

# DSC190: Machine Learning with Few Labels

## Unsupervised Learning

**Zhiting Hu**

Lecture 18, November 8, 2024

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Outline

Unsupervised learning: Variational Auto-Encoders

Presentations

- **Bobby Zhu:** Visualizing Data using t-SNE
- **Vivian Zhao:** Exploiting Audio-Visual Features with Pretrained AV-HuBERT for Multi-Modal Dysarthric Speech Reconstruction
- **Abhinav Sanisetty:** VideoPoet: A Large Language Model for Zero-Shot Video Generation
- **Zhiqing Wang:** Divide and Conquer: Leveraging Intermediate Feature Representations for Quantized Training of Neural Networks
- **Feiyang Jiang:** LoRA: Low-Rank Adaptation of Large Language Models

# Recall: Black-box Variational Inference (BBVI)

- Probabilistic model:  $\mathbf{x}$  -- observed variables,  $\mathbf{z}$  -- latent variables
- Variational distribution  $q_\lambda(\mathbf{z}|\mathbf{x})$  with parameters  $\lambda$ , e.g.,
  - Gaussian mixture distribution:
    - “A Gaussian mixture model is a universal approximator of densities, in the sense that any smooth density can be approximated with any specific nonzero amount of error by a Gaussian mixture model with enough components.” (Deep Learning book, pp.65)
  - Deep neural networks

$$\mathcal{L}(\lambda) \triangleq \mathbb{E}_{q_\lambda(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]$$

- ELBO to be maximized:

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

- Want to compute the gradient w.r.t variational parameters  $\lambda$

## BBVI with the score gradient

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- ELBO:

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

- **Question:** what's the score gradient w.r.t.  $\lambda$  ?

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_q[\nabla_\lambda \log q(\mathbf{z}|\lambda)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\lambda))]$$

## BBVI with the score gradient

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- ELBO:

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

- **Question:** what's the score gradient w.r.t.  $\lambda$  ?

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_q[\nabla_\lambda \log q(\mathbf{z}|\lambda)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\lambda))]$$

- Compute noisy unbiased gradients of the ELBO with Monte Carlo samples from the variational distribution

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S \nabla_\lambda \log q(z_s|\lambda)(\log p(\mathbf{x}, z_s) - \log q(z_s|\lambda)),$$

where  $z_s \sim q(\mathbf{z}|\lambda)$ .

## BBVI with the reparameterization gradient

$$\mathcal{L} = \mathbb{E}_{q_{\lambda}(\mathbf{z})}[f_{\lambda}(\mathbf{z})]$$
$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}}f_{\lambda}(\mathbf{z}) \nabla_{\lambda}t(\epsilon, \lambda)]$$

- ELBO:

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

- **Question:** what's the reparameterization gradient w.r.t.  $\lambda$  ?

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ \mathbf{z} &= t(\epsilon, \lambda) \end{aligned} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \nabla_{\lambda}t(\epsilon, \lambda)]$$

# Variational Autoencoders (VAEs)

# Variational Auto-Encoders (VAEs)

VAEs are a combination of the following ideas:

- Variational Inference
  - ELBO
- Variational distribution parametrized as neural networks
- Reparameterization trick



# Variational Auto-Encoders (VAEs)

- Model  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ 
  - $p_{\theta}(\mathbf{x}|\mathbf{z})$ : a.k.a., generative model, generator, (probabilistic) decoder, ...
  - $p(\mathbf{z})$ : prior, e.g., Gaussian
- Assume variational distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$ 
  - E.g., a Gaussian distribution parameterized as **deep neural networks**
  - a.k.a, recognition model, inference network, (probabilistic) encoder, ...
- ELBO:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$

Reconstruction

Divergence from prior  
(KL divergence between two Gaussians has  
an analytic form)

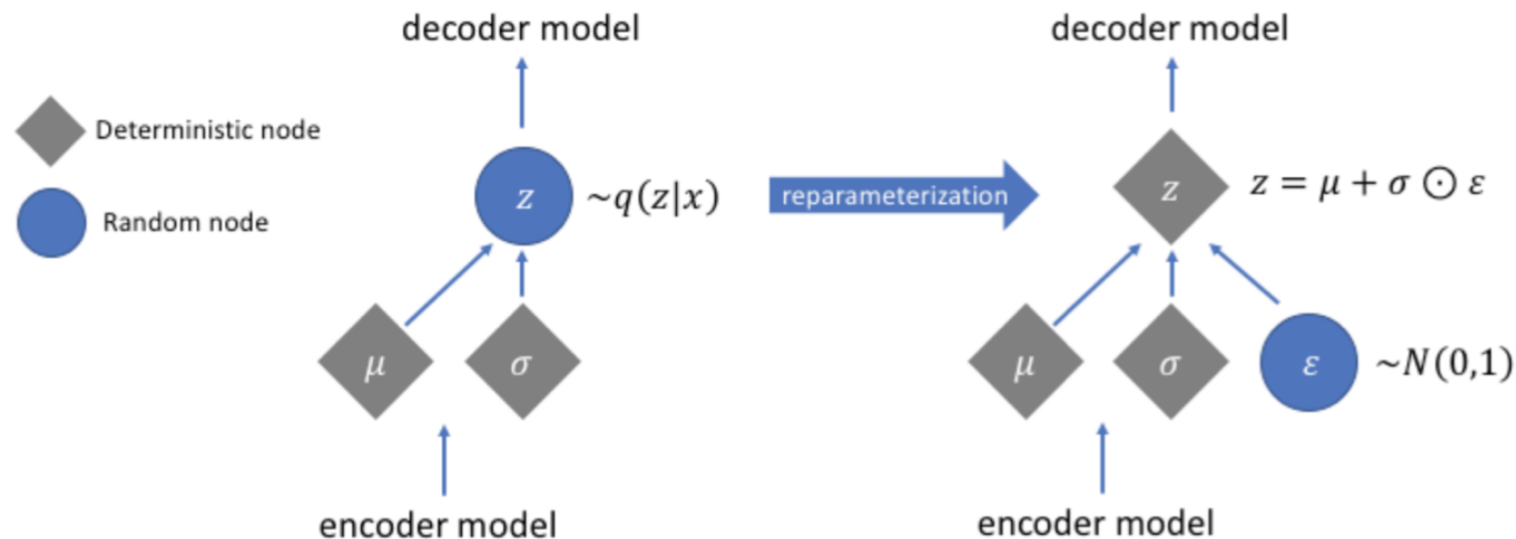
# Variational Auto-Encoders (VAEs)

