



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

Lecture 16: Generative Models [part 2]

Logistics

Project Proposals due Monday 11:59pm

- They are graded for completeness / development of the proposal
- They are **not** graded for quality of the idea

I will provide feedback on the proposals, but **don't wait for it** to work on the projects

I will hold additional office hours for feedback on **Friday** ... tentatively
4:30-5:30pm

PixelRNN and PixelCNN

Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

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

The diagram illustrates the decomposition of the likelihood of an image \mathbf{x} into a product of 1-d distributions. A green box labeled "Likelihood of image \mathbf{x} " points to $p(\mathbf{x})$. A blue box labeled "Probability of i 'th pixel value given all previous pixels" points to $p(x_i | x_1, \dots, x_{i-1})$.

then maximize likelihood of training data

Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

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

The diagram illustrates the decomposition of the likelihood of an image \mathbf{x} into a product of conditional probabilities for each pixel. The equation $p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$ is shown. A green box labeled $p(\mathbf{x})$ is connected by an upward arrow to a green box labeled "Likelihood of image \mathbf{x} ". A blue box labeled $p(x_i | x_1, \dots, x_{i-1})$ is connected by an upward arrow to a blue box labeled "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**

Optional subtitle

Explicit Density model

Use chain rule to decompose likelihood of an image \mathbf{x} into product of (many) 1-d distributions

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

Likelihood of image \mathbf{x}

Probability of i 'th pixel value given all previous pixels

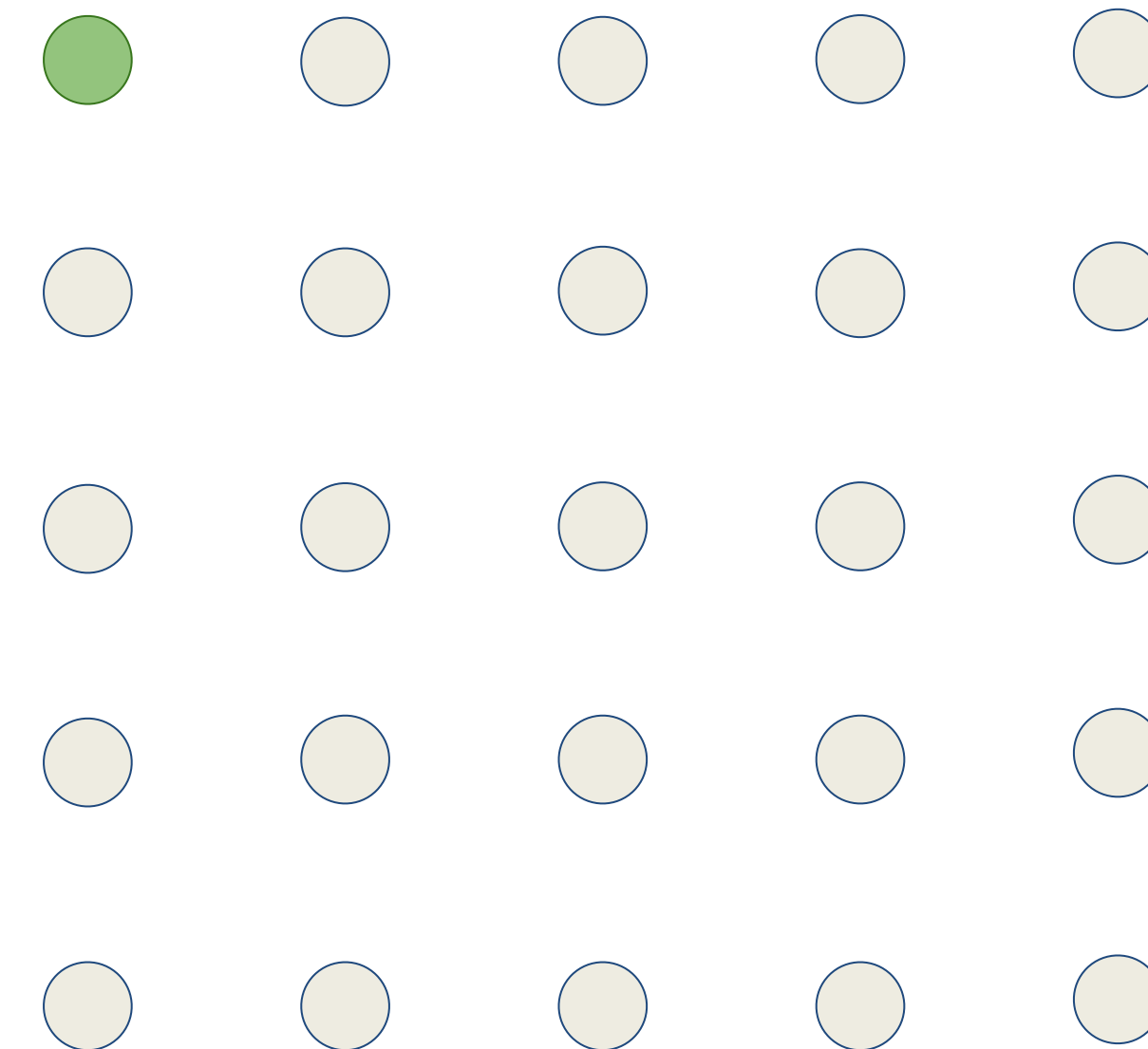
then maximize likelihood of training data

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

Also requires defining **ordering** of “previous pixels”

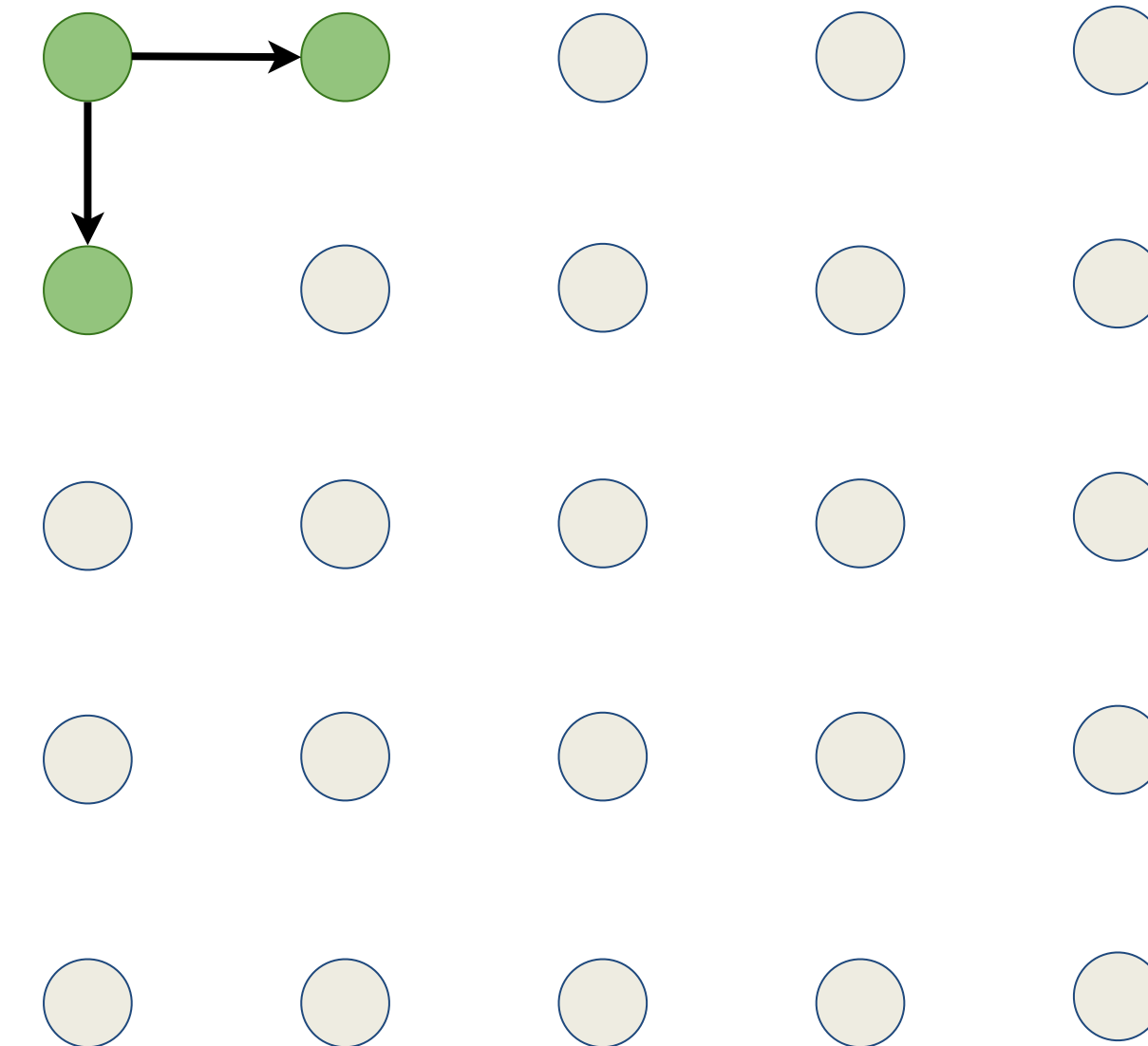
Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)



Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

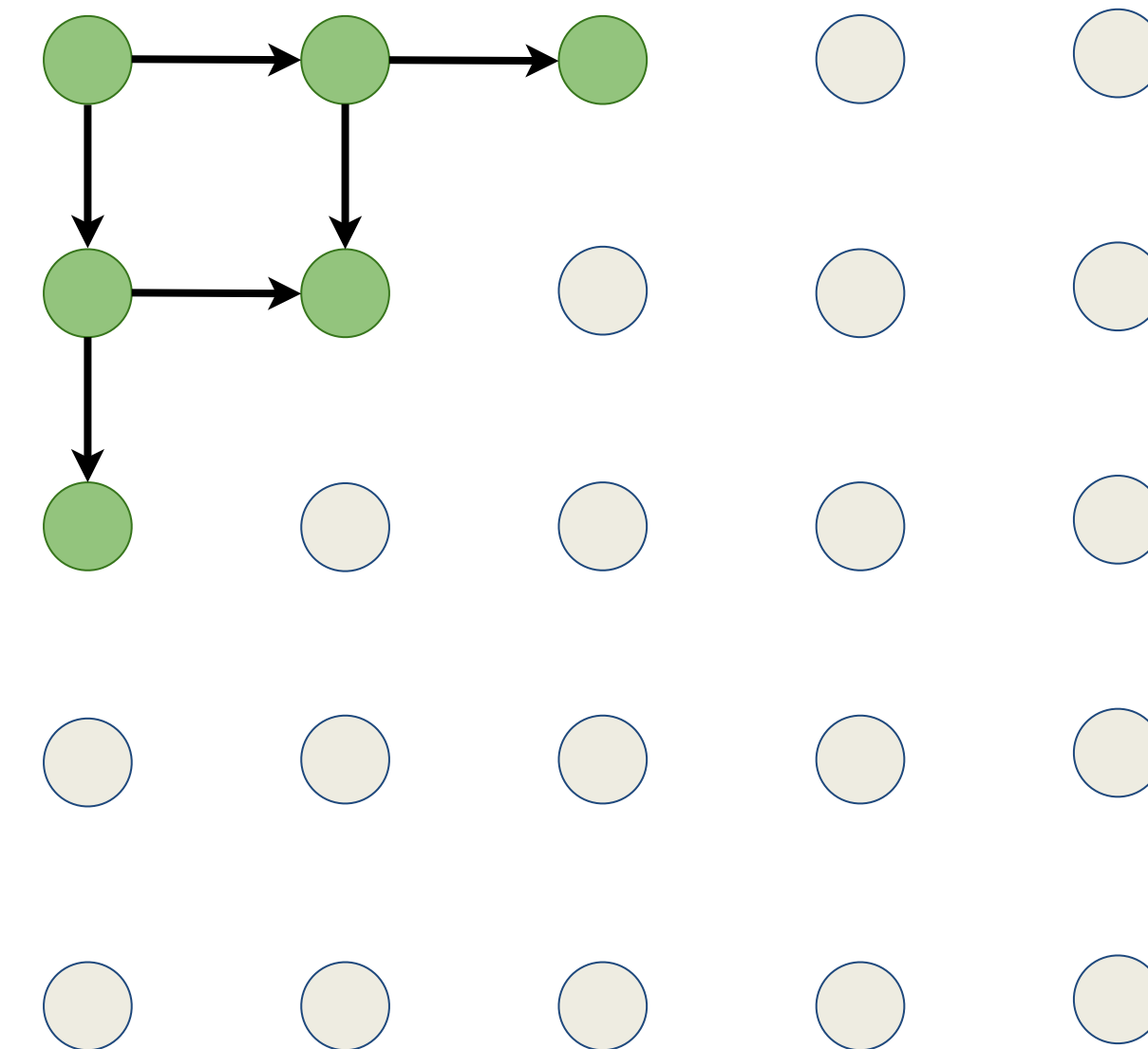


PixelRNN

[van der Oord et al., 2016]

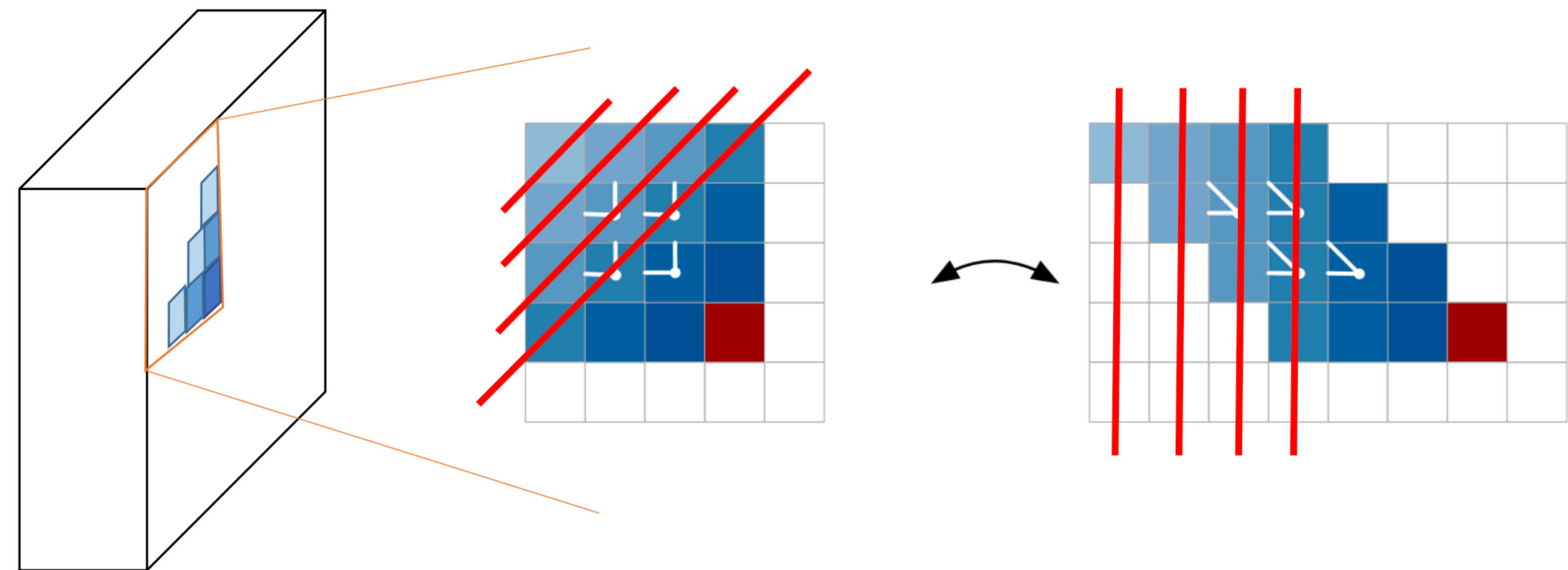
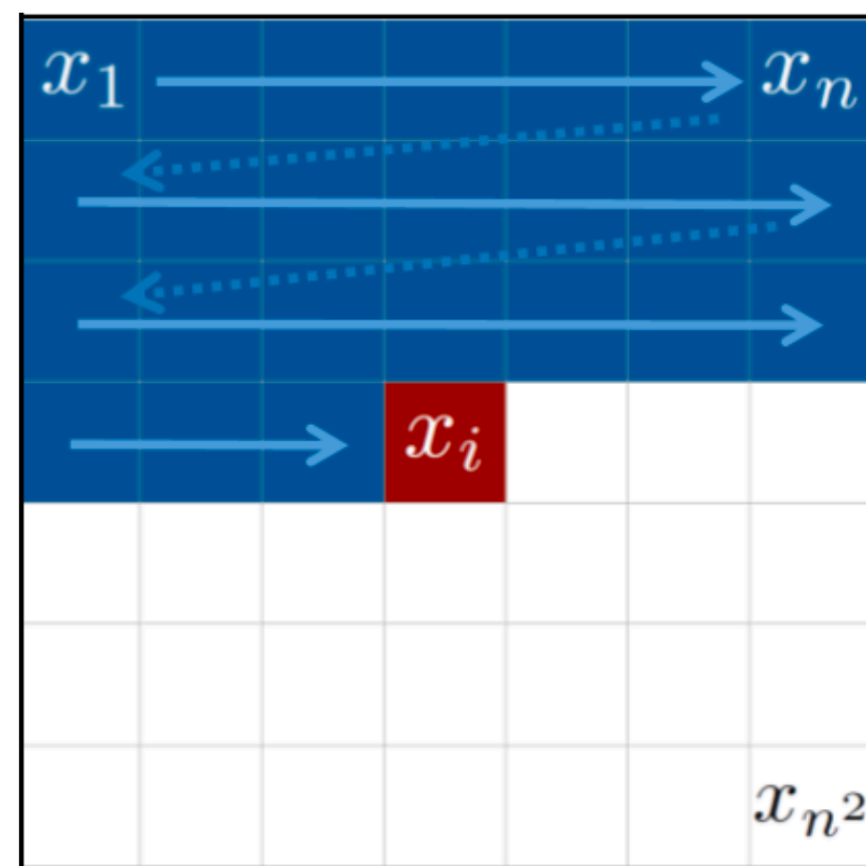
Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)



PixelRNN

[van der Oord et al., 2016]

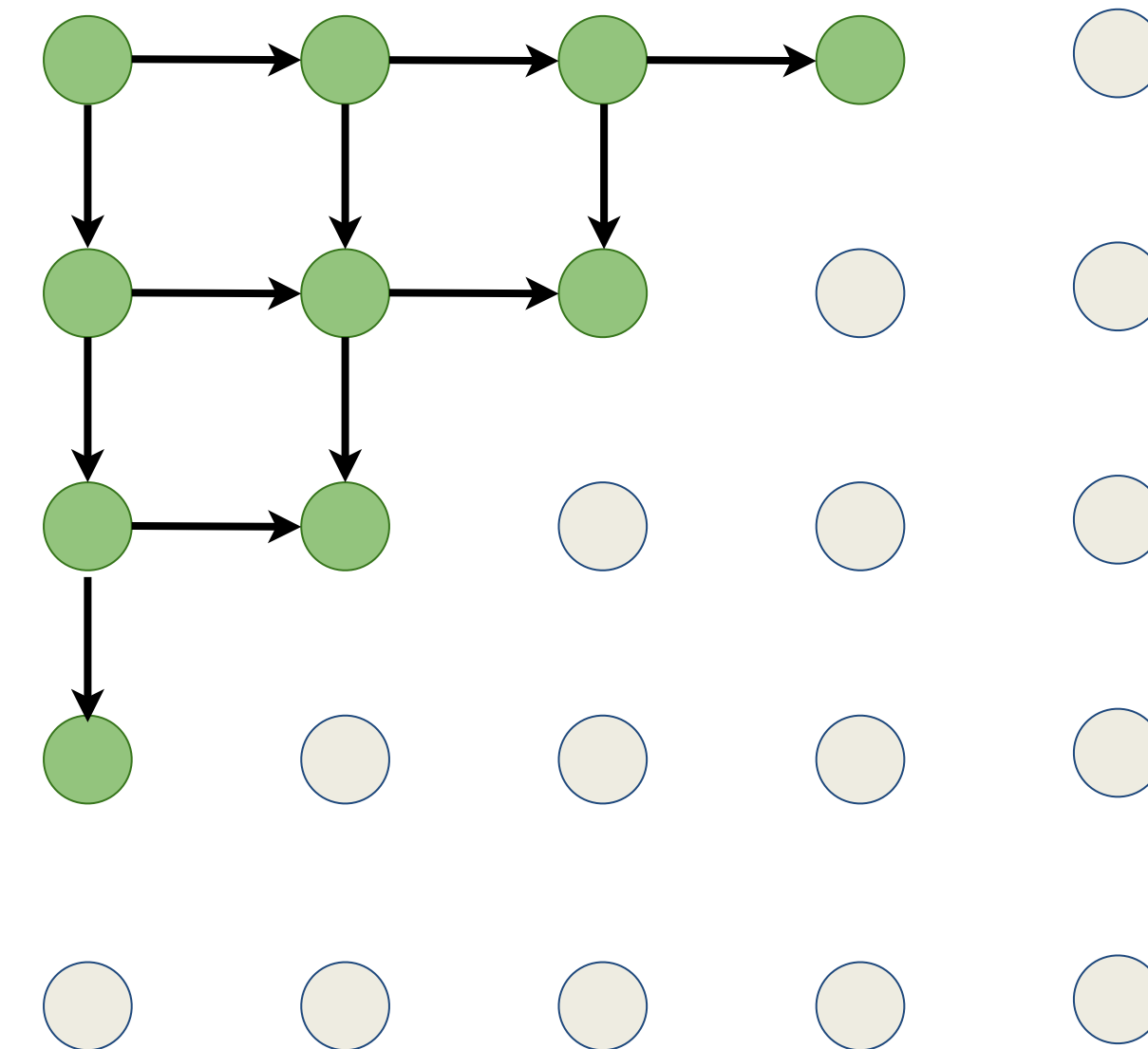


PixelRNN

[van der Oord et al., 2016]

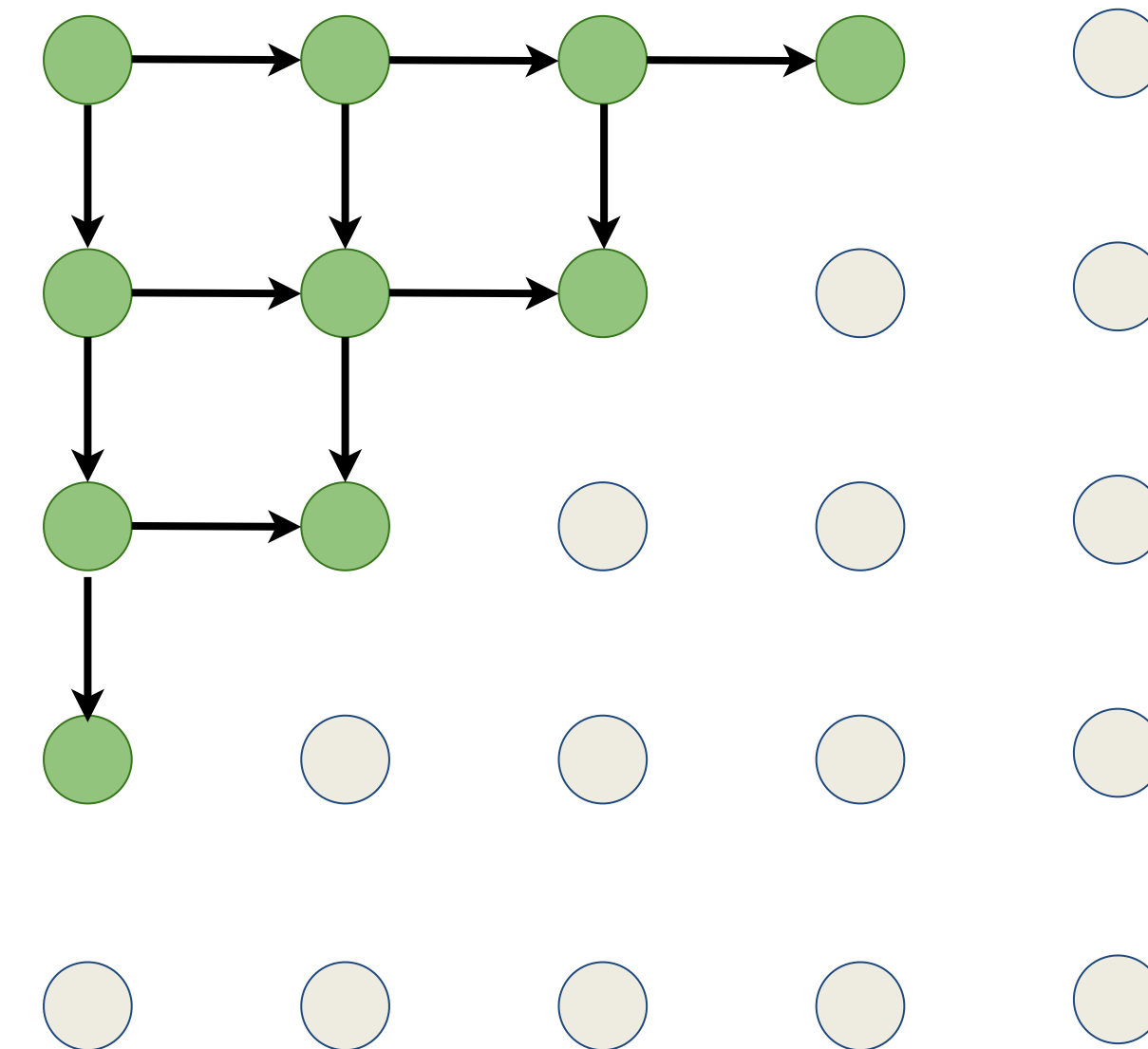
Generate image pixels starting from the corner

Dependency on previous pixels
model using an RNN (LSTM)



Generate image pixels starting from the corner

Dependency on previous pixels
model using an RNN (LSTM)



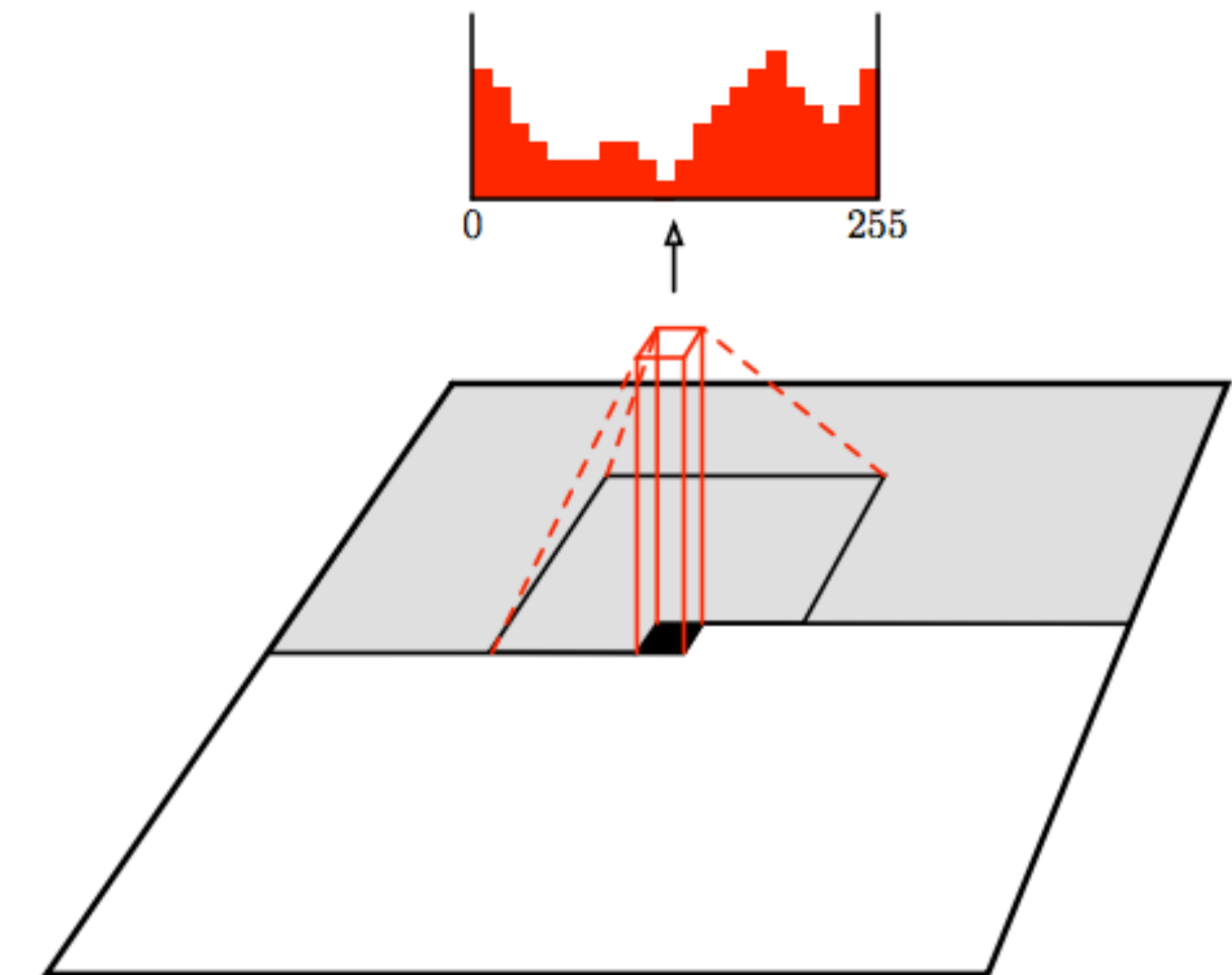
Problem: sequential generation is slow

PixelCNN

[van der Oord et al., 2016]

Still generate image pixels starting from the corner

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



PixelCNN

[van der Oord et al., 2016]

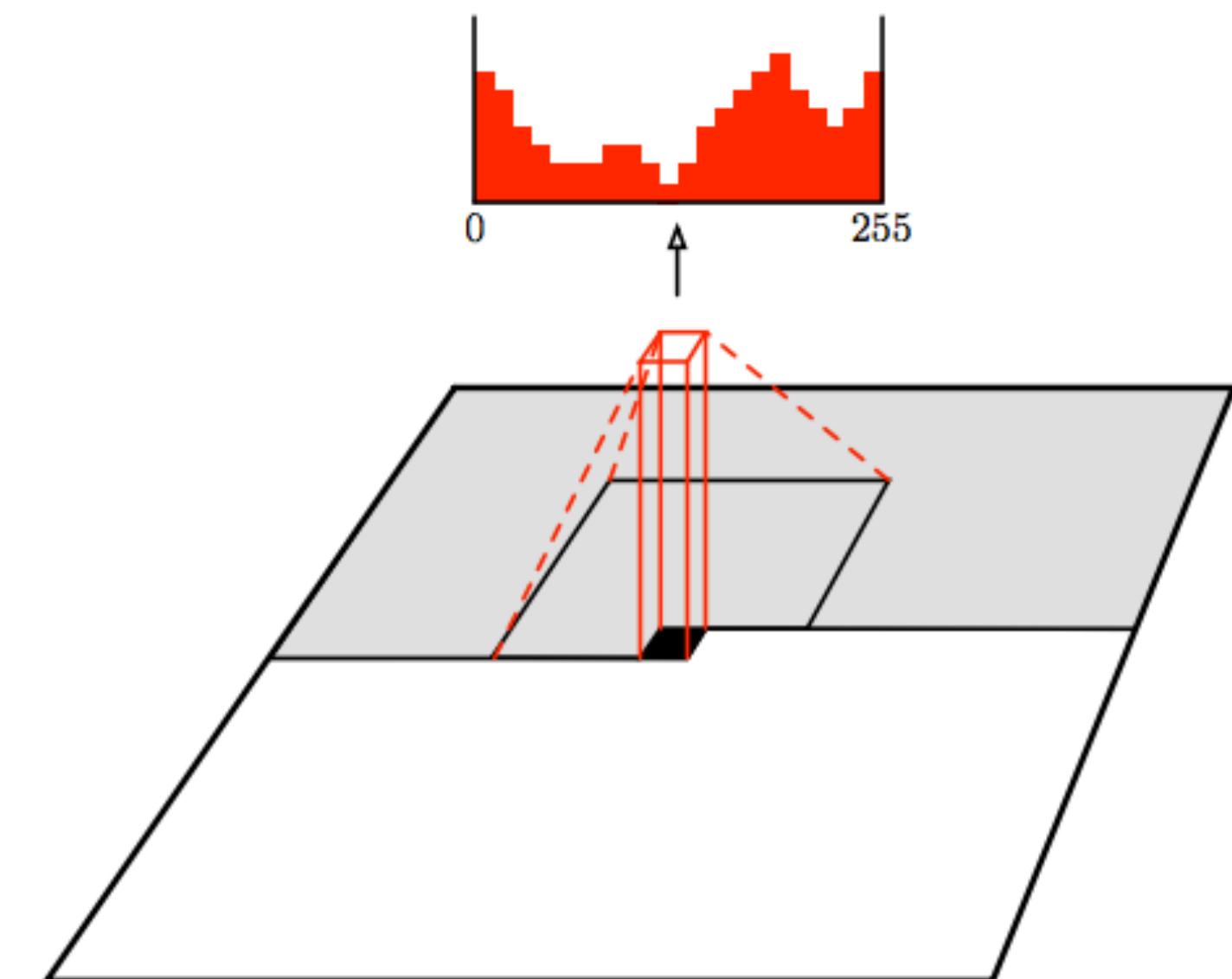
Still generate image pixels starting from the corner

