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

**Lecture 16: Generative Models [part 2]**

# Logistics

**Project Proposals Presentation Slides** due **Today 11:59pm**

— Grades and comments by Monday (sorry!)

**Project Proposal Document** due **Tuesday, November 15th**

**Assignment 4** is due today **Today 11:59pm**

**Leftovers until the end of term:**

— Assignment 5

— Project

— Paper presentation

# 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 $x$

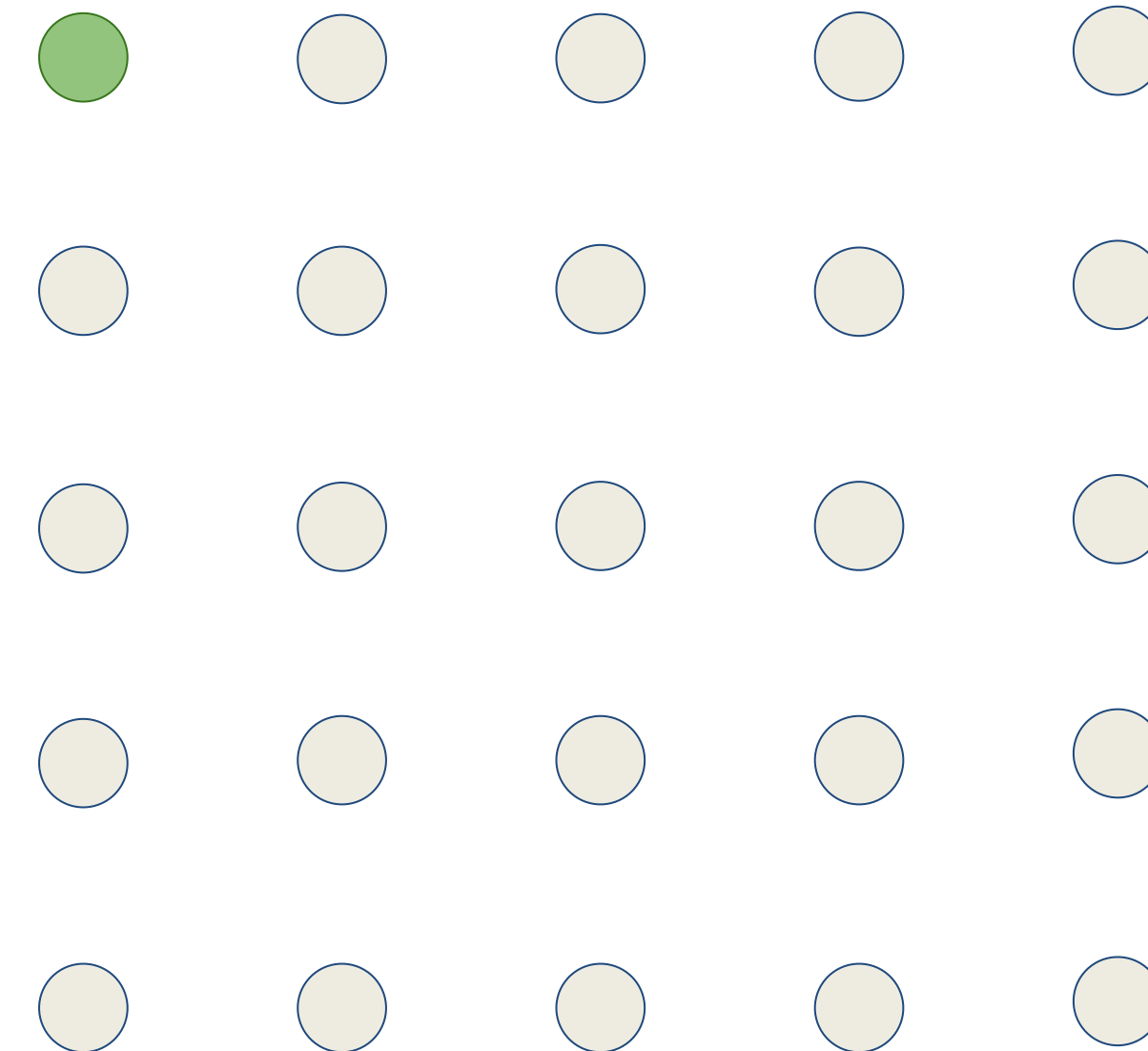Probability of i'th pixel value given all previous pixels

then maximize likelihood of training data

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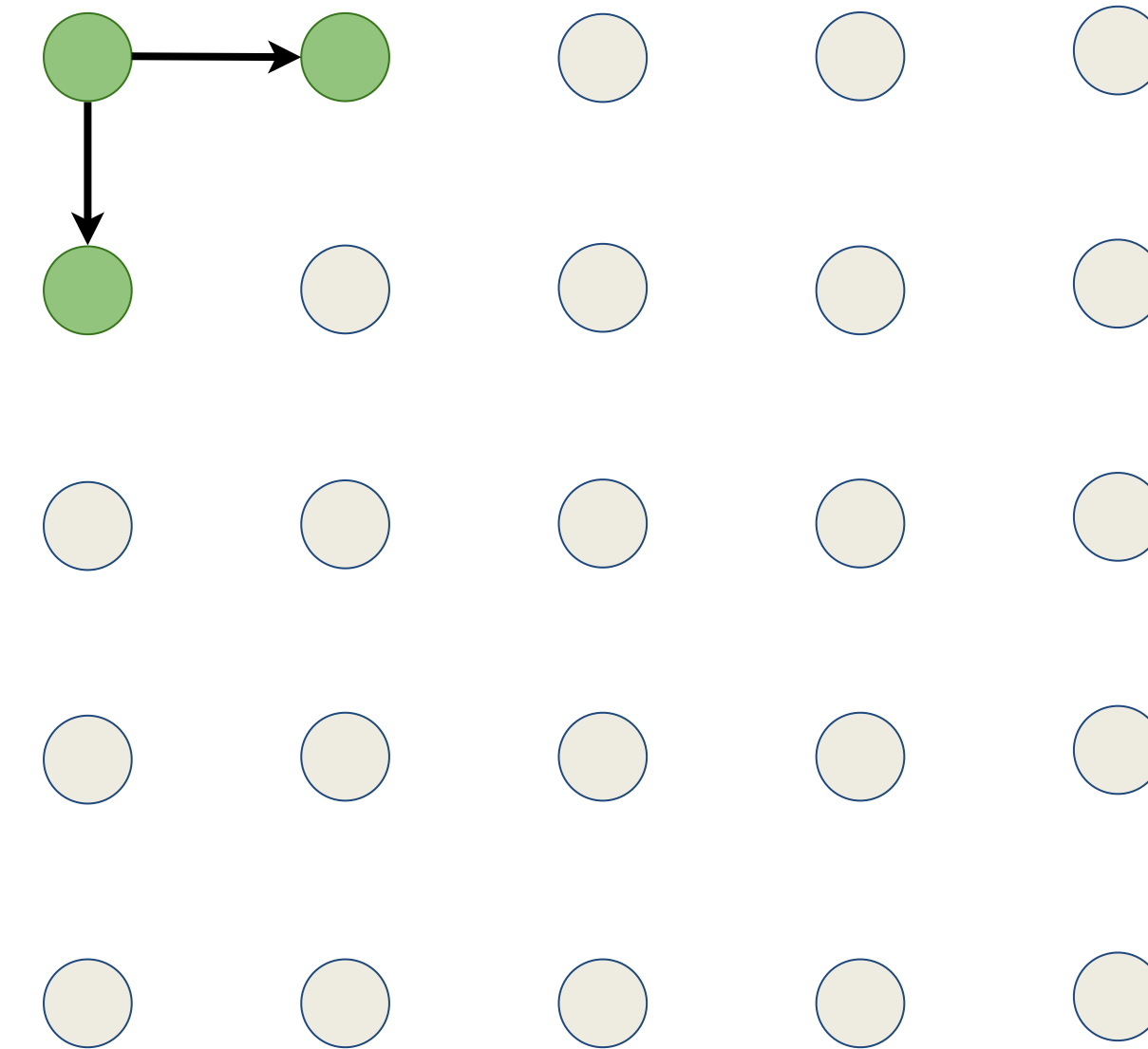
Also requires defining **ordering** of "previous pixels"

* slide from Fei-Fei 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
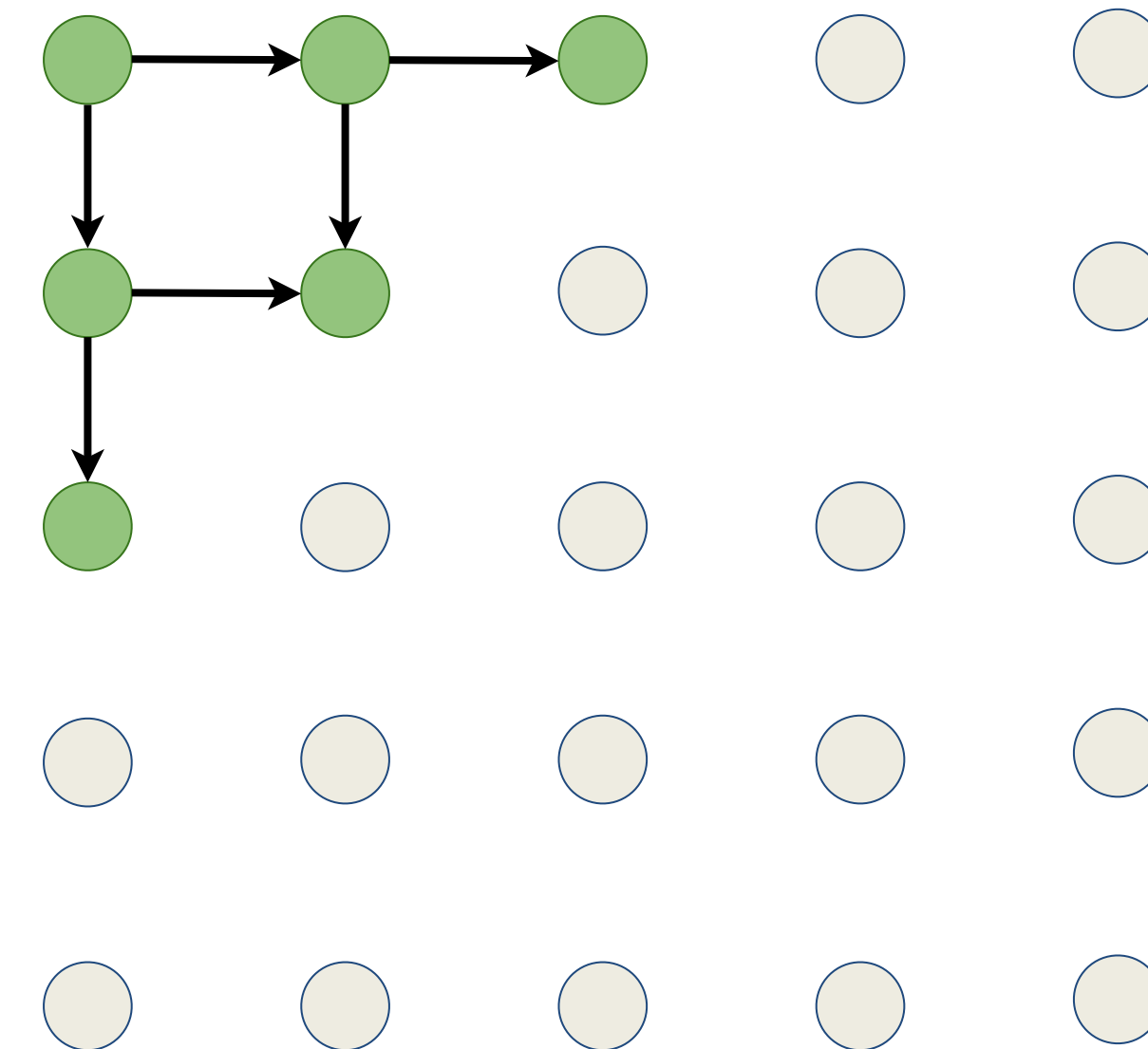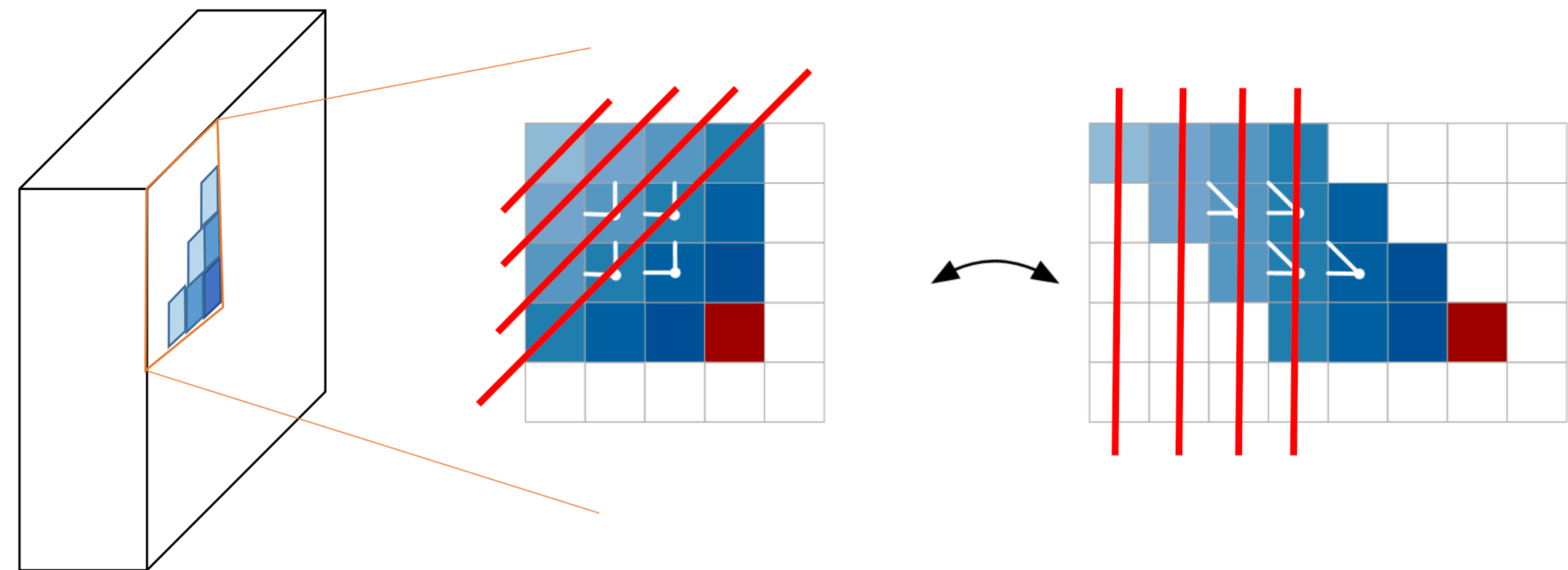
Dependency on previous pixels model using an RNN (LSTM)

