



Deep Learning for Graphics

Unsupervised Learning

Niloy Mitra

UCL

Iasonas Kokkinos

UCL/Facebook

Paul Guerrero

UCL

Vladimir Kim

Adobe Research

Kostas Rematas

U Washington

Tobias Ritschel

UCL



Timetable

	Niloy	Iasonas	Paul	Vova	Kostas	Tobias
Introduction	X	X	X			X
Theory	X					
NN Basics		X				X
Supervised Applications		X	X			
Data						X
Unsupervised Applications			X			
Beyond 2D	X			X		
Outlook	X	X	X	X	X	X

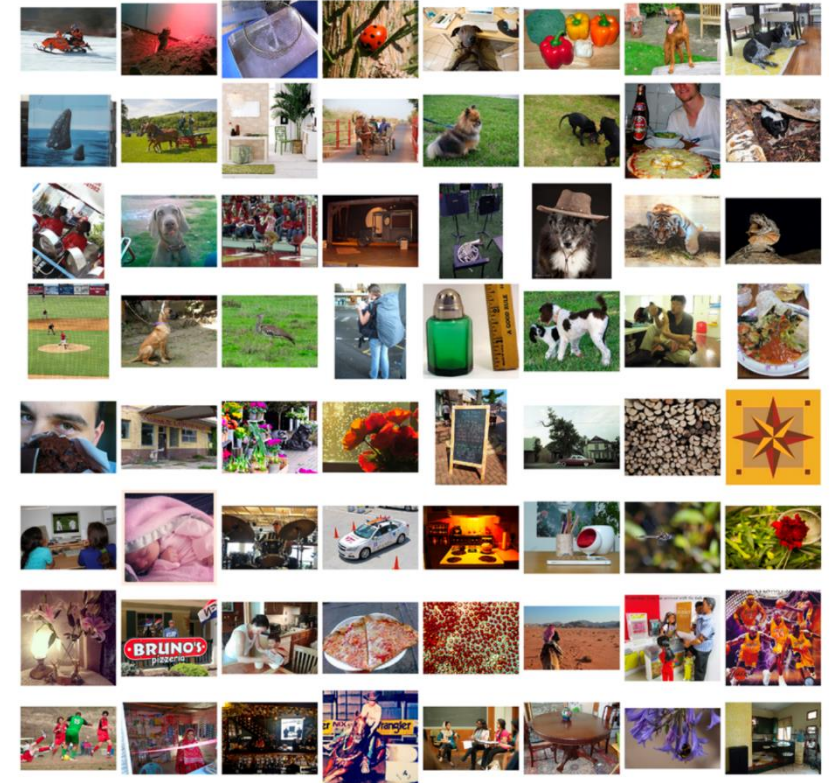
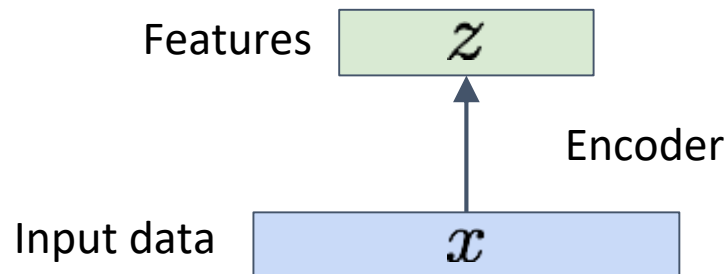
Unsupervised Learning

- 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

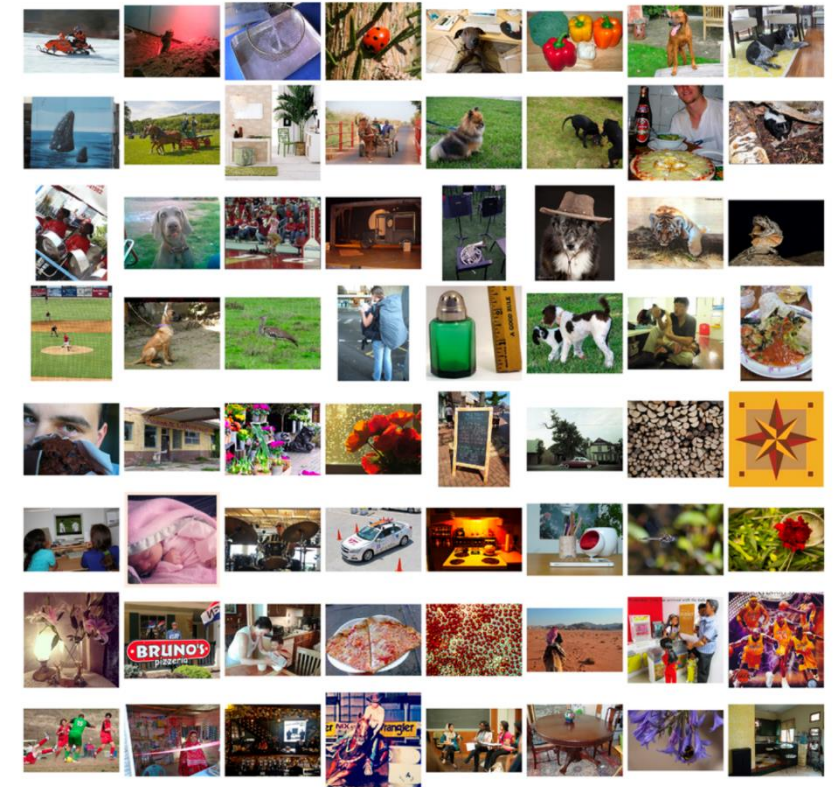
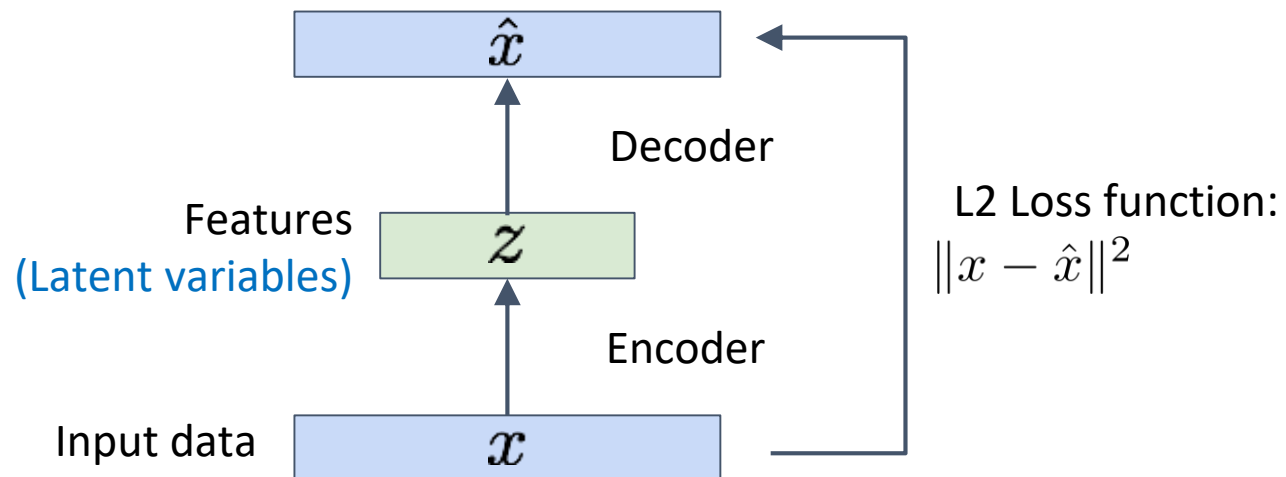
- These are good for classification, clustering, exploration, generation, ...
- We have no ground truth for them



Autoencoders

Goal: Meaningful features that capture the main factors of variation

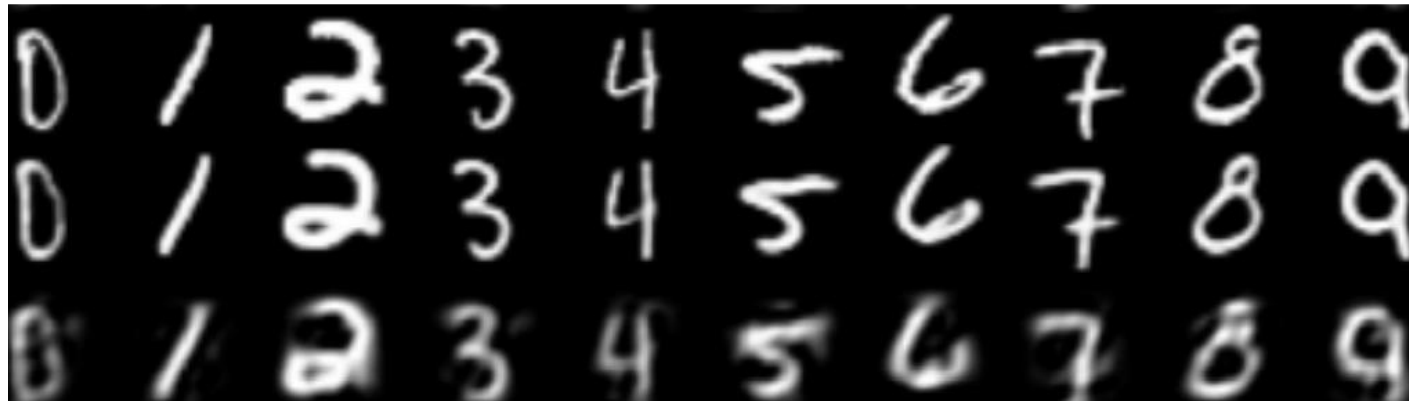
Features that can be used to reconstruct the image



Autoencoders

Linear Transformation for Encoder and Decoder give result close to PCA

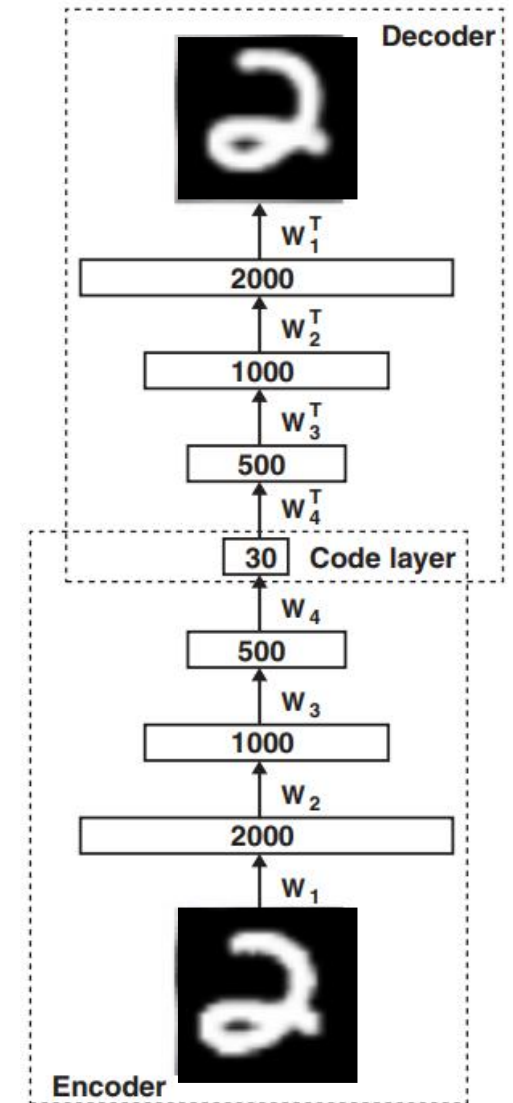
Deeper networks give better reconstructions, since basis can be non-linear



Original

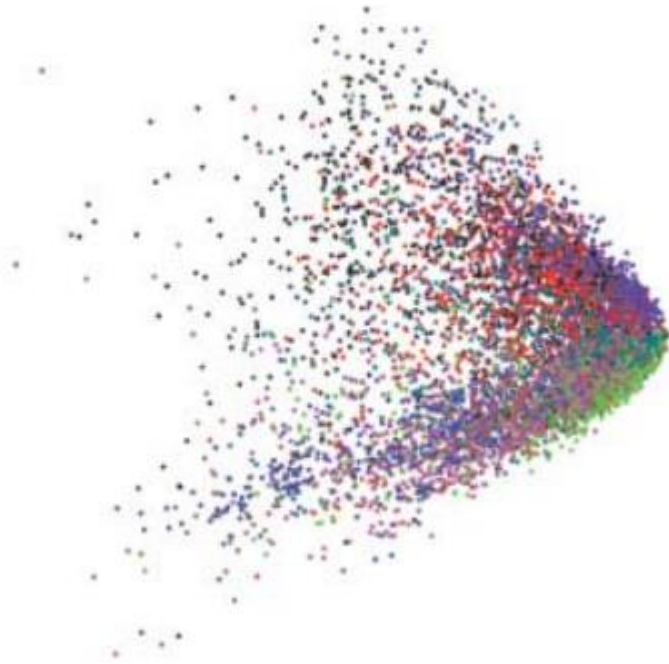
Autoencoder

PCA

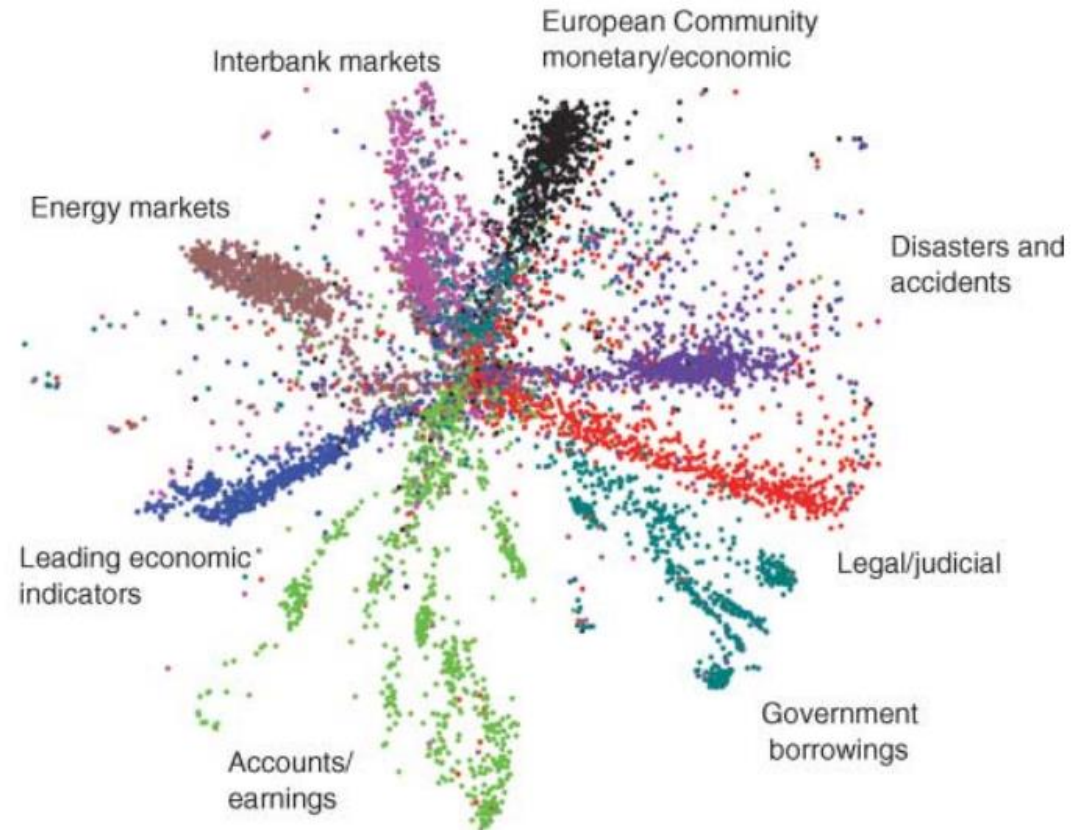


Example: Document Word Prob. → 2D Code

LSA (based on PCA)



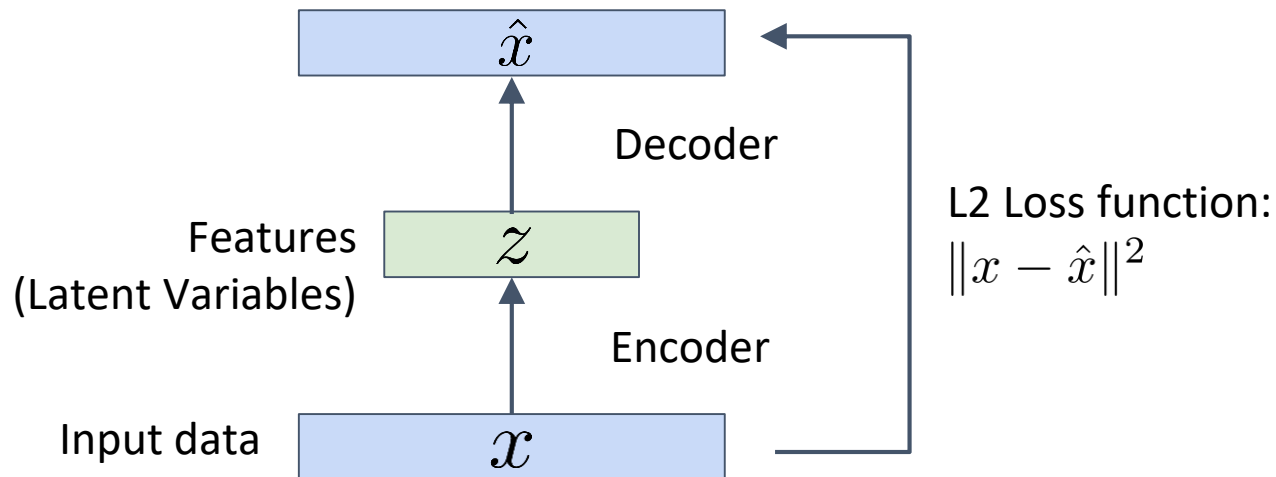
Autoencoder



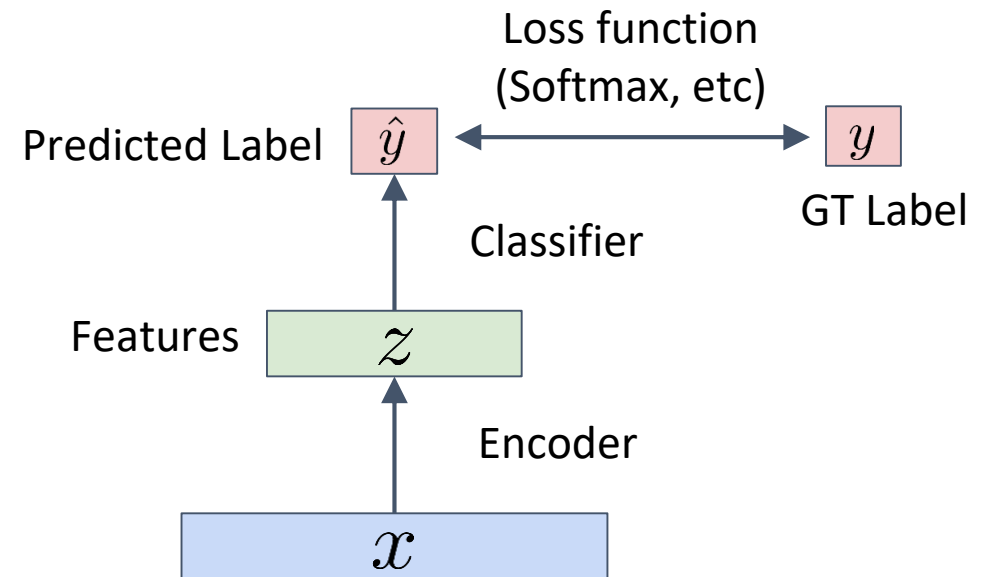
Example: Semi-Supervised Classification

- Many images, but few ground truth labels

start unsupervised
train autoencoder on many images



supervised fine-tuning
train classification network on labeled images

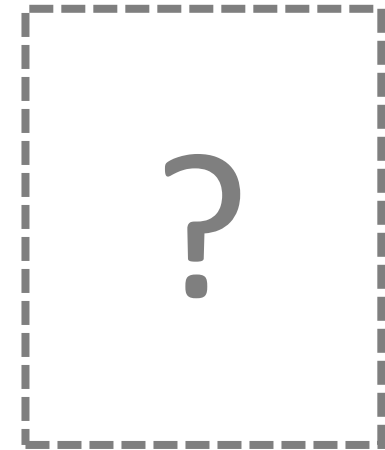


Code example

```
Autoencoder  
(autoencoder.ipynb)
```

Generative Models

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



Dataset

Generated

Generative Models

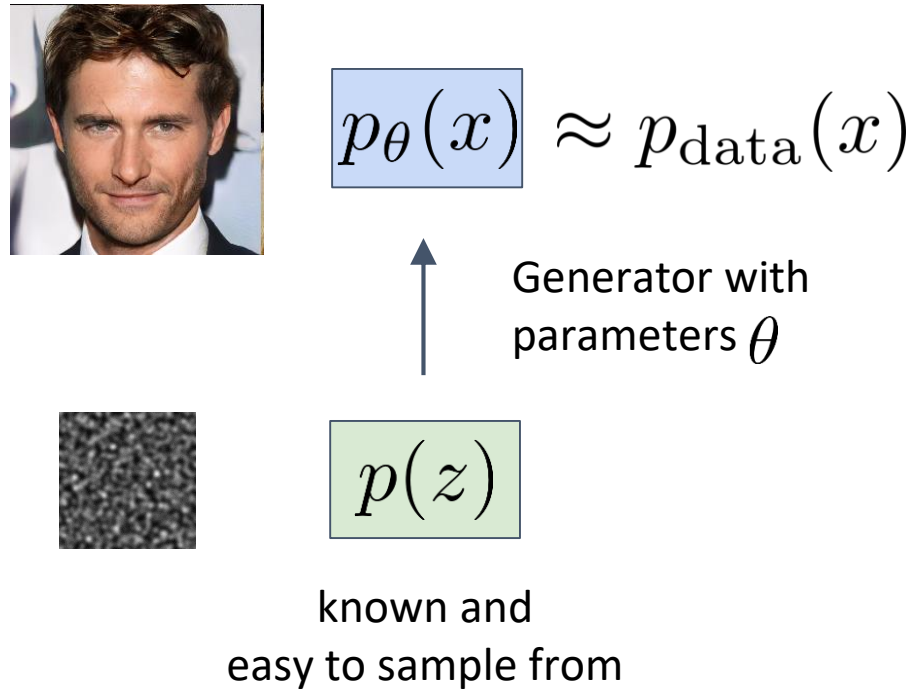
- Assumption: the dataset are samples from an unknown distribution $p_{\text{data}}(x)$
- Goal: create a new sample from $p_{\text{data}}(x)$ that is not in the dataset



