



Autoencoders

Deep Learning & Applications

Learning Problems



Supervised Learning:

- ▶ Learn model using dataset with **data-label pairs** $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$
- ▶ Examples: Classification, regression, structured prediction

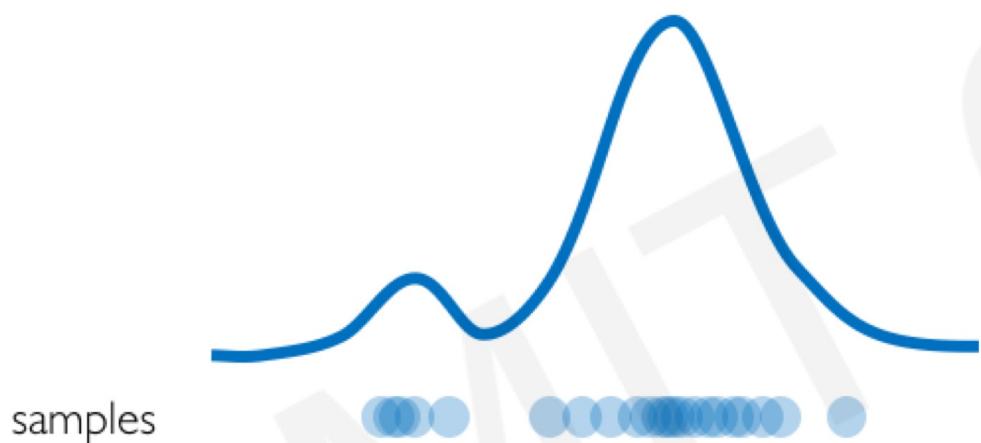
Unsupervised Learning:

- ▶ Learn model using dataset **without labels** $\{\mathbf{x}_i\}_{i=1}^N$
- ▶ Examples: Clustering, dimensionality reduction, generative models

Generative modeling

Goal: Take as input training samples from some distribution and learn a model that represents that distribution

Density Estimation

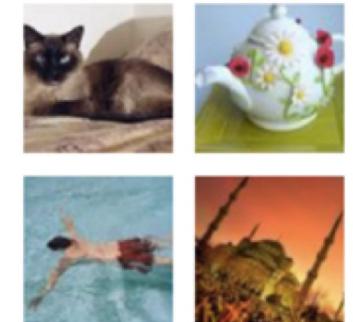


Sample Generation



Input samples

Training data $\sim P_{data}(x)$



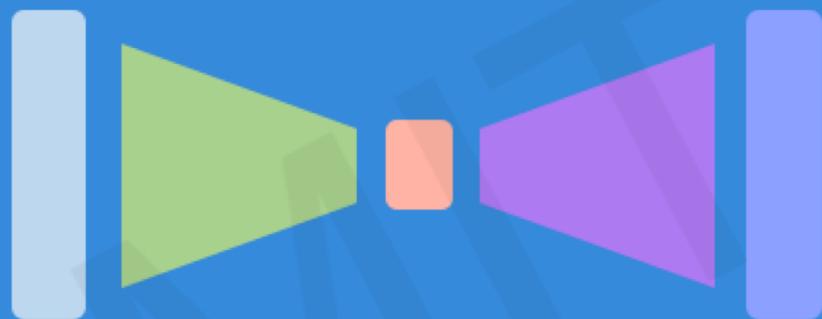
Generated samples

Generated $\sim P_{model}(x)$

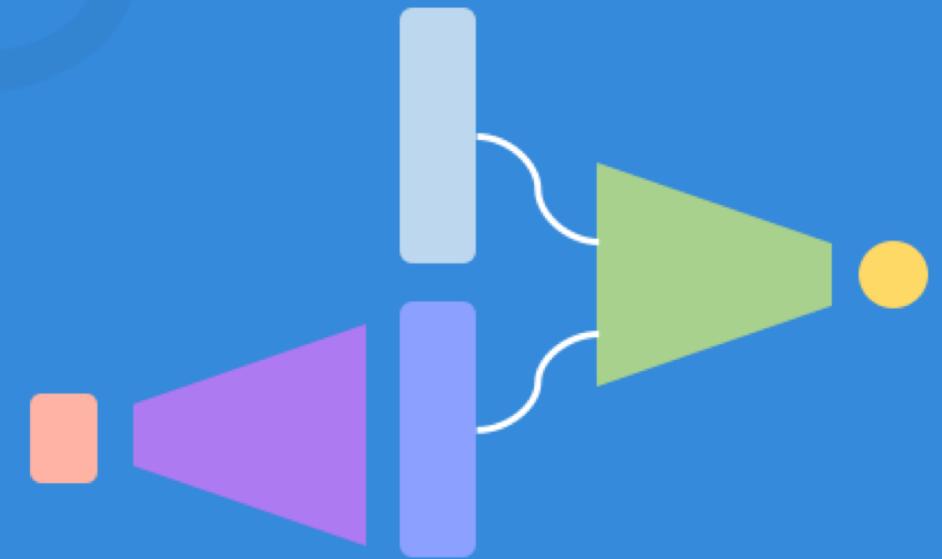
How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

Latent variable models

Autoencoders and Variational
Autoencoders (VAEs)



Generative Adversarial
Networks (GANs)



What is a latent variable?



Myth of the Cave

What is a latent variable?

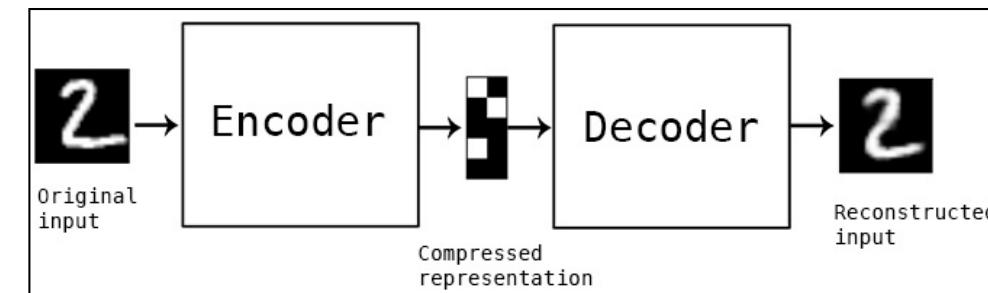
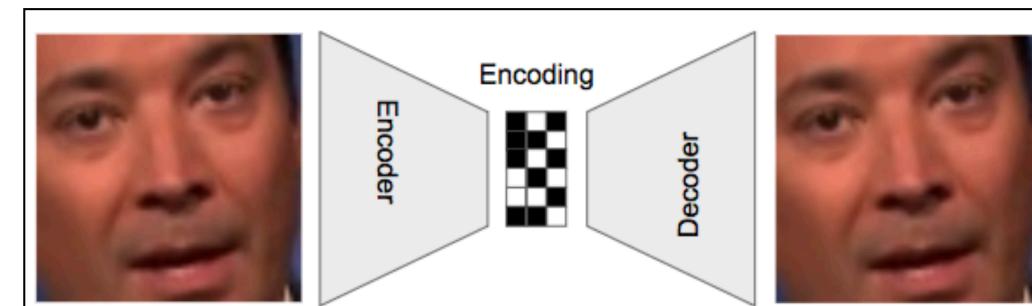
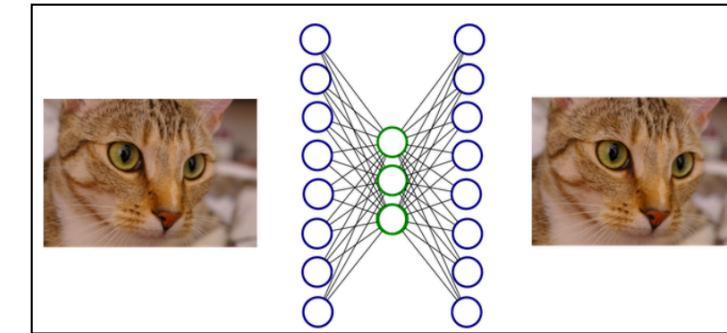


Can we learn the **true explanatory factors**, e.g. latent variables, from only observed data?

What is an Autoencoder?

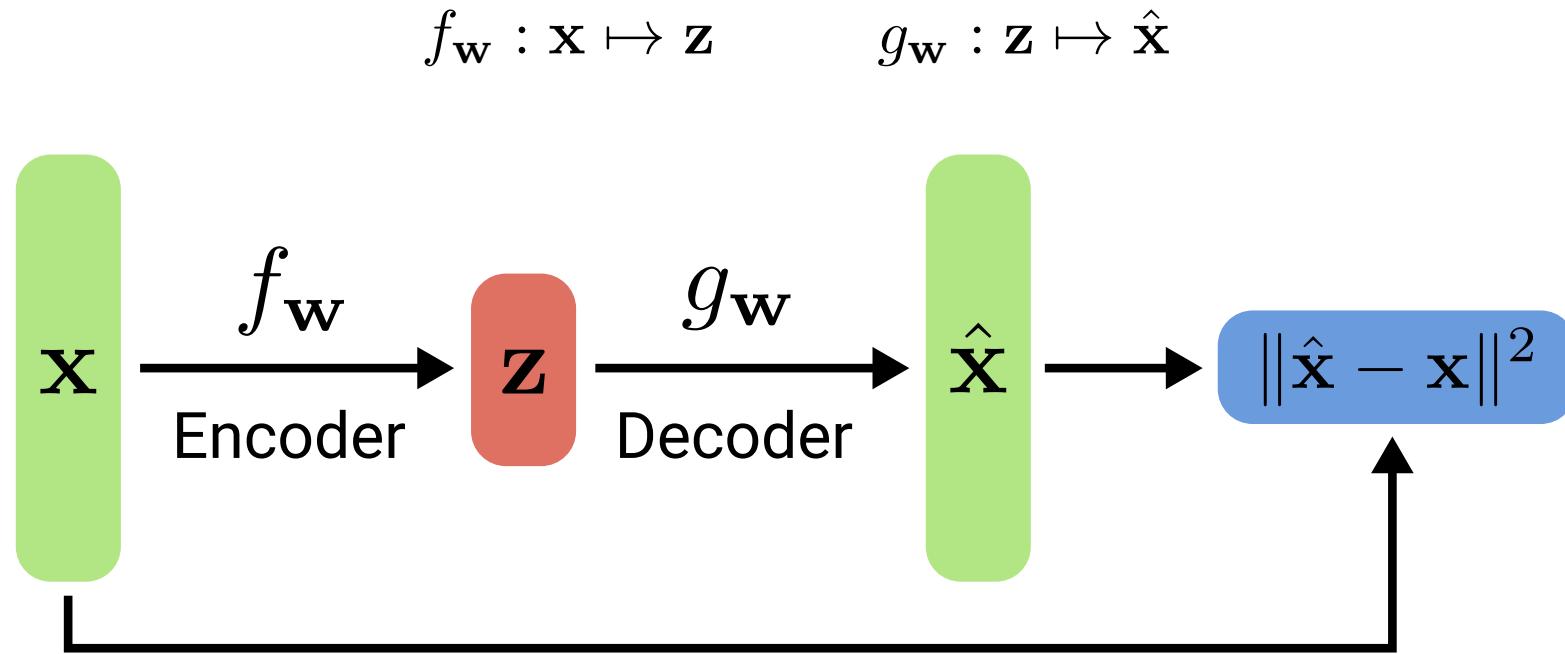
A neural network trained using unsupervised learning

- Trained to copy its input to its output



What is an Autoencoder?