# **Pixel**RNN

Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)
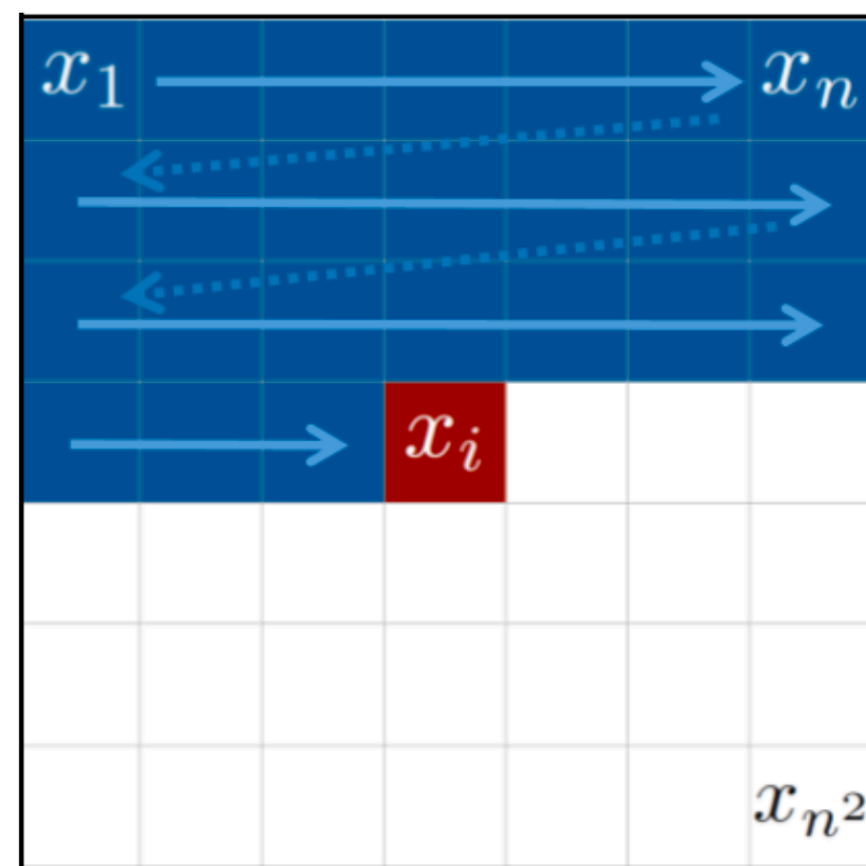
# **Pixel**RNN

$x_1$ $x_n$

$x_i$

$x_{n^2}$

# **Pixel**RNN

Generate image pixels starting
from the corner

Dependency on previous pixels
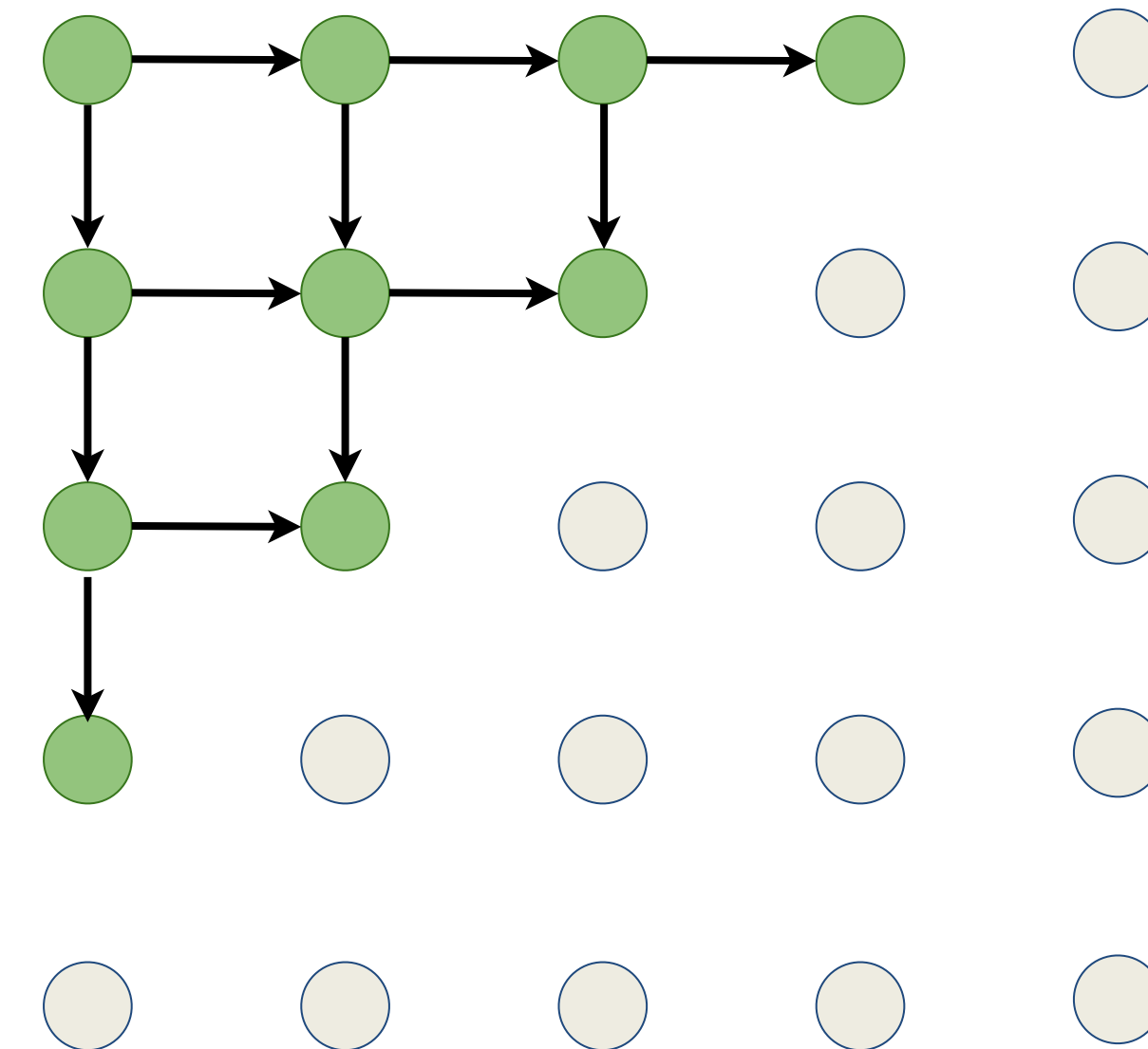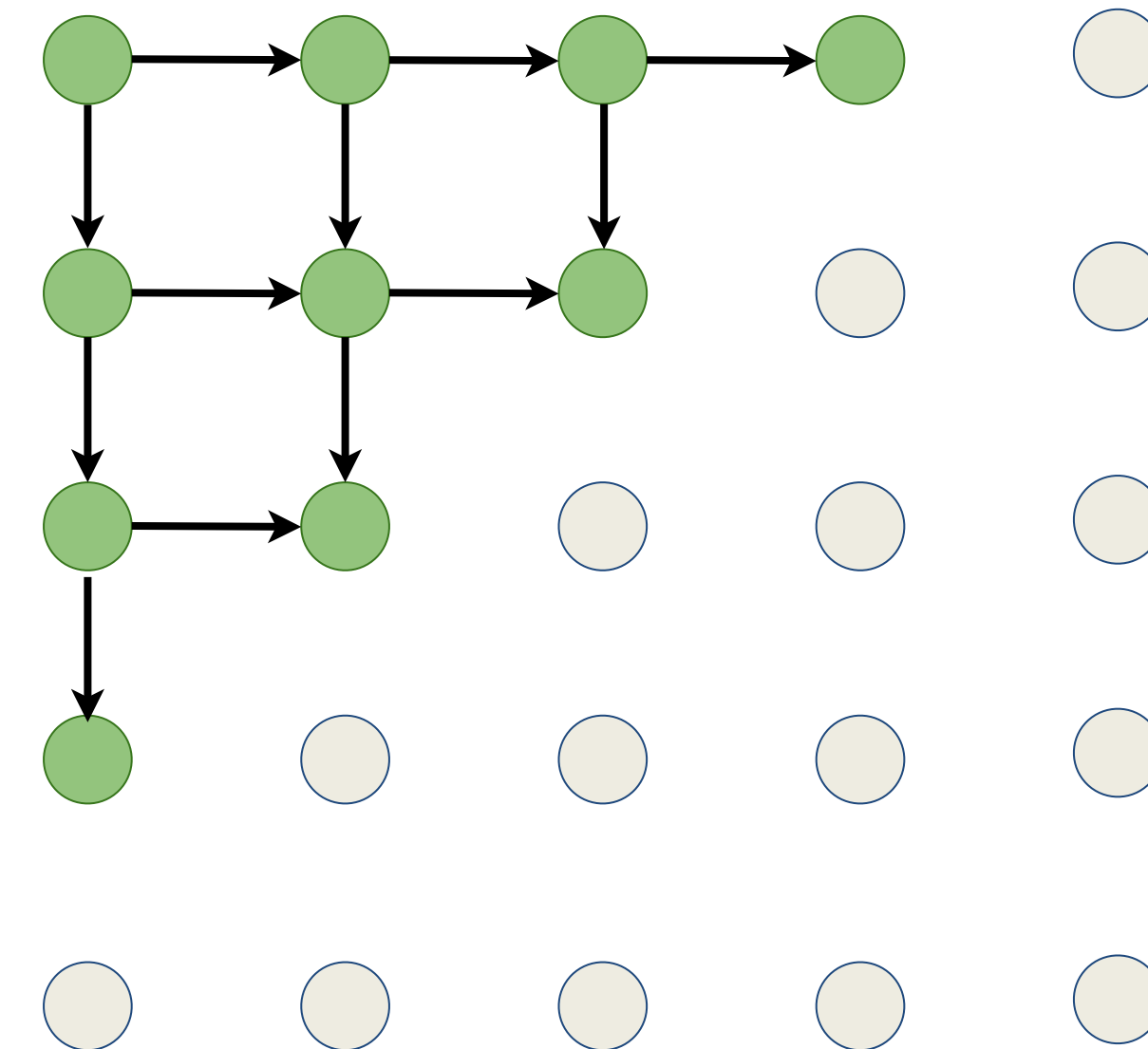model using an RNN (LSTM)

# **Pixel**RNN

Generate image pixels starting
from the corner

Dependency on previous pixels
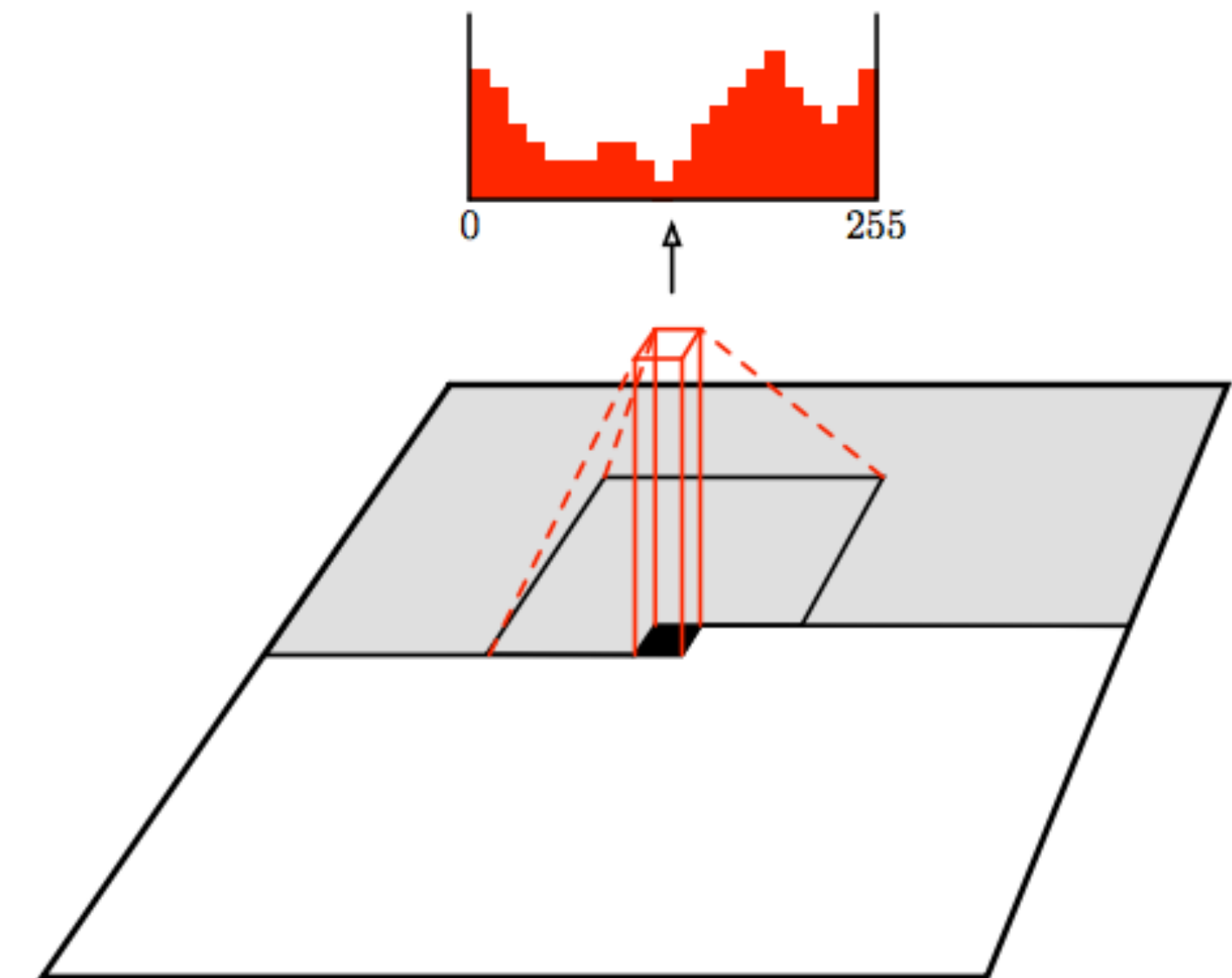model using an RNN (LSTM)



