

# COMP541

## DEEP LEARNING

Lecture #10 – Generative Adversarial Networks  
and Flow-Based Models



KOÇ  
UNIVERSITY

Aykut Erdem // Koç University // Fall 2023

# Previously on COMP541

- supervised vs. unsupervised learning
- generative modeling
- sparse coding
- autoencoders
- autoregressive generative models

Video: Samples from "cooking" subset of Kinetics, Weissenborn et al.



# Lecture overview

- supervised vs unsupervised learning
- generative modeling
- basic foundations
  - sparse coding
  - autoencoders
- generative adversarial networks (GANs)

**Disclaimer:** Some of the material and slides for this lecture were borrowed from

- Justin Johnson's EECS 498/598 class
- Ruslan Salakhutdinov's talk titled "Unsupervised Learning: Learning Deep Generative Models"
- Ian Goodfellow's tutorial on "Generative Adversarial Networks"
- Aaron Courville's IFT6135 class

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression,  
object detection, semantic  
segmentation, image captioning,  
sentiment analysis, etc.

Classification



Cat

# Supervised vs Unsupervised Learning

## Supervised Learning

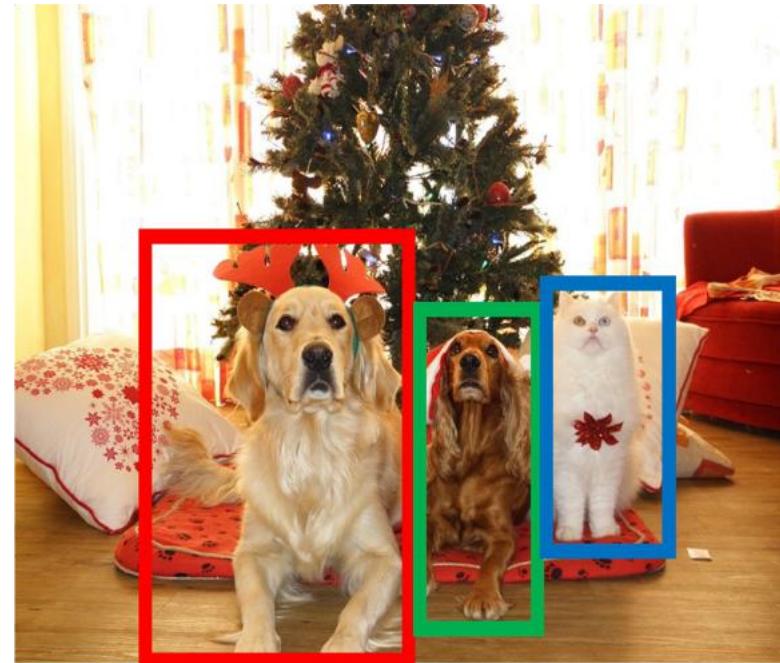
**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Object Detection



**DOG , DOG , CAT**

# Supervised vs Unsupervised Learning

## Supervised Learning

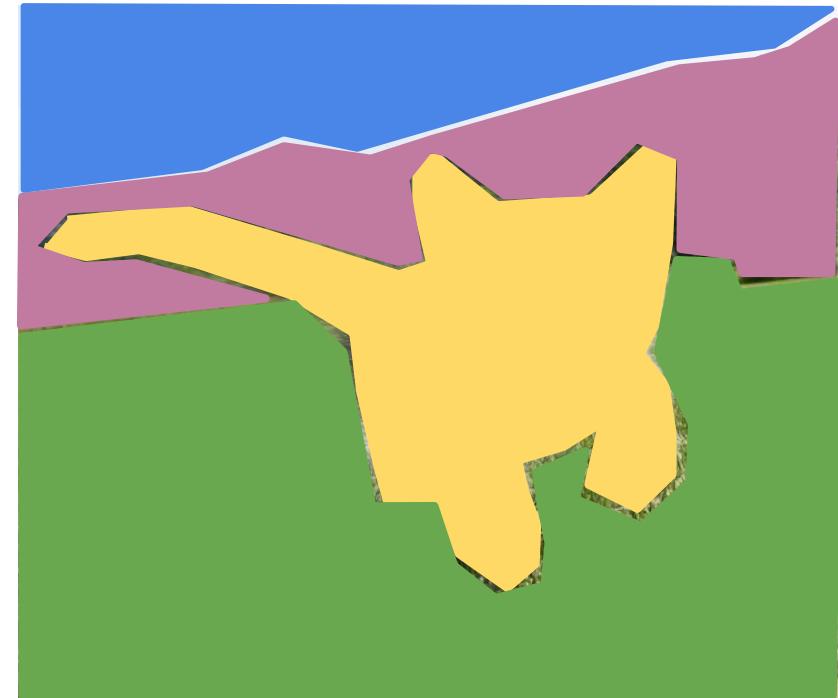
**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Semantic Segmentation



GRASS, CAT, TREE, SKY

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Image captioning



*A cat sitting on a suitcase on the floor*

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Sentiment Analysis

“This Movie is amazing. It has a great plot and talented actors, and the supporting cast is really good as well.”



# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Unsupervised Learning

**Data:**  $x$

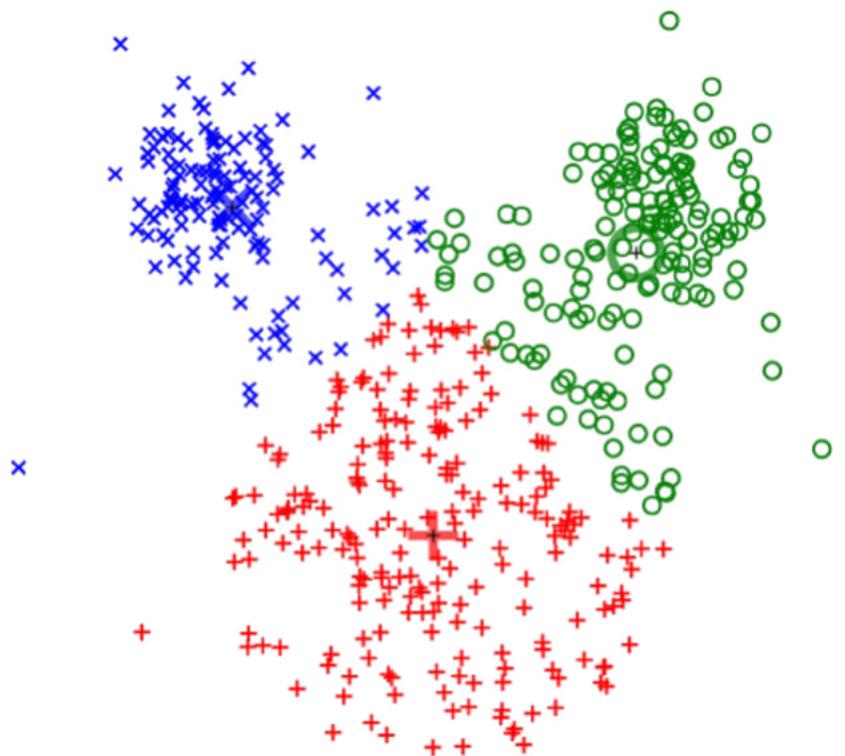
Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

Clustering  
(e.g. K-Means)



**Unsupervised Learning**

**Data:** x

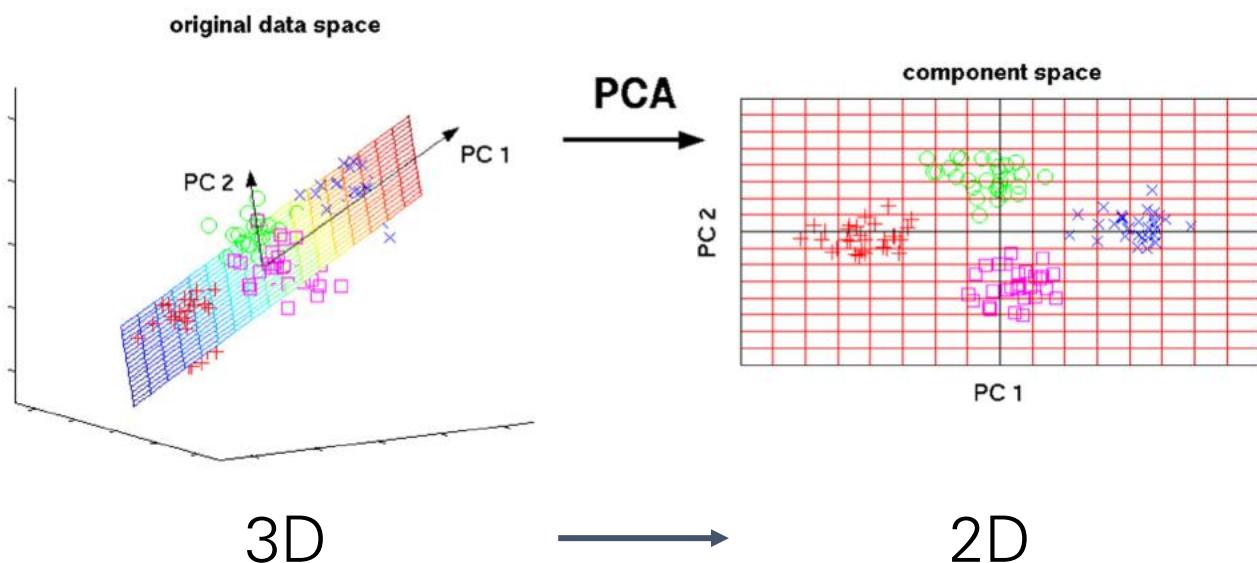
Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

Dimensionality Reduction  
(e.g. Principal Components Analysis)



## Unsupervised Learning

**Data:**  $x$

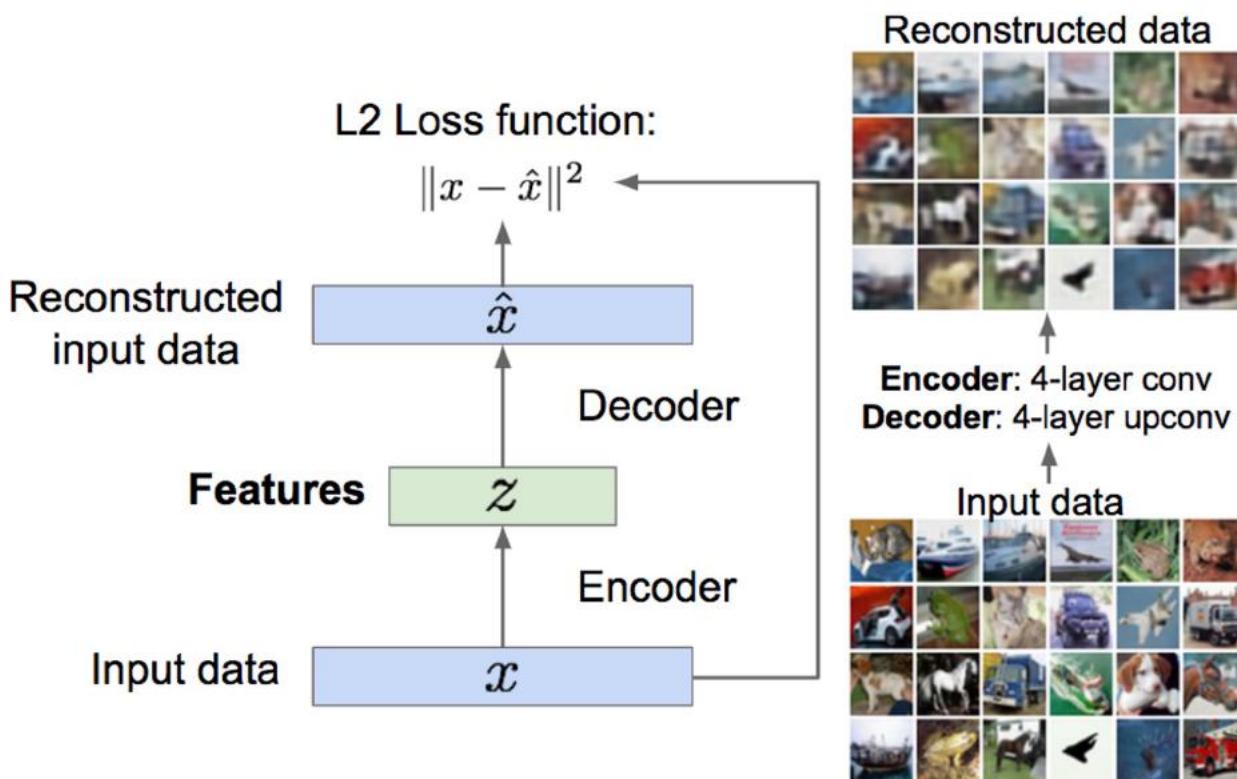
Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

Feature Learning  
(e.g. autoencoders)



## Unsupervised Learning

**Data:**  $x$

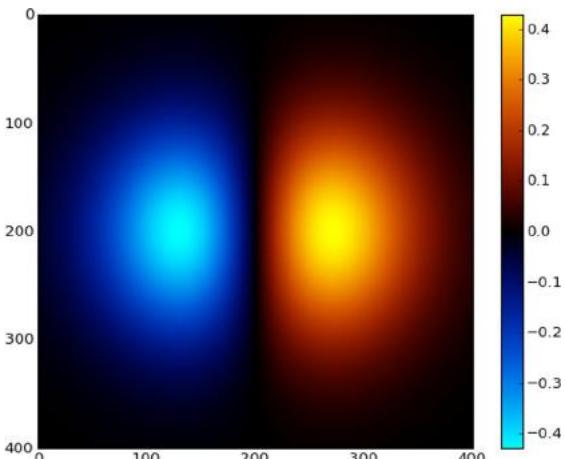
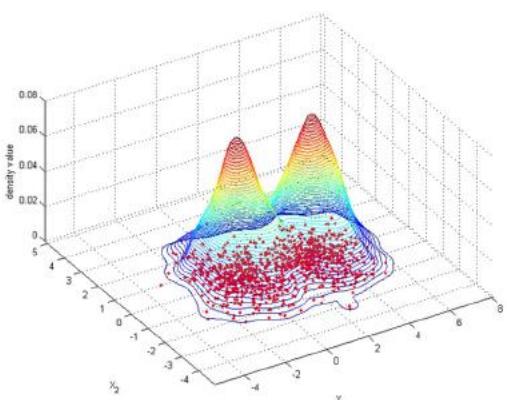
Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering,  
dimensionality reduction,  
feature learning, density  
estimation, etc.

# Supervised vs Unsupervised Learning

## Density Estimation



## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**  
Learn  $p(x|y)$

Data:  $x$



Label:  $y$   
Cat

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**  
Learn  $p(x|y)$

Data:  $x$



Label:  $y$   
Cat

**Probability Recap:**

**Density Function**

$p(x)$  assigns a positive number to each possible  $x$ ; higher numbers mean  $x$  is more likely

Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of  $x$  **compete** for density

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**  
Learn  $p(x|y)$

Data:  $x$



$P(\text{cat} | \text{kitten})$



$P(\text{dog} | \text{kitten})$



Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of  $x$  **compete** for density

**Density Function**

$p(x)$  assigns a positive number to each possible  $x$ ; higher numbers mean  $x$  is more likely

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

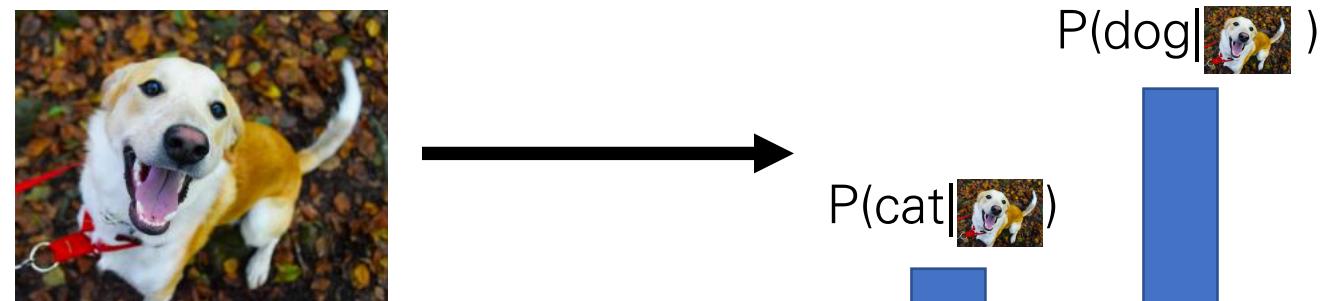
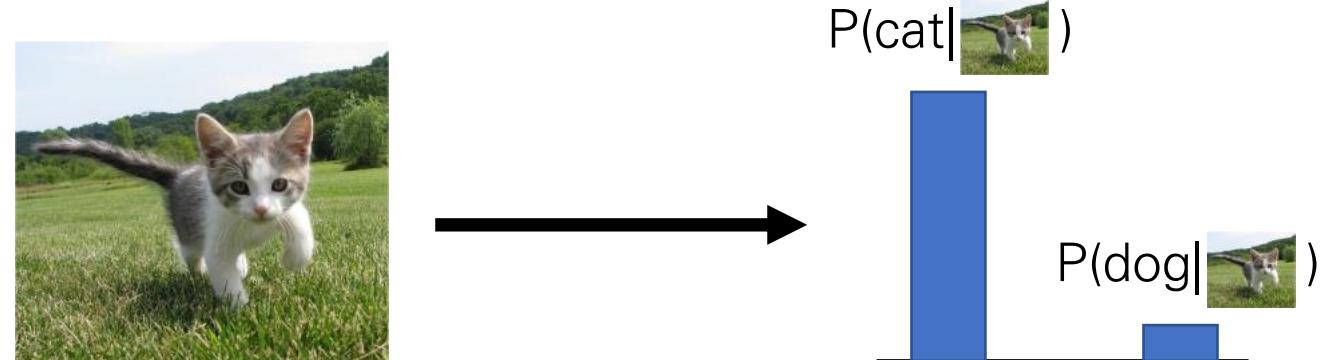
**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional**

**Generative Model:**

Learn  $p(x|y)$



Discriminative model: the possible labels for each input “compete” for probability mass.  
But no competition between **images**

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional**

**Generative Model:**

Learn  $p(x|y)$



$P(\text{cat} | \text{monkey})$

$P(\text{dog} | \text{monkey})$



$P(\text{dog} | \text{dog})$

$P(\text{cat} | \text{dog})$

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

# Discriminative vs Generative Models

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional**

**Generative Model:**

Learn  $p(x|y)$



$$P(\text{cat} | \text{monkey})$$

$$P(\text{dog} | \text{monkey})$$



$$P(\text{dog} | \text{colorful pattern})$$

$$P(\text{cat} | \text{colorful pattern})$$

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

# Discriminative vs Generative Models

**Discriminative Model:**

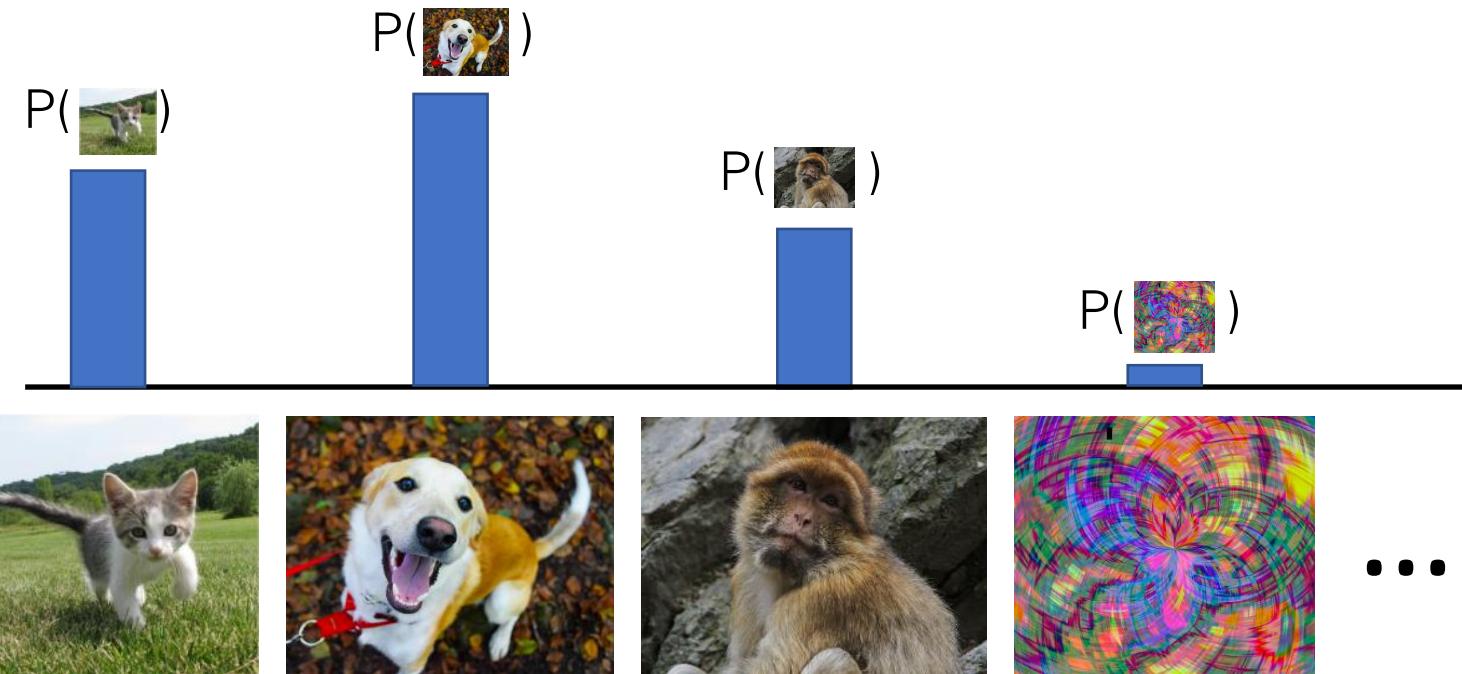
Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**

Learn  $p(x|y)$



Generative model: All possible images compete with each other for probability mass

# Discriminative vs Generative Models

**Discriminative Model:**

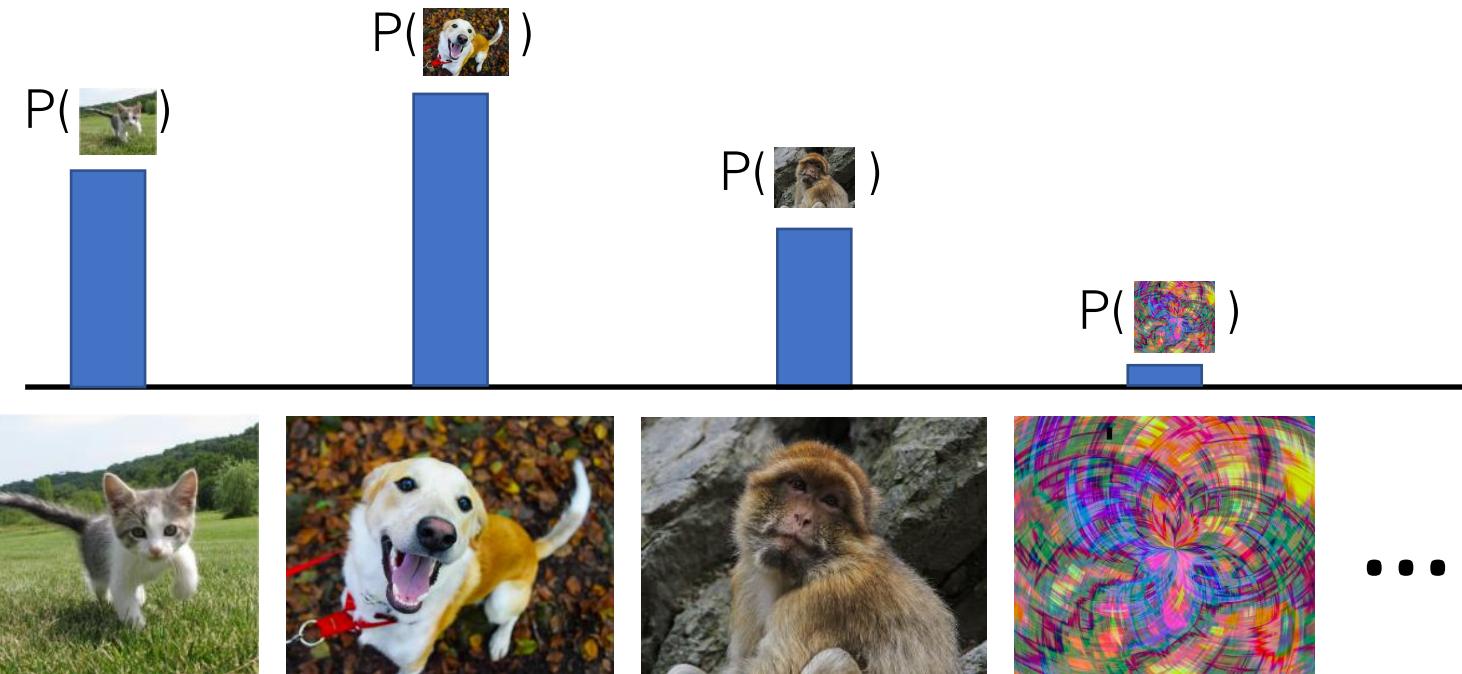
Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**

Learn  $p(x|y)$



Generative model: All possible images compete with each other for probability mass

Requires deep image understanding! Is a dog more likely to sit or stand? How about 3-legged dog vs 3-armed monkey?

