

# VAE: Variational Autoencoder

Rowel Atienza  
[github.com/roatienza](https://github.com/roatienza)  
*University of the Philippines*  
*2023*

# Generative Models Landscape (Vision)

GAN

VAE,  
AutoEncoders

Diffusion Models,  
Consistency Models

Generative Models

# GAN vs VAE

GAN focuses on modeling the input distribution,  $P(\mathbf{x})$

VAE focuses on the continuous latent space code to indirectly model the input distribution,  $P(\mathbf{x})$

$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \\ \theta^* &= \operatorname{argmax}_{\theta} \log \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \\ \theta^* &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta)\end{aligned}$$

GAN

$$P_{\theta}(\mathbf{x}) = \int P_{\theta}(\mathbf{z}|\mathbf{x})P(\mathbf{x})d\mathbf{z}$$

$$Q_{\phi}(\mathbf{z}|\mathbf{x}) \approx P_{\theta}(\mathbf{z}|\mathbf{x})$$

VAE

# GAN vs VAE

GAN tends to generate more realistic synthetic signals

~~VAE generated signals appear more blurry~~

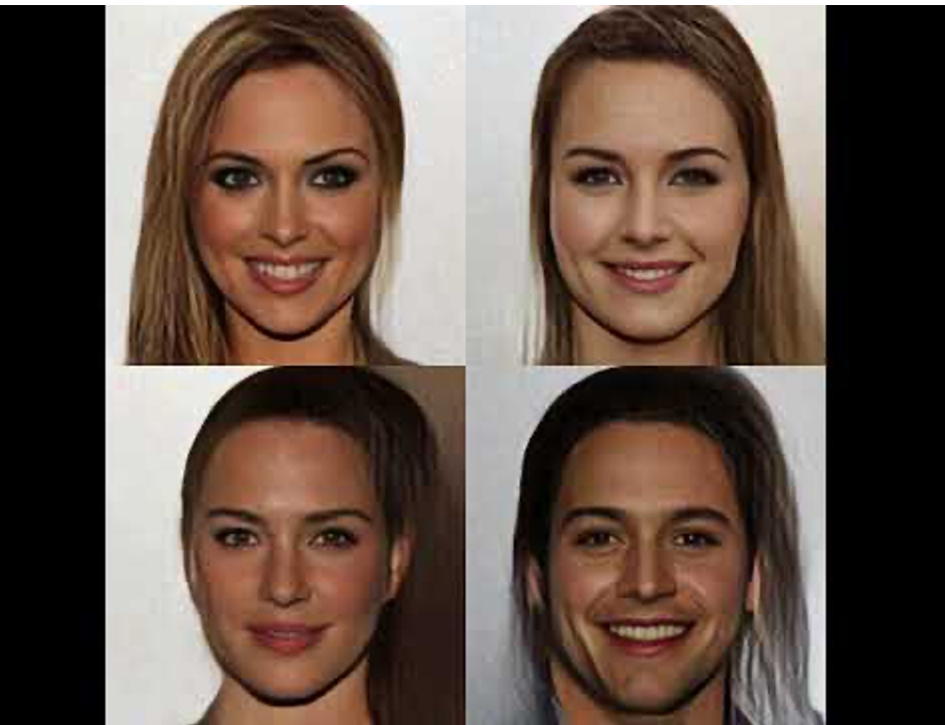
Very Deep VAEs Generalize Autoregressive Models and can Outperform them on Images (ICLR 2021)



Figure 1: **Selected samples from our very deep VAE on FFHQ-256, and a demonstration of the learned generative process.** VAEs can learn to first generate global features at low resolution,



# VAE



VAE



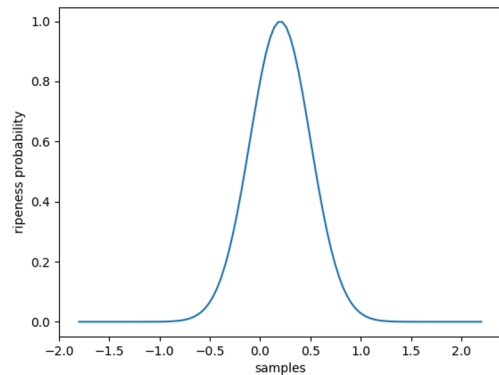
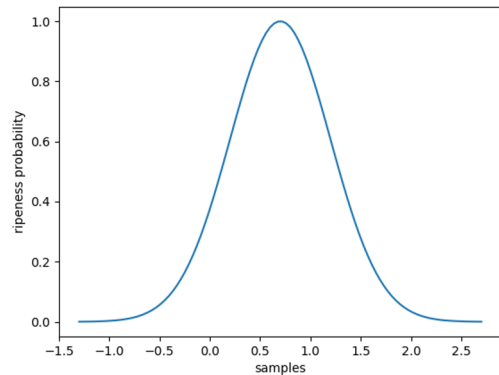
GAN

# GAN vs VAE

Unlike GAN, VAE provides a suitable variational framework for Bayesian inference and learning with latent variables

Robot Inference:

If the mango is ripe, I'll harvest it.



Input distribution  $P_{\theta}(\mathbf{x})$  and Joint distribution  $P_{\theta}(\mathbf{x}, \mathbf{z})$

$$\mathbf{x} \sim P_{\theta}(\mathbf{x})$$

In machine learning, we are interested in finding  $P_{\theta}(\mathbf{x}, \mathbf{z})$ , the joint distribution between  $\mathbf{x}$  and latent code  $\mathbf{z}$

$$P_{\theta}(\mathbf{x}) = \int P_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

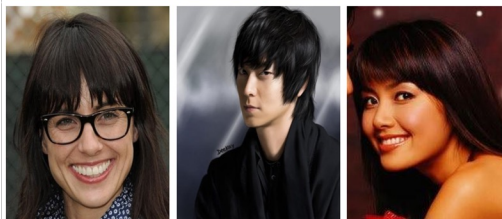
However,  $P_{\theta}(\mathbf{x}, \mathbf{z})$  is intractable or it does not have an analytic form or a good estimator.

# z as vector of attributes

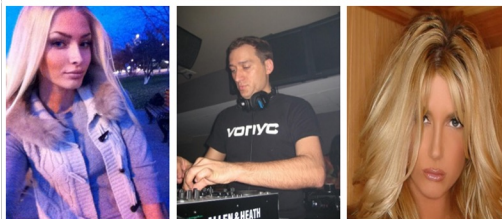
Eyeglasses



Bangs



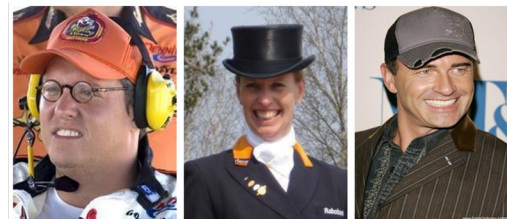
Pointy  
Nose



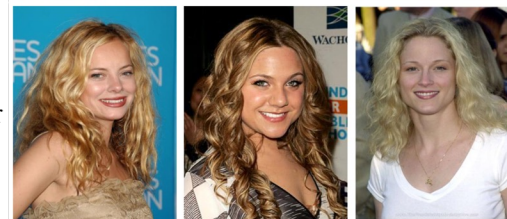
Oval Face



Wearing  
Hat



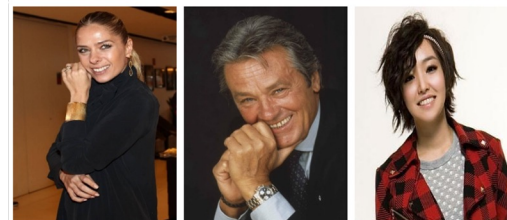
Wavy Hair



Mustache



Smiling



## Using Bayes Theorem

$$P_{\theta}(\mathbf{x}) = \int P_{\theta}(\mathbf{x}|\mathbf{z})P(\mathbf{z})d\mathbf{z}$$

A neural network can easily ignore  $P_{\theta}(\mathbf{x}|\mathbf{z})$  to come up with a trivial solution  $P_{\theta}(\mathbf{x}|\mathbf{z}) = P_{\theta}(\mathbf{x})$ .

$$P_{\theta}(\mathbf{x}) = \int P_{\theta}(\mathbf{z}|\mathbf{x})P(\mathbf{x})d\mathbf{z}$$

However,  $P_{\theta}(\mathbf{z}|\mathbf{x})$  is also intractable.

