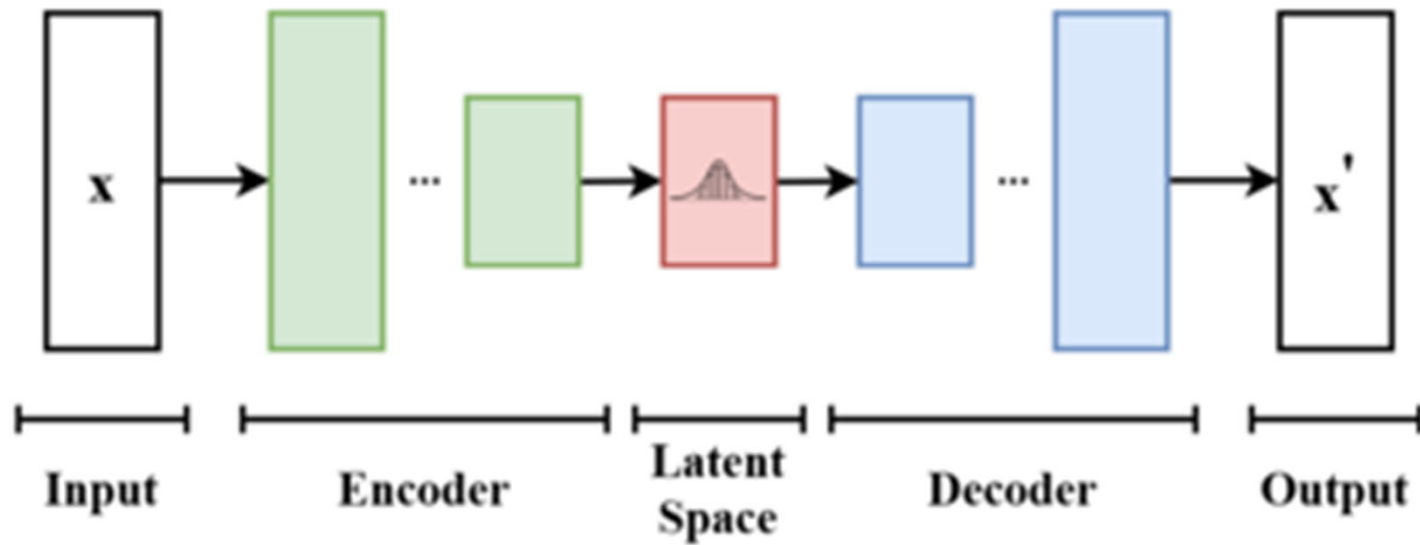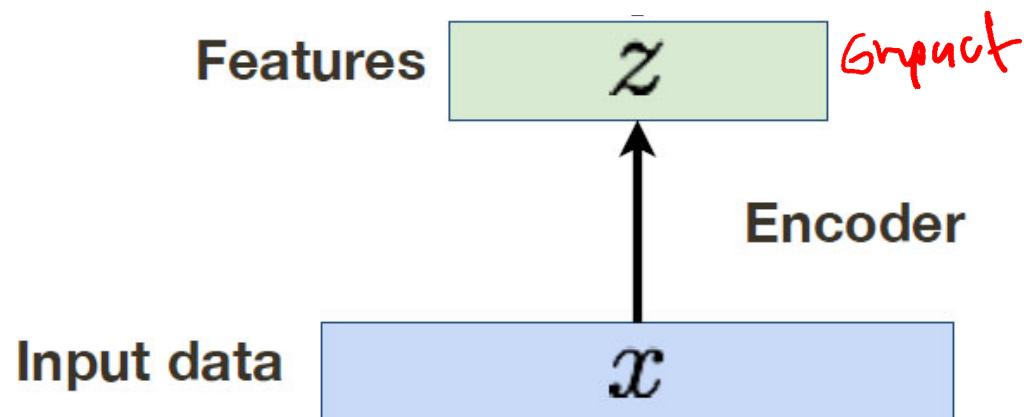# Variational Autoencoders



Evangelos Kalogerakis

# Generative models
(discussed in this course)

- Autoregressive models [done]
- **Variational Autoencoders**
- Diffusion Models [TODO next week]
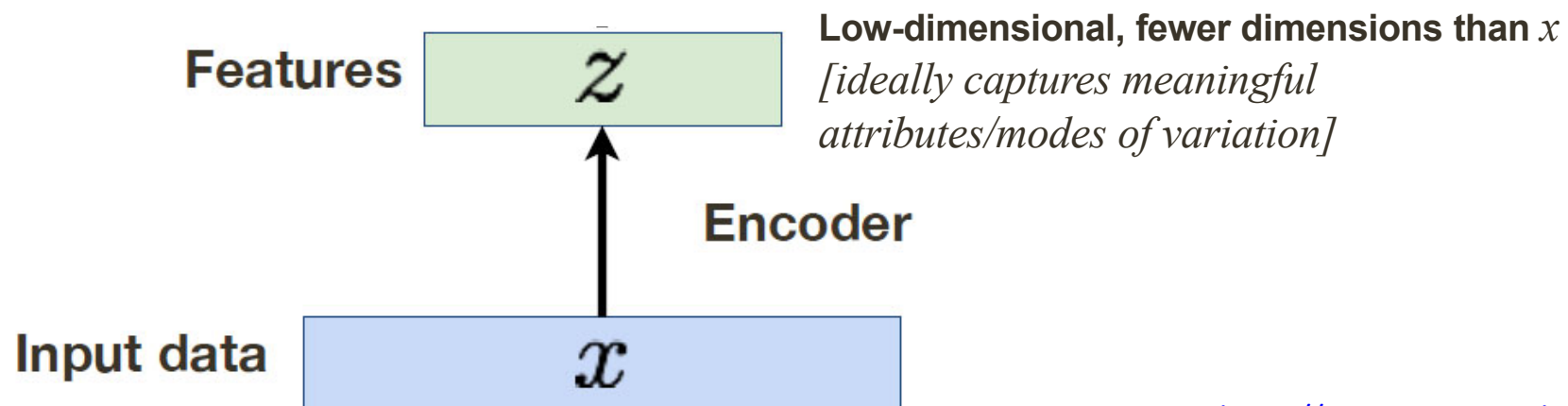- Generative Adversarial Networks [lower priority, after Easter]

# ~~Variational~~ Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from **unlabeled** training data

Features | $z$ | Output

Encoder

Input data | $x$

# ~~Variational~~ Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from **unlabeled** training data

Features $z$

Low-dimensional, fewer dimensions than $x$
*[ideally captures meaningful attributes/modes of variation]*