**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
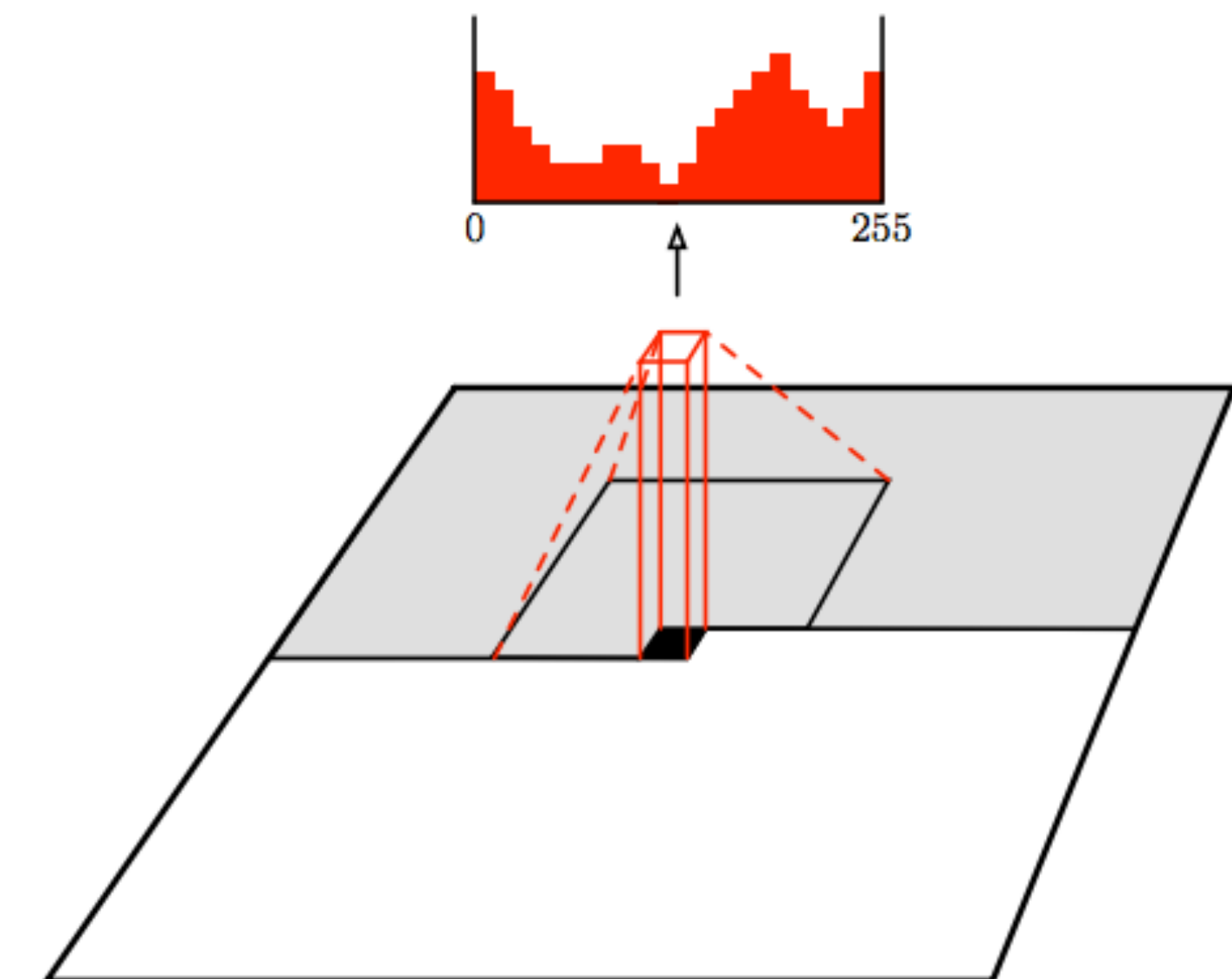context region

# **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



0                    255

# 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…

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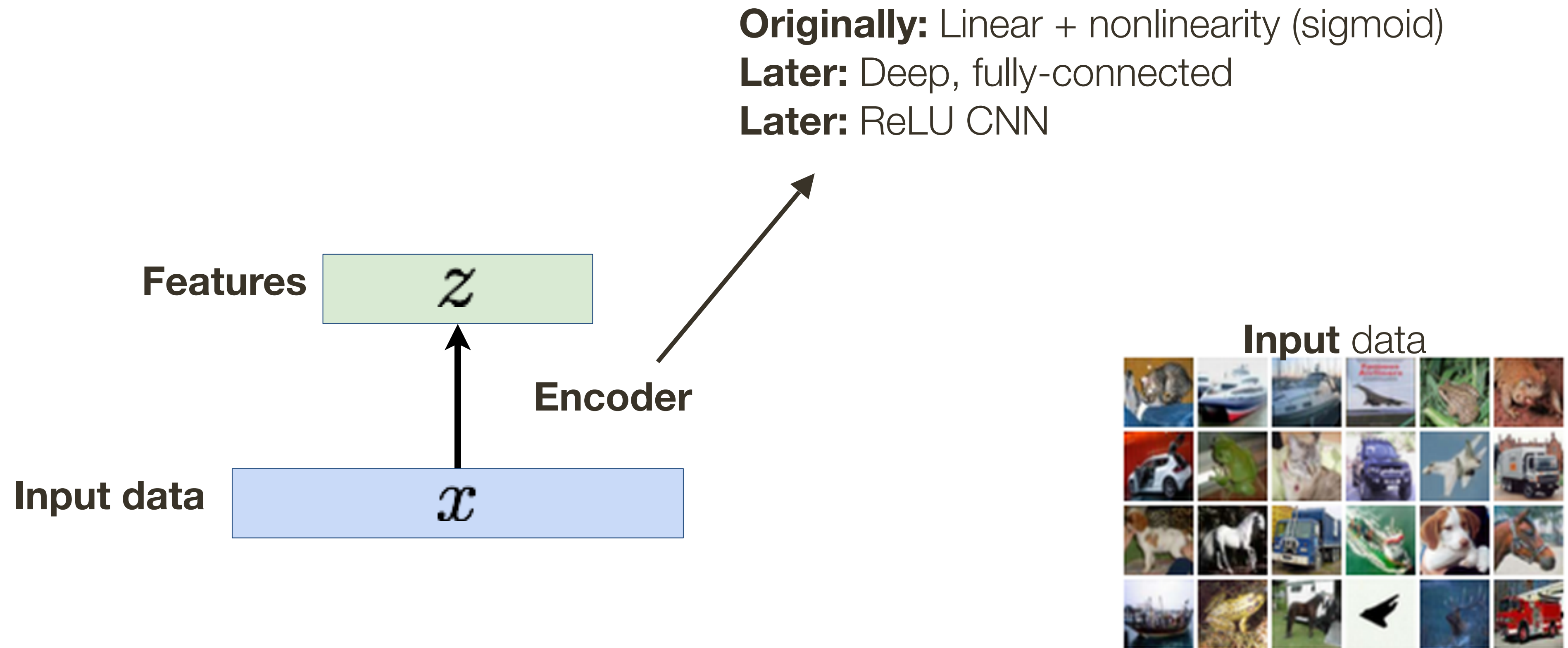
$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

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** $z$

**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)

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

**Features** $\quad z$

**Encoder**

**Input** data

**Input data** $\quad x$

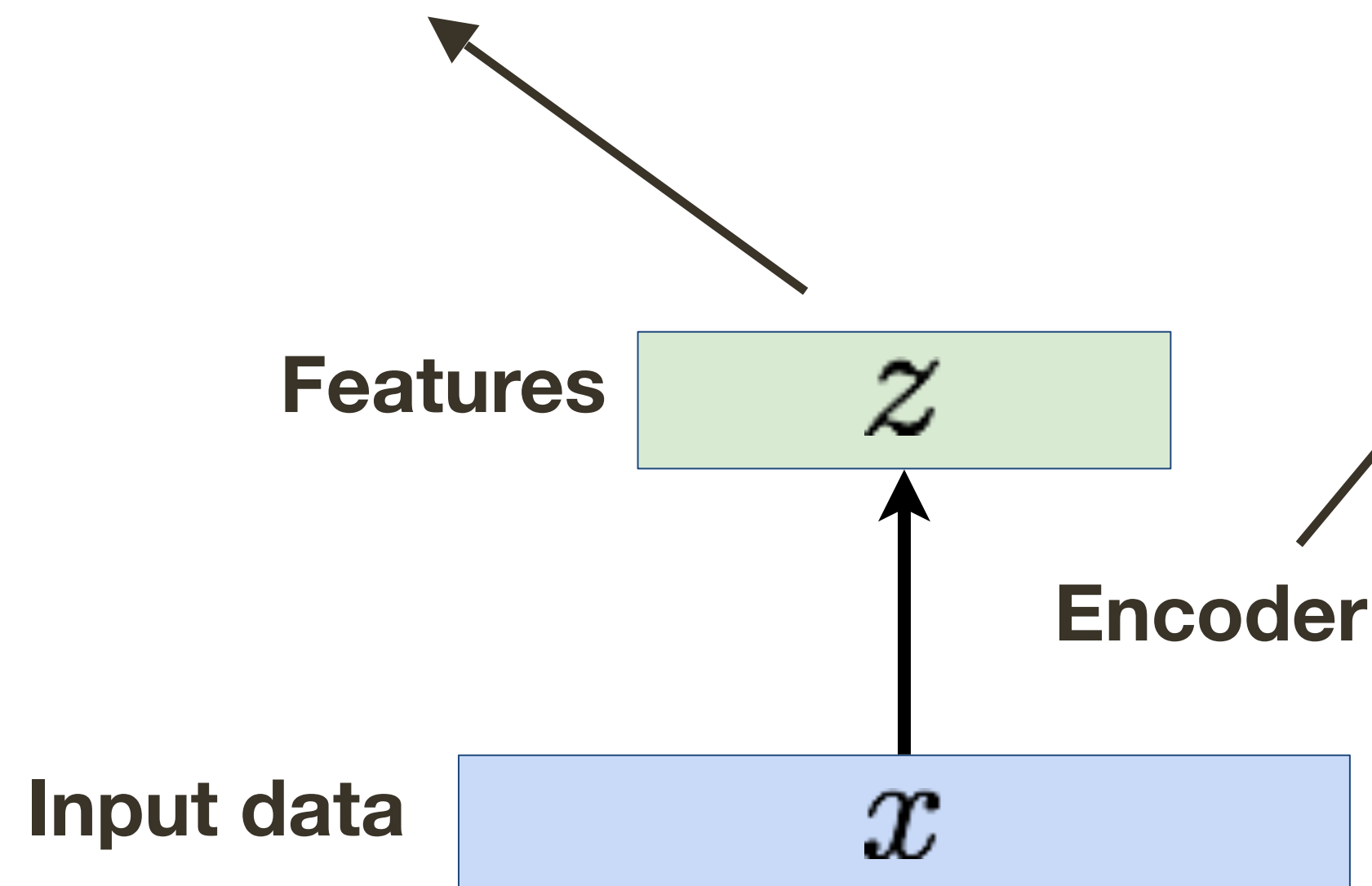# **Autoencoders** Reminder …

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
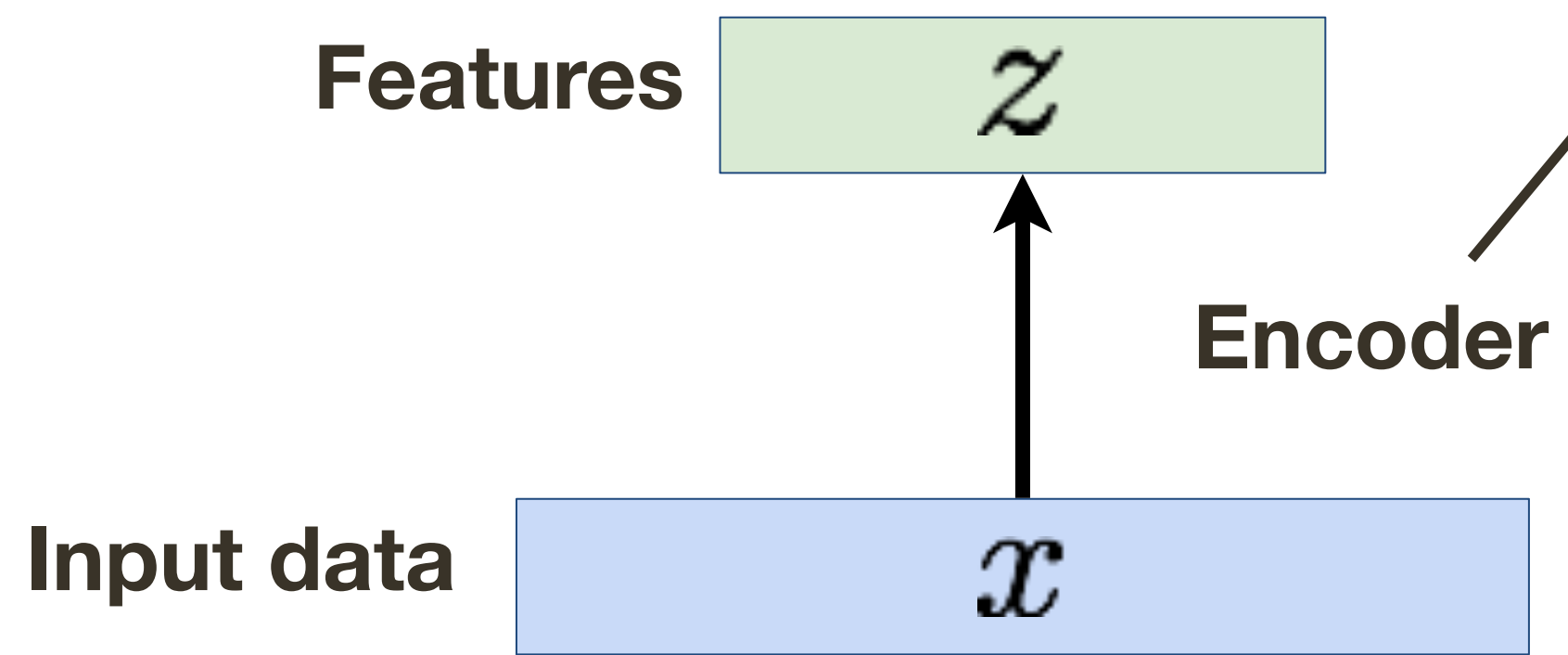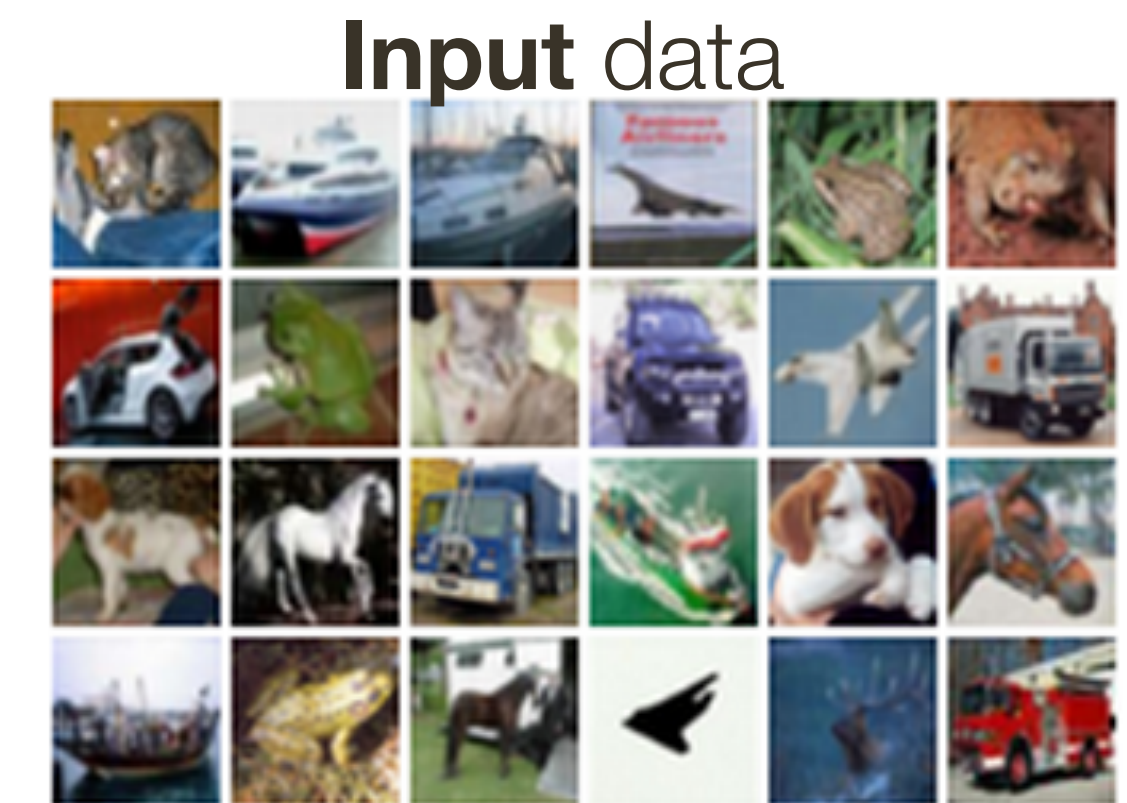**Later:** ReLU CNN