Autoencoders comprise an **encoder** f_w as well as a **decoder** g_w :



- Models of this type are called **autoencoders** as they predict their input as output

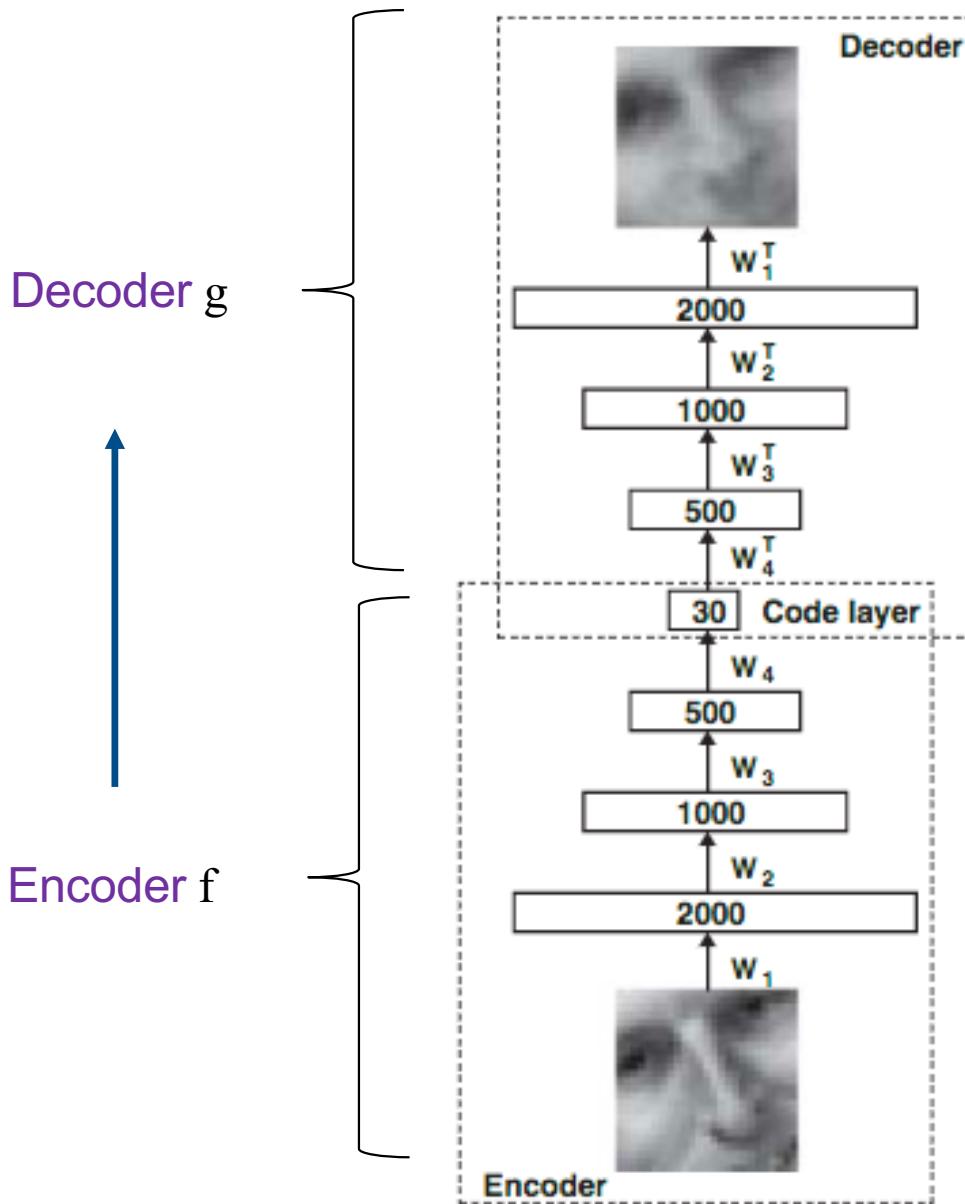
Auto encoders differ from General Data Compression

- Autoencoders are data-specific
 - i.e., only able to compress data similar to what they have been trained on
- This is different from, say, MP3 or JPEG compression algorithm
 - Which make general assumptions about "sound/images", but not about specific types of sounds/images
 - Autoencoder for pictures of cats would do poorly in compressing pictures of trees
 - Because features it would learn would be cat-specific
- Autoencoders are lossy
 - which means that the decompressed outputs will be degraded compared to the original inputs (similar to MP3 or JPEG compression).
 - This differs from lossless arithmetic compression
- Autoencoders are learnt

What does an Autoencoder Learn?

- Learning $g(f(x))=x$ everywhere is not useful
- Autoencoders are designed to be unable to copy perfectly
 - Restricted to copy only approximately
- Autoencoders learn useful properties of the data
 - Being forced to prioritize which aspects of input should be copied
- Maps an input x to an output r (called reconstruction) through an internal representation code h
 - It has a hidden layer h that describes a code used to represent the input
- The network has two parts
 - The encoder function $h=f(x)$
 - A decoder that produces a reconstruction $r=g(h)$

An autoencoder architecture



Undercomplete autoencoder

- Constrain h to have lower dimension than x
- Force it to capture most salient features of training data

Weights W are learnt using:

1. Training samples, and
2. a loss function

Autoencoder is a feed-forward non-recurrent neural net

- With an input layer, an output layer and one or more hidden layers
- Can be trained using the same techniques
 - Compute gradients using back-propagation
 - Followed by minibatch gradient descent

Autoencoder training using a loss function

- Encoder f and decoder g

$$f : X \rightarrow h$$

$$g : h \rightarrow X$$

$$\arg \min_{f,g} \|X - (f \circ g)X\|^2$$

- One hidden layer

- Non-linear encoder

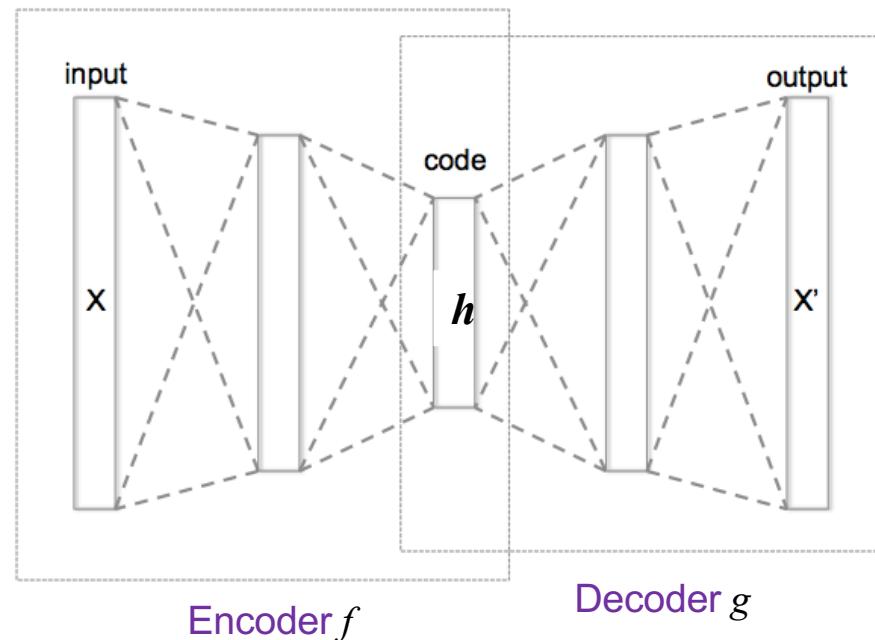
- Takes input $x \in R^d$

- Maps into output $h \in R^p$

$$h = \sigma_1(Wx + b)$$

$$x' = \sigma_2(W'h + b') \quad \sigma \text{ is an element-wise activation function such as sigmoid or Relu}$$

Autoencoder with 3 fully connected hidden layers



Trained to minimize reconstruction error (such as sum of squared errors)

$$L(x, x') = \|x - x'\|^2 = \|x - \sigma_2(W^t(\sigma_1(Wx + b)) + b')\|^2$$

Provides a compressed representation of the input x

Cases when Autoencoder Learning Fails

Where autoencoders fail to learn anything useful:

1. Capacity of encoder/decoder f/g is too high
 - Capacity controlled by depth
2. Hidden code h has dimension equal to input x
3. *Overcomplete* case: where hidden code h has dimension greater than input x
 - Even a linear encoder/decoder can learn to copy input to output without learning anything useful about data distribution

Right Autoencoder Design: Use regularization

Encoder/Decoder Capacity

If encoder f and decoder g are allowed too much capacity

- autoencoder can learn to perform the copying task without learning any useful information about distribution of data

Autoencoder with a one-dimensional code and a very powerful nonlinear encoder can learn to map $x^{(i)}$ to code i .