Encoder

Input data $x$

# ~~Variational~~ Autoencoders

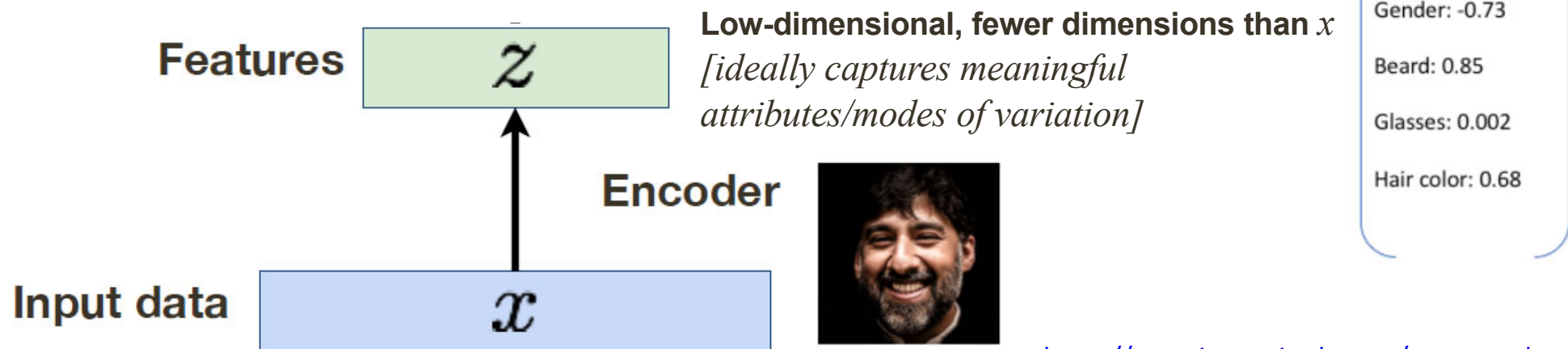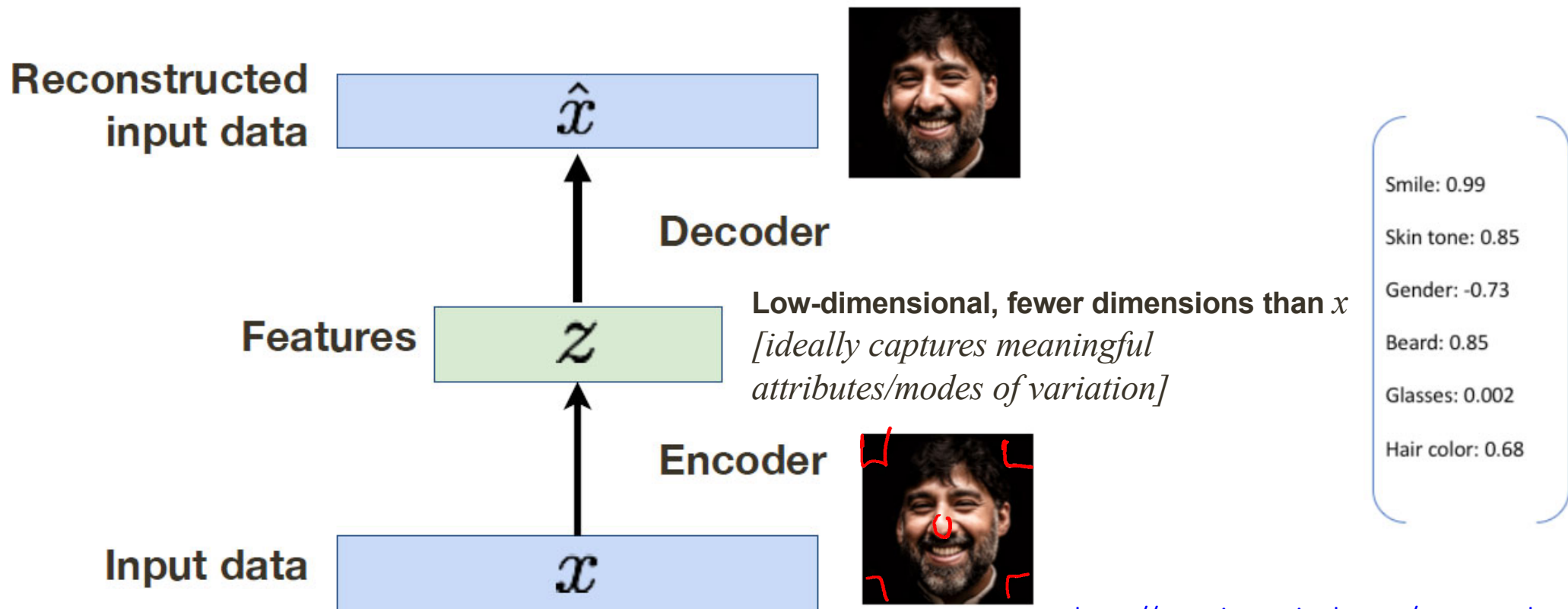Unsupervised approach for learning a lower-dimensional feature representation from **unlabeled** training data

**Features**

$z$

**Low-dimensional, fewer dimensions than** $x$
*[ideally captures meaningful attributes/modes of variation]*

**Encoder**

**Input data**

$x$

Smile: 0.99

Skin tone: 0.85

Gender: -0.73

Beard: 0.85

Glasses: 0.002

Hair color: 0.68

https://www.jeremyjordan.me/autoencoders/

# ~~Variational~~ Autoencoders

Train such that features can reconstruct original data best they can!



Reconstructed input data — $\hat{x}$

**Decoder**

Features — $z$ — **Low-dimensional, fewer dimensions than** $x$ *[ideally captures meaningful attributes/modes of variation]*

**Encoder**

Input data — $x$

Smile: 0.99

Skin tone: 0.85

Gender: -0.73

Beard: 0.85

Glasses: 0.002

Hair color: 0.68

https://www.jeremyjordan.me/autoencoders/

# Simplest Autoencoder



Input layer     Hidden layer     Output layer

$(0,0)$

$(1,1)$

$(0,2)$

$\vdots$

$\hat{x} = W' \cdot a$

"bottleneck"

$a = W \cdot x$

$256$        $256 \cdot (64 \cdot 64 \cdot 3)$

$64 \cdot 64 \cdot 3$

$\hat{x}_1 = x_1$

$\hat{x}_2 = x_2$

$\hat{x}_3 = x_3$

$\hat{x}_4 = x_4$

# Without a compact bottleneck, the network may learn an identity transformation!



Input layer

Hidden layer

Output layer

$a_1 = x_1$

$a_2 = x_2$

$\hat{x}_1 = a_1$

$\hat{x}_2 = a_2$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$

$a_1$ $a_2$ $a_3$ $a_4$ $a_5$ $a_6$

$\hat{x}_1$ $\hat{x}_2$ $\hat{x}_3$ $\hat{x}_4$ $\hat{x}_5$ $\hat{x}_6$

# ~~Variational~~ Autoencoders

Train such that features can reconstruct original data best they can!

L2 Loss function:  *Self-supervision*

$$\|\vec{x} - \hat{\vec{x}}\|^2$$

**Reconstructed input data**  $\hat{x}$

**Decoder**

**Features**  $z$  **Low-dimensional, fewer dimensions than** $x$
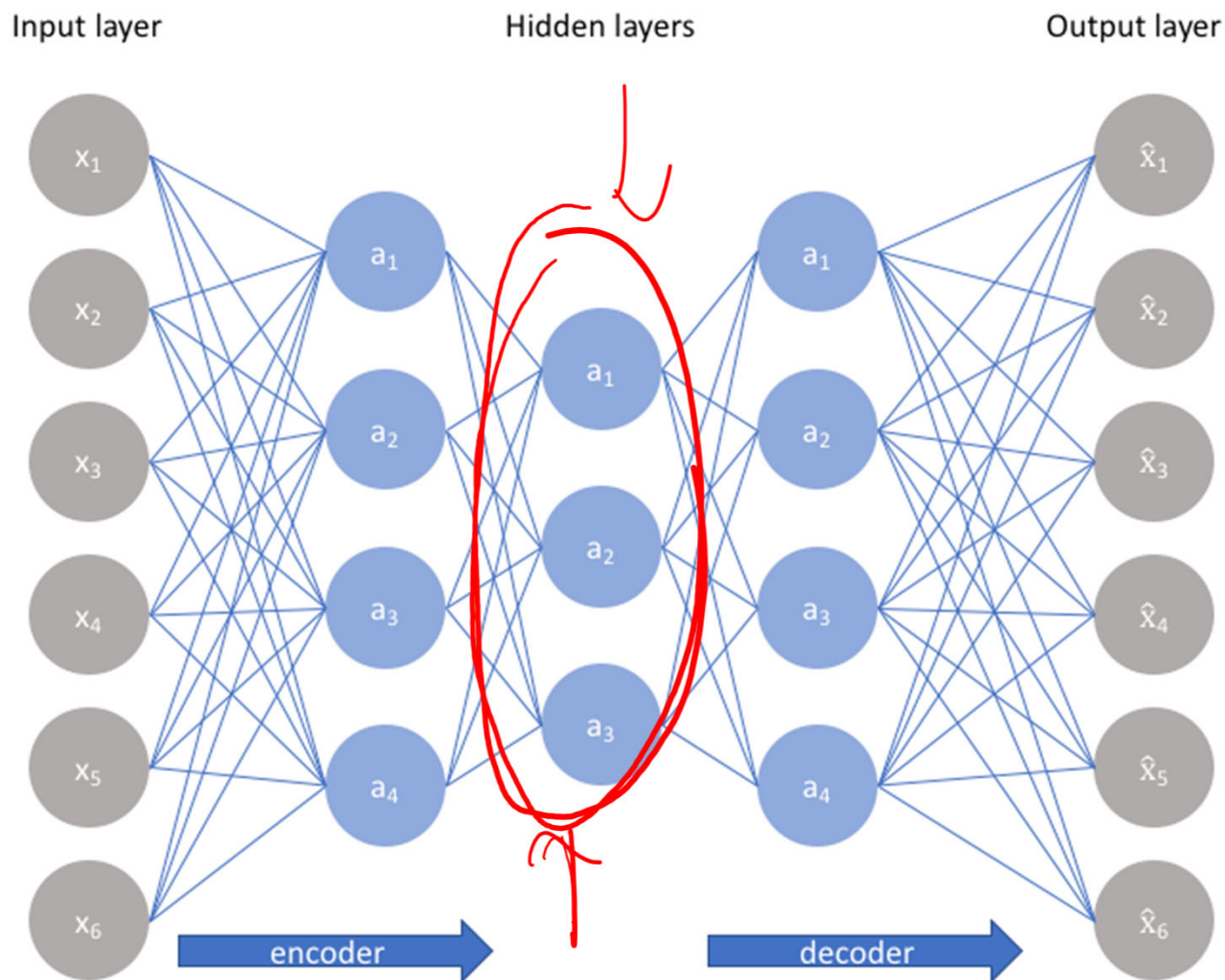*[ideally captures meaningful attributes/modes of variation]*

