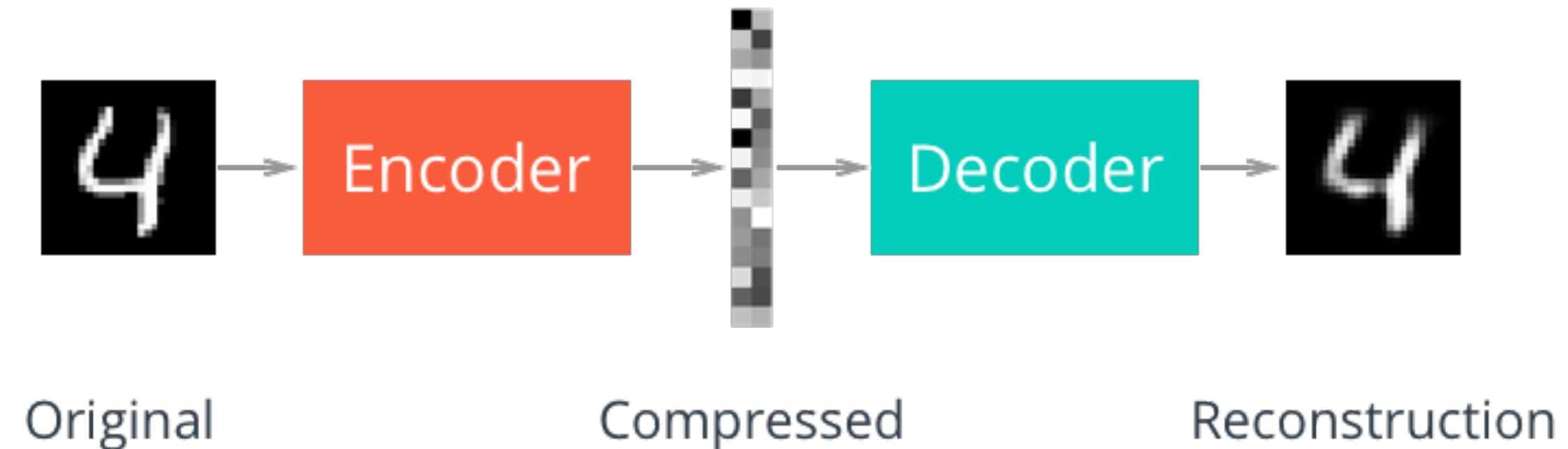


Generative Models

Variational AutoEncoders (VAEs)

Autoencoder



- Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data
- Autoencoding is training a network to **replicate** its input to its output
- Applications:
 - Data compression
 - Learning **embeddings** to support information retrieval
 - Unlabeled pre-training for semi-supervised learning
 - Generation of new instances similar to those in the training set

Variational Autoencoders (part 2)

Gaussian Distribution

In [statistics](#), a **normal distribution** or **Gaussian distribution** is a type of [continuous probability distribution](#) for a [real-valued random variable](#). The general form of its [probability density function](#) is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

The parameter μ is the [mean](#) or [expectation](#) of the distribution (and also its [median](#) and [mode](#)), while the parameter σ is its [standard deviation](#). The [variance](#) of the distribution is σ^2 . A random variable with a Gaussian distribution is said to be **normally distributed**, and is called a **normal deviate**.

Multivariate Normal Distribution

Non-degenerate case [\[edit \]](#)

The multivariate normal distribution is said to be "non-degenerate" when the symmetric covariance matrix Σ is [positive definite](#). In this case the distribution has [density](#)^[5]

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

where \mathbf{x} is a real k -dimensional column vector and $|\boldsymbol{\Sigma}| \equiv \det \boldsymbol{\Sigma}$ is the [determinant](#) of $\boldsymbol{\Sigma}$, also known as the [generalized variance](#). The equation above reduces to that of the univariate normal distribution if $\boldsymbol{\Sigma}$ is a 1×1 matrix (i.e. a single real number).

Kullback-Leibler divergence

(1) KL (Kullback–Leibler) divergence measures how one probability distribution p diverges from a second expected probability distribution q .

$$D_{KL}(p\|q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx$$

D_{KL} achieves the minimum zero when $p(x) == q(x)$ everywhere.

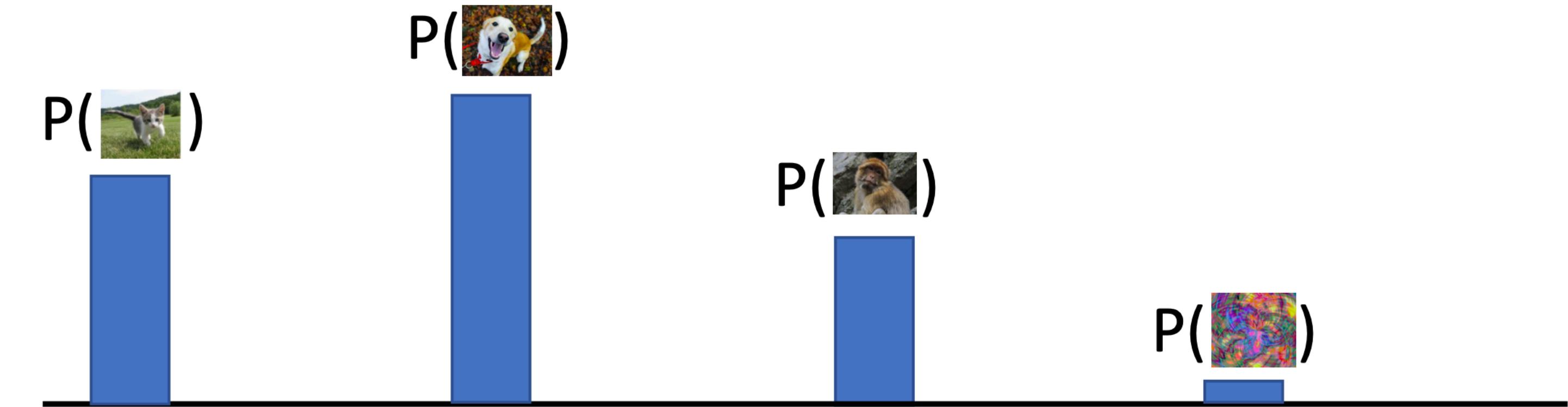
It is noticeable according to the formula that KL divergence is asymmetric. In cases where $p(x)$ is close to zero, but $q(x)$ is significantly non-zero, the q 's effect is disregarded. It could cause buggy results when we just want to measure the similarity between two equally important distributions.

Variational Autoencoders

**following “Lecture 19.pdf” (see SKOS)
slides 65-127**

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$



Generative model: All possible images compete with each other for probability mass

Conditional Generative Model: Learn $p(x|y)$

Given dataset $x^{(1)}, x^{(2)}, \dots x^{(N)}$, train the model by solving:

$$W^* = \arg \max_W \prod_i p(x^{(i)})$$

Maximize probability of training data
(Maximum likelihood estimation)

$$= \arg \max_W \sum_i \log p(x^{(i)})$$

Log trick to exchange product for sum

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data

Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

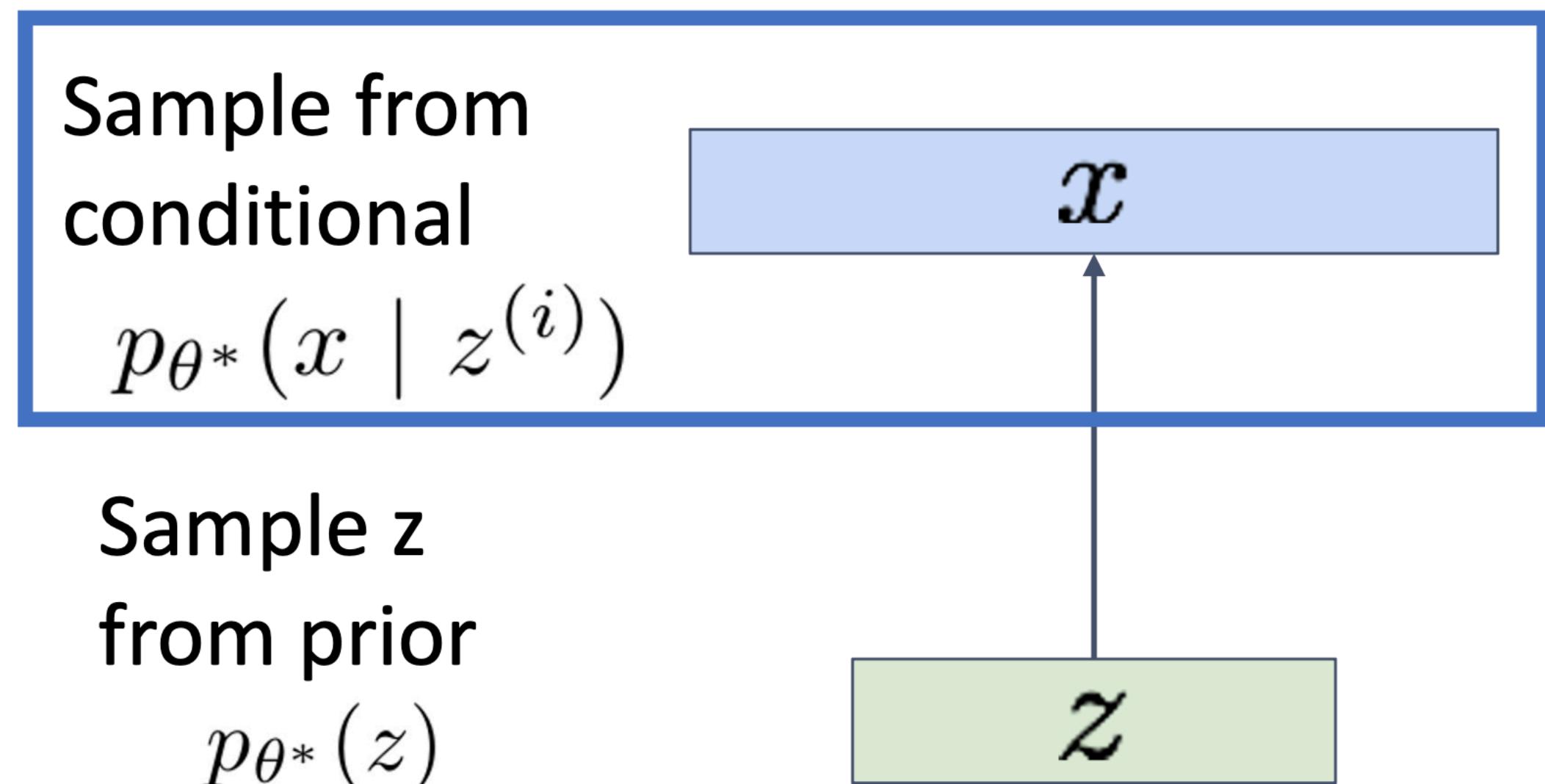
Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data

After training, sample new data like this:



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

Assume simple prior $p(z)$, e.g. Gaussian

Represent $p(x|z)$ with a neural network
(Similar to **decoder** from autencoder)

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

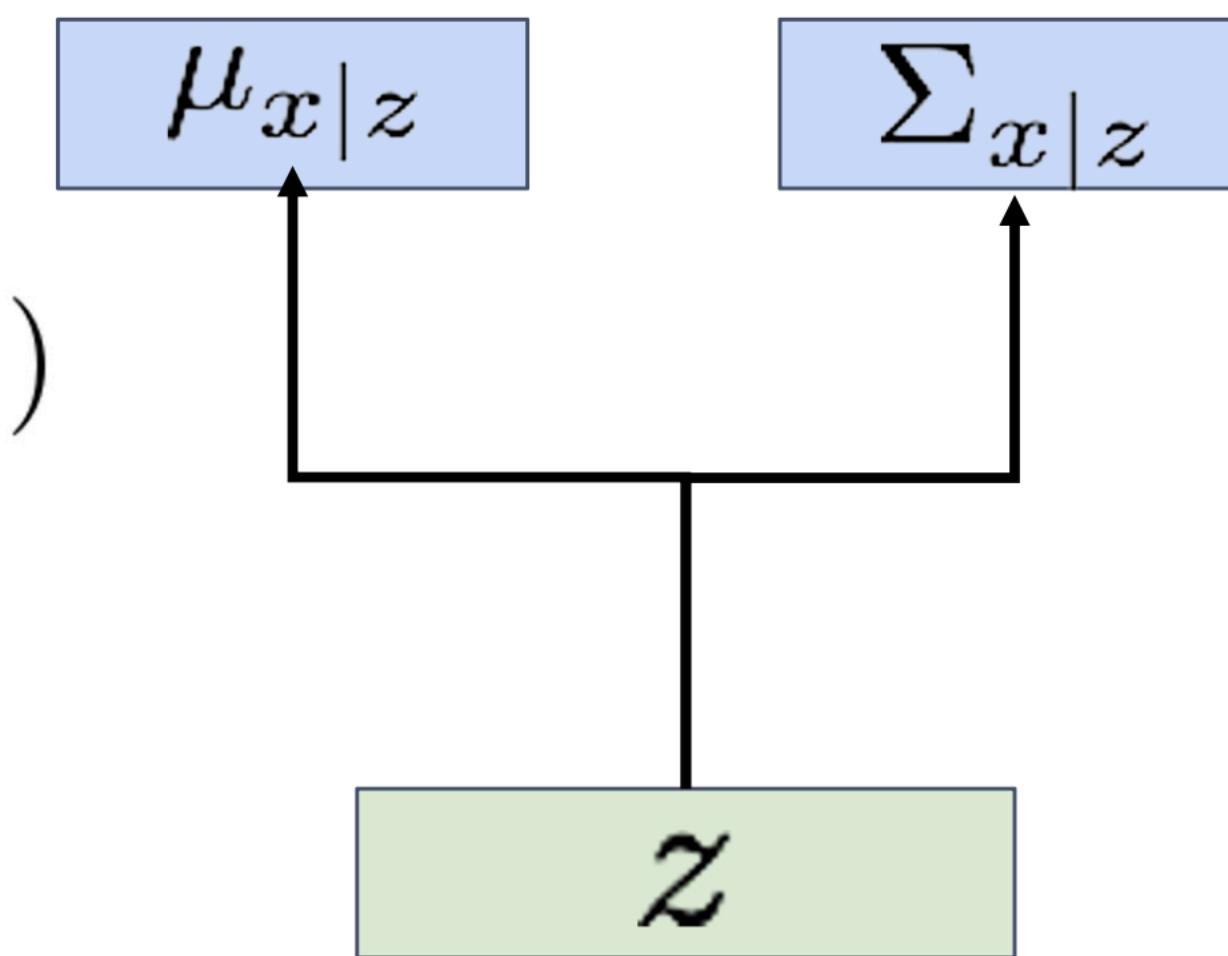
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

Assume simple prior $p(z)$, e.g. Gaussian

Represent $p(x|z)$ with a neural network
(Similar to **decoder** from autencoder)

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

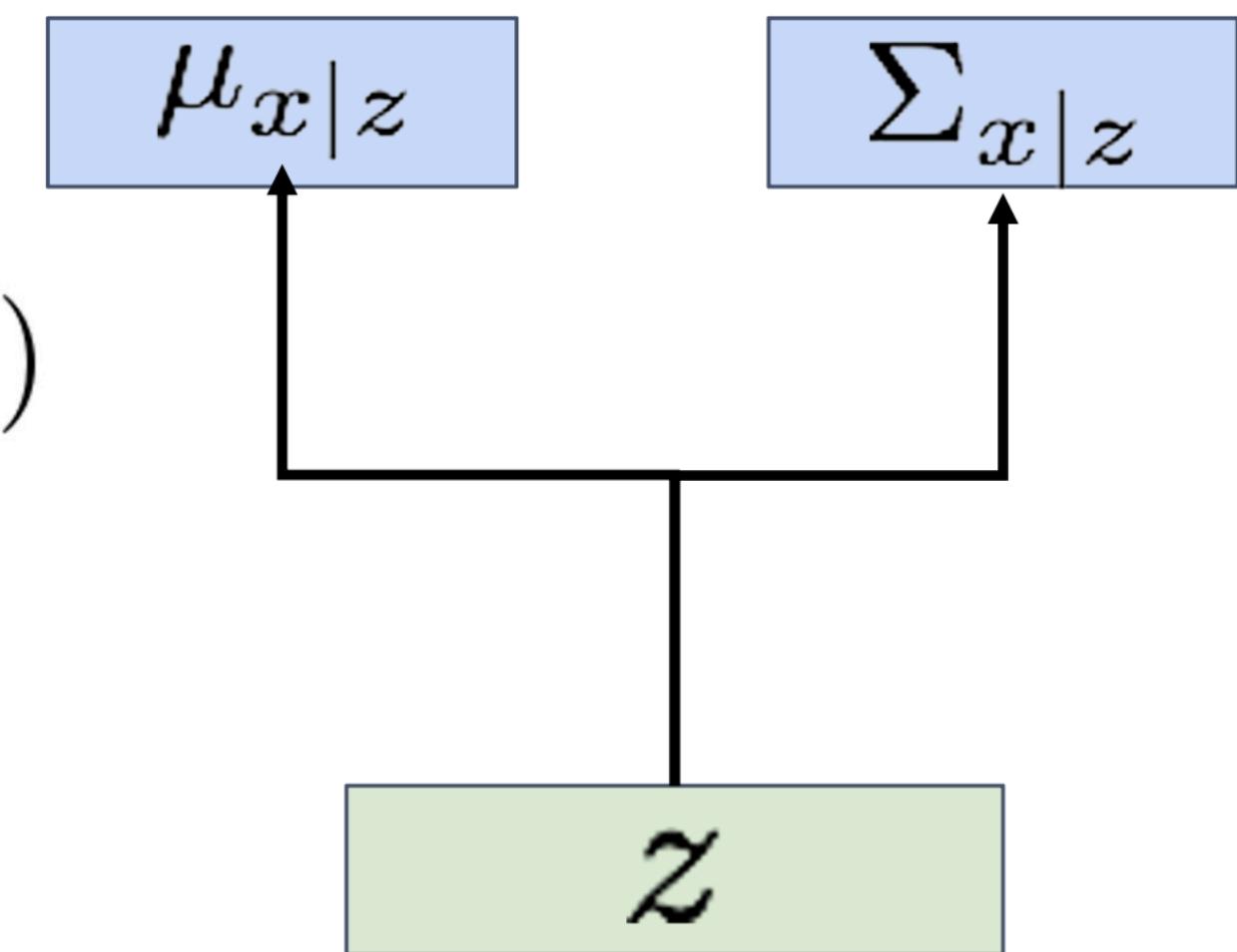
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

If we could observe the z for each x , then could train a *conditional generative model* $p(x|z)$

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

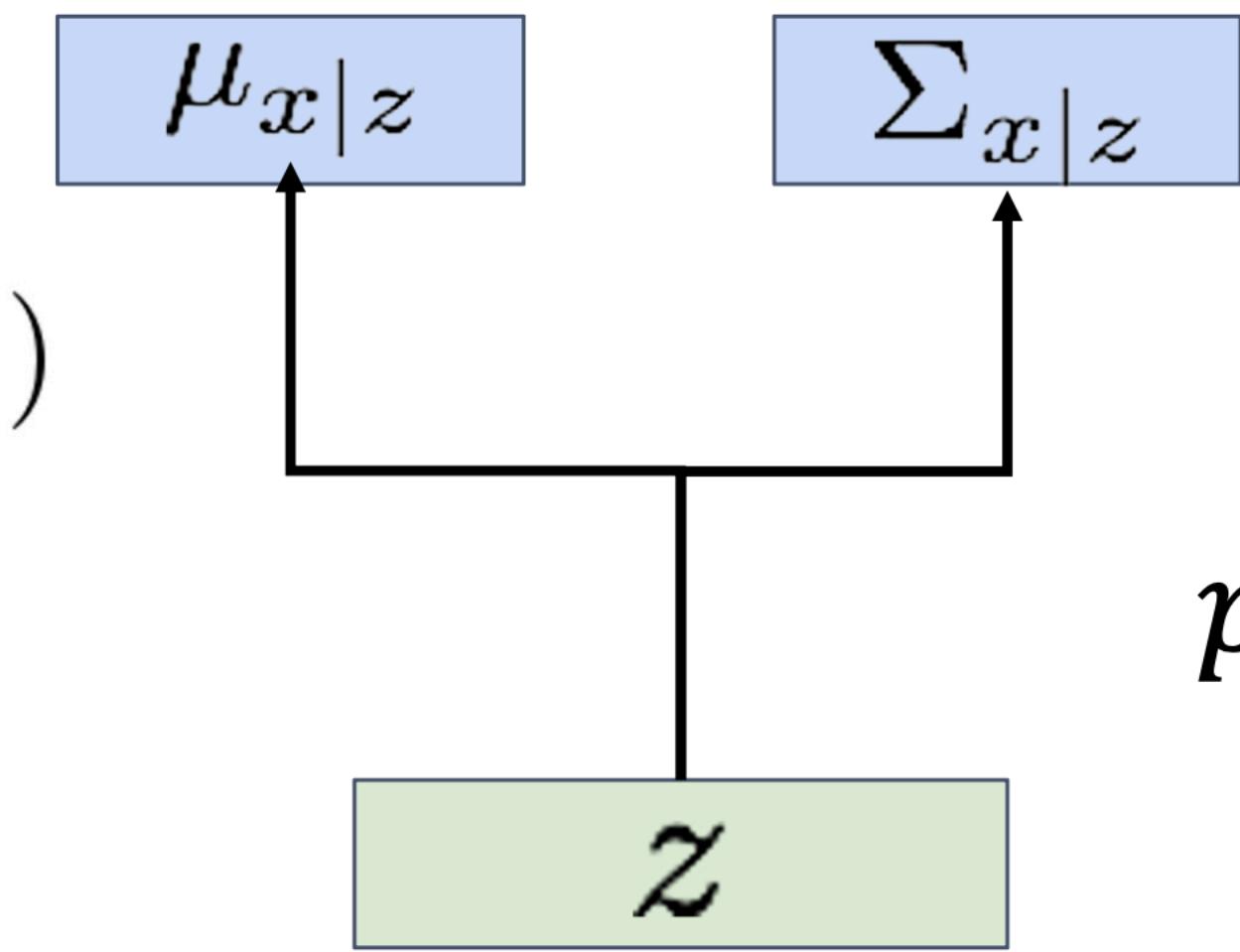
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