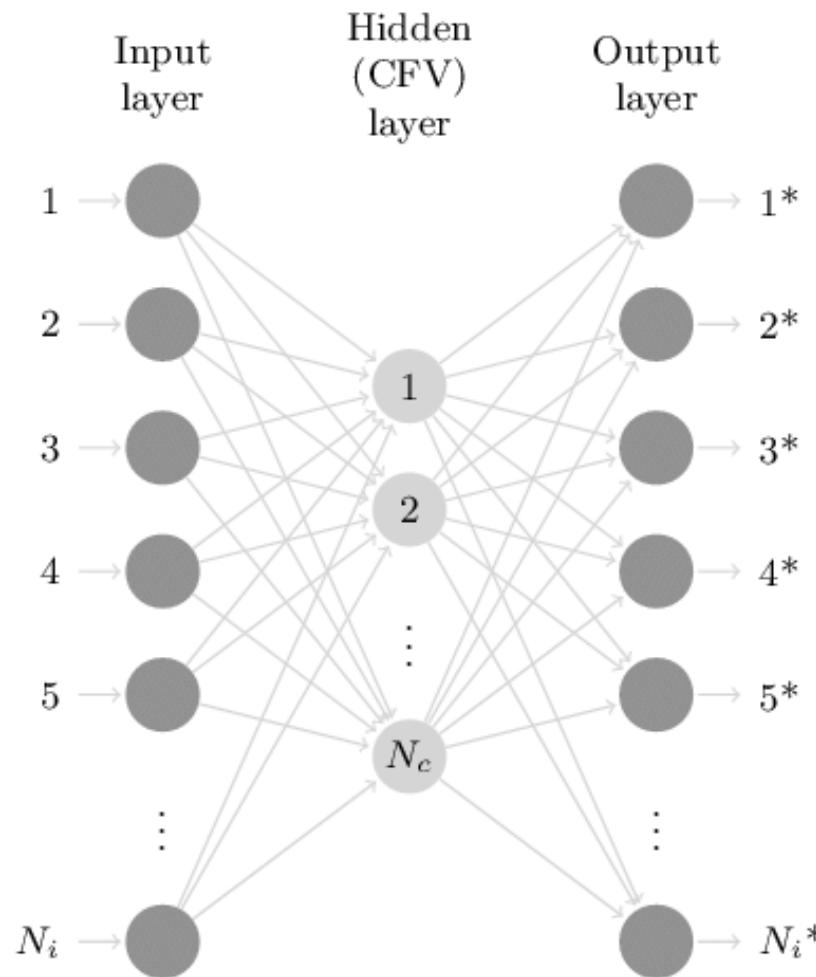
- The decoder can learn to map these integer indices back to the values of specific training examples

Autoencoder trained for copying task fails to learn anything useful if f/g capacity is too great

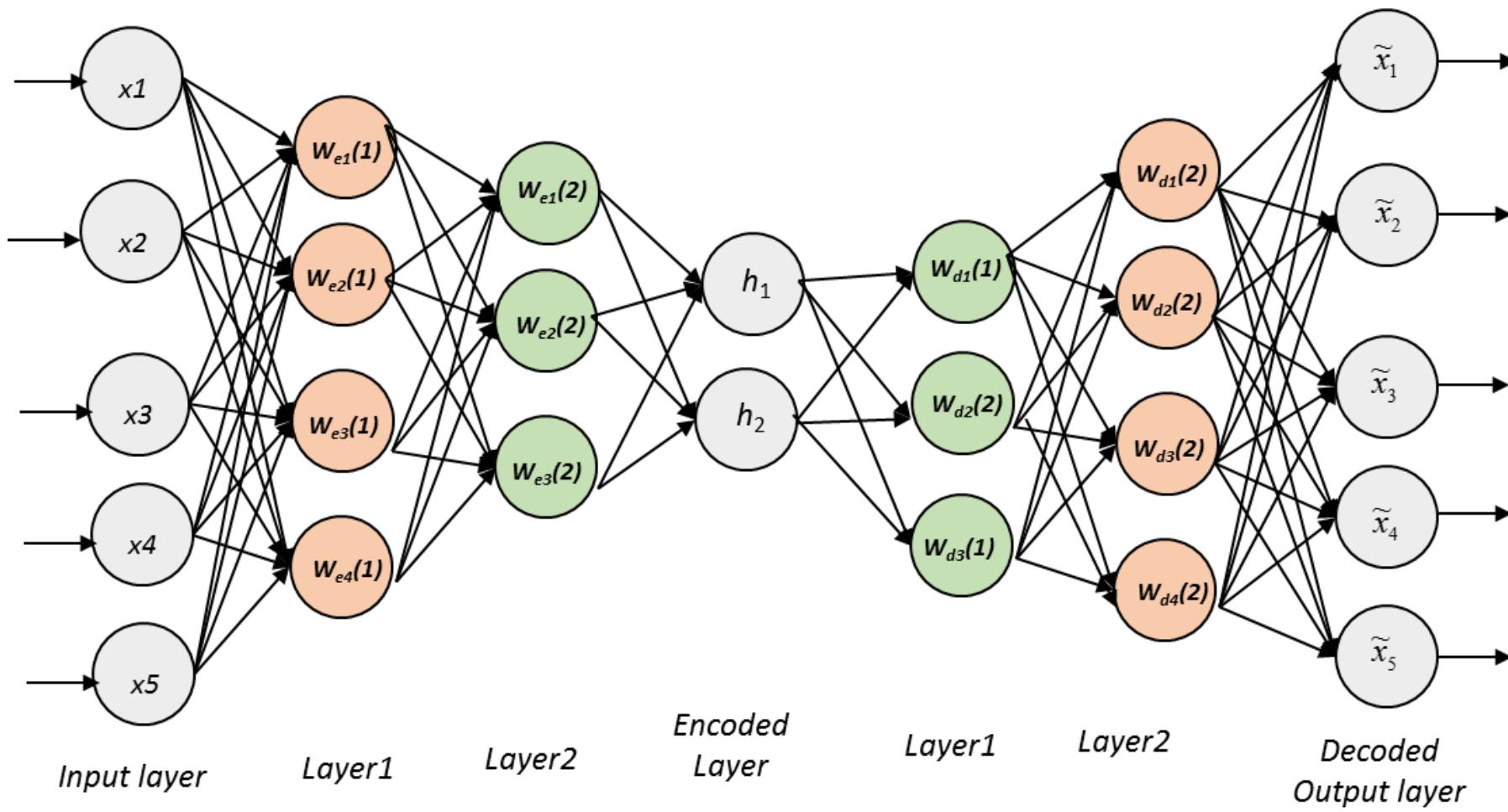
Types of Autoencoders

- Vanilla autoencoder
- Multilayer autoencoder
- Stacked autoencoder
- Deep autoencoder
- Convolutional autoencoder
- Regularized autoencoder