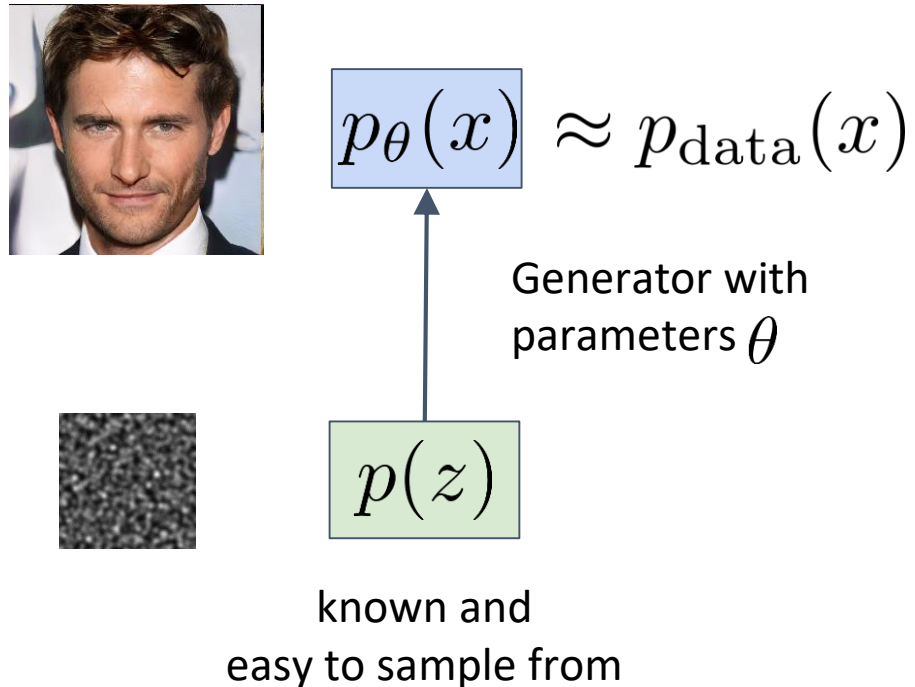
Dataset

Generated

Generative Models



Generative Models



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

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

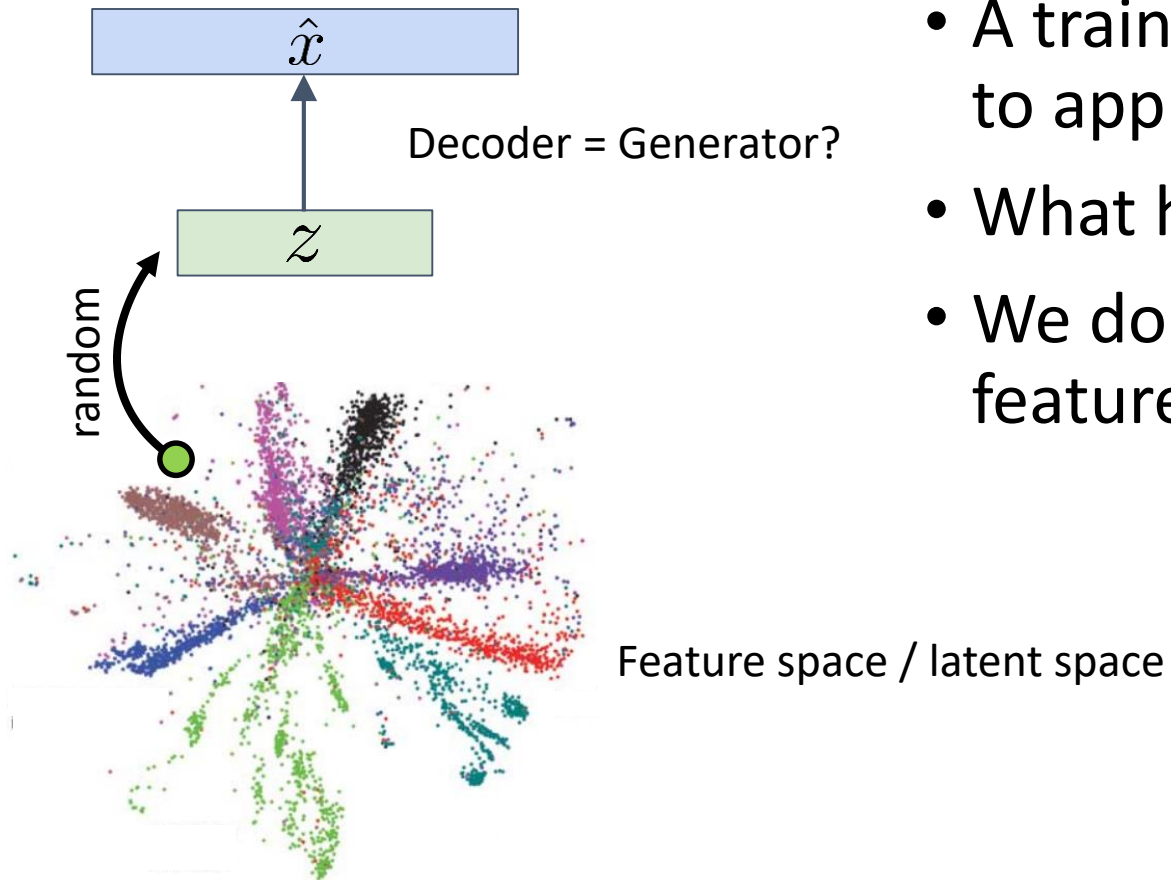
Variational Autoencoders (VAEs)

- 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 z to approximate samples from $p_{\text{data}}(x)$
- What happens if we pick a random z ?
- We do not know the distribution $p(z)$ of features that decode to likely samples

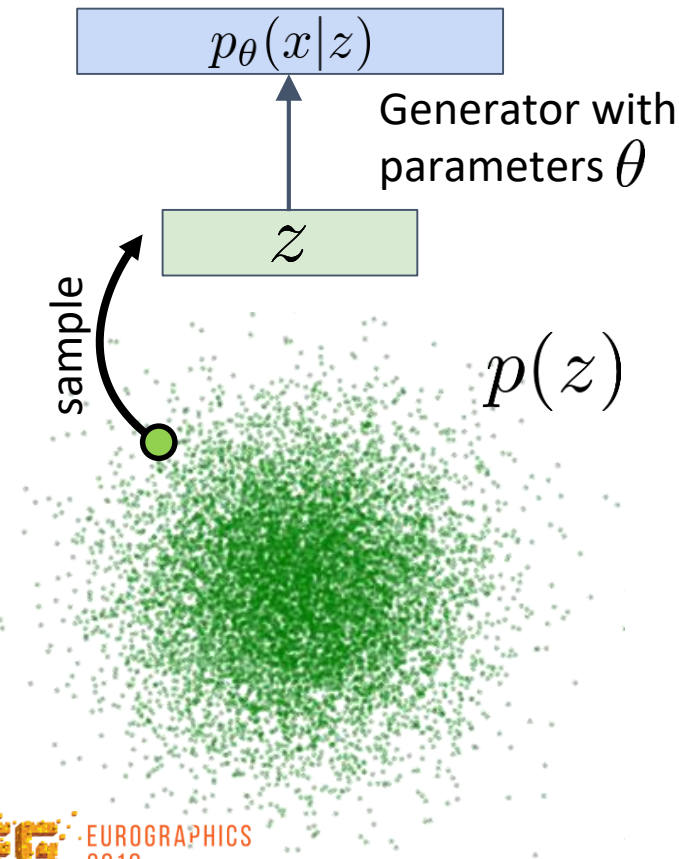
Variational Autoencoders (VAEs)

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

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

- Train the generator to maximize the likelihood of the data in $p_\theta(x)$:

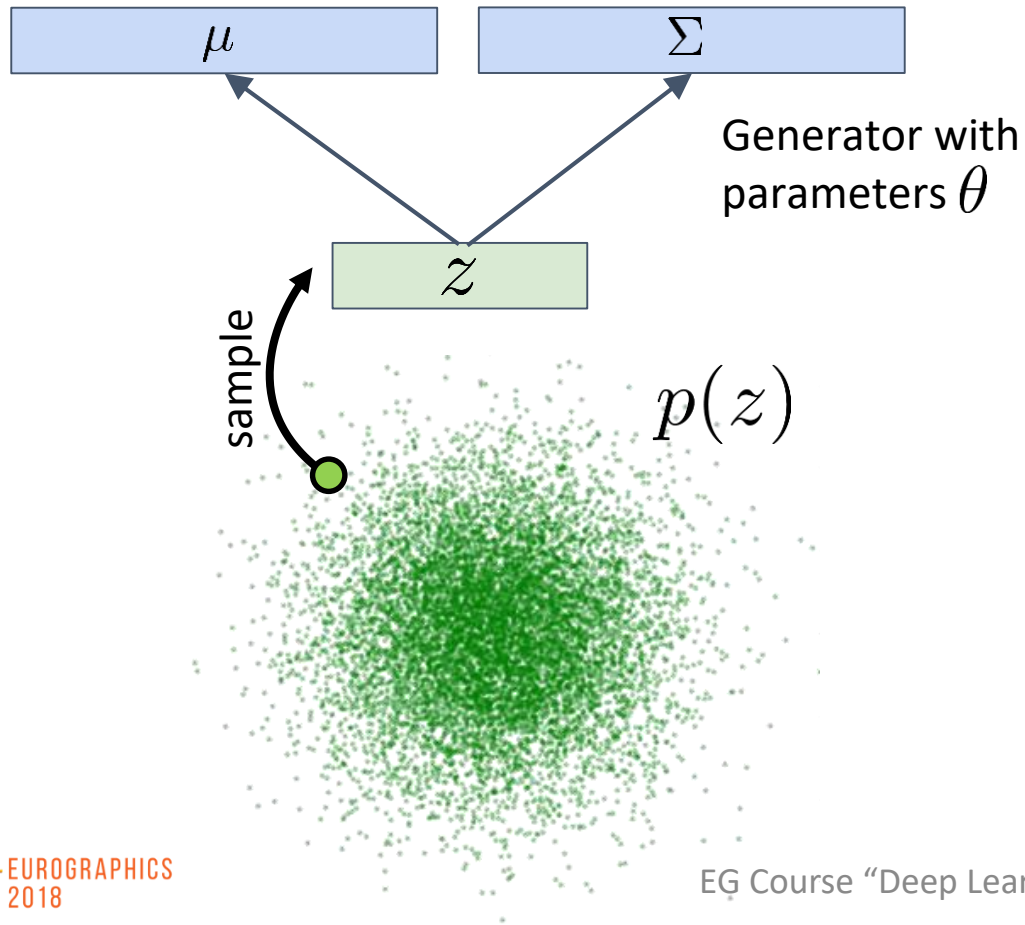
$$\max_{\theta} \sum_{x \in \text{data}} \log p_\theta(x)$$



Outputting a Distribution

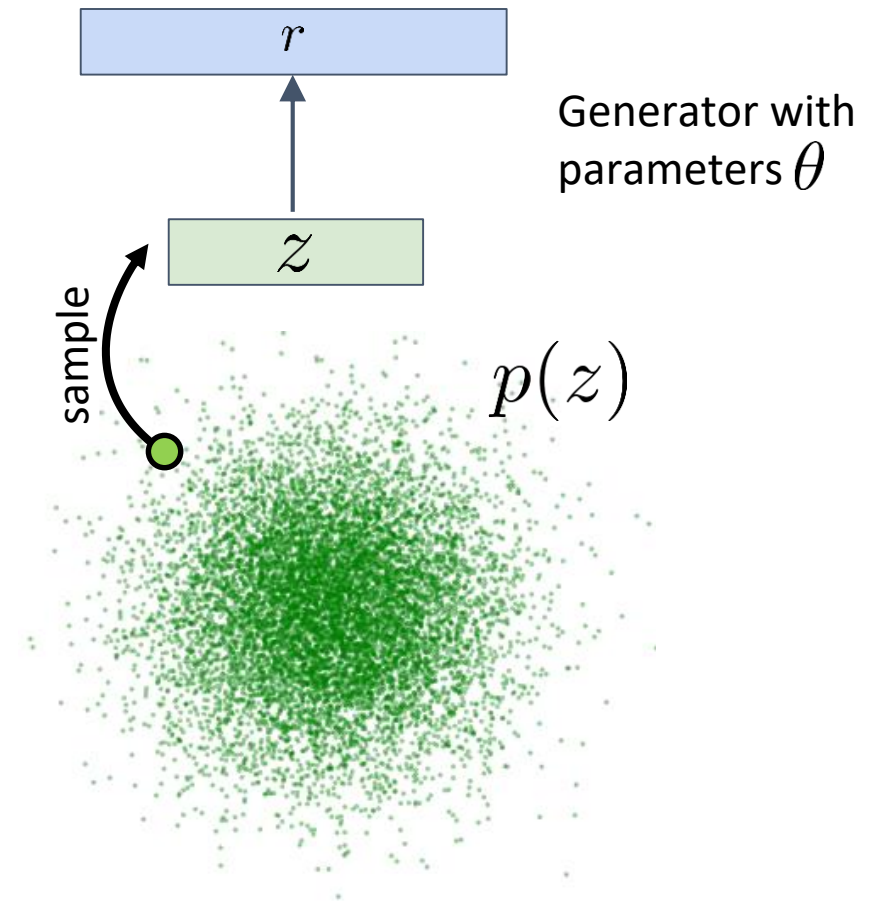
Normal distribution

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



Bernoulli distribution

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



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

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

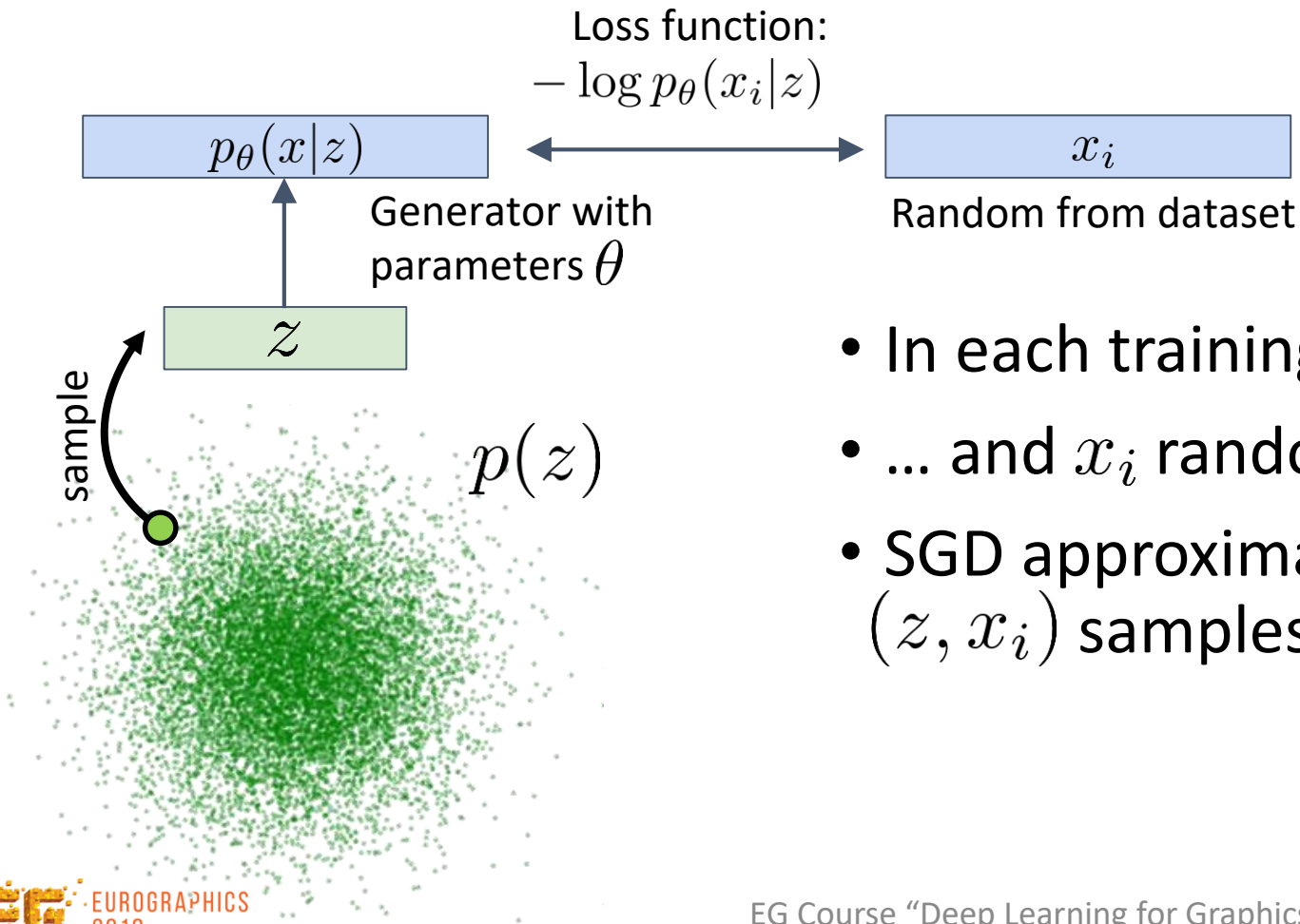
$$\theta^* = \arg \max_{\theta} \sum_{x \in \text{data}} \log \int p_{\theta}(x|z) p(z) dz$$

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

- SGD approximates the expected values over (z, x_i) samples
- In each training iteration, sample z from $p(z)$...
- ... and x_i randomly from the dataset, and maximize:

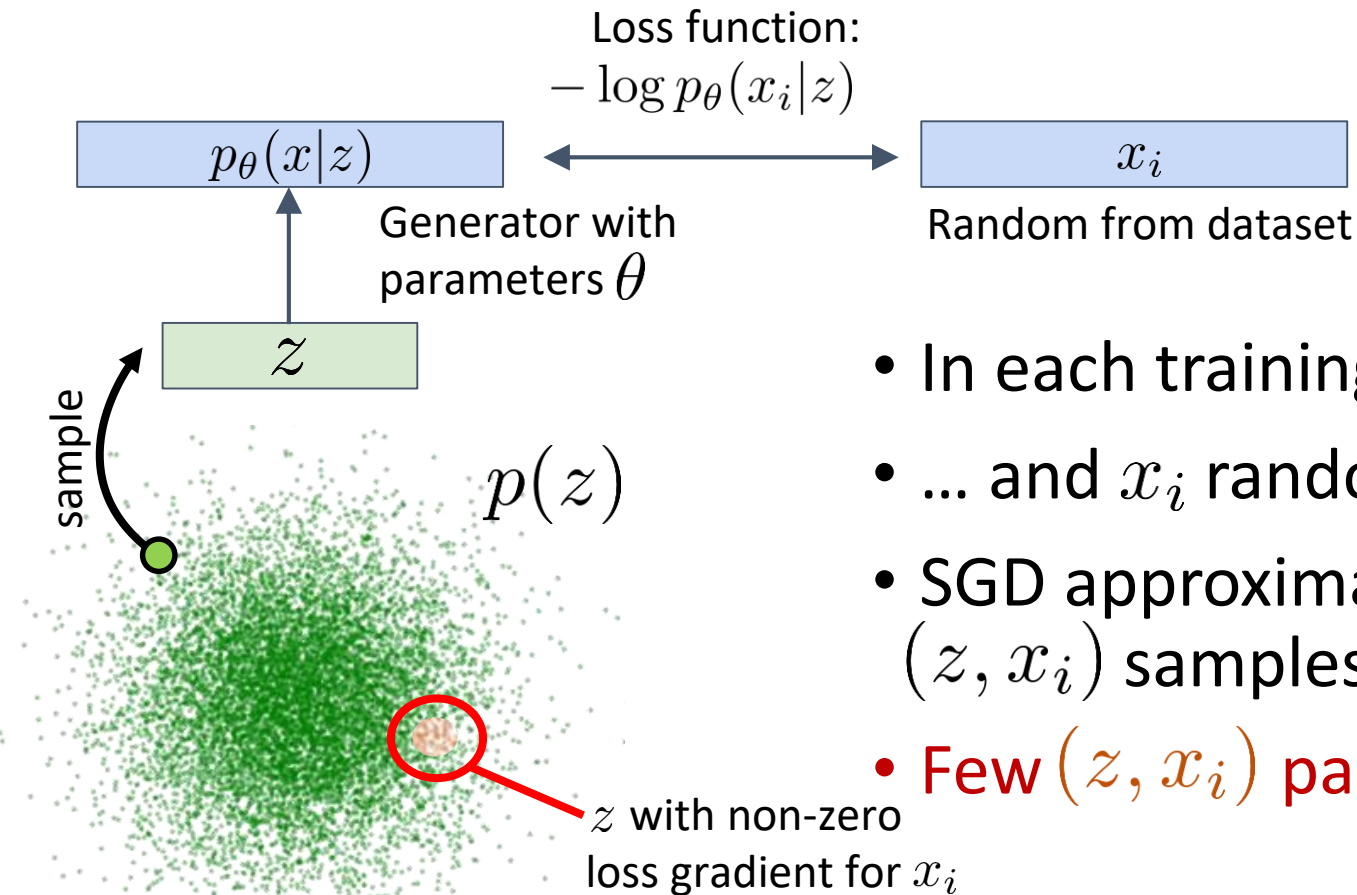
$$\max_{\theta} \log p_{\theta}(x_i|z)$$

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



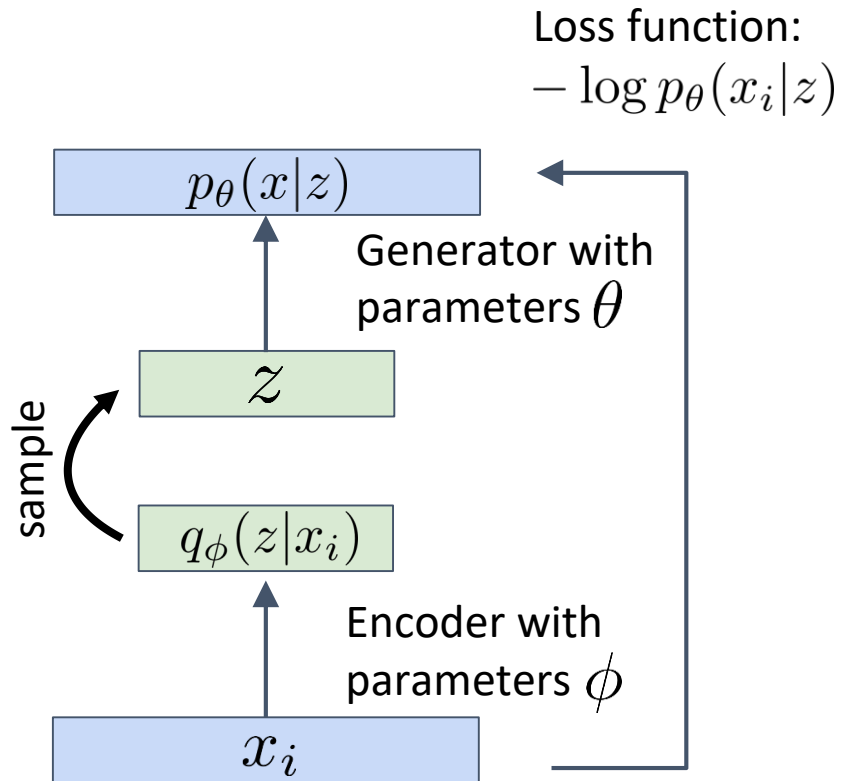
- In each training iteration, sample z from $p(z)$...
- ... and x_i randomly from the dataset
- SGD approximates the expected values over (z, x_i) samples

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



- In each training iteration, sample z from $p(z)$...
- ... and x_i randomly from the dataset
- SGD approximates the expected values over (z, x_i) samples
- **Few (z, x_i) pairs have non-zero gradients**

Variational Autoencoders (VAEs): The Encoder



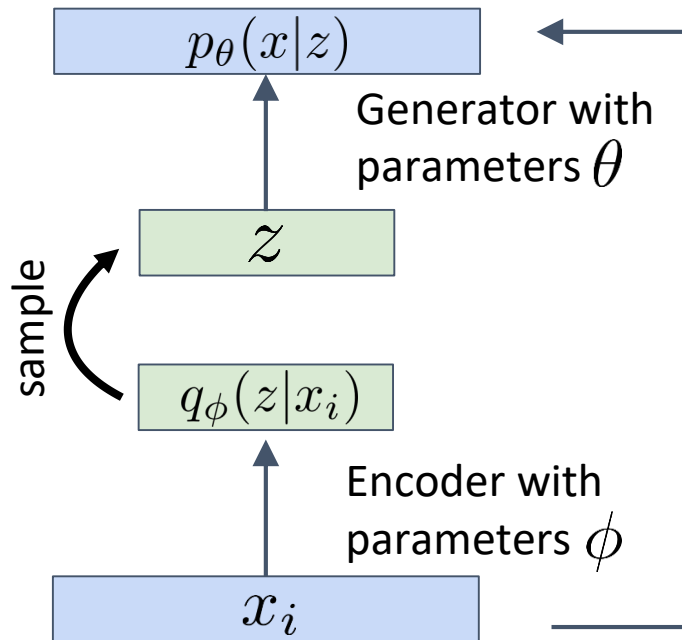
$$p_{\theta}(x) = \int p_{\theta}(x|z) p(z) dz$$

- During training, another network can guess a good z for a given x_i
- $q_{\phi}(z|x_i)$ should be much smaller than $p(z)$
- This also gives us the data point x_i

Variational Autoencoders (VAEs): The Encoder

Loss function:

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



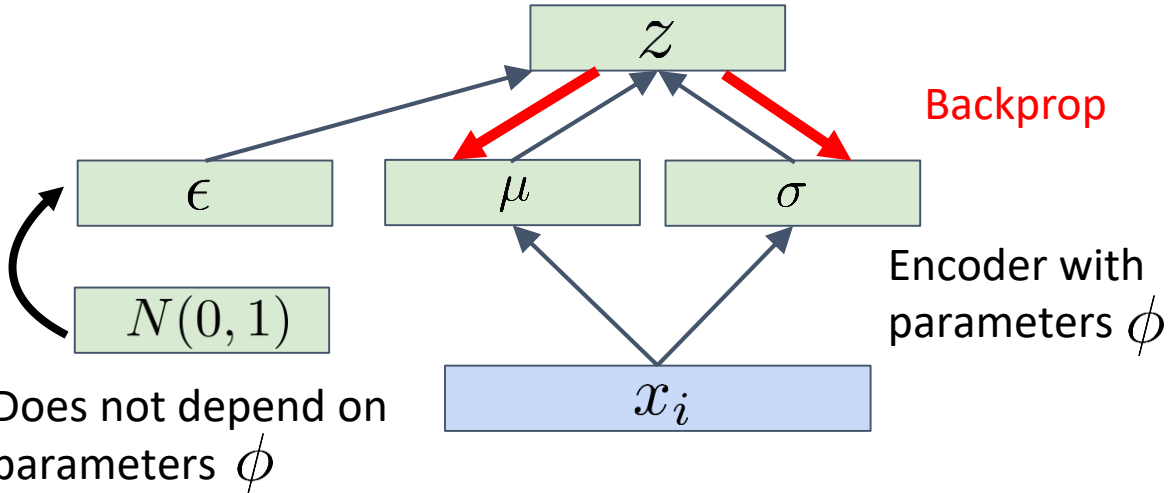
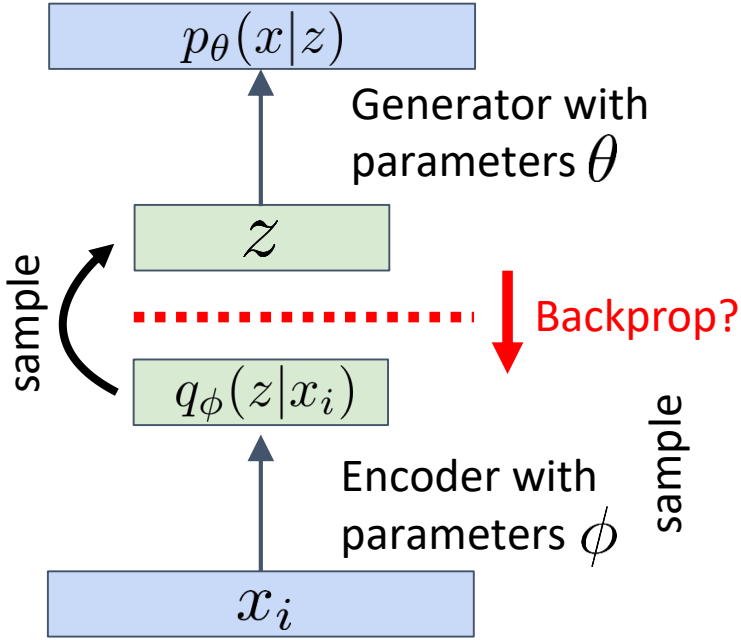
- Can we still easily sample a new z ?
- Need to make sure $q_{\phi}(z|x_i)$ approximates $p(z)$
- Regularize with **KL-divergence**
- Negative loss can be shown to be a lower bound for the likelihood, and equivalent if $q_{\phi}(z|x) = p_{\theta}(z|x)$

Reparameterization Trick

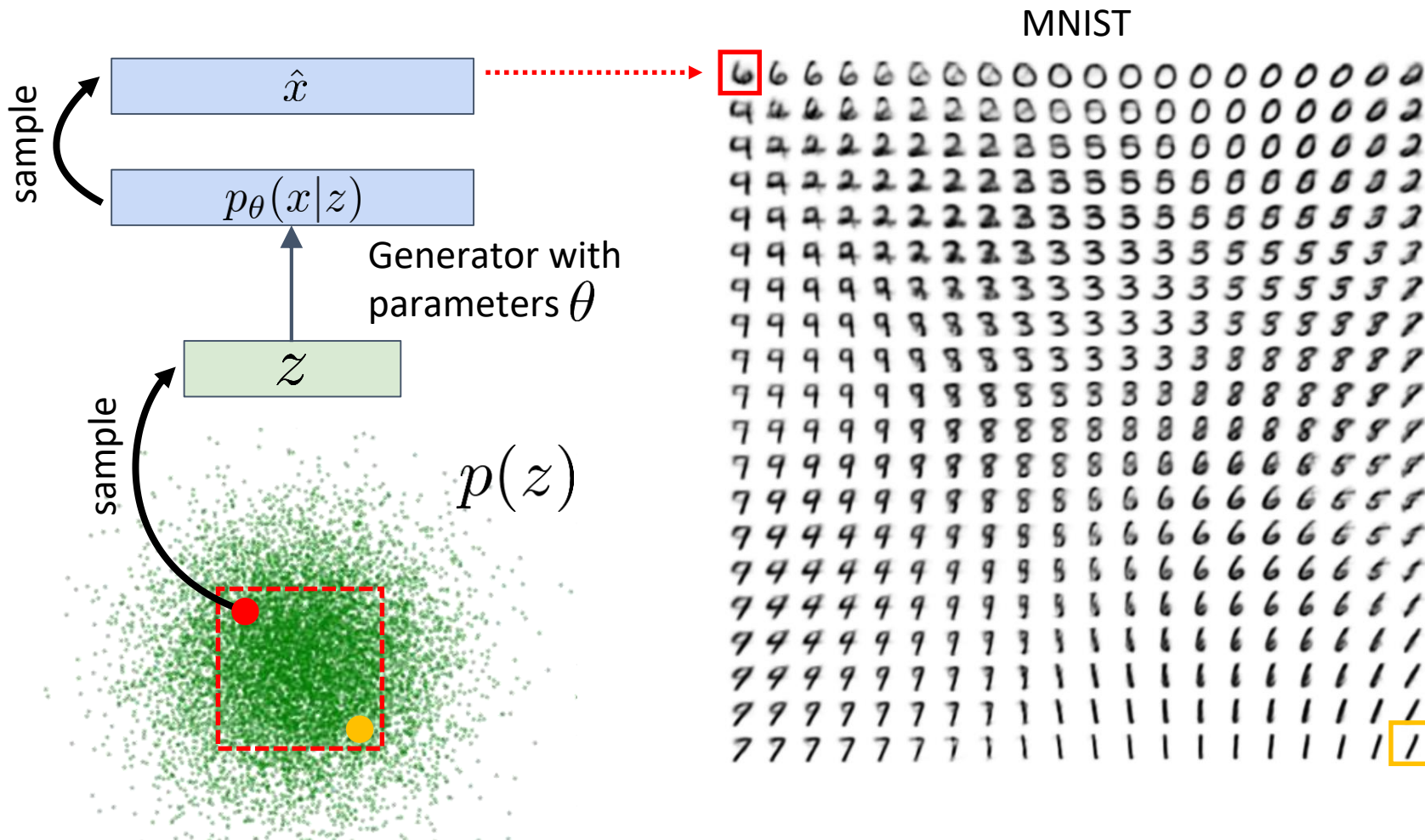
Example when $q_\phi(z|x_i) = N(z; \mu(x_i), \sigma(x_i))$:

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

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



Generating Data



MNIST

Frey Faces



Demos

VAE on MNIST

http://dpkingma.com/sgvb_mnist_demo/demo.html

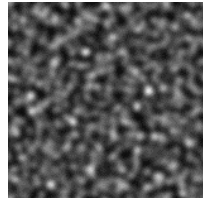
VAE on Faces

http://vdumoulin.github.io/morphing_faces/online_demo.html

Code example

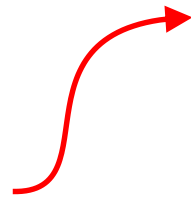
Variational Autoencoder
(variational_autoencoder.ipynb)

Generative Adversarial Networks

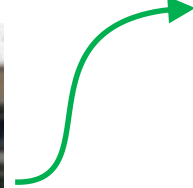


Player 1: generator

Scores if discriminator
can't distinguish output
from real image



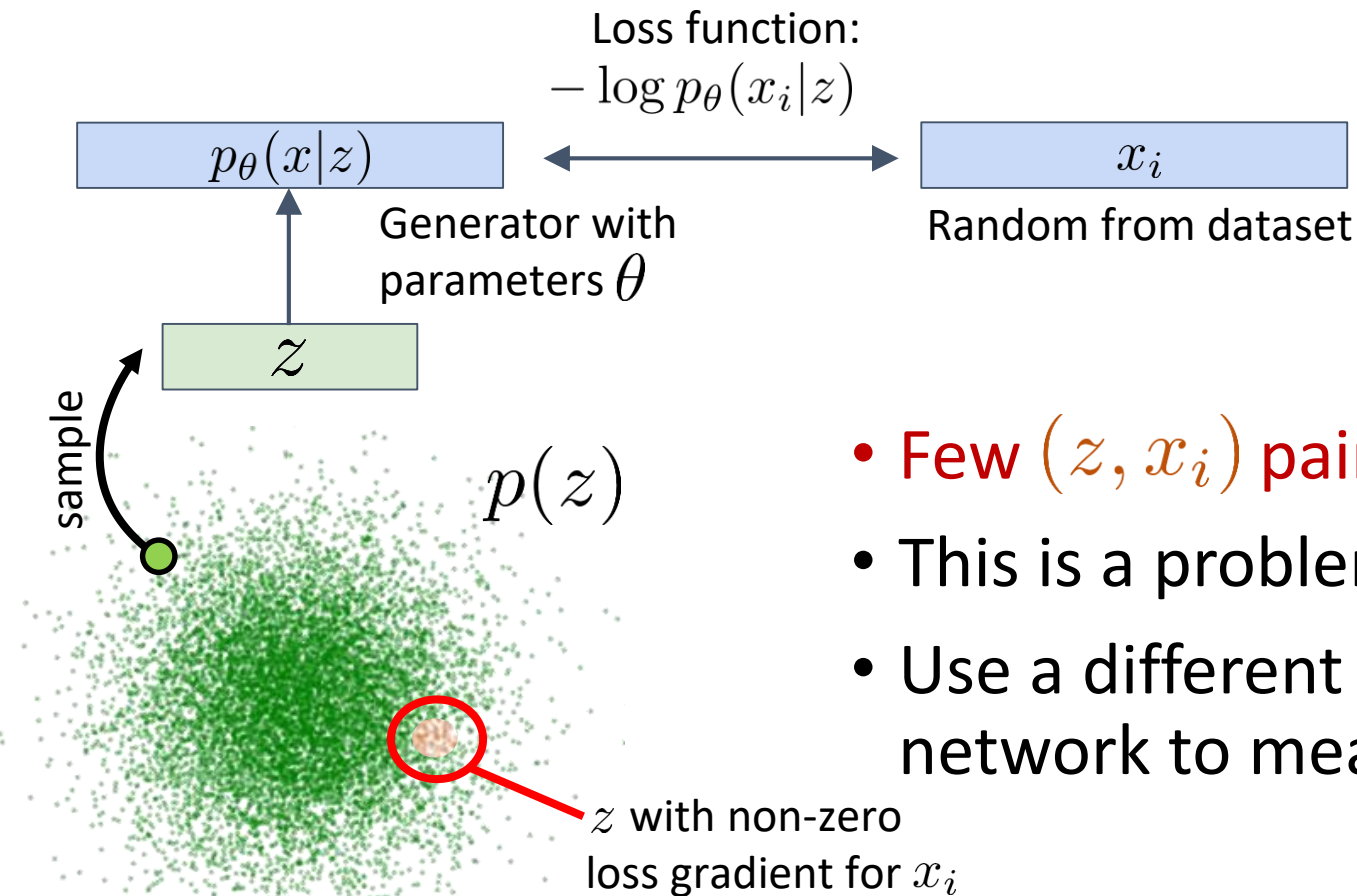
from dataset



Player 2: discriminator → real/fake

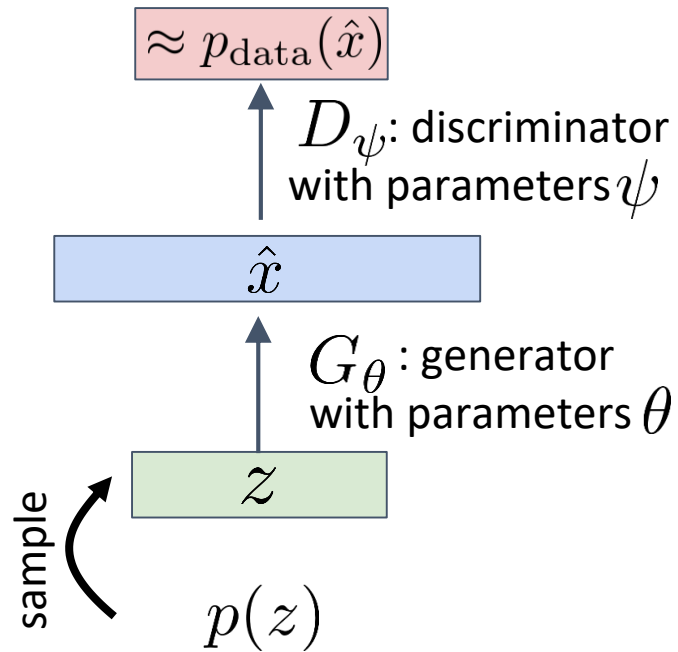
Scores if it can distinguish
between real and fake