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

Training: maximize likelihood of training images

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

Softmax loss at each pixel



PixelCNN

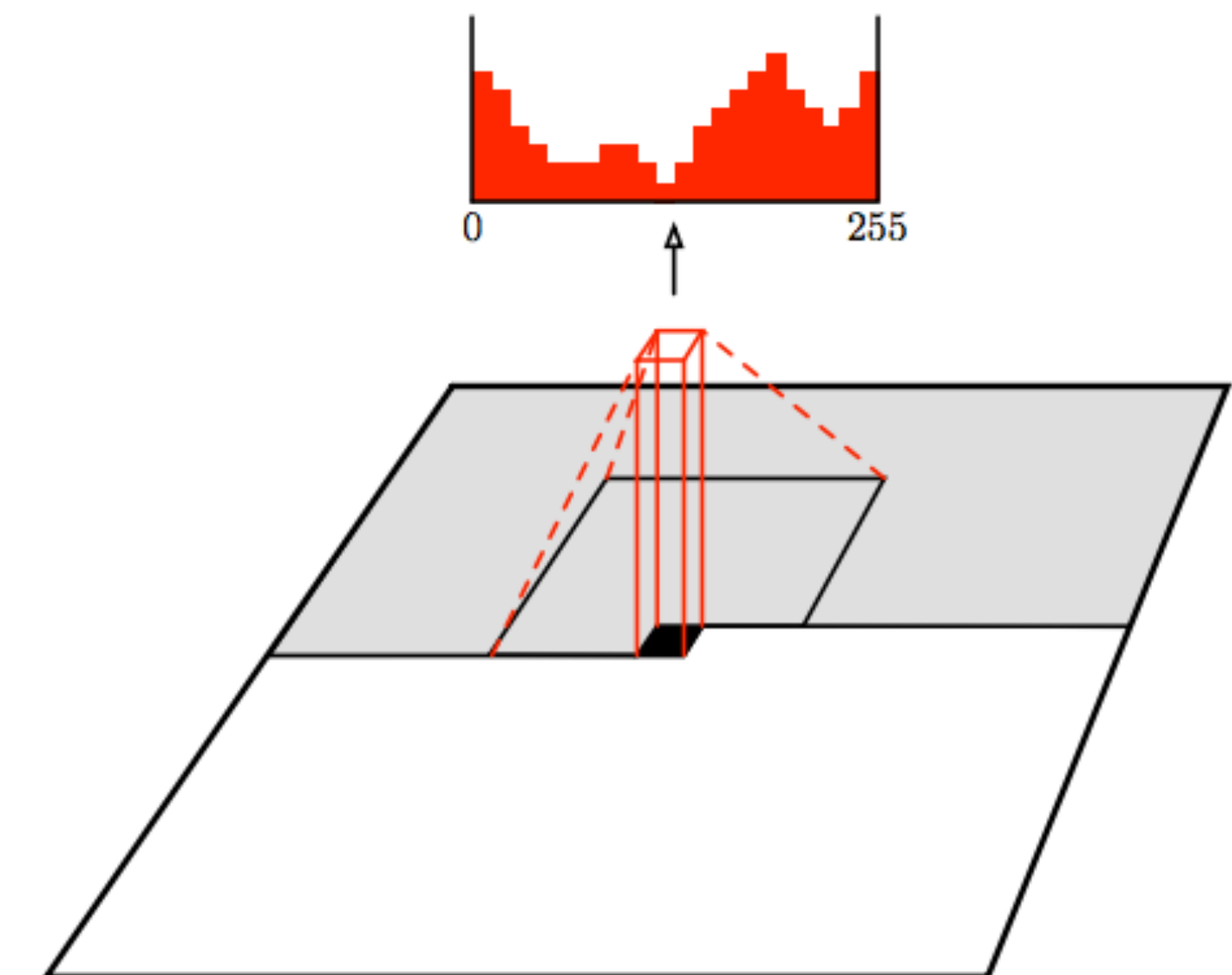
[van der Oord et al., 2016]

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

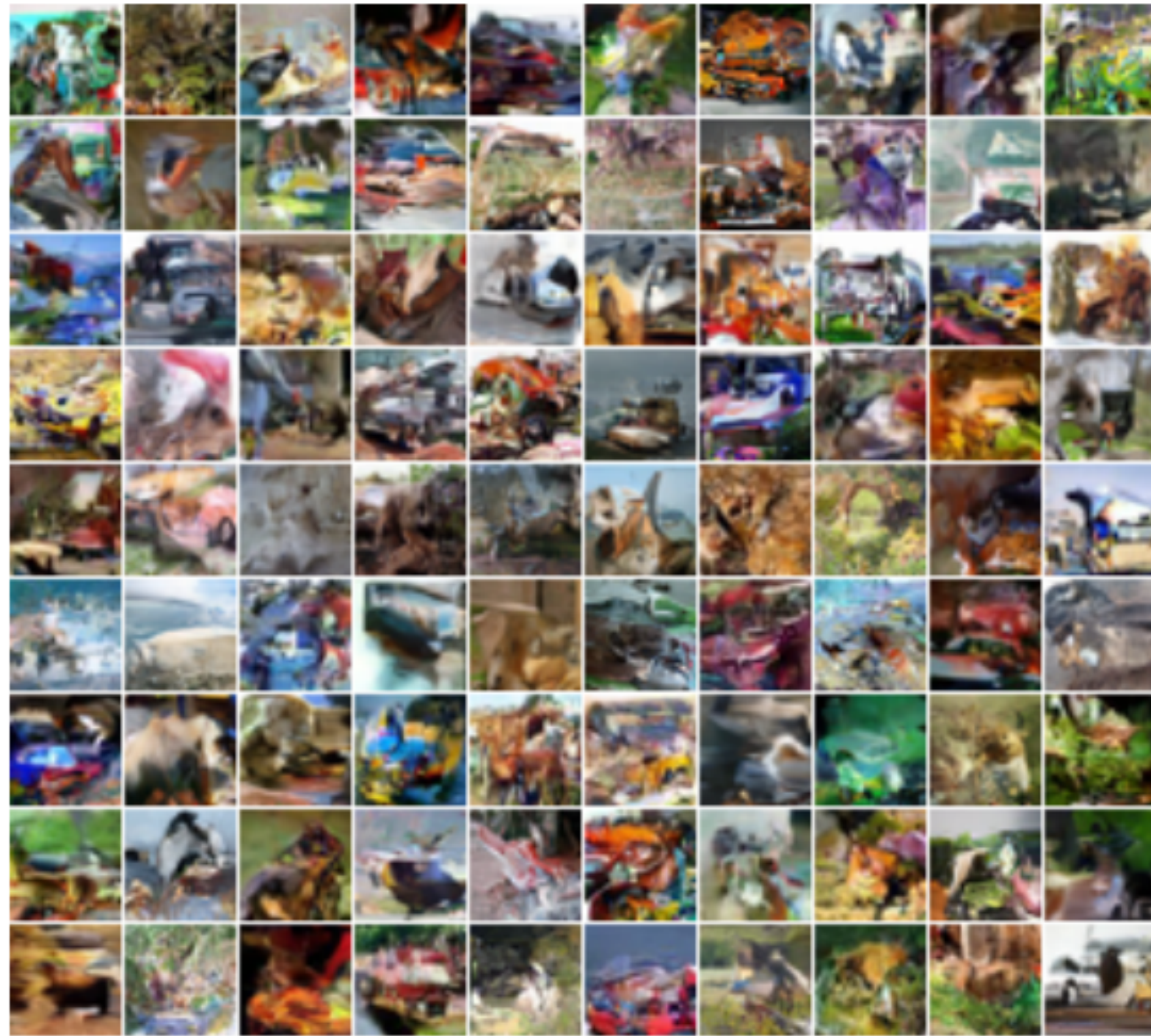
$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$



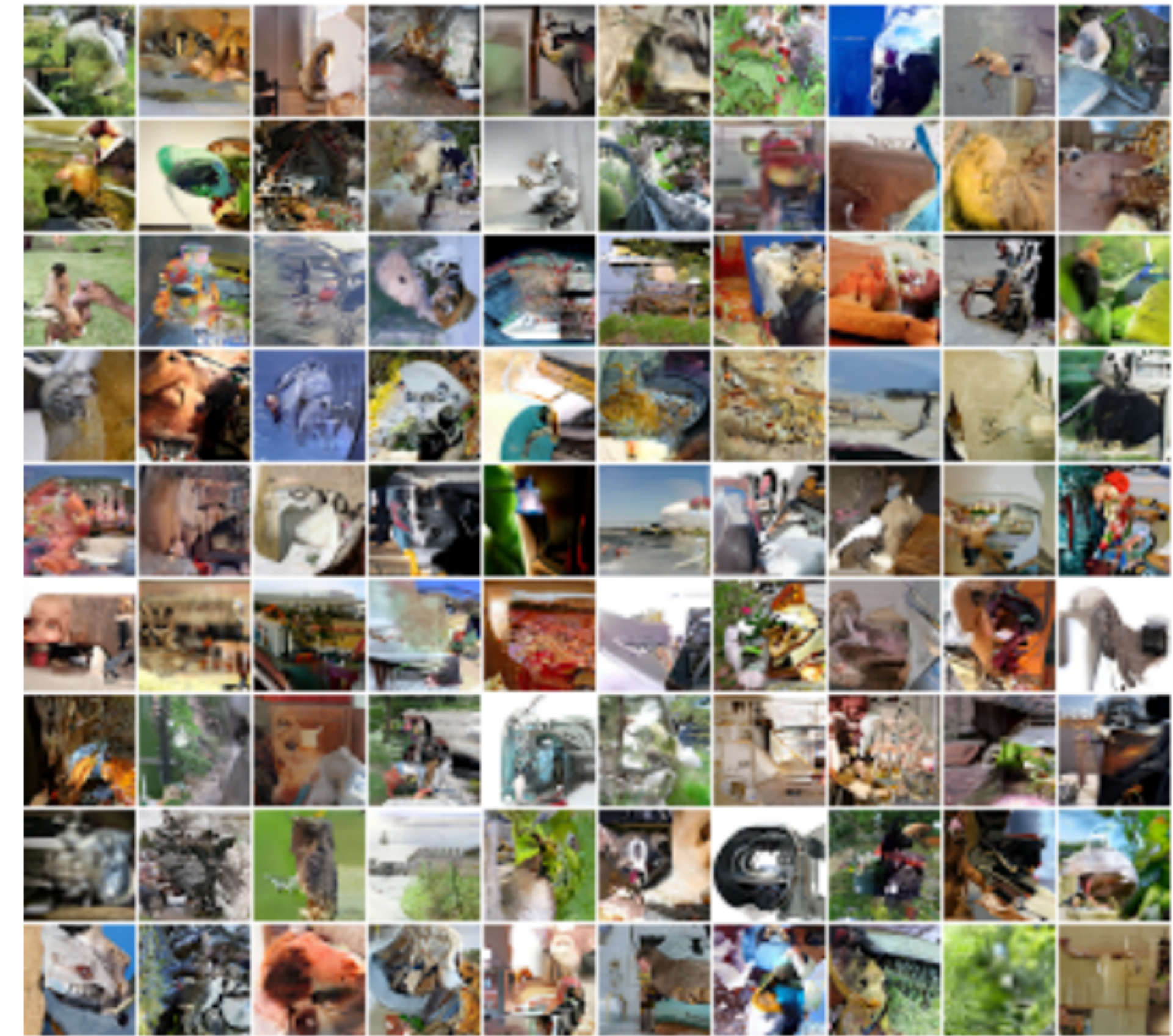
Generation is still slow (sequential), but learning is faster

Generated Samples

[van der Oord et al., 2016]



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

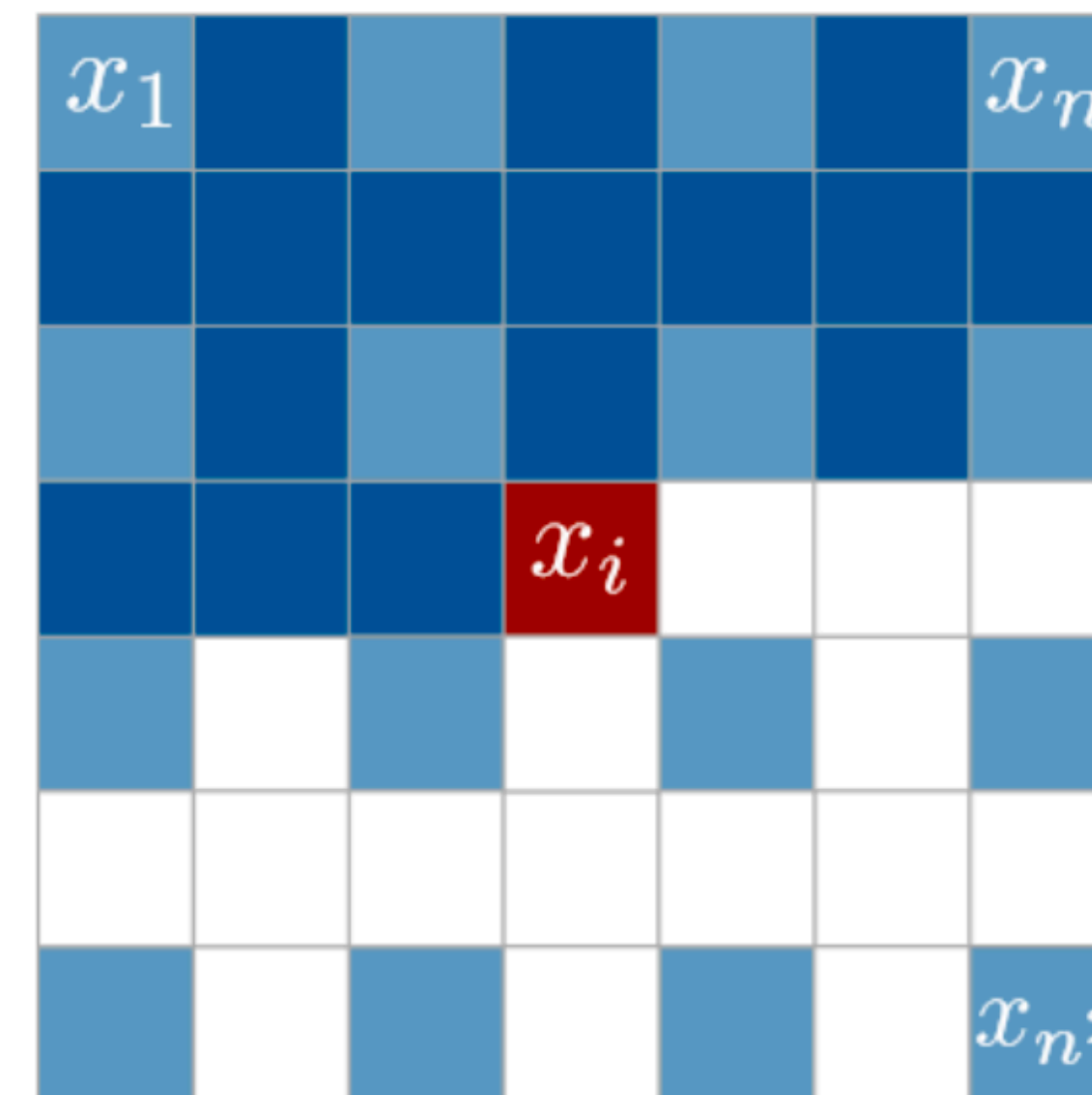
- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

Multi-scale PixelRNN

[van der Oord et al., 2016]

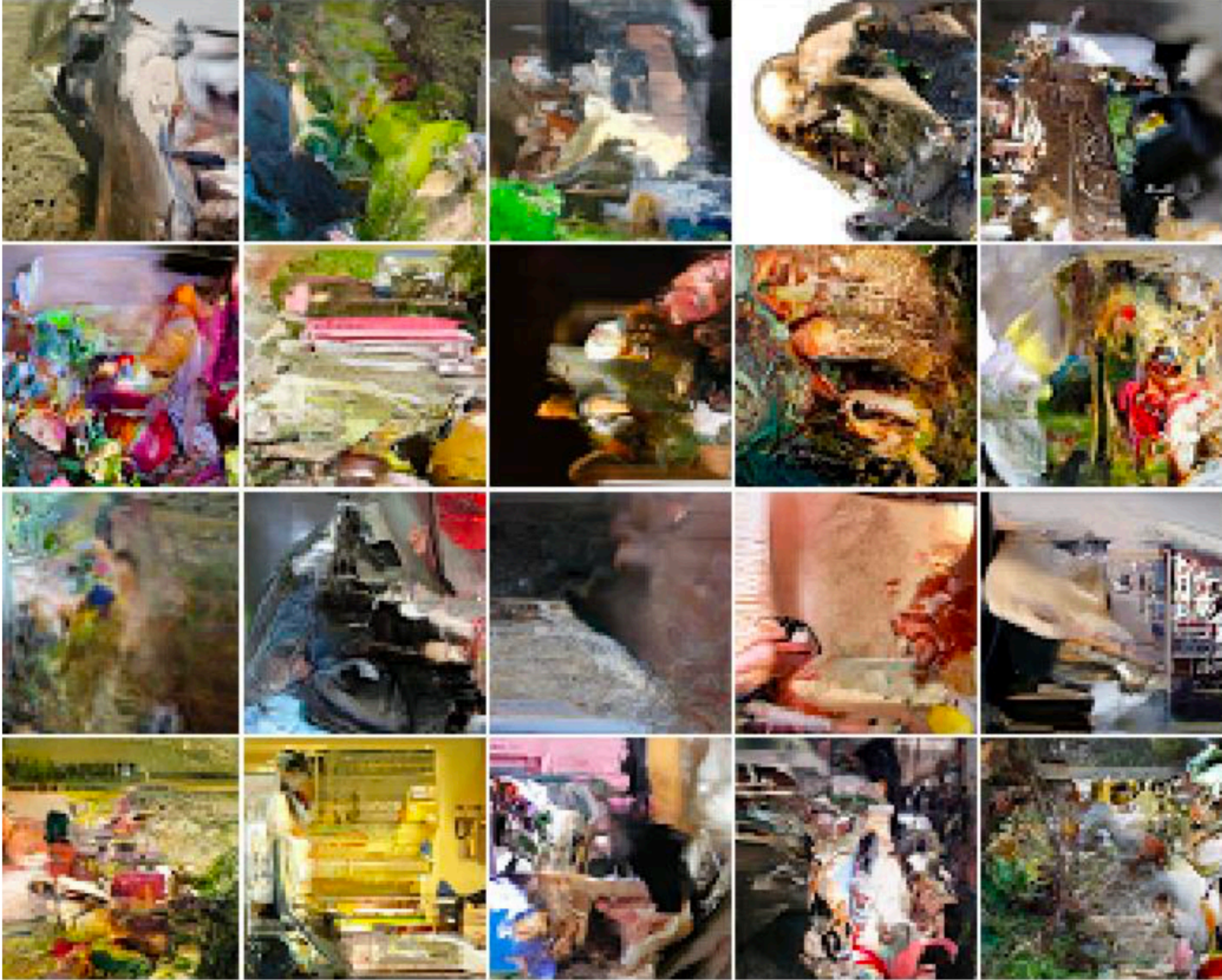
Take sub-sampled pixels as additional input pixels

Can capture better global information (more visually coherent)



Multi-scale PixelRNN

[van der Oord et al., 2016]



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

Conditional Image Generation

[van der Oord et al., 2016]

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

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



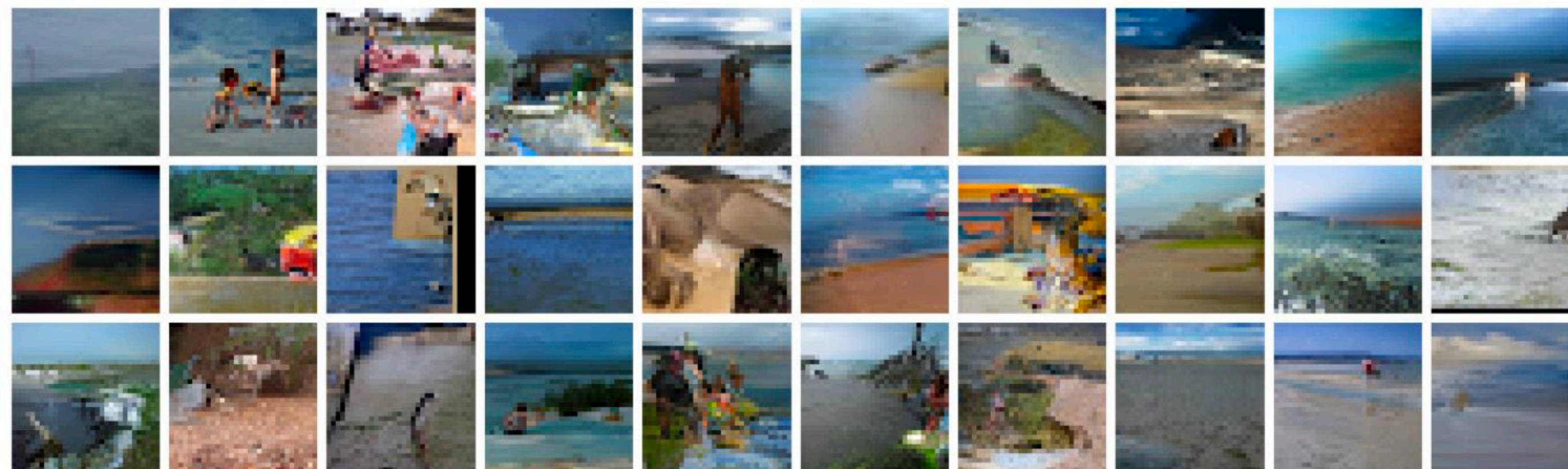
$$p(\mathbf{x}|\mathbf{h}) = p(x_1, x_2, \dots, x_{n^2}|\mathbf{h})$$

Conditional Image Generation

[van der Oord et al., 2016]

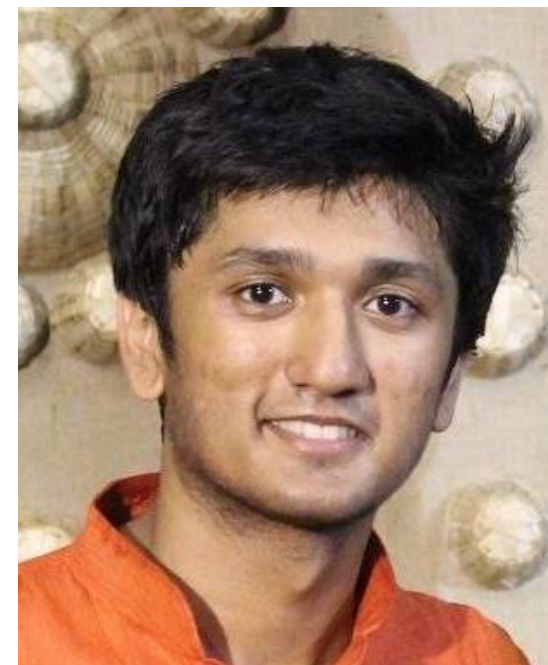


African elephant



Sandbar

Attention RNN: Structured Spatial Attention Mechanism



Siddhesh Khandelwal

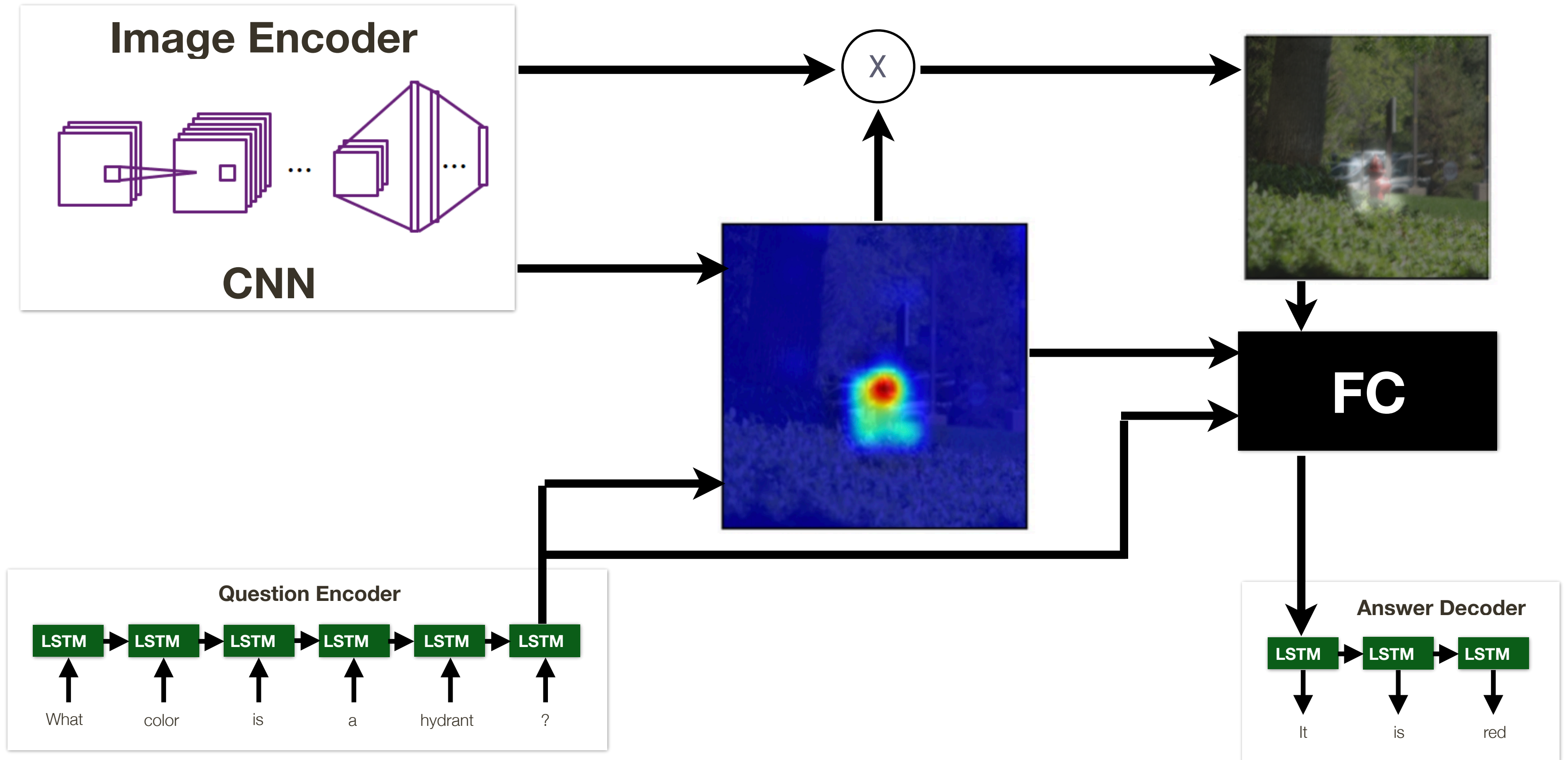
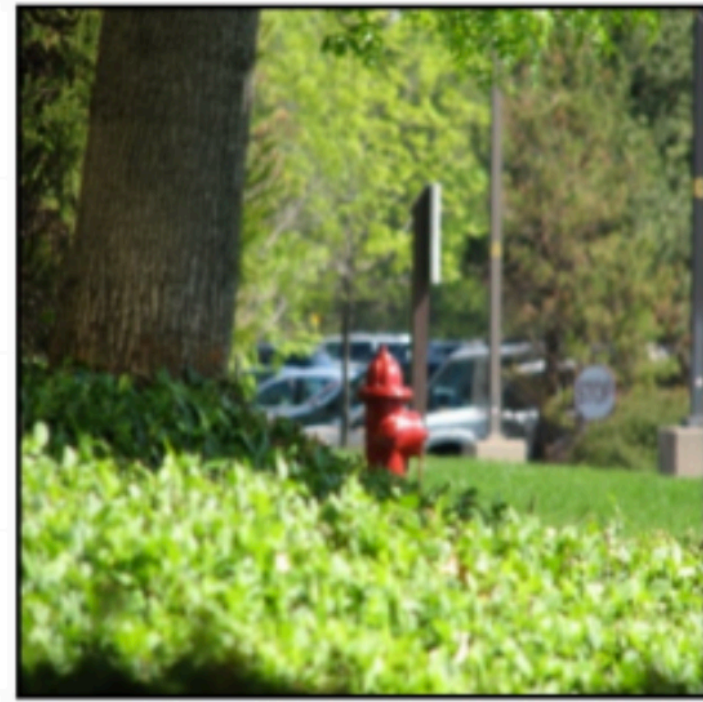


Leonid Sigal



Motivation

Attention is widely used in vision: helps identify relevant regions of the image



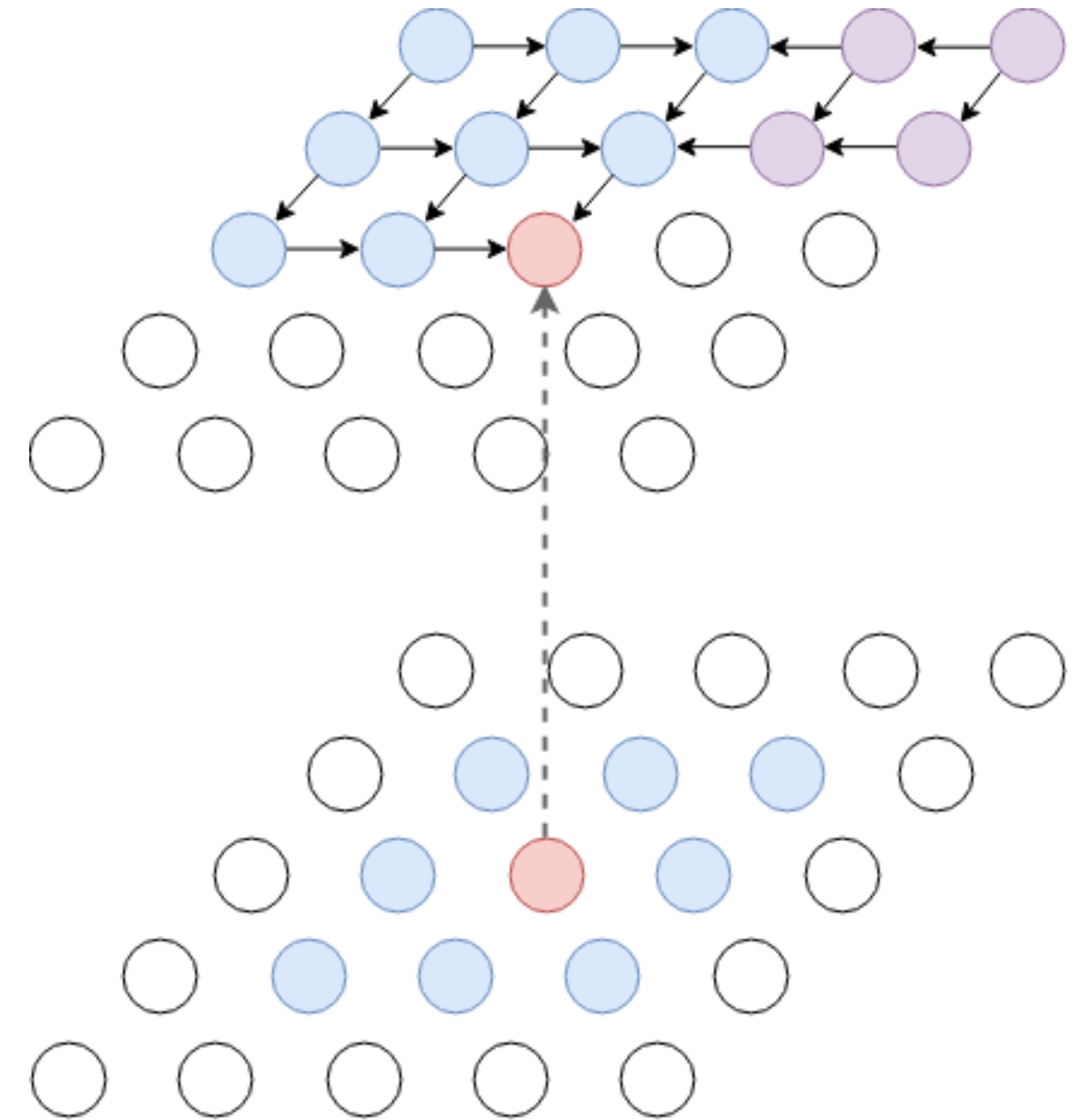
Q: What color is a hydrant?