# Discriminative vs Generative Models

**Discriminative Model:**

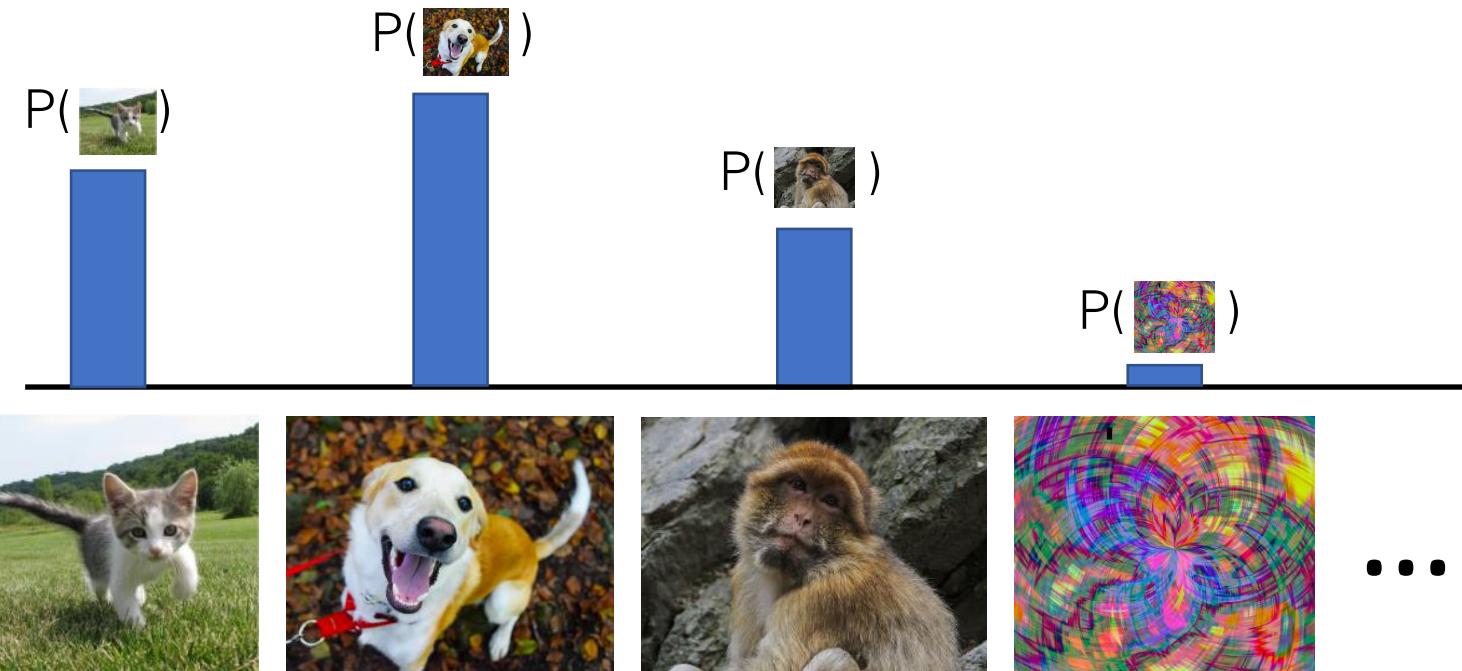
Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**

Learn  $p(x|y)$



Generative model: All possible images compete with each other for probability mass

Model can “reject” unreasonable inputs by assigning them small values

# Discriminative vs Generative Models

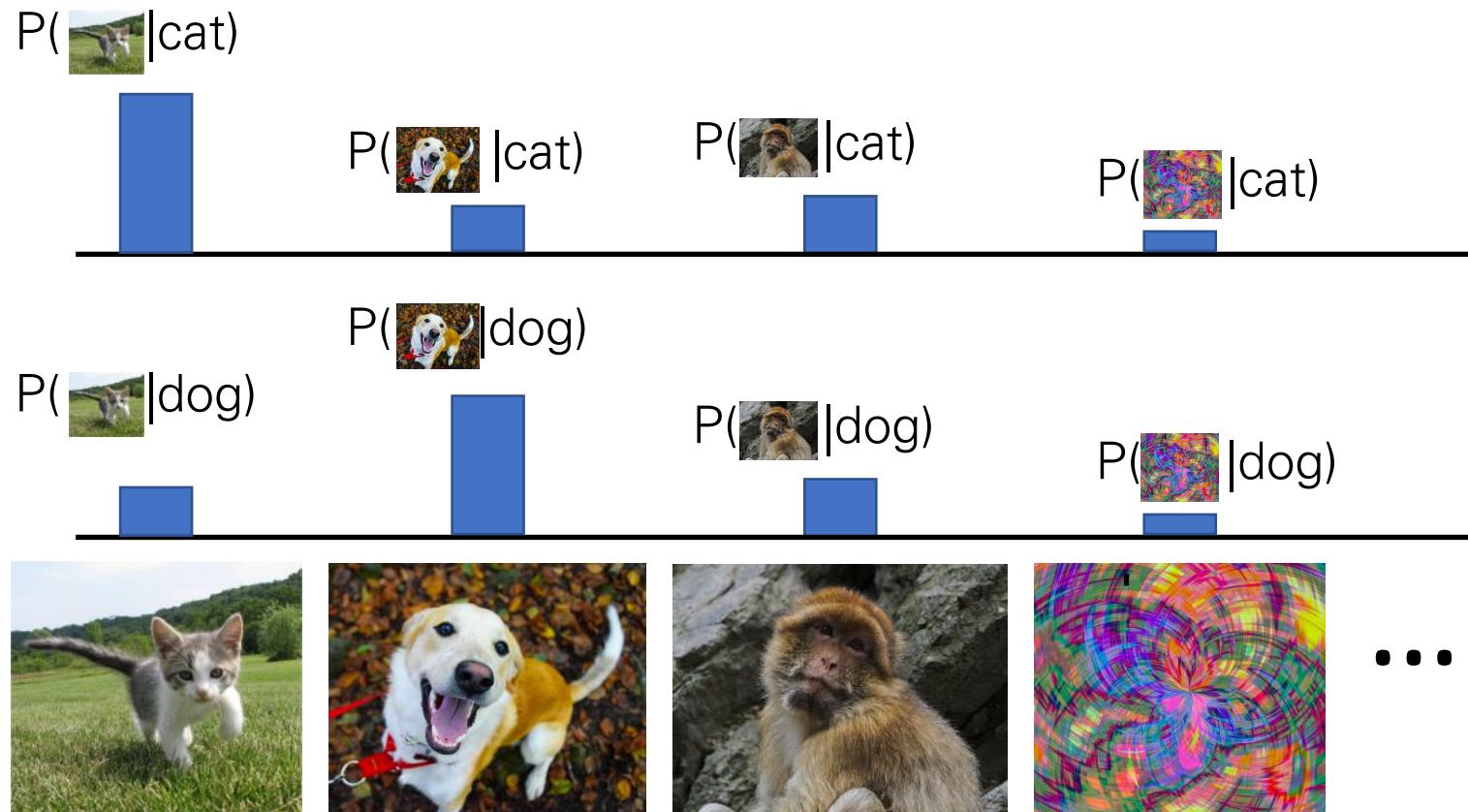
**Discriminative Model:**

Learn a probability distribution  $p(y|x)$

**Generative Model:**

Learn a probability distribution  $p(x)$

**Conditional Generative Model:**  
Learn  $p(x|y)$



Conditional Generative Model: Each possible label induces a competition among all images

# Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution  $p(y|x)$

Generative Model:

Learn a probability distribution  $p(x)$

Conditional  
Generative Model:  
Learn  $p(x|y)$

Recall **Bayes' Rule**:

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

# Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution  $p(y|x)$

Generative Model:

Learn a probability distribution  $p(x)$

Conditional Generative Model:  
Learn  $p(x|y)$

Recall **Bayes' Rule**:

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

Conditional Generative Model

Discriminative Model

(Unconditional) Generative Model

Prior over labels

We can build a conditional generative model from other components!

# What can we do with a discriminative model?

**Discriminative Model:**

Learn a probability  
distribution  $p(y|x)$



Assign labels to data  
Feature learning (with labels)

**Generative Model:**

Learn a probability  
distribution  $p(x)$

**Conditional  
Generative Model:**  
Learn  $p(x|y)$

# What can we do with a discriminative model?

**Discriminative Model:**

Learn a probability distribution  $p(y|x)$



Assign labels to data  
Feature learning (with labels)

**Generative Model:**

Learn a probability distribution  $p(x)$



Detect outliers  
Feature learning (without labels)  
Sample to **generate** new data

**Conditional**

**Generative Model:**

Learn  $p(x|y)$

# What can we do with a discriminative model?

## Discriminative Model:

Learn a probability distribution  $p(y|x)$



Assign labels to data  
Feature learning (with labels)

## Generative Model:

Learn a probability distribution  $p(x)$



Detect outliers  
Feature learning (without labels)  
Sample to **generate** new data

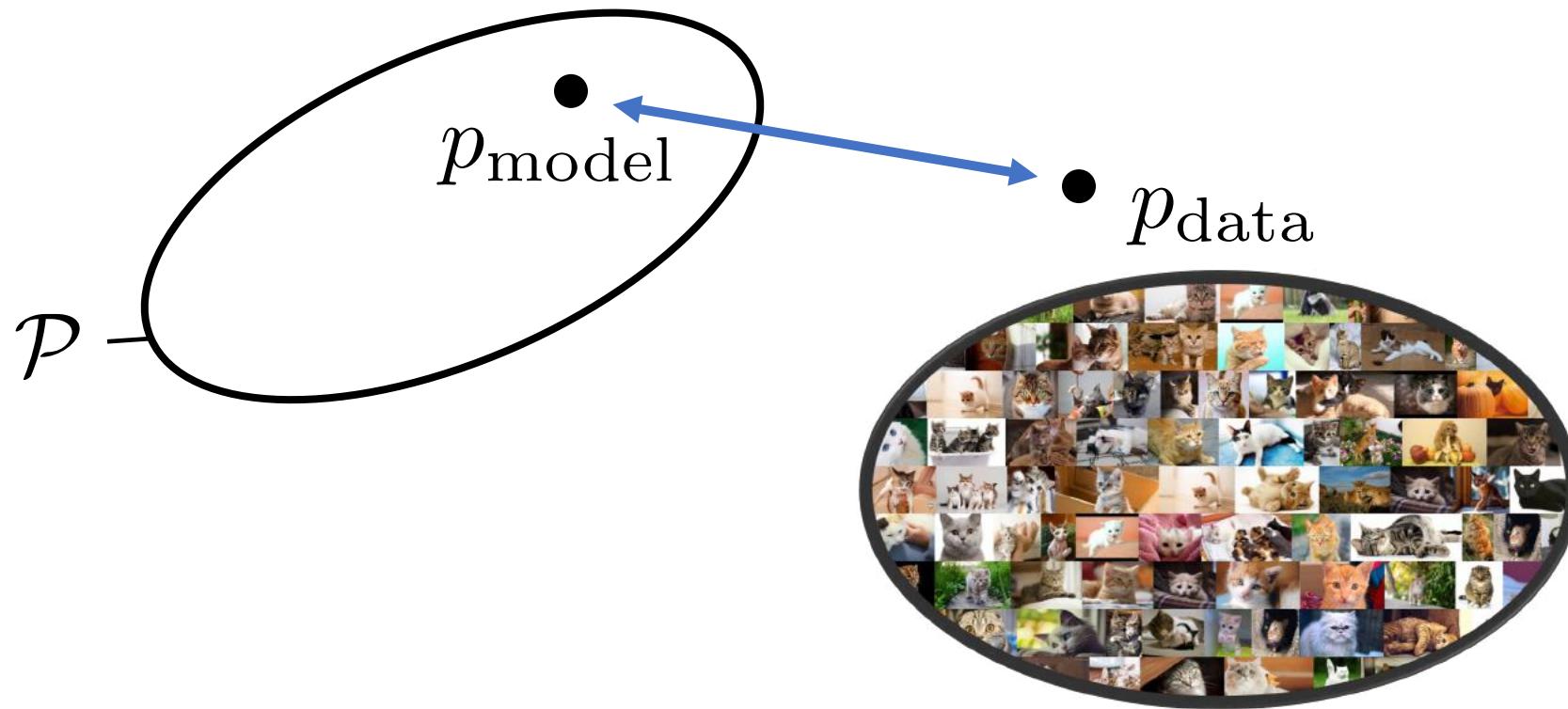
## Conditional Generative Model:

Learn  $p(x|y)$



Assign labels, while rejecting outliers!  
Generate new data conditioned on input labels

# Generative Modeling



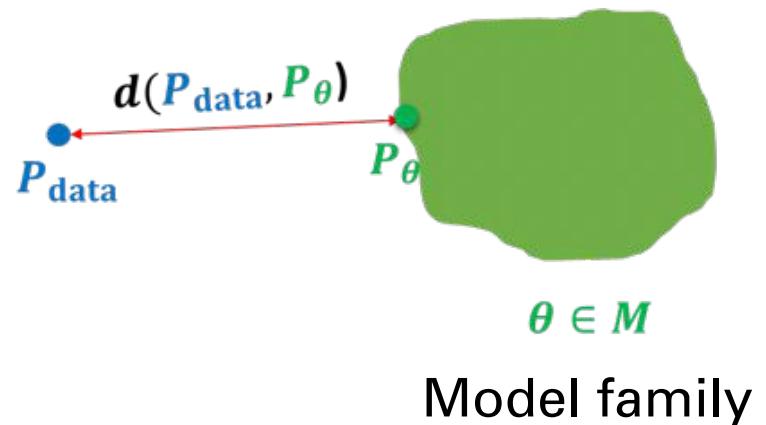
- Goal: Learn some underlying hidden structure of the training samples to generate novel samples from same data distribution

# Learning a generative model

- We are given a training set of examples, e.g., images of dogs



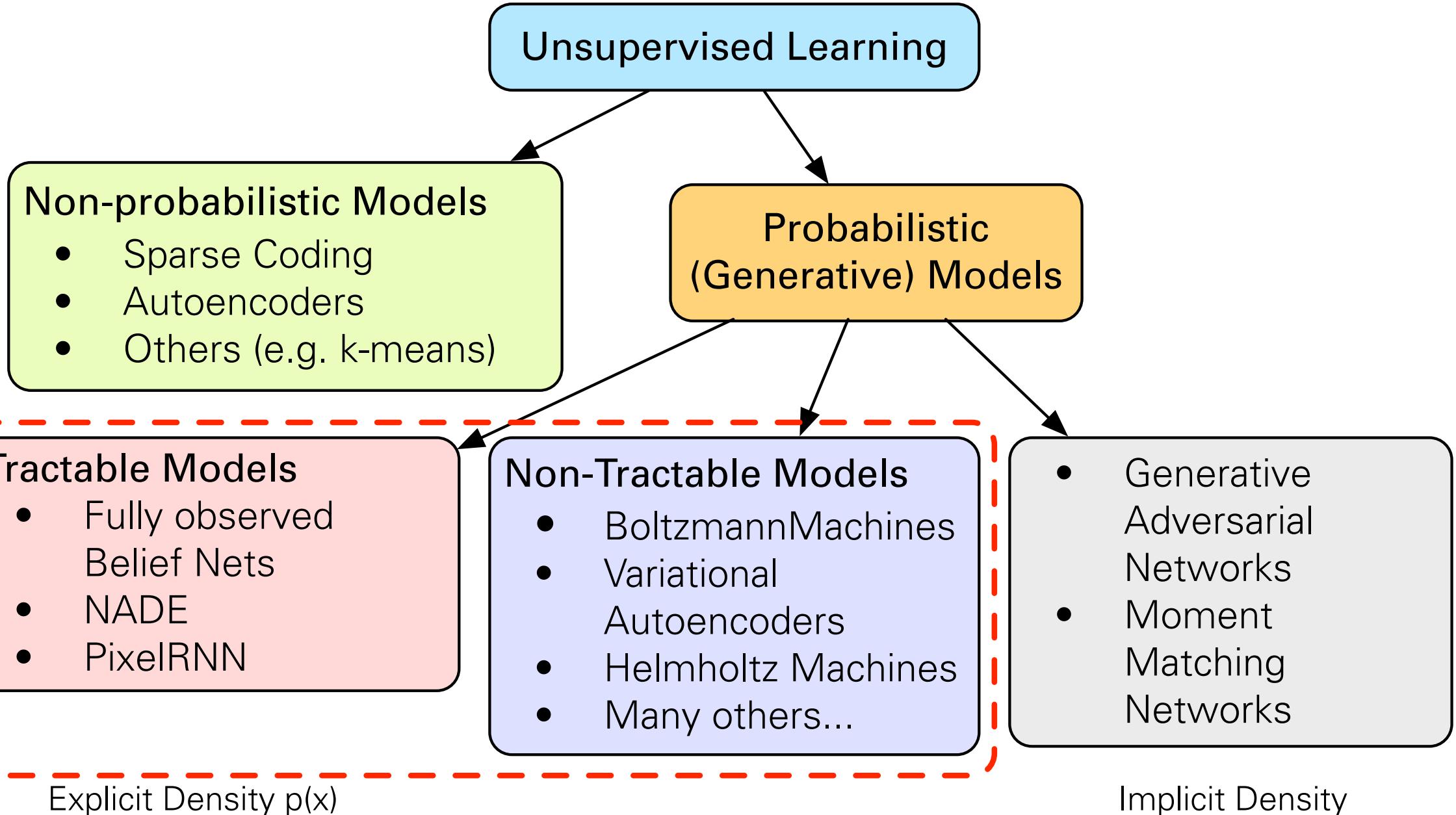
$$\begin{aligned} \mathbf{x}_i &\sim P_{\text{data}} \\ i &= 1, 2, \dots, n \end{aligned}$$



- We want to learn a probability distribution  $p(x)$  over images  $x$  s.t.
  - **Generation:** If we sample  $x_{\text{new}} \sim p(x)$ ,  $x_{\text{new}}$  should look like a dog (sampling)
  - **Density estimation:**  $p(x)$  should be high if  $x$  looks like a dog, and low otherwise (anomaly detection)
  - **Unsupervised representation learning:** We should be able to learn what these images have in common, e.g., ears, tail, etc. (features)

# Why Unsupervised Learning?

- Given high-dimensional data  $X = (x_1, \dots, x_n)$  we want to find a low-dimensional model characterizing the population.
- Recent progress mostly in supervised DL
- Real challenges for unsupervised DL
- Potential benefits:
  - **Exploit tons of unlabeled data**
  - Answer new questions about the variables observed
  - Regularizer – transfer learning – domain adaptation
  - Easier optimization (divide and conquer)
  - Joint (structured) outputs



# Unsupervised Learning

- Basic Building Blocks:
  - Sparse Coding
  - Autoencoders
- Autoregressive Generative Models
- Generative Adversarial Networks
- Variational Autoencoders
- Normalizing Flow Models

# Sparse Coding

- Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).
- **Objective:** Given a set of input data vectors  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ , learn a dictionary of bases, such that:

$$\mathbf{x}_n = \sum_{k=1}^K a_{nk} \phi_k$$

Sparse: mostly zeros

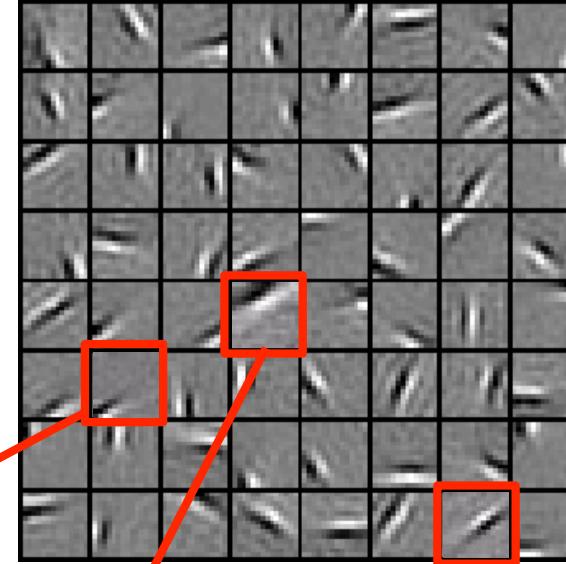
- Each data vector is represented as a sparse linear combination of bases.

# Sparse Coding

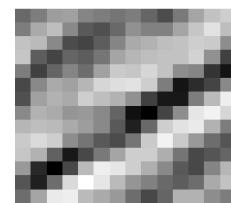
Natural Images



Learned bases: "Edges"



New example



$$x = 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{65}$$

[0.0, 0.0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Sparse Coding: Training

- Input image patches:  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^D$
- Learn dictionary of bases:  $\phi_1, \phi_2, \dots, \phi_K \in \mathbb{R}^D$

$$\min_{\mathbf{a}, \phi} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$



Reconstruction error      Sparsity penalty

- Alternating Optimization:
  1. Fix dictionary of bases and solve for activations  $\mathbf{a}$  (a standard Lasso problem).
  2. Fix activations  $\mathbf{a}$ , optimize the dictionary of bases (convex QP problem).

# Sparse Coding: Testing Time

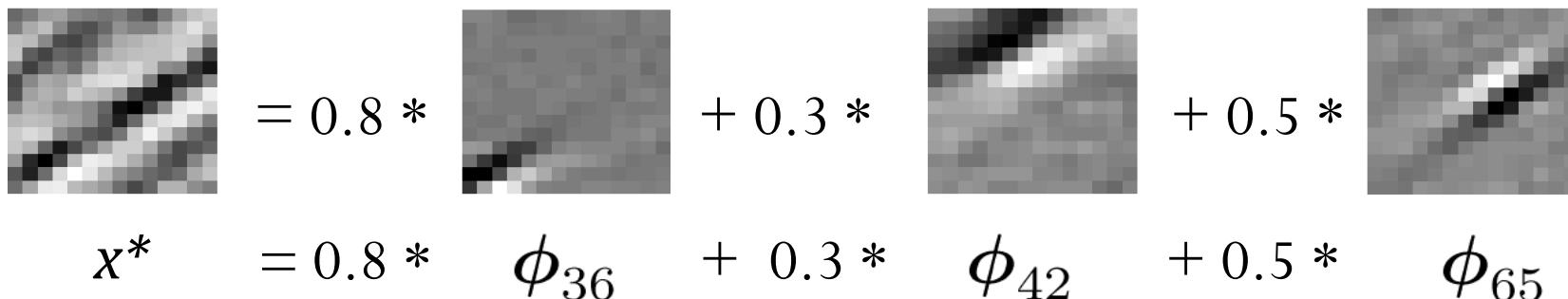
- Input: a new image patch  $\mathbf{x}^*$ , and  $K$  learned bases  $\phi_1, \phi_2, \dots, \phi_K$
- Output: sparse representation  $\mathbf{a}$  of an image patch  $\mathbf{x}^*$ .

$$\min_{\mathbf{a}} \left\| \mathbf{x}^* - \sum_{k=1}^K a_k \phi_k \right\|_2^2 + \lambda \sum_{k=1}^K |a_k|$$

# Sparse Coding: Testing Time

- Input: a new image patch  $x^*$ , and K learned bases  $\phi_1, \phi_2, \dots, \phi_K$
- Output: sparse representation  $\mathbf{a}$  of an image patch  $x^*$ .

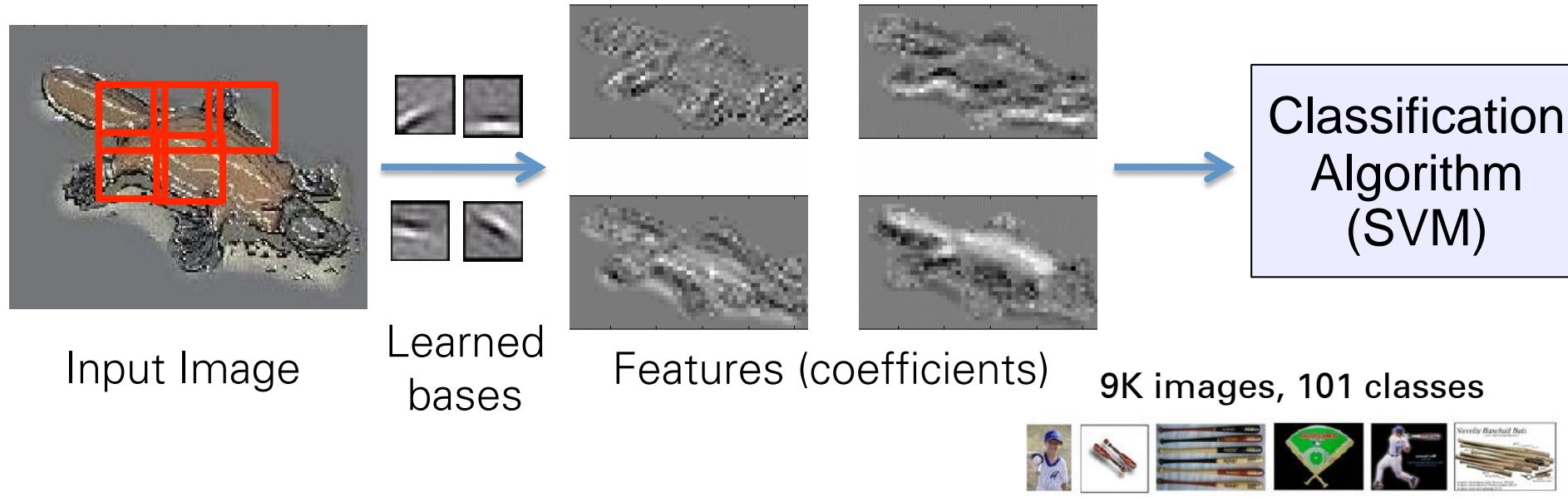
$$\min_{\mathbf{a}} \left\| \mathbf{x}^* - \sum_{k=1}^K a_k \phi_k \right\|_2^2 + \lambda \sum_{k=1}^K |a_k|$$

$$x^* = 0.8 * \phi_{36} + 0.3 * \phi_{42} + 0.5 * \phi_{65}$$


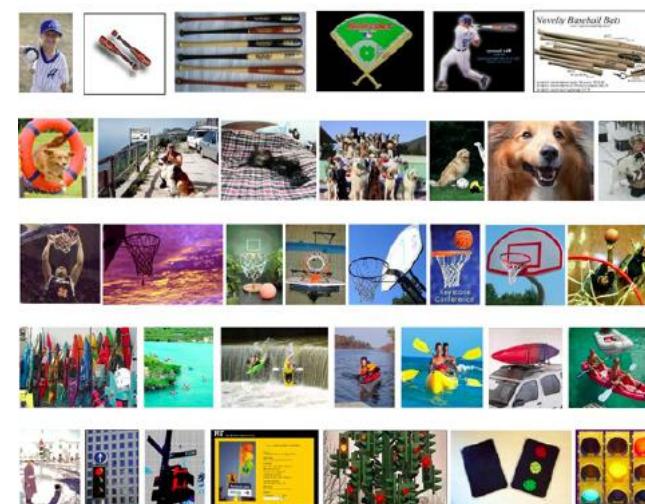
[0.0, 0.0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Image Classification

- Evaluated on Caltech101 object category dataset.