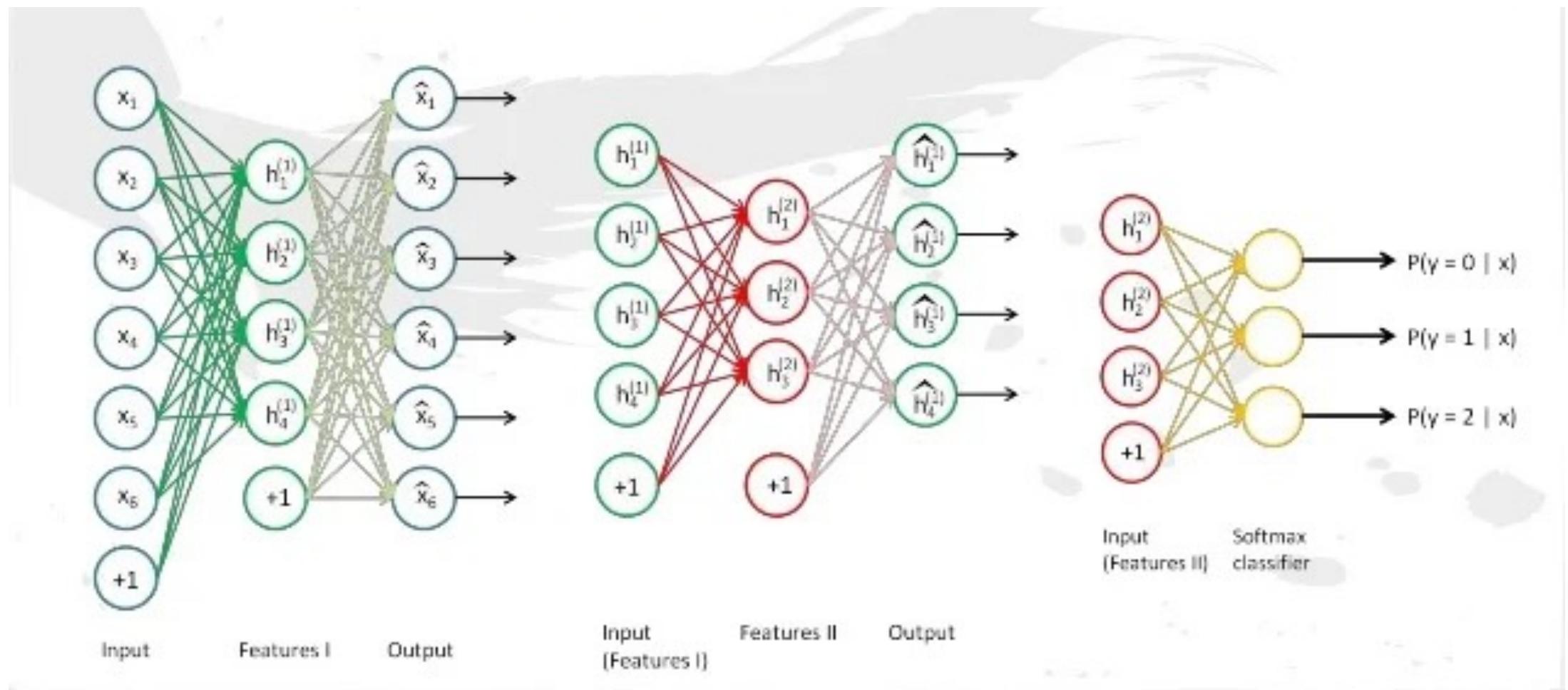
Vanilla autoencoder



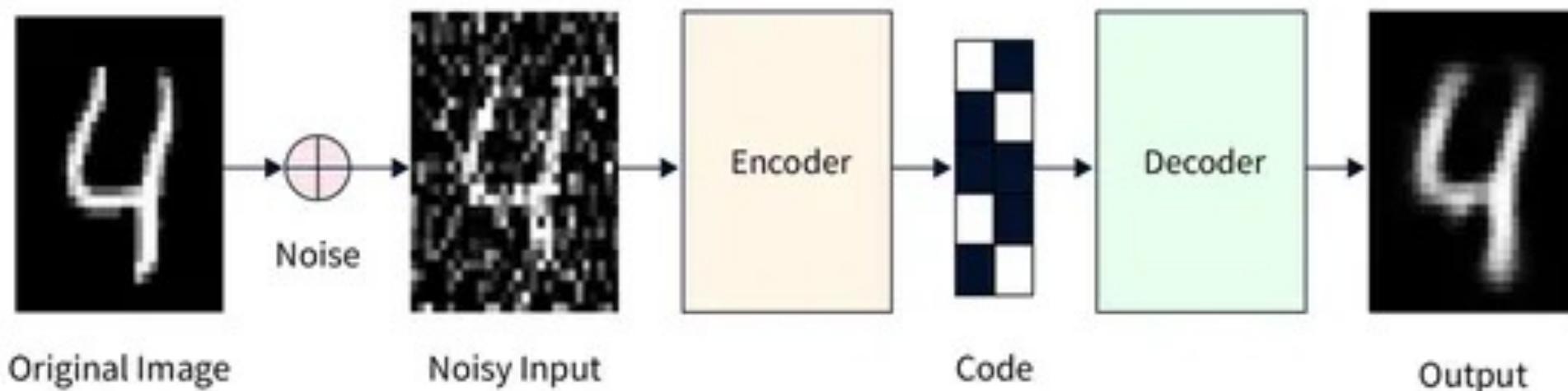
Multilayer Autoencoder



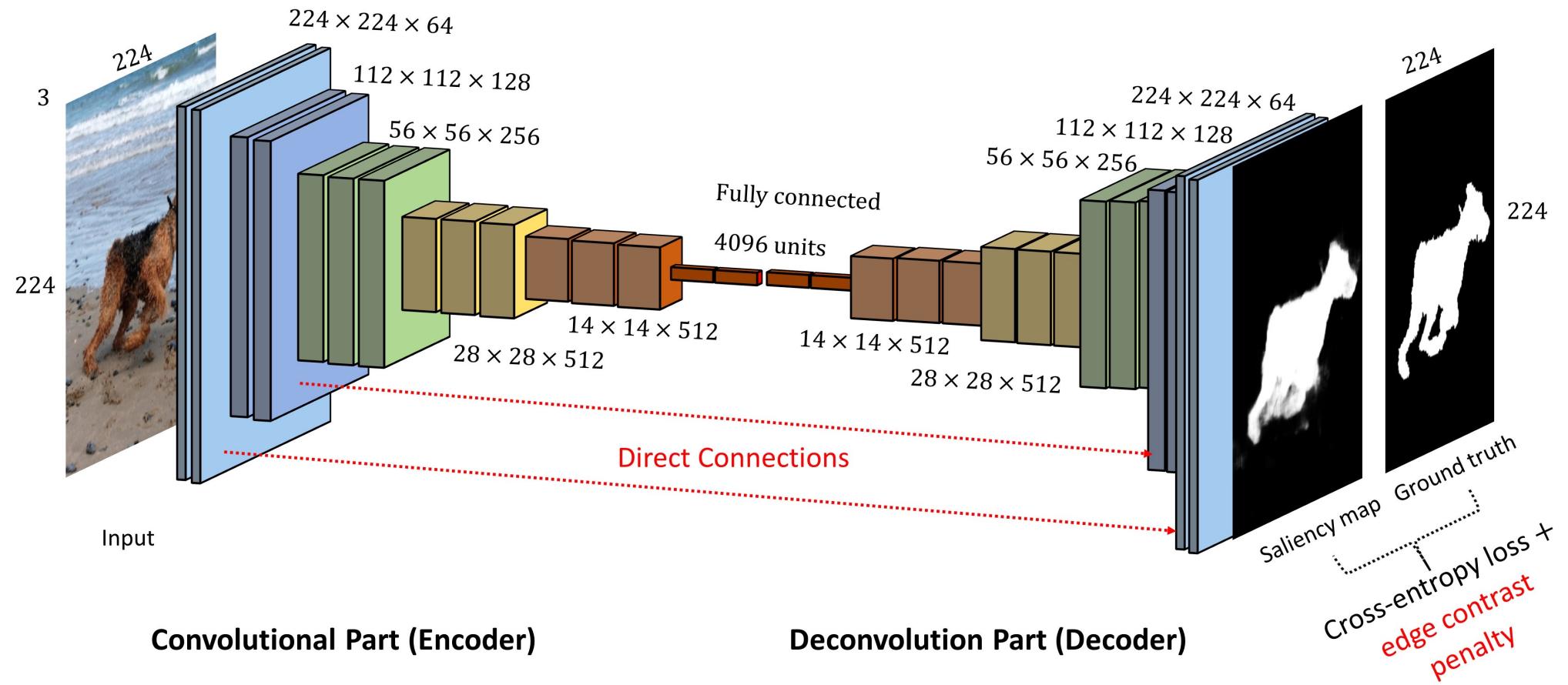
Stacked Autoencoder



Denoising Autoencoder

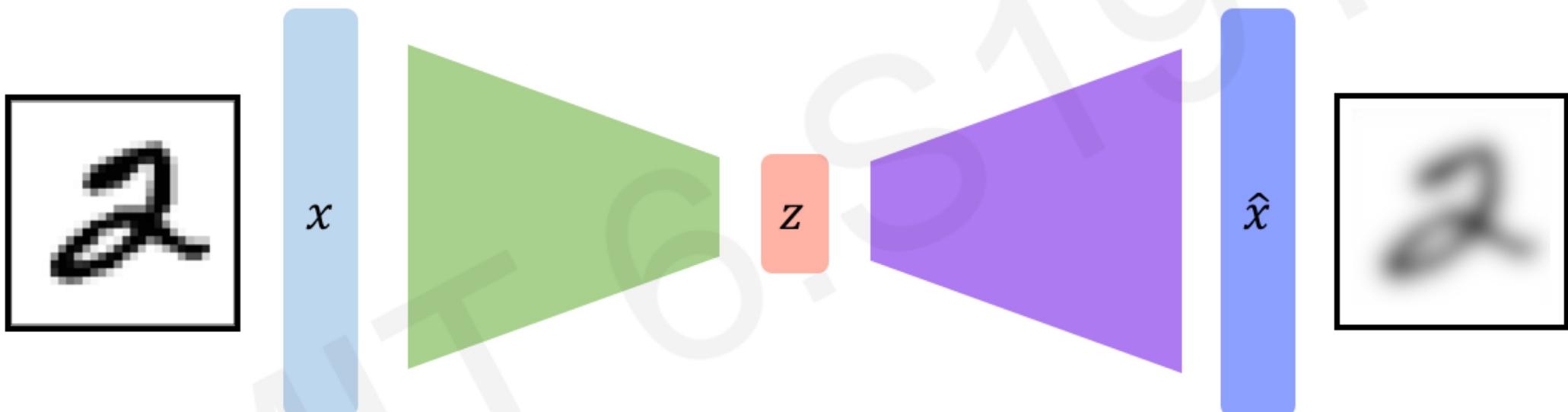


Convolutional Autoencoder

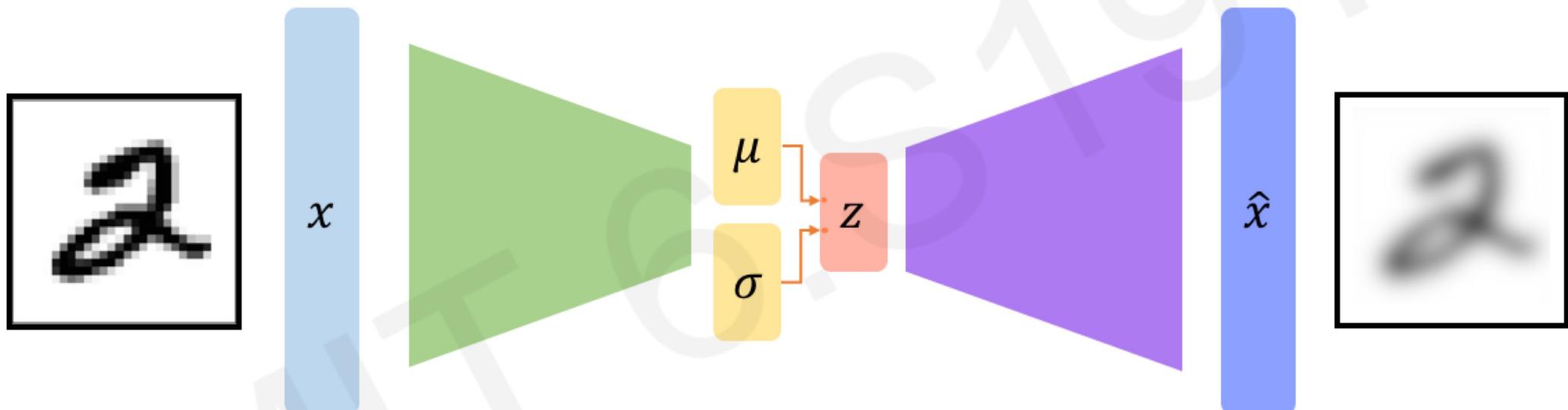


Variational Autoencoders (VAEs)