A: It is red

AttentionRNN: Structured Spatial Attention

Novel **autoregressive attention mechanism** that can encode structural dependencies among attention values

- Inspired by diagonal Bi-LSTM architecture from PixelRNN
- Spatial attention values are generated sequentially
- Image is traversed diagonally from top-left to bottom-right



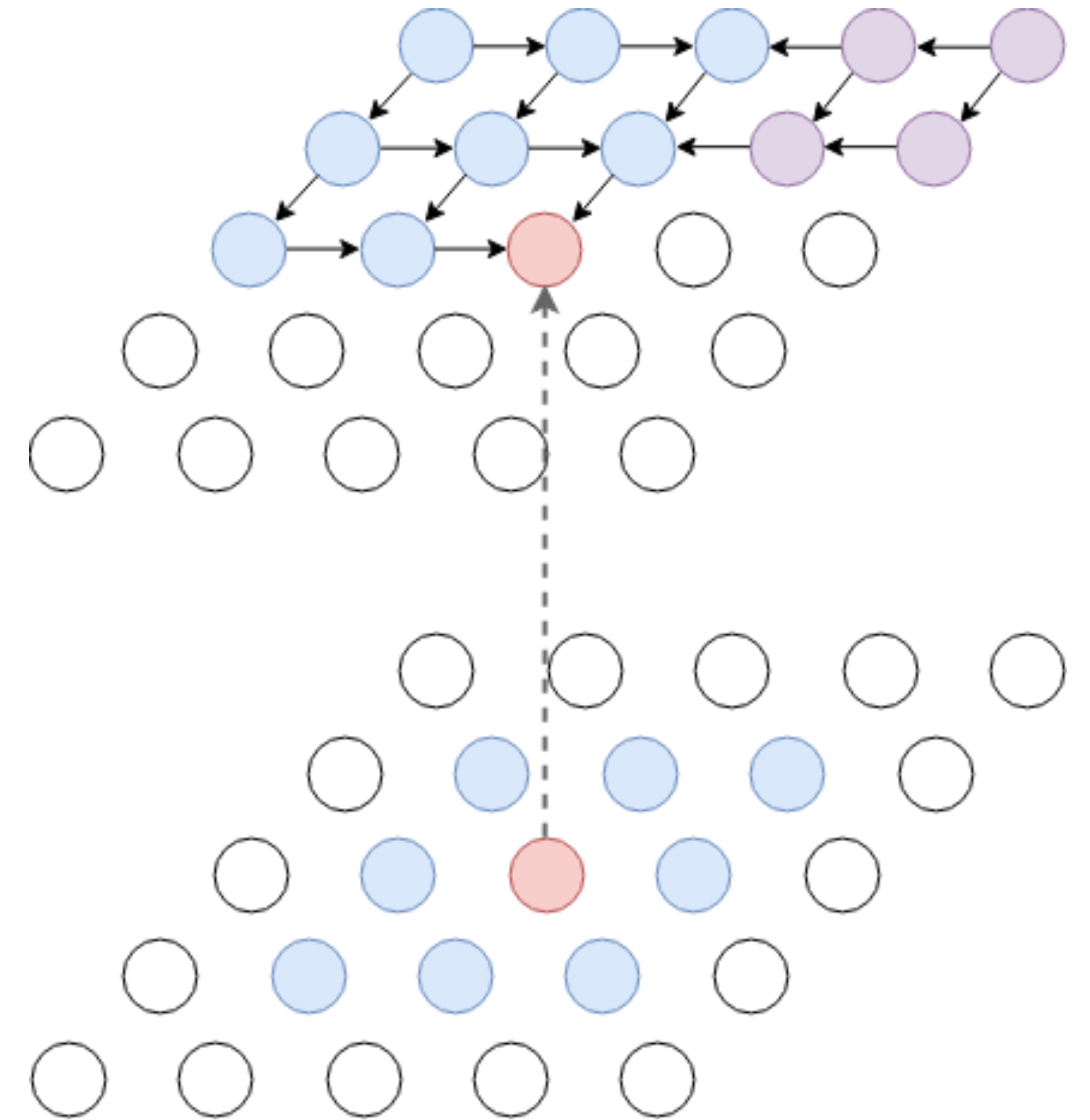
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- Previously generated attention values

Novel **autoregressive attention mechanism** that can encode structural dependencies among attention values

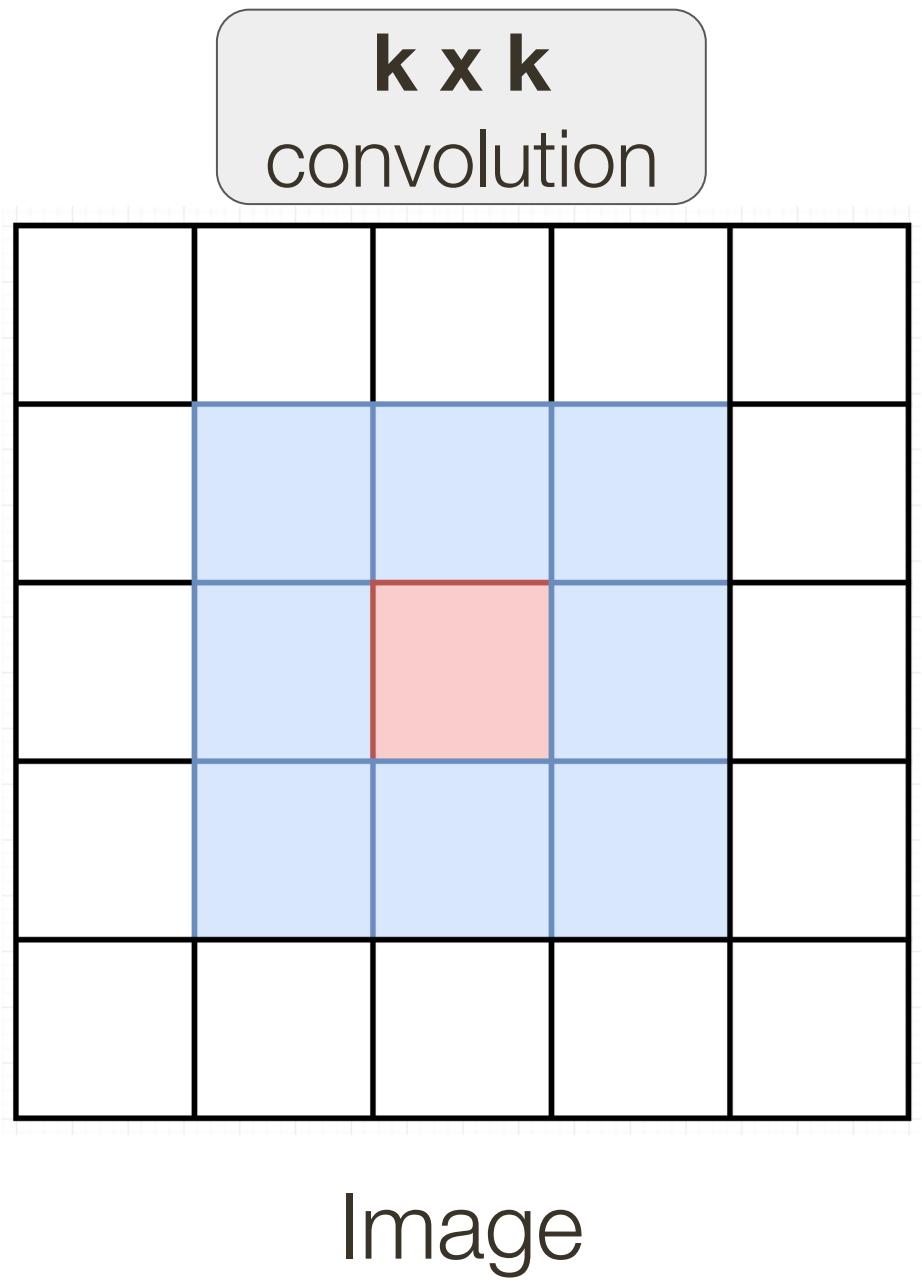
- Inspired by diagonal Bi-LSTM architecture from PixelRNN
- Spatial attention values are generated sequentially
- Image is traversed diagonally from top-left to bottom-right



AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

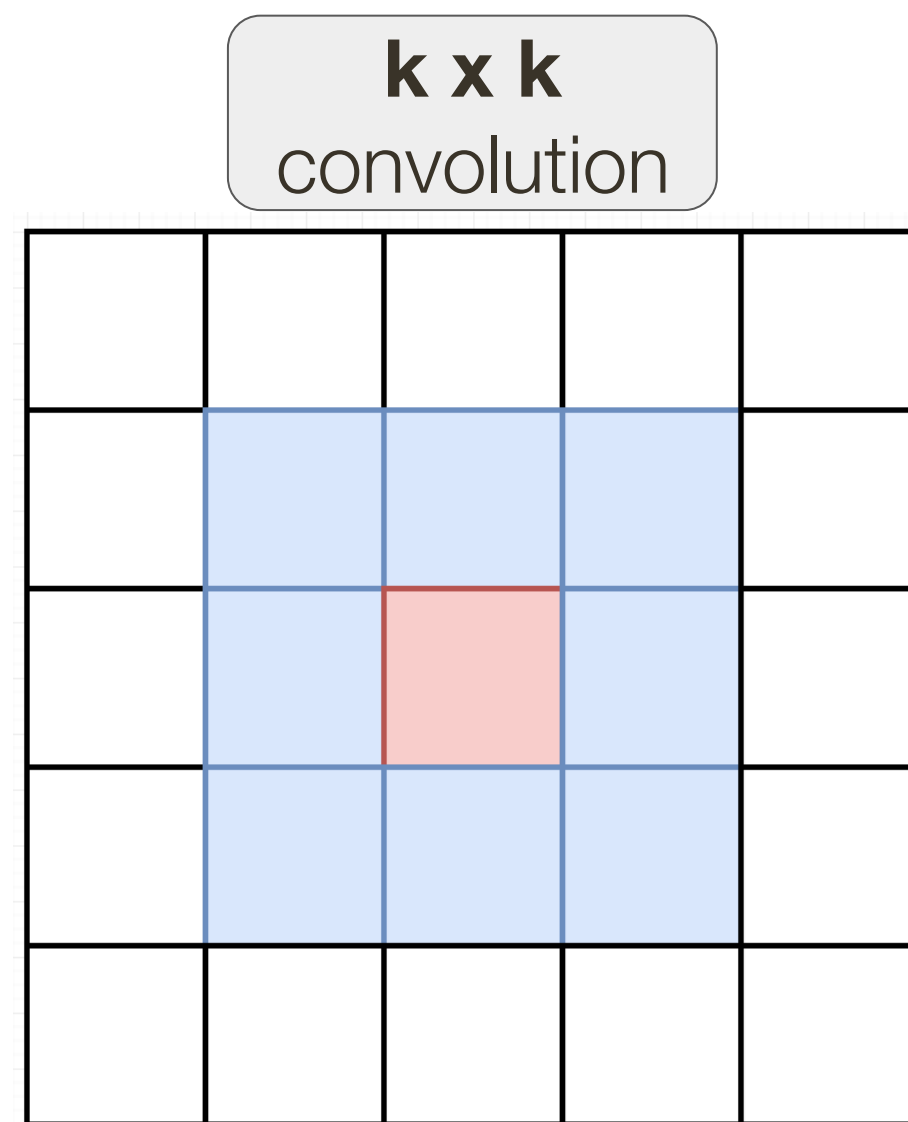
- **Local image context**
- Previously generated attention values



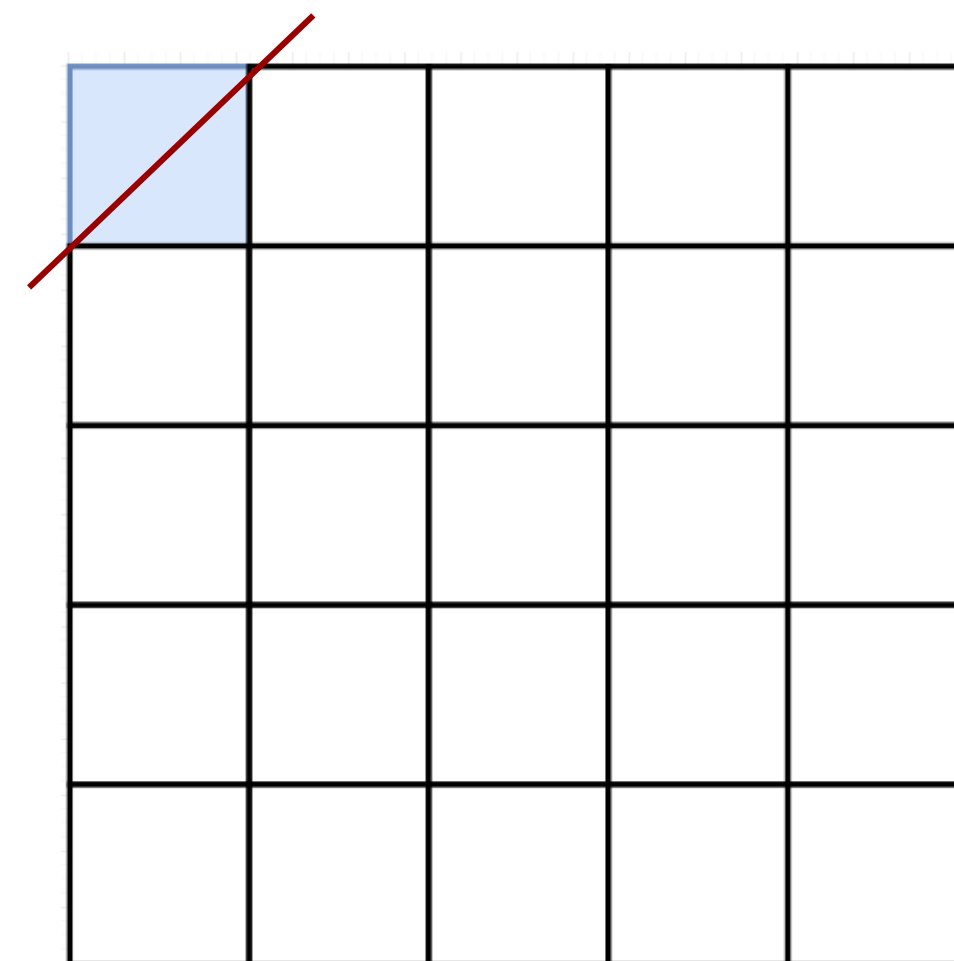
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- **Previously generated attention values**



Image

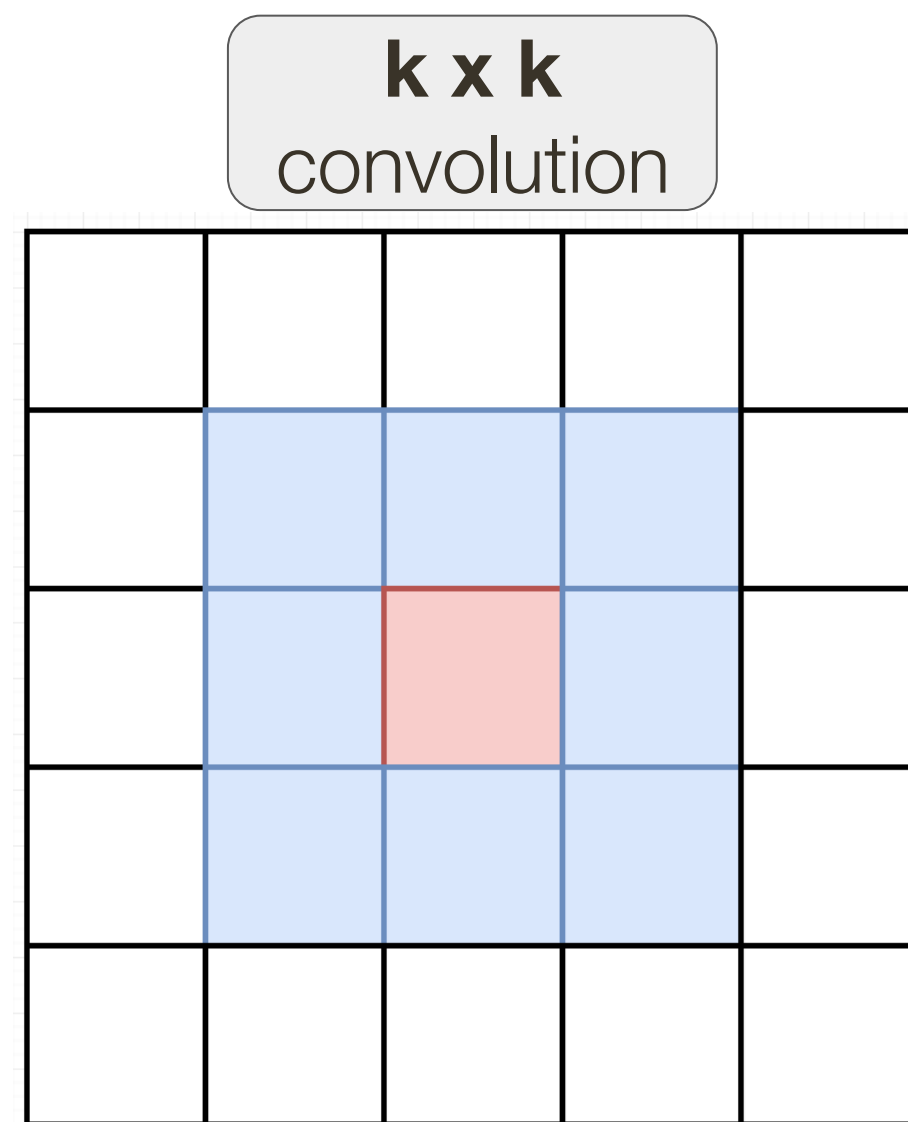


Attention Mask

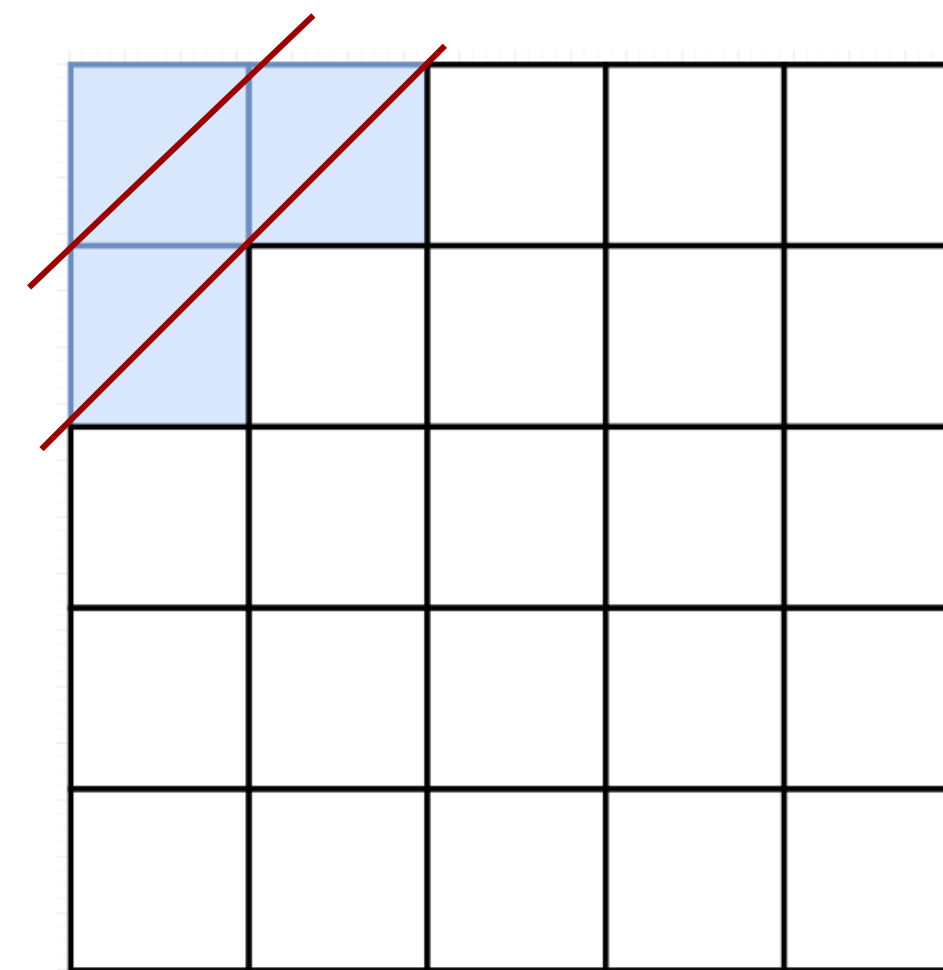
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- **Previously generated attention values**



Image

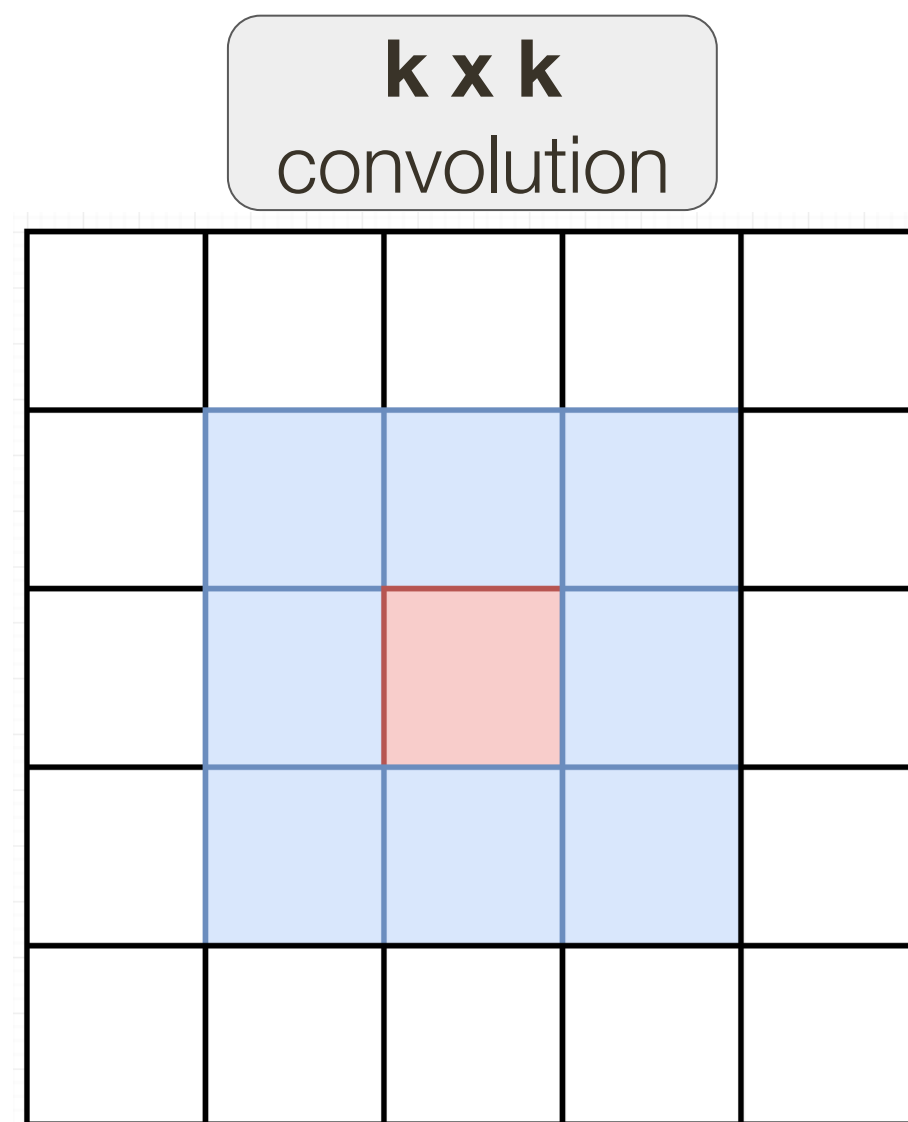


Attention Mask

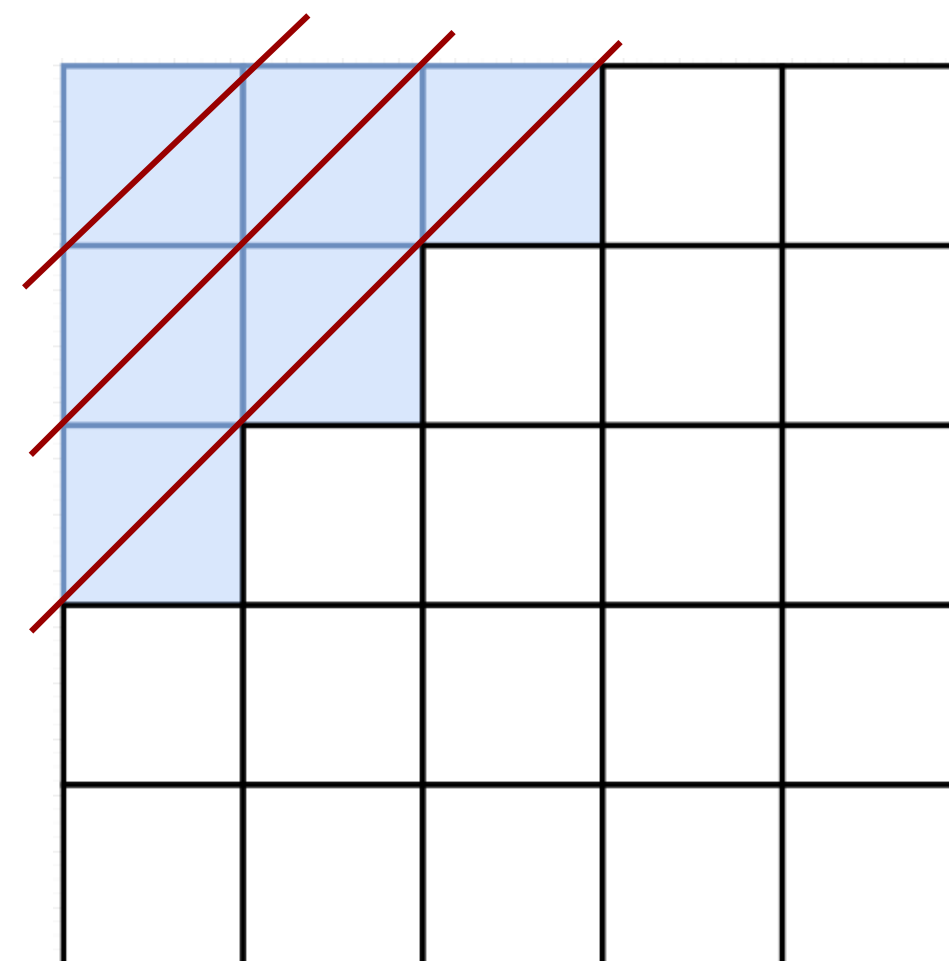
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- **Previously generated attention values**



Image

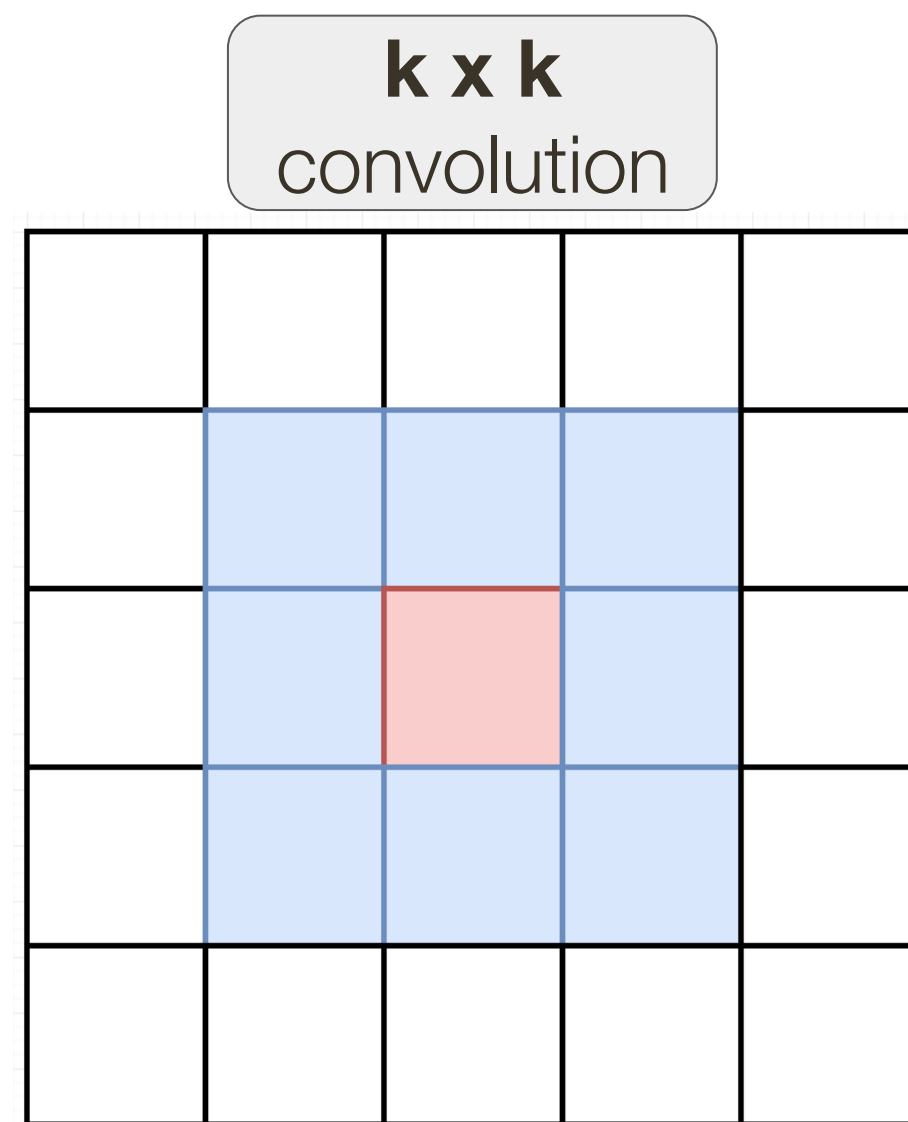


Attention Mask

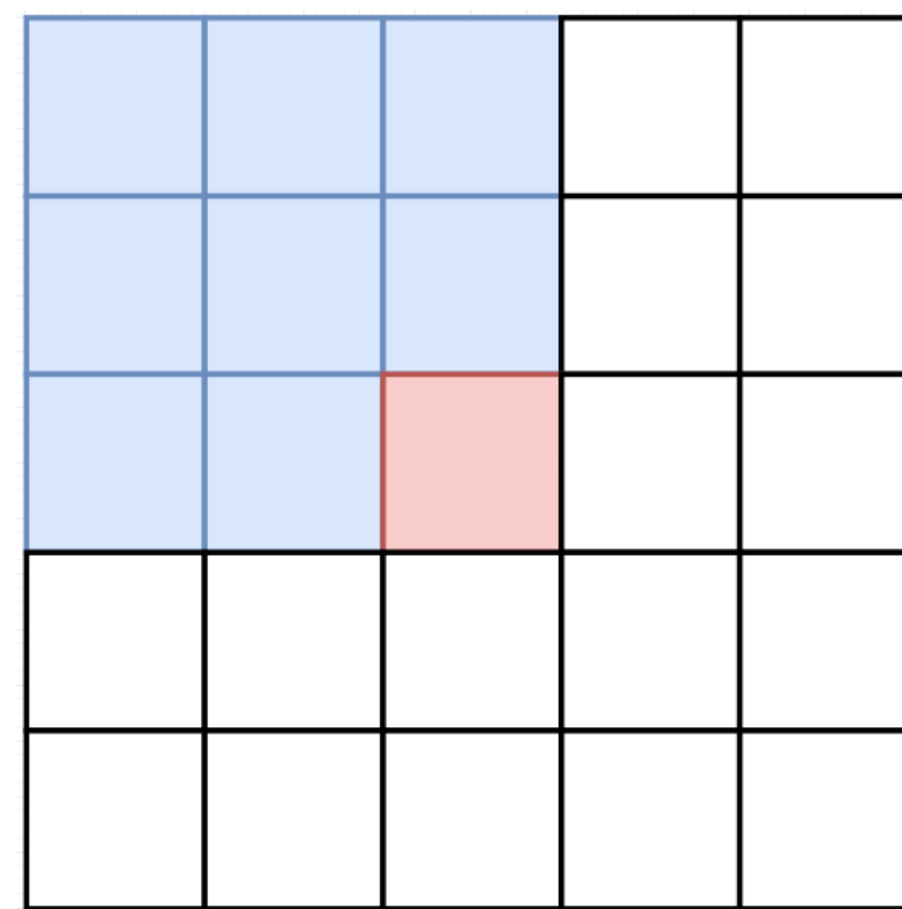
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- **Previously generated attention values**