**Encoder**

**Input data**  $x$

Smile: 0.99

Skin tone: 0.85

Gender: -0.73

Beard: 0.85

Glasses: 0.002

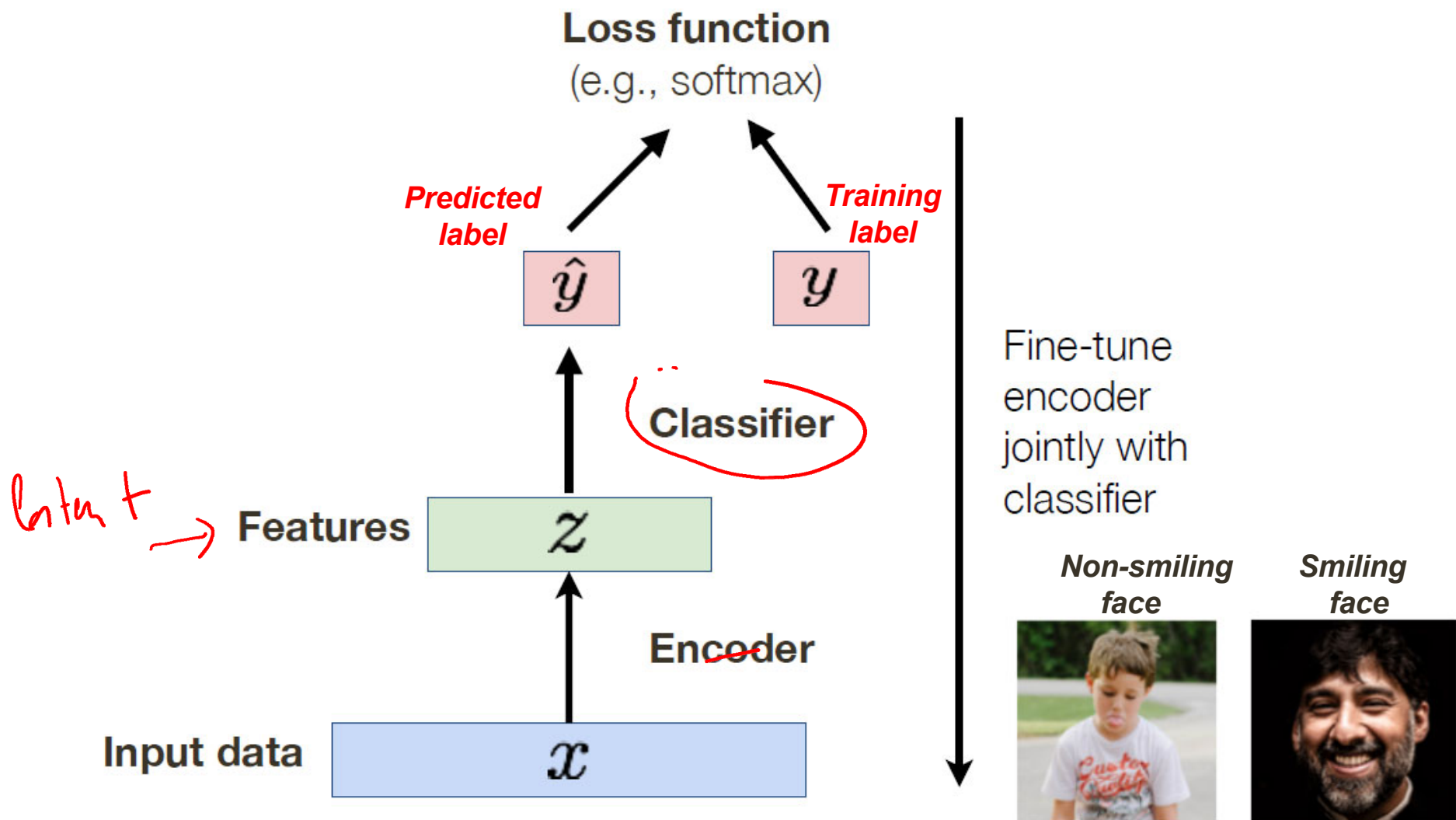Hair color: 0.68

https://www.jeremyjordan.me/autoencoders/

By penalizing the network according to the reconstruction error, our model can learn the most important attributes of the input data and how to best reconstruct the original input from the encoded bottleneck.
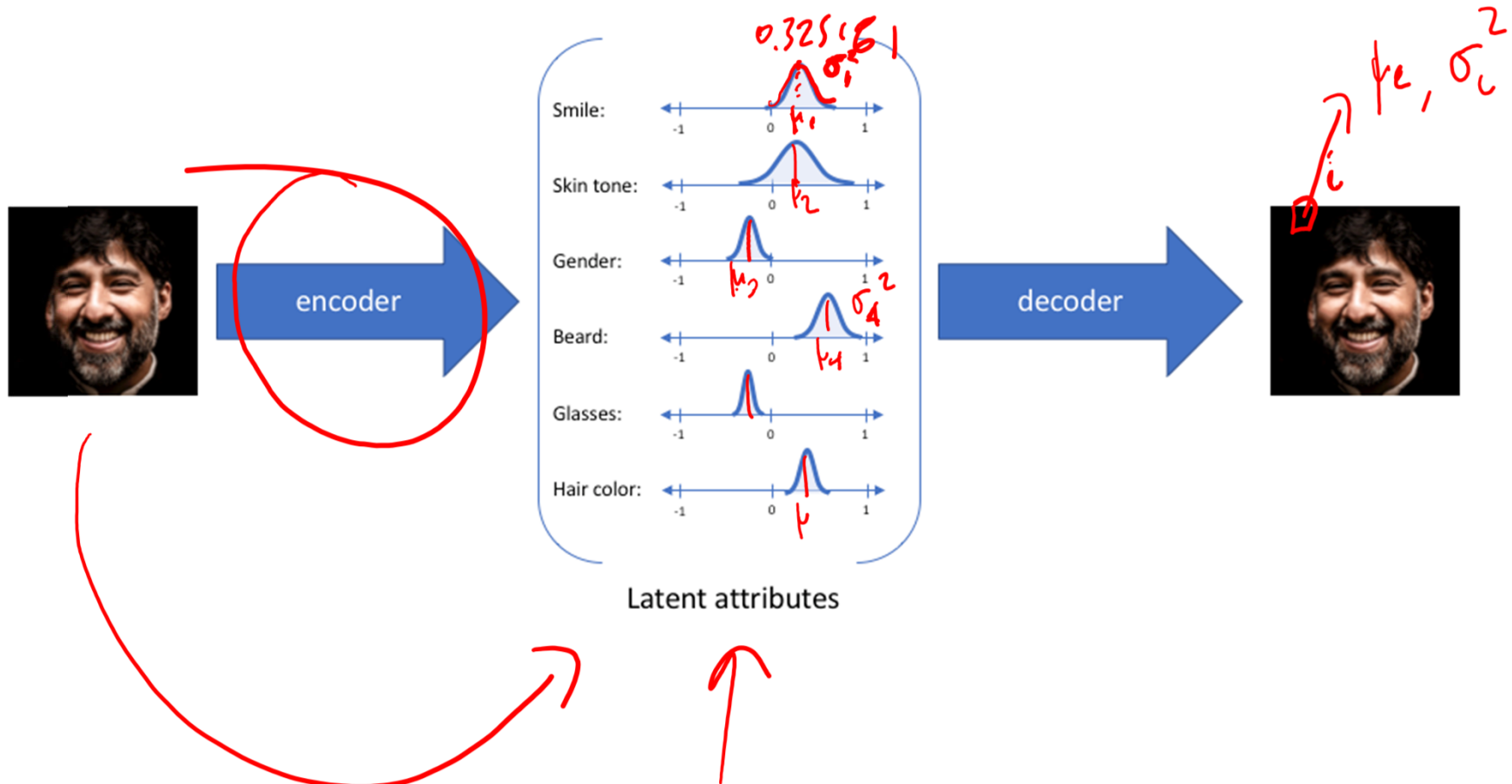
# What autoencoders are useful for?