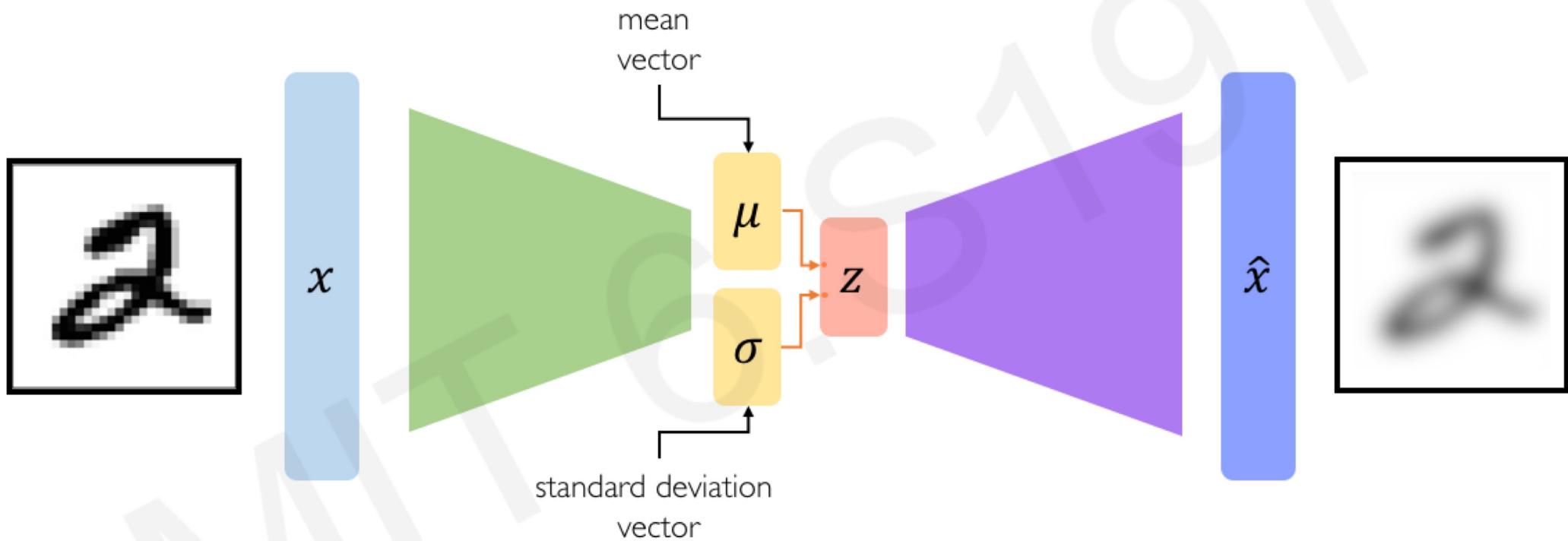
Traditional autoencoders



VAEs: key difference with traditional autoencoder



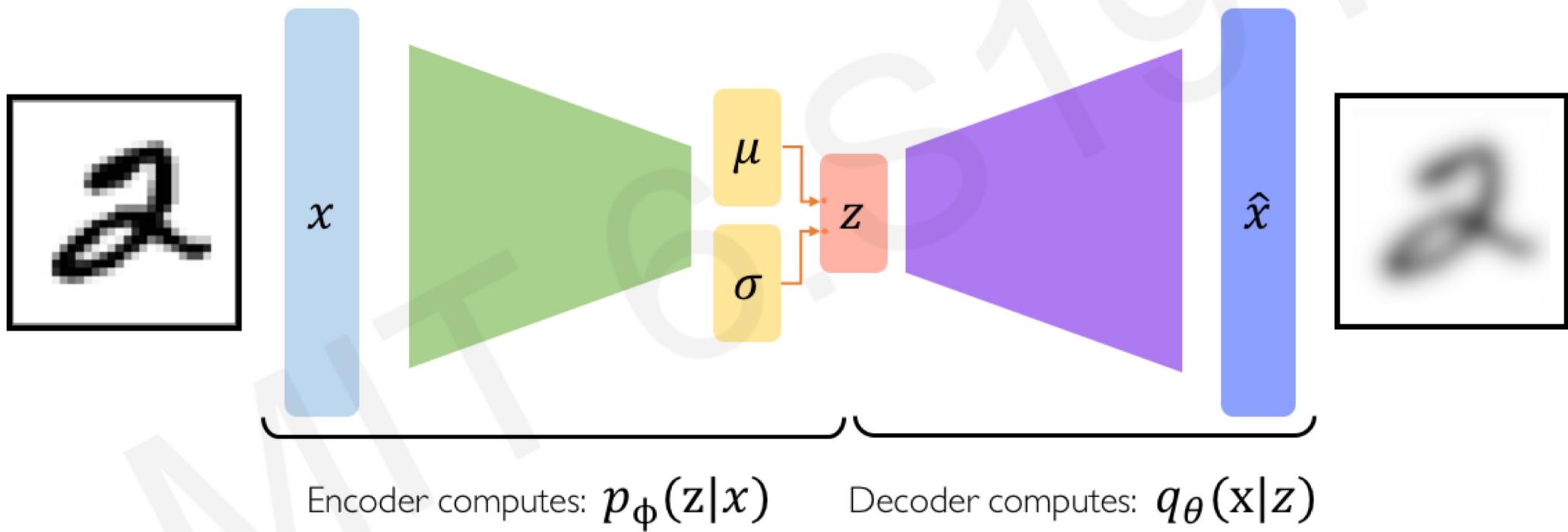
VAEs: key difference with traditional autoencoder



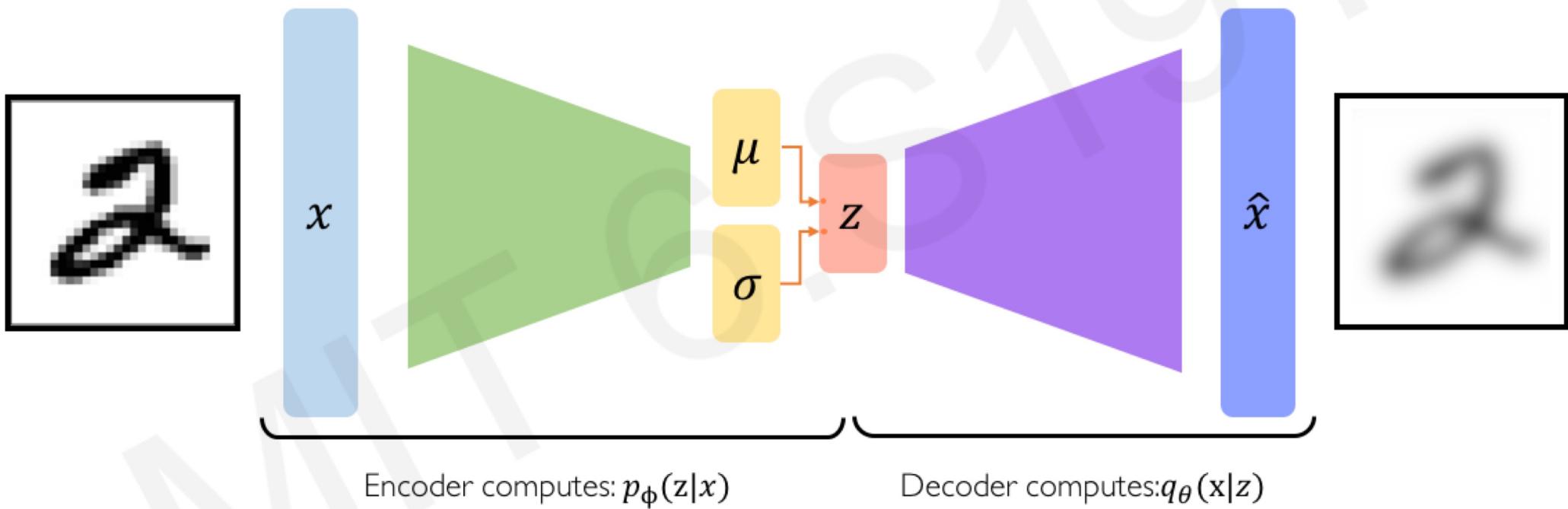
Variational autoencoders are a probabilistic twist on autoencoders!