# Variational Inference Model (Encoder)

$$Q_{\phi}(\mathbf{z}|\mathbf{x}) \approx P_{\theta}(\mathbf{z}|\mathbf{x})$$

$Q_{\phi}(\mathbf{z}|\mathbf{x})$  provides a good estimate of  $P_{\theta}(\mathbf{z}|\mathbf{x})$

$Q_{\phi}(\mathbf{z}|\mathbf{x})$  is parametric and tractable.

$Q_{\phi}(\mathbf{z}|\mathbf{x})$  can be approximated by deep neural networks by optimizing the parameters  $\phi$ .

# Variational Inference Model (Encoder)

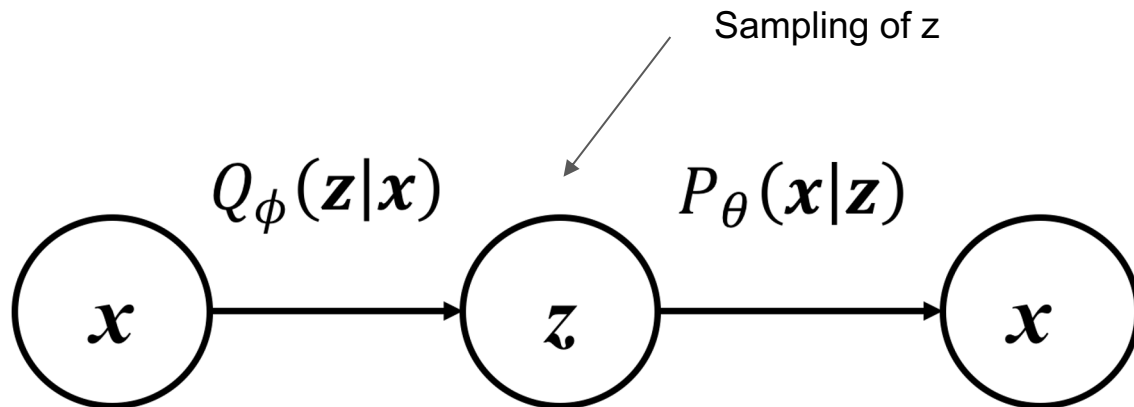
Typically,  $Q_{\phi}(\mathbf{z}|\mathbf{x})$  is chosen to be a multivariate Gaussian:

$$Q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}\left(\mathbf{z}; \boldsymbol{\mu}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}(\mathbf{x}))\right)$$

Both mean,  $\boldsymbol{\mu}(\mathbf{x})$ , and standard deviation,  $\boldsymbol{\sigma}(\mathbf{x})$ , are computed by the encoder neural network from the input data points.

The diagonal matrix implies that the elements of  $\mathbf{z}$  are independent.

# Probabilistic Graphical Model of VAE



To estimate  $P_{\theta}(x)$ , we must identify its relationship with  $Q_{\phi}(z|x)$  and  $P_{\theta}(x|z)$ .



# Core Equation of VAE

If  $Q_\phi(\mathbf{z}|\mathbf{x})$  is an estimate of  $P_\theta(\mathbf{z}|\mathbf{x})$ , the Kullback-Leibler (KL) divergence determines the distance between these two conditional densities:

$$D_{KL} \left( Q_\phi(\mathbf{z}|\mathbf{x}) \| P_\theta(\mathbf{z}|\mathbf{x}) \right) = \mathbb{E}_{\mathbf{z} \sim Q} [\log Q_\phi(\mathbf{z}|\mathbf{x}) - \log P_\theta(\mathbf{z}|\mathbf{x})]$$

Using Bayes Theorem:

$$\begin{aligned} & D_{KL} \left( Q_\phi(\mathbf{z}|\mathbf{x}) \| P_\theta(\mathbf{z}|\mathbf{x}) \right) \\ &= \mathbb{E}_{\mathbf{z} \sim Q} [\log Q_\phi(\mathbf{z}|\mathbf{x}) - \log P_\theta(\mathbf{x}|\mathbf{z}) - \log P_\theta(\mathbf{z})] + \log P_\theta(\mathbf{x}) \end{aligned}$$

# Core Equation of VAE

Recognizing  $\log Q_\phi(\mathbf{z}|\mathbf{x}) - \log P_\theta(\mathbf{z}) = D_{KL} \left( Q_\phi(\mathbf{z}|\mathbf{x}) \| P_\theta(\mathbf{z}) \right)$ :

$$\log P_\theta(\mathbf{x}) - D_{KL} \left( Q_\phi(\mathbf{z}|\mathbf{x}) \| P_\theta(\mathbf{z}|\mathbf{x}) \right) = \mathbb{E}_{\mathbf{z} \sim Q} [\log P_\theta(\mathbf{x}|\mathbf{z})] - D_{KL} \left( Q_\phi(\mathbf{z}|\mathbf{x}) \| P_\theta(\mathbf{z}) \right)$$

↑  
Term being  
maximized

↑  
Distance between Q and  
true P: Approx Zero

↑  
Decoder  
(Reconstruction)

↑  
Distance between Q  
and Prior P(z)

Left Side is known as Evidence  
Lower Bound (ELBO) of  $\log_\theta P(\mathbf{x})$

# Core Equation of VAE

Maximizing ELBO by optimizing the parameters  $\phi$  and  $\theta$  of the neural network means:

- $D_{KL} \left( Q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel P_{\theta}(\mathbf{z}|\mathbf{x}) \right) \rightarrow 0$  or the inference model is getting better in encoding the attributes of  $\mathbf{x}$  in  $\mathbf{z}$ .
- $\log P_{\theta}(\mathbf{x}|\mathbf{z})$  on the right-hand side of Equation 8.1.10 is maximized or the decoder model is getting better in reconstructing  $\mathbf{x}$  from the latent vector  $\mathbf{z}$ .

# Optimization

- The decoder term  $\mathbb{E}_{\mathbf{z} \sim Q} [\log P_{\theta}(\mathbf{x}|\mathbf{z})]$  means that the generator takes  $\mathbf{z}$  samples from the output of the inference model to reconstruct the inputs.
- Maximizing this term implies that we minimize a *Reconstruction Loss*,  $\mathcal{L}_R$ .
- If the image distribution is assumed to be Gaussian, MSE can be used.
- If every pixel is considered a Bernoulli distribution, the loss function is a binary cross entropy.

# Optimization

- The second term,  $-D_{KL} \left( Q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel P_{\theta}(\mathbf{z}) \right)$ , turns out to be straightforward to evaluate.
- $Q_{\phi}$  is a Gaussian distribution.
- Typically,  $P_{\theta}(\mathbf{z}) = P(\mathbf{z}) = \mathcal{N}(0, I)$  is also a Gaussian with zero mean and standard deviation equal to 1.0. The KL term simplifies to:

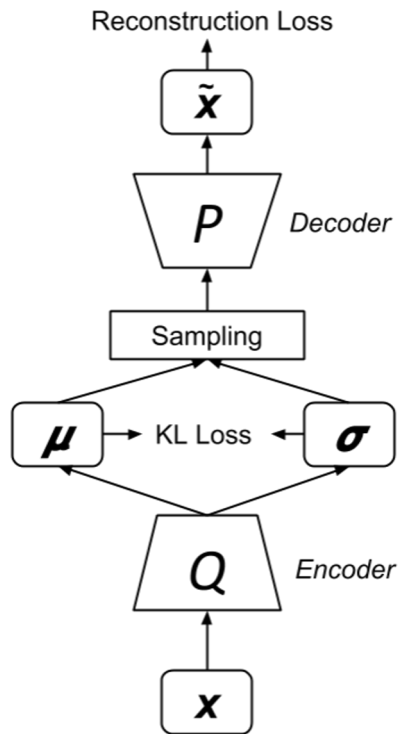
$$-\mathcal{L}_{KL} = -D_{KL} \left( Q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel P_{\theta}(\mathbf{z}) \right) = \frac{1}{2} \sum_{j=1}^J \left( 1 + \log(\sigma_j)^2 - (\mu_j)^2 - (\sigma_j)^2 \right)$$

where  $J$  is the dimensionality of  $\mathbf{z}$ . Both  $\mu_j$  and  $\sigma_j$  are functions of  $\mathbf{x}$  computed through the inference model. To maximize  $-D_{KL}$ ,  $\sigma_j \rightarrow 1$  and  $\mu_j \rightarrow 0$ .

