# Deep Generative Models

Dr. Svetlin Penkov

# Generative Modelling



$$x_i \sim P_{data}, i = 1, \ldots, N$$

$$d(P_{data}, P_\theta)$$

$$P_\theta$$

$$P_{data}$$

$$\theta \in M$$

Find **Θ\*** such that:

1. <u>Generation</u>: If $x_{new} \sim P_{\Theta^*}(x)$ then $x_{new}$ should look like a car
2. <u>Density estimation</u>: $P_{\Theta^*}(x)$ should be high if x is a car, and low otherwise
3. <u>Representation learning</u>: Learn common attributes amongst $x_i$'s

# Deep Generative Models

- Encode distributions using deep neural networks...

# Representing Distributions

- Continuous RVs e.g.

$$\mathcal{N}(\mu_\theta(x), \Sigma_\theta(x))$$

  where $\mu_\theta$ and $\Sigma_\theta$ are neural networks parameterised by $\theta$.

- Similarly, one can encode other continuous distributions (Gamma, Beta, etc…)

# Representing Distributions

- Discrete RVs e.g.

$$Cat(\theta_1, \ldots, \theta_k)$$

$$\theta_i = \frac{e^{z_i}}{\sum_{j=0}^{K} e^{z_j}} = \frac{e^{f_i(x)}}{\sum_{j=0}^{K} e^{f_j(x)}}$$

$$\theta = \mathrm{softmax}(\mathbf{f}_\phi(x))$$

where $\mathbf{f}_\phi$ is a neural network parameterised by $\phi$.

- Similarly, one can encode other discrete distributions (Bern, Mulinom, etc…)
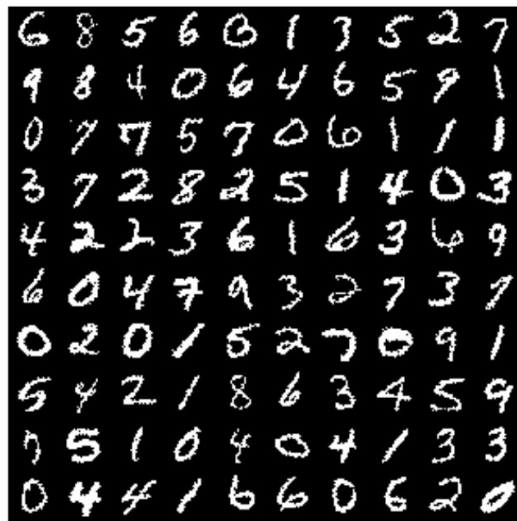
# Train on MNIST

Data

Generated Samples



- 28x28 binary images
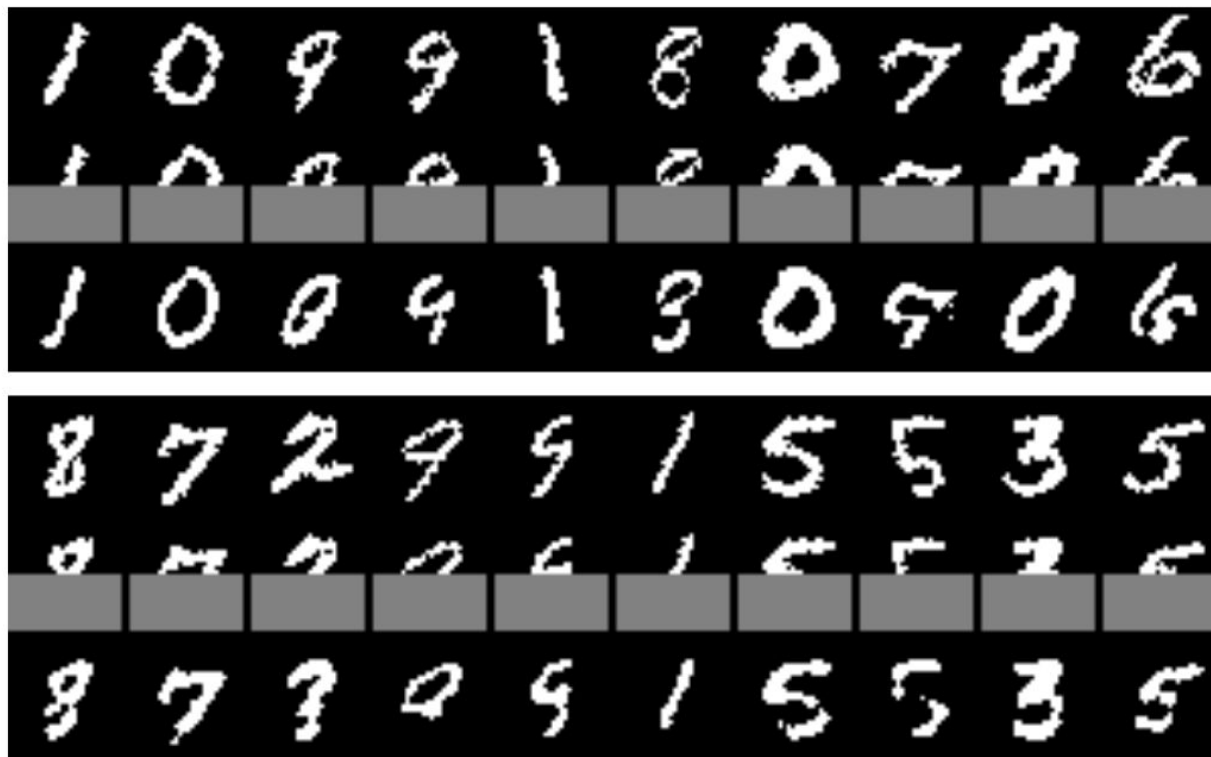
$$x \in \{0,1\}^{784}$$

