

# Variational Autoencoders: Neural Network Perspective

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in VAE

1. Generative Model
2. Black-box variational inference
3. The reparameterization trick
4. Choosing  $q$  and  $p$
5. Autoencoding Variational Bayes (AEVB)
6. Variational autoencoder (VAE)
7. VAE: The neural network perspective
8. Examples

# Topics in VAE as Neural Nets

1. What are VAEs useful for?
2. Neural Network Perspective
3. Loss Function
4. Stochastic Gradient Descent
5. Illustration of latent space with MNIST images

# What are VAEs useful for?

- To design complex generative models, e.g.,
  1. Generate fictional celebrity faces



Interpolate between points in latent space

## 2. Produce synthetic text

**“ i want to talk to you . ”**  
*“i want to be with you . ”*  
*“i do n’t want to be with you . ”*  
*i do n’t want to be with you .*  
**she did n’t want to be with him .**

## 3. Interpolate between MIDI samples



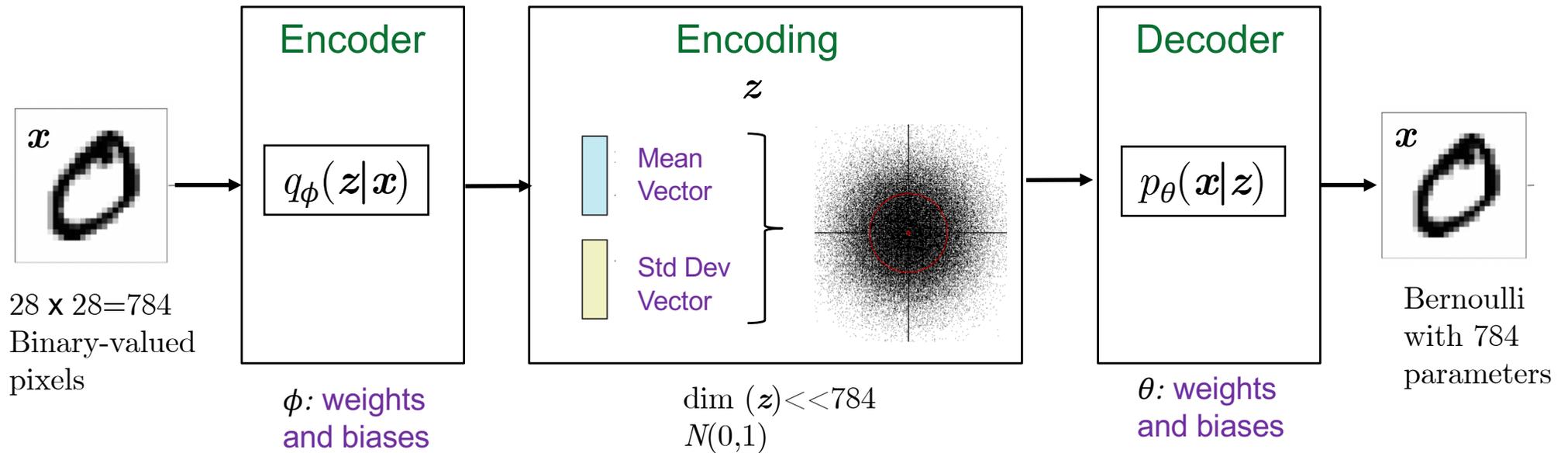
[https://www.youtube.com/watch?time\\_continue=36&v=Ir\\_AFDKOc-I](https://www.youtube.com/watch?time_continue=36&v=Ir_AFDKOc-I)

# Understanding VAEs

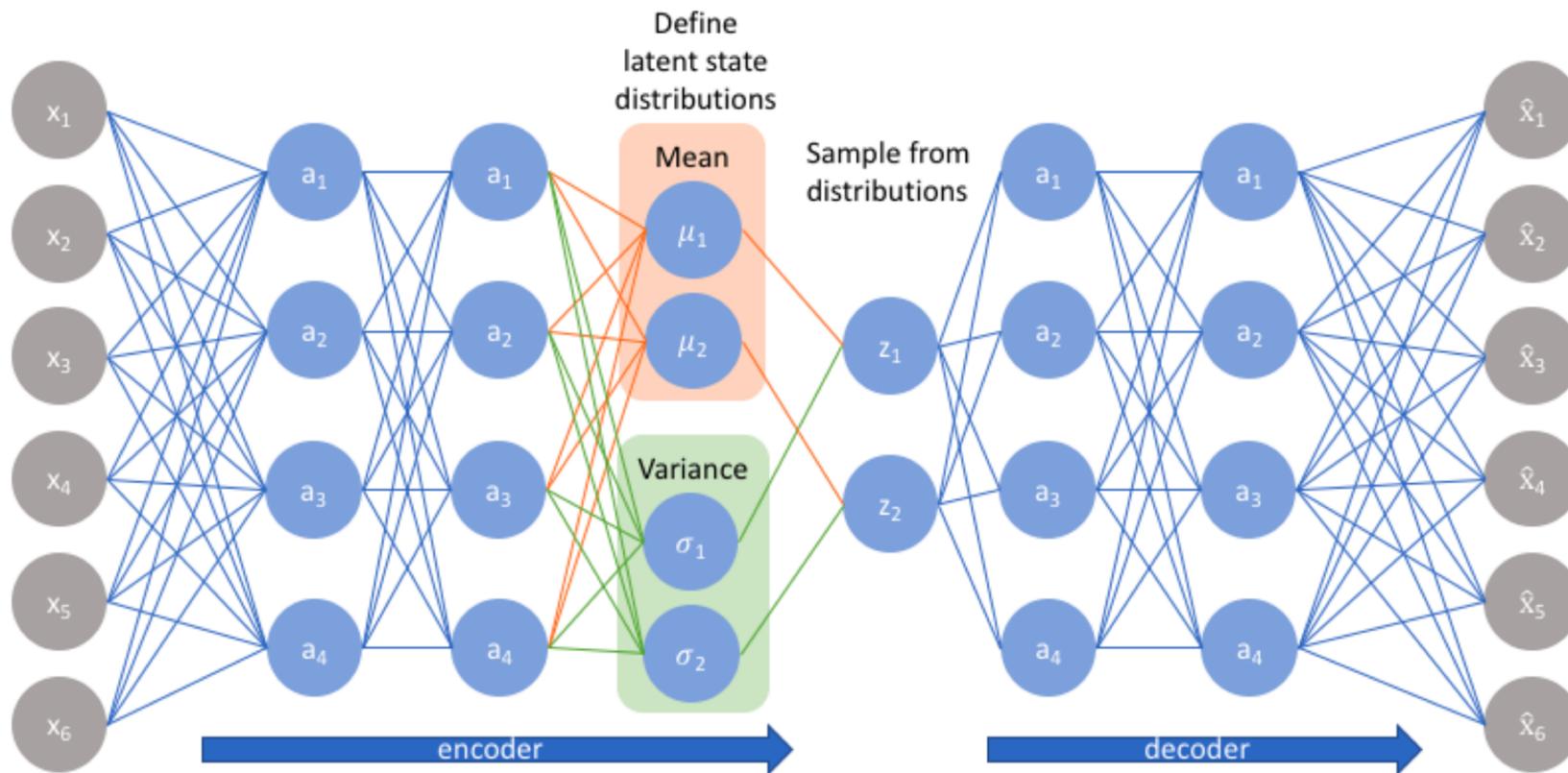
- Need two perspectives:
  1. Neural networks
  2. Probabilistic graphical models
- They do not have a shared language
- Here we discuss the neural network approach
- For the PGM approach see <https://cedar.buffalo.edu/~srihari/CSE676/21.1-VAE.pdf>