Naïve Sampling Revisited

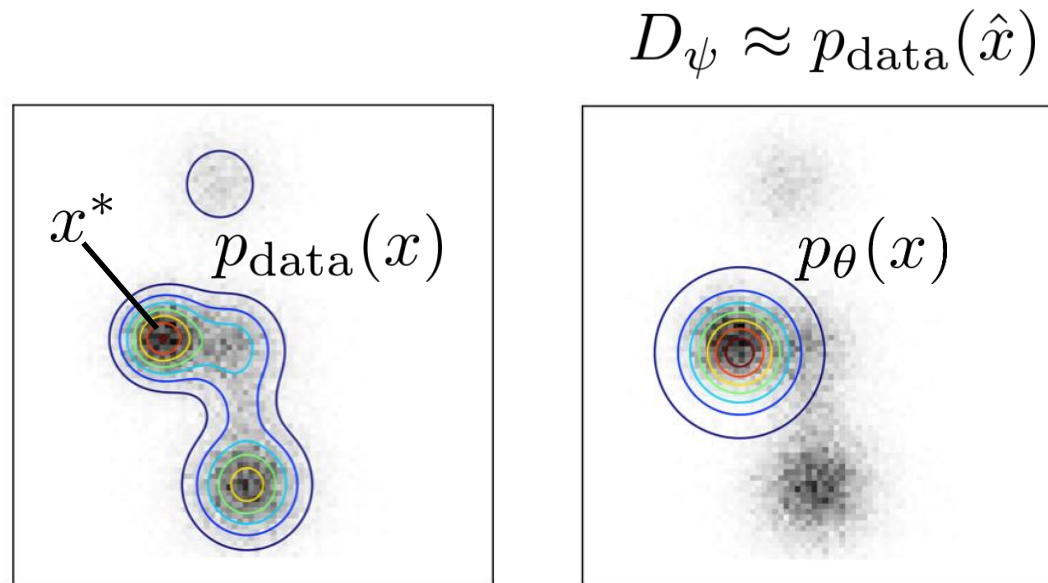


- Few (z, x_i) pairs have non-zero gradients
- This is a problem of the maximum likelihood
- Use a different loss: Train a discriminator network to measure similarity $p_{\theta}(x) \approx p_{\text{data}}(x)$

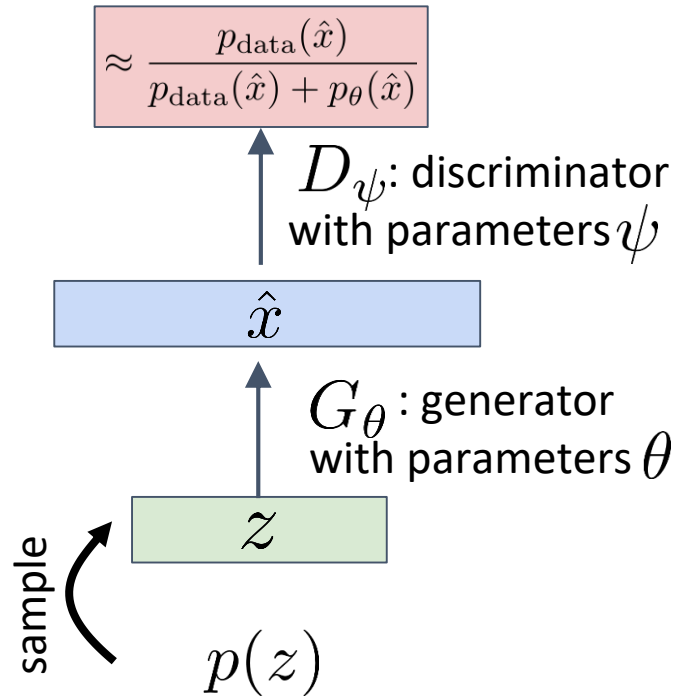
Why Adversarial?



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



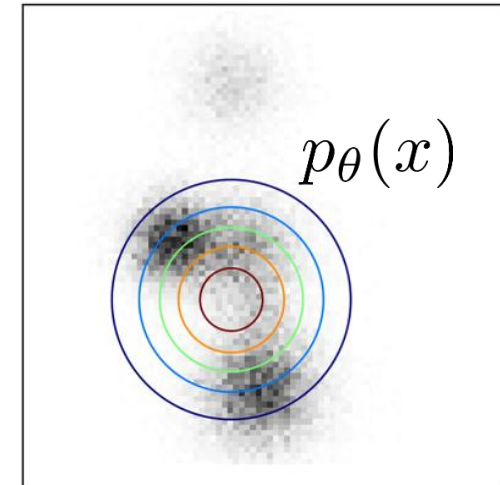
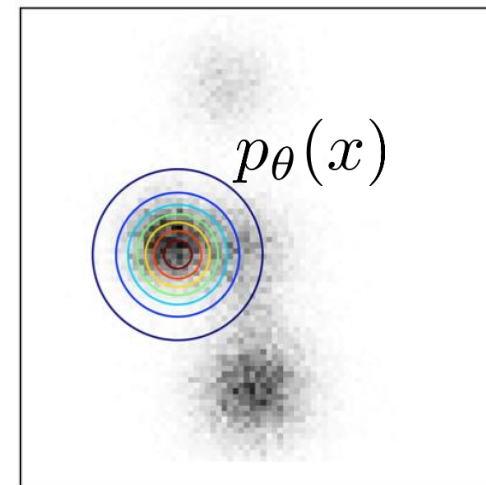
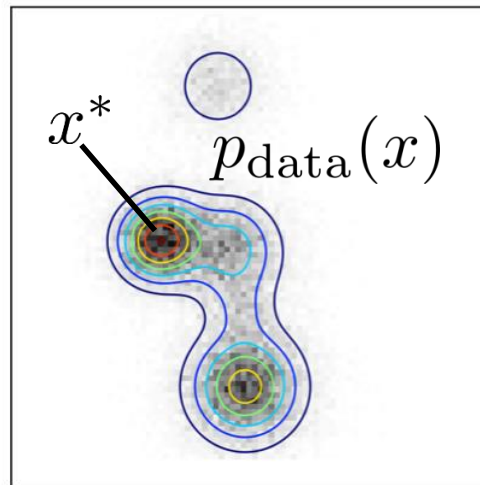
Why Adversarial?



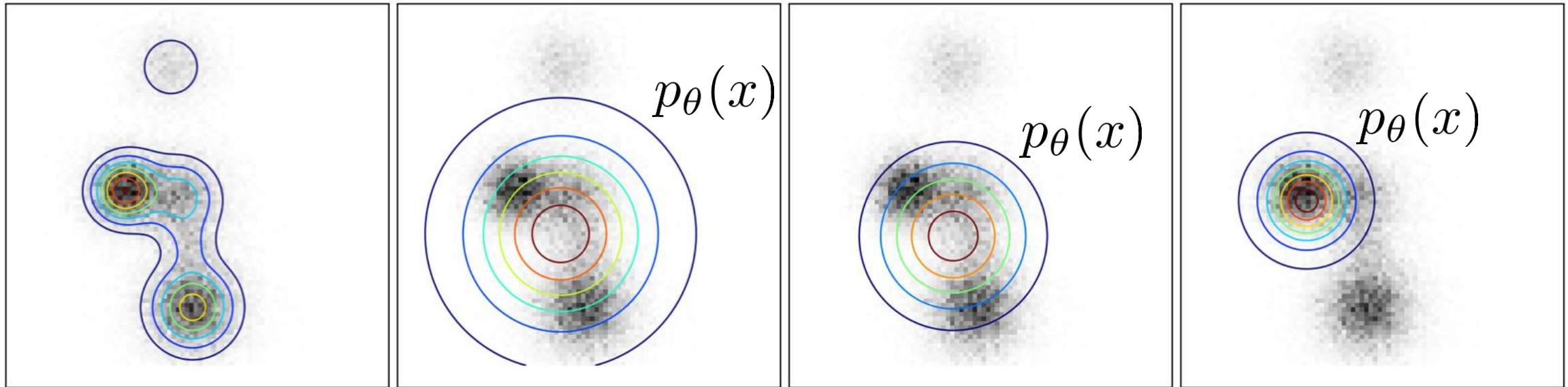
- For GANs, the discriminator instead approximates:

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

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



Why Adversarial?



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

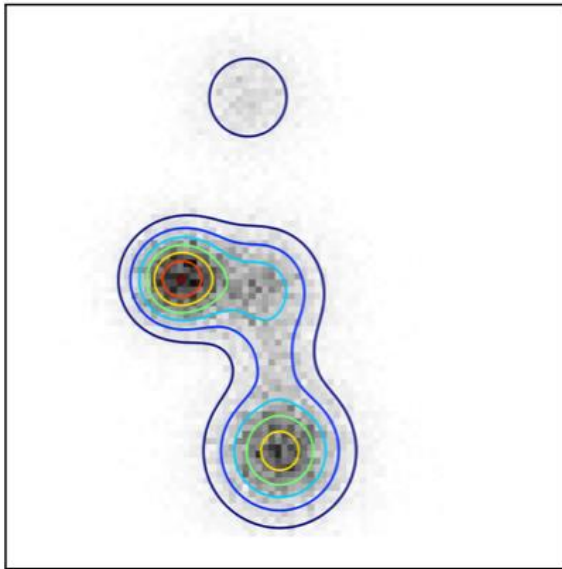
VAEs:
Maximize likelihood of
data samples in $p_{\theta}(x)$

GANs:
Adversarial game

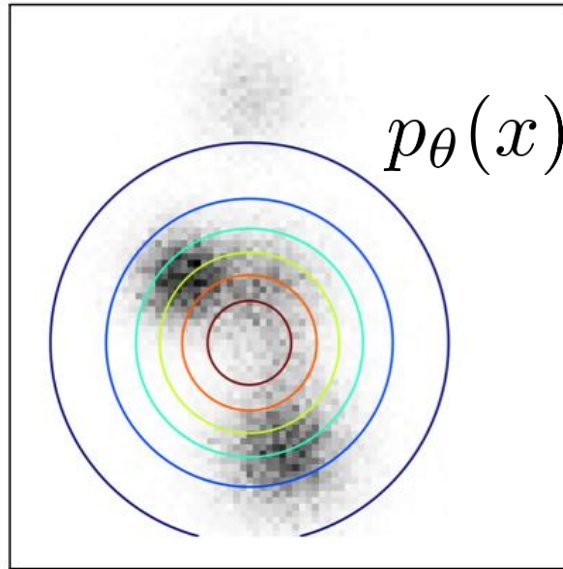
Maximize likelihood of
generator samples in
approximate $p_{\text{data}}(x)$

Why Adversarial?

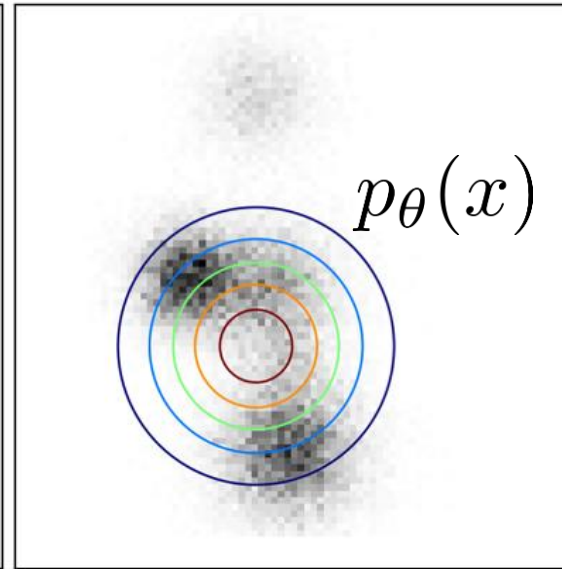
$$\approx KL(p_{\text{data}} \parallel p_{\theta}) \quad \approx JS(p_{\text{data}} \parallel p_{\theta}) \quad \approx KL(p_{\theta} \parallel p_{\text{data}})$$



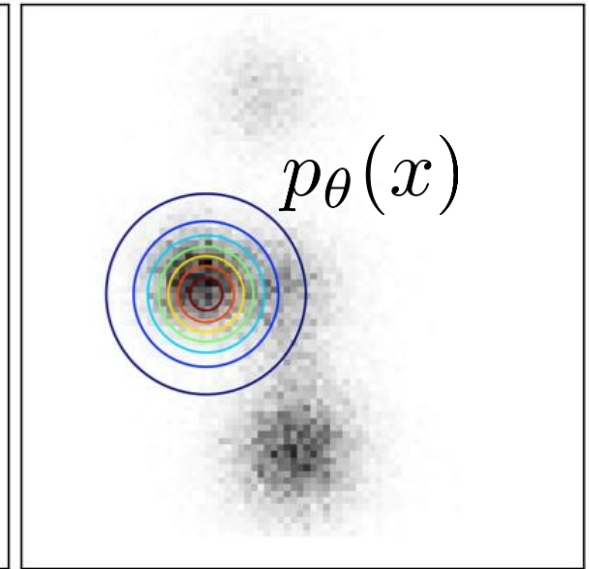
$p_{\text{data}}(x)$



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

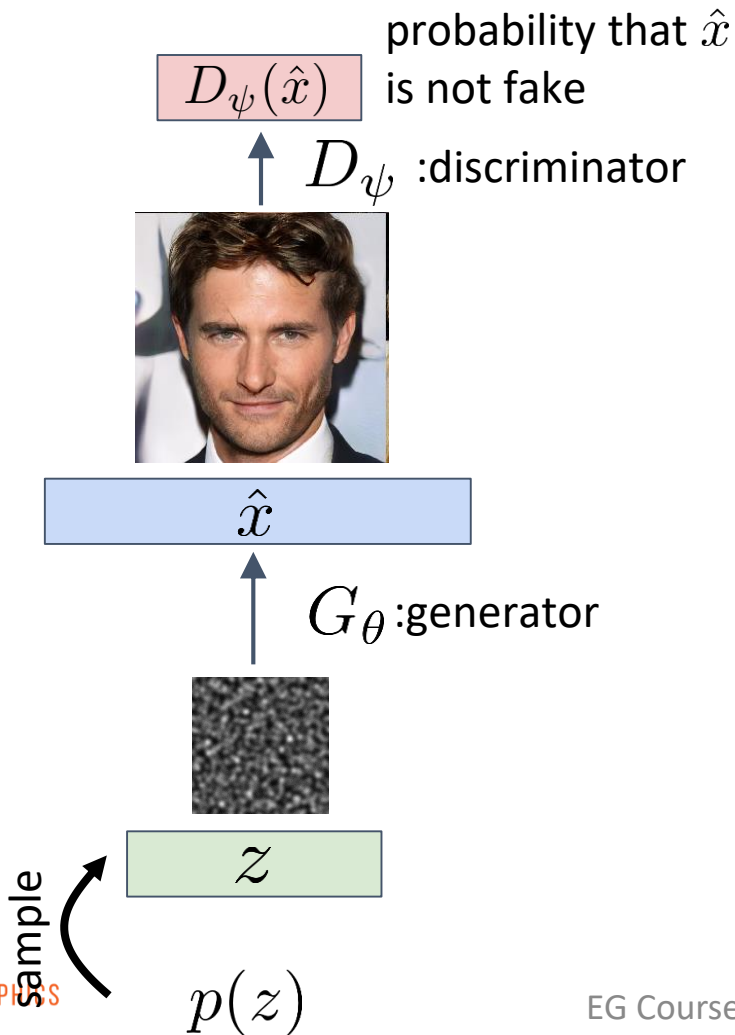


GANs:
Adversarial game



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:

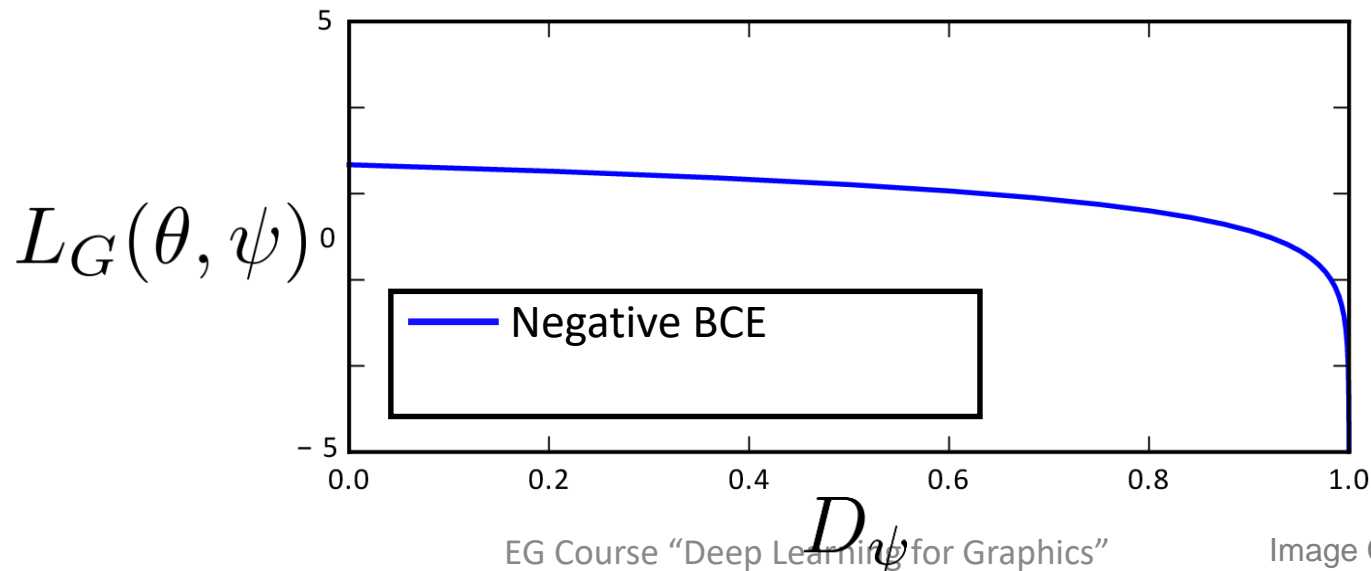
$$\max_{\theta} L(\theta, \psi)$$

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:

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



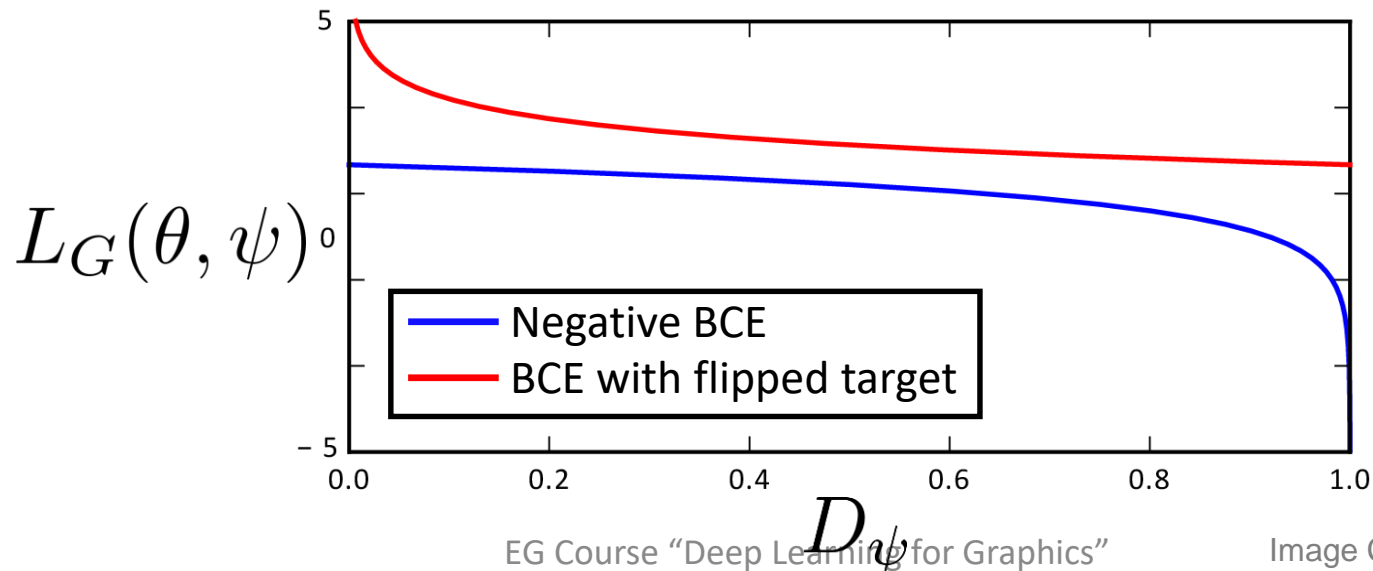
Non-saturating Heuristic

Generator loss is negative binary cross-entropy:

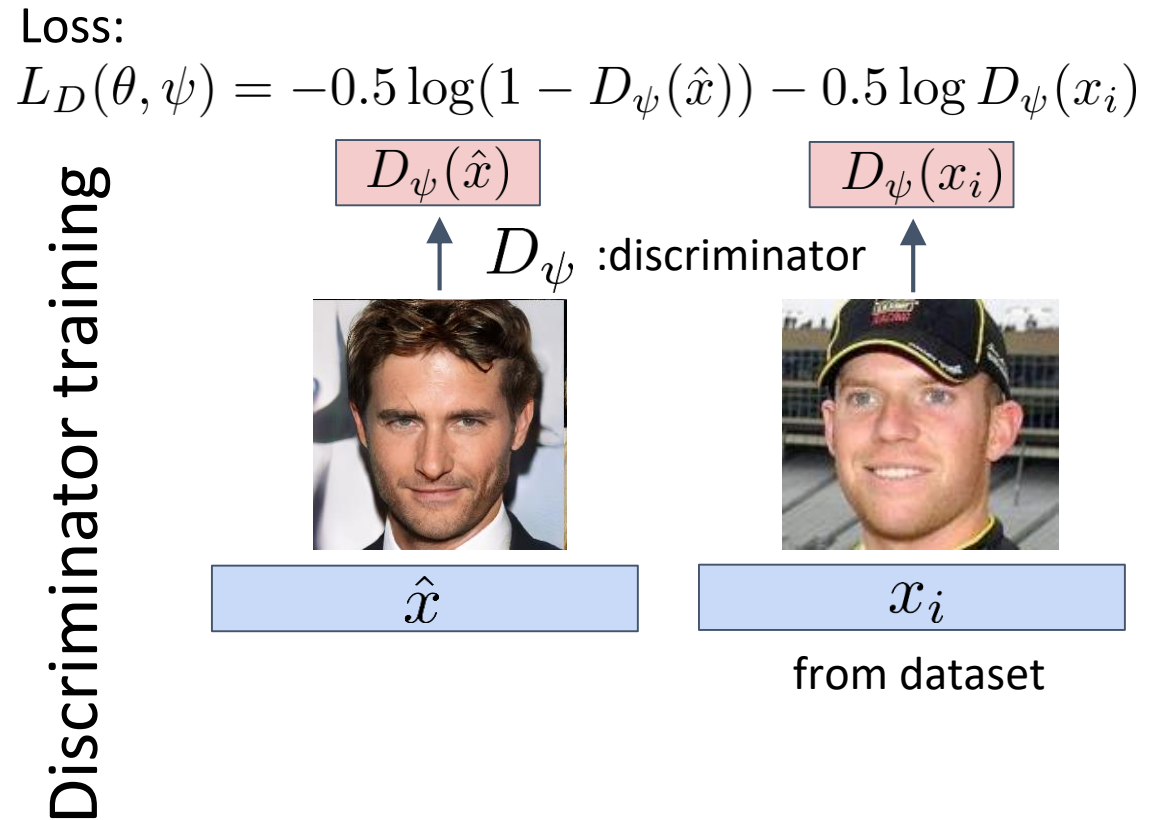
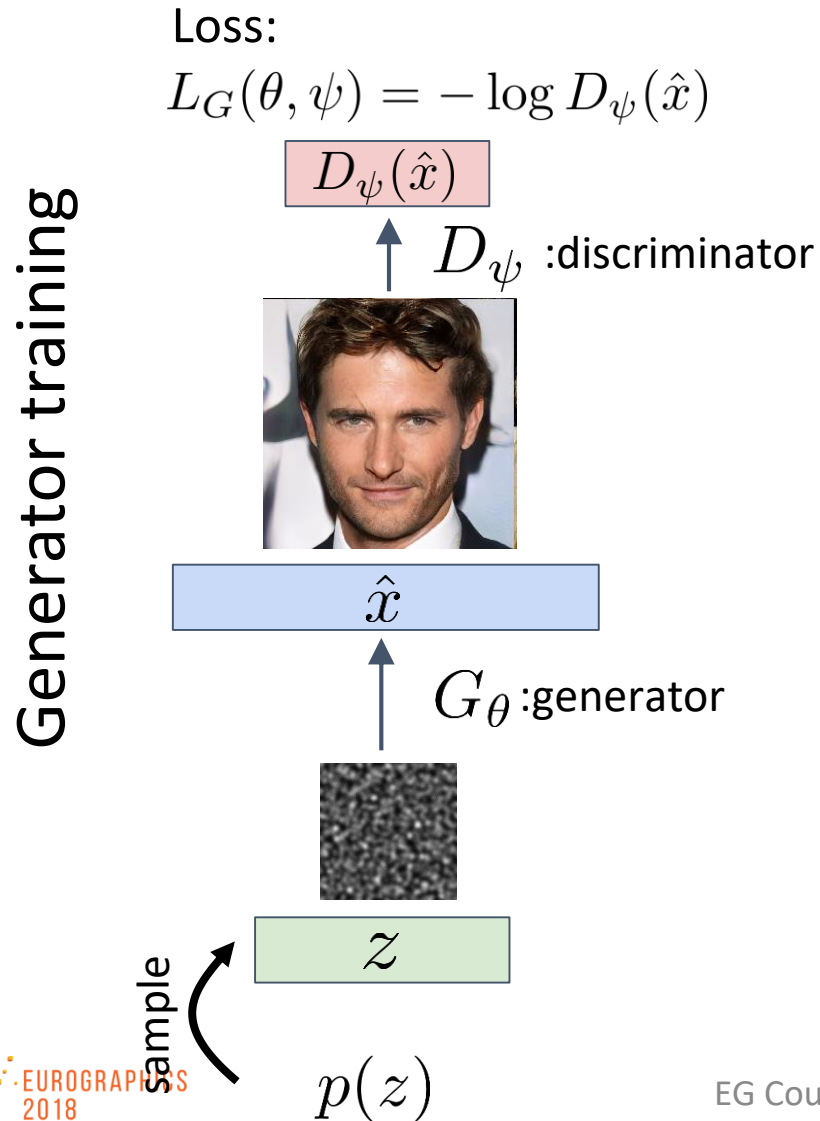
$$L_G(\theta, \psi) = 0.5 \mathbb{E}_{x \sim p_\theta} \log(1 - D_\psi(x)) \quad \text{poor convergence}$$

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

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



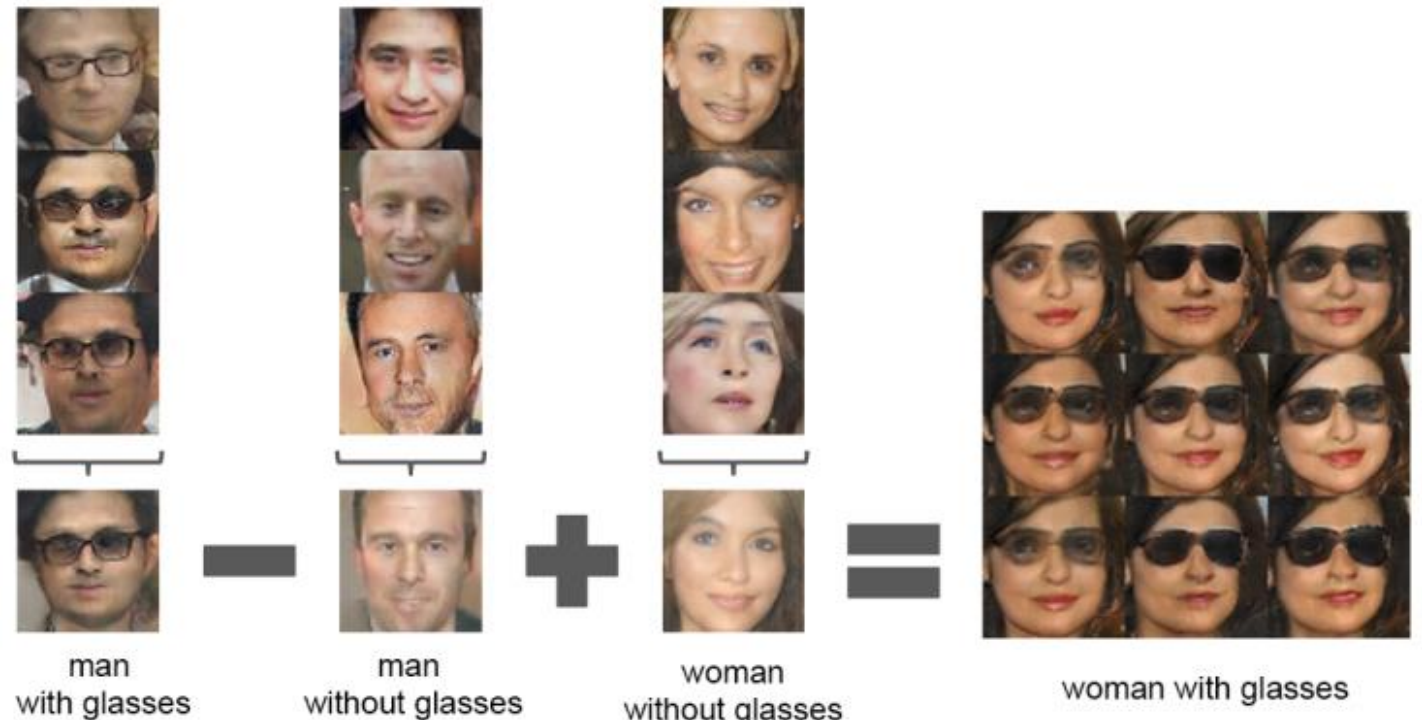
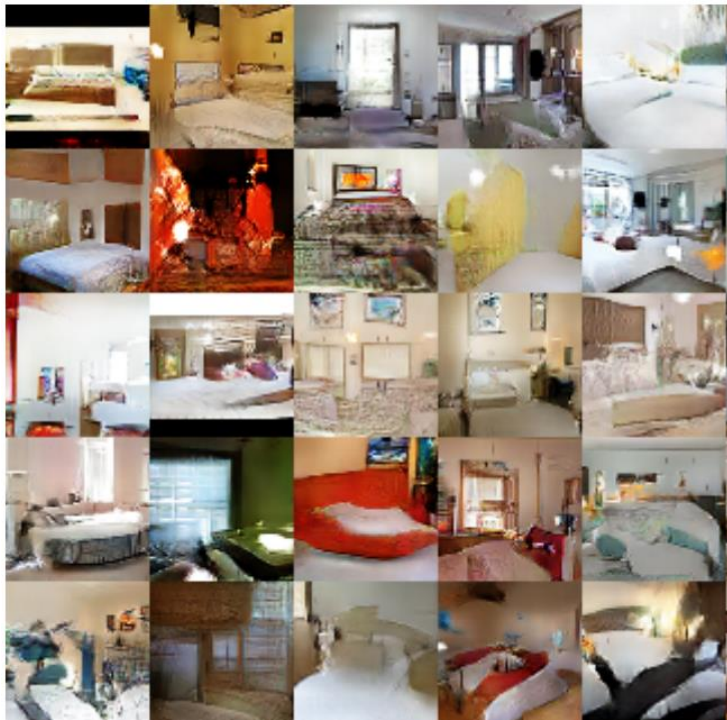
GAN Training



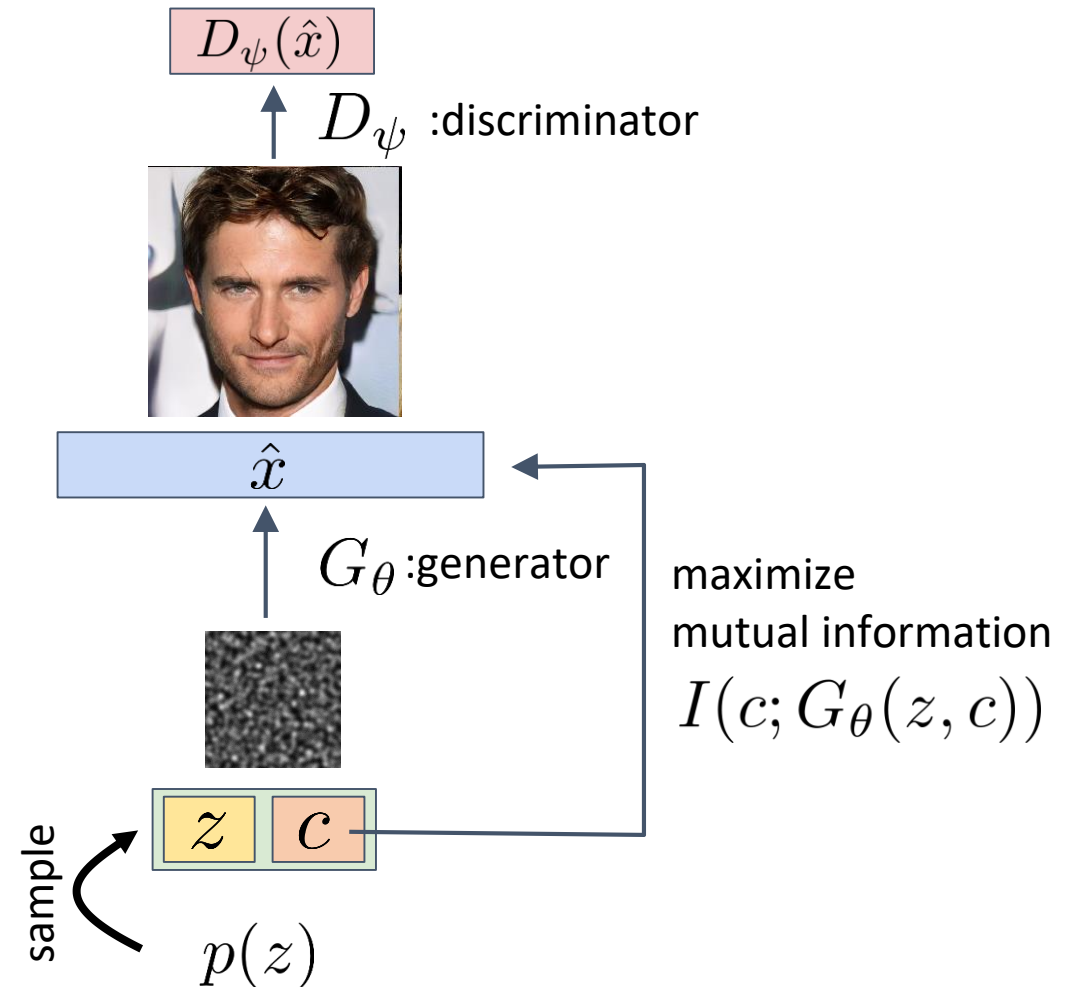
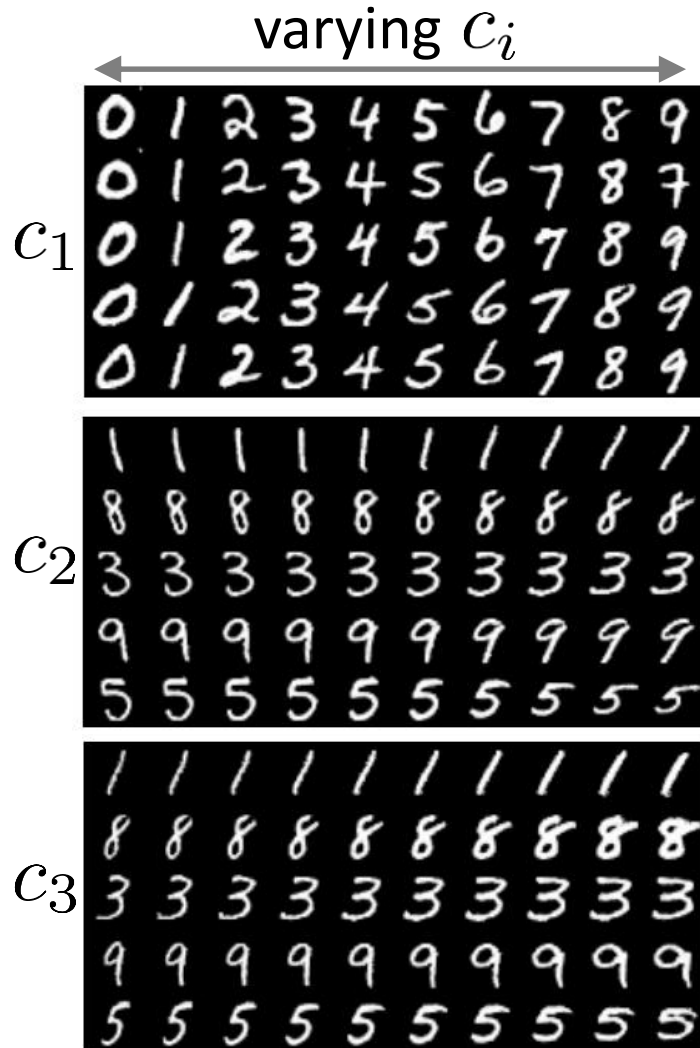
Interleave in each training step