- ELBO:

$$\begin{aligned}\mathcal{L}(\theta, \phi; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$

- Reparameterization:

- $[\mu; \sigma] = f_{\phi}(\mathbf{x})$  (a neural network)
- $\mathbf{z} = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(\mathbf{0}, \mathbf{1})$



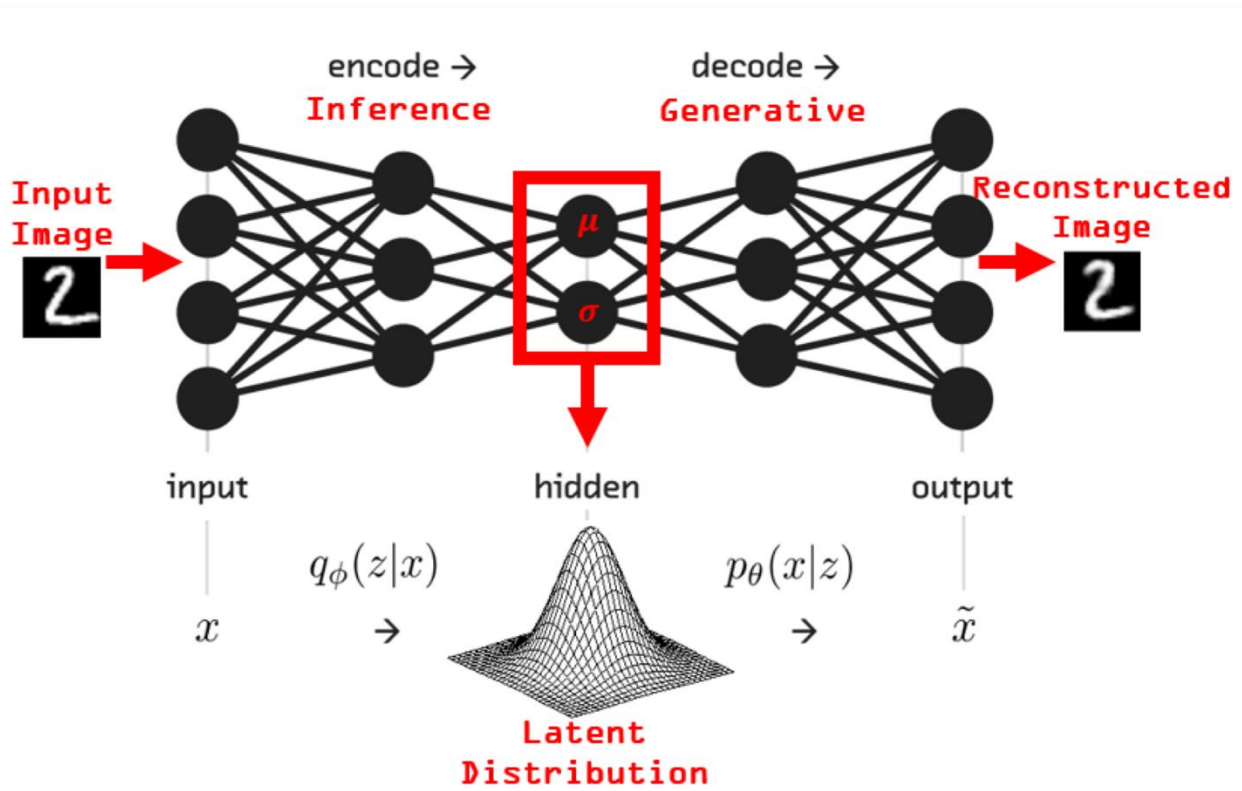
# Variational Auto-Encoders (VAEs)

- ELBO:
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] + H(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$
- Reparameterization:
  - $[\boldsymbol{\mu}; \boldsymbol{\sigma}] = f_{\boldsymbol{\phi}}(\mathbf{x})$  (a neural network)
  - $\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})$

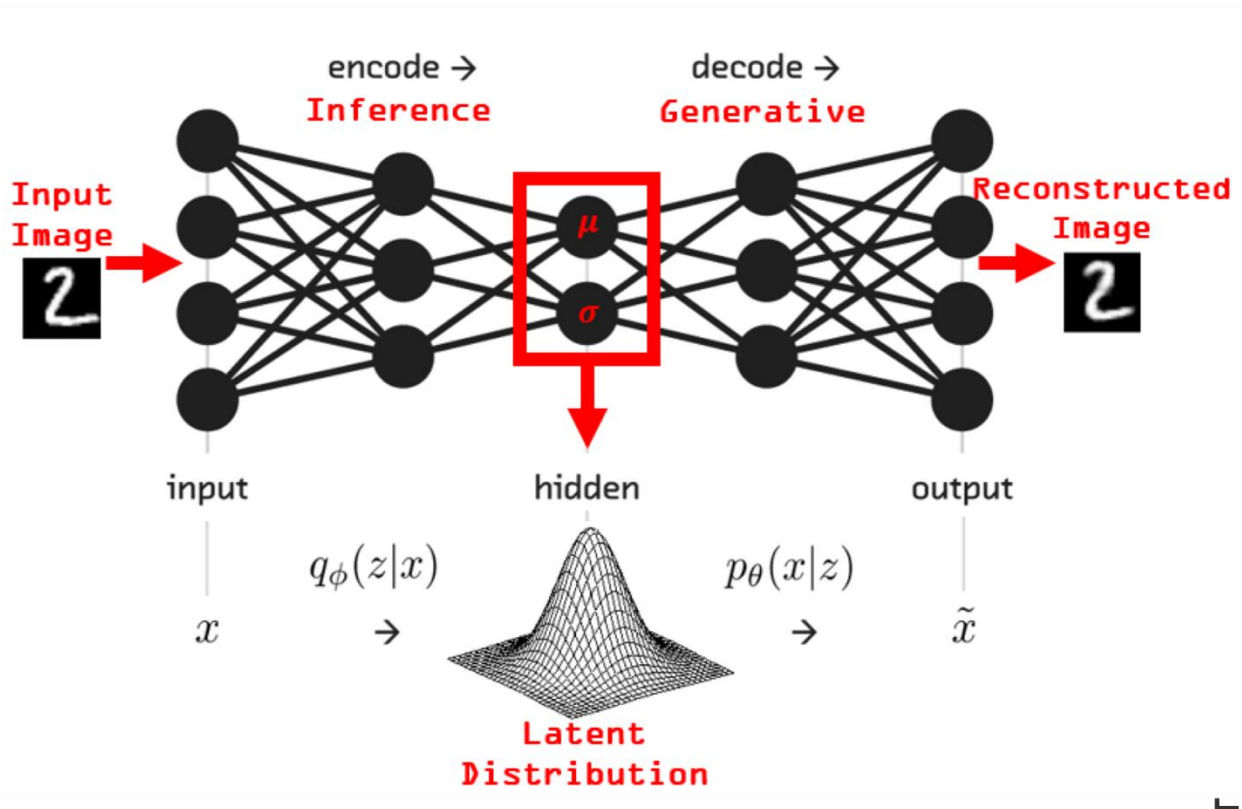
$$\nabla_{\boldsymbol{\phi}} \mathcal{L} = \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})} [\nabla_{\mathbf{z}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})] \nabla_{\boldsymbol{\phi}} \mathbf{z}(\boldsymbol{\epsilon}, \boldsymbol{\phi})]$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

