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

**Lecture 17: Generative Models Cont. (VAE, GANs)**

# Course **Logistics**

— **Great set of projects!**

    — Modes of "sub-optimality" at this stage:

            (1) not enough thought into what alterations to base-model should be tested

            (2) motivation for architectures

    — What am I expecting for the project?

    — Feedback (still working on this)

— **Proposal** and **presentation** submission (Canvas on **Friday**)

— Paper **presentations** and list

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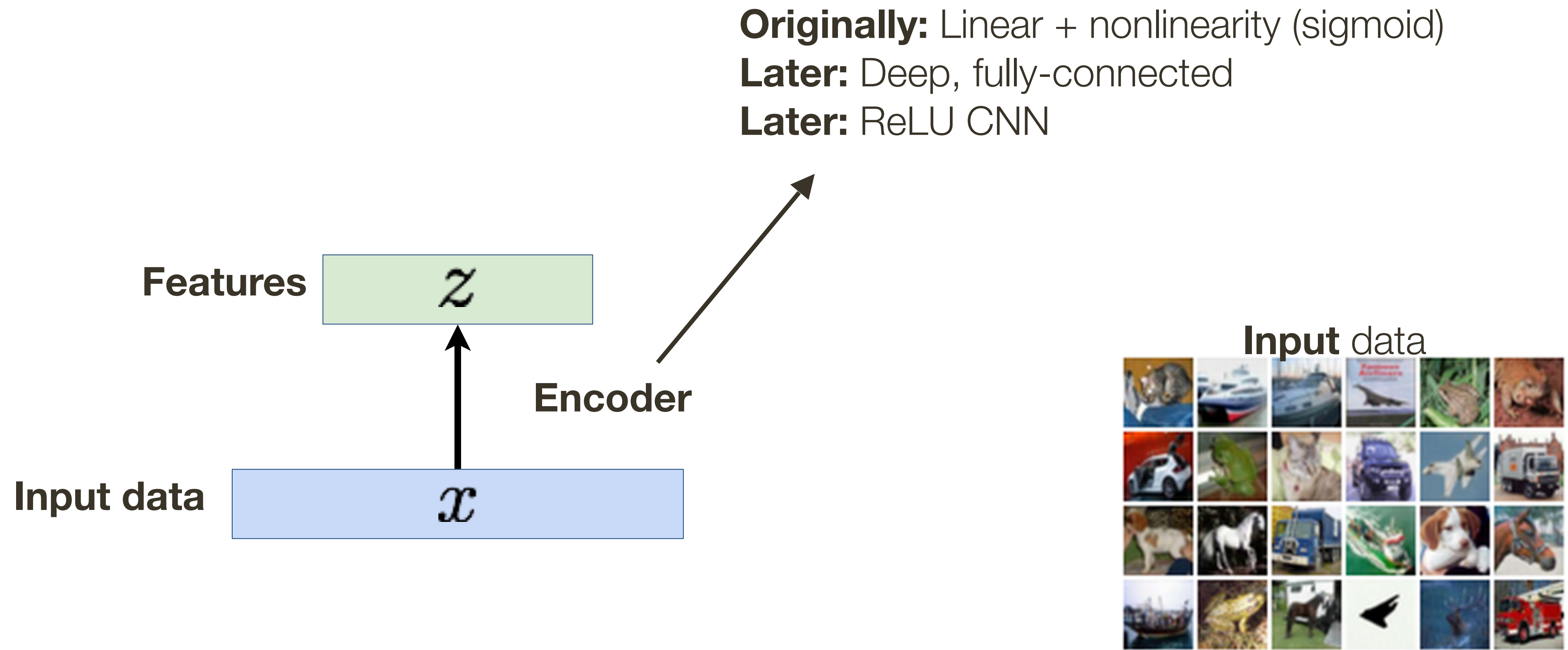
$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, ..., x_{i-1})$$

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** $x$
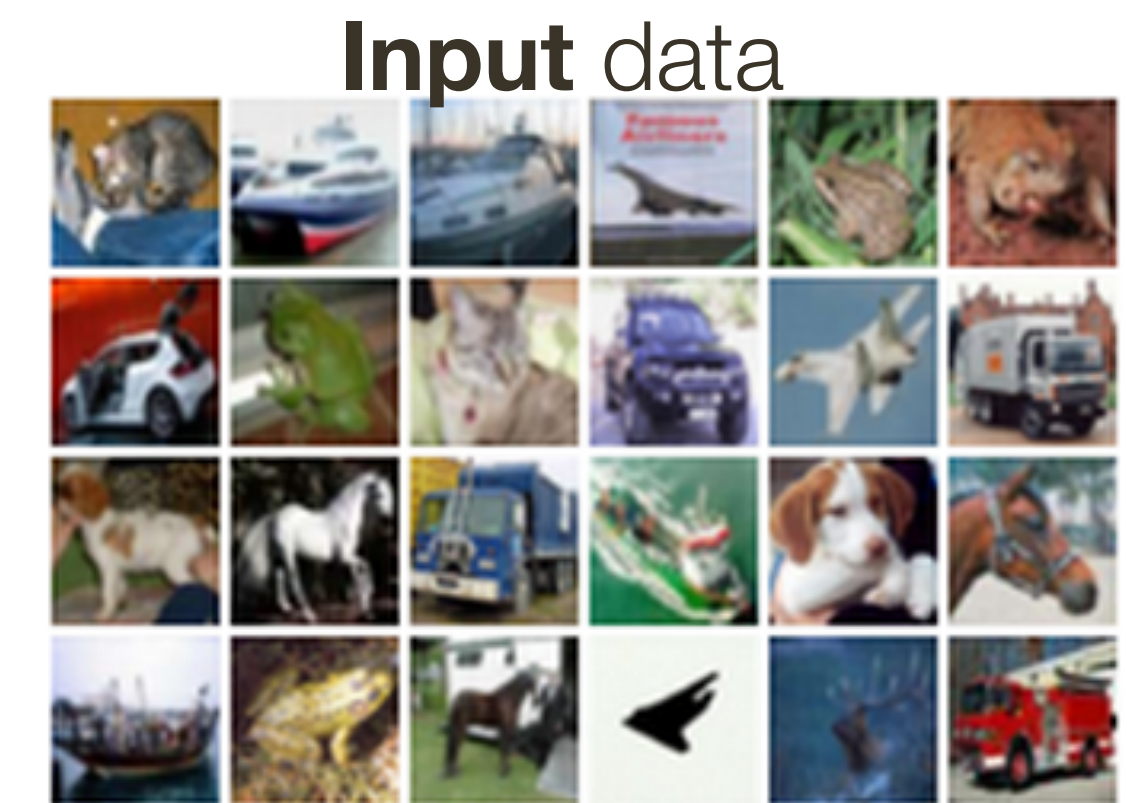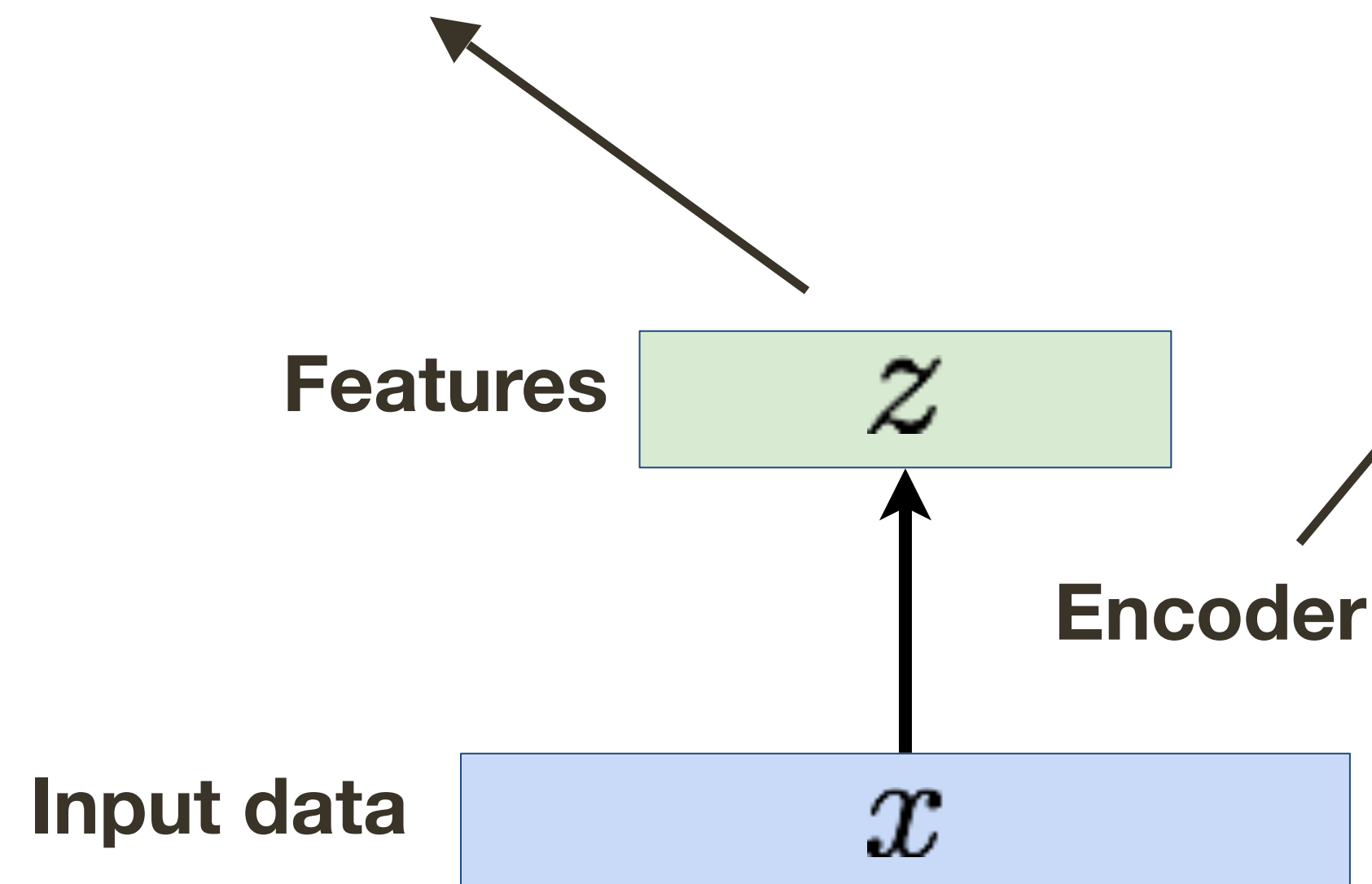