DCGAN

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



InfoGAN

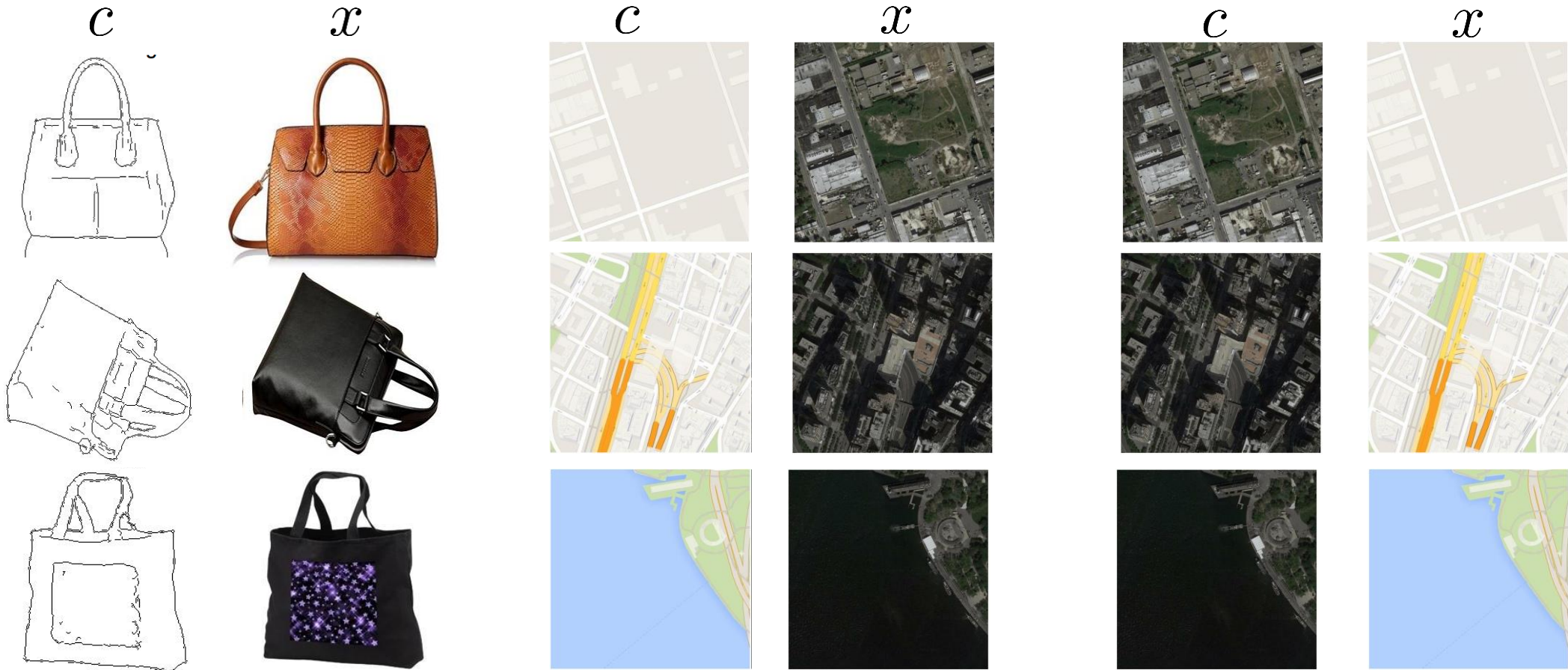


Code example

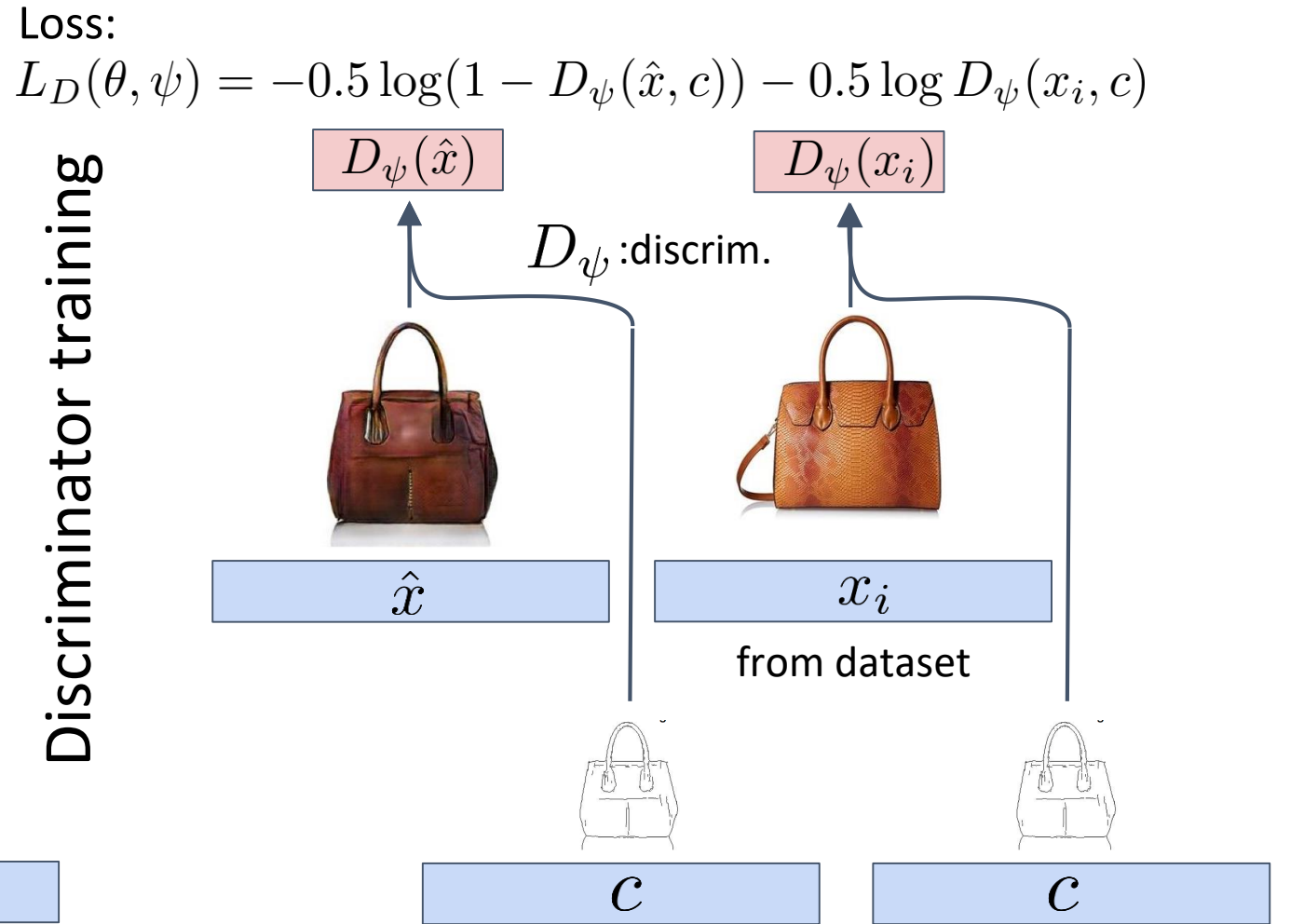
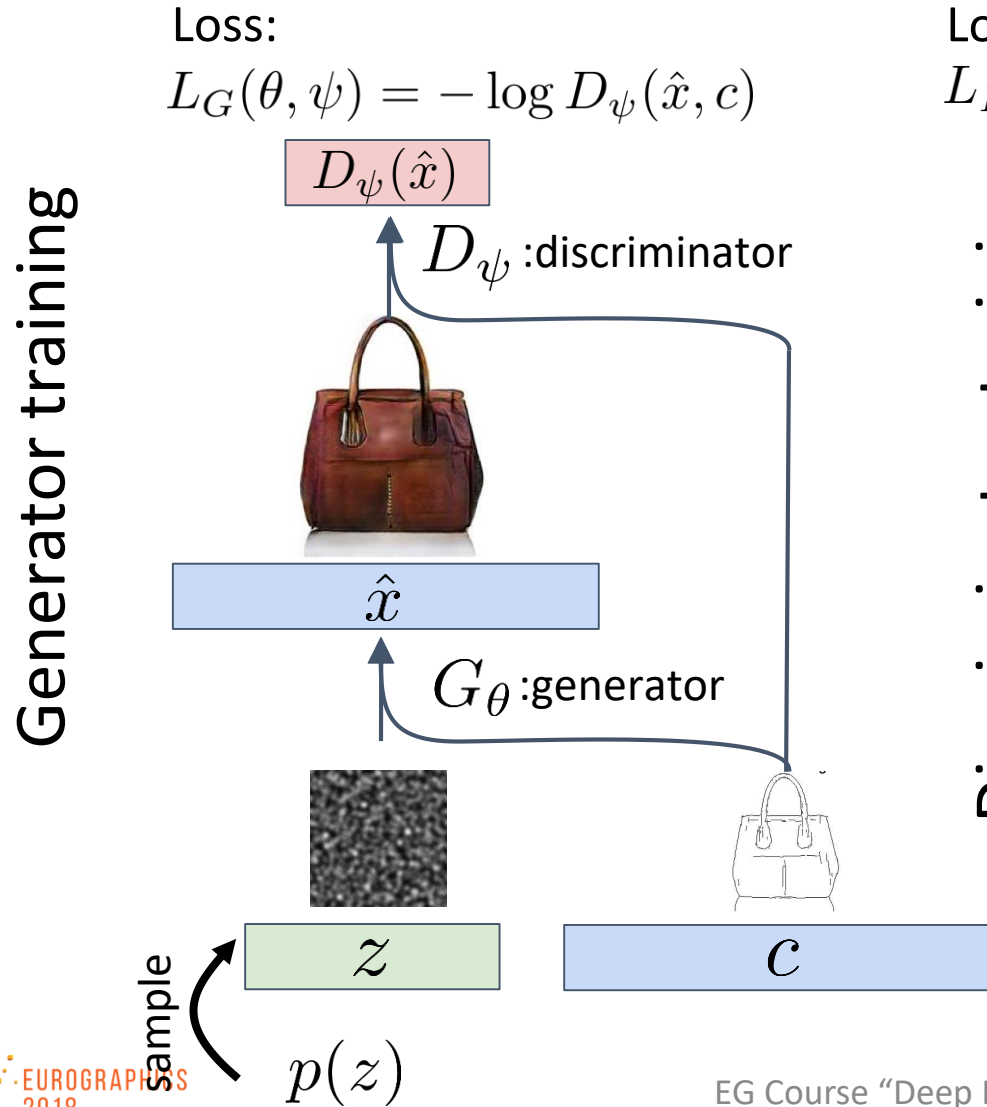
Generative Adversarial Network
(gan.ipynb)

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

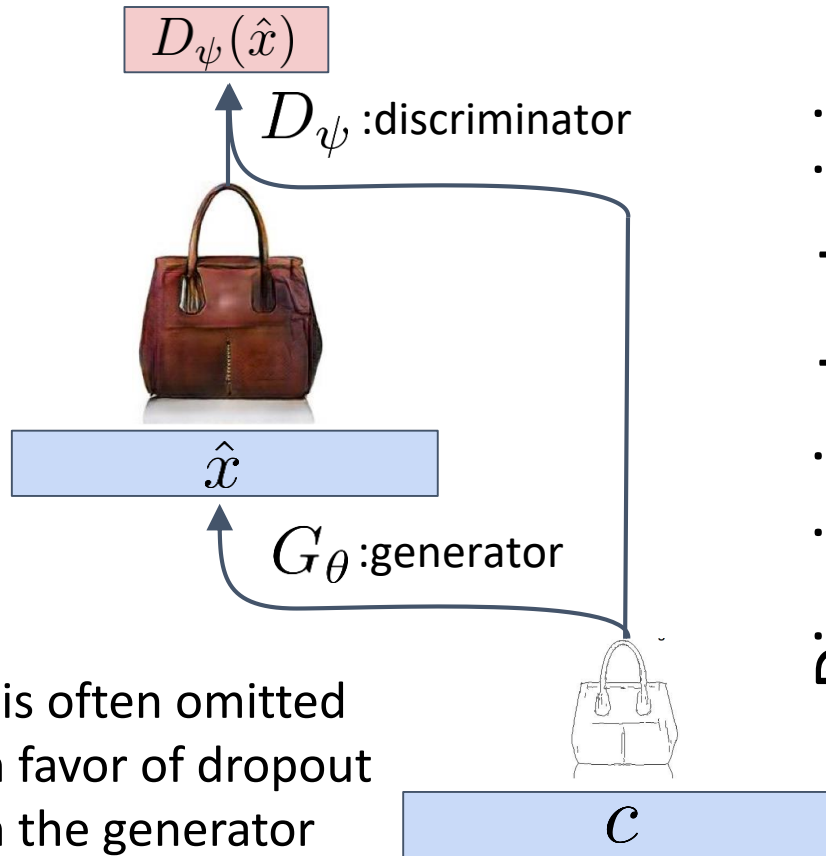


Conditional GANs: Low Variation per Condition

Loss:

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

Generator training

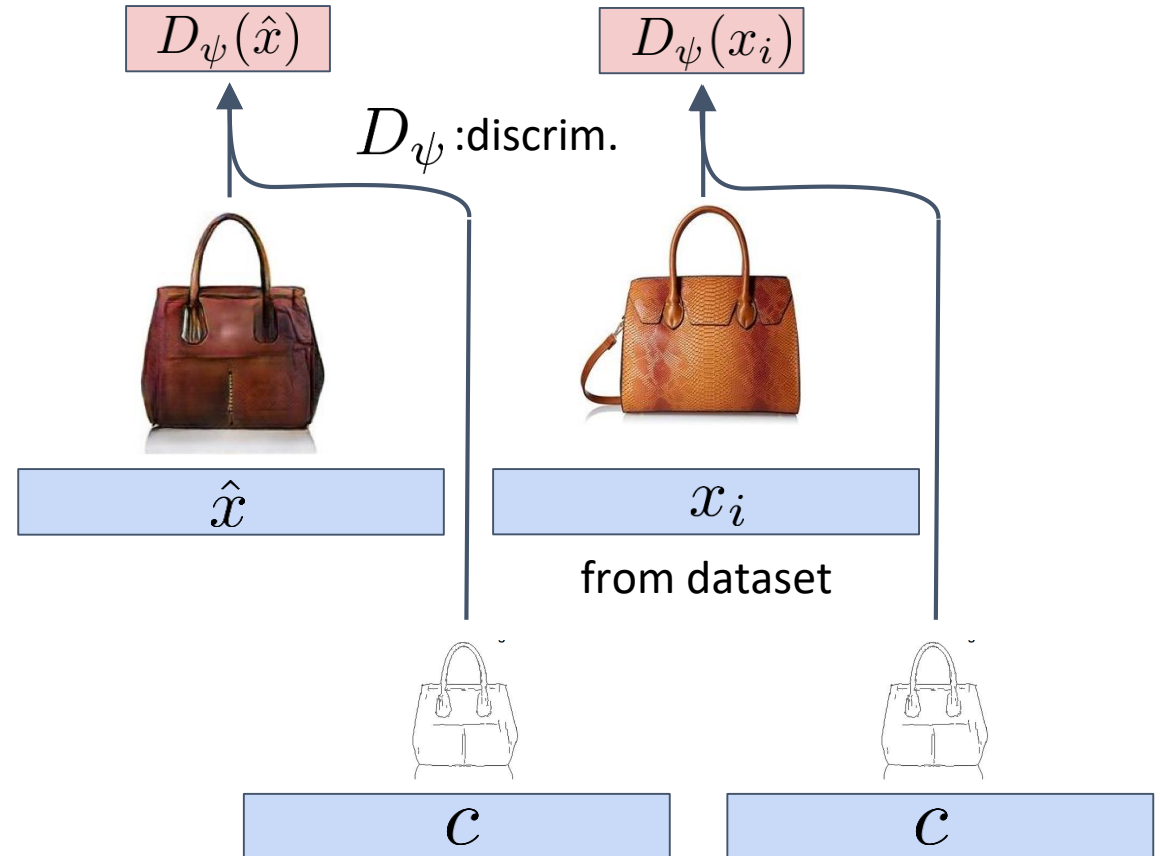


z is often omitted
in favor of dropout
in the generator

Loss:

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

Discriminator training



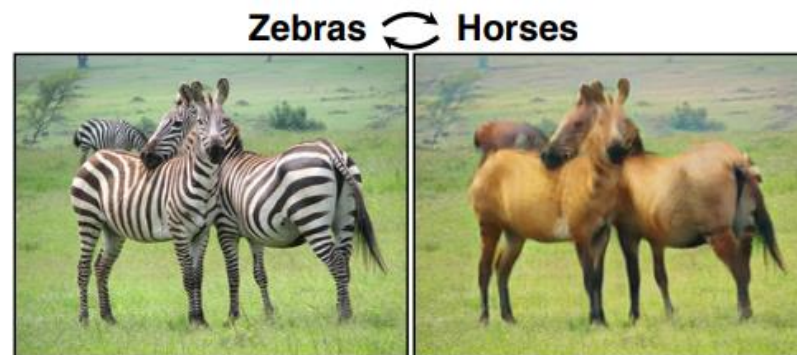
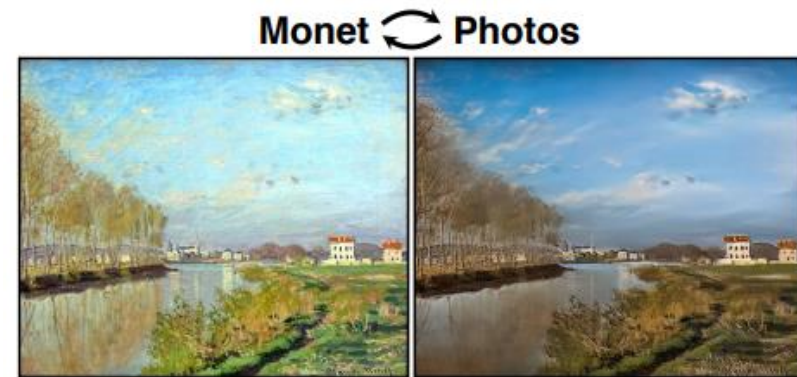
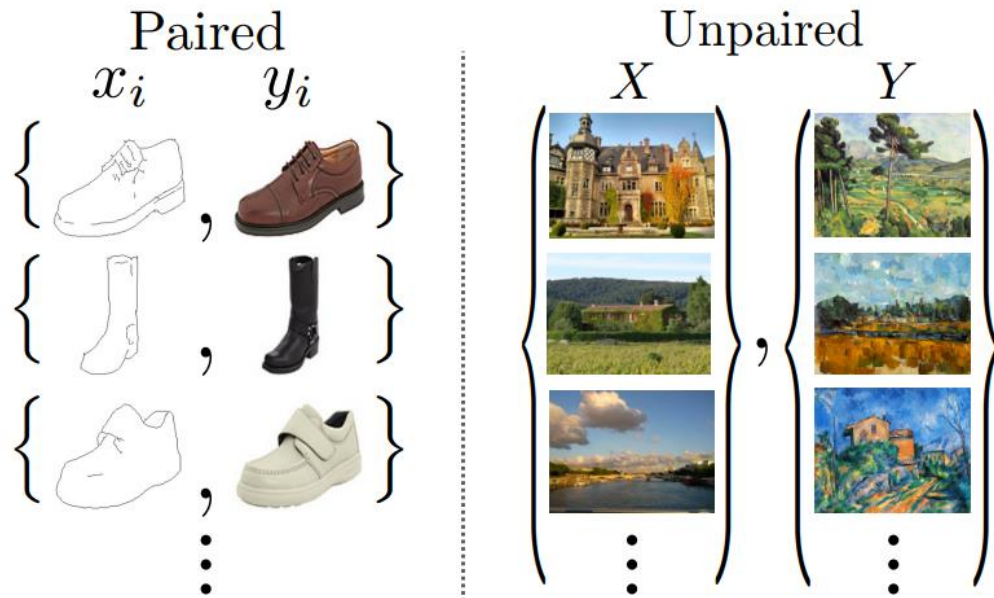
Demos

CGAN

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

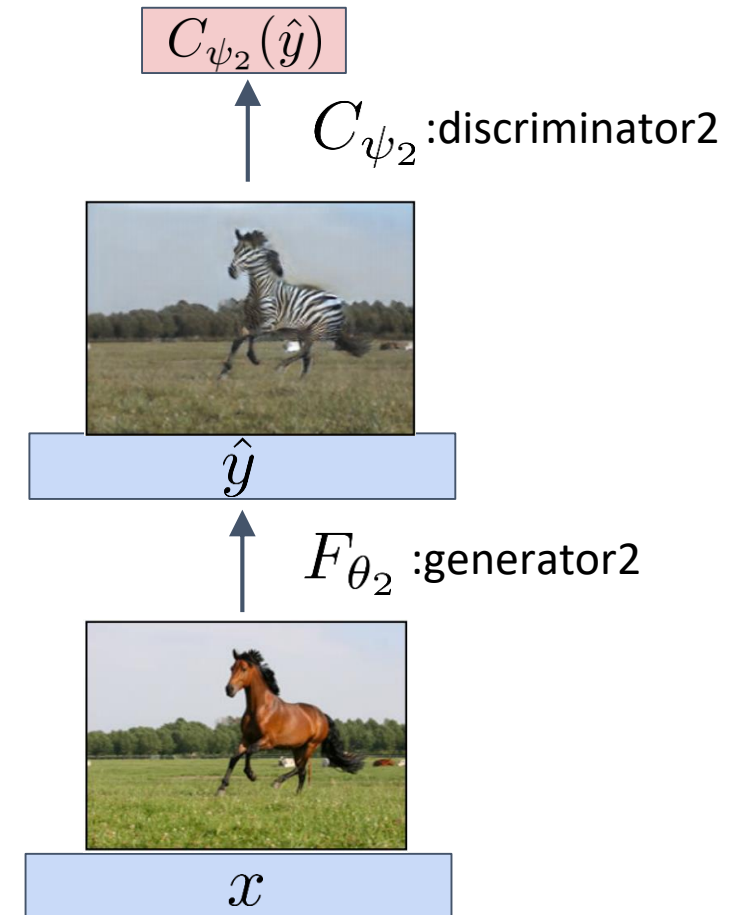
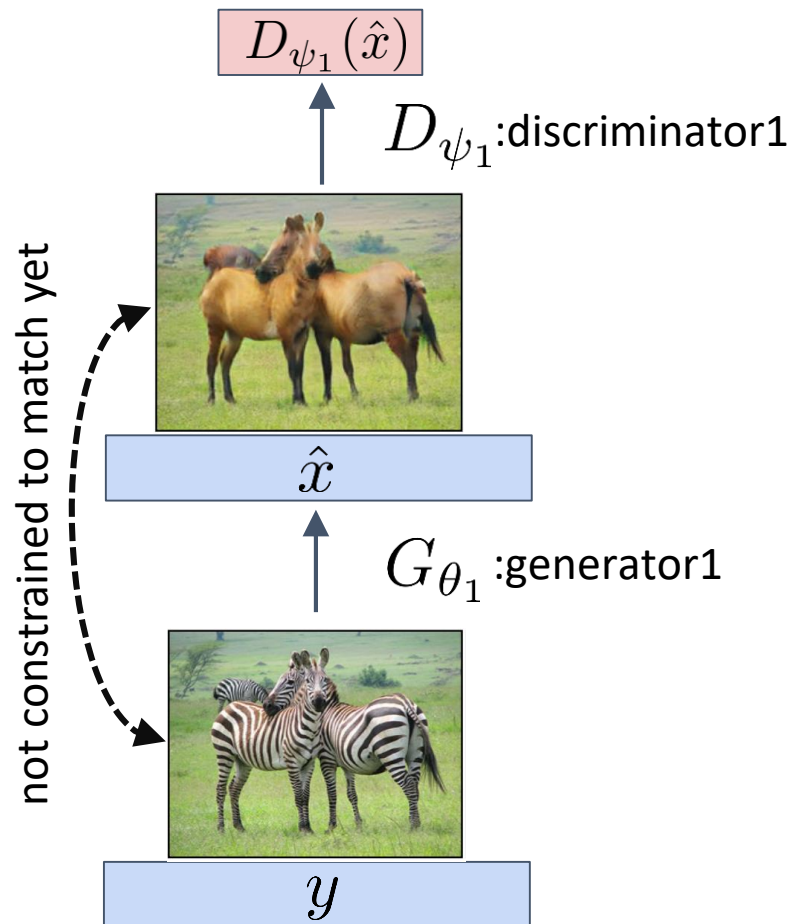
CycleGANs