- Factorise distribution

$$p(x|\alpha) = p(x_1|\alpha_1)p(x_2|x_1, \alpha_2) \ldots p(x_{784}|x_1, \ldots x_{783}, \alpha_{784})$$

where

$$p(x_i|x_{j<i}, \alpha_i) = \sigma\left(\alpha_{i,0} + \sum_{j=1}^{i} \alpha_{i,j}x_j\right)$$

G. Zhe et al., *Artificial Intelligence and Statistics*, 2015.

# Missing Data Prediction



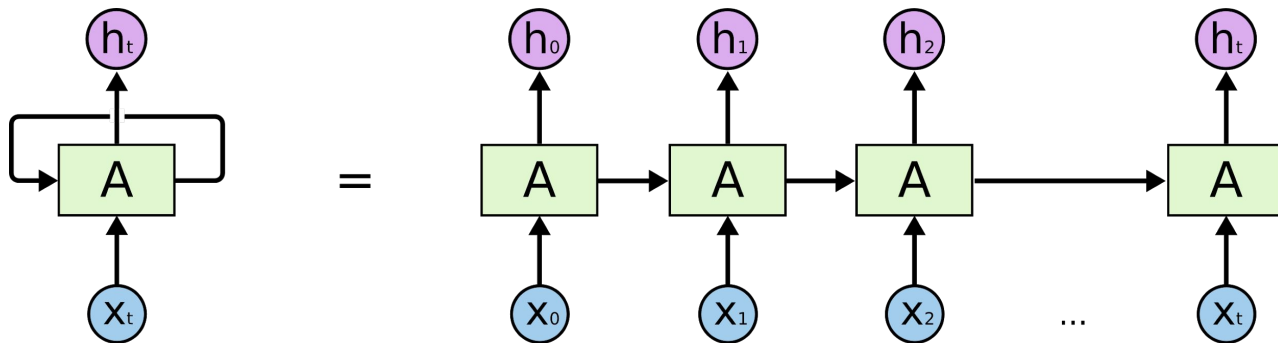G. Zhe et al., *Artificial Intelligence and Statistics*, 2015.

# PixelRNN

- Use an RNN to learn

$$p(x|\alpha) = p(x_1|\alpha_1)p(x_2|x_1, \alpha_2) \ldots p(x_{784}|x_1, \ldots x_{783}, \alpha_{784})$$
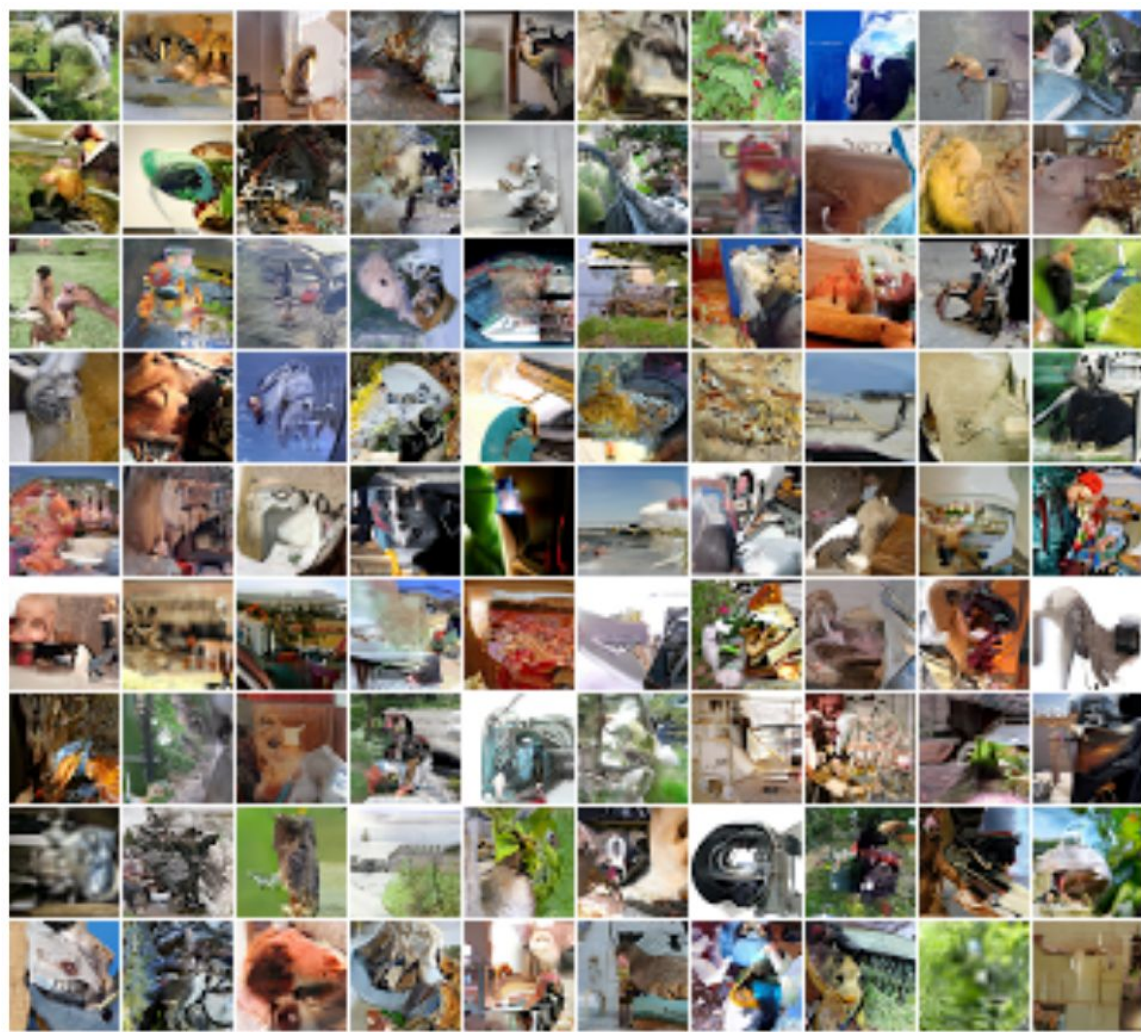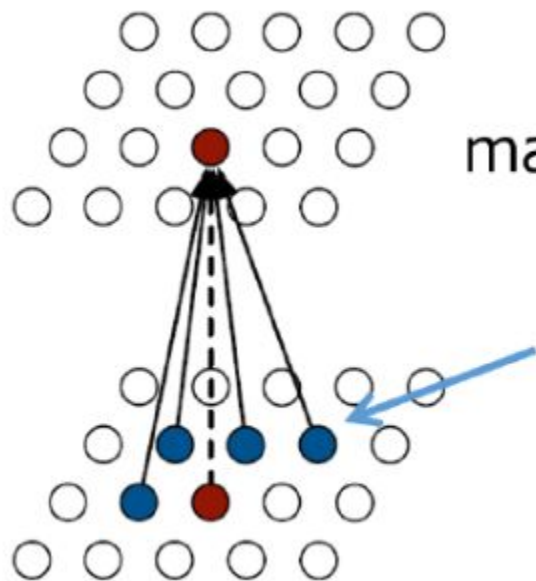
- Train RNN, such that $\quad h_t = p(x_t|x_{1:t-1})$