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

**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
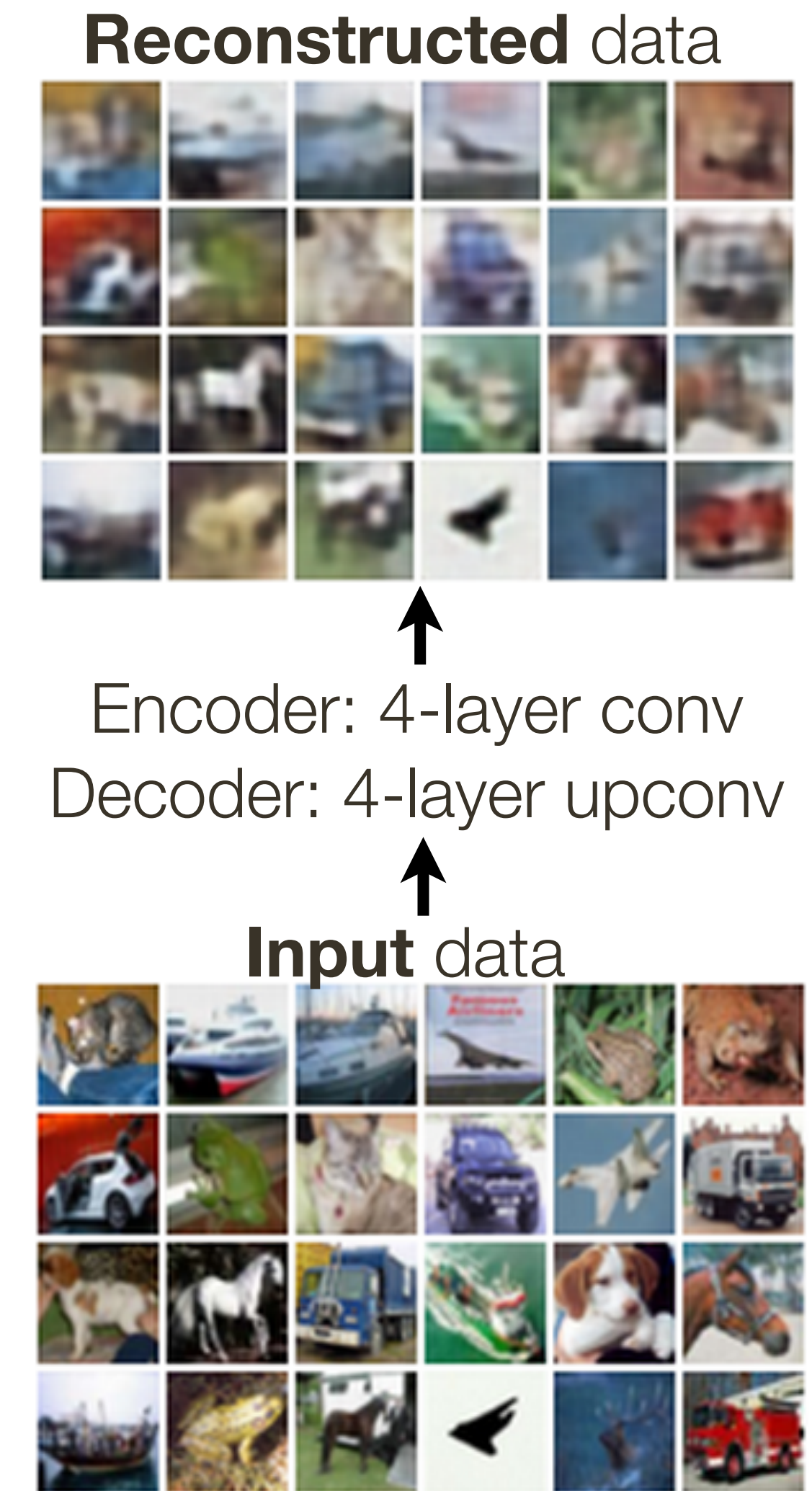Decoder: 4-layer upconv

**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

**Features**    $z$

**Encoder**

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

bird        plane

dog        deer        truck

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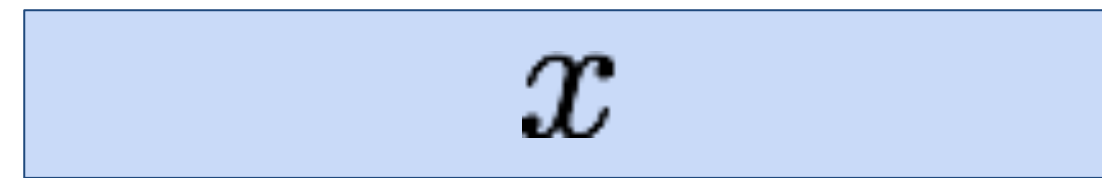
Probabilistic spin on autoencoder - will let us sample from the model to generate

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

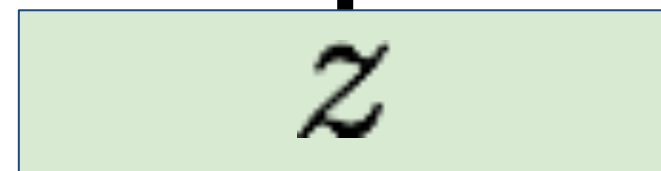Sample from
true **conditional**

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



Sample from
true **prior**

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

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

# **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**

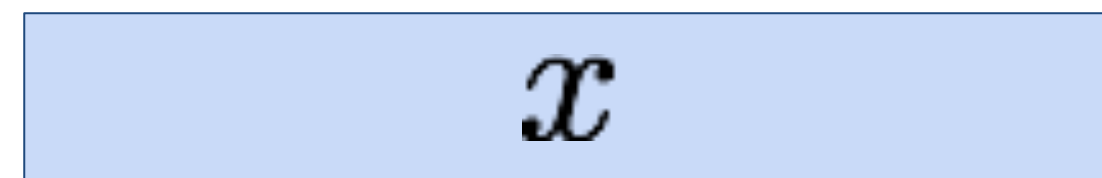$p_{\theta*}(z)$

# **Variational** Autoencoders

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

How do we **represent** this model?

Sample from
true **conditional**

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

$x$

$z$

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 $p(z)$ to be simple, e.g., Gaussian
Reasonable for latent attributes, e.g., pose, amount of smile

Conditional $p(x|z)$ 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**

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

$$\boxed{x}$$

**Decoder
Network**

$$\boxed{z}$$
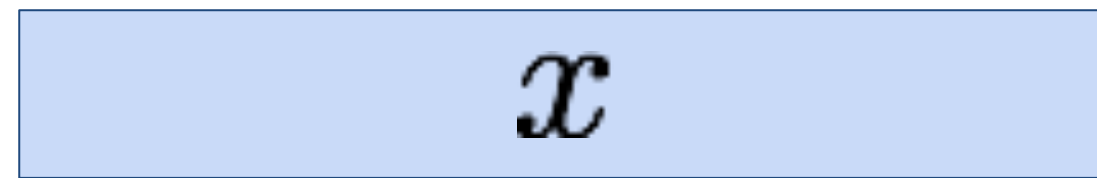
Sample from
true **prior**

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

# **Variational** Autoencoders

[ Kingma and Welling, 2014 ]

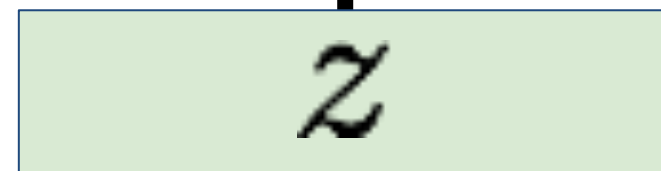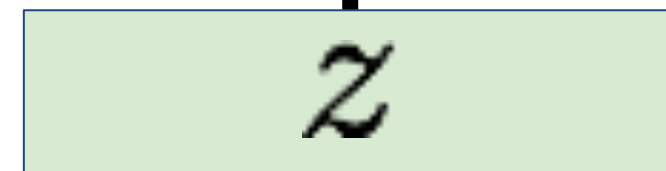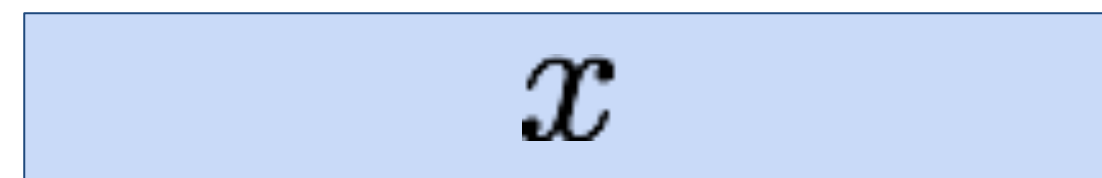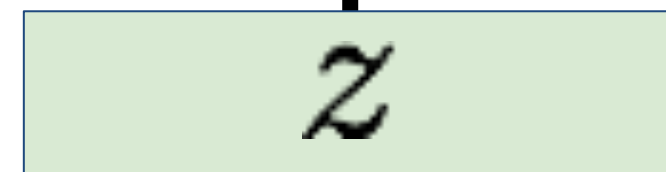We want to **estimate the true parameters** $\theta^*$ of this generative model

Sample from
true **conditional**

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

$x$

**Decoder
Network**

Sample from
true **prior**

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

$z$

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)

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

# **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)

### What is the problem with this?
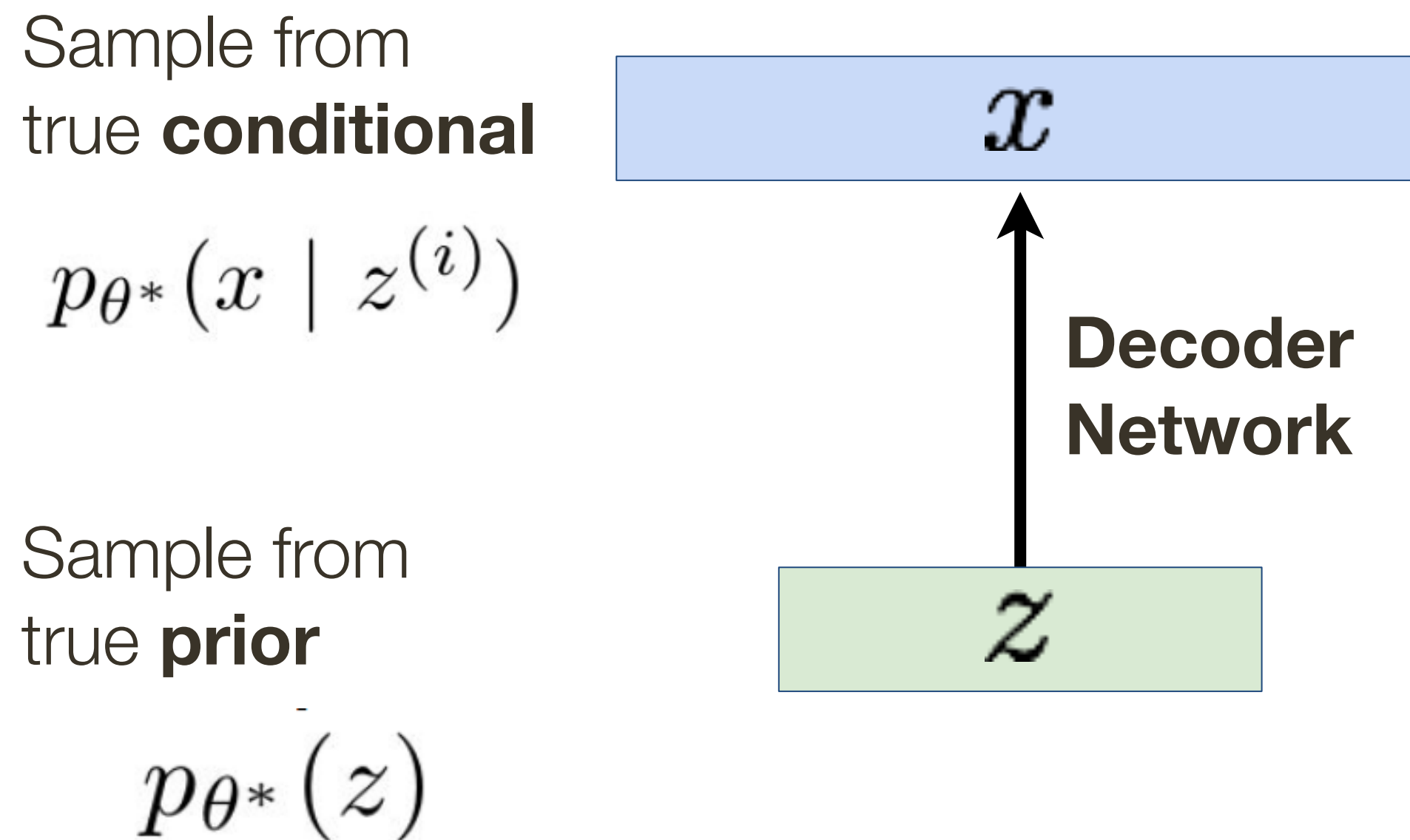
# **Variational** Autoencoders

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

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

**Decoder Network**

$$x$$

$$z$$

Sample from true **prior**

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

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

**Intractable !**

# **Intractability** in Variational Autoencoder

[ Kingma and Welling, 2014 ]

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

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

# **Intractability** in Variational Autoencoder

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

🙂

Simple **Gaussian** Prior

# **Intractability** in Variational Autoencoder

Data **likelihood**:

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

**Decoder** Neural Network 🙂

🙂

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: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

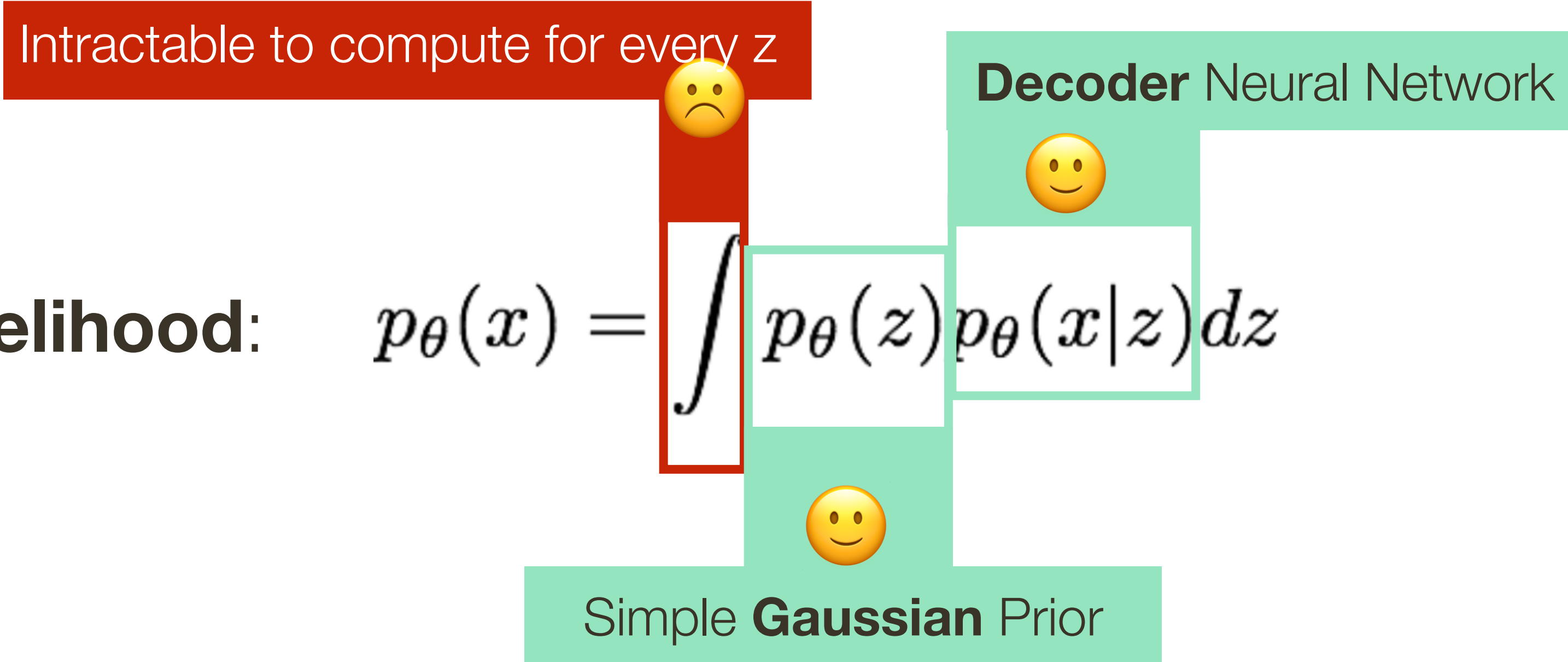**Solution:** In addition to decoder network modeling $p_\theta(x|z)$, define additional

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 $x \mid z$

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

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

**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)

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



**Encoder** Network
$q_\phi(z|x)$
(parameters $\phi$)

**Decoder** Network
$p_\theta(x|z)$
(parameters $\theta$)

# **Variational** Autoencoder

[ Kingma and Welling, 2014 ]

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



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

**Encoder** Network

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

(parameters φ)

$$x$$

**Decoder** Network

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

(parameters θ)

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

$$z$$

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

# **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:

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

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

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

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

# **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

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_θ(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_θ(z|x) **intractable** (saw earlier), can't compute this KL term :(

But we know KL divergence always >= 0.

