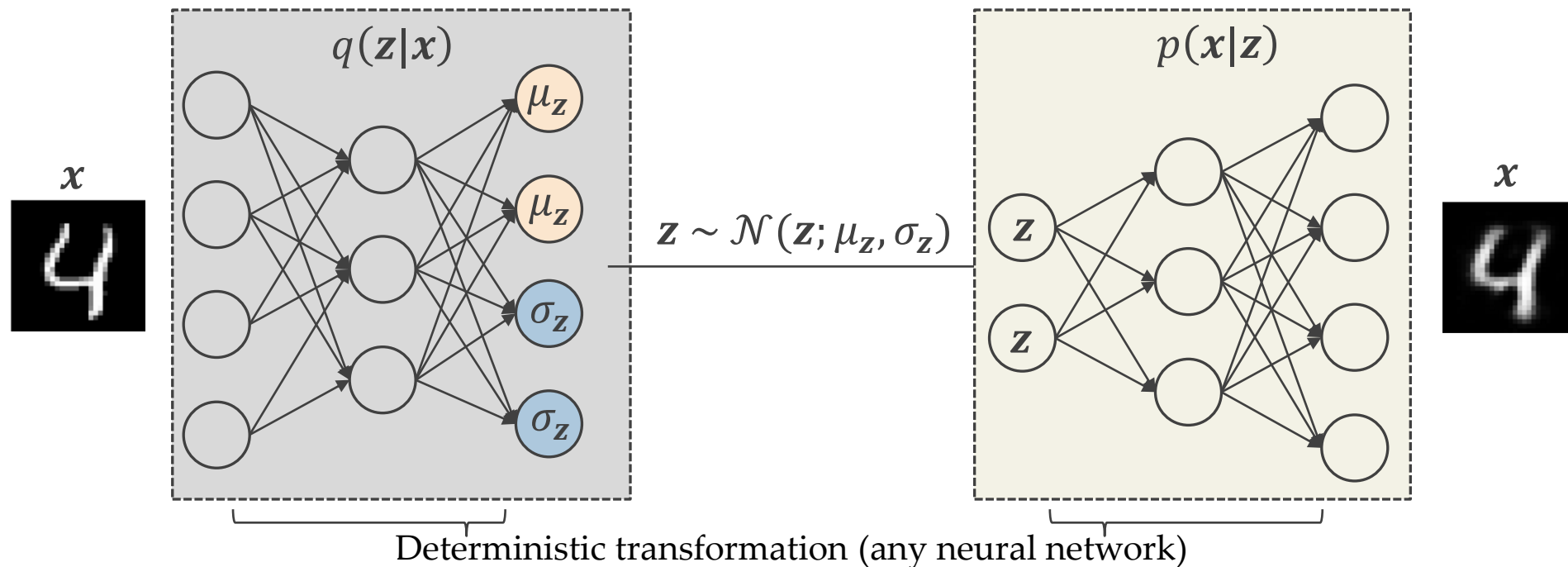
# Variational autoencoders

o Variational autoencoders is the neural network implementation of the ELBO

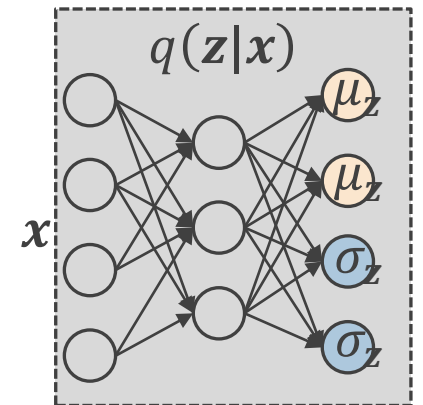$$\text{ELBO} = \mathbb{E}_{q_{\boldsymbol{\varphi}}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})] - \text{KL}\big[q_{\boldsymbol{\varphi}}(\boldsymbol{z}|\boldsymbol{x}) \parallel p(\boldsymbol{z})\big]$$

o In the standard case the approximate posterior is Gaussian $q(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}; \mu_{\boldsymbol{z}}, \sigma_{\boldsymbol{z}})$



Deterministic transformation (any neural network)

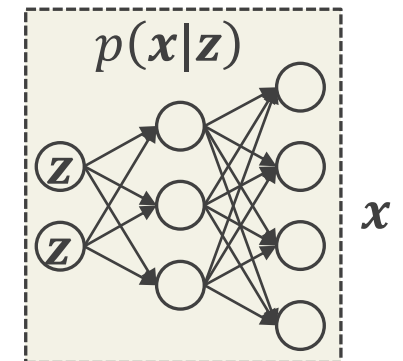# Encoder/inference network ⟺ approximate posterior

o The encoder is any standard neural network
  ◦ Modelling the approximate posterior $q(\boldsymbol{z}|\boldsymbol{x})$
  ◦ Remember: given input $\boldsymbol{x}$ we have a distribution over latent $\boldsymbol{z}$ (not single value)
  ◦ The KL term $\text{KL}[q(\boldsymbol{z}|\boldsymbol{x}) \parallel p(\boldsymbol{z})]$ encourages the posterior to not deviate too much from the prior $p(\boldsymbol{z})$

o For Gaussian $q(\boldsymbol{z}|\boldsymbol{x})$ we need two neural networks for two outputs $\mu_{\boldsymbol{z}}, \sigma_{\boldsymbol{z}}$
  ◦ The $\mu_{\boldsymbol{z}}$ is a neural net encoding the mean of $\boldsymbol{z}$ given $\boldsymbol{x}$
  ◦ The $\sigma_{\boldsymbol{z}}$ is a neural net encoding the stdev of $\boldsymbol{z}$ given $\boldsymbol{x}$
  ◦ The two neural nets can share architecture before the outputs

