

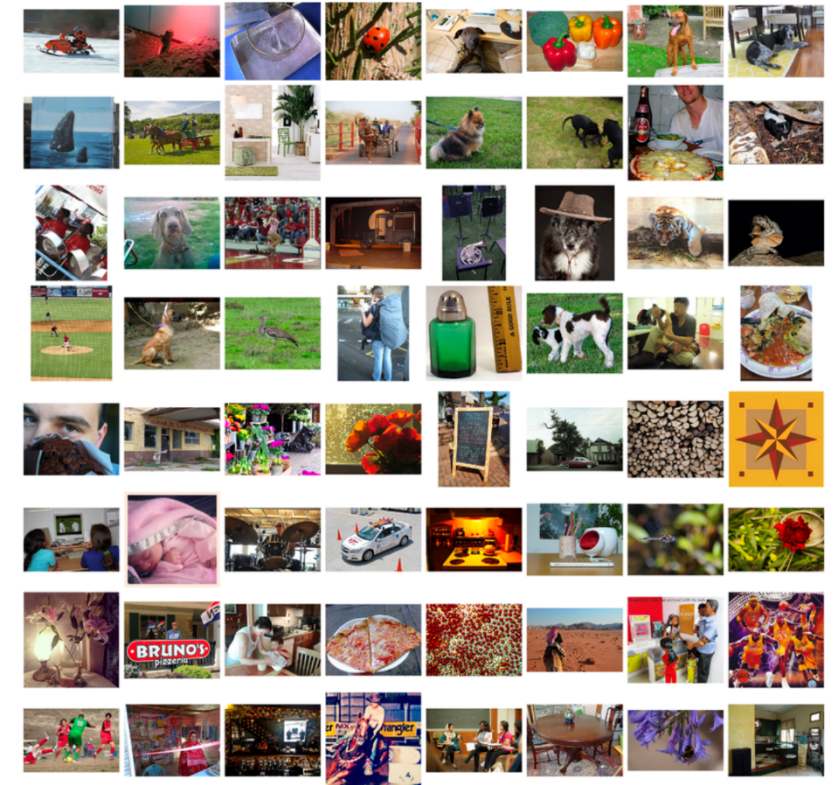
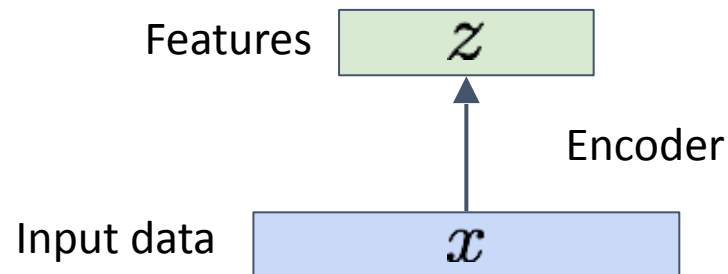
# Unsupervised Learning

- There is no direct ground truth for the quantity of interest
- Autoencoders
- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)

# Autoencoders

Goal: Meaningful features that capture the main factors of variation in the dataset

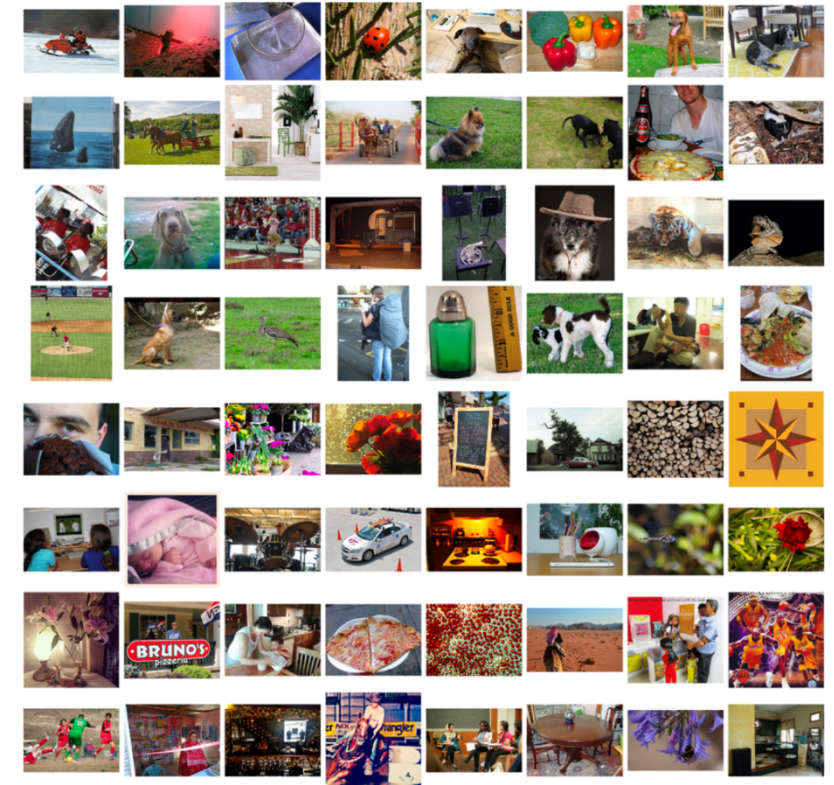
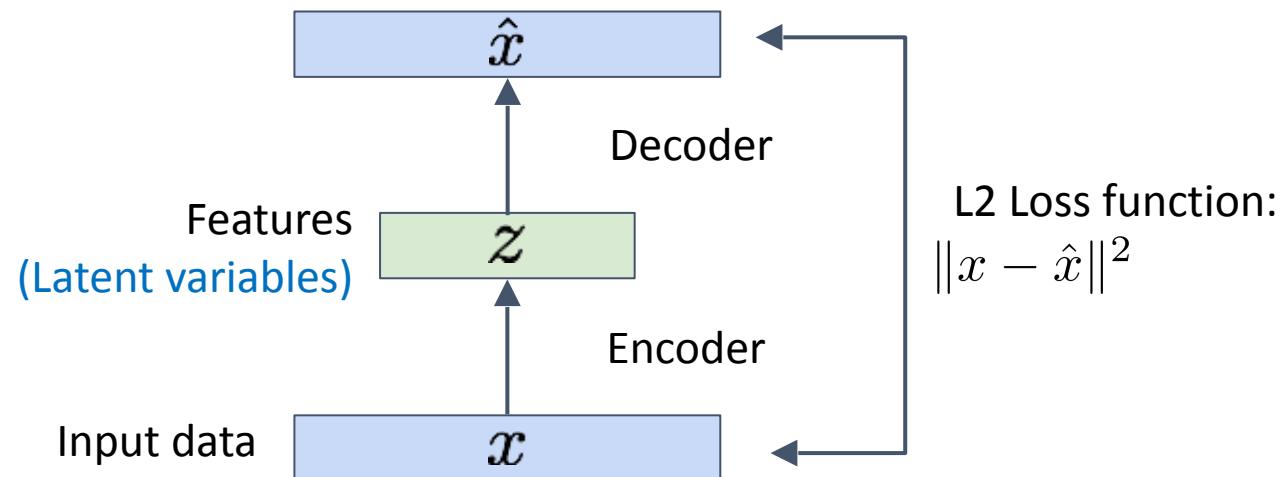
- These are good for classification, clustering, exploration, generation, ...
- We have no ground truth for them



# Autoencoders

Goal: Meaningful features that capture the main factors of variation

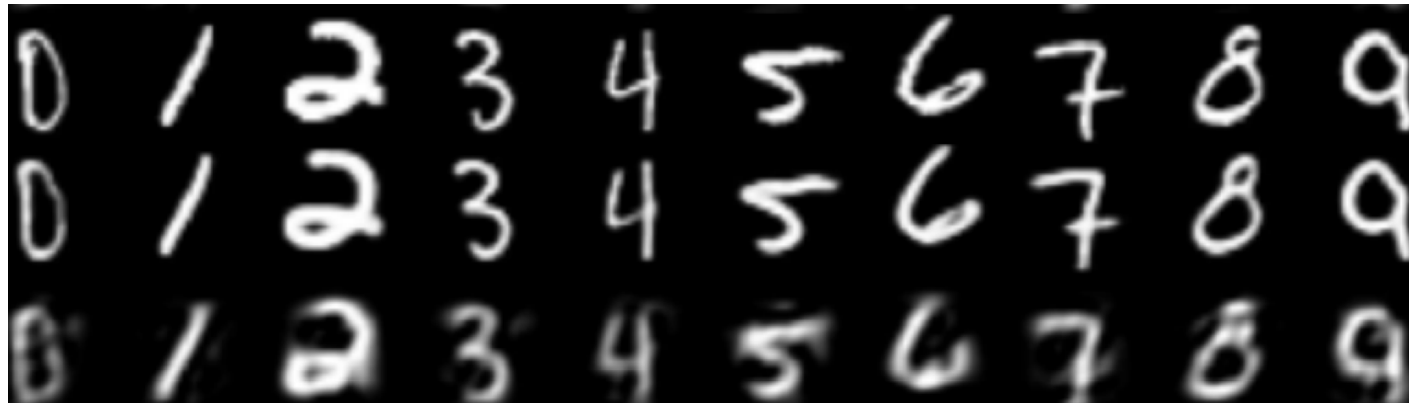
Features that can be used to reconstruct the image



# Autoencoders

Linear Transformation for Encoder and Decoder give result close to PCA

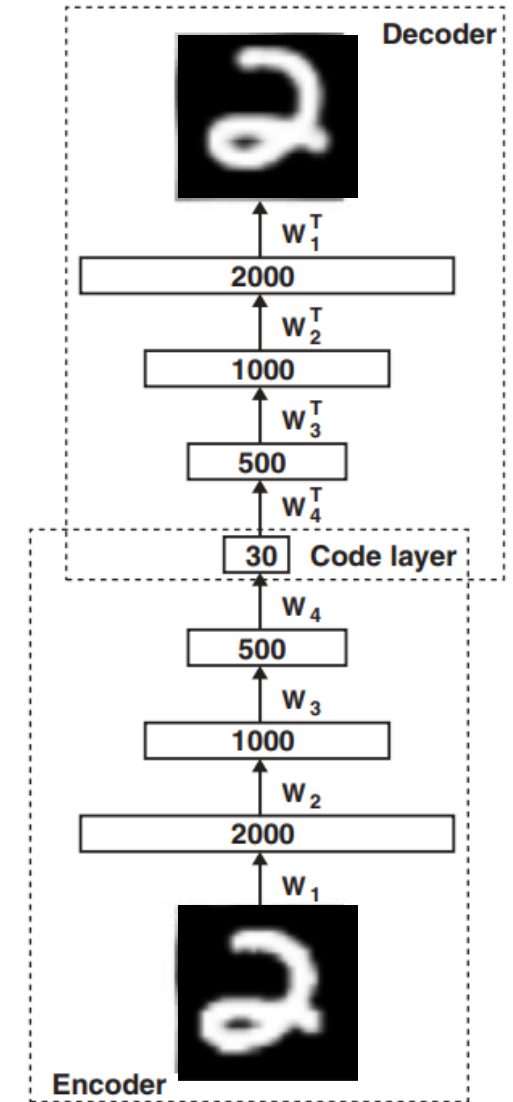
Deeper networks give better reconstructions, since basis can be non-linear