# Neural Network Perspective

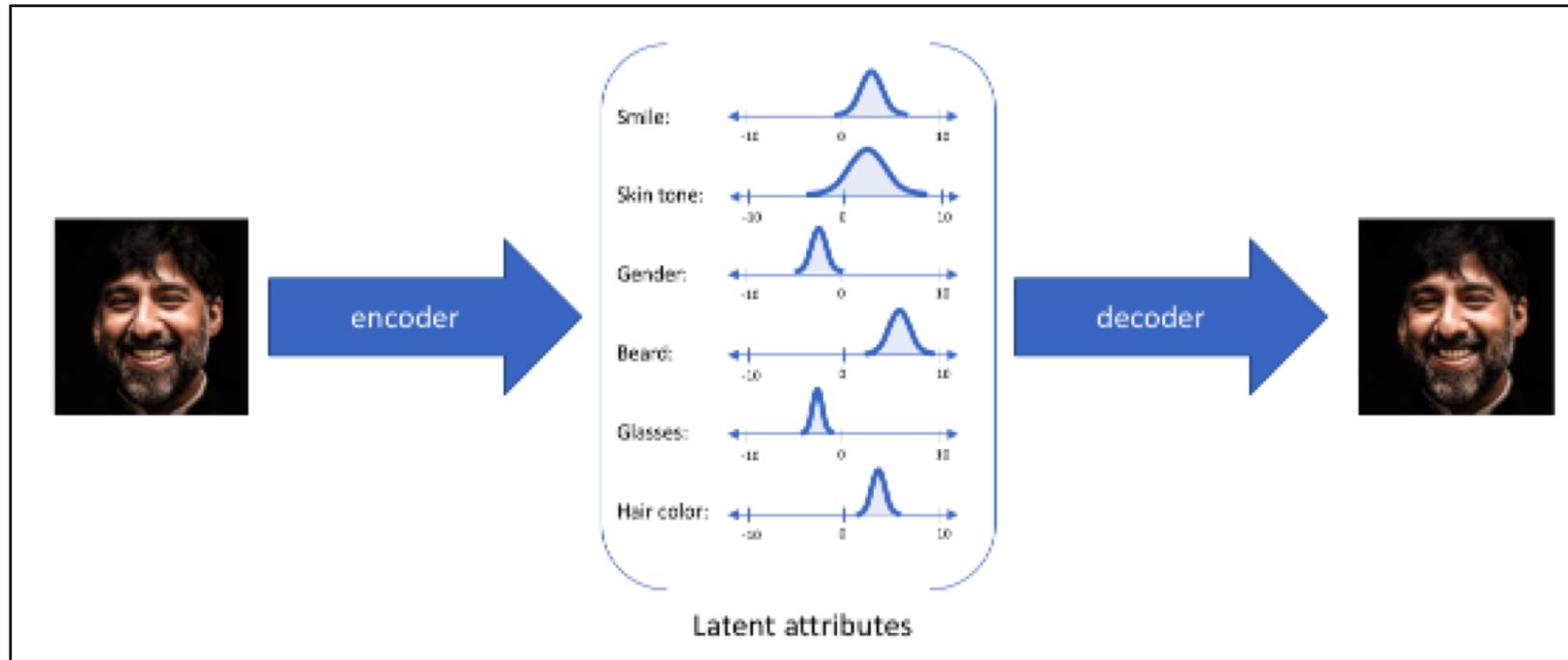
A VAE consists of an encoder, a decoder, and a loss function