## VAE Loss Function

$$\mathcal{L}_{VAE} = \mathcal{L}_R + \mathcal{L}_{KL}$$

# Training VAE



Without Reparameterization Trick

# Training VAE

- No problem with forward computation in the network
- Problem is backpropagation gradients will not pass through the stochastic **Sampling** block
- While it is fine to have stochastic inputs for neural networks, it is not possible for gradients to go through a stochastic layer.

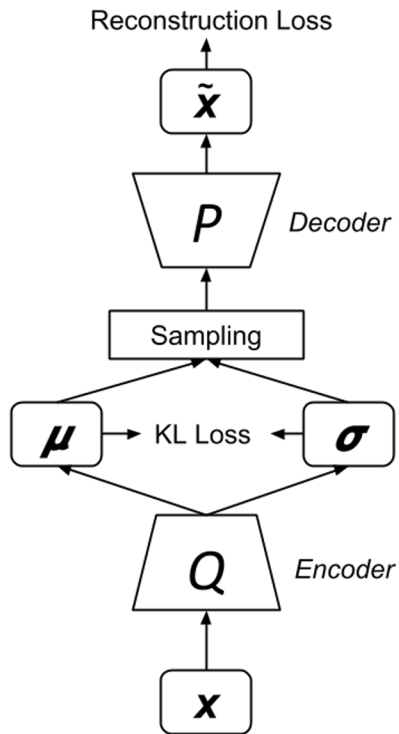


# Reparameterization Trick

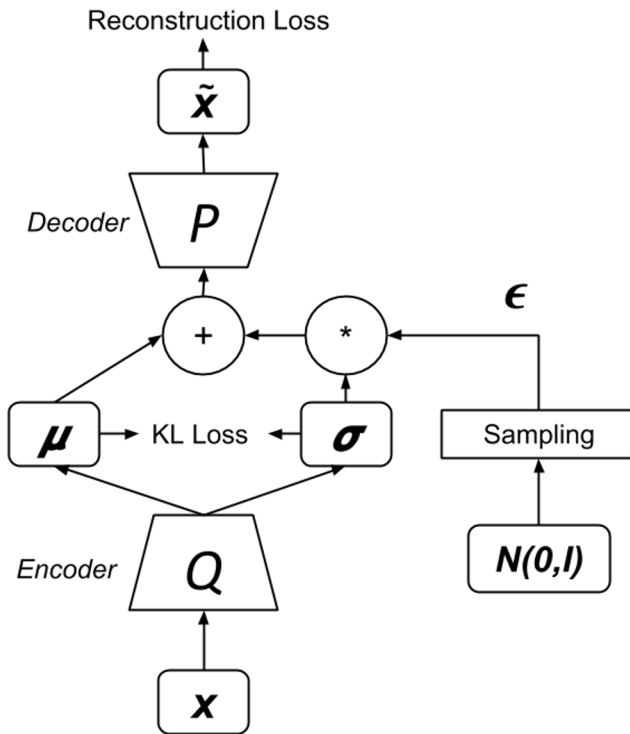
- The solution to this problem is to push out the **Sampling** process
- The solution is to sample  $\epsilon$  from an isotropic Gaussian distribution and adjust the value using the encoder predicted  $\mu$  and  $\sigma$
- Compute the sample as:

$$\text{Sample} = \mu + \epsilon\sigma$$

# Training VAE using Reparameterization Trick



Without Reparameterization Trick



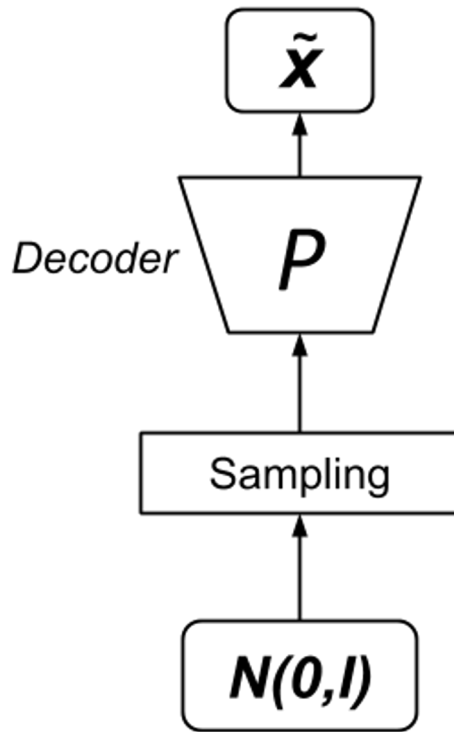
With Reparameterization Trick

# Decoder Testing

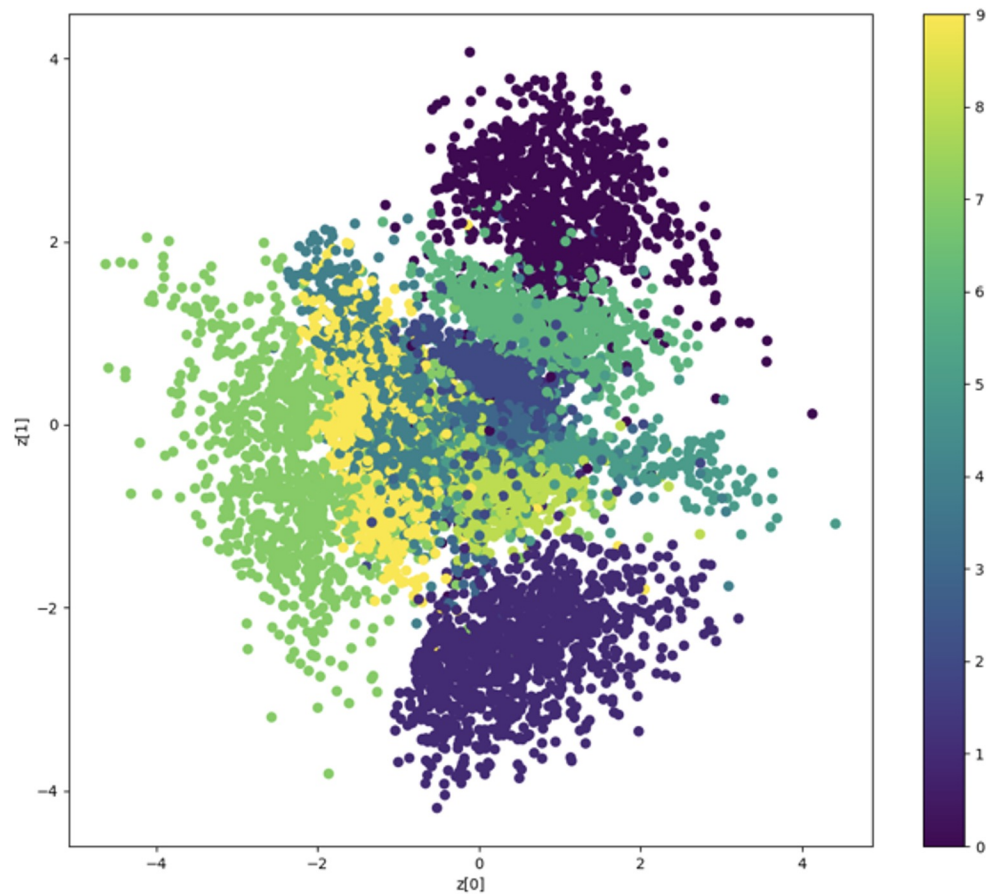
No need for the encoder

Sample directly from isotropic unit

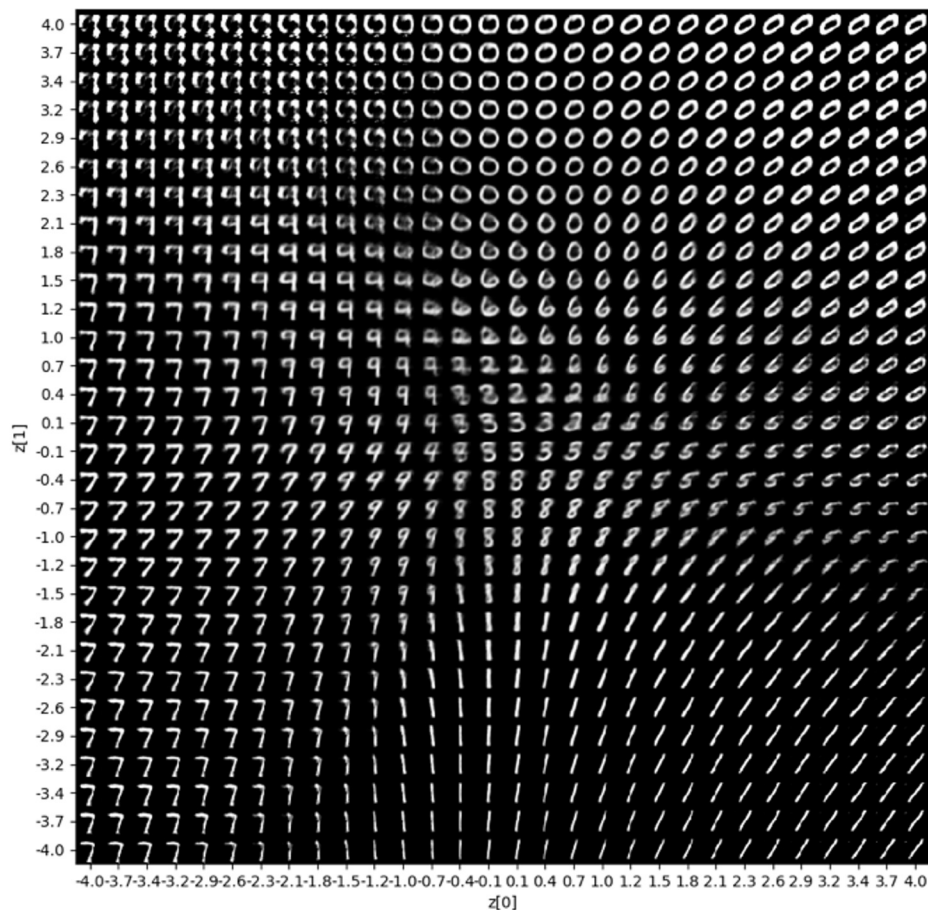
Gaussian to generate  $\mathbf{x}$



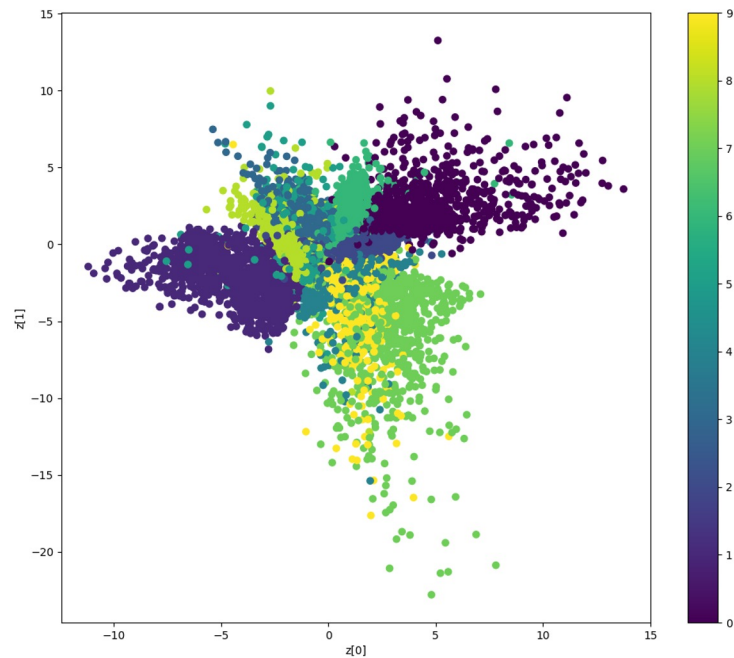
Example distribution  
for 2-dim  $z$  (VAE MLP)



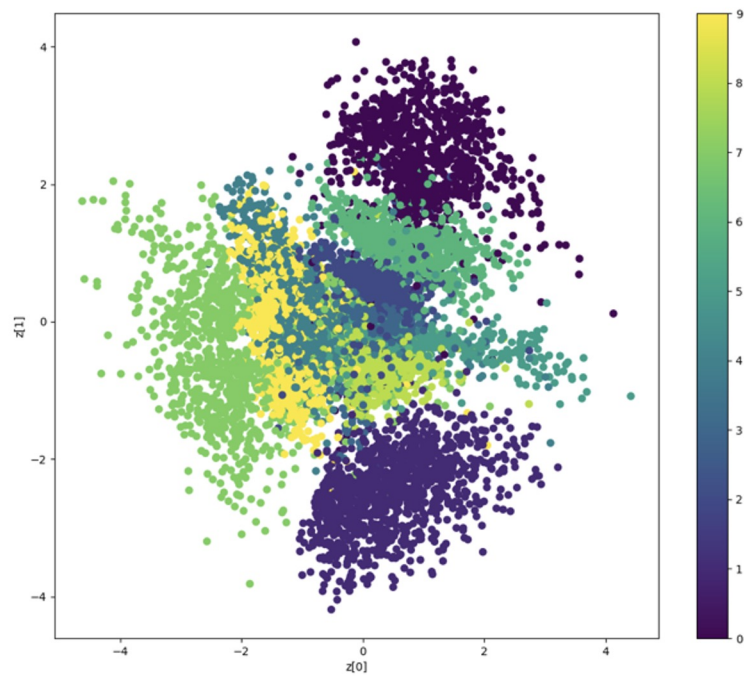
Example distribution for 2-dim  $z$   
(VAE MLP)