**Input** data

# **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
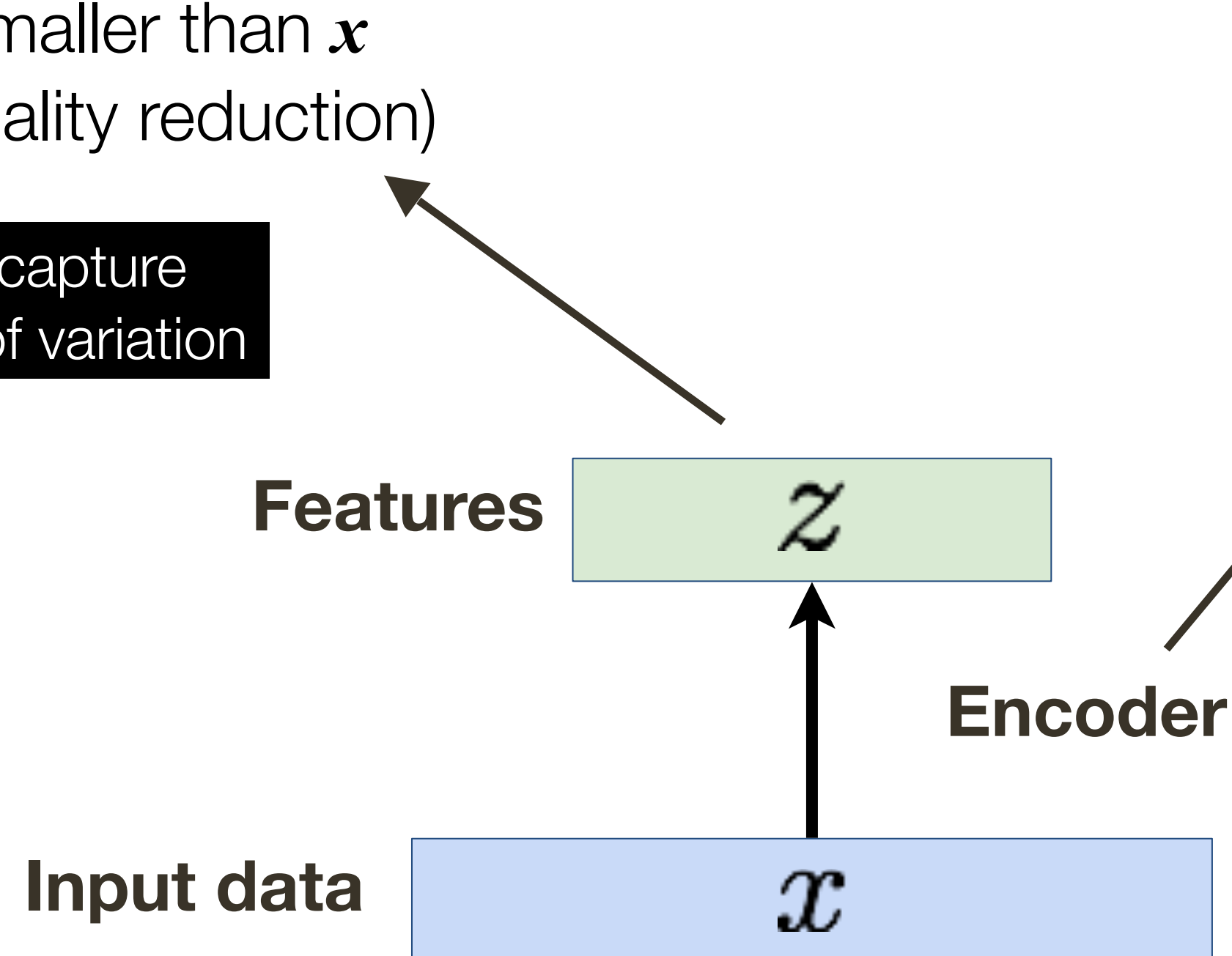**Later:** ReLU CNN

**Features** $z$

**Encoder**

**Input** data

**Input data** $x$

# **Autoencoders** Reminder …

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data
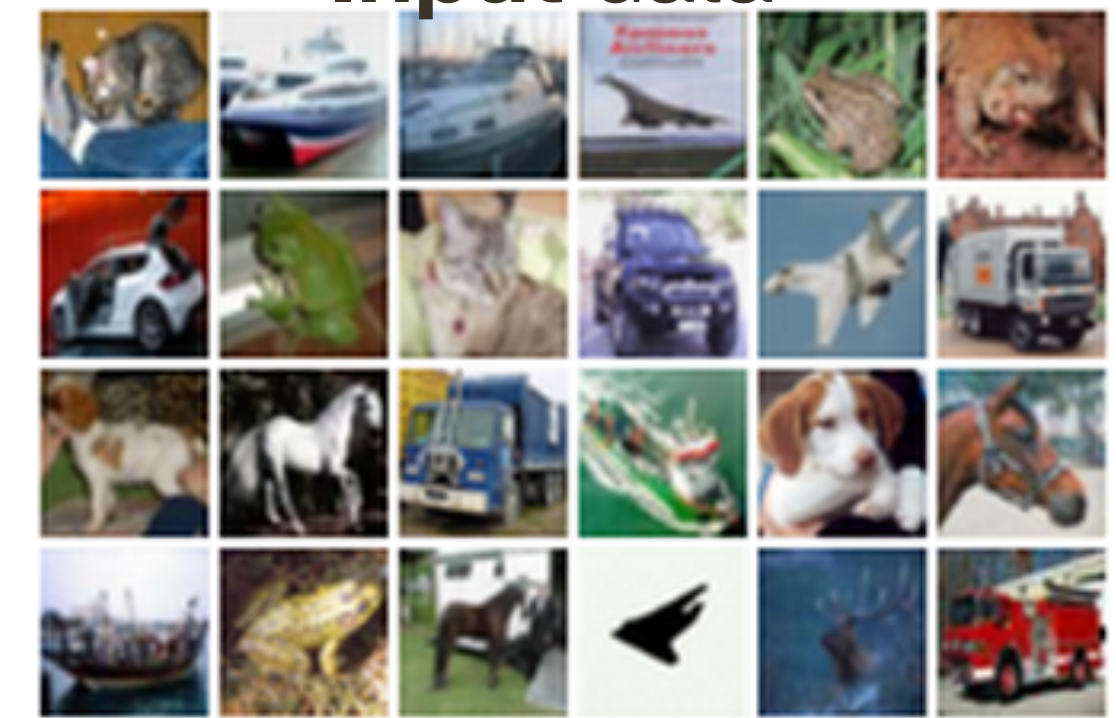
$z$ usually smaller than $x$
(dimensionality reduction)

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

Want features that capture **meaningful** factors of variation
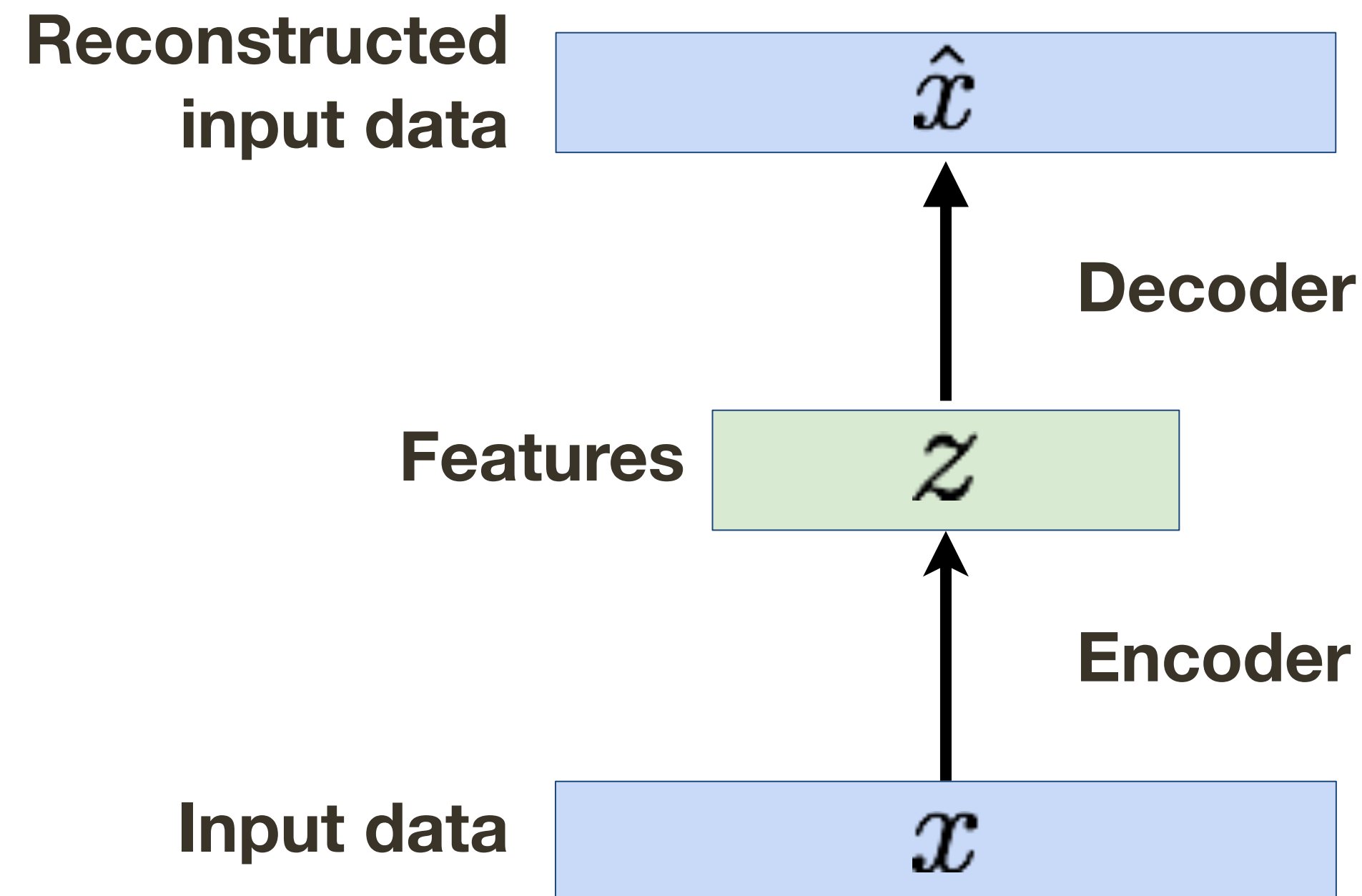
**Features** $z$

