# Topics in AI (CPSC 532L):
# Multimodal Learning with Vision, Language and Sound

**Lecture 11: Generative Models**

# Course **Logistics**

— **Assignment 3** was due yesterday

— **Assignment 2 & 3** will be posted today/tomorrow

— **Assignment 4** will be out tomorrow (Friday) and is due in a week


— Reminder: **Project presentations** next Thursday (a week from today)

   — Logistics: form will be up today

   — Send me slides to minimize laptop switching on the day

# **Supervised** vs. **Unsupervised** Learning

| **Supervised** Learning |
|---|
| **Data:** (x, y)<br>    x is data, y is label<br><br>**Goal:** Learn a *function* to map x ➞ y<br><br>**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, *etc.* |



➞ Cat

Classification

# **Supervised** vs. **Unsupervised** Learning
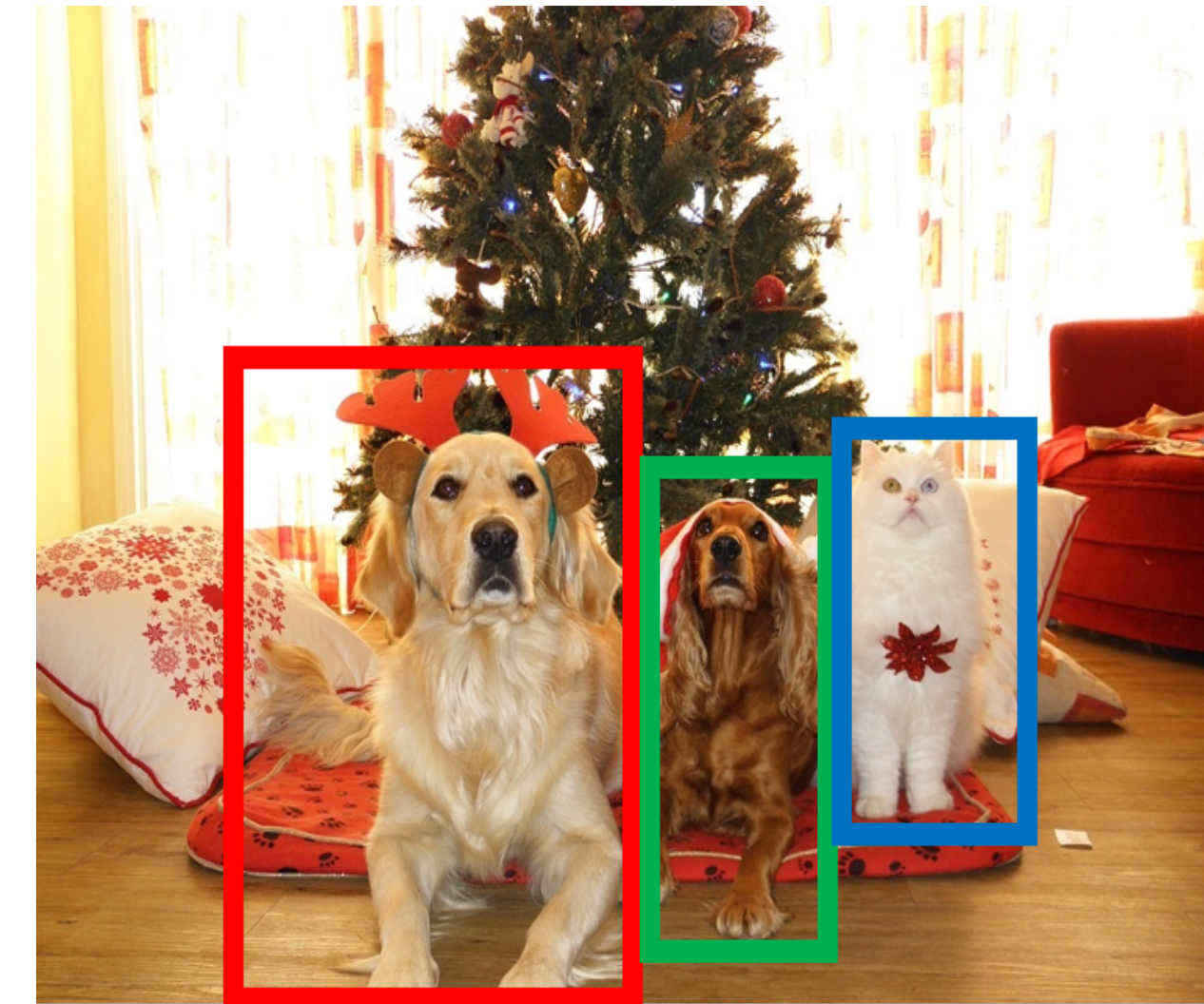


## **Supervised** Learning

**Data:** (x, y)
    x is data, y is label

**Goal:** Learn a *function* to map x ➞ y

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, *etc.*

**DOG**, **DOG**, **CAT**

Object Detection

This image is CC0 public domain

# **Supervised** vs. **Unsupervised** Learning
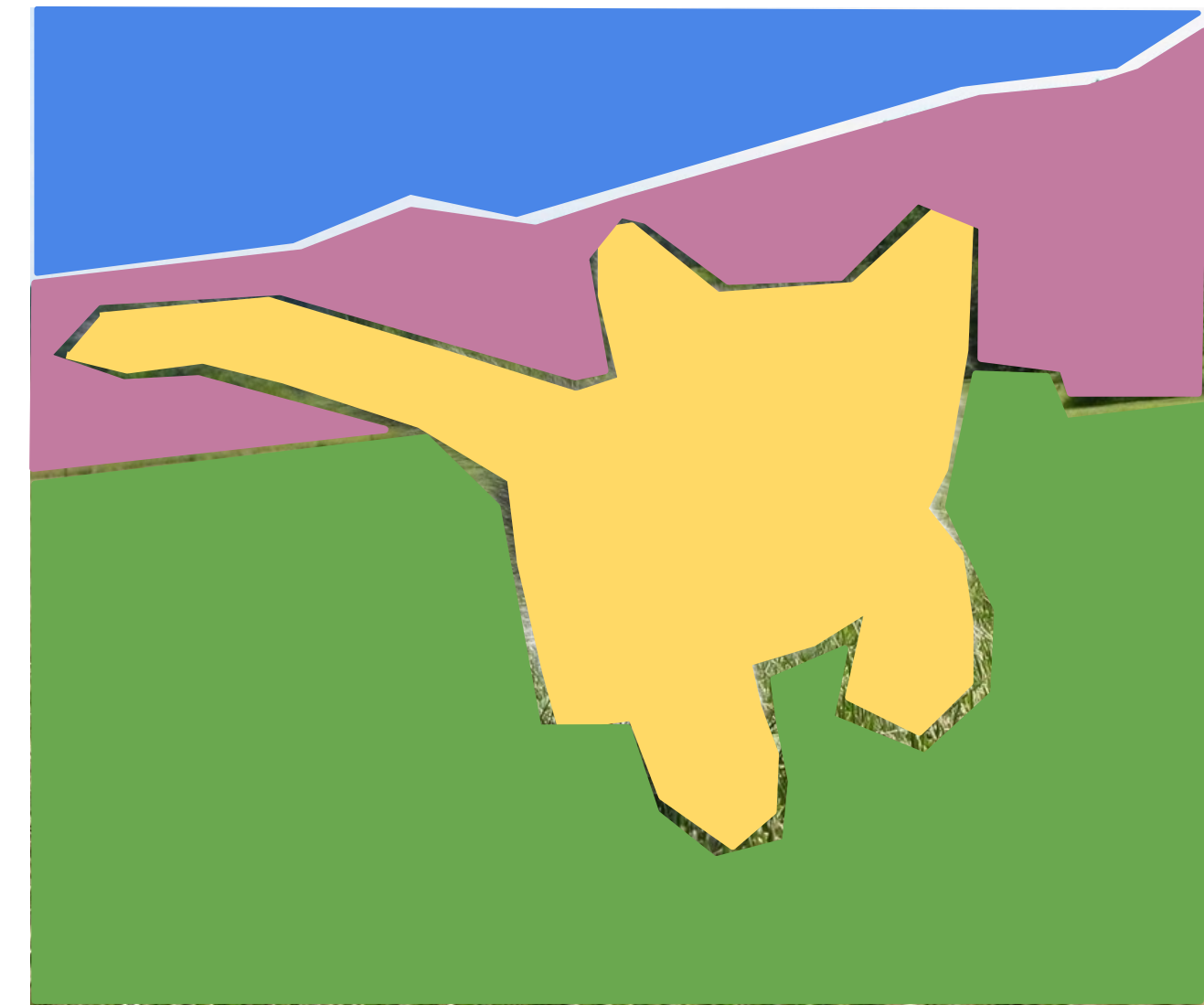
## **Supervised** Learning

**Data:** (x, y)
      x is data, y is label

**Goal:** Learn a *function* to map x➙y

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, *etc.*



**GRASS**, **CAT**, **TREE**, **SKY**

Semantic Segmentation

# **Supervised** vs. **Unsupervised** Learning



## **Supervised** Learning

**Data:** (x, y)
    x is data, y is label

**Goal:** Learn a *function* to map x➞y

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, *etc.*

A cat sitting on a suitcase on the floor

Image Captioning

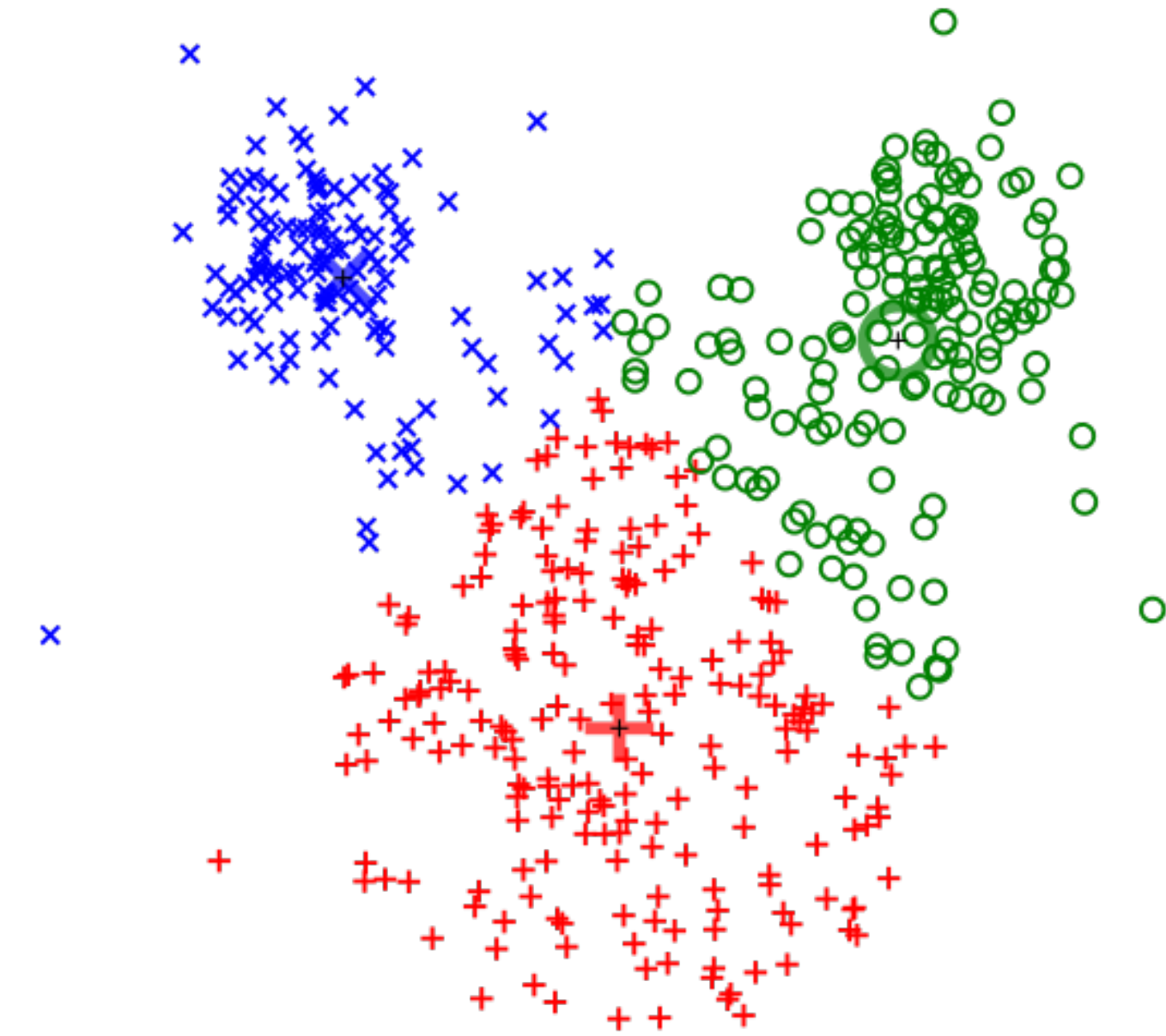# **Supervised** vs. **Unsupervised** Learning