Original

Autoencoder

PCA



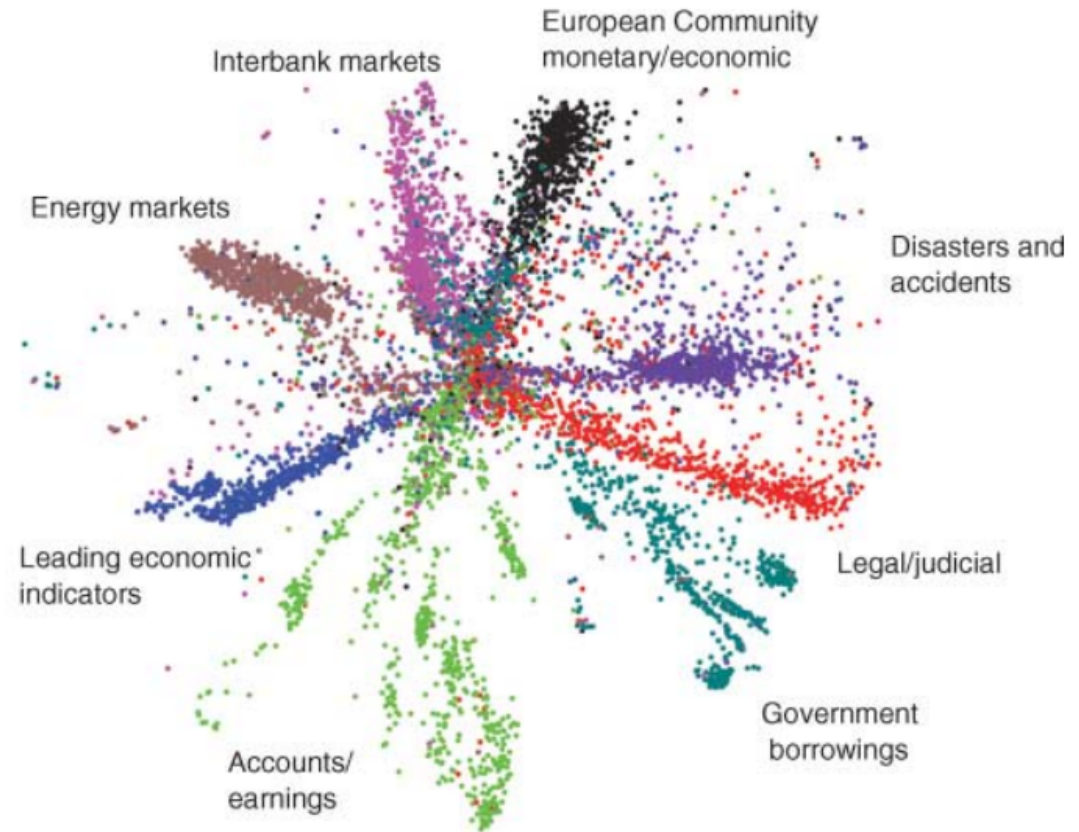


# Example: Document Word Prob. → 2D Code

PCA-based



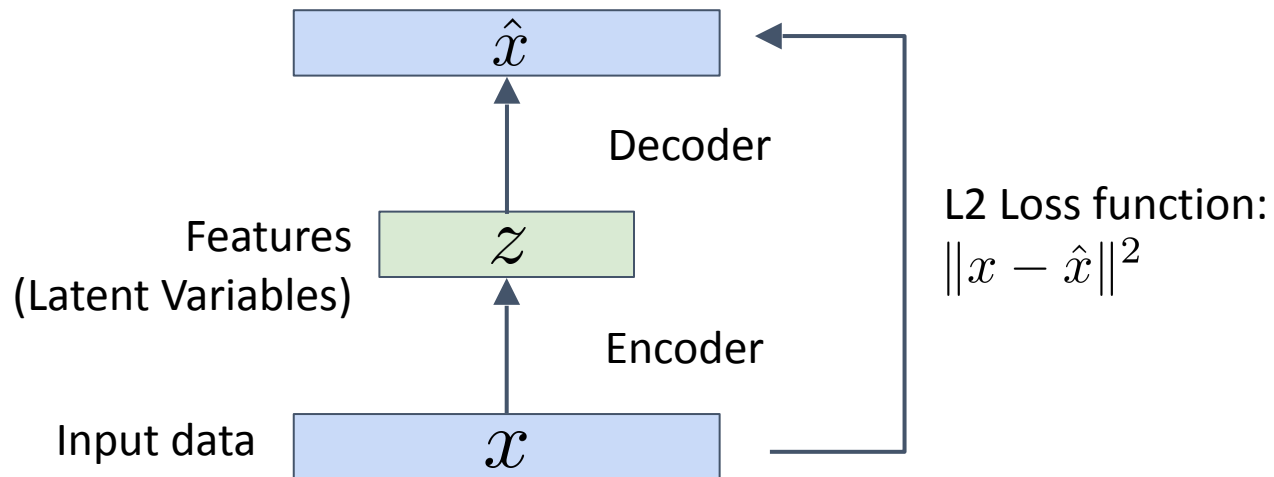
Autoencoder



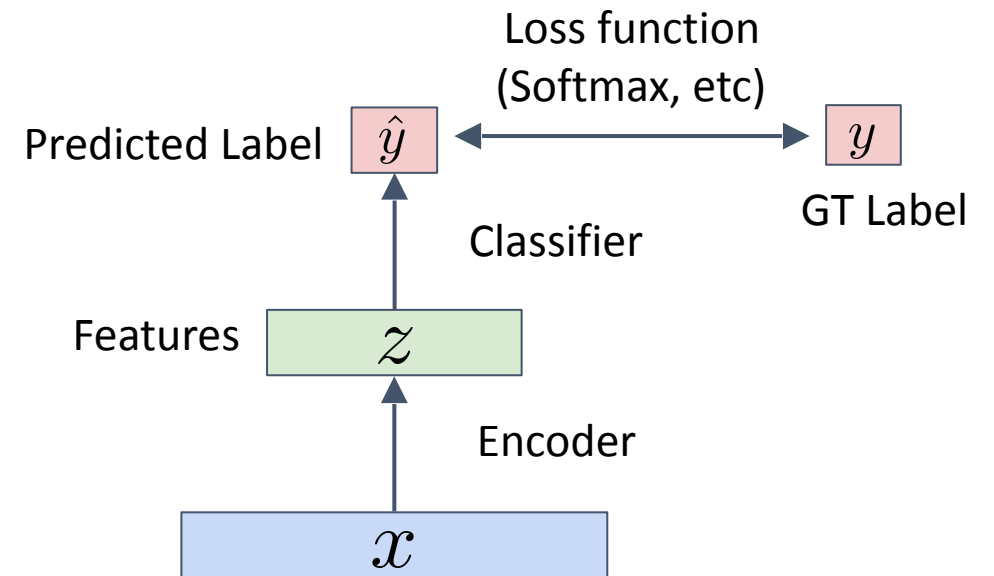
# Example: Semi-Supervised Classification

- Many images, but few ground truth labels

start unsupervised  
train autoencoder on many images



supervised fine-tuning  
train classification network on labeled images



# Autoencoder

[geometry.cs.ucl.ac.uk/creativeai](http://geometry.cs.ucl.ac.uk/creativeai)

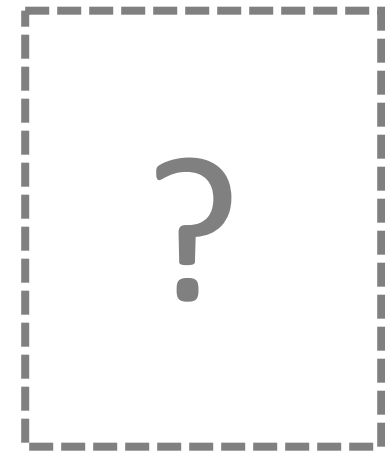


# Generative Models

- Assumption: the dataset are samples from an unknown distribution
- Goal: create a new sample from  $p_{\text{data}}(x)$  that is not in the dataset



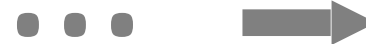
Dataset



Generated

# Generative Models

- Assumption: the dataset are samples from an unknown distribution
- Goal: create a new sample from  $p_{\text{data}}(x)$  that is not in the dataset

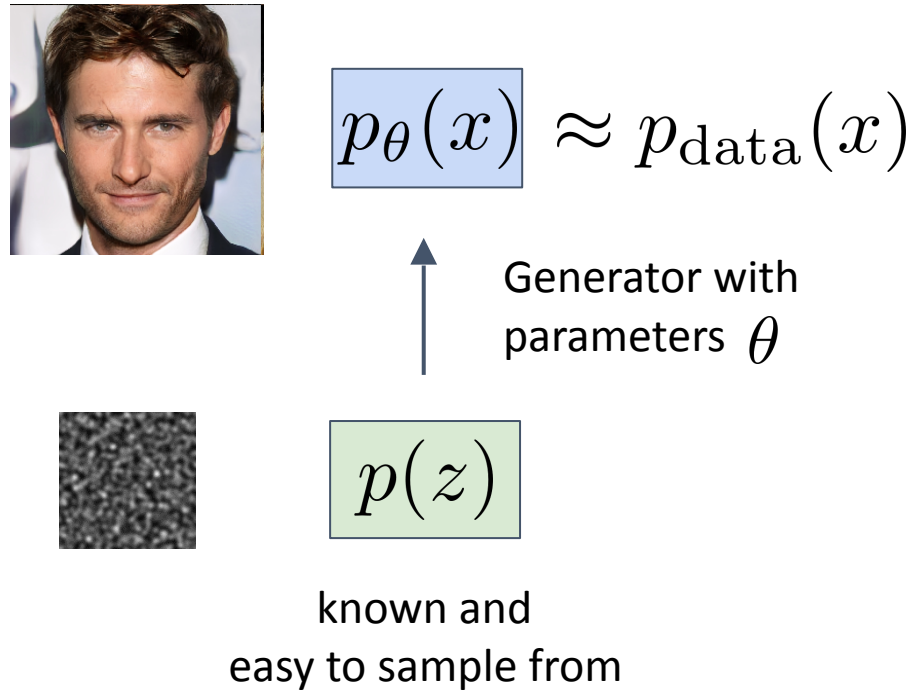


Dataset

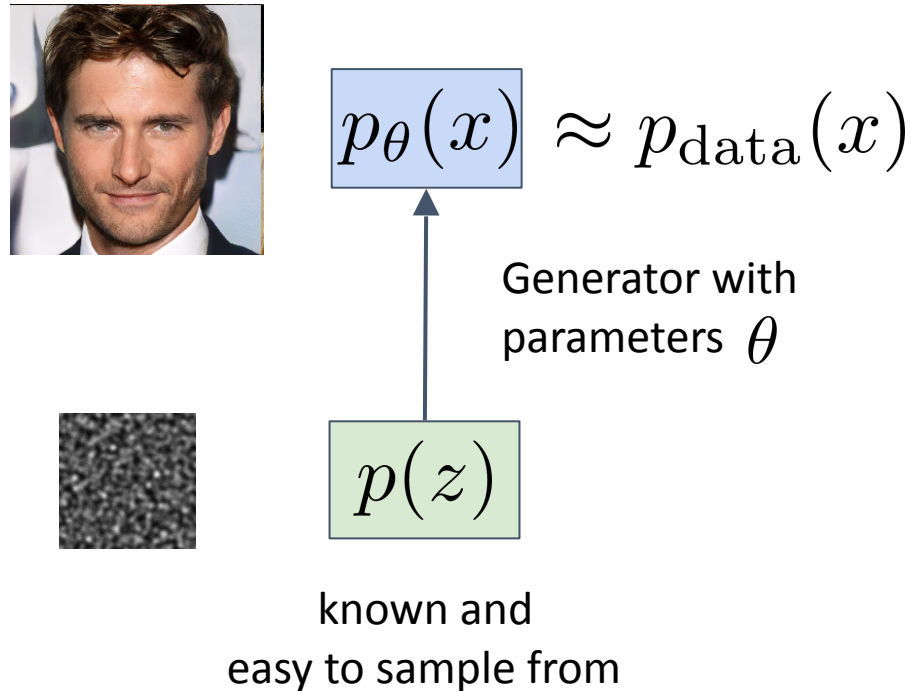
Generated



# Generative Models



# Generative Models



How to measure similarity of  $p_{\theta}(x)$  and  $p_{\text{data}}(x)$ ?

1) Likelihood of data in  $p_{\theta}(x)$

## Variational Autoencoders (VAEs)

2) Adversarial game:

*Discriminator* distinguishes  $p_{\theta}(x)$  and  $p_{\text{data}}(x)$  vs *Generator* makes it hard to distinguish

## Generative Adversarial Networks (GANs)

# Autoencoders as Generative Models?

- A trained decoder transforms some features to approximate samples from

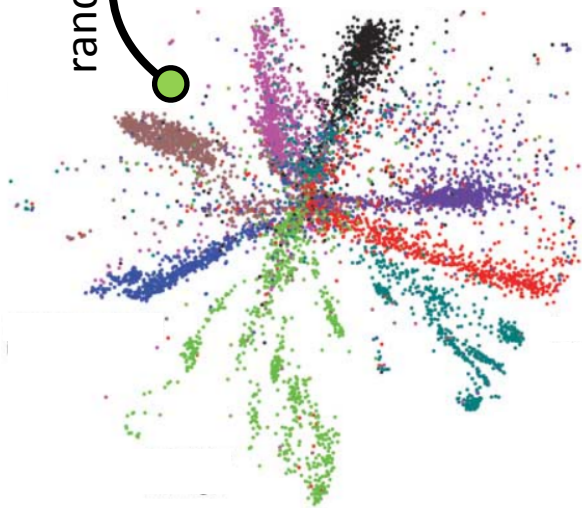
$\hat{x}$

- What happens if we pick a random ?

$p_{\text{data}}(x)$

- We do not know the distribution of features that decode to likely samples

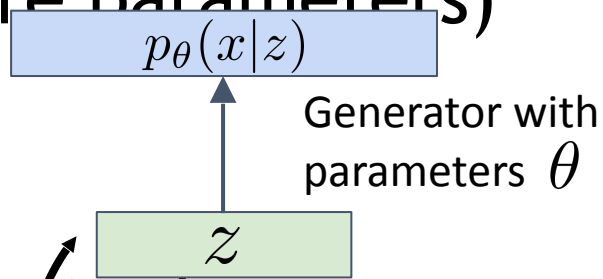
$p(z)$



Feature space / latent space

# Variational Autoencoders (VAEs)

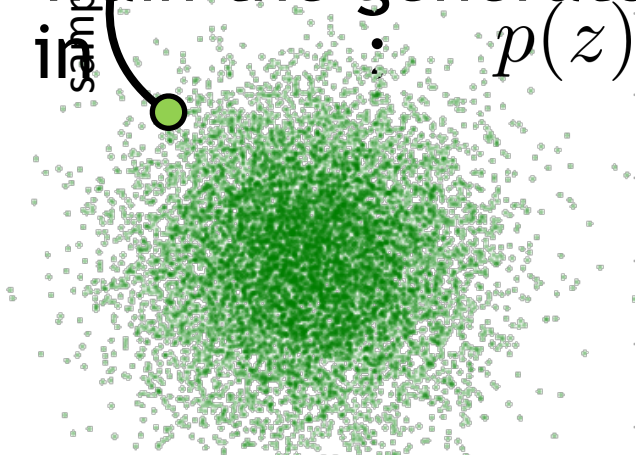
- Pick a parametric distribution  $p_\theta(x|z)$  for features
- The generator maps latent variables  $z$  to an image distribution  $p_\theta(x|z)$  (where  $\theta$  are parameters)



- Train the generator to maximize the likelihood of the data

$$p_\theta(x) = \int p_\theta(x|z) p(z) dz$$

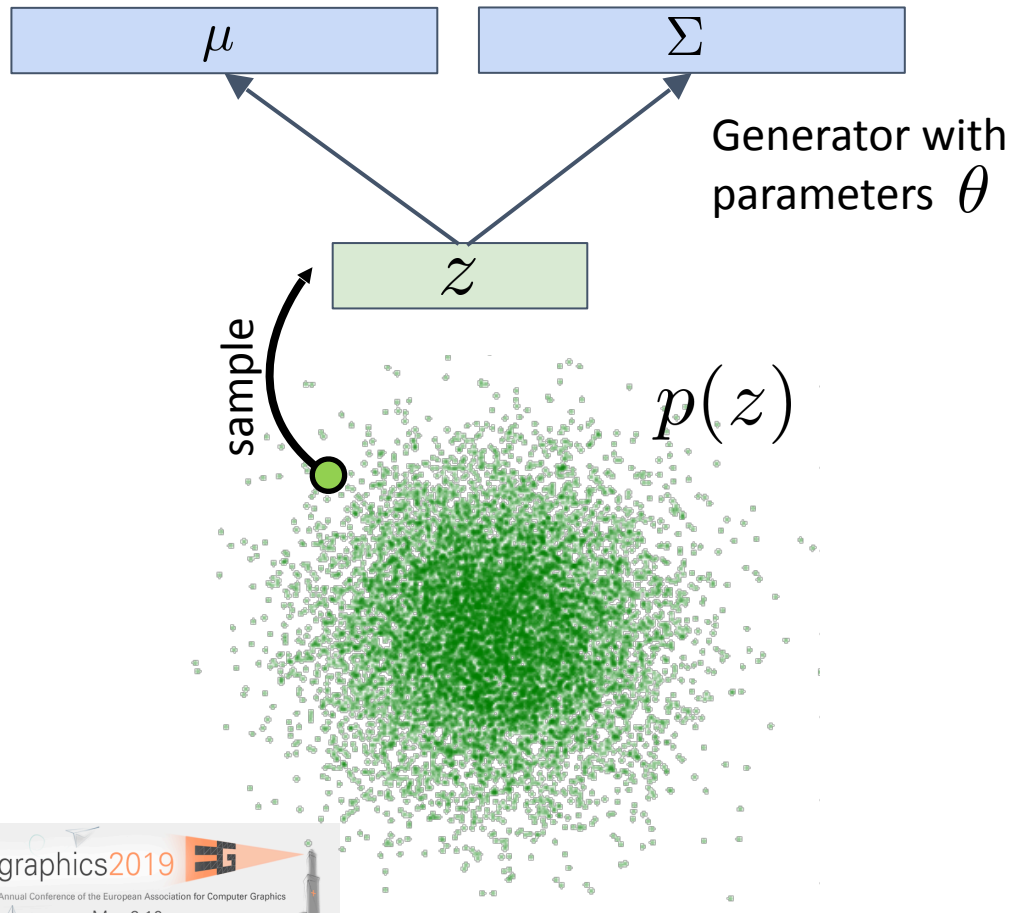
$$\max_{\theta} \sum_{x_i \in \text{data}} \log p_\theta(x_i)$$