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

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

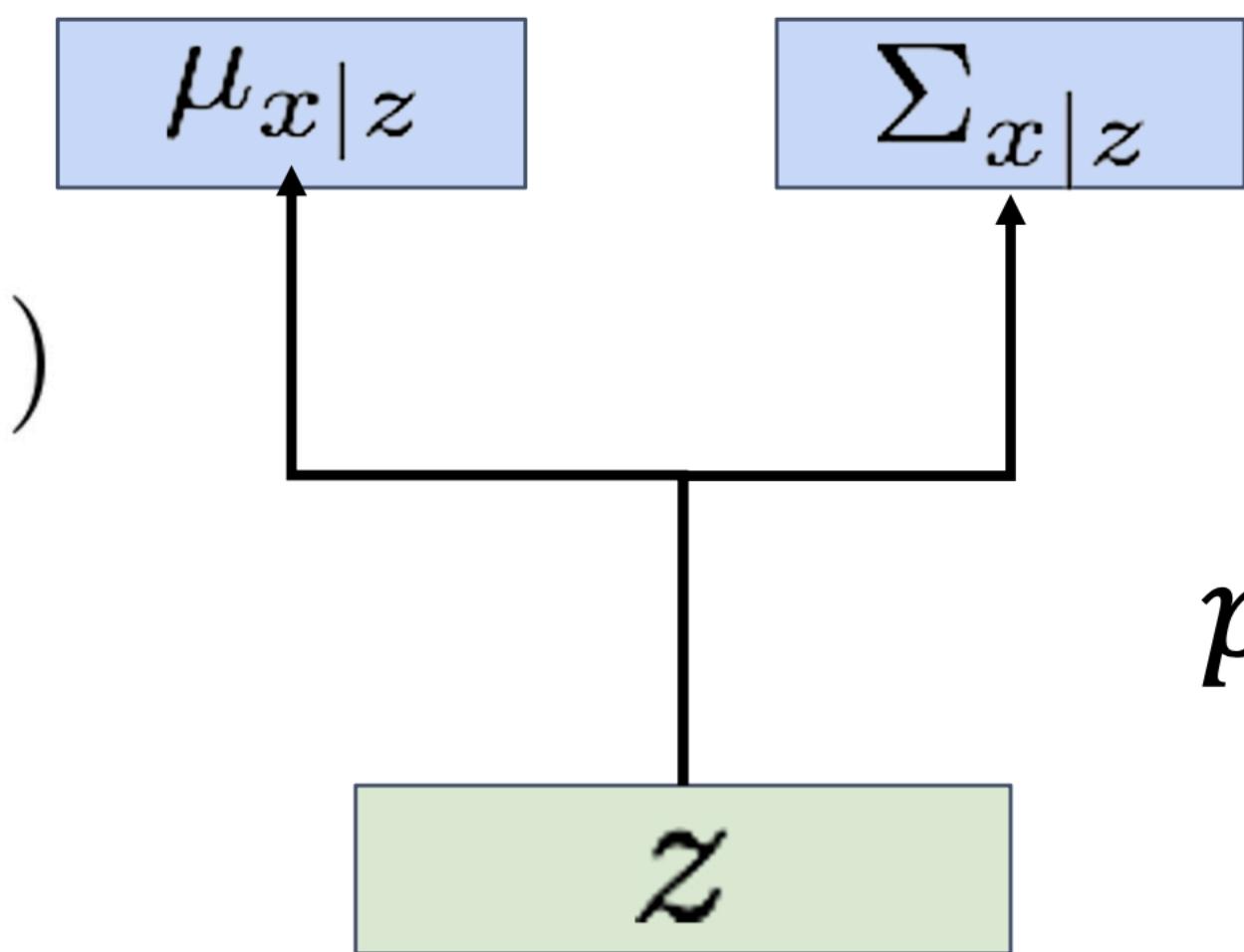
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int [p_{\theta}(x|z)p_{\theta}(z)] dz$$

Ok, can compute this with decoder network

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

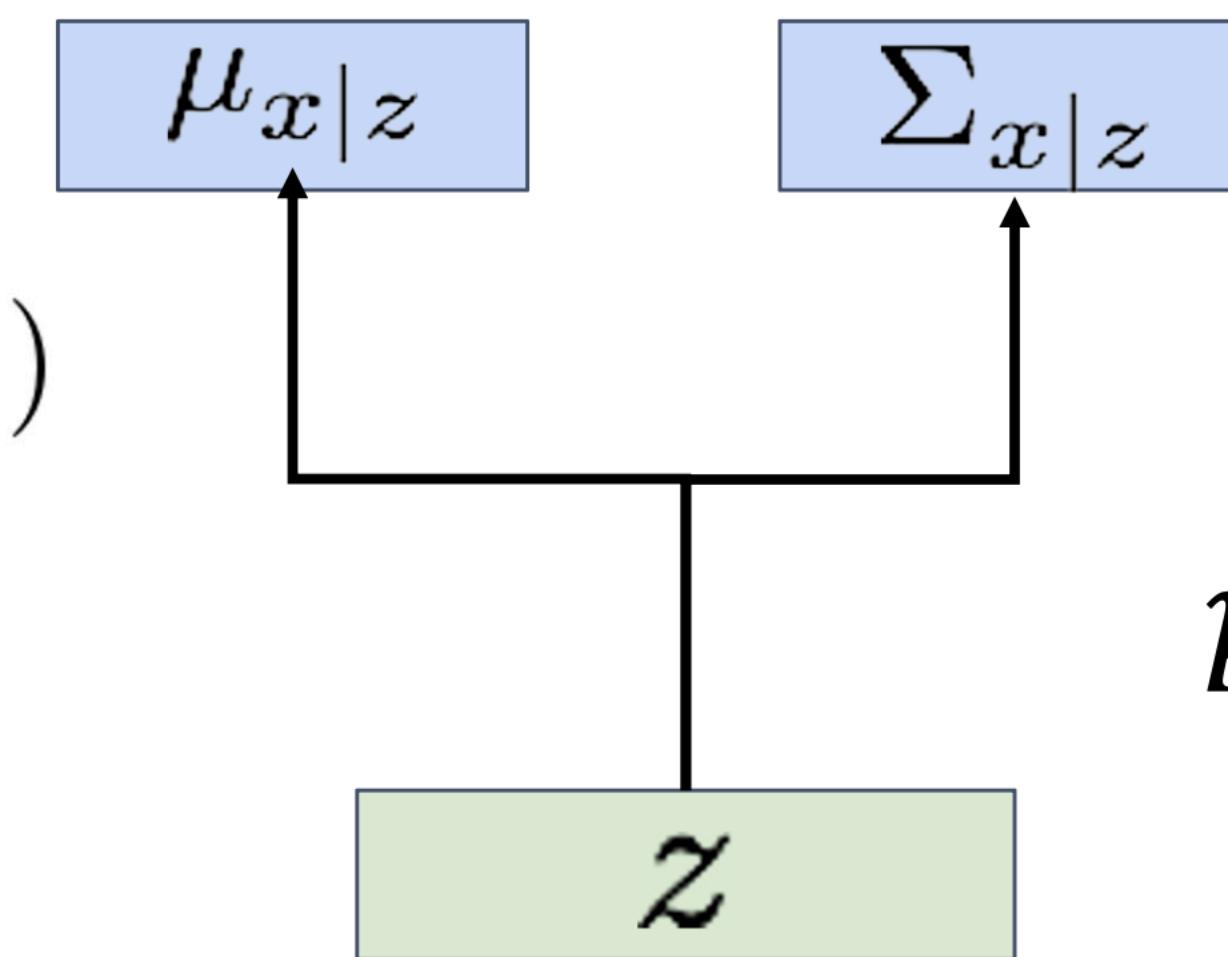
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

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

Ok, we assumed Gaussian prior for z

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

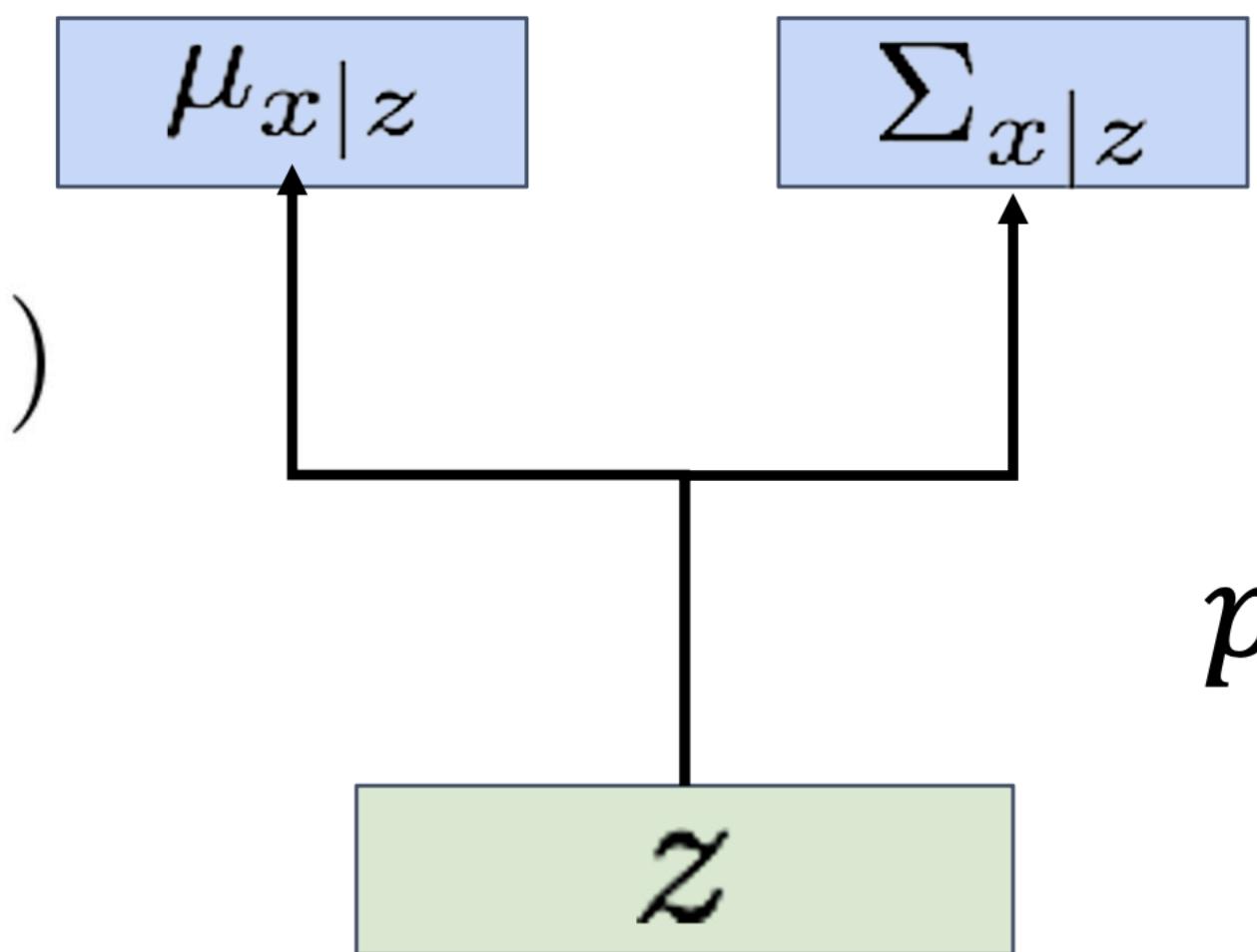
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

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

Problem: Impossible to integrate over all z !

Variational Autoencoders

Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

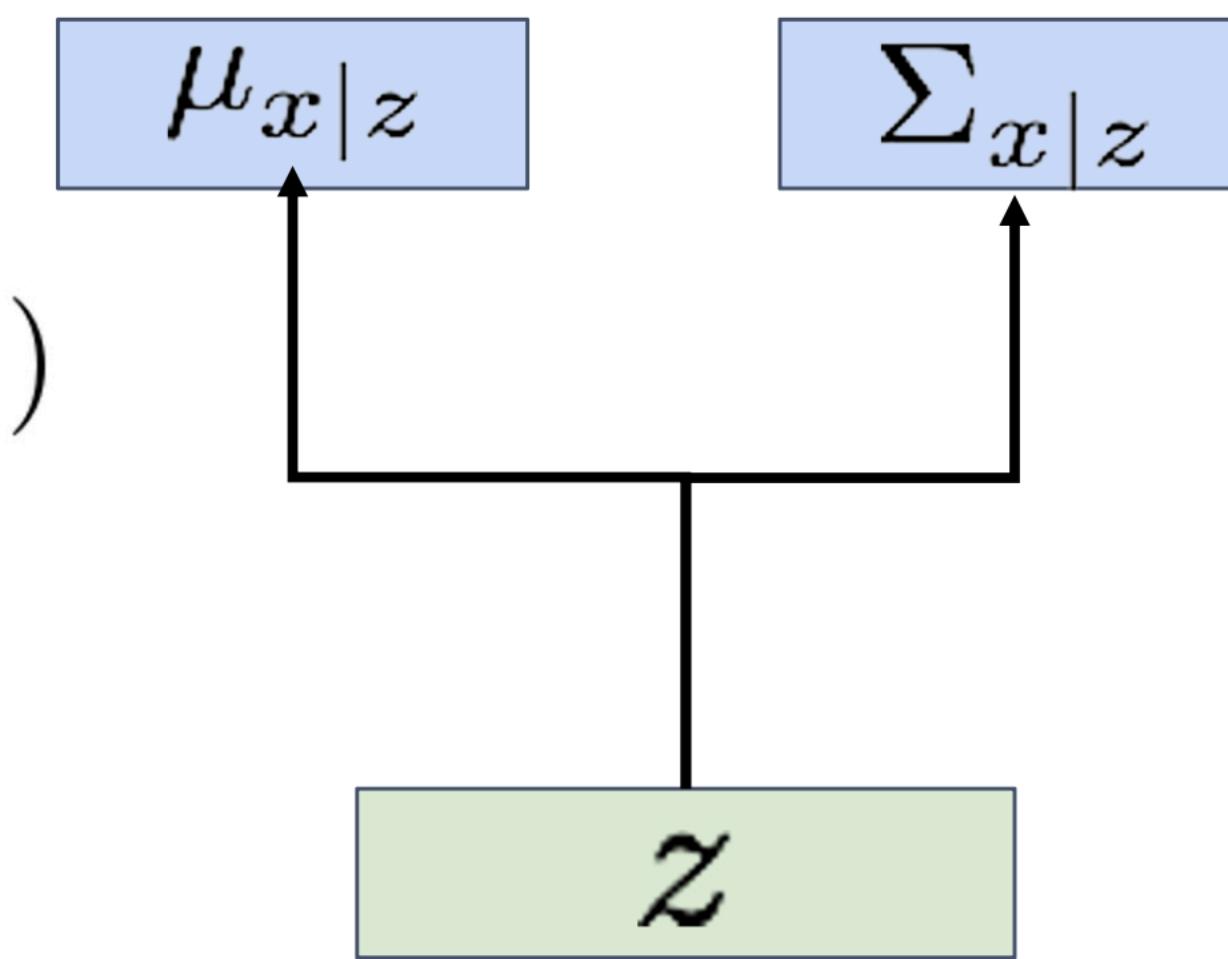
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x | z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Variational Autoencoders

Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

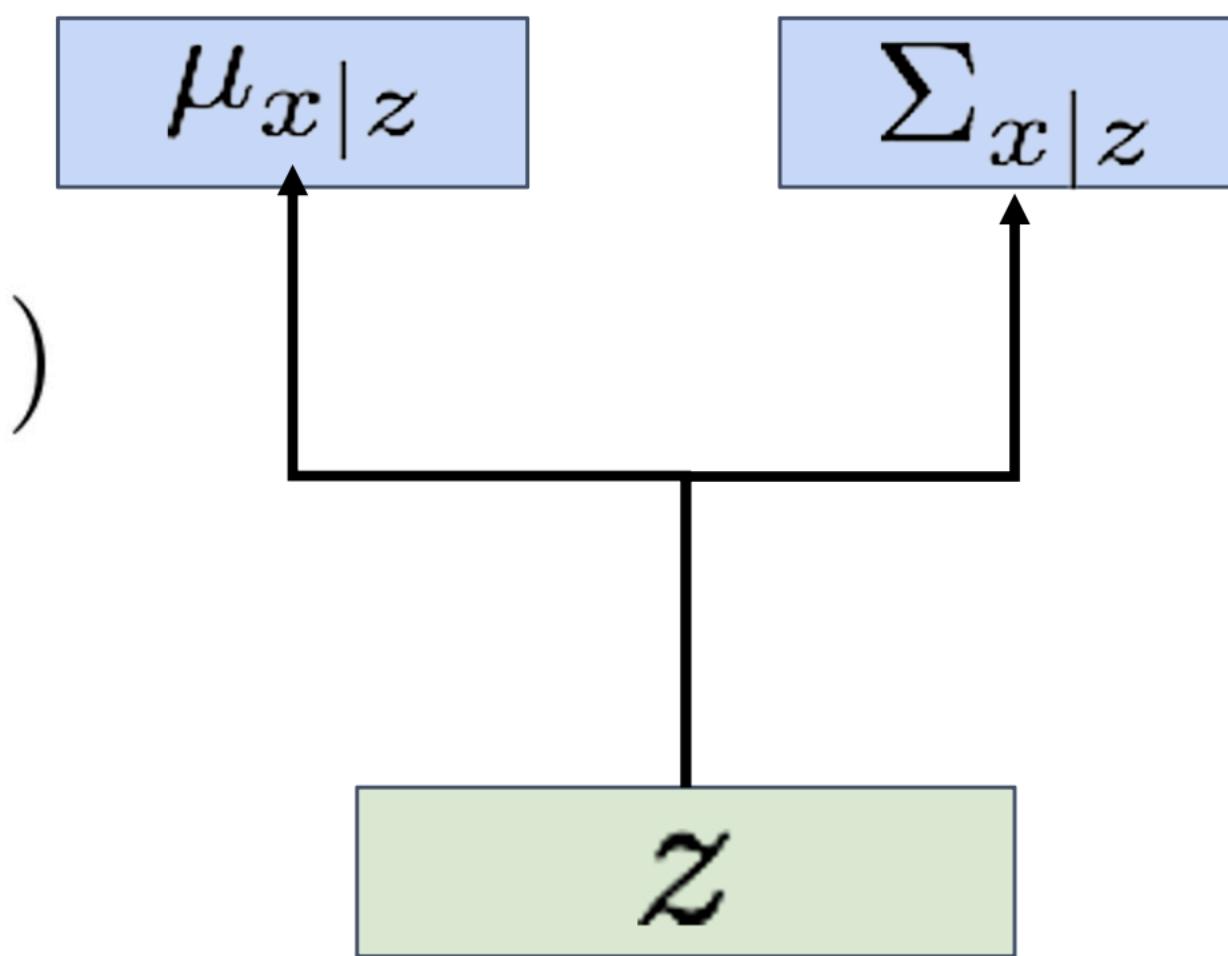
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x | z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Ok, compute with decoder network

Variational Autoencoders

Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

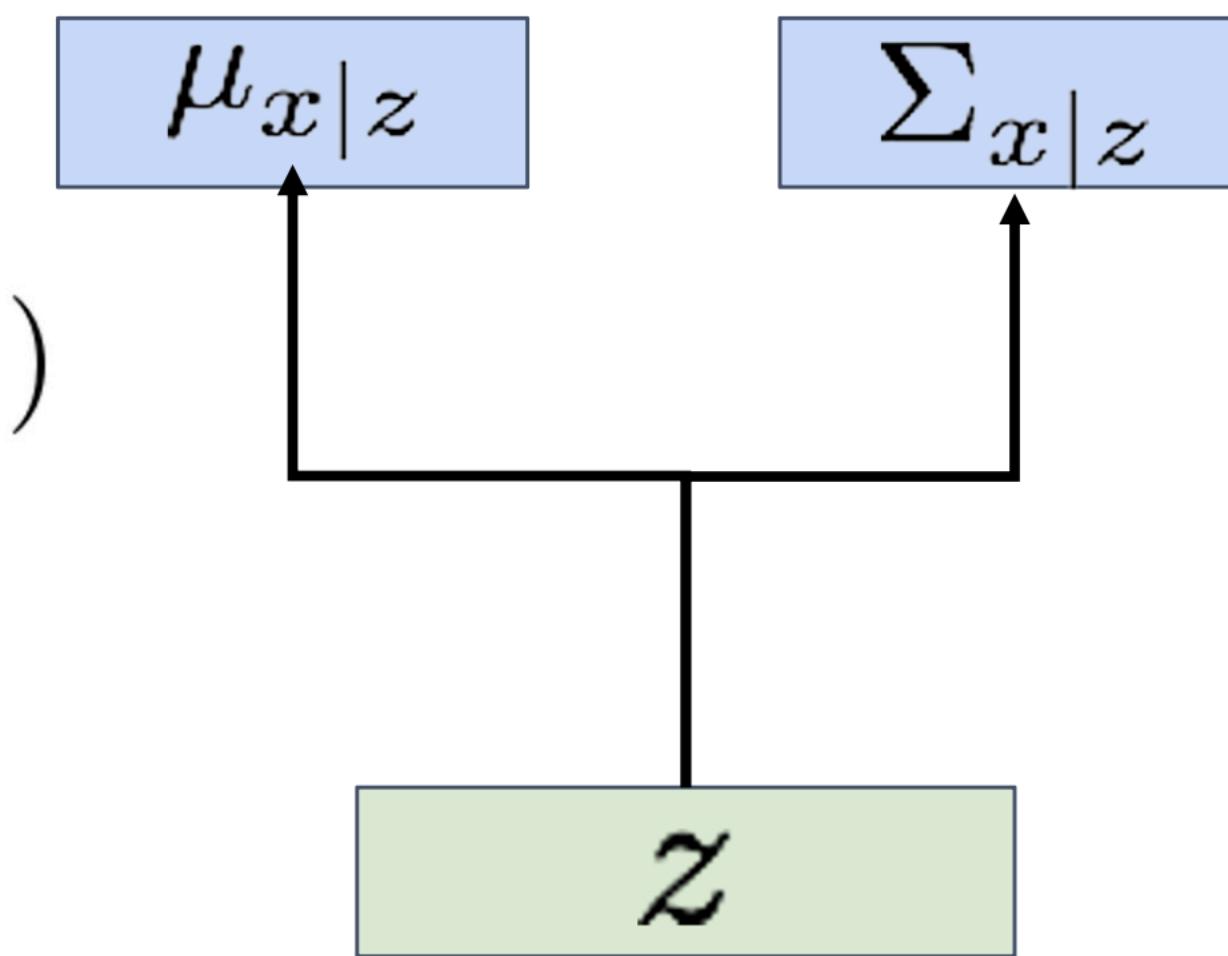
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x | z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Ok, we assumed Gaussian prior