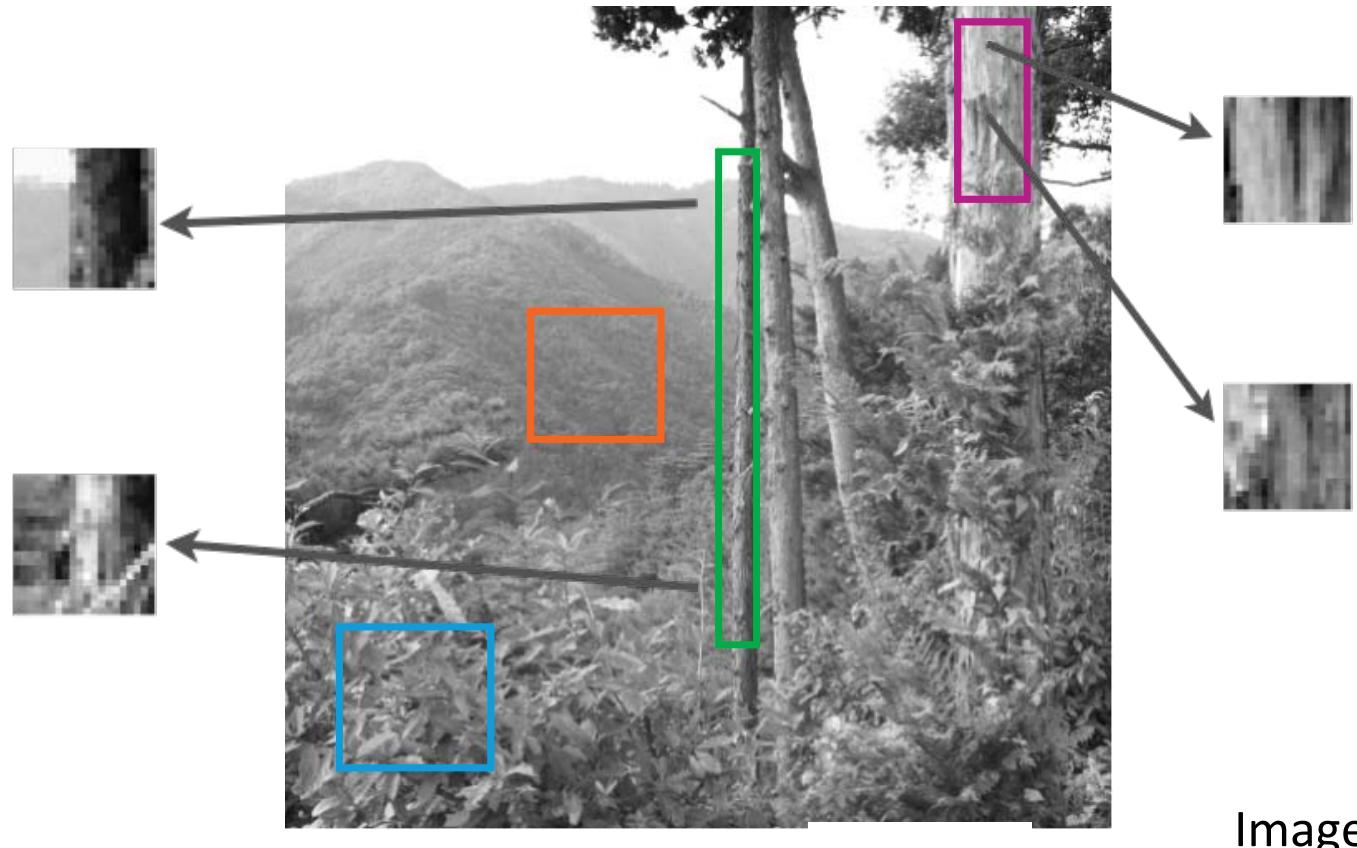
Algorithm	Accuracy
Baseline (Fei-Fei et al., 2004)	16%
PCA	37%
<b>Sparse Coding</b>	<b>47%</b>



(Lee, Battle, Raina, Ng, NIPS 2007) 40

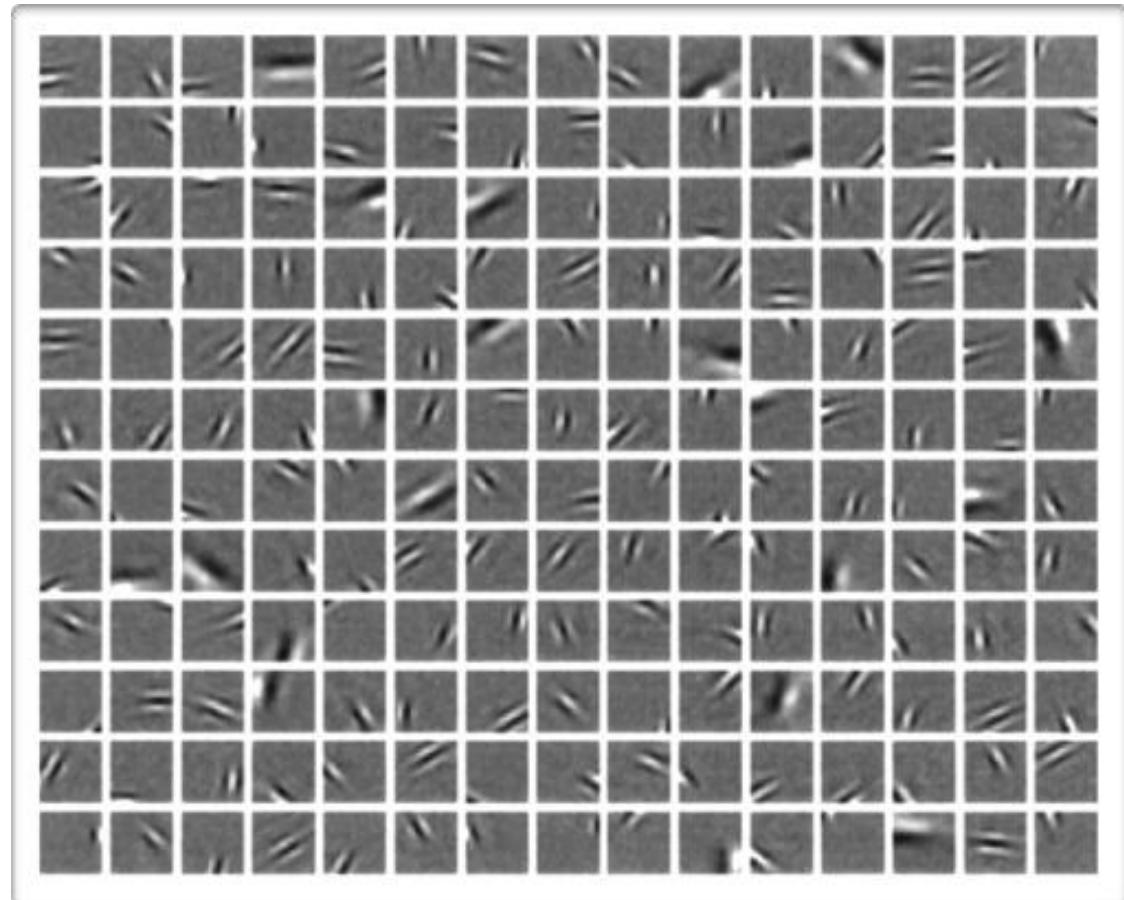
# Modeling Image Patches

- Natural image patches:
  - small **image regions** extracted from an image of nature (forest, grass, ...)



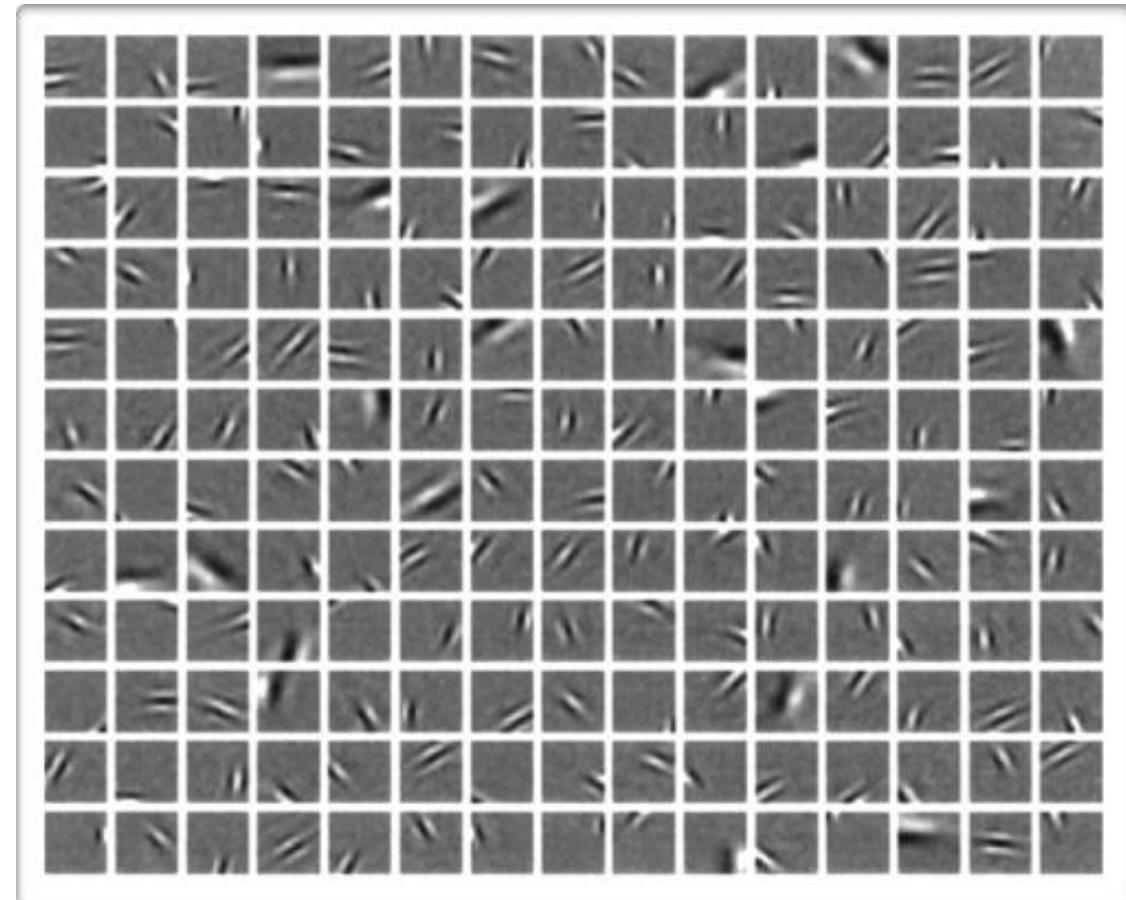
# Relationship to V1

- When trained on natural image patches
  - the dictionary columns ("atoms") look like **edge detectors**
  - each atom is tuned to a particular **position, orientation** and **spatial frequency**
  - V1 neurons in the mammalian brain have a similar behavior



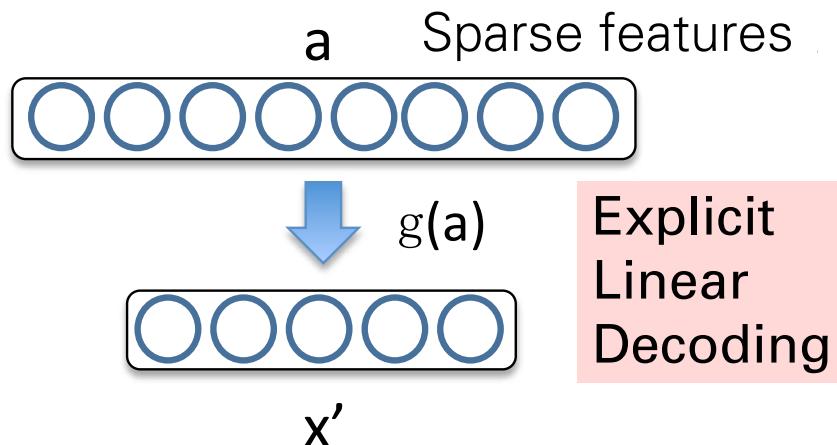
# Relationship to V1

- Suggests that the brain might be learning a sparse code of visual stimulus
  - Since then, many other models have been shown to learn similar features
  - they usually all incorporate a notion of sparsity



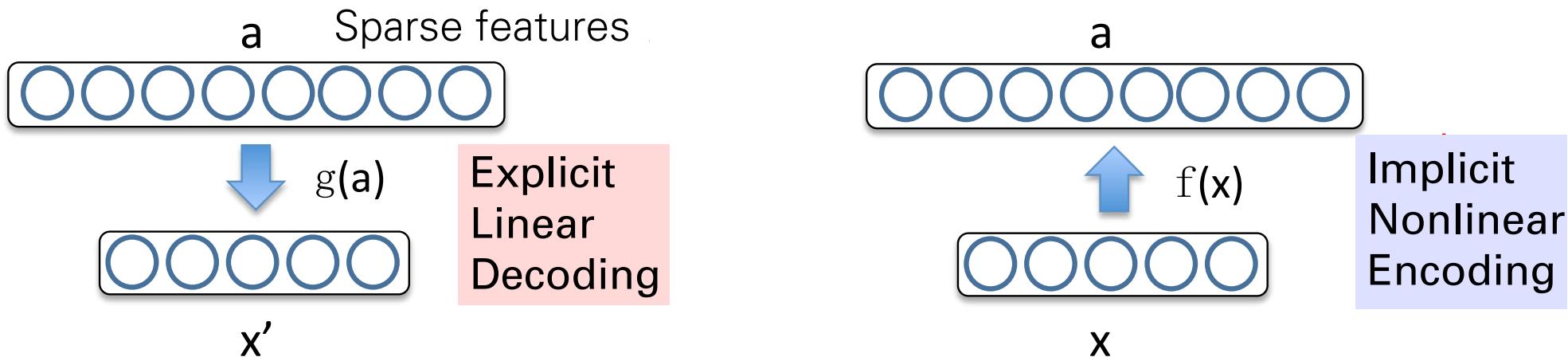
# Interpreting Sparse Coding

$$\min_{\mathbf{a}, \boldsymbol{\phi}} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \boldsymbol{\phi}_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$



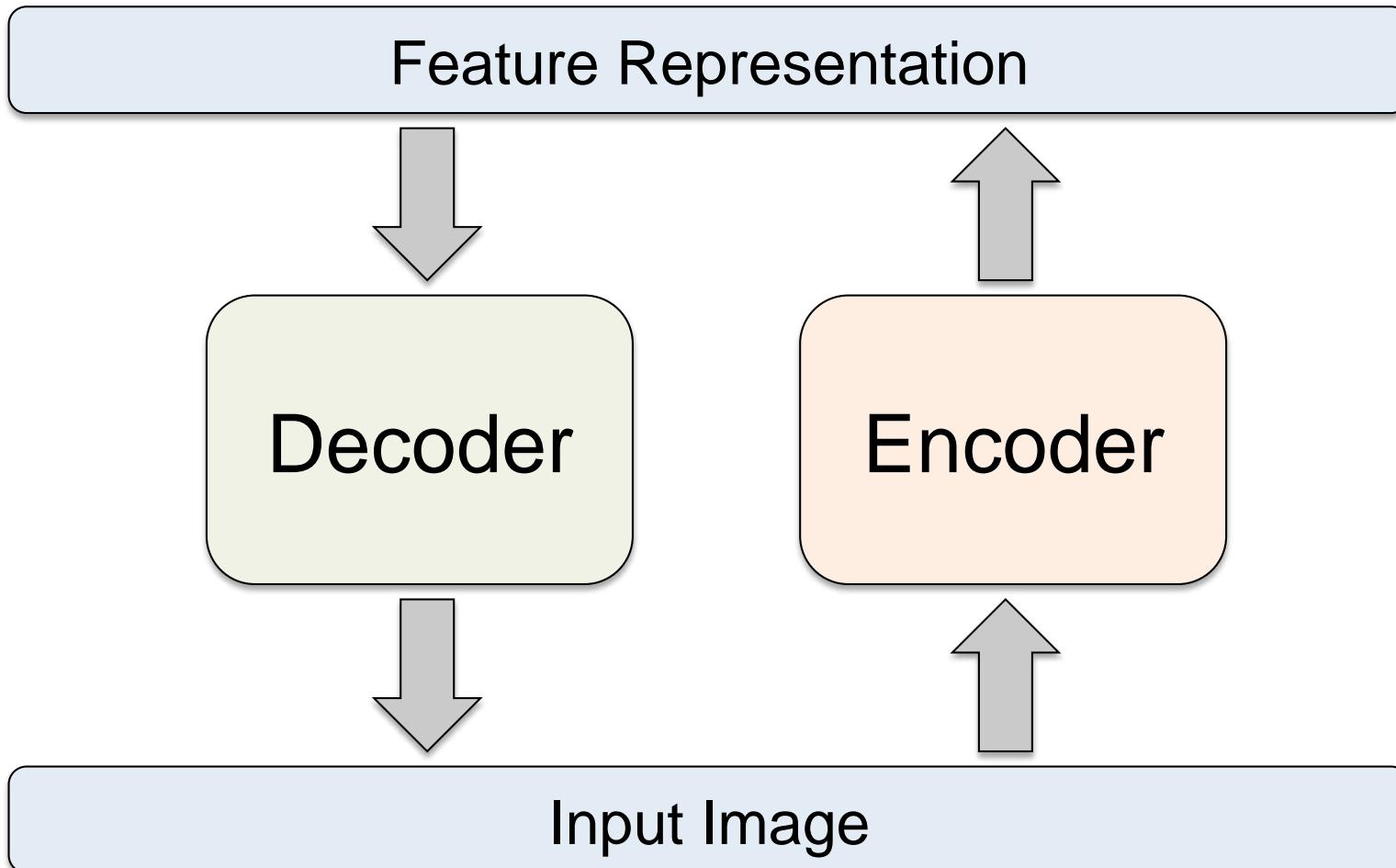
# Interpreting Sparse Coding

$$\min_{\mathbf{a}, \boldsymbol{\phi}} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \boldsymbol{\phi}_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$

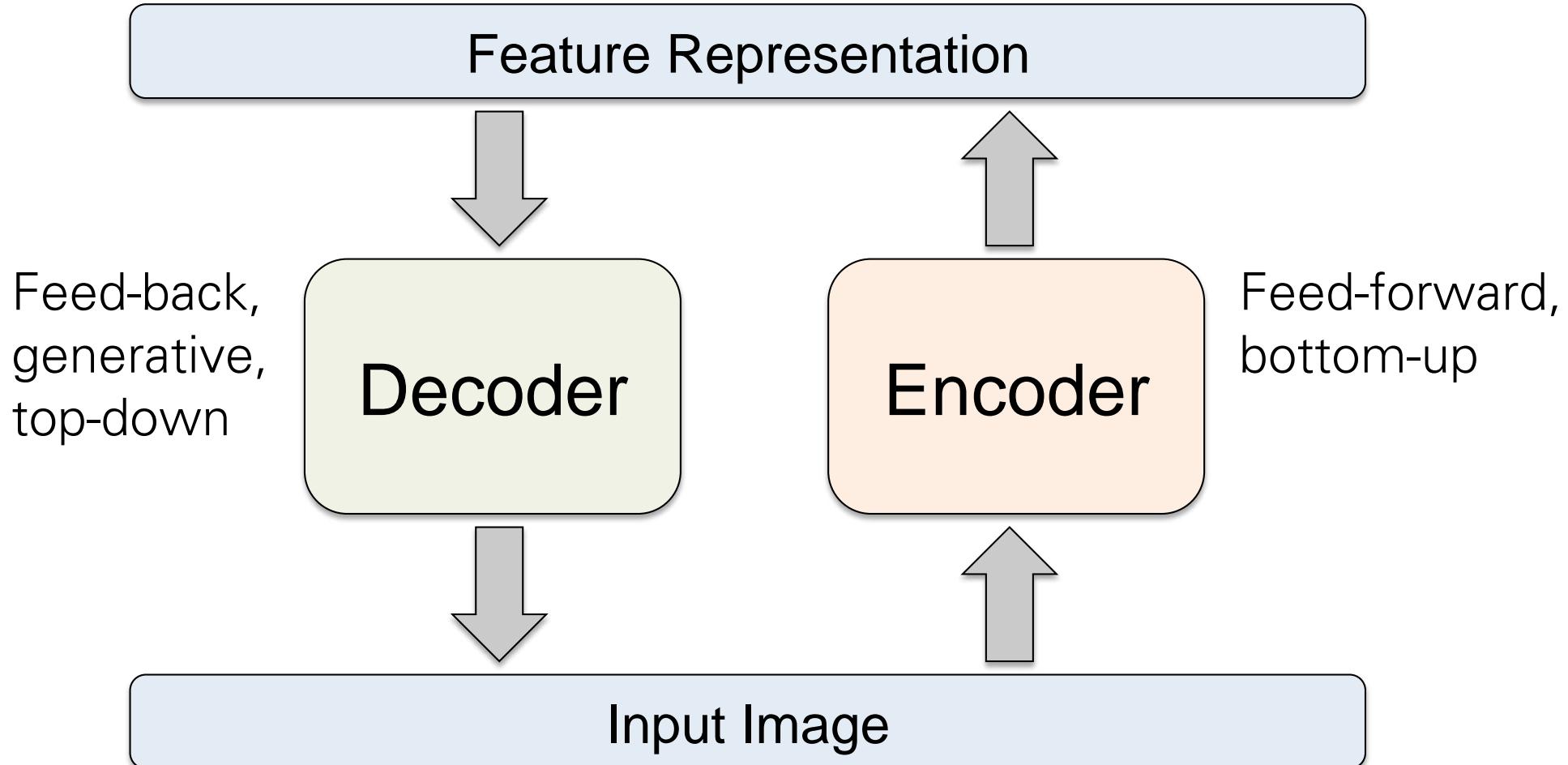


- Sparse, over-complete representation **a**.
- **Encoding**  $\mathbf{a} = f(\mathbf{x})$  is implicit and nonlinear function of  $\mathbf{x}$ .
- **Reconstruction** (or decoding)  $\mathbf{x}' = g(\mathbf{a})$  is linear and explicit.