**Encoder**

**Input** data

**Input data** $x$

# **Autoencoders** Reminder …

Train such that features can reconstruct original data best they can

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

**Decoder**

**Features**   $z$
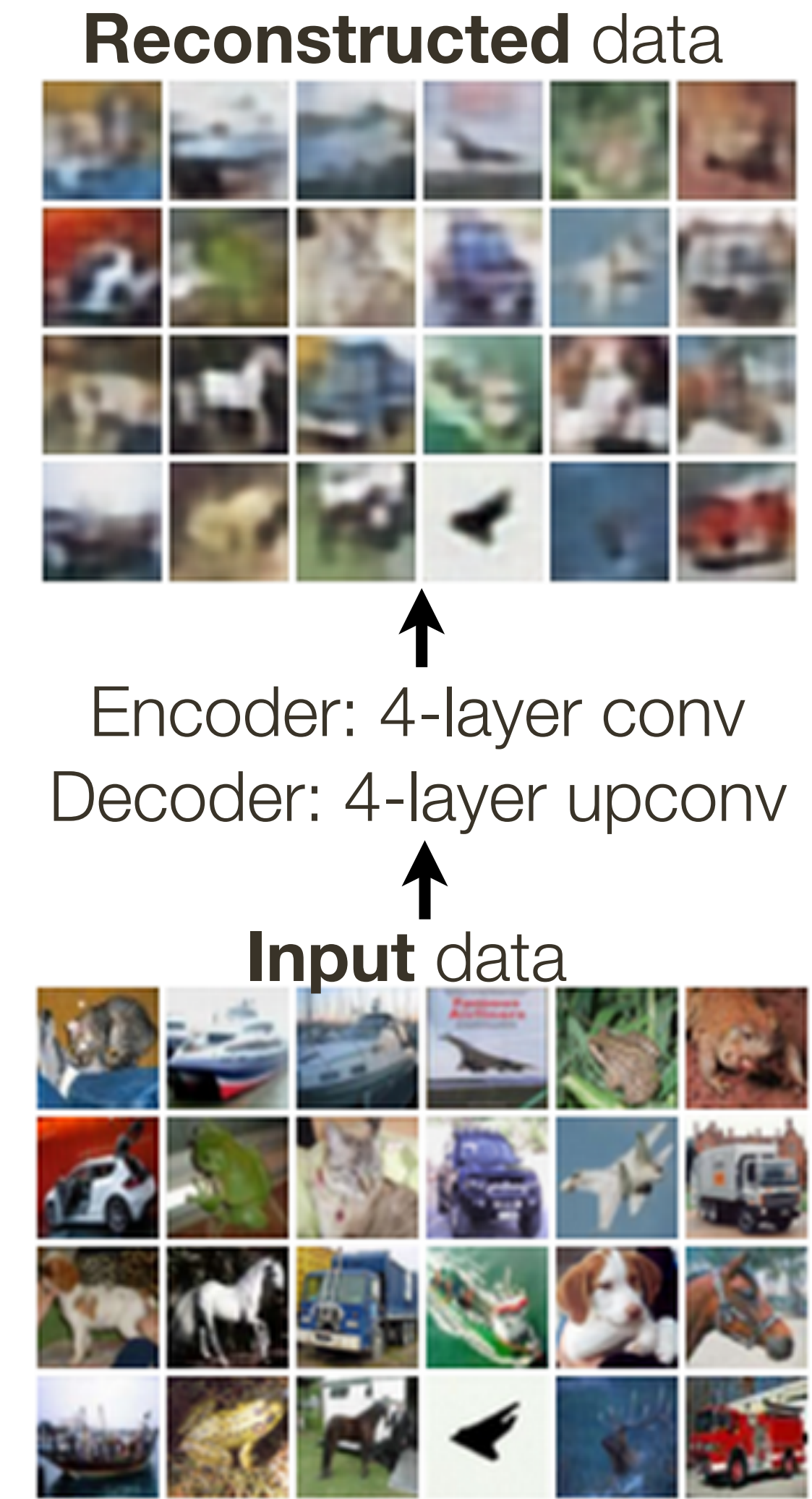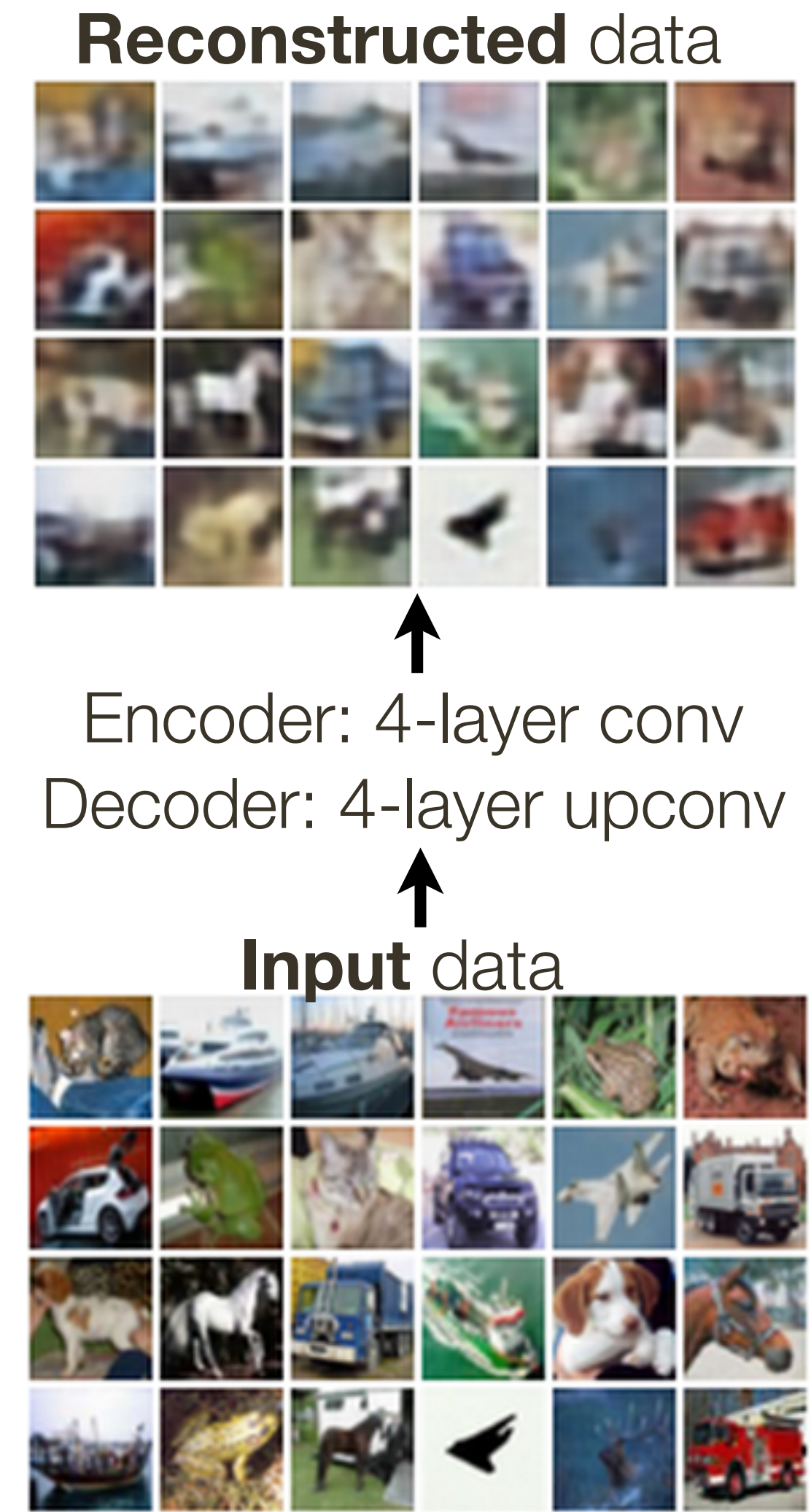
**Encoder**

**Input data**   $x$

**Input** data

# **Autoencoders** Reminder …

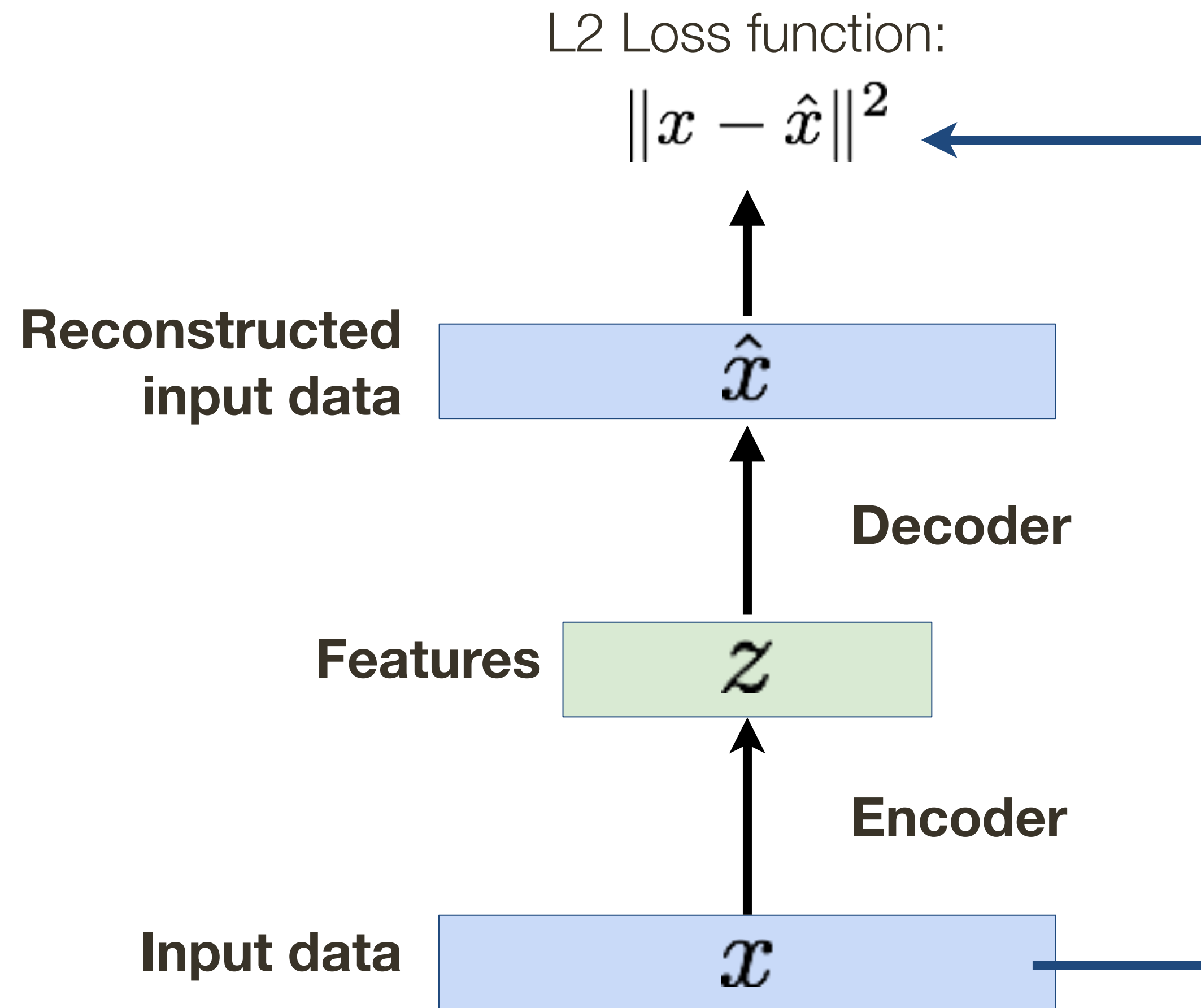Train such that features can reconstruct original data
best they can

**Reconstructed data**



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

**Decoder**

**Features**  $z$

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

**Encoder**

**Input** data



**Input data**  $x$

# **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**

**Features**    $z$

**Encoder**

**Input data**    $x$

Fine-tune
encoder
jointly with
classifier

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

bird        plane

dog      deer      truck

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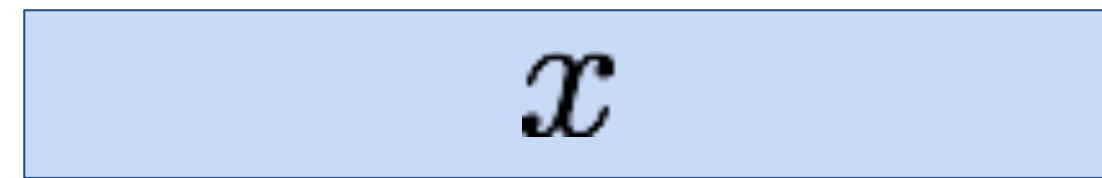
Probabilistic spin on autoencoder - will let us sample from the model to generate

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

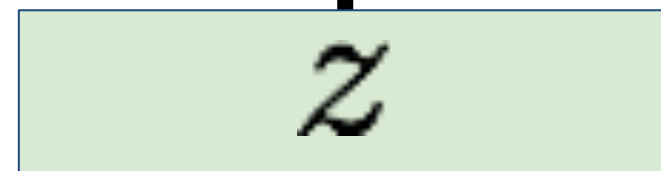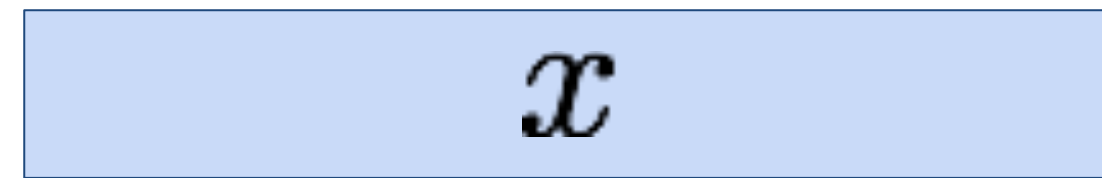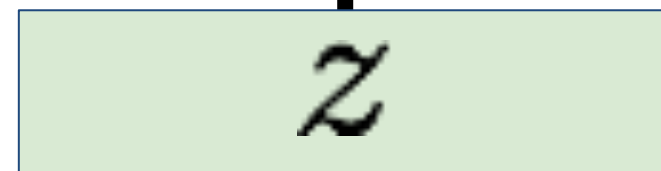Sample from
true **conditional**

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

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**

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



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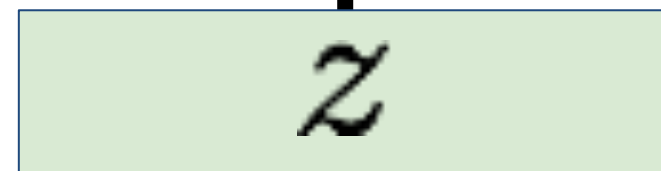
$$p_{\theta*}(x \mid z^{(i)})$$



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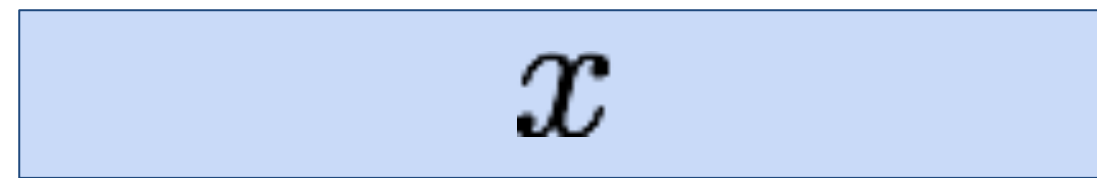
We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **represent** this model?

Sample from
true **conditional**
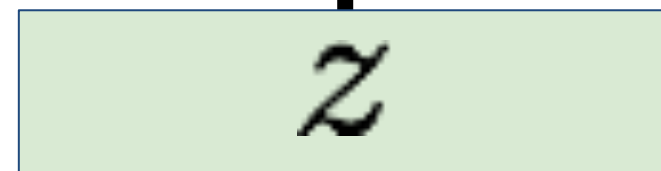
$$p_{\theta*}(x \mid z^{(i)})$$



**Decoder
Network**

Sample from
true **prior**

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

Choose prior $p(z)$ 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
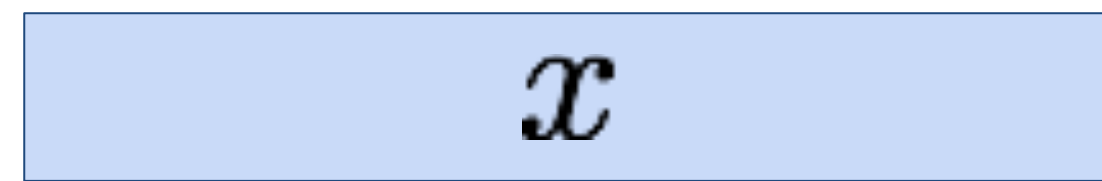
We want to **estimate the true parameters** $\theta*$ of this generative model

How do we **train** this model?

Sample from
true **conditional**

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



**Decoder
Network**

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**

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



**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)

# **Variational** Autoencoders

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

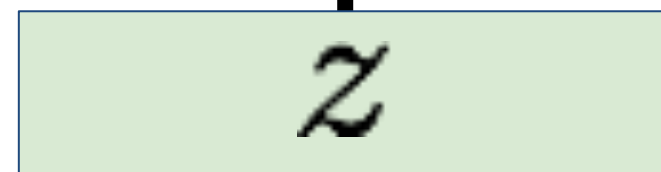How do we **train** this model?