# VAE using neural networks



# Ex: Mean and Std Dev Vectors



Source: *Auto-Encoding Variational Bayes*, Diederik P Kingma, Max Welling

# The encoder neural network

- Input:  $\mathbf{x}$ , output: latent  $z$ , & weights/biases:  $\theta$ 
  - E.g.,  $\mathbf{x}$  is a 28 x 28 handwritten digit image
  - Maps 784 space into much smaller latent space  $z$ 
    - A ‘bottleneck’ because encoder must learn efficient compression into lower-dimensional space.
  - Denote the encoder  $q_{\phi}(z|\mathbf{x})$
- The lower-dimensional space is stochastic
  - Encoder outputs *parameters* to  $q_{\phi}(z|\mathbf{x})$ , a Gaussian
- We can sample from this distribution to get noisy values of the representations  $z$

# The decoder neural network

- The decoder gets as input the latent representation  $z$
- It outputs parameters of distribution of the data
- Decoder is denoted by  $p_{\theta}(x|z)$ 
  - Weights and biases:  $\theta$
- With handwritten digits with binary pixels: 0 or 1
  - The distribution of a single pixel is a Bernoulli
  - Output  $x$  consists of 784 Bernoulli parameters, one per pixel
    - real-valued numbers between 0 and 1

# Information loss

- Information is lost because it goes from a smaller to a larger dimensionality.
- How much information is lost?
  - We measure this using the reconstruction log-likelihood  $\log p_{\theta}(\mathbf{x}|\mathbf{z})$  whose units are nats
- This measure tells us how effectively the decoder has learned to reconstruct an input image  $\mathbf{x}$  given its latent representation  $\mathbf{z}$

# Loss Function

- The *loss function* of the VAE is the negative log-likelihood with a regularizer
- Since no global representations are shared by datapoints, the total loss is  $\sum_i l_i$  for  $i=1,..N$

– where loss  $l_i$  for datapoint  $\mathbf{x}_i$  is:

$$l_i(\theta, \phi) = -E_{z \sim q_\phi(z|\mathbf{x}_i)} \left[ \log p_\theta(\mathbf{x}_i | z) \right] + KL(q_\phi(z | \mathbf{x}_i) || p(z))$$

– The first term is the reconstruction loss

- Encourages decoder to learn to reconstruct the data

– The second term is the regularizer term

- It is K-L divergence between  $q_\phi(z|\mathbf{x})$  and  $p(z)$
- Measures information lost when using  $q$  to represent  $p$

# Specifying $p$

- In VAE,  $p$  is specified as a standard Normal distribution or  $p(\mathbf{z}) \sim N(0,1)$
- If the encoder outputs representations  $\mathbf{z}$  that are different than those from a standard normal distribution, it will receive a penalty in the loss
- This regularizer term means 'keep the representations  $\mathbf{z}$  of each digit sufficiently diverse'

# Stochastic Gradient Descent

- We train the variational autoencoder using gradient descent to optimize the loss with respect to the parameters of the encoder and decoder  $\phi$  and  $\theta$
- For SGD with step size  $\rho$ , the encoder parameters are updated using
$$\phi \leftarrow \phi - \rho \partial l / \partial \phi$$
- and the decoder is updated similarly

# Taking gradients

- Finally, how to take derivatives w.r.t. the parameters of a stochastic variable
- If we are given  $z$  that is drawn from a distribution  $q_\phi(z|\mathbf{x})$ , and we want to take derivatives of a function of  $z$  with respect to  $\phi$ , how do we do that?
- The  $z$  sample is fixed, but intuitively its derivative should be nonzero

# The reparameterization trick

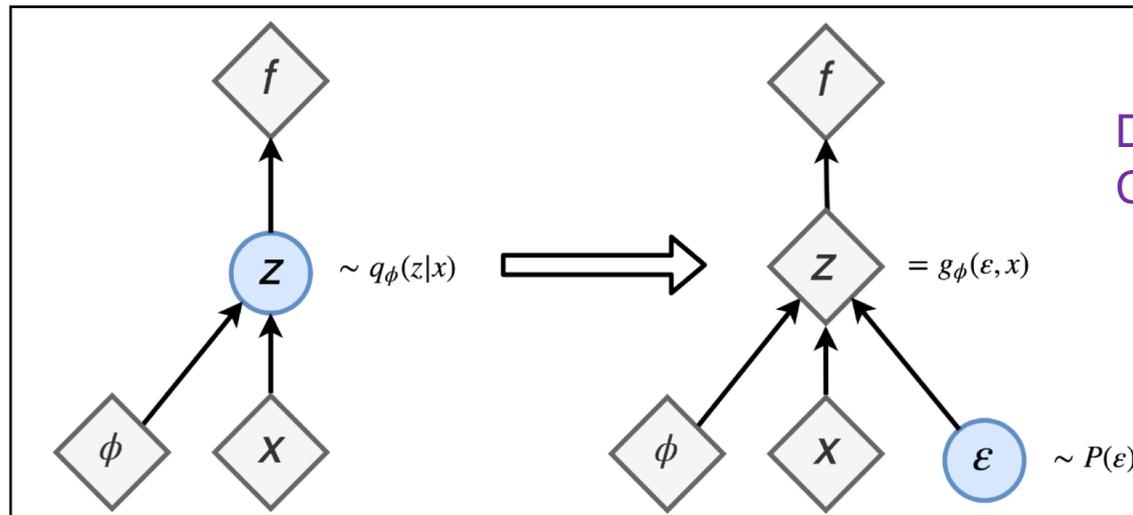
- For some distributions, it is possible to reparametrize samples such that the stochasticity is independent of the parameters
- We want samples to depend deterministically on the parameters of the distribution
- E.g., in a normally-distributed variable with mean  $\mu$  and standard deviation  $\sigma$  we can sample from it like this:

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

$\odot$  is element-wise (Hadamard) multiplication

# Reparameterization

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim N(0,1)$$



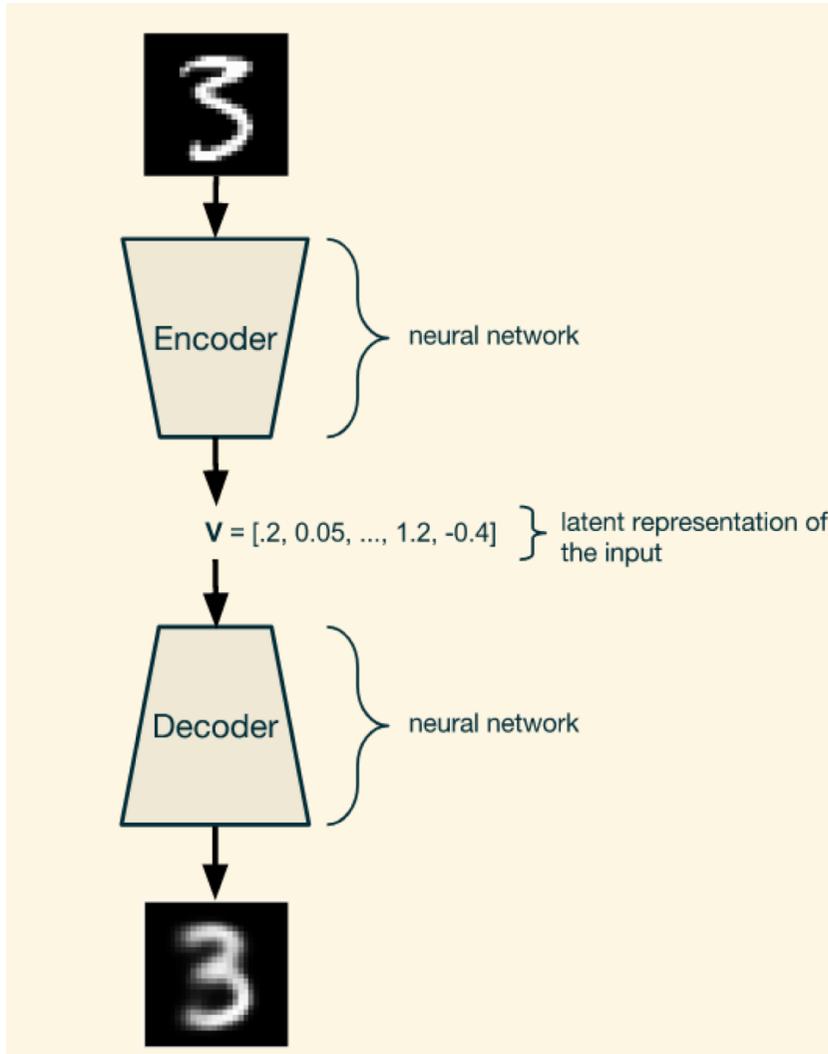
Diamonds: Deterministic dependencies  
Circles: Random variables

- Going from  $\sim$  denoting a draw to the equals sign  $=$  is the crucial step
  - We now have a function that depends on parameters deterministically
- Can take derivatives of functions involving  $\mathbf{z}$ ,  $f(\mathbf{z})$  w.r.t. the parameters of its distribution  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$

# Role of reparameterization

- In the VAE, the mean and variance are output by an inference network with parameters  $\theta$  that we optimize
- The reparameterization trick lets us backpropagate (take derivatives using the chain rule) with respect to  $\theta$  through the objective (the ELBO) which is a function of samples of the latent variables  $z$

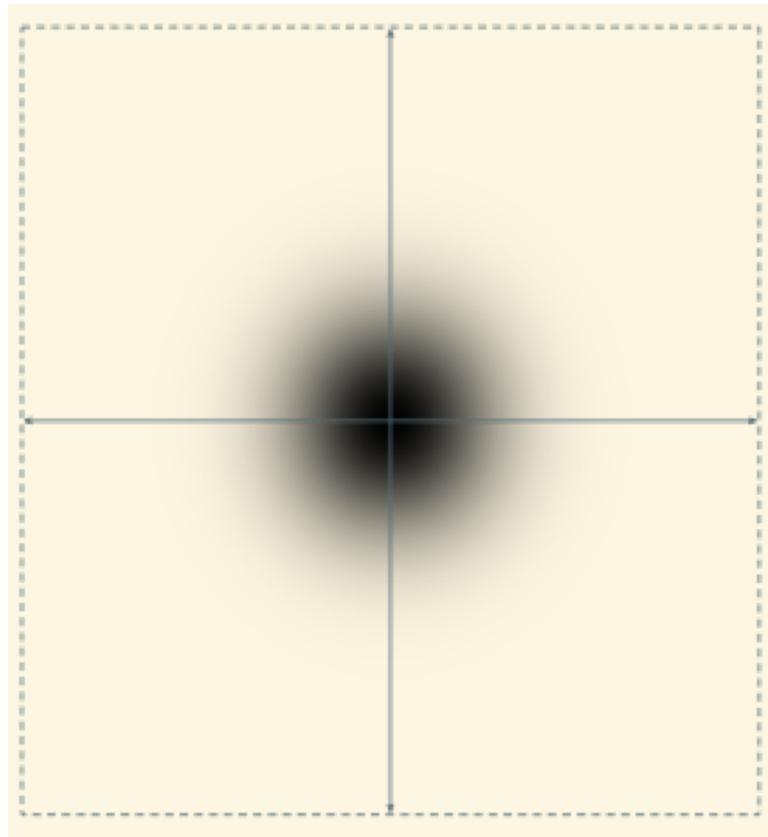
# Standard Autoencoder



A standard autoencoder trained on MNIST digits may not provide a reasonable output when a  $V$  image is input