# Autoencoder

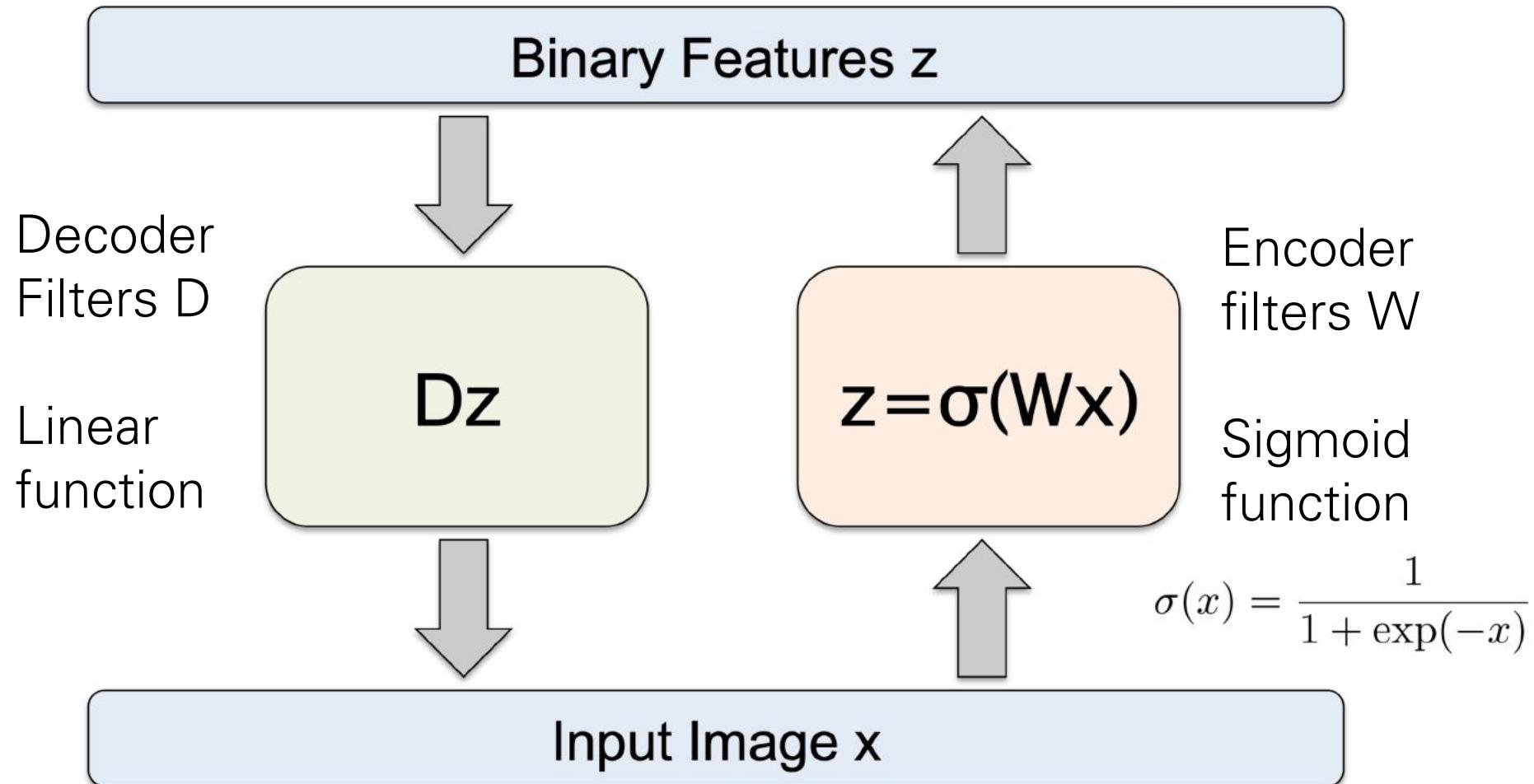


# Autoencoder

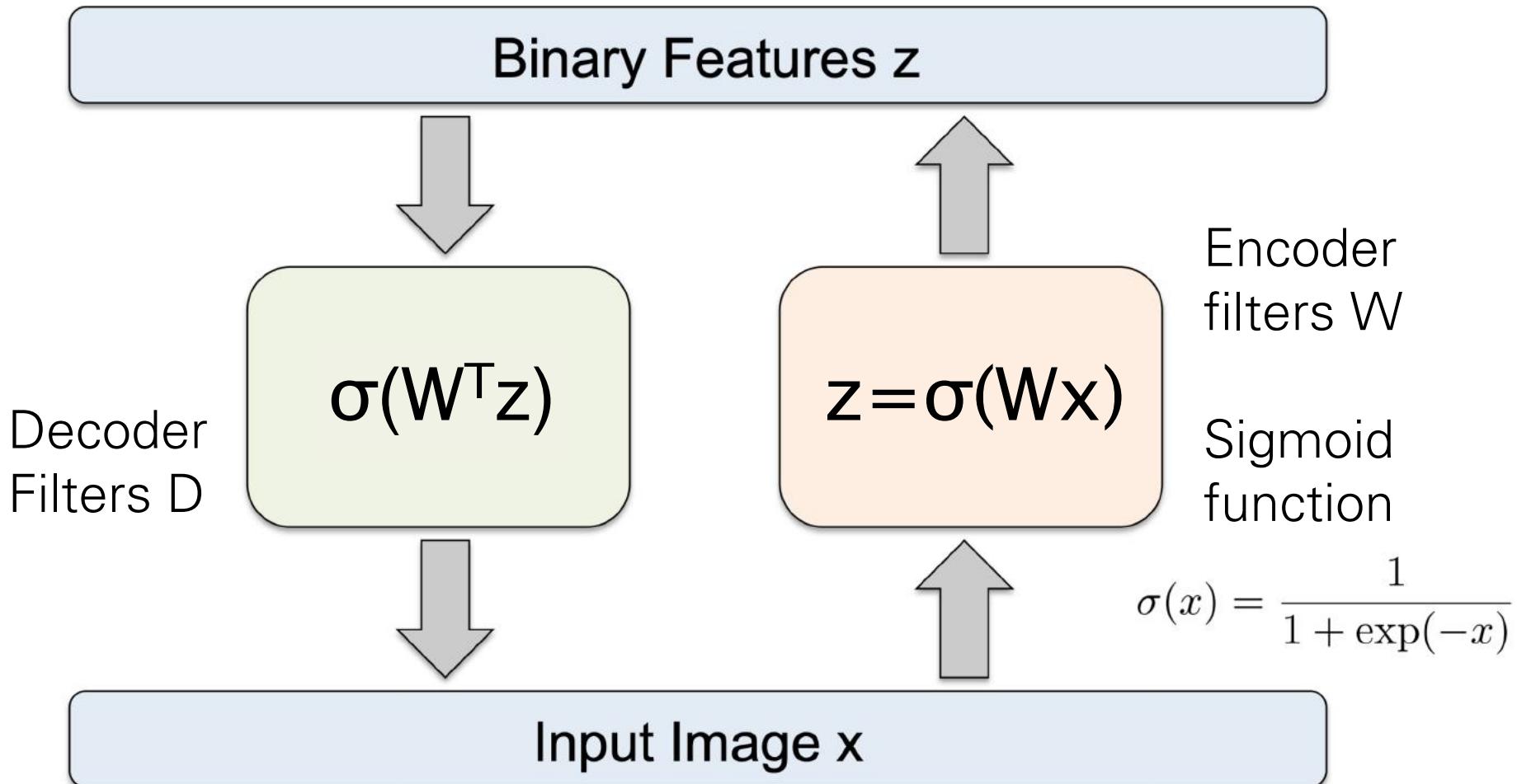


- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

# Autoencoder



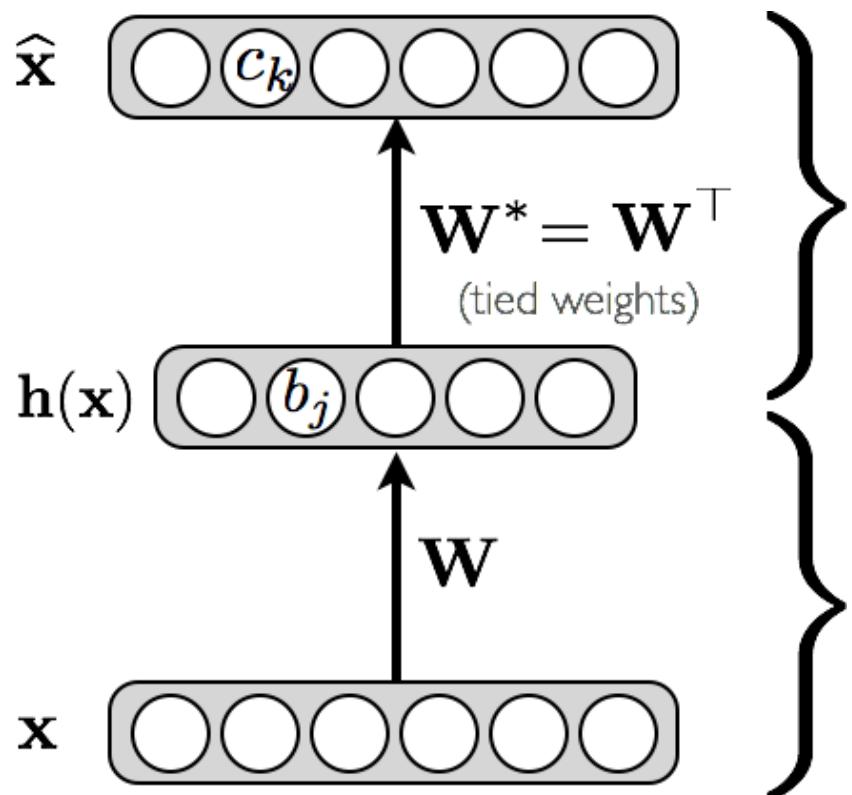
# Autoencoder



- Need additional constraints to avoid learning an identity.
- Relates to Restricted Boltzmann Machines (later).

# Autoencoder

- Feed-forward neural network trained to reproduce its input at the output layer



## Decoder

$$\begin{aligned}\hat{x} &= o(\hat{a}(x)) \\ &= \text{sigm}(\mathbf{c} + \mathbf{W}^* \mathbf{h}(x))\end{aligned}$$

for binary units

## Encoder

$$\begin{aligned}h(x) &= g(a(x)) \\ &= \text{sigm}(\mathbf{b} + \mathbf{W}x)\end{aligned}$$

# Loss Function

- **Loss function for binary inputs**

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

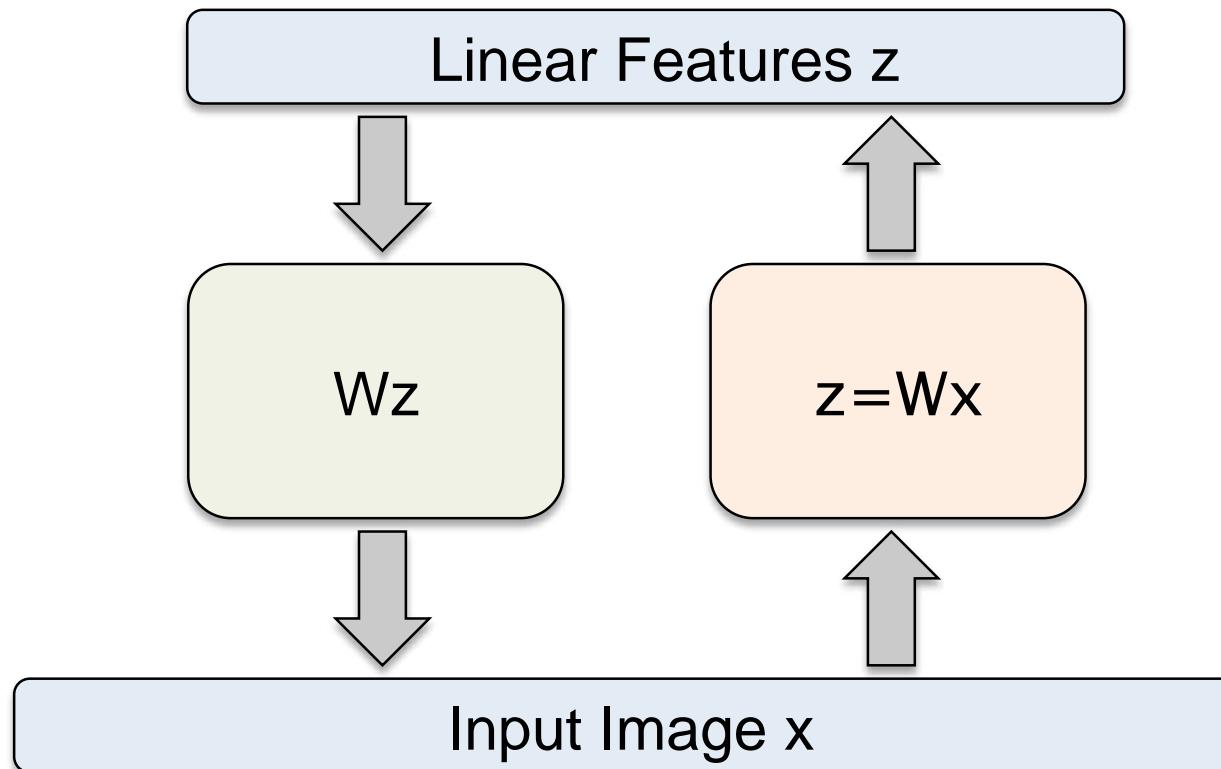
- Cross-entropy error function (reconstruction loss)  $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$

- **Loss function for real-valued inputs**

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

- reconstruction loss)
- we use a linear activation function at the output

# Autoencoder



- With nonlinear hidden units, we have a nonlinear generalization of PCA.
- If the **hidden and output layers are linear**, it will learn hidden units that are a linear function of the data and minimize the squared error.
- The K hidden units will span the same space as the first k principal components. The weight vectors may not be orthogonal.

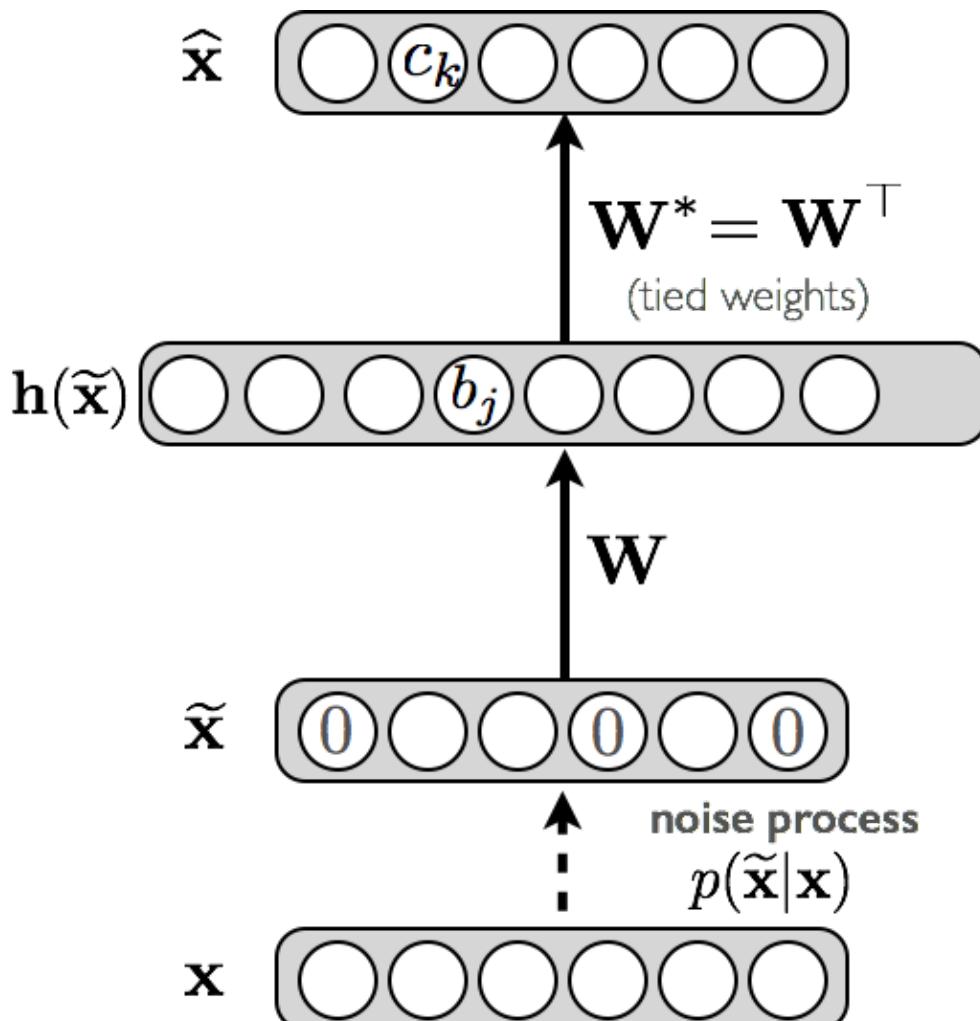
# Denoising Autoencoder

- **Idea:** Representation should be robust to introduction of noise:

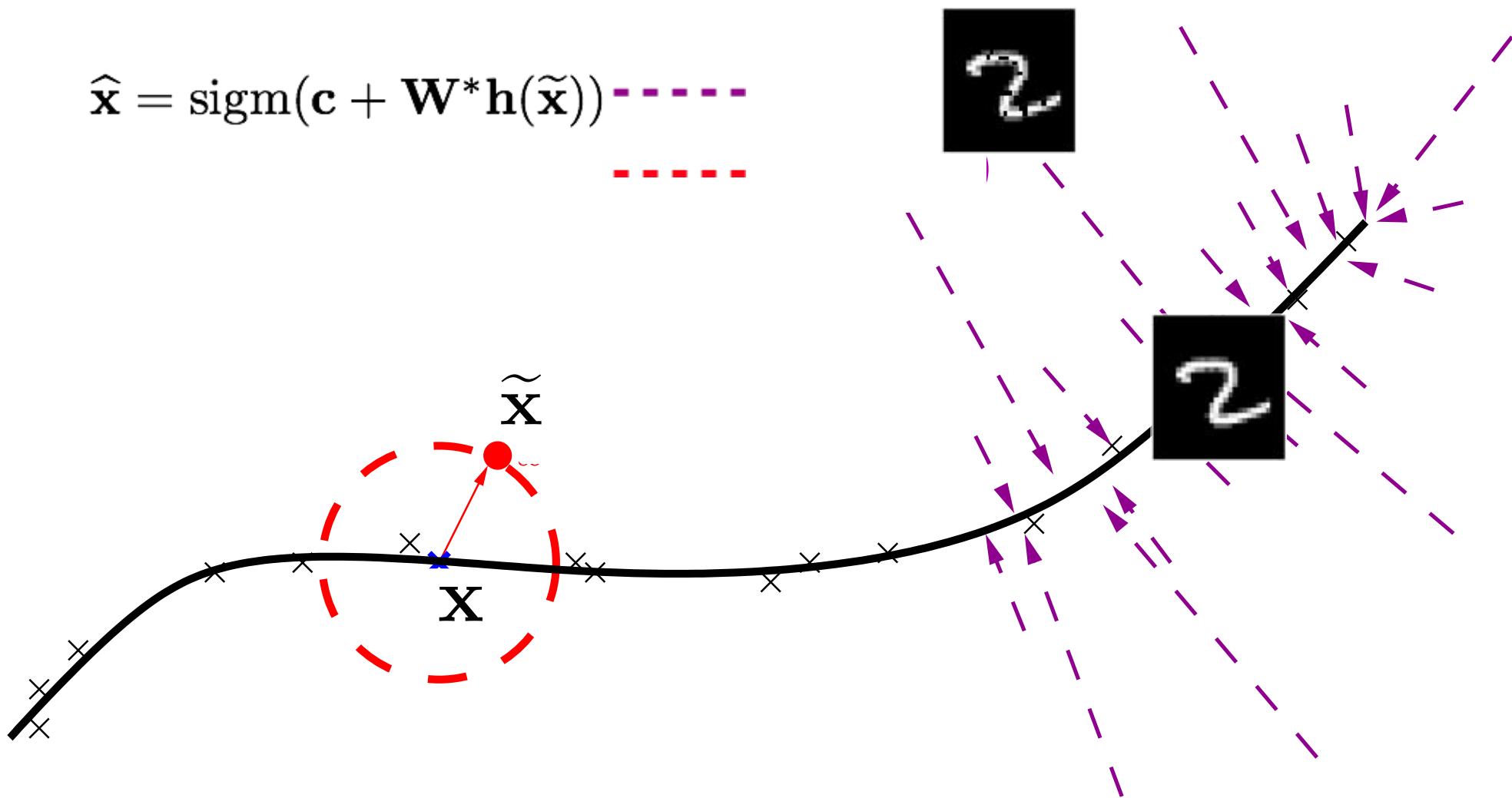
- random assignment of subset of inputs to 0, with probability  $\nu$
- Similar to dropouts on the input layer
- Gaussian additive noise

- **Reconstruction**  $\hat{\mathbf{x}}$  computed from the corrupted input  $\tilde{\mathbf{x}}$

- **Loss function** compares  $\hat{\mathbf{x}}$  reconstruction with the noiseless input  $\mathbf{x}$

