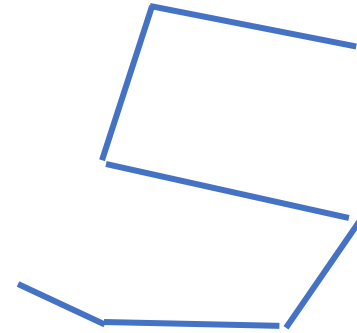
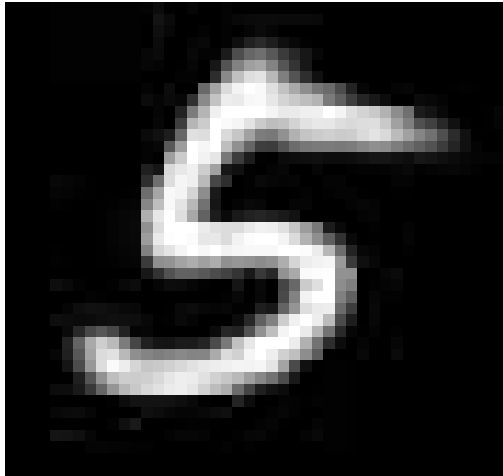


# Unsupervised image vectorisation

Attila Kun

# Vectorisation

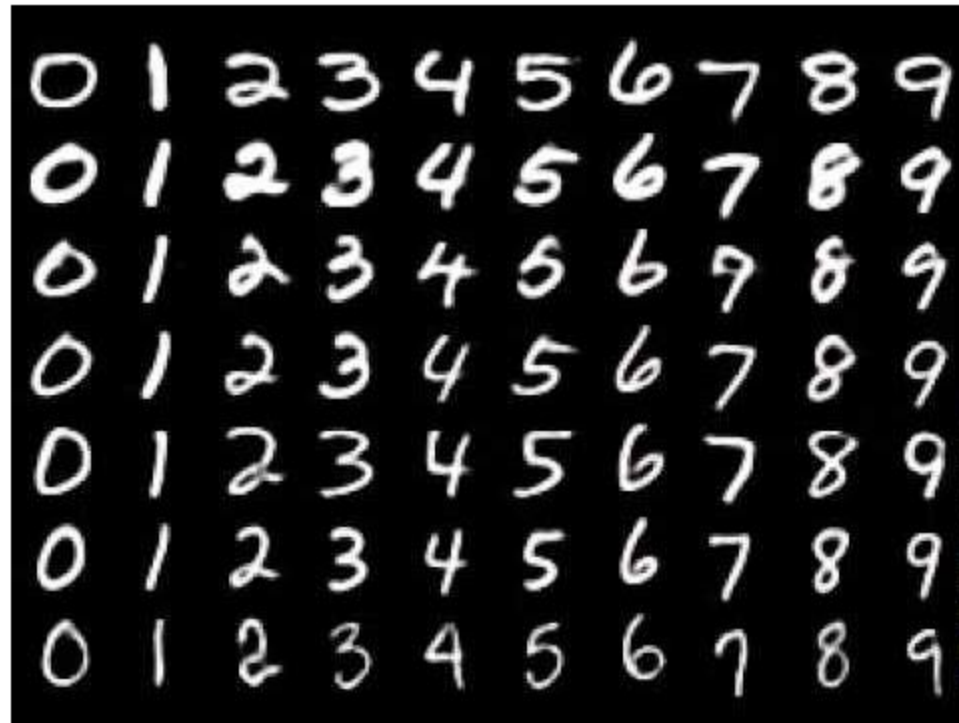


Abstract language (example: SVG)

Usefulness:

1. Scale agnostic rendering
2. Dimensionality reduction

# MNIST database



- Collected from Census Bureau employees
- 70,000 images

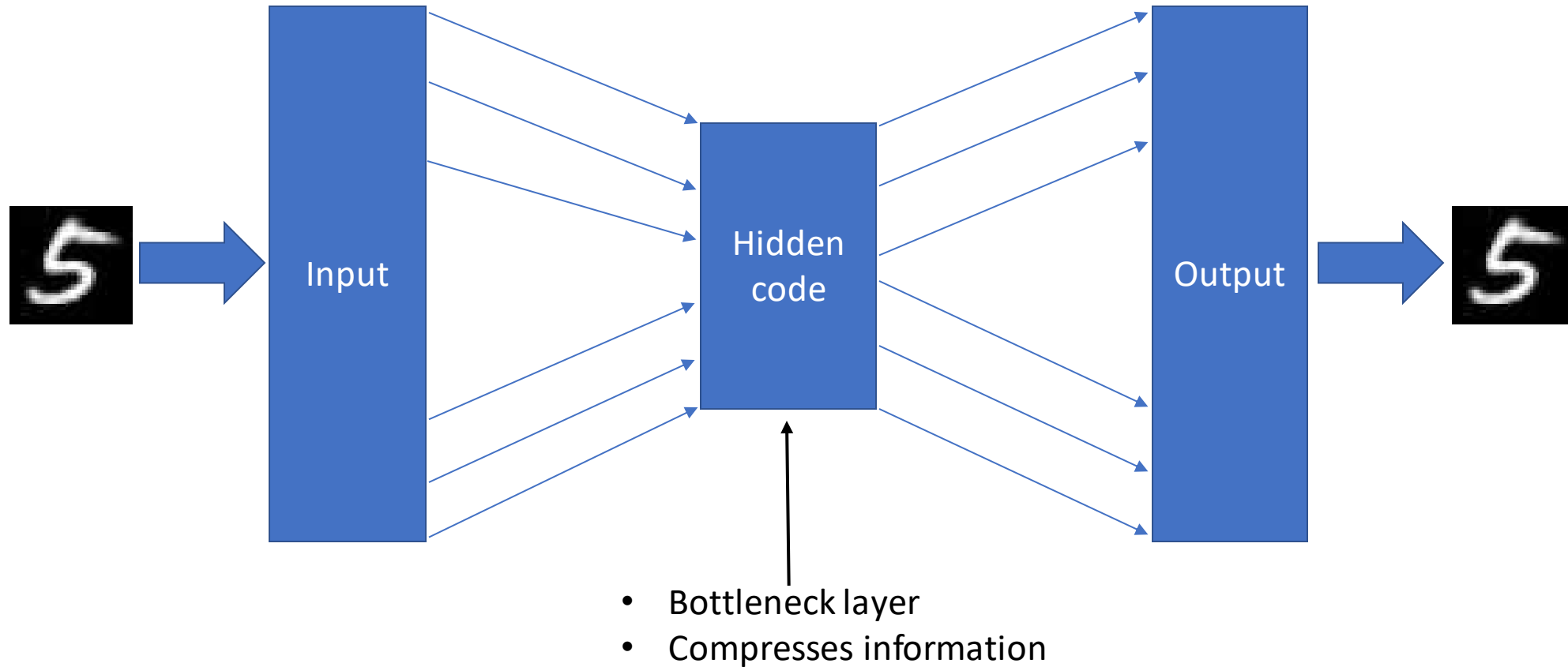
# Goals

- Vectorise MNIST digits using fix number of lines
- Use no labels
- Make use of deep learning as opposed to classical techniques
  - PyTorch (automatic differentiation, stochastic gradient descent)
  - Rust