Sample from the mean and standard dev. to compute latent sample

VAE optimization

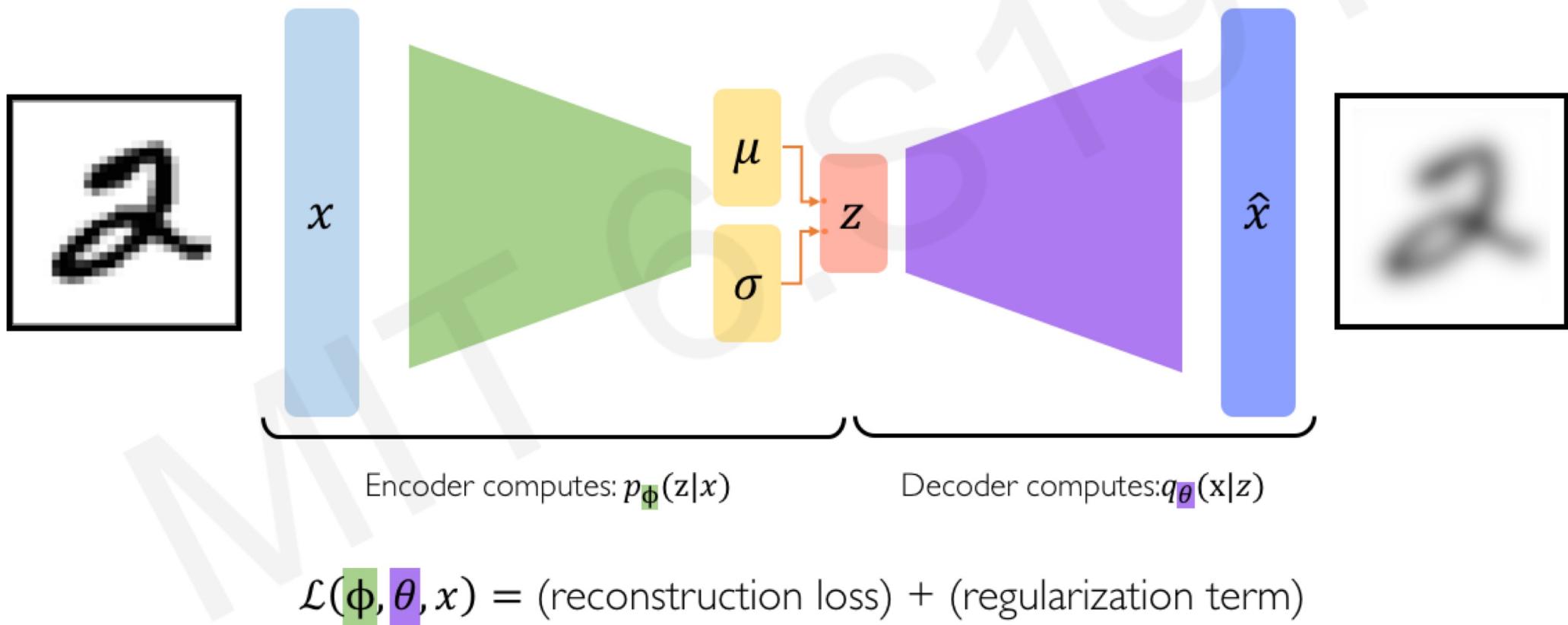


VAE optimization

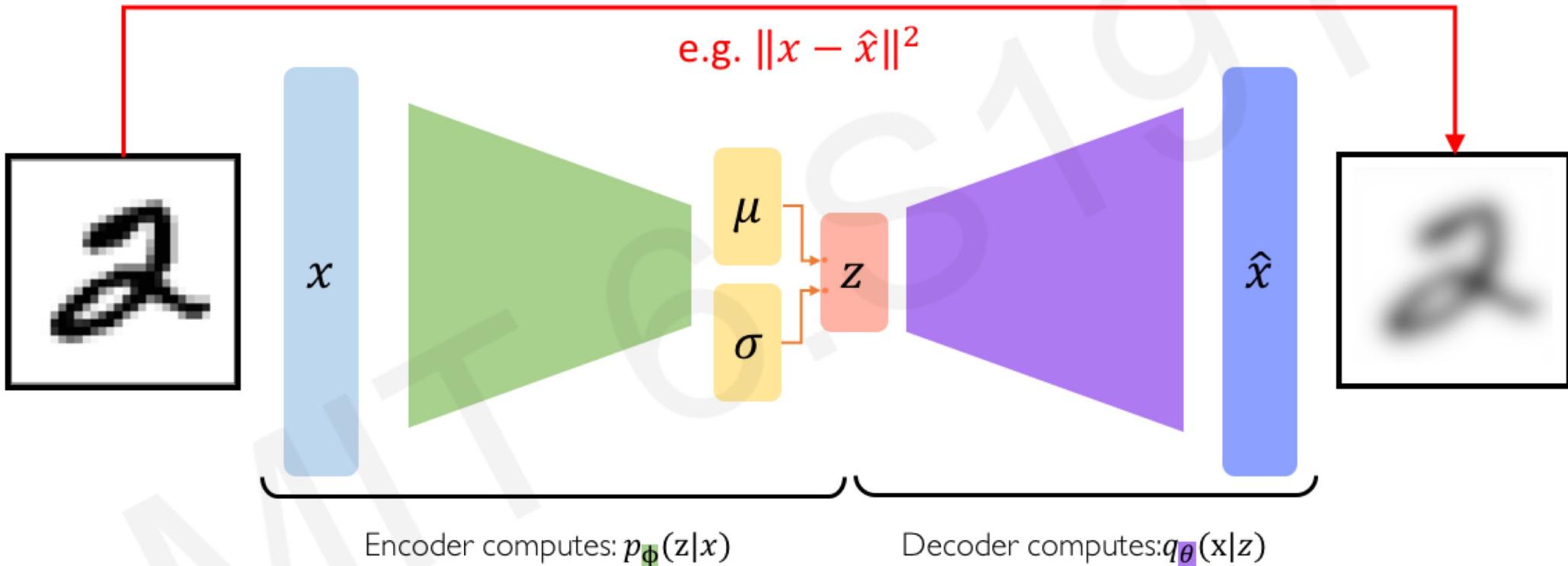


$$\mathcal{L}(\phi, \theta) = (\text{reconstruction loss}) + (\text{regularization term})$$

VAE optimization

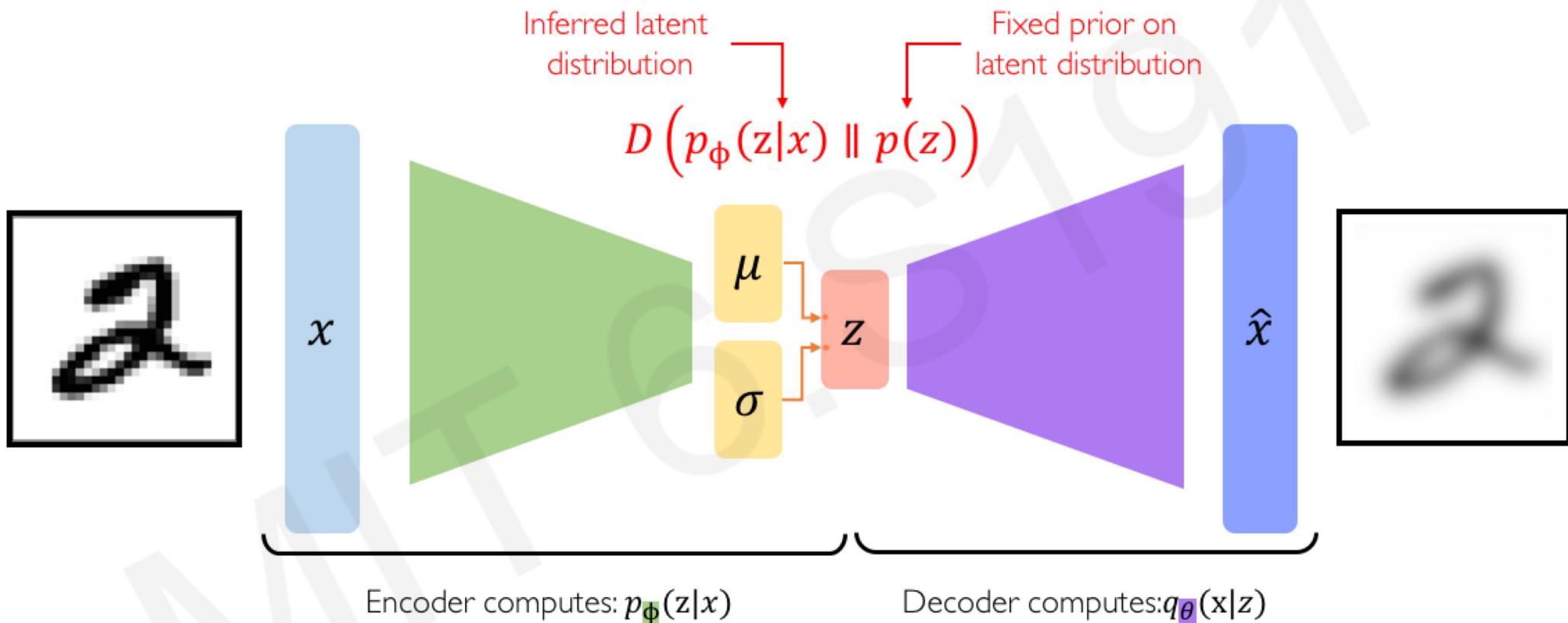


VAE optimization



$$\mathcal{L}(\phi, \theta, x) = \text{(reconstruction loss)} + \text{(regularization term)}$$

VAE optimization

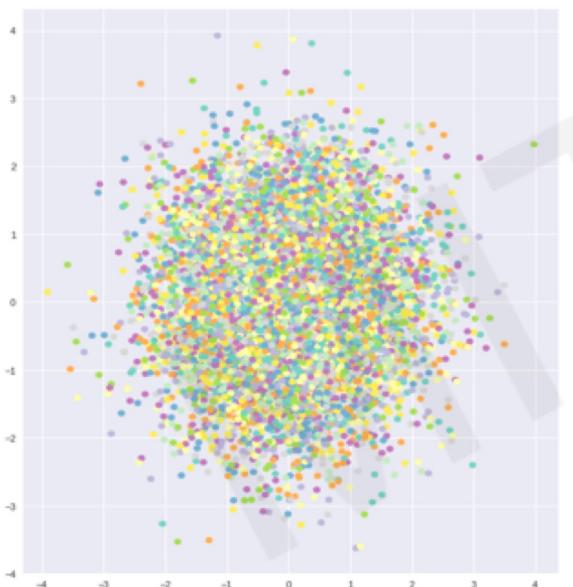


$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + \boxed{(\text{regularization term})}$$

Priors on the latent distribution

$$D(p_{\phi}(z|x) \parallel p(z))$$

Inferred latent distribution ↑
Fixed prior on latent distribution ↑



Common choice of prior:

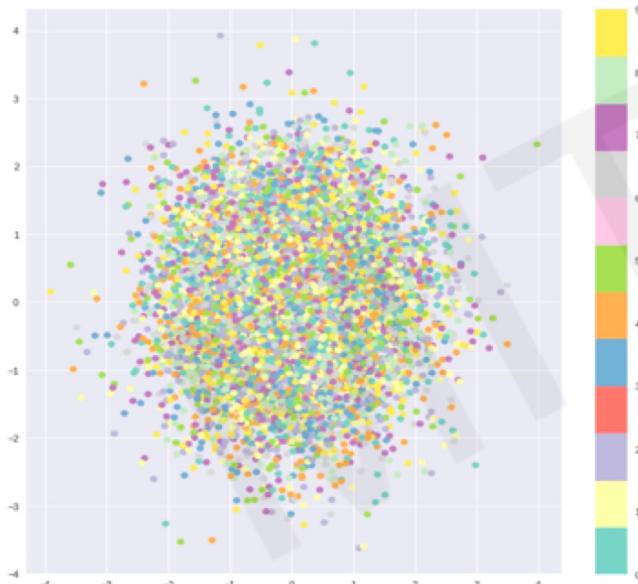
$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to “cheat” by clustering points in specific regions (ie. memorizing the data)

Priors on the latent distribution

$$D(p_{\phi}(z|x) \parallel p(z)) \\ = -\frac{1}{2} \sum_{j=0}^{k-1} (\sigma_j + \mu_j^2 - 1 - \log \sigma_j)$$

KL-divergence between
the two distributions

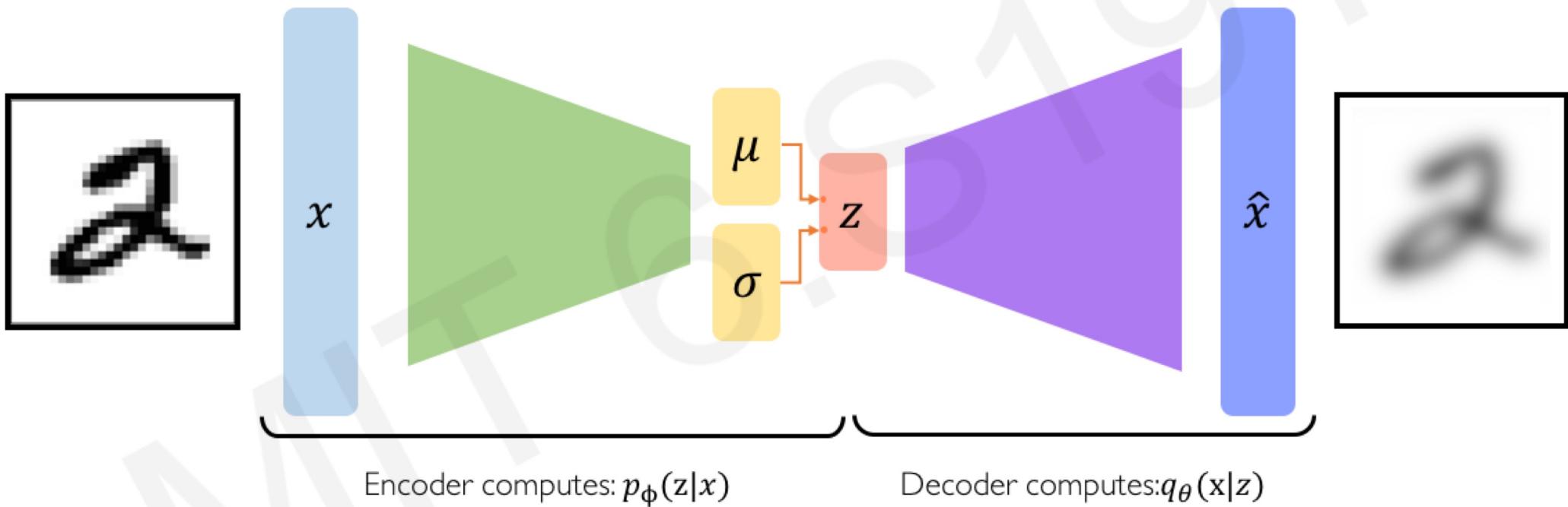


Common choice of prior:

$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to “cheat” by clustering points in specific regions (ie. memorizing the data)

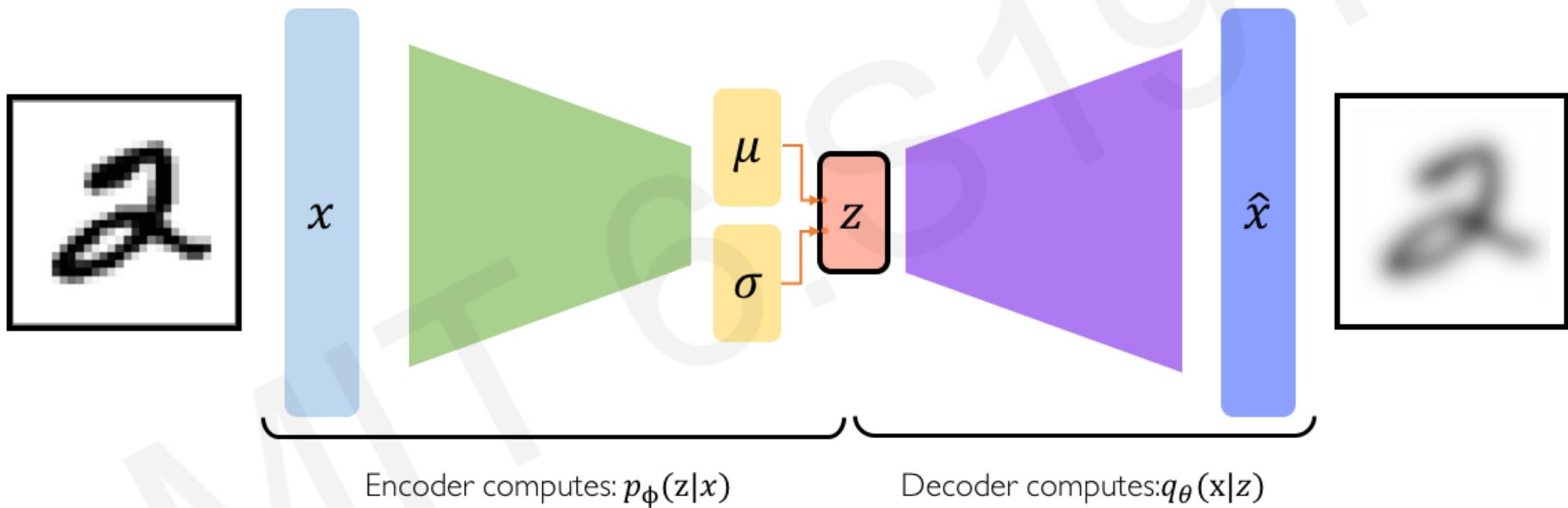
VAEs computation graph



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

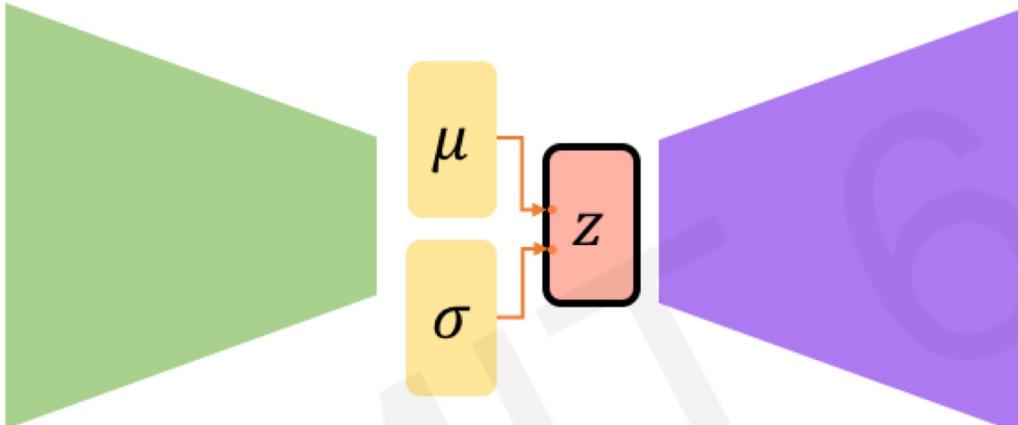
VAEs computation graph

Problem: We cannot backpropagate gradients through sampling layers!



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

Reparametrizing the sampling layer



Key Idea:

$$z \sim \mathcal{N}(\mu, \sigma^2)$$

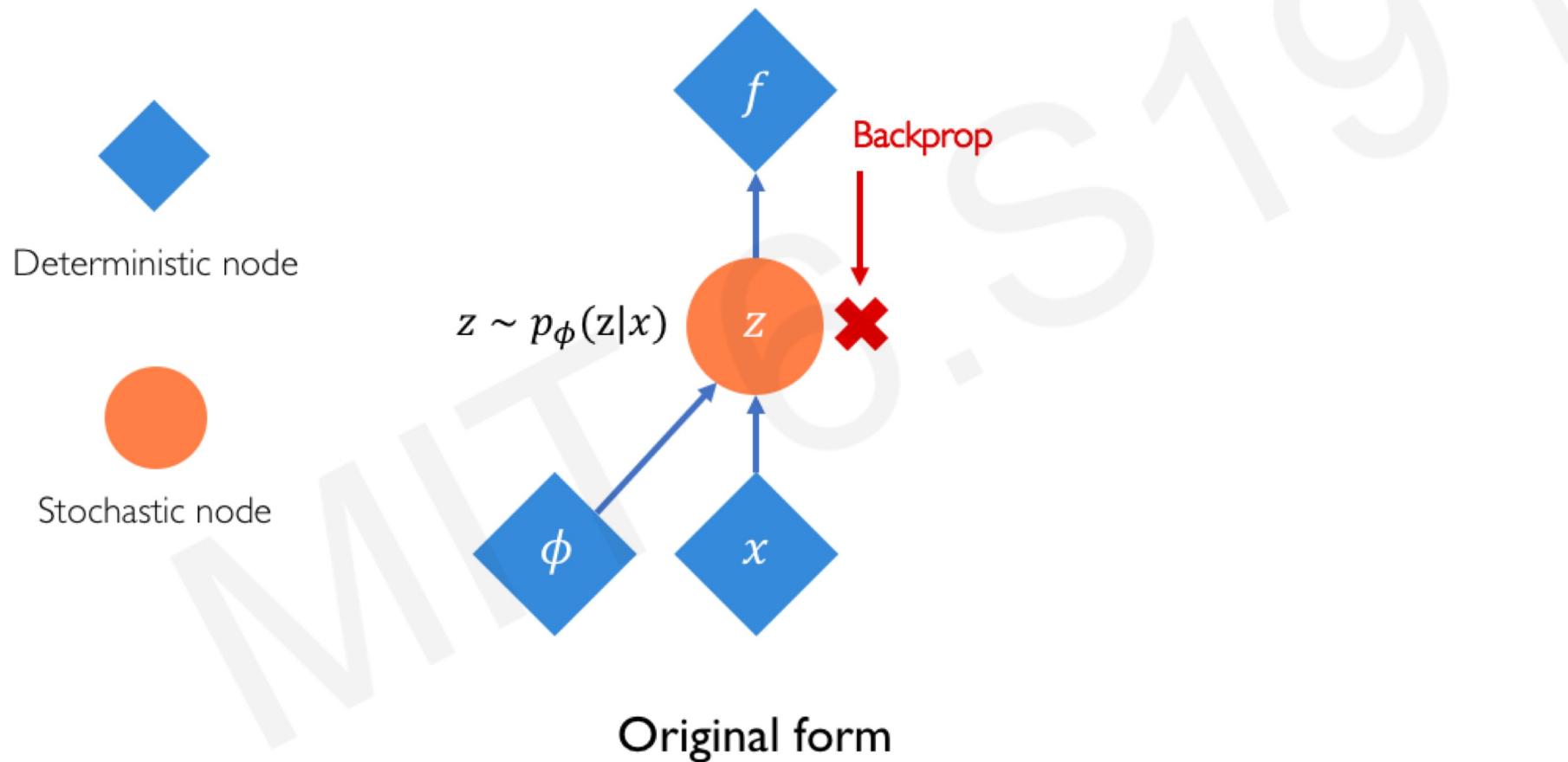
Consider the sampled latent vector z as a sum of

- a fixed μ vector,
- and fixed σ vector scaled by random constants drawn from the prior distribution

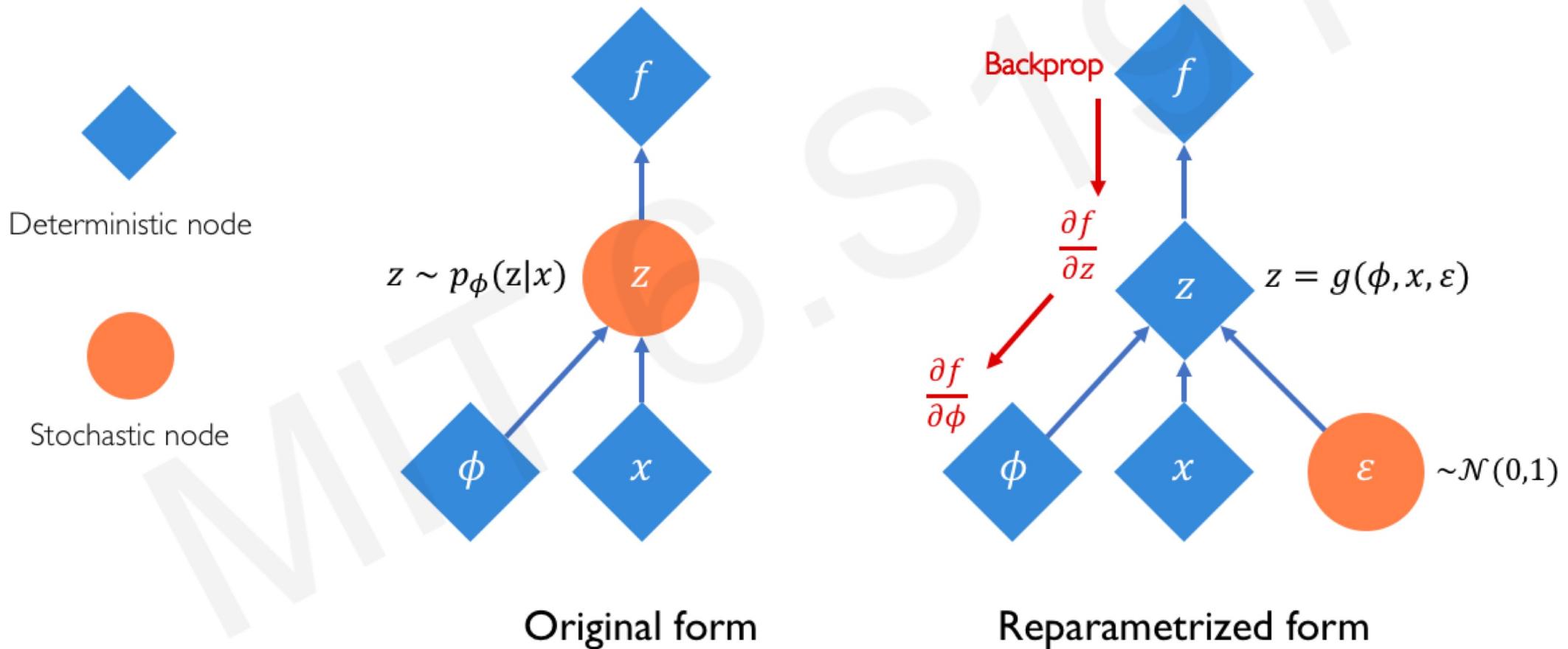
$$\Rightarrow z = \mu + \sigma \odot \varepsilon$$