# PixelRNN

- Deal with RGB images

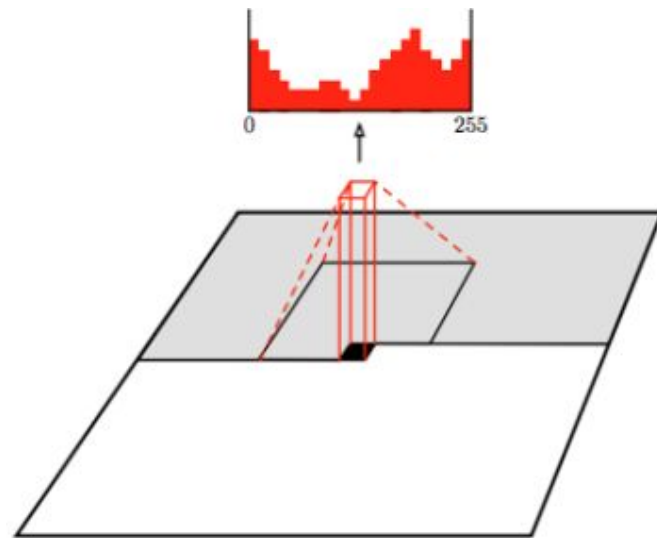$$p(x_t|x_{1:t-1}) = p(x_t^R|x_{1:t-1})p(x_t^G|x_{1:t-1})p(x_t^B|x_{1:t-1})$$
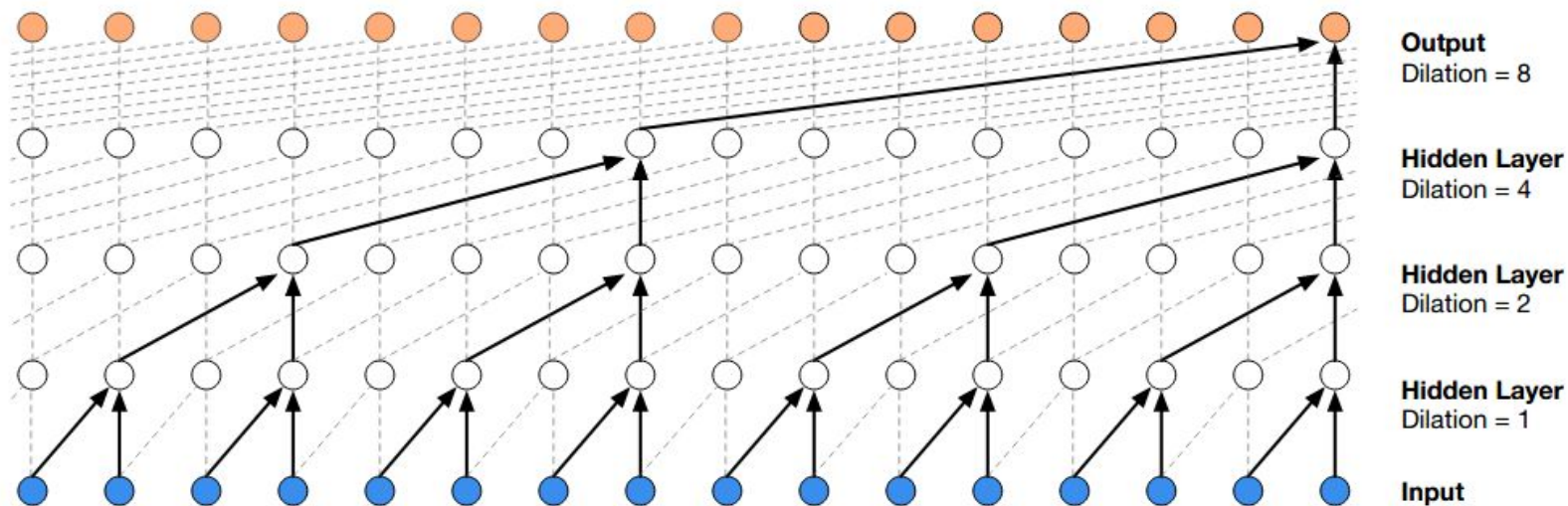
Oord et al., ICML, 2016