**Unsupervised** Learning

**Data:** x
   Just data, no labels!

**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, *etc.*

k-means clustering

This image is CC0 public domain
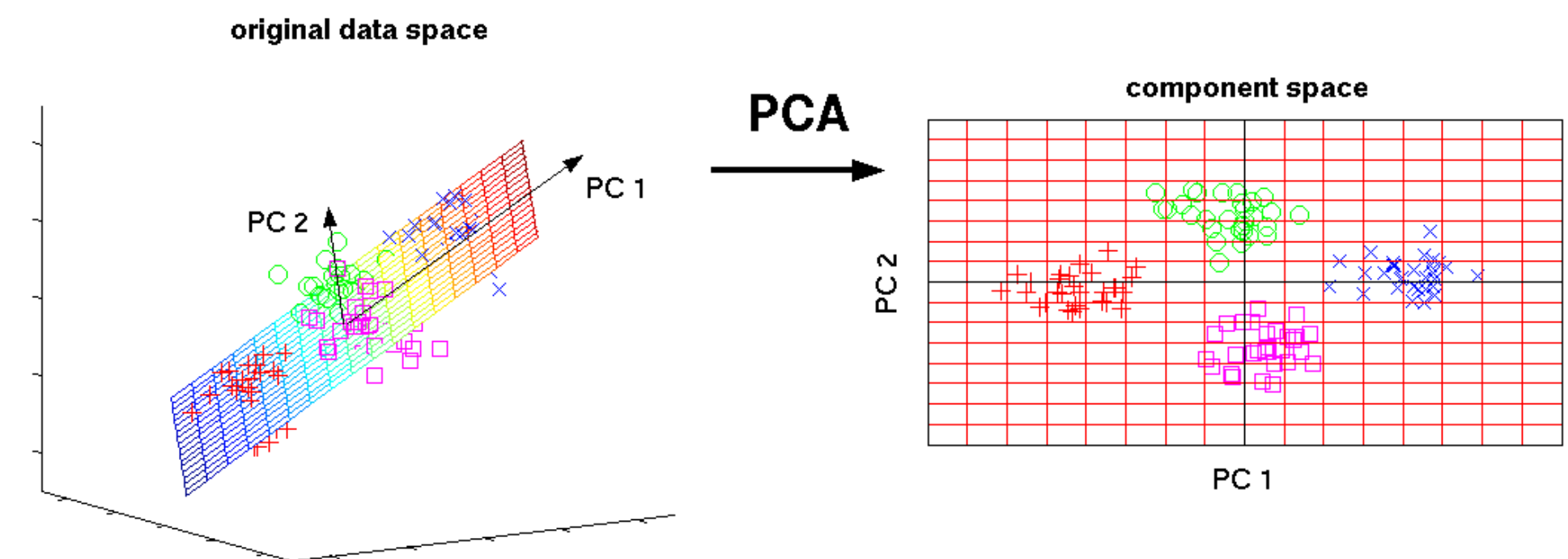
# **Supervised** vs. **Unsupervised** Learning



**Unsupervised** Learning

**Data:** x
   Just data, no labels!

**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, *etc.*

dimensionality reduction

# **Supervised** vs. **Unsupervised** Learning

## **Unsupervised** Learning
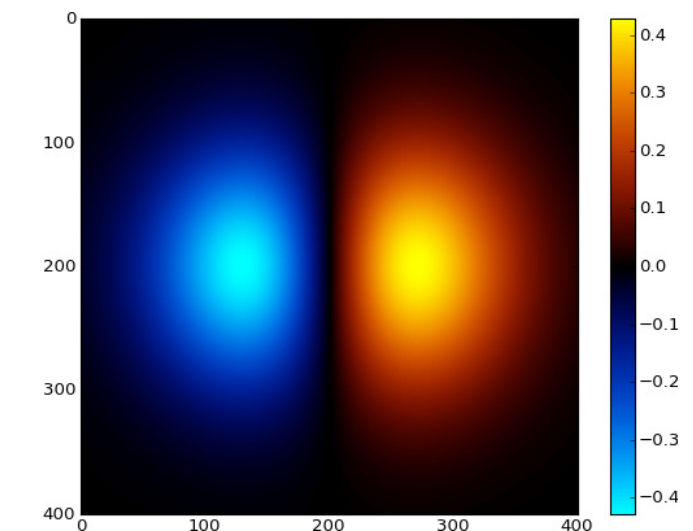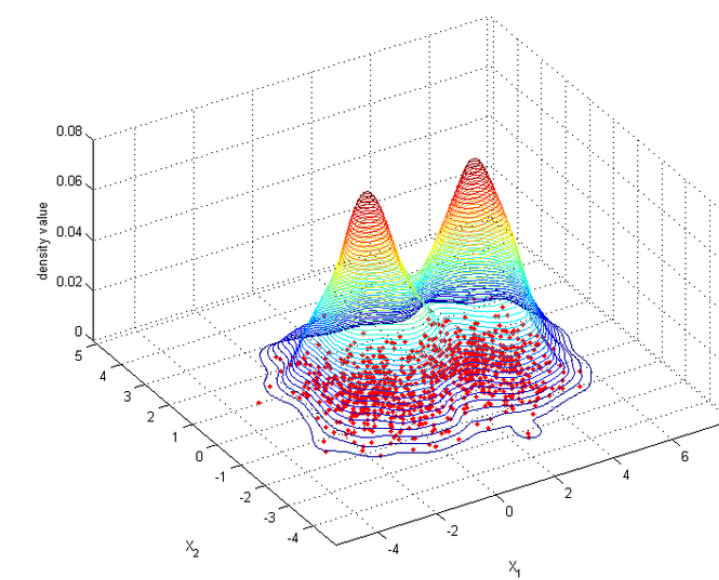
**Data:** x
    Just data, no labels!

**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, *etc.*



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-dim density estimation



2-dim density estimation

2-d density images left and right are CC0 public domain

# **Supervised** vs. **Unsupervised** Learning

| **Supervised** Learning | **Unsupervised** Learning |
|---|---|
| **Data:** (x, y)<br><br>    x is data, y is label<br><br>**Goal:** Learn a *function* to map x➞y<br><br>**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, *etc.* | **Data:** x<br><br>    Just data, no labels!<br><br>**Goal:** Learn some underlying hidden *structure* of the data<br><br>**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, *etc.* |

# **Generative** Models

Given training data, generate new samples from the same distribution



Training data ~ $p_{\text{data}}(\boldsymbol{x})$

Generated samples ~ $p_{\text{model}}(\boldsymbol{x})$

Want to learn $p_{\text{model}}(\boldsymbol{x})$ similar to $p_{\text{data}}(\boldsymbol{x})$

# **Generative** Models

Given training data, generate new samples from the same distribution



Training data ~ $p_{\text{data}}(\boldsymbol{x})$



Generated samples ~ $p_{\text{model}}(\boldsymbol{x})$

Want to learn $p_{\text{model}}(\boldsymbol{x})$ similar to $p_{\text{data}}(\boldsymbol{x})$

Addresses **density estimation**, a core problem in unsupervised learning

— **Explicit** density estimation: explicitly define and solve for $p_{\text{model}}(\boldsymbol{x})$

— **Implicit** density estimation: learn model that can sample from $p_{\text{model}}(\boldsymbol{x})$ w/o explicitly defining it
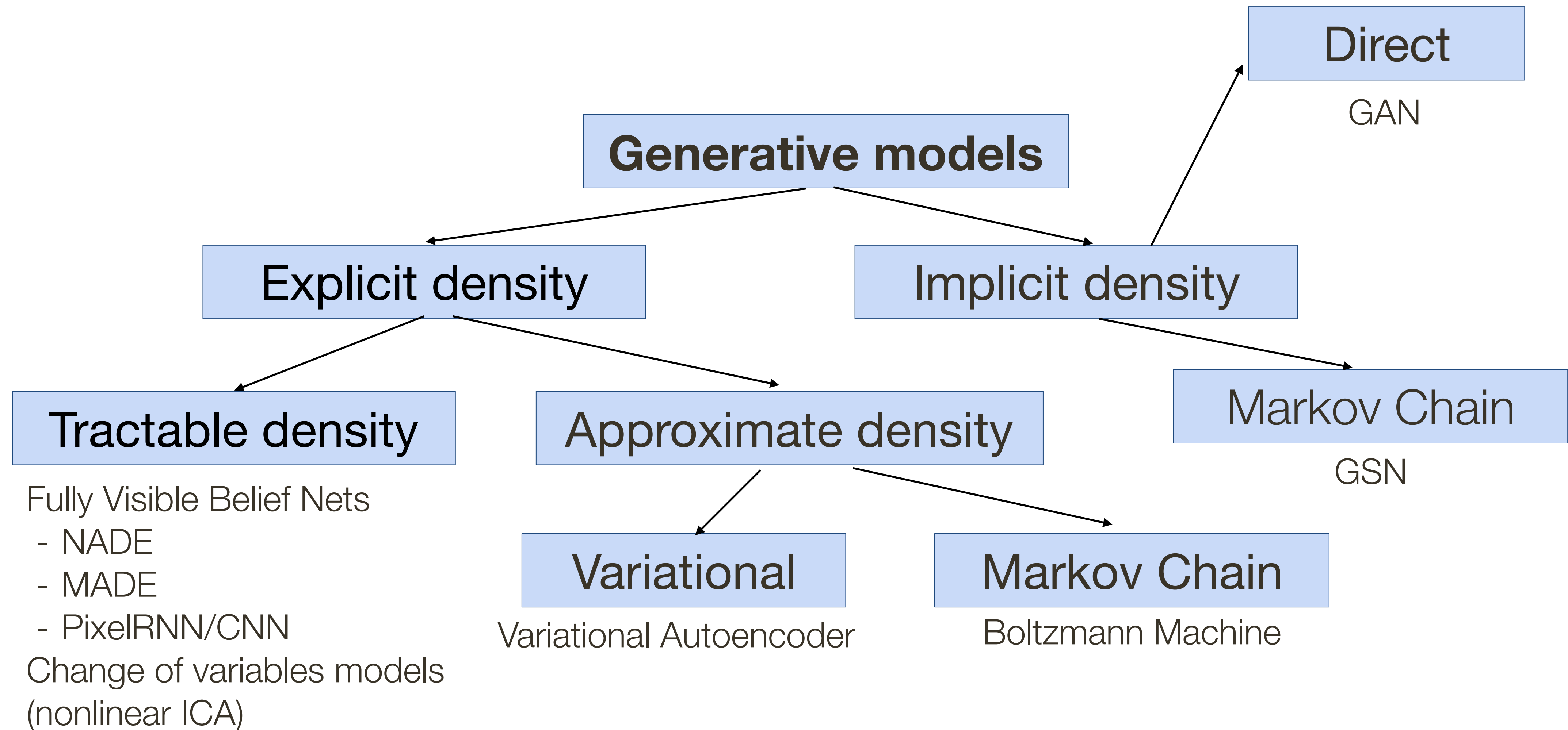
# **Taxonomy** of Generative Models



Direct

GAN

**Generative models**

Explicit density

Implicit density

Tractable density

Approximate density

Markov Chain

GSN

Fully Visible Belief Nets
- NADE
- MADE
- PixelRNN/CNN
Change of variables models
(nonlinear ICA)