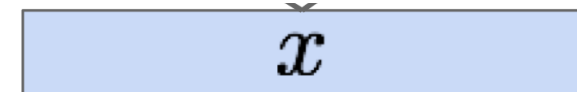
# Example: VAEs for images



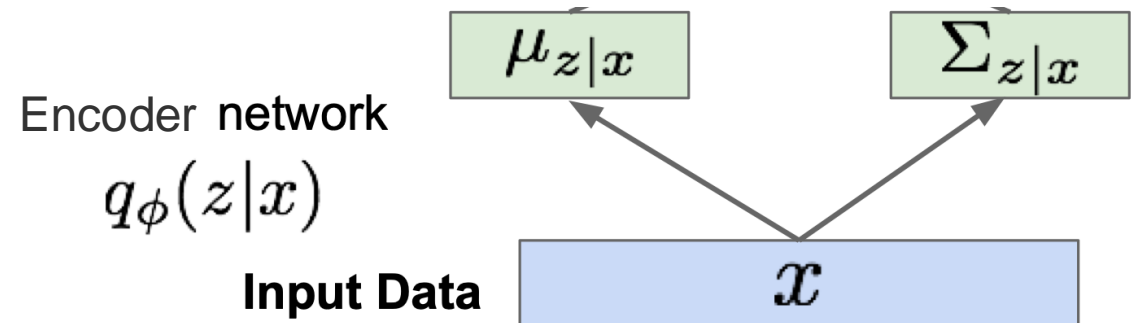
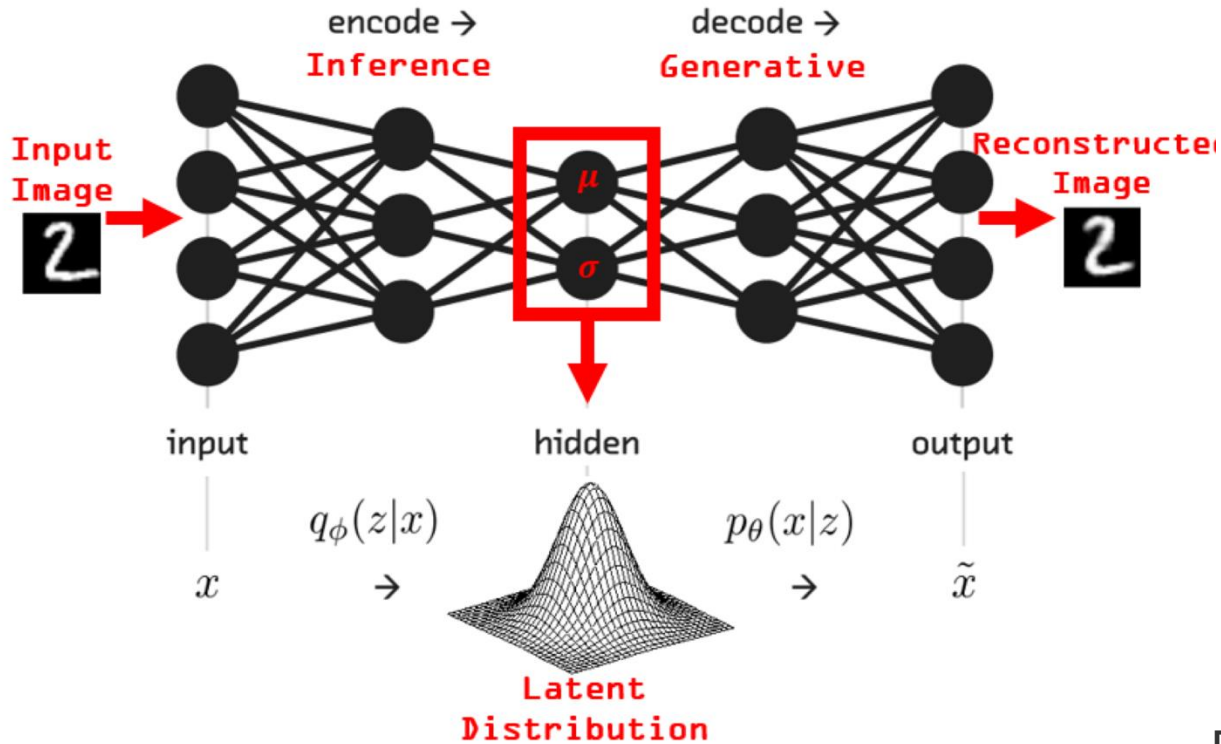
# Example: VAEs for images



