# Deep Learning

Image Generation: Autoencoders and VAEs

Eric Ewing Tuesday, 10/28

**CSCI 1470** 



Source: ChatGPT image generator, prompt "Ok, hear me out. I need a friendly cartoon mole person with a hat and magnifying glass looking at the earth's core all in one image."

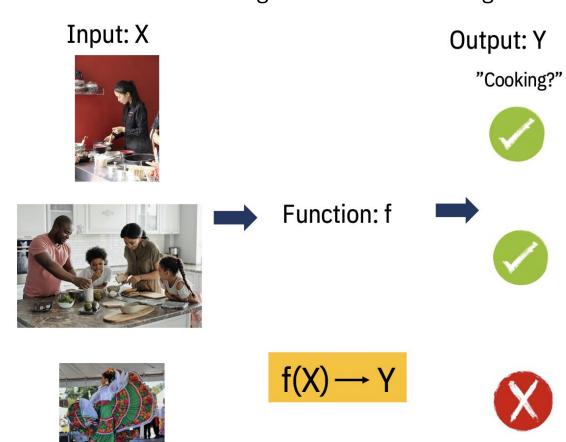
## Logistics

- Final Project groups will come out this afternoon (or evening)
  - If you haven't filled out the form, you should fill it out even if you don't have group-mates in mind
- Weekly Quiz is out

## Recap: Supervised Learning

- Supervised learning requires labels
- Learn a function that takes in input features and outputs labels

Is it an image of someone cooking?



What are some pros and cons of supervised learning?

## Recap: Supervised Learning

#### Pros:

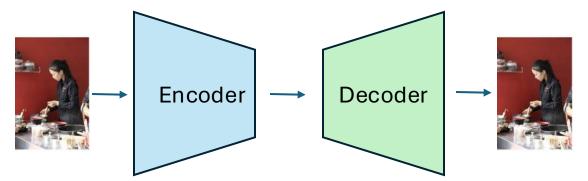
- Can produce very high quality models with sufficient data
- Performance is easy to measure (e.g., accuracy)

#### Cons:

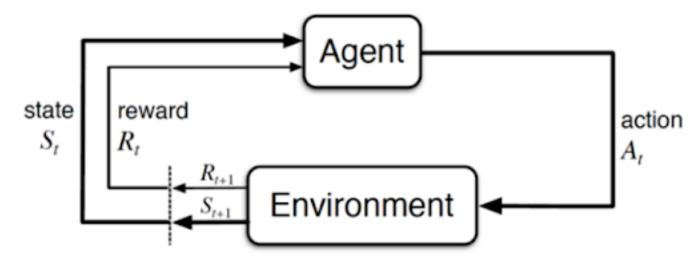
- Reliant on availability of labels
- Reliant on quality of labels

#### What's Left?

Unsupervised Learning: Learning without labels



Reinforcement Learning: Learning from experiences



## Image Capabilities in LLMs

## How might we incorporate other capabilities into our LLMs?

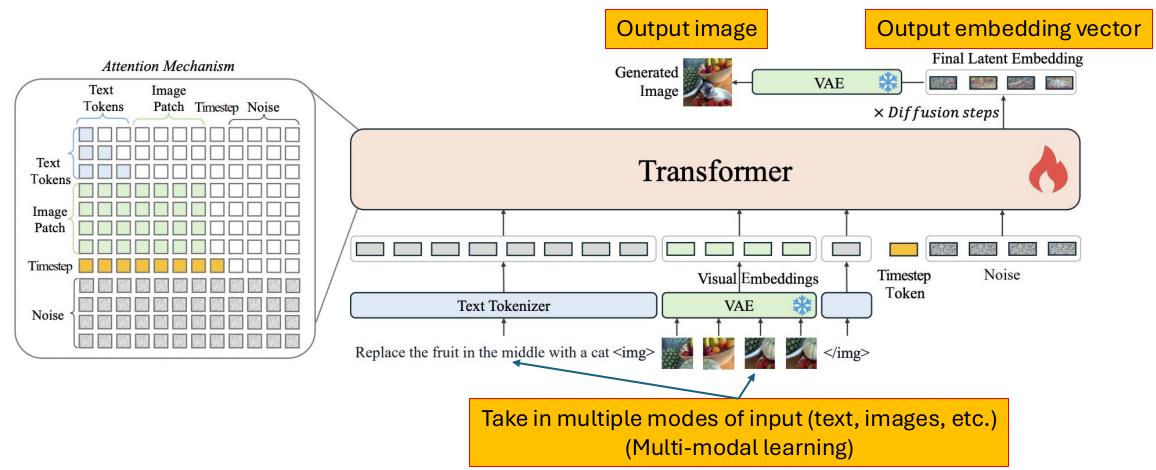
- 1. Generate Images (call other model)
- 2. Take images as input and generate text

```
from openai import OpenAI
import base64
client = OpenAI()
response = client.responses.create(
    model="gpt-5",
    input="Generate an image of gray tabby cat hugging an otter with an orange scarf"
    tools=[{"type": "image_generation"}],
# Save the image to a file
image_data = [
    output.result
    for output in response.output
    if output.type == "image_generation_call"
if image data:
    image_base64 = image_data[0]
    with open("otter.png", "wb") as f:
        f.write(base64.b64decode(image base64))
```

Give this cat a detective hat and a monocle



## What *might* this look like?



Source: OmniGen: Unified Image Generation: https://arxiv.org/pdf/2409.11340

## What *might* this look like?

VAE: Variational Auto-Encoder Noise: for "Diffusion" process What's this? Final Latent Embedding Attention Mechanism Generated VAE Text Image **Image Tokens** Patch Timestep Noise  $\times$  Diffusion steps Transformer Text Tokens Image Patch Timestep Noise Timestep Visual Embeddings Token **Text Tokenizer** And why is VAE Noise there noise? Replace the fruit in the middle with a cat <img> And This?

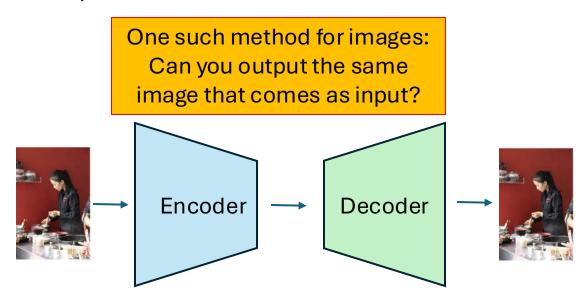
Source: OmniGen: Unified Image Generation: https://arxiv.org/pdf/2409.11340

### Foundation Models

Key Question: What is the equivalent of Language Modeling for data other than natural language? (From last class)

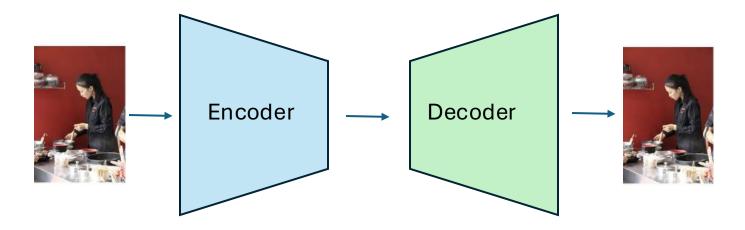
#### **Desired Properties:**

- 1. No need for human labeling
- 2. Large amount of data available
- 3. Ability to learn a "general" representation of the data