Variational Autoencoders

Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

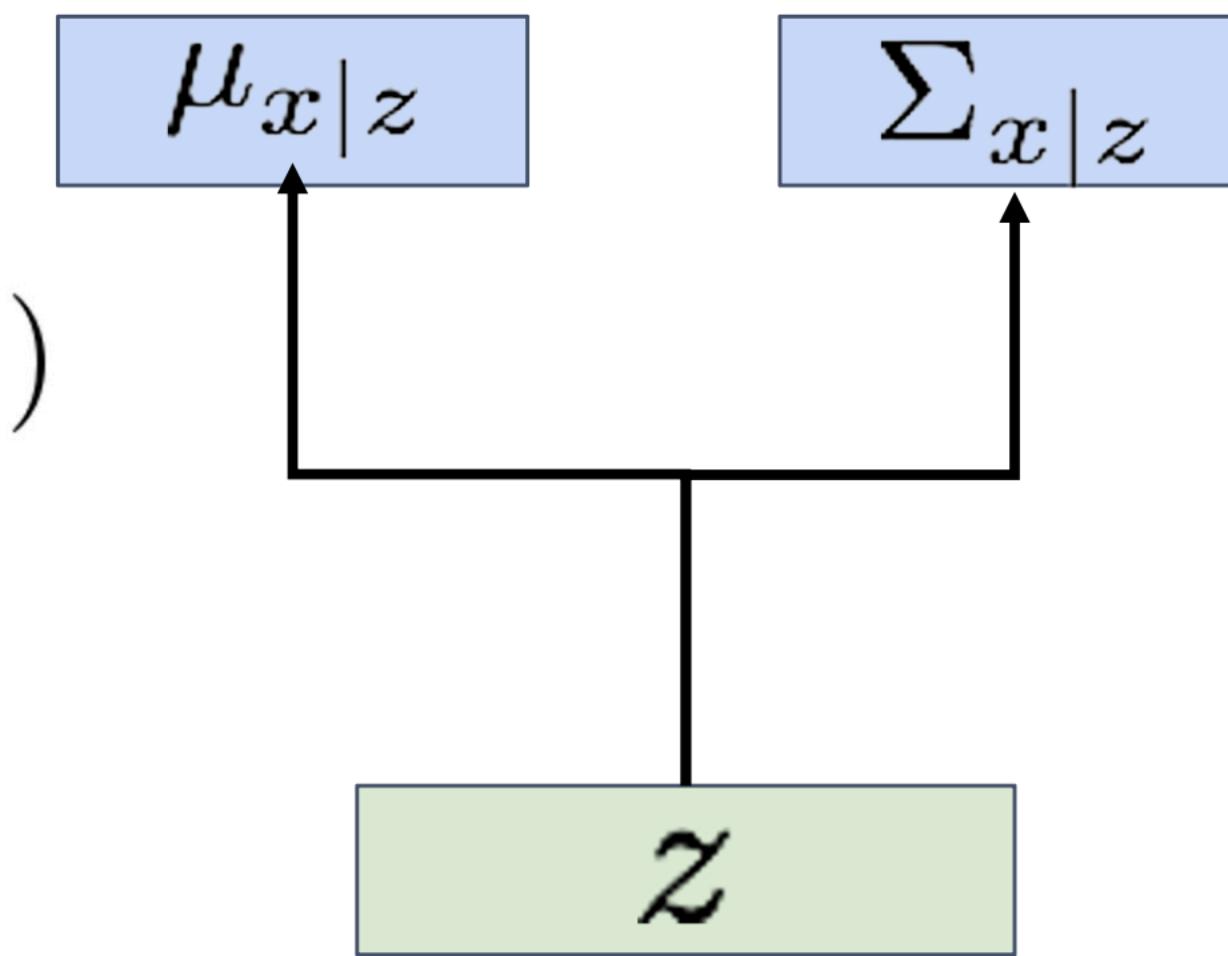
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x | z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Problem: No way to compute this!

Variational Autoencoders

Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

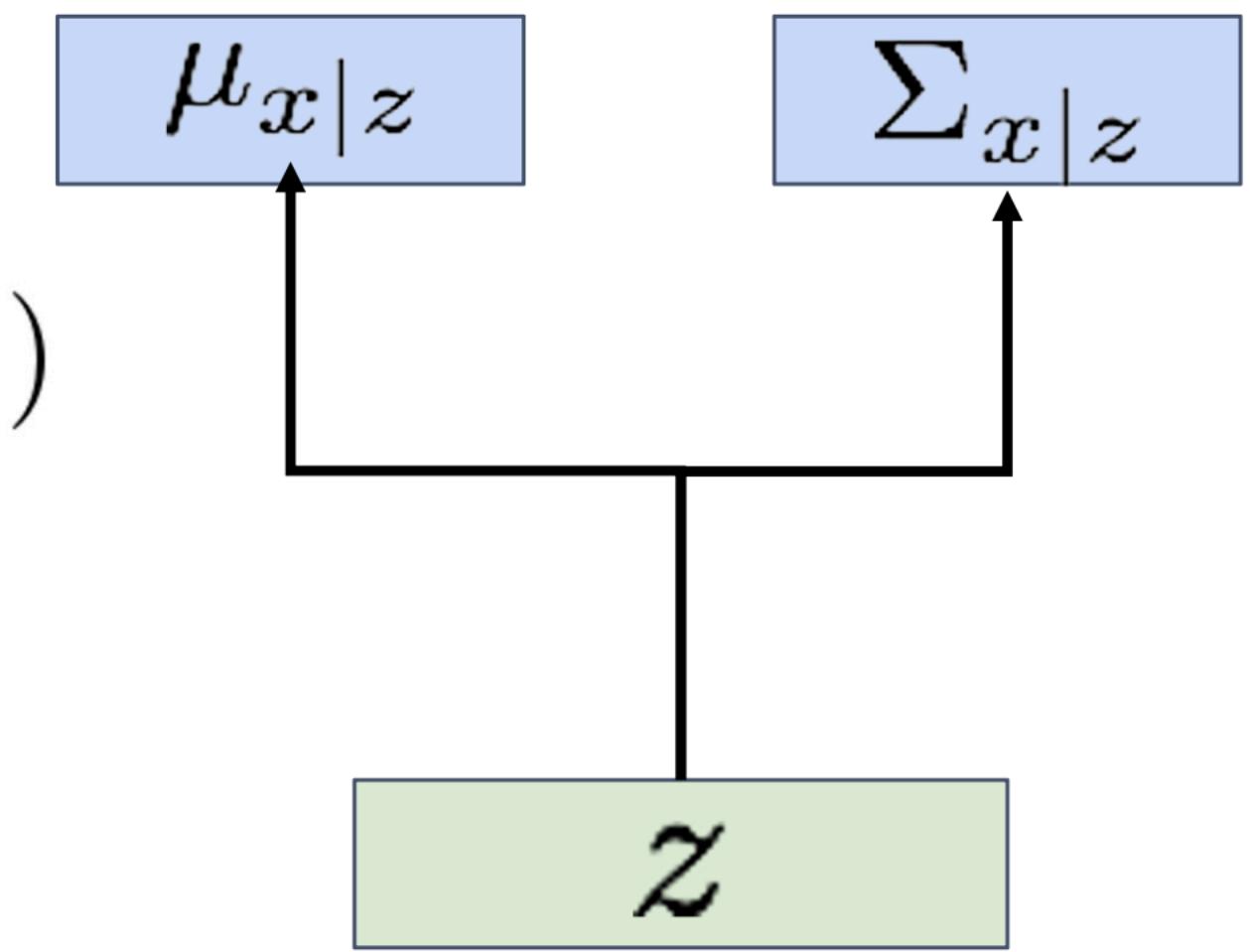
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x | z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(z | x)}$$

Solution: Train another network (**encoder**) that learns $q_{\phi}(z | x) \approx p_{\theta}(z | x)$

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

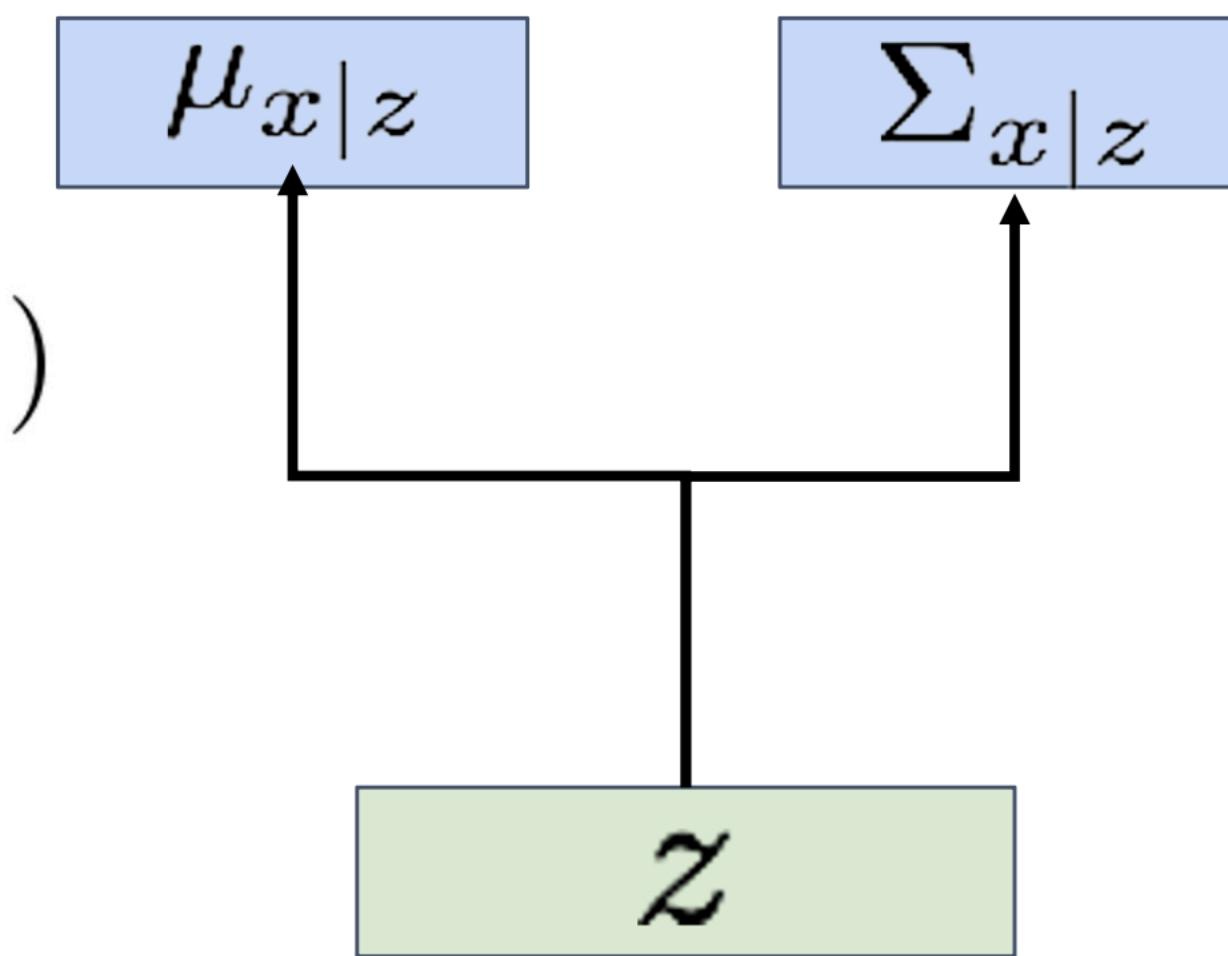
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample z from prior

$$p_{\theta^*}(z)$$



Recall $p(x, z) = p(x \mid z)p(z) = p(z \mid x)p(x)$

Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

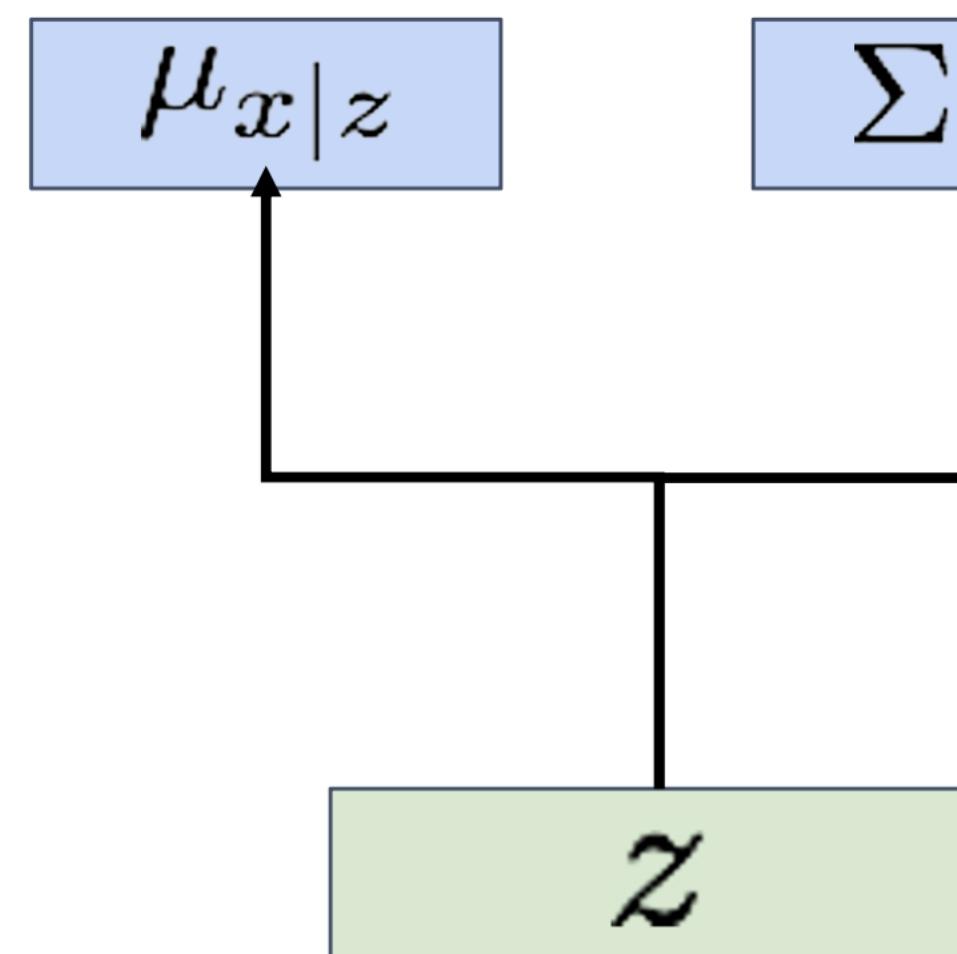
$$p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)} \approx \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{q_{\phi}(z \mid x)}$$

Use **encoder** to compute $q_{\phi}(z \mid x) \approx p_{\theta}(z \mid x)$

Variational Autoencoders

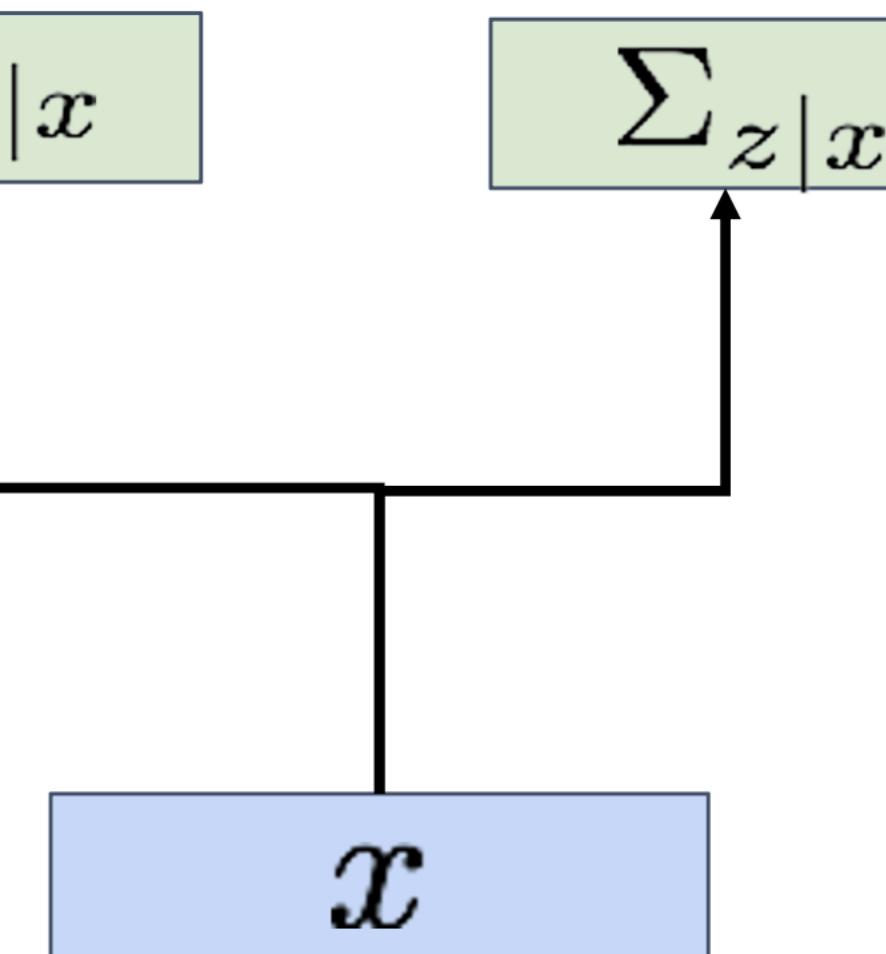
Decoder network inputs
latent code z , gives
distribution over data x

$$p_\theta(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



Encoder network inputs
data x , gives distribution
over latent codes z

$$q_\phi(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$



If we can ensure that
 $q_\phi(z | x) \approx p_\theta(z | x)$,
then we can approximate

$$p_\theta(x) \approx \frac{p_\theta(x | z)p(z)}{q_\phi(z | x)}$$

Idea: Jointly train both
encoder and decoder

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= \log p_\theta(x|z) - \log \frac{q_\phi(z|x)}{p(z)} + \log \frac{q_\phi(z|x)}{p_\theta(z|x)}\end{aligned}$$

Split up using rules for logarithms

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= \log p_\theta(x|z) - \log \frac{q_\phi(z|x)}{p(z)} + \log \frac{q_\phi(z|x)}{p_\theta(z|x)}\end{aligned}$$

$$\log p_\theta(x) = E_{z \sim q_\phi(z|x)} [\log p_\theta(x)]$$

We can wrap in an expectation since it doesn't depend on z

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right]\end{aligned}$$

$$\log p_\theta(x) = E_{z \sim q_\phi(z|x)}[\log p_\theta(x)]$$

We can wrap in an expectation since it doesn't depend on z

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)}$$

$$\begin{aligned} &= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x)) \end{aligned}$$

Data reconstruction

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)}$$

$$= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right]$$
$$= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x))$$

KL divergence between prior, and
samples from the encoder network

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)}$$

$$= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right]$$
$$= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x))$$

KL divergence between encoder
and posterior of decoder

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)}$$

$$\begin{aligned} &= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x)) \end{aligned}$$

KL is ≥ 0 , so dropping this term gives a **lower bound** on the data likelihood:

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)}$$

$$= E_z[\log p_\theta(x|z)] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right]$$

$$= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x))$$

$$\log p_\theta(x) \geq E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right)$$

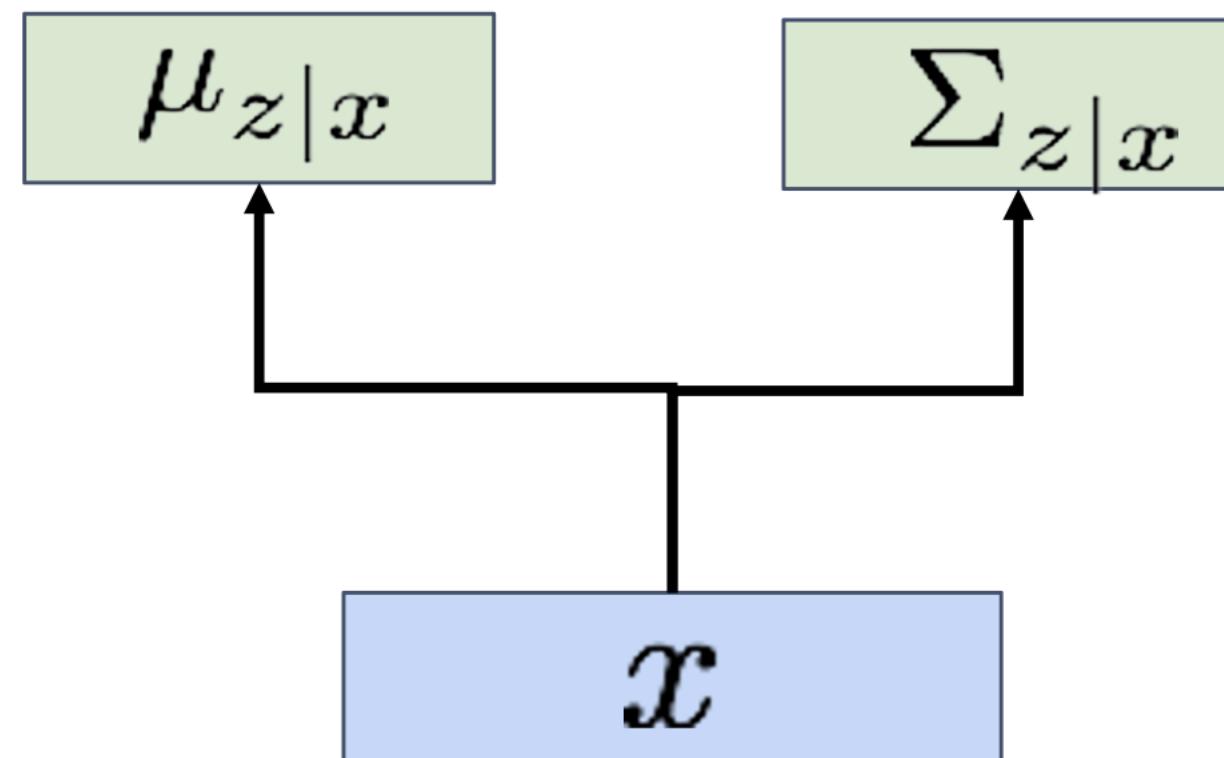
Variational Autoencoders

Jointly train **encoder q** and **decoder p** to maximize
the **variational lower bound** on the data likelihood
Also called **Evidence Lower Bound (ELBo)**

$$\log p_\theta(x) \geq E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))$$