# Outputting a Distribution

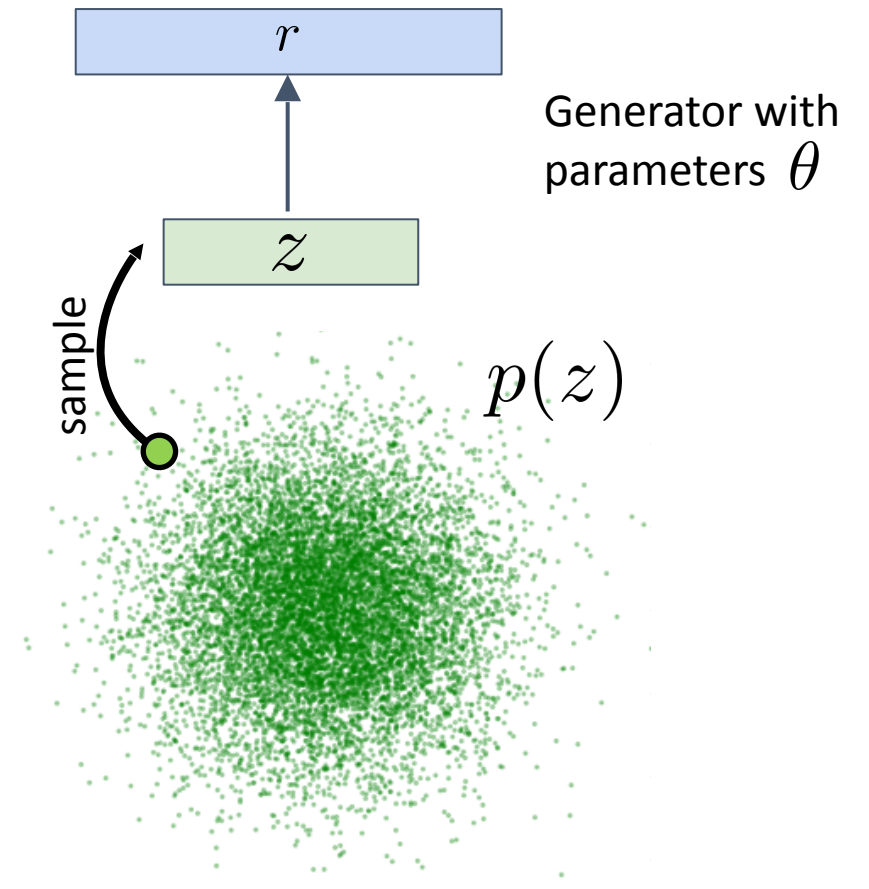
Normal distribution

$$p_{\theta}(x|z) = N(x; \mu(z), \Sigma(z))$$



Bernoulli distribution

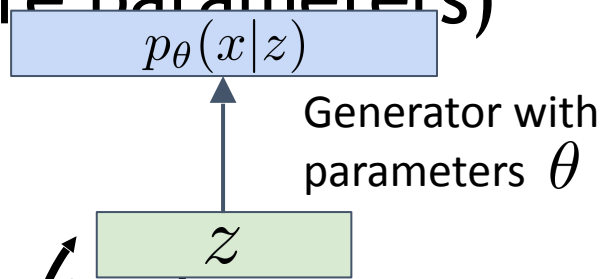
$$p_{\theta}(x|z) = \text{Bern}(x; r(z))$$





# Variational Autoencoders (VAEs)

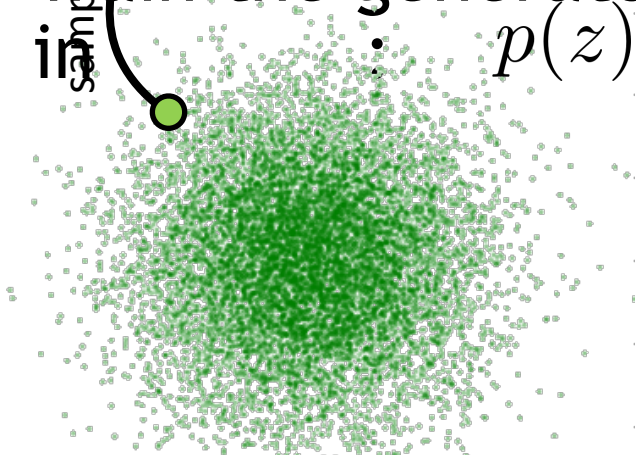
- Pick a parametric distribution  $p_\theta(x|z)$  for features
- The generator maps latent variables  $z$  to an image distribution  $p_\theta(x|z)$  (where  $\theta$  are parameters)



- Train the generator to maximize the likelihood of the data

$$p_\theta(x) = \int p_\theta(x|z) p(z) dz$$

$$\max_{\theta} \sum_{x_i \in \text{data}} \log p_\theta(x_i)$$



# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

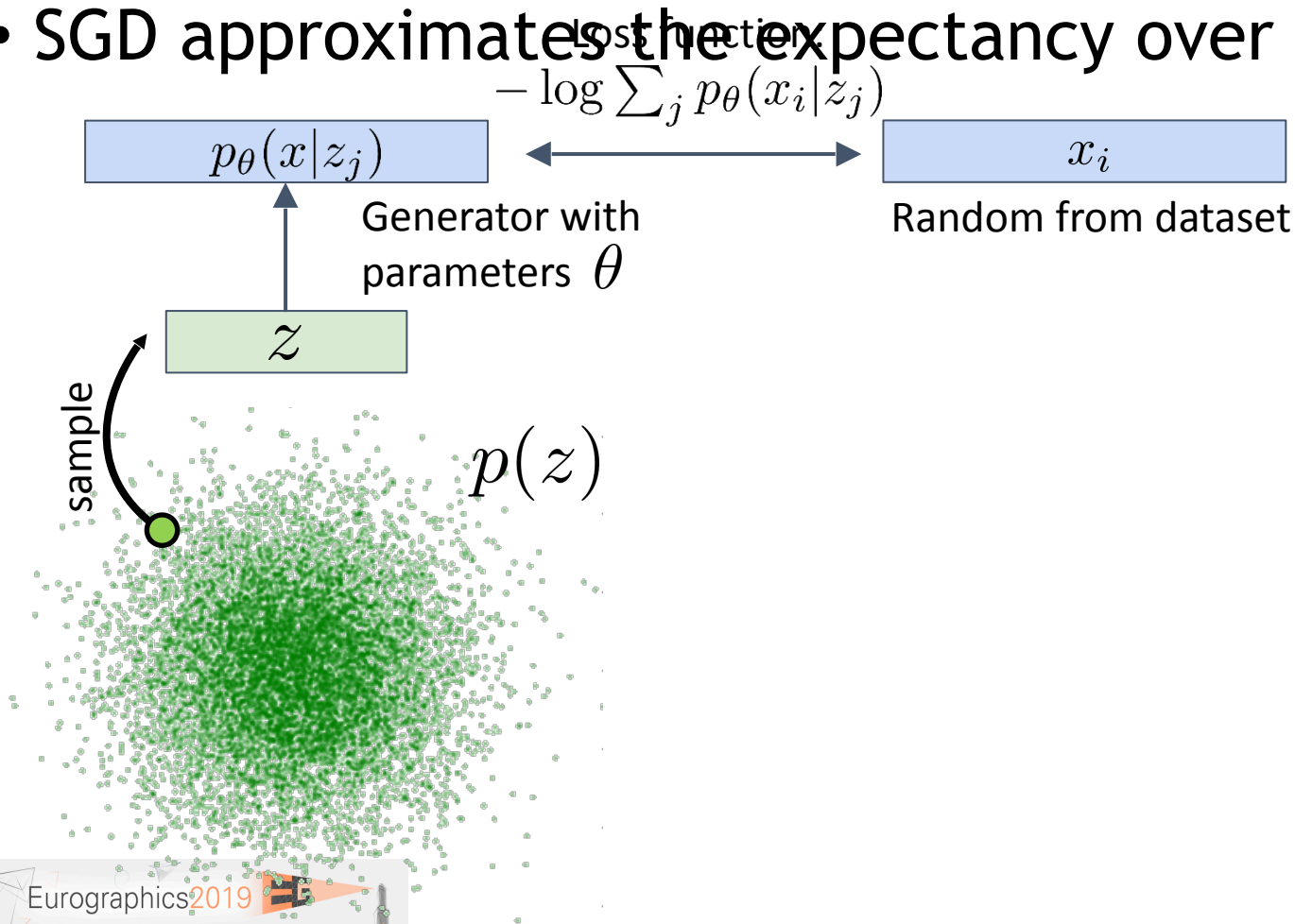
- Approximate Integral with Monte-Carlo in each iteration
- SGD approximates the sum over data generated distribution:

$$\theta^* = \arg \max_{\theta} \sum_{x_i \in \text{data}} \log \int p_{\theta}(x_i|z) p(z) dz$$

$$\theta^* \approx \arg \max_{\theta} \mathbb{E}_{x_i \sim p_{\text{data}}(x)} \log \mathbb{E}_{z \sim p(z)} p_{\theta}(x_i|z)$$

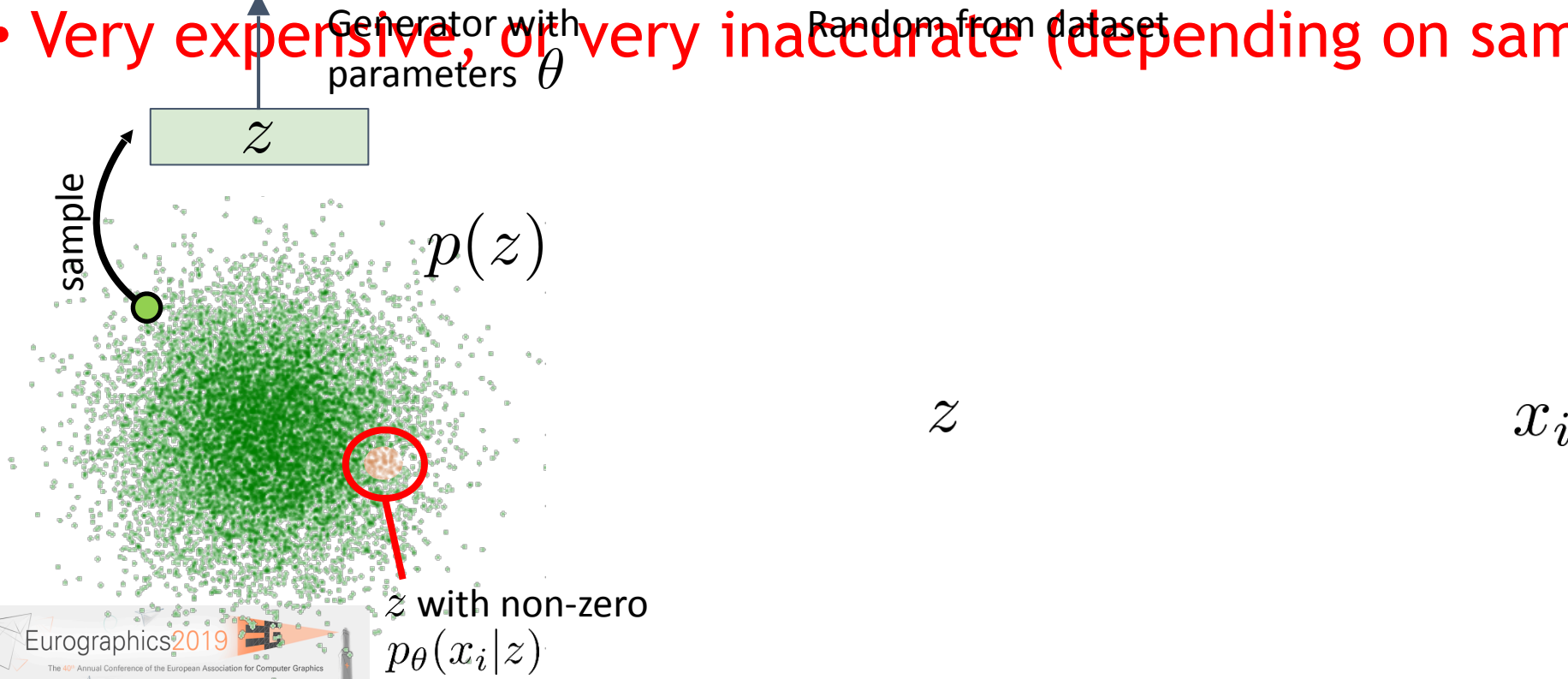
# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

- Approximate Integral with Monte-Carlo in each iteration
- SGD approximates the expectancy over data



# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

- Approximate Integral with Monte-Carlo in each iteration
- SGD approximates the <sup>Loss function</sup> expectancy over data
- $0 < p_{\theta}(x_i | z_j) < 1$   $\leftarrow$   $p$  close to a given  $x_i$
- **Very expensive, or very inaccurate (depending on sample count)**



# Variational Autoencoders (VAEs): The Encoder

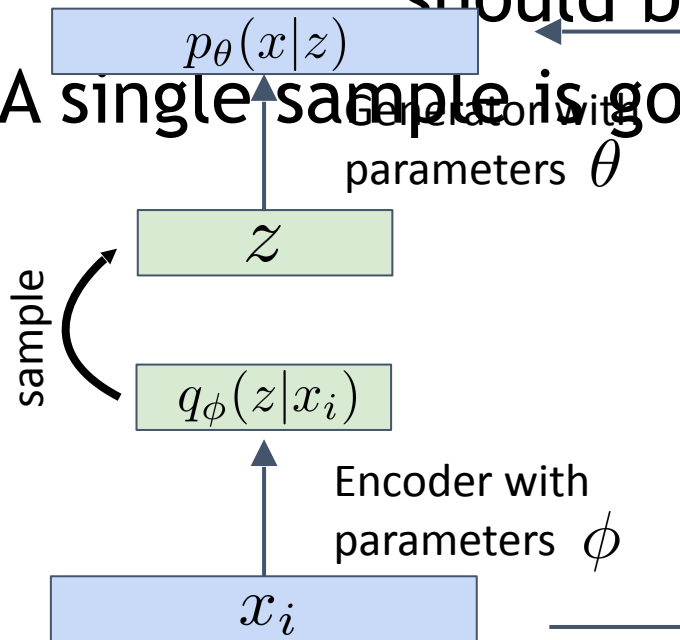
- During training, another network can learn a distribution of good for a given

Loss function:

$$-\log p_{\theta}(x_i|z)$$

$$p_{\theta}(x) = \int p_{\theta}(x|z) p(z) dz$$

- should be much smaller than
- A single sample is good enough



$z$                        $x_i$

$q_{\phi}(z|x_i)$                        $p(z)$

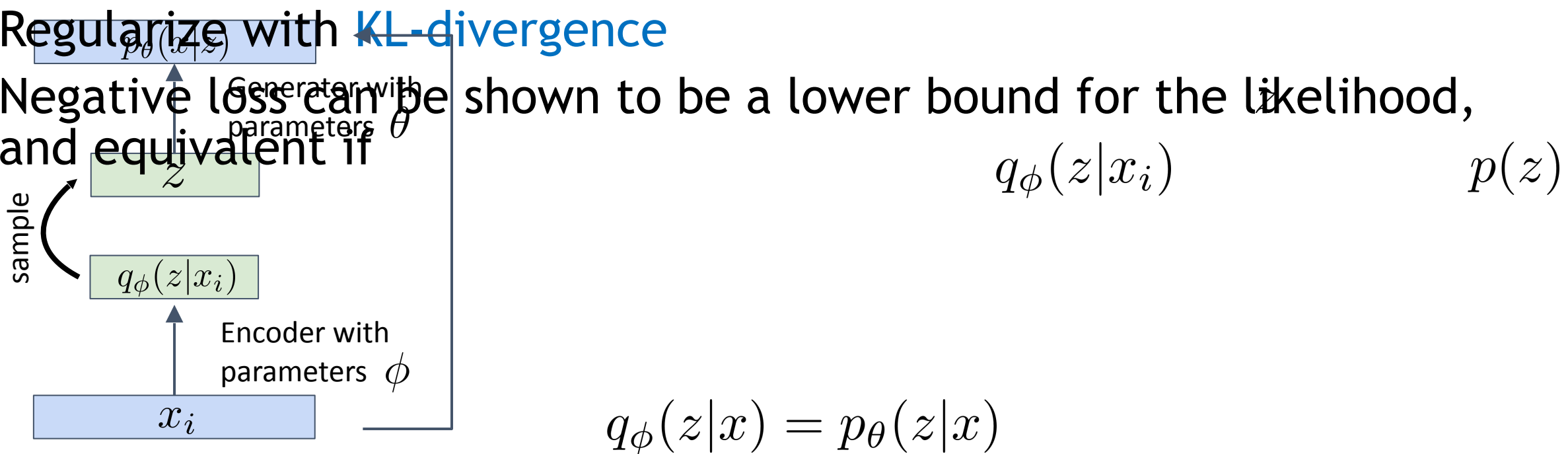


# Variational Autoencoders (VAEs): The Encoder

- Can we still easily sample a new  $x$  ?
- Need to make sure  $q_\phi(z|x_i)$  approximates  $p_\theta(x_i|z)$
- Regularize with **KL-divergence**
- Negative loss can be shown to be a lower bound for the likelihood, and equivalent if