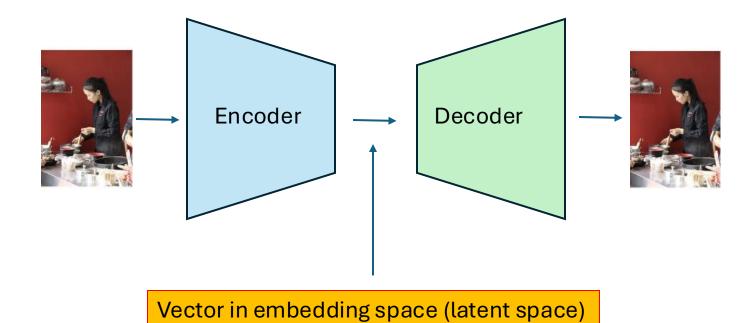
### Autoencoders

Autoencoder: an encoder-decoder architecture that tries to produce its own input



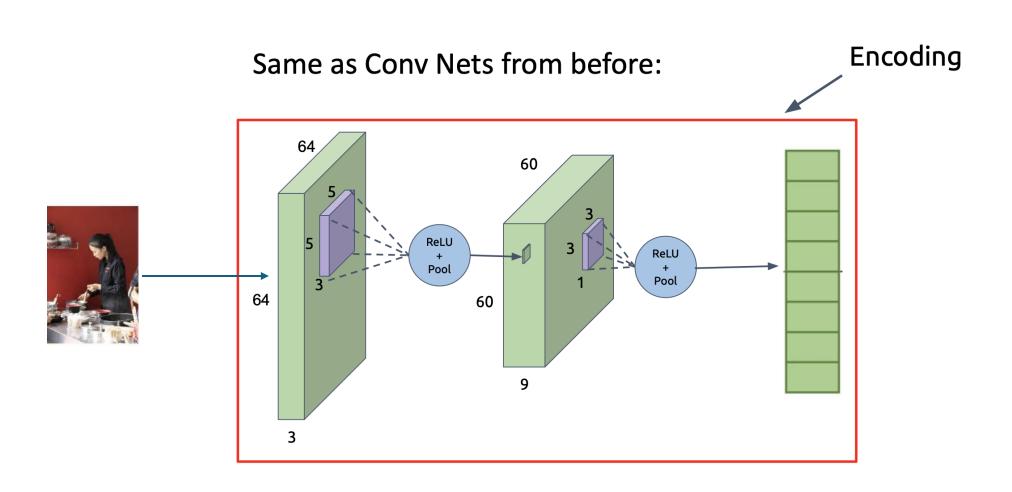
But what's hard about this? It's very easy to learn a function that outputs the input to the function (i.e., the identity function)

## Autoencoders



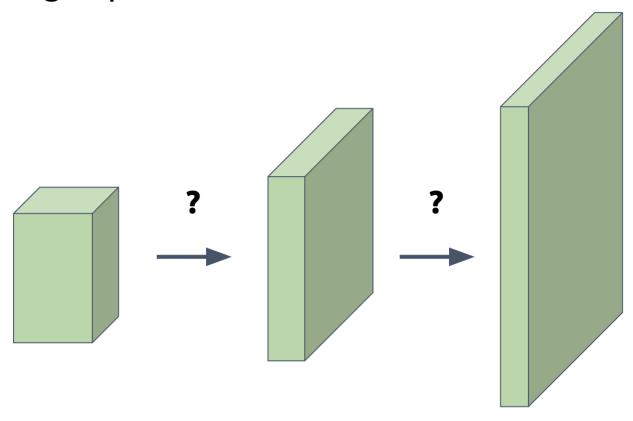
Much smaller than the original input size

## Convolutional Autoencoders: Encoding



## Autoencoders: Decoding

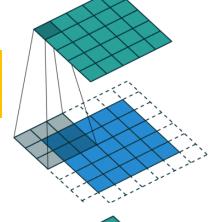
- Convolution as we know it only keeps resolution same or decreases it
- How do we go up in resolution?



## **Transposed Convolution**

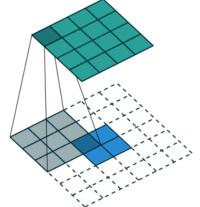
Standard Convolution: Kernel slides along input

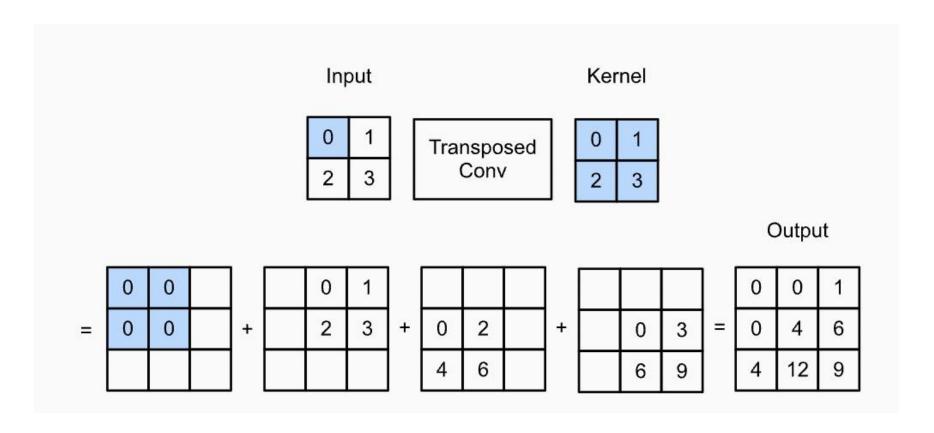
Green: output Blue: input



Transpose Convolution: Single input is multiplied by entire kernel, summed

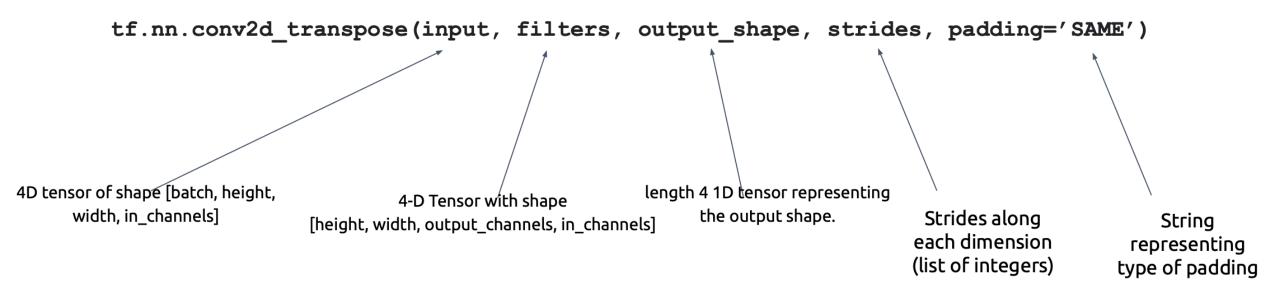
Green: Input Blue: Output





(stride of 1)

## **Transpose Convolution in Tensorflow**



# **Specifying Output Size**

 An image can be the result of the same convolution on images of different resolution

We need to specify which one we want.

57	60	
66	61	

1	2	3
4	5	6
7	8	9

Kernel

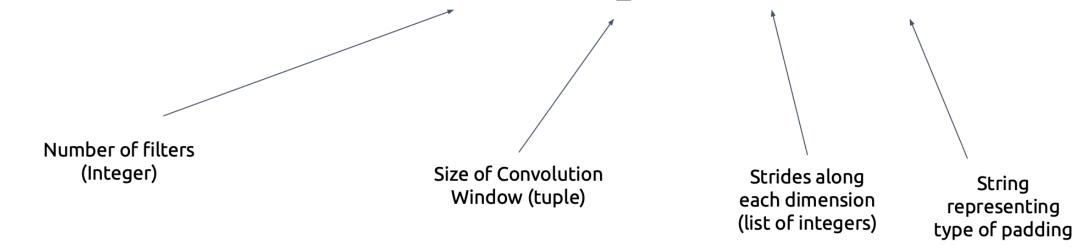
2	1	0	3
0	0	1	2
3	1	2	0
0	2	2	1

2	1	0	3	0
0	0	1	2	0
3	1	2	0	0
0	2	2	1	0
0	0	0	0	0



## **Transpose Convolution in Keras**

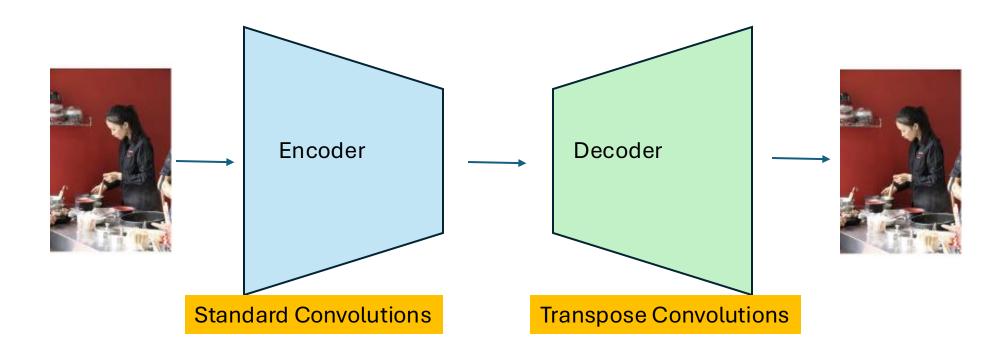
tf.keras.layers.Conv2DTranspose(filters, kernel\_size, strides, padding='SAME'



Note: Output Shape is inferred, but can be specified via the "output\_padding" parameter

Documentation here: <a href="https://www.tensorflow.org/api\_docs/python/tf/keras/layers/Conv2DTranspose">https://www.tensorflow.org/api\_docs/python/tf/keras/layers/Conv2DTranspose</a>

## Convolutional Autoencoder



#### **Loss Function**

What do you think is an appropriate loss function?