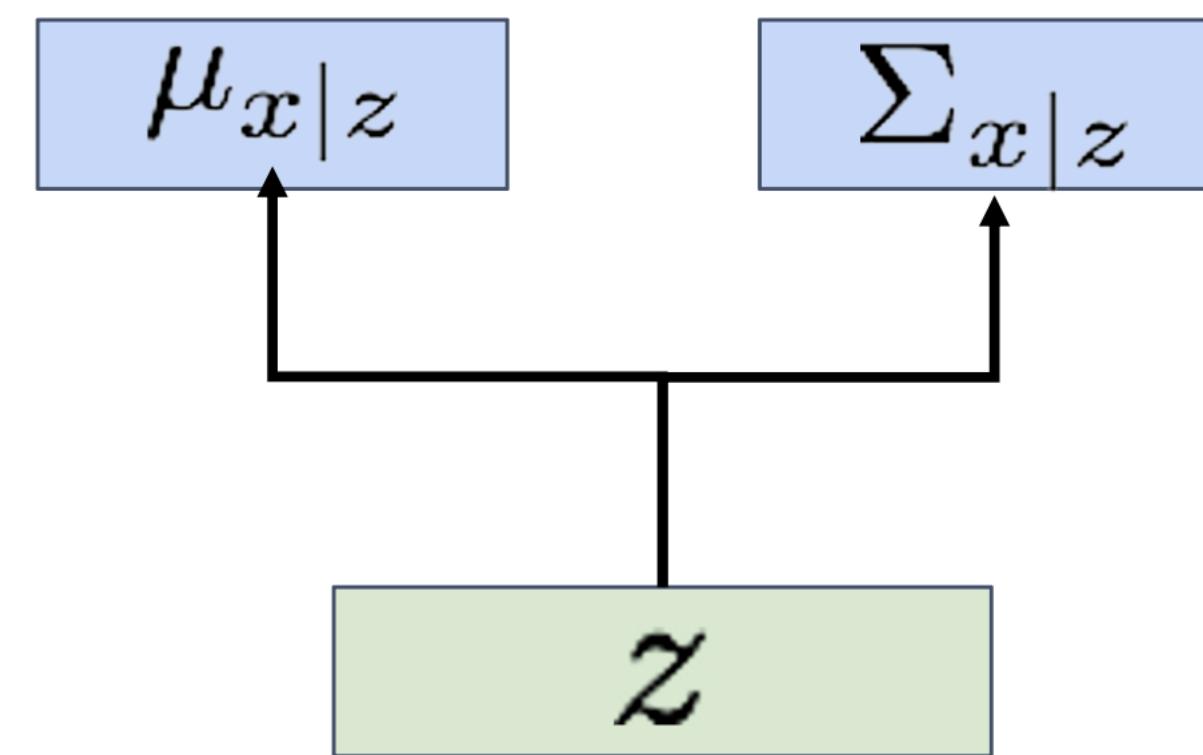
Encoder Network

$$q_\phi(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$



Decoder Network

$$p_\theta(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



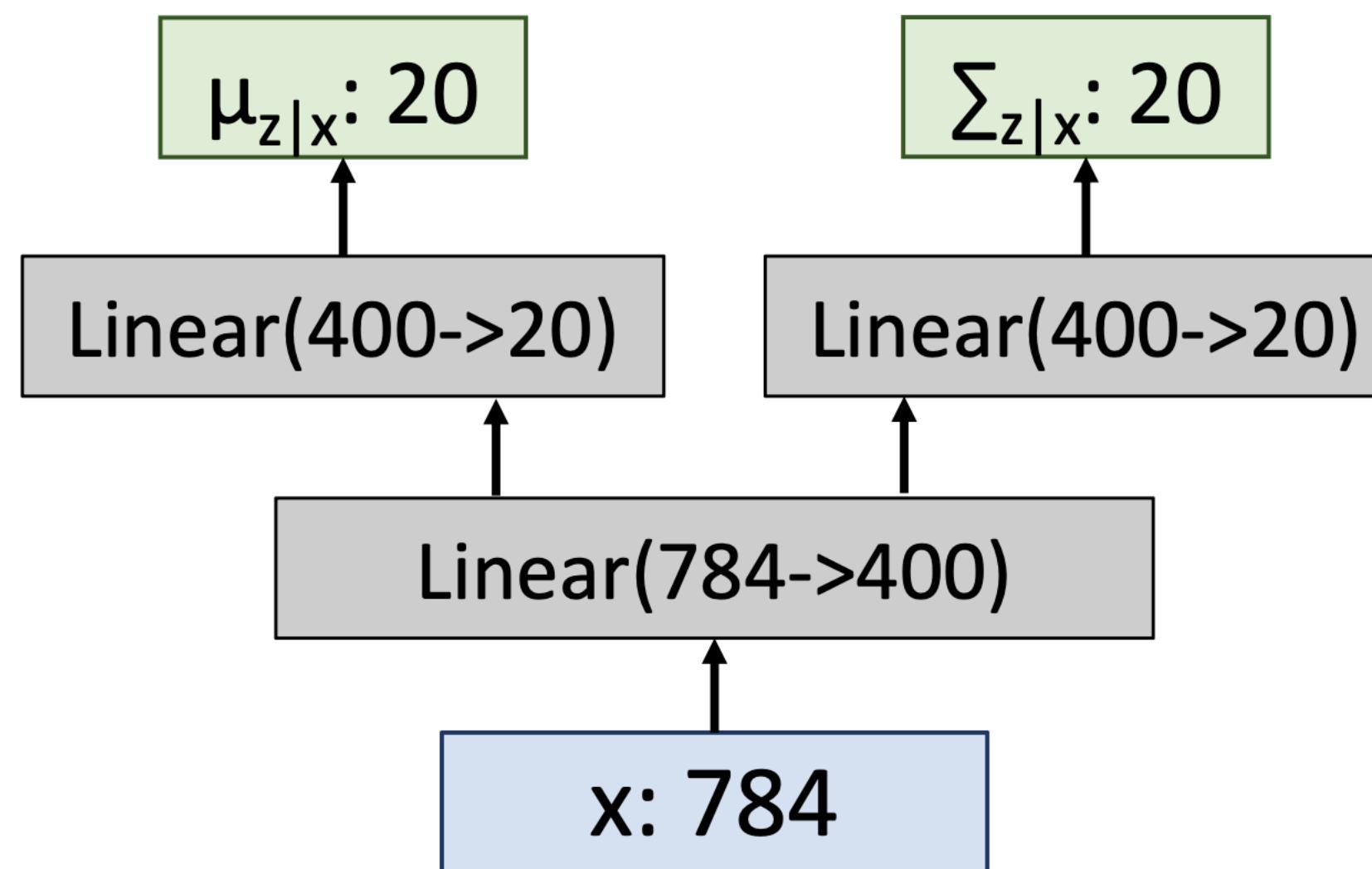
Example: Fully-Connected VAE

x : 28x28 image, flattened to 784-dim vector

z : 20-dim vector

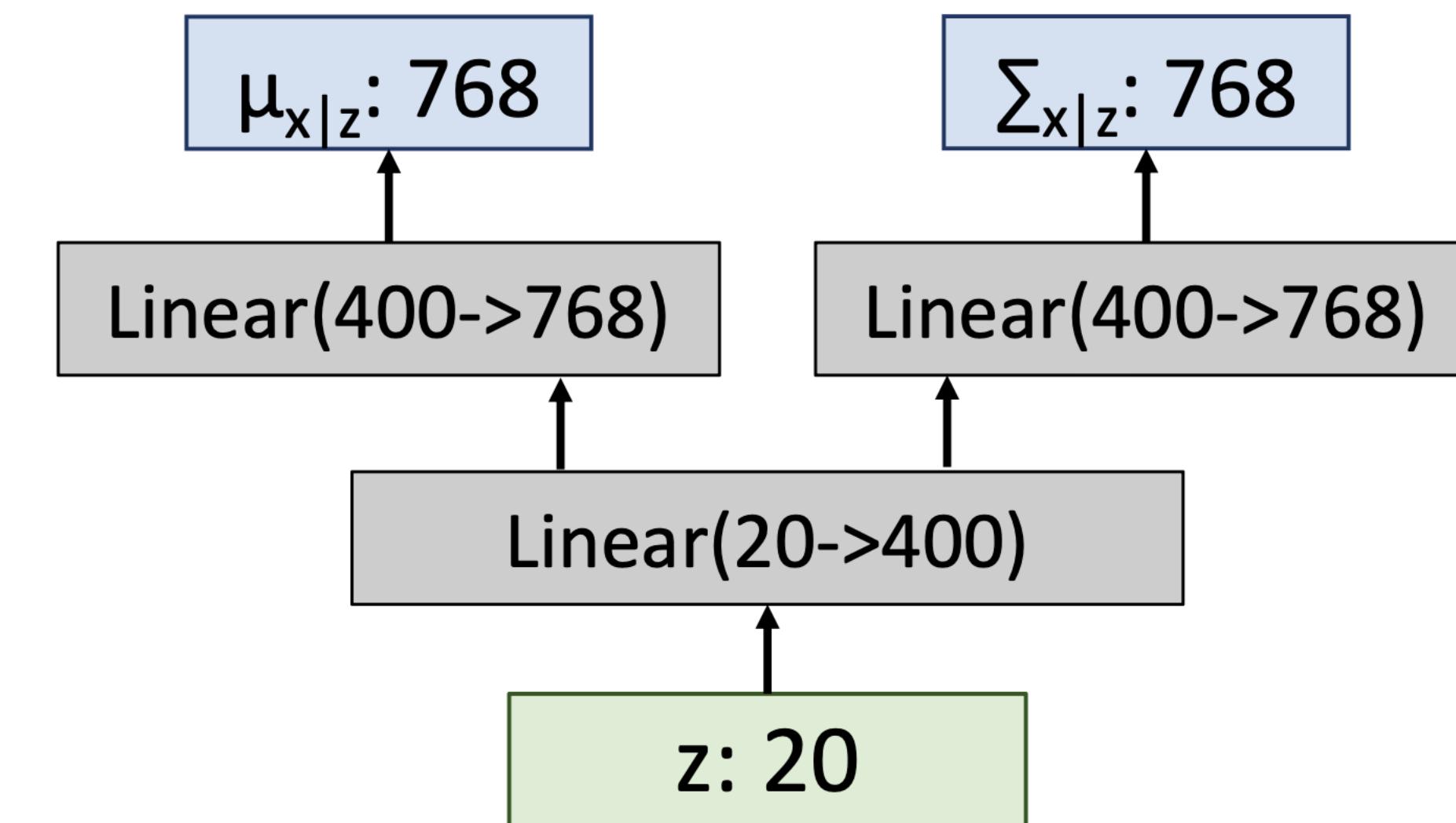
Encoder Network

$$q_{\phi}(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$



Decoder Network

$$p_{\theta}(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

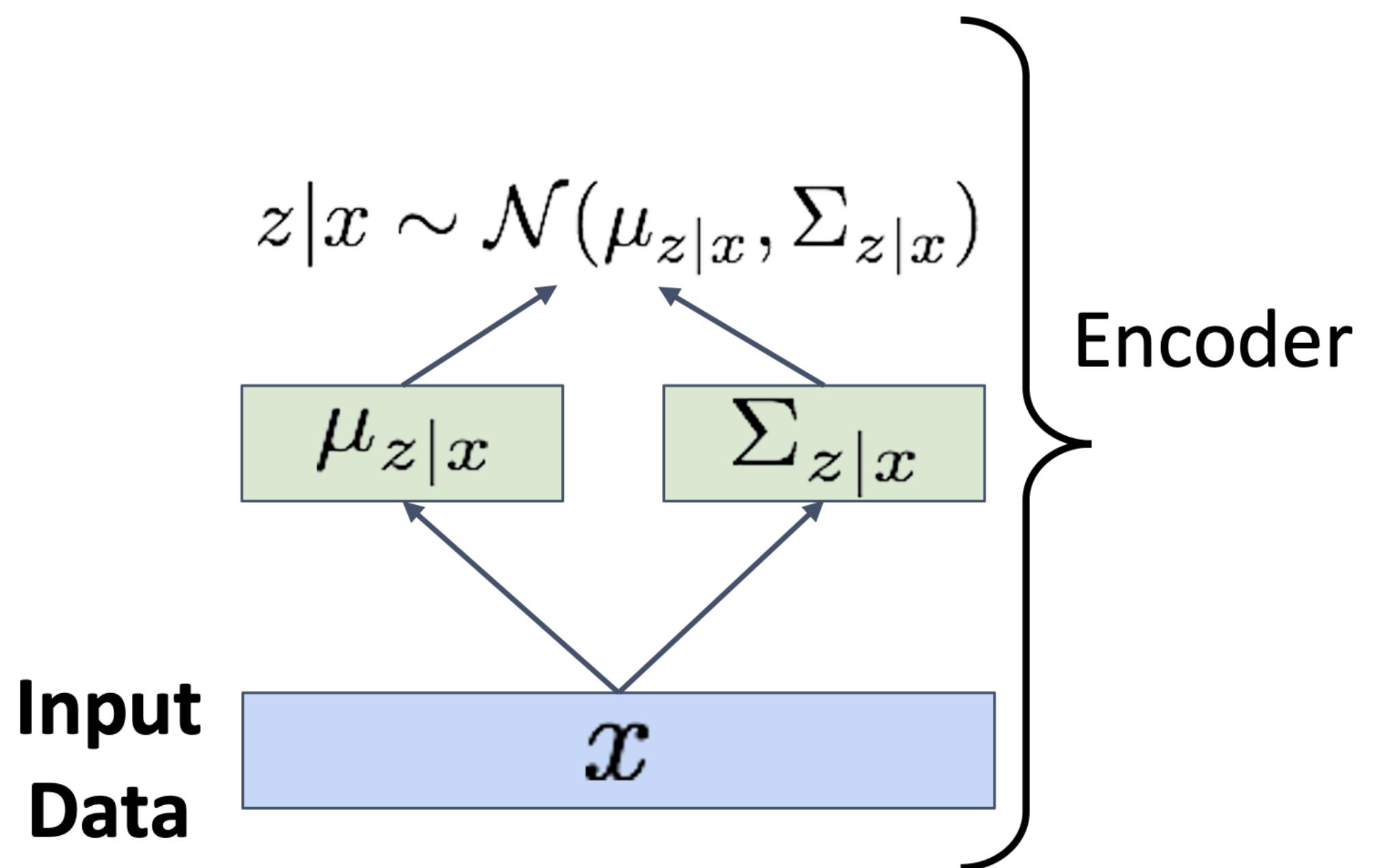
Input
Data

Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes

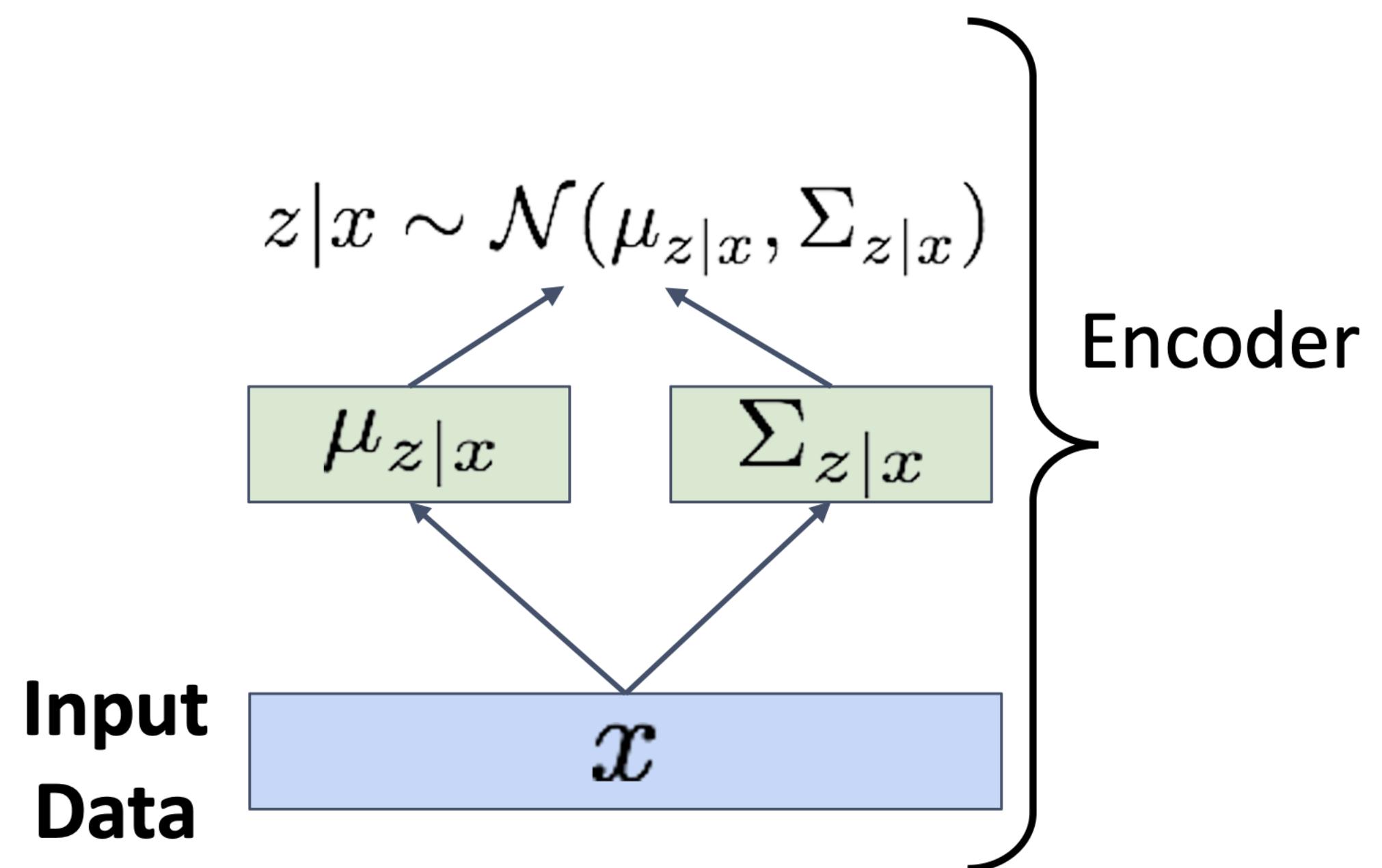


Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**



Variational Autoencoders

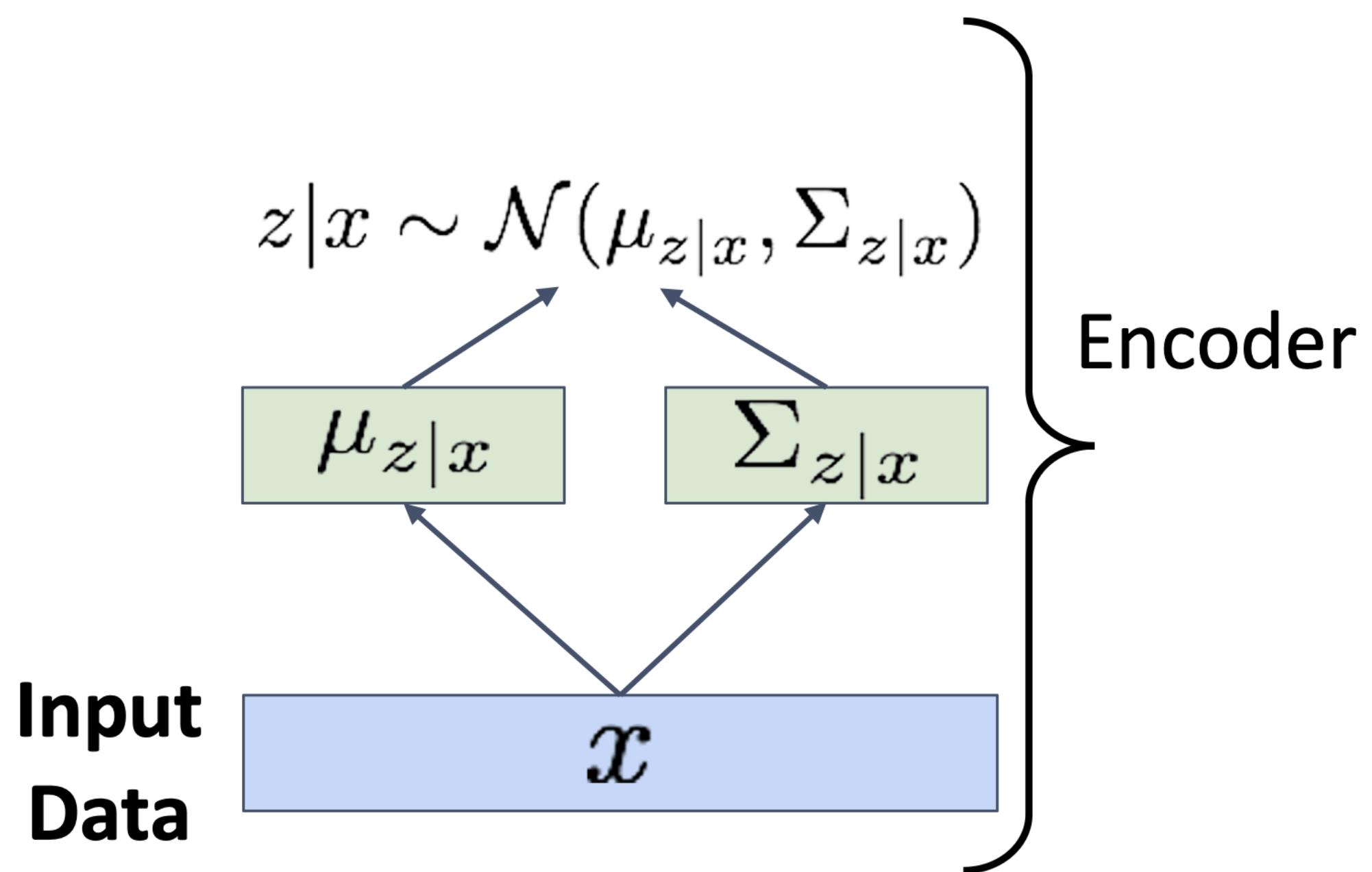
Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**

$$\begin{aligned} -D_{KL} (q_\phi(z|x), p(z)) &= \int_Z q_\phi(z|x) \log \frac{p(z)}{q_\phi(z|x)} dz \\ &= \int_Z N(z; \mu_{z|x}, \Sigma_{z|x}) \log \frac{N(z; 0, I)}{N(z; \mu_{z|x}, \Sigma_{z|x})} dz \\ &= \frac{1}{2} \sum_{j=1}^J \left(1 + \log \left((\Sigma_{z|x})_j^2 \right) - (\mu_{z|x})_j^2 - (\Sigma_{z|x})_j^2 \right) \end{aligned}$$