(Non-variational) Autoencoders cannot be used to sample new data!

For recognition: after pre-training with a reconstruction loss, fine-tune encoder for a recognition task with **few amounts of data!**

# Variational Autoencoders

The encoder instead outputs a range of possible values (a prob distribution) from which we'll randomly sample to feed into our decoder model.
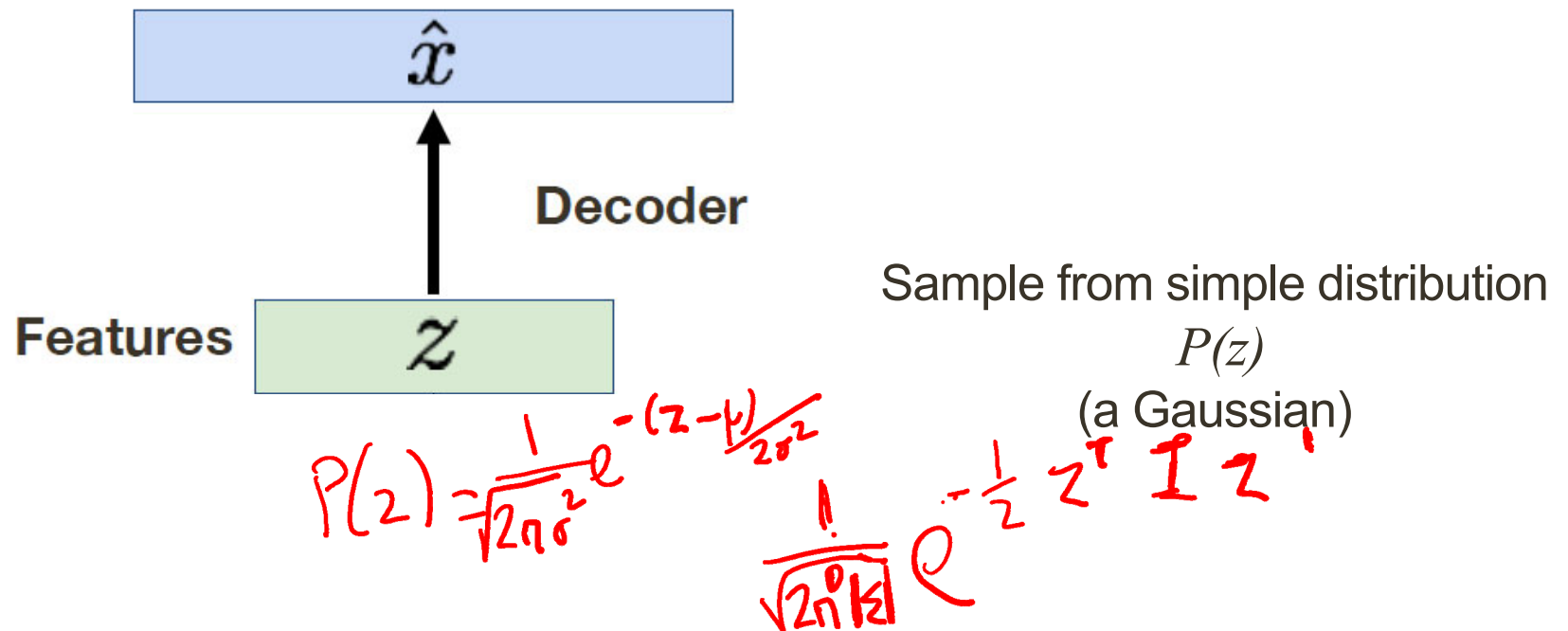**=> Enforce a continuous, smooth latent space representation.**

# Variational Autoencoders
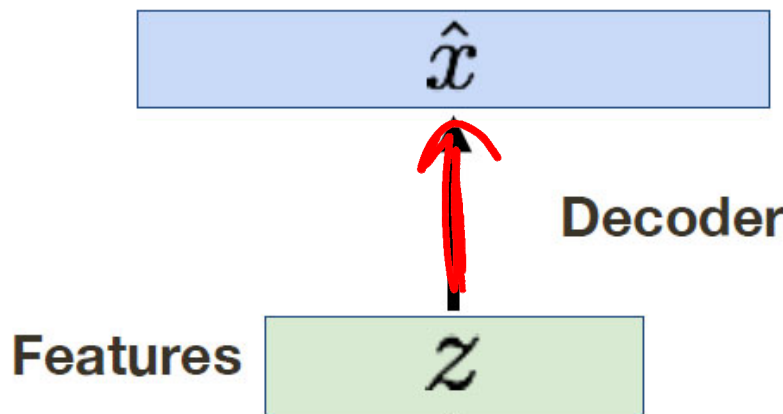
**Allow us to generate data!**

Assume training data is generated from underlying unobserved **latent representation z**

**At test time:**



$\hat{x}$

Decoder

Features $\mathcal{z}$

Sample from simple distribution
$P(z)$
(a Gaussian)

$$P(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(z-\mu)}{2\sigma^2}}$$

$$\frac{1}{\sqrt{2\pi|z|}} e^{-\frac{1}{2} z^T I z}$$

$\mu(z)$ $\sigma^2(z)$ $\max \log \left( P(x_1) \cdot P(x_2) \cdot P(x_3) \cdot \ldots \right)$

# Variational Autoencoders

$z = [0.0, -1.5, 0.9, \ldots]$

$(\mu, \sigma^2)$

**Allow us to generate data!**

Assume training data is generated from underlying unobserved

**latent representation z**

**At test time:**

$P(x), P(z|x)$



$\hat{x}$

Decoder

Features    $z$

Sample from complex cond. distribution

$P(x \mid z)$

(a neural network with learned param $\theta$)
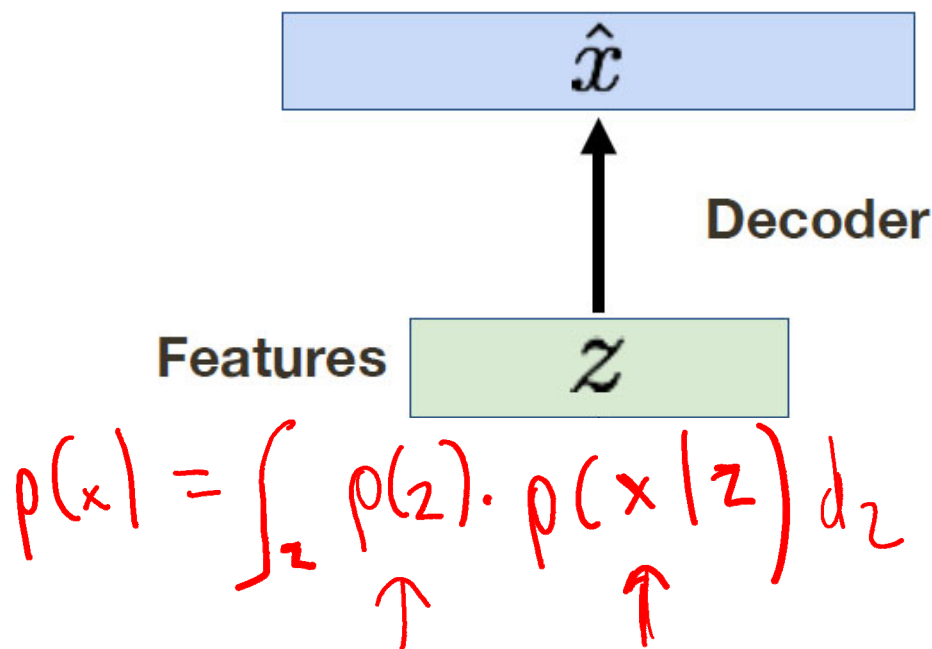
Sample from simple distribution

$P(z)$

(a Gaussian)

$\rightarrow P(x, z) = P(z) \cdot P(x|z)$

Gaussian      Gaussian

# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x_i) = \int_z p_\theta(z)p_\theta(x|z)dz$$

$\hat{x}$

Decoder

Features | $z$

Sample from complex cond. distribution
*P(x | z)*
(a neural network with learned param *θ*)