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

# **Variational** Autoencoder

Derivation of lower bound of the data likelihood

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

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

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

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

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

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

**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)
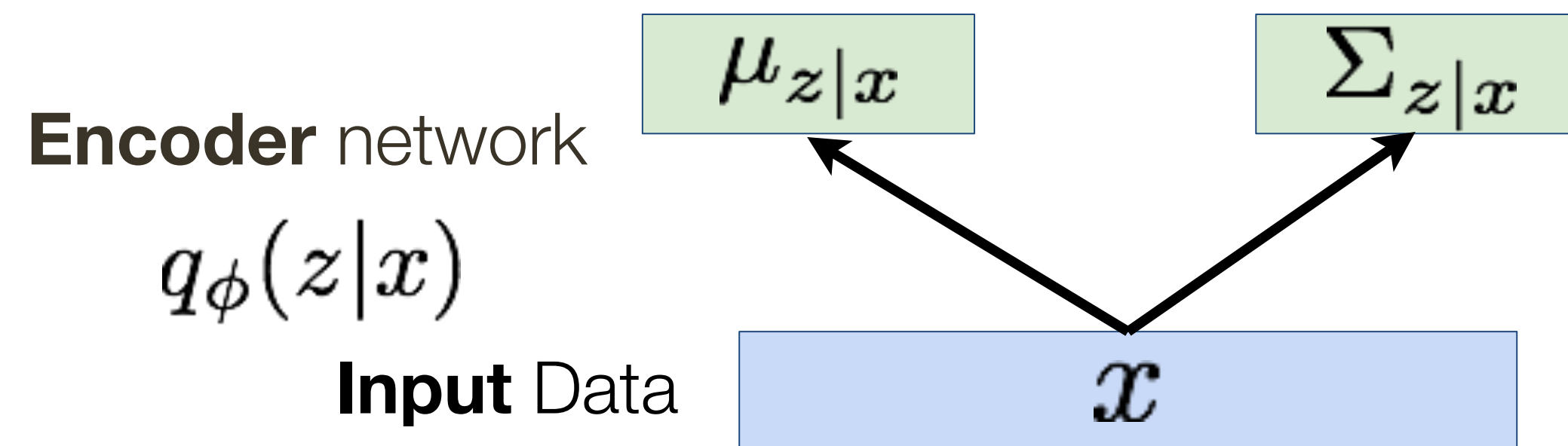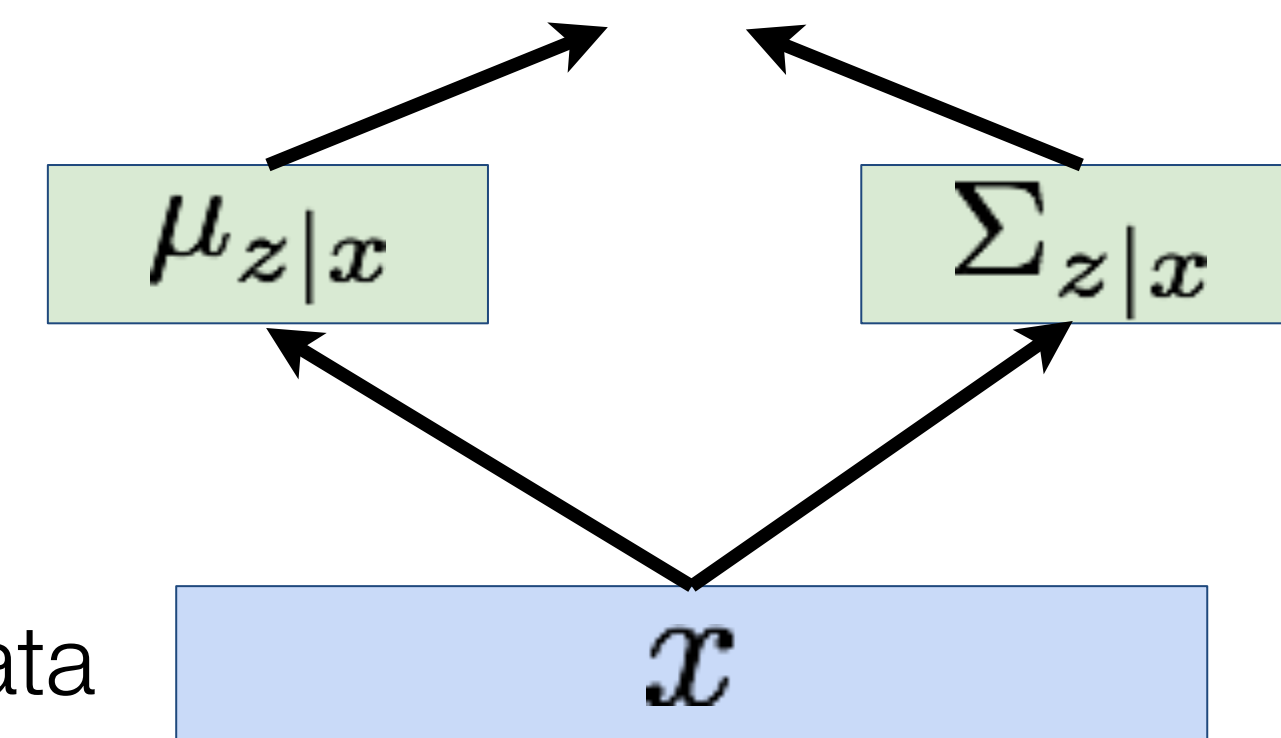for a given mini batch of input data

**Input** Data    $x$

# **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

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

Make approximate
posterior distribution
close to prior
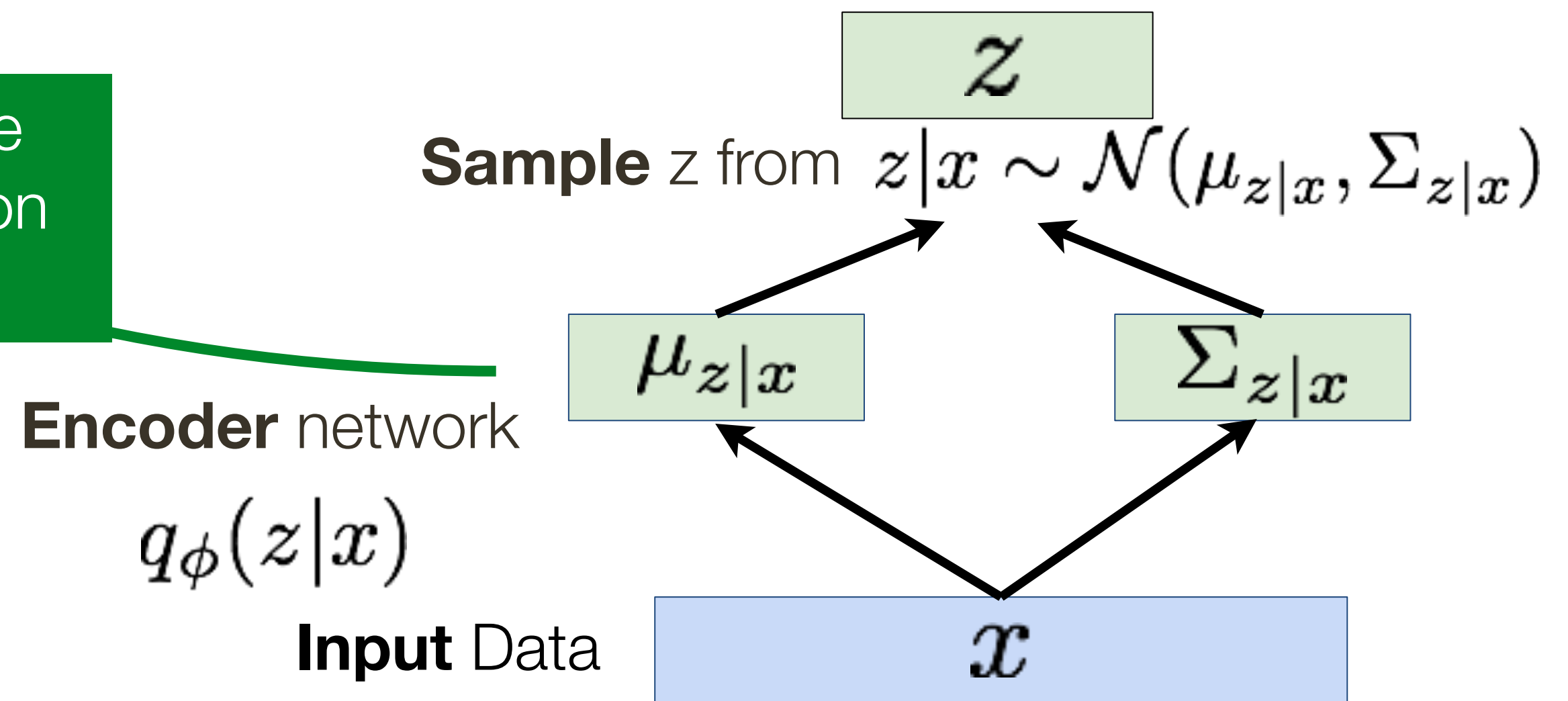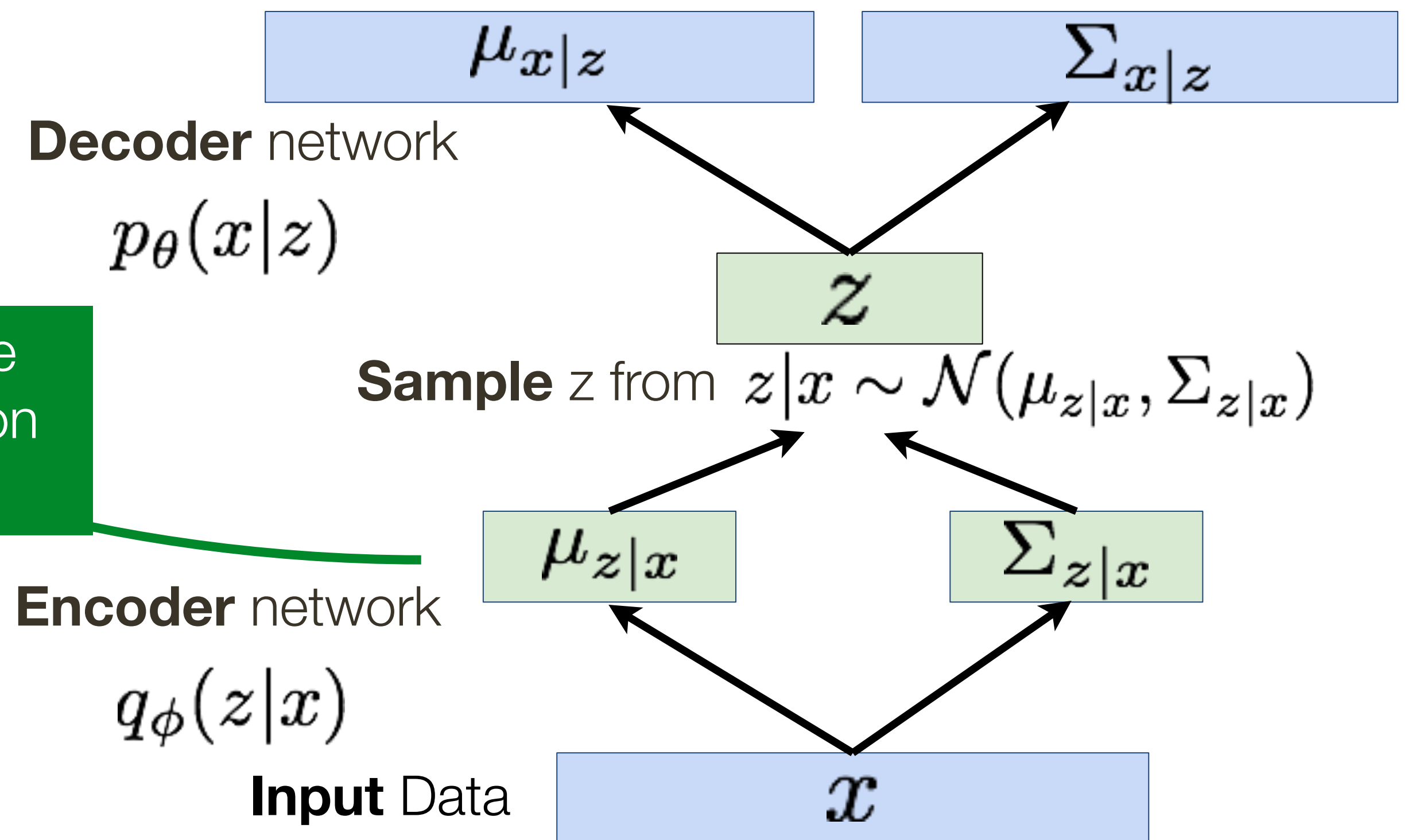
**Encoder** network