Reconstruction loss (MSE): How far is each output pixel from the corresponding input pixel?

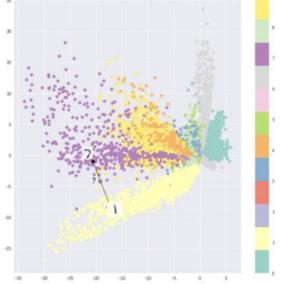
## Autoencoders

What do Autoencoders actually learn?

- 1. Encoder learns a dimensionality reduction (from image to vector)
- 2. Decoder learns an image generation function (from vector to image)

Encodings of MNIST data points with a trained autoencoder (dimensionality reduced further by PCA)

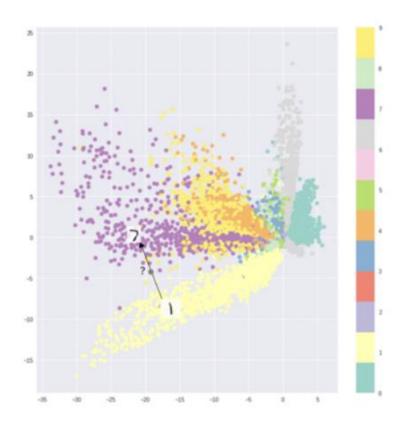
Encoders can be used to learn insights into structure of data

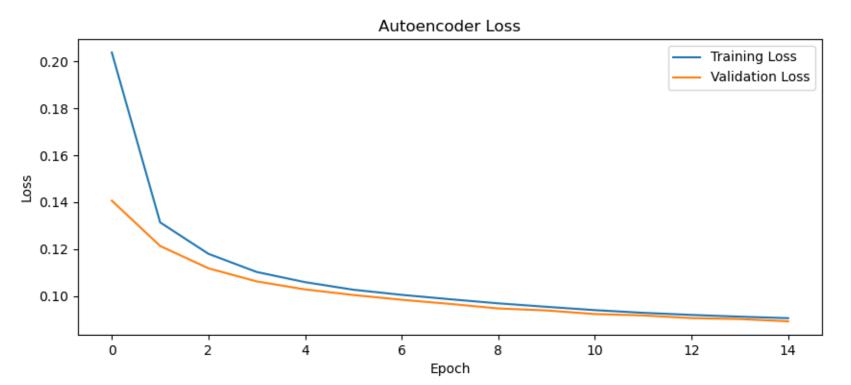


Decoders can be used to generate "new" images

## Generating Images

- How can we generate a "new" image using a decoder?
- Sample a vector in latent space and send it to the decoder...
- But how do you choose which vector?
- What if you wanted to generate a specific image? How would you find the right vector?

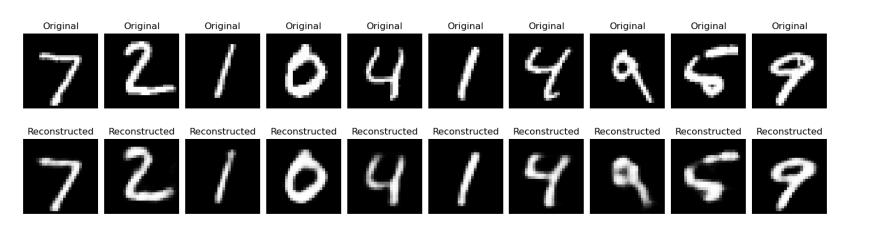


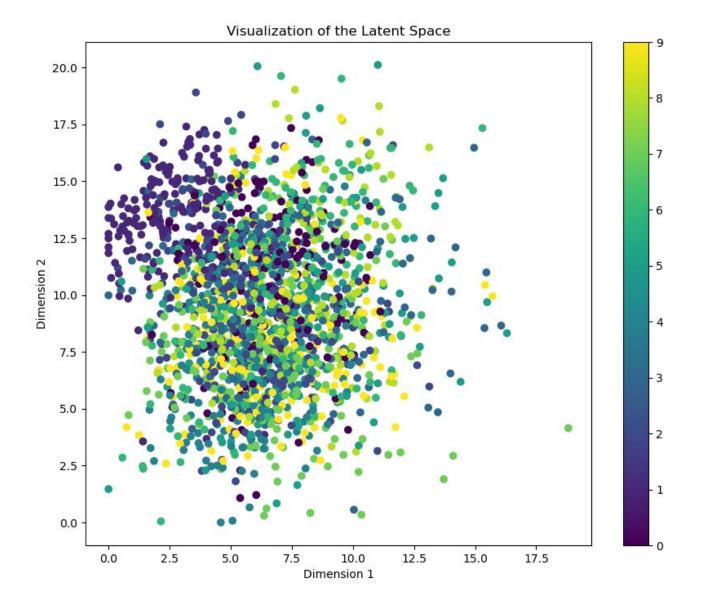


(Encoding size = 32)

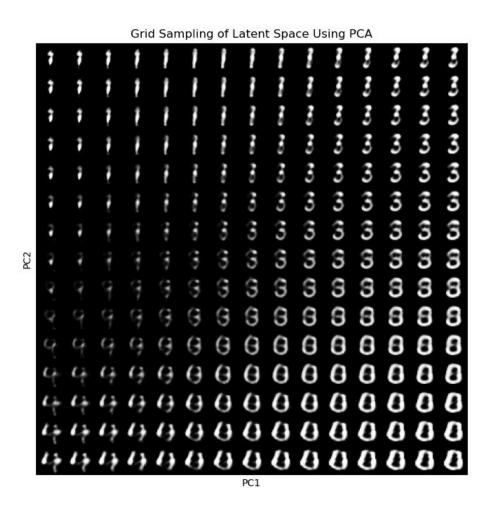
Is this a good autoencoder?

Why is loss alone (even with validation loss) not enough to tell us?





## Grid sample latent space and pass to encoder



Explained variance: PC1=0.32, PC2=0.11

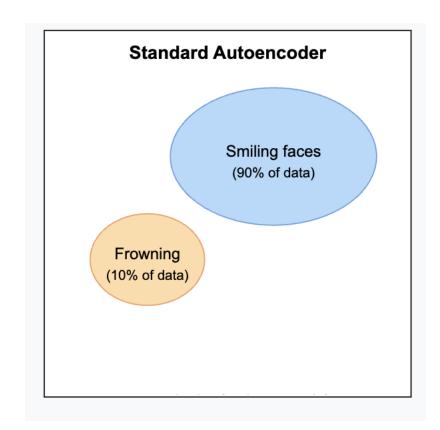
# Autoencoders can generate, but are not generative

#### Recall:

- Discriminative models learn P(y | X)
- Generative models learn P(X)

When we randomly sample, we may get some "invalid" outputs. A generative model could assign these invalid outputs a low probability P(X)

Nothing constrains the latent space of an autoencoder to represent probability distributions

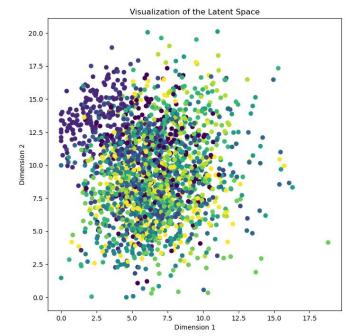


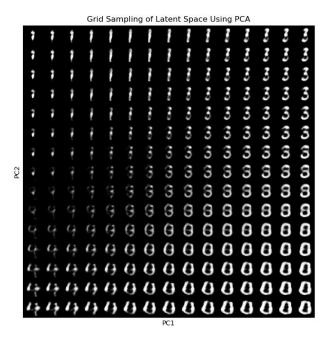
### Issues with Autoencoders

 Vectors close together in latent space may not produce similar outputs

 Tend to overfit data (struggle to produce "new" outputs)

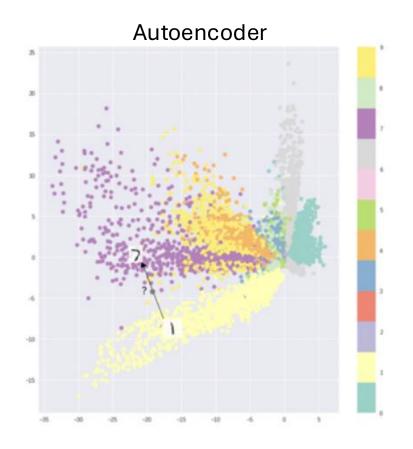
How to address issues with overfitting outputs? Try to learn more *variation* in outputs.