Closed form solution when
 q_ϕ is diagonal Gaussian and
 p is unit Gaussian!
(Assume z has dimension J)

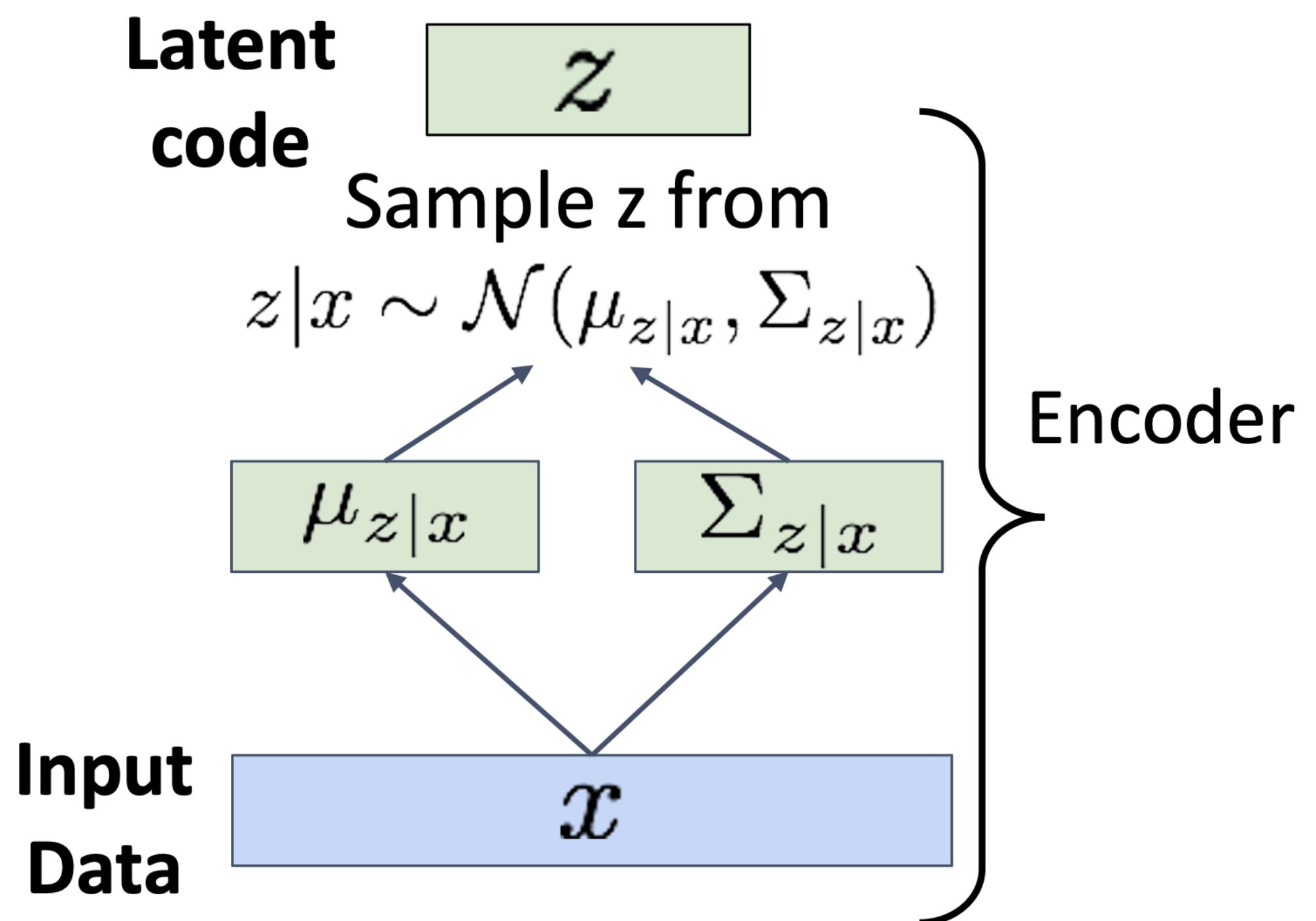


Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**
3. Sample code z from encoder output

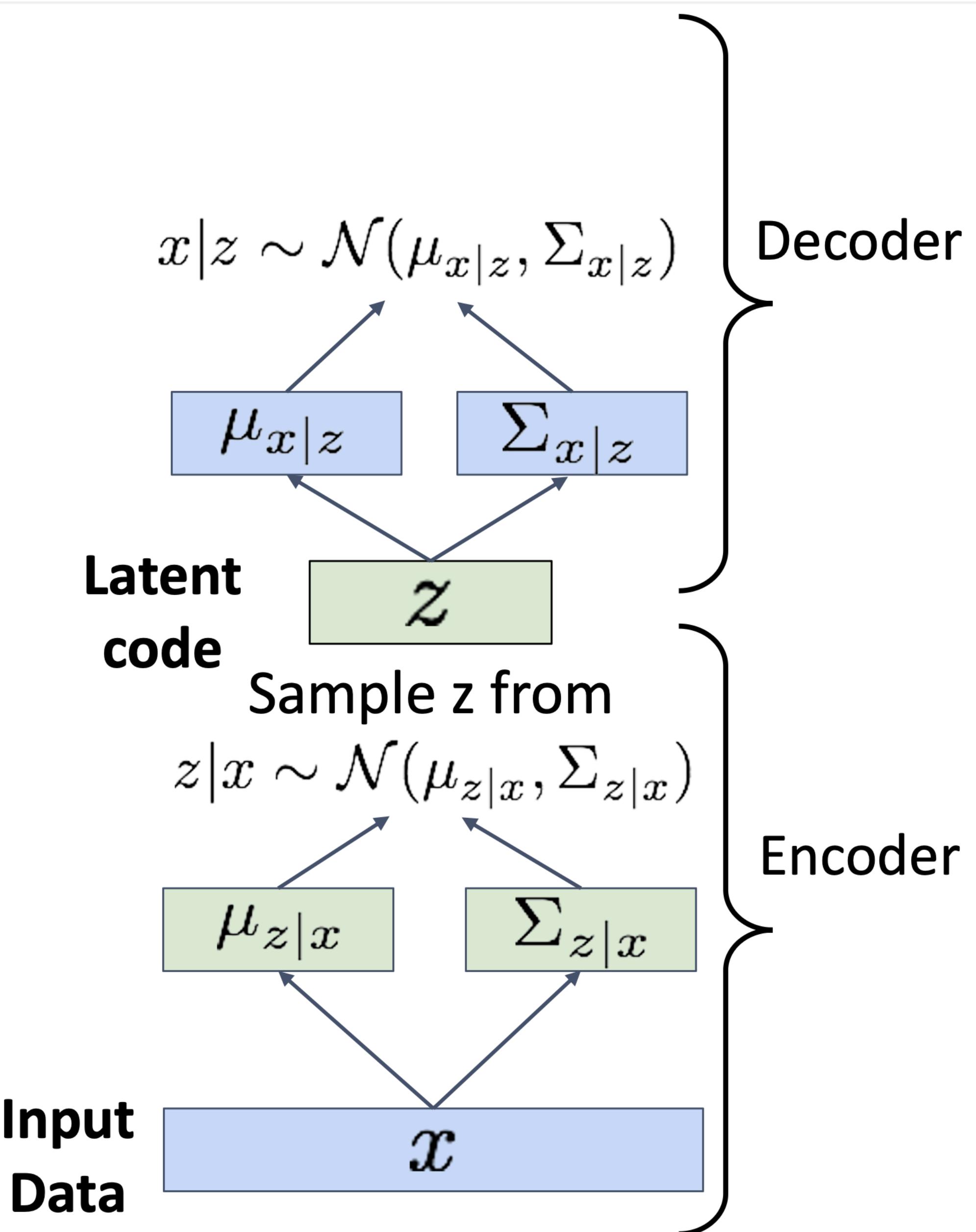


Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**
3. Sample code z from encoder output
4. Run sampled code through **decoder** to get a distribution over data samples

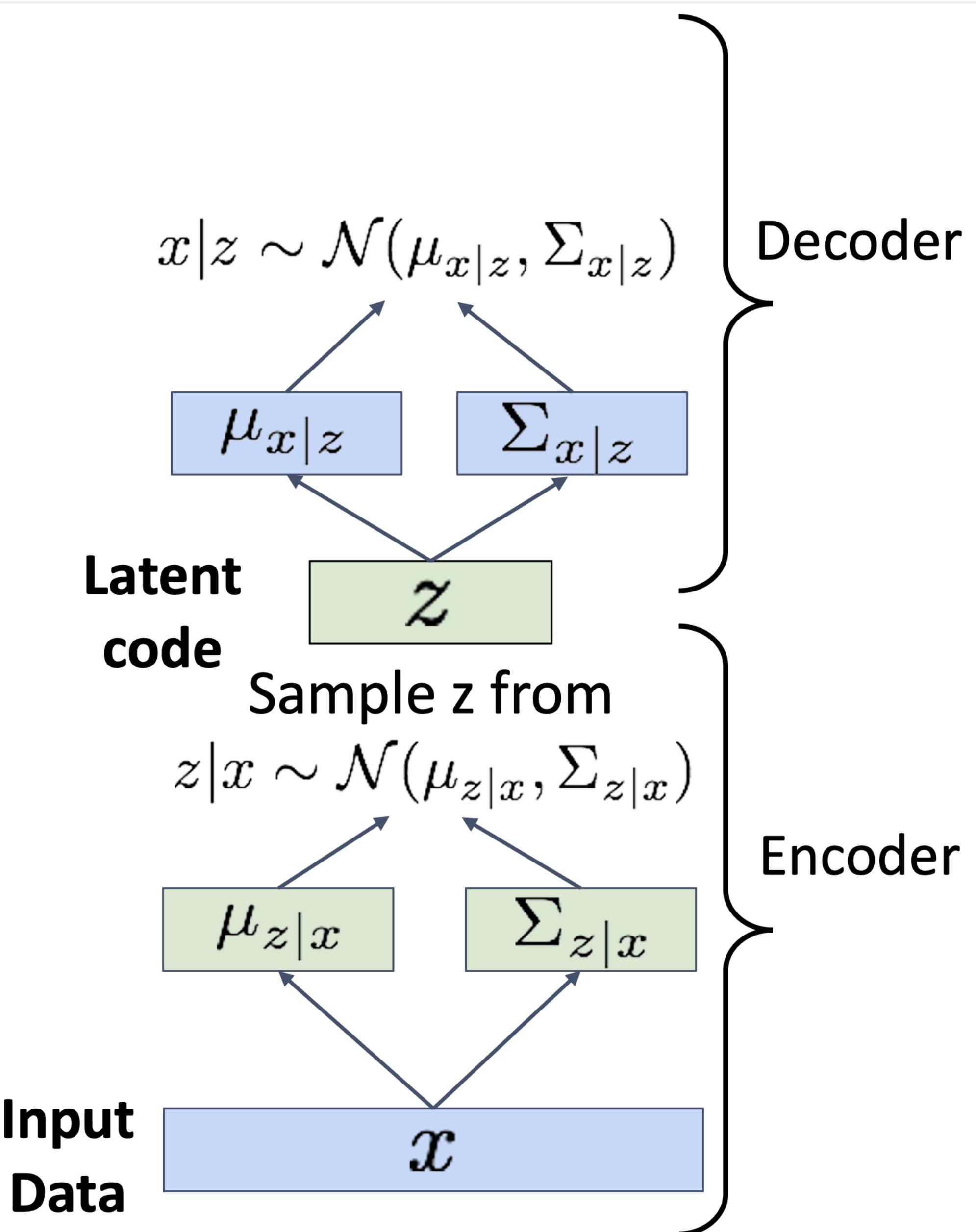


Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**
3. Sample code z from encoder output
4. Run sampled code through **decoder** to get a distribution over data samples
5. **Original input data should be likely under the distribution output from (4)!**

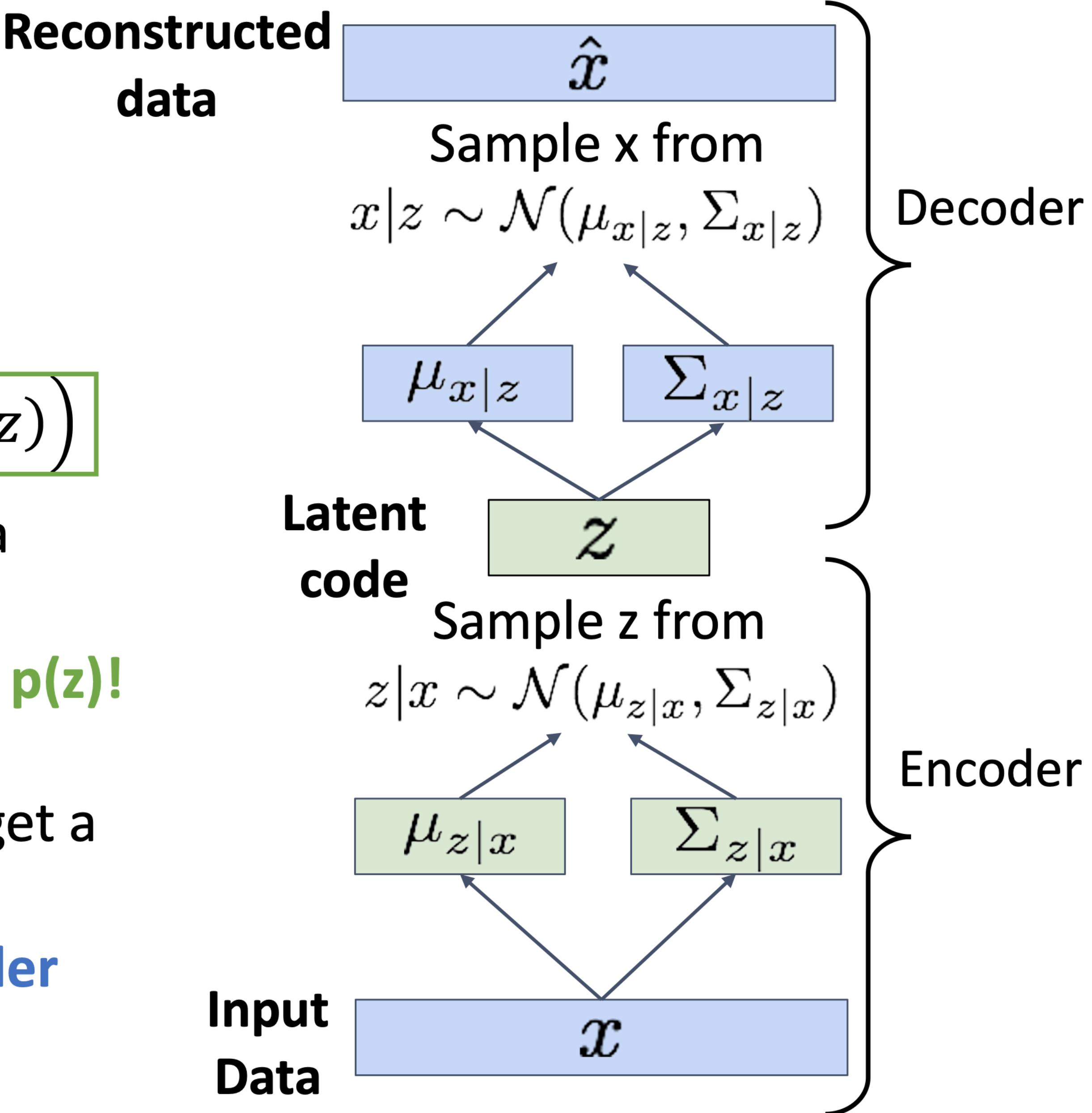


Variational Autoencoders

Train by maximizing the
variational lower bound

$$E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))$$

1. Run input data through **encoder** to get a distribution over latent codes
2. **Encoder output should match the prior $p(z)$!**
3. Sample code z from encoder output
4. Run sampled code through **decoder** to get a distribution over data samples
5. **Original input data should be likely under the distribution output from (4)!**
6. Can sample a reconstruction from (4)



Variational Autoencoders: Generating Data

After training we can
generate new data!

1. Sample z from prior $p(z)$

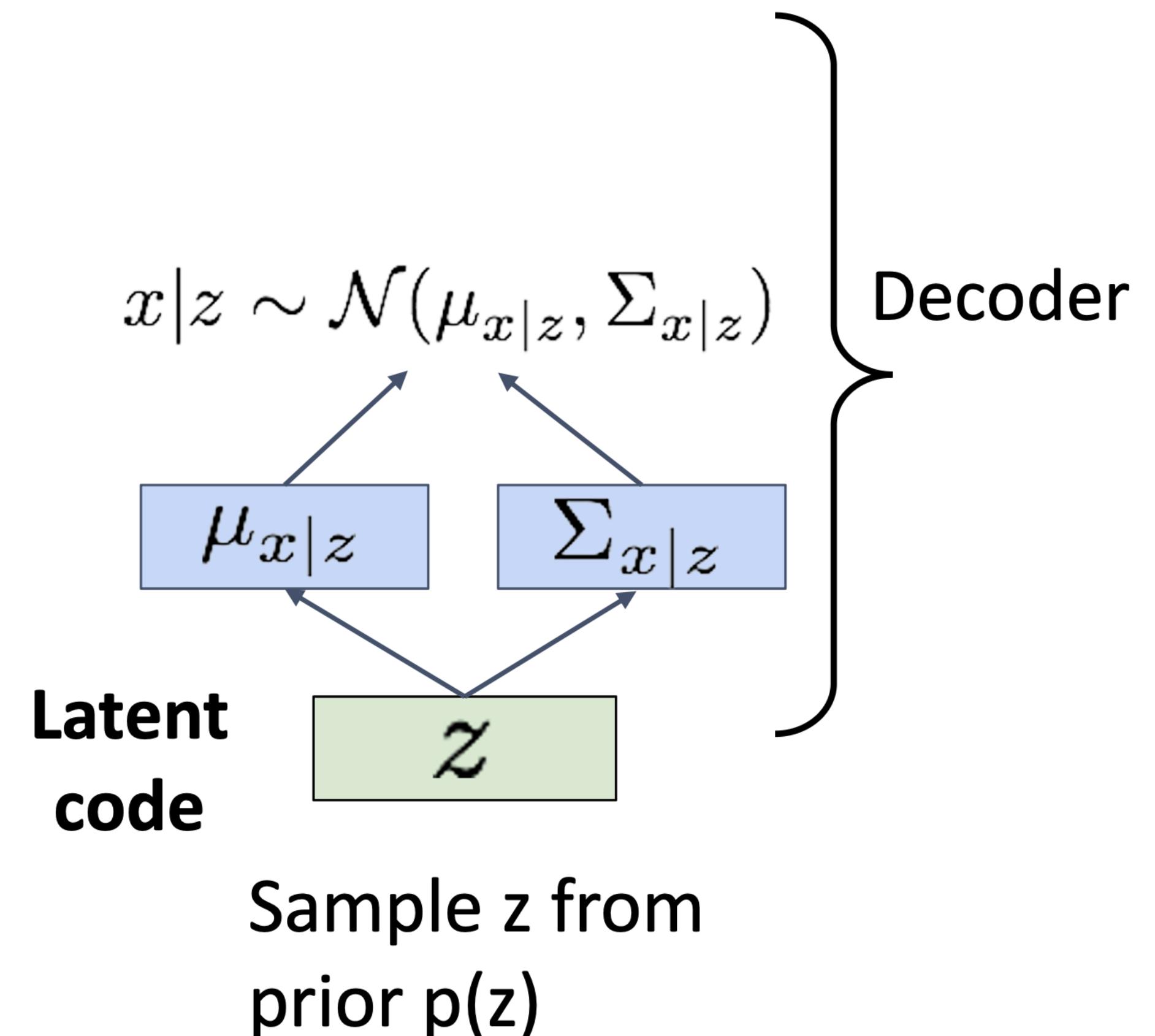
**Latent
code** 

Sample z from
prior $p(z)$

Variational Autoencoders: Generating Data

After training we can
generate new data!