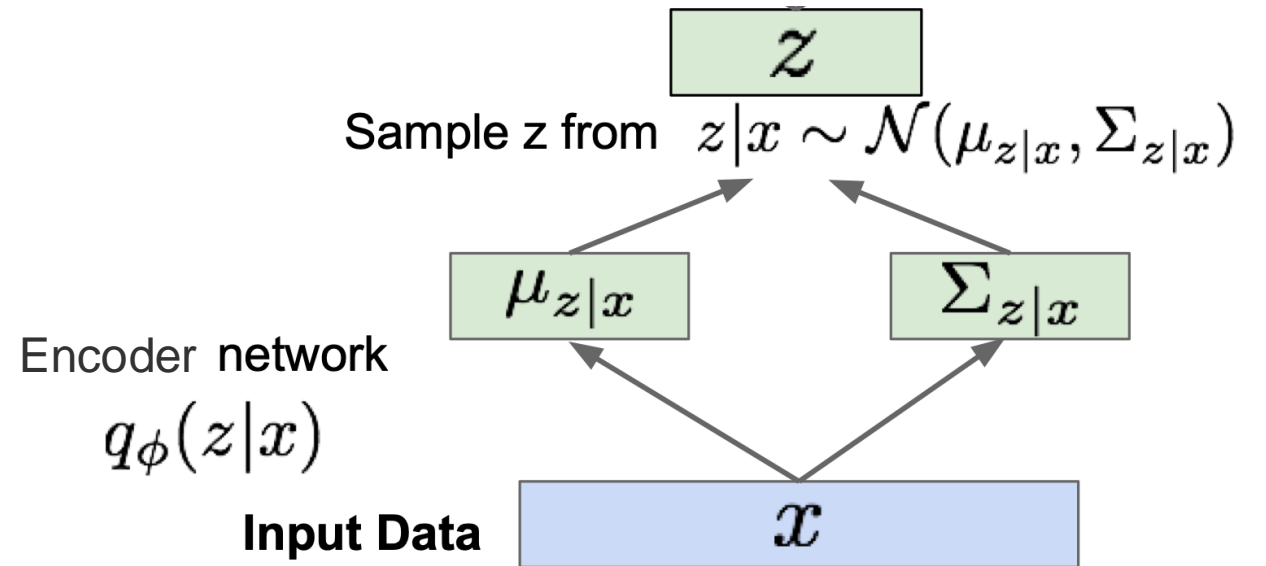
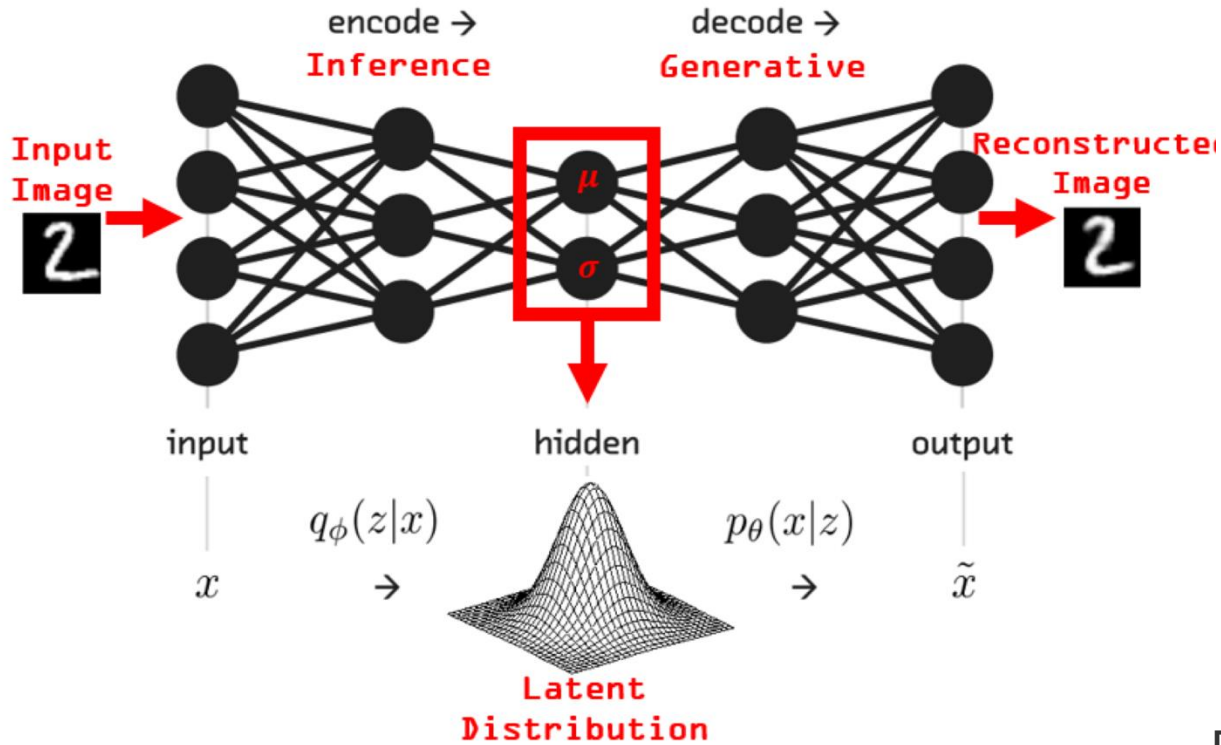
Input Data