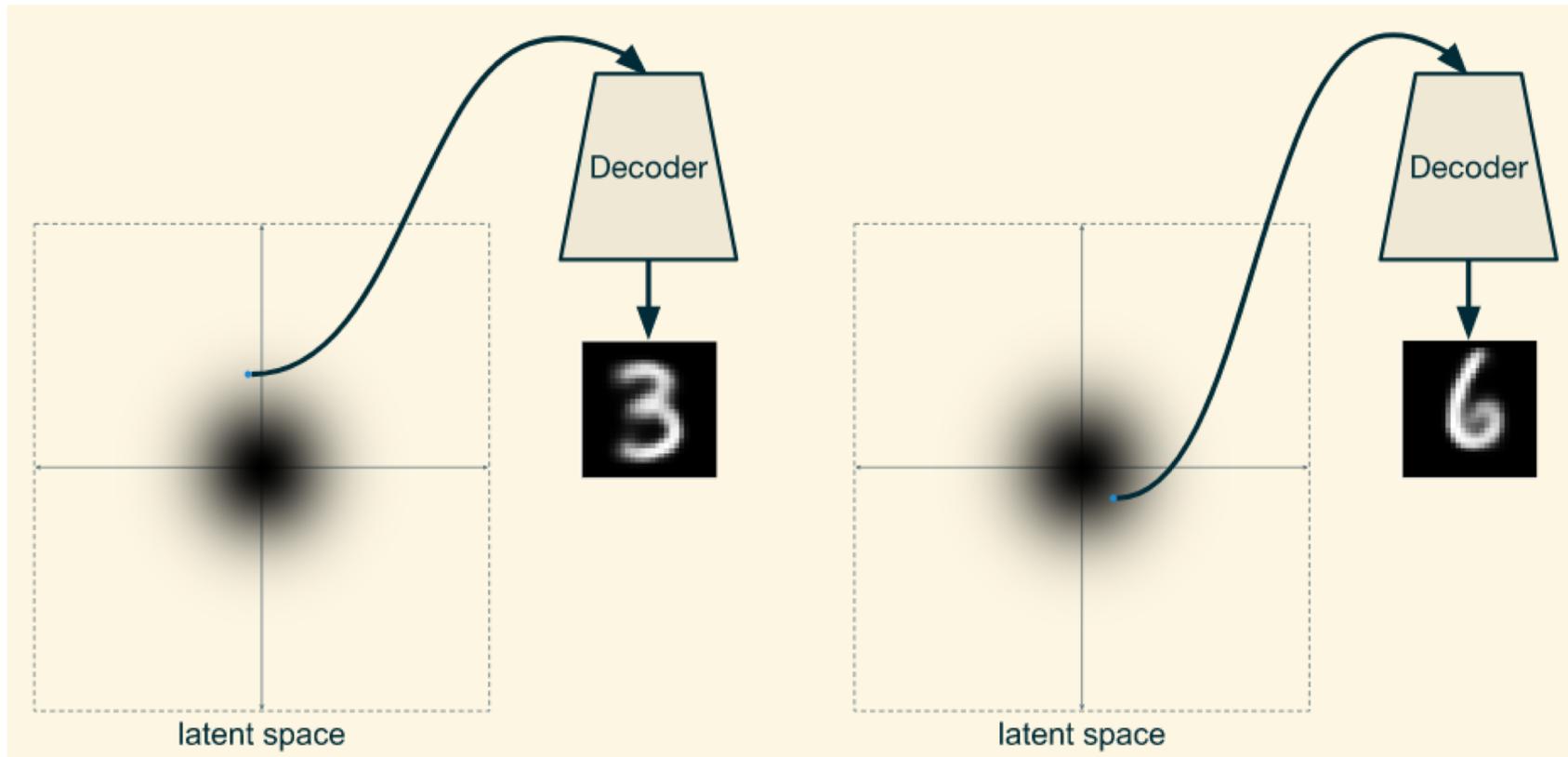
# Normal Distribution of MNIST

A standard normal distribution



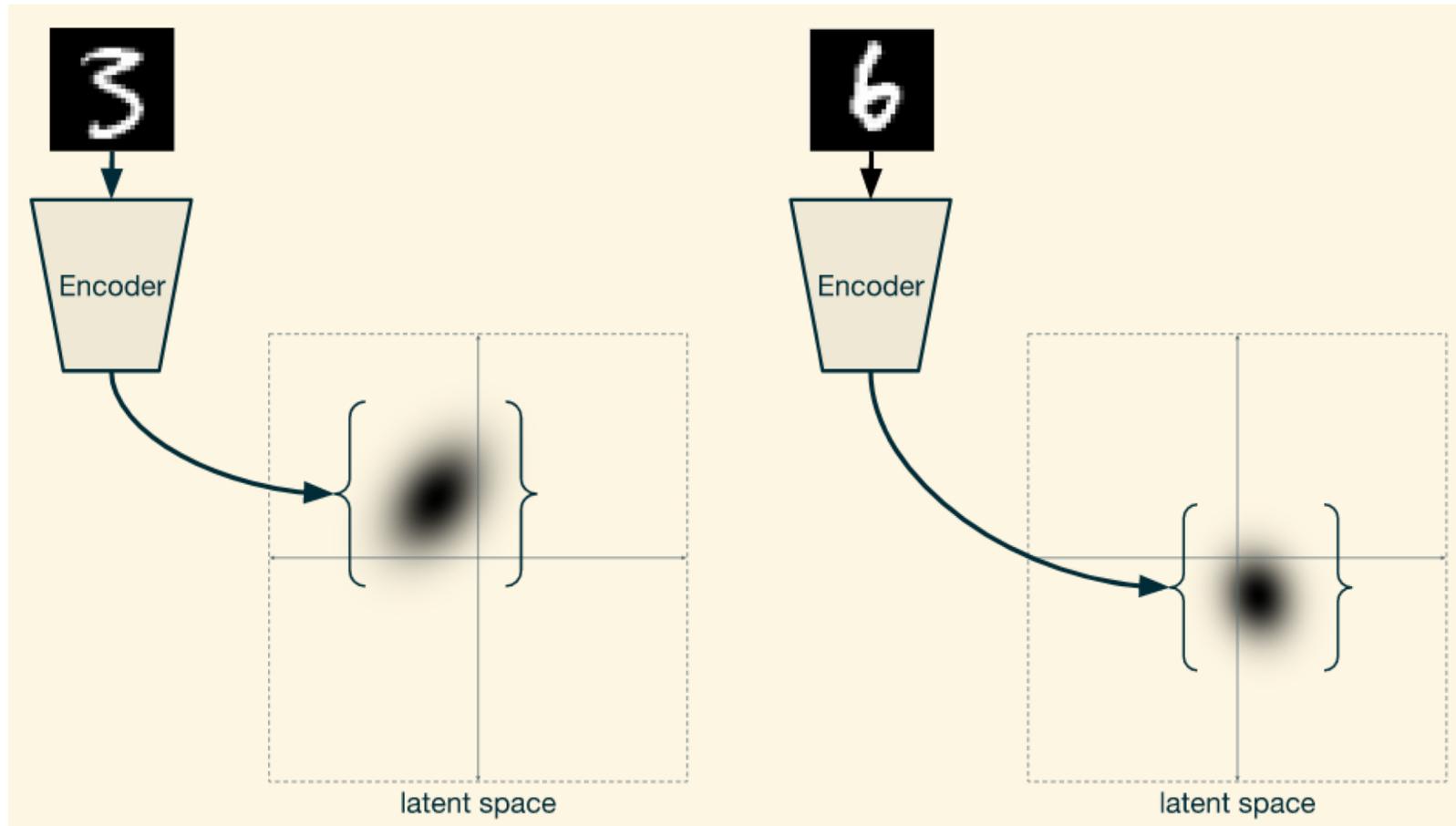
This is how we would like points corresponding to MNIST digit images to be distributed in the latent space