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**

Sample from true **prior**

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

(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
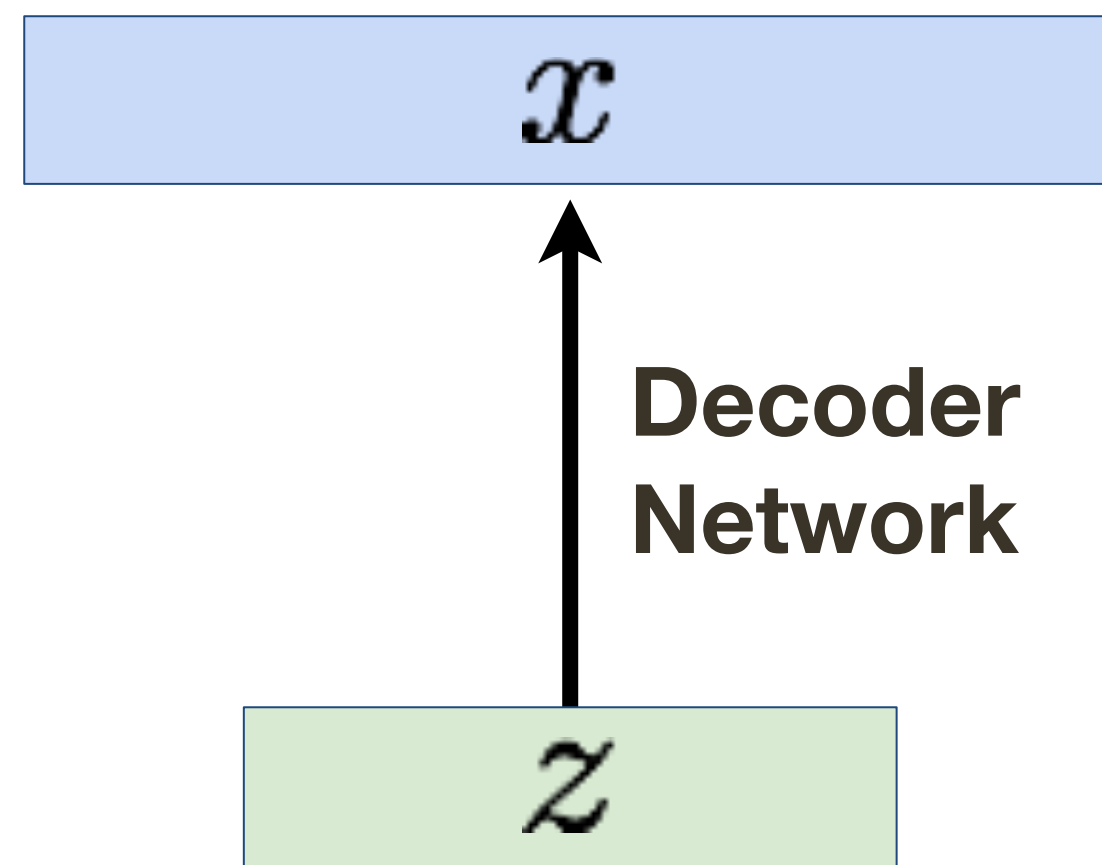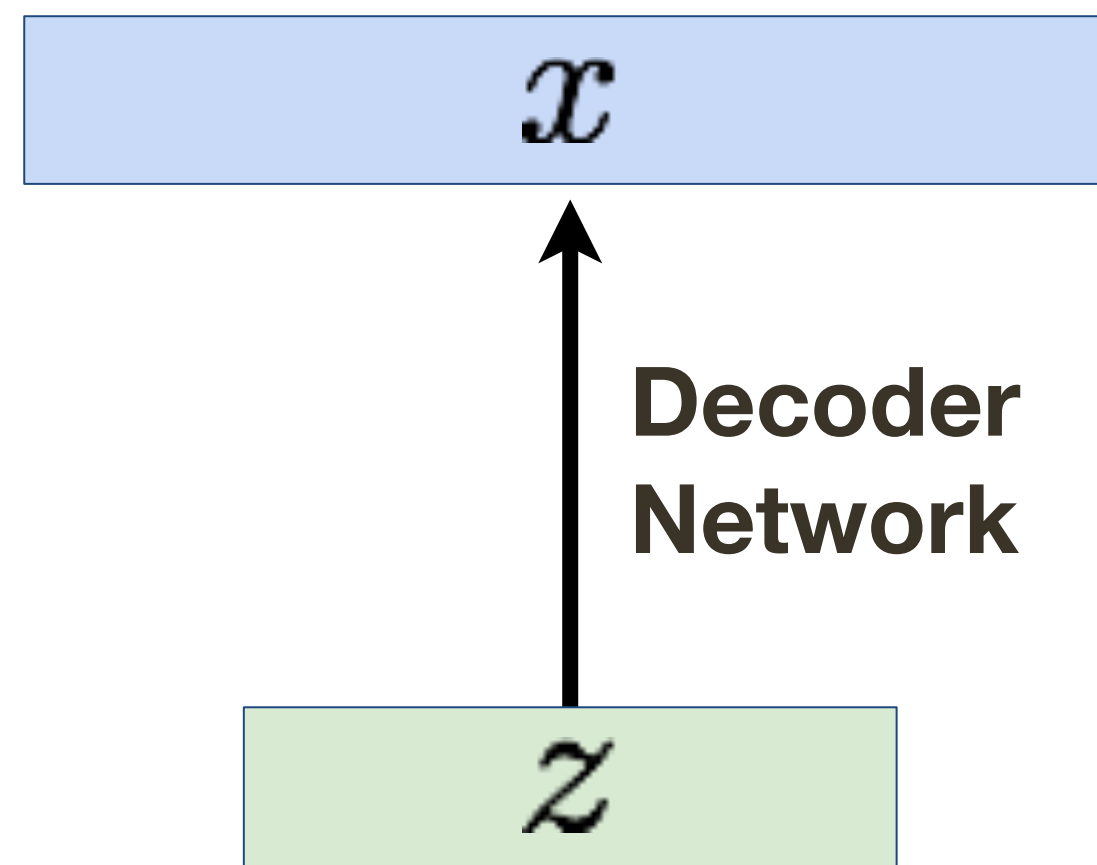
Sample from
true **conditional**

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

Sample from
true **prior**

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

$x$

**Decoder
Network**

$z$

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

🙂

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

Simple **Gaussian** Prior

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

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

Simple **Gaussian** Prior

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

**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 $z \,|\, x$

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

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

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

**Encoder** Network

$q_\phi(z|x)$

(parameters $\phi$)

$x$

**Decoder** Network

$p_\theta(x|z)$

(parameters $\theta$)

$z$

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

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

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

**Encoder** Network

$q_\phi(z|x)$

(parameters φ)

**Decoder** Network

$p_\theta(x|z)$

(parameters θ)

$\boxed{x}$

$\boxed{z}$

# **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:

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

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

# **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

[ 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:

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

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

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

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

* 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:

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

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

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

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

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

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:

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

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

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

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

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

**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)
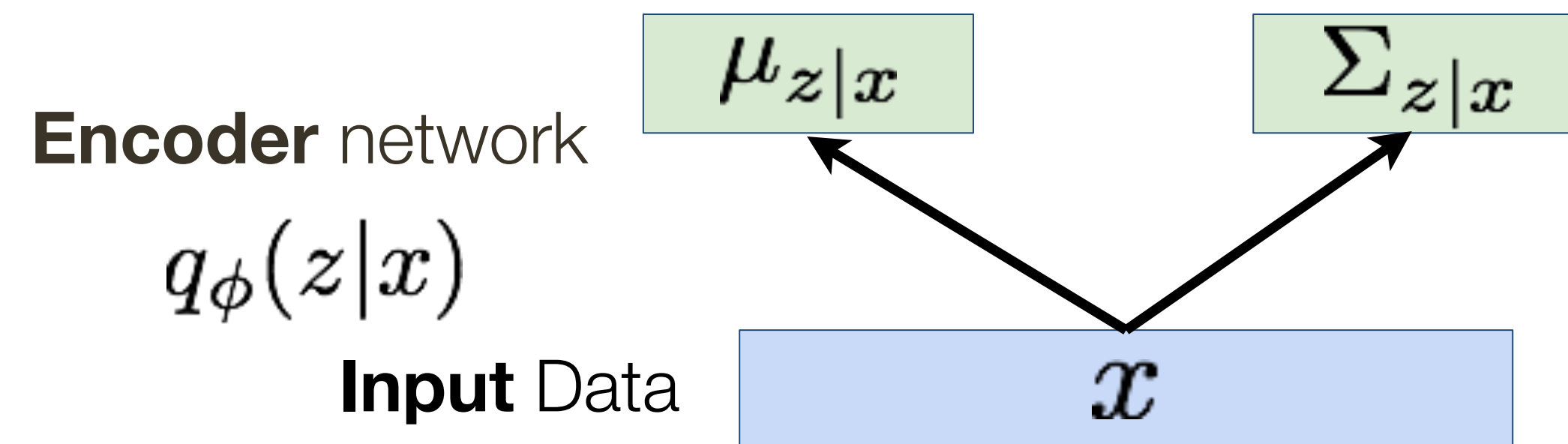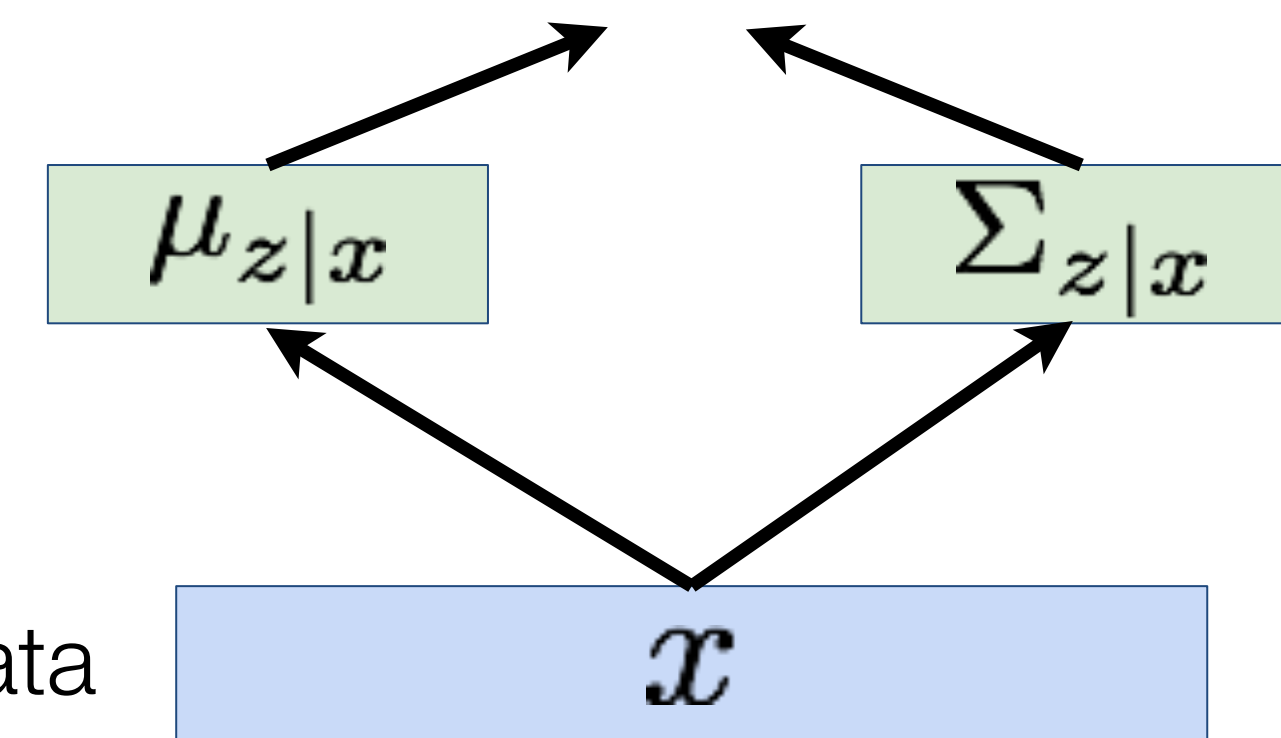for a given mini batch of input data

**Input** Data    $x$

# **Variational** Autoencoder: Learning

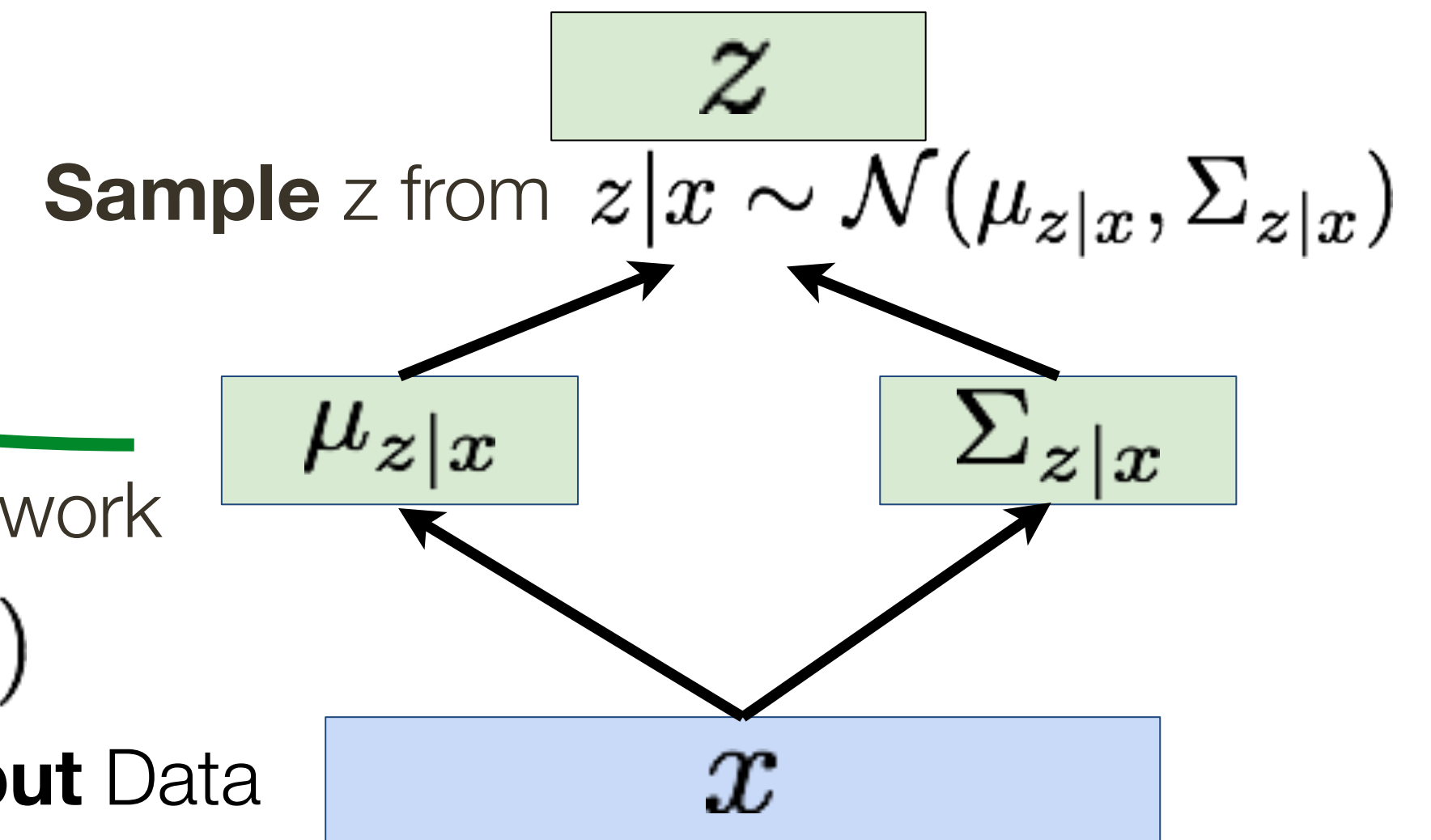**Putting it all together**:
maximizing the likelihood lower
bound