# Decoder network ⇔ generative model

o The decoder model is also a neural network
  ◦ It receives a stochastic input **z** and returns as output a generation

o The output modelled with a distribution according to the data type
  ◦ For continuous values could be a Gaussian
  ◦ For binary values Bernoulli distribution

o With generative models often convenient to think of the generation process
  ◦ Then the encoder is the variational approximation to ensure tractability

o Check the graphical model
  ◦ Sample $z \sim p(z)$ from the prior
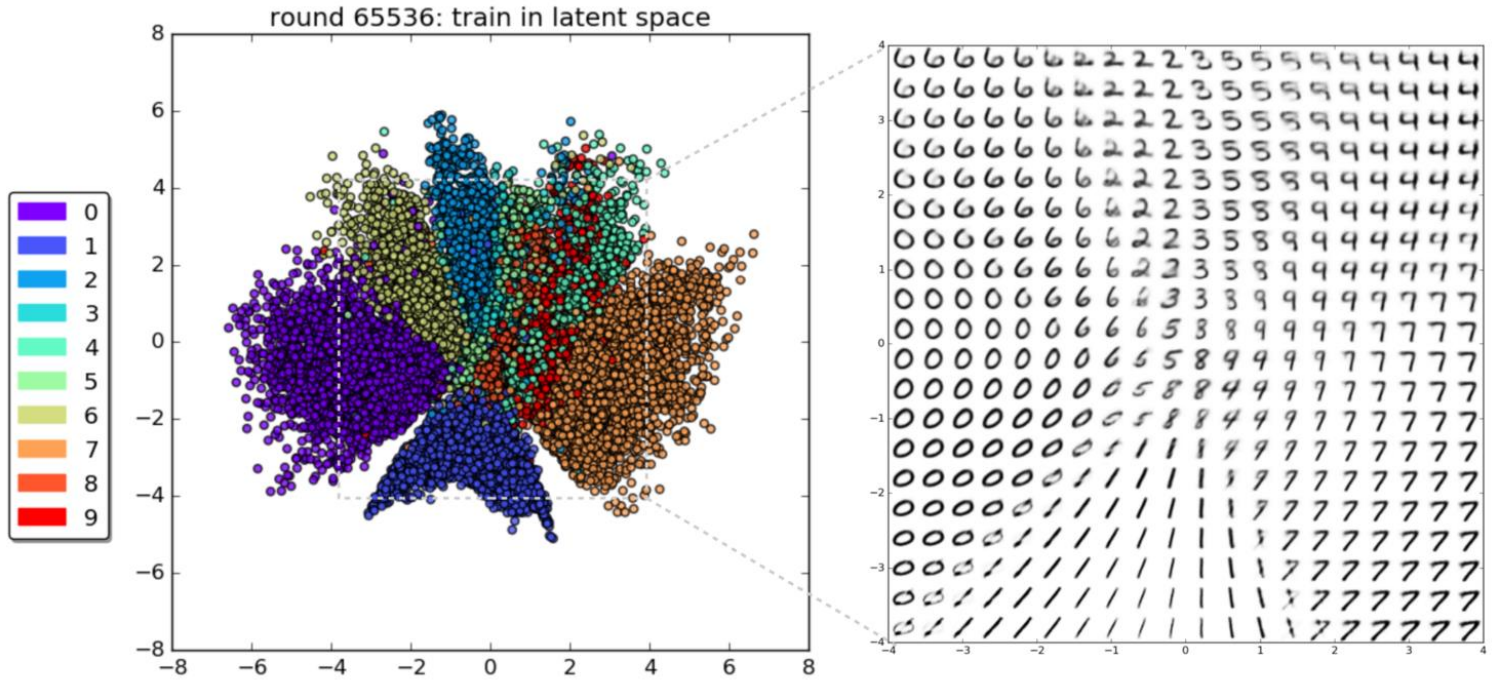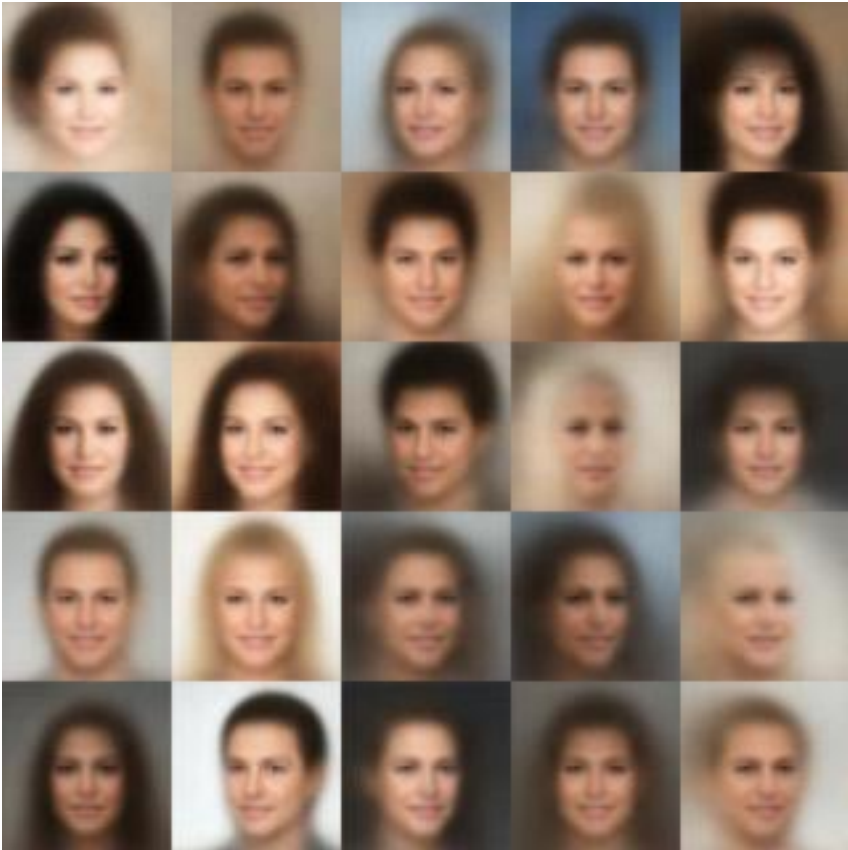  ◦ Given $z$ generate $x \sim p(x|z)$