- Less supervision than CGANs: mapping between unpaired datasets
- Two GANs + cycle consistency

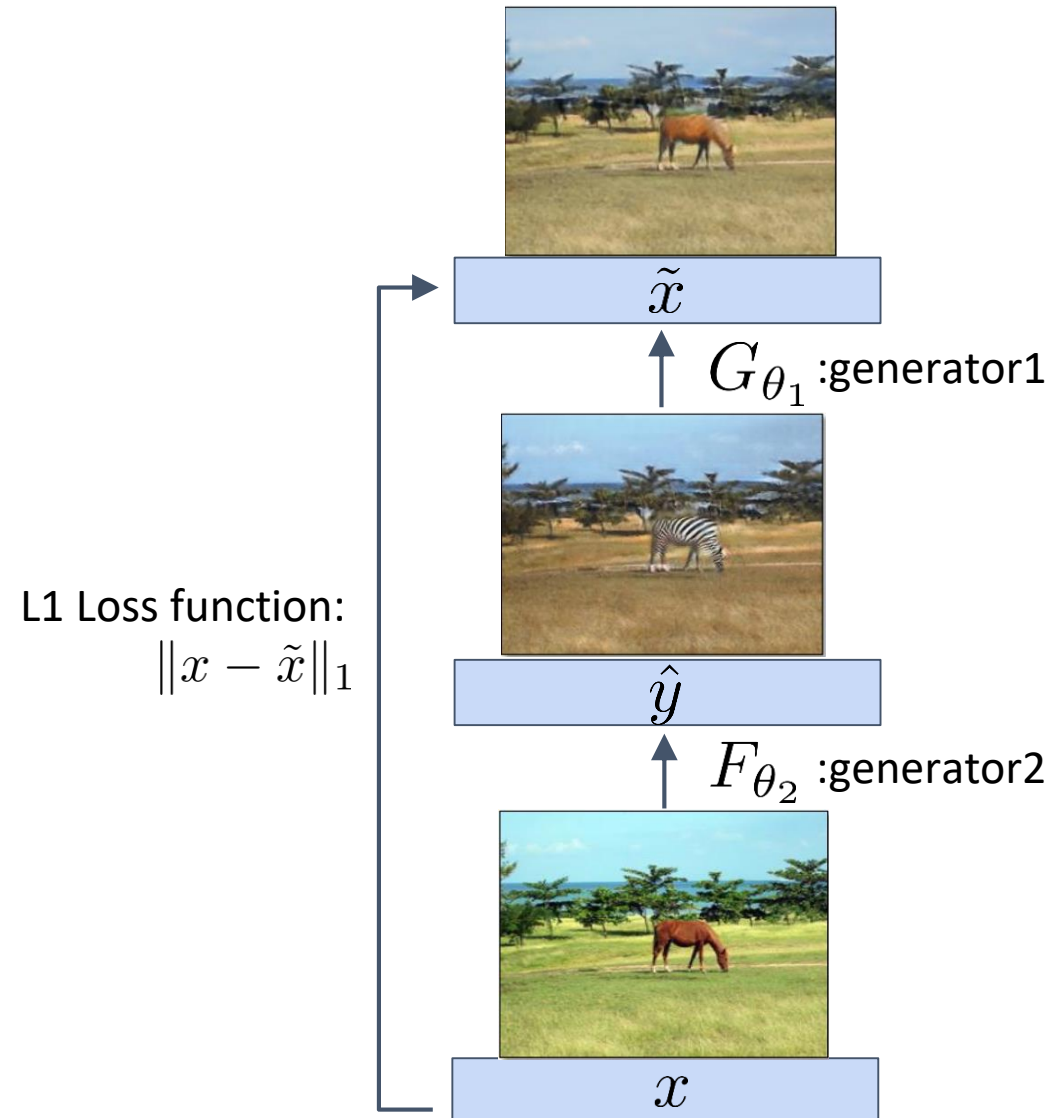
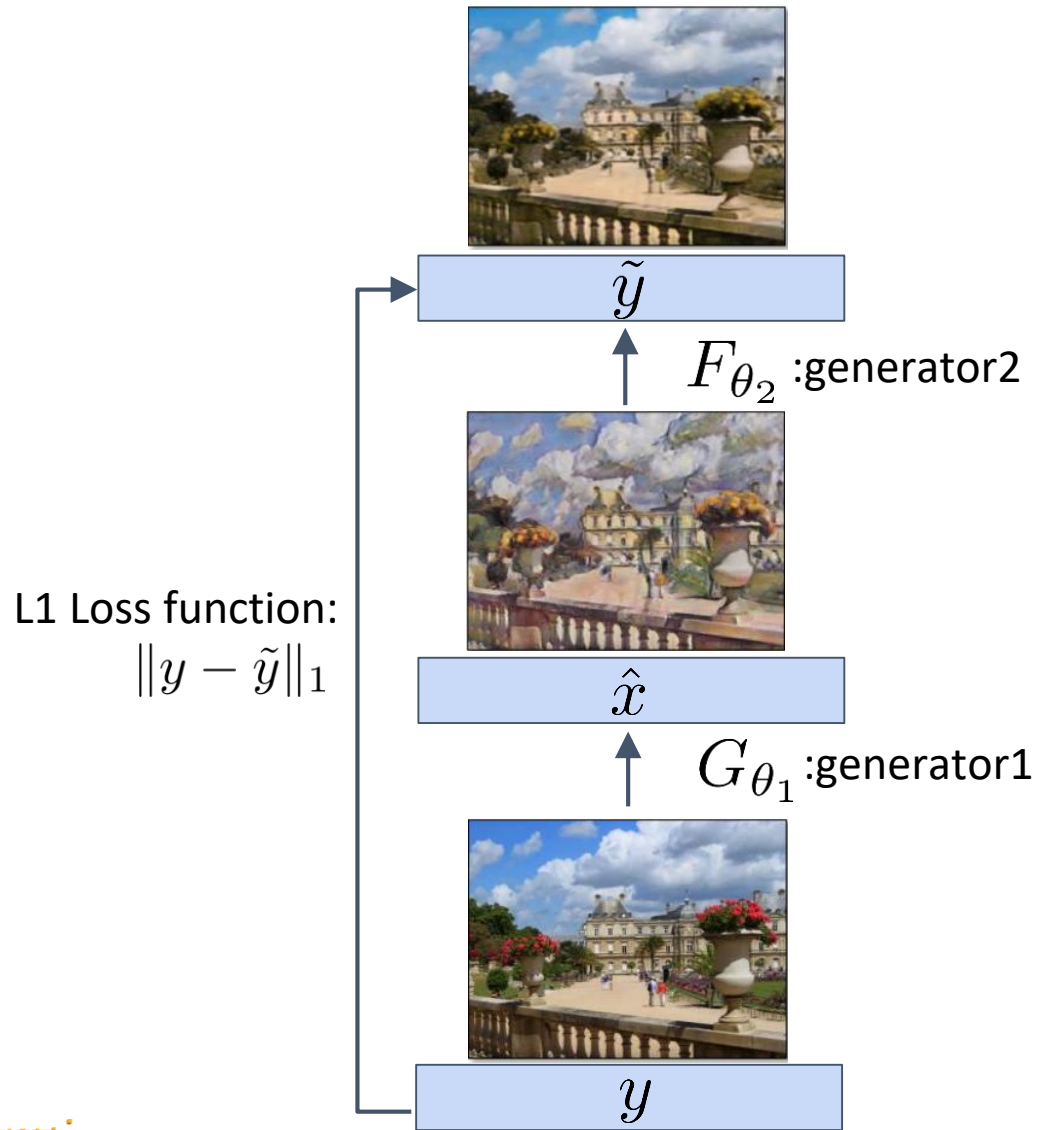


CycleGAN: Two GANs ...

- Not conditional, so this alone does not constrain generator input and output to match



CycleGAN: ... and Cycle Consistency



Unstable Training

GAN training can be unstable

Three current research problems (may be related):

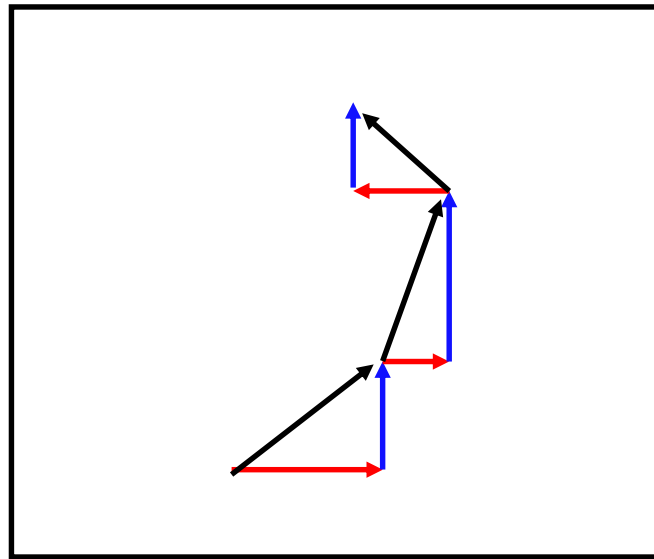
- Reaching a Nash equilibrium (the gradient for both L_G and L_D is 0)
- p_θ and p_{data} initially don't overlap
- Mode Collapse

GAN Training

- Vector-valued loss: $\mathbf{L}(\theta, \psi) = \begin{pmatrix} L_G(\theta, \psi) \\ L_D(\theta, \psi) \end{pmatrix}$
- In each iteration, gradient descent approximately follows this vector over the parameter space (θ, ψ) :

$$\mathbf{V}(\theta, \psi) = \begin{pmatrix} \frac{\partial}{\partial \theta} L_G(\theta, \psi) \\ \frac{\partial}{\partial \psi} L_D(\theta, \psi) \end{pmatrix}$$

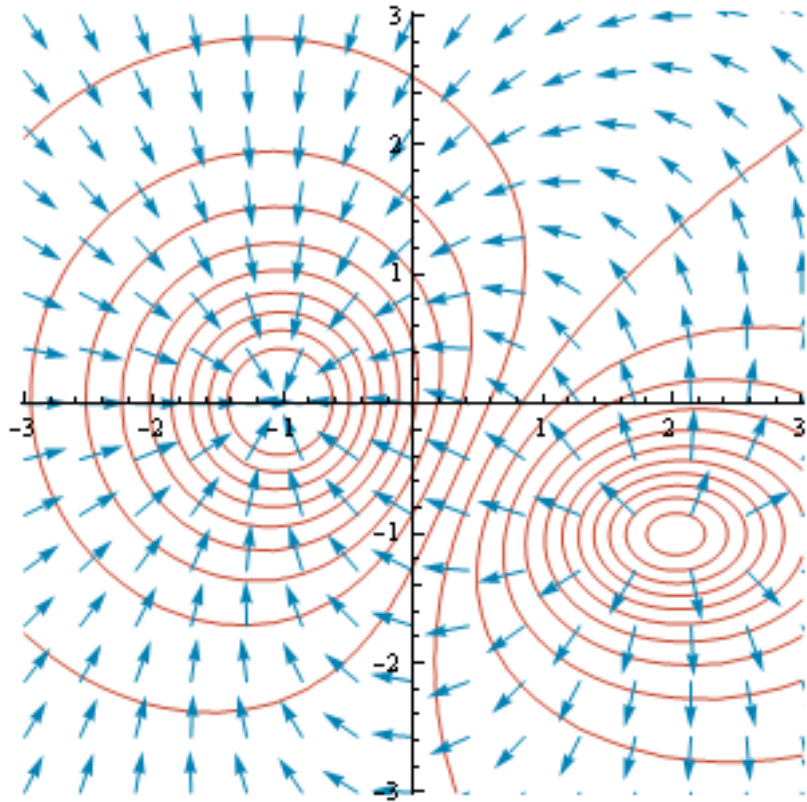
ψ



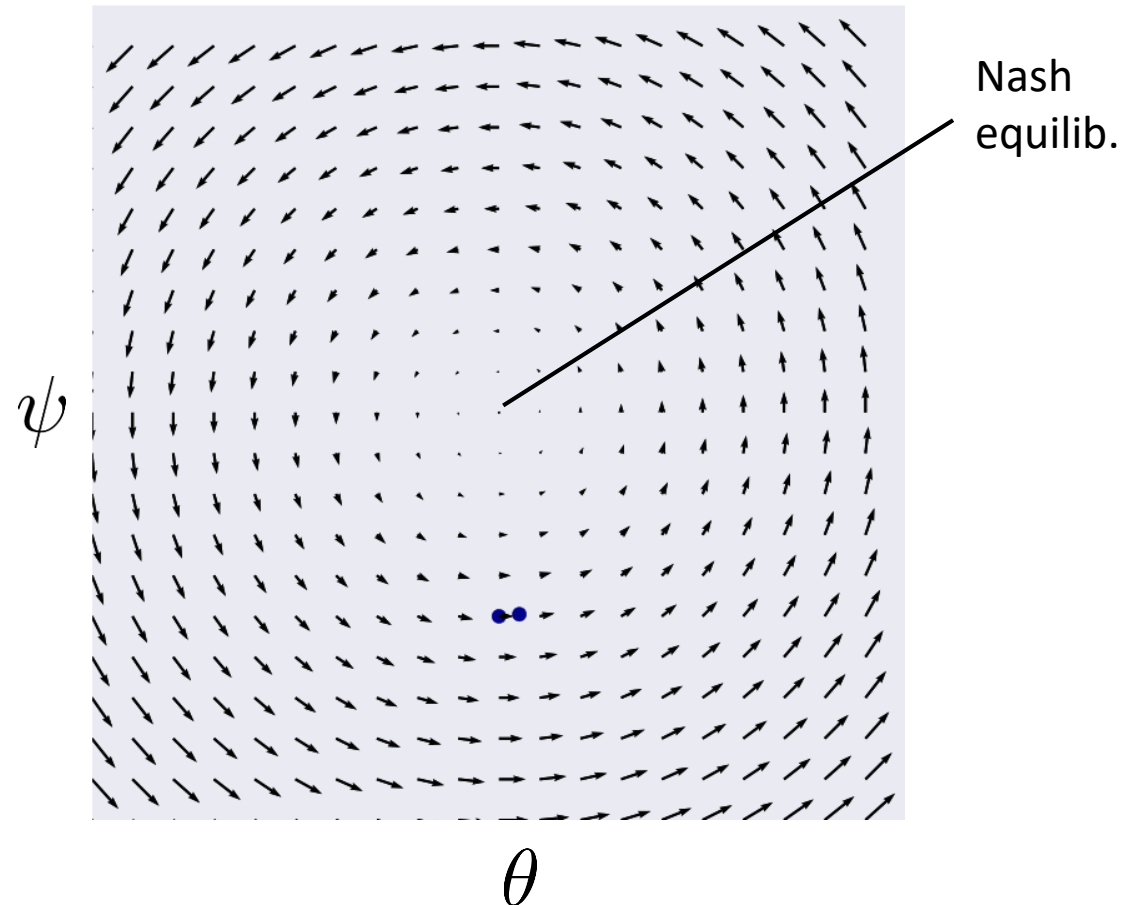
- $\frac{\partial}{\partial \theta} L_G(\theta, \psi)$
- $\frac{\partial}{\partial \psi} L_D(\theta, \psi)$
- $\mathbf{V}(\theta, \psi)$

Reaching Nash Equilibrium

Gradient field example



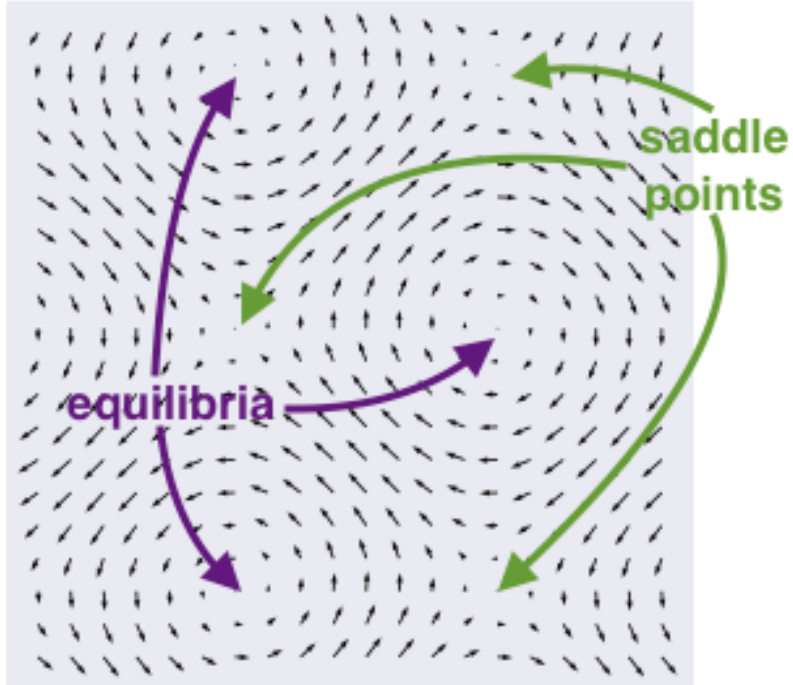
$V(\theta, \psi)$ Example



Reaching Nash Equilibrium

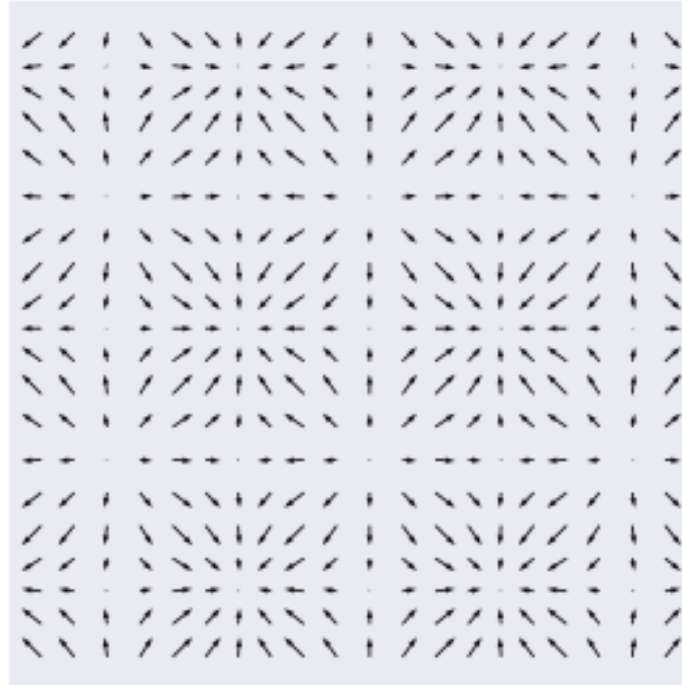
Solution attempt: relaxation with term: $-\nabla L = \frac{\partial}{\partial \theta} \|\mathbf{V}(\theta, \psi)\|_2^2$