Sample from simple distribution
*P(z)*
(a Gaussian)

$$p(x) = \int_z p(z) \cdot p(x|z)\, dz$$

# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \int_z p_\theta(z) p_\theta(x|z) dz$$

🙂

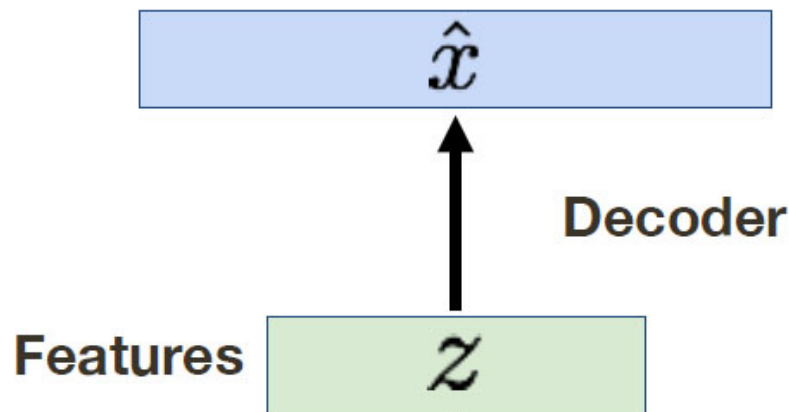Simple **Gaussian** Prior

Sample from complex cond. distribution
*P(x | z)*
(a neural network with learned param *θ*)
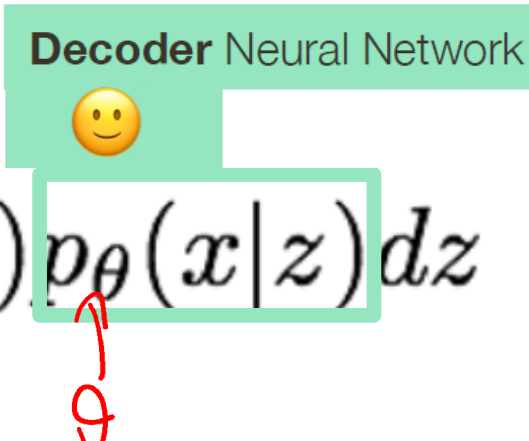
Sample from simple distribution
*P(z)*
(a Gaussian)

$\hat{x}$

Decoder

Features    $z$

# Variational Autoencoders

🙂

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \int_z p_\theta(z) p_\theta(x|z) dz$$

$\theta$



Sample from complex cond. distribution
*P(x | z)*
(a neural network with learned param $\theta$)
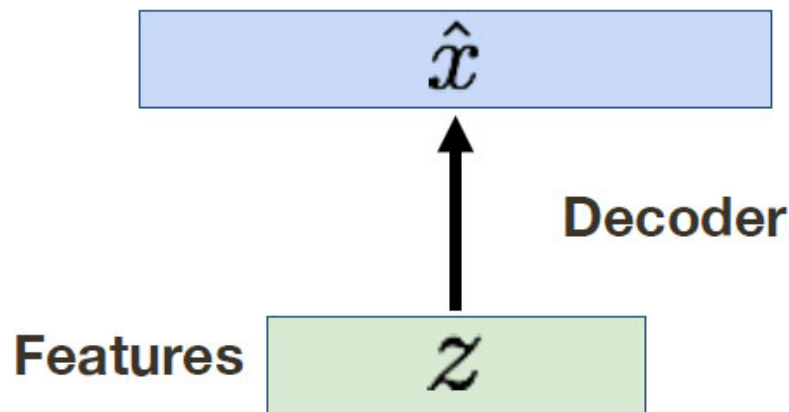
Sample from simple distribution
*P(z)*
(a Gaussian)

# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \int_z p_\theta(z)p_\theta(x|z)dz$$

Intractable to compute for every z 🙁

Sample from complex cond. distribution
$P(x \mid z)$
(a neural network with learned param $\theta$)

Sample from simple distribution
$P(z)$
(a Gaussian)
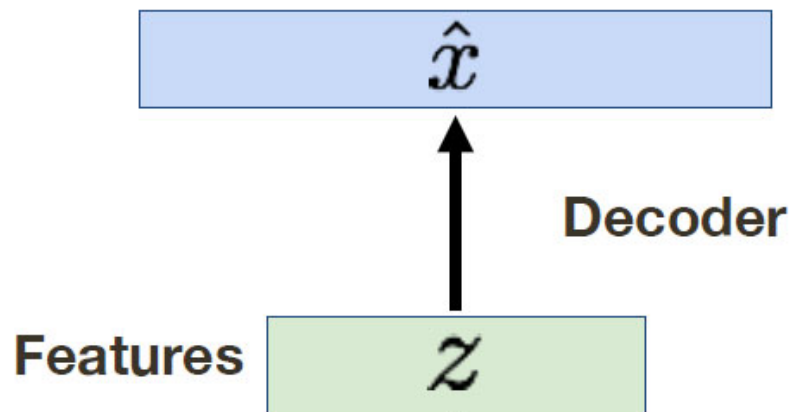
$\hat{x}$

Decoder

Features $z$

# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \boxed{\int_z} p_\theta(z) p_\theta(x|z) dz$$

Intractable to compute for every z

☹️

$\hat{x}$

**Decoder**

Features | $z$

Sample from complex cond. distribution
$P(x \mid z)$
(a neural network with learned param $\theta$)

Sample from simple distribution
$P(z)$
(a Gaussian)

latent | data

**Posterior** density is also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / \boxed{p_\theta(x)}$
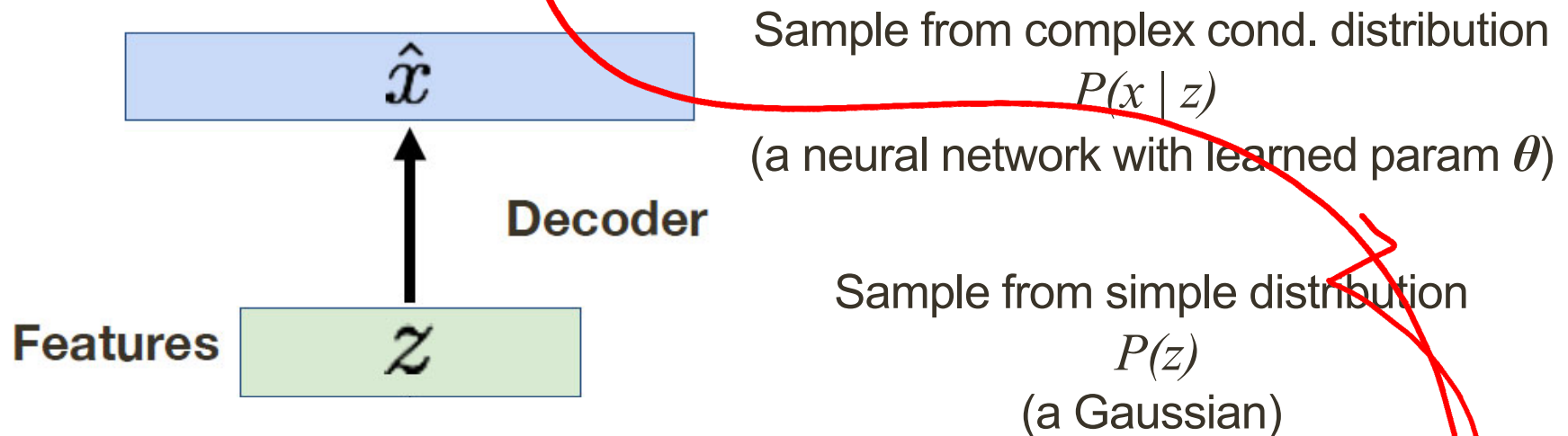
# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \boxed{\int_z} p_\theta(z) p_\theta(x|z) dz$$

Intractable to compute for every z 😦

Sample from complex cond. distribution
$P(x \mid z)$
(a neural network with learned param $\theta$)

$\hat{x}$

**Decoder**

**Features** $\quad z$

Sample from simple distribution
$P(z)$
(a Gaussian)

Solution: approximate $p_\theta(z \mid x)$ with a tractable distribution $q_\phi(z \mid x)$