Image

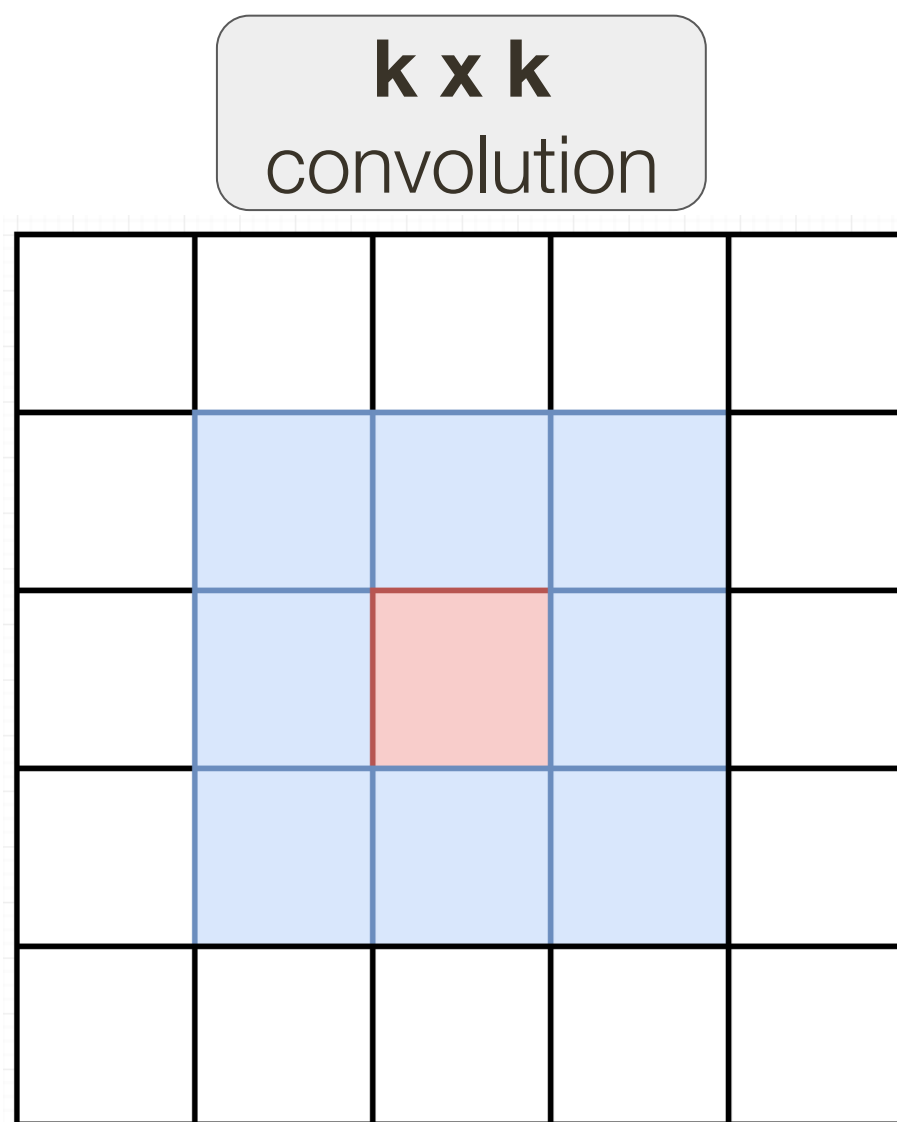


Attention Mask

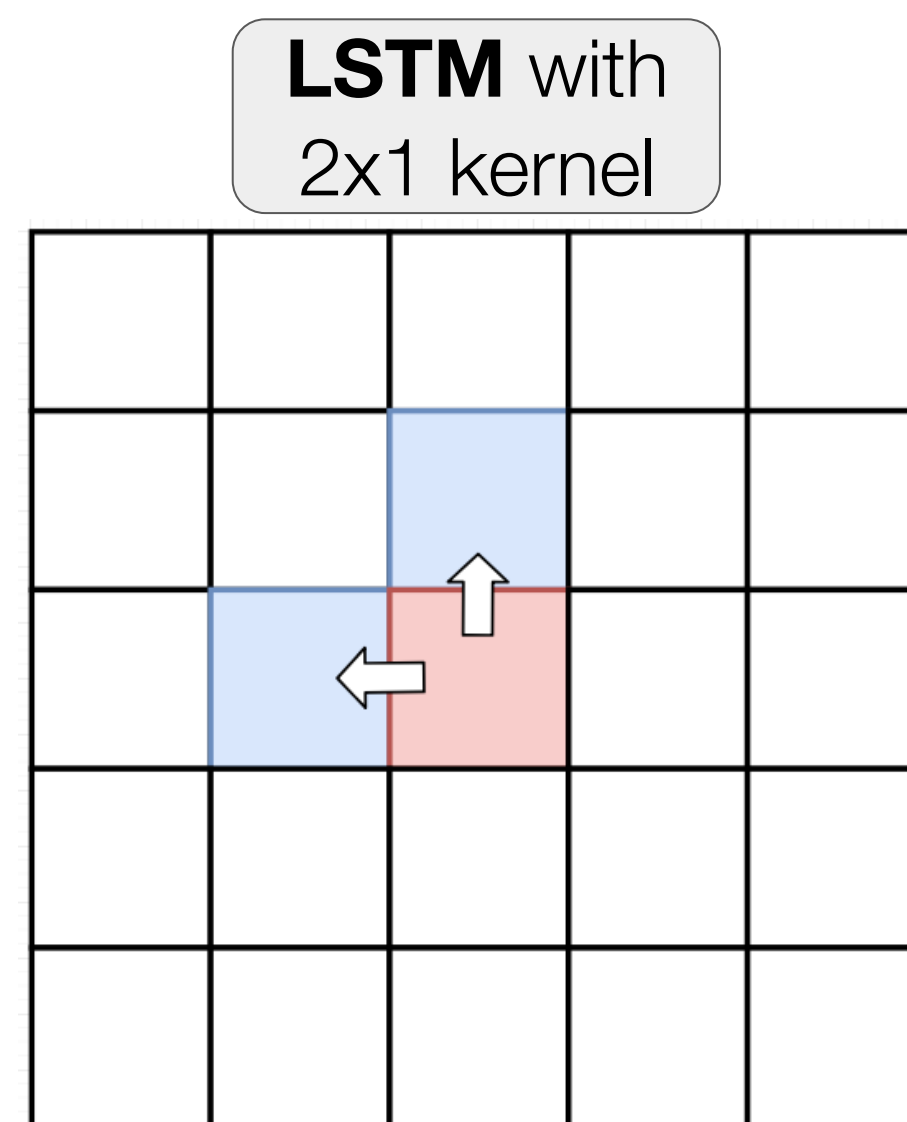
AttentionRNN: Structured Spatial Attention

Each **attention value** depends on

- Local image context
- Previously generated attention values



Image

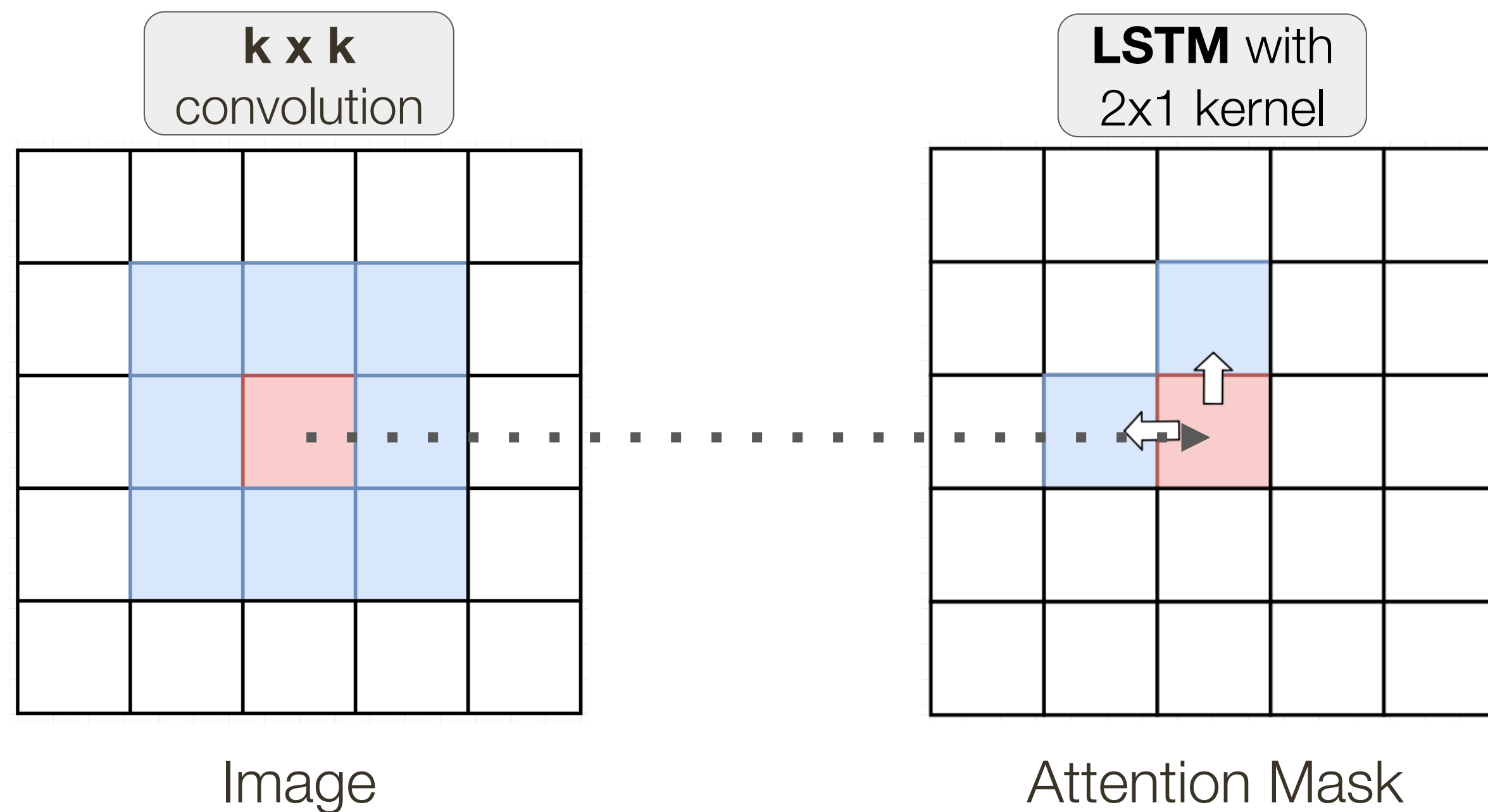


Attention Mask

AttentionRNN: Structured Spatial Attention

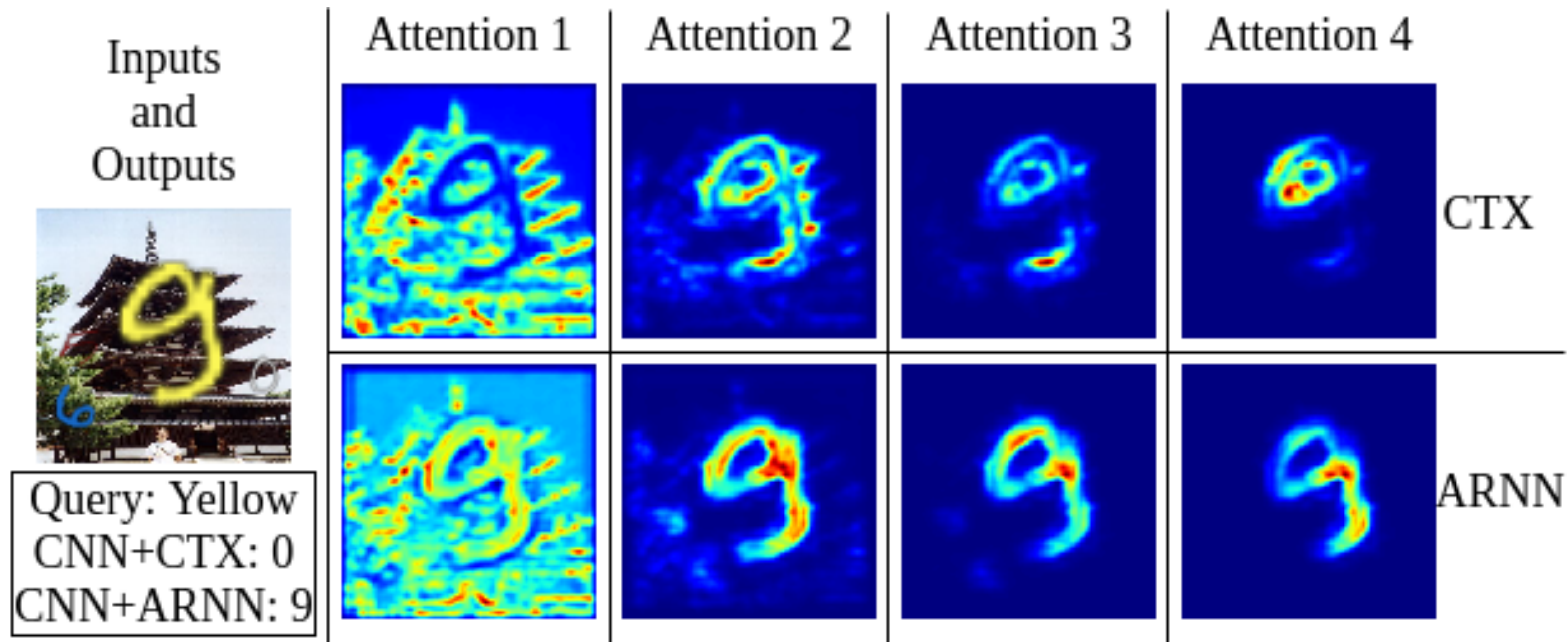
Each **attention value** depends on

- Local image context
- Previously generated attention values



Experiments: Visual Digit Prediction

Task: Given an image, predict a digit number specified by a query color



	SAN	\neg CTX	CTX	ARNN
Correctness	0.1258	0.2017	0.2835	0.3729

Variational Autoencoders

(VAE)

So far ...

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

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

So far ...

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

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

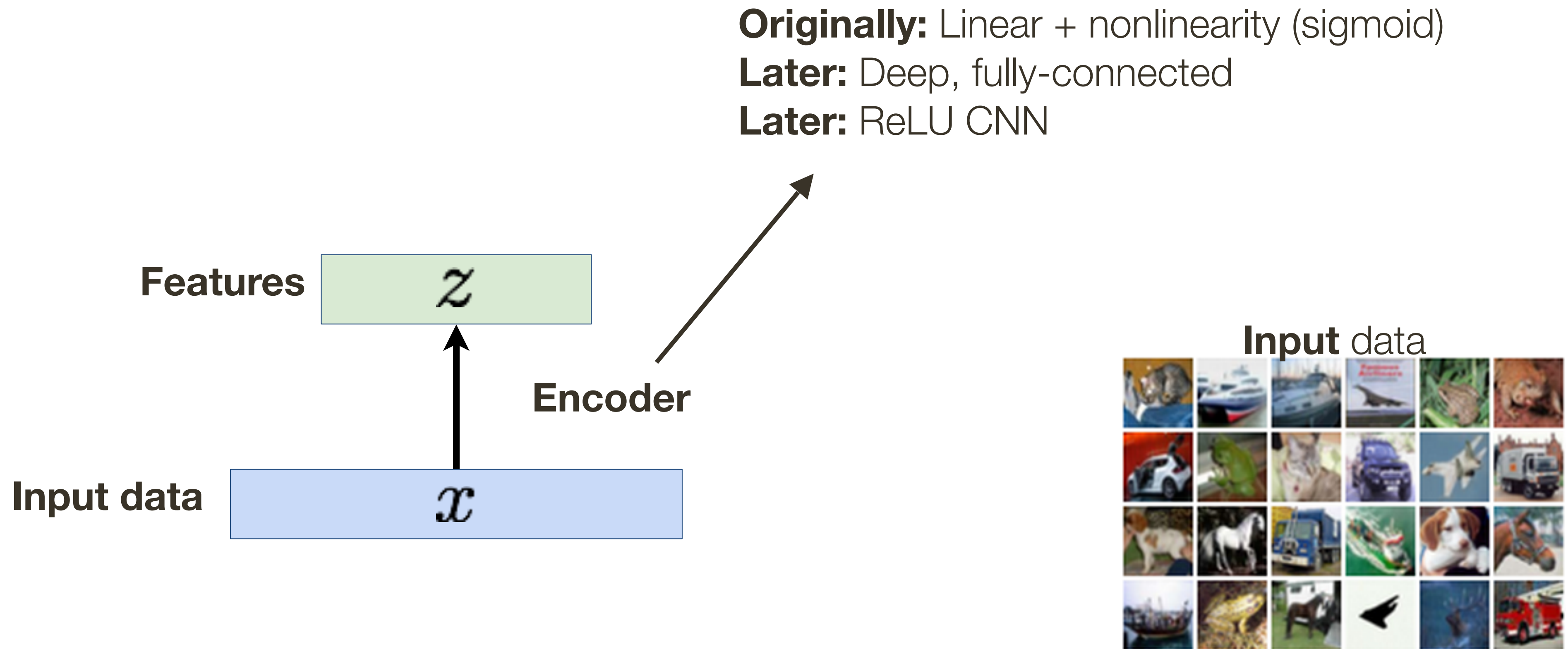
VAEs define intractable density function with latent variables z (that we need to marginalize):

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

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

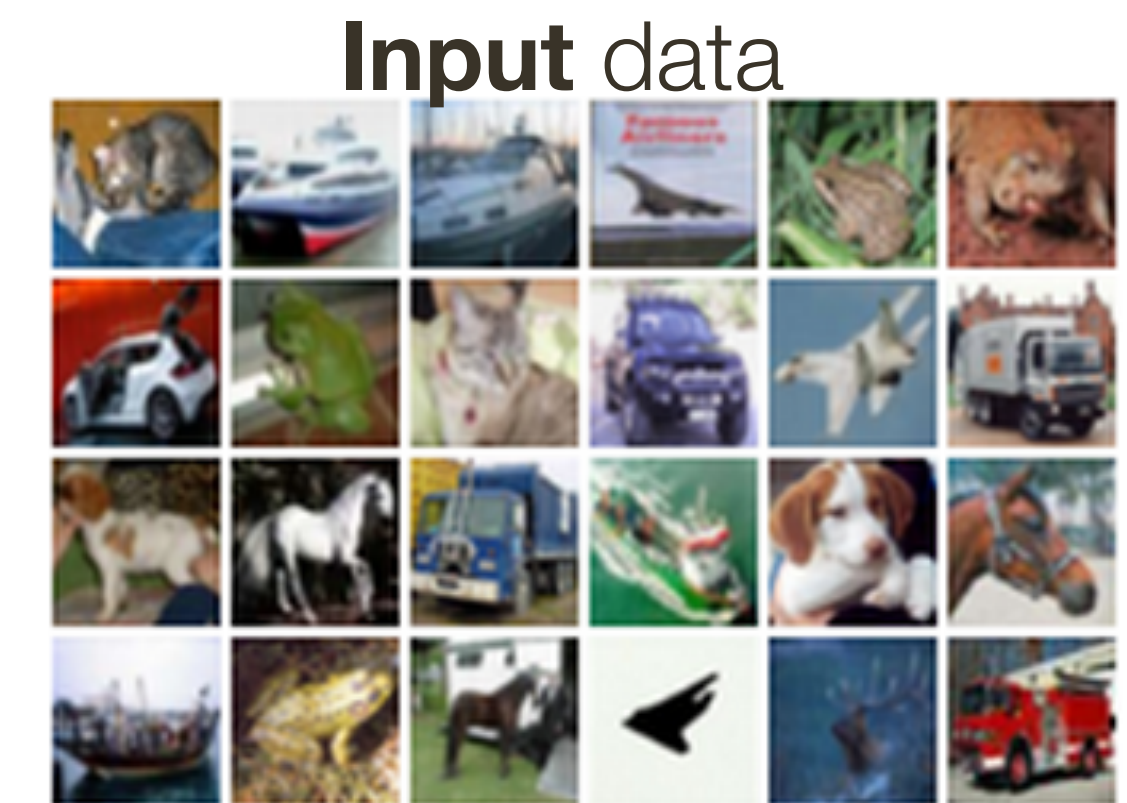
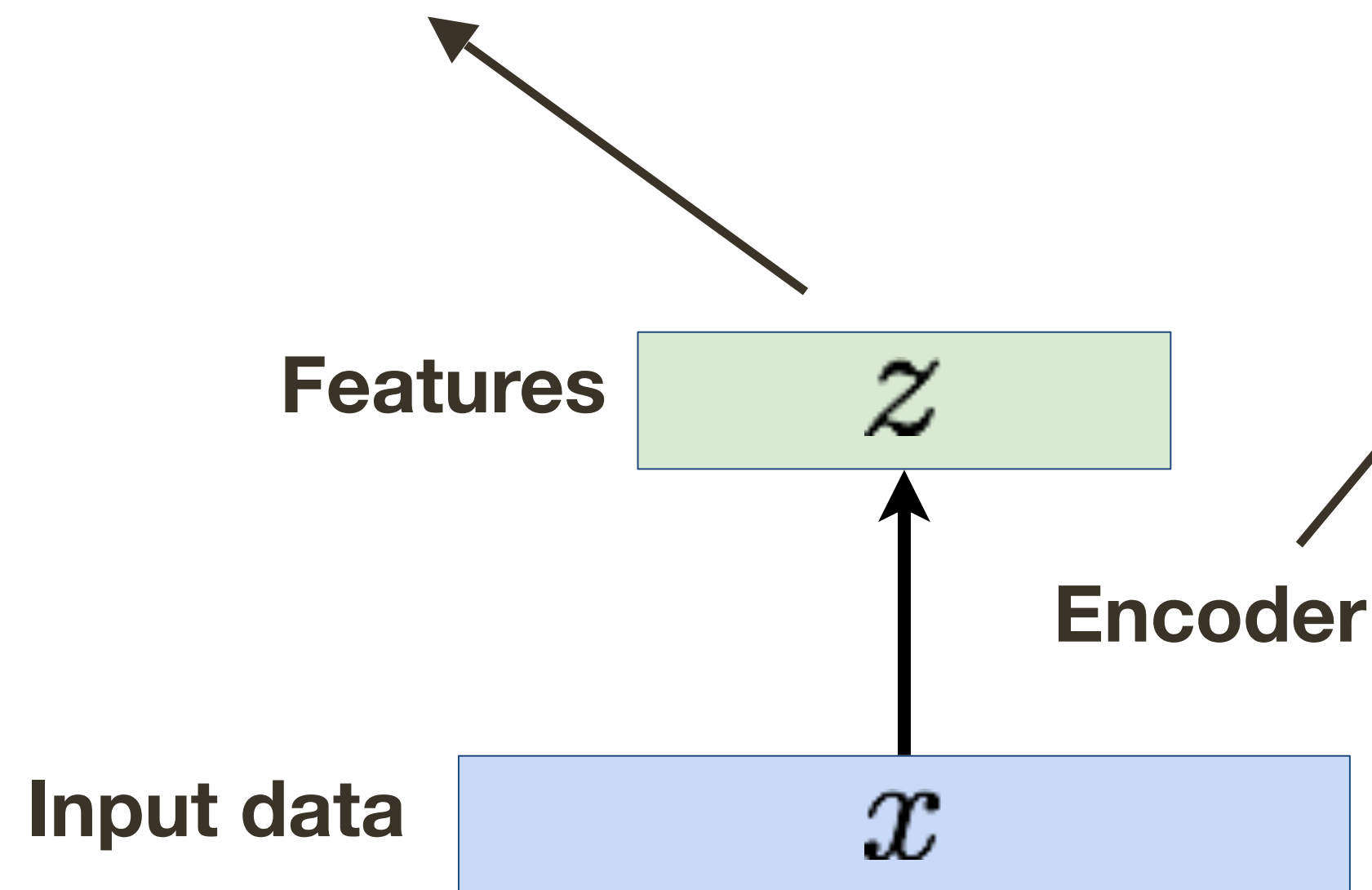


Autoencoders Reminder ...

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

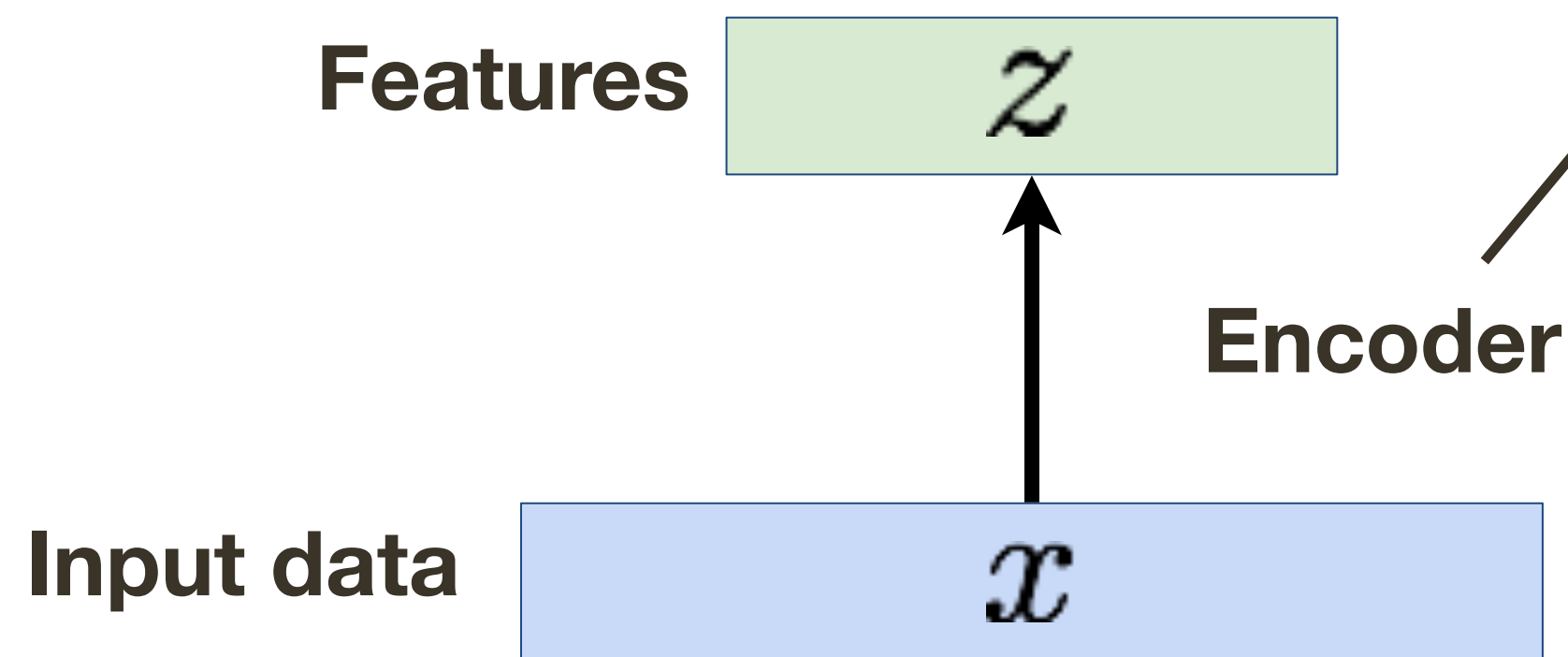


Autoencoders Reminder ...

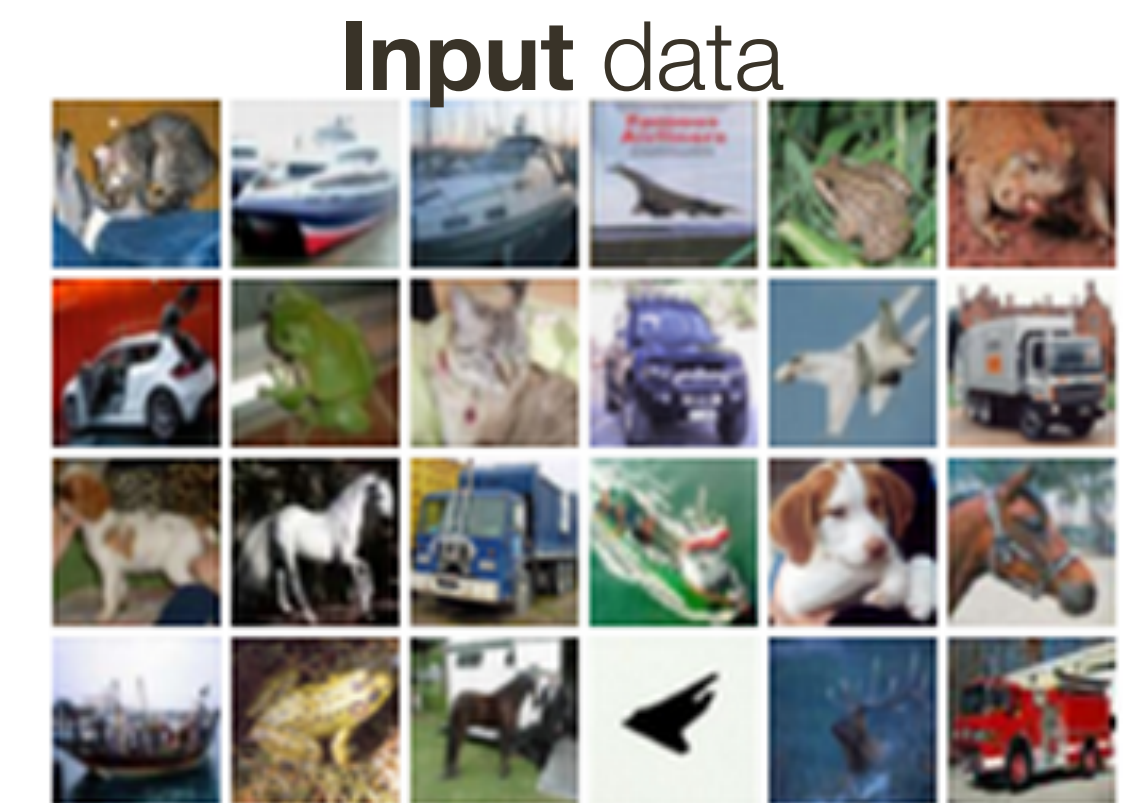
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