Non-conservative field v



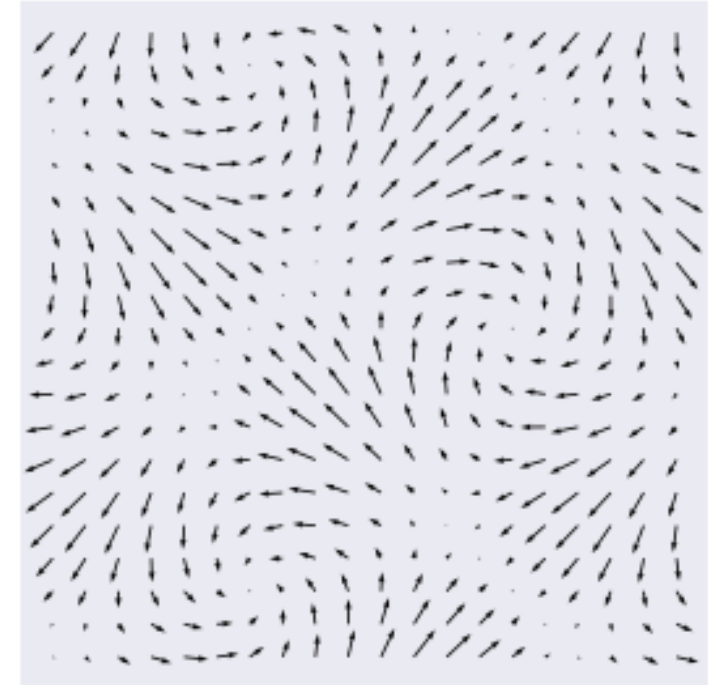
no relaxation has cycles

Conservative field $-\nabla L$



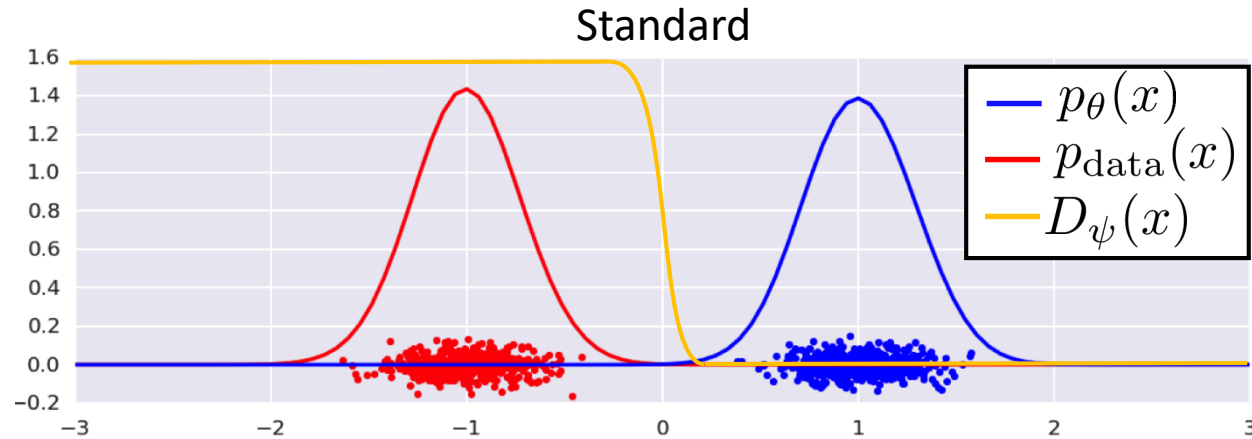
full relaxation introduces
bad Nash equilibria

Combined field $v - 0.6\nabla L$

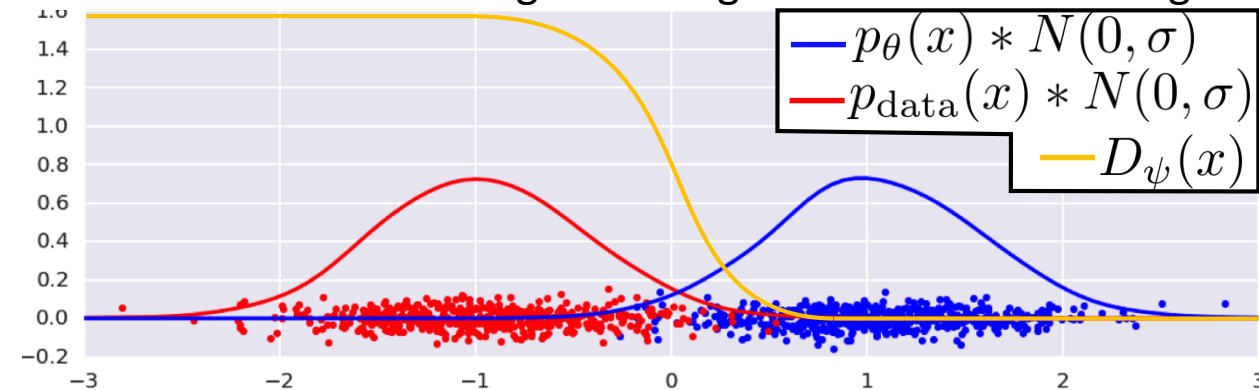


mixture works sometimes

Generator and Data Distribution Don't Overlap



Instance noise: adding noise to generated and real images

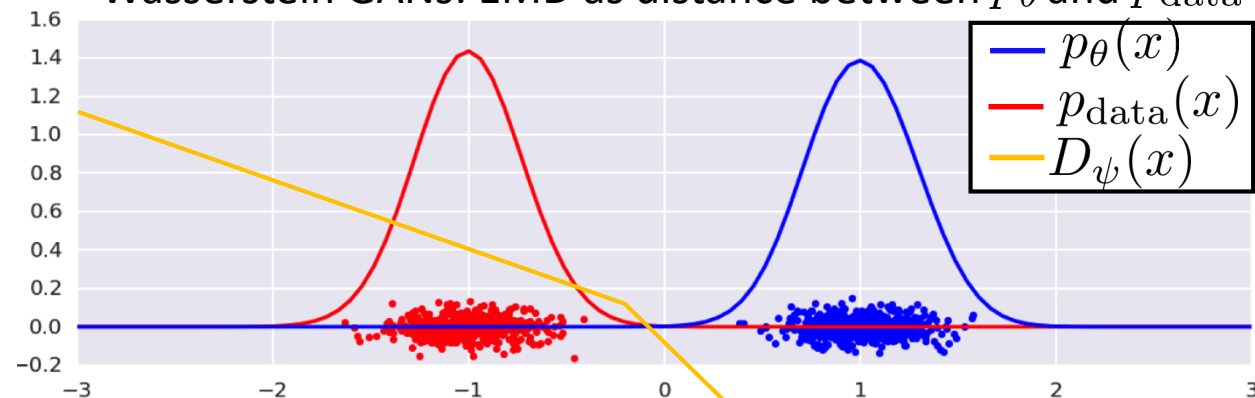


Roth et al. suggest an analytic convolution with a gaussian:

Stabilizing Training of Generative Adversarial Networks

through Regularization, Roth et al. 2017

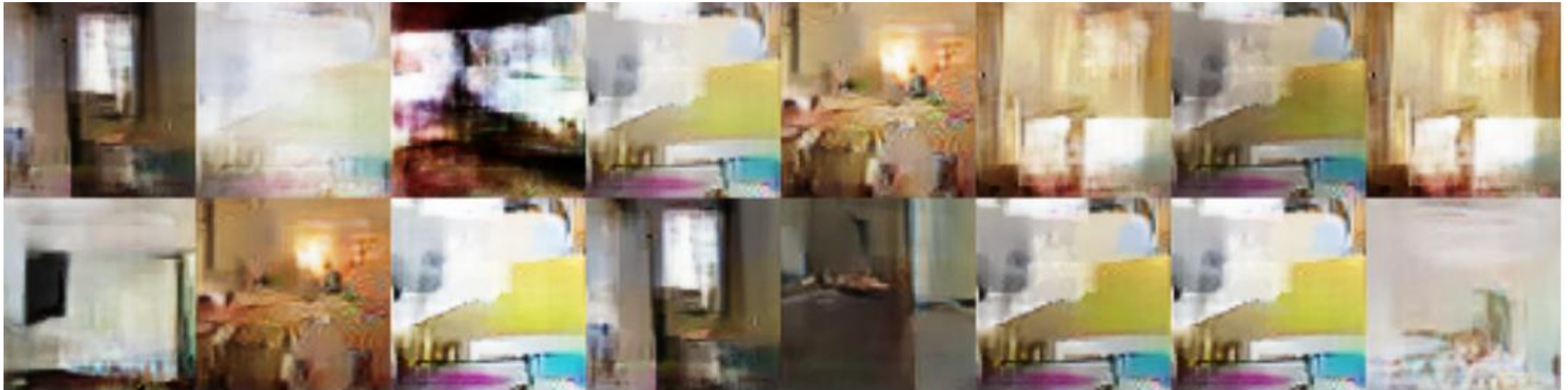
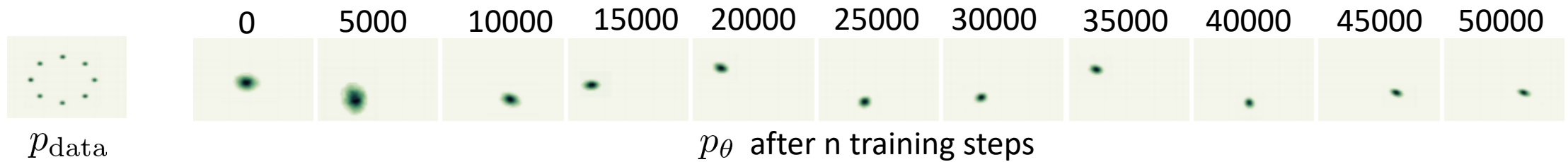
Wasserstein GANs: EMD as distance between p_{θ} and p_{data}



Mode Collapse

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

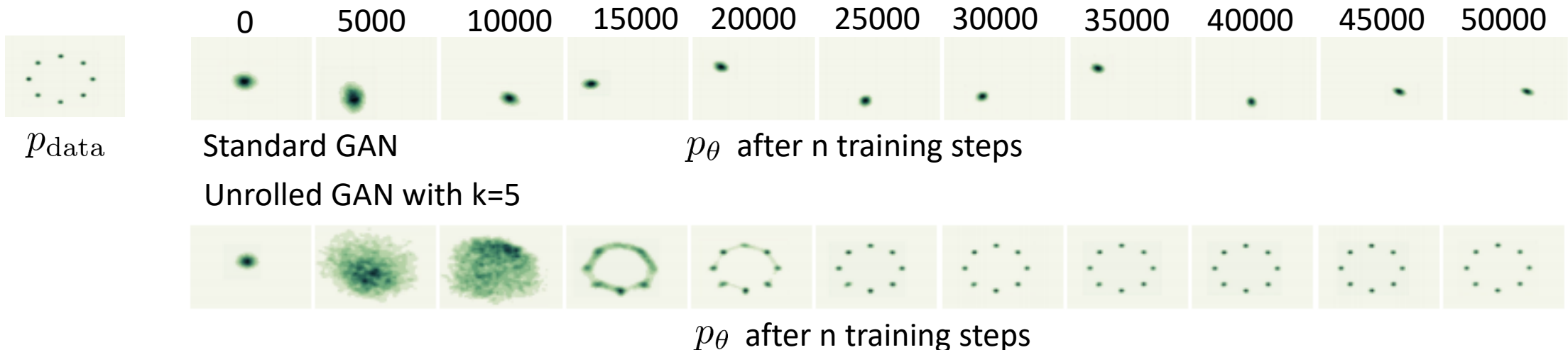
p_{θ} only covers one or a few modes of p_{data}



Mode Collapse

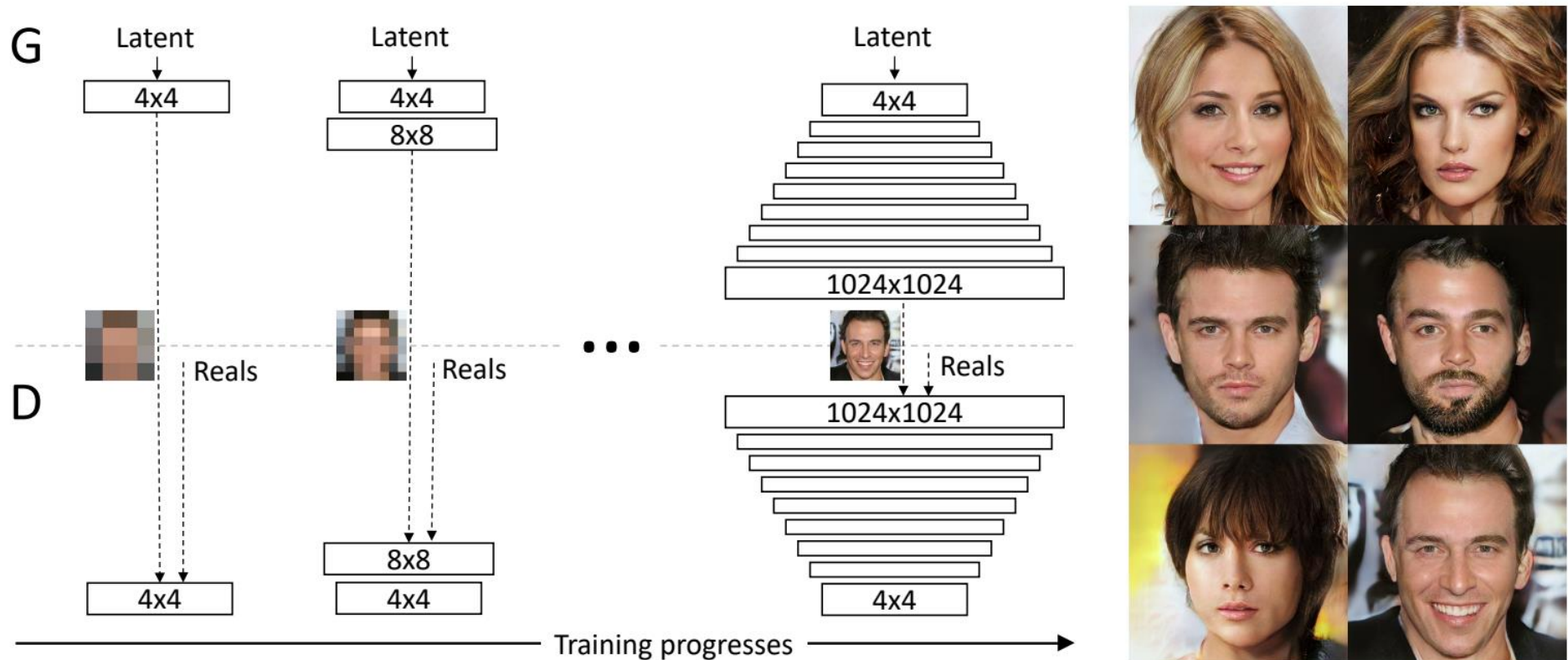
Solution attempts:

- 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



Progressive GANs

- Resolution is increased progressively during training
- Also other tricks like using minibatch statistics and normalizing feature vectors



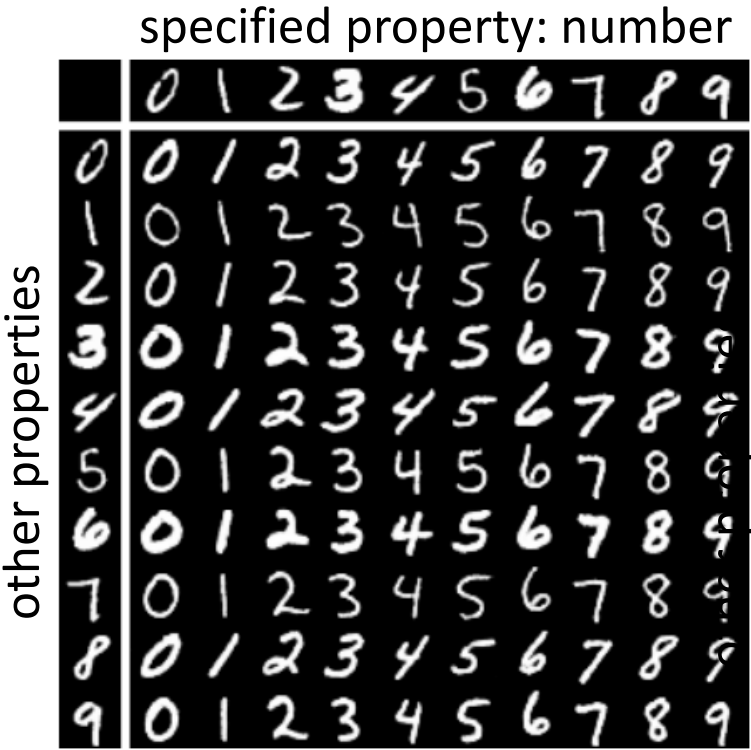
Disentanglement

z

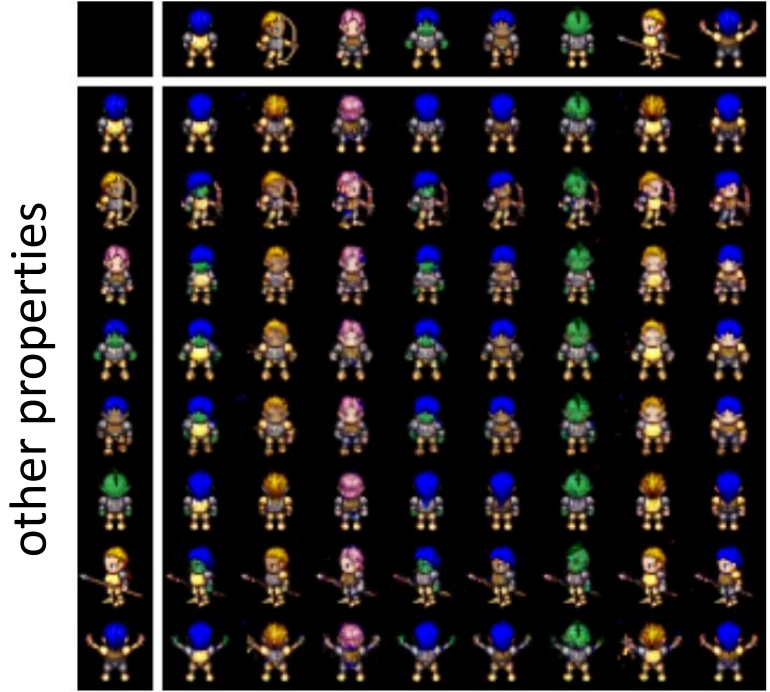
Entangled: different properties may be mixed up over all dimensions

z_a z_b ...

Disentangled: different properties are in different dimensions



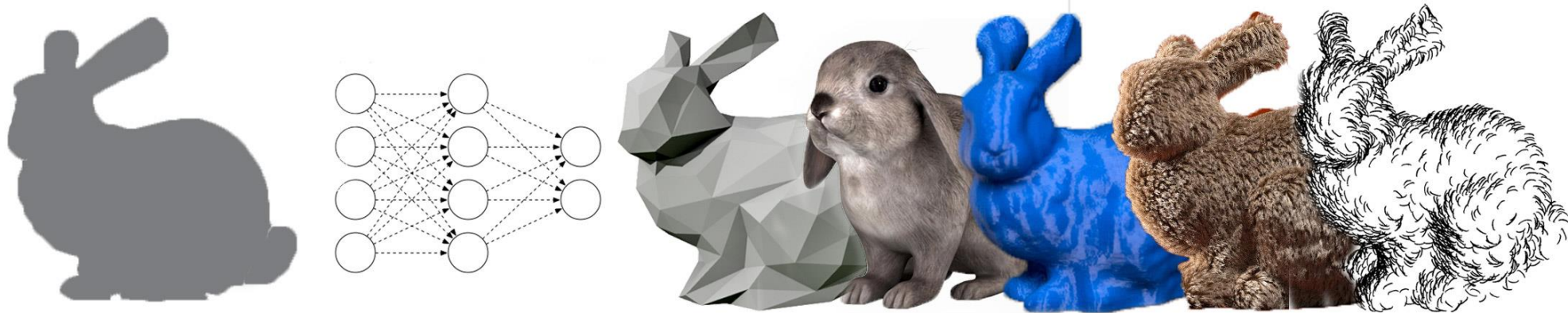
specified property: character



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)

Thank you!



<http://geometry.cs.ucl.ac.uk/dl4g/>