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

**Input** Data

$$\mu_{z|x}$$

$$\Sigma_{z|x}$$

$$x$$

# **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

$$z$$

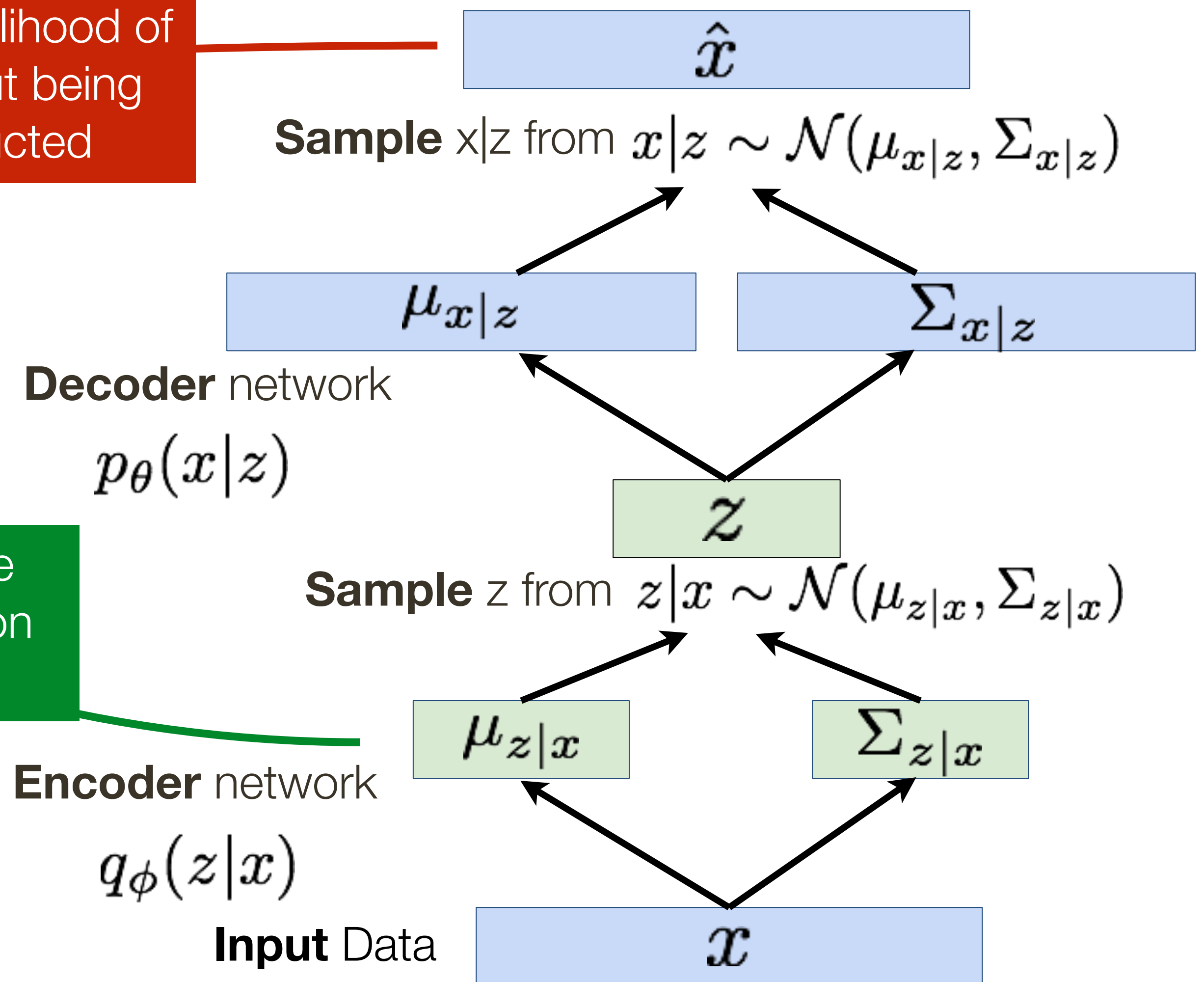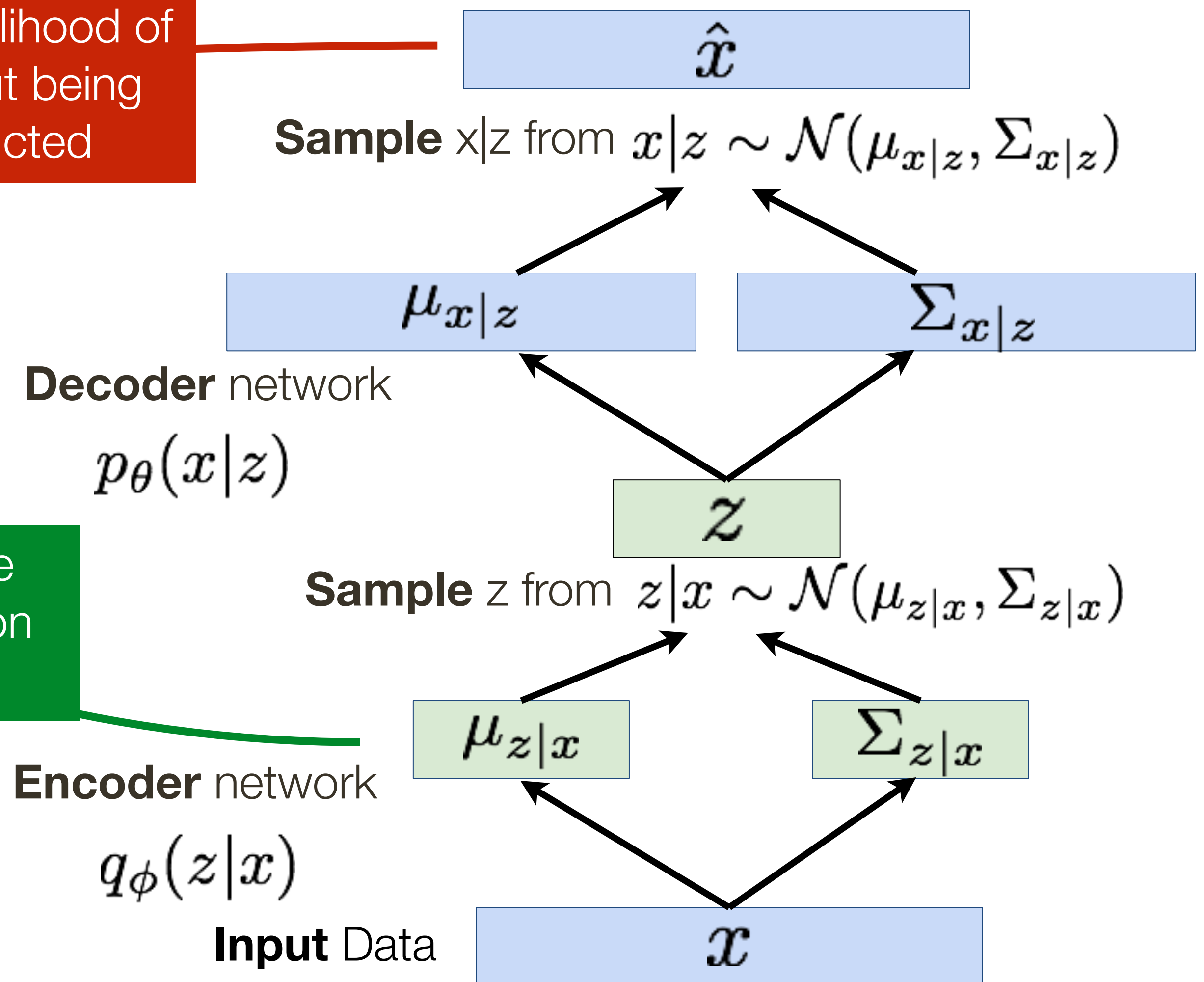**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

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

**Encoder** network

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

**Input** Data $\qquad x$

# **Variational** Autoencoder: Learning

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

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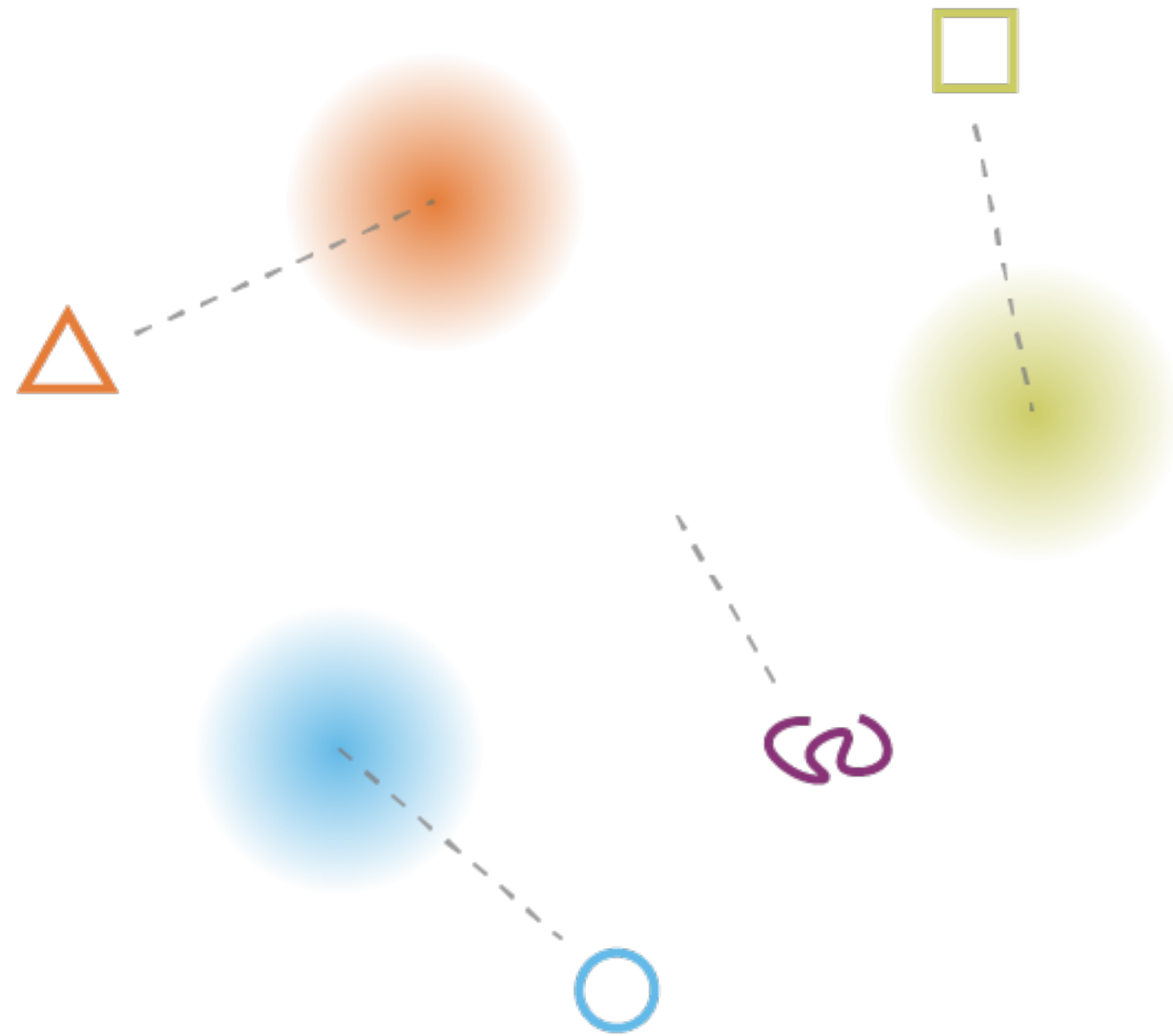
Make approximate posterior distribution close to prior

$\mu_{x|z}$

$\Sigma_{x|z}$

**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
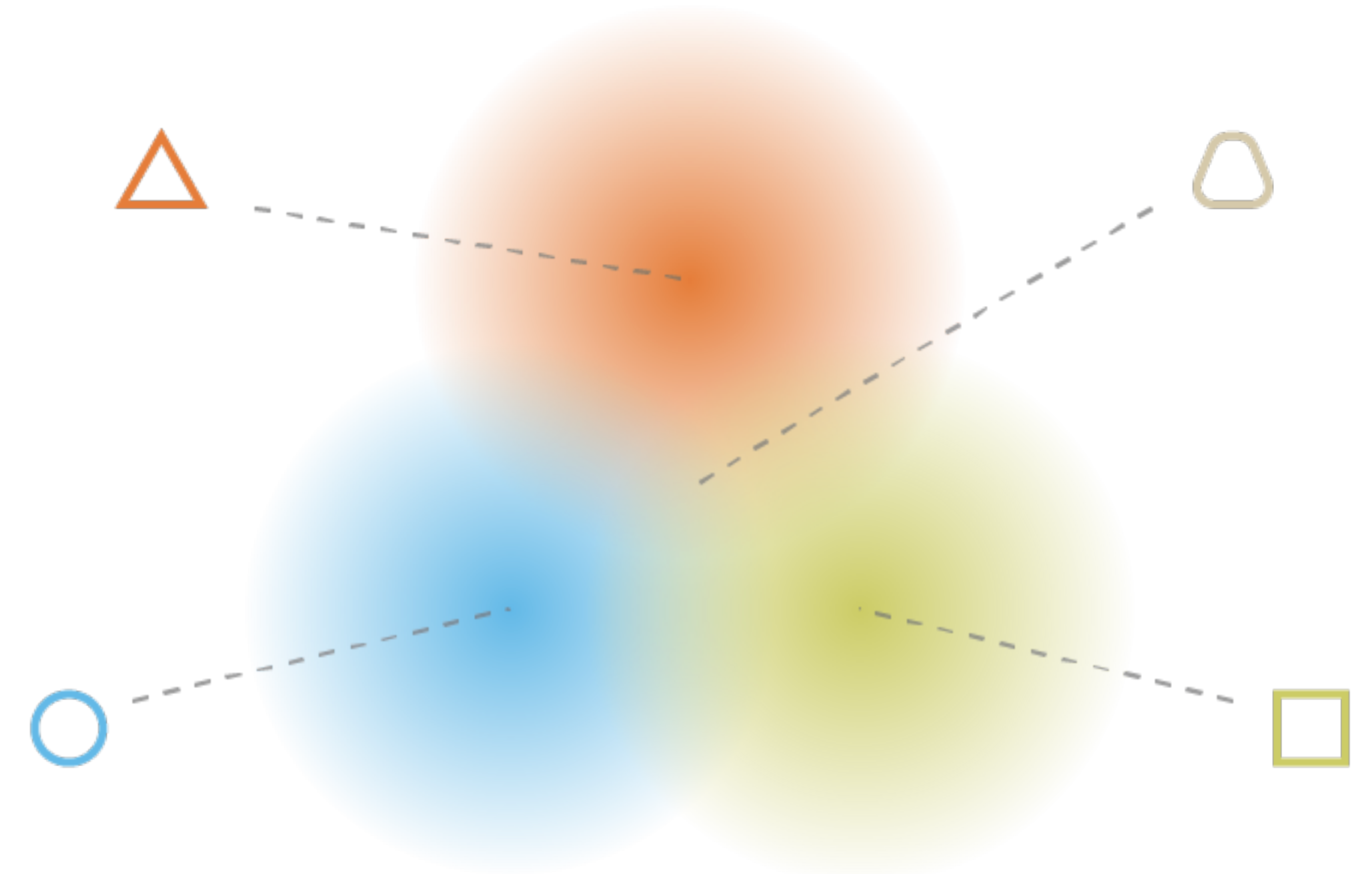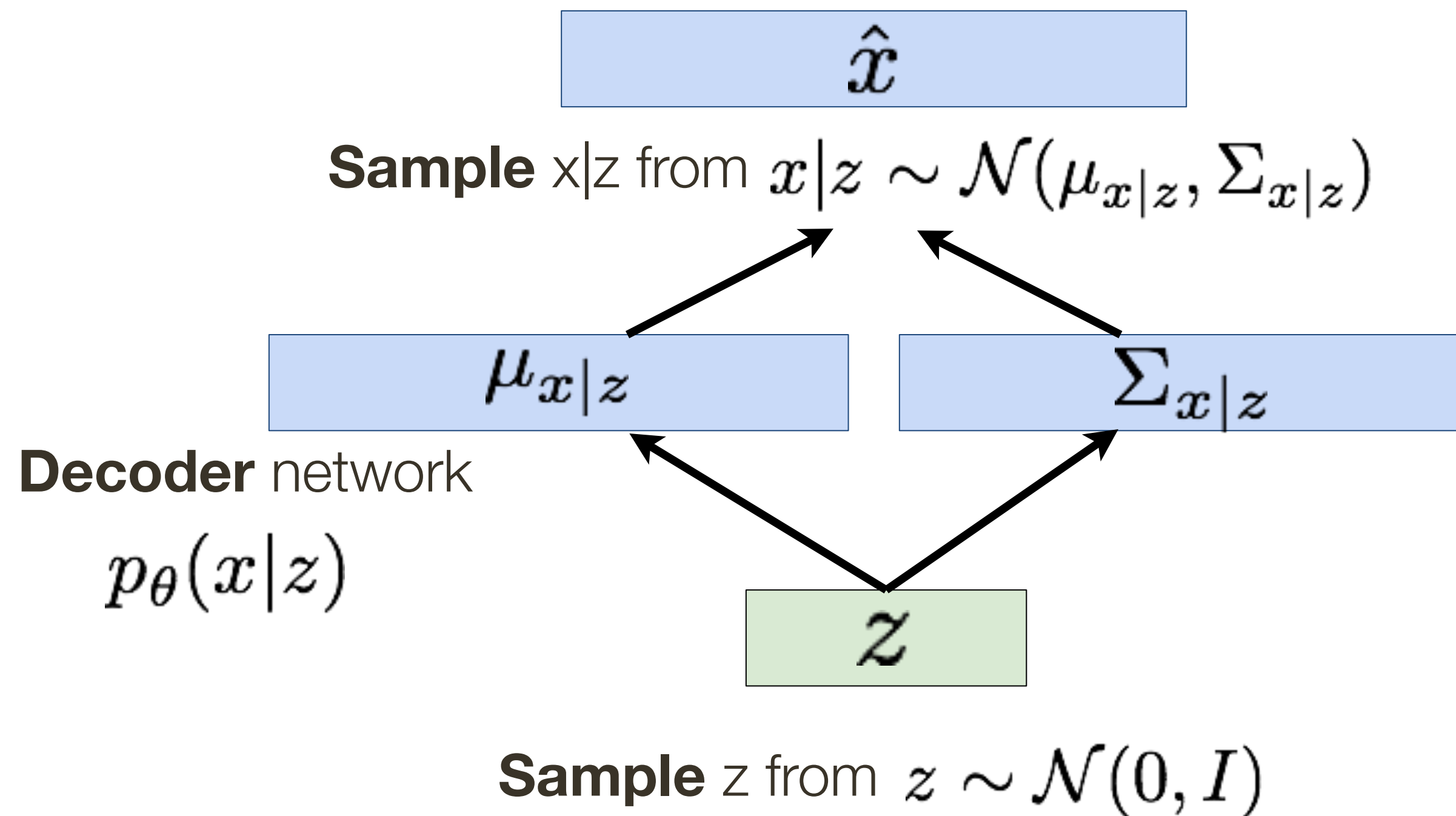posterior distribution
close to prior