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

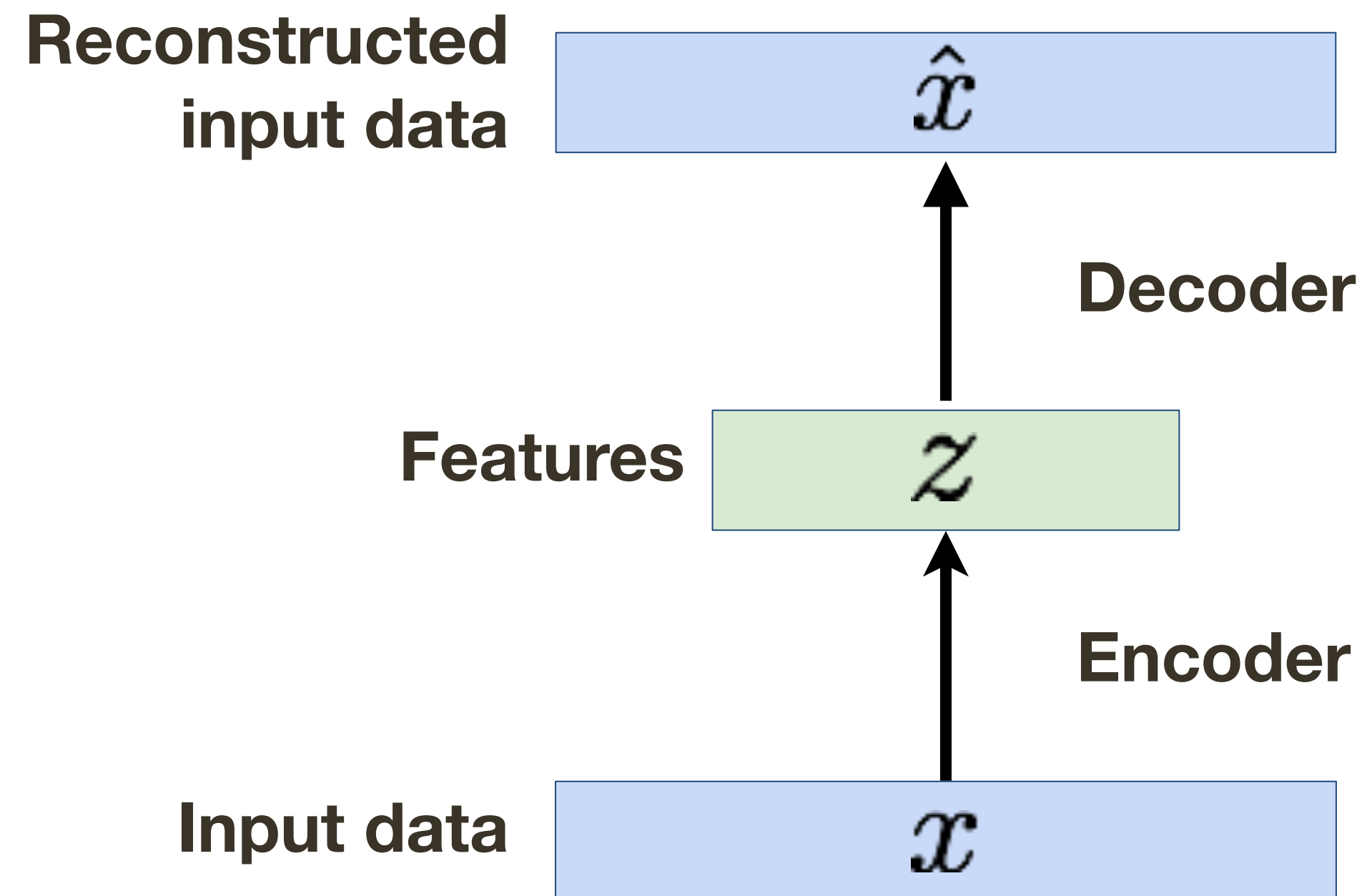


Originally: Linear + nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN



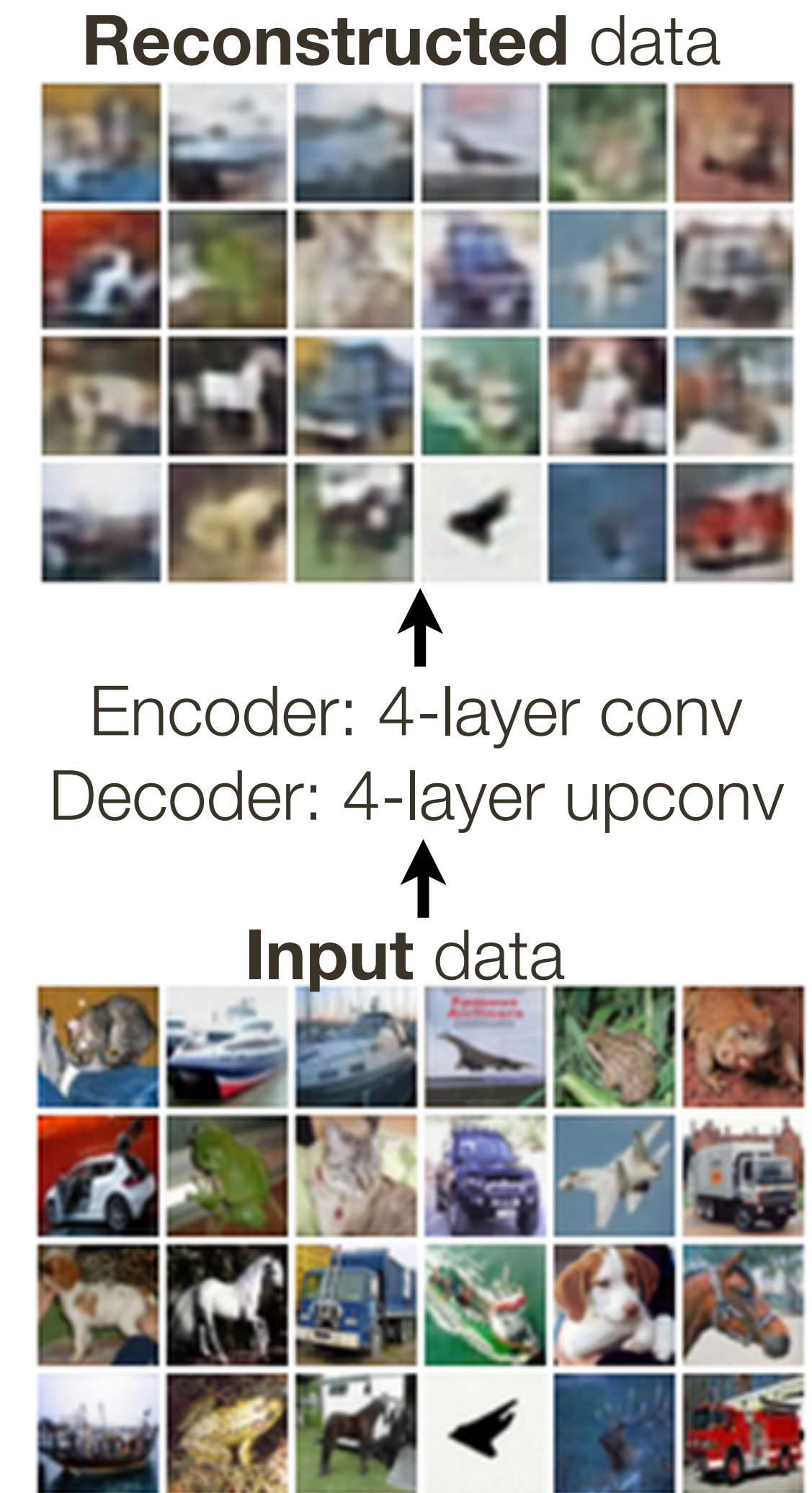
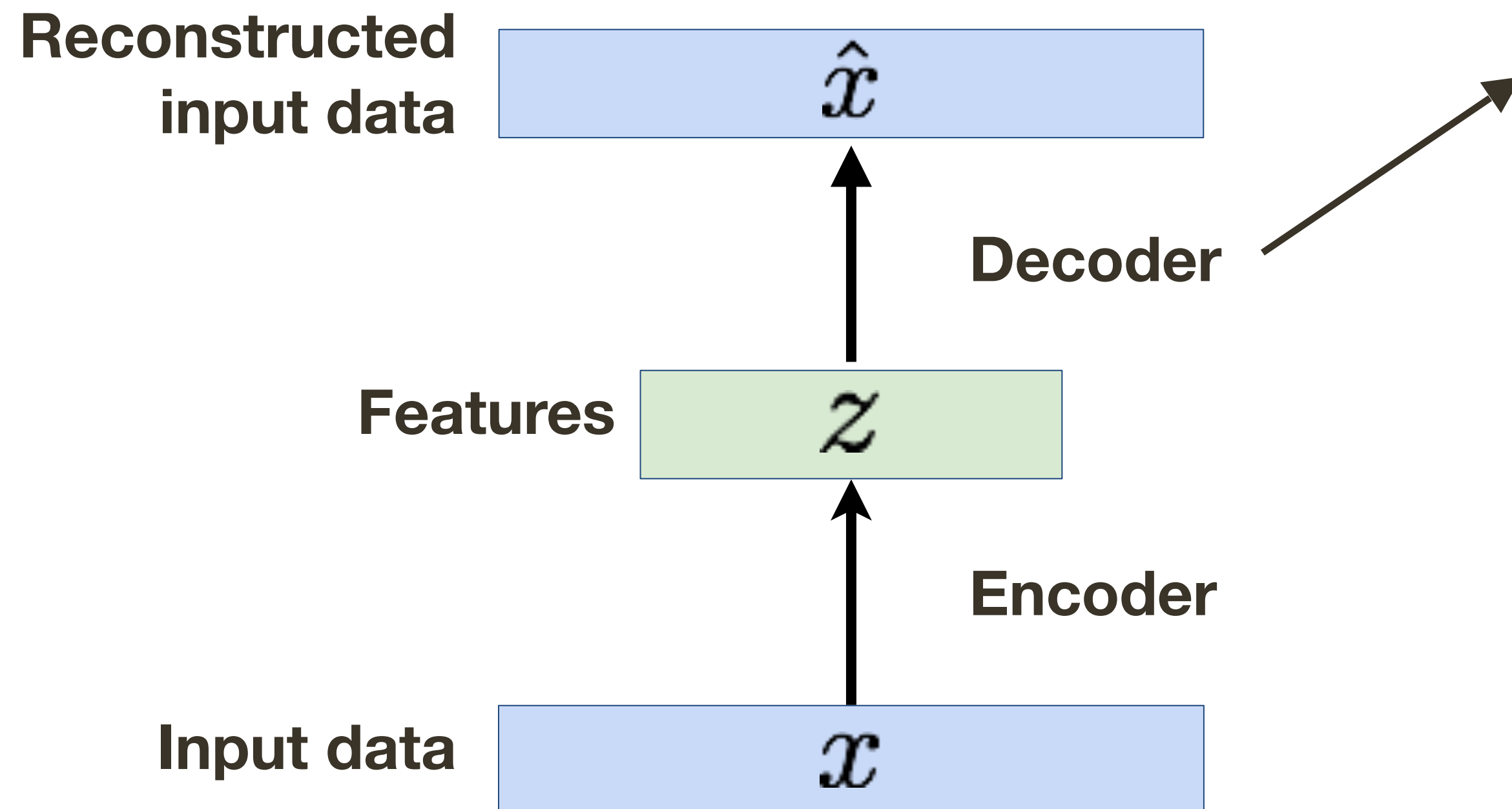
Autoencoders Reminder ...

Train such that features can reconstruct original data best they can

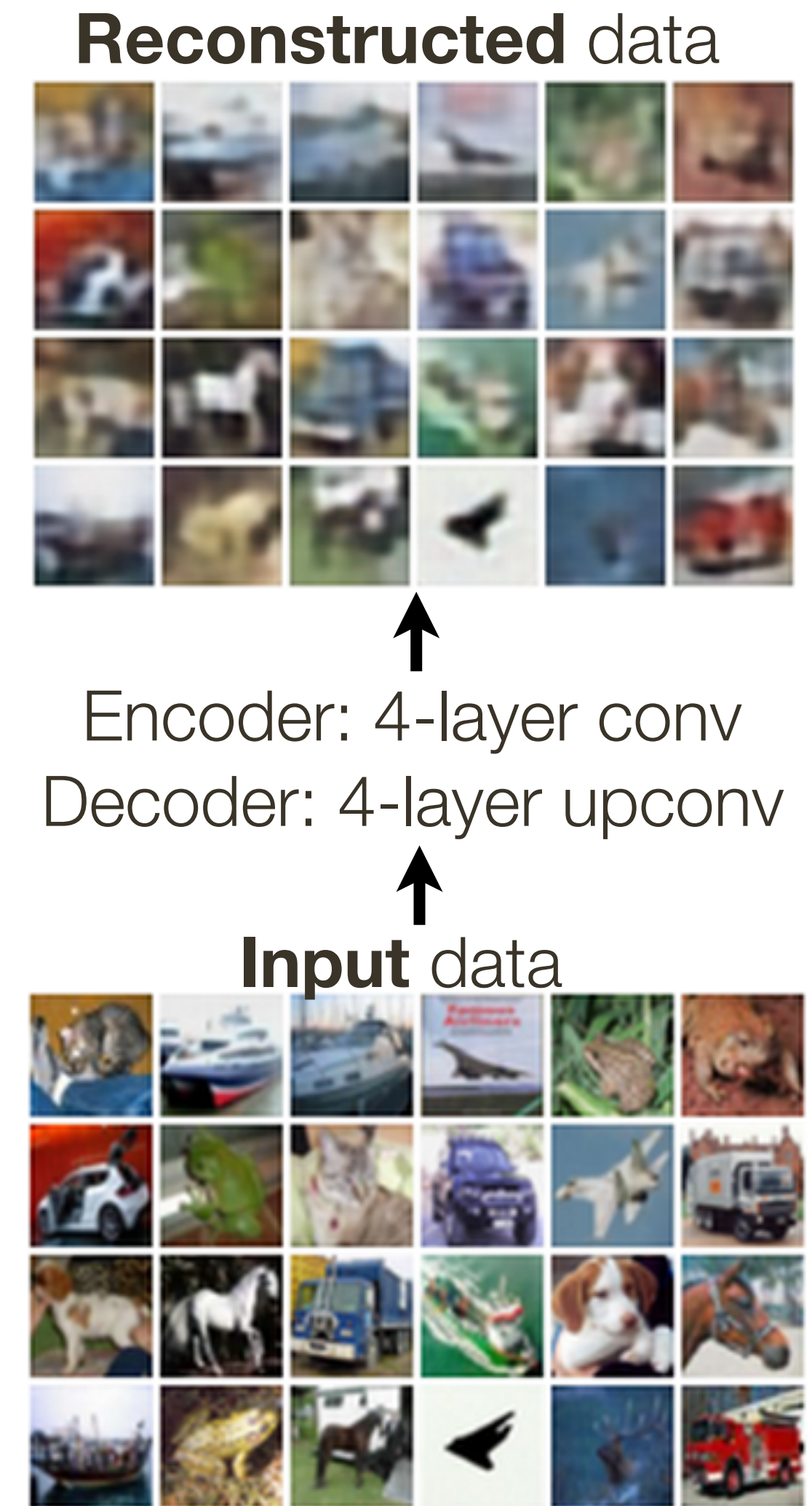
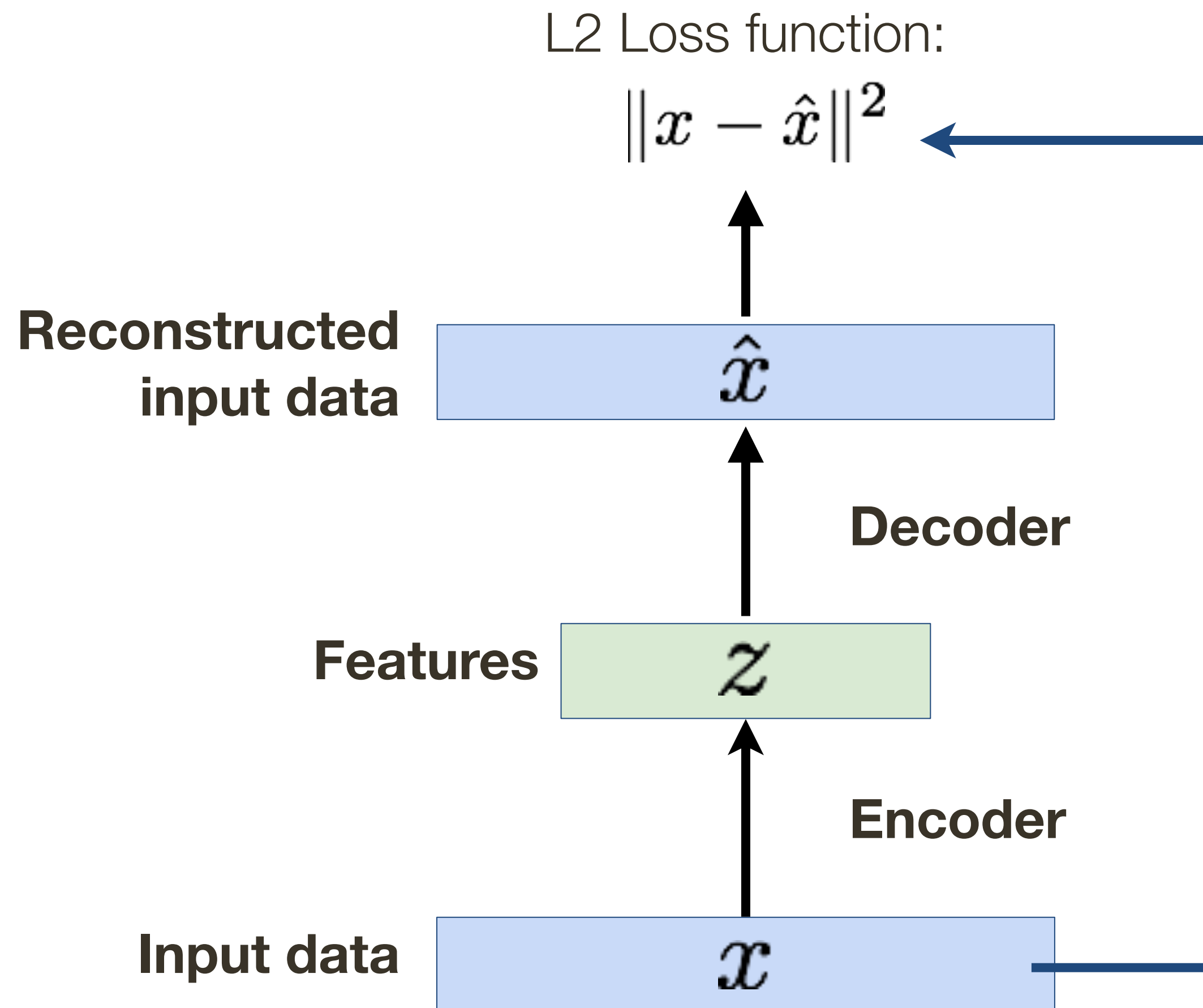


Autoencoders Reminder ...

Train such that features can reconstruct original data best they can

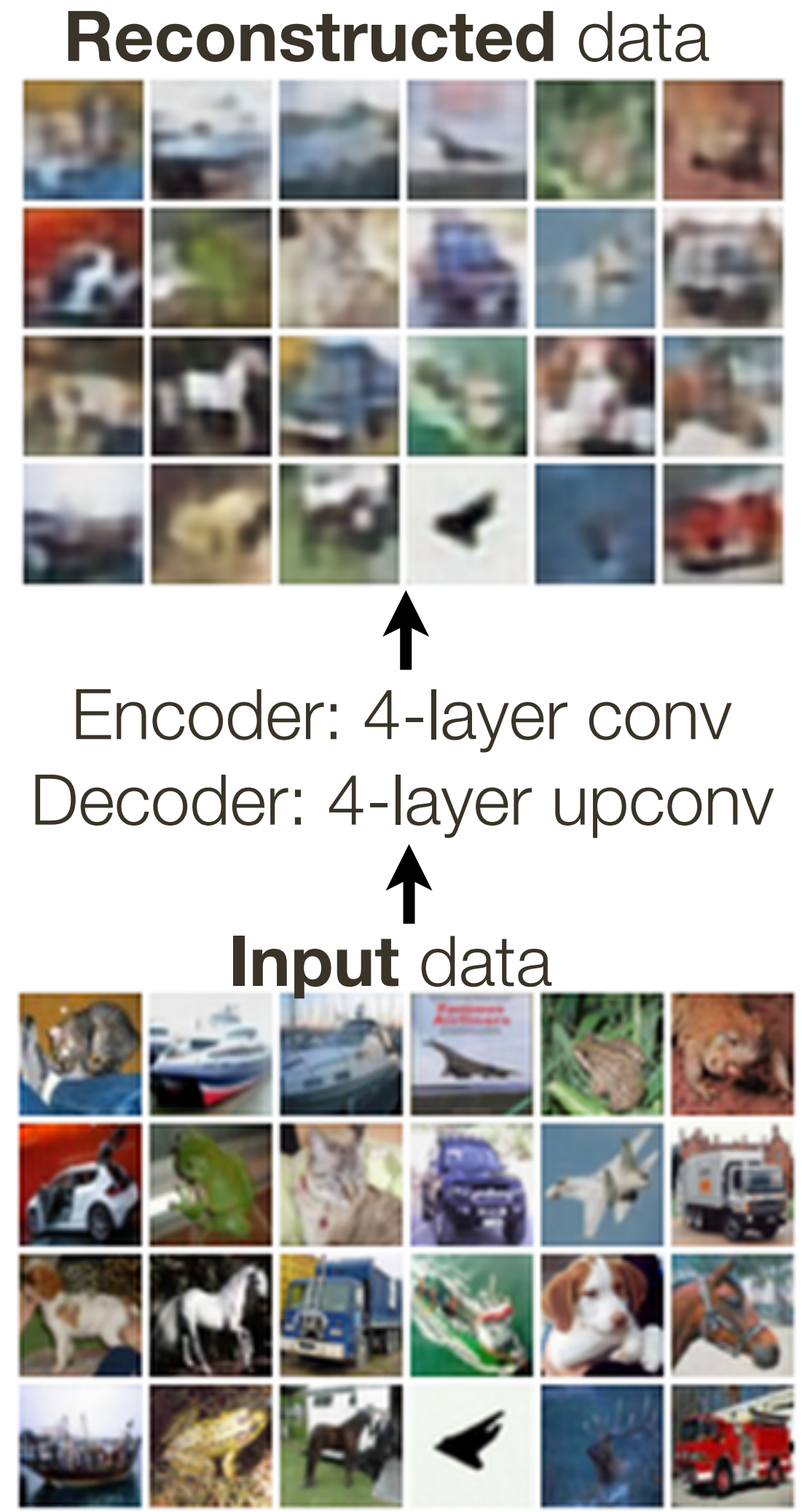
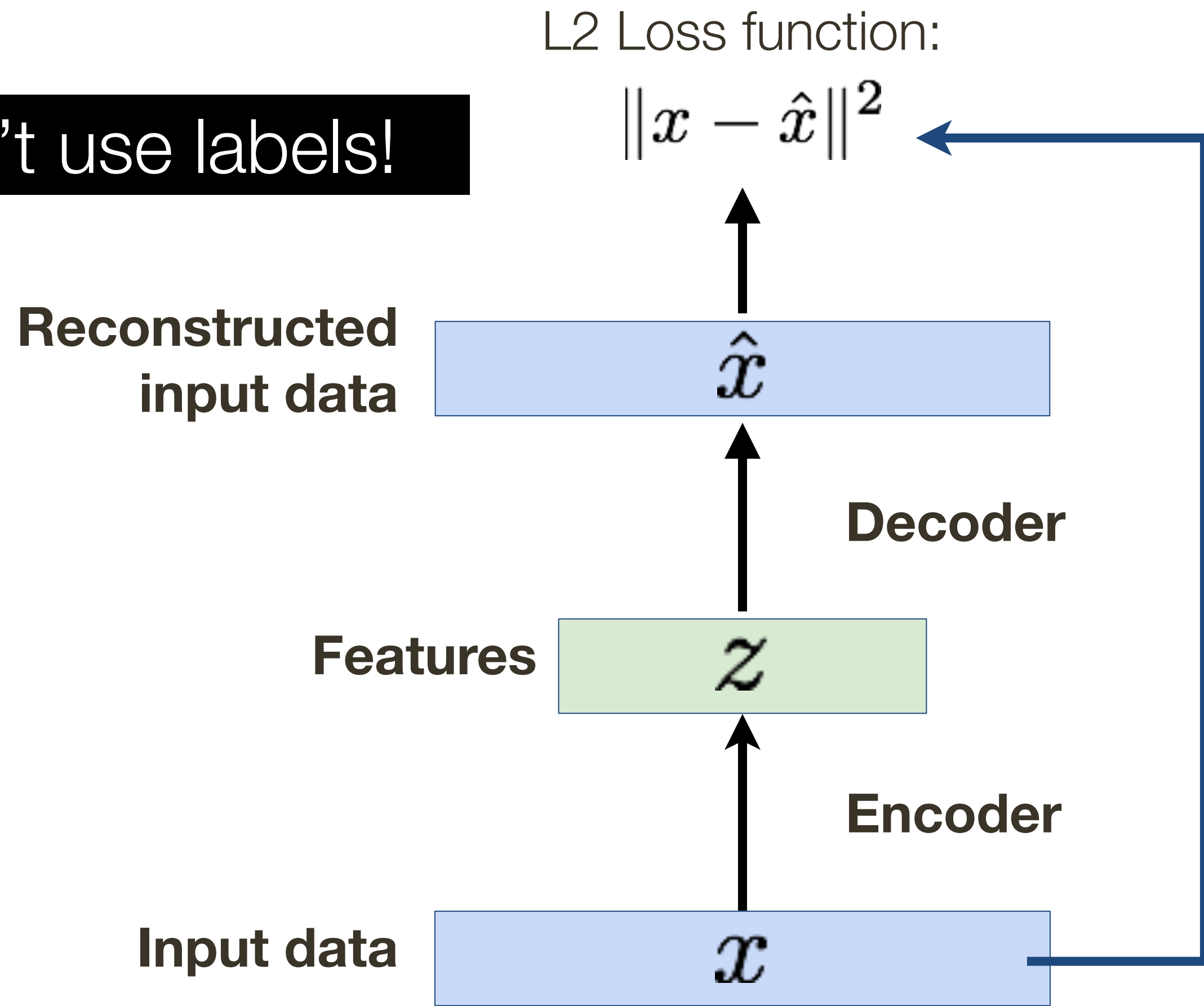


Autoencoders Reminder ...



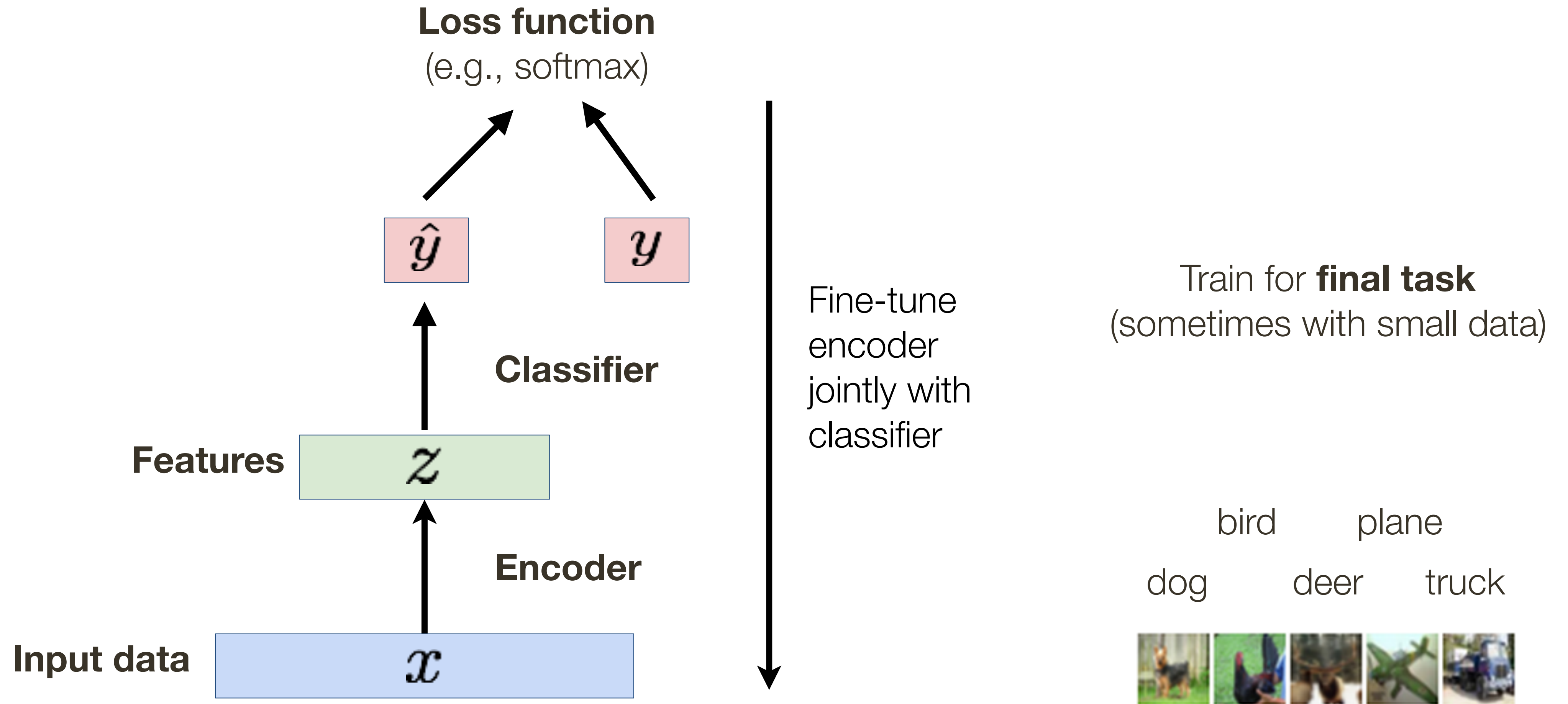
Autoencoders Reminder ...

Doesn't use labels!



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

Autoencoders Reminder ...



Variational Autoencoders

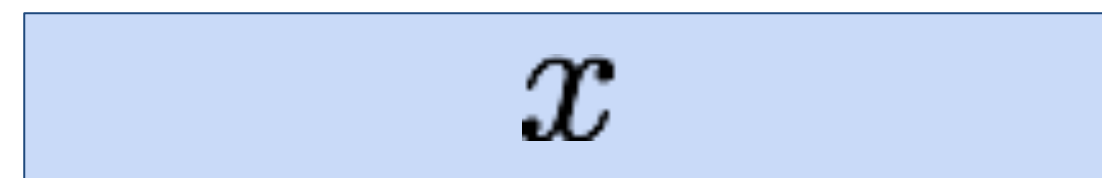
[Kingma and Welling, 2014]

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

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

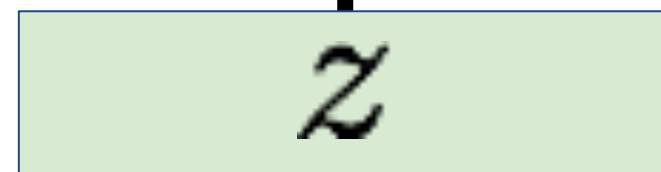
Sample from
true **conditional**

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



Sample from
true **prior**

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



Variational Autoencoders

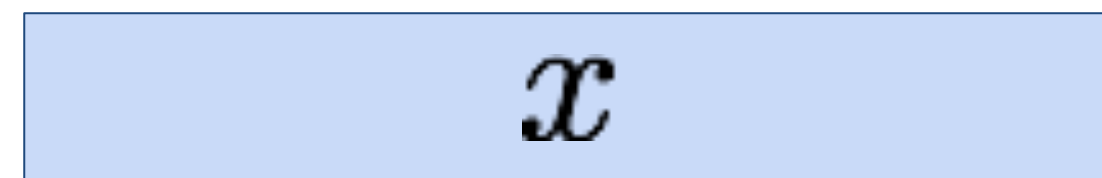
[Kingma and Welling, 2014]

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**

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



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.*)

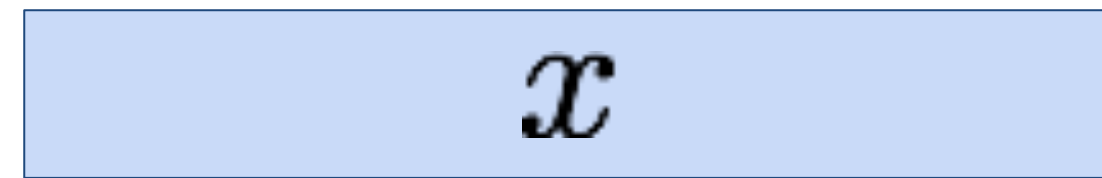
Variational Autoencoders

[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

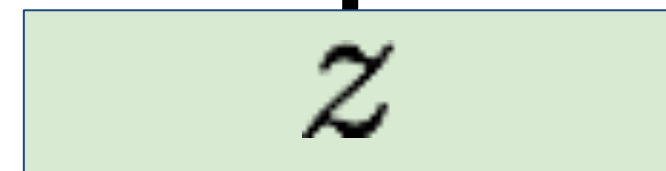
Sample from
true **conditional**

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



Sample from
true **prior**

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



Variational Autoencoders

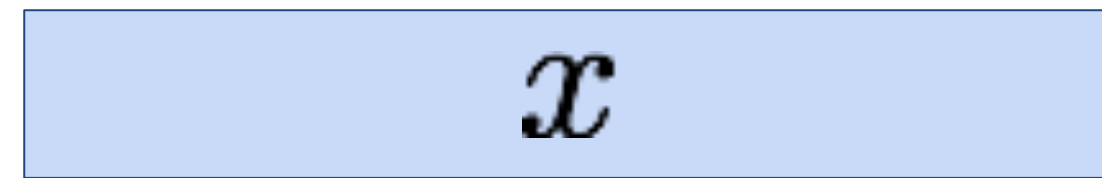
[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

How do we **represent** this model?

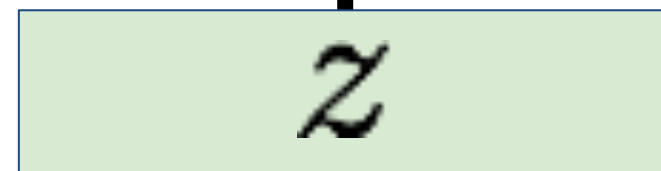
Sample from
true **conditional**

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



Sample from
true **prior**

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



Variational Autoencoders

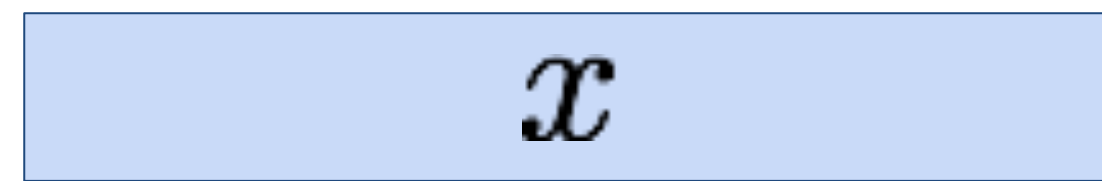
[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

How do we **represent** this model?

Sample from
true **conditional**

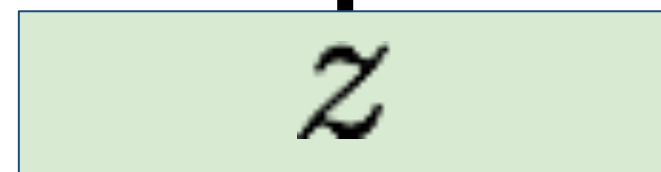
$$p_{\theta^*}(x \mid z^{(i)})$$



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

[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

How do we **represent** this model?

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

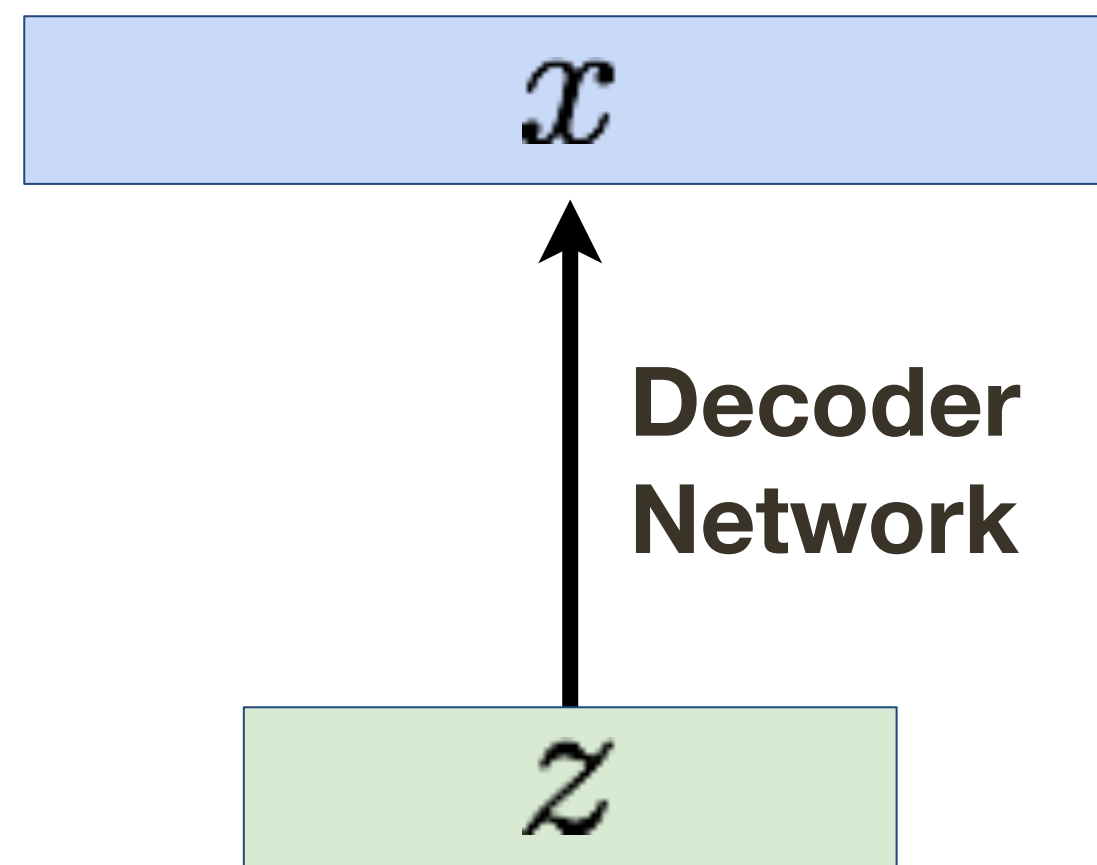
Conditional $p(\mathbf{x}|\mathbf{z})$ is complex (generates image)
Represent with Neural Network

Sample from
true **conditional**

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

Sample from
true **prior**

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



Variational Autoencoders

[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

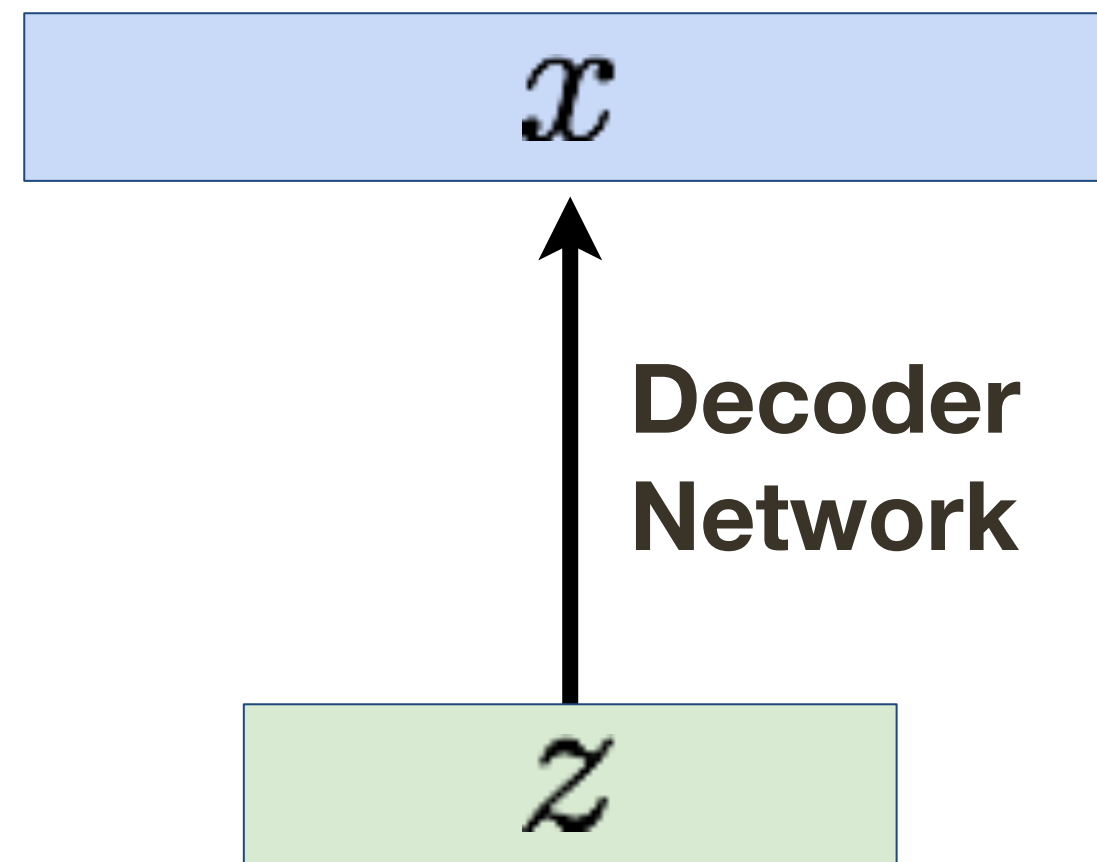
How do we **train** this model?