Variational

Markov Chain
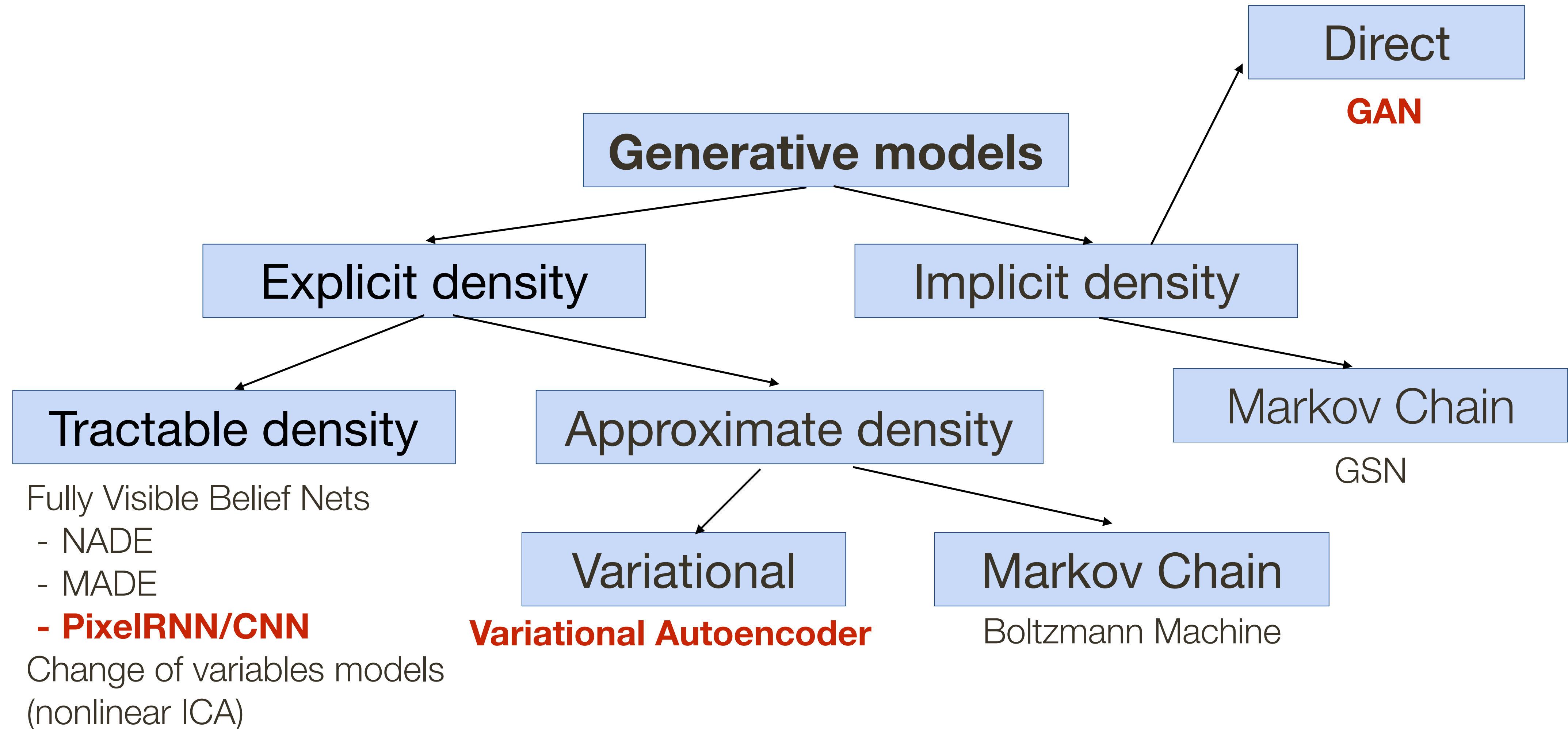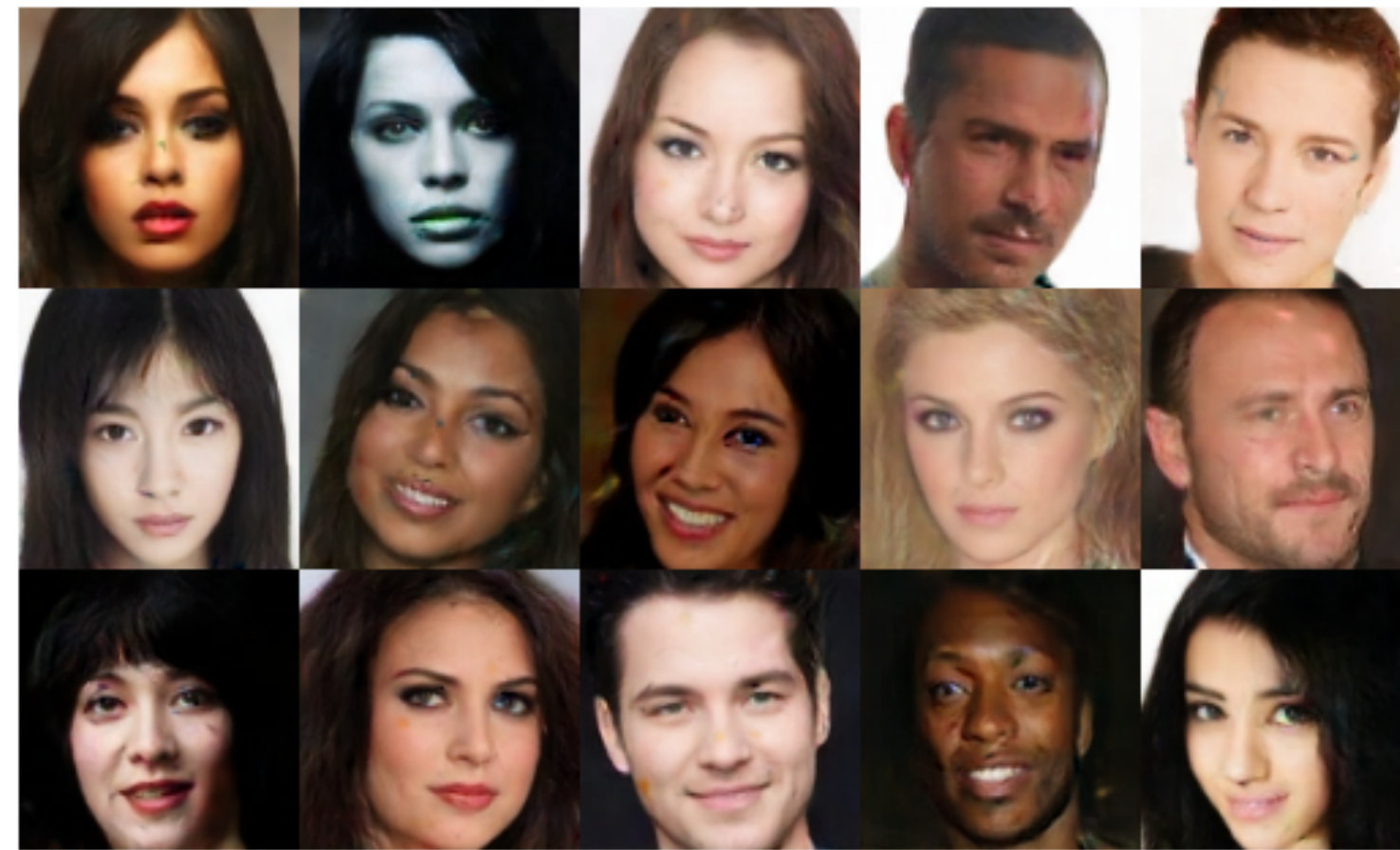
Variational Autoencoder

Boltzmann Machine

Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.
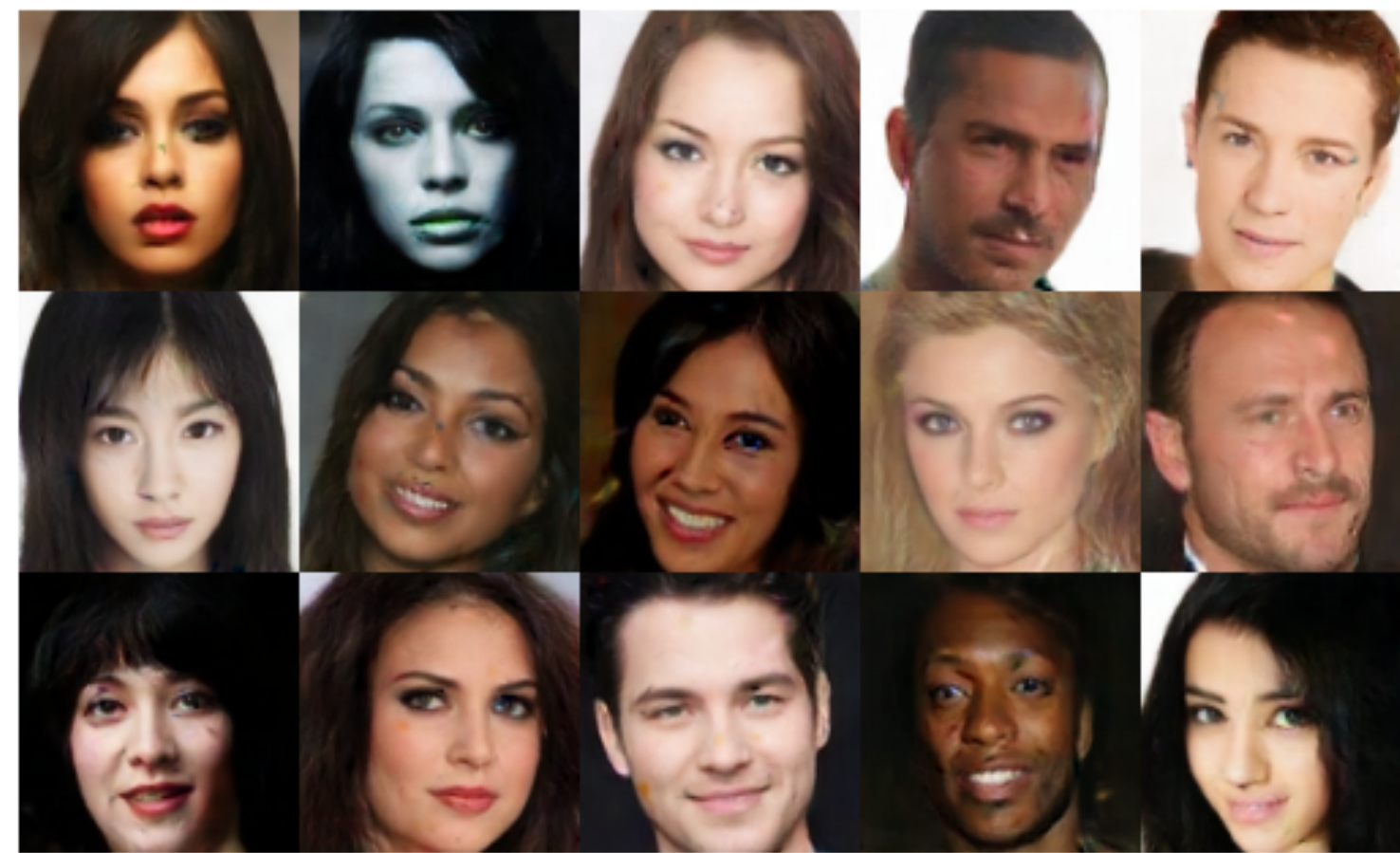
* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# **Taxonomy** of Generative Models



Direct

**GAN**

**Generative models**

Explicit density

Implicit density

Tractable density

Approximate density

Markov Chain

GSN

Fully Visible Belief Nets
  - NADE
  - MADE
  - **PixelRNN/CNN**
Change of variables models
(nonlinear ICA)

Variational

**Variational Autoencoder**

Markov Chain

Boltzmann Machine

Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.
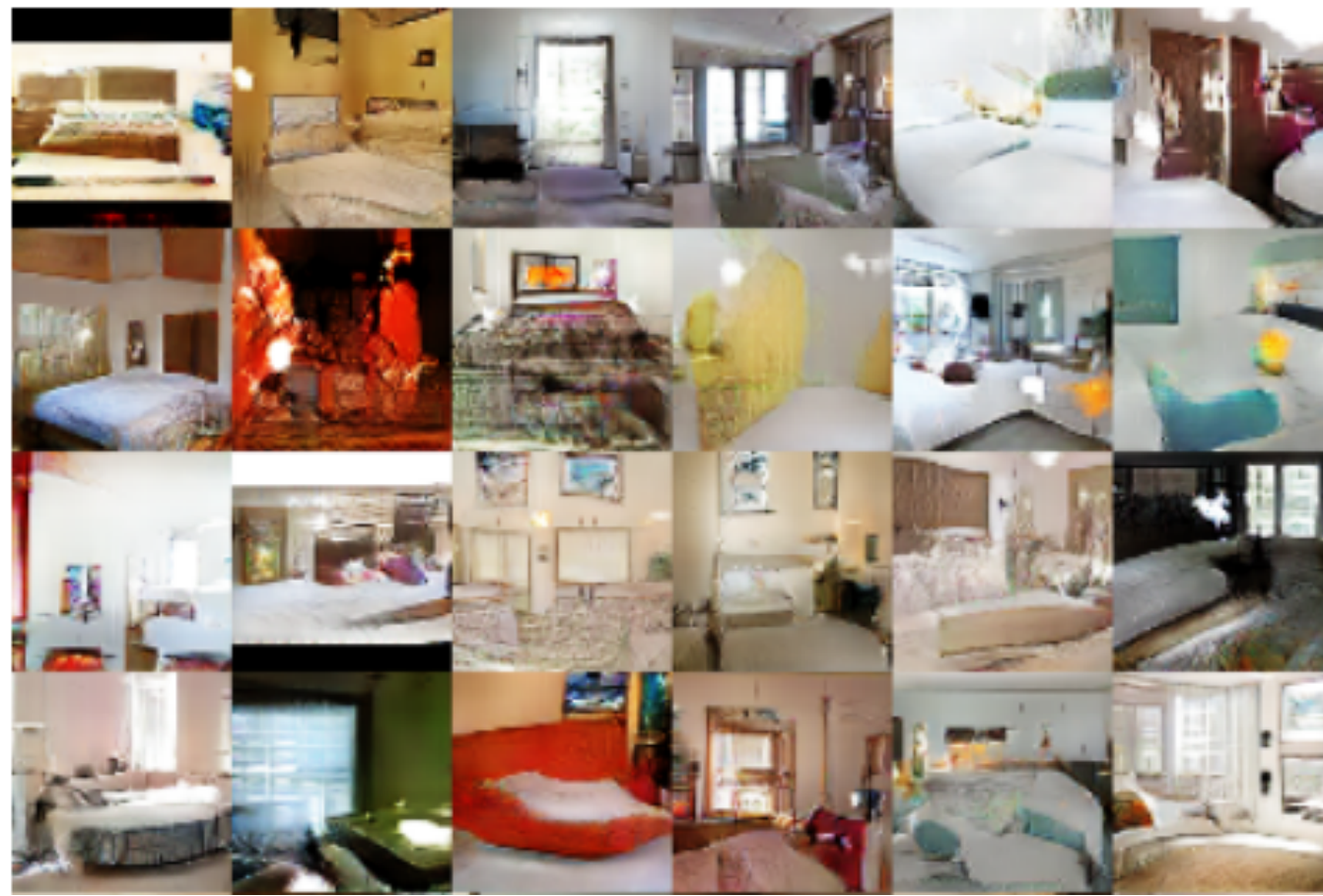
# Why **Generative** Models?