$$\hat{x}$$

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

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

**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

$$x$$

# **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: Learning

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



**what can happen without regularisation** ❌    ✔ **what we want to obtain with regularisation**

https://towardsdatascience.com/understanding-variational-autoencoders-vaes-f70510919f73

# **Variational** Autoencoder: Generating Data

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



$\hat{x}$

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

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

**Decoder** network

$p_\theta(x|z)$

$z$

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

# **Variational** Autoencoder: Generating Data

# **Variational** Autoencoder: Generating Data

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

**Data manifold** for 2-d z

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

$\hat{x}$

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

**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

# **Conditional** VAE: Diverse Image Colorization



[ Deshpande et al., 2017 ]

# **Conditional** VAE: Temporal Predictions



(a) Frame 1  (b) Frame 2 (ground truth)  (c) Frame 2 (Sample 1)  (d) Frame 2 (Sample 2)

# **Variational** Autoencoder (VAE)

# **Variational** Autoencoder (VAE) **+ LSTM**

# **VAE + LSTM** with Structured Latent Space  [ He et al., 2018 ]

# Results: Chair CAD dataset

(a) Partial control.

(b) Full control.

## Ablation

| | Bound | Static | −C | | +C | |
|---|---|---|---|---|---|---|
| | | | −S | +S | −S | +S |
| Intra-E ↓ | 1.98 | 40.33 | 17.64 | 7.79 | 14.81 | **5.50** |
| Inter-E ↑ | 1.39 | 0.42 | 0.73 | 1.35 | 1.02 | **1.37** |
| I-Score ↑ | 4.01 | 1.28 | 1.83 | 3.63 | 2.56 | **3.94** |

## Quantitative

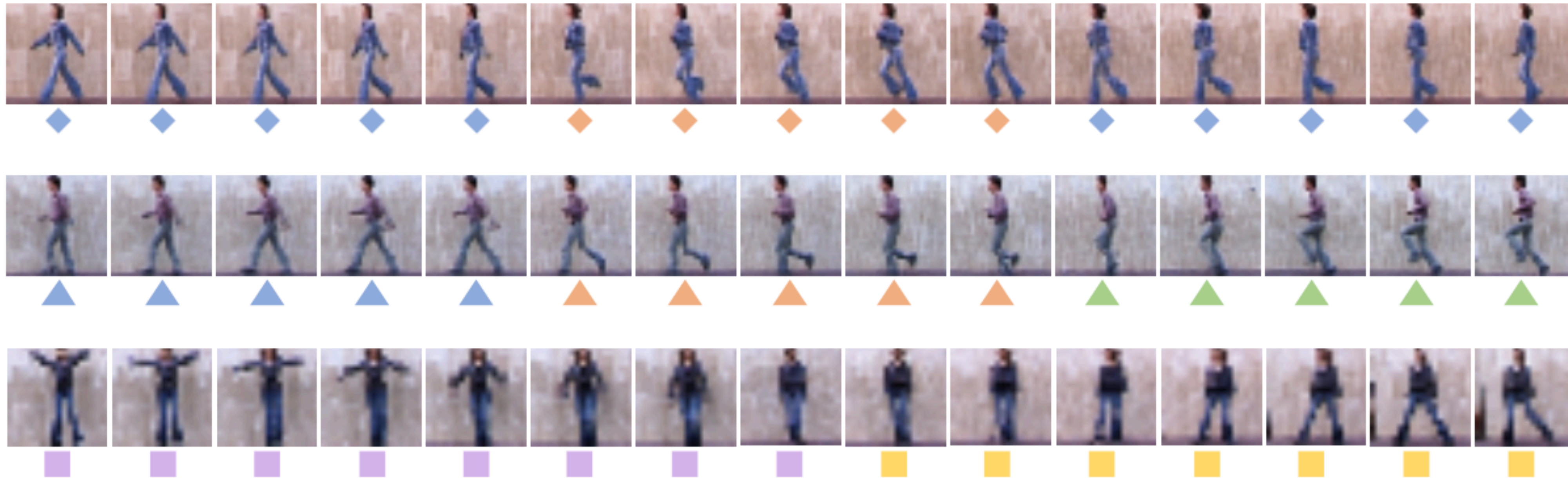| | | Chair CAD [1, 40] | |
|---|---|---|---|
| | Bound | Deep Rot. [40] | VideoVAE (ours) |
| | | ◉ | ◉ |
| Intra-E ↓ | 1.98 | 14.68 | **5.50** |
| Inter-E ↑ | 1.39 | 1.34 | **1.37** |
| I-Score ↑ | 4.01 | 3.39 | **3.94** |

# **Results:** Weizmann Human Action dataset

✓ Identity = ◆ | ▲ | ■    ✓ Action = ● walking | ● running | ● skipping | ● jumping jack | ● side step
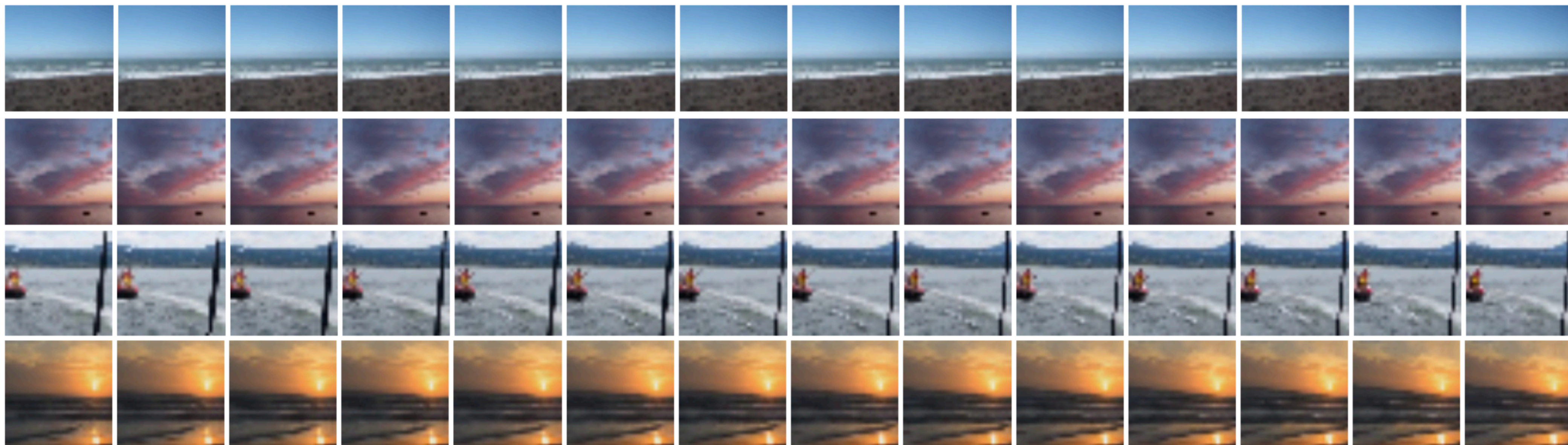
← generate →

| Weizmann Human Action [2] | | | |
|---|---|---|---|
| | Bound | MoCoGAN [32] | VideoVAE (ours) |
| | | ○ | ○ ◉ |
| Intra-E ↓ | 0.63 | 23.58 | 9.53 **9.44** |
| Inter-E ↑ | 4.49 | 2.91 | **4.37** **4.37** |
| I-Score ↑ | 89.12 | 13.87 | 69.55 **70.10** |

# **Results:** MIT Flickr

[ He et al., 2018 ]



| | YFCC [31] — MIT Flickr [34] | | |
|---|---|---|---|
| Bound | VGAN [34] | VideoVAE (ours) | |
| | ○ | ○ | ● |
| Intra-E ↓ 30.34 | 46.96 | 44.03 | **38.20** |
| Inter-E ↑ 0.693 | **0.692** | 0.691 | **0.692** |
| I-Score ↑ 1.87 | 1.58 | 1.62 | **1.81** |

# **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)

# VAE /w (powerful) PixelCNN Decoder

**Problem**: If the decoder is too powerful, it may just ignore the latent variables (i.e. posterior collapse). This happens when the decoder can make the reconstruction loss incredibly small, such that the regularization term dominates the loss function. In such a case, the encoder will learn to reduce the regularization term, and produce meaningless latents to match $p(z) = N(0, 1)$.
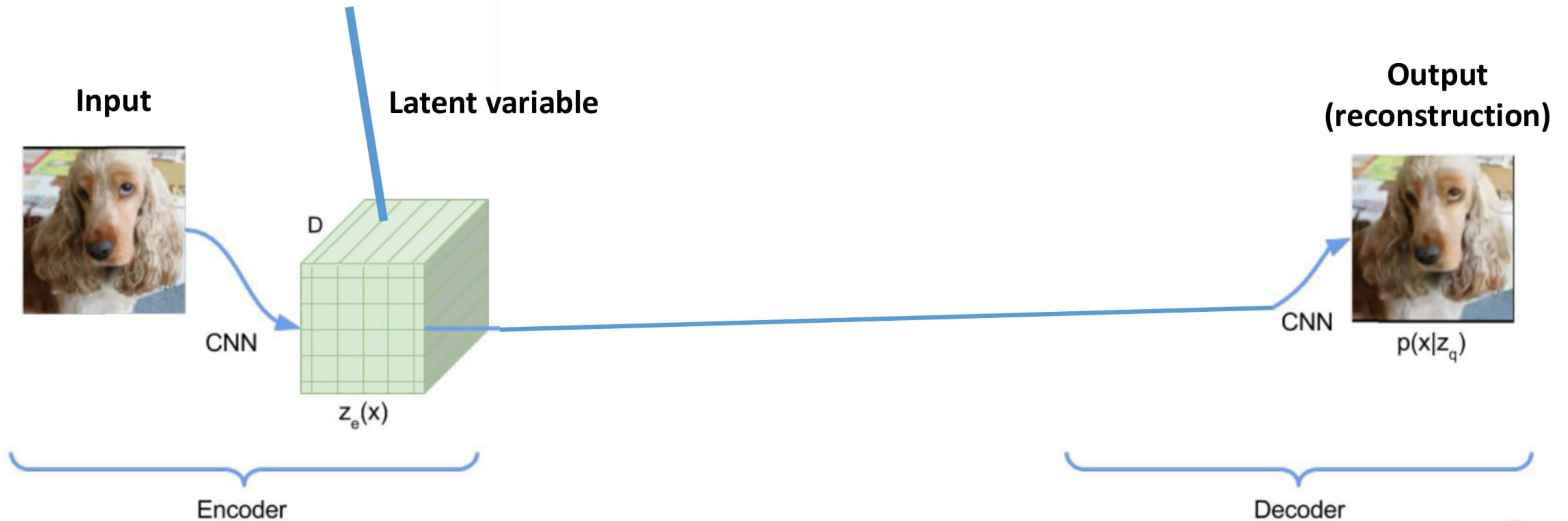
# Vector Quantized Variational Autoencoders (VQ-VAE)

# **Autoencoders** Reminder …

How to discretize?

For the example:
We take this to be a 4 x 4 image
with 2 channels.

**Input**

**Latent variable**

**Output (reconstruction)**

D

CNN

CNN

$z_e(x)$

$p(x|z_q)$

Encoder

Decoder

We can train this system end-to-end
using MSE (reconstruction loss)

# **Autoencoders** Reminder …

How to discretize?