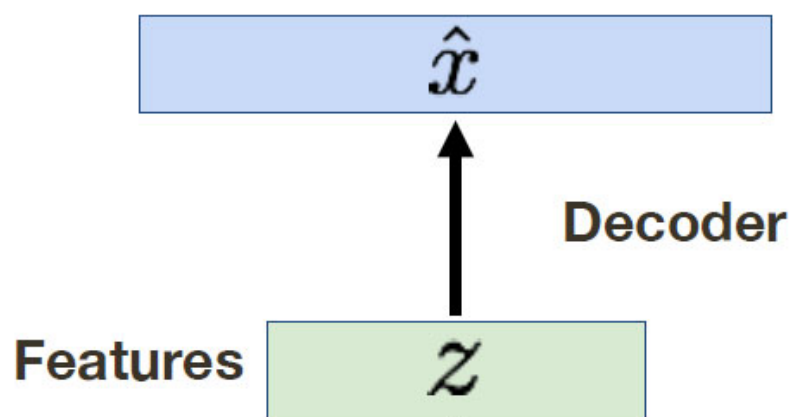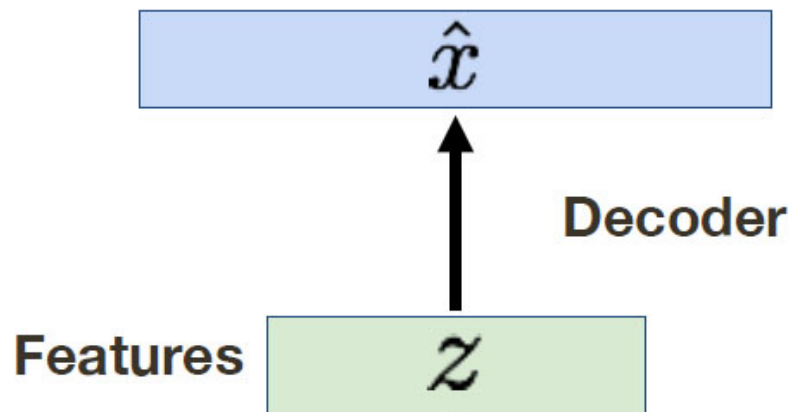**Posterior** density is also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$

# Variational Autoencoders

**How to train this model?**

Maximum likelihood:

$$p_\theta(x) = \int_z p_\theta(z) p_\theta(x|z) dz$$

Intractable to compute for every z ☹



Sample from complex cond. distribution
$P(x \mid z)$
(a neural network with learned param $\theta$)

Sample from simple distribution
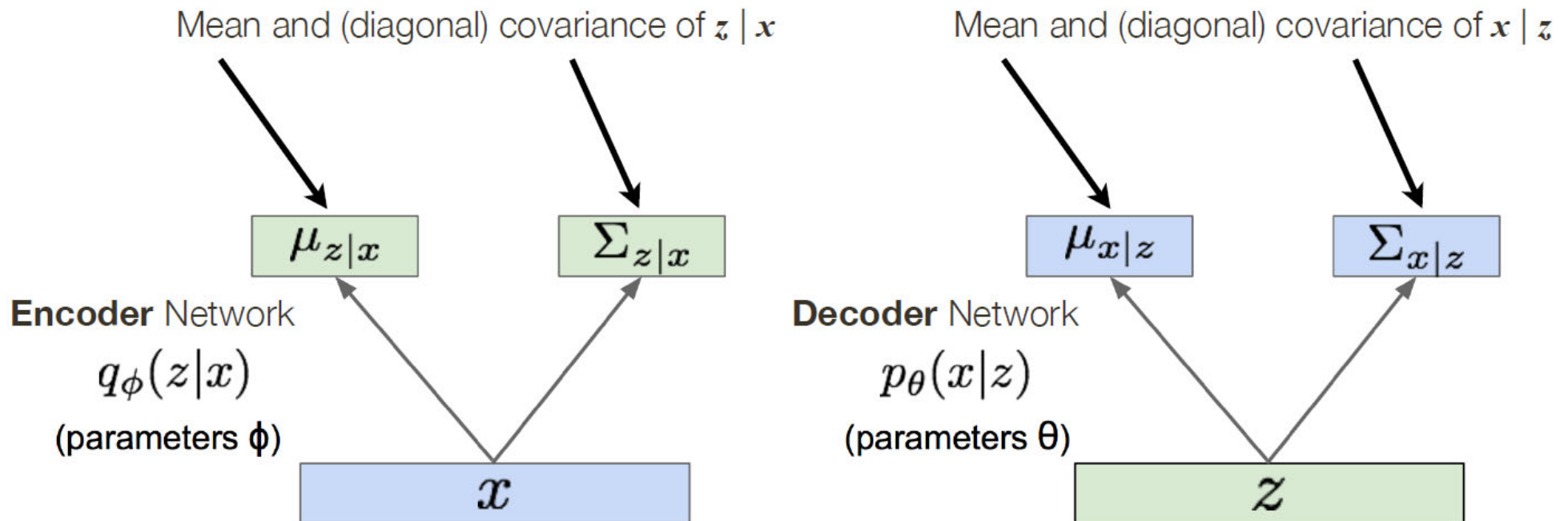$P(z)$
(a Gaussian)

Solution: approximate $p_\theta(z \mid x)$ with a **neural network** $q_\phi(z \mid x)$ [encoder]

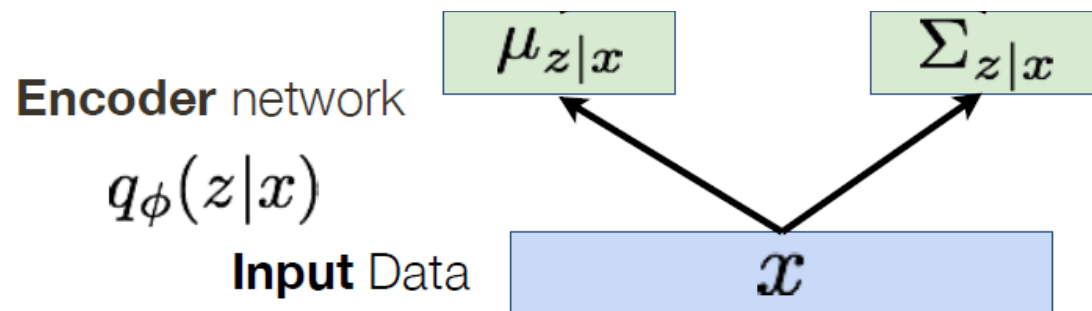**Posterior** density is also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$

# Variational Autoencoders

Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic (they model Gaussian distributions)

# Forward pass during training

**Encoder** network

$q_\phi(z|x)$

$\mu_{z|x}$

$\Sigma_{z|x}$

**Input** Data

$x$

# Forward pass during training



**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$z$

**Encoder** network

$q_\phi(z|x)$

$\mu_{z|x}$

$\Sigma_{z|x}$

**Input** Data

$x$

# Forward pass during training

# VAE Loss

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log\left( p_\theta\left(x|z\right)\right) \ + \ KL\left(q_\phi\left(z|x\right)\middle\| p_\theta(z)\right)$$

| $\mu_{x|z}$ | $\Sigma_{x|z}$ |
|:---:|:---:|

**Decoder** network

$p_\theta(x|z)$

| $z$ |
|:---:|

**Sample** z from $z|x \sim \mathcal{N}\left(\mu_{z|x}, \Sigma_{z|x}\right)$

| $\mu_{z|x}$ | $\Sigma_{z|x}$ |
|:---:|:---:|

**Encoder** network

$q_\phi(z|x)$

**Input** Data | $x$ |

# VAE Loss

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log\left(p_\theta\left(x|z\right)\right) + KL\left(q_\phi\left(z|x\right)\middle\| p_\theta(z)\right)$$