Loss function:

$$-\log p_\theta(x_i|z) + KL(q_\phi(z|x_i) || p(z))$$

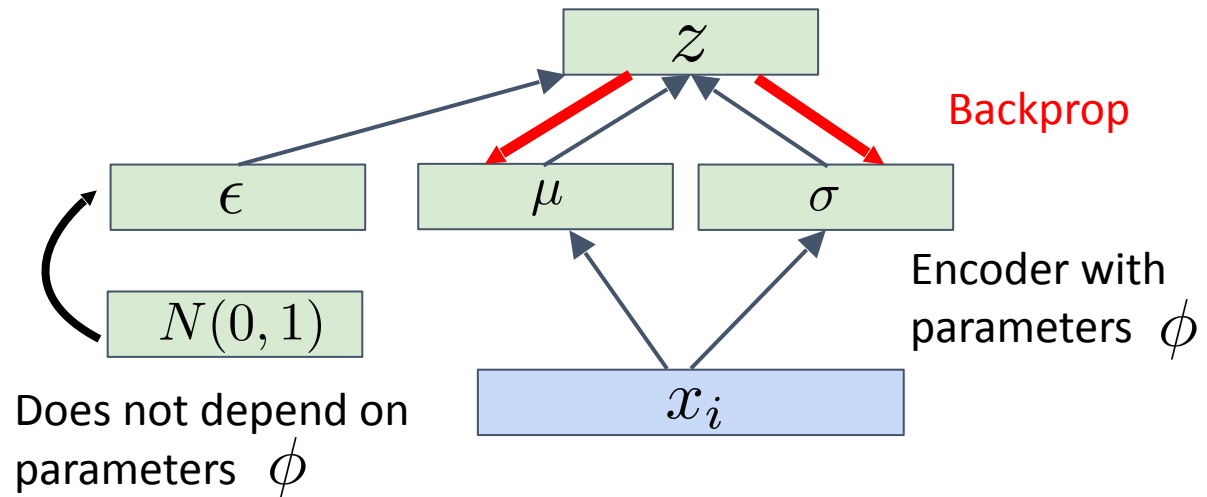
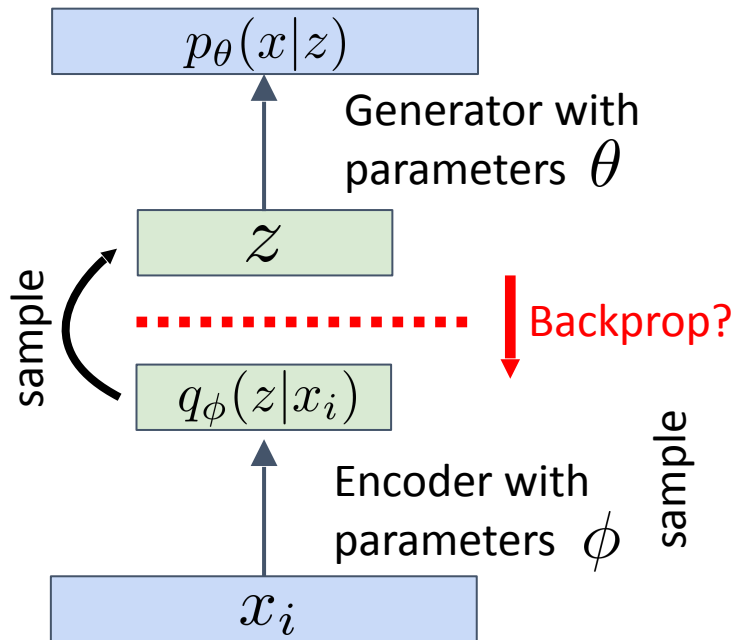


# Reparameterization Trick

Example when  $q_\phi(z|x_i) = N(z; \mu(x_i), \sigma(x_i))$

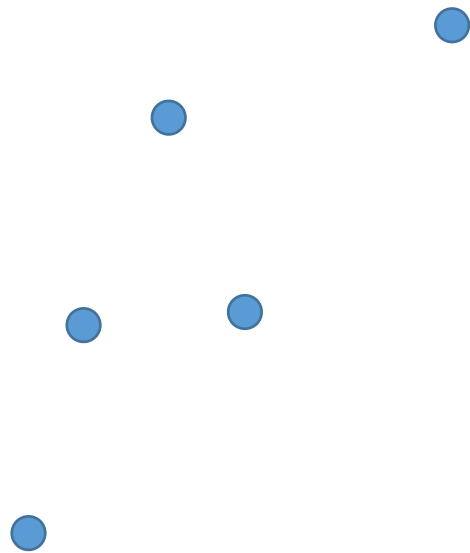
$z = \sigma + \mu \cdot \epsilon$ , where  $\epsilon \sim N(0, 1)$

$$\frac{\partial z}{\partial \phi} = \frac{\partial \mu}{\partial \phi} + \frac{\partial \sigma}{\partial \phi} \cdot \epsilon$$

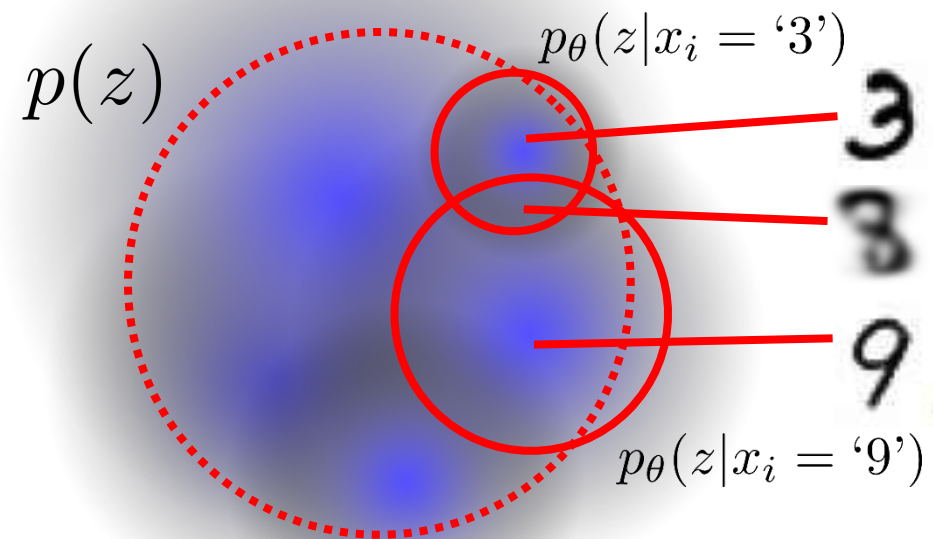


# Feature Space of Autoencoders vs. VAEs

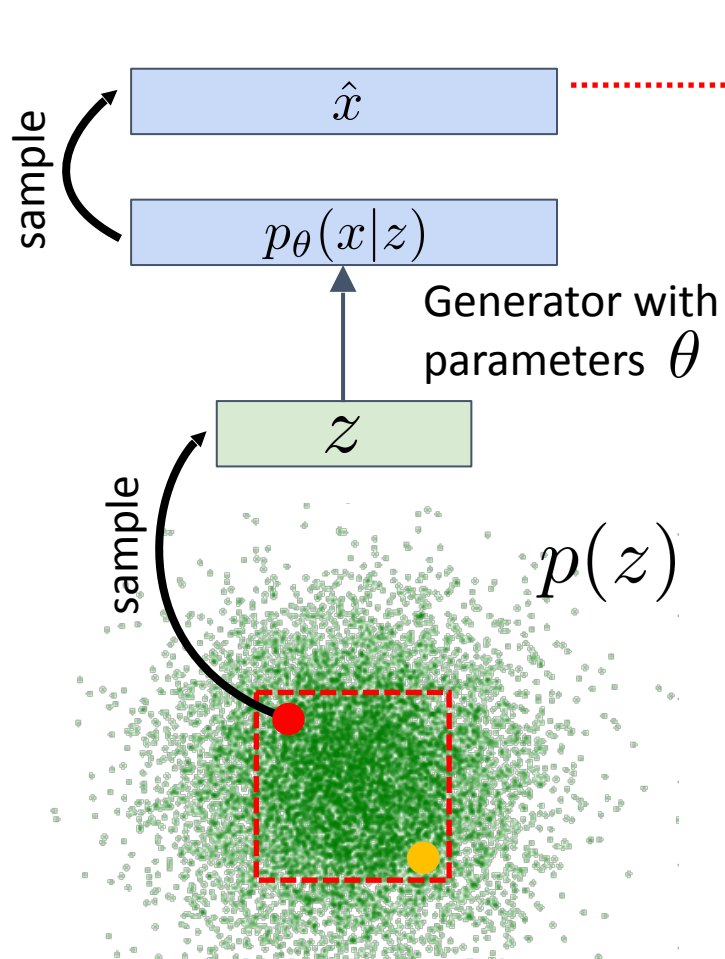
Autoencoder



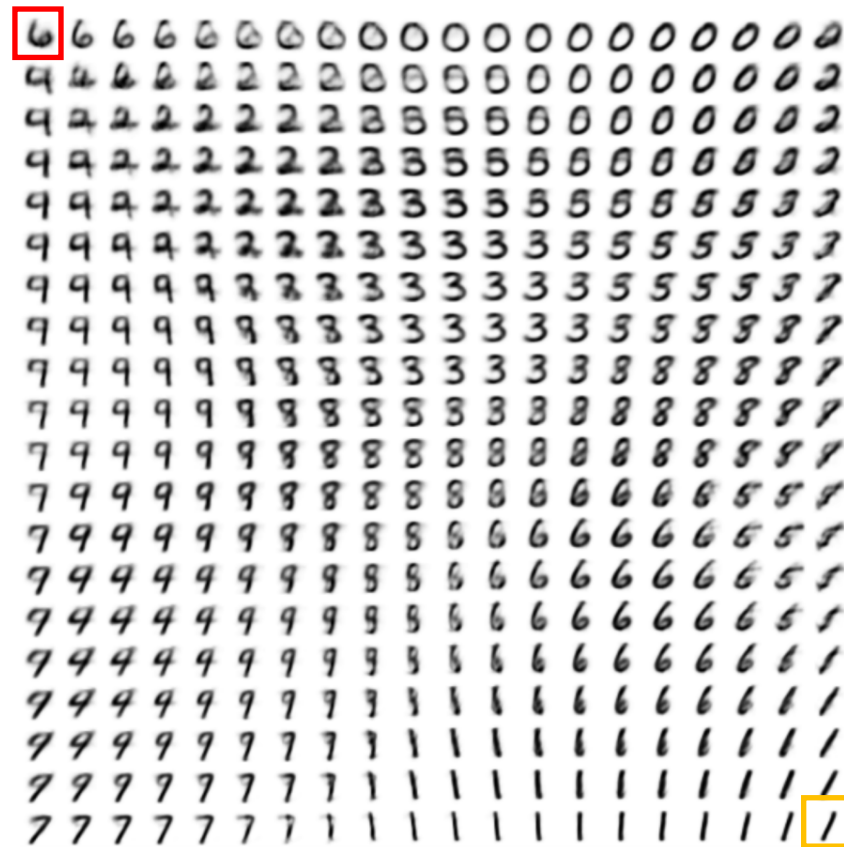
VAE



# Generating Data



MNIST



Frey Faces



## VAE on MNIST

<https://www.siares.com/projects/variational-autoencoder>



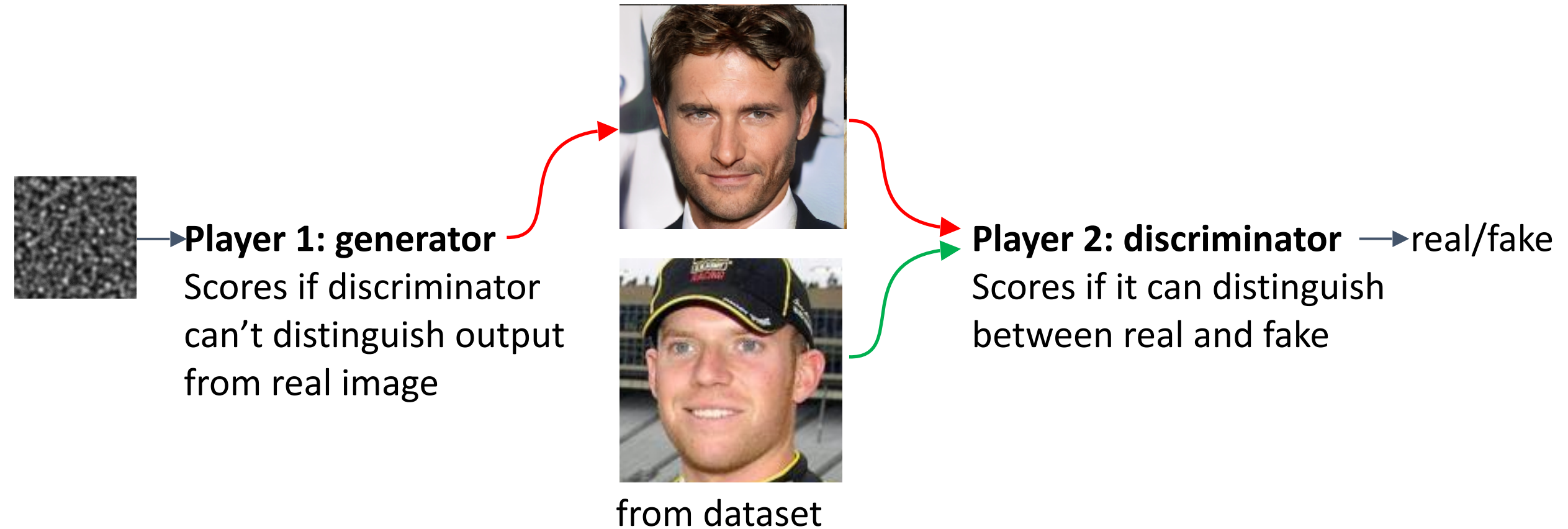


# Variational Autoencoder

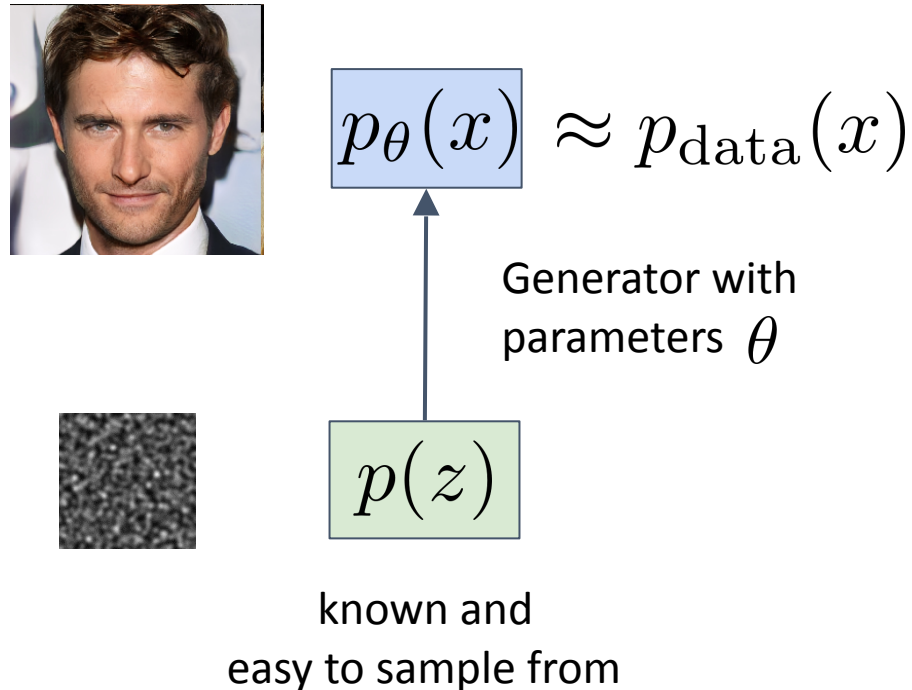
[geometry.cs.ucl.ac.uk/creativeai](http://geometry.cs.ucl.ac.uk/creativeai)



# Generative Adversarial Networks



# Generative Models



How to measure similarity of  $p_{\theta}(x)$  and  $p_{\text{data}}(x)$ ?