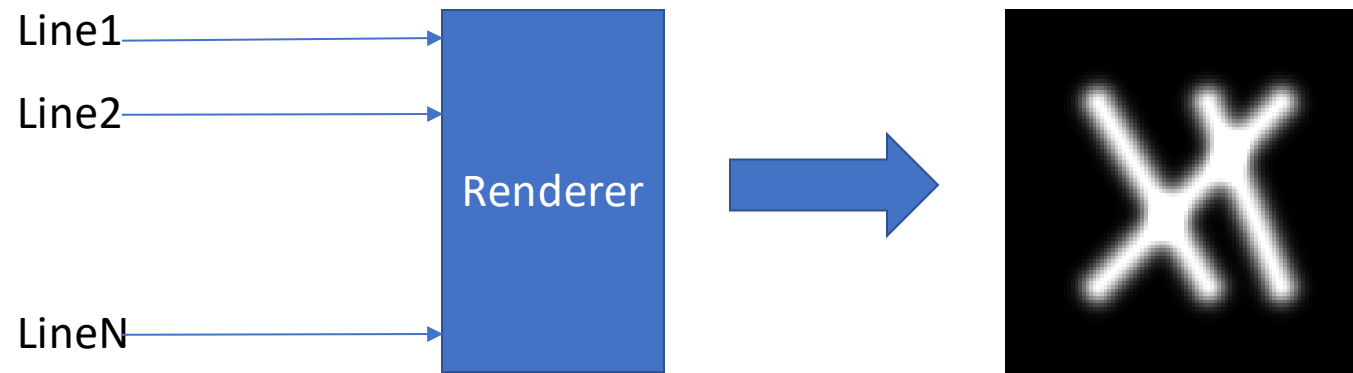
# 3 approaches

1. Traditional autoencoder
2. Variational autoencoder
3. Likelihood-free inference

# Traditional autoencoder

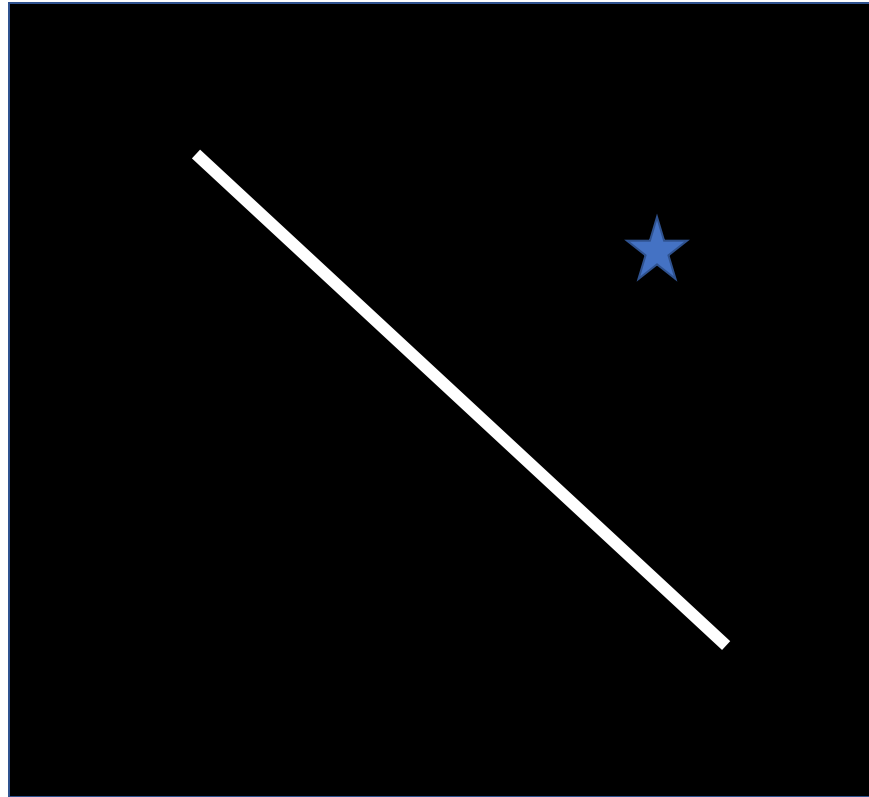


# Differentiable renderer



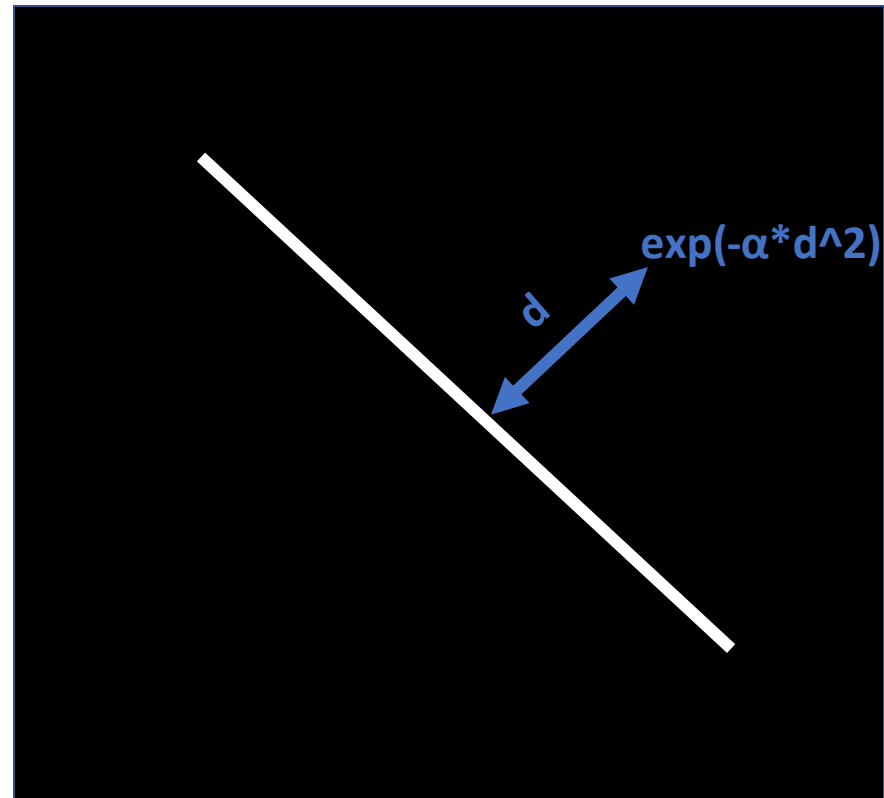
- Heat map
- Continuous function
- Soft renderer