Reconstruction Loss

(outputs should be as close as possible to input)

reduces to $(x - \mu_{x|z})^2$ for fixed output covariance

$\mu_{x|z}$ $\quad$ $\Sigma_{x|z}$

**Decoder** network

$p_\theta(x|z)$

$z$

**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$\mu_{z|x}$ $\quad$ $\Sigma_{z|x}$

**Encoder** network

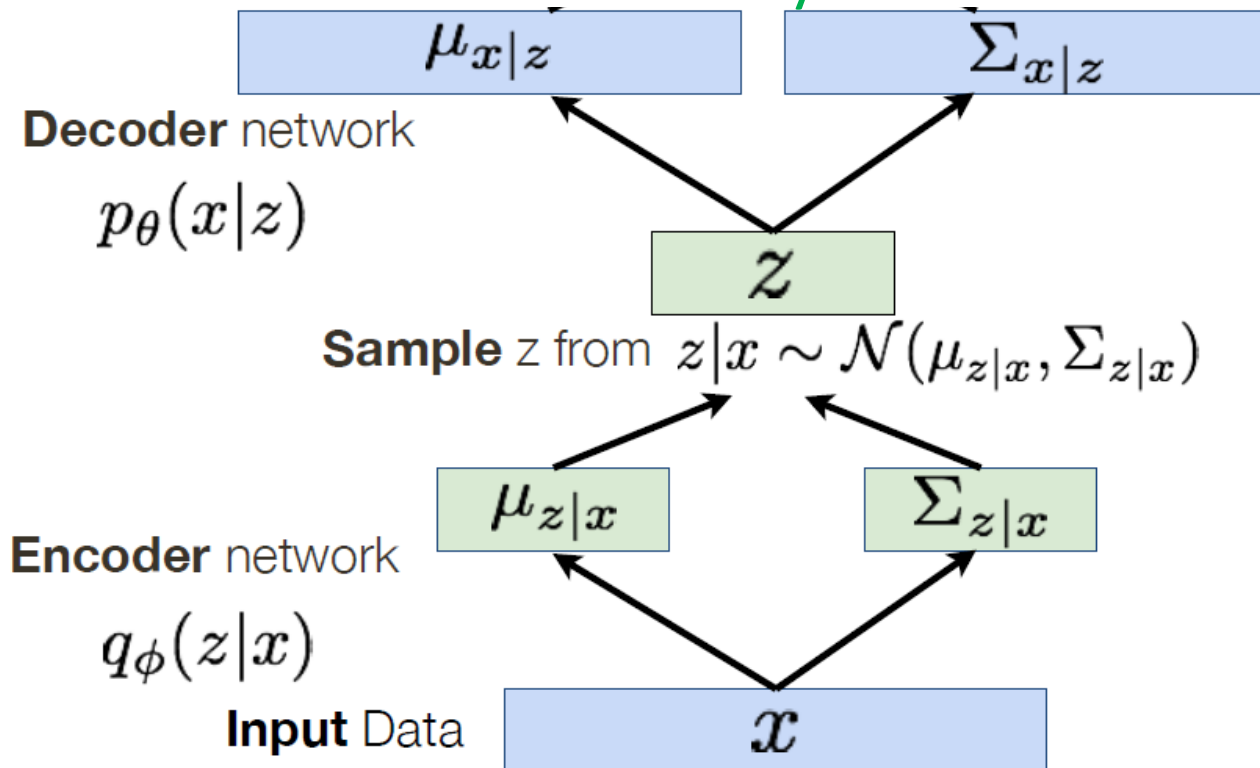$q_\phi(z|x)$

**Input** Data $\quad x$

# VAE Loss

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log\left(p_\theta\left(x|z\right)\right) + \boxed{KL\left(q_\phi\left(z|x\right)\Big\| p_\theta(z)\right)}$$

Regularization term
Make distribution of the latent space produced
by the encoder close to a standard Gaussian.

| $\mu_{x|z}$ | $\Sigma_{x|z}$ |
|---|---|

**Decoder** network

$p_\theta(x|z)$

| $z$ |
|---|

**Sample** z from $\ z|x \sim \mathcal{N}\left(\mu_{z|x}, \Sigma_{z|x}\right)$

| $\mu_{z|x}$ | $\Sigma_{z|x}$ |
|---|---|

**Encoder** network

$q_\phi(z|x)$

**Input** Data    $x$

# KL divergence

A measure of how one probability distribution is different from a second one:

$$KL\left(q_\phi\left(z|x\right)\middle\|p_\theta(z)\right) = \int_z q_\phi\left(z|x\right)\log\frac{q_\phi\left(z|x\right)}{p_\theta(z)}$$

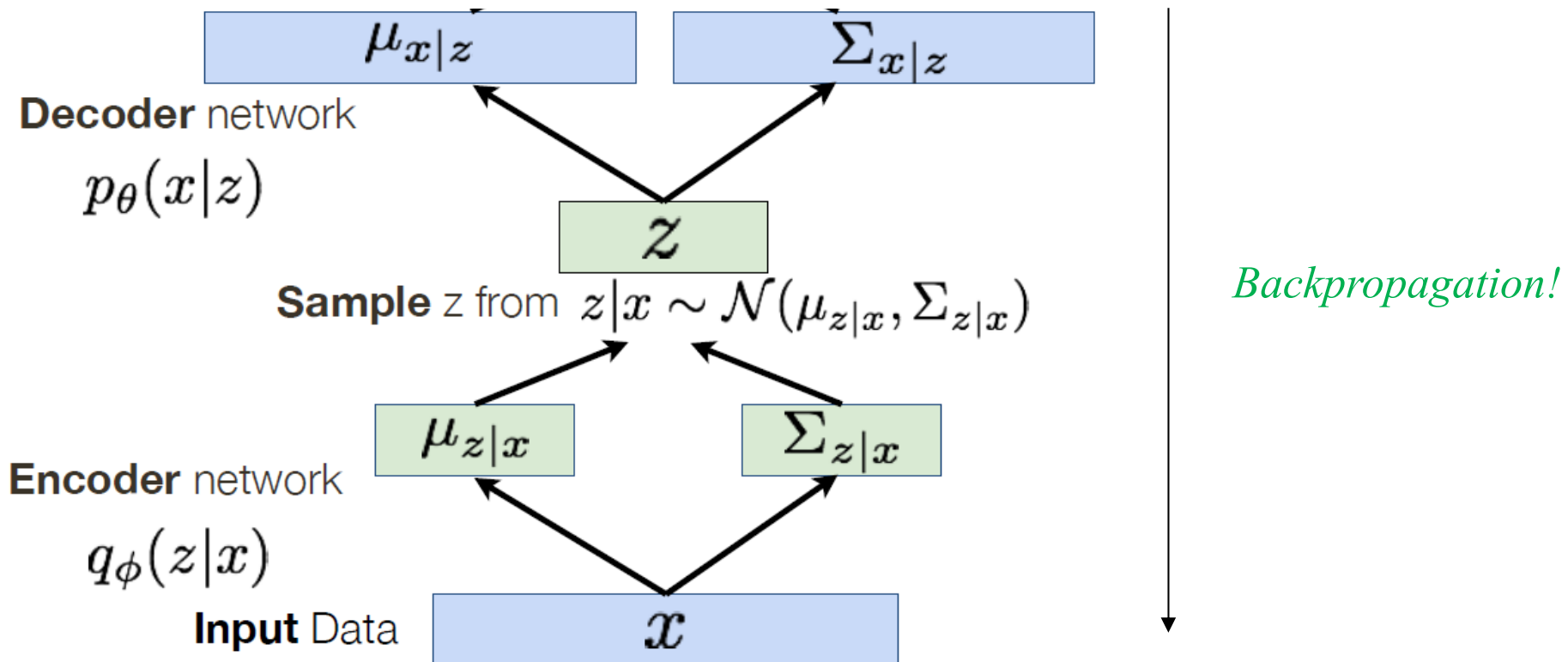In our case, we want our latent space $p_\theta(z)$ to be $N(0, I)$

# VAE Loss

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log\left(p_\theta\left(x|z\right)\right) + KL\left(q_\phi\left(z|x\right)\middle\| p_\theta(z)\right)$$