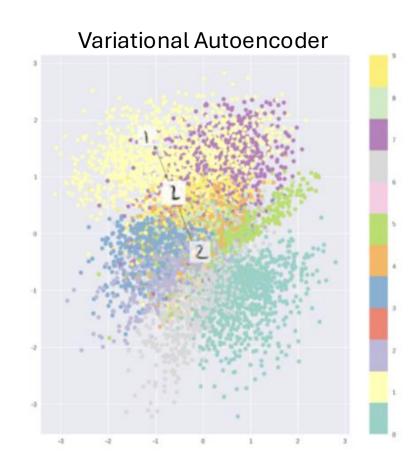




## Issues with Autoencoders

What might a better latent space look like for generation?





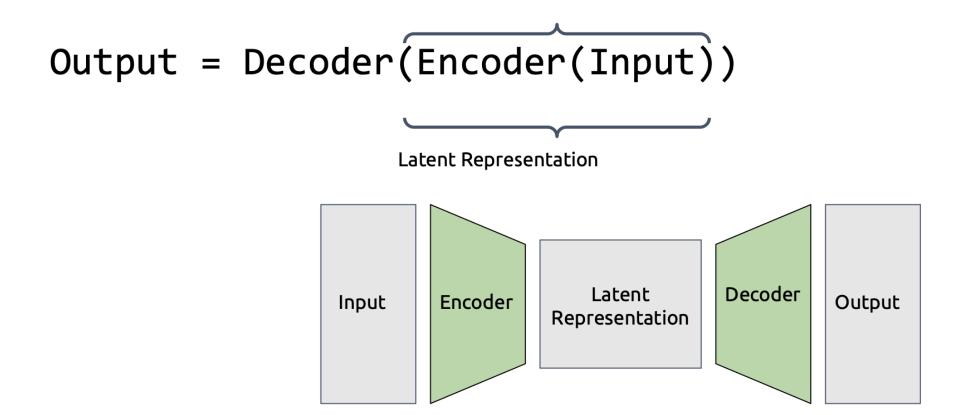
#### Variational Autoencoders

Autoencoder's goal: Reconstruct the original input

Variational Autoencoder's goal: Generate a new output that resembles the input

## Building up the VAE Architecture

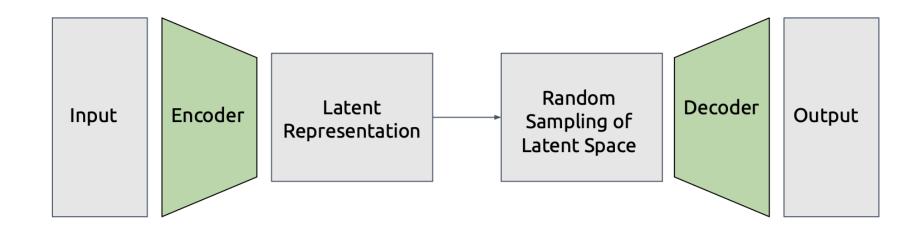
If we were to describe an autoencoder functionally:



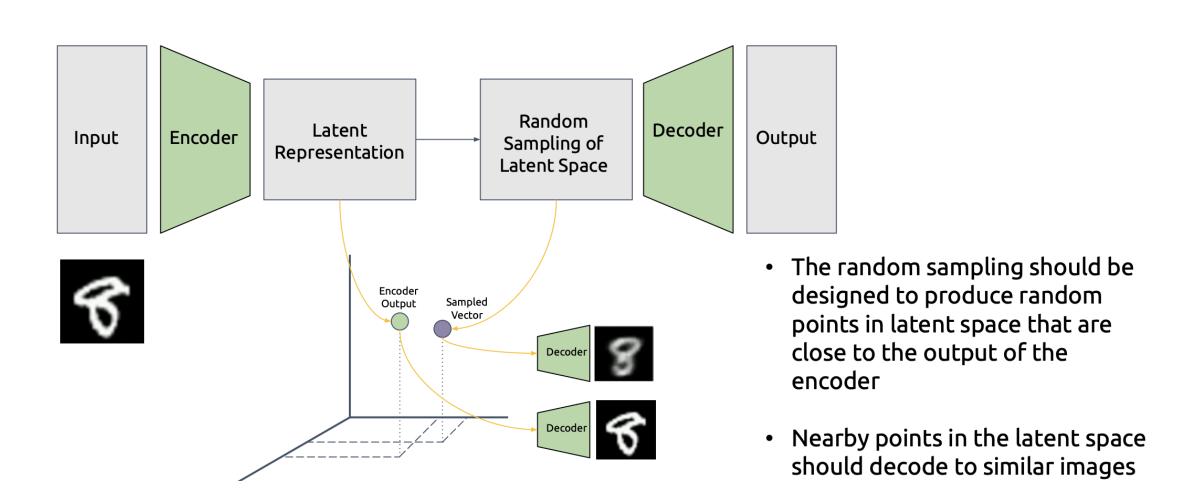
## Building up the VAE Architecture

For variational autoencoders, we also do a random sampling operation at the bottleneck

Output = Decoder(random\_sample(Encoder(Input)))



# How does random sampling in latent space lead to variation?

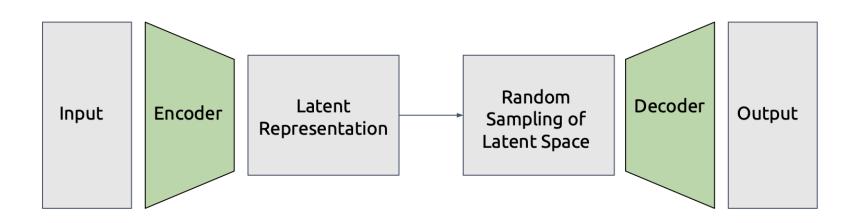


## How should random\_sample be defined?

Output = Decoder(random\_sample(Encoder(Input)))

- We want the sample to be close to the encoder output
- One option: sample from a Gaussian centered at Encoder(Input)

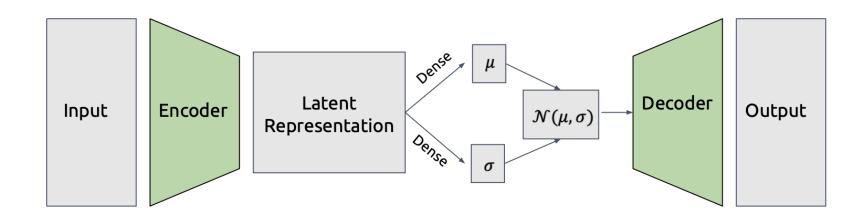
What can we modify?



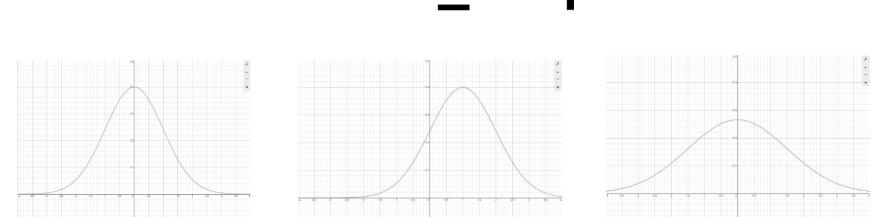
## How should random\_sample be defined?

Output = Decoder(random\_sample(Encoder(Input)))

- We want the sample to be close to the encoder output
- One option: sample from a Gaussian centered at Encoder(Input)
- Use two dense layers to convert the encoder output into the mean and standard deviation of the Gaussian



# How should random\_sample be defined?



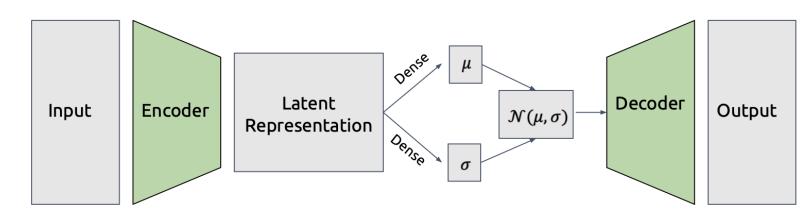




$$\mu = 1$$
 $\sigma = 1$ 

$$\mu = 0$$

$$\sigma = 1.5$$

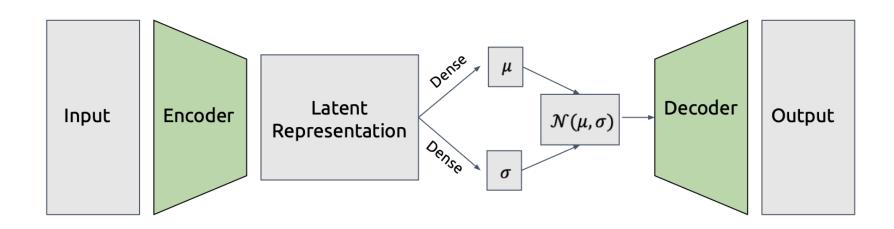


## Training a VAE

#### Two goals:

- 1. Reproduce an output similar to the input (Input ≈ Output)
- 2. Have some variation in our output (Input ≠ Output )

- Seems like two conflicting goals!
- How do we resolve these two goals?