— Realistic **samples** for artwork, super-resolution, colorization, *etc.*

# Why **Generative** Models?

— Realistic **samples** for artwork, super-resolution, colorization, *etc.*



— Generative models of time-series data can be used for **simulation**, **predictions** and planning (reinforcement learning applications)

— Training generative models can also enable inference of latent representation that can be useful as **general features**

— **Dreaming** / hypothesis visualization

# PixelRNN and PixelCNN

# **Pixel**RNN

**Explicit** Density model

Use chain rule to decompose likelihood of an image $x$ into product of (many) 1-d distributions

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

Likelihood of image $x$

Probability of i'th pixel value given all previous pixels

then maximize likelihood of training data

# **Pixel**RNN

## **Explicit** Density model

Use chain rule to decompose likelihood of an image $x$ into product of (many) 1-d distributions

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

Likelihood of image $x$

Probability of i'th pixel value given all previous pixels

Complex distribution over pixel values, so lets model using **neural network**

then maximize likelihood of training data

# **Pixel**RNN

**Explicit** Density model

Use chain rule to decompose likelihood of an image $x$ into product of (many) 1-d distributions

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

Likelihood of image $x$

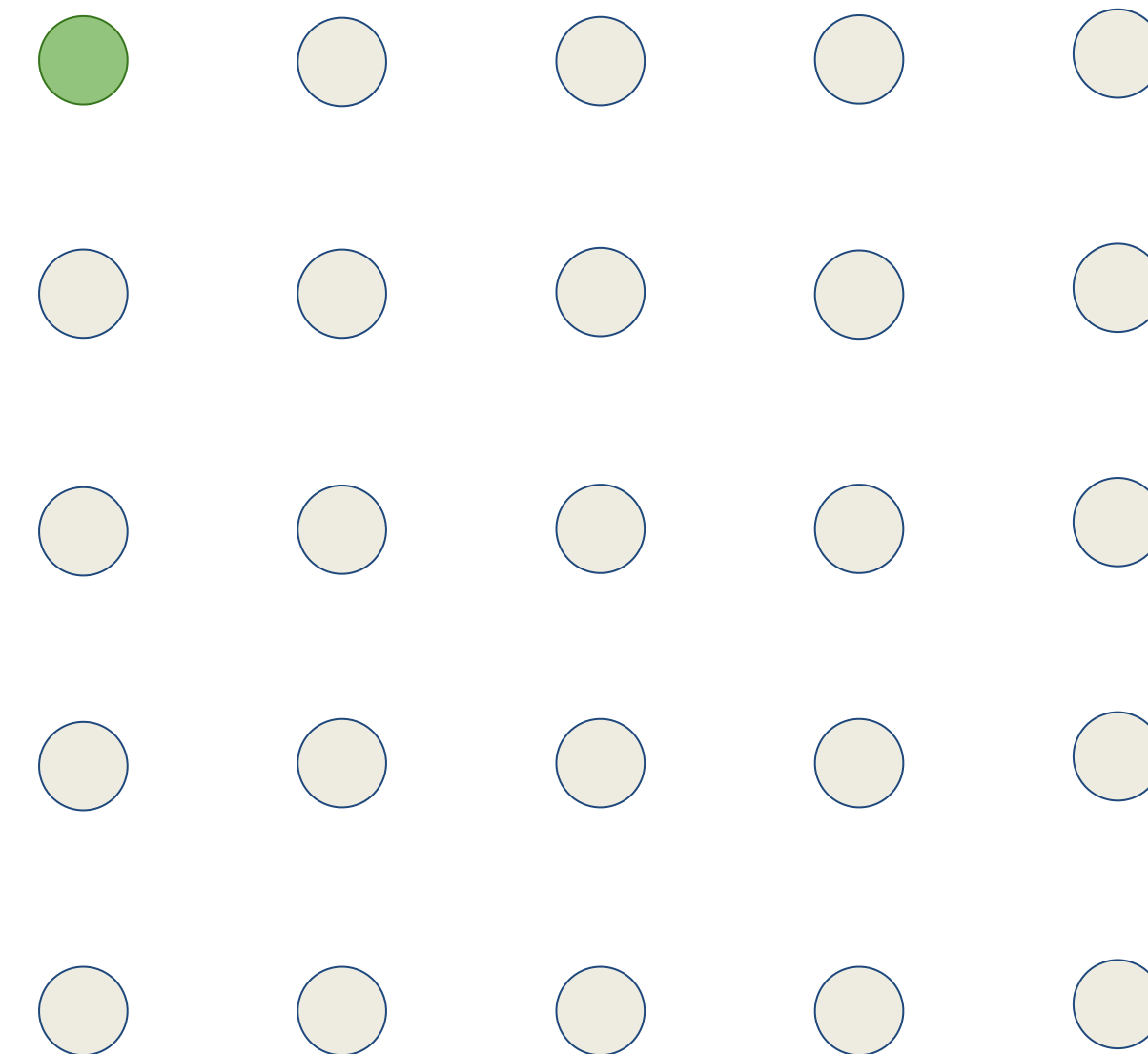Probability of i'th pixel value given all previous pixels

then maximize likelihood of training data

Complex distribution over pixel values, so lets model using **neural network**
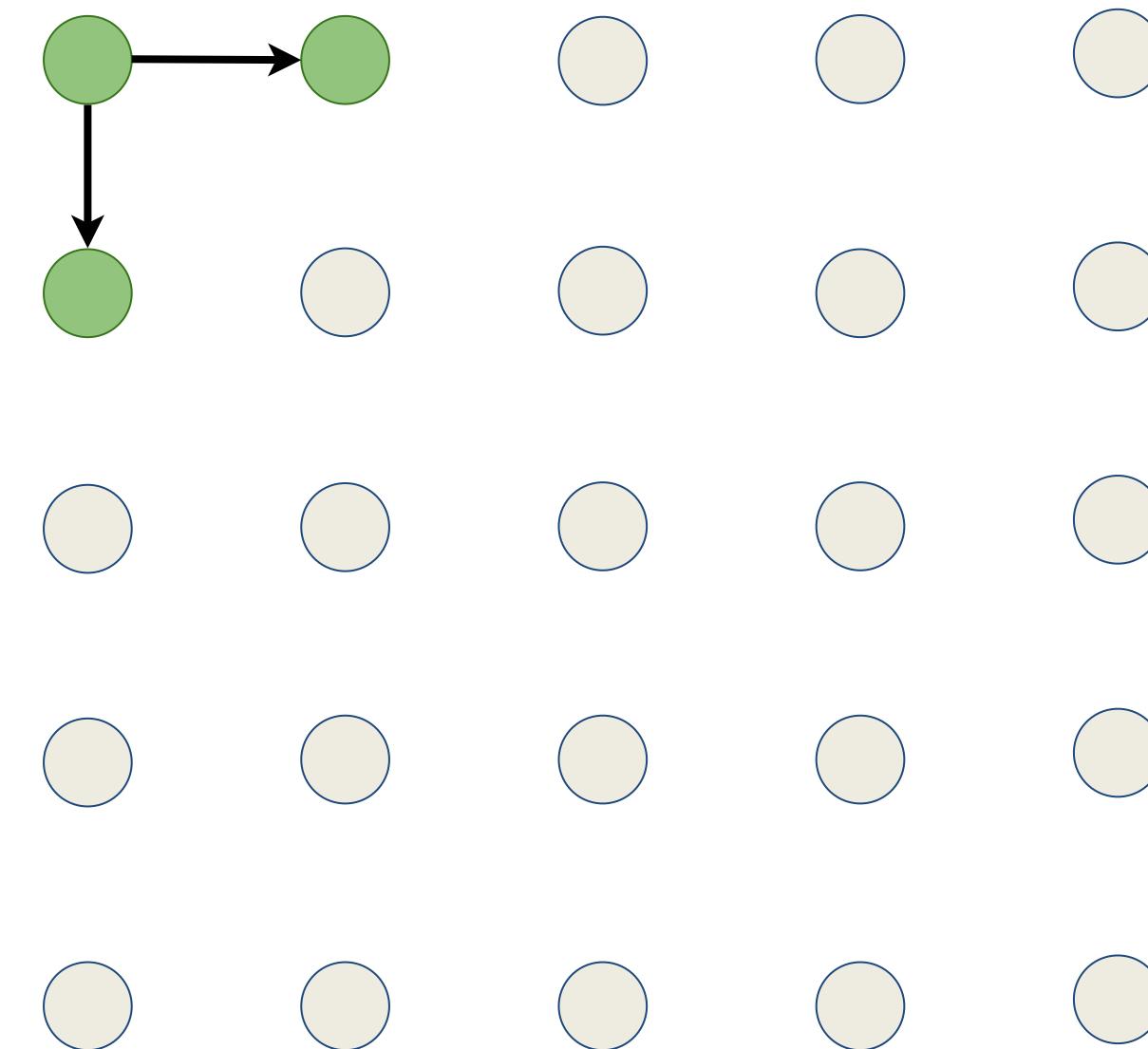
Also requires defining **ordering** of "previous pixels"

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# **Pixel**RNN

Generate image pixels starting
from the corner

Dependency on previous pixels
model using an RNN (LSTM)

# **Pixel**RNN

Generate image pixels starting
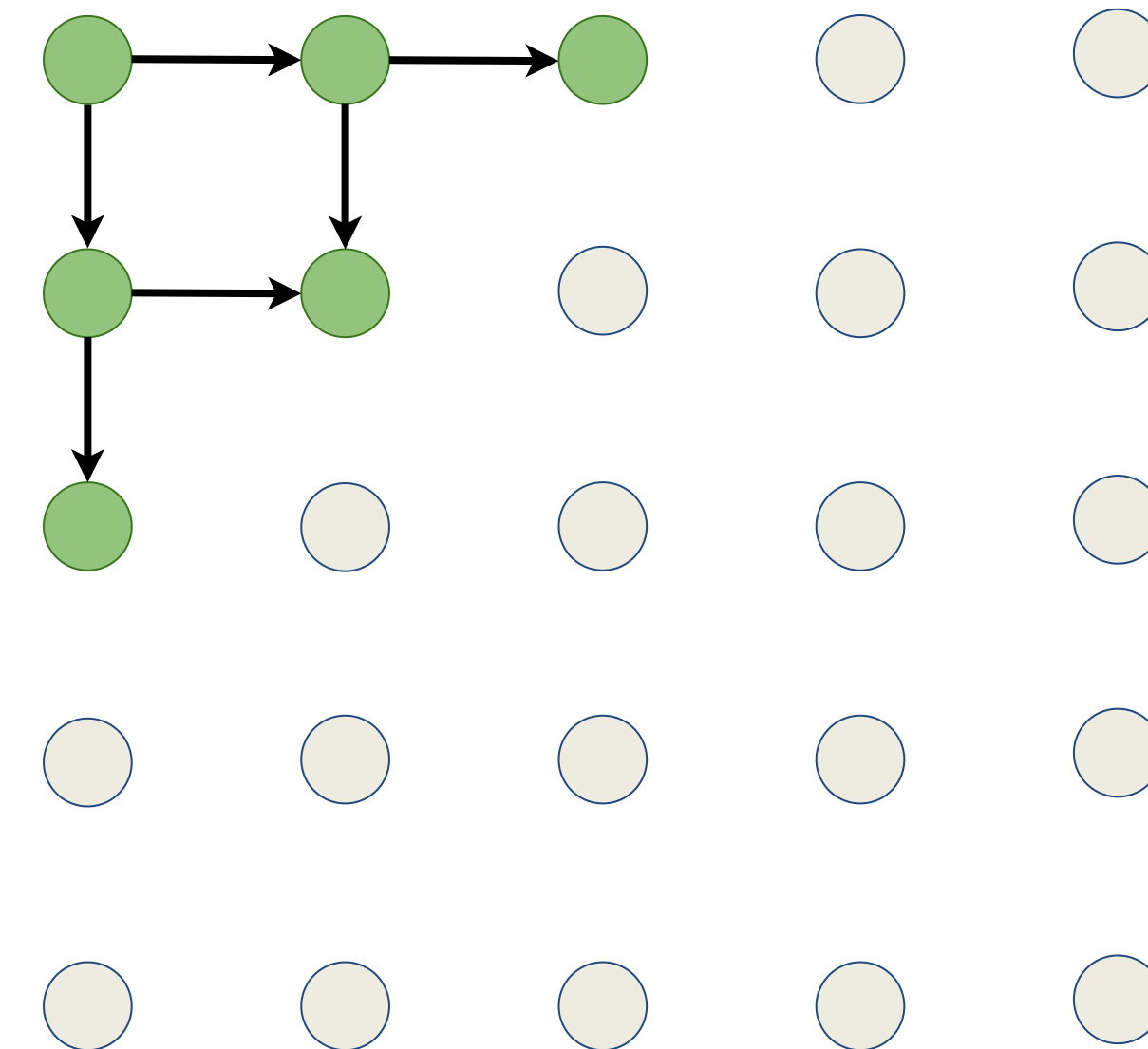from the corner


Dependency on previous pixels
model using an RNN (LSTM)

# **Pixel**RNN

Generate image pixels starting
from the corner


Dependency on previous pixels
model using an RNN (LSTM)