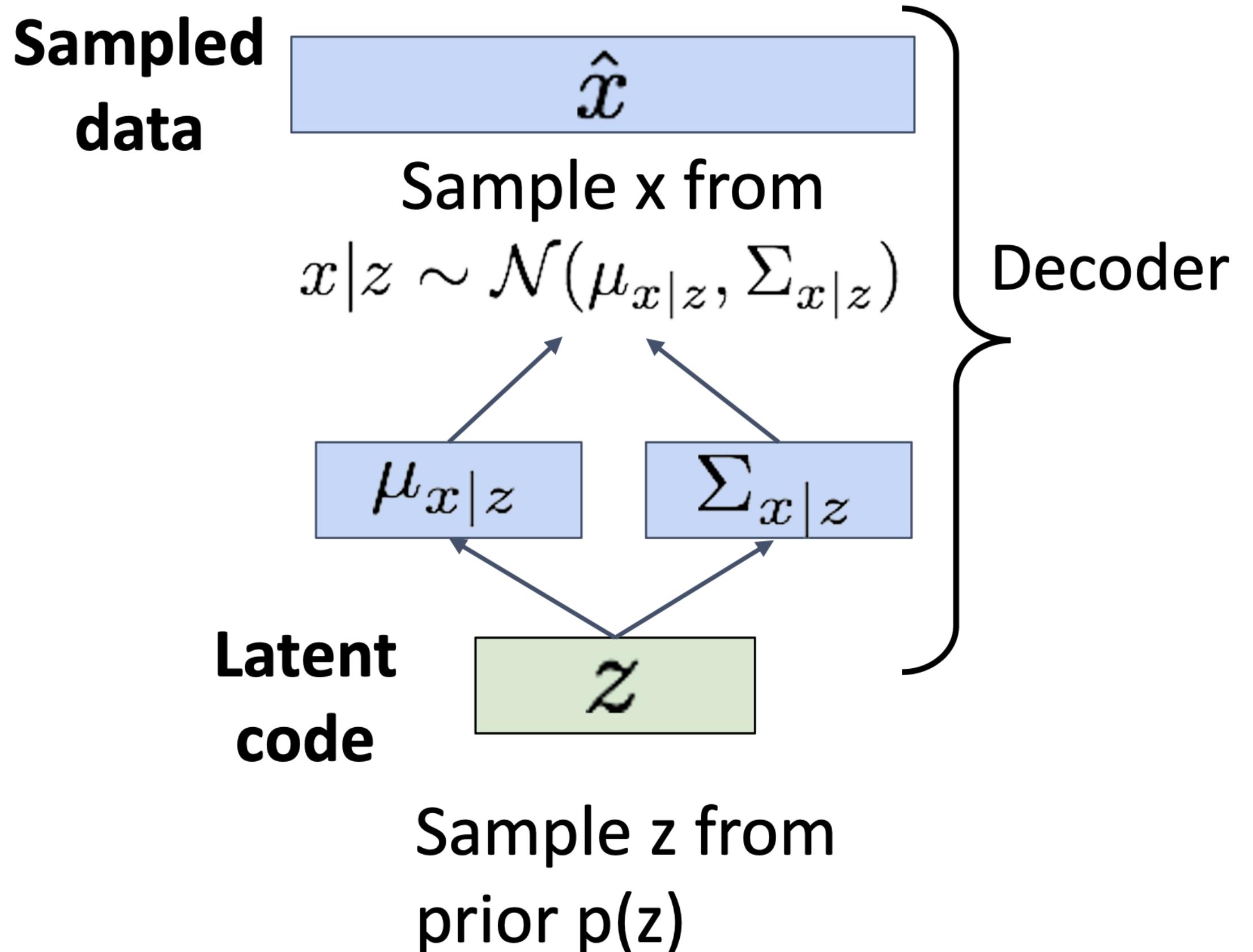
1. Sample z from prior $p(z)$
2. Run sampled z through decoder to
get distribution over data x



Variational Autoencoders: Generating Data

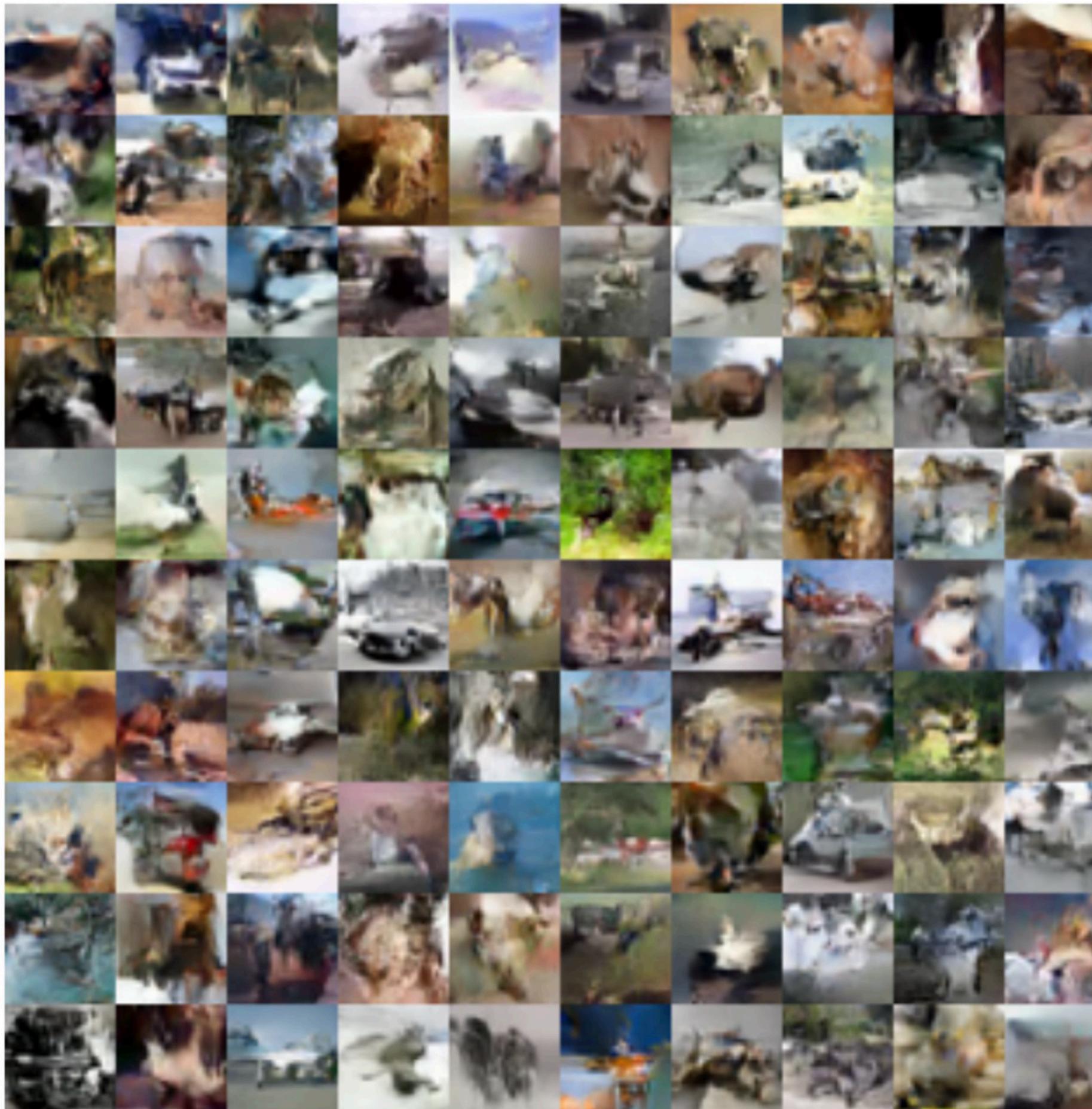
After training we can generate new data!

1. Sample z from prior $p(z)$
2. Run sampled z through decoder to get distribution over data x
3. Sample from distribution in (2) to generate data



Variational Autoencoders: Generating Data

32x32 CIFAR-10



Labeled Faces in the Wild



Figures from (L) Dirk Kingma et al. 2016; (R) Anders Larsen et al. 2017.

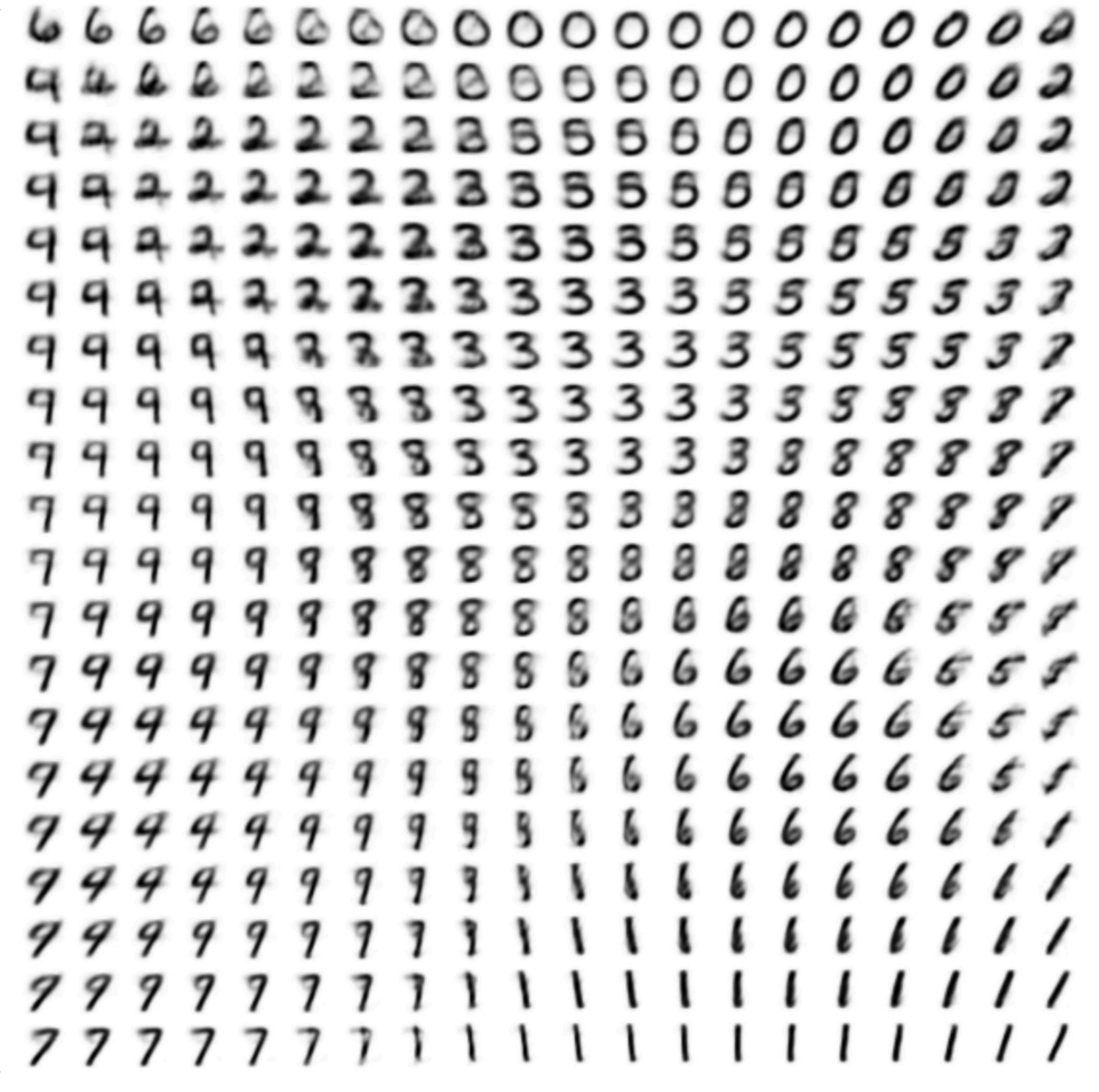
Variational Autoencoders

The diagonal prior on $p(z)$ causes dimensions of z to be independent

“Disentangling factors of variation”

Vary z_1

Vary z_2

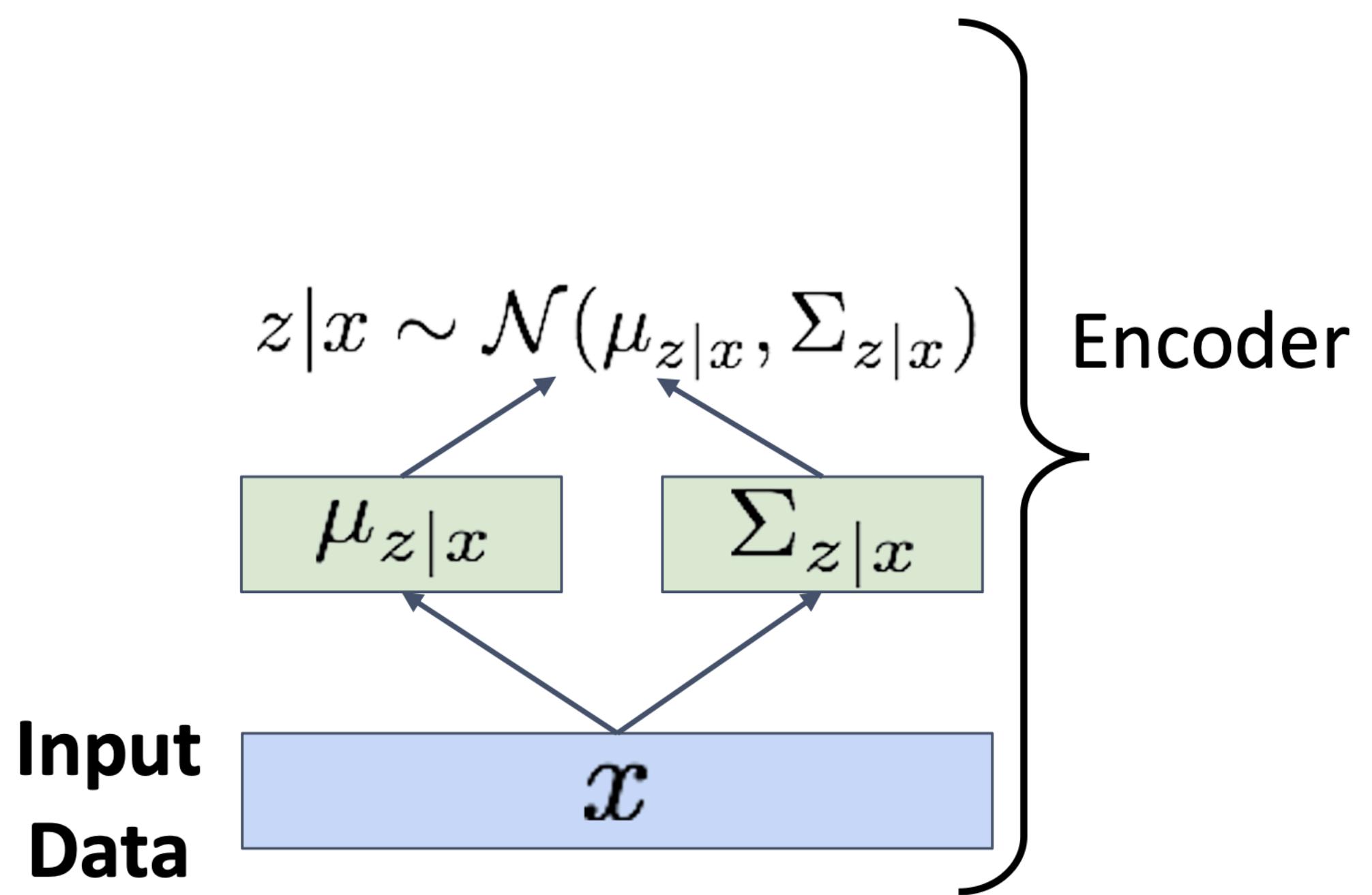


Kingma and Welling, Auto-Encoding Variational Bayes, ICLR 2014

Variational Autoencoders

After training we can **edit images**

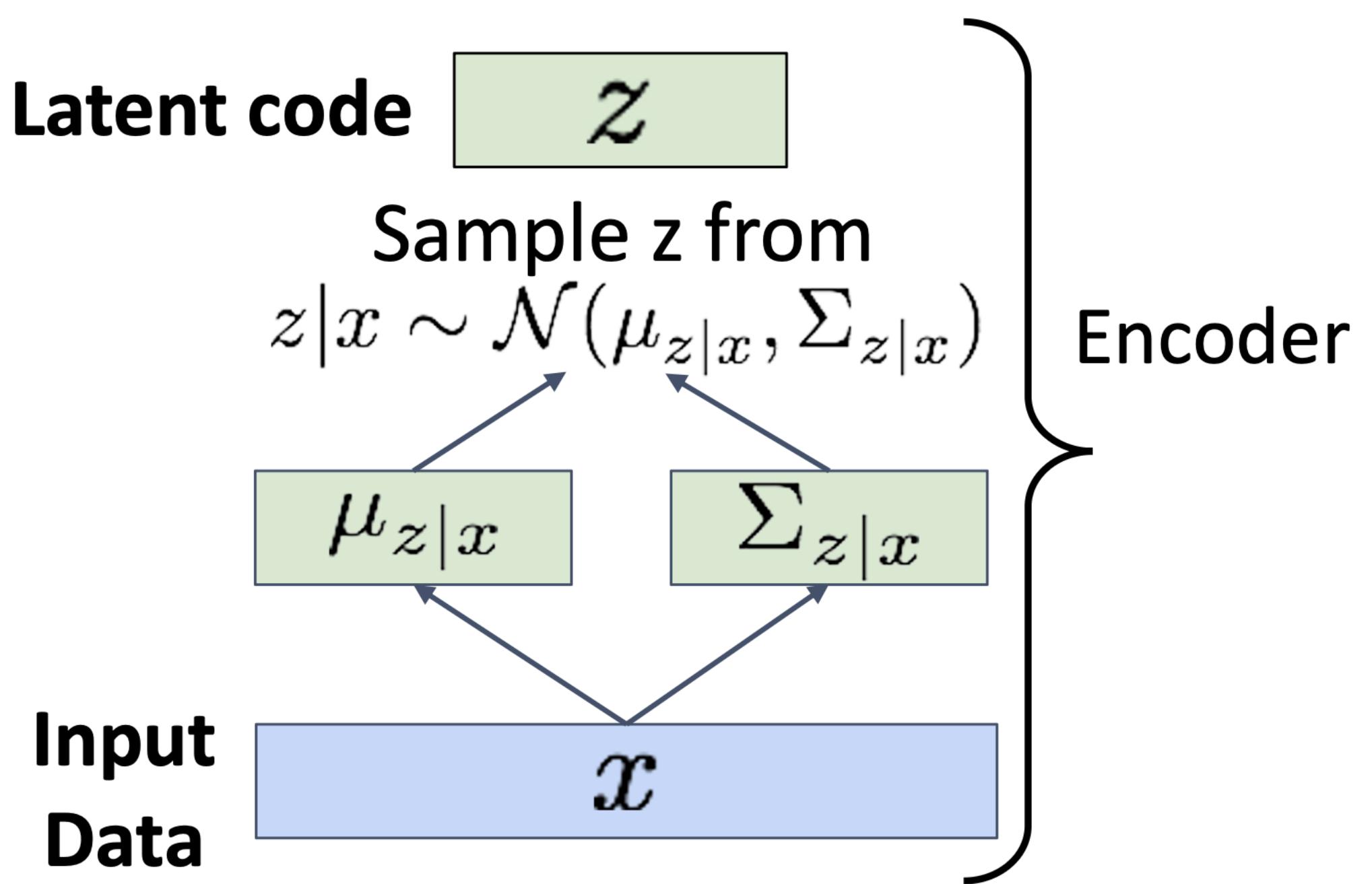
1. Run input data through **encoder** to get a distribution over latent codes



Variational Autoencoders

After training we can **edit images**

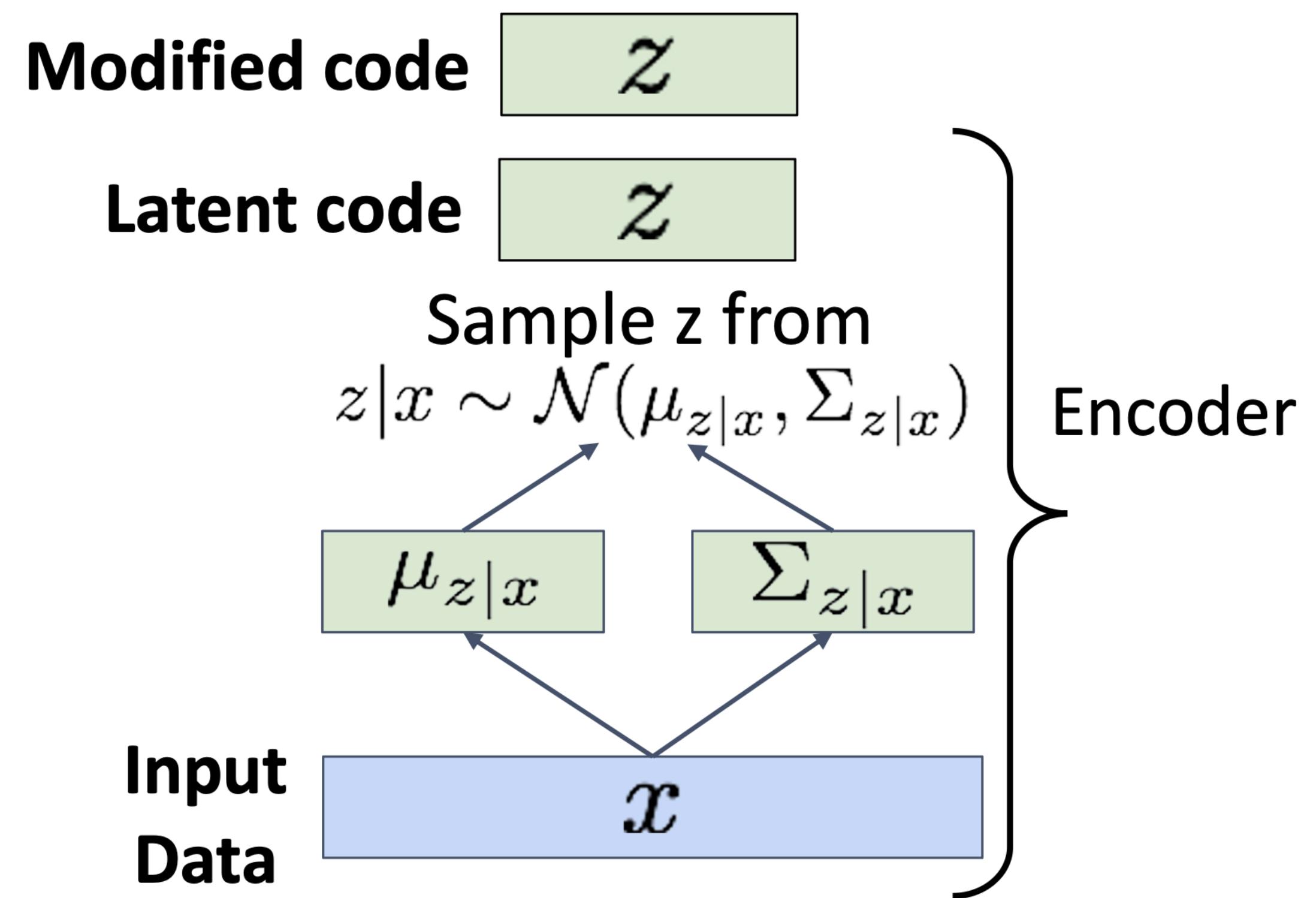
1. Run input data through **encoder** to get a distribution over latent codes
2. Sample code z from encoder output



