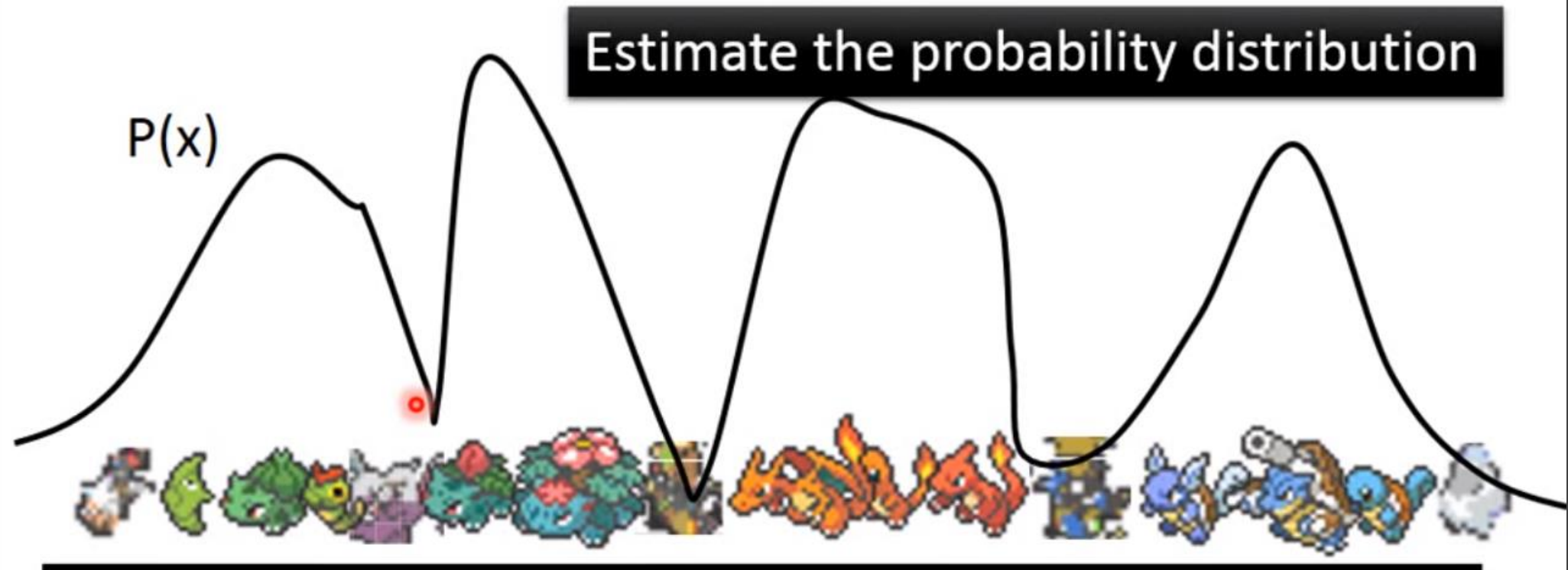


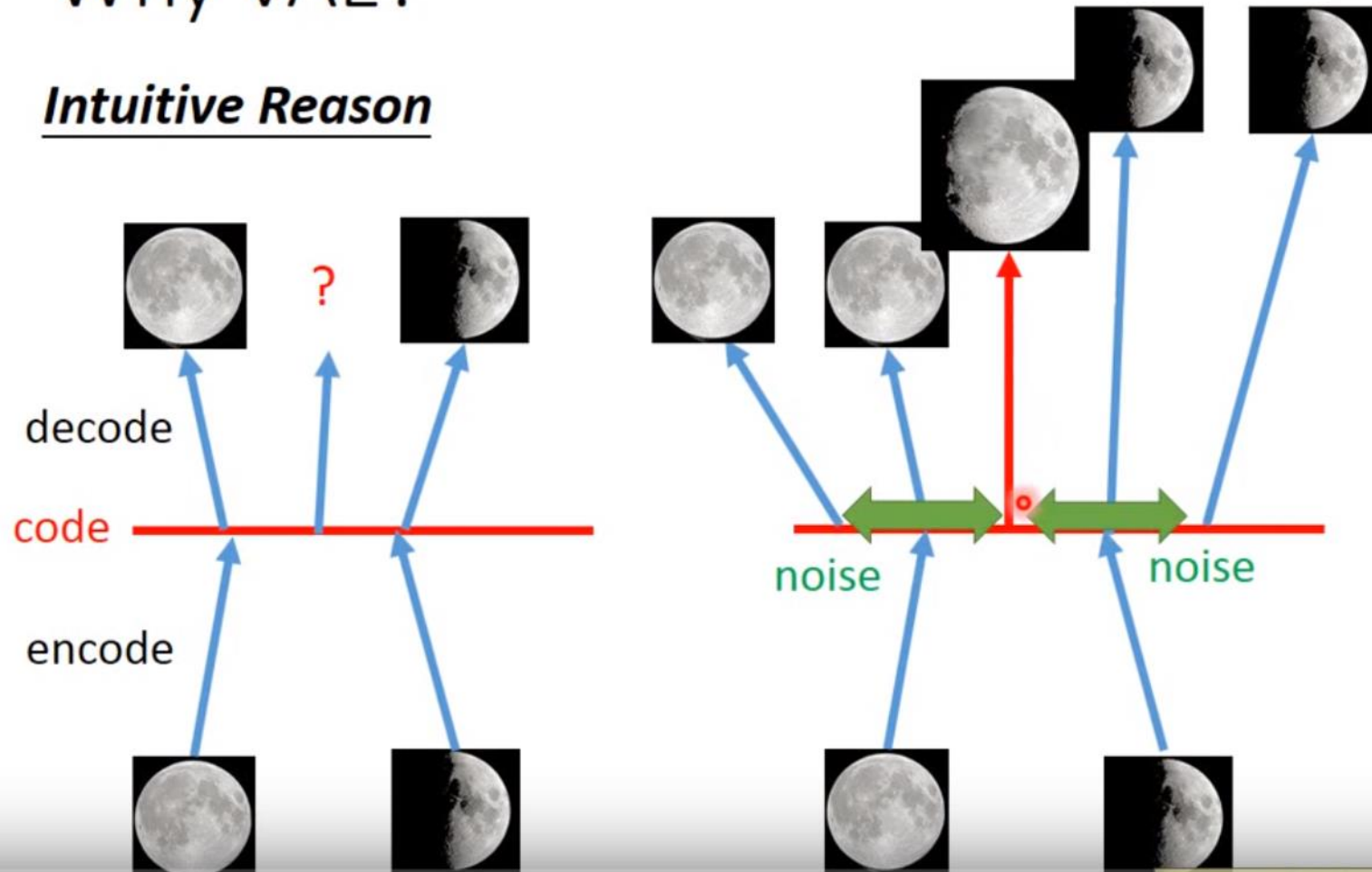
Why VAE?