$x_1$      $x_n$

$x_i$

$x_{n^2}$

# **Pixel**RNN

Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

# **Pixel**RNN

Generate image pixels starting
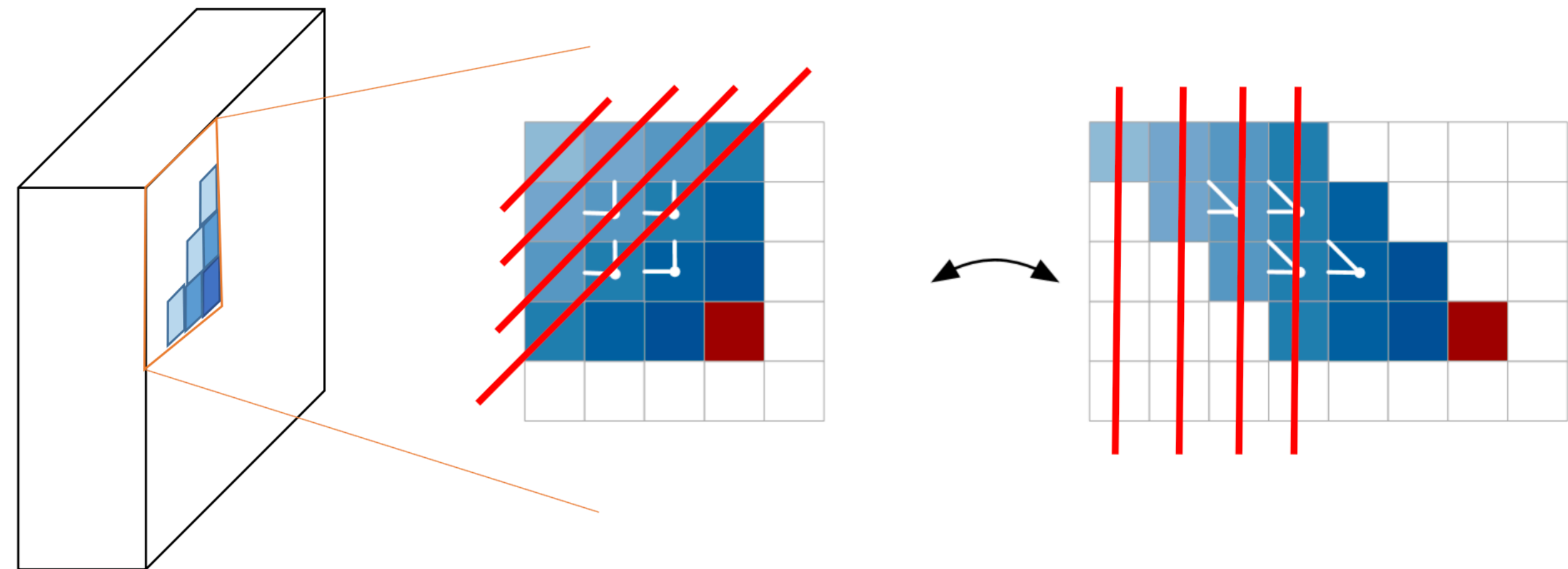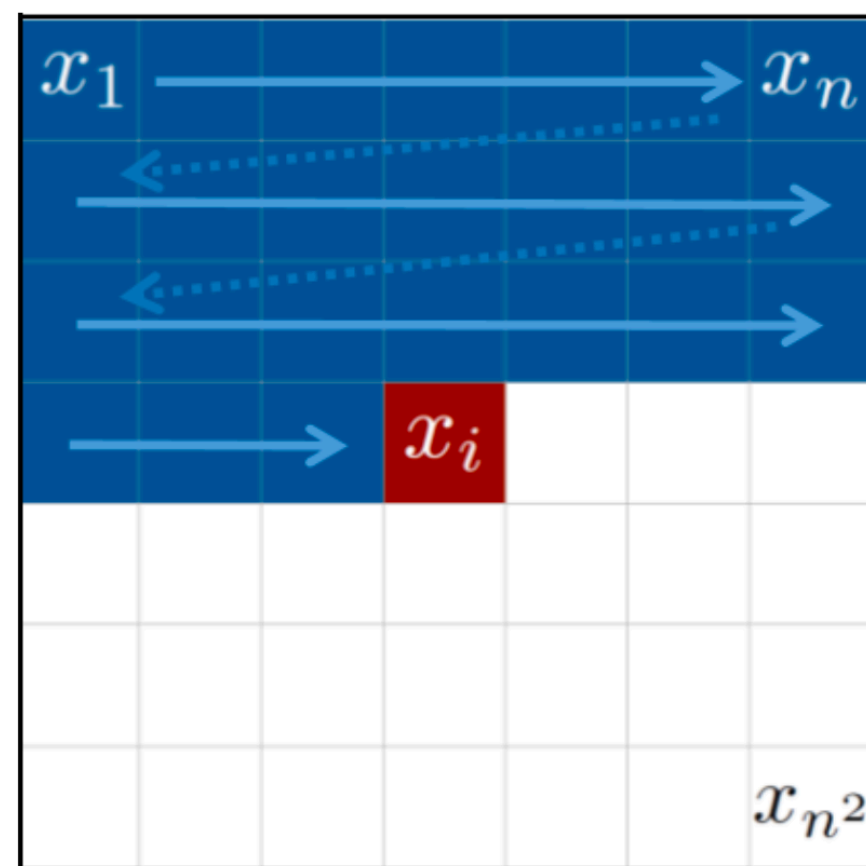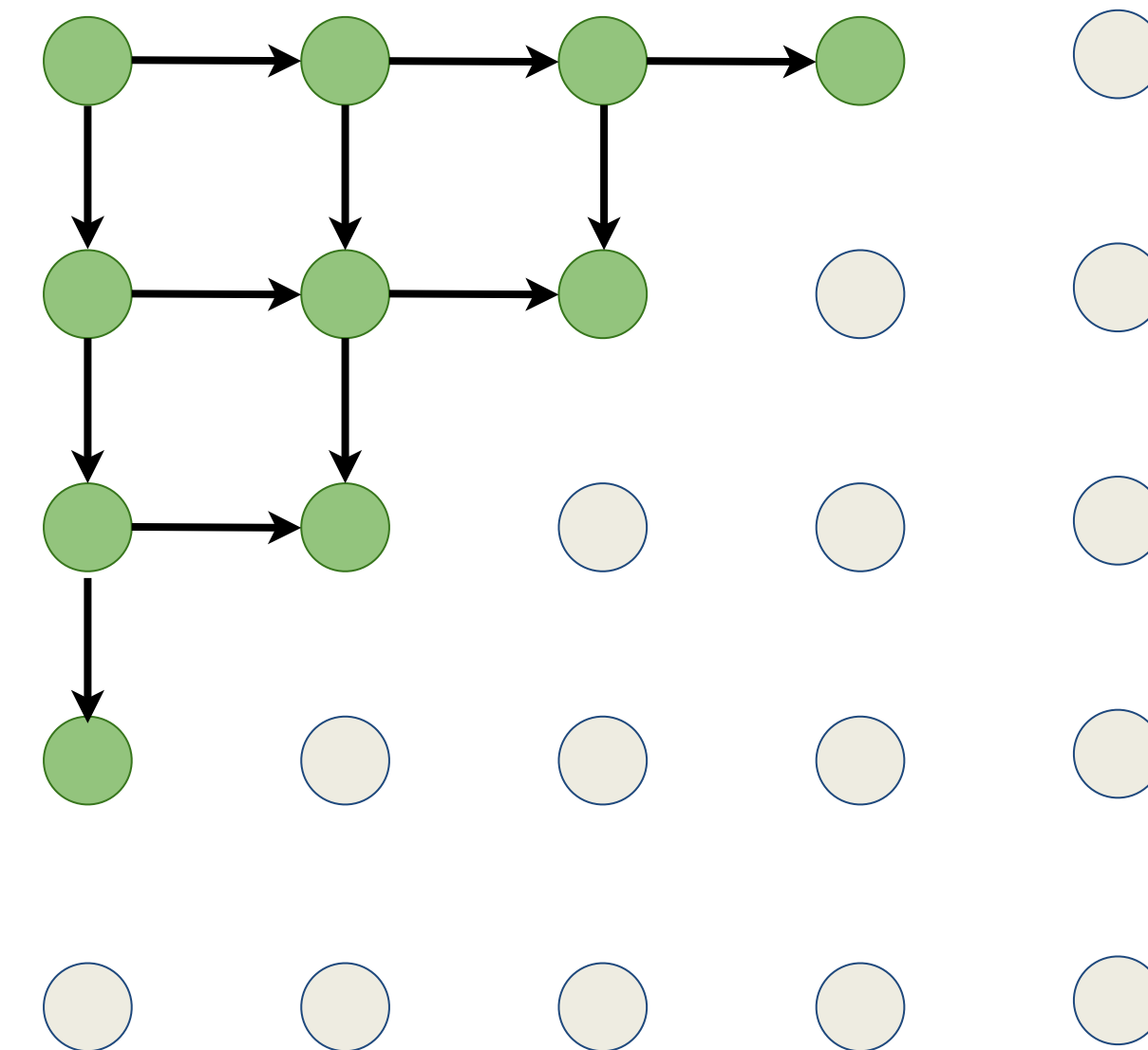from the corner

Dependency on previous pixels
model using an RNN (LSTM)



**Problem:** sequential generation is slow

# **Pixel**CNN

Still generate image pixels
starting from the corner

Dependency on previous pixels
now modeled using a CNN over
context region

# **Pixel**CNN

Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

**Training:** maximize likelihood of training images

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

**Softmax** loss at each pixel

# **Pixel**CNN

Still generate image pixels
starting from the corner

Dependency on previous pixels
now modeled using a CNN over
context region

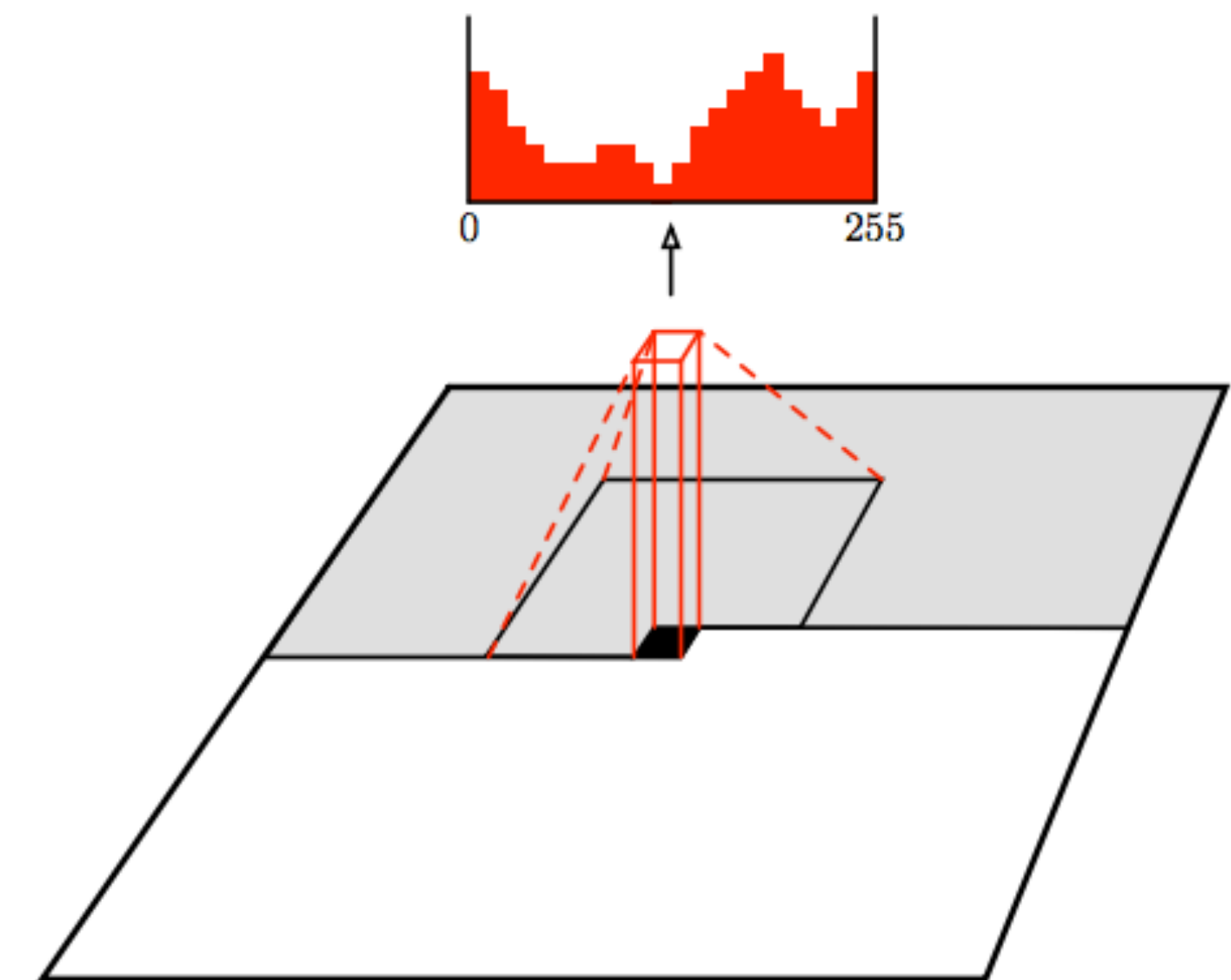**Training:** maximize likelihood of
training images

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$



Generation is still slow (sequential),
but learning is faster

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# **Generated** Samples

32x32 **CIFAR-10**



32x32 **ImageNet**

# PixelRNN and PixelCNN

**Pros:**

— Can explicitly compute likelihood p(x)
— Explicit likelihood of training data gives good
  evaluation metric
— Good samples

**Con:**

— Sequential generation => slow

**Improving** PixelCNN performance

— Gated convolutional layers
— Short-cut connections
— Discretized logistic loss
— Multi-scale
— Training tricks
— Etc…

# **Multi-scale** PixelRNN

Take sub-sampled pixels as additional input pixels

Can capture better global information (more visually coherent)

# Multi-scale PixelRNN

# **Conditional** Image Generation

Similar to PixelRNN/CNN but conditioned on a high-level image description vector $\mathbf{h}$

$$p(\mathbf{x}) = p(x_1, x_2, ..., x_{n^2})$$

$$\downarrow$$

$$p(\mathbf{x}|\mathbf{h}) = p(x_1, x_2, ..., x_{n^2}|\mathbf{h})$$

* slide from Hsiao-Ching Chang, Ameya Patil, Anand Bhattad

# **Conditional** Image Generation

African elephant



Sandbar

# Variational Autoencoders (VAE)

# So far …

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

# **So** far …

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent variables z (that we need to marginalize):

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

**cannot optimize directly**, derive and optimize lower bound of likelihood instead

# **Autoencoders** Reminder …

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

**Originally:** Linear + nonlinearity (sigmoid)
**Later:** Deep, fully-connected
**Later:** ReLU CNN

**Features** $\boxed{z}$

**Encoder**

**Input data** $\boxed{x}$

**Input** data

# **Autoencoders** Reminder …

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

$z$ usually smaller than $x$
(dimensionality reduction)

**Originally:** Linear + nonlinearity (sigmoid)
**Later:** Deep, fully-connected
**Later:** ReLU CNN

**Features** $z$

**Encoder**

**Input data** $x$

**Input** data

# **Autoencoders** Reminder ...

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

$z$ usually smaller than $x$
(dimensionality reduction)

Want features that capture
**meaningful** factors of variation

**Originally:** Linear + nonlinearity (sigmoid)
**Later:** Deep, fully-connected
**Later:** ReLU CNN

**Features** $z$

**Encoder**

**Input data** $x$

**Input** data

# **Autoencoders** Reminder …

Train such that features can reconstruct original data best they can

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

**Decoder**

**Features** $z$

**Encoder**

**Input data** $x$

**Input** data

# **Autoencoders** Reminder …

Train such that features can reconstruct original data best they can

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

**Decoder**

**Features** $z$

**Encoder**

**Input data** $x$

**Reconstructed** data



Encoder: 4-layer conv
Decoder: 4-layer upconv

**Input** data

# **Autoencoders** Reminder …

L2 Loss function:

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

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

**Decoder**

**Features** $z$

**Encoder**

**Input data** $x$

**Reconstructed** data



Encoder: 4-layer conv
Decoder: 4-layer upconv

**Input** data

# **Autoencoders** Reminder …

L2 Loss function:

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

Doesn't use labels!