Variational Autoencoders

After training we can **edit images**

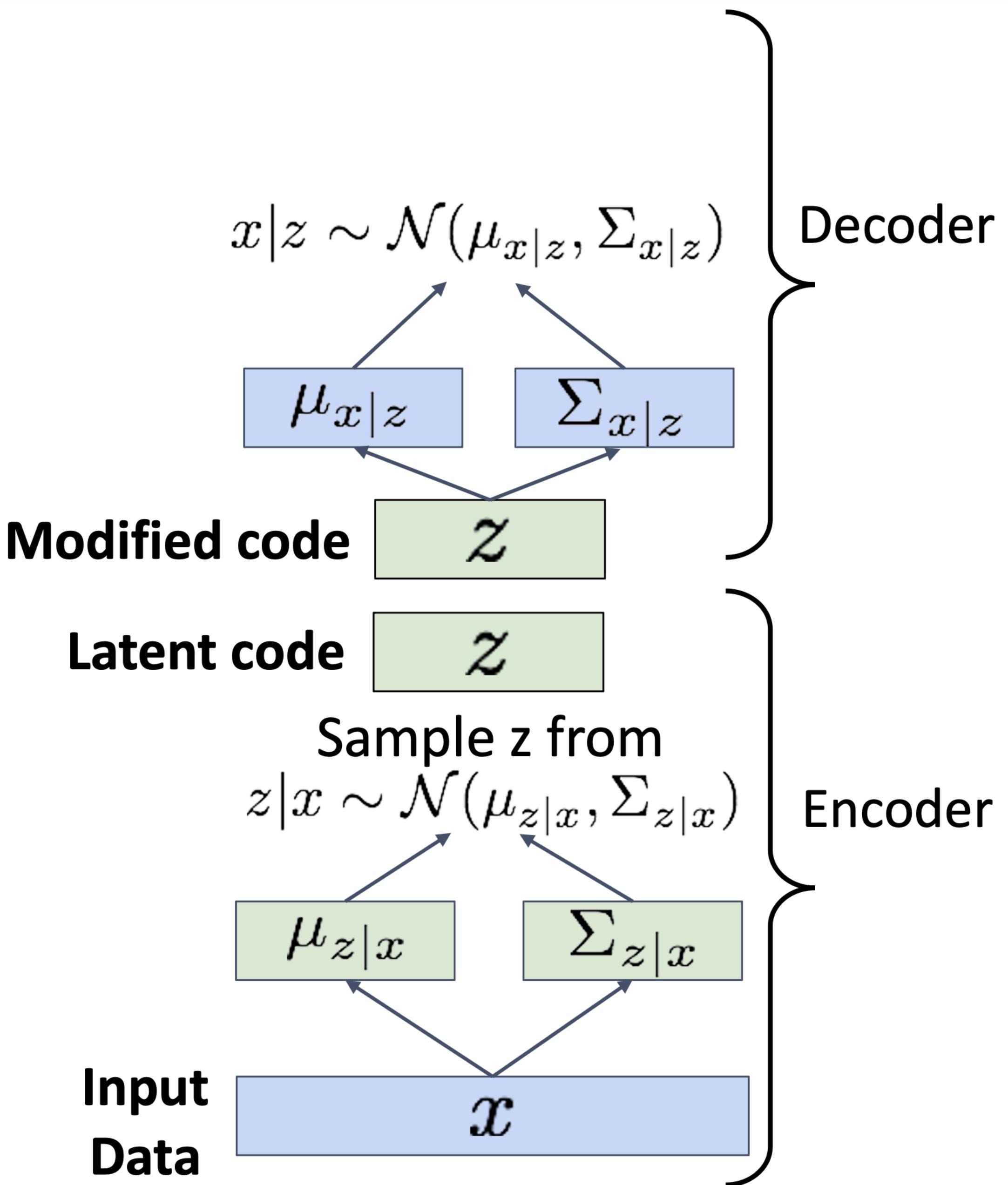
1. Run input data through **encoder** to get a distribution over latent codes
2. Sample code z from encoder output
3. Modify some dimensions of sampled code



Variational Autoencoders

After training we can **edit images**

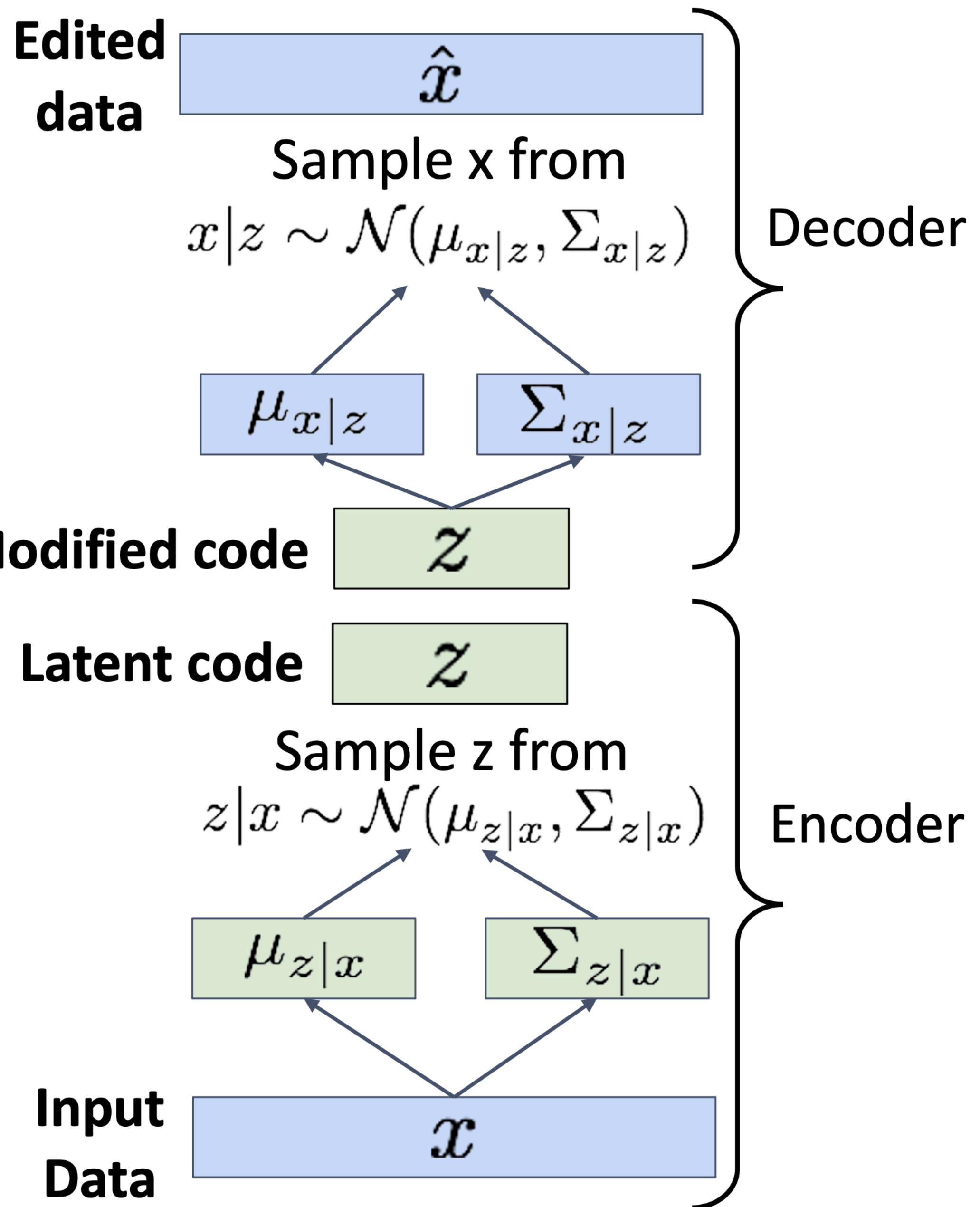
1. Run input data through **encoder** to get a distribution over latent codes
2. Sample code z from encoder output
3. Modify some dimensions of sampled code
4. Run modified z through **decoder** to get a distribution over data sample



Variational Autoencoders

After training we can **edit images**

1. Run input data through **encoder** to get a distribution over latent codes
2. Sample code z from encoder output
3. Modify some dimensions of sampled code
4. Run modified z through **decoder** to get a distribution over data samples
5. Sample new data from (4)



Variational Autoencoders

The diagonal prior on $p(z)$ causes dimensions of z to be independent

“Disentangling factors of variation”

Degree of smile

Vary z_1

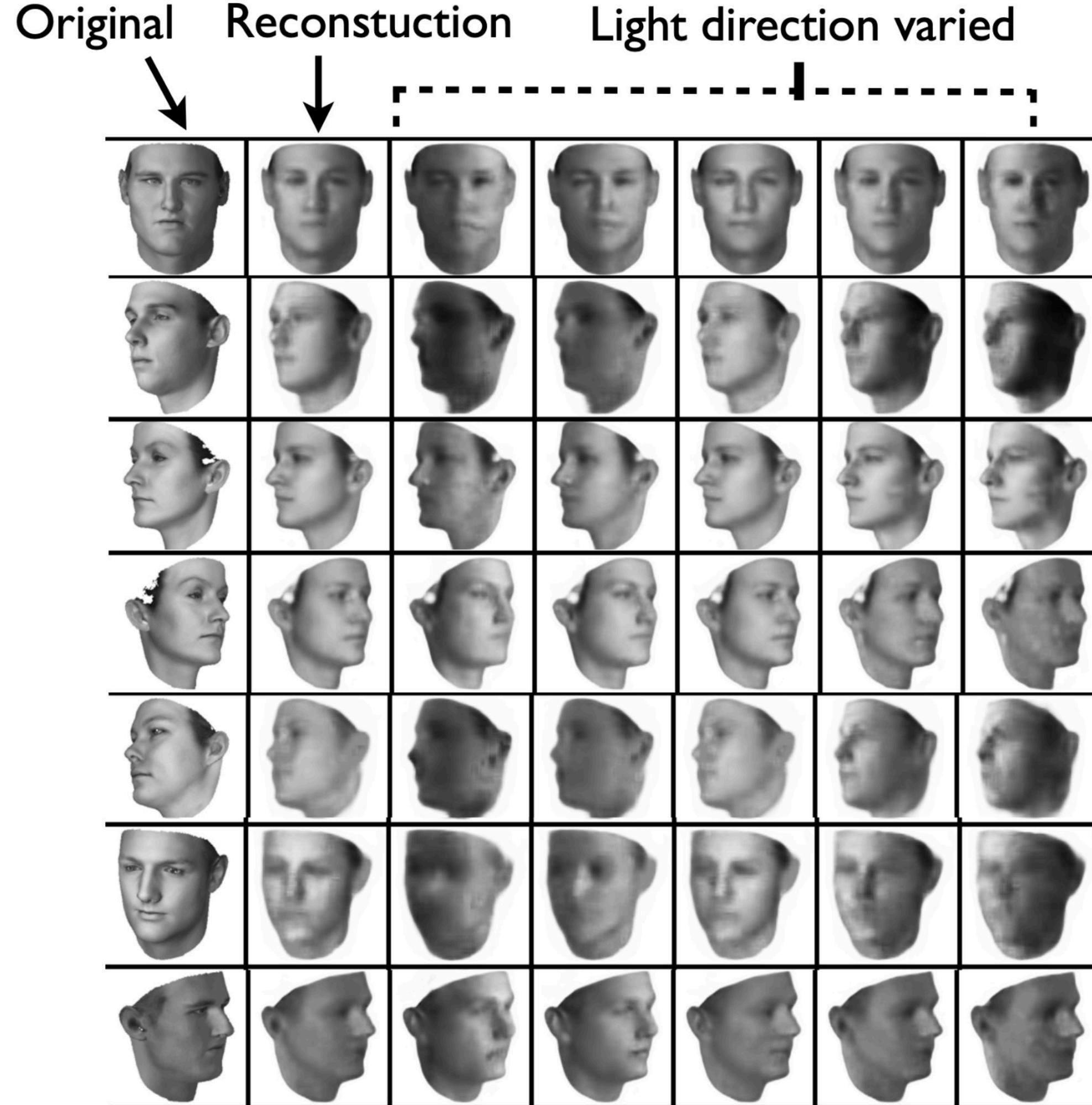
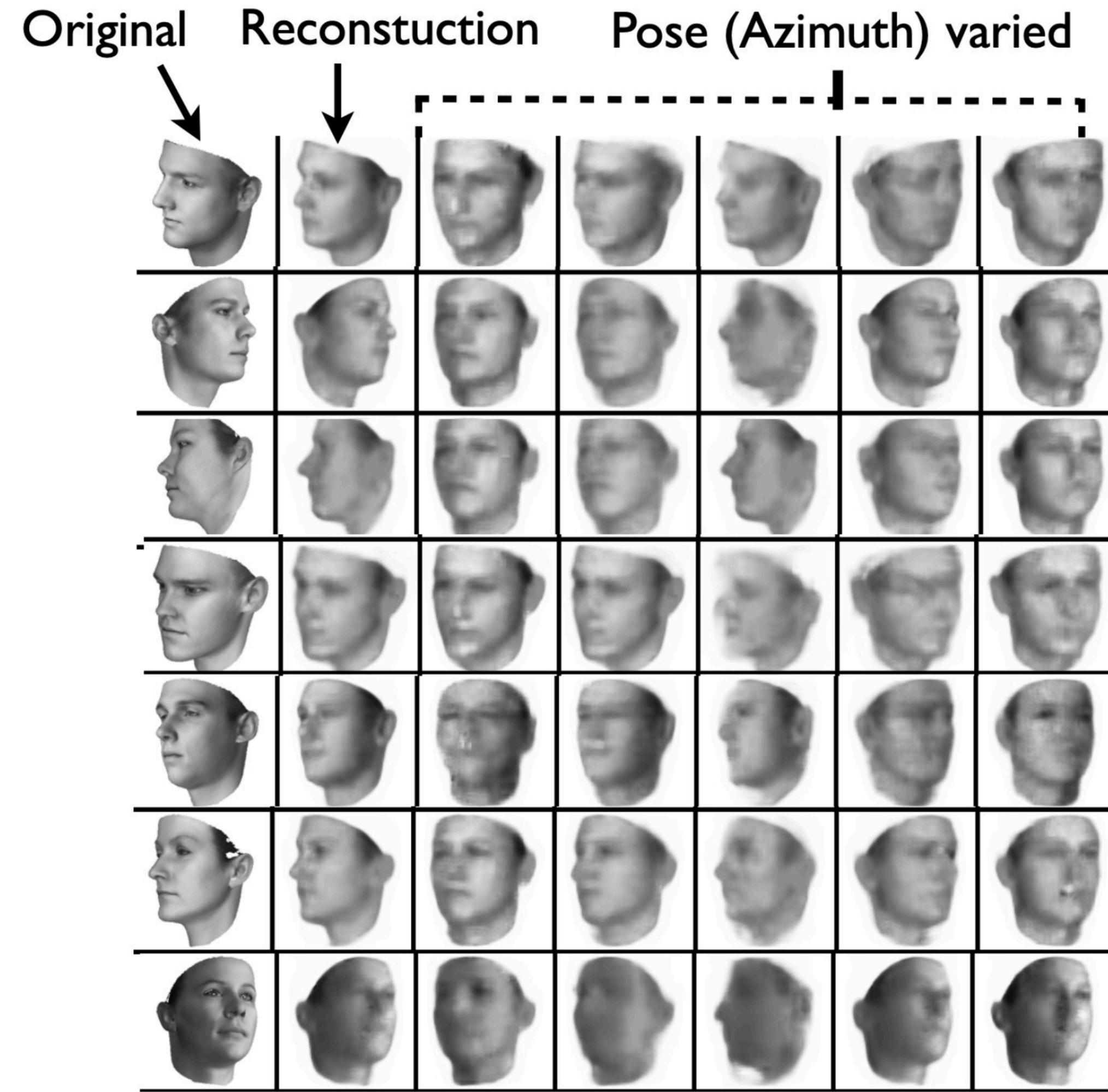
Head pose

Vary z_2



Kingma and Welling, Auto-Encoding Variational Bayes, ICLR 2014

Variational Autoencoders: Image Editing



Kulkarni et al, "Deep Convolutional Inverse Graphics Networks", NeurIPS 2014

Variational Autoencoder: Summary

Probabilistic spin to traditional autoencoders => allows generating data

Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Allows inference of $q(z|x)$, can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

Active areas of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian, e.g., Gaussian Mixture Models (GMMs)
- Incorporating structure in latent variables, e.g., Categorical Distributions

Reparametrization trick

For example, a common choice of the form of $q_\phi(\mathbf{z}|\mathbf{x})$ is a multivariate Gaussian with a diagonal covariance structure:

$$\begin{aligned}\mathbf{z} &\sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)} \mathbf{I}) \\ \mathbf{z} &= \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})\end{aligned}\quad ; \text{ Reparameterization trick.}$$

where \odot refers to element-wise product.

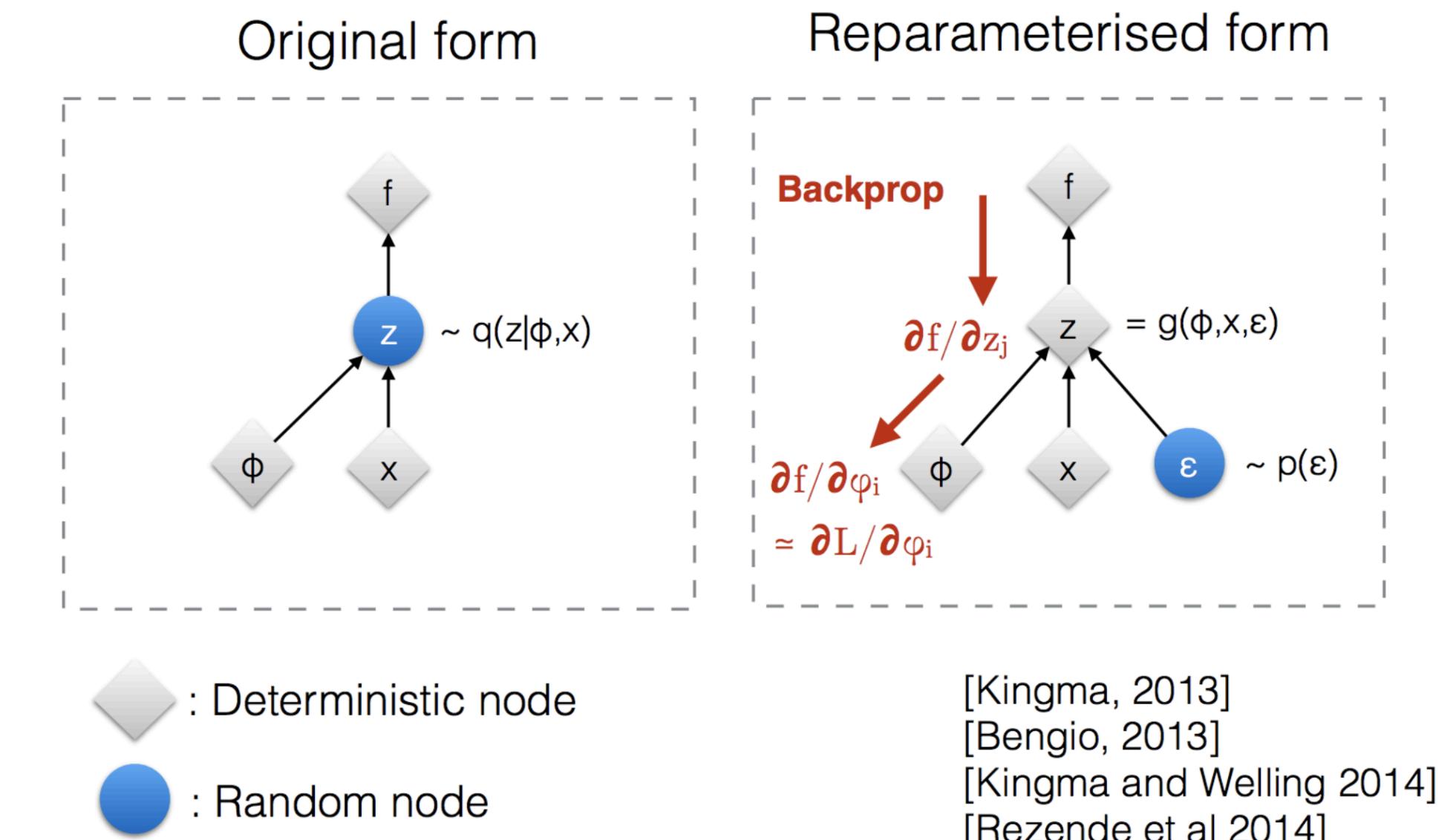


Fig. 8. Illustration of how the reparameterization trick makes the \mathbf{z} sampling process trainable.(Image source: Slide 12 in Kingma's NIPS 2015 workshop talk)