# 2-dim Autoencoder vs 2-dim VAE

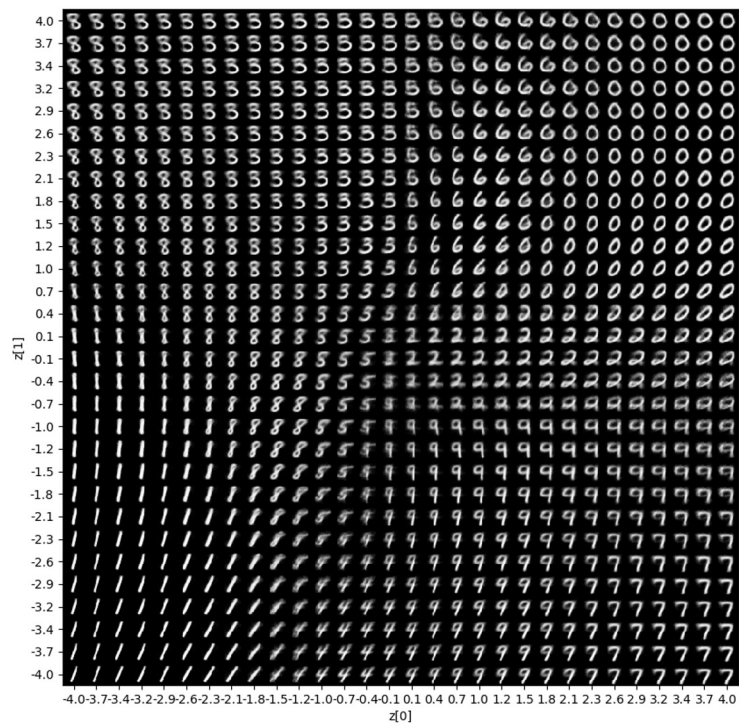


Autoencoder

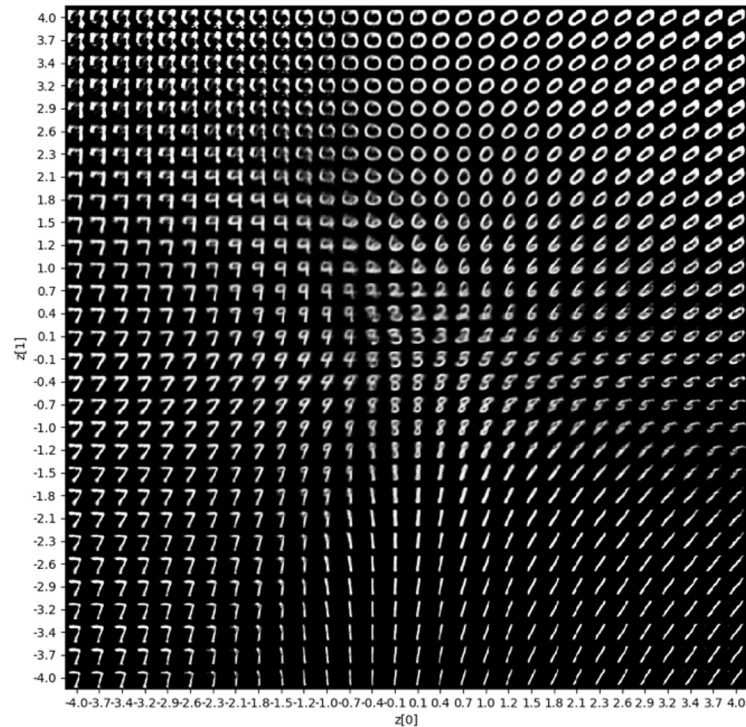


VAE

# 2-dim Autoencoder vs 2-dim VAE



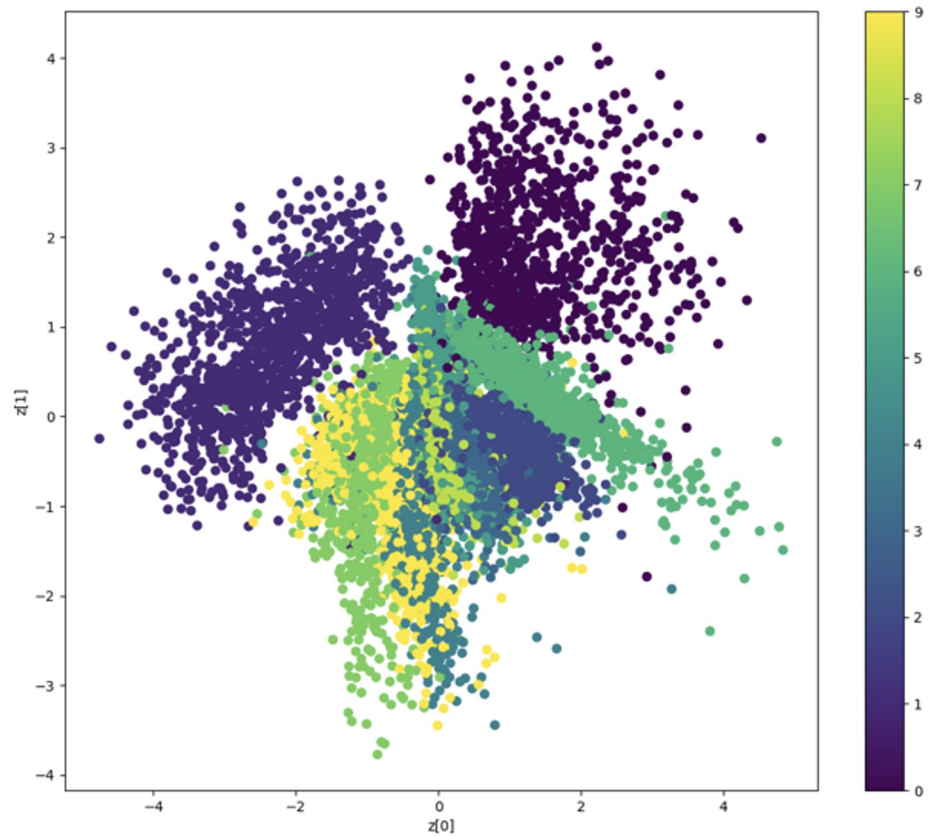
Autoencoder



VAE

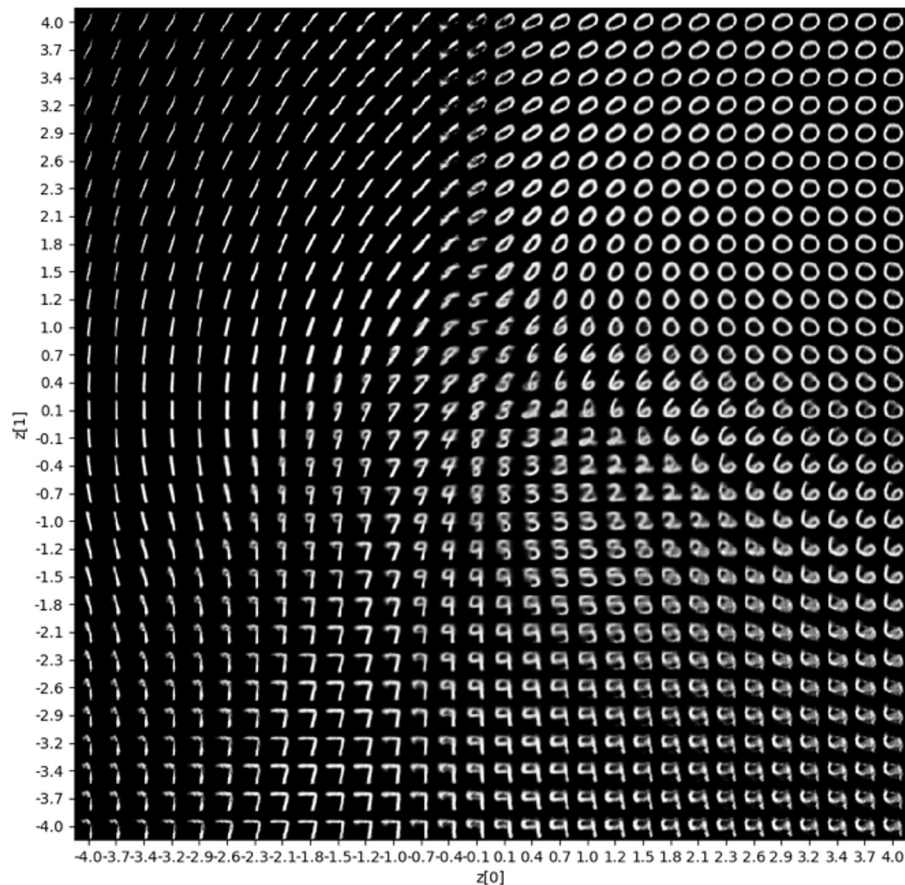
Example  
distribution for  
2-dim  $z$

(VAE CNN)





Example distribution for 2-dim  $z$   
(VAE CNN)