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

**Decoder**

**Features** $z$

**Encoder**

**Input data** $x$

**Reconstructed** data

Encoder: 4-layer conv
Decoder: 4-layer upconv

**Input** data

# **Autoencoders** Reminder ...

**Loss function**
(e.g., softmax)

$\hat{y}$    $y$

**Classifier**

Fine-tune
encoder
jointly with
classifier

Train for **final task**
(sometimes with small data)

**Features**    $z$

bird        plane

dog        deer        truck

**Encoder**

**Input data**    $x$

# **Variational** Autoencoders

Probabilistic spin on autoencoder - will let us sample from the model to generate

Assume training data is generated from underlying unobserved (latent) representation z

Sample from
true **conditional**

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



Sample from
true **prior**

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

# **Variational** Autoencoders

Probabilistic spin on autoencoder - will let us sample from the model to generate

Assume training data is generated from underlying unobserved (latent) representation z

Sample from
true **conditional**

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



**Intuition:** $x$ is an image, $z$ is latent factors used to generate $x$ (e.g., attributes, orientation, *etc.*)

Sample from
true **prior**

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

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

Sample from
true **conditional**

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

$x$

$z$

Sample from
true **prior**

$p_{\theta*}(z)$

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta^*$ of this generative model

How do we **represent** this model?

Sample from
true **conditional**

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



$x$

$z$

Sample from
true **prior**

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

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **represent** this model?

Sample from
true **conditional**

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



Choose prior $p(z)$ to be simple, e.g., Gaussian

Reasonable for latent attributes, e.g., pose, amount of smile

Sample from
true **prior**

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

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta^*$ of this generative model

How do we **represent** this model?

Sample from
true **conditional**

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



**Decoder
Network**

Sample from
true **prior**

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

Choose prior $p(\boldsymbol{z})$ to be simple, e.g., Gaussian
Reasonable for latent attributes, e.g., pose, amount of smile

Conditional $p(\boldsymbol{x}|\boldsymbol{z})$ is complex (generates image)
Represent with Neural Network

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **train** this model?

Sample from
true **conditional**

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



**Decoder
Network**

Sample from
true **prior**

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

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **train** this model?

Sample from true **conditional**

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

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data



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

Sample from true **prior**

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

(now with latent $z$ that we need to marginalize)

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **train** this model?
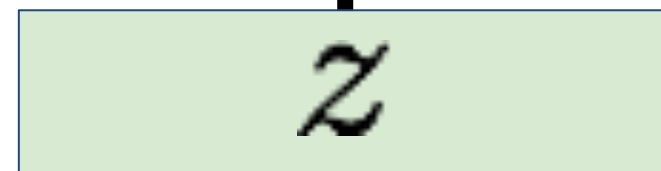
Sample from true **conditional**

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

Sample from true **prior**

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

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data

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

(now with latent $z$ that we need to marginalize)



**Decoder Network**

$x$

$z$

## What is the problem with this?

# **Variational** Autoencoders

We want to **estimate the true parameters** $\theta*$ of this generative model

Sample from true **conditional**

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



**Decoder Network**
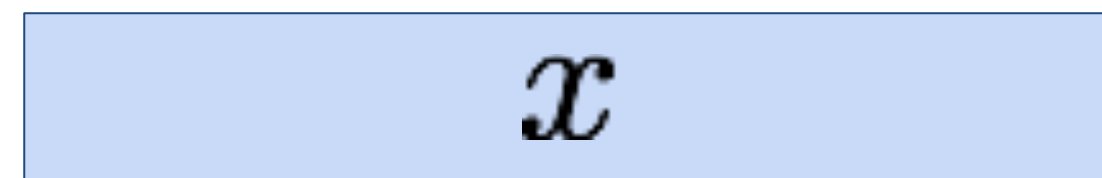
Sample from true **prior**

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

How do we **train** this model?

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data

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

(now with latent $z$ that we need to marginalize)

**Intractable !**

# **Intractability** in Variational Autoencoder

Data **likelihood**:

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

# **Intractability** in Variational Autoencoder

Data **likelihood**: $\qquad p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$

🙂

Simple **Gaussian** Prior

# **Intractability** in Variational Autoencoder

**Decoder** Neural Network

🙂

Data **likelihood**:  $p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$

🙂

Simple **Gaussian** Prior

# **Intractability** in Variational Autoencoder

Intractable to compute for every z

**Decoder** Neural Network

🙂

☹️

Data **likelihood**:
$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

🙂

Simple **Gaussian** Prior

# **Intractability** in Variational Autoencoder

Intractable to compute for every z

**Decoder** Neural Network

Data **likelihood**: $p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$

Simple **Gaussian** Prior

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

# **Intractability** in Variational Autoencoder

Intractable to compute for every z

**Decoder** Neural Network

Data **likelihood**:

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

Simple **Gaussian** Prior

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

# **Intractability** in Variational Autoencoder

Intractable to compute for every z

**Decoder** Neural Network

☹️

🙂

Data **likelihood**: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

🙂

Simple **Gaussian** Prior

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

**Solution:** In addition to decoder network modeling $p_\theta(x|z)$, define additional

encoder network $q_\phi(z|x)$ that approximates $p_\theta(z|x)$

— Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

# **Variational** Autoencoder

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



Mean and (diagonal) covariance of $z \mid x$

Mean and (diagonal) covariance of $x \mid z$

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

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

**Encoder** Network

**Decoder** Network

$$q_\phi(z|x)$$

$$p_\theta(x|z)$$

(parameters φ)

(parameters θ)

$$x$$

$$z$$

# **Variational** Autoencoder

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



Why?   Mean and (**diagonal**) covariance of $z \mid x$

Mean and (**diagonal**) covariance of $x \mid z$

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

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

**Encoder** Network

$q_\phi(z|x)$

(parameters φ)

$x$

**Decoder** Network

$p_\theta(x|z)$

(parameters θ)

$z$

# **Variational** Autoencoder

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



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

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

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

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

**Encoder** Network

$q_\phi(z|x)$

(parameters φ)

**Decoder** Network

$p_\theta(x|z)$

(parameters θ)

$x$

$z$

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

Taking expectation with respect to z
(using encoder network) will come in
handy later

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \qquad \text{(Bayes' Rule)}$$

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$
\begin{aligned}
\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
&= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)} \\
&= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)} \\
&= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}
\end{aligned}
$$

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})}\right] \quad (\text{Logarithms})$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))$$

Expectation with respect to z
(using encoder network) leads to nice KL terms

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))$$

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through **reparam. trick**, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice **closed-form solution**!

$p_\theta(z|x)$ **intractable** (saw earlier), can't compute this KL term :(

But we know KL divergence always >= 0.

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \mid\mid p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + D_{KL}(q_\phi(z \mid x^{(i)}) \mid\mid p_\theta(z \mid x^{(i)}))$$

**Tractable lower bound** which we can take gradient of
and optimize! (pθ(x|z) differentiable, KL term differentiable)

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z \mid x^{(i)}))$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("**ELBO**")

**Training:** Maximize lower bound

$$\theta^*, \phi^* = \arg\max_{\theta, \phi} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$$

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

**Reconstruct Input Data**

**Make approximate posterior close to the prior**

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("**ELBO**")

**Training:** Maximize lower bound

$$\theta^*, \phi^* = \arg\max_{\theta,\phi} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$$

# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower
bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Lets look at **computing the bound** (forward pass)
for a given mini batch of input data

**Input** Data $\quad\boxed{x}$

# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower
bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

**Encoder** network

$q_\phi(z|x)$

**Input** Data

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

$x$

# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower
bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate
posterior distribution
close to prior

**Encoder** network

$$q_\phi(z|x)$$

**Input** Data

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

$$x$$

# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower
bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate
posterior distribution
close to prior

$z$

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

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

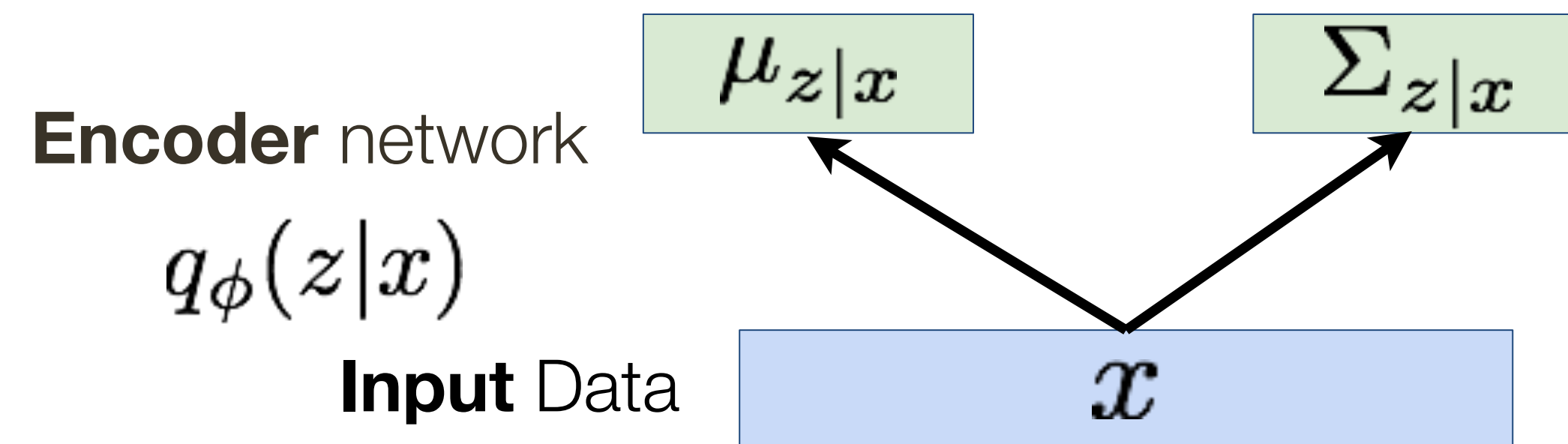**Encoder** network

$q_\phi(z|x)$

**Input** Data    $x$

# **Variational** Autoencoder: Learning

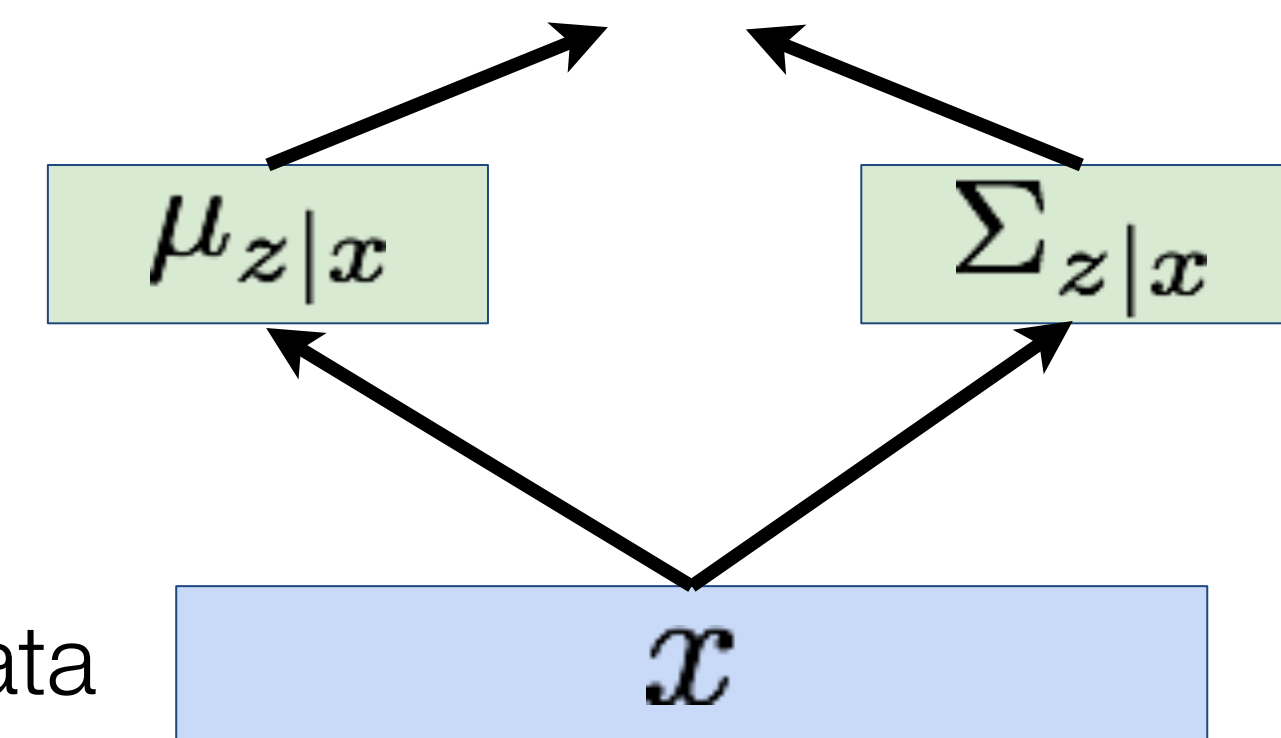**Putting it all together**:
maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

