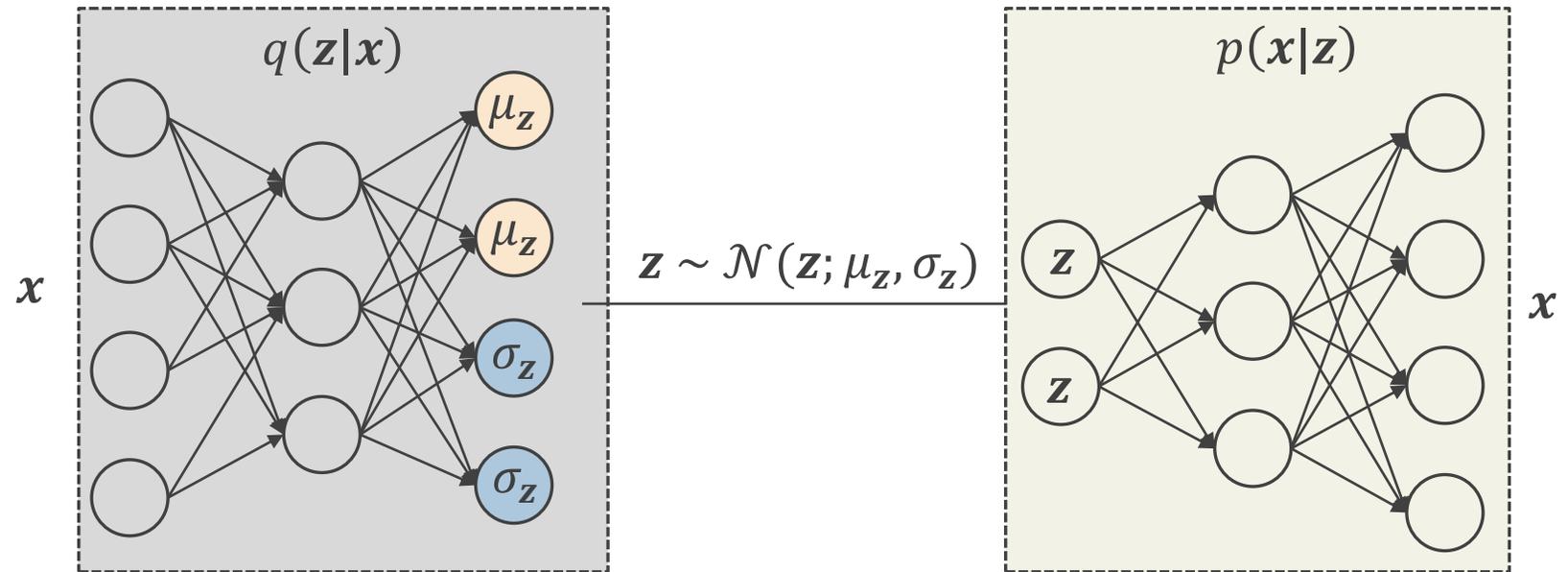
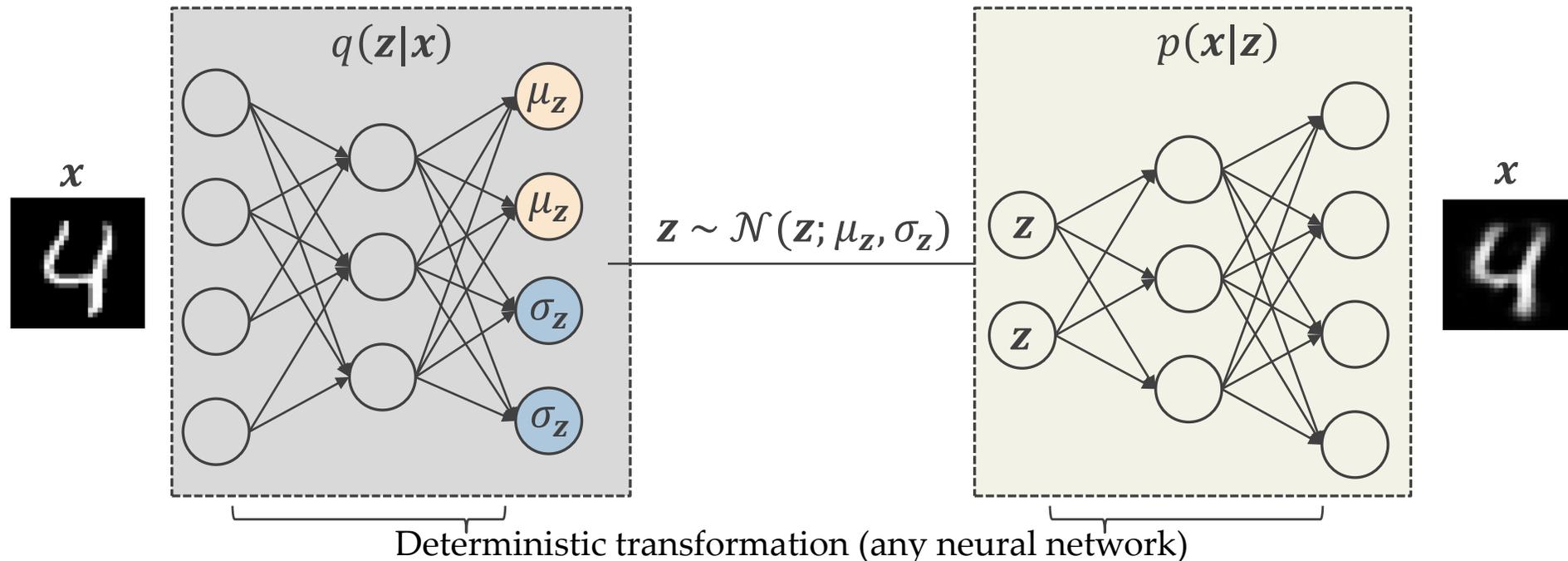