# Decoder of a VAE

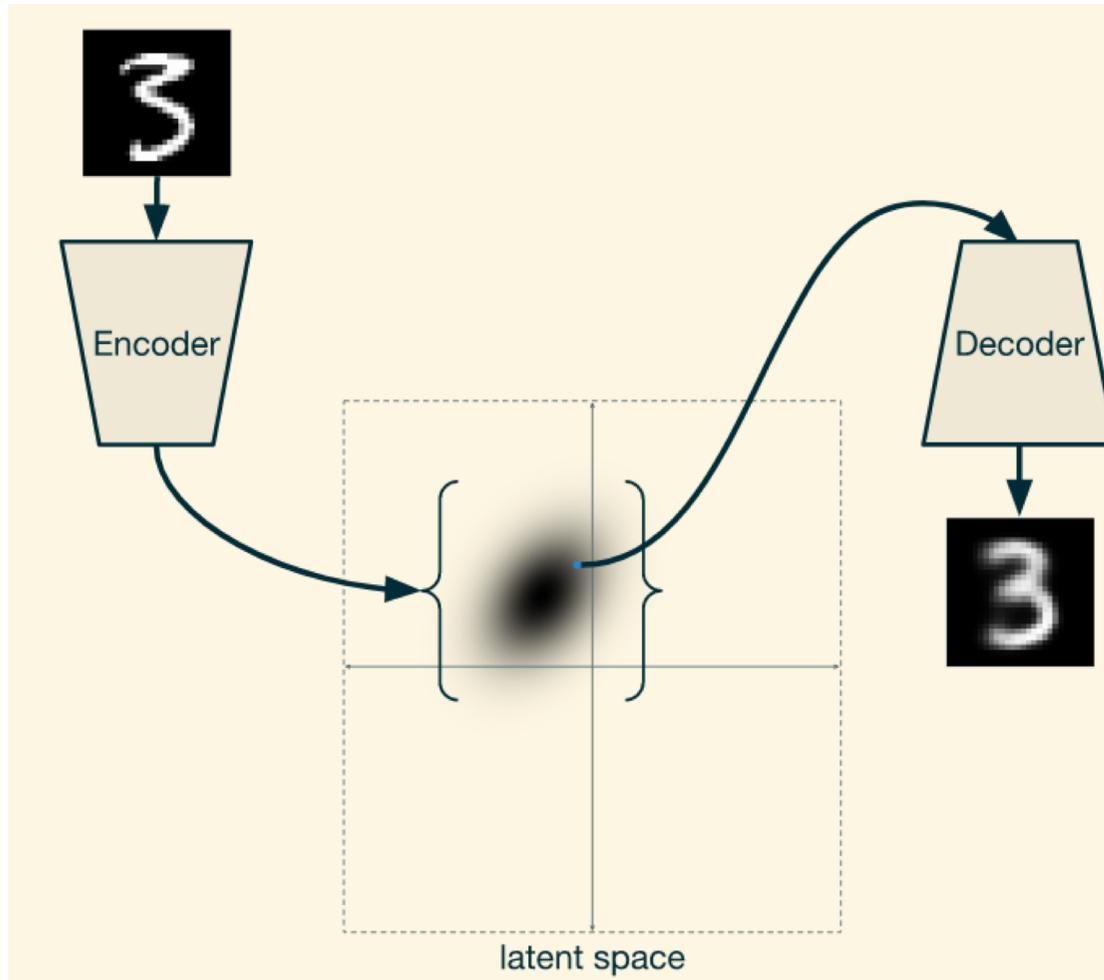


3 s are in first quadrant, 6 s are in third quadrant

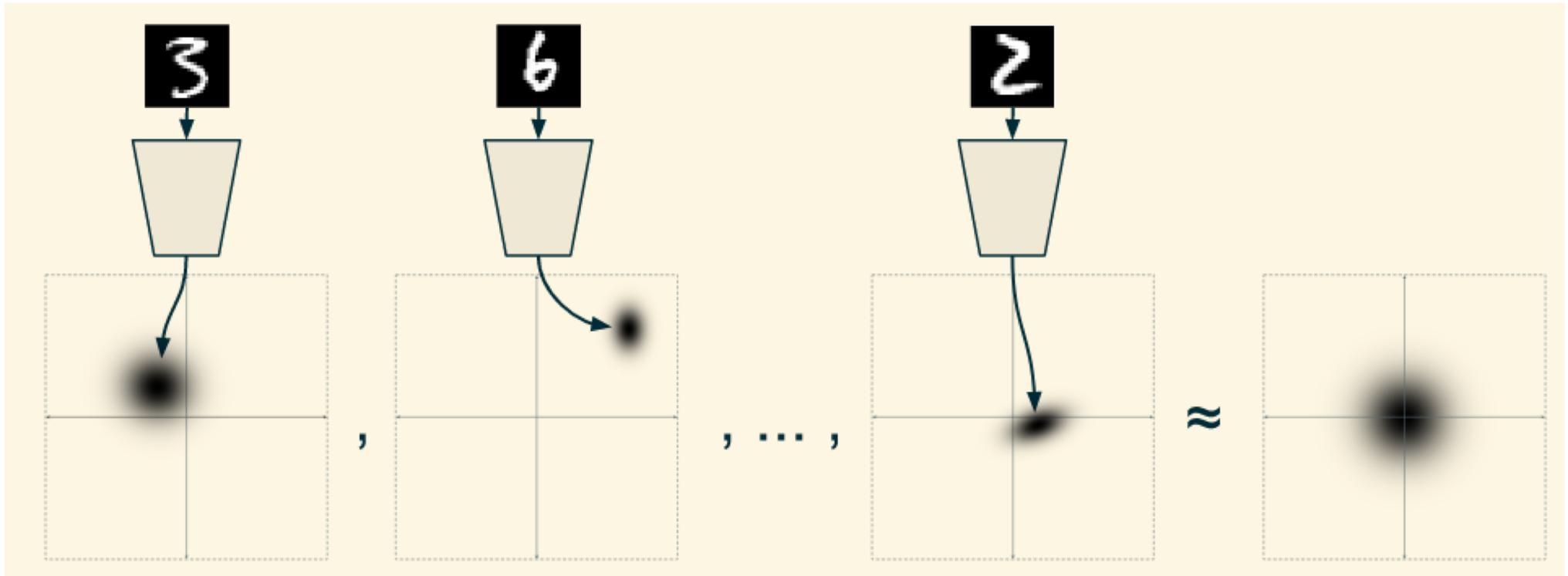
# Encoder of a VAE



# MNIST Variational Autoencoder

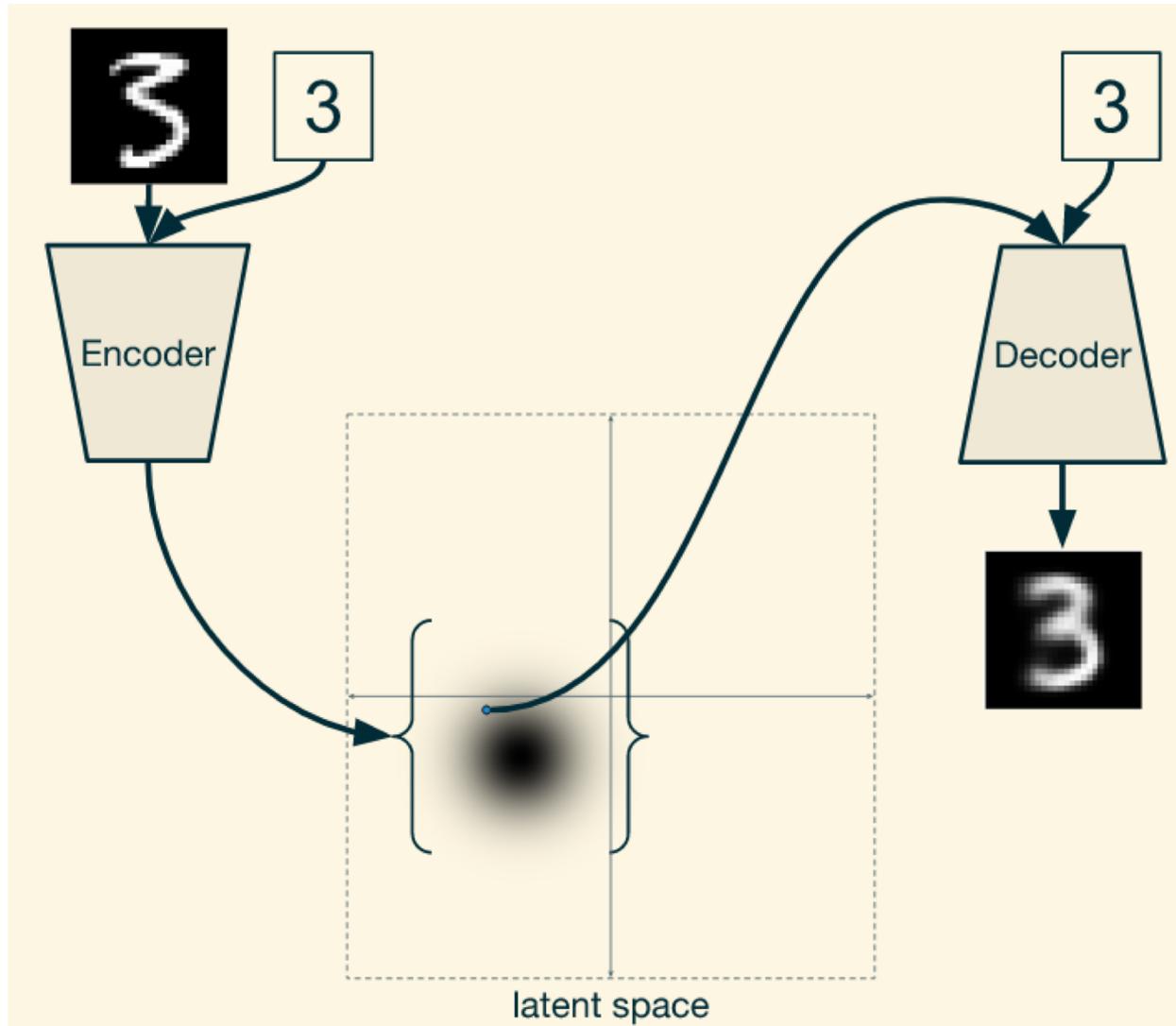