Sample from
true **conditional**

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

Sample from
true **prior**

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



Variational Autoencoders

[Kingma and Welling, 2014]

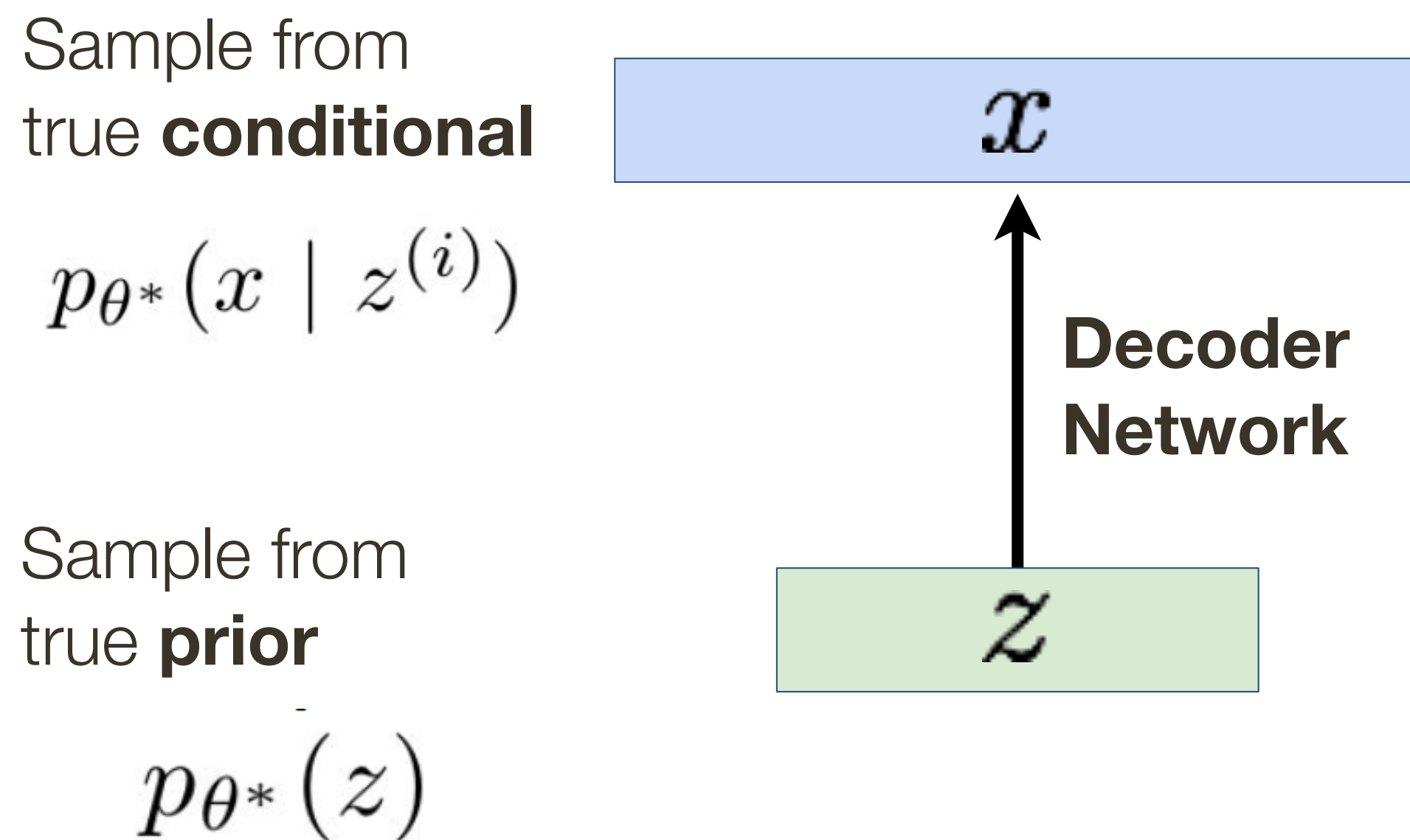
We want to **estimate the true parameters** θ^* of this generative model

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)



Variational Autoencoders

[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

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)

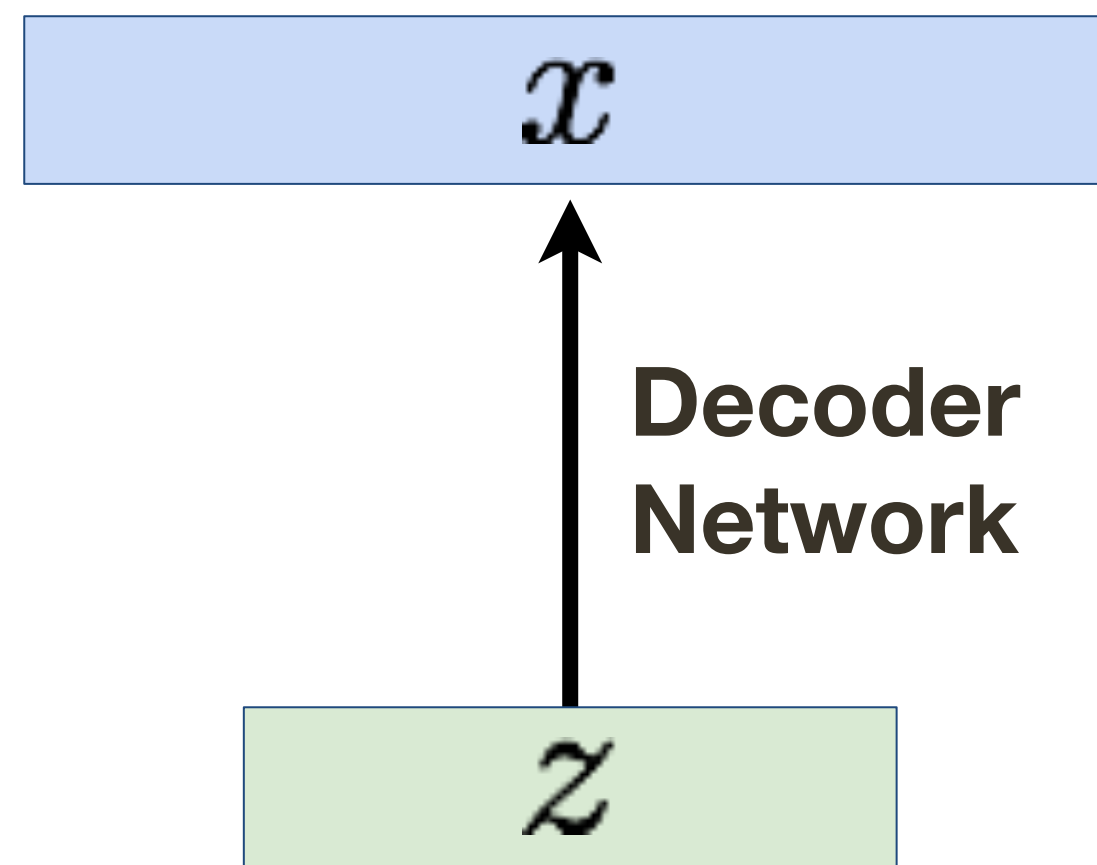
What is the problem with this?

Sample from
true **conditional**

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

Sample from
true **prior**

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



Variational Autoencoders

[Kingma and Welling, 2014]

We want to **estimate the true parameters** θ^* of this generative model

How do we **train** this model?

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

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

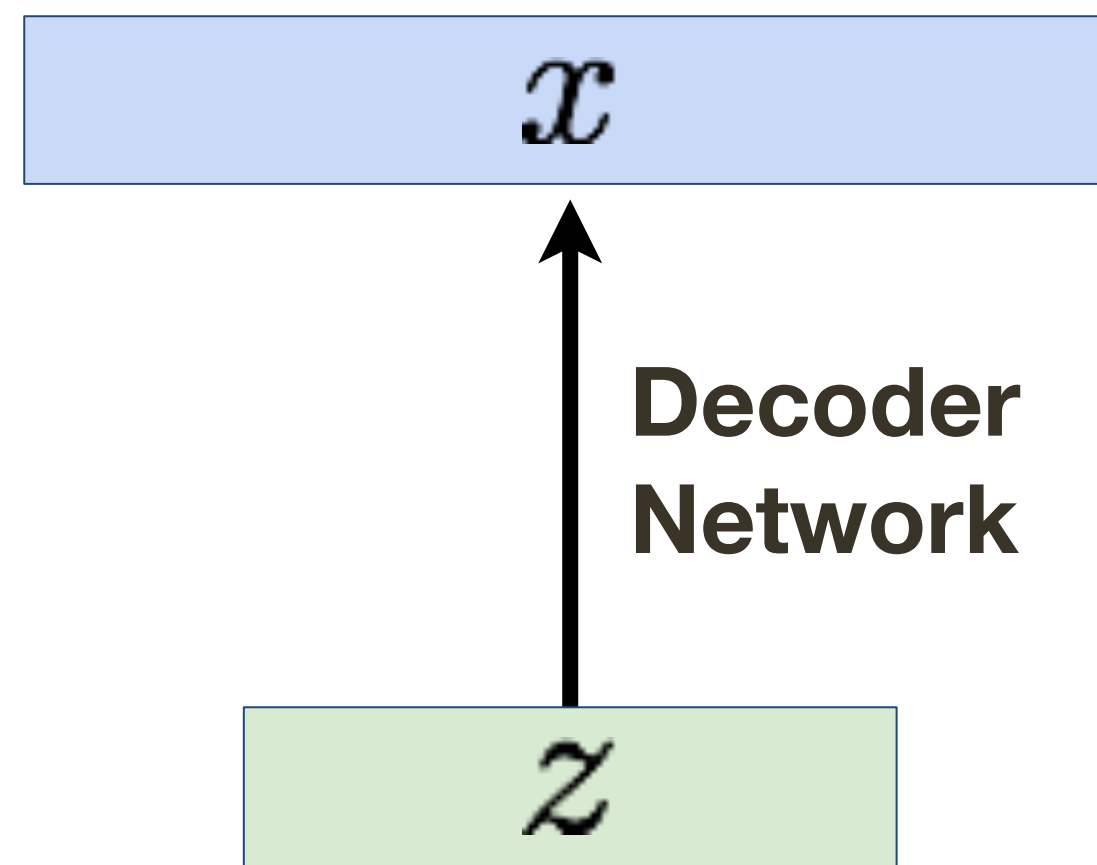
(now with latent z that we need to marginalize)

Sample from true **conditional**

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

Sample from true **prior**

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



Intractable !

Intractability in Variational Autoencoder

[Kingma and Welling, 2014]

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

Intractability in Variational Autoencoder

[Kingma and Welling, 2014]

Data **likelihood**:

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



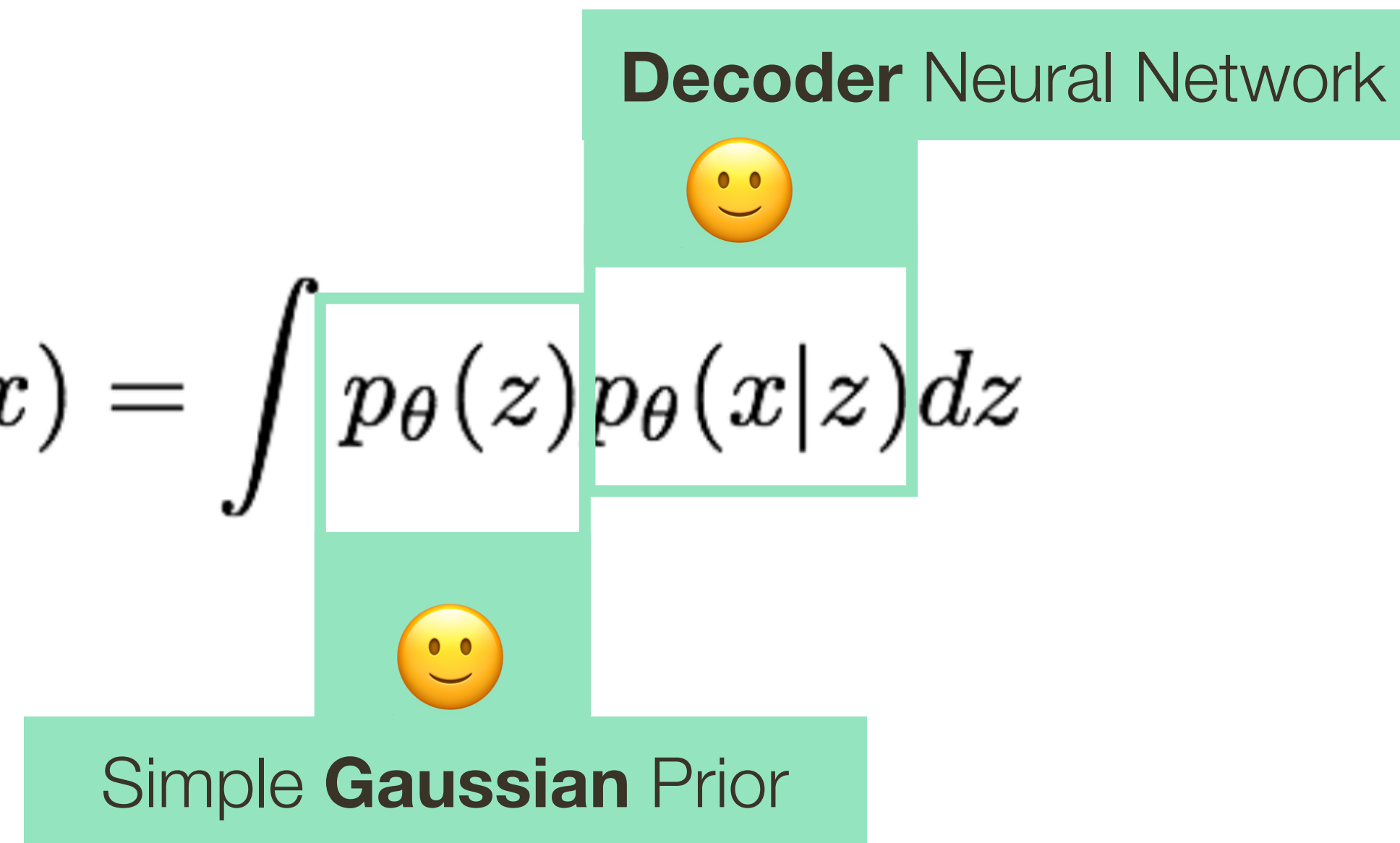
Simple **Gaussian** Prior

Intractability in Variational Autoencoder

[Kingma and Welling, 2014]

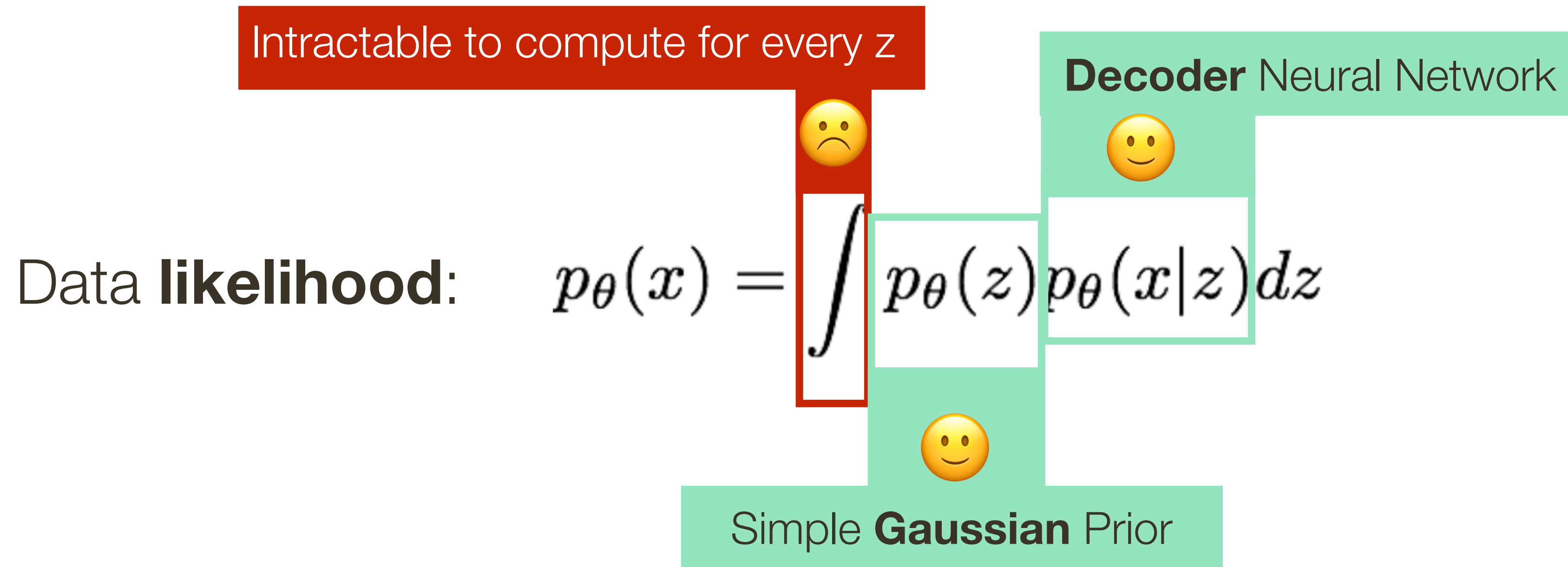
Data **likelihood**:

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



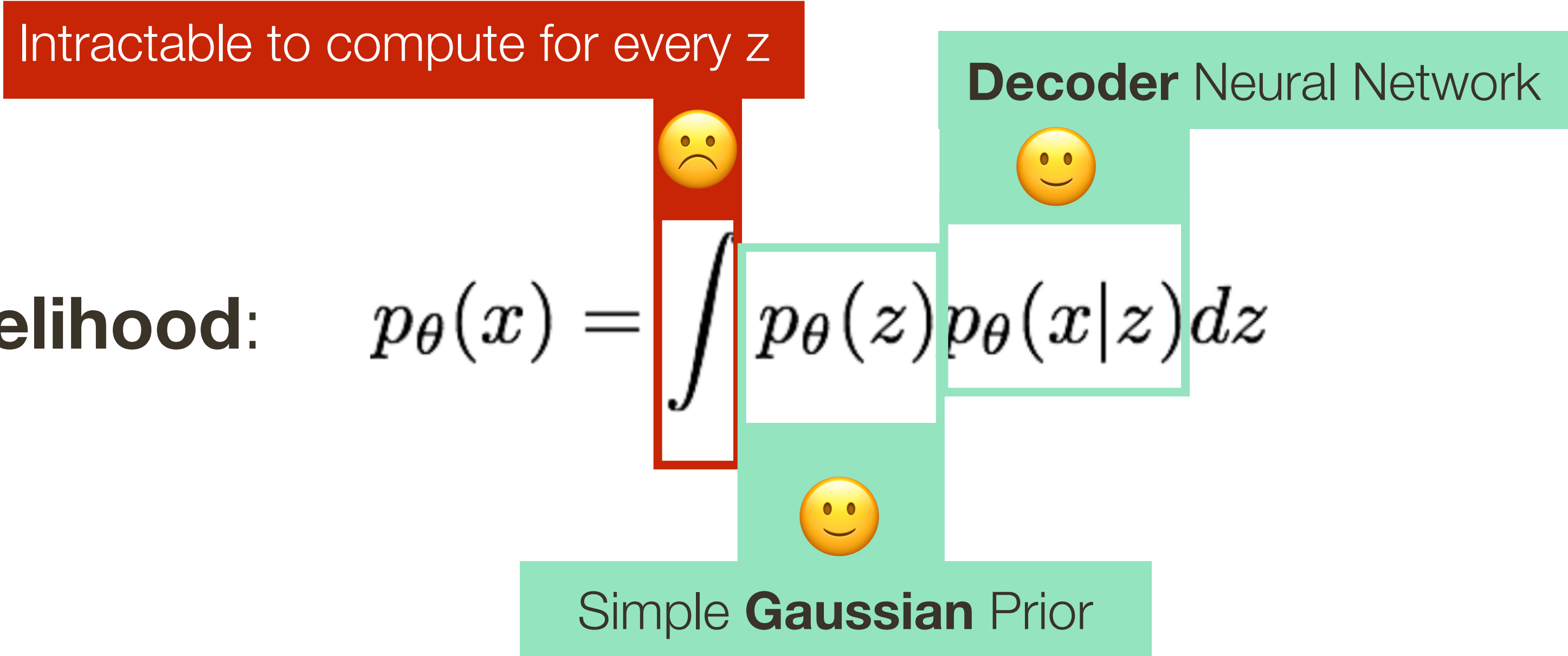
Intractability in Variational Autoencoder

[Kingma and Welling, 2014]



Intractability in Variational Autoencoder

[Kingma and Welling, 2014]



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

Intractability in Variational Autoencoder

[Kingma and Welling, 2014]

Intractable to compute for every z

Data **likelihood**:

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

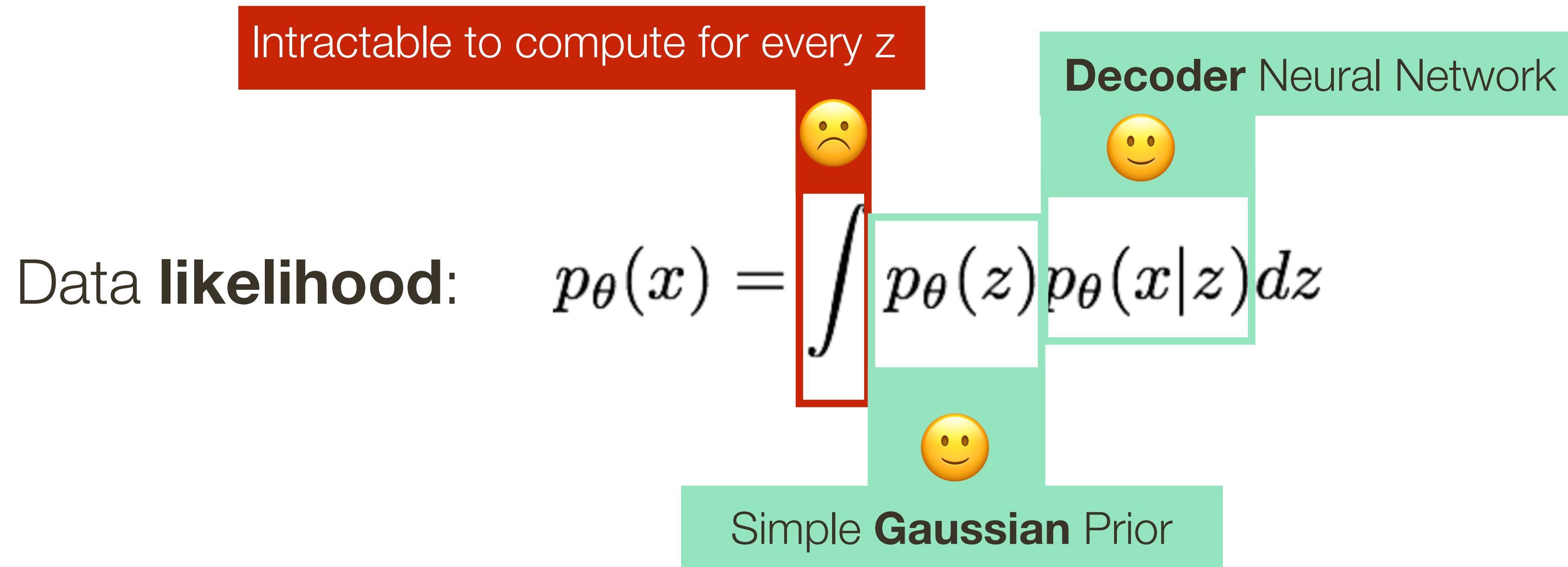
Decoder Neural Network

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

[Kingma and Welling, 2014]



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

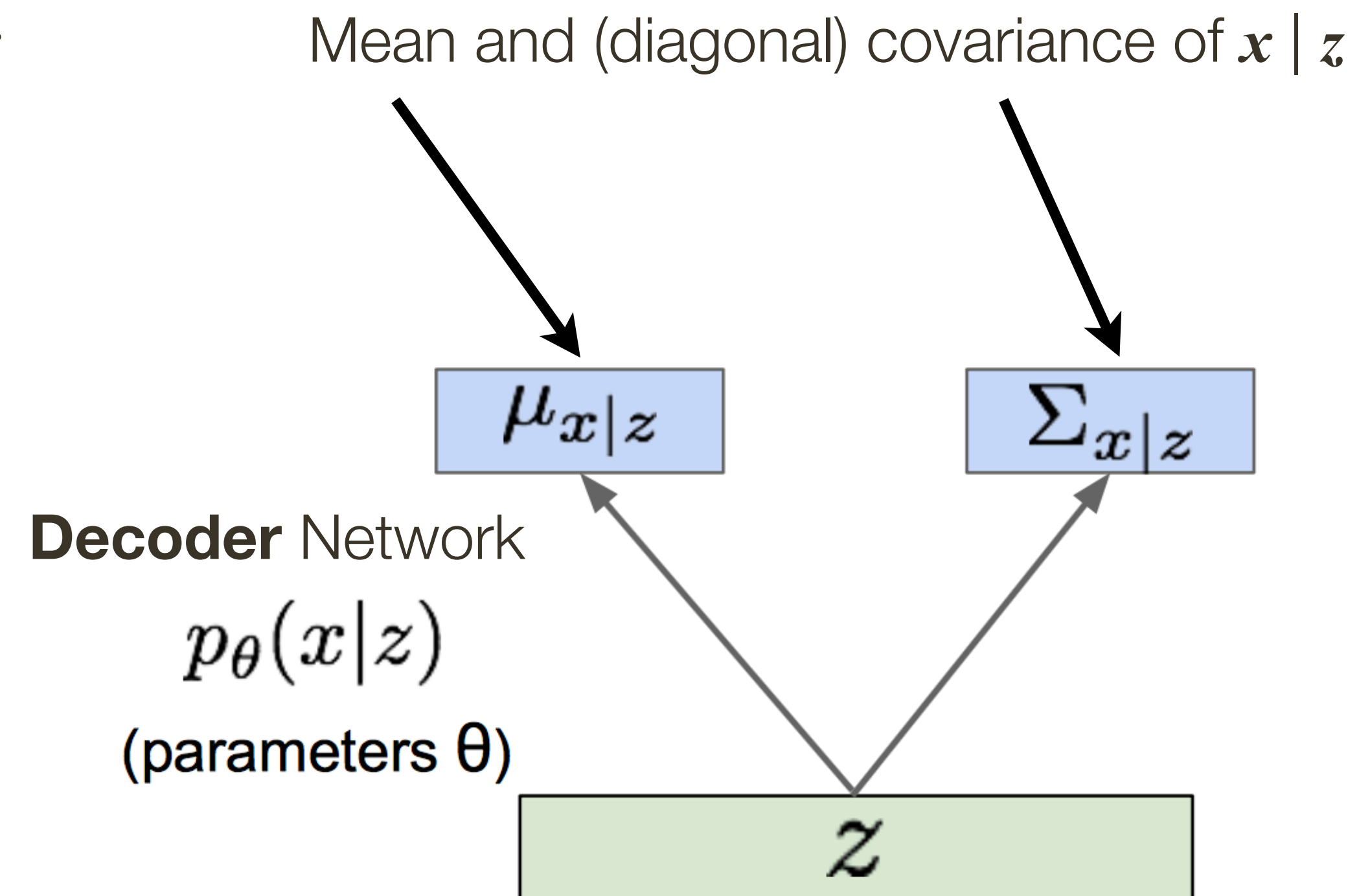
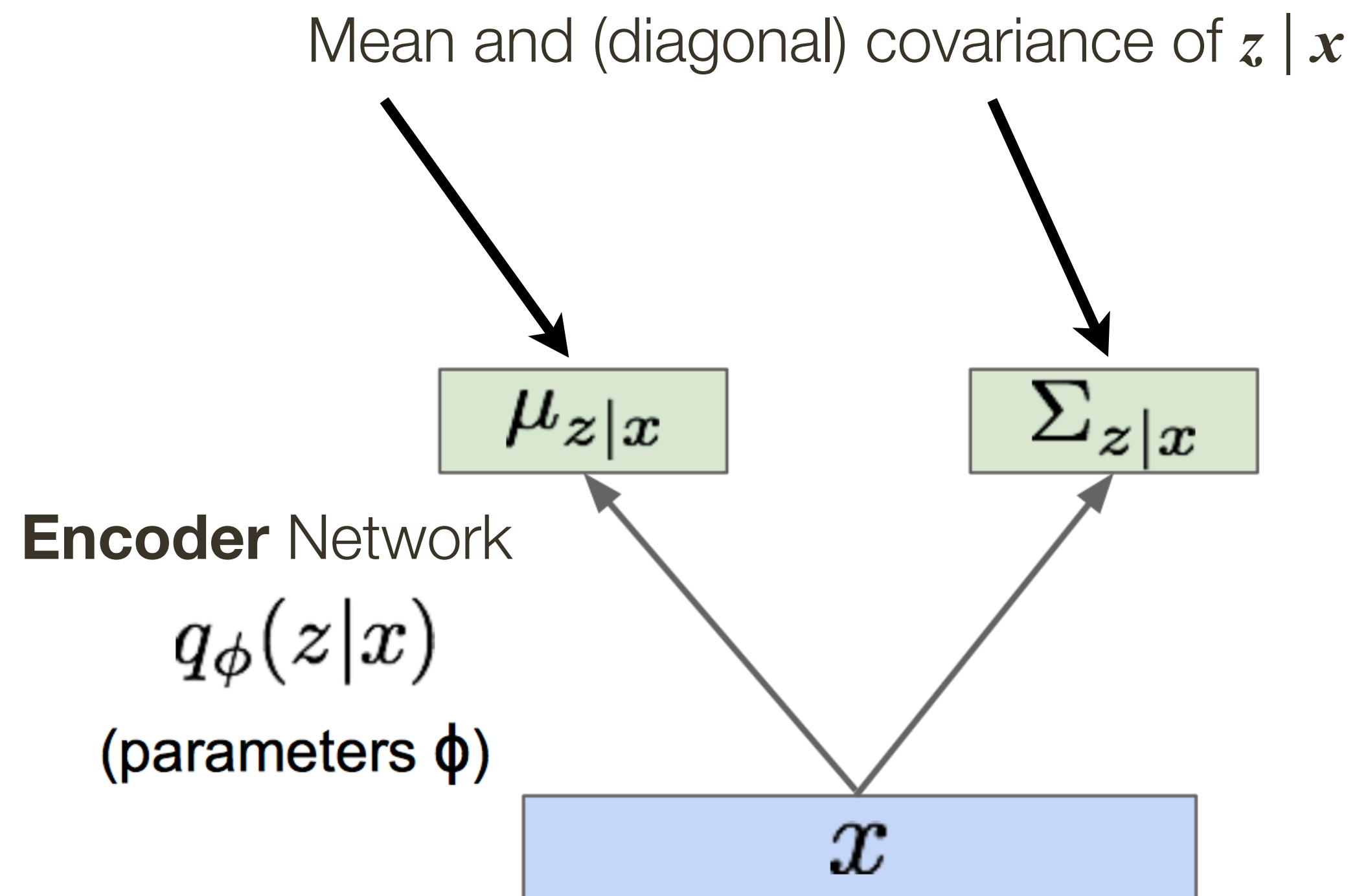
Solution: In addition to decoder network modeling $p_{\theta}(x|z)$, define additional encoder network $q_{\phi}(z|x)$ that approximates $p_{\theta}(z|x)$