Variational autoencoders



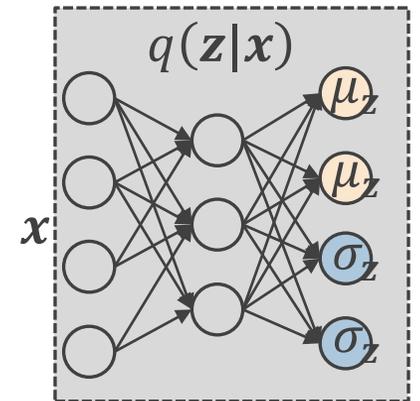
Variational autoencoders

- Variational autoencoders is the neural network implementation of the ELBO
$$\text{ELBO} = \mathbb{E}_{q_{\varphi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}[q_{\varphi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})]$$
- In the standard case the approximate posterior is Gaussian $q(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\mathbf{z}}, \sigma_{\mathbf{z}})$



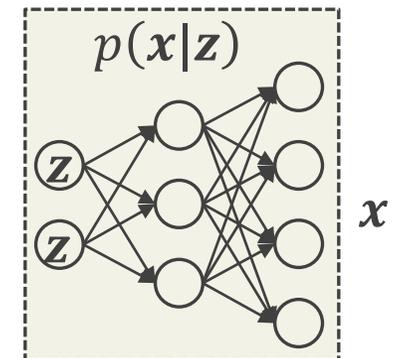
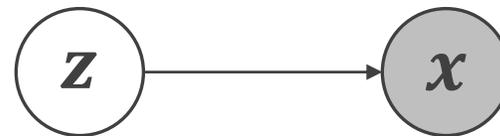
Encoder/inference network \Leftrightarrow approximate posterior

- The encoder is any standard neural network
 - Modelling the approximate posterior $q(\mathbf{z}|\mathbf{x})$
 - Remember: given input \mathbf{x} we have a distribution over latent \mathbf{z} (not single value)
 - The KL term $\text{KL}[q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})]$ encourages the posterior to not deviate too much from the prior $p(\mathbf{z})$
- For Gaussian $q(\mathbf{z}|\mathbf{x})$ we need two neural networks for two outputs $\mu_{\mathbf{z}}, \sigma_{\mathbf{z}}$
 - The $\mu_{\mathbf{z}}$ is a neural net encoding the mean of \mathbf{z} given \mathbf{x}
 - The $\sigma_{\mathbf{z}}$ is a neural net encoding the stdev of \mathbf{z} given \mathbf{x}
 - The two neural nets can share architecture before the outputs