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

**Encoder** network
$q_\phi(z|x)$
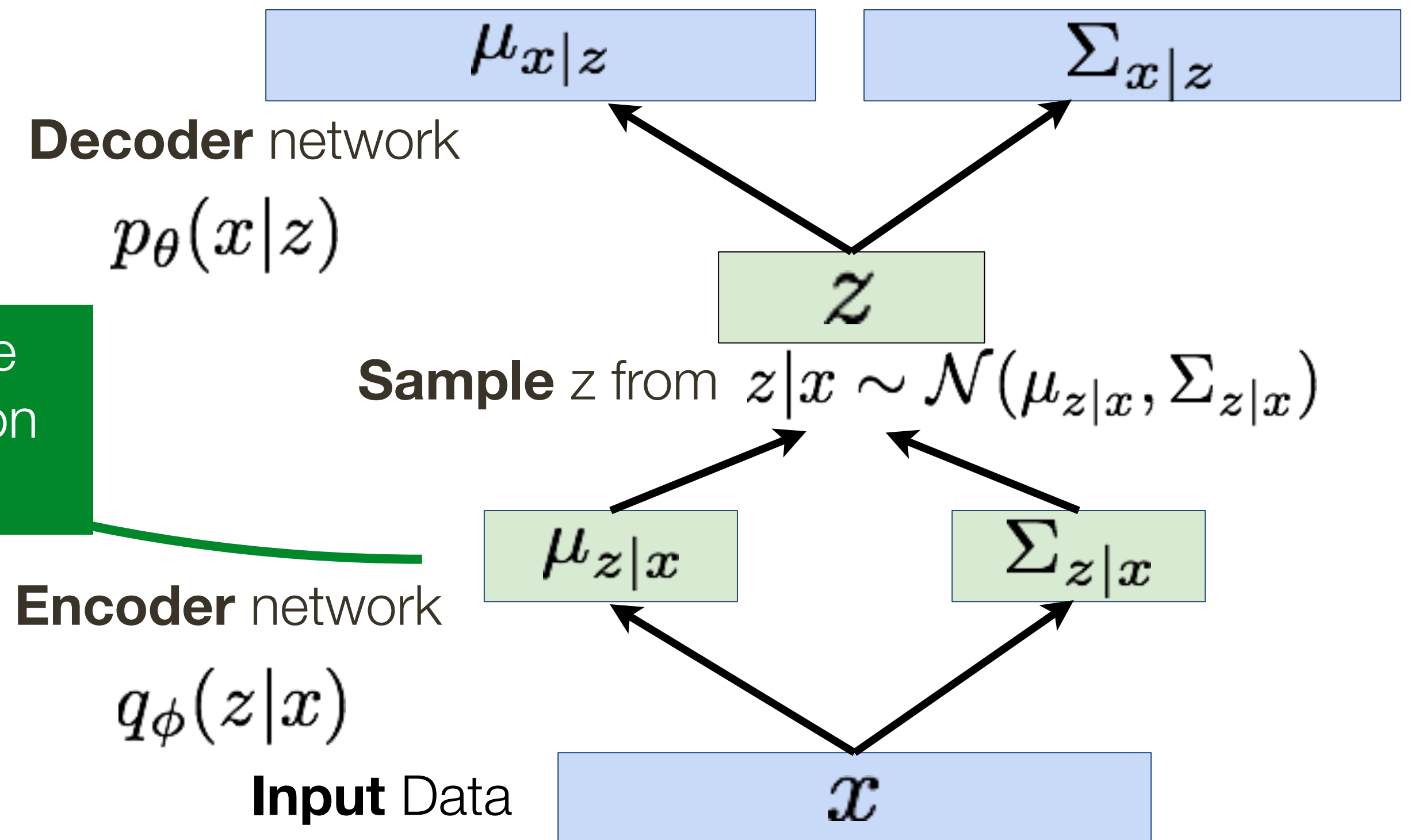
$\mu_{z|x}$   $\Sigma_{z|x}$

**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

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

**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

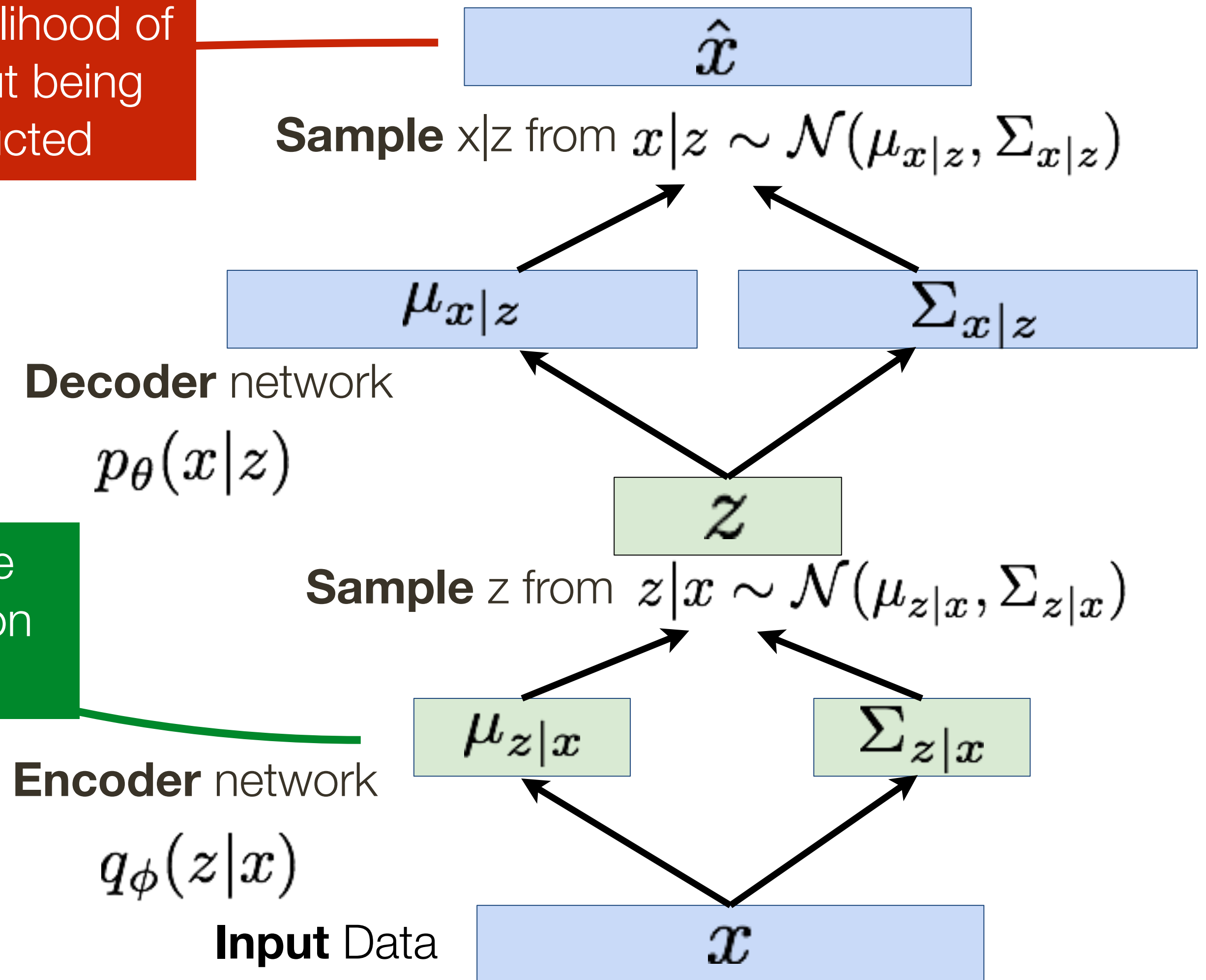**Sample** z from $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

$z$

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

**Encoder** network

$q_\phi(z|x)$

**Input** Data

$x$

# **Variational** Autoencoder: Learning

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

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

Make approximate
posterior distribution
close to prior

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

**Decoder** network

$p_\theta(x\mid z)$

$z$

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

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

**Encoder** network

$q_\phi(z\mid x)$

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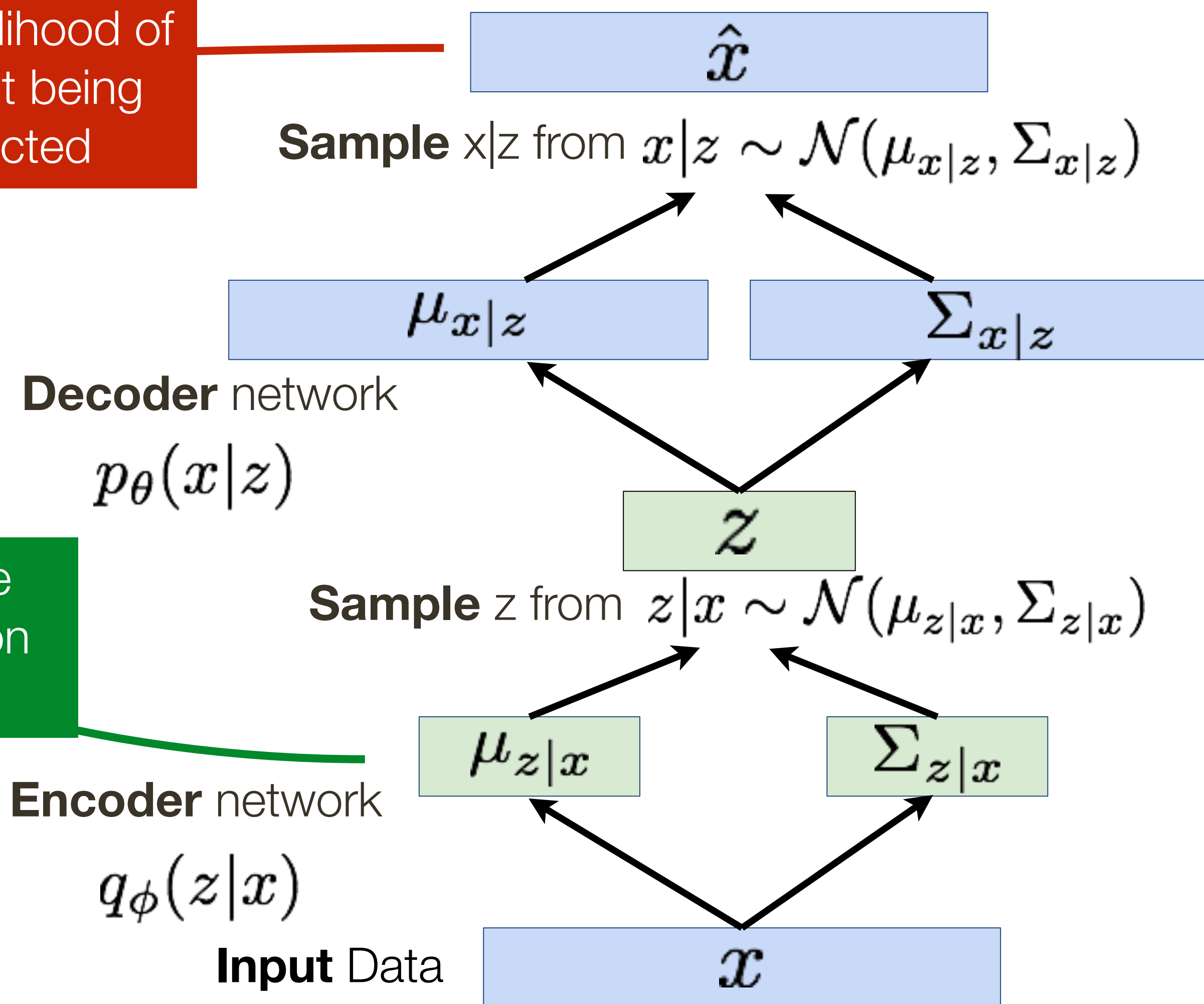
# **Variational** Autoencoder: Learning

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

Maximize likelihood of original input being reconstructed

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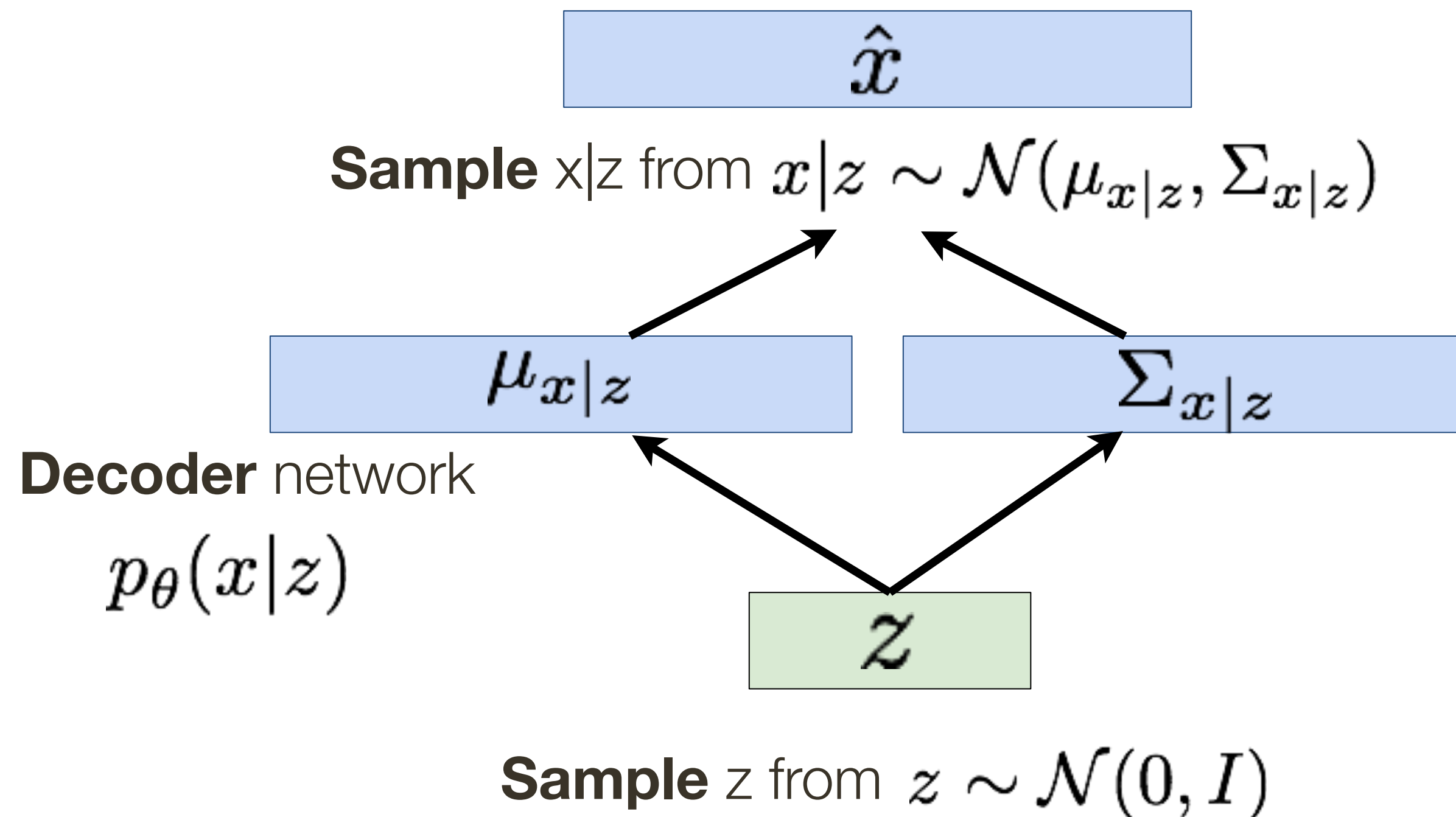
Make approximate posterior distribution close to prior