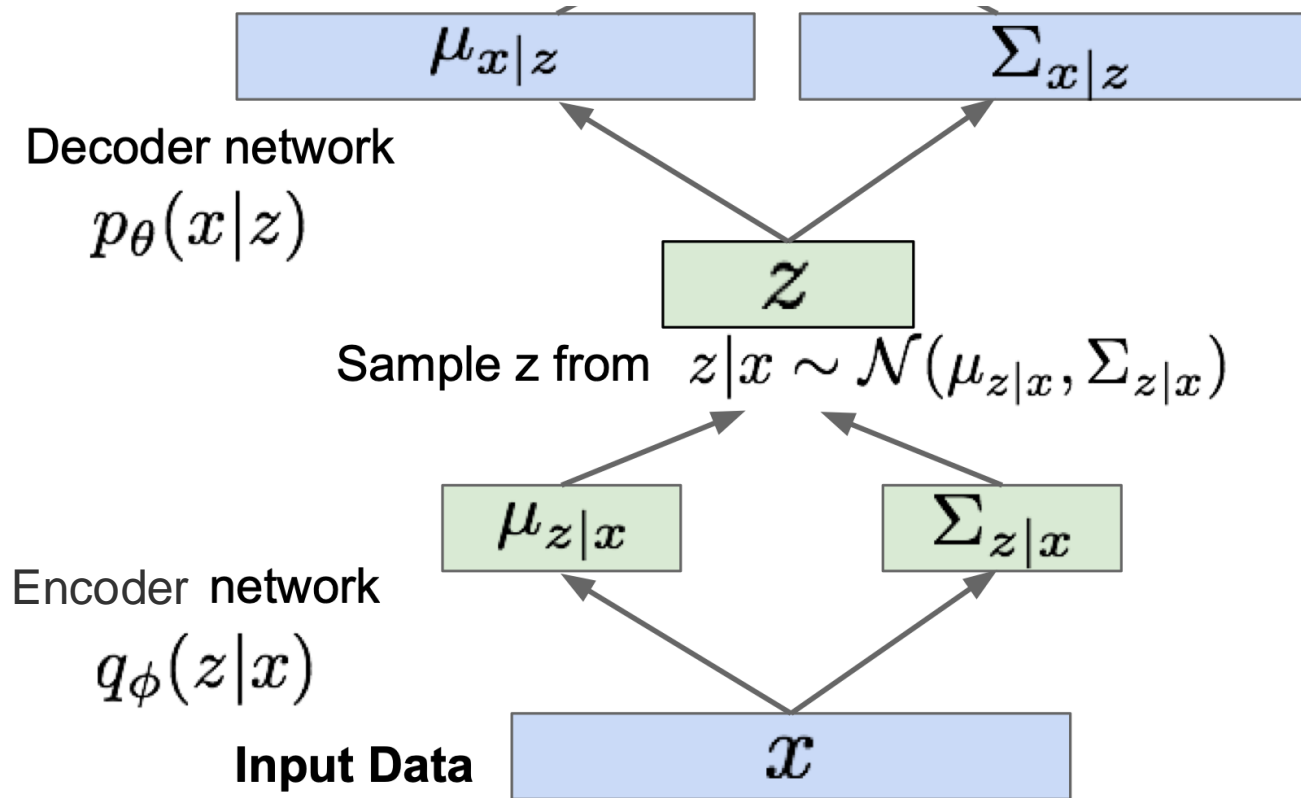
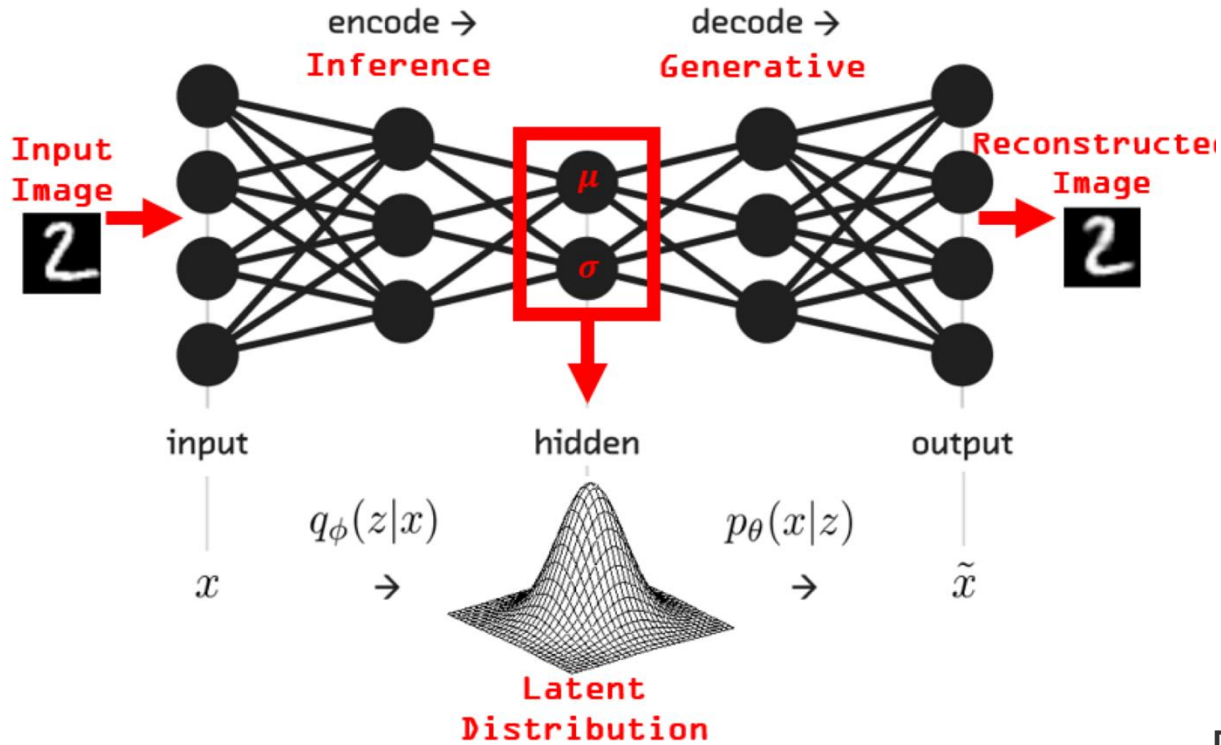
# Example: VAEs for images