Decoder network \Leftrightarrow generative model

- The decoder model is also a neural network
 - It receives a stochastic input \mathbf{z} and returns as output a generation
- The output modelled with a distribution according to the data type
 - For continuous values could be a Gaussian
 - For binary values Bernoulli distribution
- With generative models often convenient to think of the generation process
 - Then the encoder is the variational approximation to ensure tractability
- Check the graphical model
 - Sample $\mathbf{z} \sim p(\mathbf{z})$ from the prior
 - Given \mathbf{z} generate $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})$



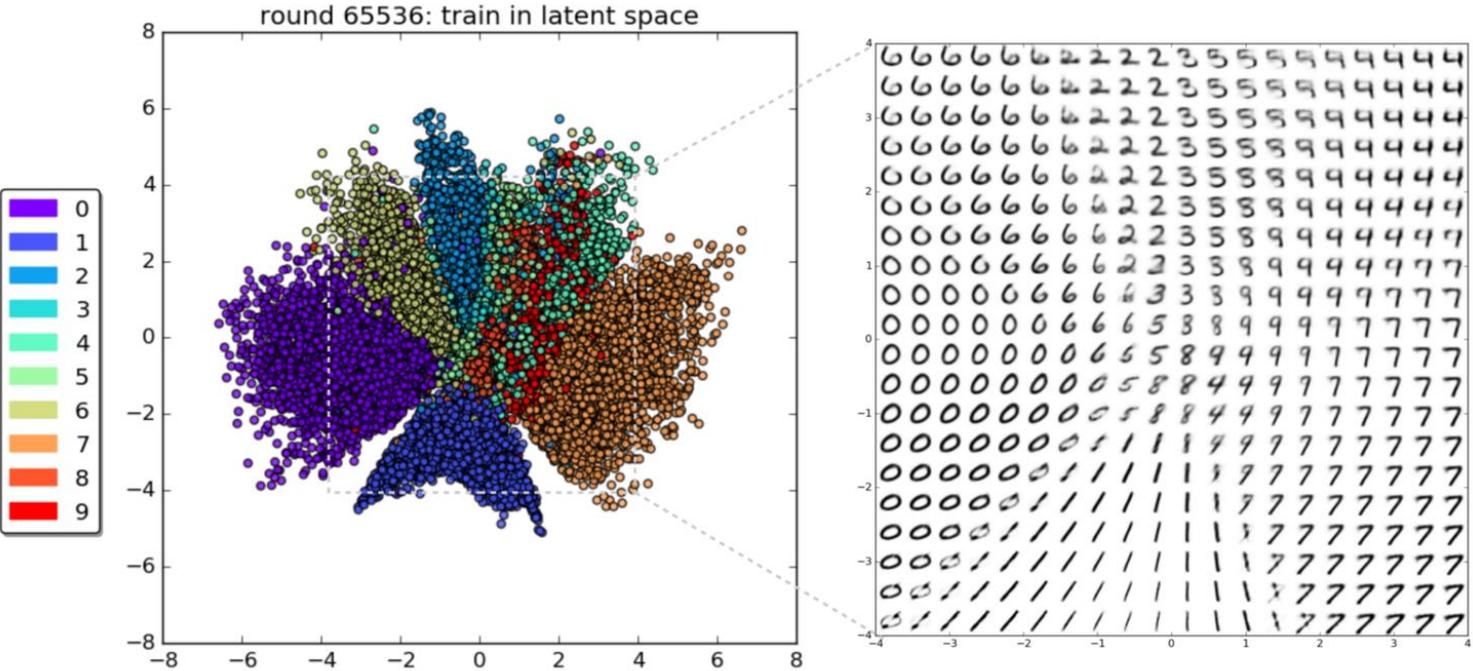
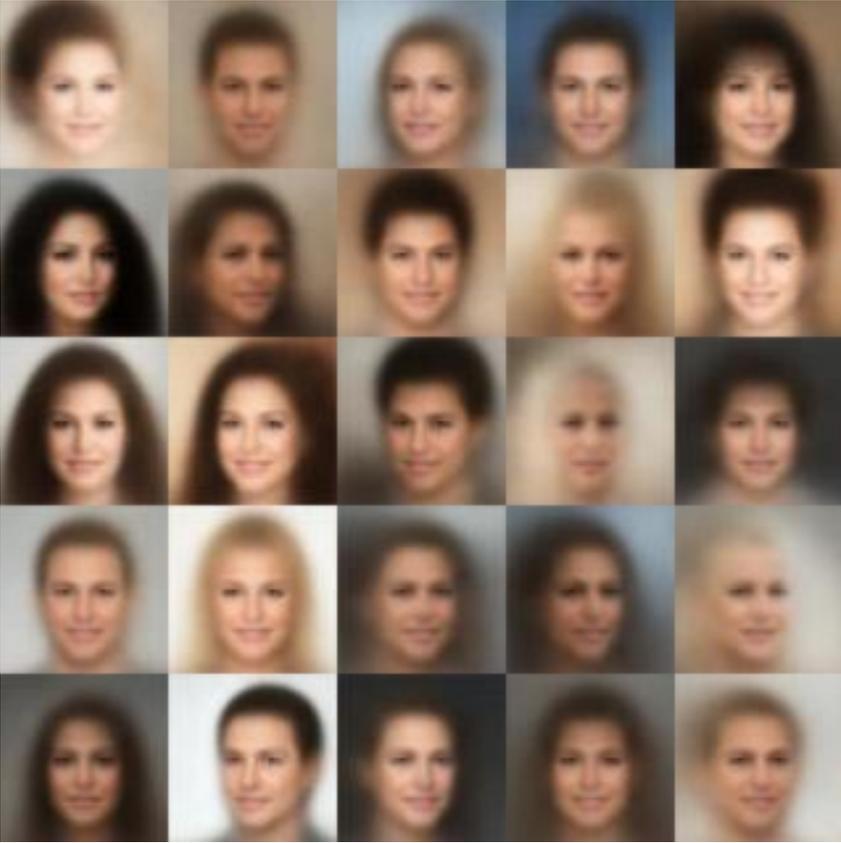
Prior distribution