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

Variational Autoencoder

[Kingma and Welling, 2014]

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

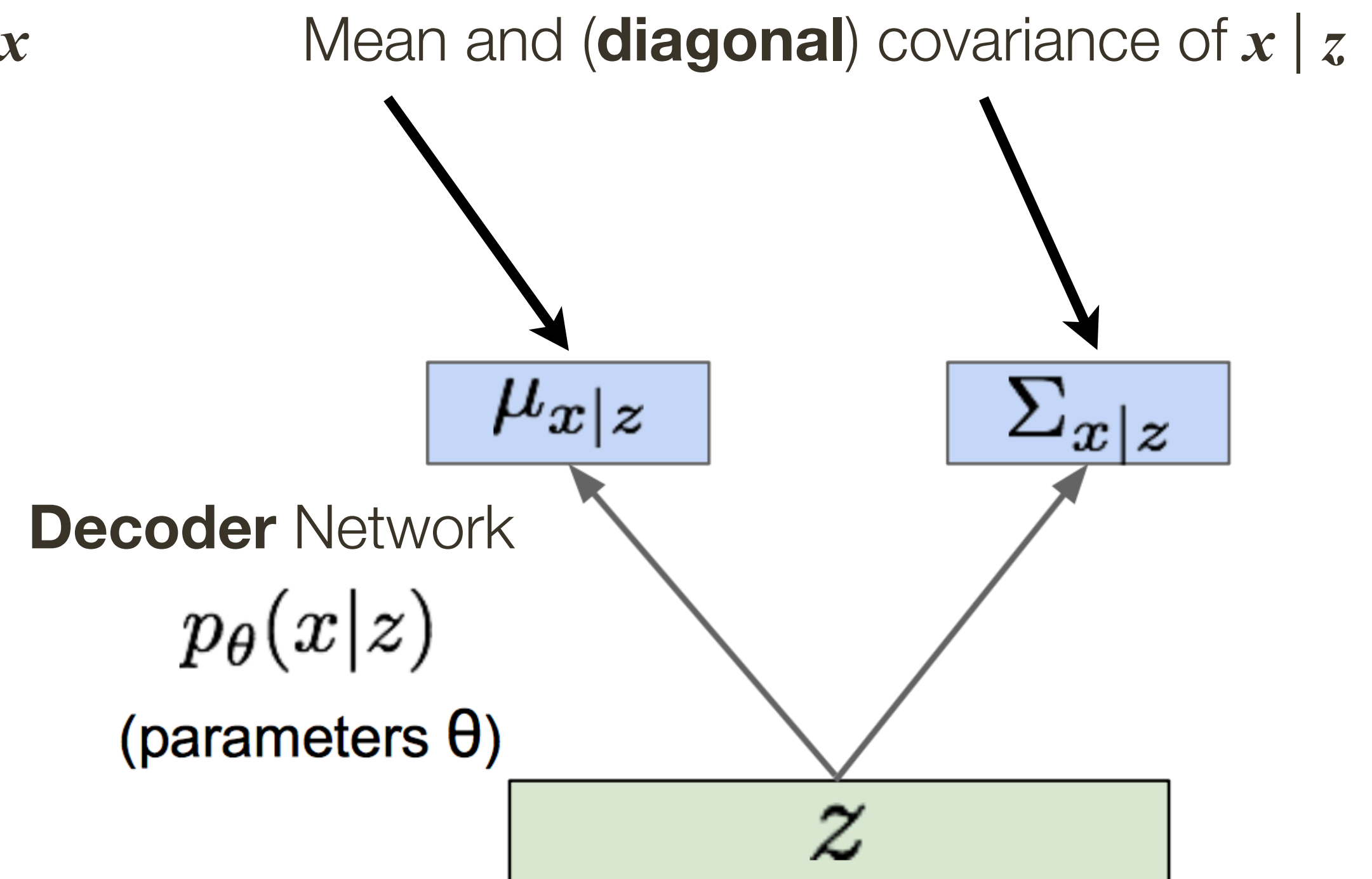
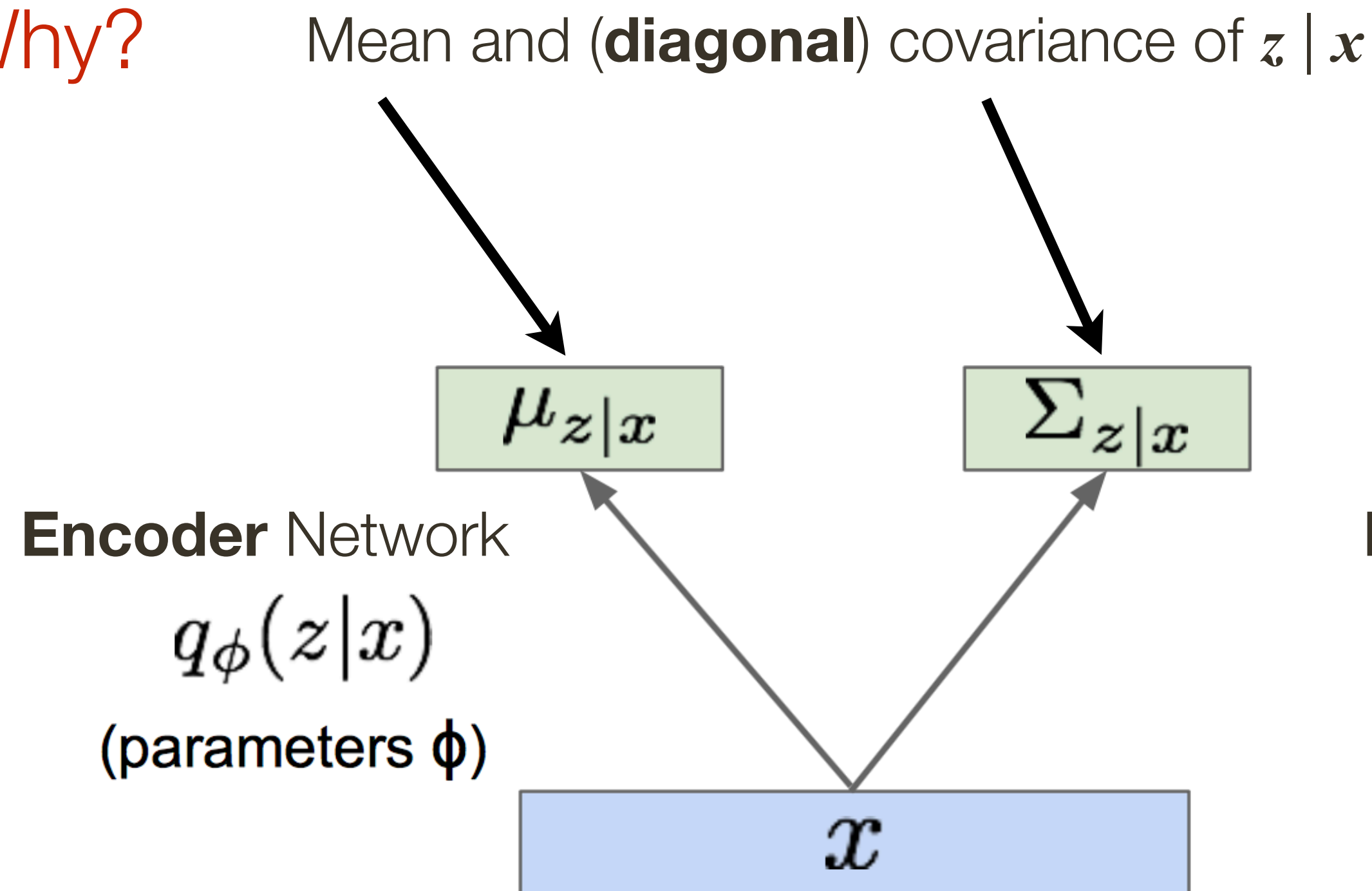


Variational Autoencoder

[Kingma and Welling, 2014]

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

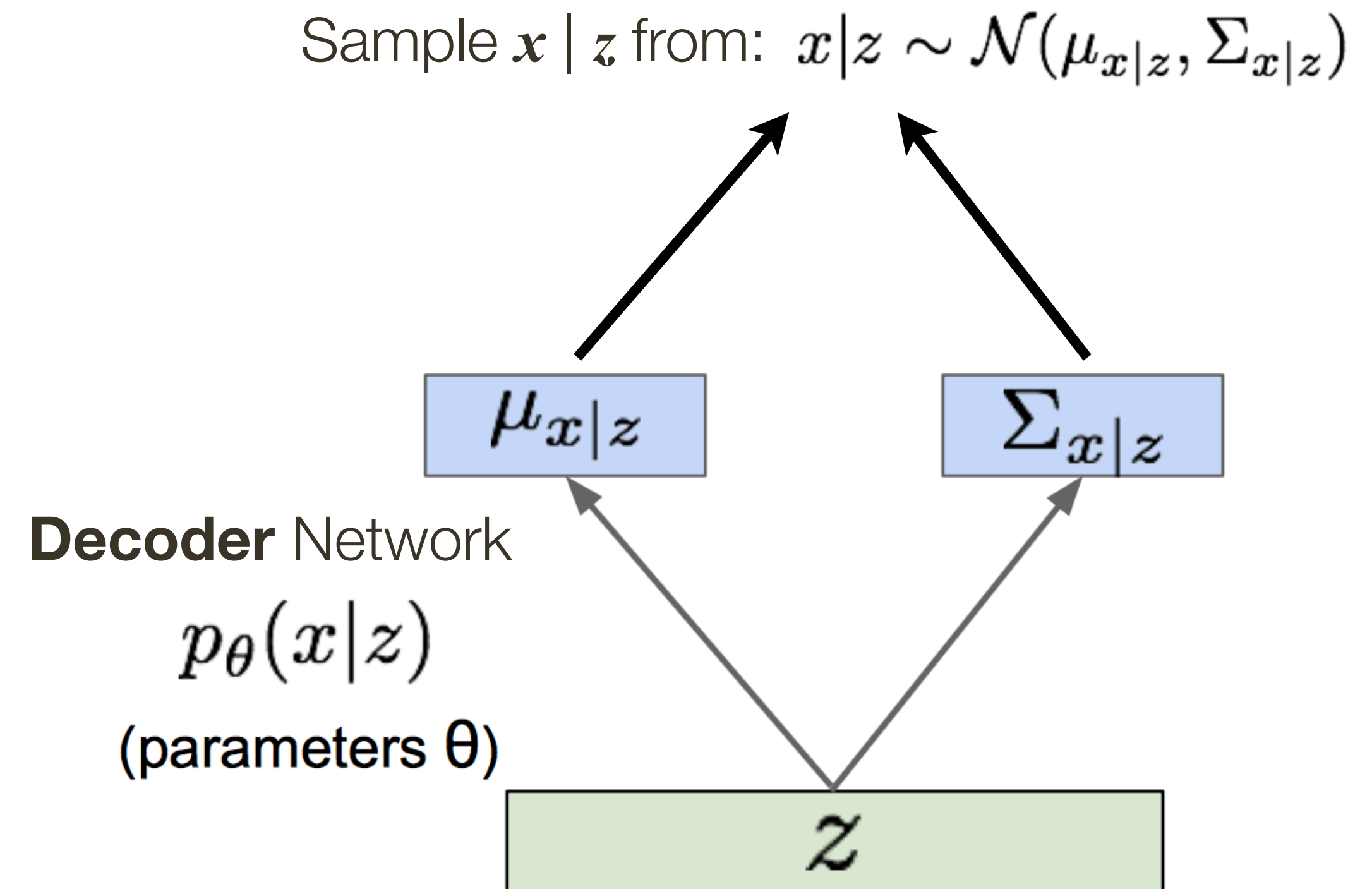
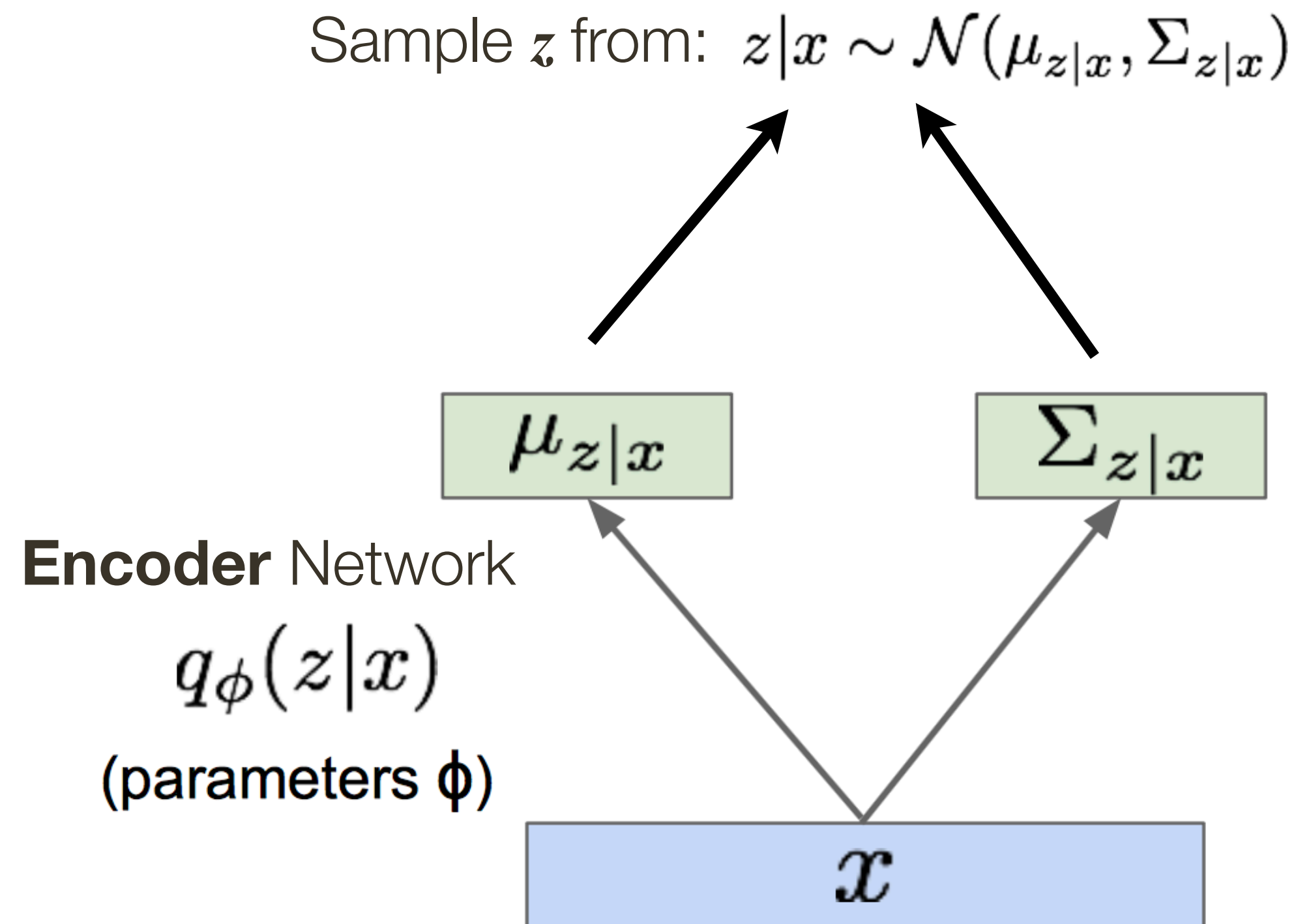
Why?



Variational Autoencoder

[Kingma and Welling, 2014]

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



Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

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

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)})) \text{ Does not depend on } z \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z) q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)}) q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z \left[\log p_{\theta}(x^{(i)} | z) \right] - \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z \left[\log p_{\theta}(x^{(i)} | z) \right] - \underbrace{D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z))}_{\text{Expectation with respect to } z} + \underbrace{D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))}_{\text{(using encoder network) leads to nice KL terms}}\end{aligned}$$

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

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

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

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

Tractable lower bound which we can take gradient of and optimize! ($p_{\theta}(x|z)$ differentiable, KL term differentiable)

Variational Autoencoder

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

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

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

[Kingma and Welling, 2014]

Derivation of lower bound of the data likelihood

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

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

**Reconstruct
Input Data**

**Make approximate posterior
close to the prior**

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

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

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

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

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

Input Data

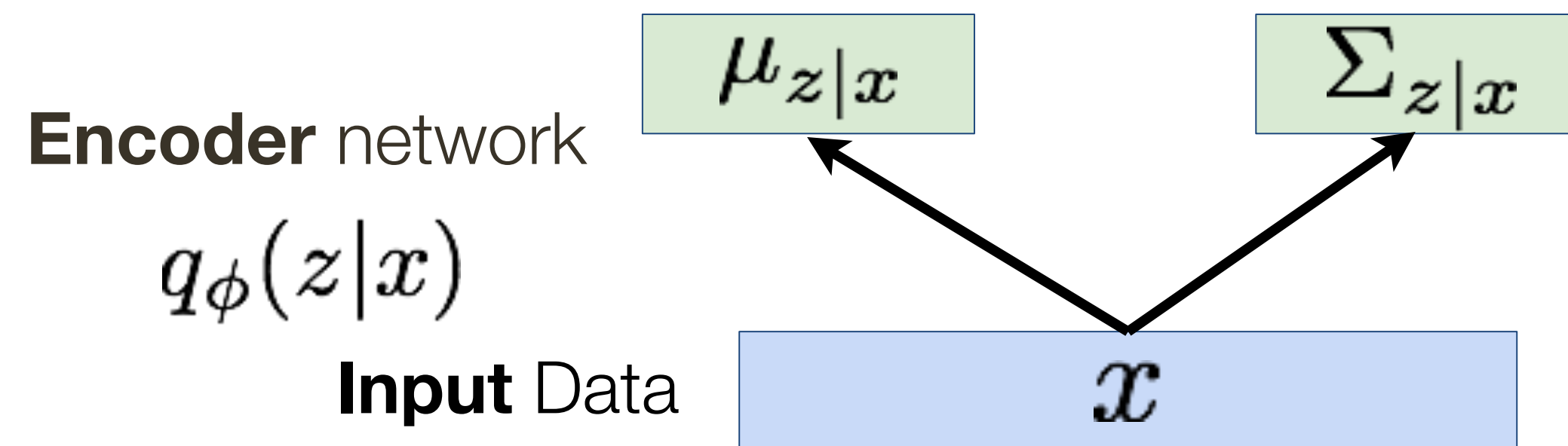
\mathcal{X}

Variational Autoencoder: Learning

Putting it all together:

maximizing the likelihood lower bound

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



Variational Autoencoder: Learning

Putting it all together:

maximizing the likelihood lower bound

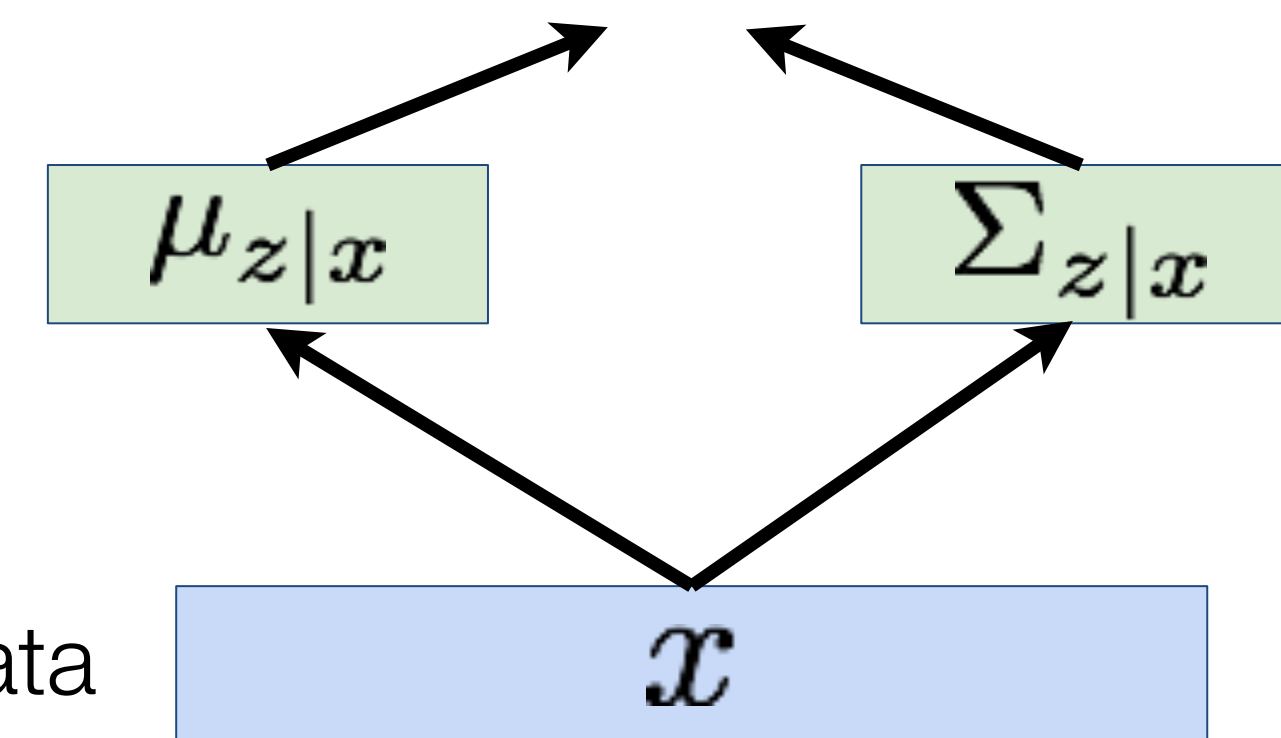
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Encoder network

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

Input Data



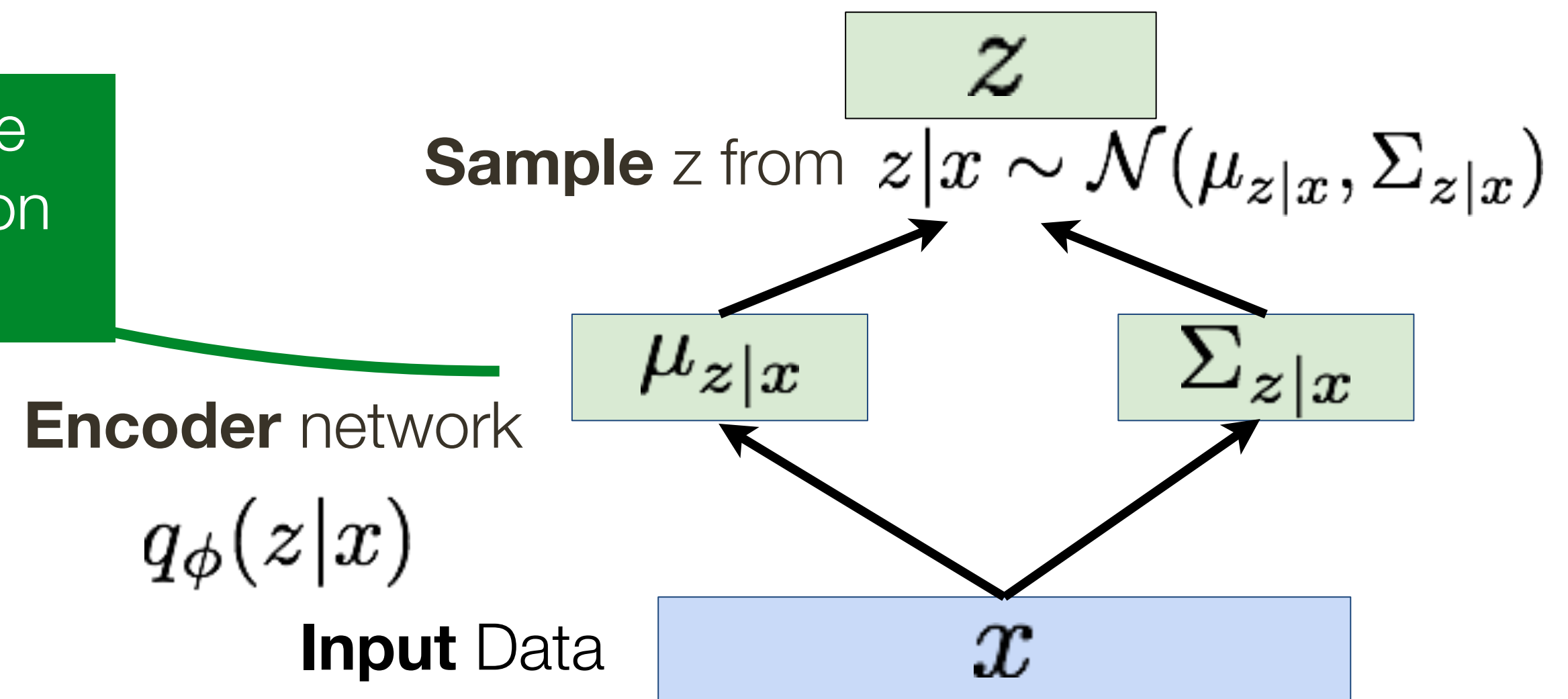
Variational Autoencoder: Learning

Putting it all together:

maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior



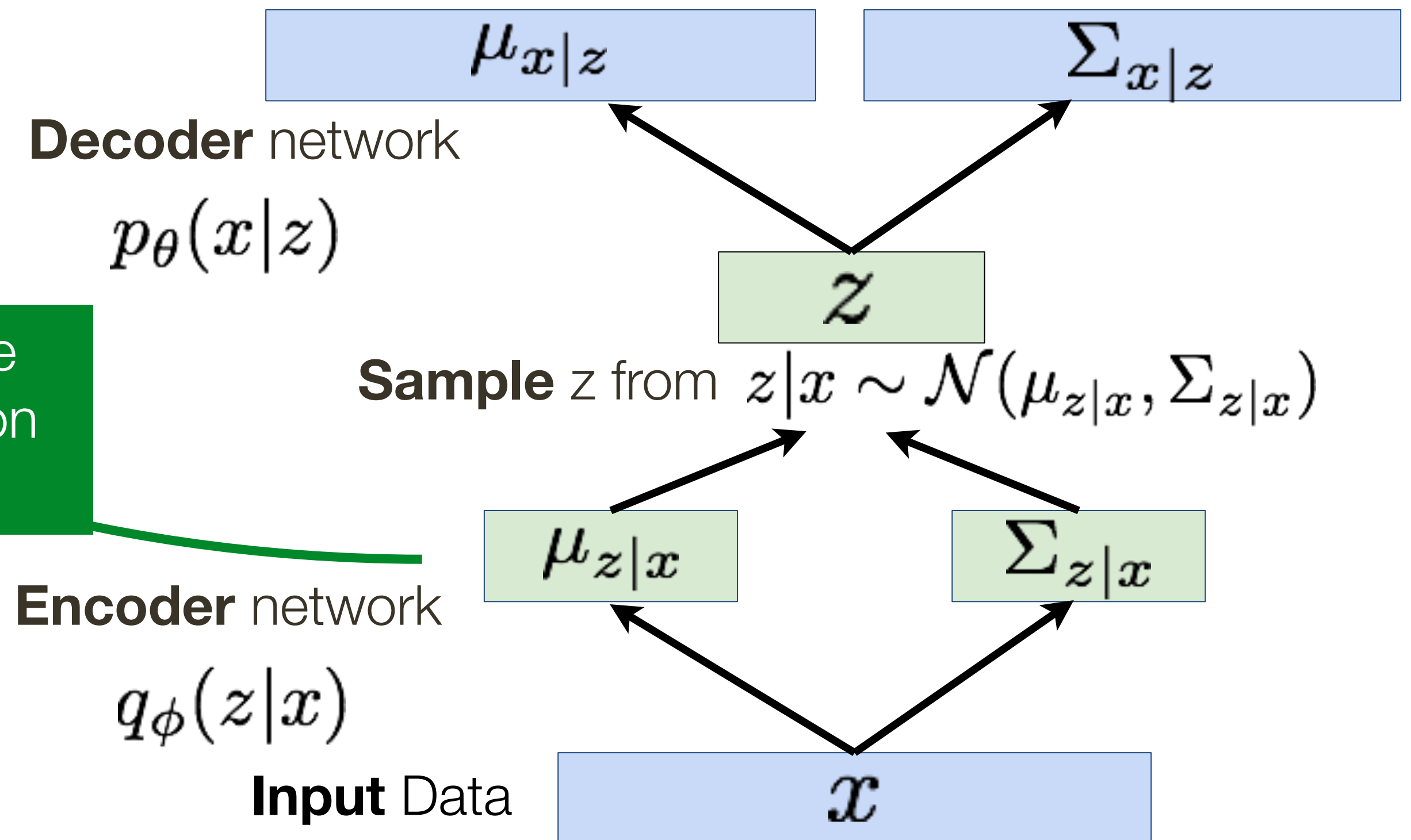
Variational Autoencoder: Learning

Putting it all together:

maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior



Variational Autoencoder: Learning

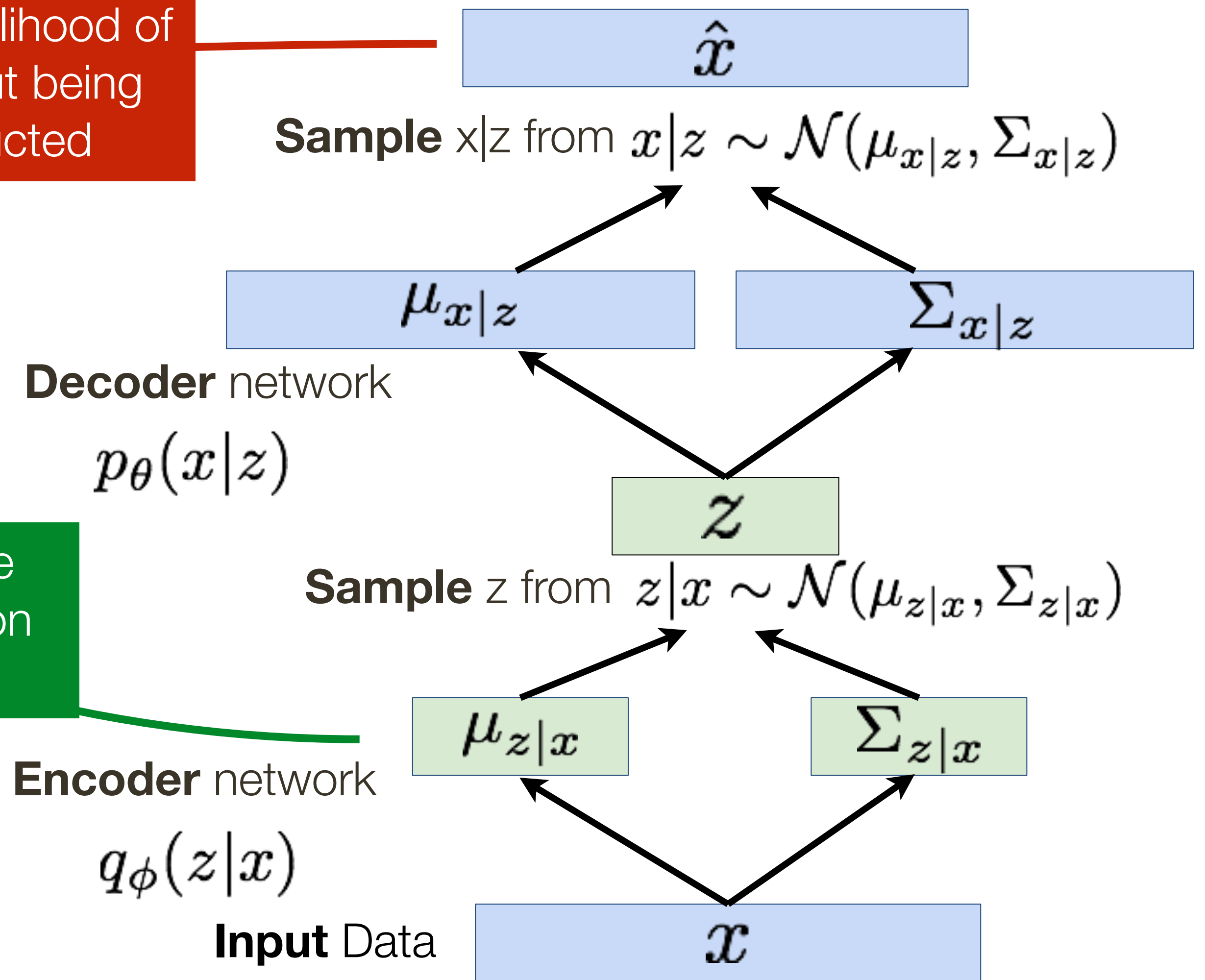
Putting it all together:

maximizing the likelihood lower bound

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

Maximize likelihood of original input being reconstructed

Make approximate posterior distribution close to prior



Variational Autoencoder: Learning

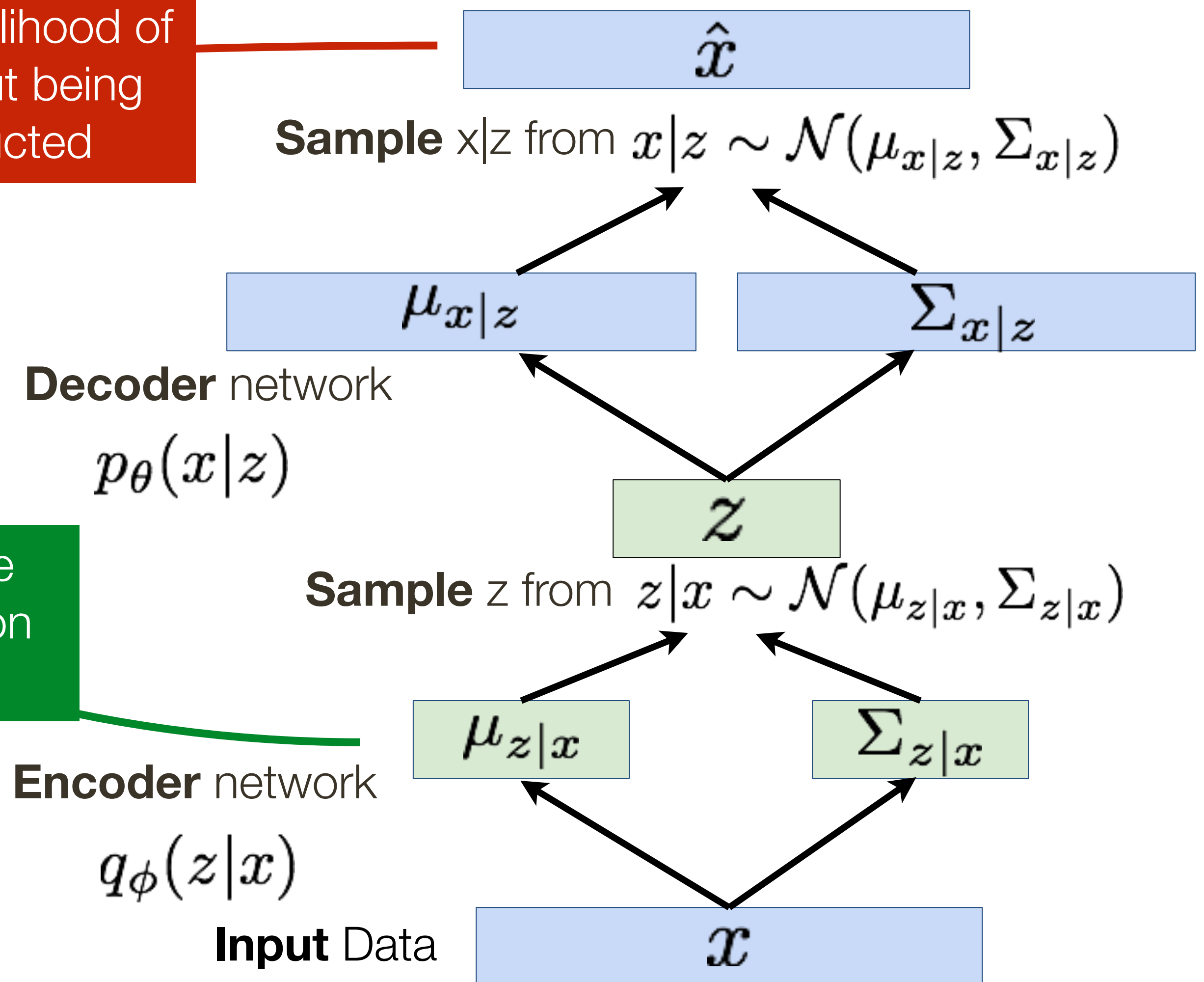
Putting it all together:

maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

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

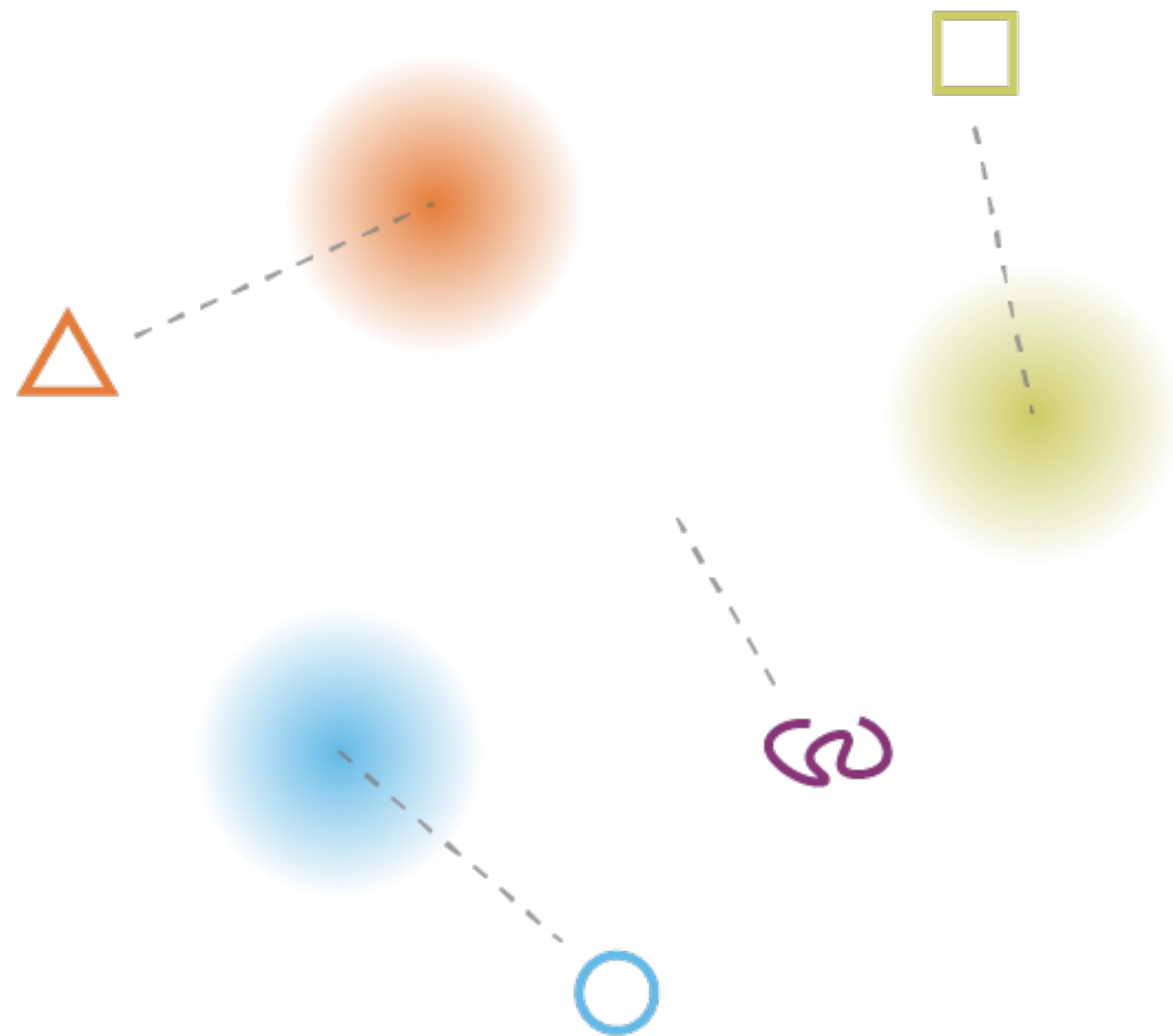
Make approximate posterior distribution close to prior



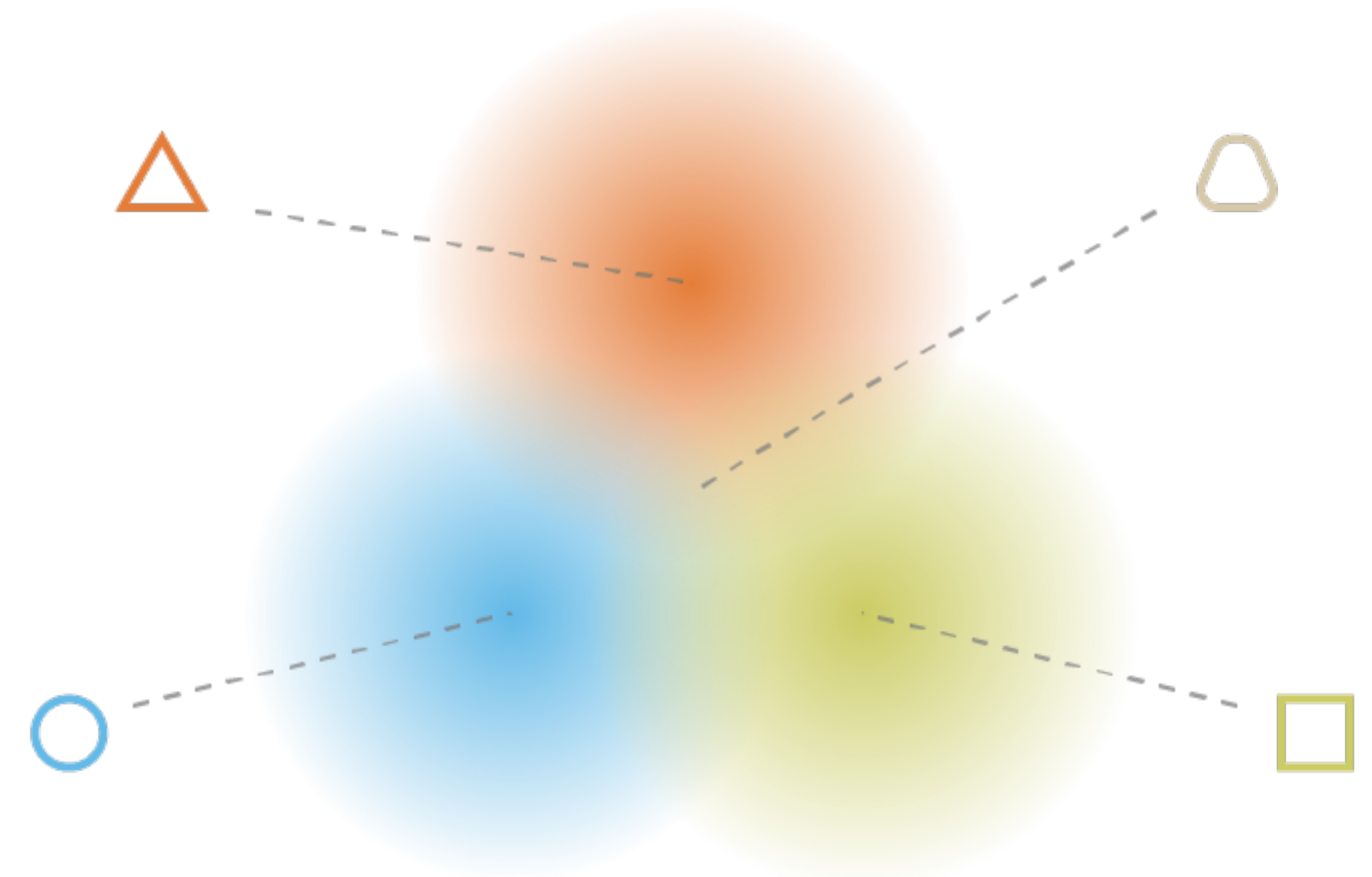
For every minibatch of input data: compute this forward pass, and then backprop!

Variational Autoencoder: Learning

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



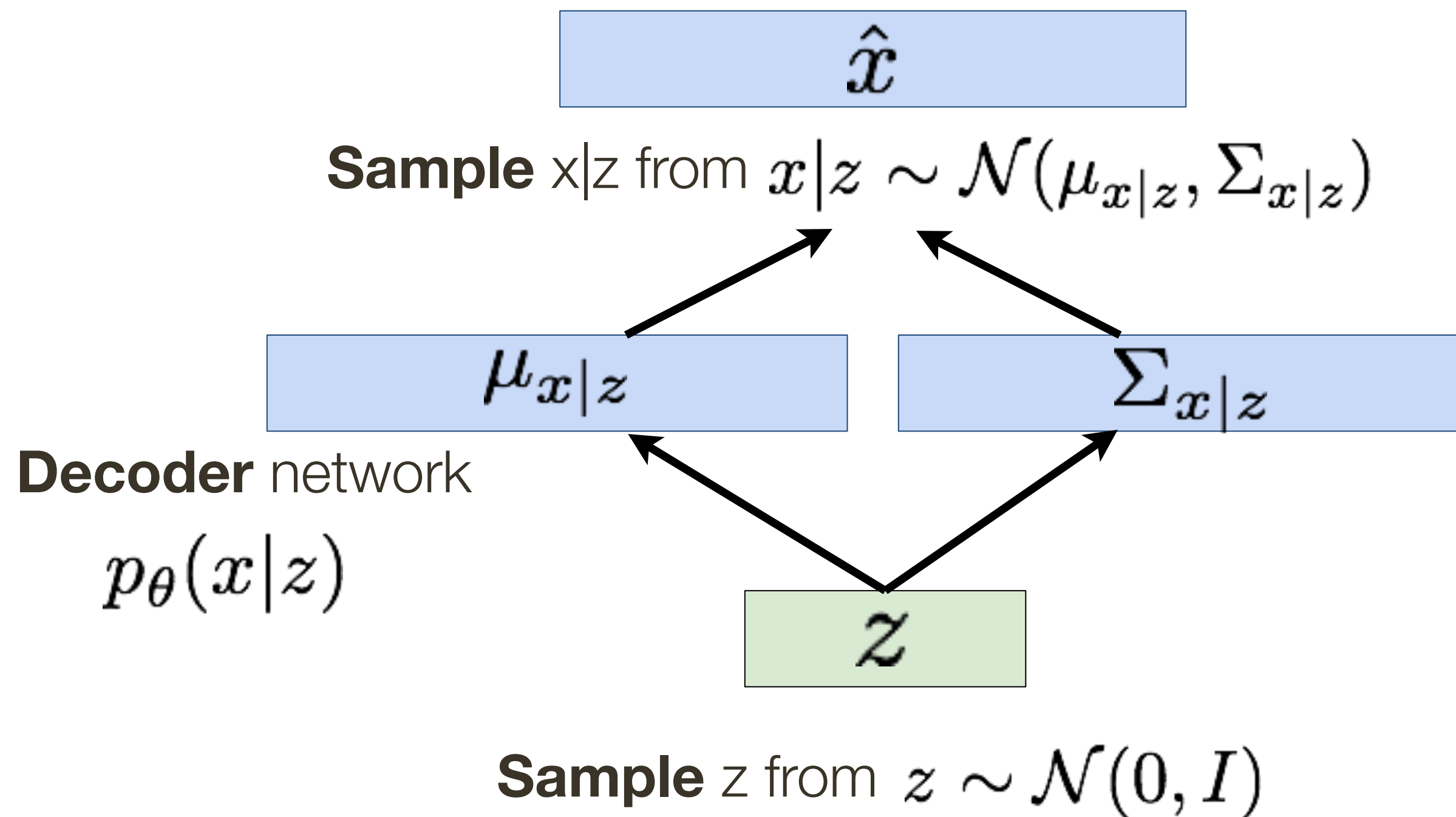
what can happen without regularisation



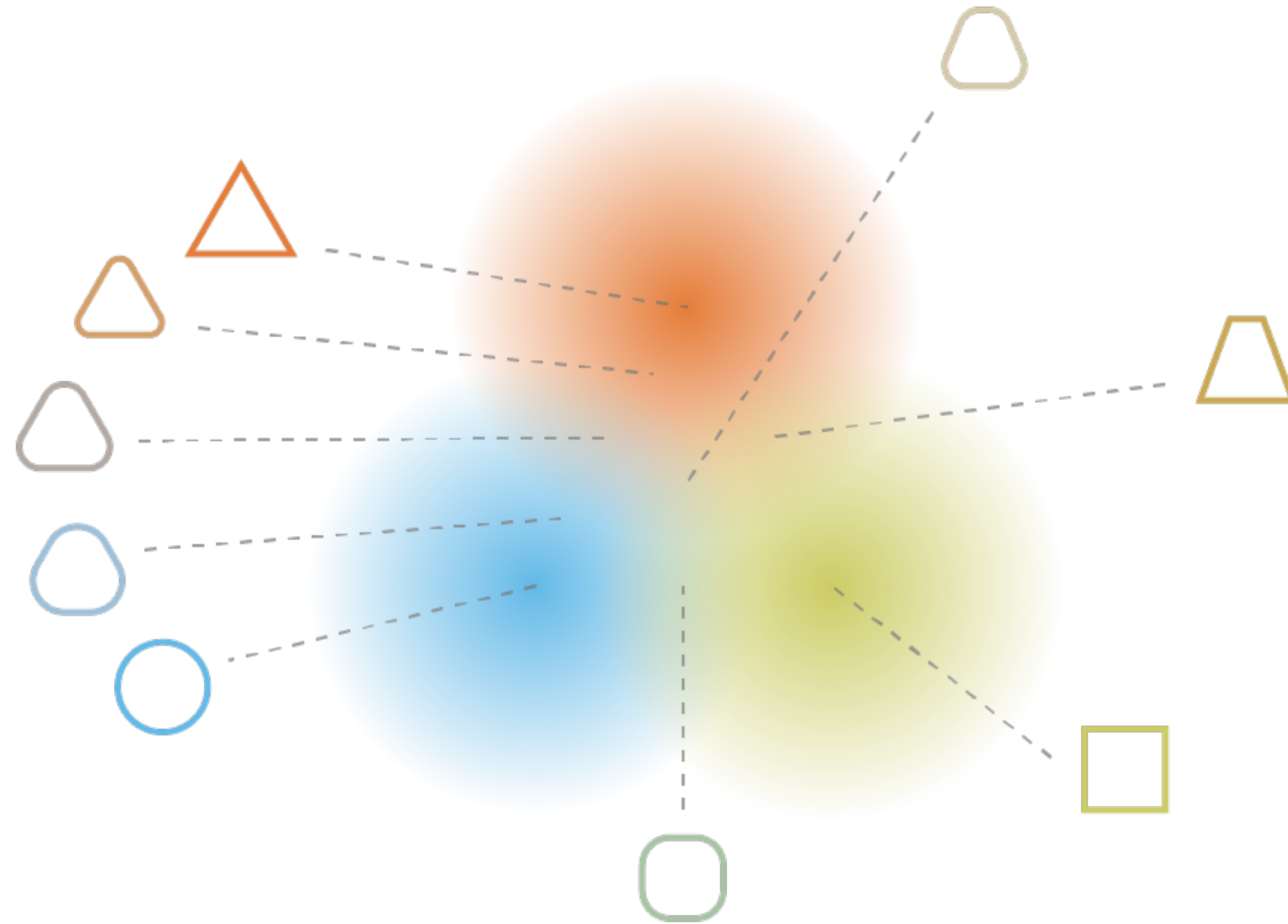
what we want to obtain with regularisation

Variational Autoencoder: Generating Data

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

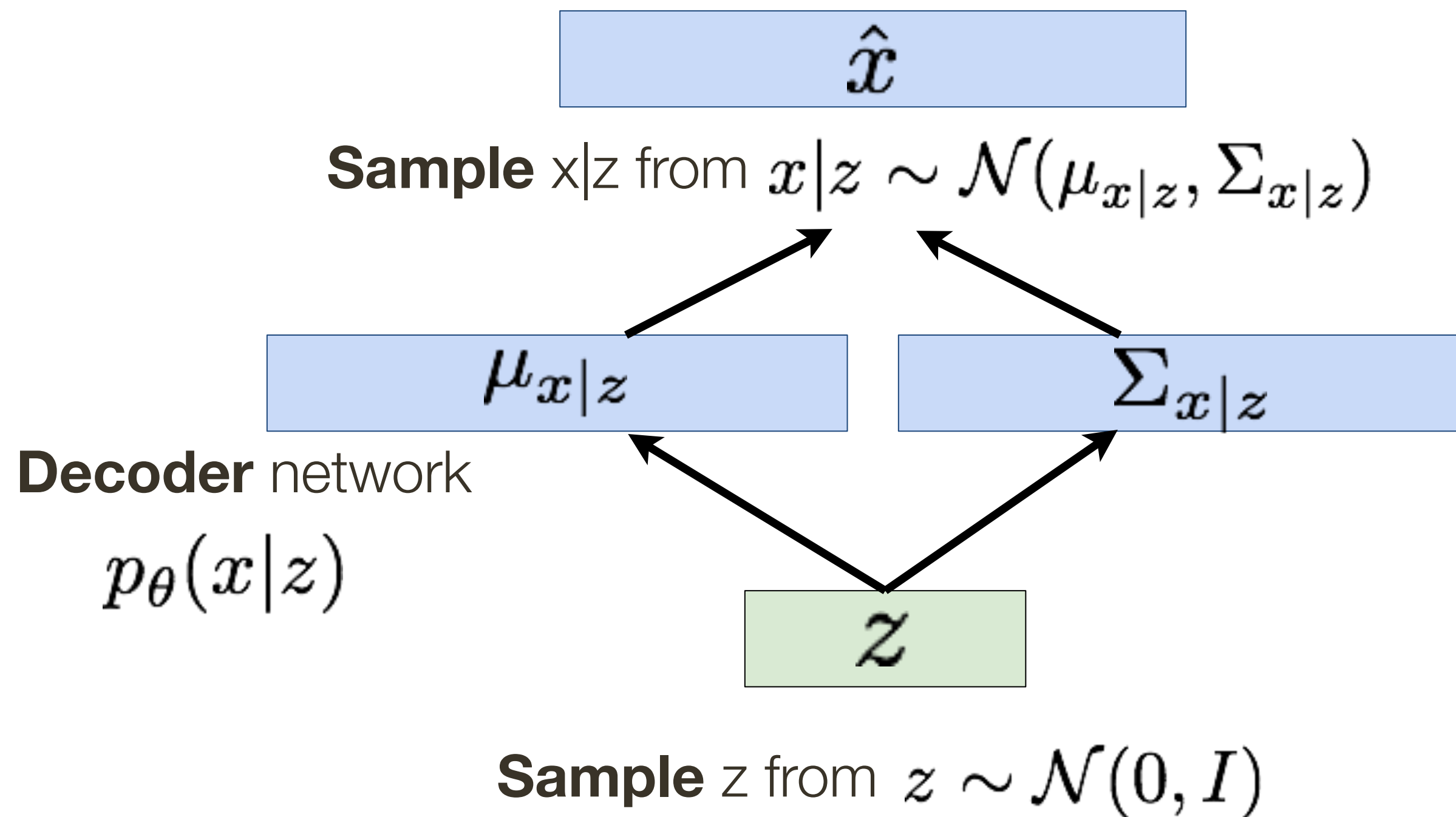


Variational Autoencoder: Generating Data

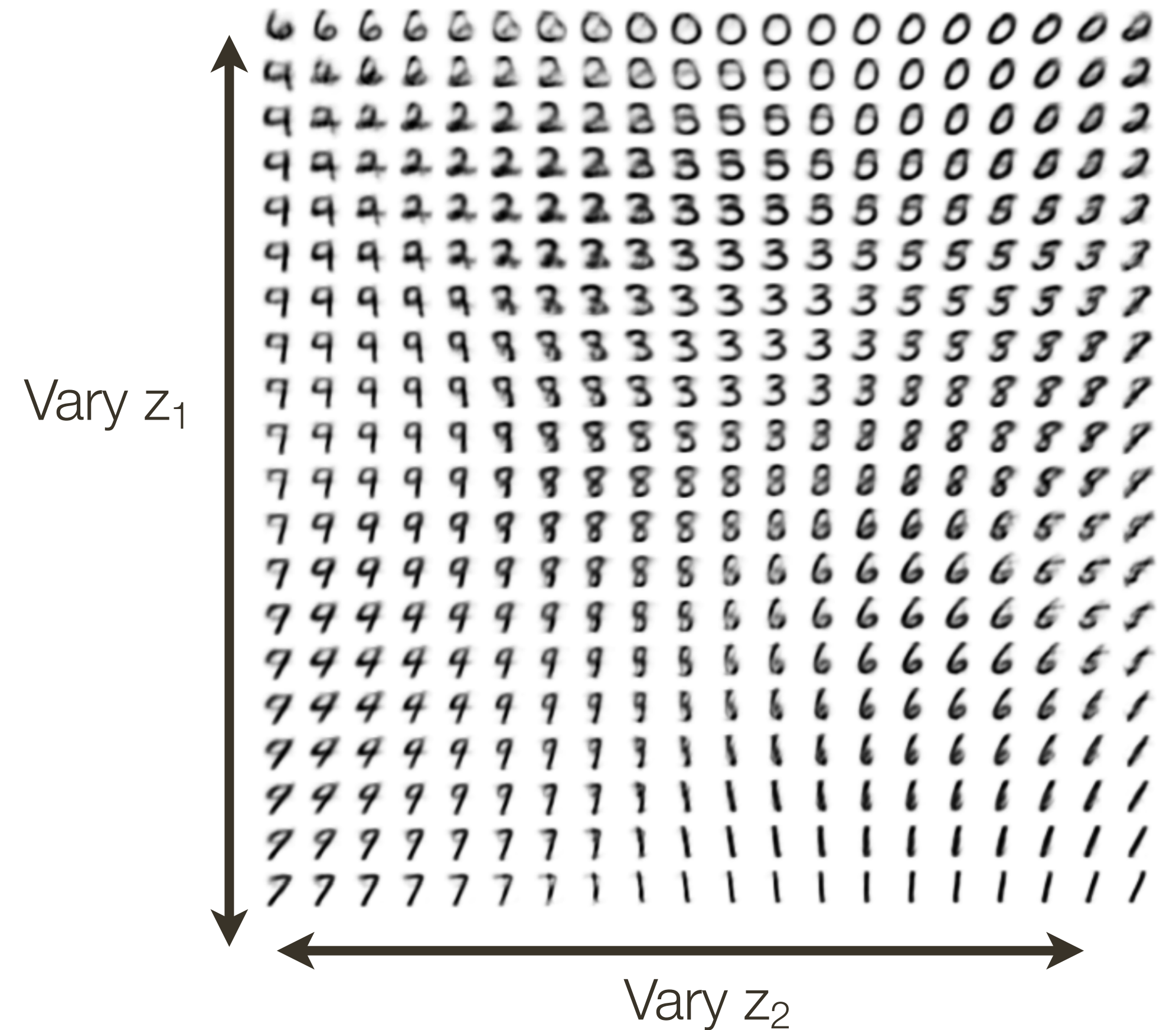


Variational Autoencoder: Generating Data

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



Data manifold for 2-d z

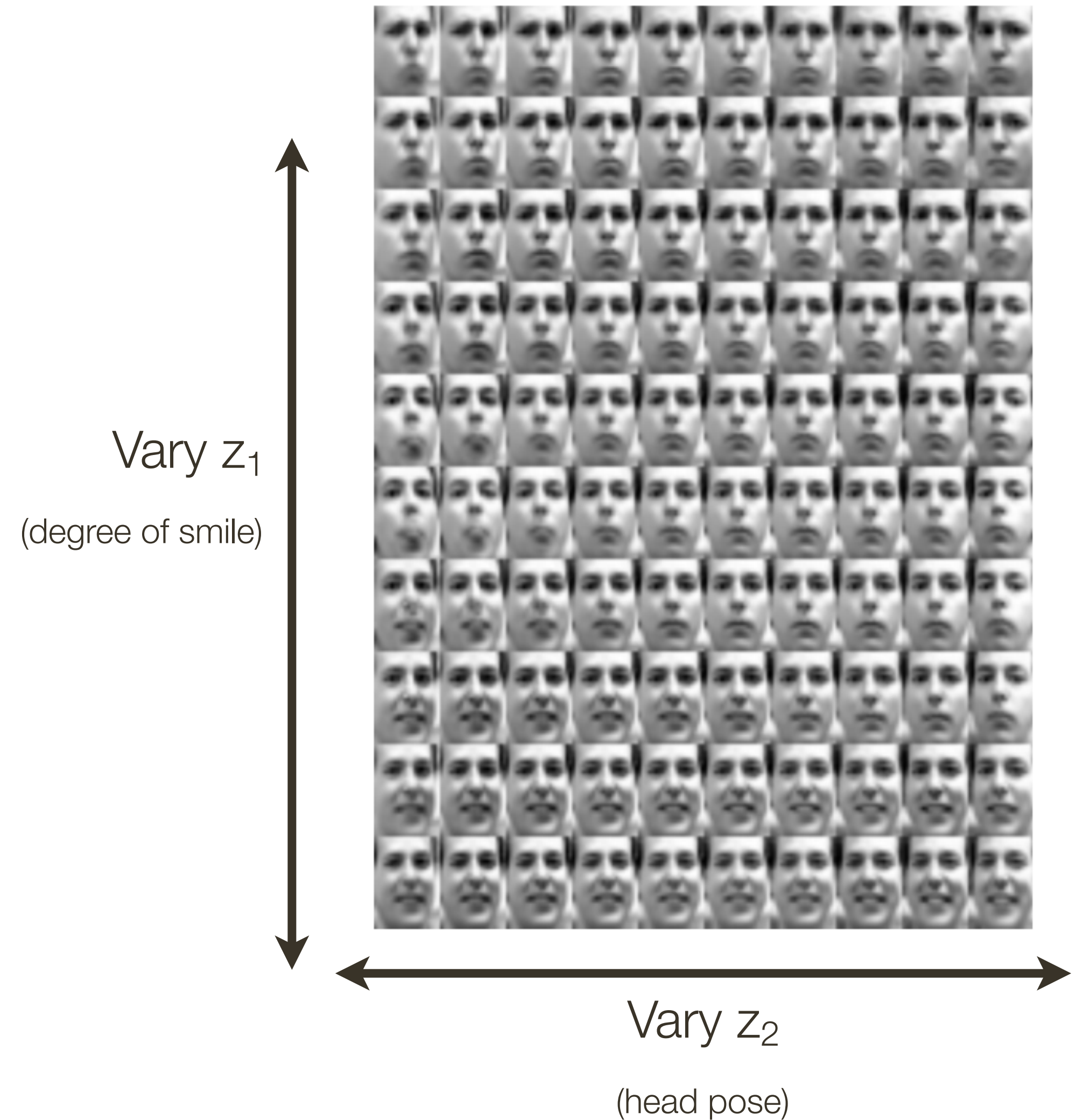


Variational Autoencoder: Generating Data

Diagonal prior on $z \Rightarrow$
independent latent variables

Different dimensions of z encode
interpretable factors of variation

Data manifold for 2-d z



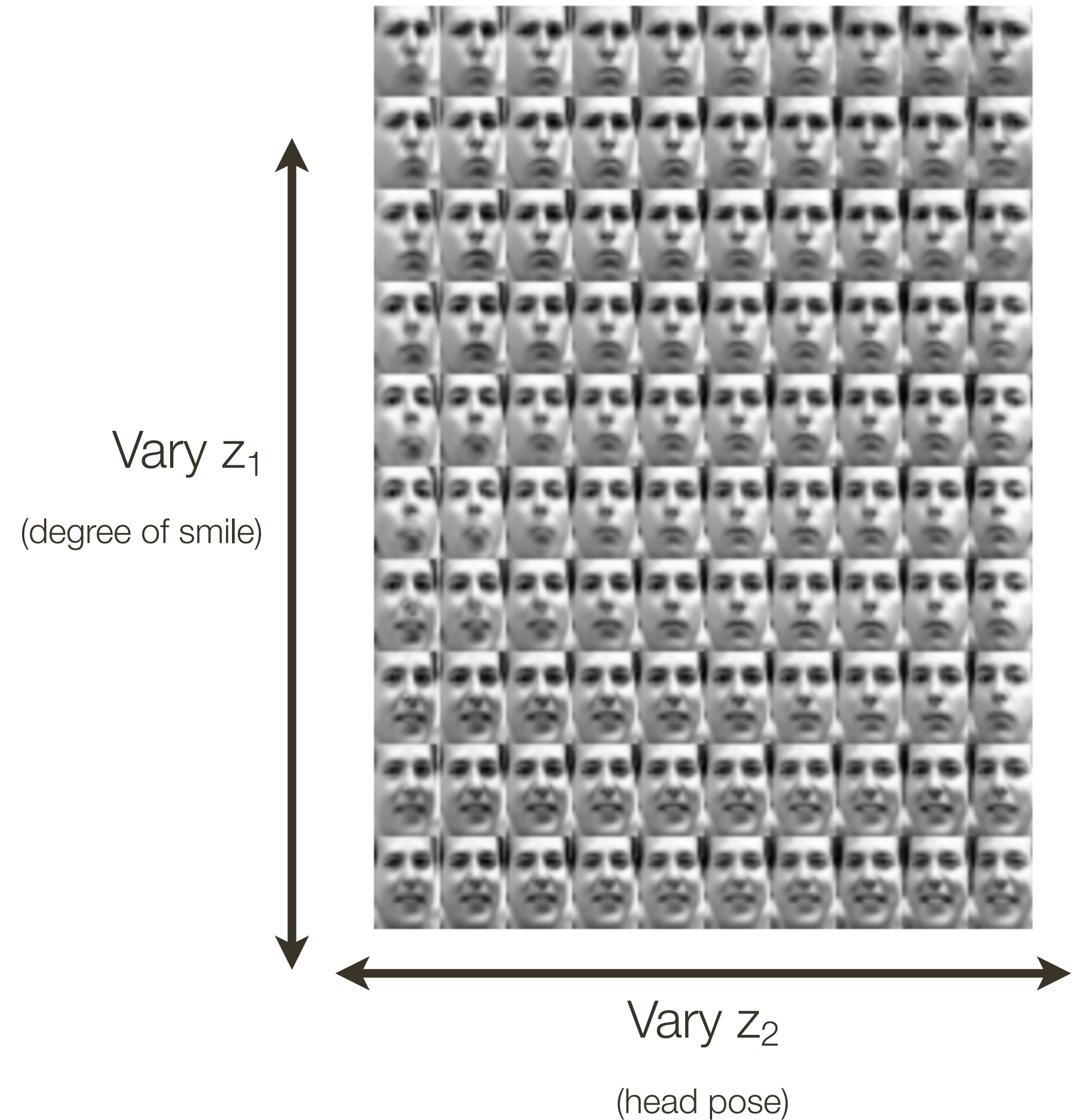
Variational Autoencoder: Generating Data

Diagonal prior on $z \Rightarrow$
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



Variational Autoencoder: Generating Data



32x32 CIFAR-10



Labeled Faces in the Wild