$\hat{x}$

**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    $x$

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

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

**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**



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

**Data manifold** for 2-d z



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

(a) Frame 1

(b) Frame 2 (ground truth)

(c) Frame 2 (Sample 1)

(d) Frame 2 (Sample 2)

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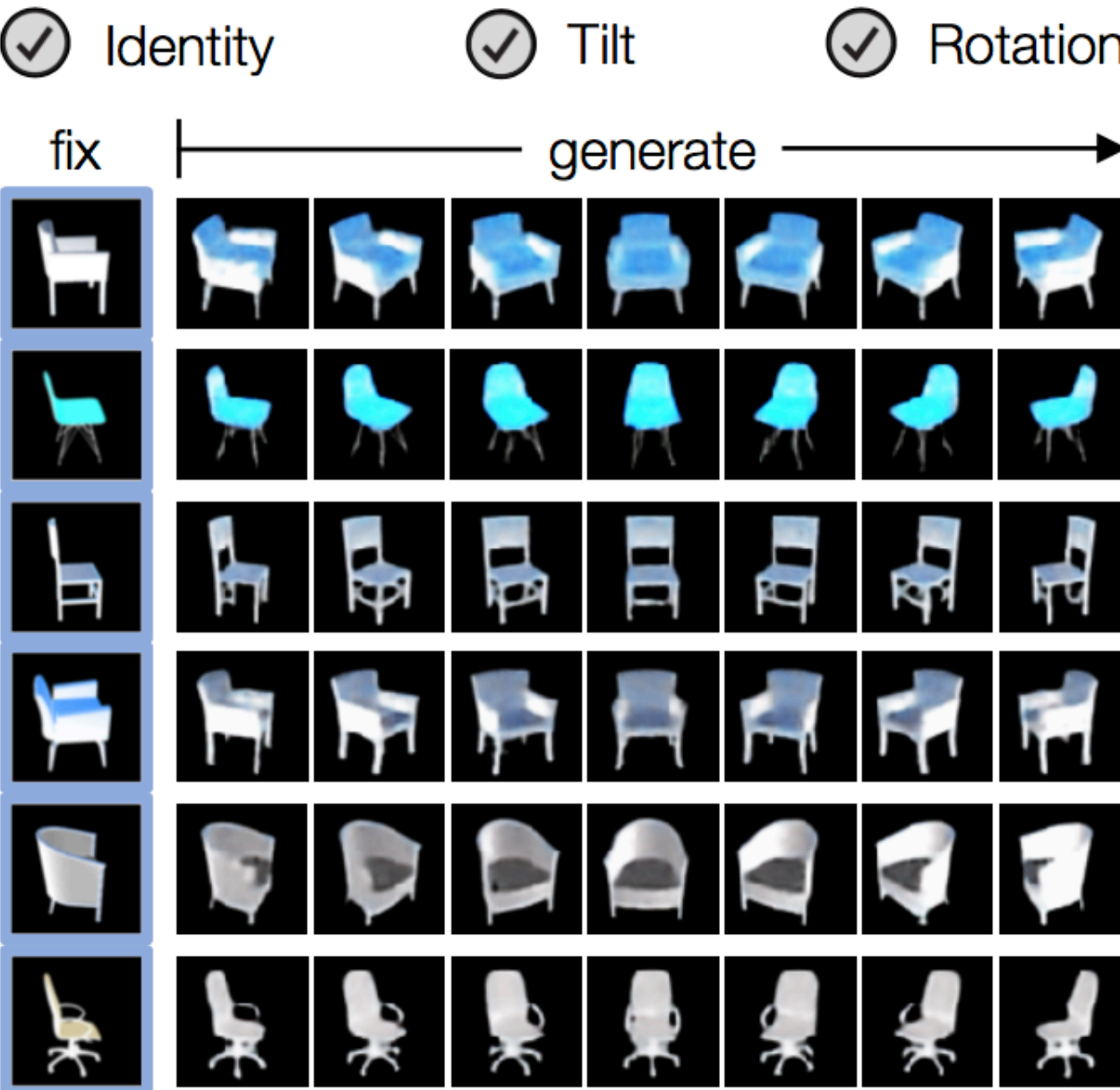
# **Variational** Autoencoder (VAE) **+ LSTM**

# VAE + LSTM with Structured Latent Space

# Results: Chair CAD dataset



(a) Partial control.

(b) Full control.

## Ablation

| | Bound | Static | $-C$ | | $+C$ | |
|---|---|---|---|---|---|---|
| | | | $-S$ | $+S$ | $-S$ | $+S$ |
| Intra-E $\downarrow$ | 1.98 | 40.33 | 17.64 | 7.79 | 14.81 | **5.50** |
| Inter-E $\uparrow$ | 1.39 | 0.42 | 0.73 | 1.35 | 1.02 | **1.37** |
| I-Score $\uparrow$ | 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 $\downarrow$ | 1.98 | 14.68 | **5.50** |
| Inter-E $\uparrow$ | 1.39 | 1.34 | **1.37** |
| I-Score $\uparrow$ | 4.01 | 3.39 | **3.94** |

# Results: Weizmann Human Action dataset



| | 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



| | 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)

# **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

# Generative Adversarial Networks (GANs)

# **Generative** Adversarial Networks

**Problem:** Want to sample from complex, high-dimensional training distribution. There is no direct way to do this!

# **Generative** Adversarial Networks

**Problem:** Want to sample from complex, high-dimensional training distribution. There is no direct way to do this!

**Solution:** Sample from a simple distributions, e.g., random noise. Learn transformation to the training distribution
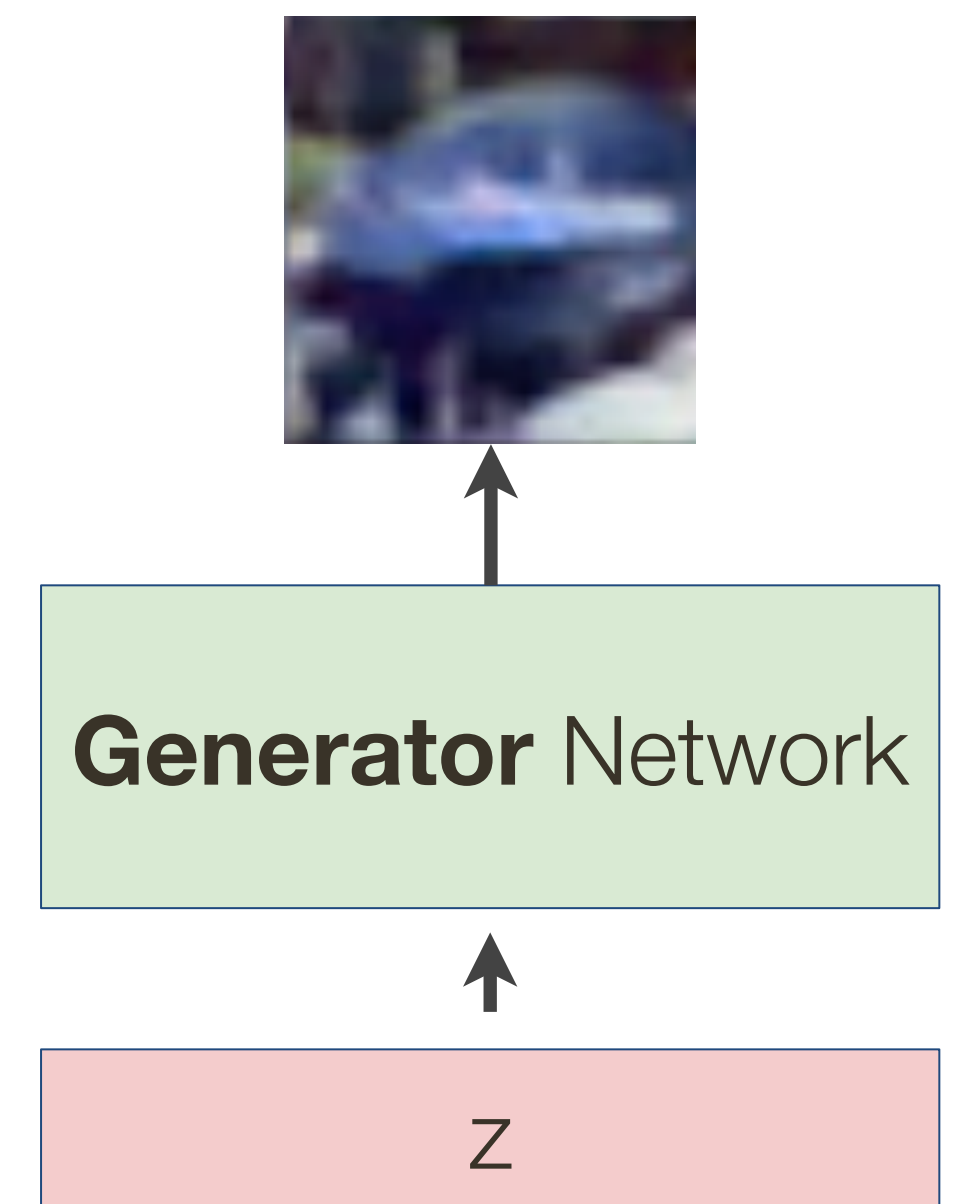
# **Generative** Adversarial Networks

**Problem:** Want to sample from complex, high-dimensional training distribution. There is no direct way to do this!

**Solution:** Sample from a simple distributions, e.g., random noise. Learn transformation to the training distribution

**Question:** What can we use to represent complex transformation function?

# **Generative** Adversarial Networks

[ Goodfellow et al., 2014 ]

**Problem:** Want to sample from complex, high-dimensional training distribution. There is no direct way to do this!

**Output**: Sample from training distribution

**Solution:** Sample from a simple distributions, e.g., random noise. Learn transformation to the training distribution

**Question:** What can we use to represent complex transformation function?

**Generator** Network

**Input**: Random noise

z

# Training GANs: Two-player Game

**Generator** network: try to fool the discriminator by generating real-looking images
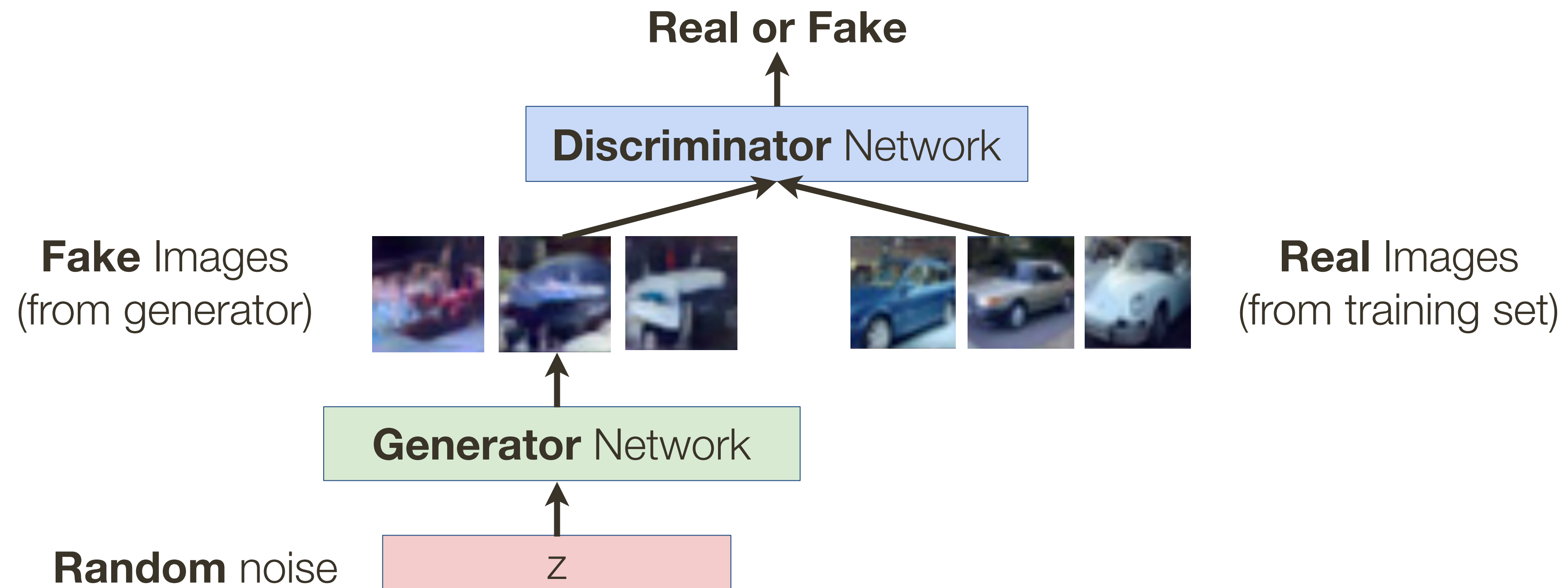**Discriminator** network: try to distinguish between real and fake images