## Weighted Combination of Losses

 $L_1$  = loss associated with producing output similar to input

 ${\cal L}_2$  = loss associated with producing output with some variation to input

Total Loss: 
$$\lambda \in [0,\infty]$$
 Input Encoder Latent Representation Representation Output

### VAE Losses

 $L_1$  is easy, we've seen this before

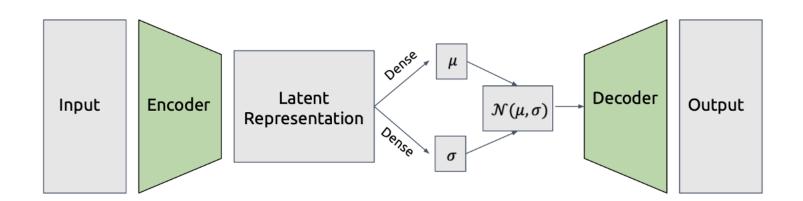
$$L_1(x, \hat{x}) = \left| |x - \hat{x}| \right|_2$$
(MSE)

But what is  $L_2$ ? How do we measure how much variation our output has?

$$L_2(??,??) = ????$$

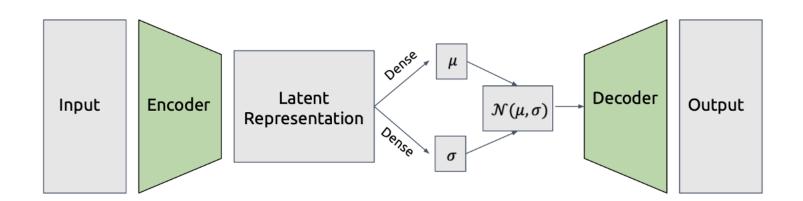
Whatever our loss function, it needs to encourage  $\sigma > 0$ , or else the model will force  $\sigma$  to 0 in an effort to create the best recreations possible.

If  $\sigma = 0$ , then the VAE will behave the same as an autoencoder!



Whatever our loss function, it needs to encourage  $\sigma > 0$ , or else the model will force  $\sigma$  to 0 in an effort to create the best recreations possible.

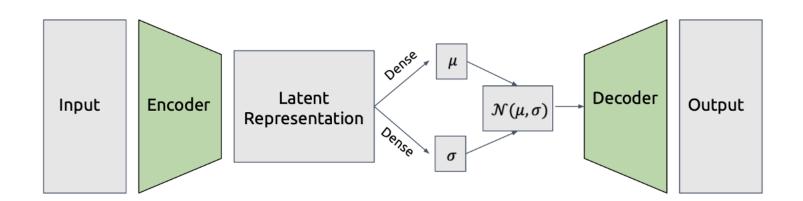
But it can't be too big... because too much variation will create poor reconstructions.



Whatever our loss function, it needs to encourage  $\sigma > 0$ , or else the model will force  $\sigma$  to 0 in an effort to create the best recreations possible.

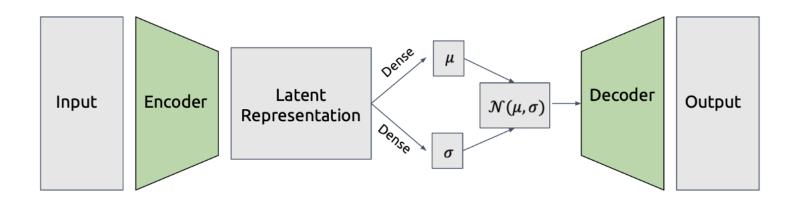
But it can't be too big... because too much variation will create poor reconstructions.

Also, what should  $\mu$  be?



### The idea:

- Introduce a prior probability function we we want our latent space to look like.
- Encourage  $N(\mu, \sigma)$  close to N(0, 1)
- (This will have beneficial properties we'll see later)



# How do we measure distance between probabilities?

Kullback-Leibler (KL) Divergence

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} \frac{p(x)}{p(x)} \log\left(\frac{p(x)}{q(x)}\right) dx$$

### What this says:

- "Everywhere that p has probability density..."
- "...the difference between p and q should be small"
  - Difference in log probabilities (remember that  $\log\left(\frac{a}{b}\right) = \log(a) \log(b)$ )

### KL Divergence

- Expensive to compute, in general (no closed form, have to numerically approximate the integral)
- But! There is a closed form for Gaussians:

$$D_{KL}(\mathcal{N}(\mu, \sigma^2) || \mathcal{N}(0, 1)) = \frac{1}{2} \sum_{i=1}^{k} (\mu_i^2 + \sigma_i^2 - \ln \sigma_i^2 - 1)$$

K is the dimensionality of  $\vec{\mu}$ ,  $\vec{\sigma}$  (i.e., the size of the encoding)

### The Final VAE Loss Function

We now have all the tools necessary to construct our loss function.

$$L = L_1 + \lambda L_2 \qquad \qquad \lambda \in [0, \infty]$$

Which turns into this:

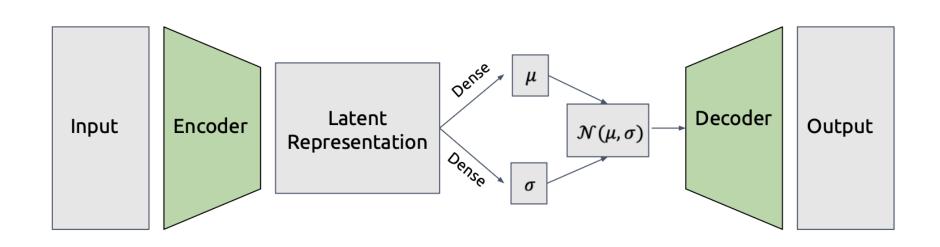
$$L = ||x - \hat{x}||_{2}^{2} + \lambda D_{KL}(\mathcal{N}(\mu, \sigma), \mathcal{N}(0, 1))$$

### Any questions?



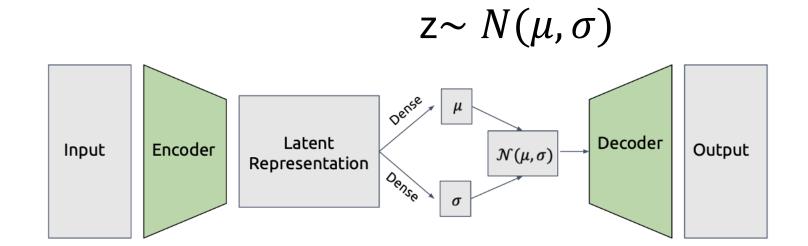
# Putting it all together

$$L = ||x - \hat{x}||_2^2 + \lambda D_{KL}(\mathcal{N}(\mu, \sigma), \mathcal{N}(0, 1))$$



### There's just one issue

How do we take the gradient of a sampling operation?



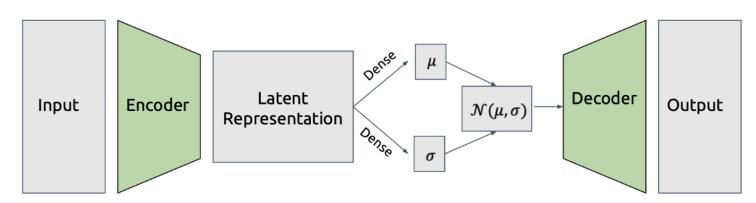
### Reparametization Trick

$$z \sim N(\mu, \sigma)$$

Can be rewritten as:

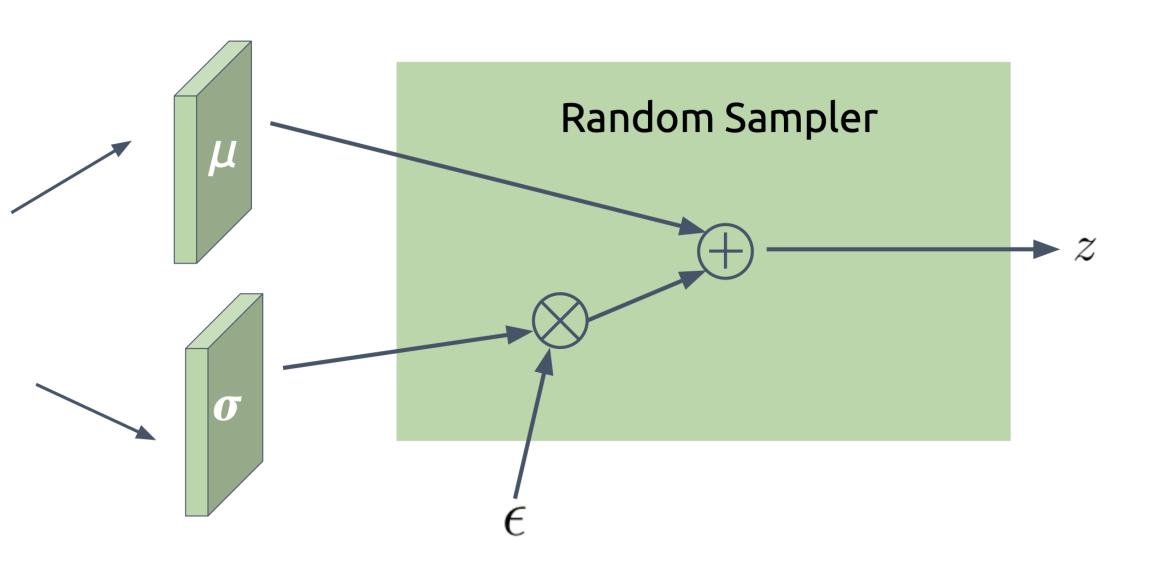
$$z = \mu + \epsilon \sigma$$
, where  $\epsilon \sim N(0, 1)$ 

Random sampling operation  $(\epsilon)$  no longer depends on learnable parameters

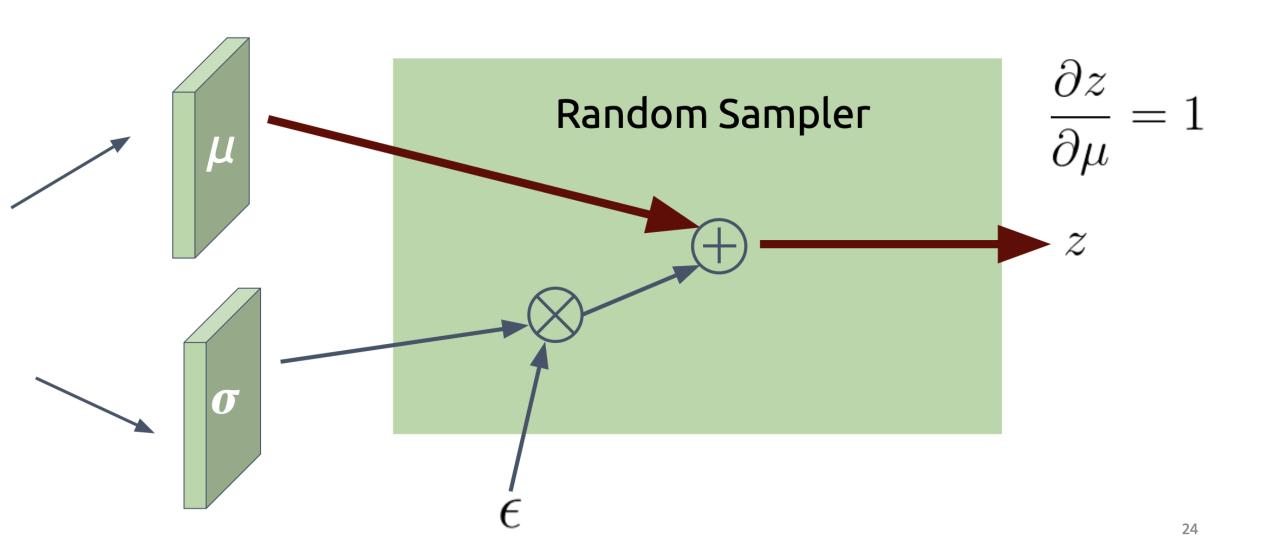


Another explanation of why this is needed: <a href="https://gregorygundersen.com/blog/2018/04/29/reparameterization/">https://gregorygundersen.com/blog/2018/04/29/reparameterization/</a>

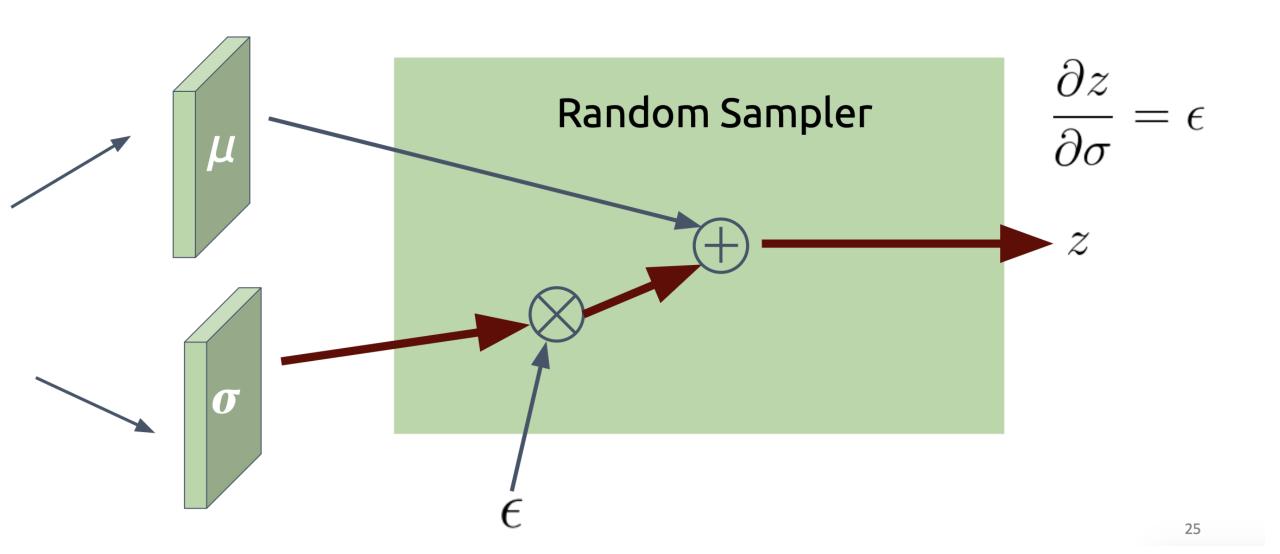
### Random Sampler with Reparameterization Trick



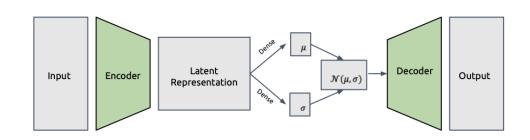
# Random Sampler with Reparameterization Trick



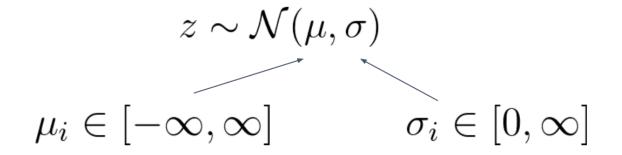
## Random Sampler with Reparameterization Trick



### One more practical detail

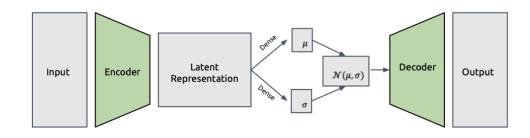


Let's again consider our sampling operation

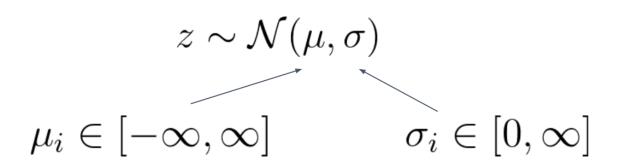


- Nothing prevents the neural network from outputting negative values for the standard deviation.
- Instead of predicting  $\sigma$ , we will instead predict  $\log(\sigma^2)$  . This ensures that every  $\sigma_i \in [0,\infty]$ 
  - ullet i.e. just treat the output of the Dense layer as if it is  $\log(\sigma^2)$

### One more practical detail



Let's again consider our sampling operation



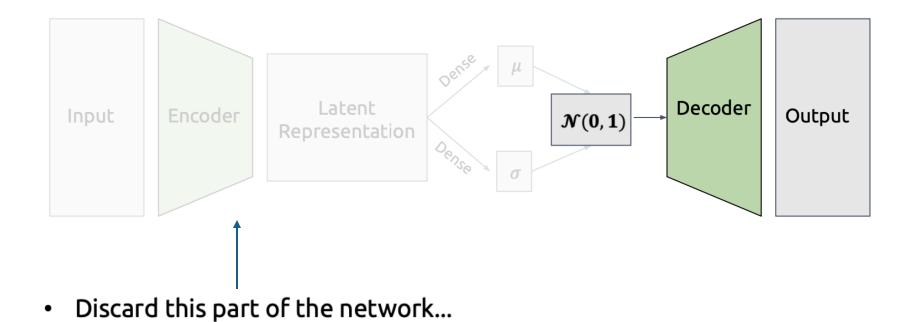


- Instead of predicting  $\sigma$ , we will instead predict  $\log(\sigma^2)$  . This ensures that every  $\sigma_i \in [0,\infty]$ 
  - i.e. just treat the output of the Dense layer as if it is  $\log(\sigma^2)$

$$D_{KL}(\mathcal{N}(\mu, \sigma^2) || \mathcal{N}(0, 1)) = \frac{1}{2} \sum_{i=1}^{k} (\mu_i^2 + \sigma_i^2 - \ln \sigma_i^2 - 1)$$