# Differentiable renderer



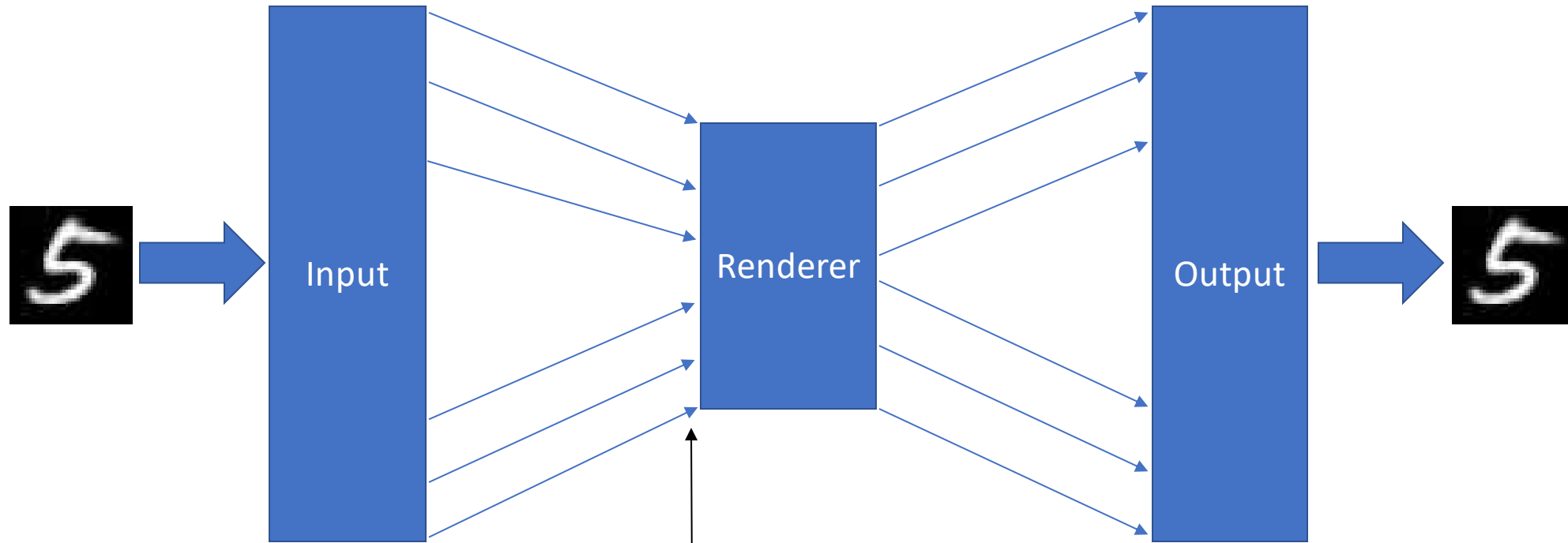


# Differentiable renderer



Repeat for every pixel on GPU

# Modified autoencoder



- Network learns how to produce useful inputs
- Well defined interface
- Inputs are the inferred vectors

# Training

- Whole architecture is differentiable
- We can optimise using stochastic gradient descent
- Source of randomness: minibatch training