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

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

$$z$$

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

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

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

**Input** Data $\qquad x$

# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower bound

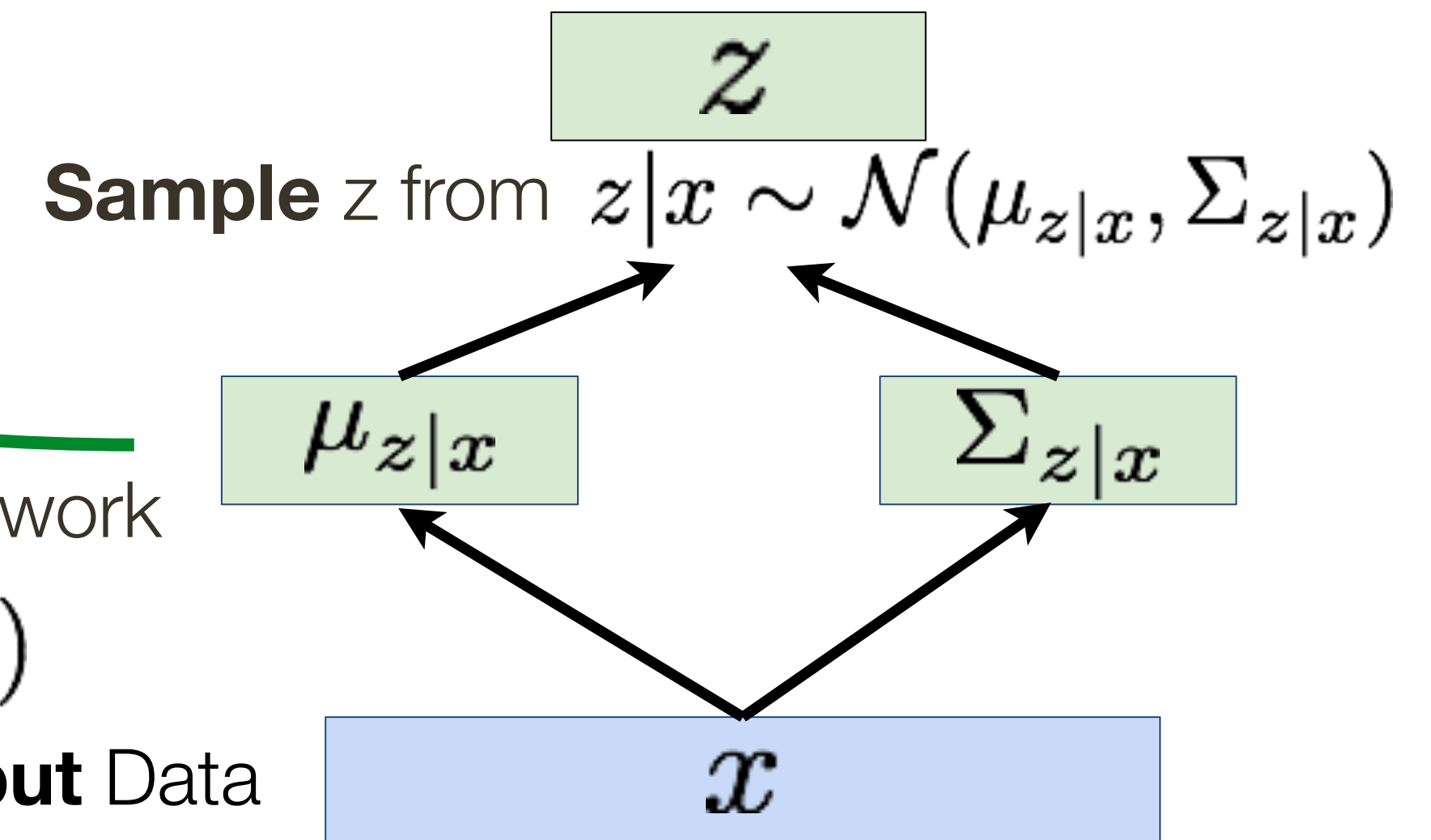Maximize likelihood of original input being reconstructed

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

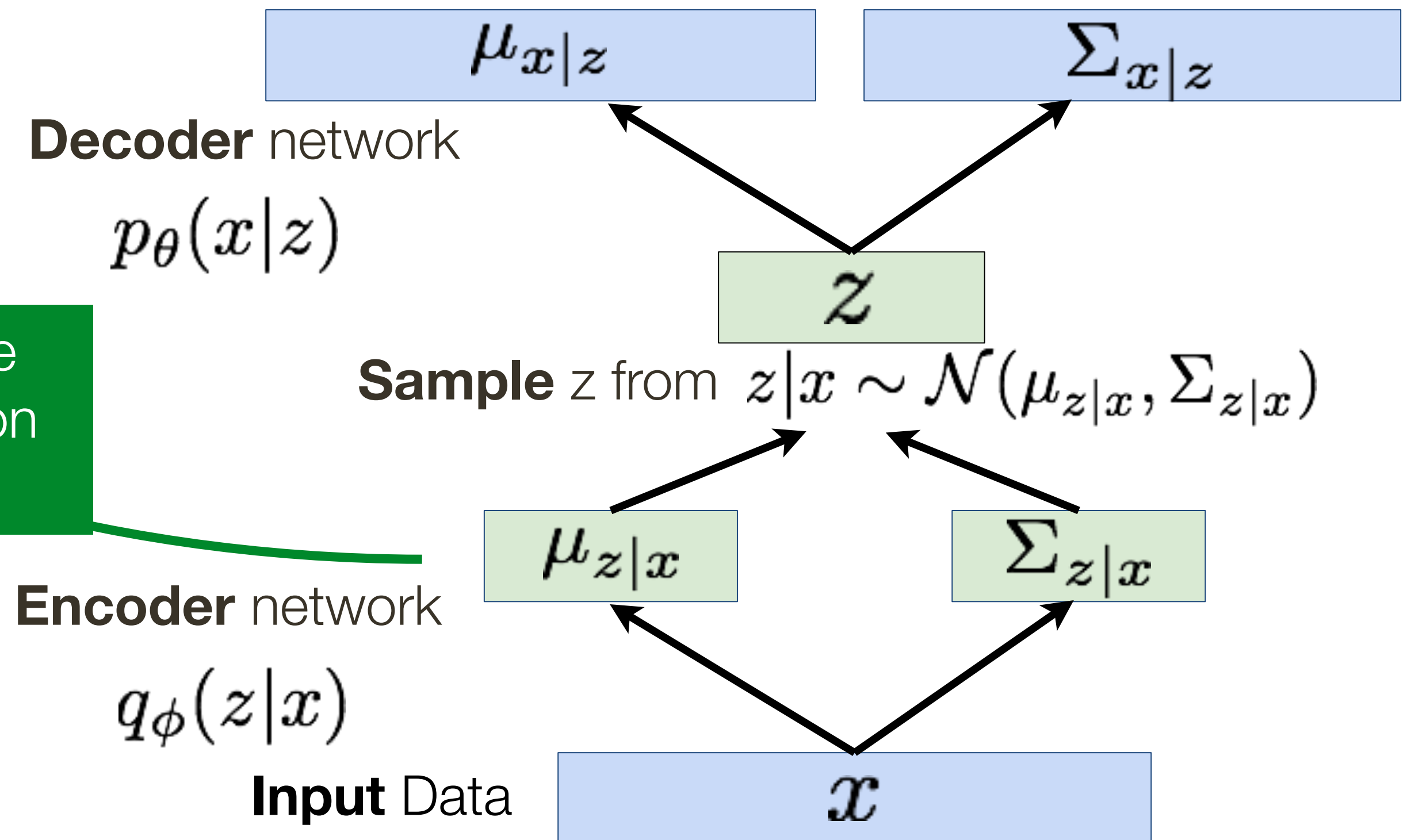Make approximate posterior distribution close to prior

$\hat{x}$

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

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

**Encoder** network

$q_\phi(z|x)$

**Input** Data     $x$
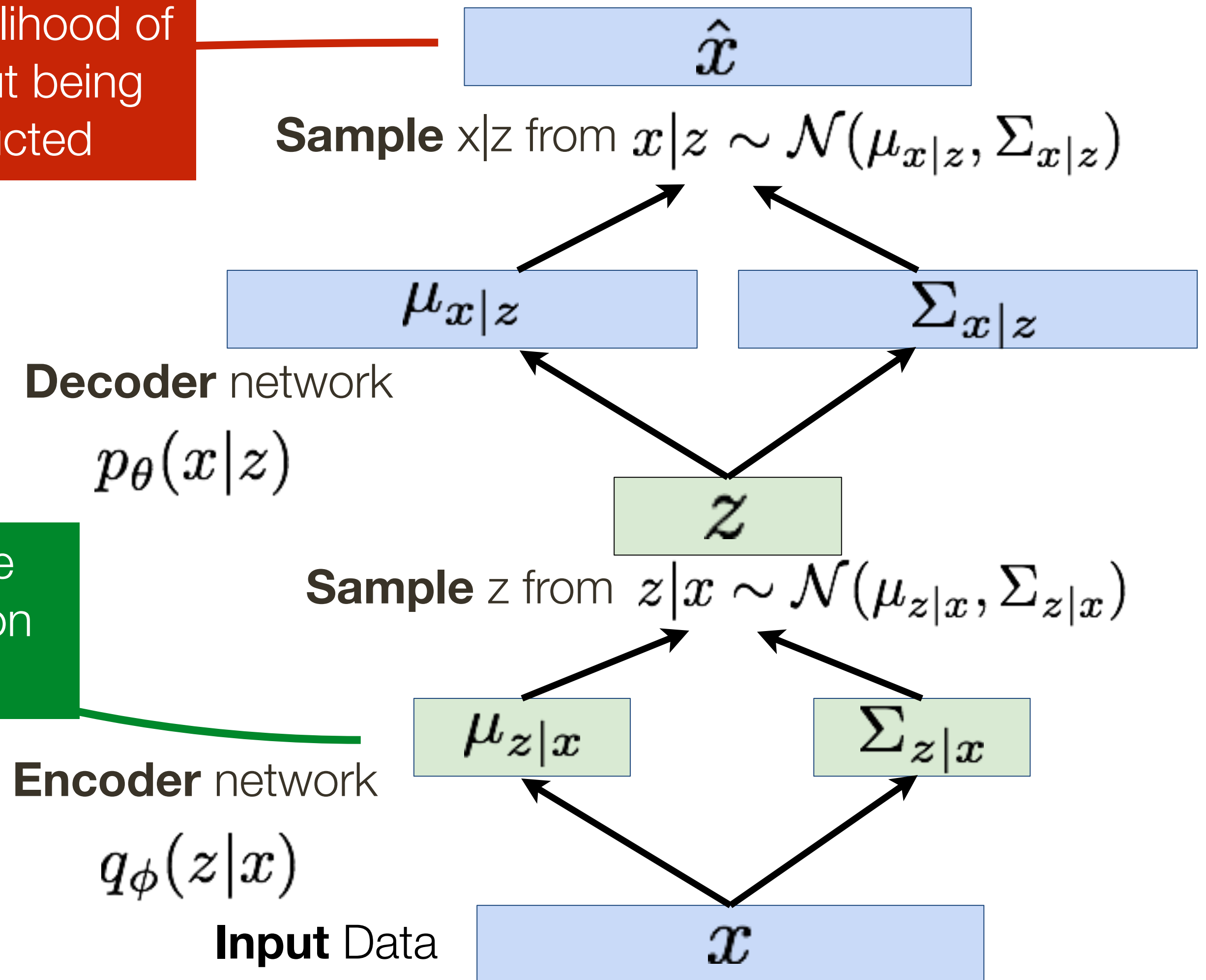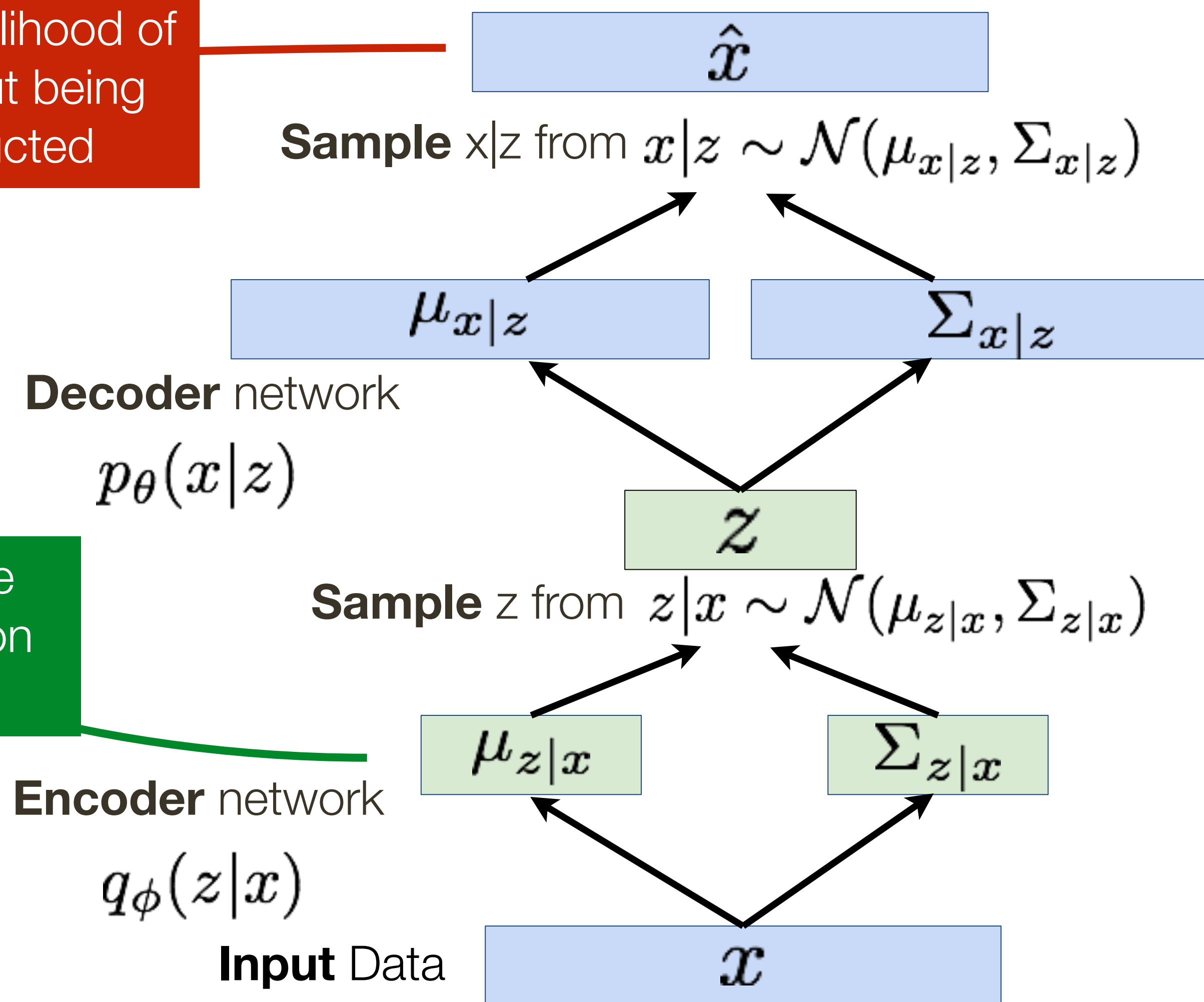
# **Variational** Autoencoder: Learning

**Putting it all together**:
maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

$\hat{x}$

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

$\mu_{z|x}$　$\Sigma_{z|x}$

**Encoder** network

$q_\phi(z|x)$

**Input** Data　$x$

For every minibatch of input data: compute this forward pass, and then backprop!