# Conditional VAE

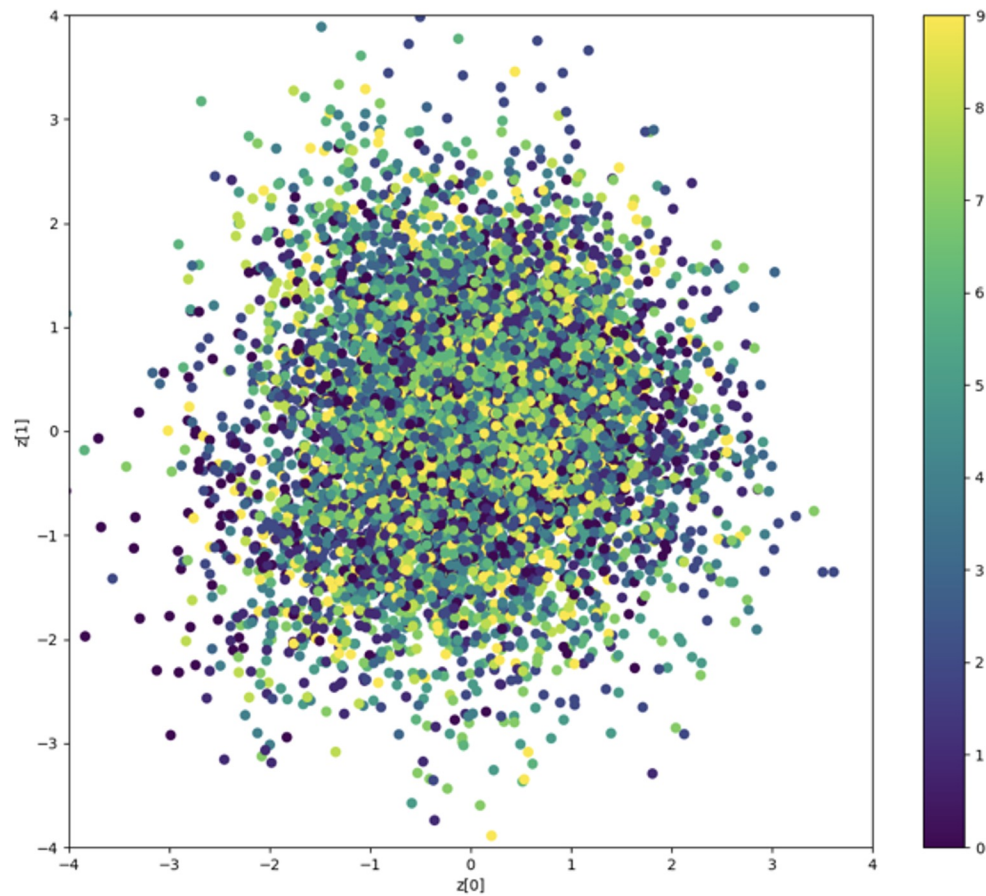
$$\log P_{\theta}(\mathbf{x}|\mathbf{c}) - D_{KL} \left( Q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{c}) \| P_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{c}) \right) = \mathbb{E}_{\mathbf{z} \sim Q} [\log P_{\theta}(\mathbf{x}|\mathbf{z}, \mathbf{c})] - D_{KL} \left( Q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{c}) \| P_{\theta}(\mathbf{z}|\mathbf{c}) \right)$$

## No change in the loss function

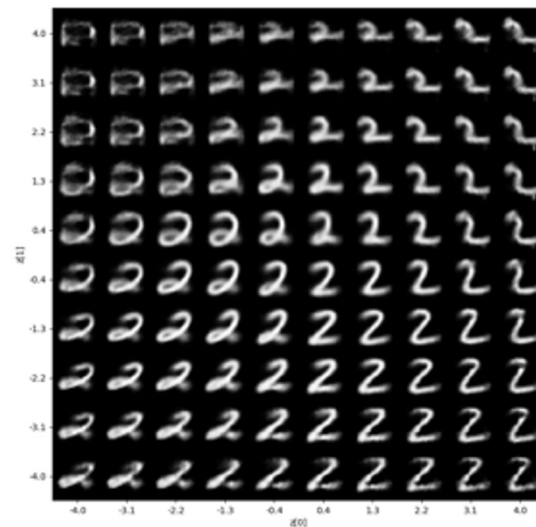
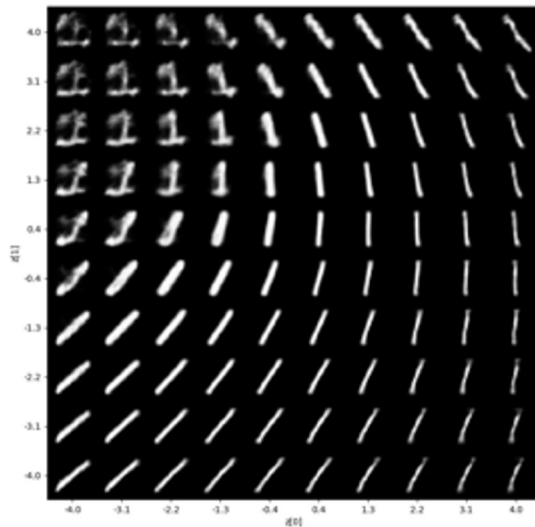
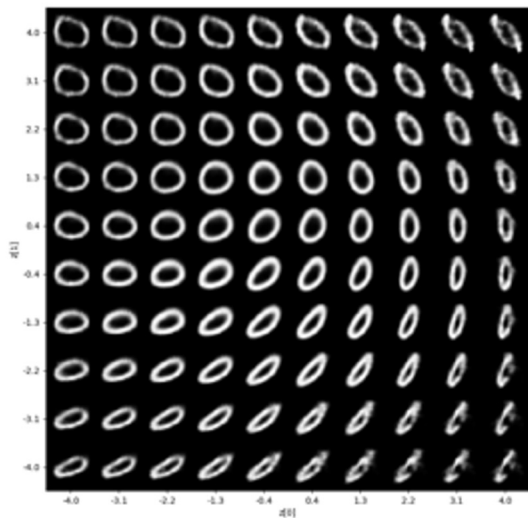
- *Reconstruction Loss* of the decoder given both the latent vector and the condition
- *KL Loss* between the encoder given both the latent vector and the condition and the prior distribution given the condition. Similar to VAE, typically we choose  $P_{\theta}(\mathbf{z}|\mathbf{c}) = P(\mathbf{z}|\mathbf{c}) = \mathcal{N}(0, I)$ .

Example distribution for  
2-dim  $z$

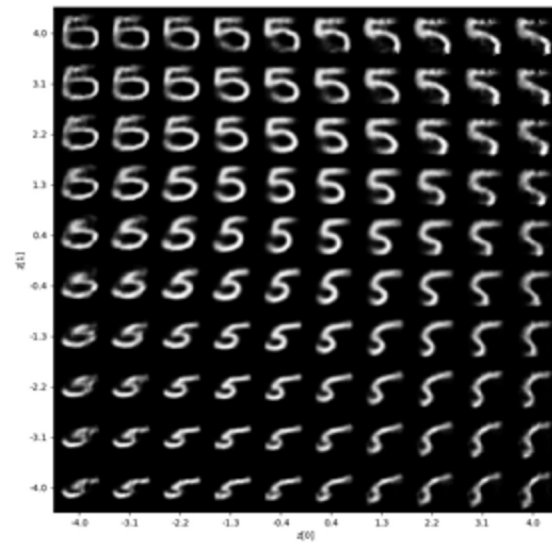
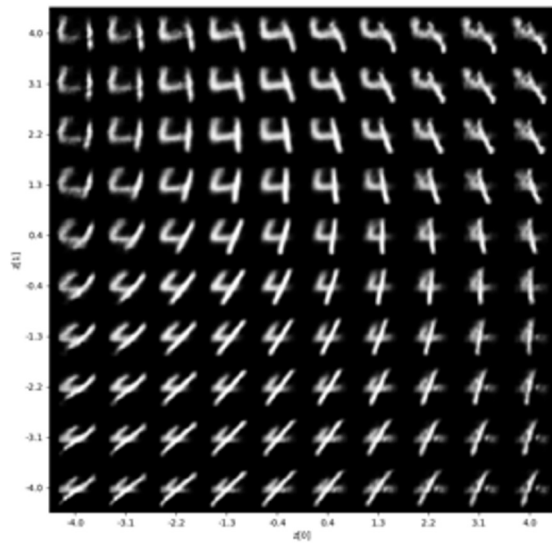
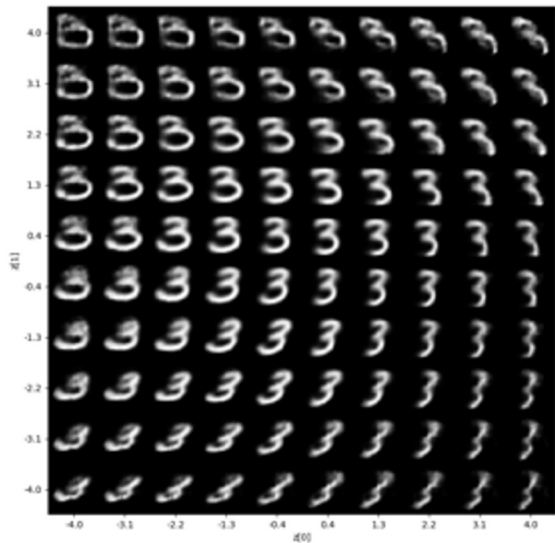
(VAE CNN)