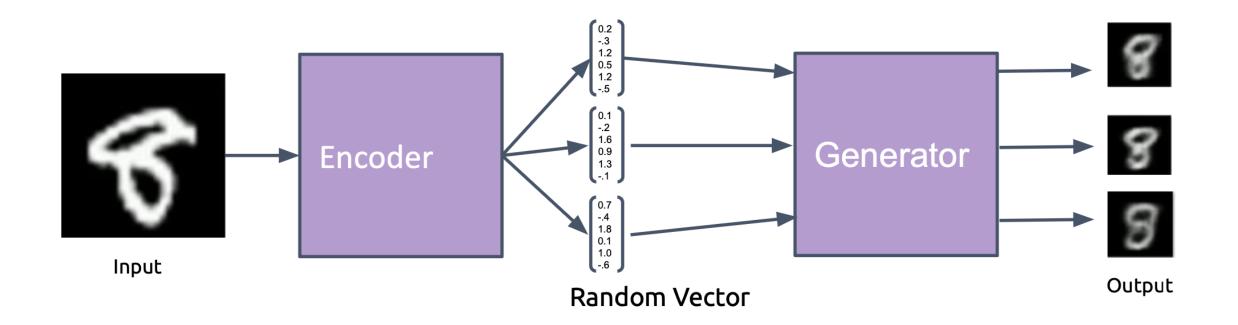
# Sampling from a VAE

...and set  $(\mu, \sigma) = (0, 1)$ 



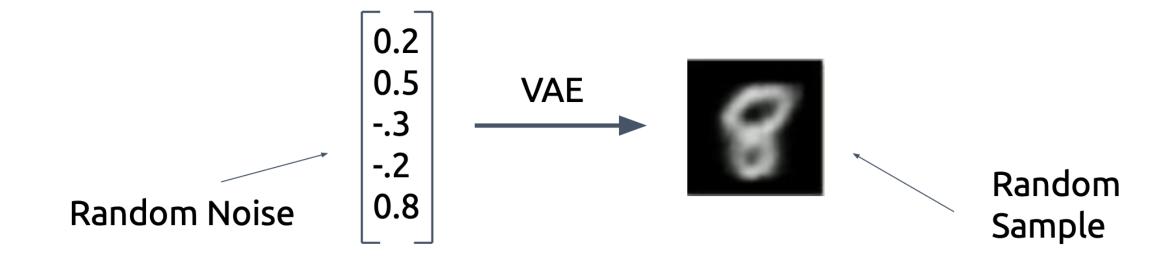
# Sampling from a VAE

 We can use a trained VAE to generate random variants of an input data point...



# Sampling from a VAE

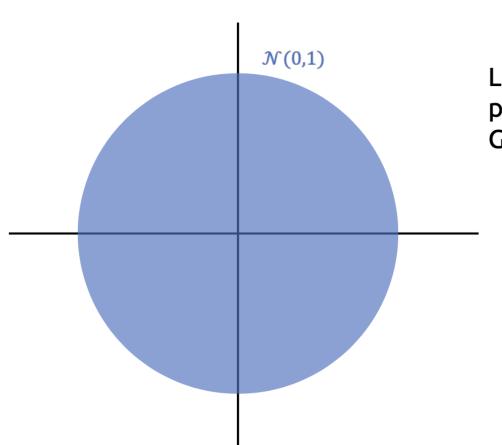
... But ultimately, we want to draw random samples from a VAE



### How can we do this?

This is where our particular choice of training loss will pay off

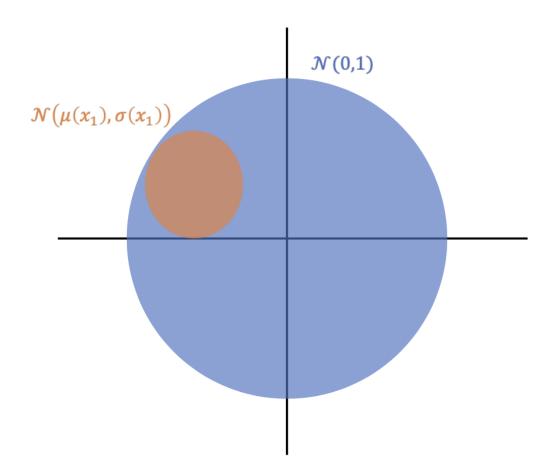
# Encoding different points into latent space



Let this circle represent the probability density of a unit Gaussian in latent space

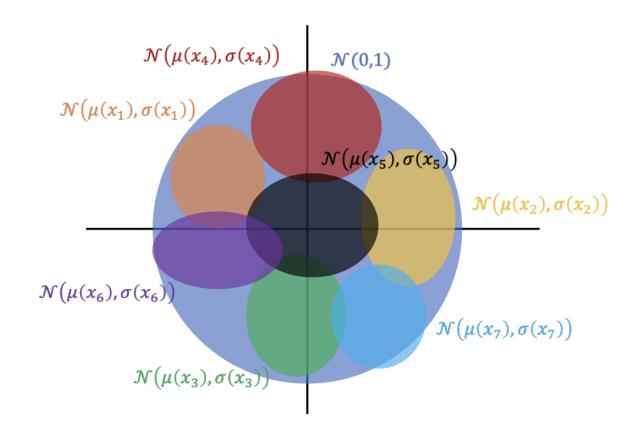
# Encoding different points into latent space

Let this circle represent the probability density of the  $\mathcal{N}(\mu, \sigma)$  distribution that the encoder predicts given an input data point  $x_1$ 



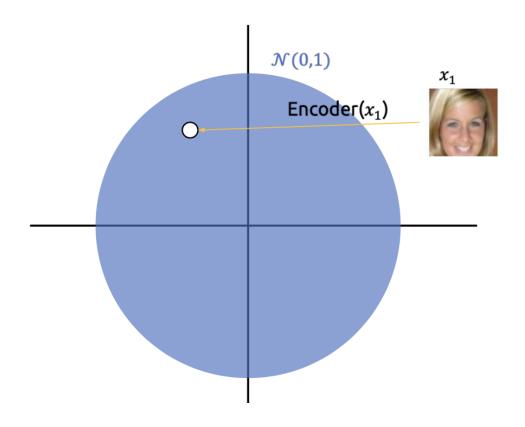
### Encoding different points into latent space

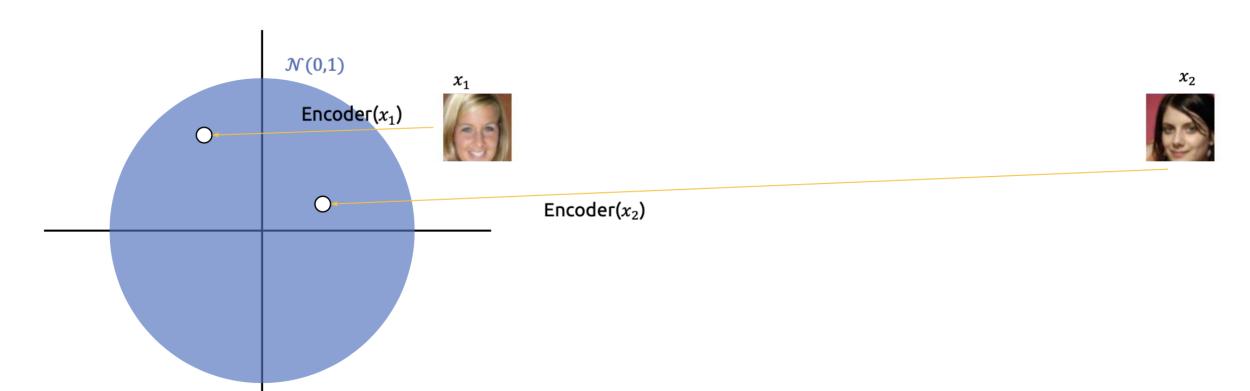
$$L = ||x - \hat{x}||_{2}^{2} + \lambda D_{KL}(\mathcal{N}(\mu, \sigma), \mathcal{N}(0, 1))|$$