# PixelCNN

masked convolution

# WaveNet

# Variational Autoencoder

- Find common features amongst $x_i$'s

- Model the features as latent random variables

- Expected data likelihood per point $x_i$

$$l(x_i, \theta) = -\mathbb{E}_{p(z|x_i)}[\log p_\theta(x_i|z)]$$

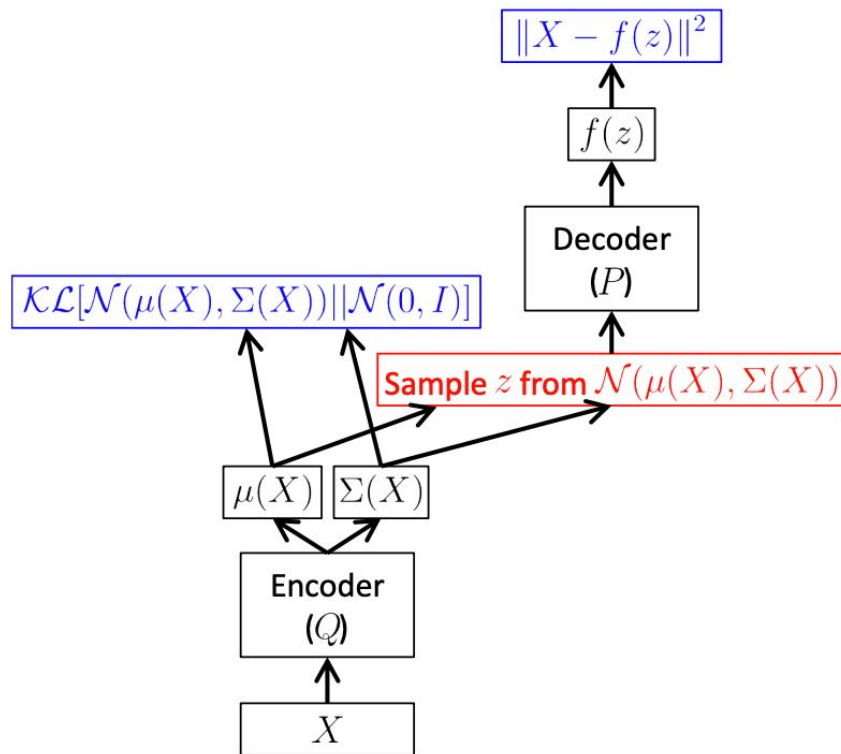# Variational Autoencoder (VAE)

- Find common features amongst $x_i$'s

- Model the features as latent random variables

- Expected data likelihood per point $x_i$

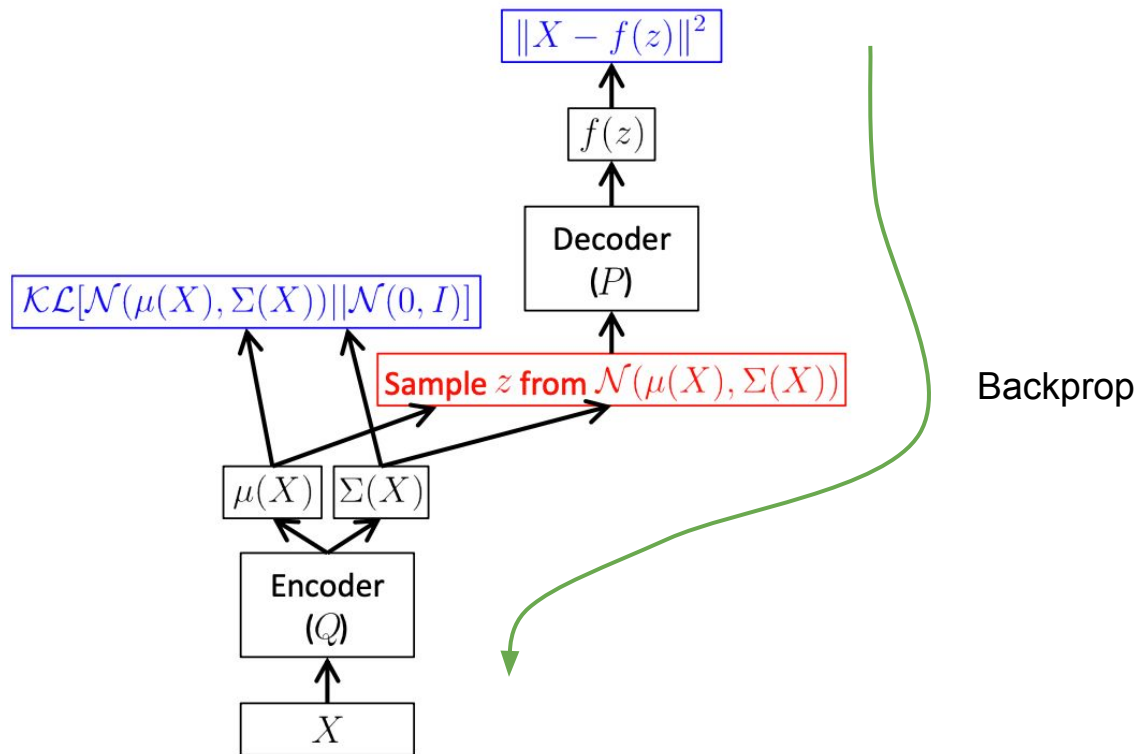$$l(x_i, \theta) = -\mathbb{E}_{p(z|x_i)}[\log p_\theta(x_i|z)]$$