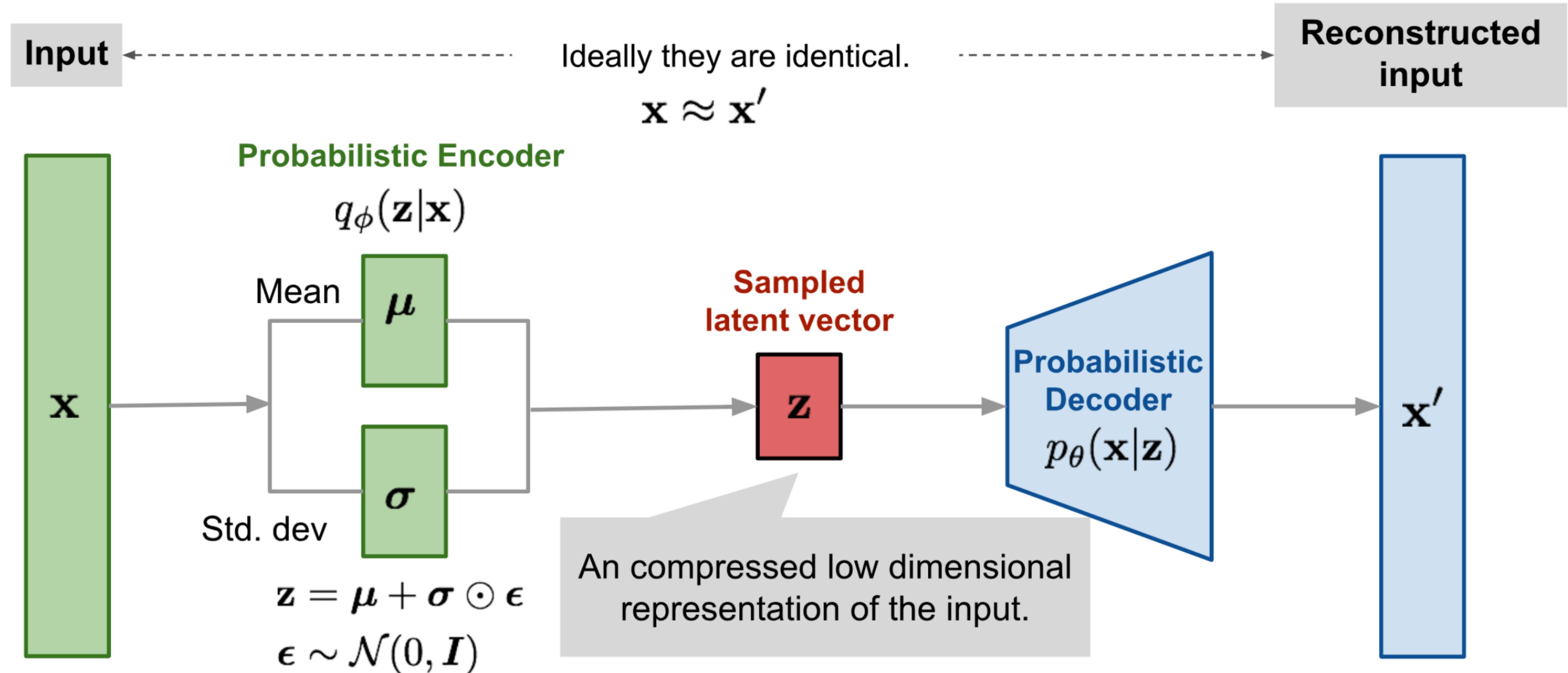
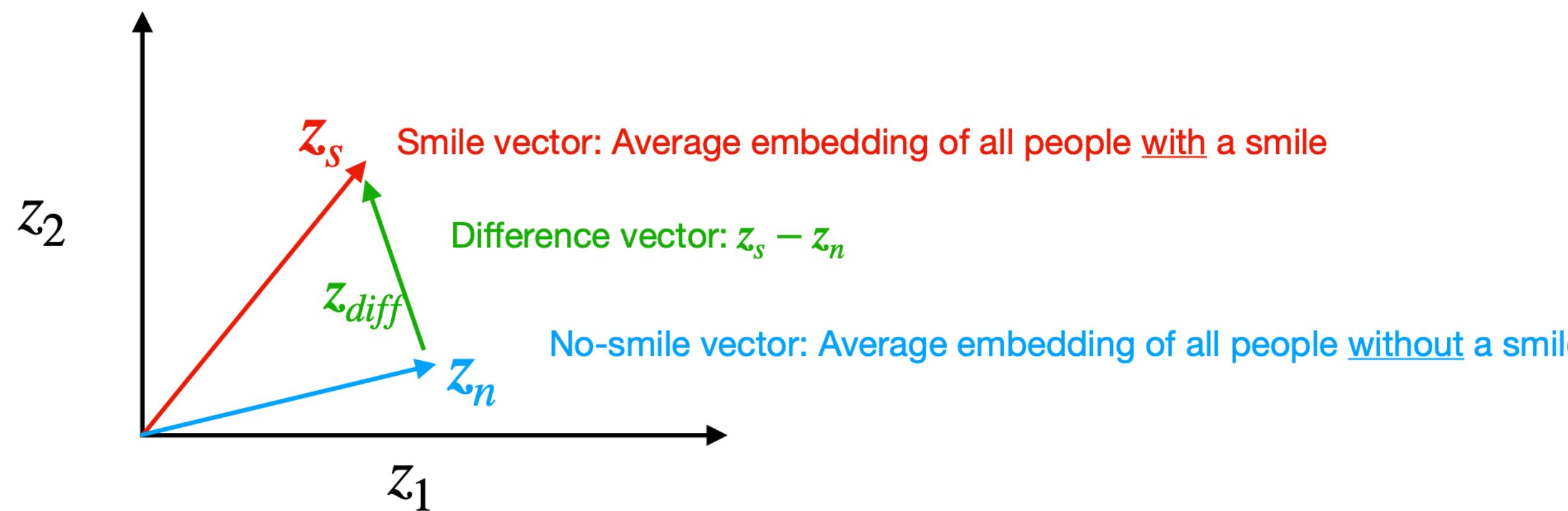


Fig. 9. Illustration of variational autoencoder model with the multivariate Gaussian assumption.

Latent space arithmetic

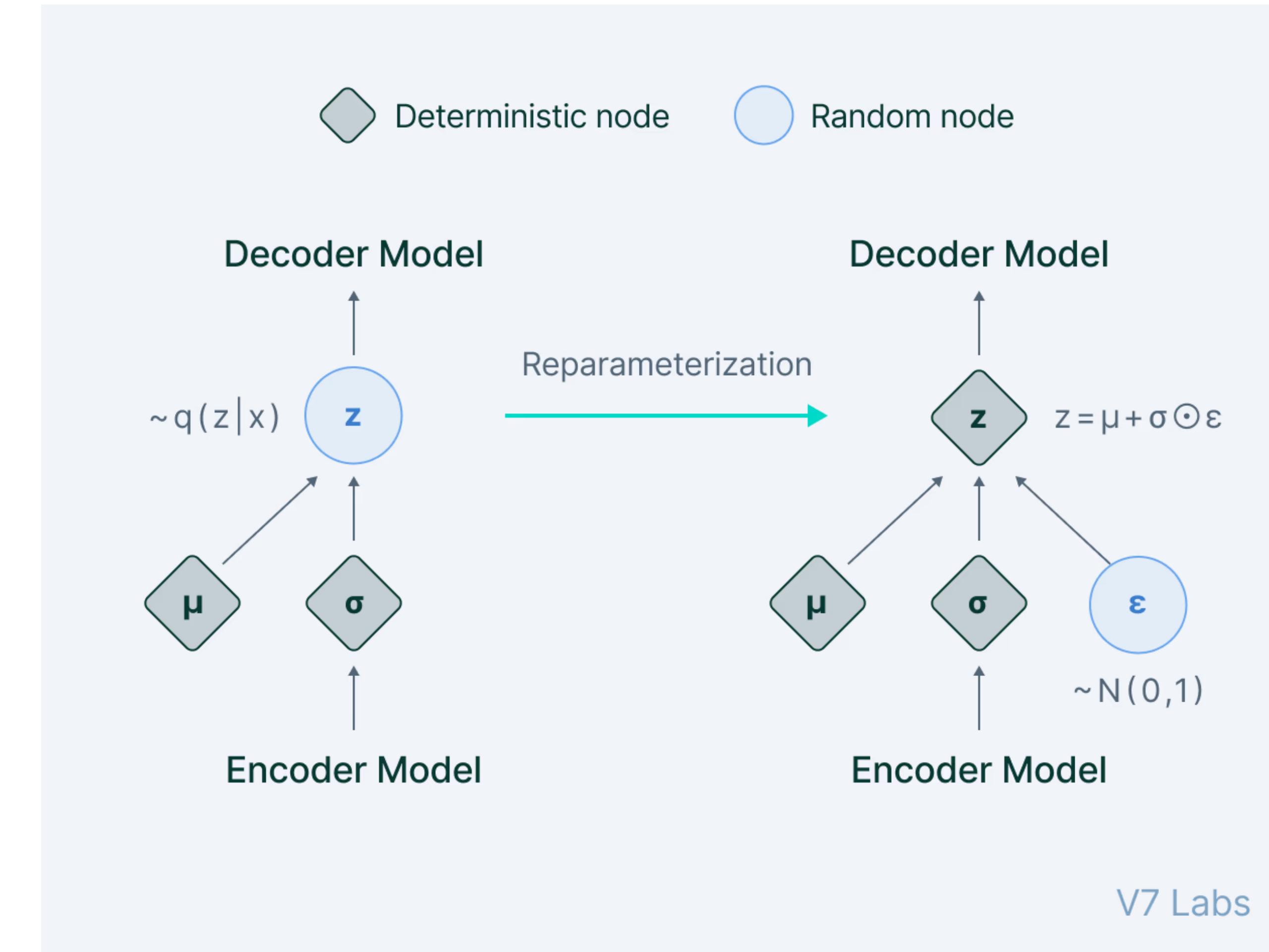


E.g., we can give a sad person a smile by

- $z_{new} = z_{orig} + \alpha \cdot z_{diff}$

Variational Autoencoders

Reparametrization trick

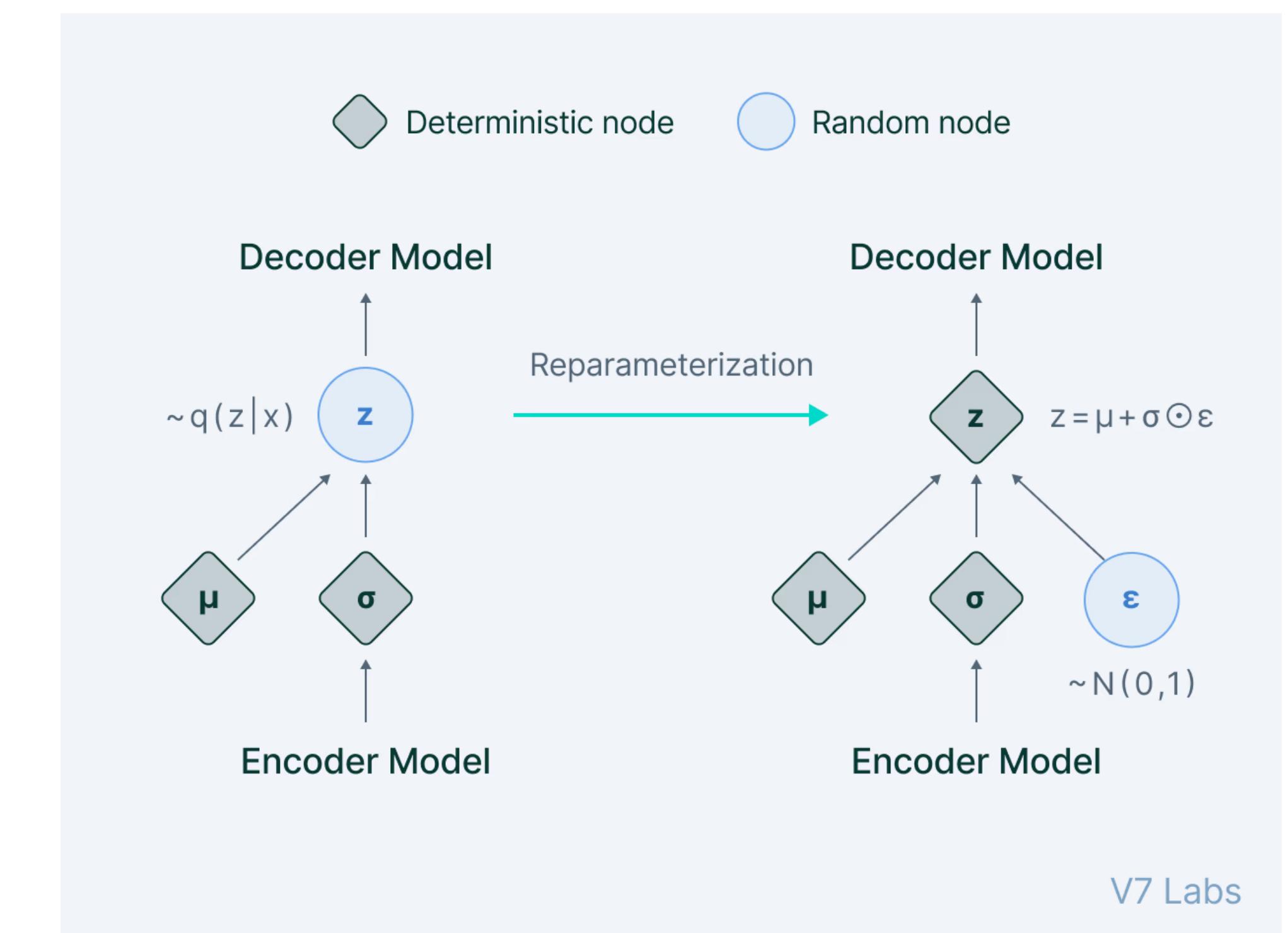


Variational Autoencoders (part 2)

Reparametrization trick

```
class VariationalEncoder(nn.Module):
    def __init__(self, latent_dims):
        super(VariationalEncoder, self).__init__()
        self.linear1 = nn.Linear(784, 512)
        self.mu_layer = nn.Linear(512, latent_dims)
        self.logvar_layer = nn.Linear(512, latent_dims)

    def forward(self, x):
        x = torch.flatten(x, start_dim=1)
        x = F.relu(self.linear1(x))
        mu = self.mu_layer(x)
        logvar = self.logvar_layer(x)
        std = torch.exp(0.5 * logvar)
        eps = torch.randn_like(std)
        z = mu + std * eps
        return z, mu, logvar
```



Variational Autoencoders

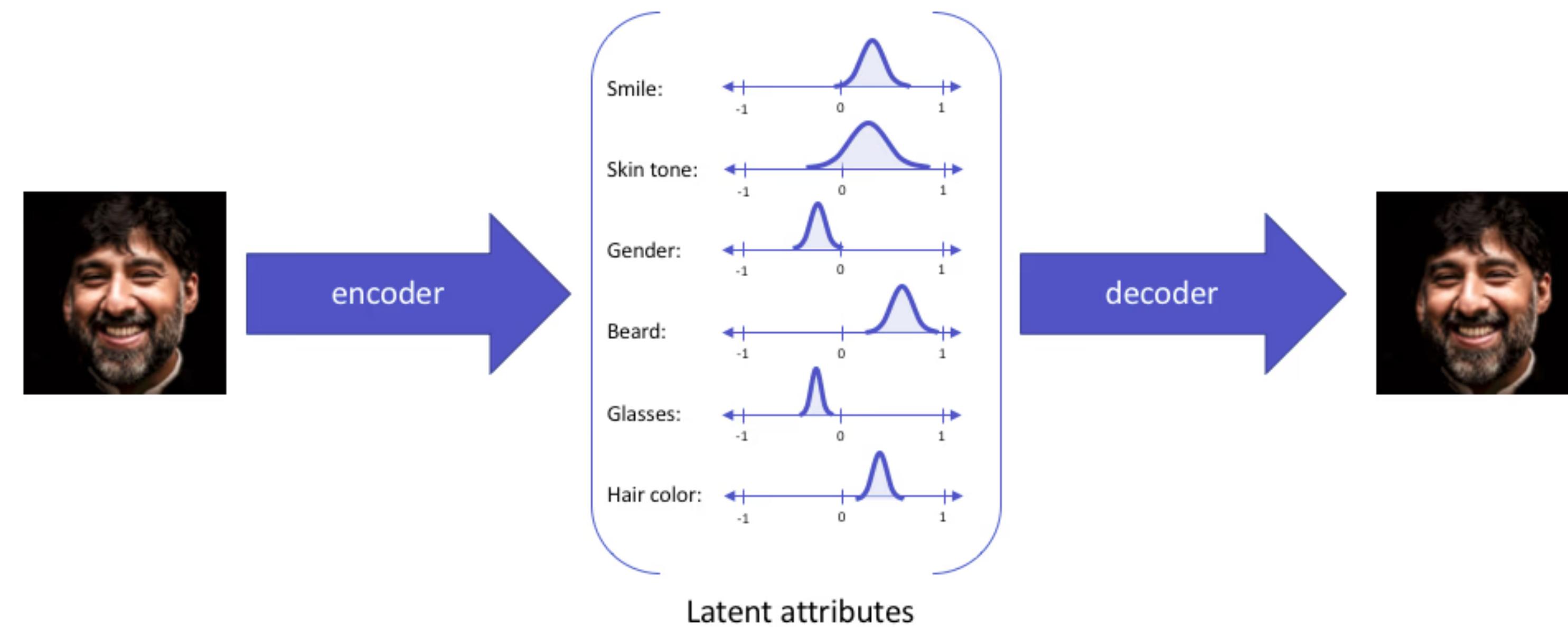
Demo code

see <https://avandekleut.github.io/vae/>

Who is right?

<https://chatgpt.com/share/6834a21f-7e84-8007-bf80-a1b61c6d94ef>

Disentangled Representation Learning



source: <https://www.v7labs.com/blog/autoencoders-guide>

β -VAE

Published as a conference paper at ICLR 2017

β -VAE: LEARNING BASIC VISUAL CONCEPTS WITH A CONSTRAINED VARIATIONAL FRAMEWORK

Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot,
Matthew Botvinick, Shakir Mohamed, Alexander Lerchner
Google DeepMind
`{irinah, lmatthey, arkap, cpburgess, glorotx,
botvinick, shakir, lerchner}@google.com`

$$\mathcal{F}(\theta, \phi, \beta; \mathbf{x}, \mathbf{z}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

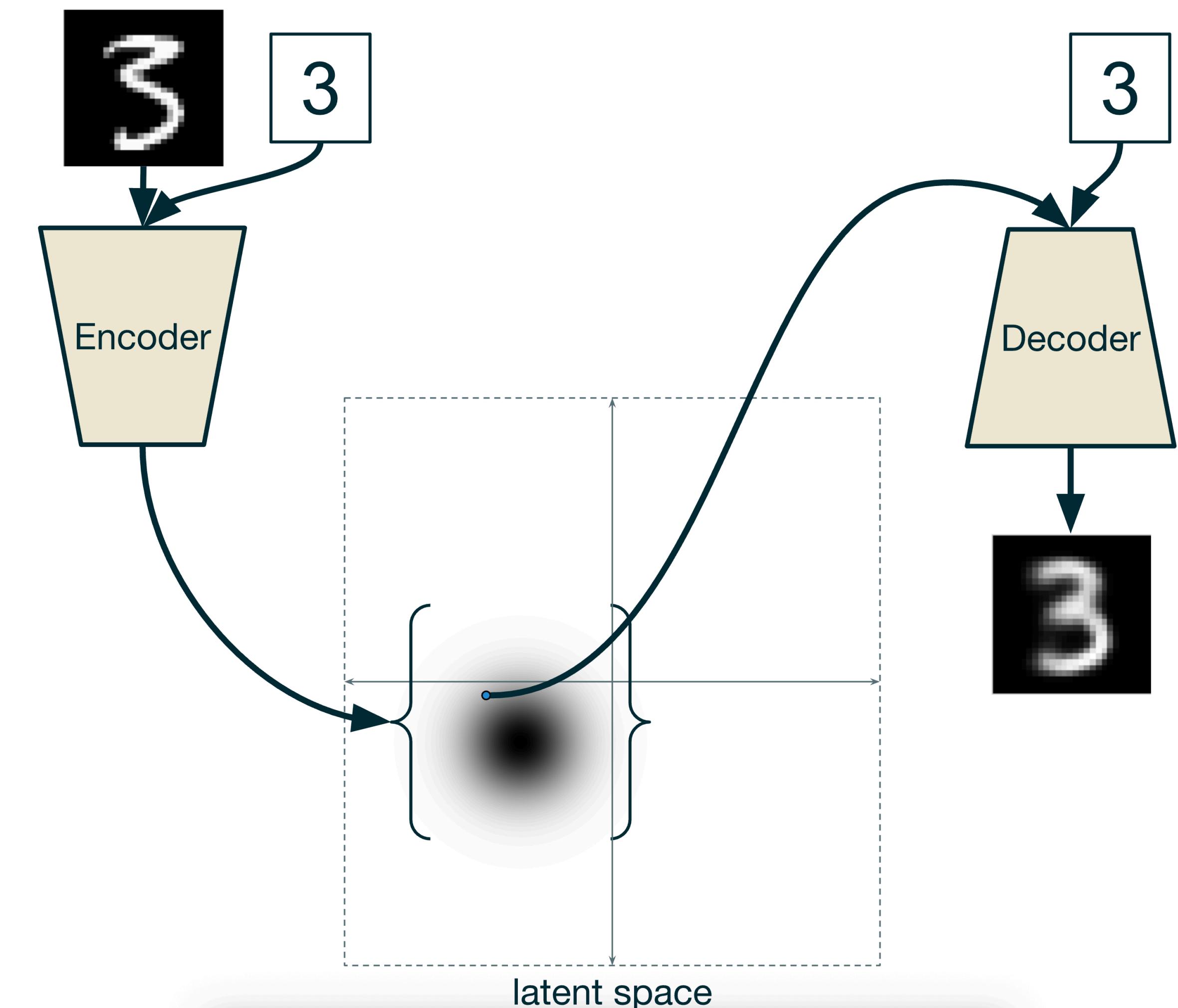
ABSTRACT

Learning an interpretable factorised representation of the independent data generative factors of the world without supervision is an important precursor for the development of artificial intelligence that is able to learn and reason in the same way that humans do. We introduce β -VAE, a new state-of-the-art framework for automated discovery of interpretable factorised latent representations from raw image data in a completely unsupervised manner. Our approach is a modification of the variational autoencoder (VAE) framework. We introduce an adjustable hyperparameter β that balances latent channel capacity and independence constraints with reconstruction accuracy. We demonstrate that β -VAE with appropriately tuned $\beta > 1$ qualitatively outperforms VAE ($\beta = 1$), as well as state of the art unsupervised (InfoGAN) and semi-supervised (DC-IGN) approaches to disentangled factor learning on a variety of datasets (*celeba*, *faces* and *chairs*). Furthermore, we devise a protocol to quantitatively compare the degree of disentanglement learnt by different models, and show that our approach also significantly outperforms all baselines quantitatively. Unlike InfoGAN, β -VAE is stable to train, makes few assumptions about the data and relies on tuning a single hyperparameter β , which can be directly optimised through a hyperparameter search using weakly labelled data or through heuristic visual inspection for purely unsupervised data.

See:

<https://lilianweng.github.io/posts/2018-08-12-vae/#beta-vae>

Conditional VAE



Conditional VAE

* demo code