- Back to what we want to do



Why VAE?

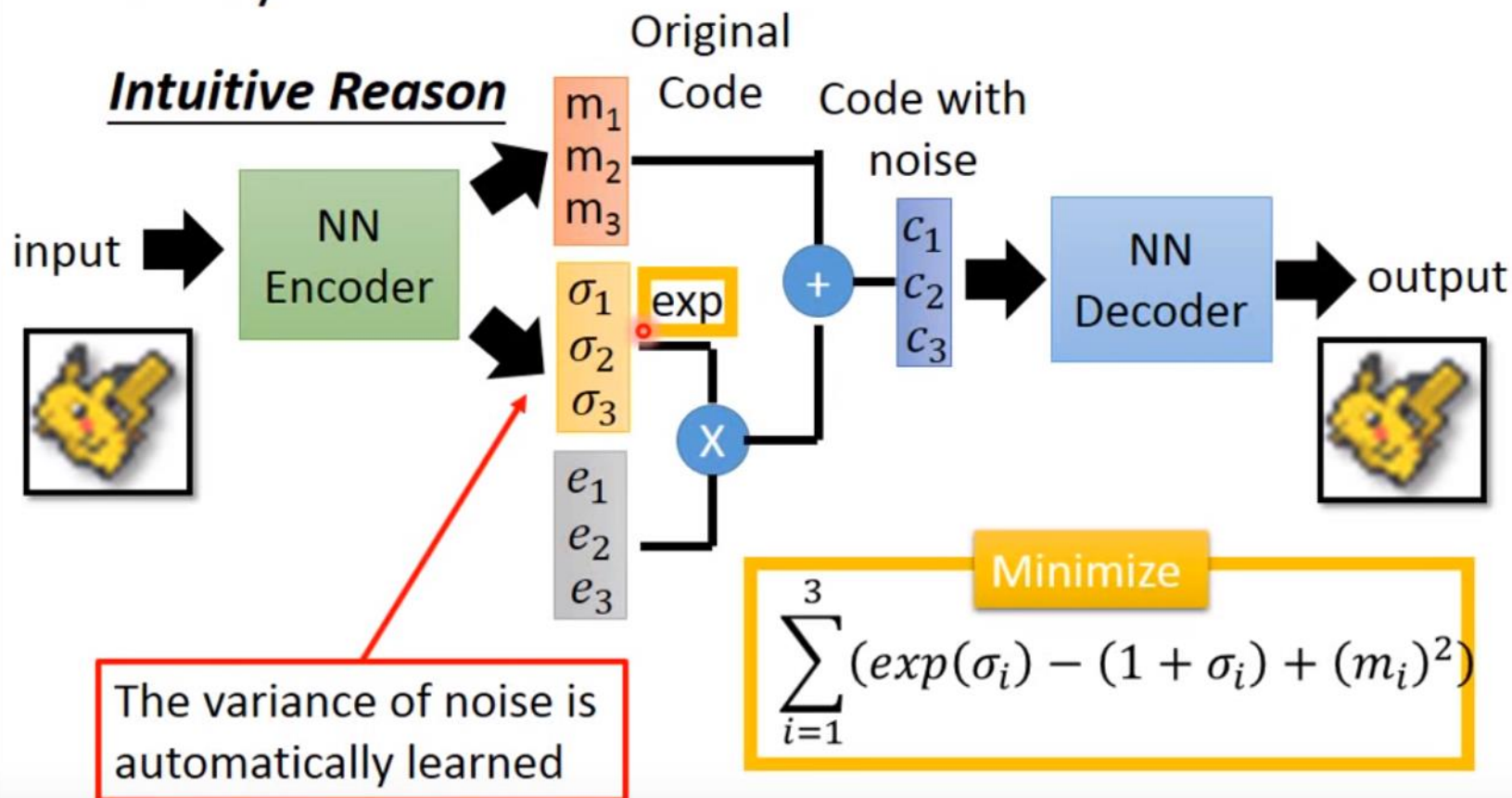
Intuitive Reason



Why VAE?

What will happen if we only minimize reconstruction error?