# Training GANs: Two-player Game

**Generator** network: try to fool the discriminator by generating real-looking images
**Discriminator** network: try to distinguish between real and fake images



Real or Fake

Discriminator Network

Fake Images
(from generator)

Real Images
(from training set)

Generator Network

Random noise       z

# **Training** GANs: Two-player Game

**Generator** network: try to fool the discriminator by generating real-looking images
**Discriminator** network: try to distinguish between real and fake images

Train jointly in **minimax game**
Minimax objective function:

Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output
for real data x

Discriminator output for
generated fake data G(z)

- **Discriminator** ($\theta_d$) wants to maximize objective such that D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)
- **Generator** ($\theta_g$) wants to minimize objective such that D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)

# Training GANs: Two-player Game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient **ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient **descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

# Training GANs: Two-player Game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient **ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$
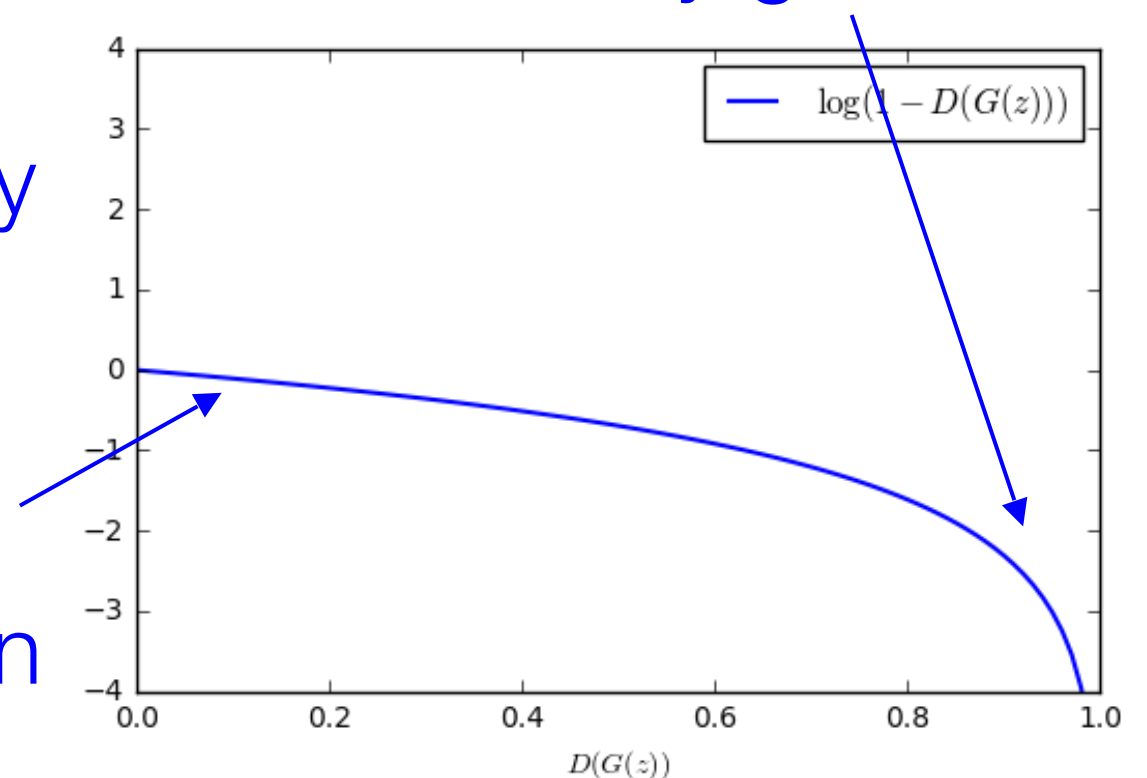
2. Gradient **descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

In practice, optimizing this generator
objective does not work well!

# **Training** GANs: Two-player Game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient **ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient **descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

In practice, optimizing this generator objective does not work well!

Gradient signal dominated by region where sample is already good

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!



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

# **Training** GANs: Two-player Game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:
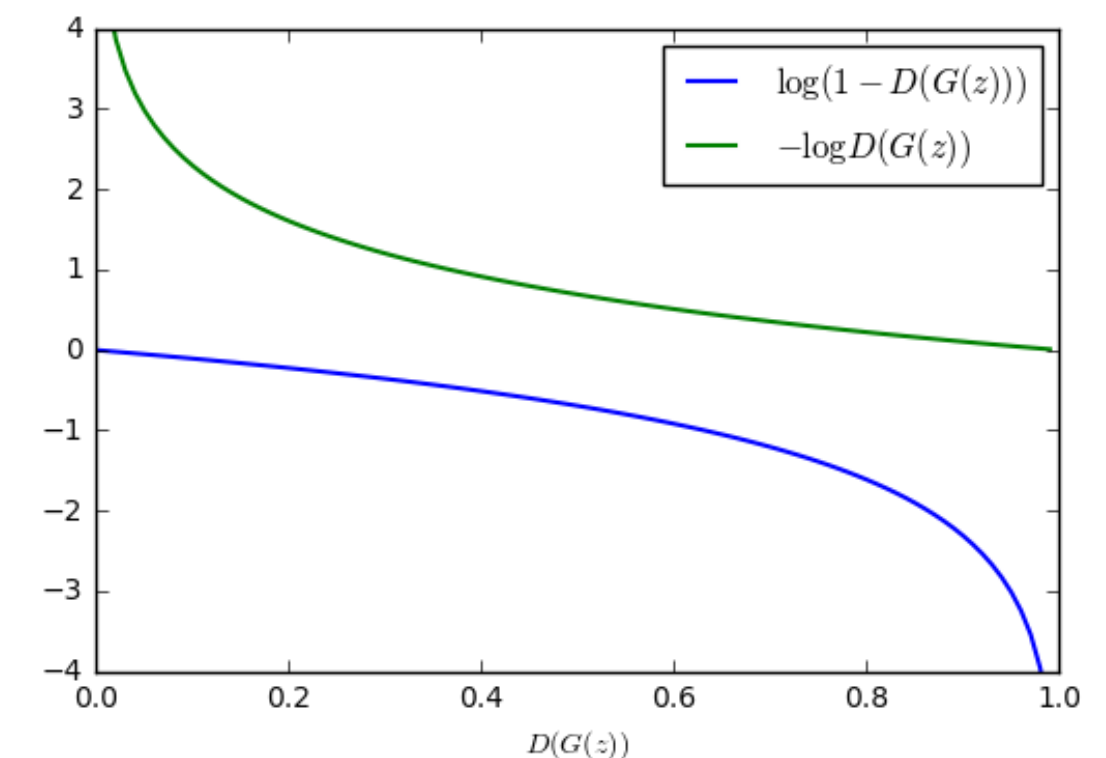
1. Gradient **ascent** on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

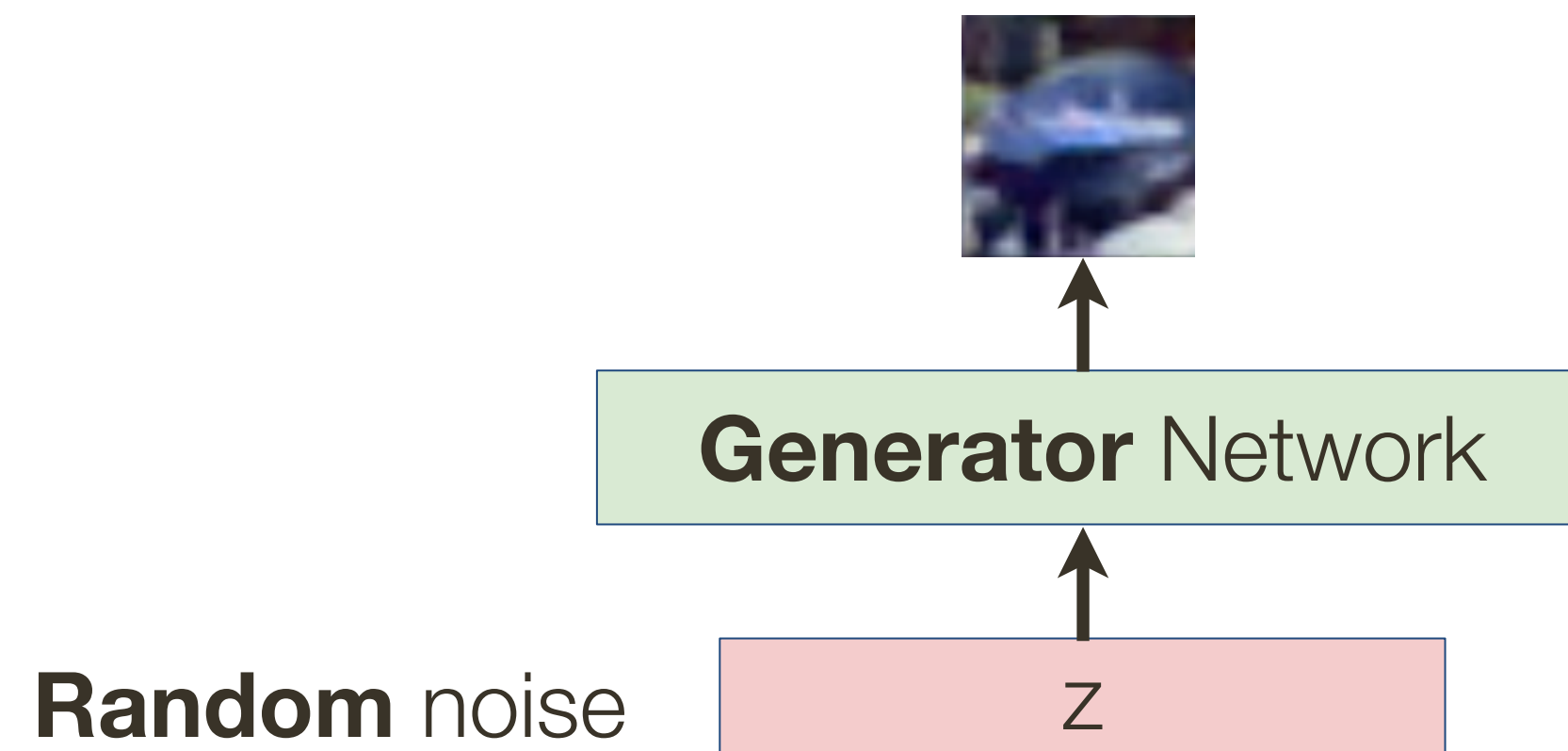2. Instead, gradient **ascent** on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.
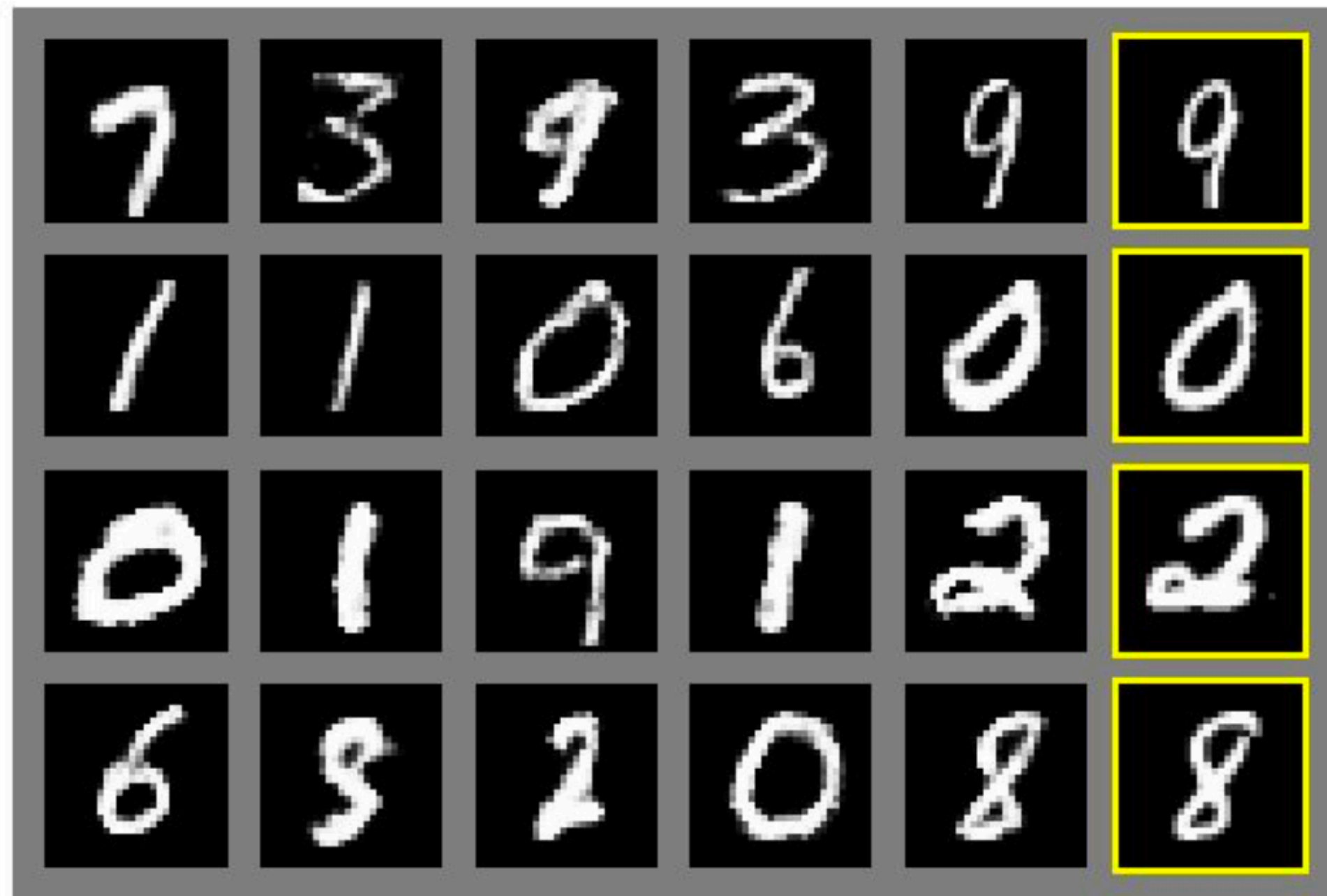
# Sampling **GAN**s