# Results

Original MNIST



Vectorised output

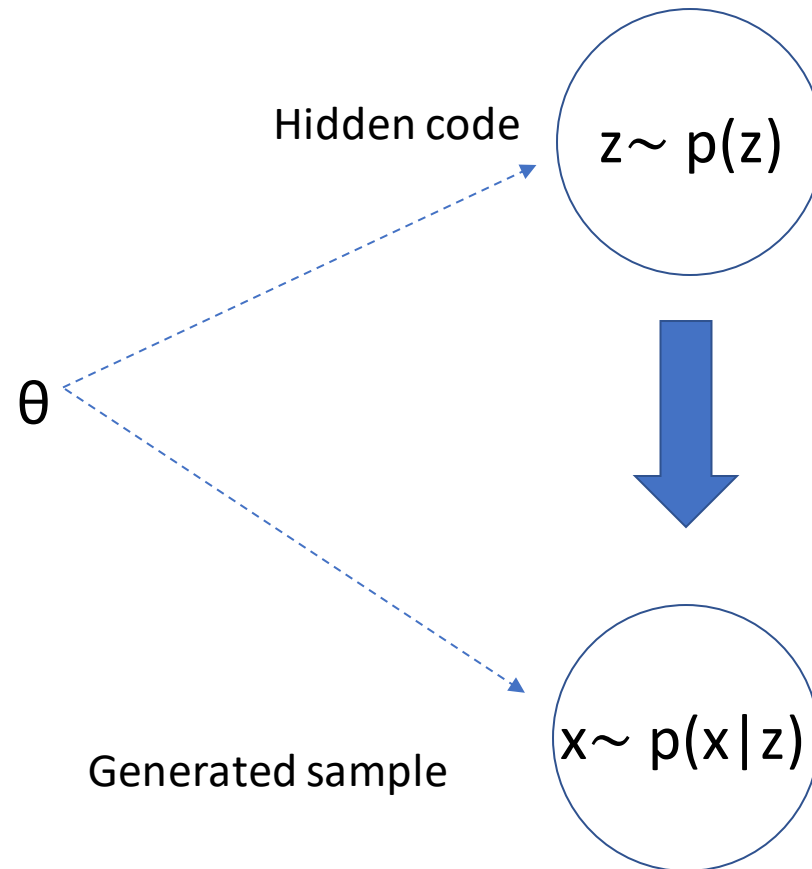
Pro: fast training

Con: weak results, cumbersome soft rendering implementation

Learns something but not great. How to improve?

# Variational autoencoder

Generative model:



$p(z|x)?$

Approximate the posterior  
using  $q(z|x)$  parameterised  
by  $\phi$

# Variational autoencoder

In the image vectorisation context:

- $p(z)$ : Vectors
- $p(x|z)$ : Generated images
- $p(z|x) \approx q(z|x)$ : Inferred vectors of an image. Finding this is our goal.
- $p(x|z)$ ,  $q(z|x)$  are neural networks
- $p(z)$  is multivariate standard normal

# Variational lower bound

$$\log p_{\theta}(x_i) \geq \underbrace{-D_{KL}(q_{\phi}(z_i|x_i) || p_{\theta}(z_i))}_{\text{regularisation}} + \underbrace{\mathbb{E}_{q_{\phi}(z_i|x_i)} [\log p_{\theta}(x_i|z_i)]}_{\text{reconstruction}}$$

If the lower bound is maximal w.r.t.  $\theta$  and  $\phi$ , then  $q(z|x)=p(z|x)$

# Implementing $p(x|z)$

- How to turn vectors into an image?
- We need a renderer



# Non-differentiable renderer

- Differentiable programming is annoying
- Non-differentiable programming is more general

Solutions:

1. Approximating gradients with a neural network
2. Native implementation (Rust) => big speedup

# Putting it together

- Approximate renderer + variational objective
- Training via stochastic gradient descent