1) Likelihood of data in  $p_{\theta}(x)$

## Variational Autoencoders (VAEs)

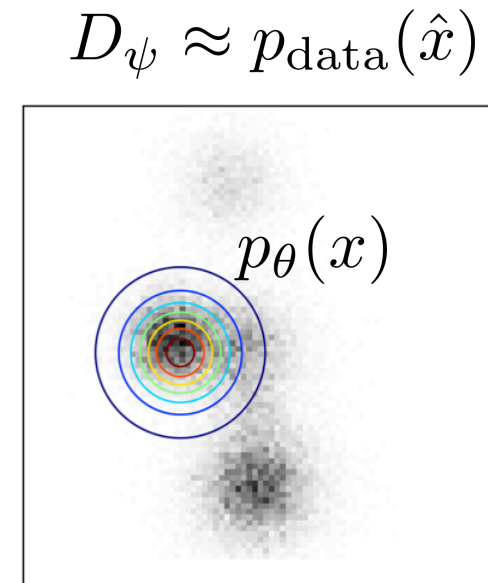
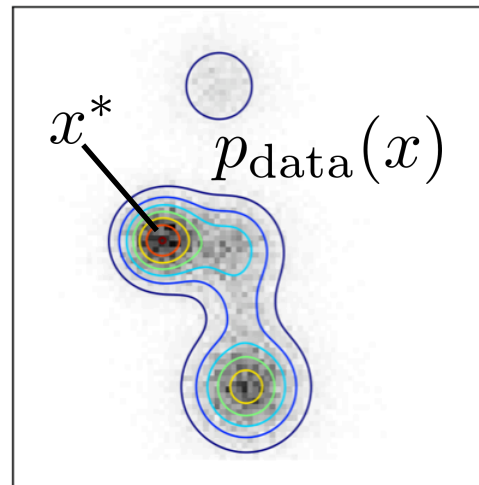
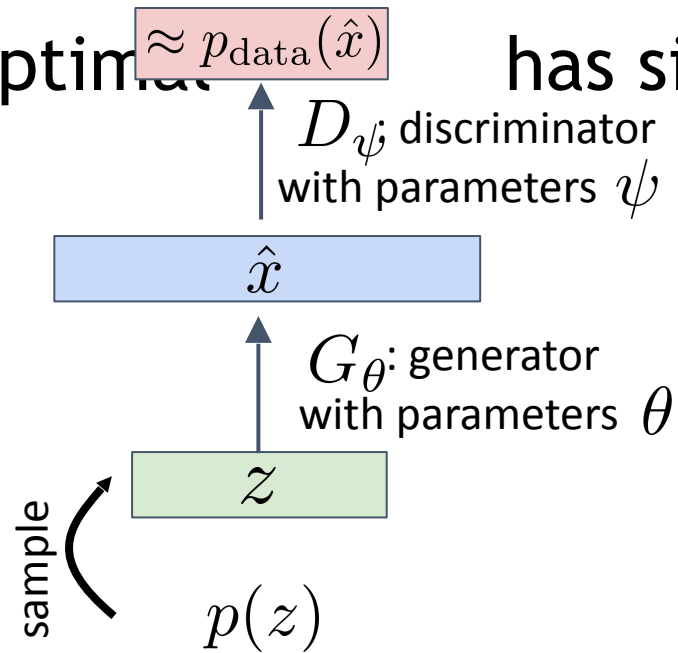
2) Adversarial game:

*Discriminator* distinguishes  $p_{\theta}(x)$  and  $p_{\text{data}}(x)$  vs *Generator* makes it hard to distinguish

## Generative Adversarial Networks (GANs)

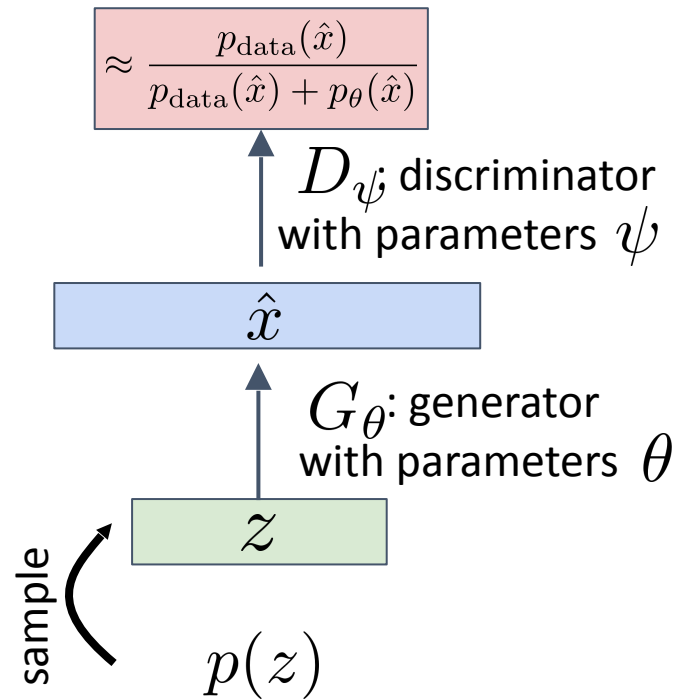
# Why Adversarial?

- If discriminator approximates  $p_{\text{data}}(x)$  :
- at maximum of  $\int p_{\text{data}}(x) p_{\theta}(x) dx$  has lowest loss
- Optimal  $p_{\theta}(x)$  has single mode at  $x^*$ , small variance



# Why Adversarial?

- For GANs, the discriminator instead approximates:



$$\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)} \rightarrow \text{depends on the generator}$$

$$D_\psi \approx p_{\text{data}}(\hat{x}) \quad D_\psi \approx \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)}$$

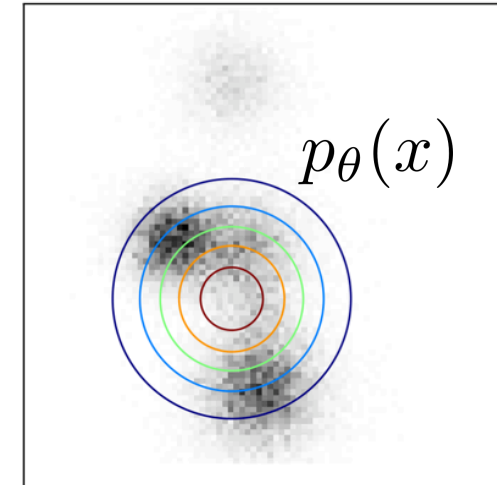
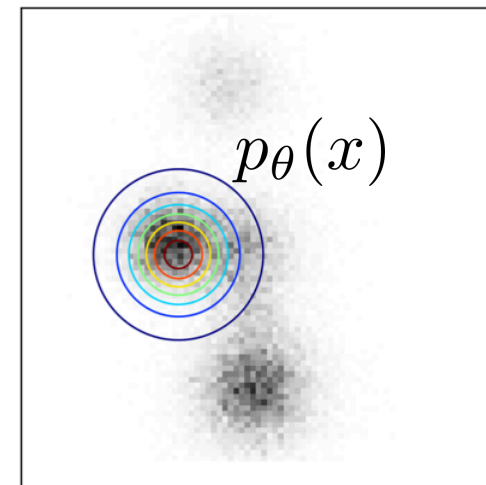
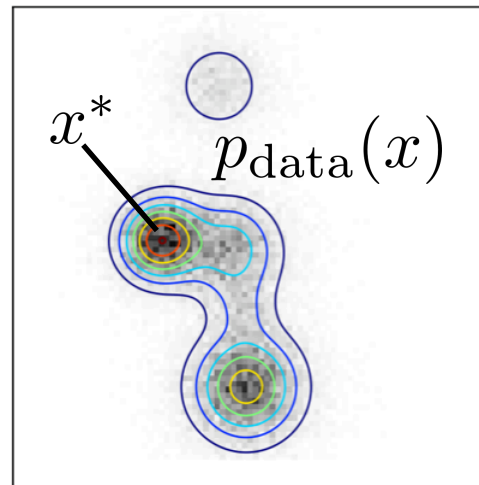
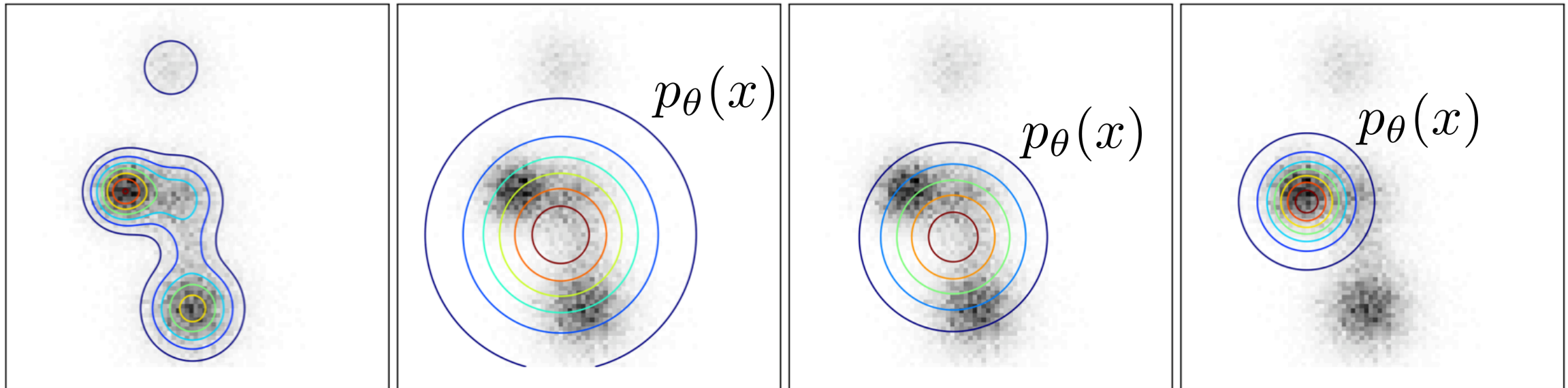


Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

# Why Adversarial?



$p_{\text{data}}(x)$

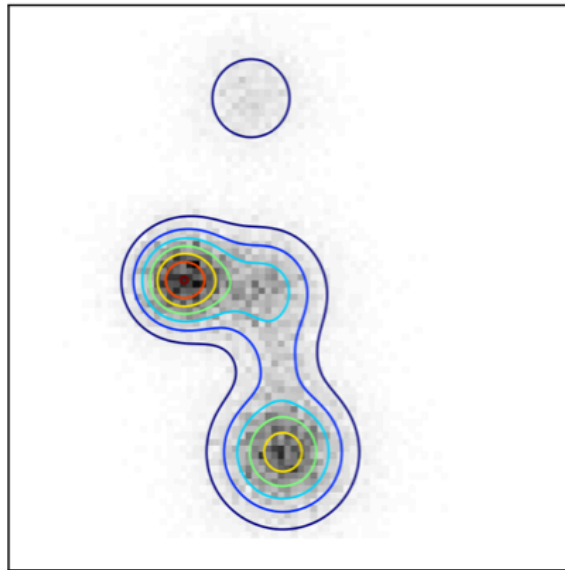
VAEs:  
Maximize likelihood of  
**data samples** in  $p_{\theta}(x)$

GANs:  
Adversarial game

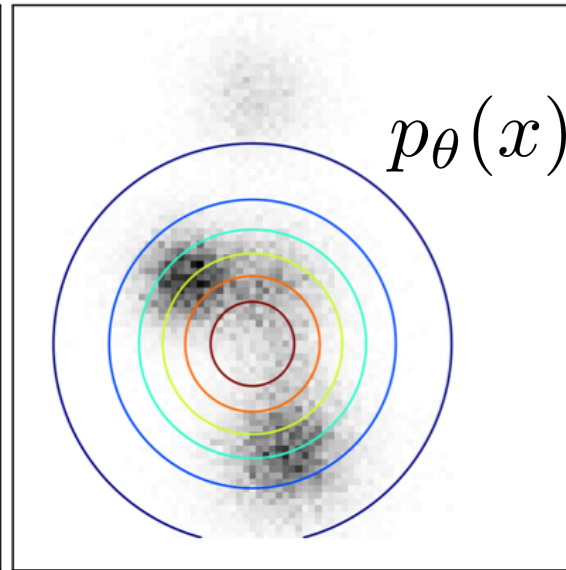
Maximize likelihood of  
**generator samples** in  
approximate  $p_{\text{data}}(x)$

# Why Adversarial?

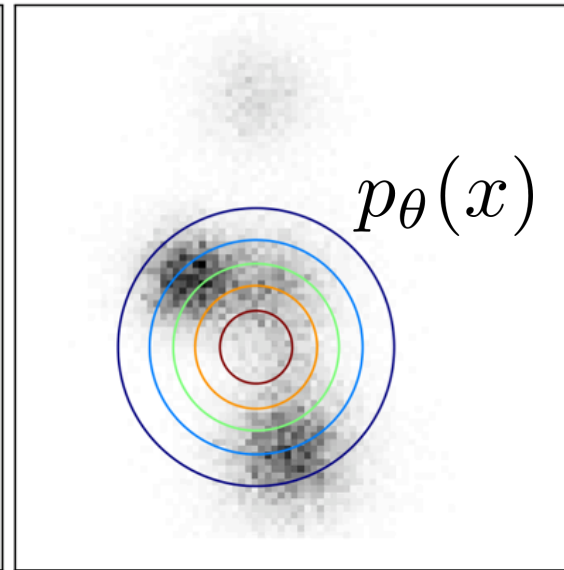
$$\approx KL(p_{\text{data}} \parallel p_{\theta}) \quad \approx JS(p_{\text{data}} \parallel p_{\theta}) \quad \approx KL(p_{\theta} \parallel p_{\text{data}})$$