**Random** noise    z

# Generative Adversarial Nets

Generated Samples

# GANs with Convolutional Architectures

# GANs with Convolutional Architectures

Interpolating between points in latent space

# GANs: Interpretable Vector Math

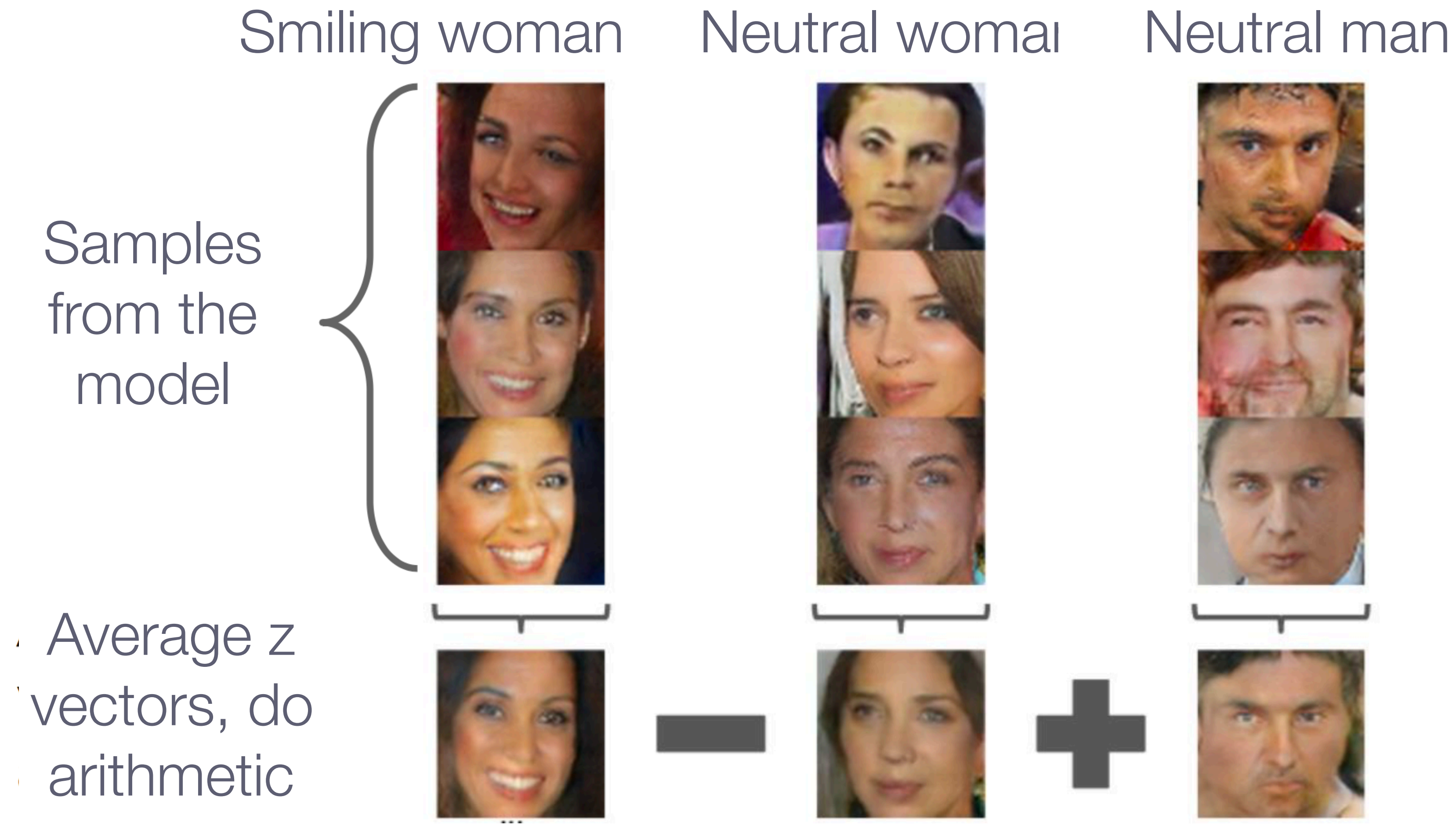Smiling woman        Neutral woman        Neutral man

Samples from the model

# GANs: Interpretable Vector Math

Smiling woman    Neutral woman    Neutral man

Samples from the model

Average z vectors, do arithmetic

# GANs: Interpretable Vector Math

Smiling woman    Neutral woman    Neutral man

Samples from the model

Smiling man

Average z vectors, do arithmetic

# GANs: Interpretable Vector Math

Glasses Man        No Glasses Man        No Glasses Woman

Samples
from the
model

# GANs: Interpretable Vector Math

Glasses Man    No Glasses Man    No Glasses Woman

Samples from the model

Average z vectors, do arithmetic

# GANs: Interpretable Vector Math

Glasses Man     No Glasses Man     No Glasses Woman

Samples from the model

Woman with Glasses

Average z vectors, do arithmetic

# Year of the **GAN**

**Better training and generation**


(a) Church outdoor.  (b) Dining room.
(c) Kitchen.  (d) Conference room.
LSGAN. Mao et al. 2017.


BEGAN. Bertholet et al. 2017.

**Source->Target domain transfer**


Input    Output
horse → zebra
zebra → horse
apple → orange
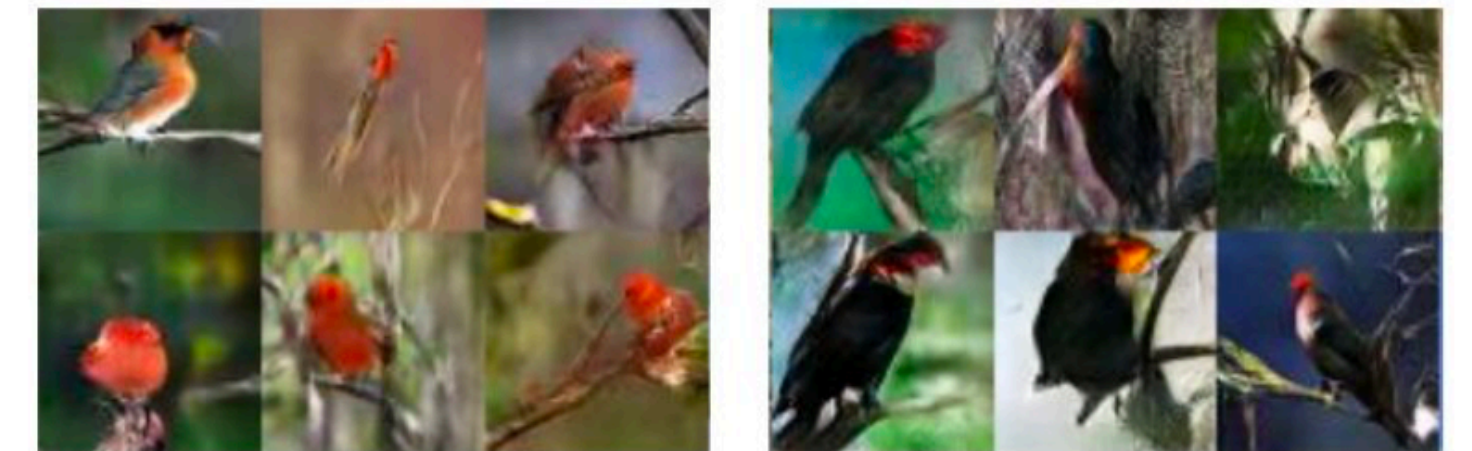CycleGAN. Zhu et al. 2017.


Input    Output
→ summer Yosemite
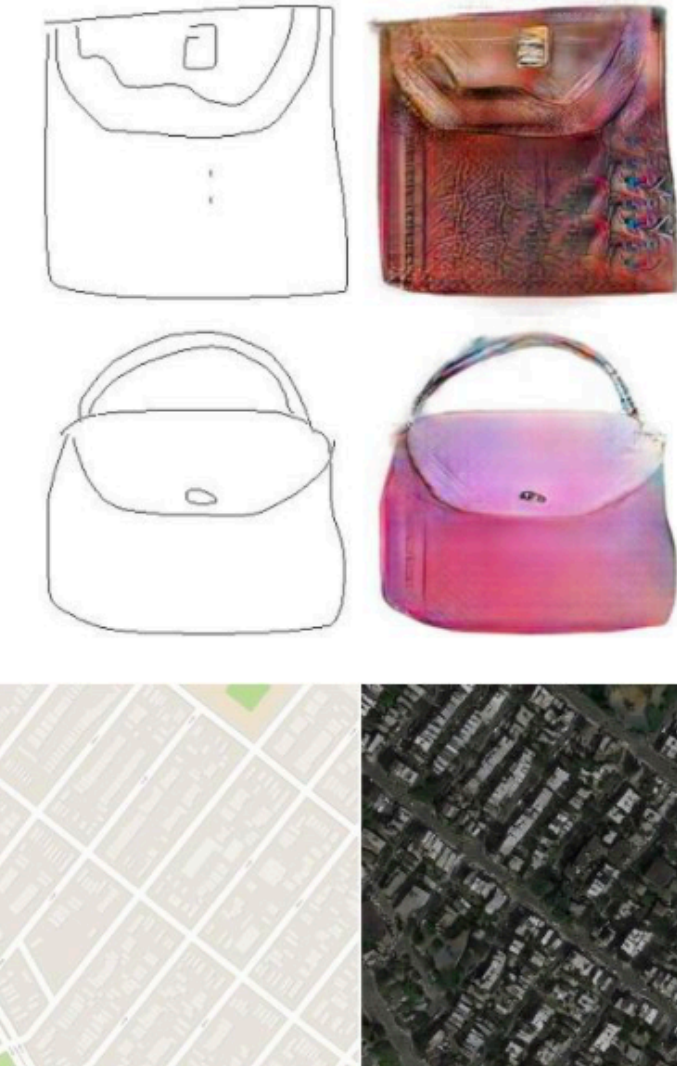→ winter Yosemite

**Text -> Image Synthesis**

this small bird has a pink breast and crown, and black primaries and secondaries.

this magnificent fellow is almost all black with a red crest, and white cheek patch.


Reed et al. 2017.

**Many GAN applications**


Pix2pix. Isola 2017. Many examples at
https://phillipi.github.io/pix2pix/

# Year of the **GAN**

- GAN - Generative Adversarial Networks
- 3D-GAN - Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling
- acGAN - Face Aging With Conditional Generative Adversarial Networks
- AC-GAN - Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN - AdaGAN: Boosting Generative Models
- AEGAN - Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN - Amortised MAP Inference for Image Super-resolution
- AL-CGAN - Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI - Adversarially Learned Inference
- AM-GAN - Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN - Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery
- ArtGAN - ArtGAN: Artwork Synthesis with Conditional Categorial GANs
- b-GAN - b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN - Deep and Hierarchical Implicit Models
- BEGAN - BEGAN: Boundary Equilibrium Generative Adversarial Networks
- BiGAN - Adversarial Feature Learning
- BS-GAN - Boundary-Seeking Generative Adversarial Networks
- CGAN - Conditional Generative Adversarial Nets
- CaloGAN - CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks
- CCGAN - Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- CatGAN - Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks
- CoGAN - Coupled Generative Adversarial Networks

- Context-RNN-GAN - Contextual RNN-GANs for Abstract Reasoning Diagram Generation
- C-RNN-GAN - C-RNN-GAN: Continuous recurrent neural networks with adversarial training
- CS-GAN - Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets
- CVAE-GAN - CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training
- CycleGAN - Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
- DTN - Unsupervised Cross-Domain Image Generation
- DCGAN - Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
- DiscoGAN - Learning to Discover Cross-Domain Relations with Generative Adversarial Networks
- DR-GAN - Disentangled Representation Learning GAN for Pose-Invariant Face Recognition
- DualGAN - DualGAN: Unsupervised Dual Learning for Image-to-Image Translation
- EBGAN - Energy-based Generative Adversarial Network
- f-GAN - f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
- FF-GAN - Towards Large-Pose Face Frontalization in the Wild
- GAWWN - Learning What and Where to Draw
- GeneGAN - GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data
- Geometric GAN - Geometric GAN
- GoGAN - Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking
- GP-GAN - GP-GAN: Towards Realistic High-Resolution Image Blending
- IAN - Neural Photo Editing with Introspective Adversarial Networks
- iGAN - Generative Visual Manipulation on the Natural Image Manifold
- IcGAN - Invertible Conditional GANs for image editing
- ID-CGAN - Image De-raining Using a Conditional Generative Adversarial Network
- Improved GAN - Improved Techniques for Training GANs
- InfoGAN - InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- LAGAN - Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis
- LAPGAN - Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks

# GANs

Don't work with an explicit density function

Take game-theoretic approach: learn to generate from training distribution through 2-player game

## Pros:
— Beautiful, state-of-the-art samples!

## Cons:

— Trickier / more unstable to train

— Can't solve inference queries such as p(x), p(z|x)

## **Active area** of research:

— Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)

— Conditional GANs, GANs for all kinds of applications