Intuitive Reason

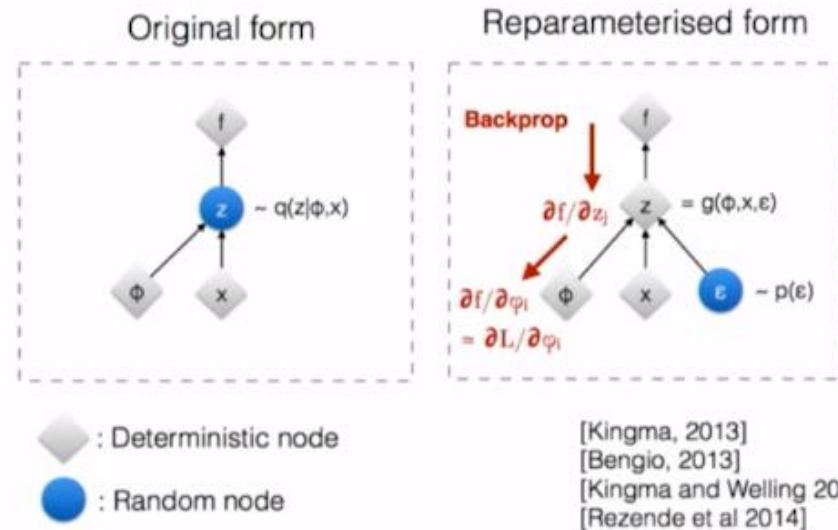


Re-Parameterization Trick - (The great thing about VAE)

Backpropagation not possible through random sampling!

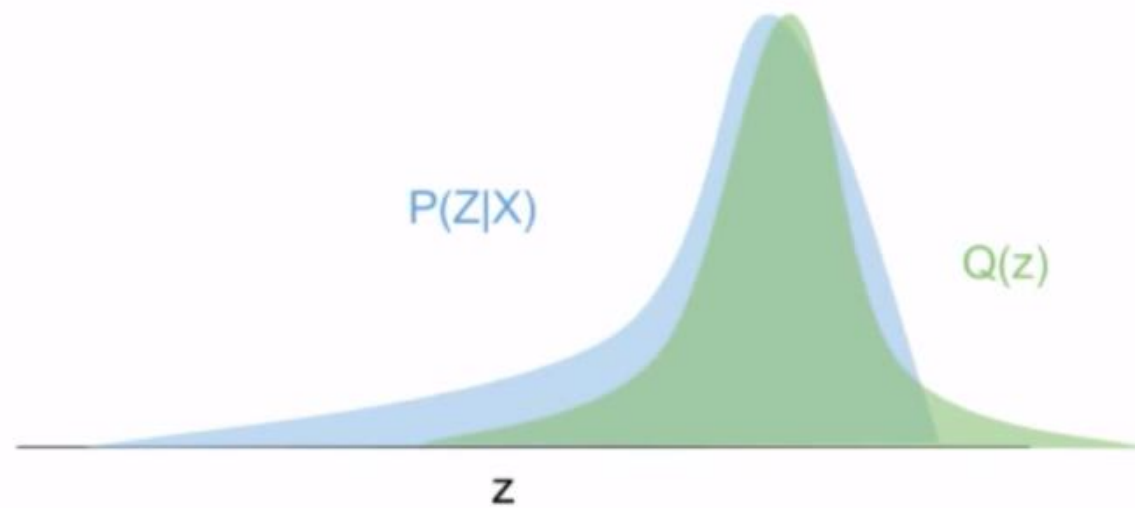
$$z^{(i,l)} = \mu^{(i)} + \sigma^{(i)} \odot \varepsilon_i$$

$$\varepsilon_i \sim N(0,1)$$



[Sampling Generative Networks <https://arxiv.org/abs/1609.04468>]