# Example: VAEs for images

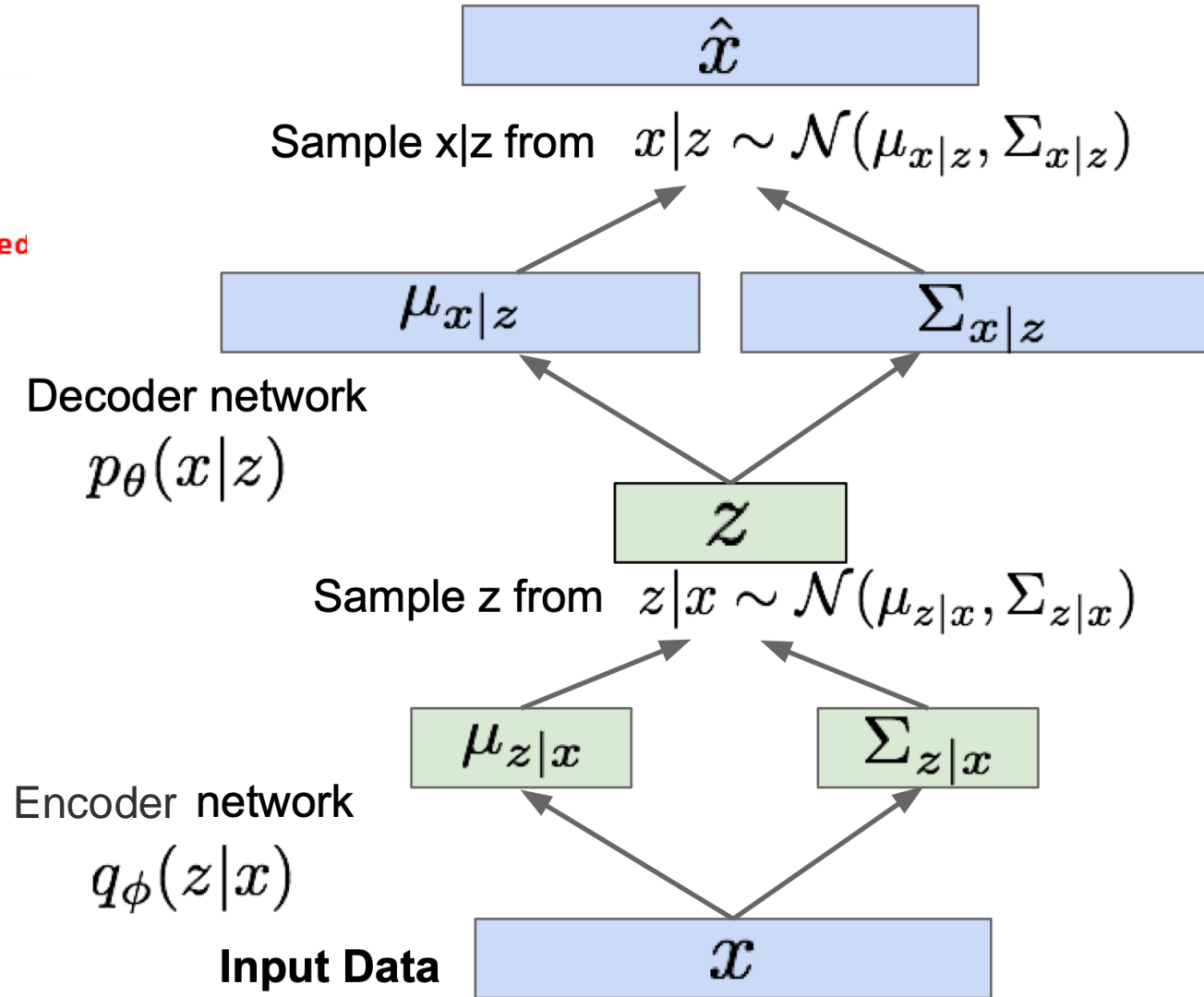
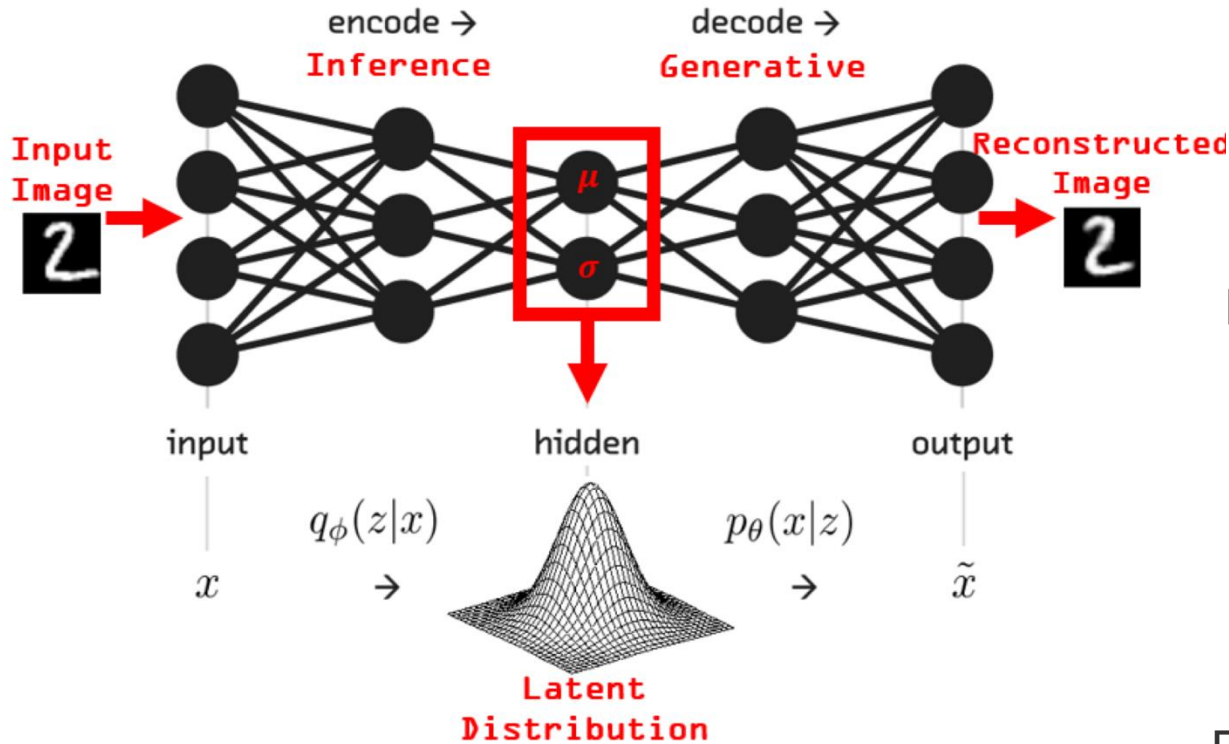