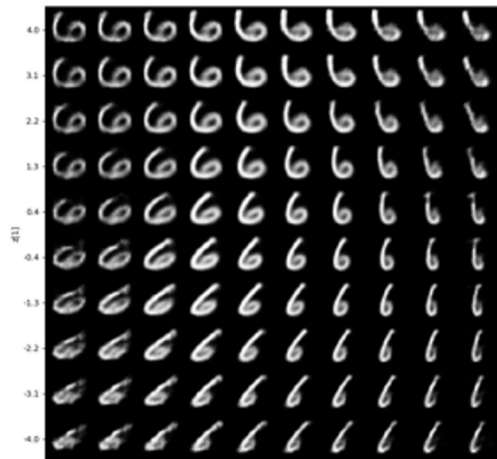


# CVAE Sample Outputs

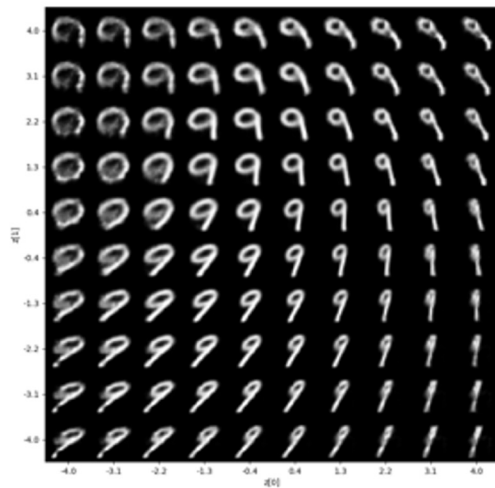
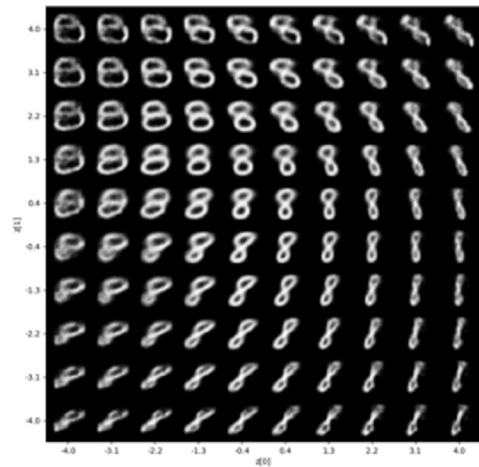


# CVAE Sample Outputs





# CVAE Sample Outputs



# Beta-VAE - Disentangled Representation VAE

- Intuitively, we define a vector representation as disentangled, if it can be decomposed into a number of subspaces, each one of which is compatible with, and can be transformed independently by a unique symmetry transformation. [1]
- A disentangled code or representation is a tensor that can change a specific feature or attribute of the output data while not affecting the other attributes. [2]
- Capture the independent features of a given data in such a way that if one feature changes, the others remain unaffected [3]

1. Higgins, Irina, et al. "Towards a Definition of Disentangled Representations." arXiv preprint arXiv:1812.02230 (2018).
2. Atienza, R. Advanced Deep Learning with Keras. Packt Pub (2018)
3. Locatello et al. "Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations." ICML 2019

## $\beta$ -VAE - Change in Loss Function

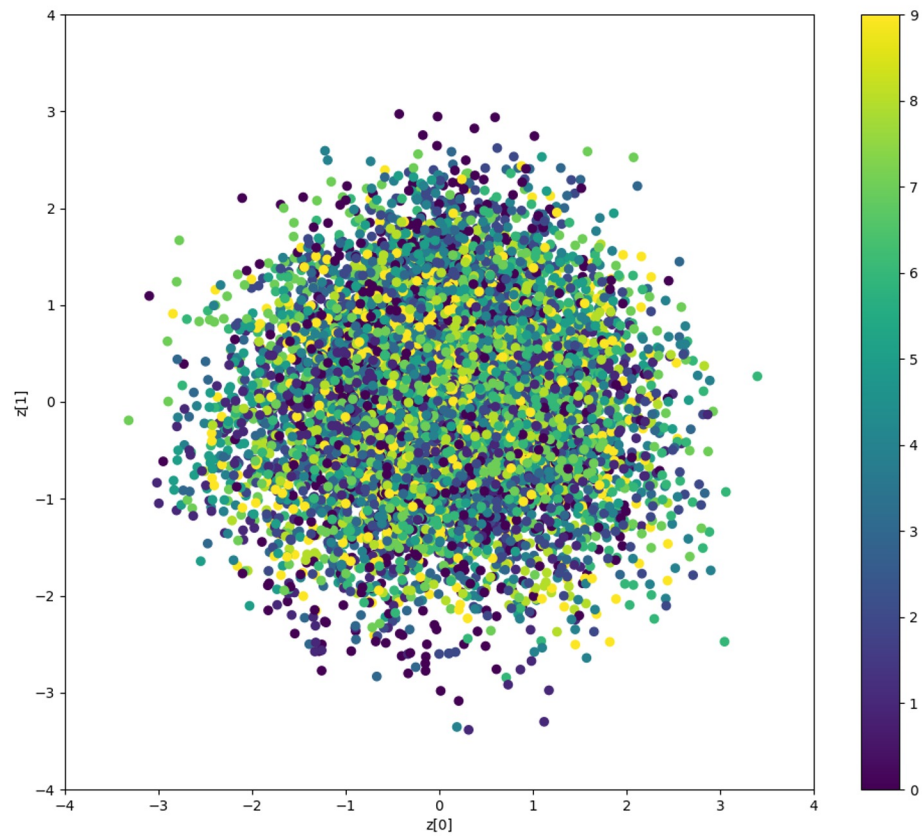
$$\mathcal{L}_{\beta\text{-VAE}} = \mathcal{L}_R + \beta \mathcal{L}_{KL}$$

$\beta > 1$  acts as a loss regularizer. The effect is to enforce tighter standard deviation resulting to disentanglement of latent codes. CVAE and VAE are special cases where  $\beta = 1$ .