KL-Divergence - Background



Gaussian Mixture Model

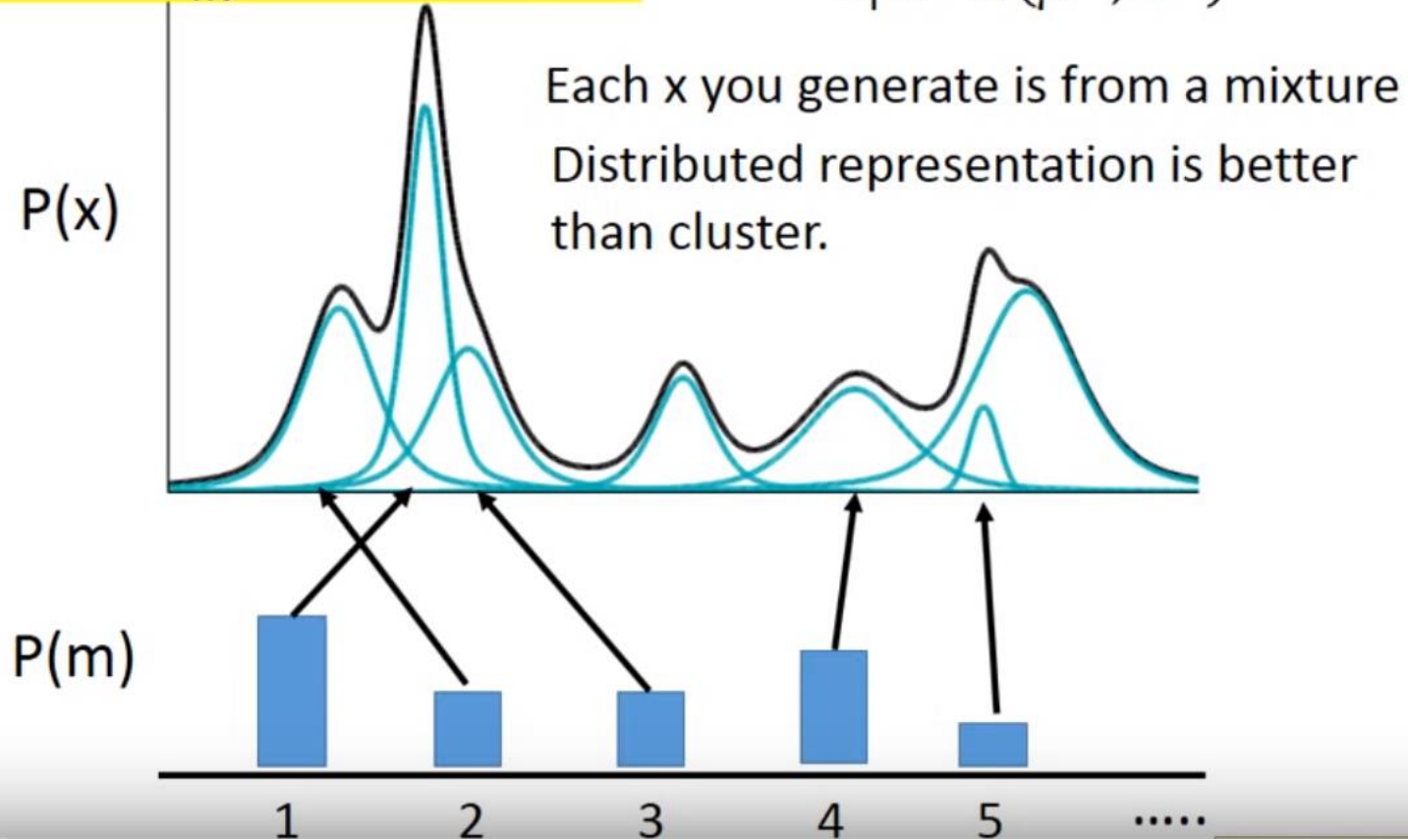
$$P(x) = \sum_m P(m)P(x|m)$$

How to sample?

$m \sim P(m)$ (multinomial)