- Posterior is intractable so approximate

$$l(x_i, \theta, \phi) = -\underbrace{\mathbb{E}_{q_\phi(z|x_i)}[\log p_\theta(x_i|z)]}_{\substack{\text{Expected} \\ \text{Reconstruction} \\ \text{Error}}} + \underbrace{\mathbb{KL}(q_\phi(z|x_i)||p(z))}_{\text{Regulariser}}$$
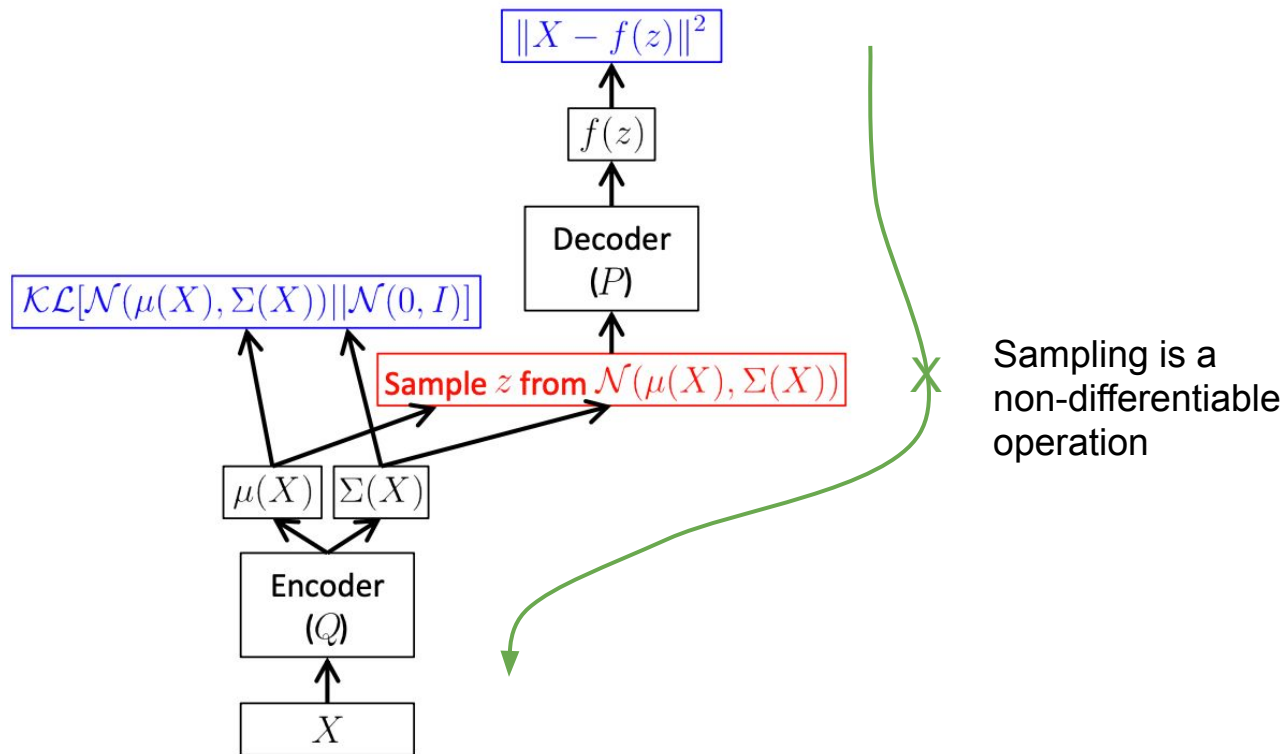
Kingma et al, arXiv:1312.6114, 2014