# Example: VAEs for images





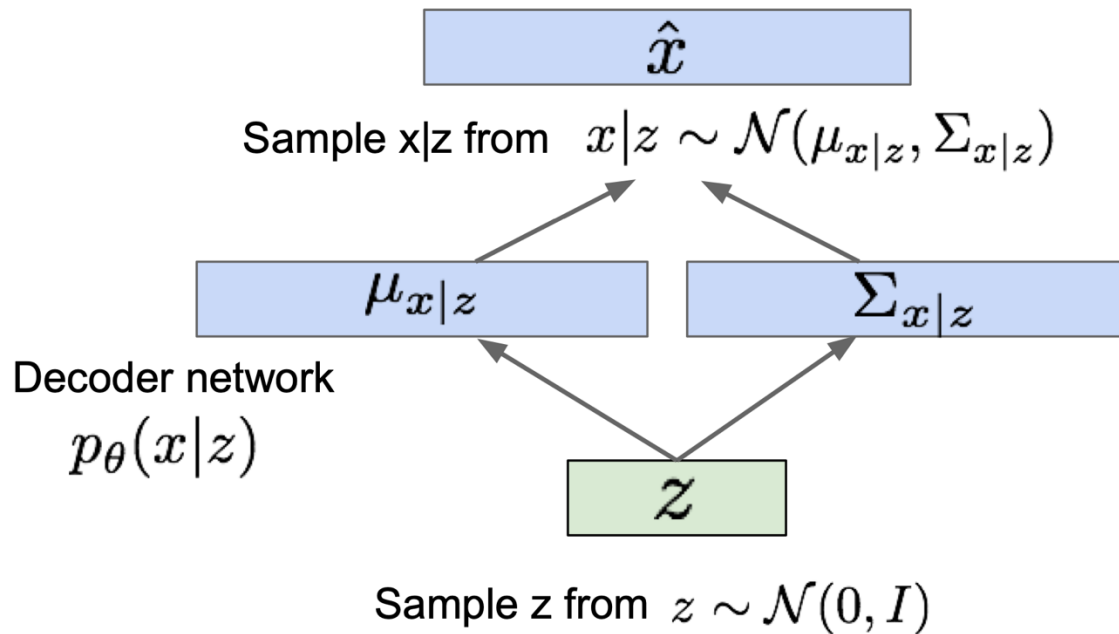
# Example: VAEs for images



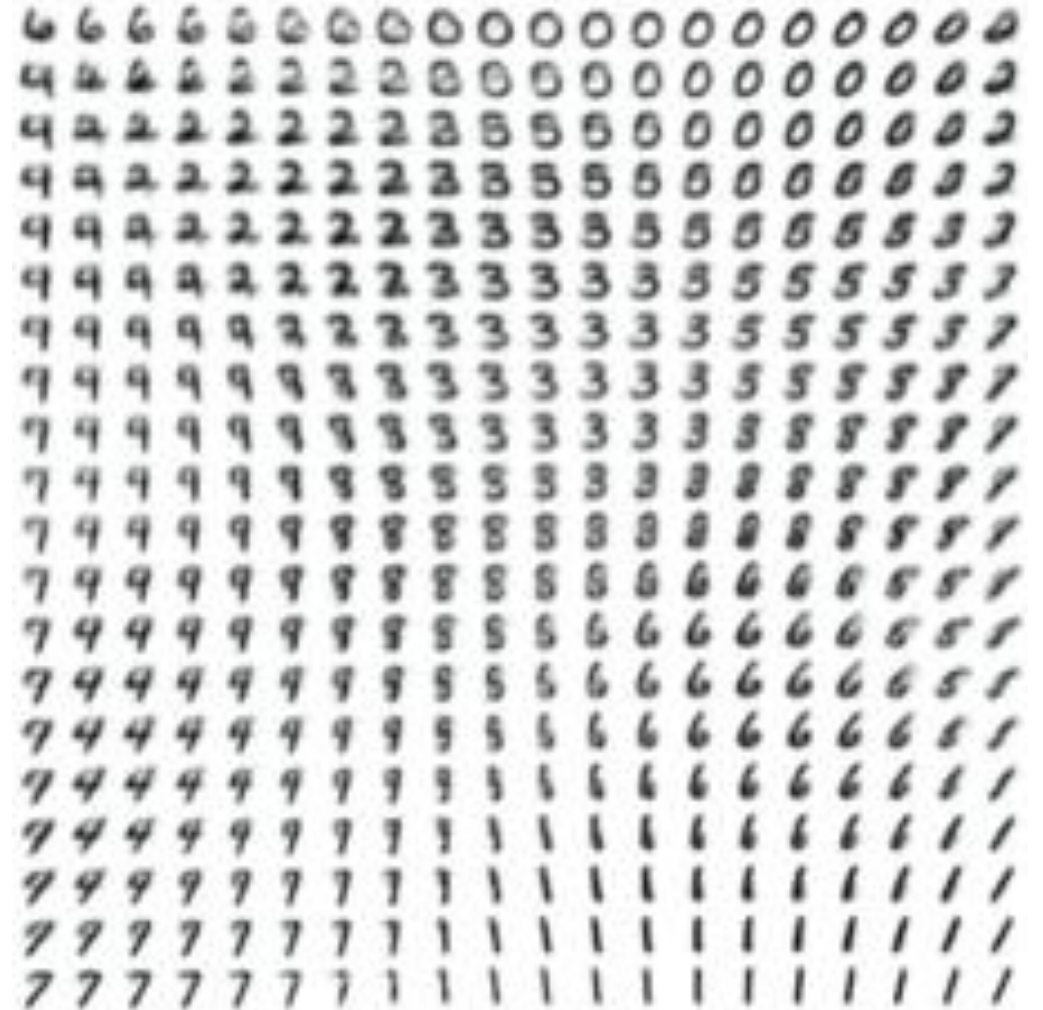
# Example: VAEs for images

Generating samples:

- Use decoder network. Now sample  $z$  from prior!