# Results

3 7 5 6 0 6 2 6  
3 7 5 6 0 6 2 6

7 3 4 9 5 7 3 1  
7 3 4 9 5 7 3 1

0 5 4 2 8 7 6 0  
0 5 4 2 8 7 6 0

6 6 9 5 0 8 9 1  
6 6 9 5 0 8 9 1

5 9 6 0 7 5 7 9  
5 9 6 0 7 5 7 9

1 0 4 6 7 9 5 0  
1 0 4 6 7 9 5 0

# Sampling from prior

Regularised



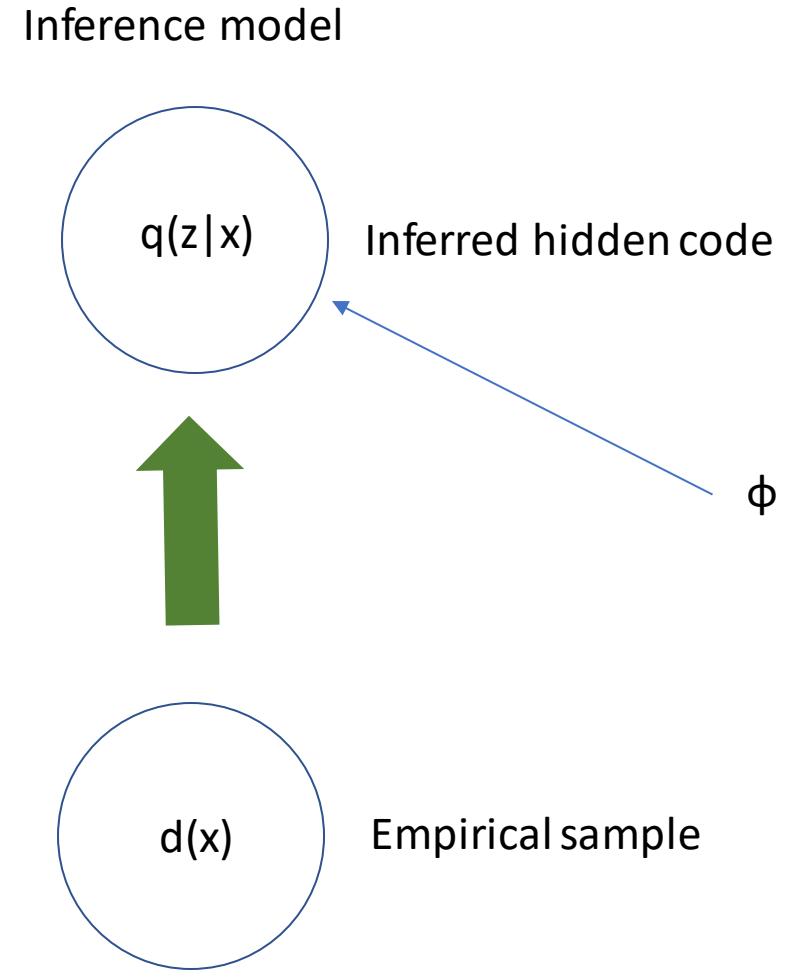
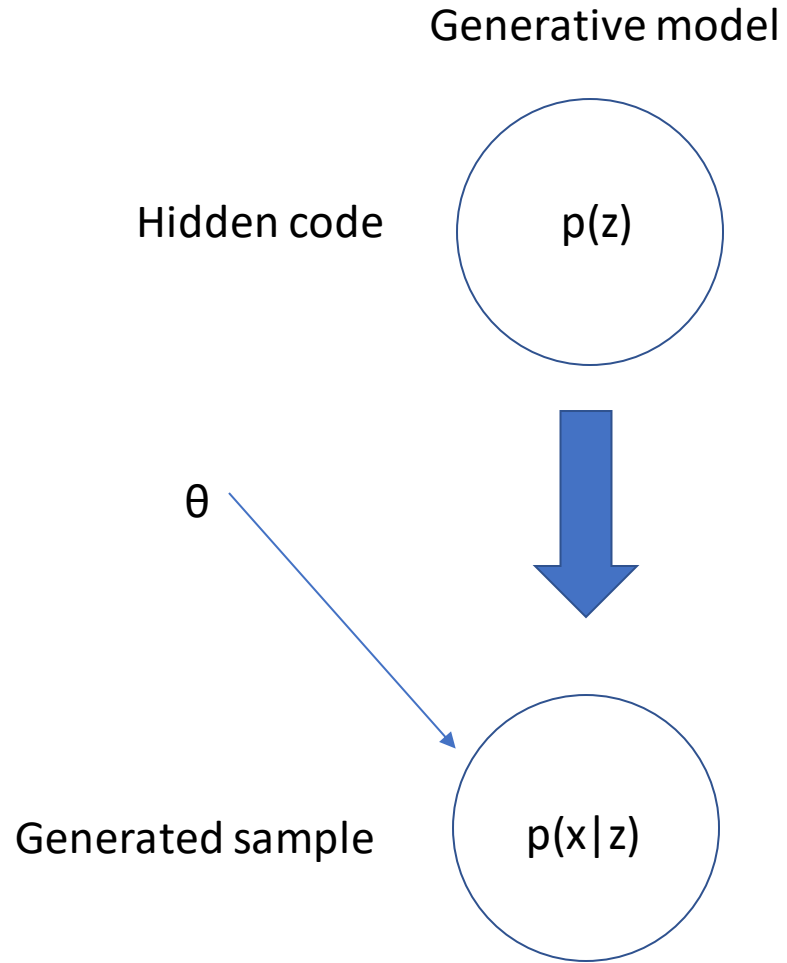
Unregularised