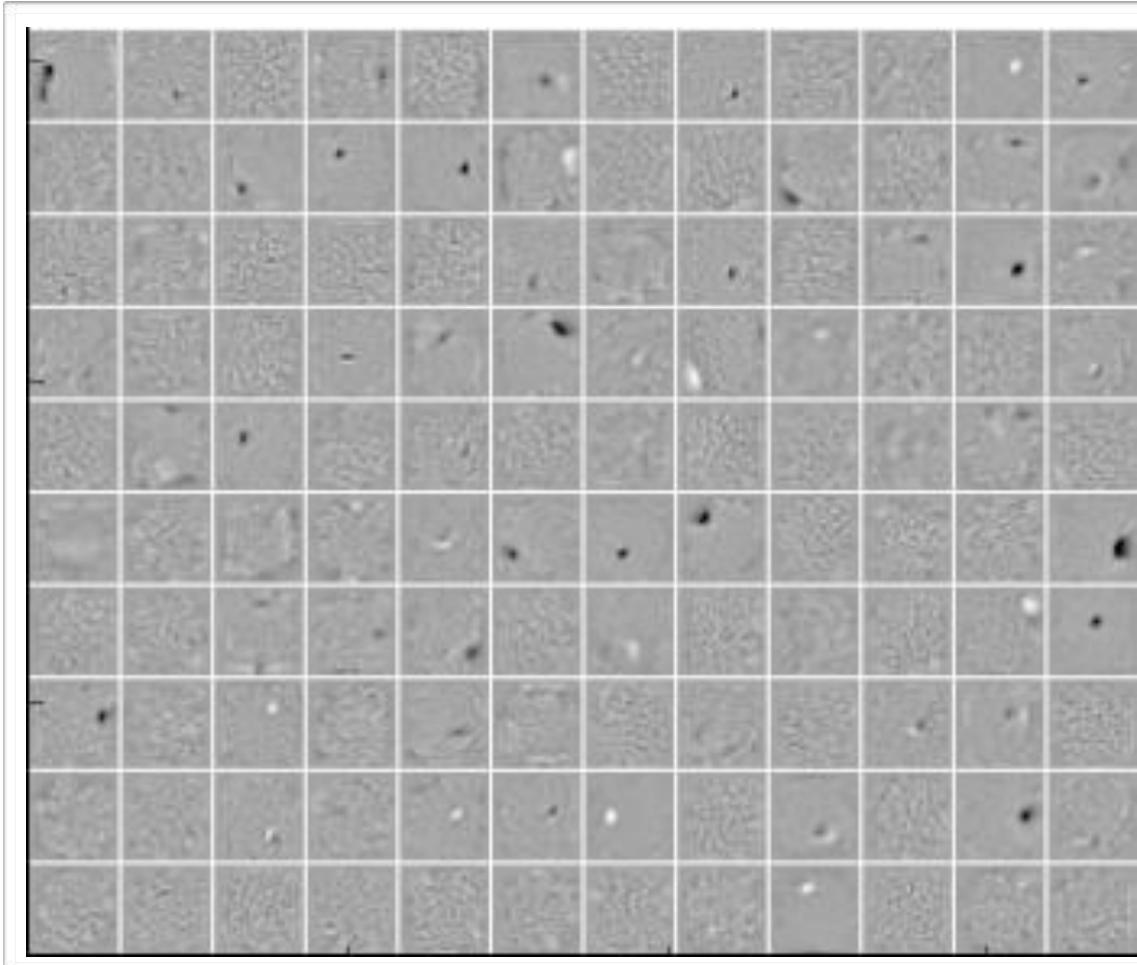


# Denoising Autoencoder

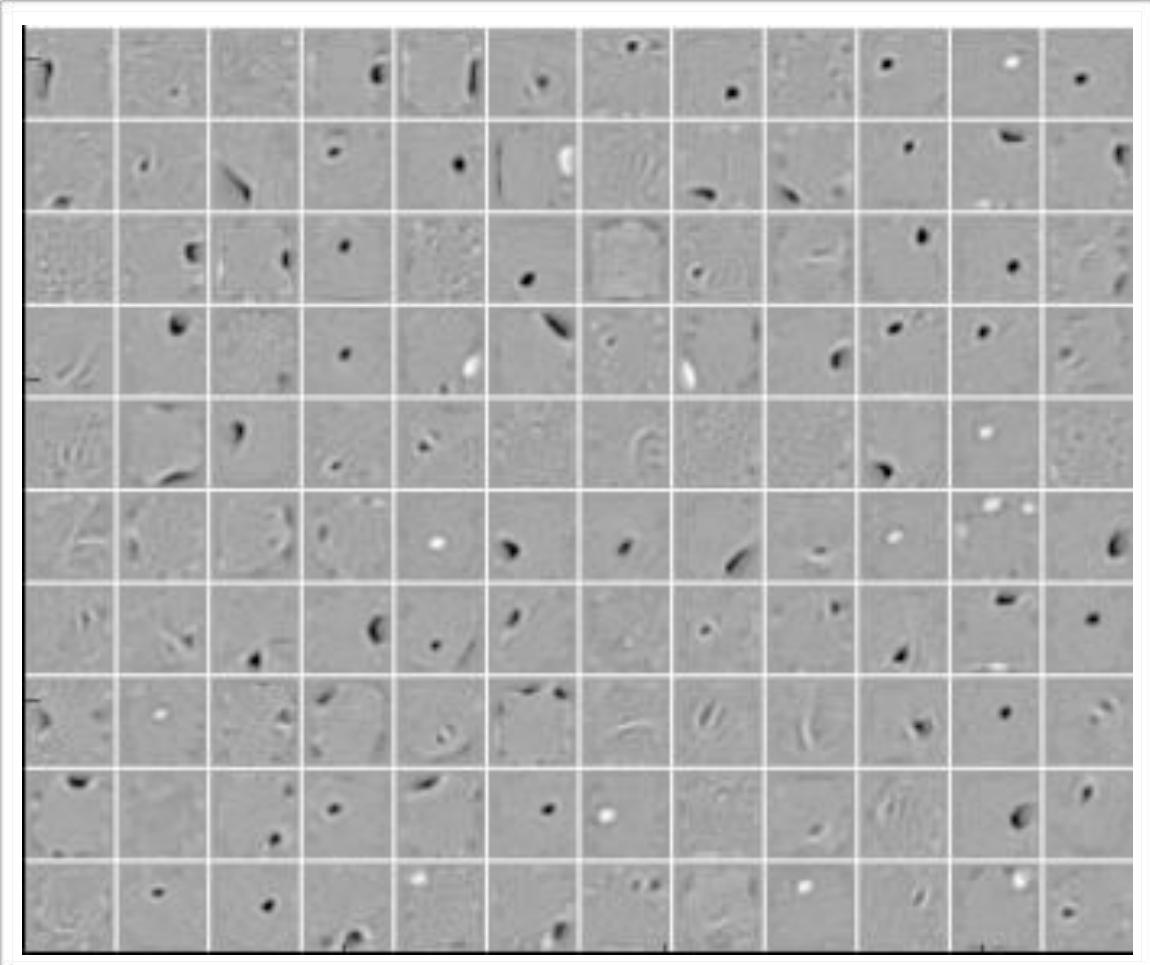


# Learned Filters

Non-corrupted

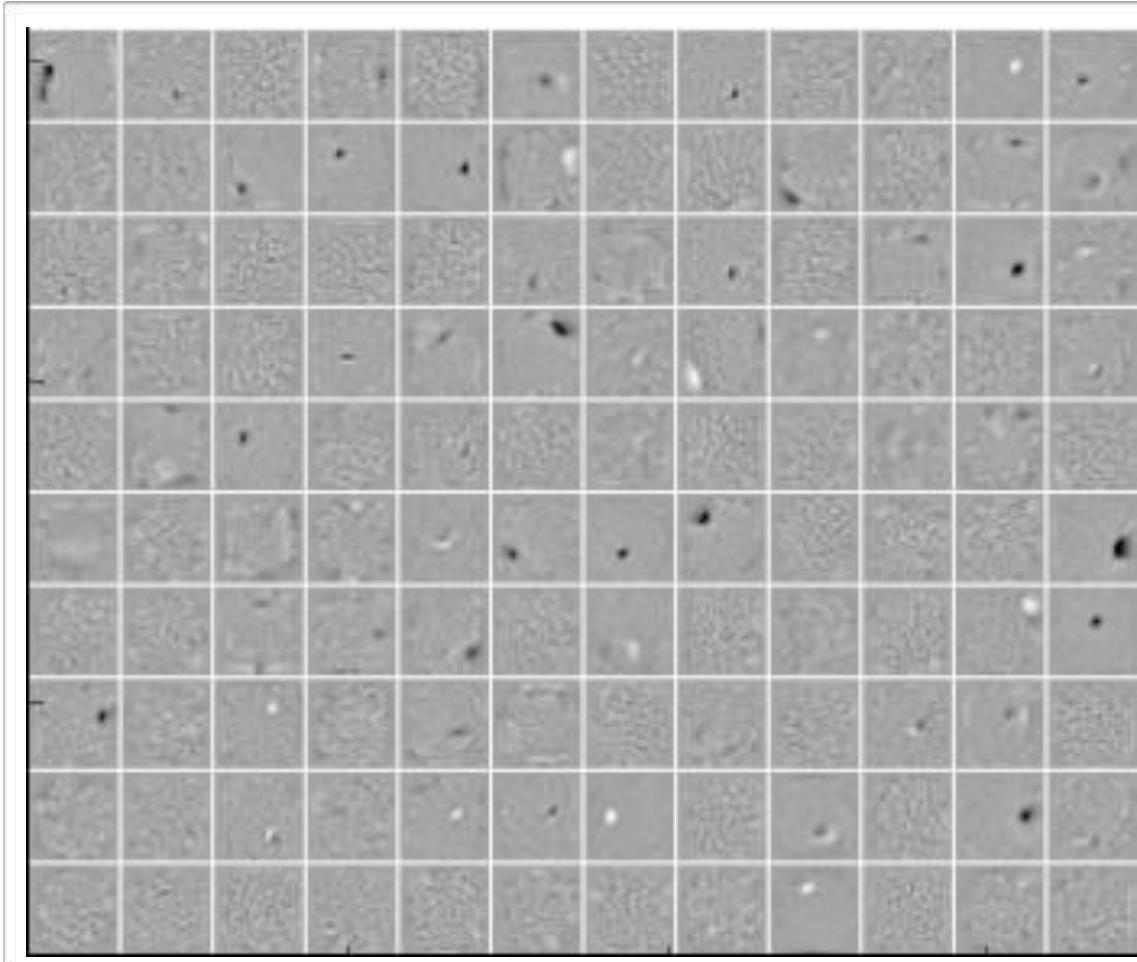


25% corrupted input

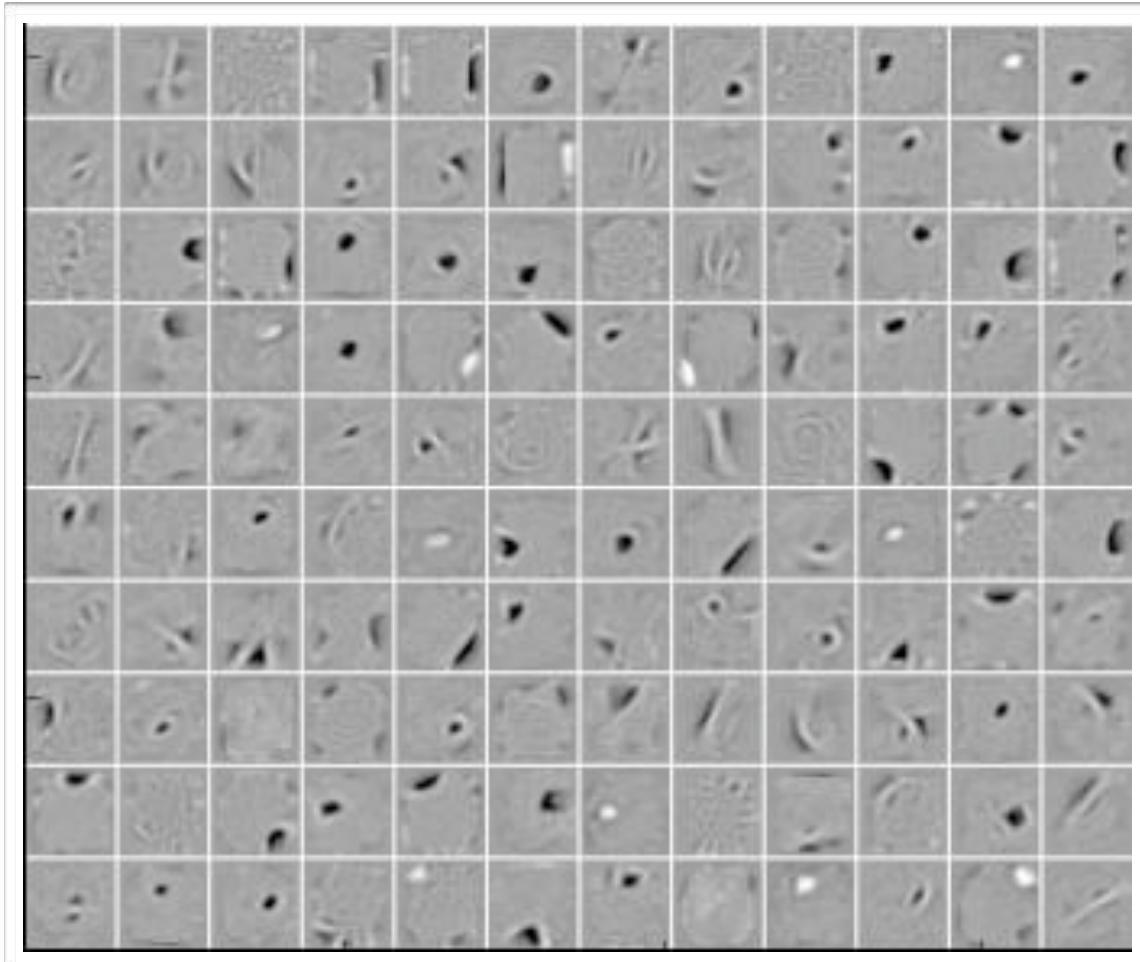


# Learned Filters

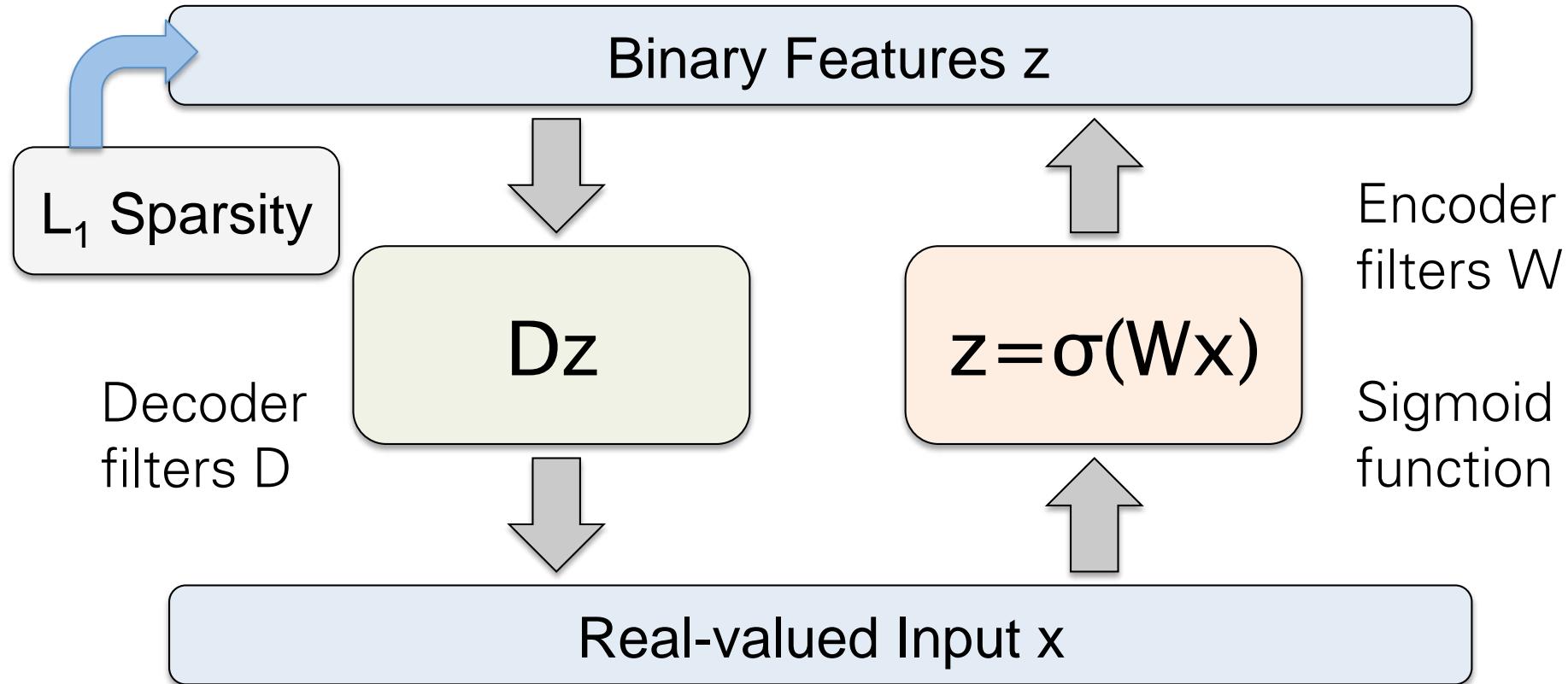
Non-corrupted



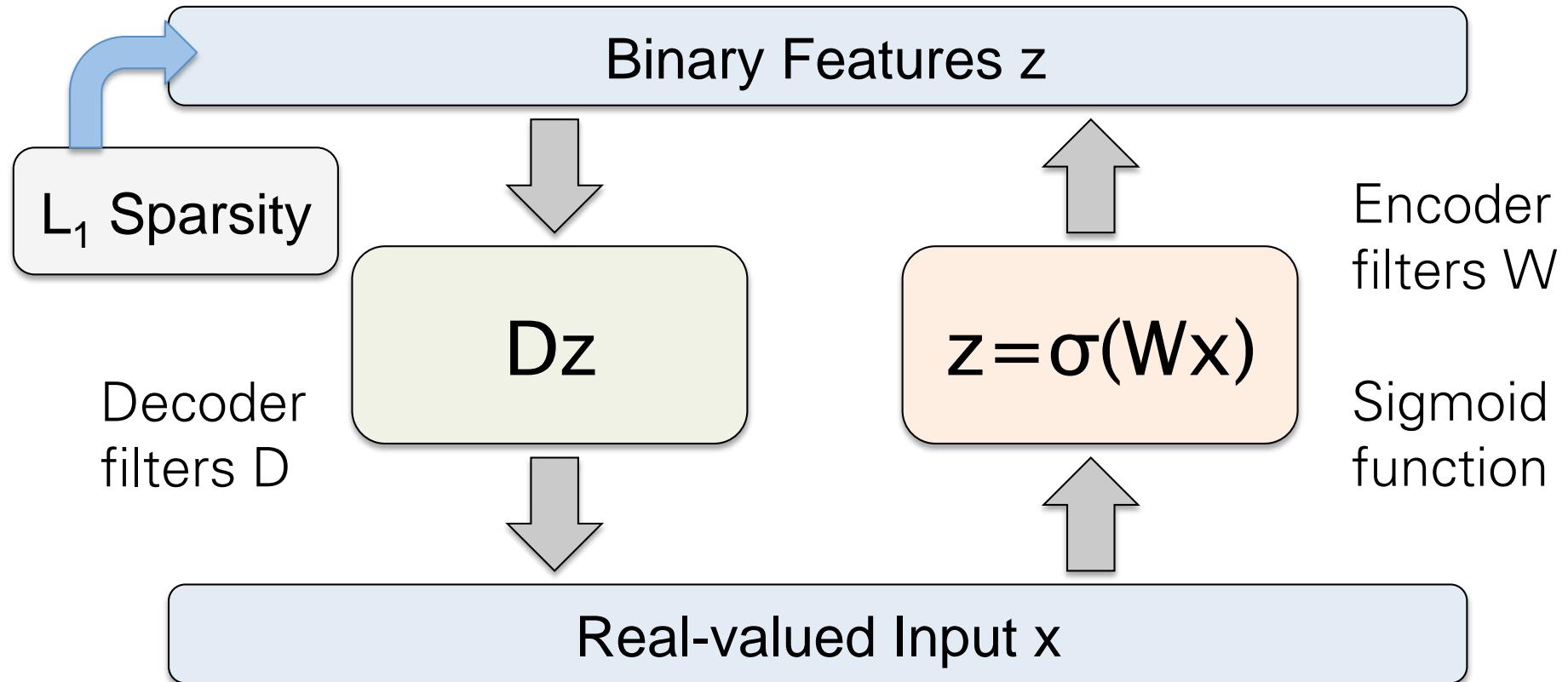
50% corrupted input



# Predictive Sparse Decomposition



# Predictive Sparse Decomposition



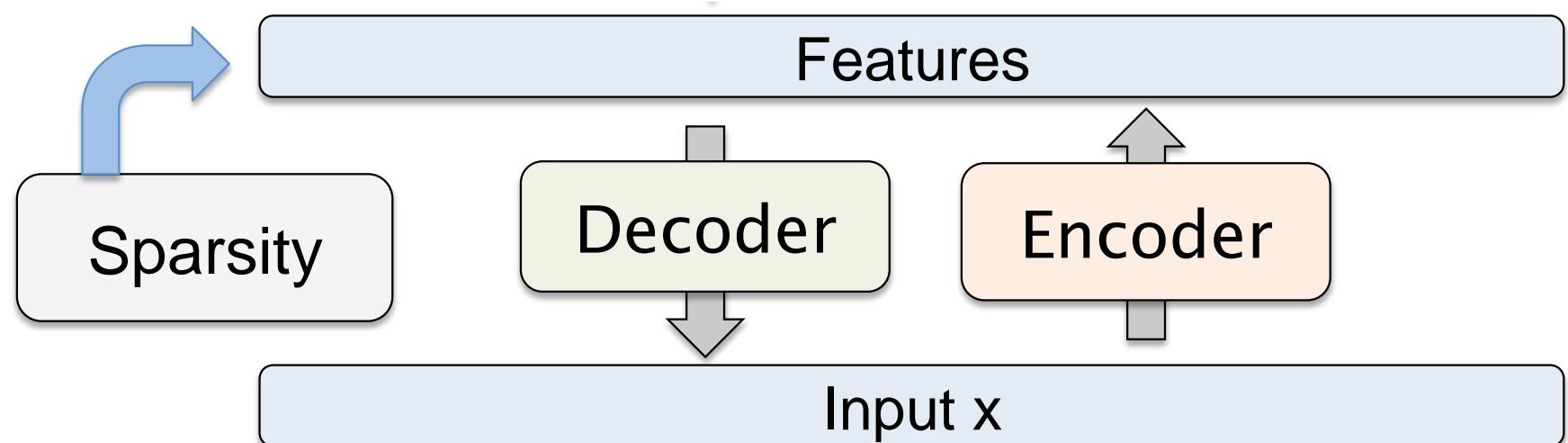
At training  
time

$$\min_{D, W, \mathbf{z}} ||D\mathbf{z} - \mathbf{x}||_2^2 + \lambda |\mathbf{z}|_1 + ||\sigma(W\mathbf{x}) - \mathbf{z}||_2^2$$