# **Variational** Autoencoder: Generating Data

Use decoder network and sample z from **prior**



$\hat{x}$

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

$\mu_{x|z}$ $\quad$ $\Sigma_{x|z}$
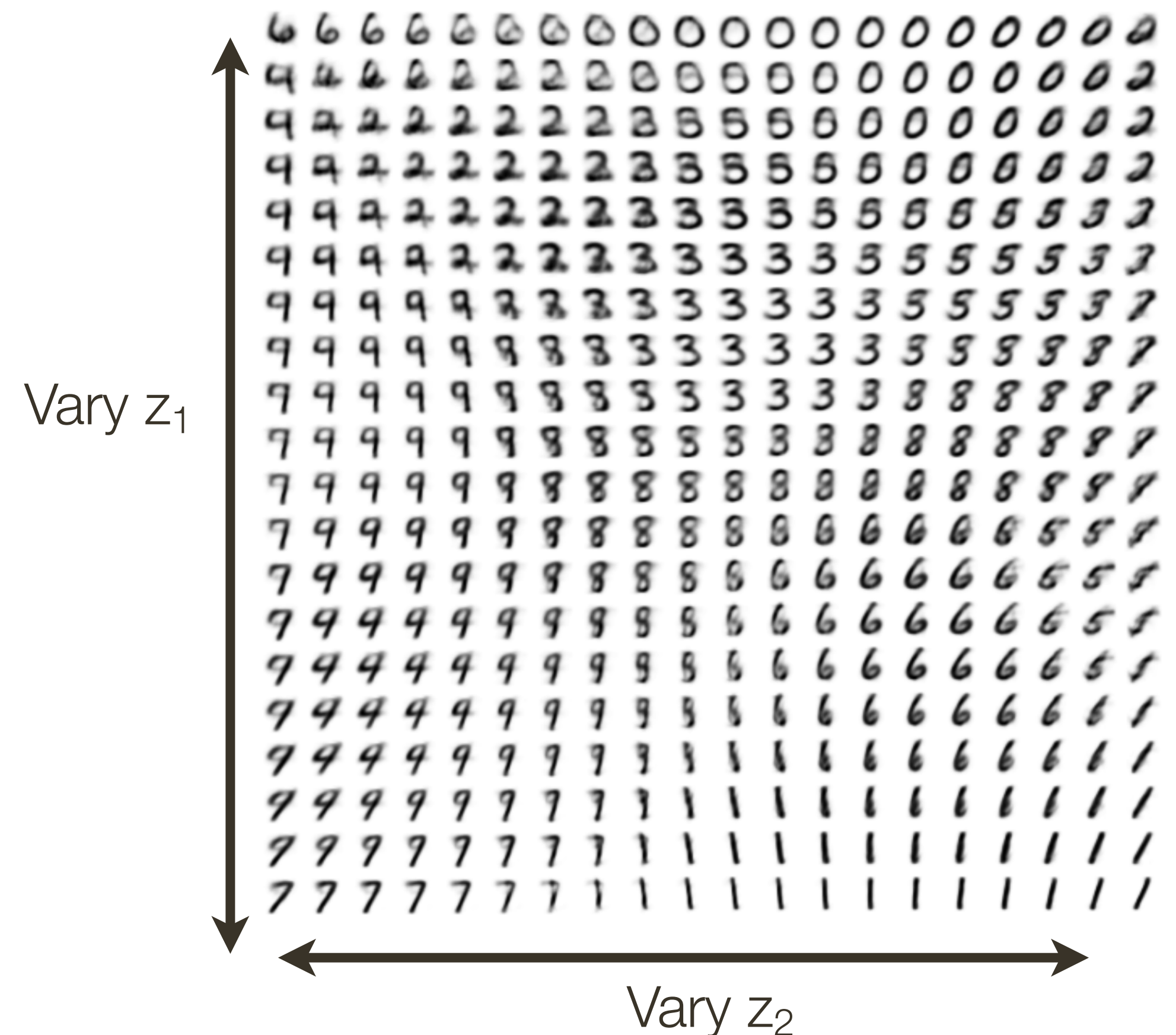
**Decoder** network
$p_\theta(x|z)$

$z$

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

# **Variational** Autoencoder: Generating Data

Use decoder network and sample z from **prior**



**Data manifold** for 2-d z

$\hat{x}$

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

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

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

$z$

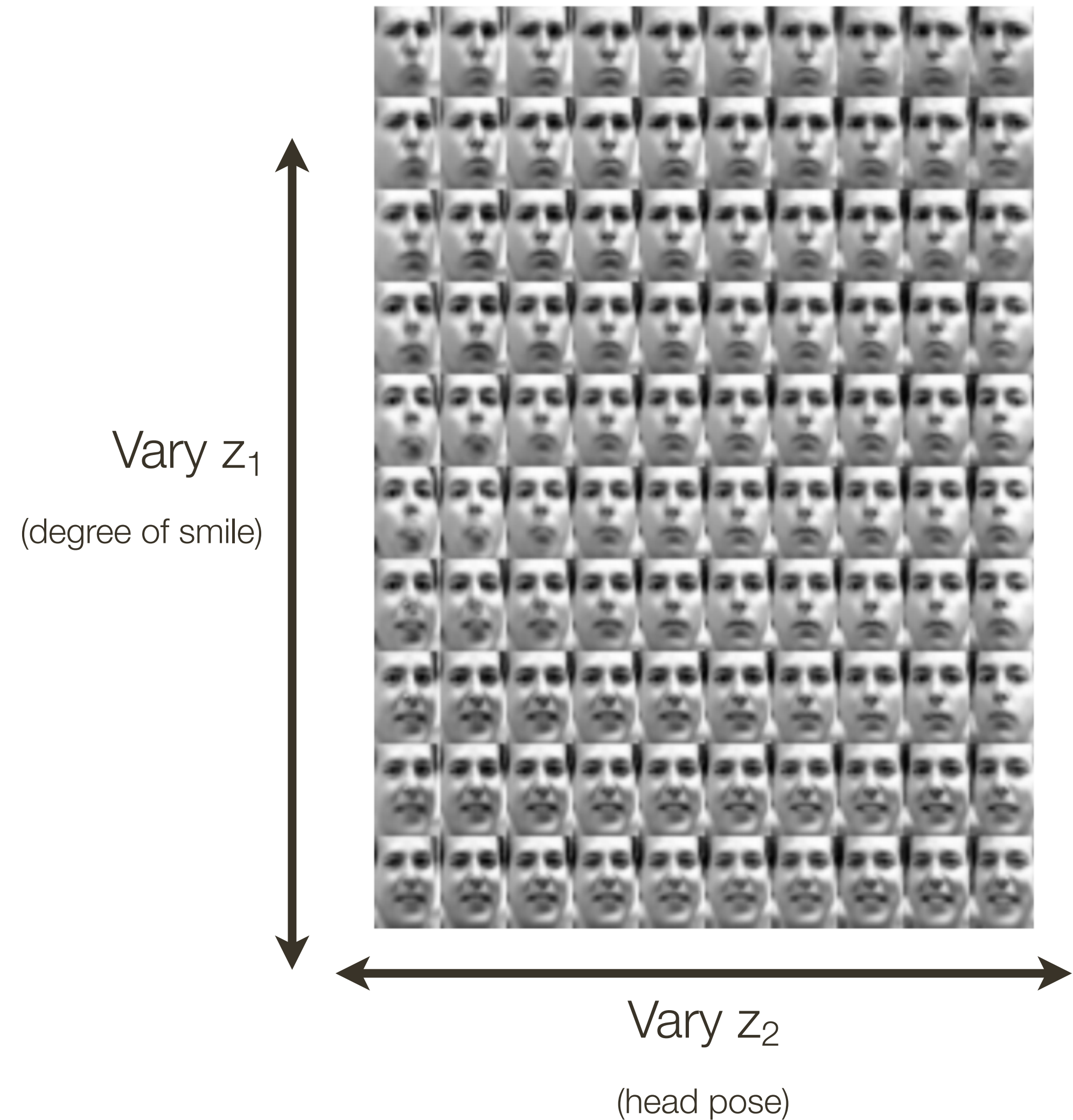**Sample** z from $z \sim \mathcal{N}(0, I)$

Vary $z_1$

Vary $z_2$

# **Variational** Autoencoder: Generating Data

Diagonal prior on z =>
independent latent variables

Different dimensions of z encode
interpretable factors of variation

**Data manifold** for 2-d z
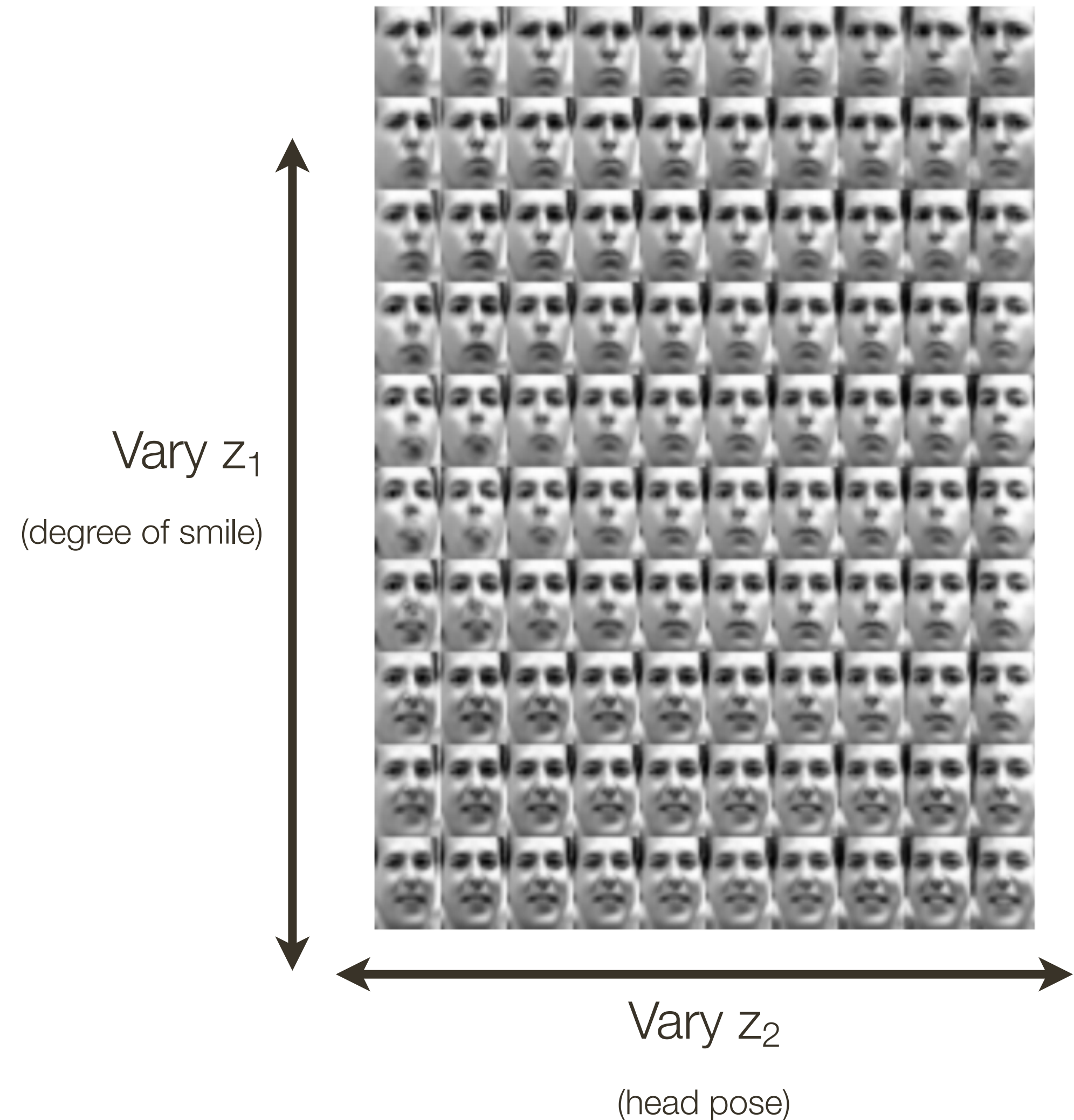


Vary $z_1$

(degree of smile)

Vary $z_2$

(head pose)

# **Variational** Autoencoder: Generating Data

Diagonal prior on z =>
independent latent variables

Different dimensions of z encode
interpretable factors of variation

Also good feature representation that can
be computed using $q_\phi(z|x)$!

**Data manifold** for 2-d z



Vary $z_1$

(degree of smile)

Vary $z_2$

(head pose)

# **Variational** Autoencoder: Generating Data



32x32 CIFAR-10



Labeled Faces in the Wild

# Conditional VAEs

(a) Frame 1

(b) Frame 2 (ground truth)

(c) Frame 2 (Sample 1)

(d) Frame 2 (Sample 2)

# **Variational** Autoencoders

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 area** of research:
- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
- Incorporating structure in latent variables (our submission to CVPR)