- The prior distribution acts as a regularizer
- The prior $p(\mathbf{z})$ is often the unit Gaussian $p(\mathbf{z}) \sim \mathcal{N}(0, 1)$
- 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 θ and φ for the decoder
- Objectives:
 - We want to predict to good distributions for \mathbf{z} for (seen & unseen) inputs \mathbf{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' → good extrapolations

Interpolation in the latent VAE space



Training Variational Autoencoders

- Maximize the ELBO

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\varphi}) &= \mathbb{E}_{q_{\boldsymbol{\varphi}}(Z|x)}[\log p_{\boldsymbol{\theta}}(x|Z)] - \text{KL}(q_{\boldsymbol{\varphi}}(Z|x) || p_{\lambda}(Z)) \\ &= \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|x) \log p_{\boldsymbol{\theta}}(x|\mathbf{z}) d\mathbf{z} - \int_{\mathbf{z}} q_{\boldsymbol{\varphi}}(\mathbf{z}|x) \log \frac{q_{\boldsymbol{\varphi}}(\mathbf{z}|x)}{p(\mathbf{z})} d\mathbf{z}\end{aligned}$$

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

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

- The first term is an integral (expectation) that we cannot solve analytically
 - Sample from the approximate posterior $q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})$ instead and do MC average
 - Pick a $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ that couples well with log
- With a ‘low variance estimator’ a single sample \mathbf{z} is enough
 - Stochasticity is desirable → reduces overfitting
- Reparameterization trick for low variance estimation

Training KL regularization term

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

- The second term is an integral which corresponds to KL distance
- 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 \rightarrow very convenient
 - *E.g.*, compute the KL divergence between a centered $N(0, 1)$ and a non-centered $N(\mu, \sigma)$ gaussian
- If closed formula not easy, MC averaging with sampling from $q_{\boldsymbol{\varphi}}(\mathbf{z}|\mathbf{x})$ is possible