# VAE Computational Graph

# VAE Computational Graph
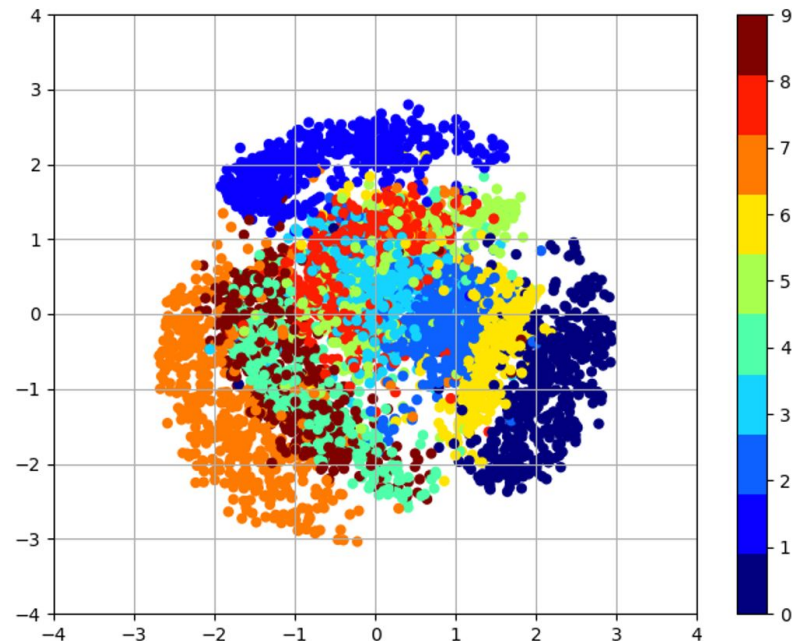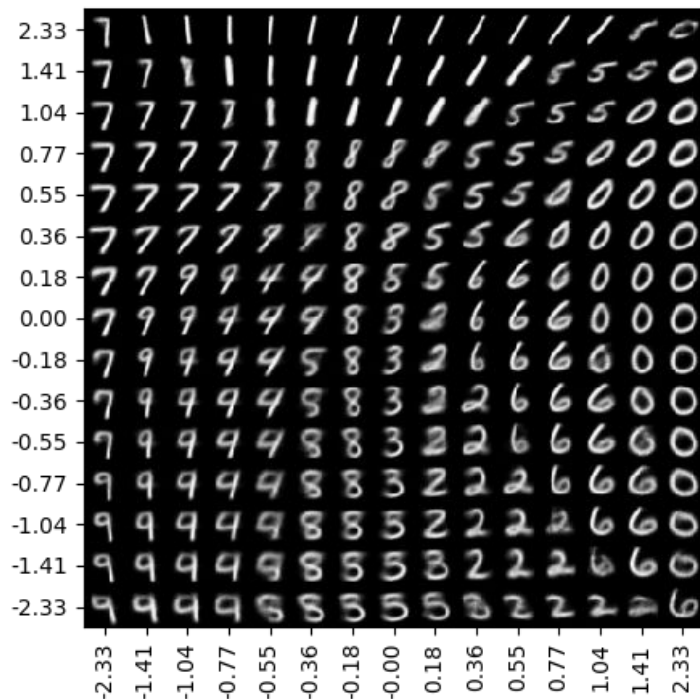
# VAE Computational Graph



Sampling is a non-differentiable operation

# VAE Reparametrisation Trick

# VAE



Input   ←- - - - - - - - - - - - - - - Ideally they are identical. - - - - - - - - - - - -→ Reconstructed input

$$\mathbf{x} \approx \mathbf{x}'$$

**Probabilistic Encoder**

$$q_\phi(\mathbf{z}|\mathbf{x})$$

Mean $\boldsymbol{\mu}$

Std. dev $\boldsymbol{\sigma}$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I})$$

**Sampled latent vector**

$\mathbf{z}$

An compressed low dimensional representation of the input.

**Probabilistic Decoder**

$$p_\theta(\mathbf{x}|\mathbf{z})$$

$\mathbf{x}$

$\mathbf{x}'$

# Latent Space Manifolds

# Representation Learning



**Probabilistic Encoder**

$q_\phi(\mathbf{z}|\mathbf{x})$

Mean $\boldsymbol{\mu}$

Std. dev $\boldsymbol{\sigma}$

**Sampled latent vector**

$\mathbf{z}$

- Maximise independence
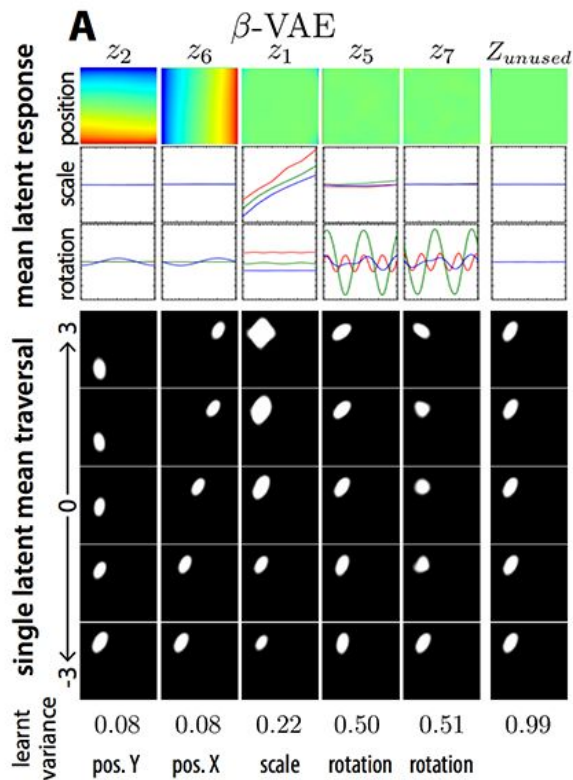- Minimise mutual information
- Use structured proposal $q$
- .
- .