m is an integer

$x|m \sim N(\mu^m, \Sigma^m)$



So VAE is version of Distributed representation of GMM

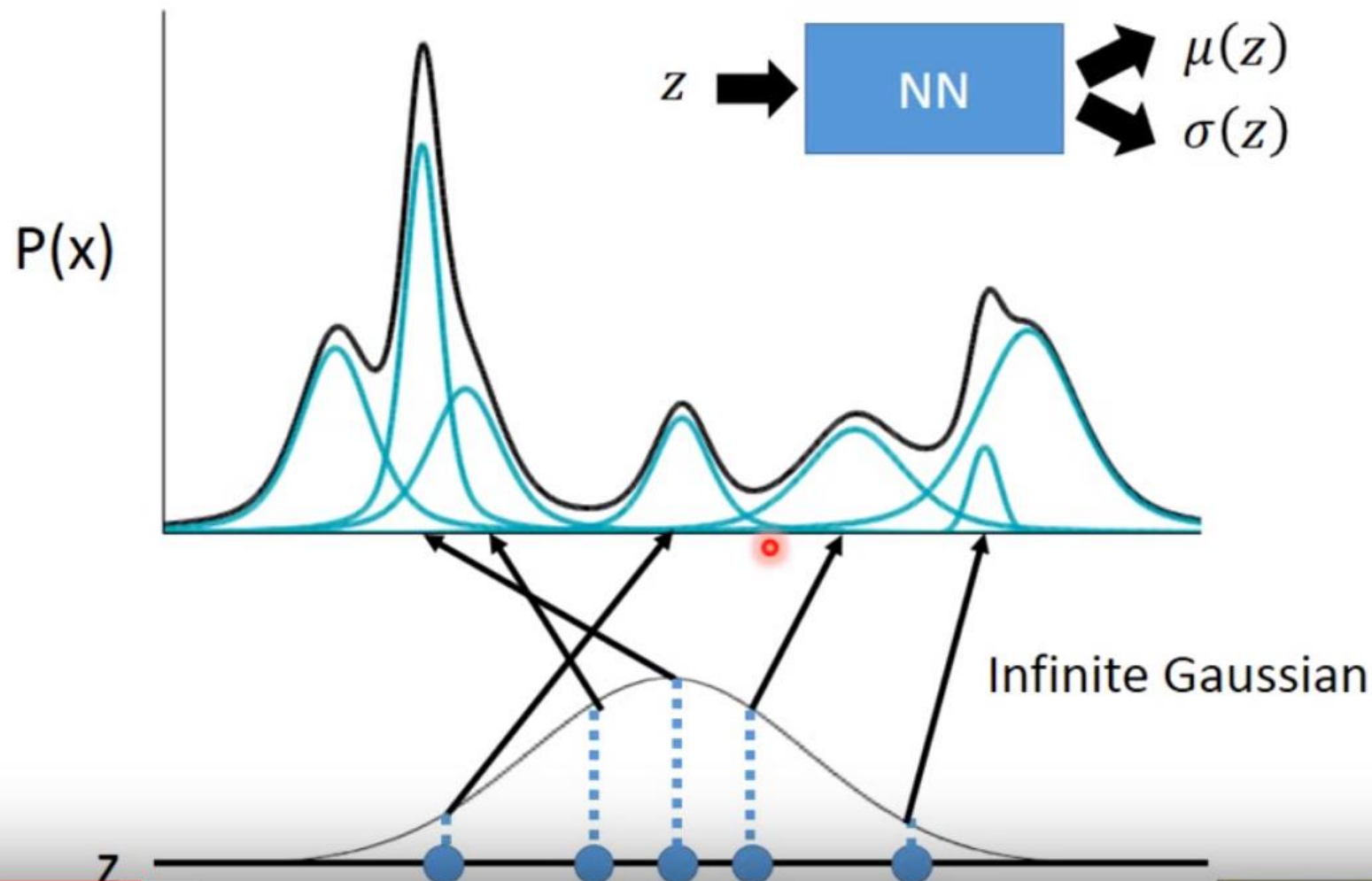
VAE

$$z \sim N(0, I)$$

z is a vector from normal distribution

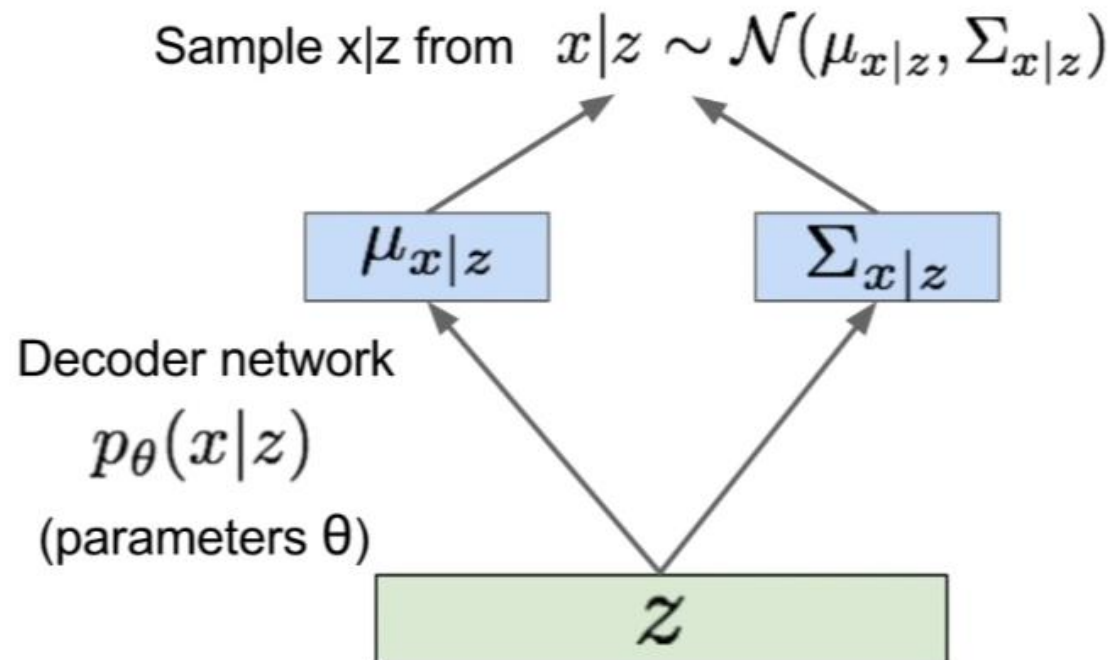
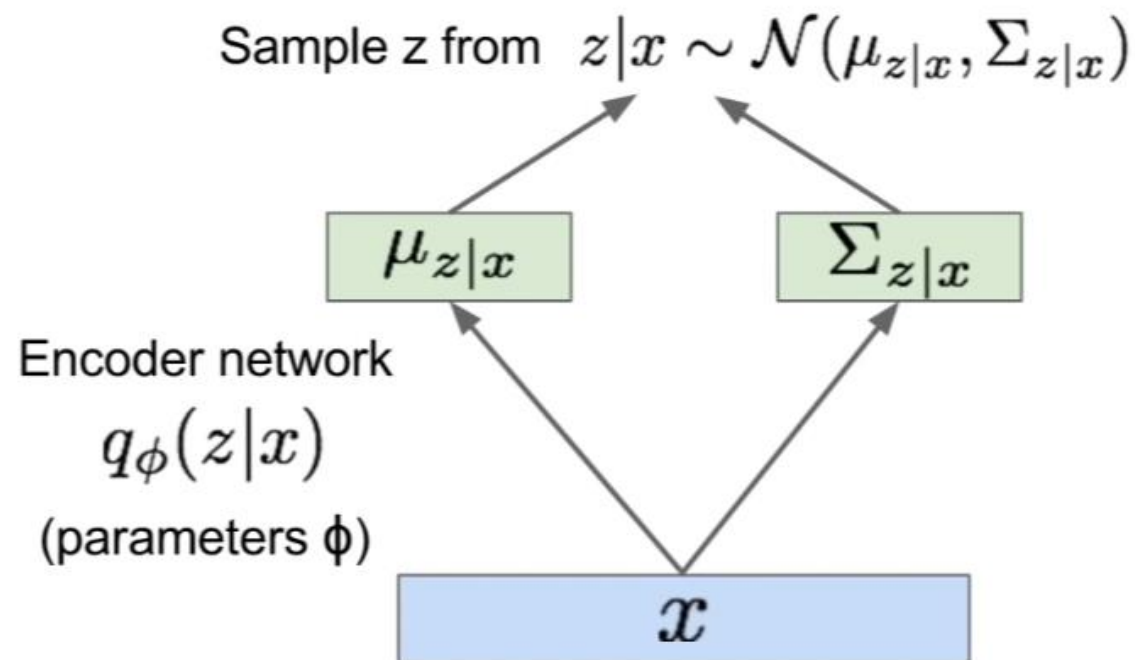
$$x|z \sim N(\mu(z), \sigma(z))$$