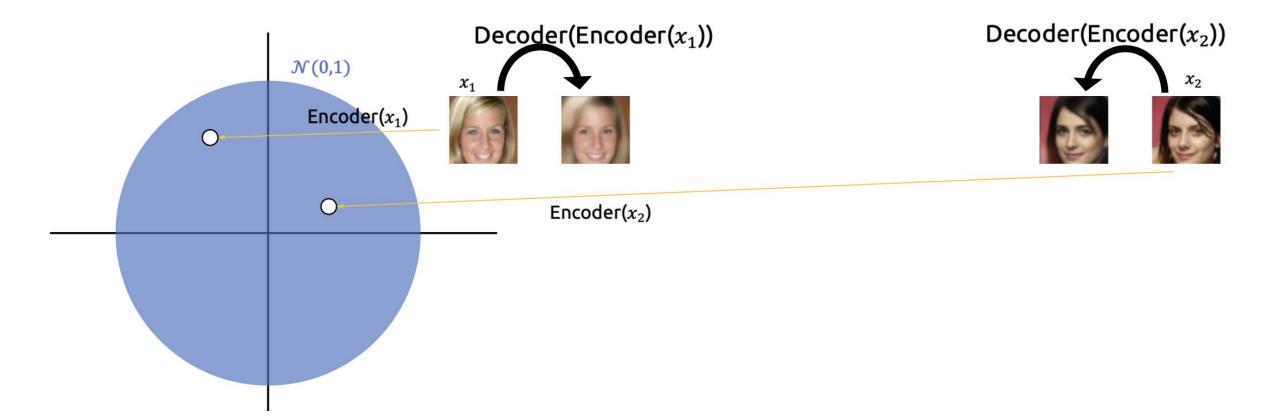


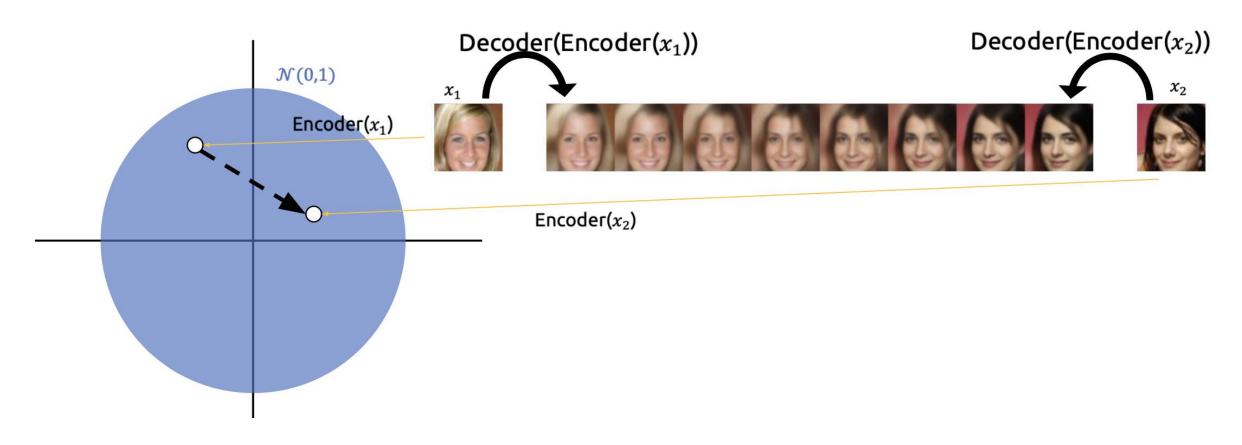
Because of our KL divergence loss, the  $\mathcal{N}(\mu, \sigma)$  for any input data point has to be somewhat similar to  $\mathcal{N}(0,1)$ 

So, if we sample a point from  $\mathcal{N}(0,1)$ , it is very likely to fall within one of these encoded







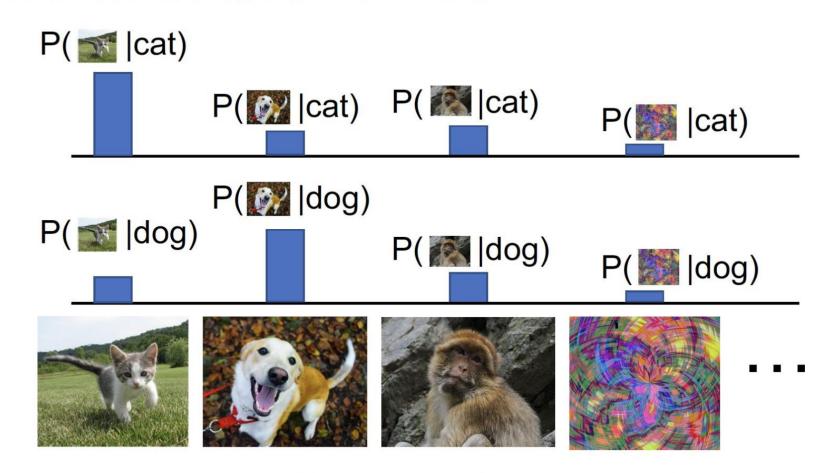


### Discriminative vs Generative Models

Discriminative Model: Learn a probability distribution p(y|x)

Generative Model: Learn a probability distribution p(x)

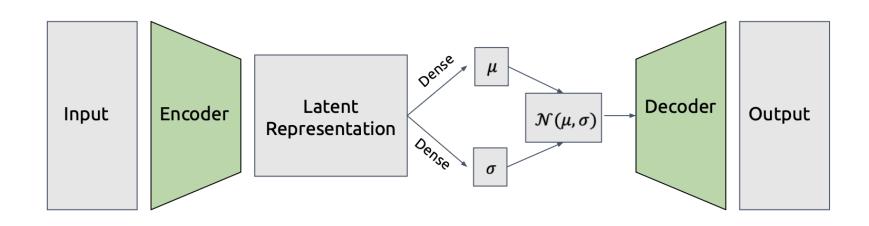
**Conditional Generative Model:** Learn p(x|y)



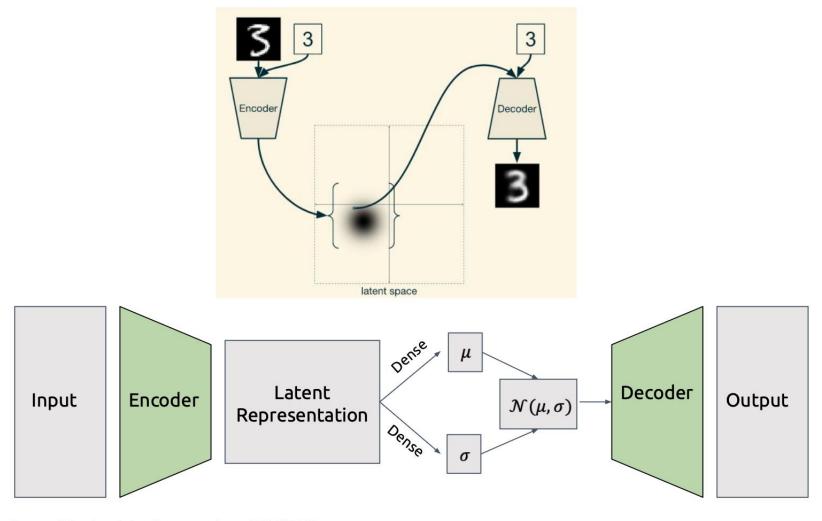
Conditional Generative Model: Each possible label induces a competition among all images

Credit: UMich EECS498

### **Conditional VAE**



### **Conditional VAE**



# VAE output

### Input



VAE reconstruction



https://towardsdatascience.com/what-the-heck-are-vae-gans-17b86023588a

What's the issue here?

Why?

### Why are VAE samples blurry?

- Our reconstruction loss is the culprit
- Mean Square Error (MSE) loss looks at each pixel in isolation
- If no pixel is too far from its target value, the loss won't be too bad
- Individual pixels look OK, but larger-scale features in the image aren't recognizable

#### Solutions?

Let's choose a different reconstruction loss!

#### Input



VAE reconstruction



https://towardsdatascience.com/what-the-heck-are-vae-gans-17b86023588a

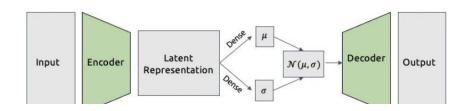
# Recap

Variational Autoencoders (VAEs)

### **Loss Function**

Reparameterization Trick

#### **Conditional VAEs**



Input



VAE reconstruction



https://towardsdatascience.com/wha t-the-heck-are-vae-gans-17b860235 88a