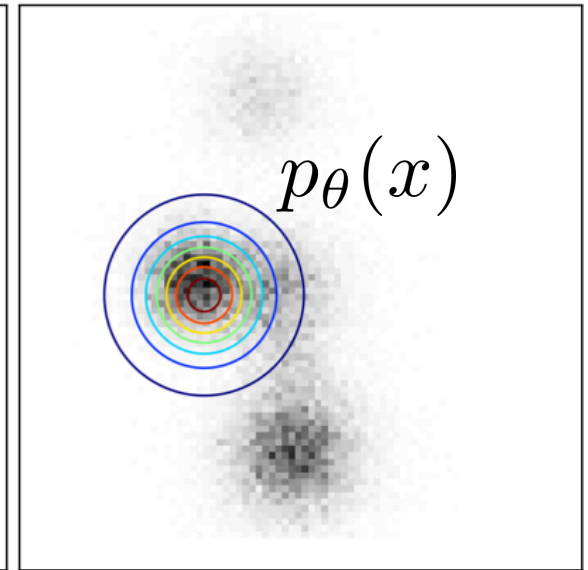
$p_{\text{data}}(x)$



VAEs:  
Maximize likelihood of  
**data samples** in  $p_{\theta}(x)$



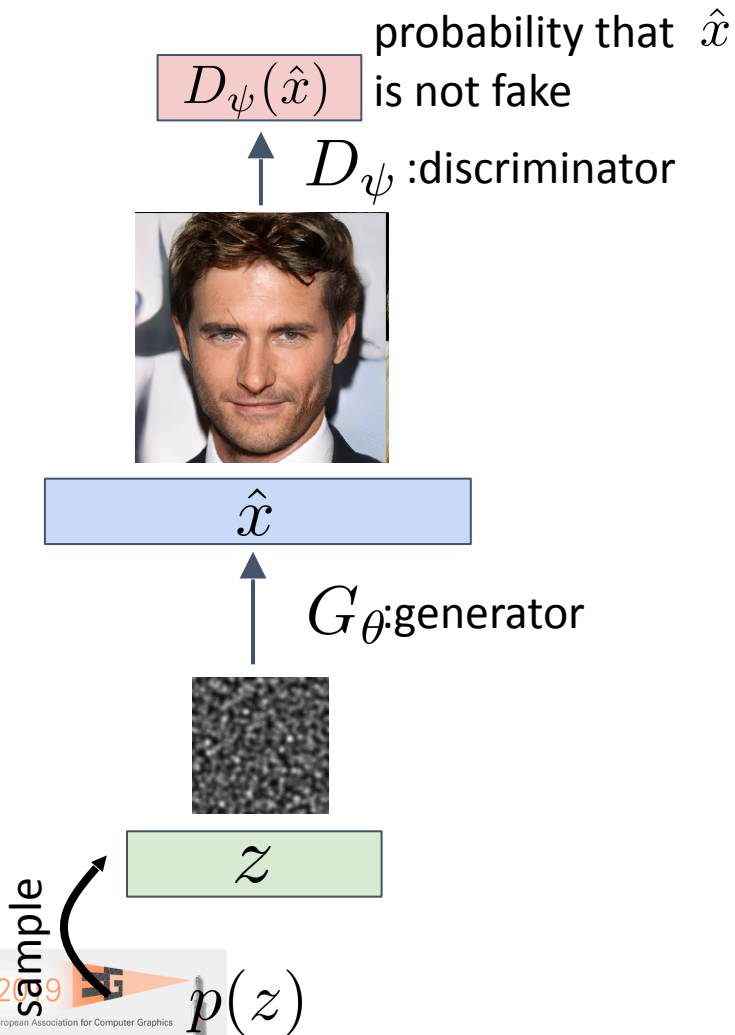
GANs:  
Adversarial game



Maximize likelihood of  
**generator samples** in  
approximate  $p_{\text{data}}(x)$



# GAN Objective



fake/real classification loss (BCE):

$$L(\theta, \psi) = -0.5 \mathbb{E}_{x \sim p_{\text{data}}} \log D_\psi(x) - 0.5 \mathbb{E}_{x \sim p_\theta} \log(1 - D_\psi(x))$$

Discriminator objective:

$$\min_{\psi} L(\theta, \psi)$$

Generator objective:

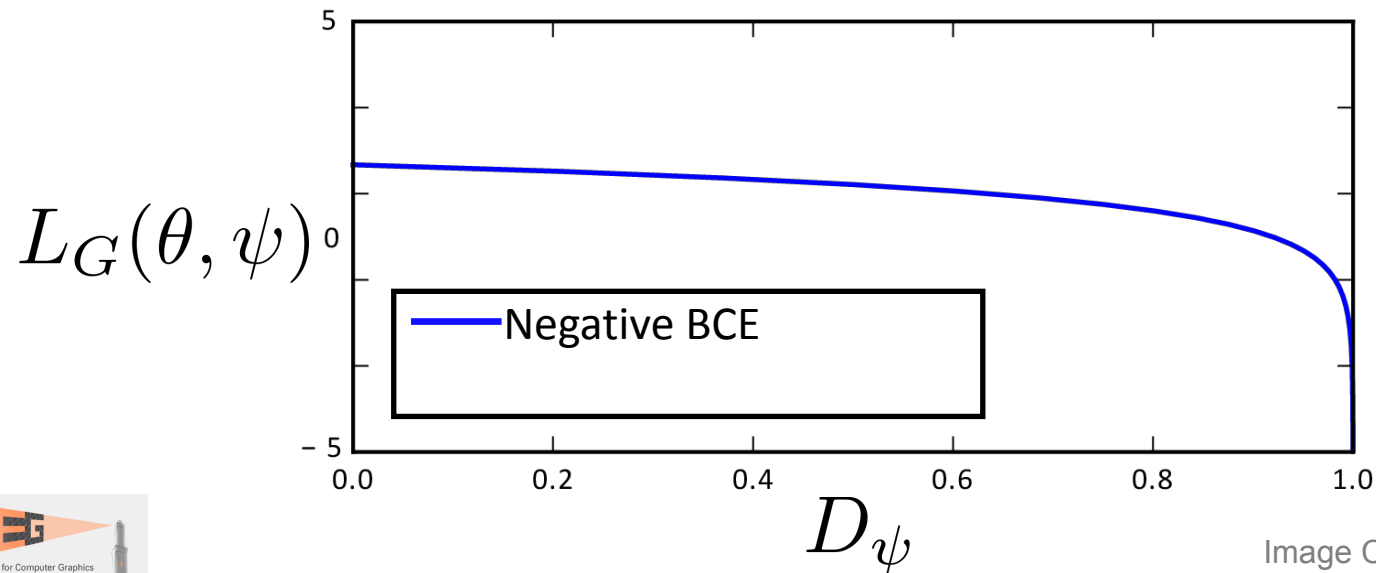
$$\max_{\theta} L(\theta, \psi)$$

# Non-saturating Heuristic

$$L(\theta, \psi) = -0.5 \mathbb{E}_{x \sim p_{\text{data}}} \log D_{\psi}(x) \\ - 0.5 \mathbb{E}_{x \sim p_{\theta}} \log(1 - D_{\psi}(x))$$

Generator loss is negative binary cross-entropy:

$$L_G(\theta, \psi) = 0.5 \mathbb{E}_{x \sim p_{\theta}} \log(1 - D_{\psi}(x)) \text{ poor convergence}$$



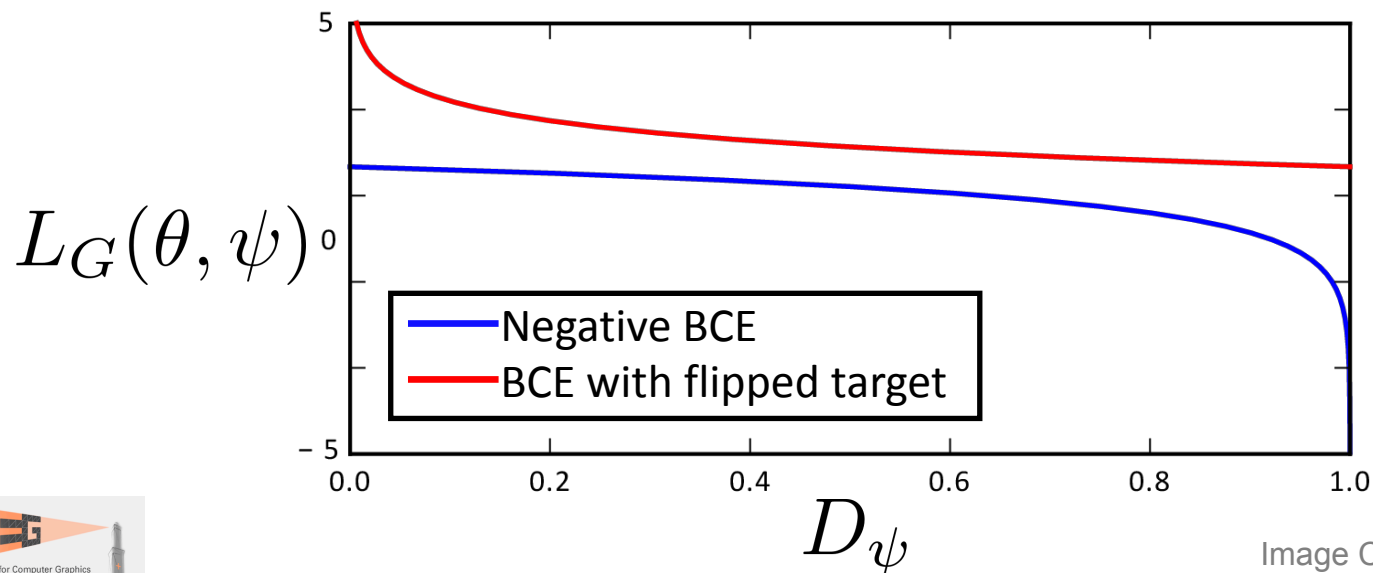
# Non-saturating Heuristic

Generator loss is negative binary cross-entropy:

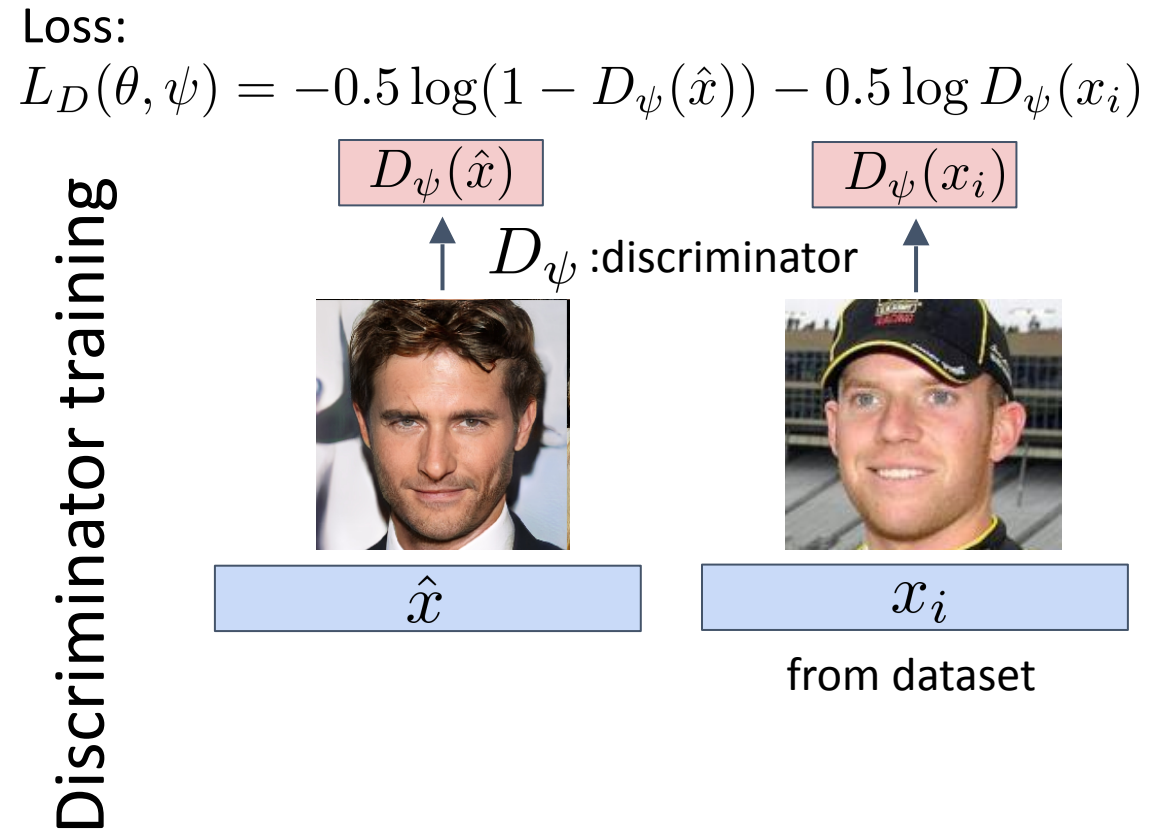
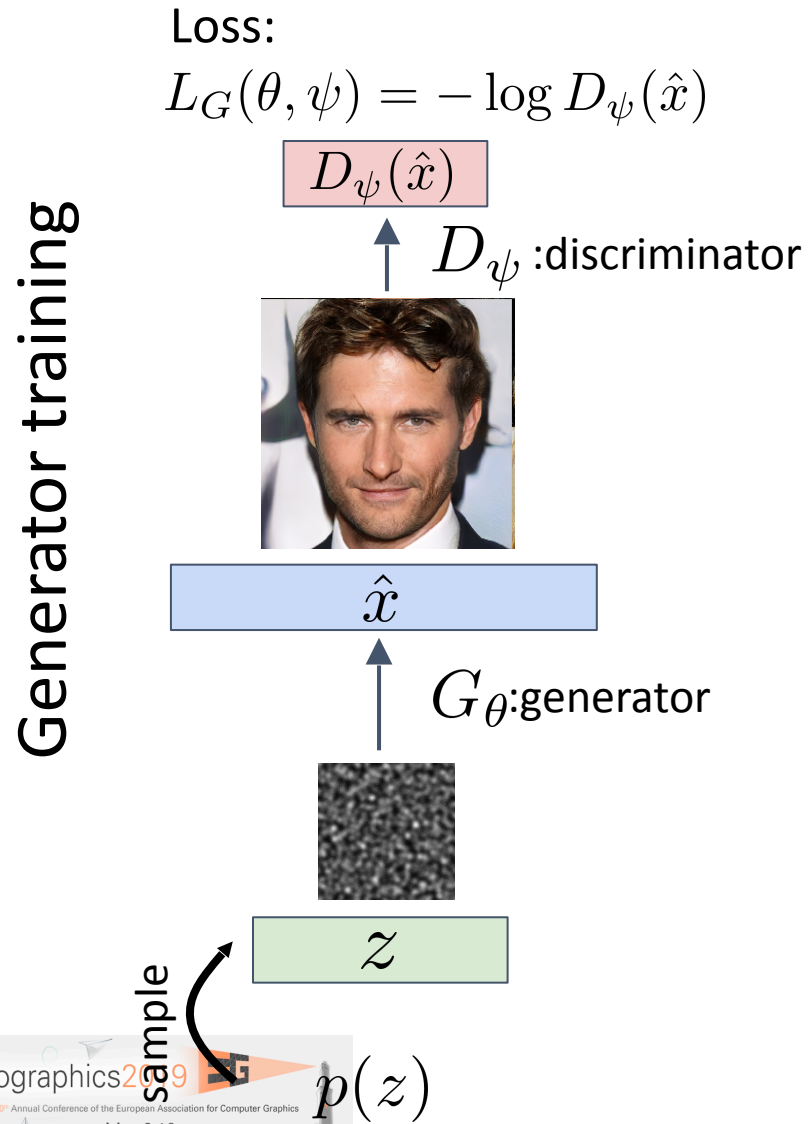
$$L_G(\theta, \psi) = 0.5 \mathbb{E}_{x \sim p_\theta} \log(1 - D_\psi(x)) \quad \text{poor convergence}$$

Flip target class instead of flipping the sign for generator loss:

$$L_G(\theta, \psi) = -0.5 \mathbb{E}_{x \sim p_\theta} \log D_\psi(x) \quad \text{good convergence – like BCE}$$



# GAN Training



Interleave in each training step

# DCGAN

- First paper to successfully use CNNs with GANs
- Due to using novel components (at that time) like batch norm., ReLUs, etc.