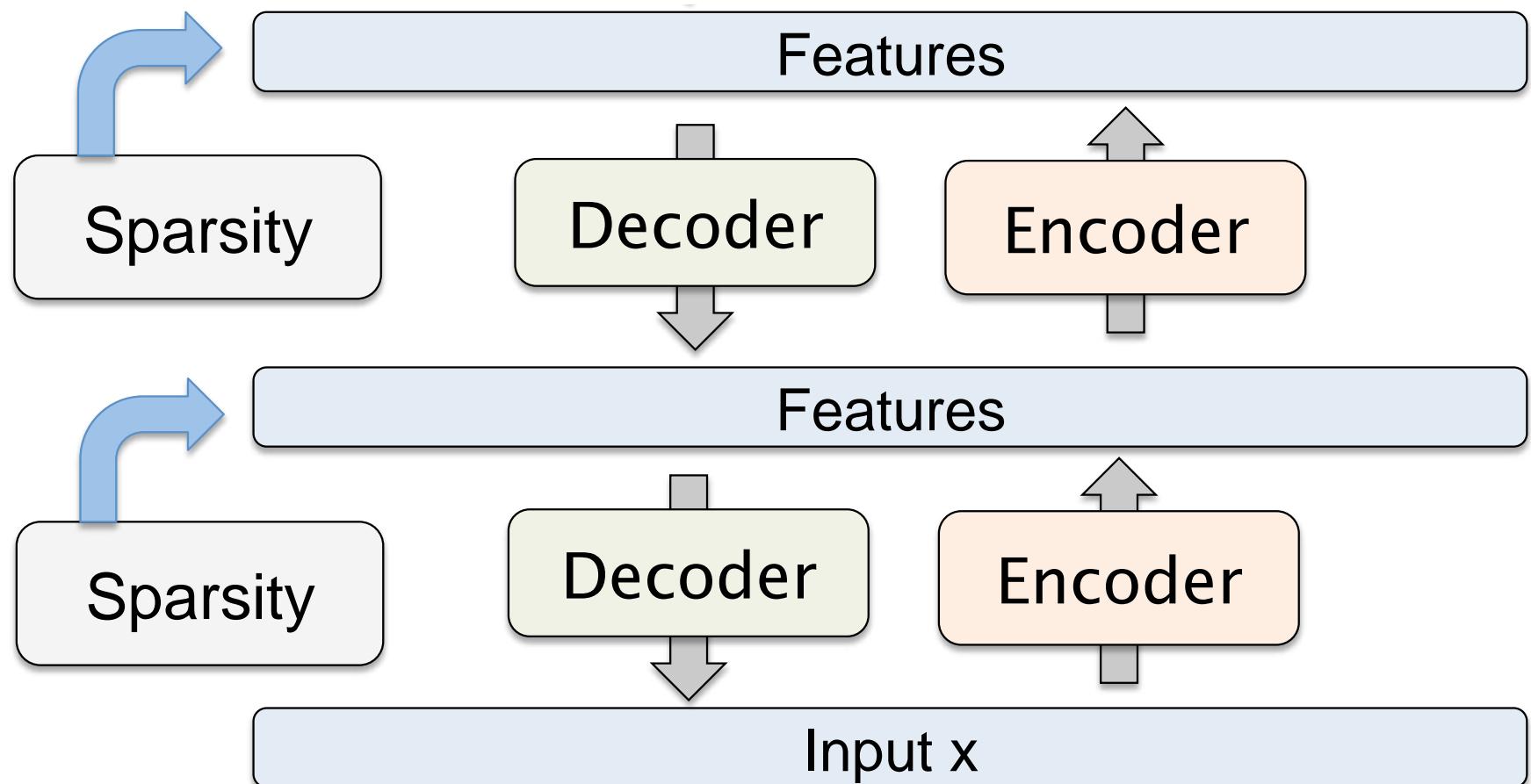
Decoder

Encoder

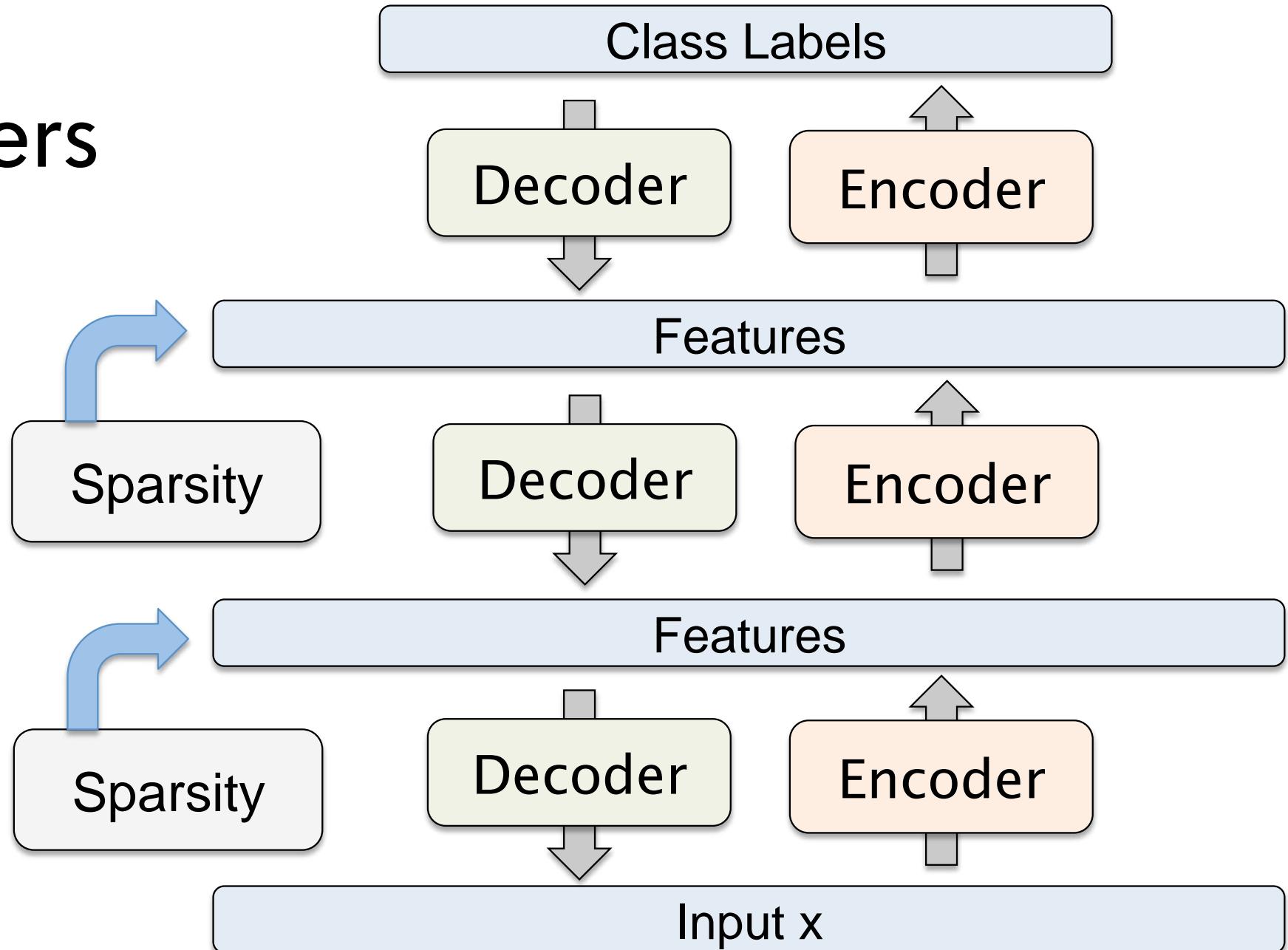
# Stacked Autoencoders



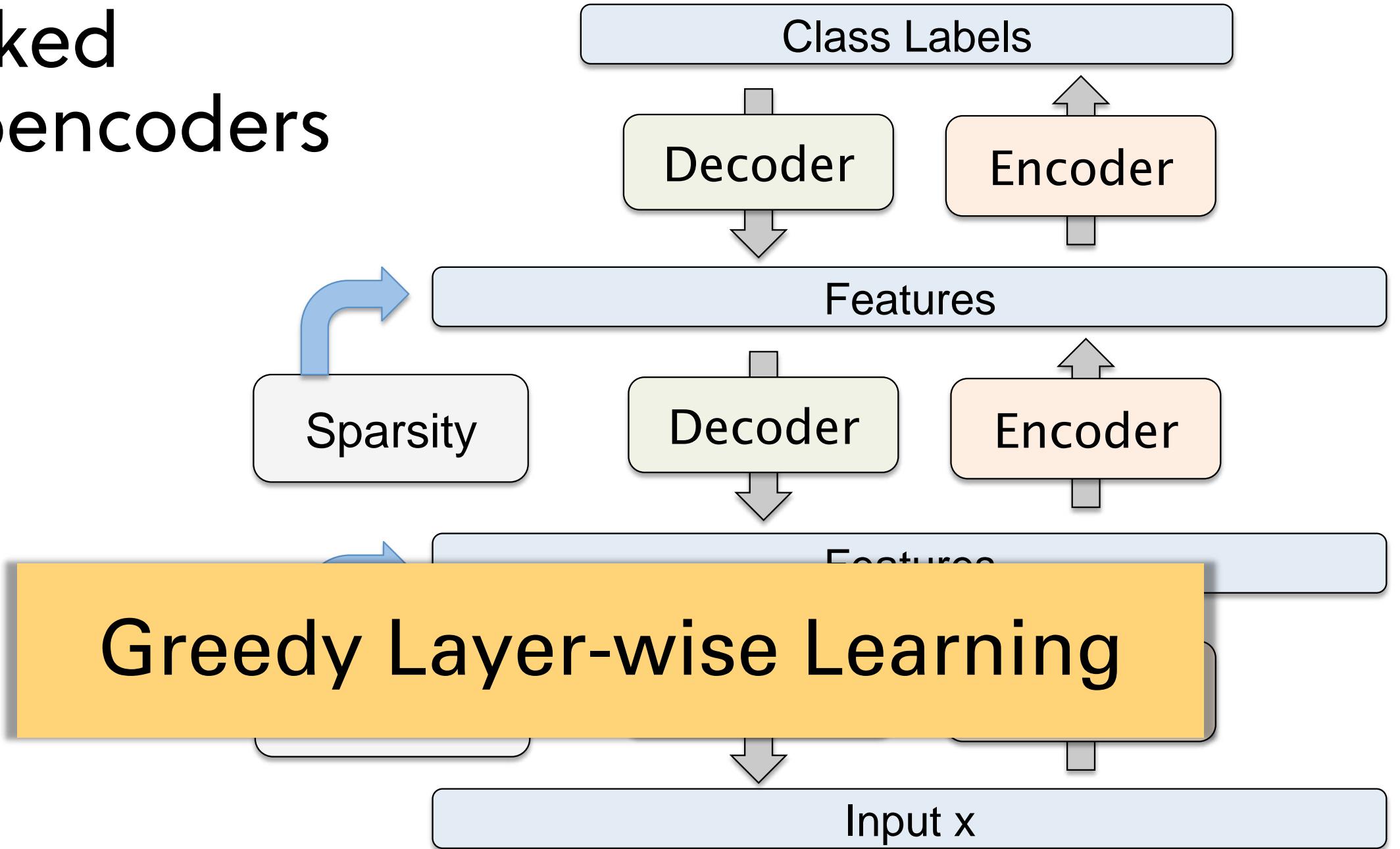
# Stacked Autoencoders



# Stacked Autoencoders

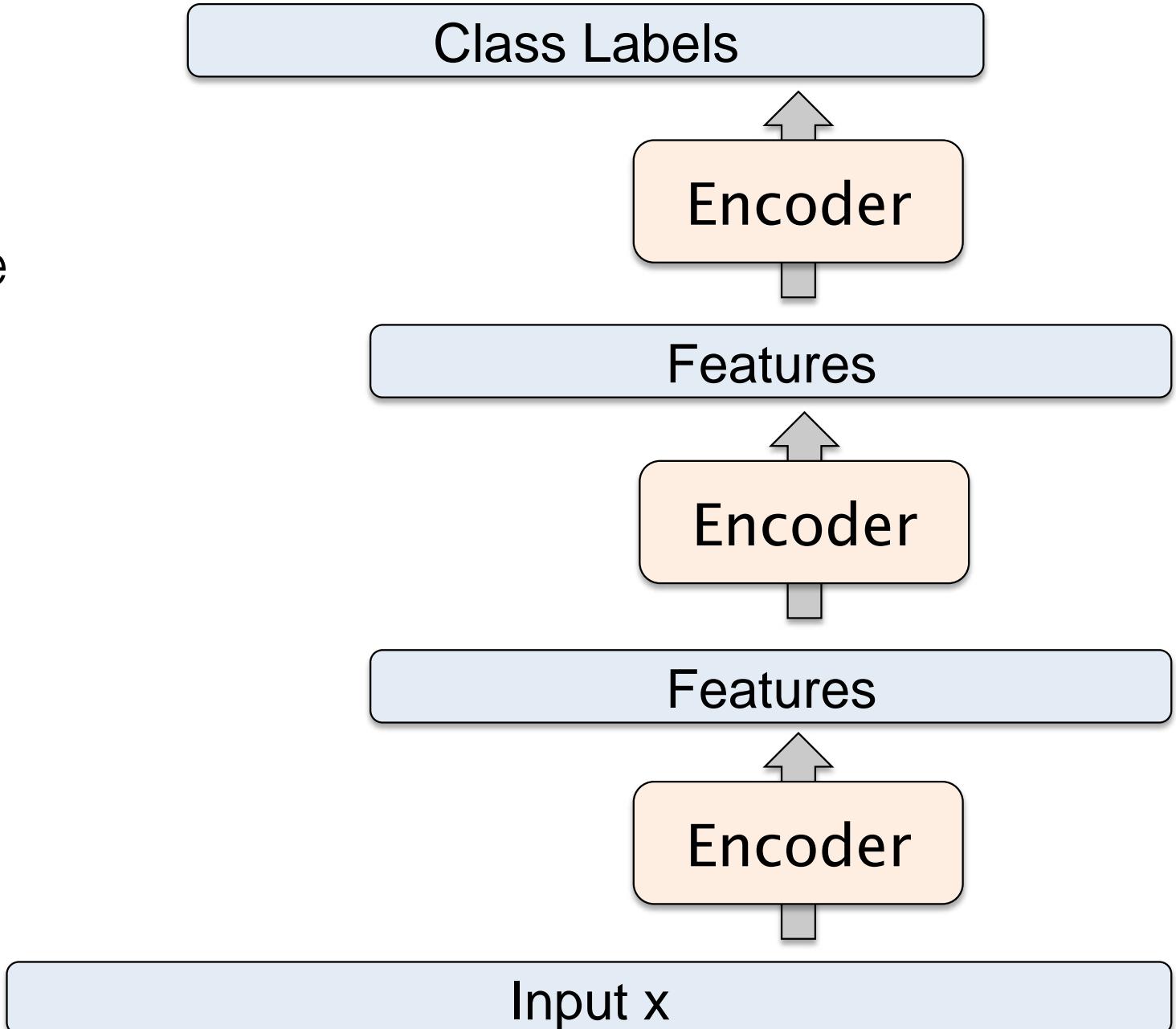


# Stacked Autoencoders



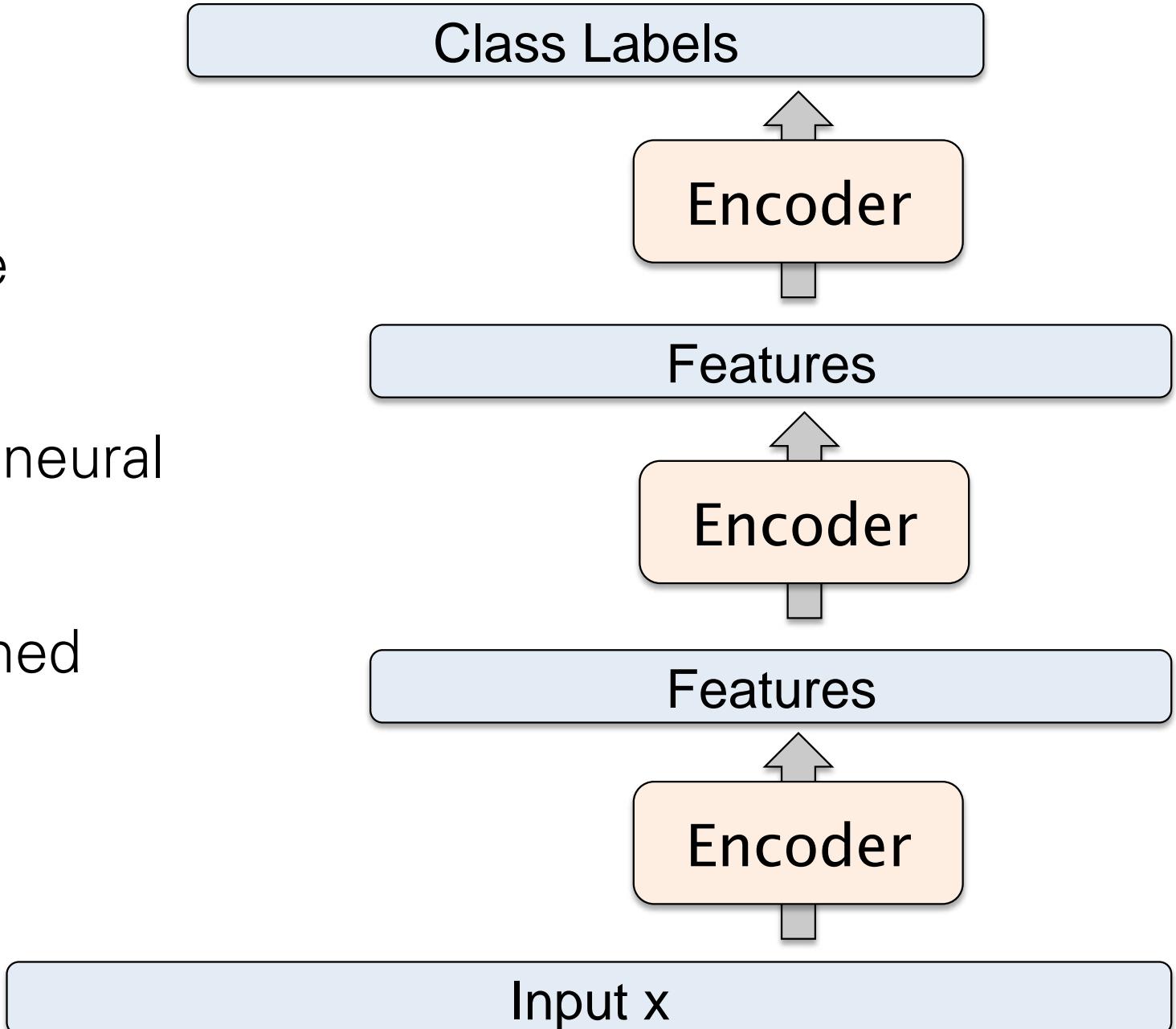
# Stacked Autoencoders

- Remove decoders and use feed-forward part.

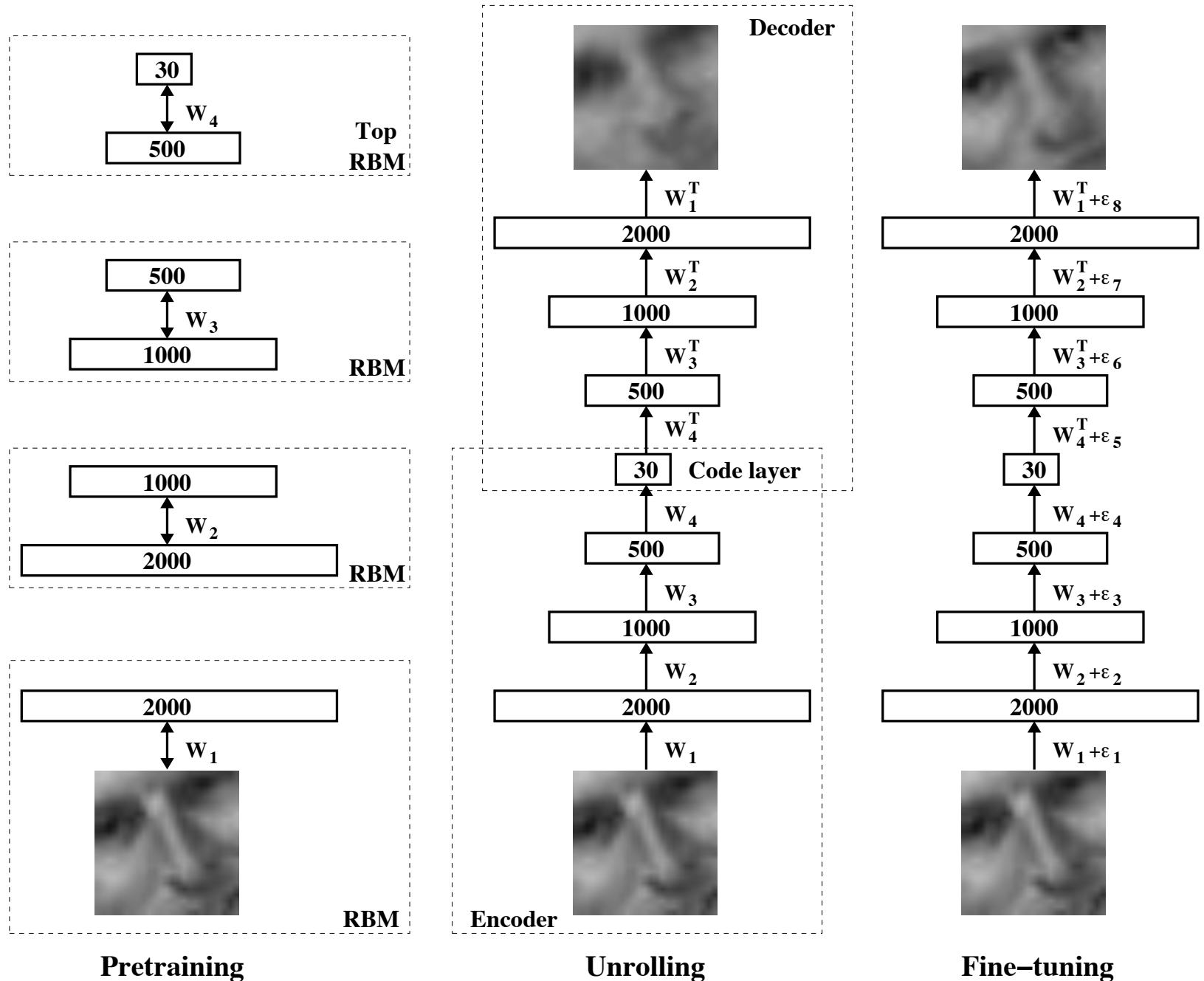


# Stacked Autoencoders

- Remove decoders and use feed-forward part.
- Standard, or convolutional neural network architecture.
- Parameters can be fine-tuned using backpropagation.

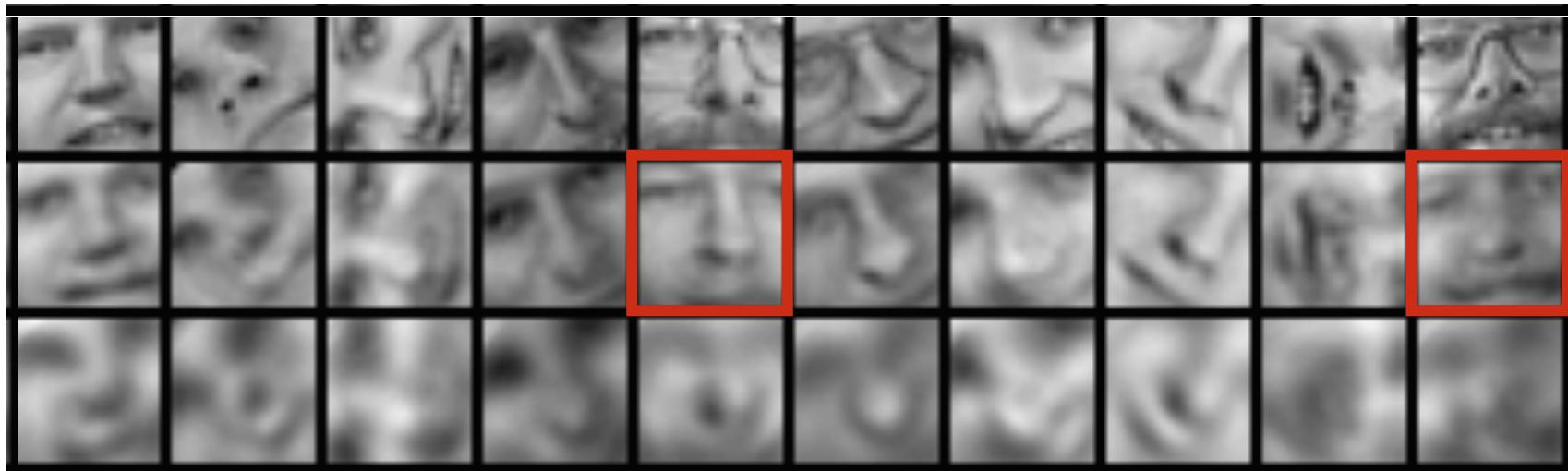


# Deep Autoencoder



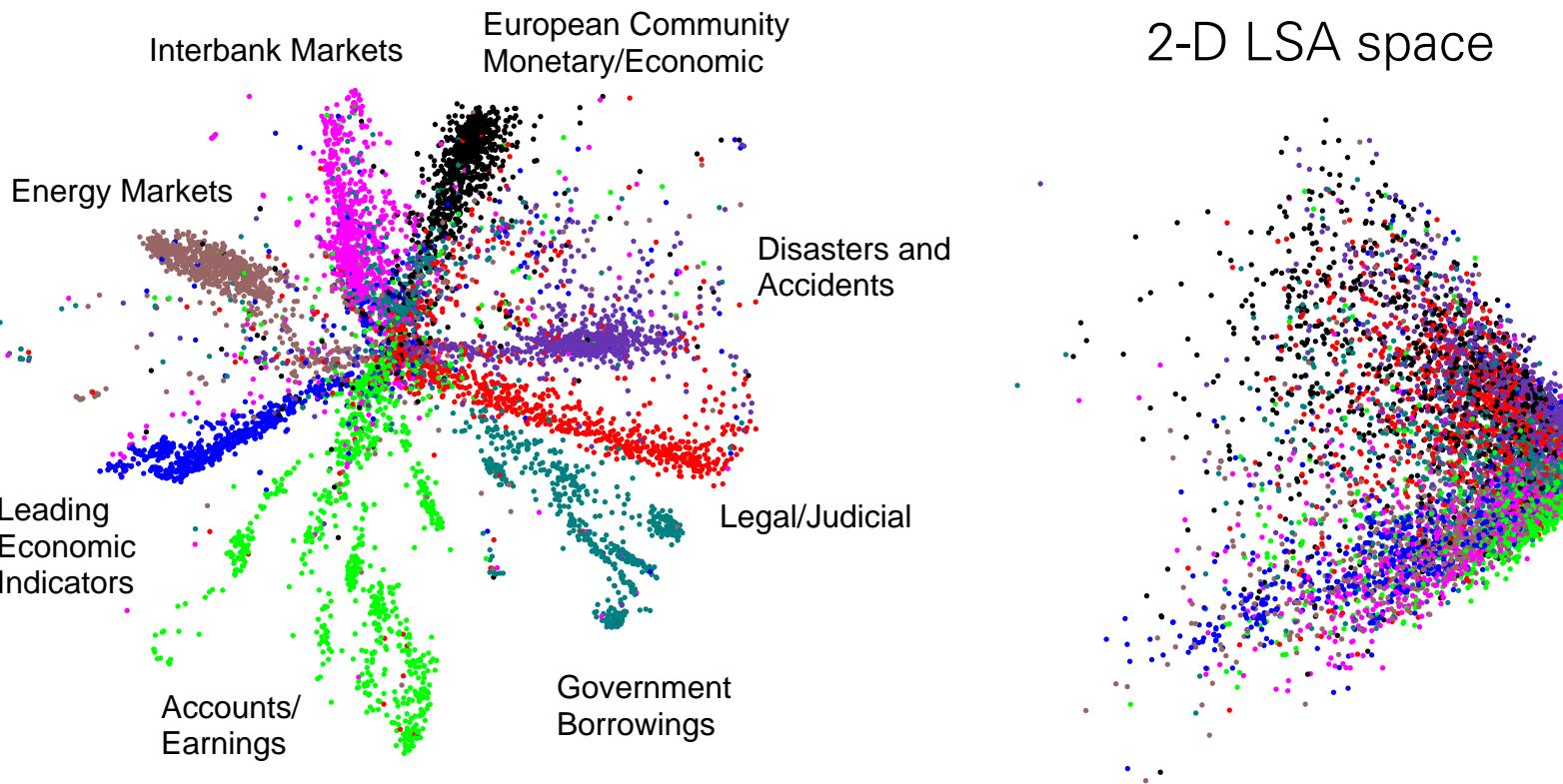
# Deep Autoencoders

- 25x25 – 2000 – 1000 – 500 – 30 autoencoder to extract 30-D real-valued codes for Oliver face patches.



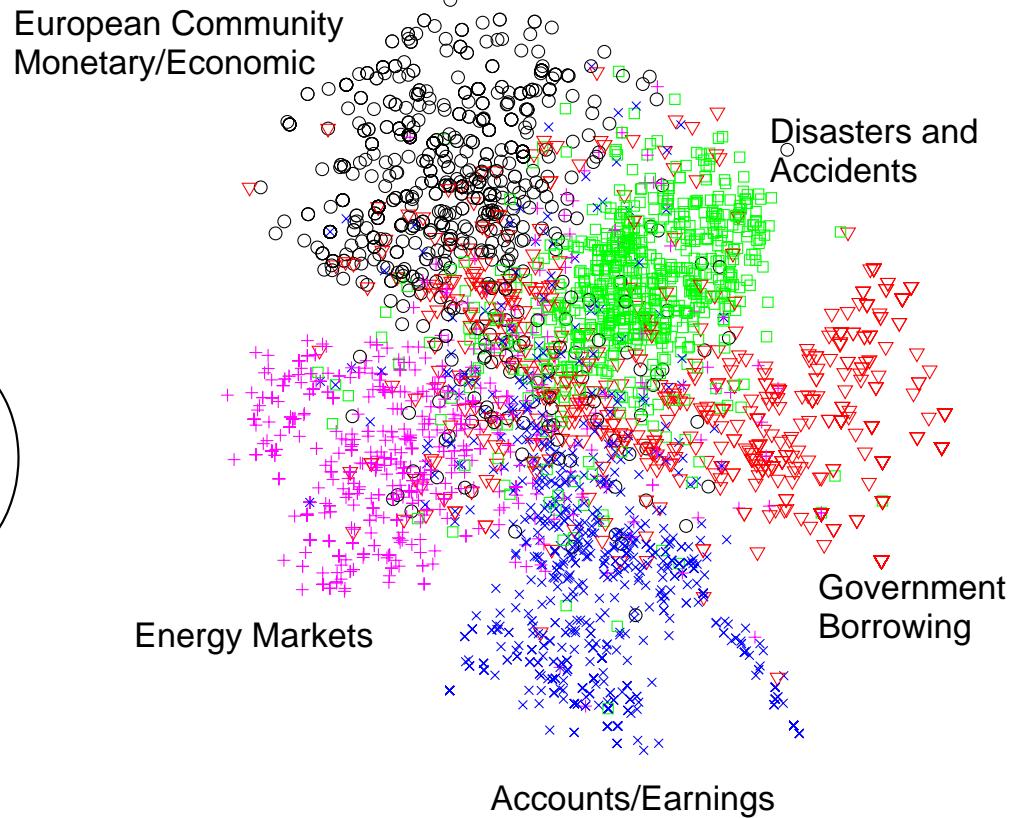
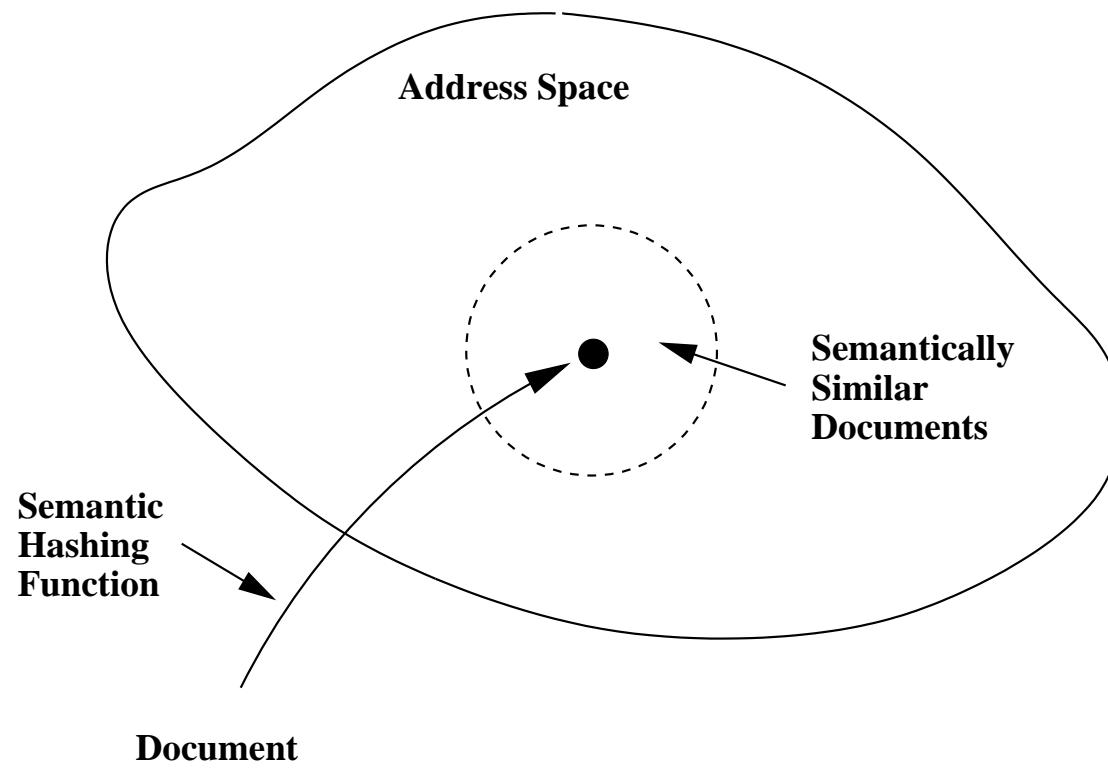
- **Top:** Random samples from the test dataset.
- **Middle:** Reconstructions by the 30-dimensional deep autoencoder.
- **Bottom:** Reconstructions by the 30-dimensional PCA.

# Information Retrieval



- The Reuters Corpus Volume II contains 804,414 newswire stories (randomly split into **402,207 training** and **402,207 test**).
- “Bag-of-words” representation: each article is represented as a vector containing the counts of the most frequently used 2000 words in the training set.