# Prior distribution

o The prior distribution acts as a regularizer

o The prior $p(\mathbf{z})$ is often the unit Gaussian $p(\mathbf{z}) \sim \mathcal{N}(0, 1)$

o If we expect/desire different nature of $\mathbf{z}$, *e.g.,* sparsity or binary latents
  ◦ → pick a different prior
  ◦ The sampled $\mathbf{z}$ will be from that prior
  ◦ The KL term will regularize the encoder to be close to the prior

# Learning the variational autoencoders

- The variational autoencoder is two neural networks with inputs or outputs that are stochastic (represented by distributions, not single values)

- We must train the neural networks
  - I.e., fit good parameters $\boldsymbol{\theta}$ and $\boldsymbol{\varphi}$ for the decoder

- Objectives:
  - We want to predict to good distributions for $\boldsymbol{z}$ for (seen & unseen) inputs $\boldsymbol{x}$
  - We want on average our approximate posterior to be close to the prior
  - We want to reconstruct inputs well
  - We want generations that look 'real' $\rightarrow$ good extrapolations

# Interpolation in the latent VAE space

# Training Variational Autoencoders

o Maximize the ELBO

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|z)] - \mathrm{KL}(q_\varphi(z|x)||p(z))$$

$$= \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}) \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z} - \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}) \log \frac{q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \, d\mathbf{z}$$

o Normally you derive the math between each integral
  ◦ Good exercise: derive the ELBO for Gaussian latents and Bernoulli outputs

o Often, the integrals make some terms intractable. How to train?
  ◦ Backpropagation with Monte Carlo (MC) averaging
  ◦ Forward propagation means evaluating the two terms
  ◦ Backpropagation → compute gradients with respect to the $\boldsymbol{\theta}$ and $\boldsymbol{\varphi}$

# Training reconstruction term

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \int_{\boldsymbol{z}} q_{\boldsymbol{\varphi}}(\boldsymbol{z}|\boldsymbol{x}) \log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z}) \, d\boldsymbol{z} - \int_{\boldsymbol{z}} q_{\boldsymbol{\varphi}}(\boldsymbol{z}|\boldsymbol{x}) \log \frac{q_{\boldsymbol{\varphi}}(\boldsymbol{z}|\boldsymbol{x})}{p(\boldsymbol{z})} \, d\boldsymbol{z}$$

o The first term is an integral (expectation) that we cannot solve analytically
  ◦ Sample from the approximate posterior $q_{\varphi}(\boldsymbol{z}|\boldsymbol{x})$ instead and do MC average
  ◦ Pick a $p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})$ that couples well with log

o With a 'low variance estimator' a single sample $\boldsymbol{z}$ is enough
  ◦ Stochasticity is desirable → reduces overfitting

o Reparameterization trick for low variance estimation

# Training KL regularization term

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}) \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z} - \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x}) \log \frac{q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \, d\mathbf{z}$$

o The second term is an integral which corresponds to KL distance

o For known distributions, *e.g.*, both $q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ Gaussians, the KL often reduces to a closed formula → very convenient
  ◦ E.g., compute the KL divergence between a centered $N(0, 1)$ and a non-centered $N(\mu, \sigma)$ gaussian

o If closed formula not easy, MC averaging with sampling from $q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})$ is possible