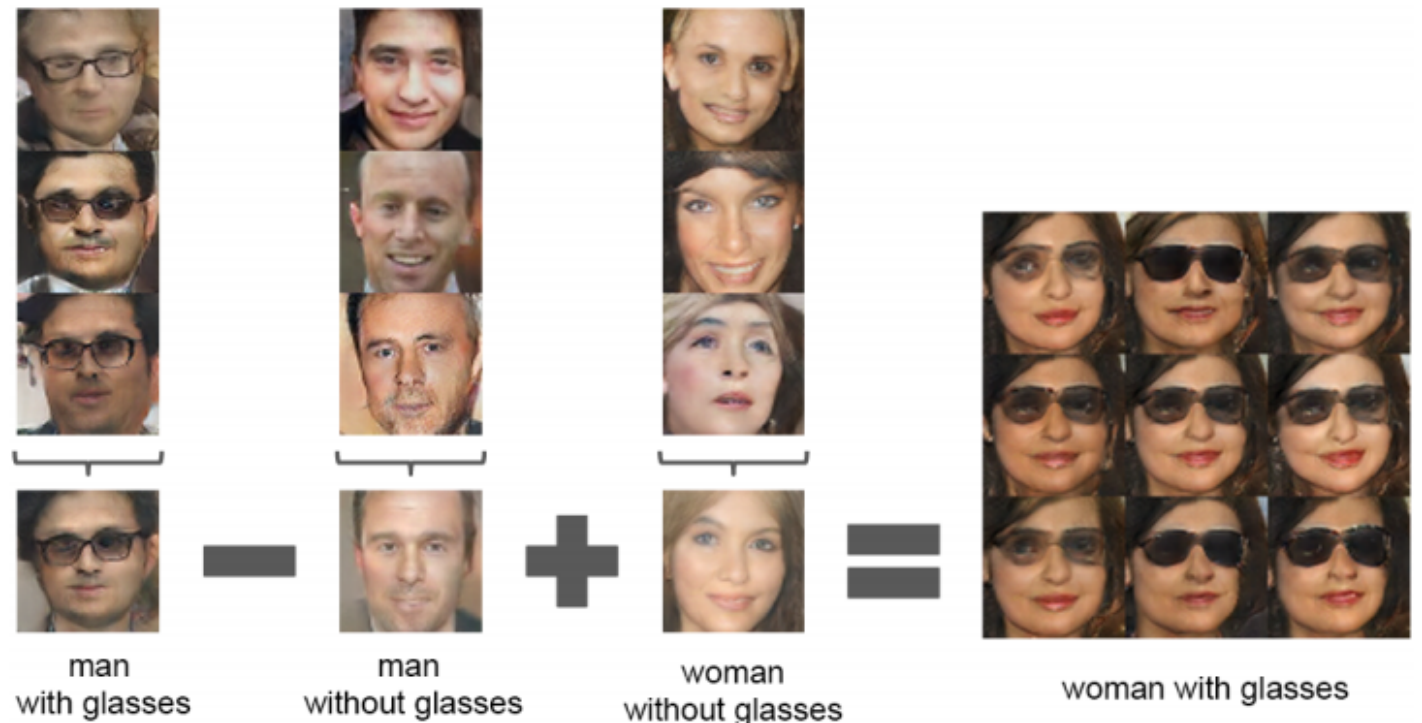
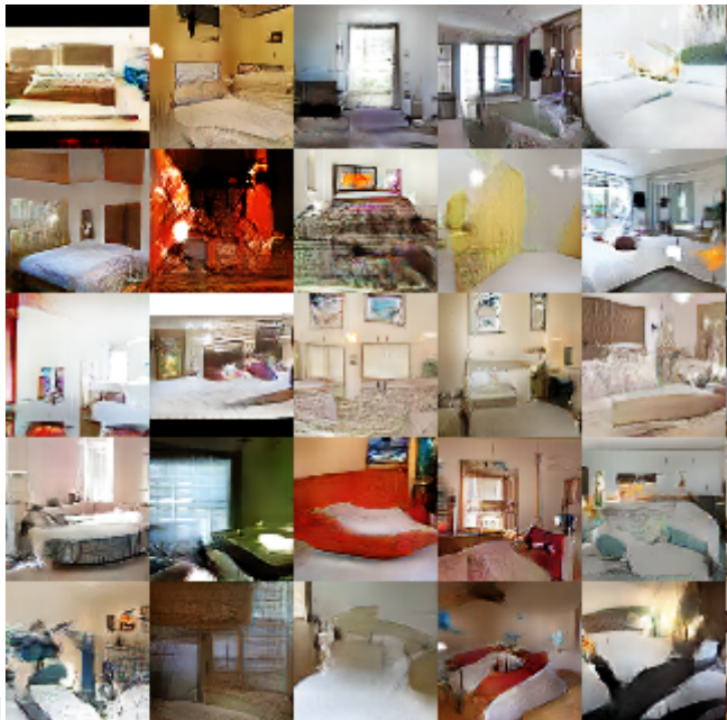


Image Credit: *Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks*, Radford et al.

# Generative Adversarial Network

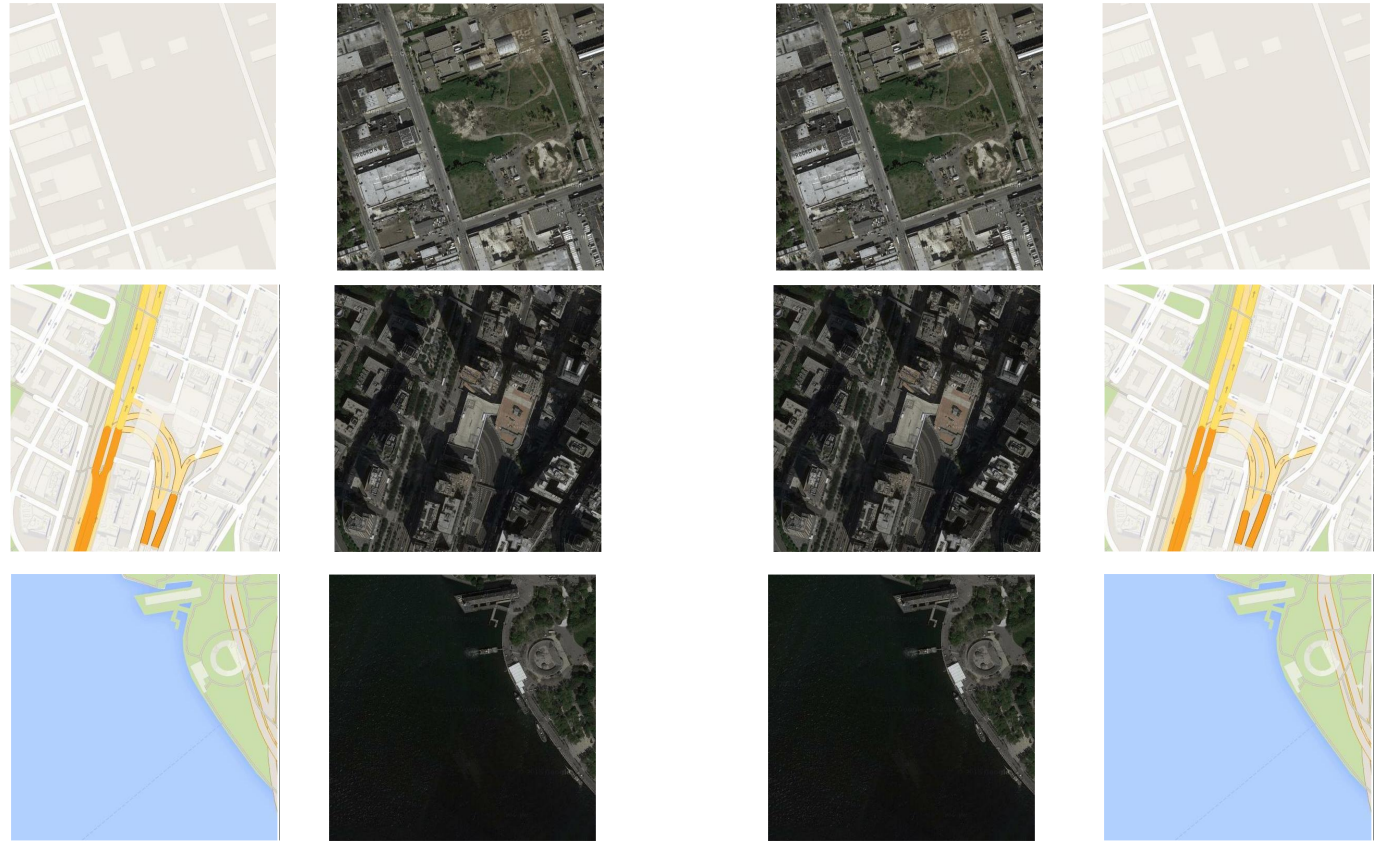
[geometry.cs.ucl.ac.uk/creativeai](http://geometry.cs.ucl.ac.uk/creativeai)





# Conditional GANs (CGANs)

- $\approx$  learn a mapping between images from example pairs
- Approximate sampling from a conditional distribution  $p_{\text{data}}(x | c)$



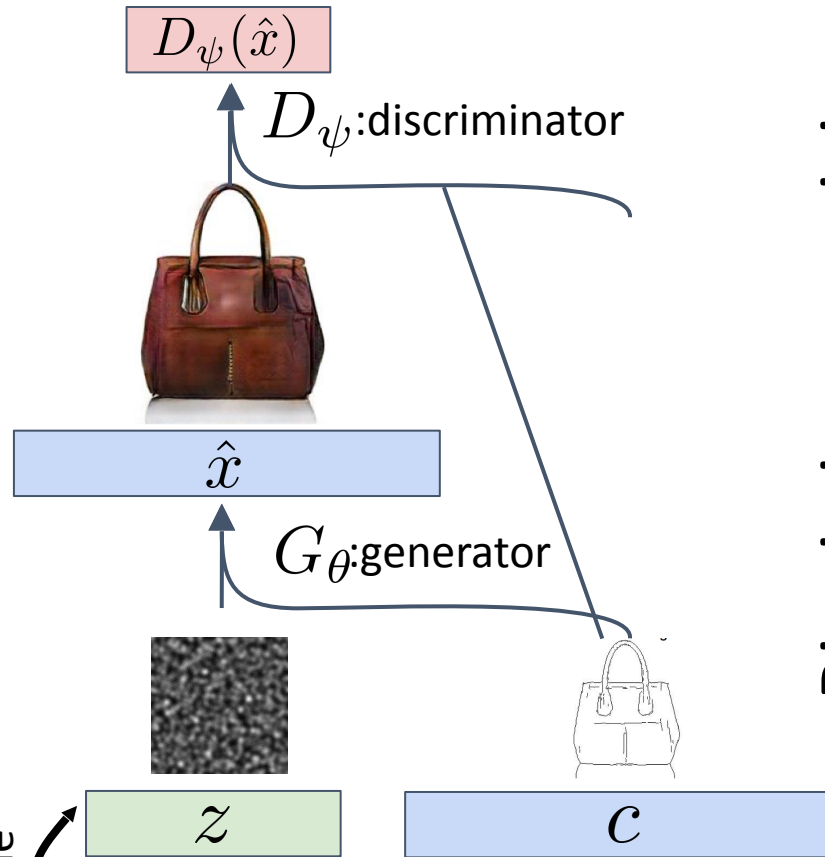


# Conditional GANs

Loss:

$$L_G(\theta, \psi) = -\log D_\psi(\hat{x}, c)$$

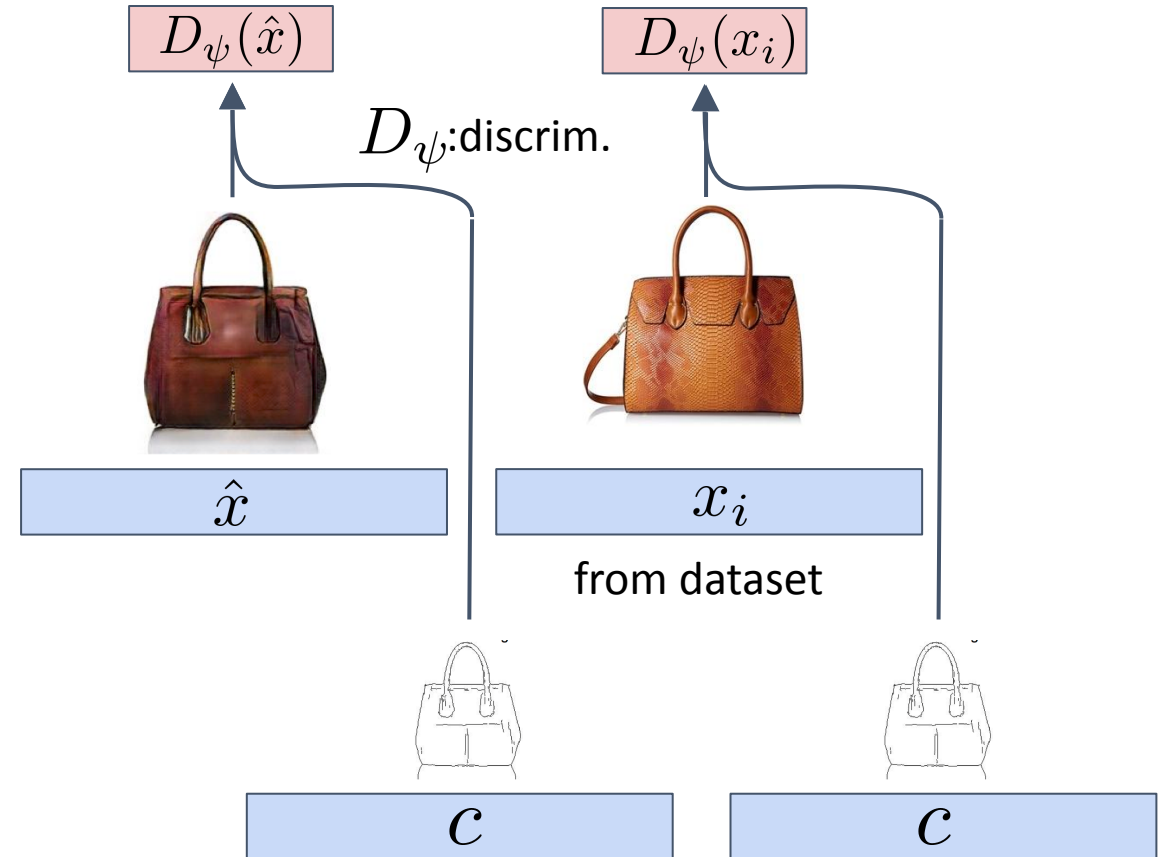
Generator training



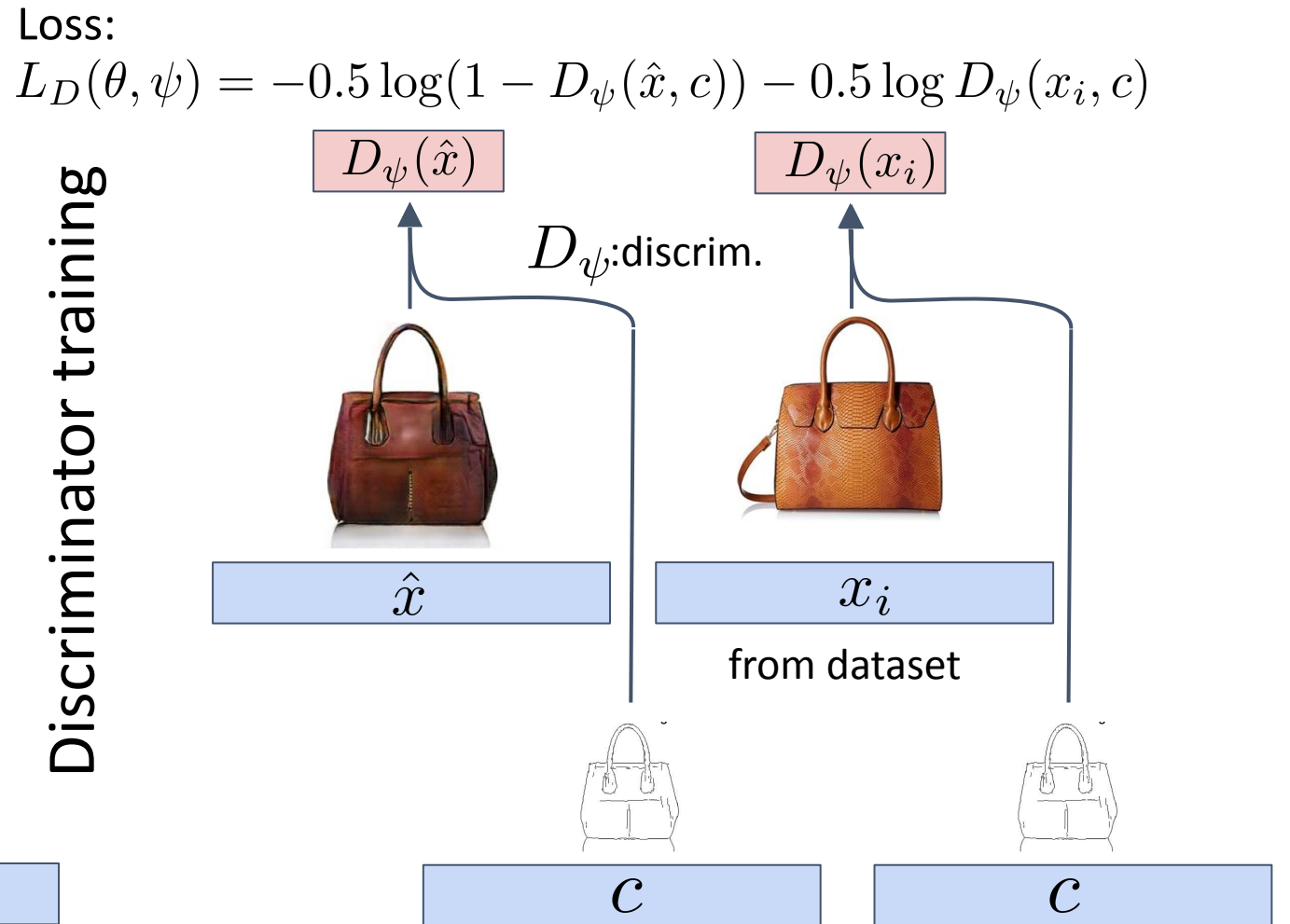
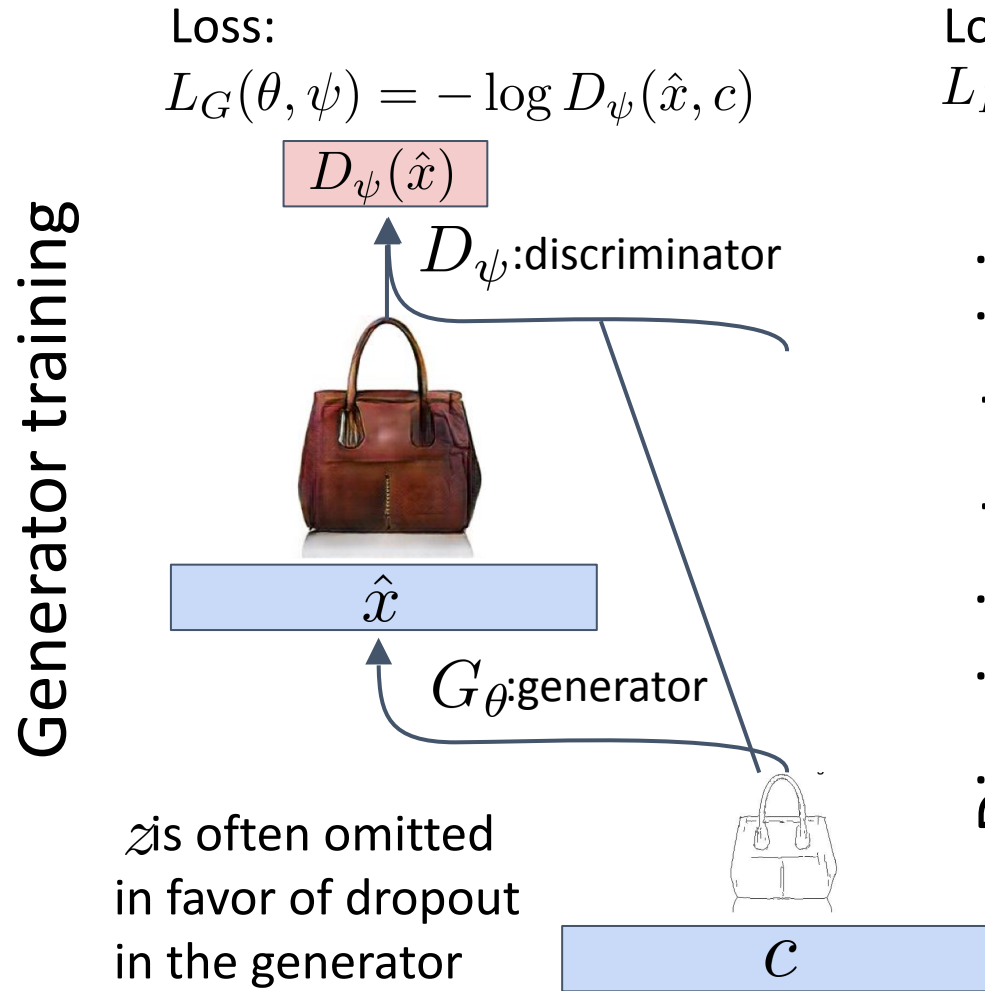
Loss:

$$L_D(\theta, \psi) = -0.5 \log(1 - D_\psi(\hat{x}, c)) - 0.5 \log D_\psi(x_i, c)$$

Discriminator training



# Conditional GANs: Low Variation per Condition



# CGAN

<https://affinelayer.com/pixsrv/index.html>



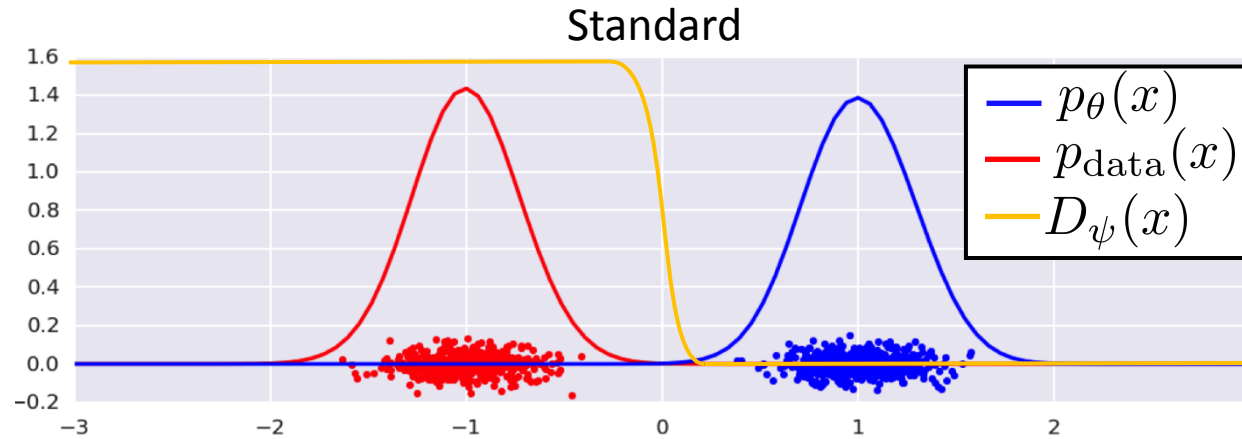
# Unstable Training

GAN training can be unstable

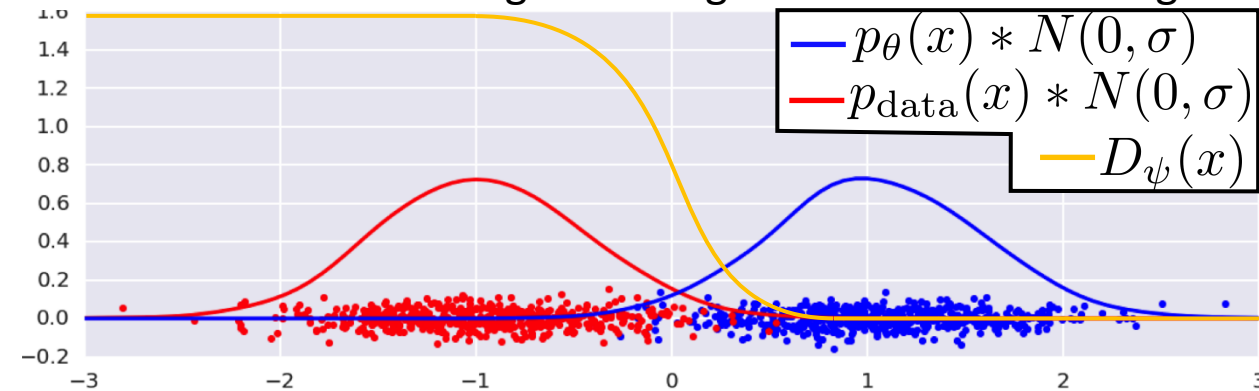
Three current research problems (may be related):

- Reaching a Nash equilibrium (the gradient for both  $L_G$  and  $L_D$  is 0)
- $p_\theta$  and  $p_{\text{data}}$  initially don't overlap
- Mode Collapse

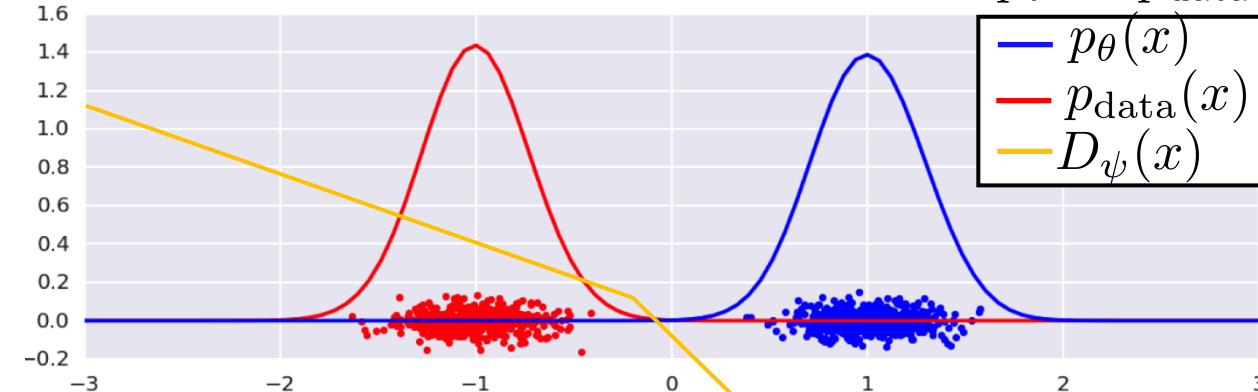
# Generator and Data Distribution Don't Overlap



Instance noise: adding noise to generated and real images



Wasserstein GANs: EMD as distance between  $p_{\theta}$  and  $p_{\text{data}}$



Roth et al. suggest an analytic convolution with a gaussian:

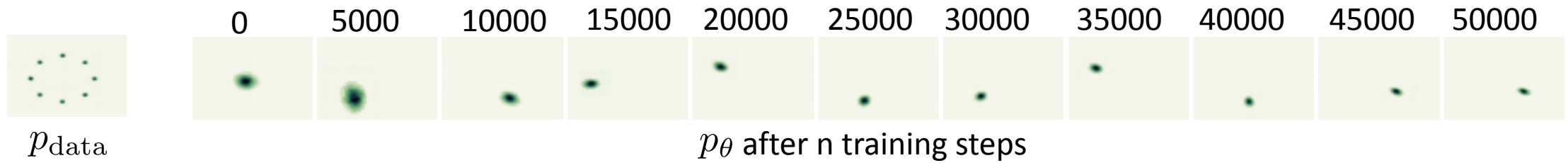
*Stabilizing Training of Generative Adversarial Networks*

*through Regularization, Roth et al. 2017*

# Mode Collapse

$$\text{Optimal } D_{\psi}(x): \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\theta}(x)}$$

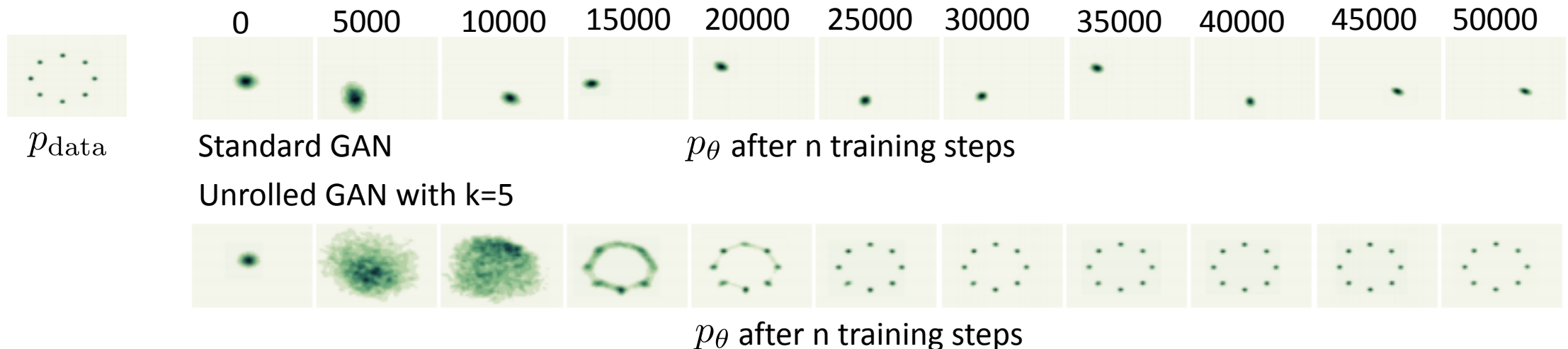
$p_{\theta}$  only covers one or a few modes of  $p_{\text{data}}$



# Mode Collapse

## Solution attempts:

- Minibatch comparisons: Discriminator can compare instances in a minibatch (*Improved Techniques for Training GANs*, Salimans et al.)
- Unrolled GANs: Take  $k$  steps with the discriminator in each iteration, and backpropagate through all of them to update the generator





# Summary

- Autoencoders
  - Can infer useful latent representation for a dataset
  - Bad generators
- VAEs
  - Can infer a useful latent representation for a dataset
  - Better generators due to latent space regularization
  - Lower quality reconstructions and generated samples (usually blurry)
- GANs
  - Can not find a latent representation for a given sample (no encoder)
  - Usually better generators than VAEs
  - Currently unstable training (active research)