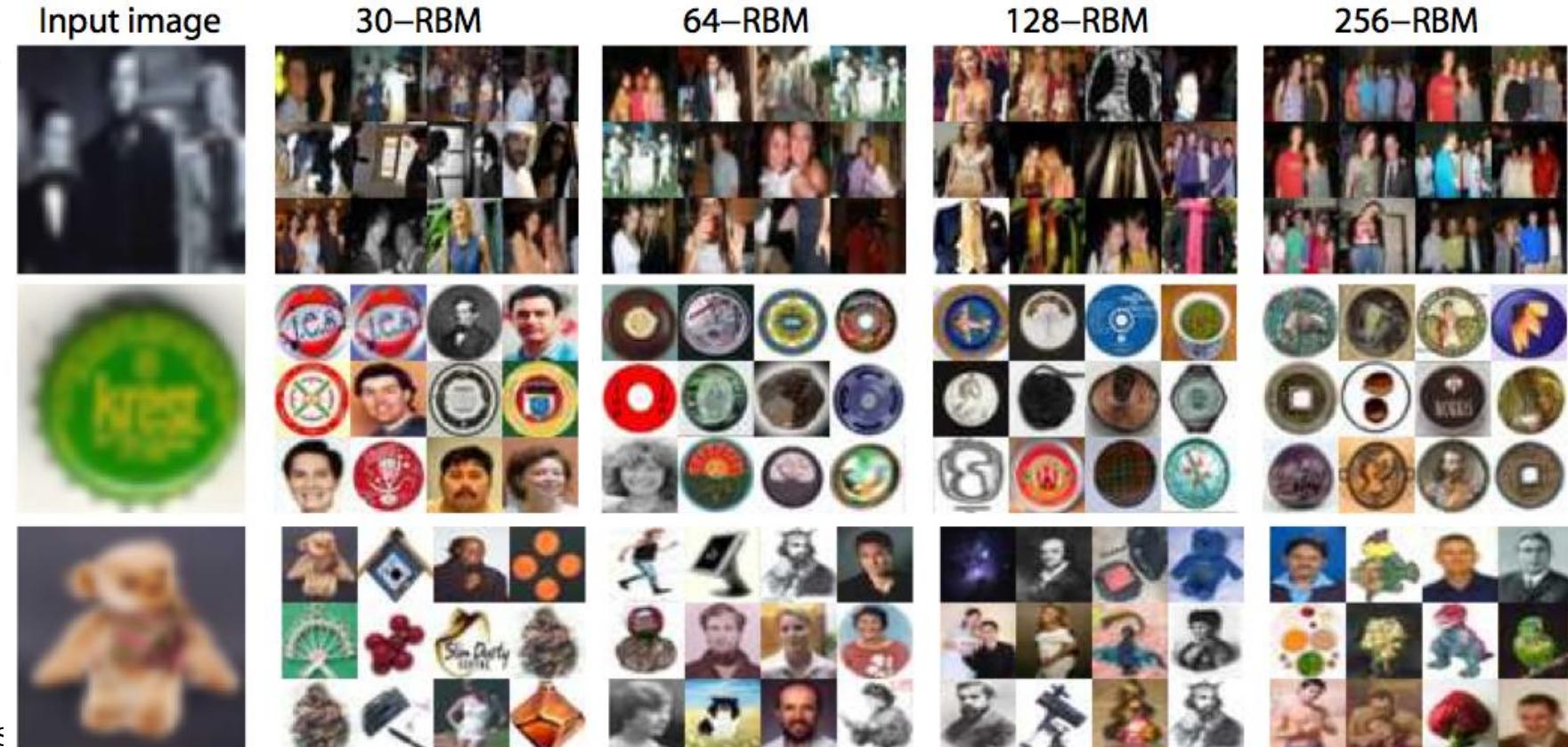
# Semantic Hashing



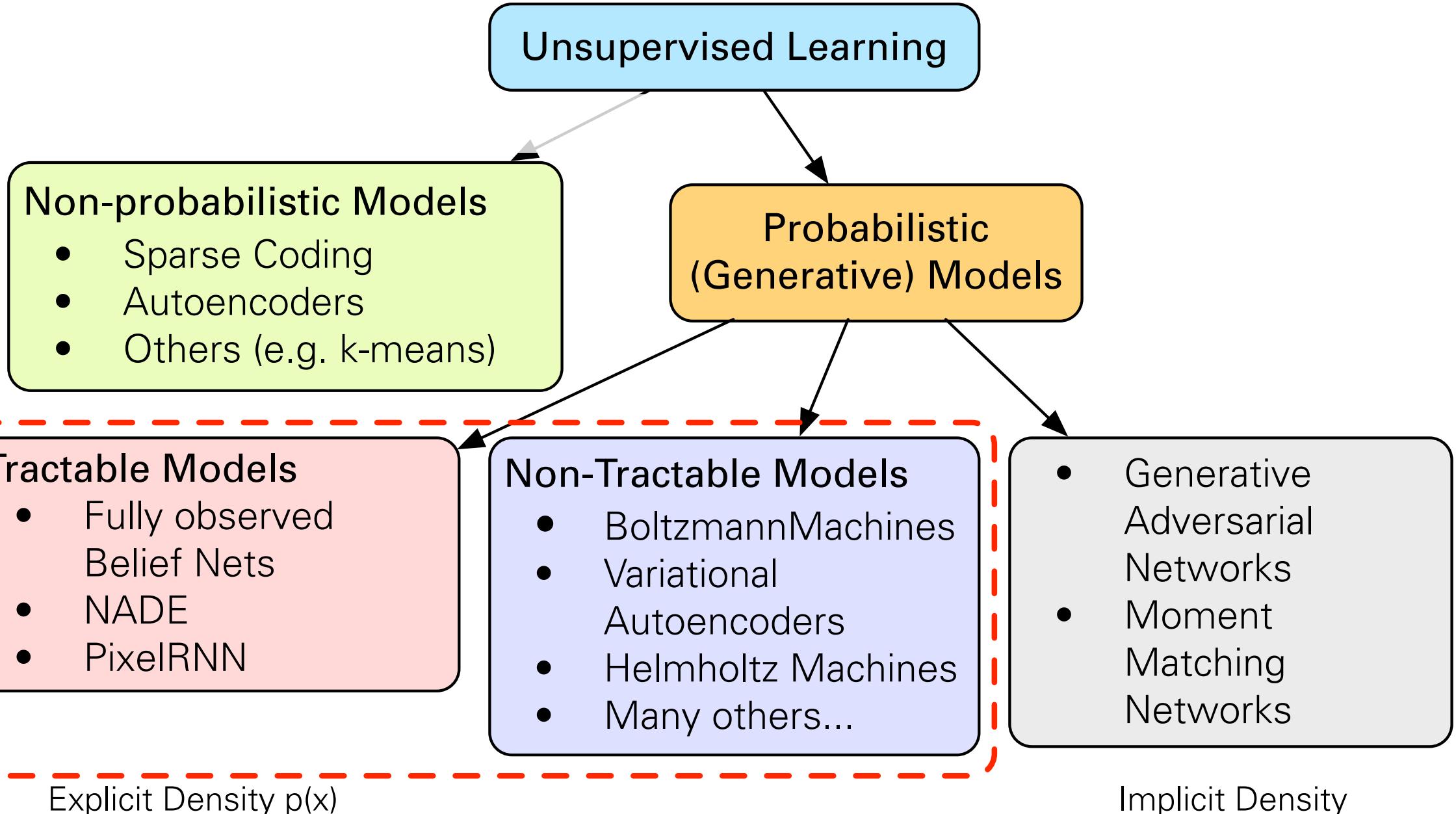
- Learn to map documents into **semantic 20-D binary codes**.
- Retrieve similar documents stored at the nearby addresses **with no search at all**.

# Searching Large Image Database using Binary Codes

- Map images into binary codes for fast retrieval.



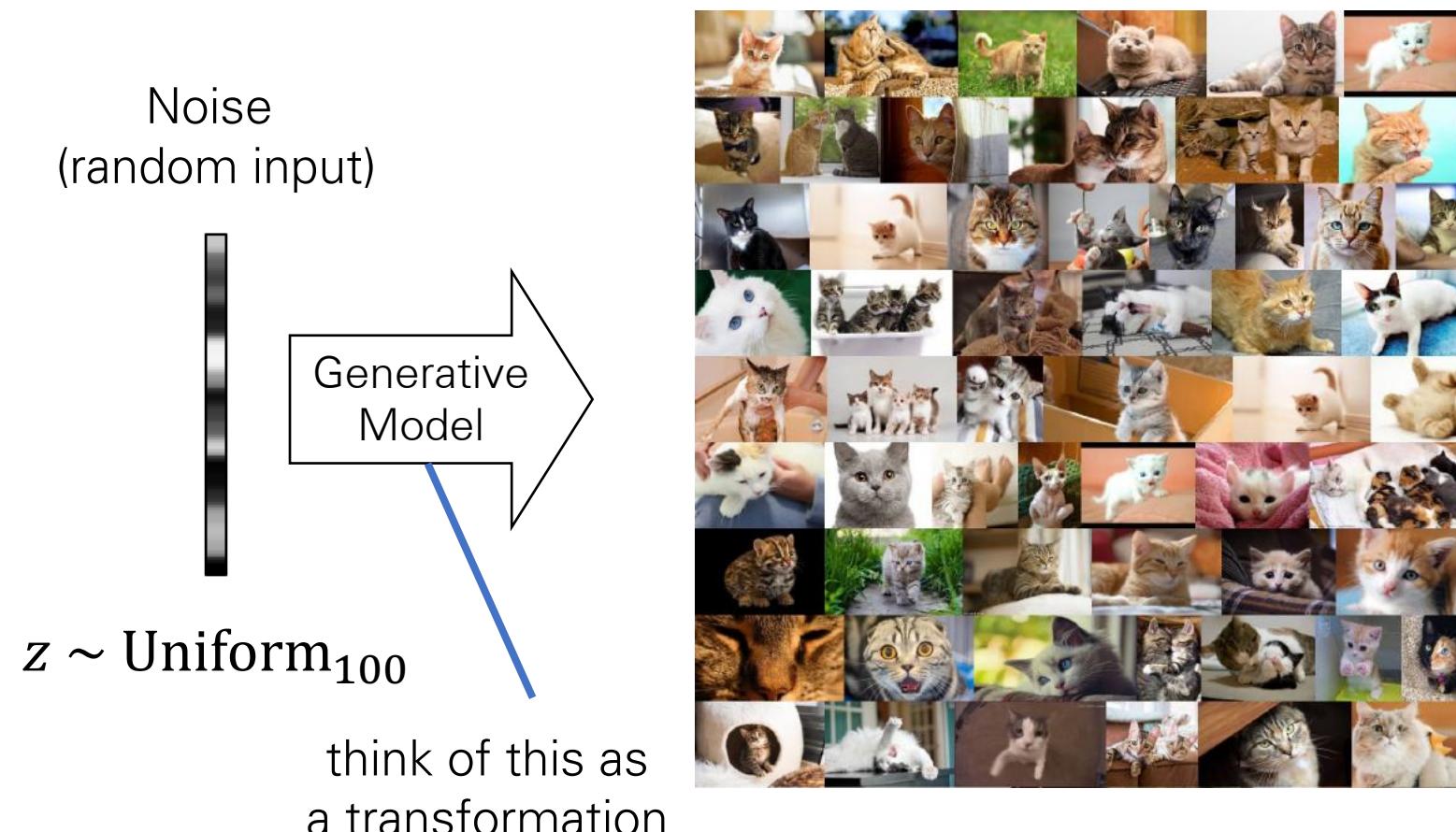
- Small Codes, Torralba, Fergus
- Spectral Hashing, Y. Weiss, A. Torralba, R. Fergus, NIPS 2008
- Kulis and Darrell, NIPS 2009, Gong and Lazebnik, CVPR 2011
- Norouzi and Fleet, ICML 2011



# Generative Adversarial Networks

# Genetive Adversarial Networks (GANs)

(Goodfellow et al., 2014)



- A game-theoretic likelihood free model

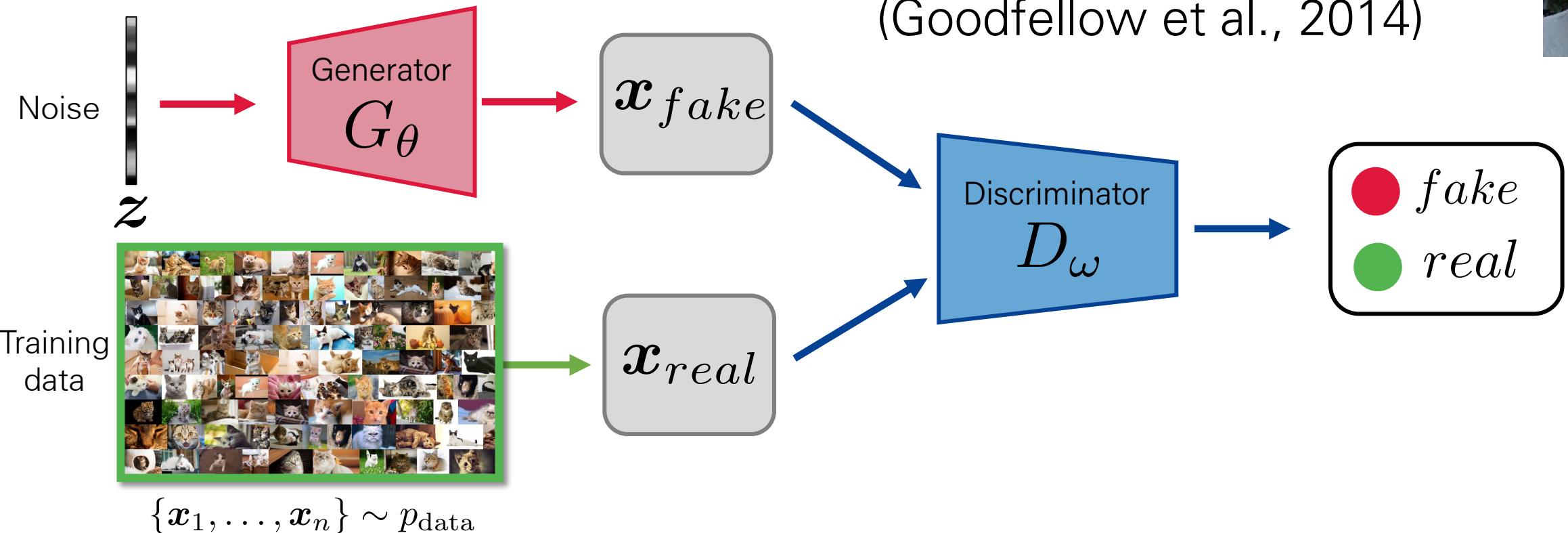
## Advantages:

- Uses a latent code
- No Markov chains needed
- Produces the best looking samples

# Genetive Adversarial Networks (GANs)

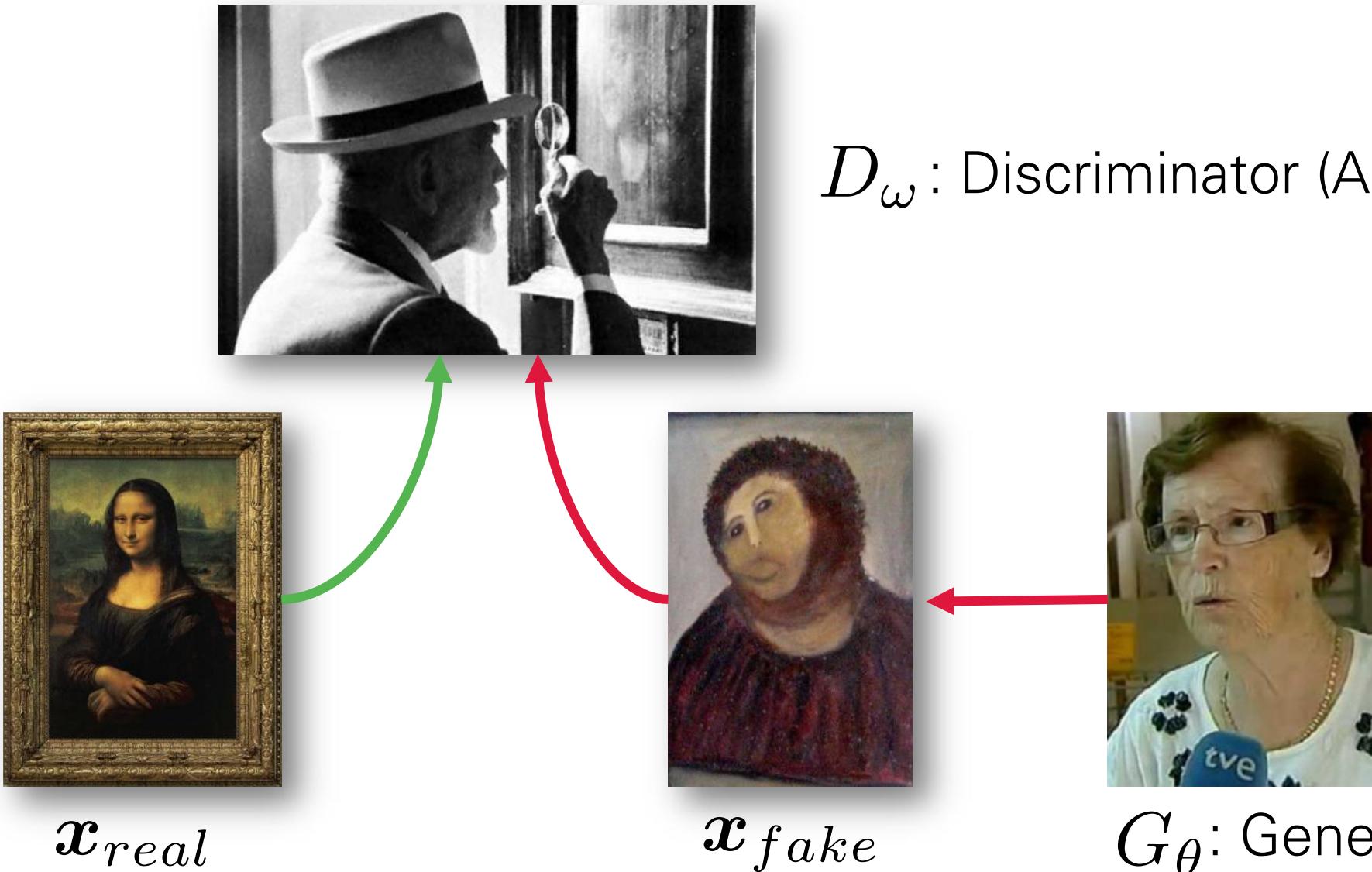


(Goodfellow et al., 2014)



- A game between a generator  $G_\theta(z)$  and a discriminator  $D_\omega(x)$ 
  - Generator tries to fool discriminator (i.e. generate realistic samples)
  - Discriminator tries to distinguish fake from real samples

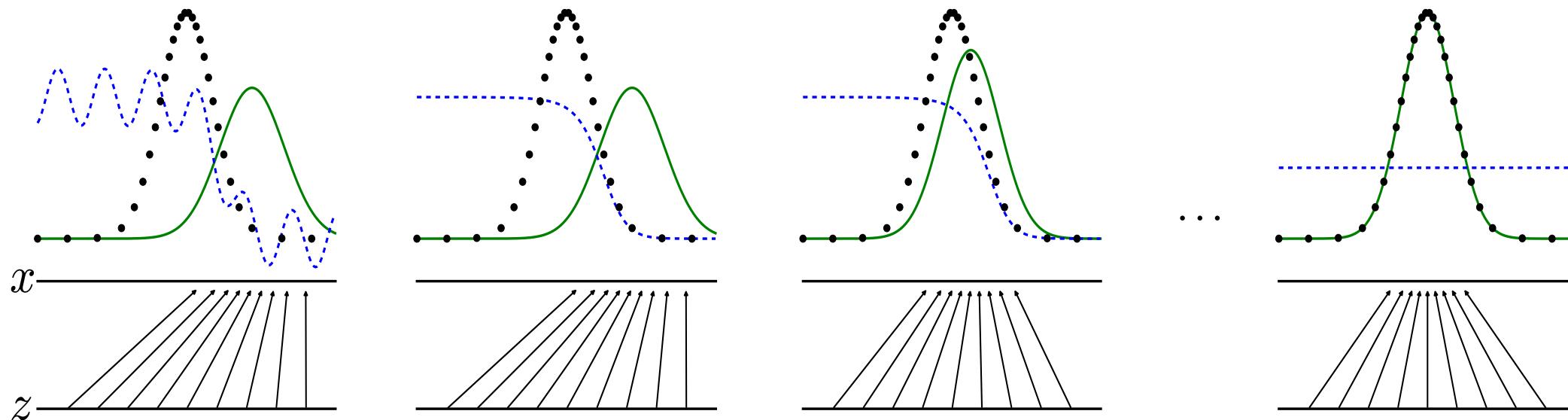
# Intuition behind GANs



# Training Procedure

(Goodfellow et al., 2014)

- Use SGD on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples



# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$



Real data



Noise vector used  
to generate data

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

Cross-entropy  
loss for binary  
classification

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Generator maximizes the log-probability  
of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence between  $p_{\text{data}}$  and  $p_x$

# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_{\omega}(x)] + \mathbb{E}_{z \sim p_z} [\log (1 - D_{\omega}(G_{\theta}(z)))]$$



Real data



Noise vector used

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{x \sim}$$

Important question is  
“Does this converge??”

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_z \mathbb{E}_{x \sim}$$

of the discriminator being mistaken

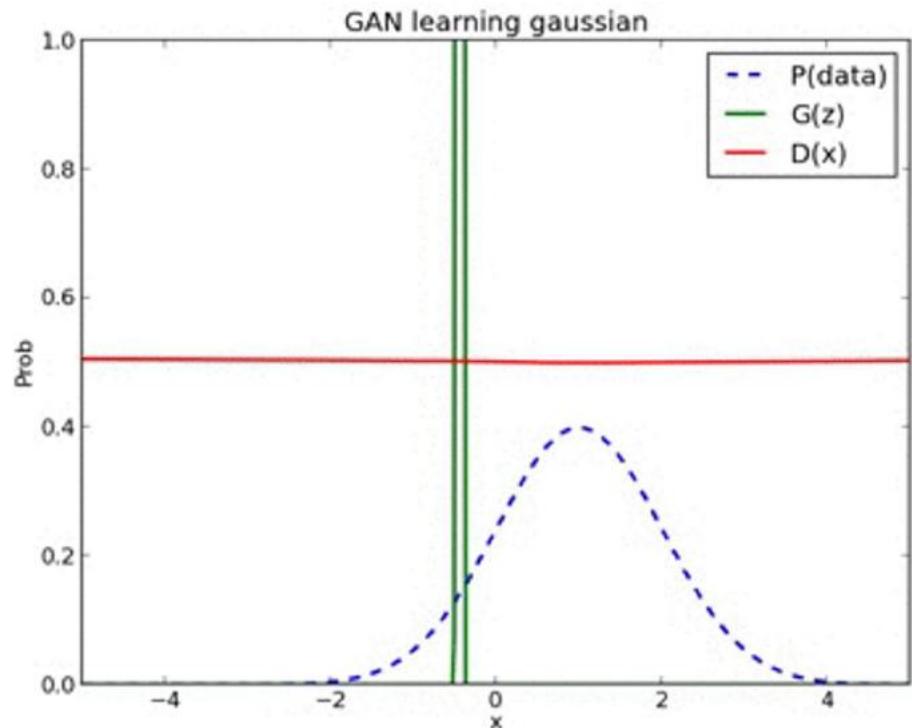
cross-entropy  
loss for binary  
classification

g-probability

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence

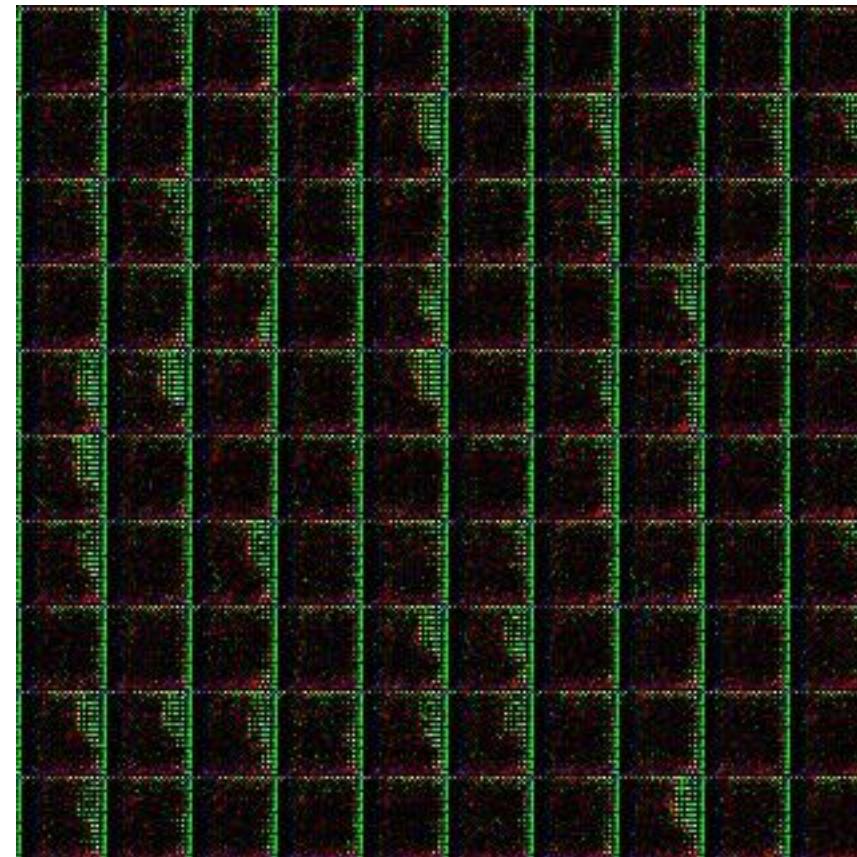
# Training Procedure

(Goodfellow et al., 2014)



Source: Alec Radford

Generating 1D points



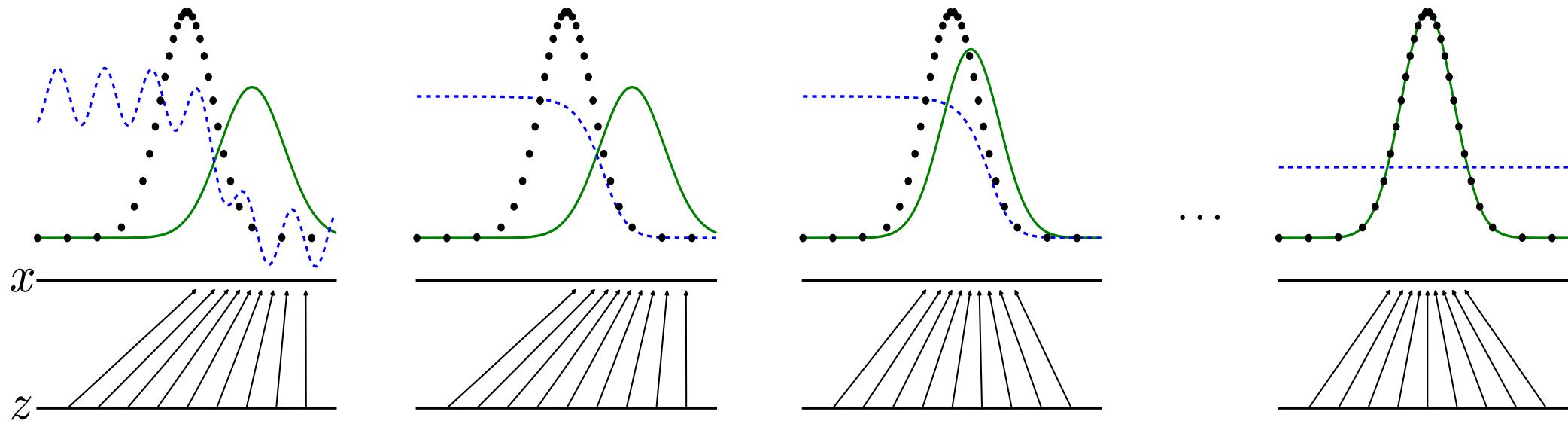
Source: OpenAI blog

Generating images

# Training Procedure

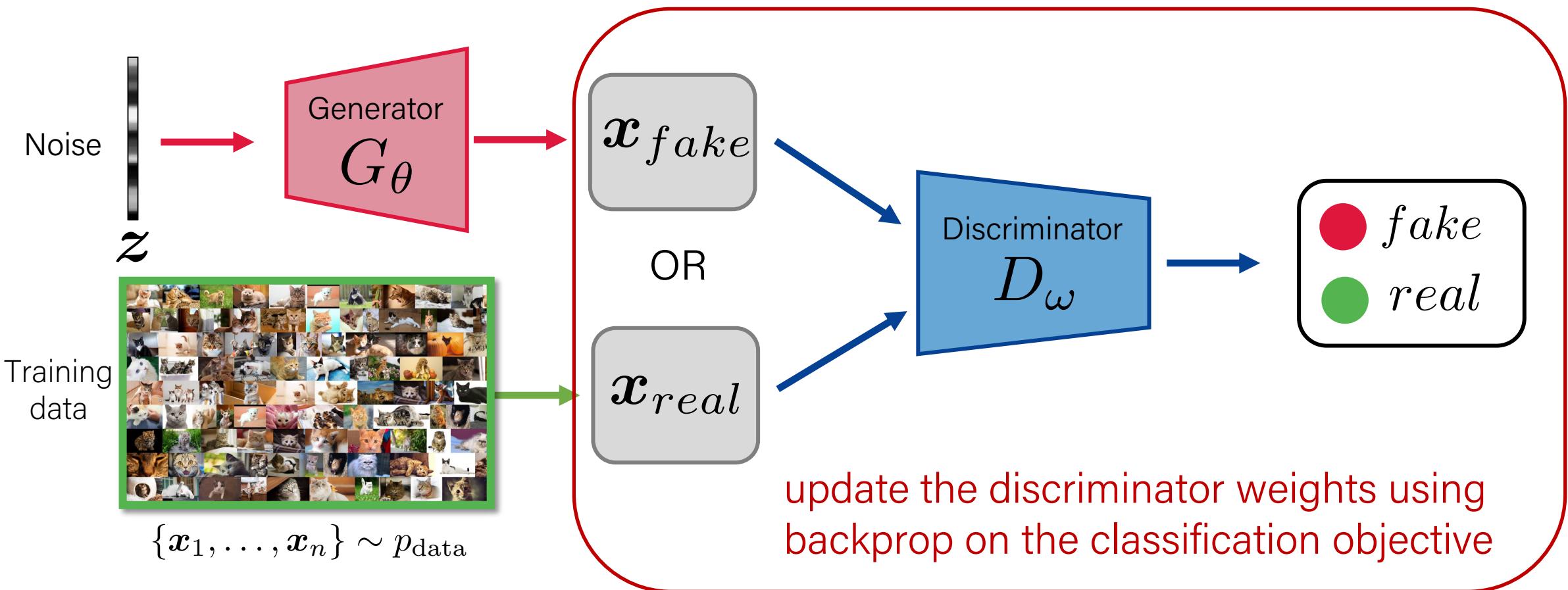
(Goodfellow et al., 2014)