Each dimension of z represents an attribute



Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Variational Autoencoders

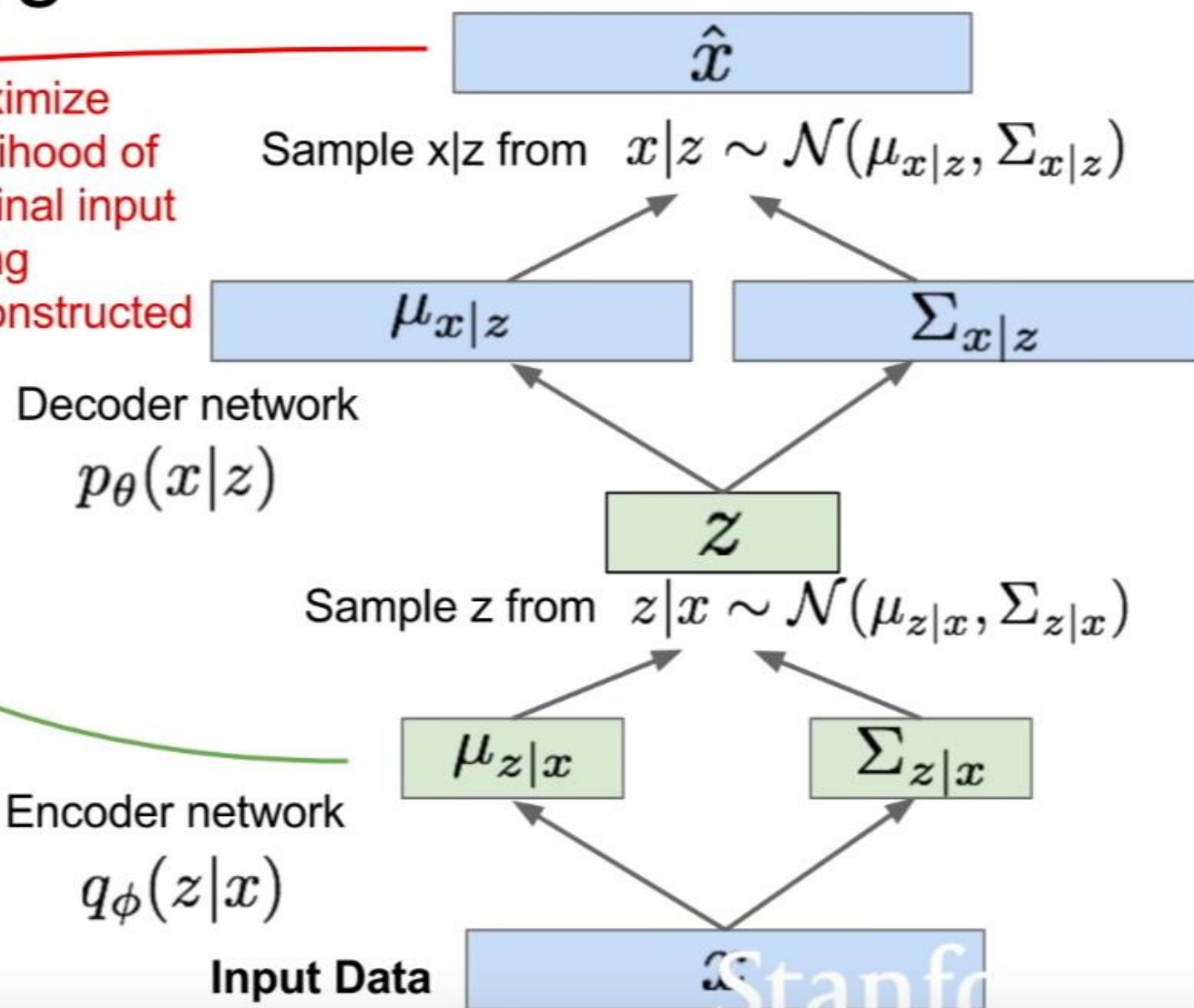
Putting it all together: maximizing the likelihood lower bound

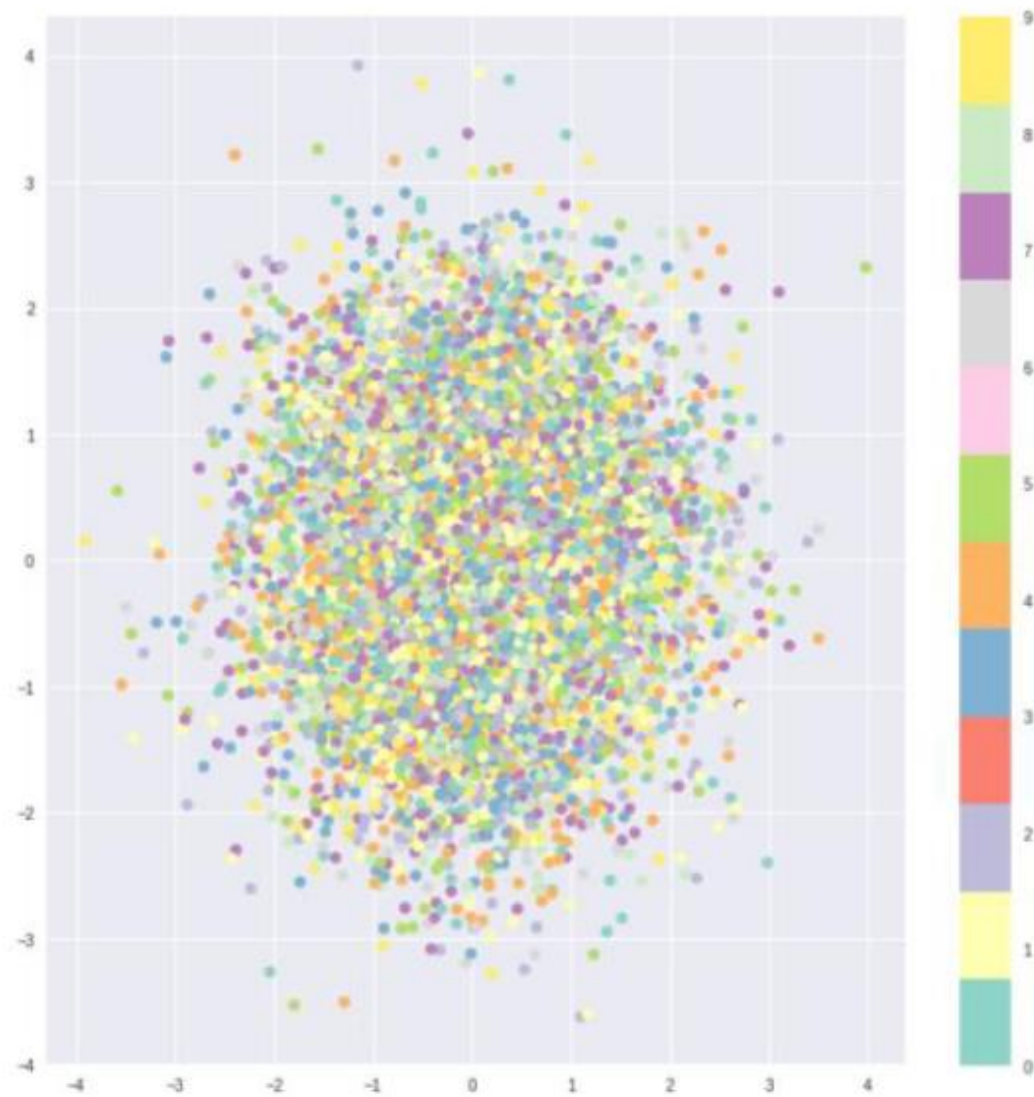
$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