# $\beta$-VAE

$$l(x_i, \theta, \phi) = -\mathbb{E}_{q_\phi(z|x_i)}[\log p_\theta(x_i|z)] + \beta\mathbb{KL}(q_\phi(z|x_i)||p(z))$$
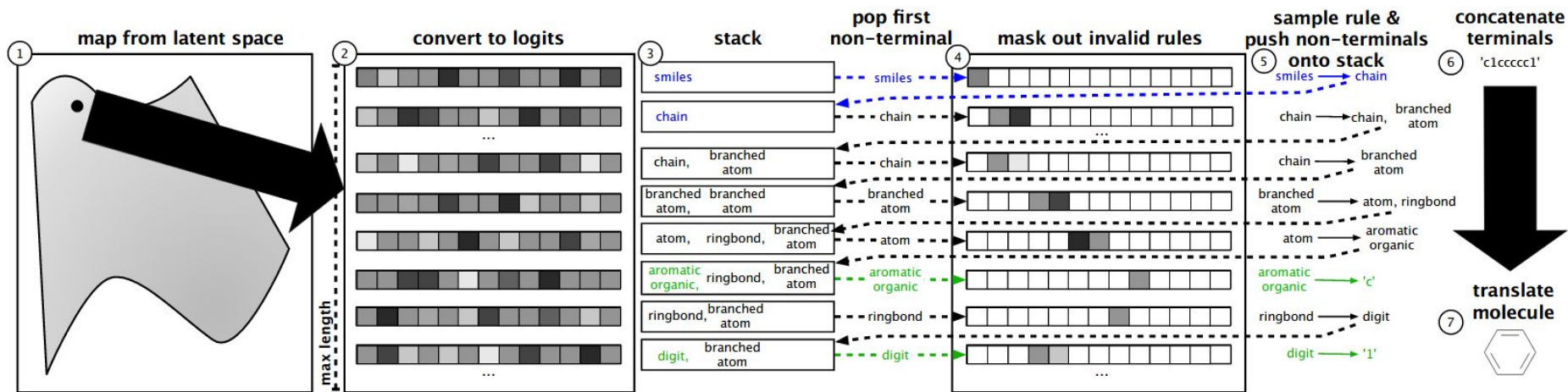


Higgins et al., *ICLR,* 2017

# Structured data

- Lots of redundancy in images
- What if we want to generate well formed strings
  - e.g. maths expressions, programs, DNA sequences, chemical molecule descriptions
- Even a single error could make the sample useless
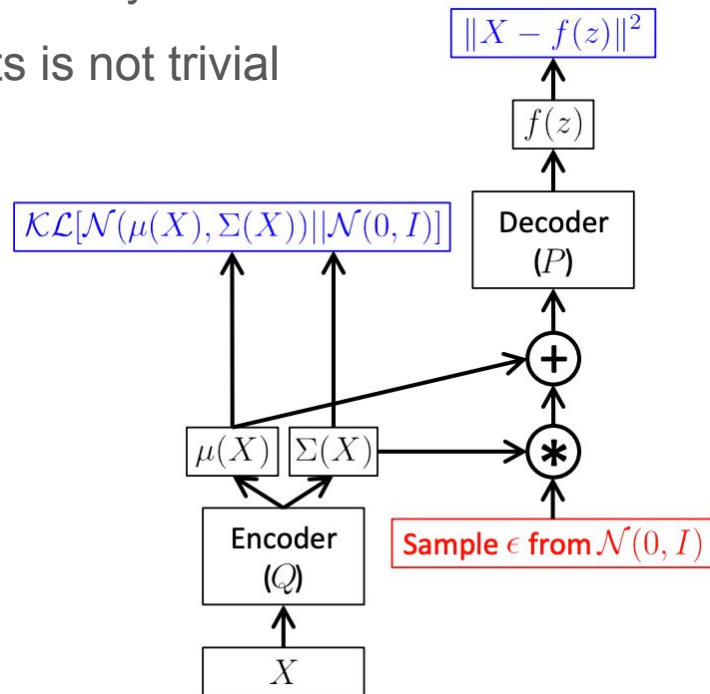
# Grammer VAE (GVAE) - Encoder
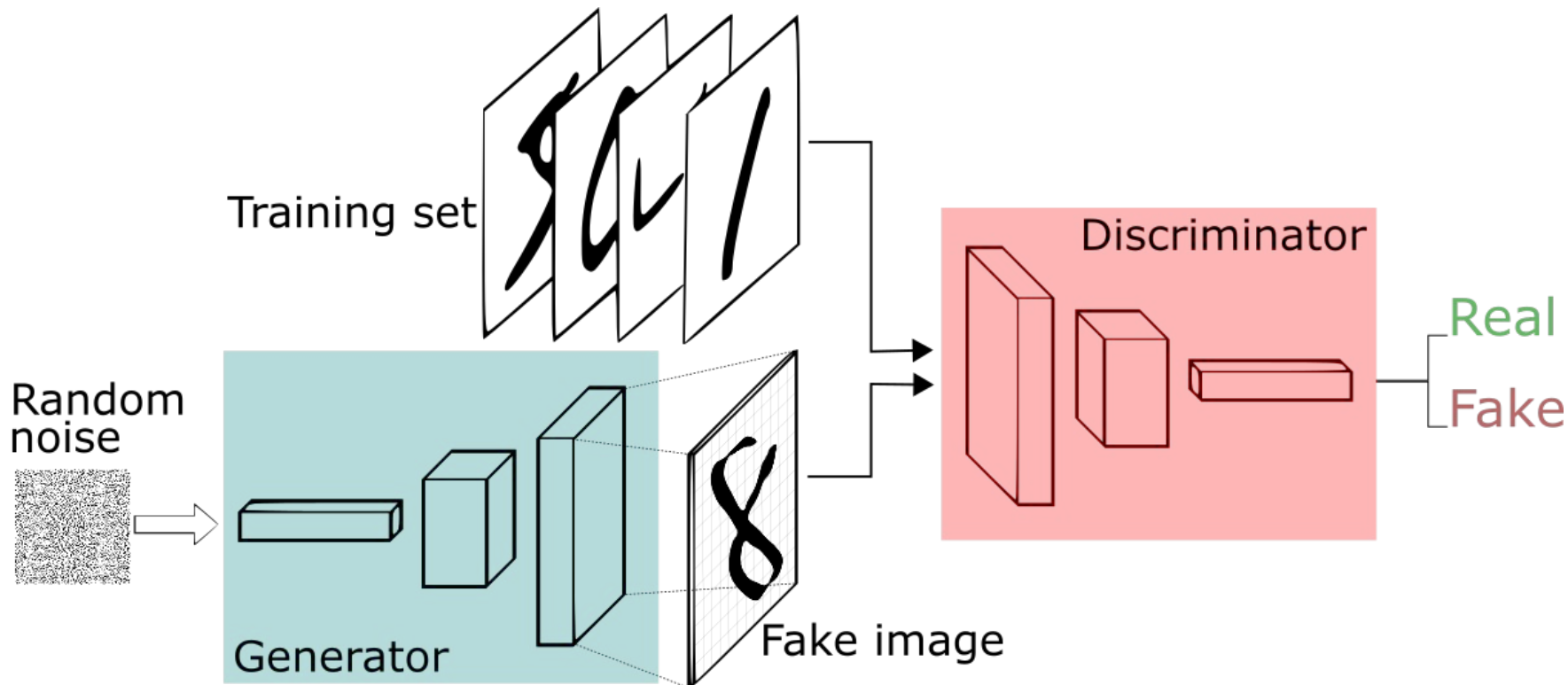
# Grammer VAE (GVAE) - Decoder