# Structure of Latent Space



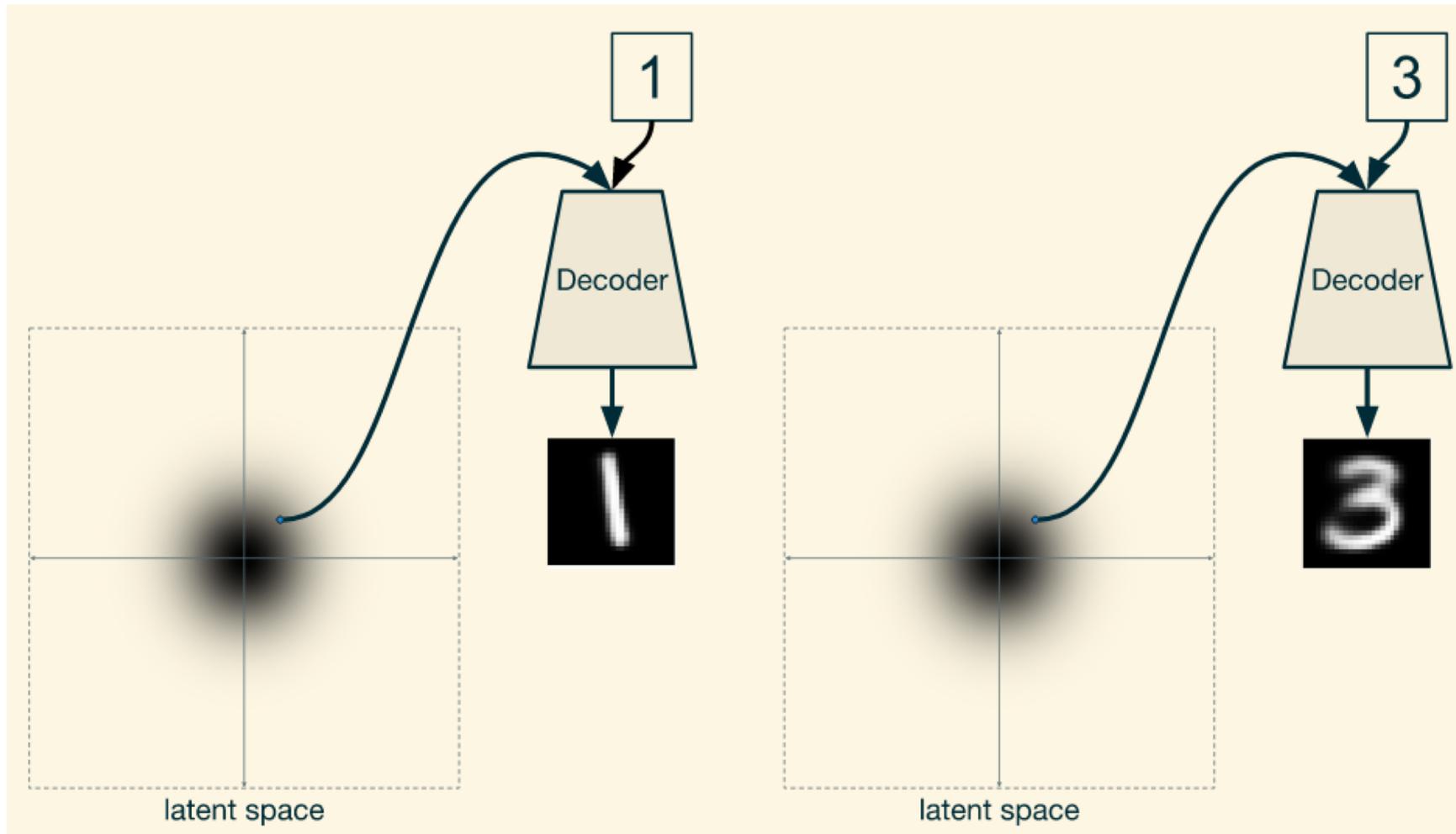
Decoder expects the latent space to be normally distributed  
Whether the sum of distributions produced by the encoder approximates a standard Normal distribution is measured by the KL divergence

# Conditional VAE



The number information is fed as a one-hot vector

# Generating Images from a VAE



Feed a random point in latent space and desired number.

Even if the same latent point is used for two different numbers, the process will work correctly since the latent space only encodes features such as stroke width or angle