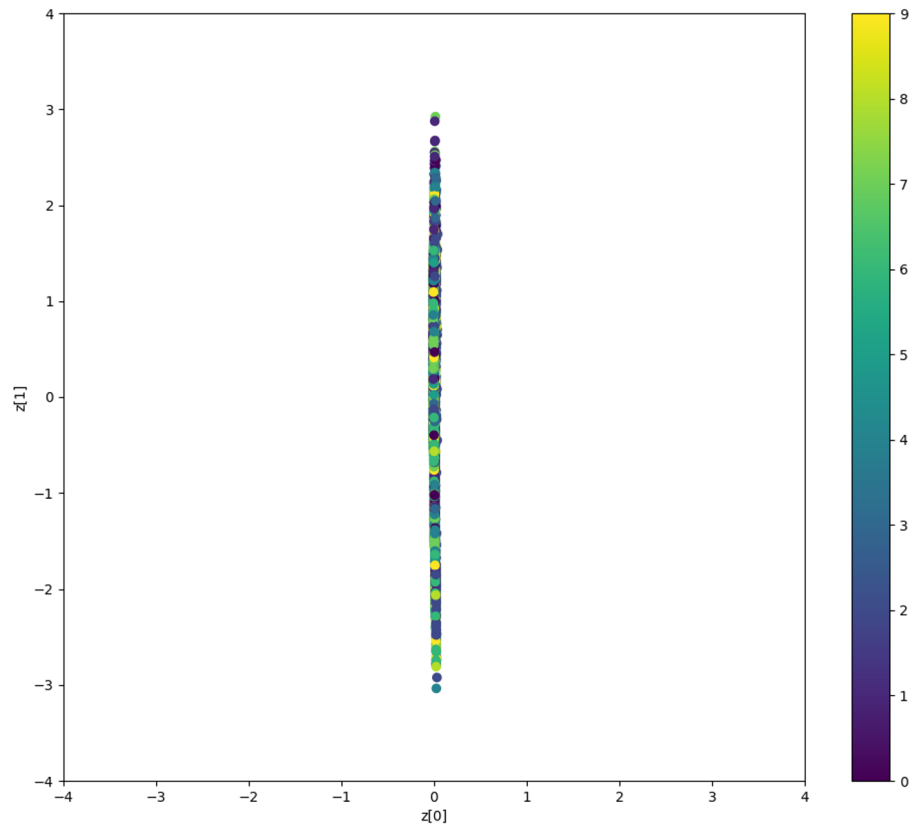


How to determine  $\beta$ ? Answer: Trial and error

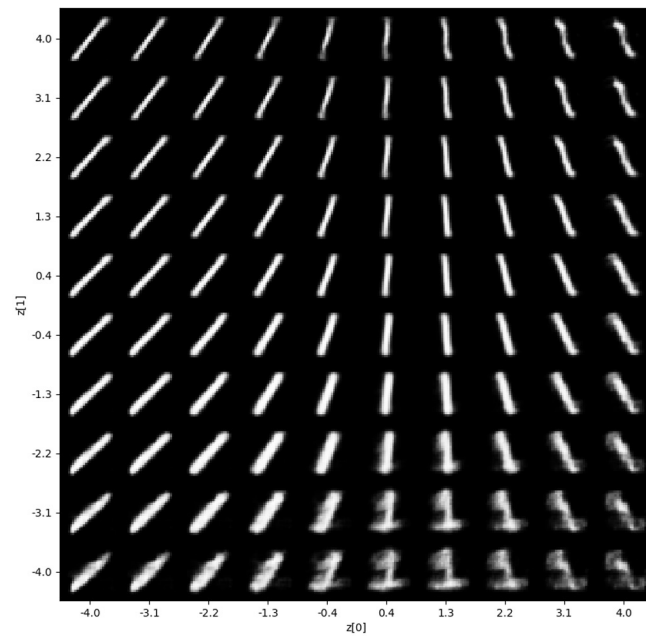
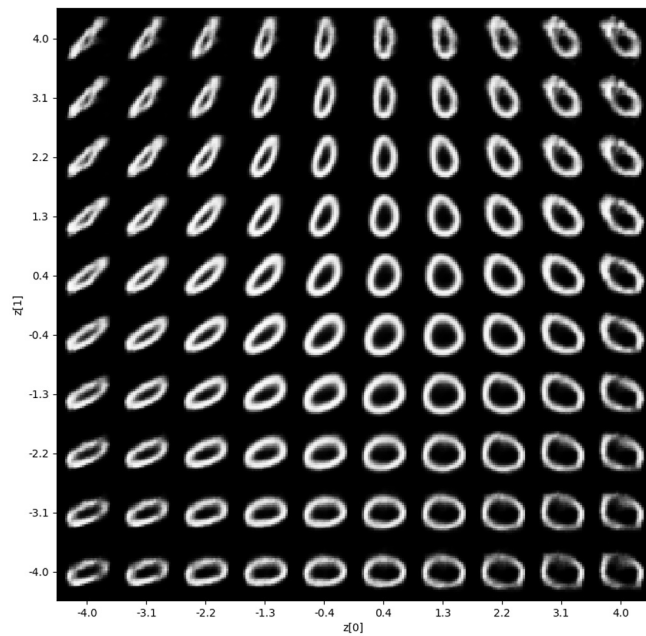
$$\beta=7$$



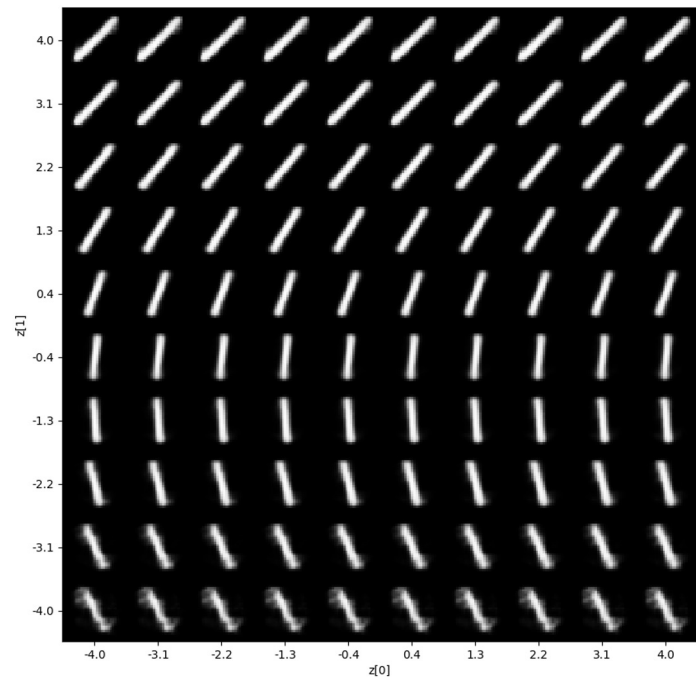
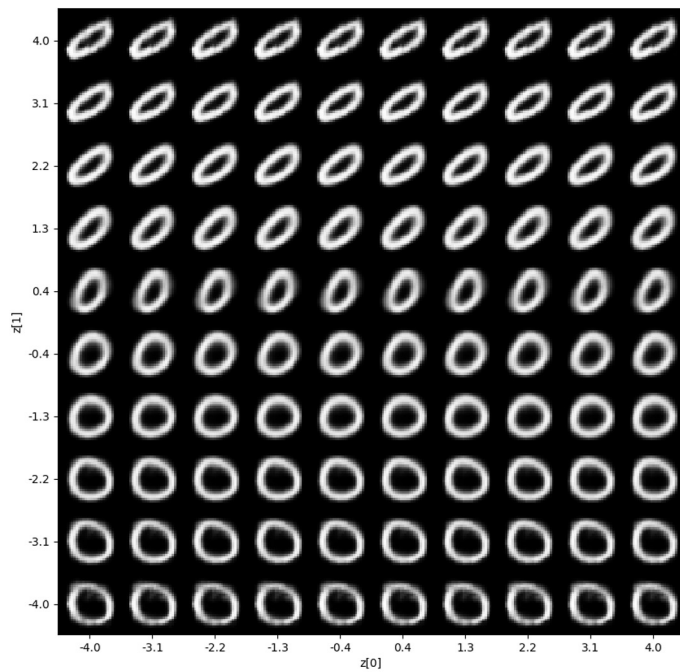
$$\beta=10$$



$$\beta=7$$



$$\beta=10$$



# References

Kingma, Diederik P., and Max Welling. "Auto-encoding Variational Bayes." arXiv preprint arXiv:1312.6114 (2013).

Sohn, Kihyuk, Honglak Lee, and Xinchun Yan. "Learning structured output representation using deep conditional generative models." Advances in Neural Information Processing Systems. 2015.

Doersch, Carl. "Tutorial on variational autoencoders." arXiv preprint arXiv:1606.05908 (2016).

Atienza, R. "Advanced Deep Learning with TF2 and Keras." Packt Pub 2020

End