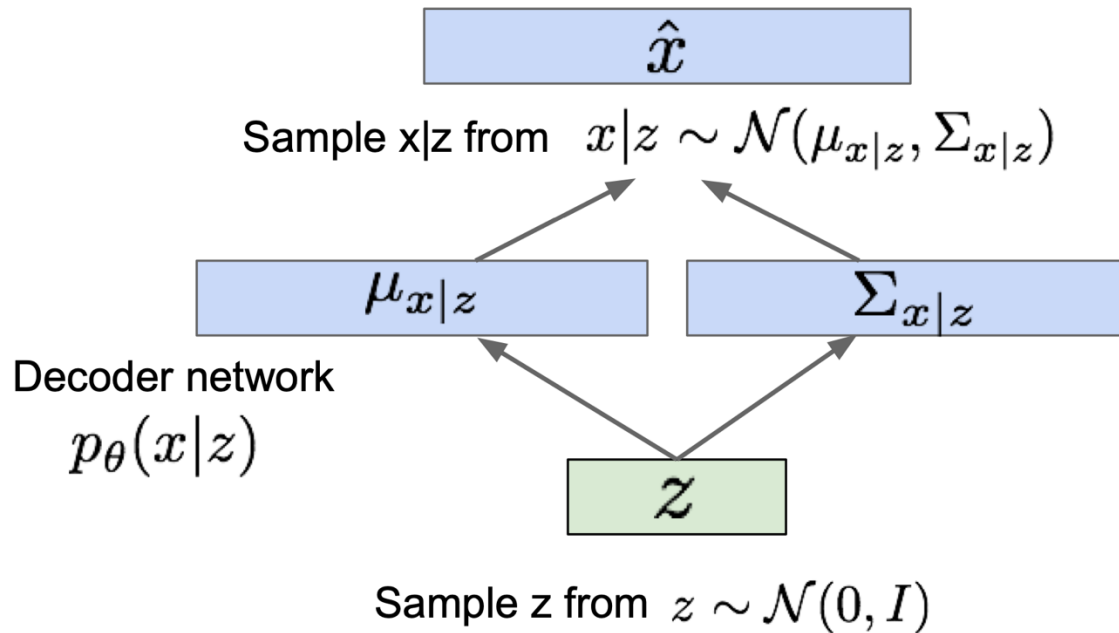
Data manifold for 2-d  $z$



# Example: VAEs for images

Generating samples:

- Use decoder network. Now sample  $z$  from prior!



Data manifold for 2-d  $z$



Vary  $z_1$   
(Degree of smile)

Vary  $z_2$  (head pose)

## Example: VAEs for text

- Latent code interpolation and sentences generation from VAEs [Bowman et al., 2015].

---

**“ i want to talk to you . ”**

*“i want to be with you . ”*

*“i do n’t want to be with you . ”*

*i do n’t want to be with you .*

**she did n’t want to be with him .**

---

## Note: Amortized Variational Inference

- Variational distribution as an **inference model**  $q_{\phi}(\mathbf{z}|\mathbf{x})$  with parameters  $\phi$  (which was traditionally factored over samples)
- Amortize the cost of inference by learning a **single** data-dependent inference model
- The trained inference model can be used for quick inference on new data

# Variational Auto-encoders: Summary

- A combination of the following ideas:
  - Variational Inference: ELBO
  - Variational distribution parametrized as neural networks
  - Reparameterization trick

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

← Reconstruction

↓ Divergence from prior

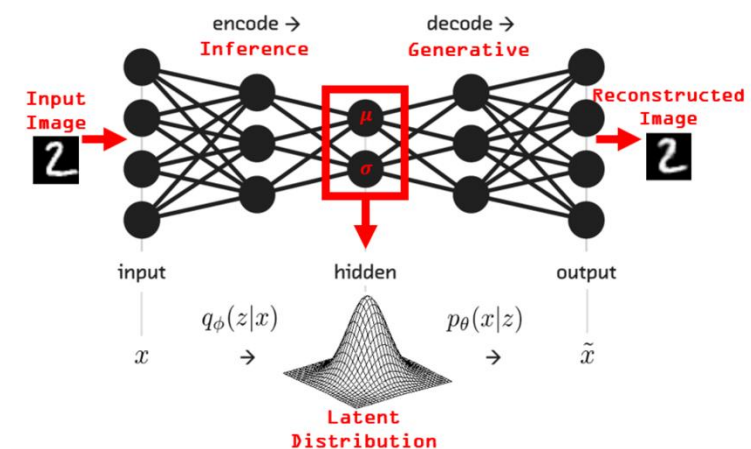
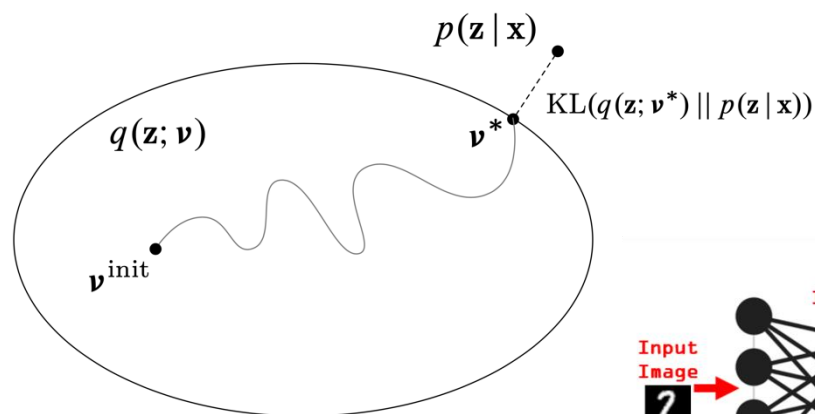


(Razavi et al., 2019)

- Pros:
  - Principled approach to generative models
  - Allows inference of  $q(\mathbf{z}|\mathbf{x})$ , can be useful feature representation for other tasks
- Cons:
  - Samples blurrier and lower quality compared to GANs
  - Tend to collapse on text data

# Summary: Supervised / Unsupervised Learning

- Supervised Learning
  - Maximum likelihood estimation (MLE)
- Unsupervised learning
  - Maximum likelihood estimation (MLE) with latent variables
    - Marginal log-likelihood
  - EM algorithm for MLE
    - ELBO / Variational free energy
  - Variational Inference
    - ELBO / Variational free energy
    - Variational distributions
      - Factorized (mean-field VI)
      - Mixture of Gaussians (Black-box VI)
      - Neural-based (VAEs)



# **Presentations**



**Questions?**