where $\varepsilon \sim \mathcal{N}(0,1)$

Reparametrizing the sampling layer



Reparametrizing the sampling layer



VAEs: Latent perturbation

Slowly increase or decrease a **single latent variable**

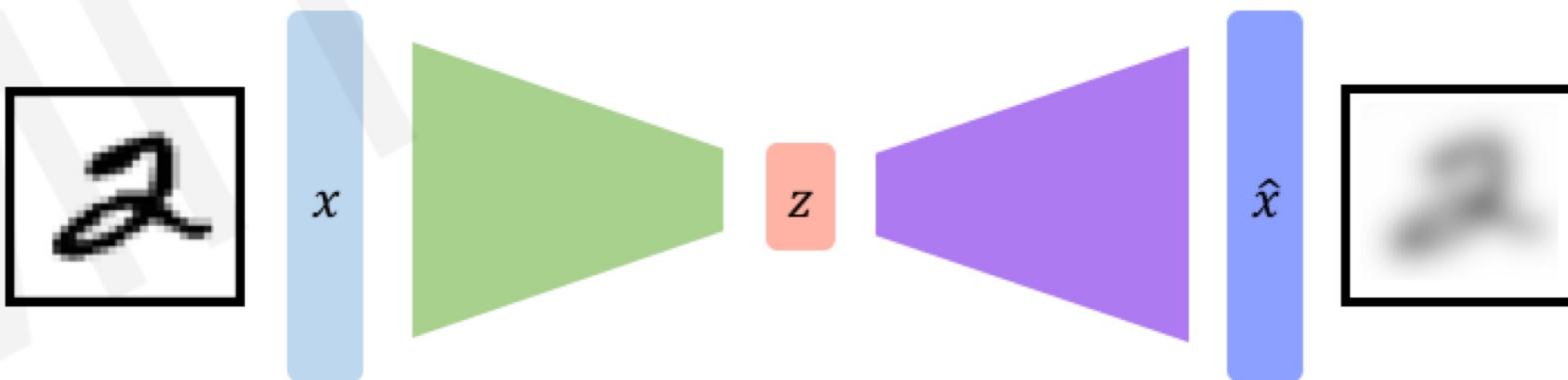
Keep all other variables fixed



Different dimensions of z encodes **different interpretable latent features**

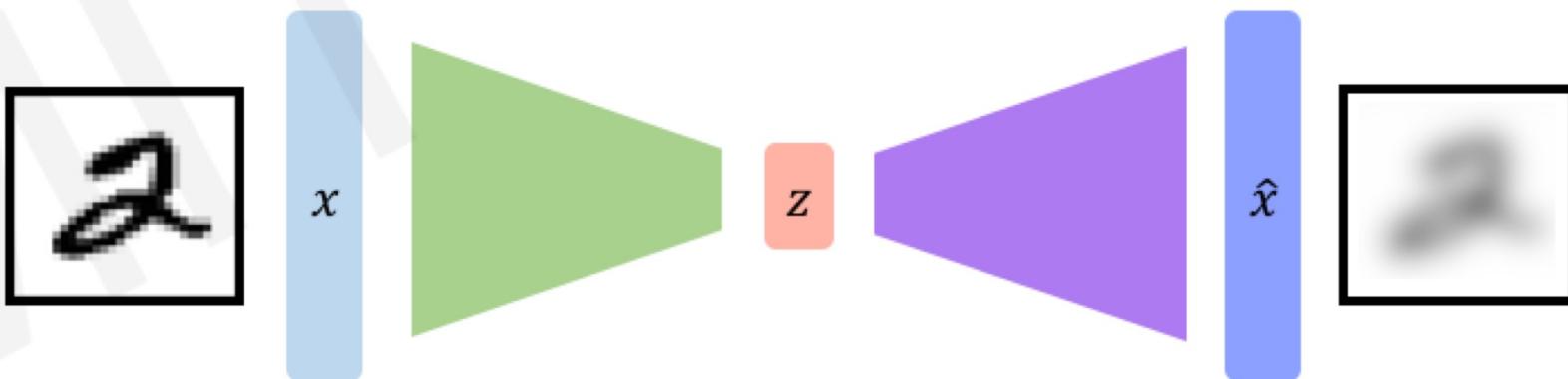
VAE summary

- Compress representation of world to something we can use to learn



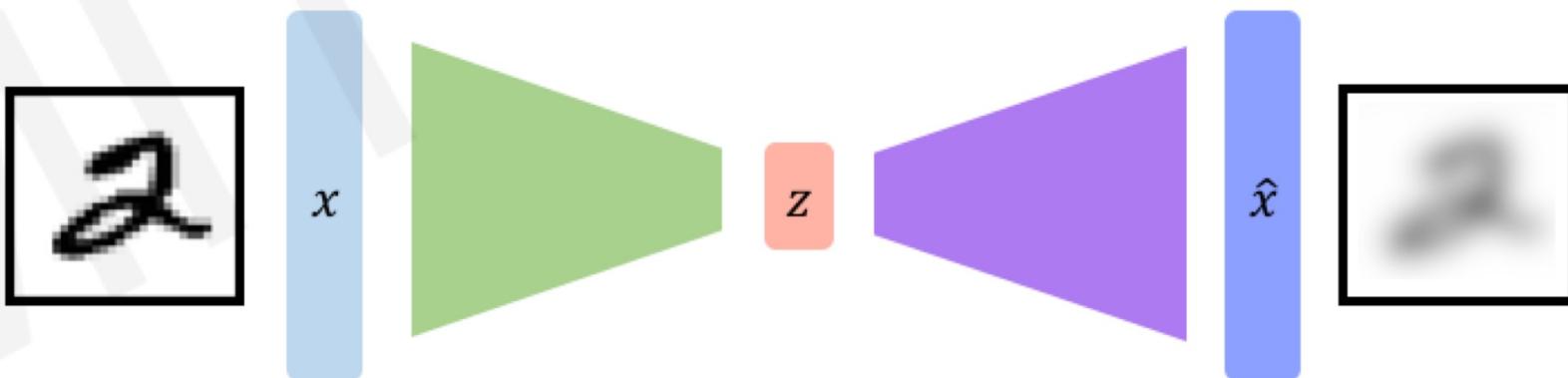
VAE summary

- . Compress representation of world to something we can use to learn
- . Reconstruction allows for unsupervised learning (no labels!)



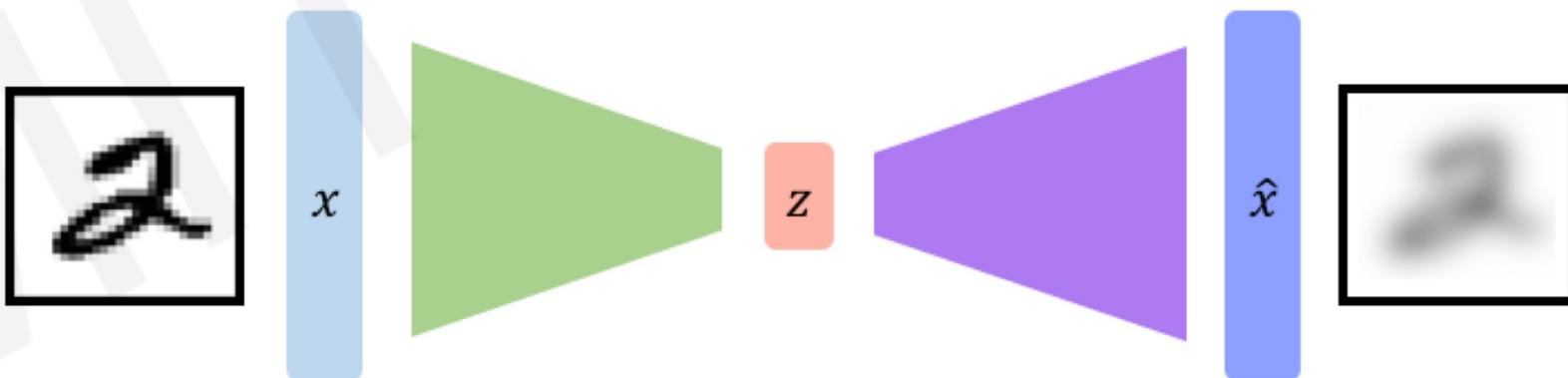
VAE summary

- . Compress representation of world to something we can use to learn
- . Reconstruction allows for unsupervised learning (no labels!)
- . Reparameterization trick to train end-to-end



VAE summary

- . Compress representation of world to something we can use to learn
- . Reconstruction allows for unsupervised learning (no labels!)
- . Reparameterization trick to train end-to-end
- . Interpret hidden latent variables using perturbation



VAE summary

- . Compress representation of world to something we can use to learn
- . Reconstruction allows for unsupervised learning (no labels!)
- . Reparameterization trick to train end-to-end
- . Interpret hidden latent variables using perturbation
- . Generating new examples