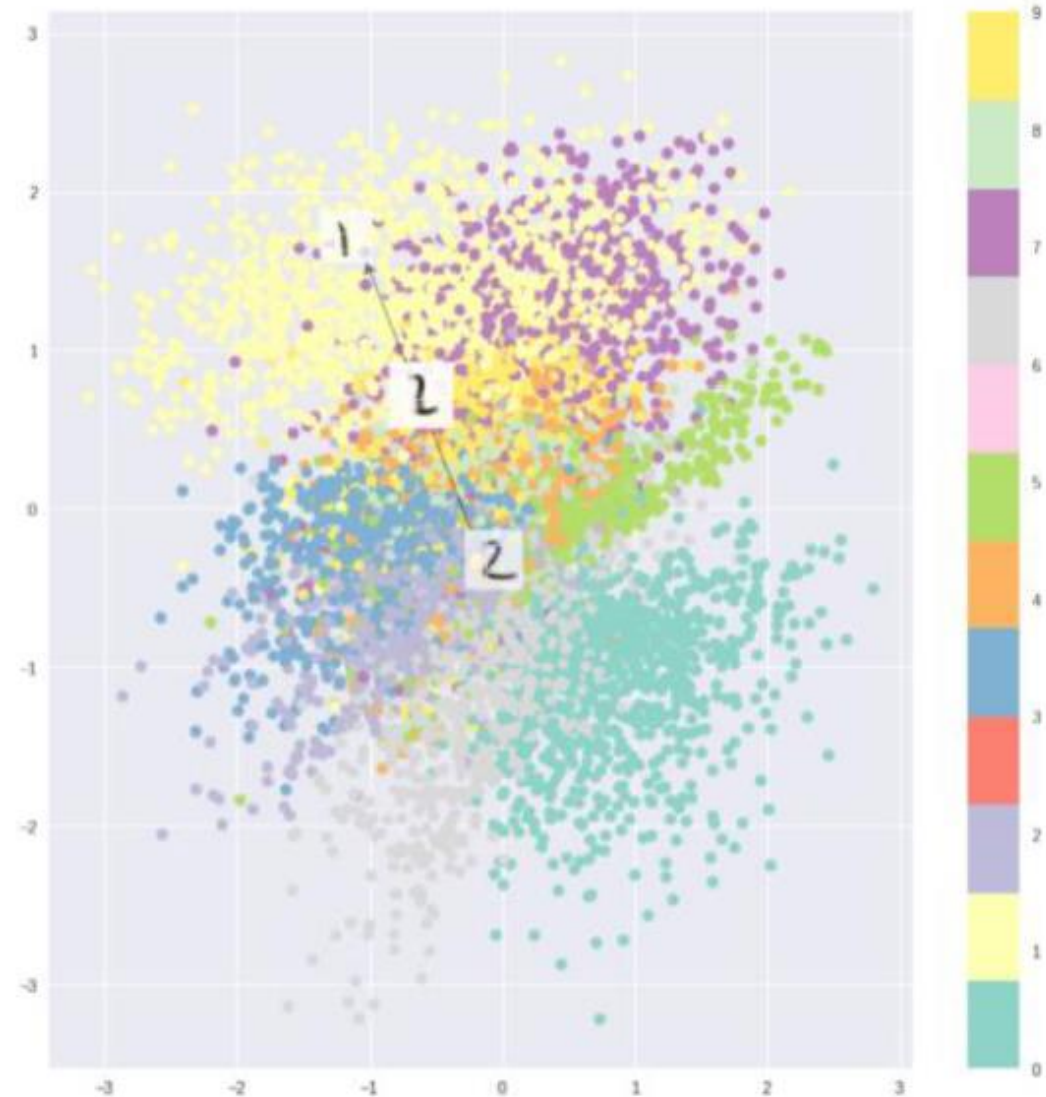
For every minibatch of input data: compute this forward pass, and then backprop!