$$\lambda\left(x - \mu_{x|z}\right)^2 + \sum_{d=1}^{D}\left(\sigma_{z|x}^2[d] + \mu_{z|x}^2[d] - \log\sigma_{z|x}[d] - 1\right)$$
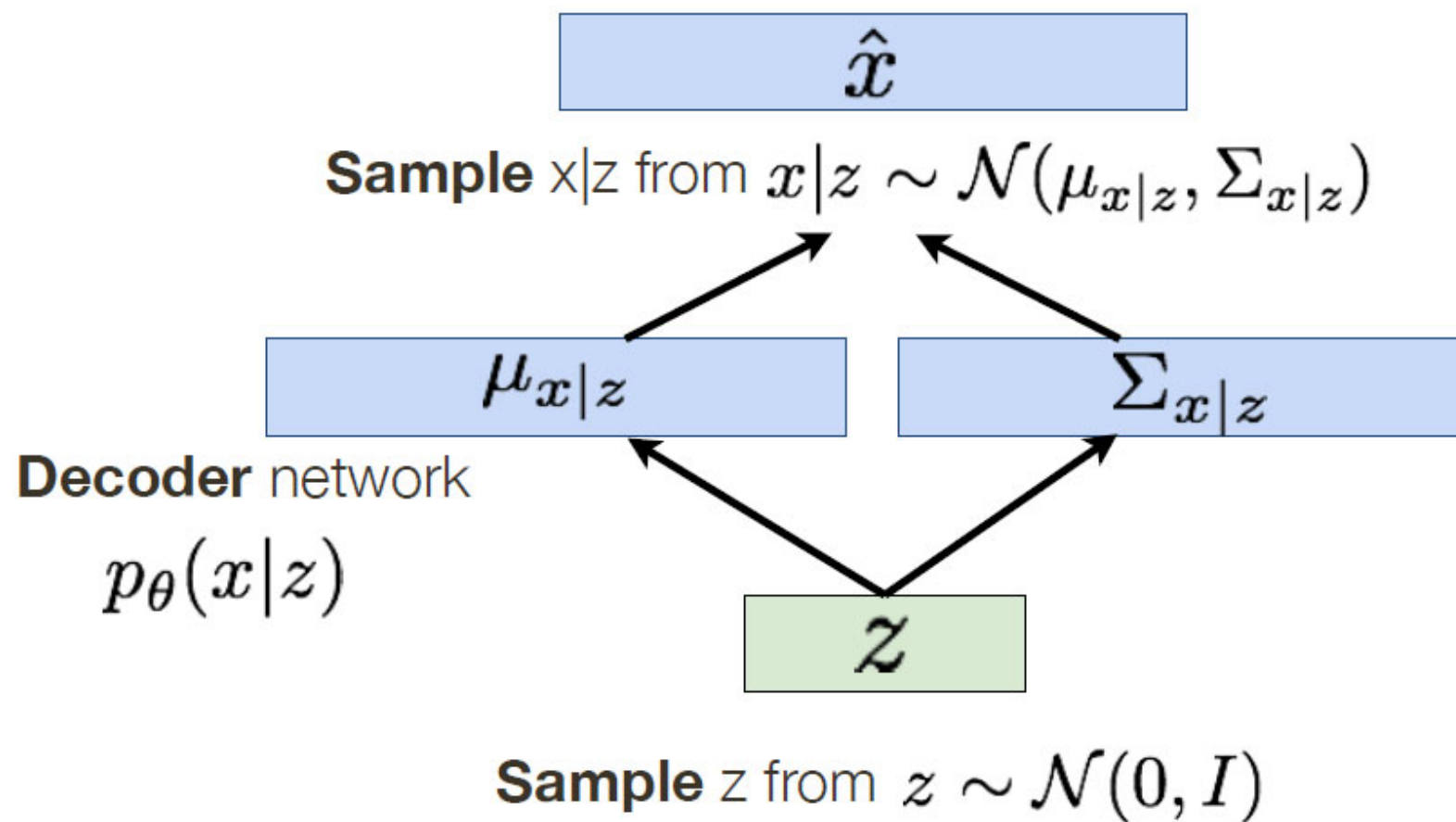
*(where λ is a weighting term)*



**Decoder** network
$p_\theta(x|z)$

**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

**Encoder** network
$q_\phi(z|x)$

**Input** Data

*Backpropagation!*

# VAE Loss (skipping proofs...)

$$-\mathbf{E}_{z \sim q_\phi(z|x)} \log\big(p_\theta(x|z)\big) \;+\; KL\big(q_\phi(z|x)\big\| p_\theta(z)\big)$$

Minimize upper bound
on loss we care about!

$$\geq -\log p_\theta(x)$$



**Decoder** network
$p_\theta(x|z)$

$\mu_{x|z}$    $\Sigma_{x|z}$

$z$

**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$\mu_{z|x}$    $\Sigma_{z|x}$

**Encoder** network
$q_\phi(z|x)$

**Input** Data    $x$

*Backpropagation!*

# Test time



$\hat{x}$

**Sample** x|z from $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{x|z}$        $\Sigma_{x|z}$

**Decoder** network
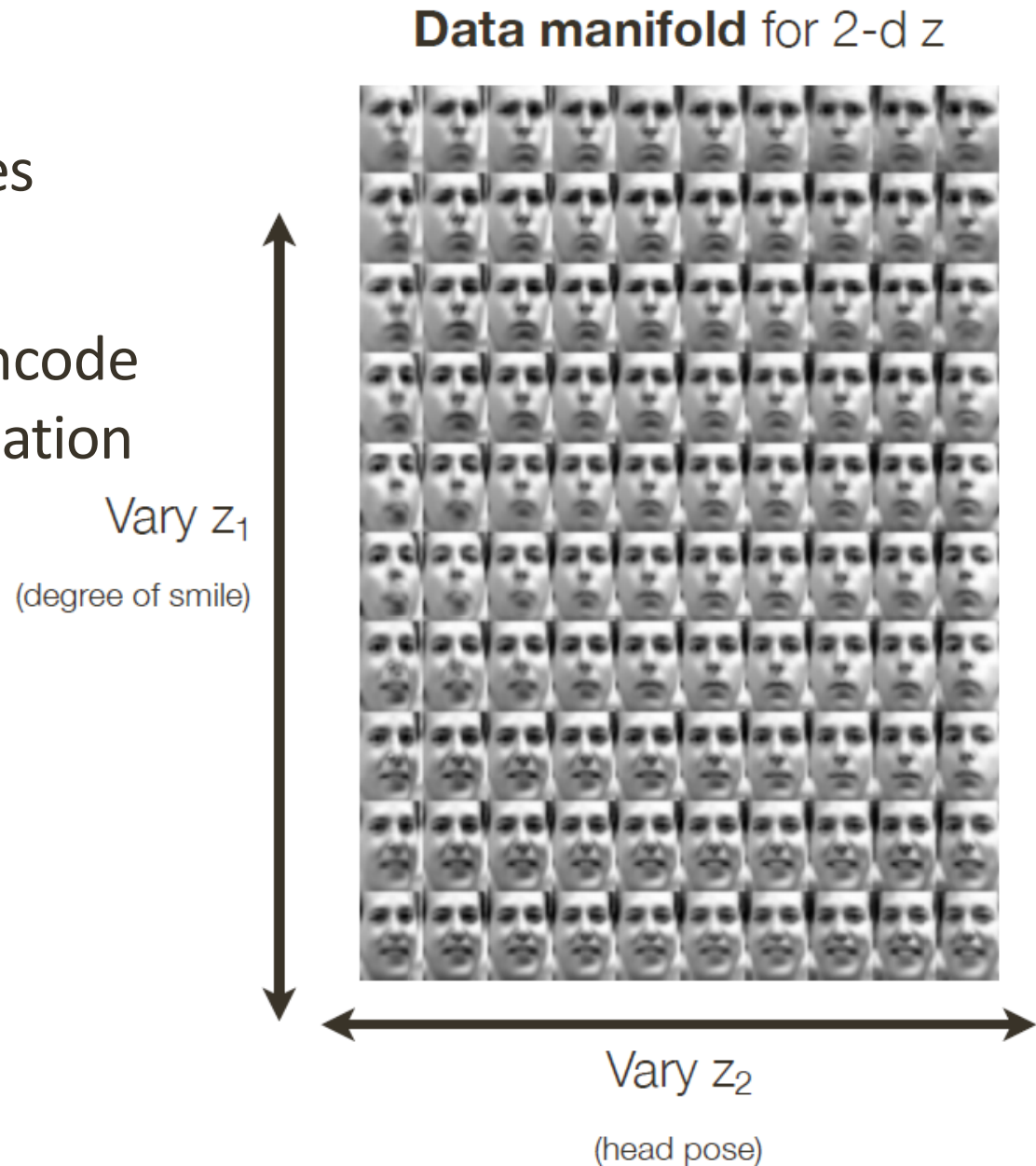
$p_\theta(x|z)$

$z$

**Sample** z from $z \sim \mathcal{N}(0, I)$

# VAE Latent space

- Diagonal prior on z => independent latent variables

- Different dimensions of z encode interpretable factors of variation



**Data manifold** for 2-d z

Vary $z_1$

(degree of smile)

Vary $z_2$

(head pose)

# VAE useful literature

- **Understanding Variational Autoencoders:**
  **https://www.jeremyjordan.me/variational-autoencoders/**

- **Tutorial on Variational Autoencoders**
  **https://arxiv.org/pdf/1606.05908.pdf**

- **Today VAEs are mostly used to produce a low-dimensional latent space of data – latent diffusion models operate on this space…**