# Likelihood free inference

- Variational autoencoders use fully specified distributions
- Sometimes we can only sample from an implicit distribution

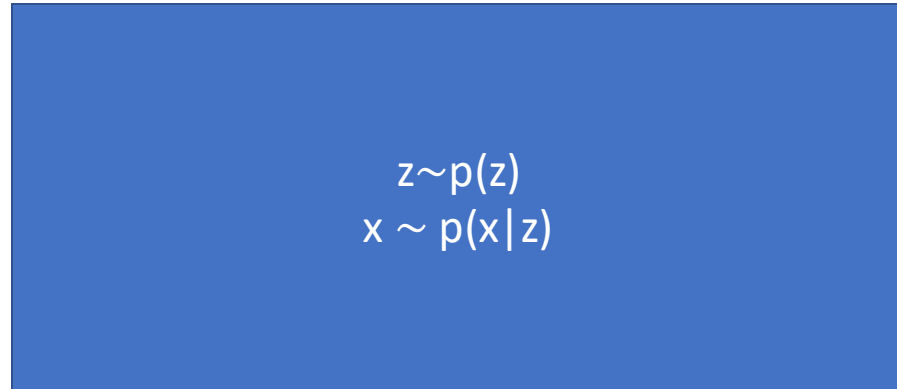
# Implicit distributions



# Discriminator

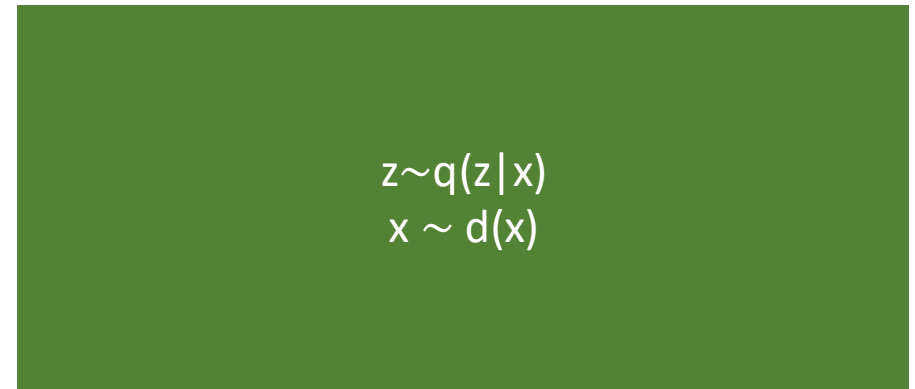
- Train a discriminator that can tell the difference between pairs

Generative model



VS

Inference model



$$\min_{\theta, \phi} \max_{\psi} \mathbb{E}_{p(x_n, z_n | \theta)} [\log D(x_n, z_n; \psi)] + \mathbb{E}_{q(x_n, z_n | \phi)} [\log(1 - D(x_n, z_n; \psi))]$$

# Problems

- Mode collapse
- Complication: approximate rendering
- Possible solution: Different objective (Wasserstein)



# Summary

1. Traditional autoencoder (differentiable renderer) 😐
2. Variational autoencoder (non-differentiable renderer) 😊
3. Likelihood-free variational inference 😞

## Conclusion:

- Classical techniques might be still better
- More complicated approaches like SPIRAL would give better results at the cost of higher complexity

Thank you. Questions?