Maximize likelihood of original input being reconstructed

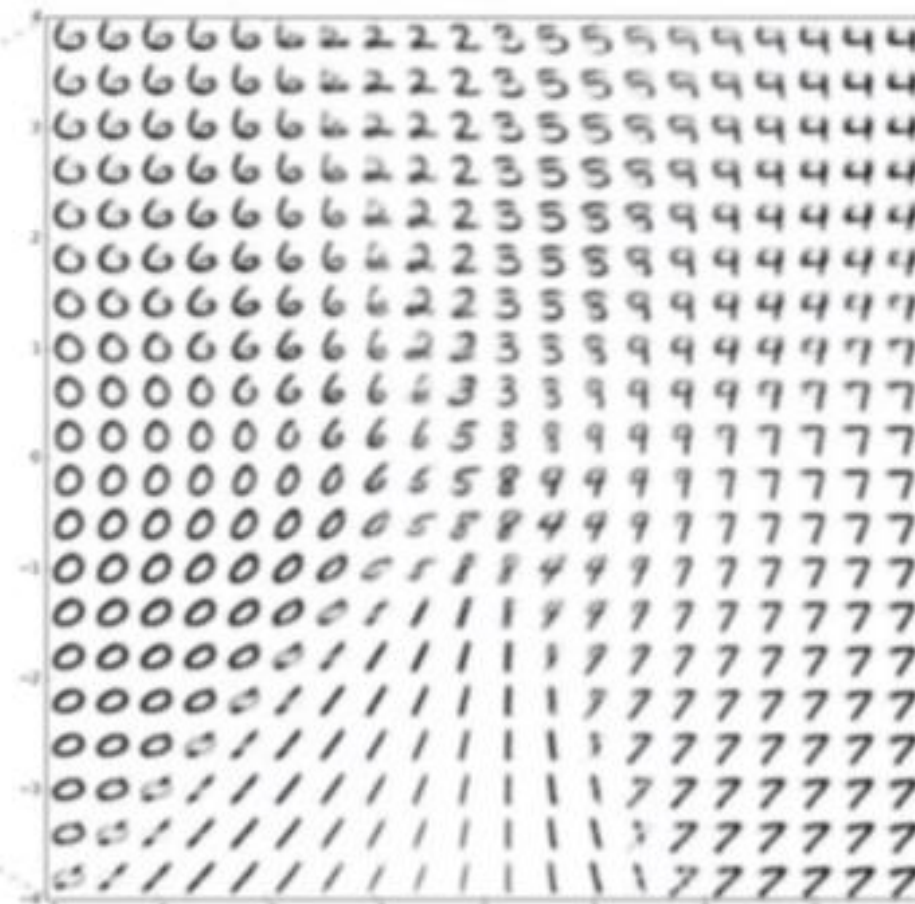
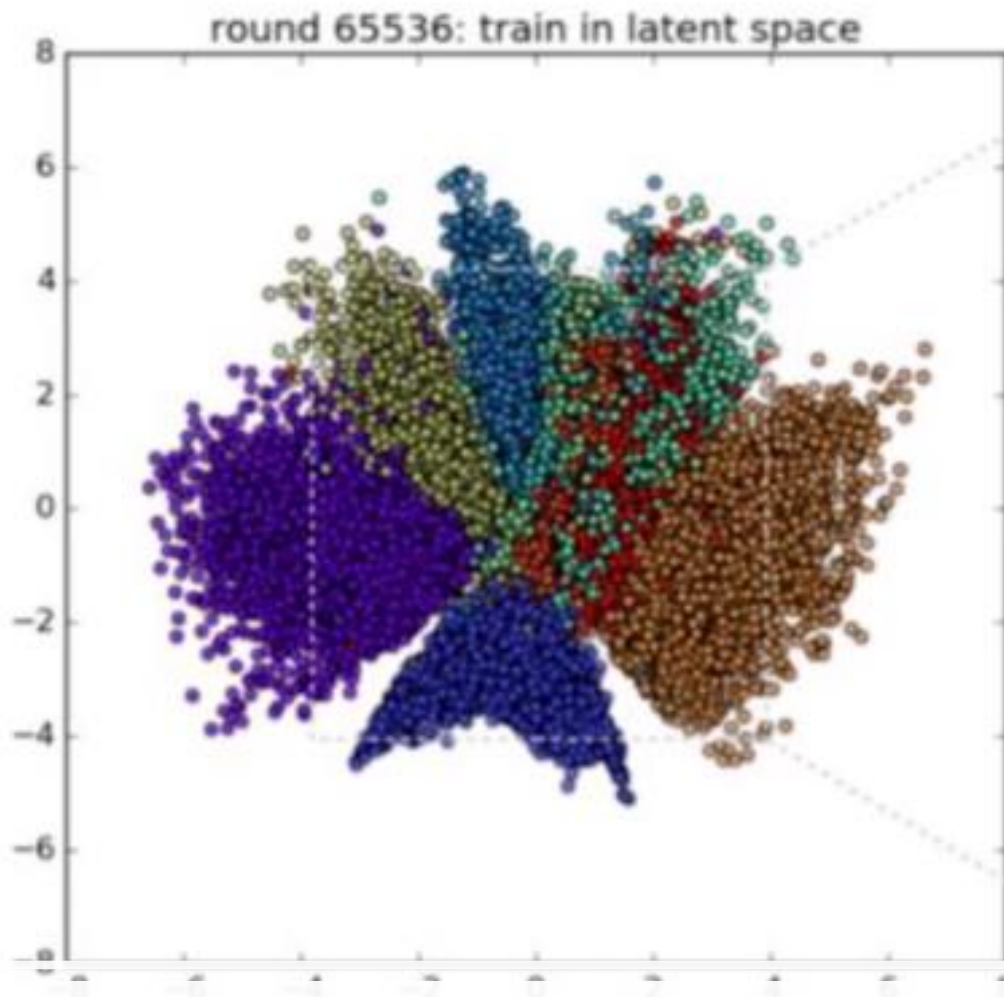
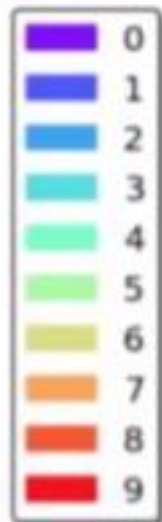
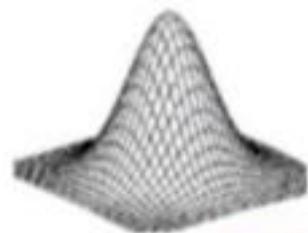




Optimizing using pure KL divergence loss



Optimizing using both reconstruction loss and KL divergence loss



Cross entropy Sampling