# Distance Metric

- Sampling more complex images tends to give blurry results

- Defining a distance metric for complex objects is not trivial

- Can we avoid / learn the distance metric?

# Generative Adversarial Networks (GANs)

Goodfellow, NeurIPS, 2014

# Generative Adversarial Networks (GANs)

- 2-player game objective function (i.e. How well are fake samples detected?)

$$\min_{\theta} \max_{\phi} V(G_{\theta}, D_{\phi}) = E_{\mathbf{x} \sim p_{\text{data}}}[\log D_{\phi}(\mathbf{x})] + E_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))]$$

Does D output 1 when data is real?

Does D output 0 when data is from generator?

- G wants to deceive D (decrease objective)

- D wants to detect generated samples (maximise objective)

- Tricky to train (mode collapse)

# 4.5 Years of Progress on Faces



2014

2015

2016

2017

2018

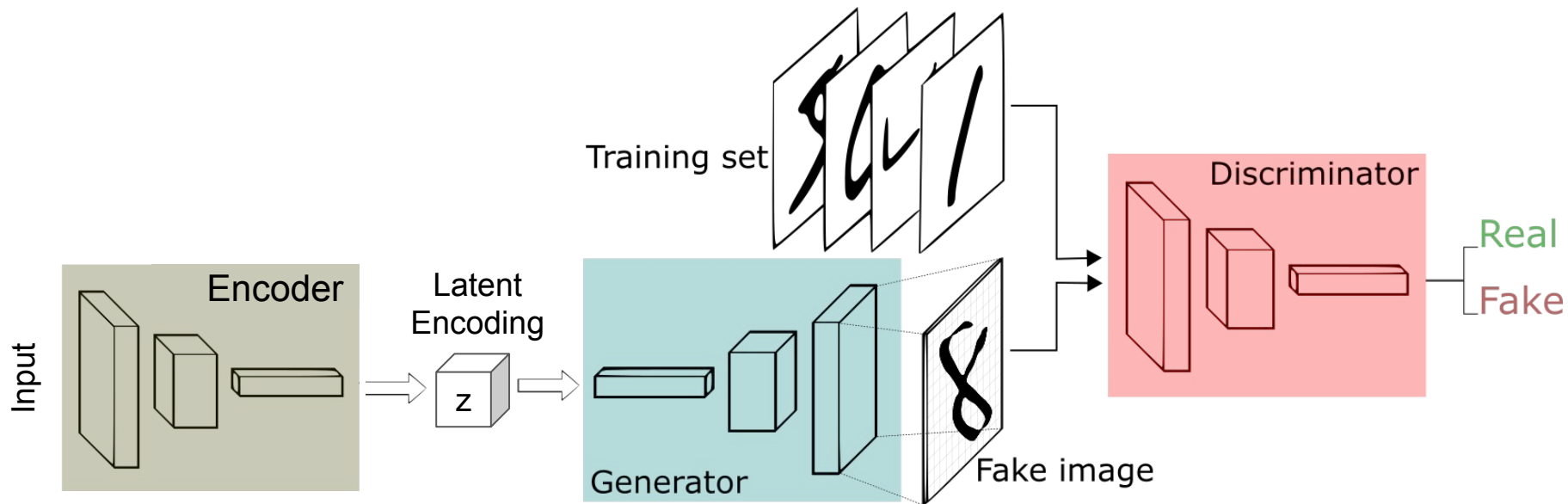# 2 Years of Progress on ImageNet



Odena et al 2016

Miyato et al 2017

Zhang et al 2018

Brock et al 2018

# Conditional GANs

# CycleGAN



Monet ⟳ Photos

Monet → photo

photo → Monet

Zebras ⟳ Horses

zebra → horse

horse → zebra

Summer ⟳ Winter

summer → winter

winter → summer

Photograph → Monet    Van Gogh    Cezanne    Ukiyo-e

Zhu et al, *ICCV*, 2017

# CycleGAN

# Everybody Dance Now



Chan et al., ICCV, 2019