4 x 4 image with 2 channels.
We plot all pixel values (16) in 2D
(since we have 2 channels)



Channel 2

D

CNN

$z_e(x)$

Encoder

Channel 1

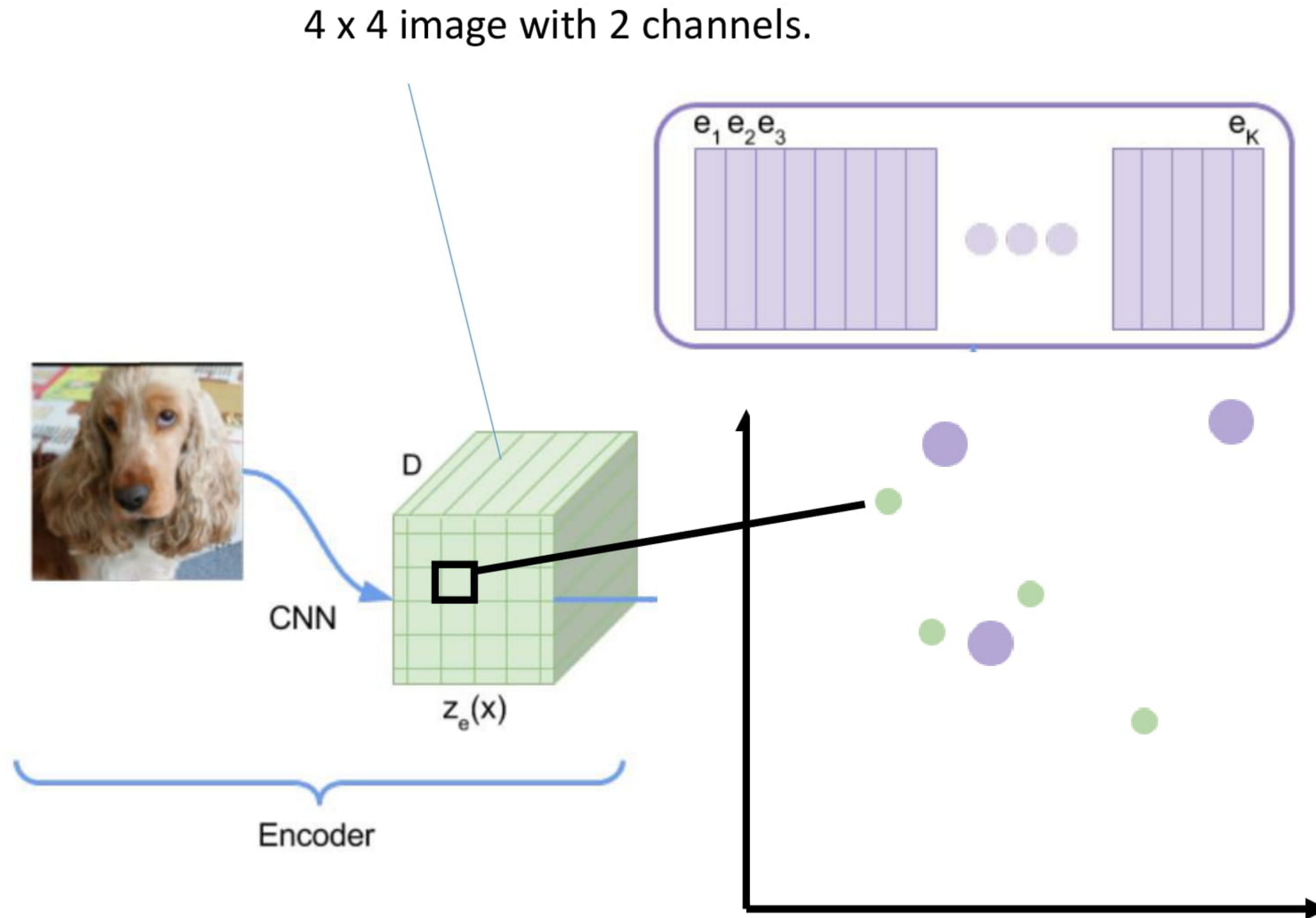# **Vector Quantized** - VAE

4 x 4 image with 2 channels.

$e_1, e_2 e_3$

$e_K$

D

CNN

$z_e(x)$

Encoder

Make dictionary of vectors
$e_1, \ldots, e_K$

Each $e_i$ has 2 dimensions.

# **Vector Quantized** - VAE

4 x 4 image with 2 channels.

Make dictionary of vectors
$e_1, \ldots, e_K$

Each $e_i$ has 2 dimensions.

$e_1 \, e_2 e_3$      $e_K$

$e_3$     $e_2$

D

CNN

$z_e(x)$

$e_1$

Encoder

For each latent pixel, look up
nearest dictionary element $e$

# **Vector Quantized** - VAE

4 x 4 image with 2 channels.

Each $e_i$ has 2 dimensions.

# **Vector Quantized** - VAE

Latent is 1 channel image and contains the id of each e for each pixel (**discrete**).

# VQ-VAE — **Training**

- How to backpropegate through the discretization?
  - Lets say a gradient is incoming to a dictionary vector
  - We do not update the dictionary vector (fixed)
  - Instead we apply the gradient of e to the non-discretized vector

# VQ-VAE — **Training**

$$L = \log p(x|z_q(x)) + \|\mathrm{sg}[z_e(x)] - e\|_2^2 + \beta\|z_e(x) - \mathrm{sg}[e]\|_2^2,$$

- How to backpropegate through the discretization?
  - Lets say a gradient is incoming to a dictionary vector
  - We do not update the dictionary vector (fixed)
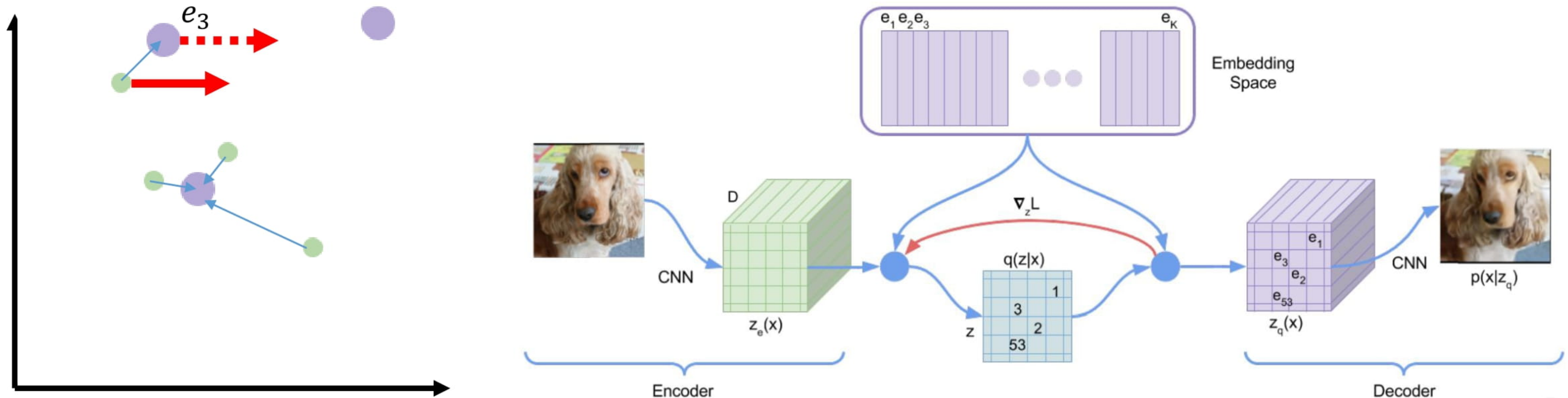  - Instead we apply the gradient of e to the non-discretized vector
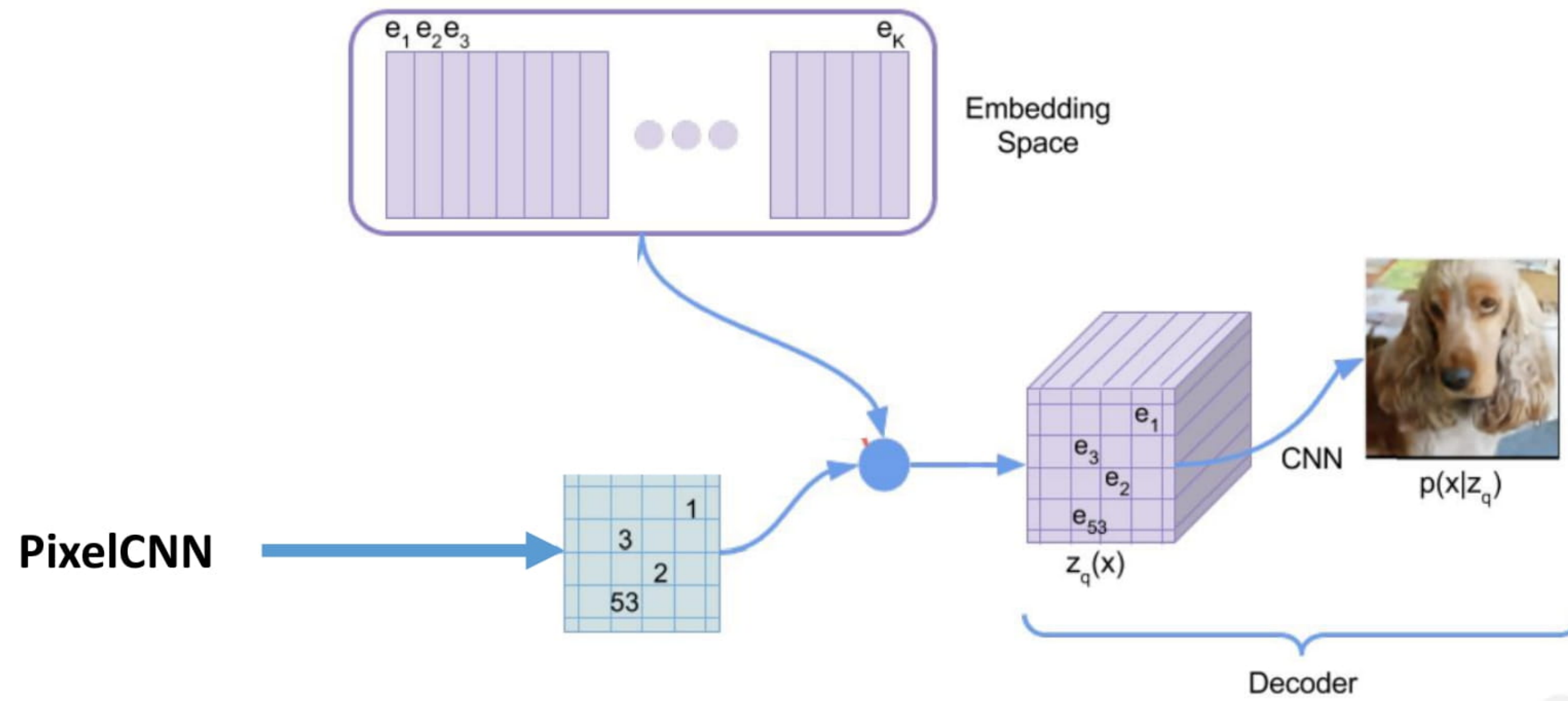
# VQ-VAE — **Training**

$$L = \log p(x|z_q(x)) + \|\mathrm{sg}[z_e(x)] - e\|_2^2 + \beta\|z_e(x) - \mathrm{sg}[e]\|_2^2,$$

**Reconstruction loss**, which optimizes both the encoder and decoder

**Regularization**, ensures that encoder does not grow arbitrarily

**VQ loss**, which moves the embedding vectors towards encoder outputs

# VQ-VAE — **Sampling / Generation**



Class: pickup

# VQ-VAE — **Sampling / Generation**

- Comparable with VAE on CIFAR-10 in terms of density estimation
- Reconstructions on ImageNet are very good



Figure 2: Left: ImageNet 128x128x3 images, right: reconstructions from a VQ-VAE with a 32x32x1 latent space, with K=512.

# VQ-VAE vs. GAN



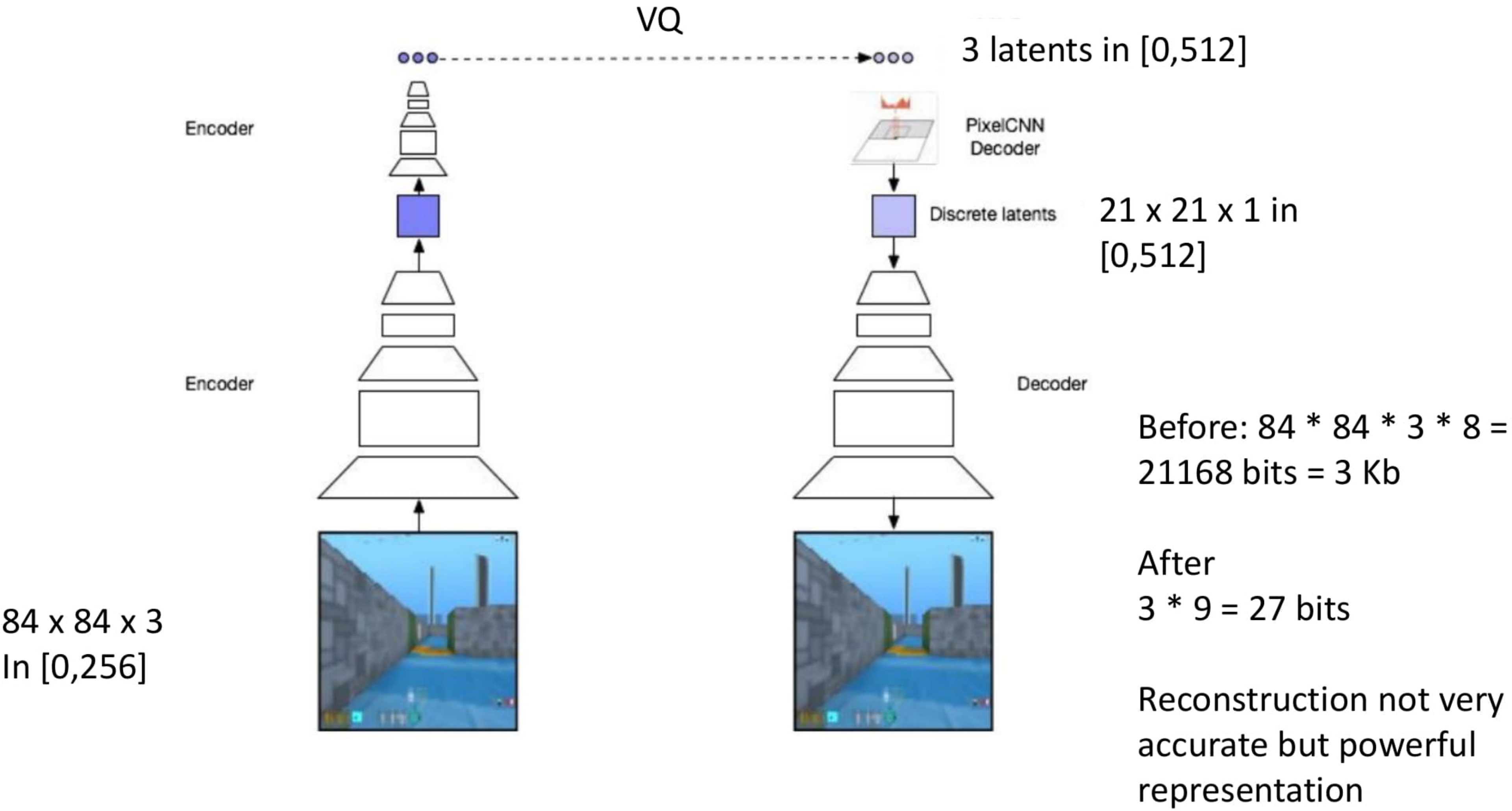VQ-VAE (Proposed)                    BigGAN deep

# Relationship of VQ-VAE to VAE

**VAE**: Assumes Gaussian prior over continuous latent space

**VQ-VAE**: Assumes uniform categorical distribution over discrete keywords (all keywords are equally likely)

# Comparison

|  | GAN | Variational Autoencoder | Pixel CNN | VQ-VAE (This talk) |
|---|---|---|---|---|
| **Compute exact likelihood p(x)** | ✘ | ✘ | ✔ | ✘ |
| **Has latent variable z** | ✔ | ✔ | ✘ | ✔ |
| **Compute latent variable z (inference)** | ✘ | ✔ | ✘ | ✔ |
| **Discrete latent variable** | ✘ | ✘ | ✘ | ✔ |
| **Stable training?** | ✘ | ✔ | ✔ | ✔ |
| **Sharp images?** | ✔ | ✘ | ✔ | ✔? |

# **Multi-stage** VQ-VAE



VQ

3 latents in [0,512]

Encoder

PixelCNN
Decoder

Discrete latents

21 x 21 x 1 in
[0,512]

Encoder

Decoder

84 x 84 x 3
In [0,256]

Before: 84 * 84 * 3 * 8 =
21168 bits = 3 Kb

After
3 * 9 = 27 bits

Reconstruction not very
accurate but powerful
representation

# 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

What if we give up on explicitly modeling density, and just want to sample?

# **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

What if we give up on explicitly modeling density, and just want to sample?

GANs: don't work with any explicit density function