- Use SGD on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples



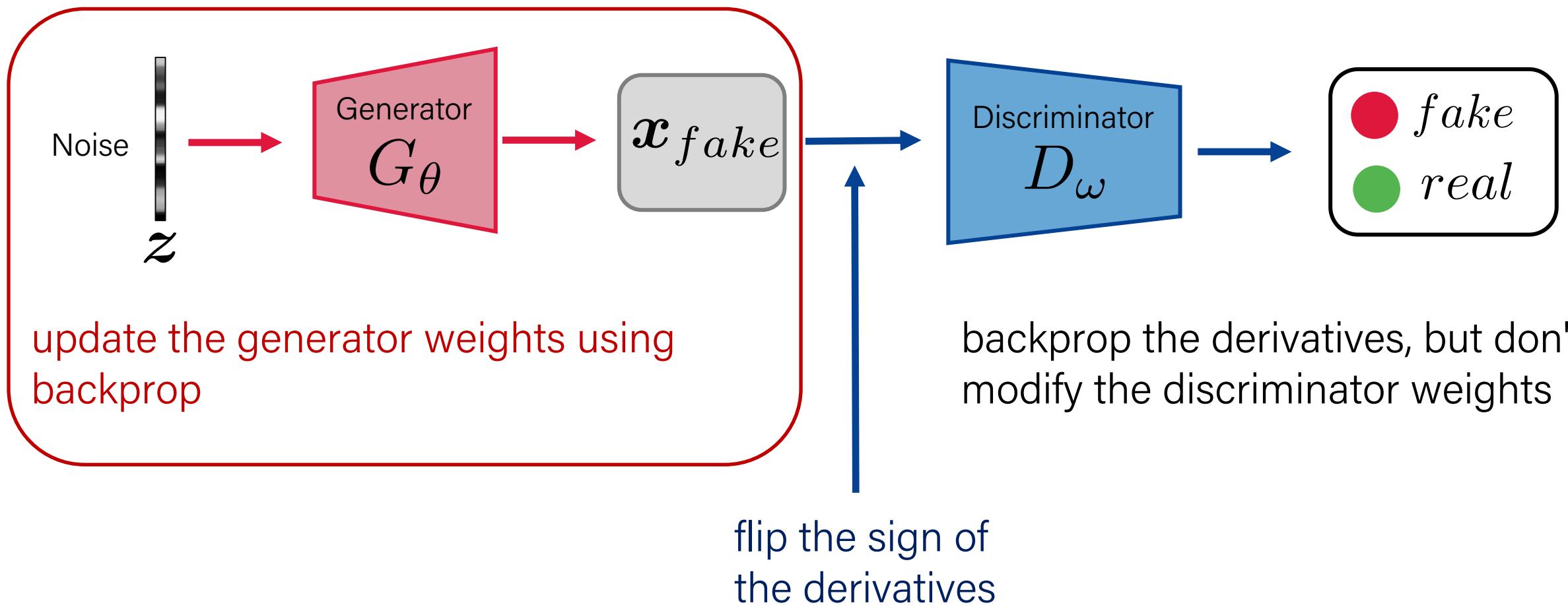
# Training Procedure

- Updating the discriminator:



# Training Procedure

- Updating the generator:



# Results

(Goodfellow et al., 2014)

- The generator uses a mixture of rectifier linear activations and/or sigmoid activations
- The discriminator net used maxout activations.



MNIST samples



TFD samples



CIFAR10 samples  
(fully-connected model)



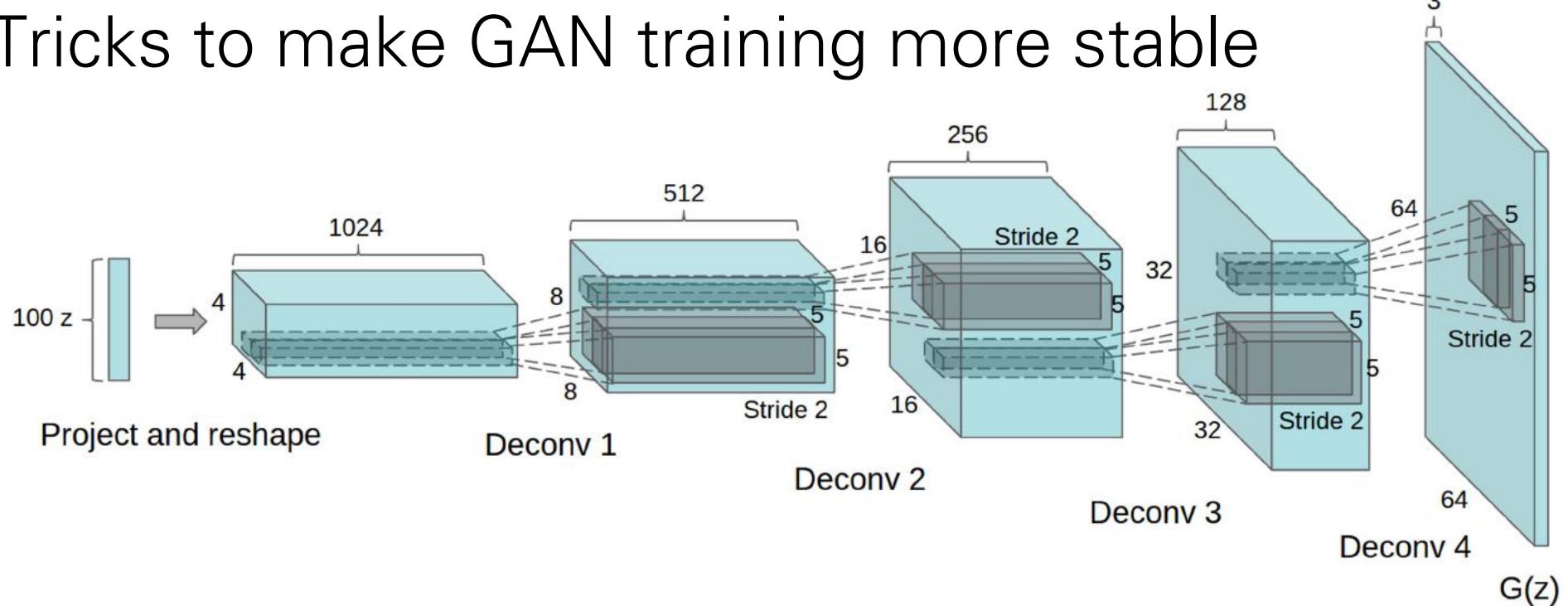
CIFAR10 samples  
(convolutional discriminator,  
deconvolutional generator)

# Deep Convolutional GANs (DCGAN)



(Radford et al., 2015)

- Idea: Tricks to make GAN training more stable



- No fully connected layers
- Batch Normalization  
(Ioffe and Szegedy, 2015)
- Leaky Rectifier in D
- Use Adam (Kingma and Ba, 2015)
- Tweak Adam hyperparameters a bit  
(lr=0.0002, b1=0.5)

# DCGAN for LSUN Bedrooms

64×64 pixels  
~3M images

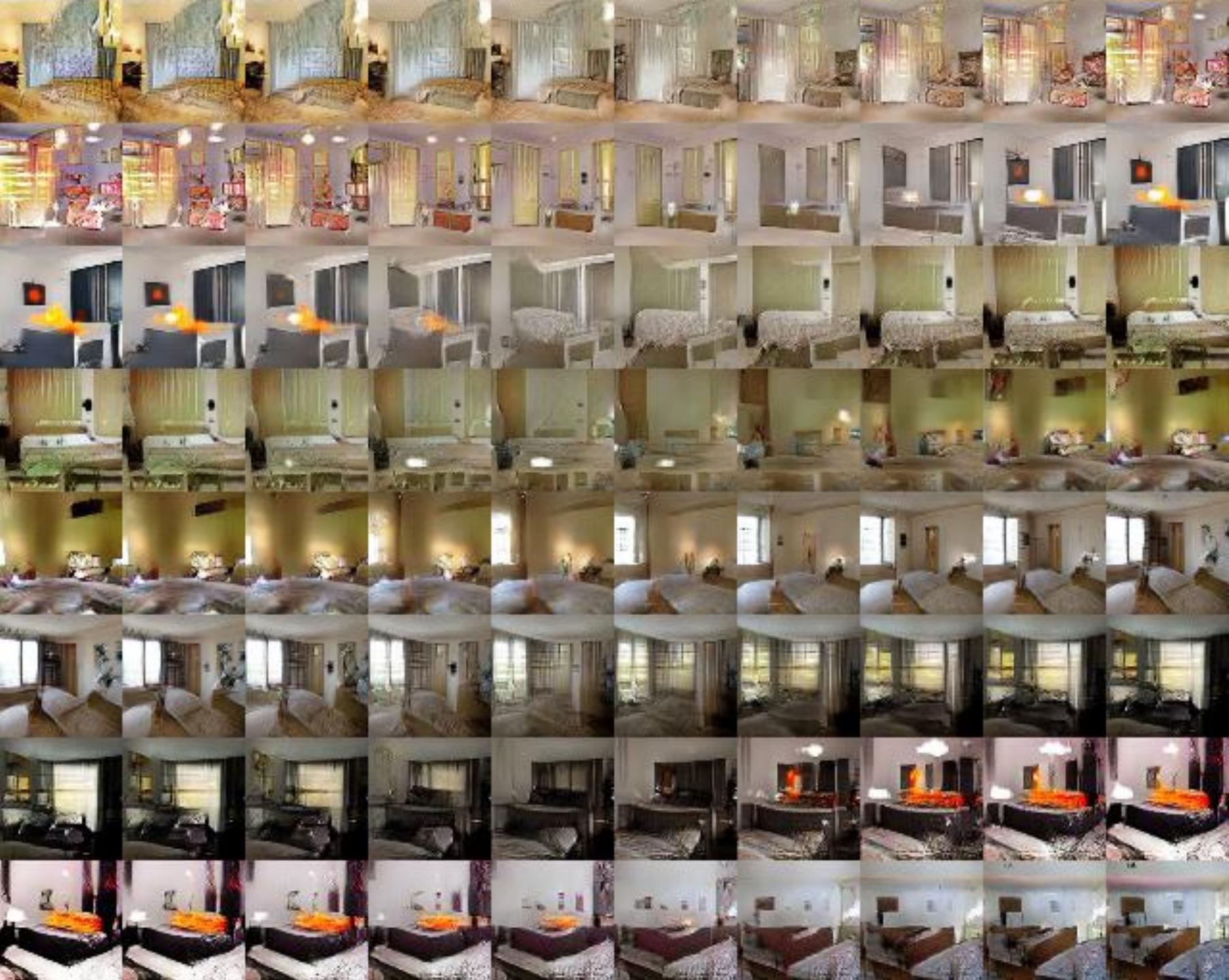
(Radford et al.,  
2015)



# Walking over the latent space

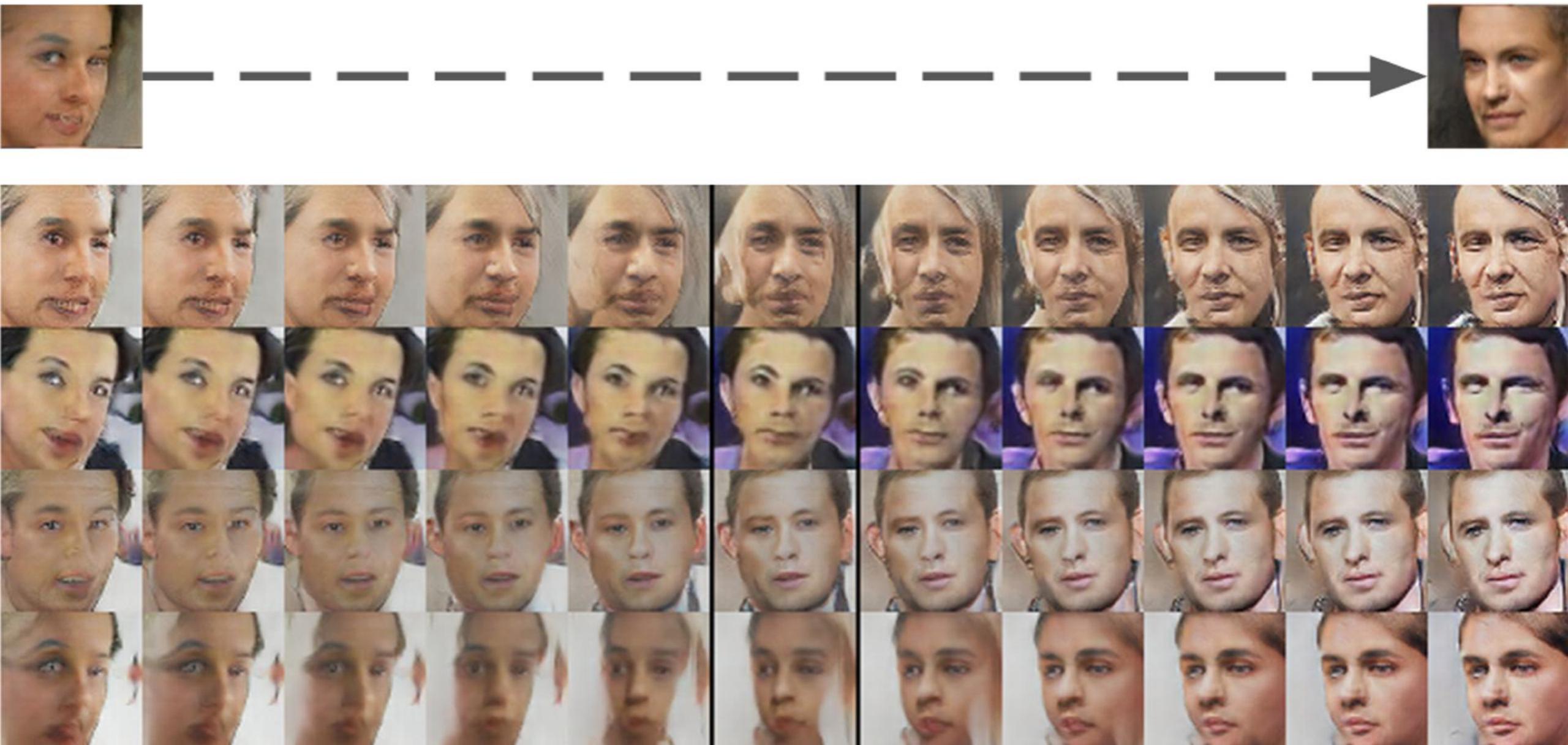
(Radford et al., 2015)

- Interpolation suggests non-overfitting behavior



# Walking over the latent space

(Radford et al., 2015)



# Vector Space Arithmetic

(Radford et al., 2015)



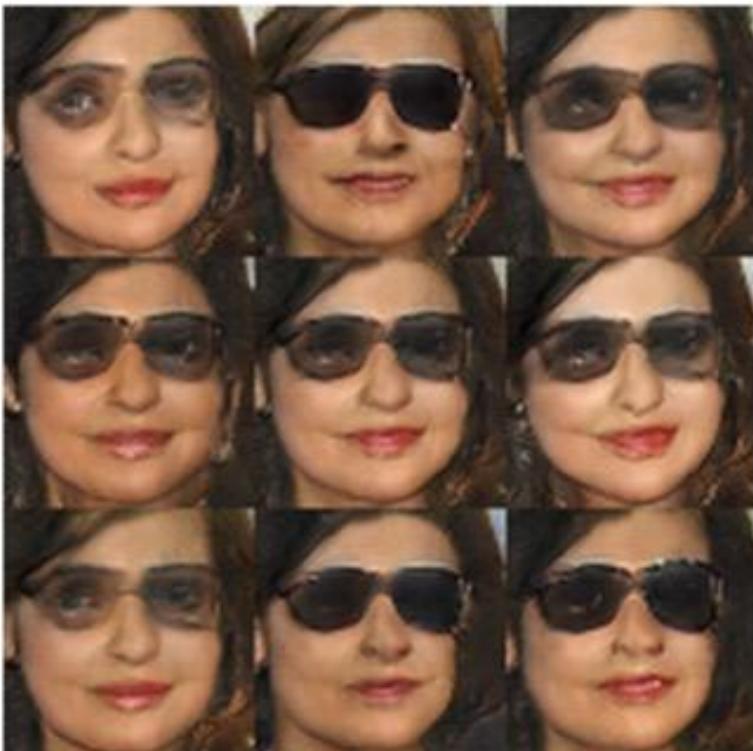
man  
with glasses



man  
without glasses



woman  
without glasses



woman with glasses

# Vector Space Arithmetic

(Radford et al., 2015)



smiling  
woman



neutral  
woman

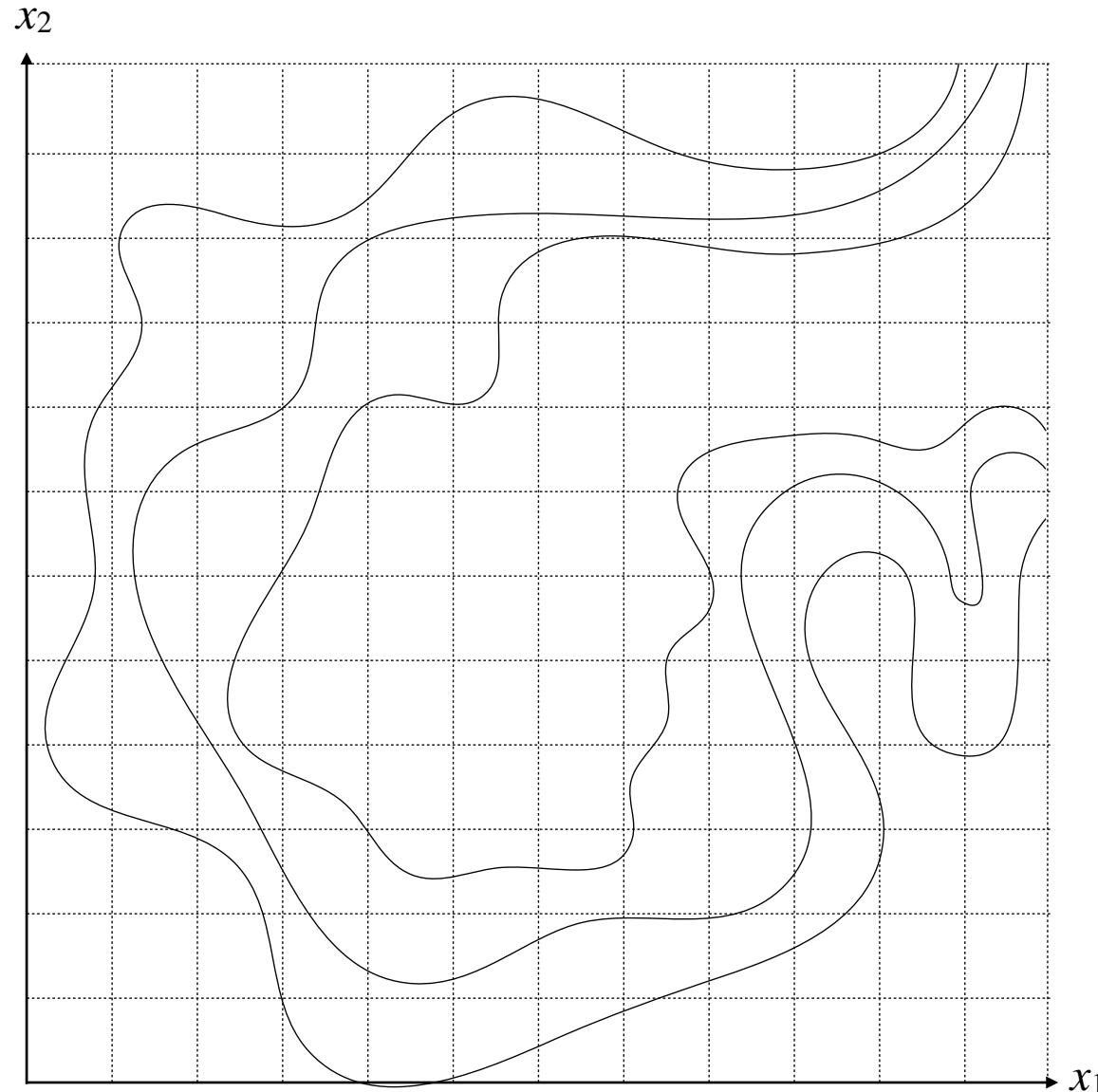


neutral  
man

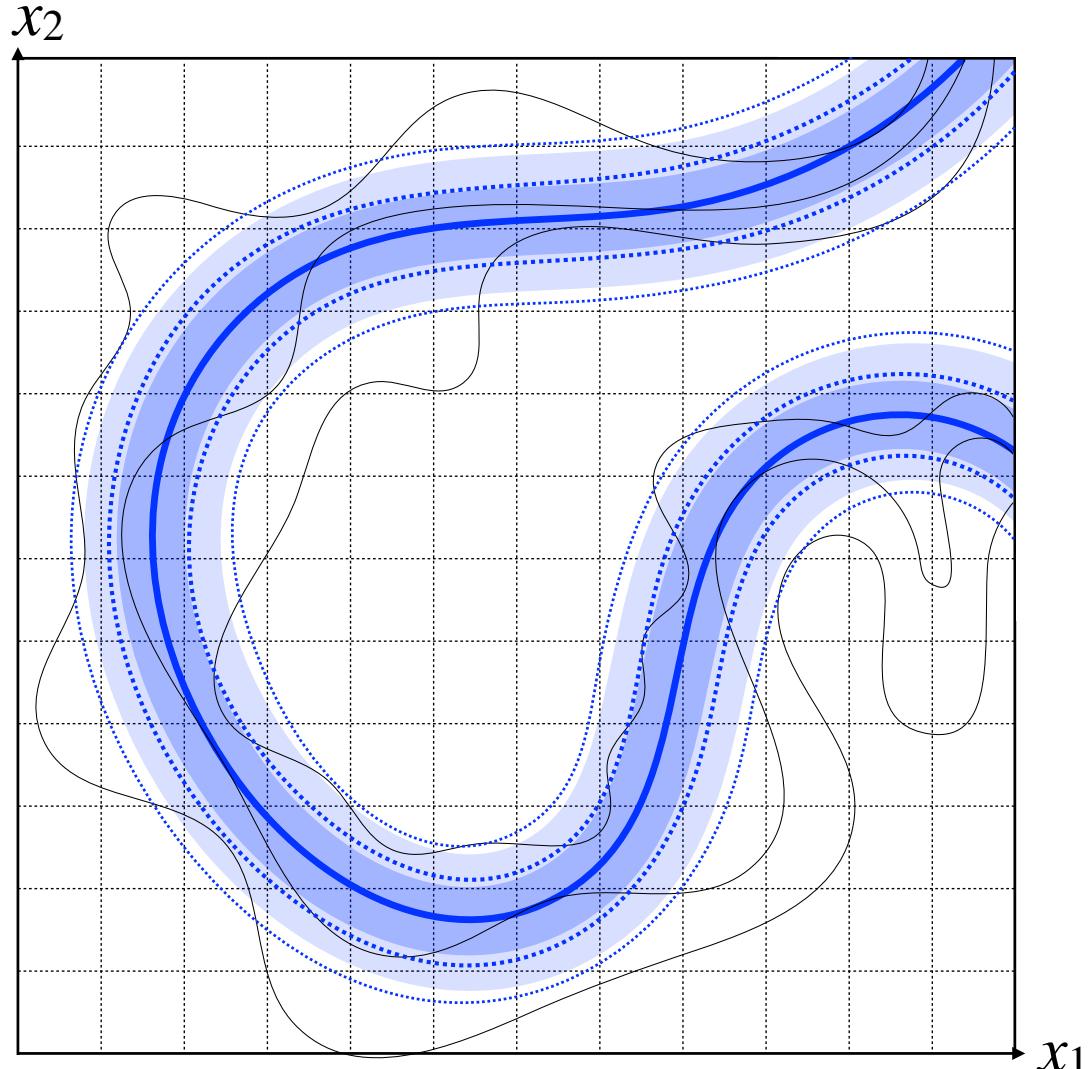


smiling man

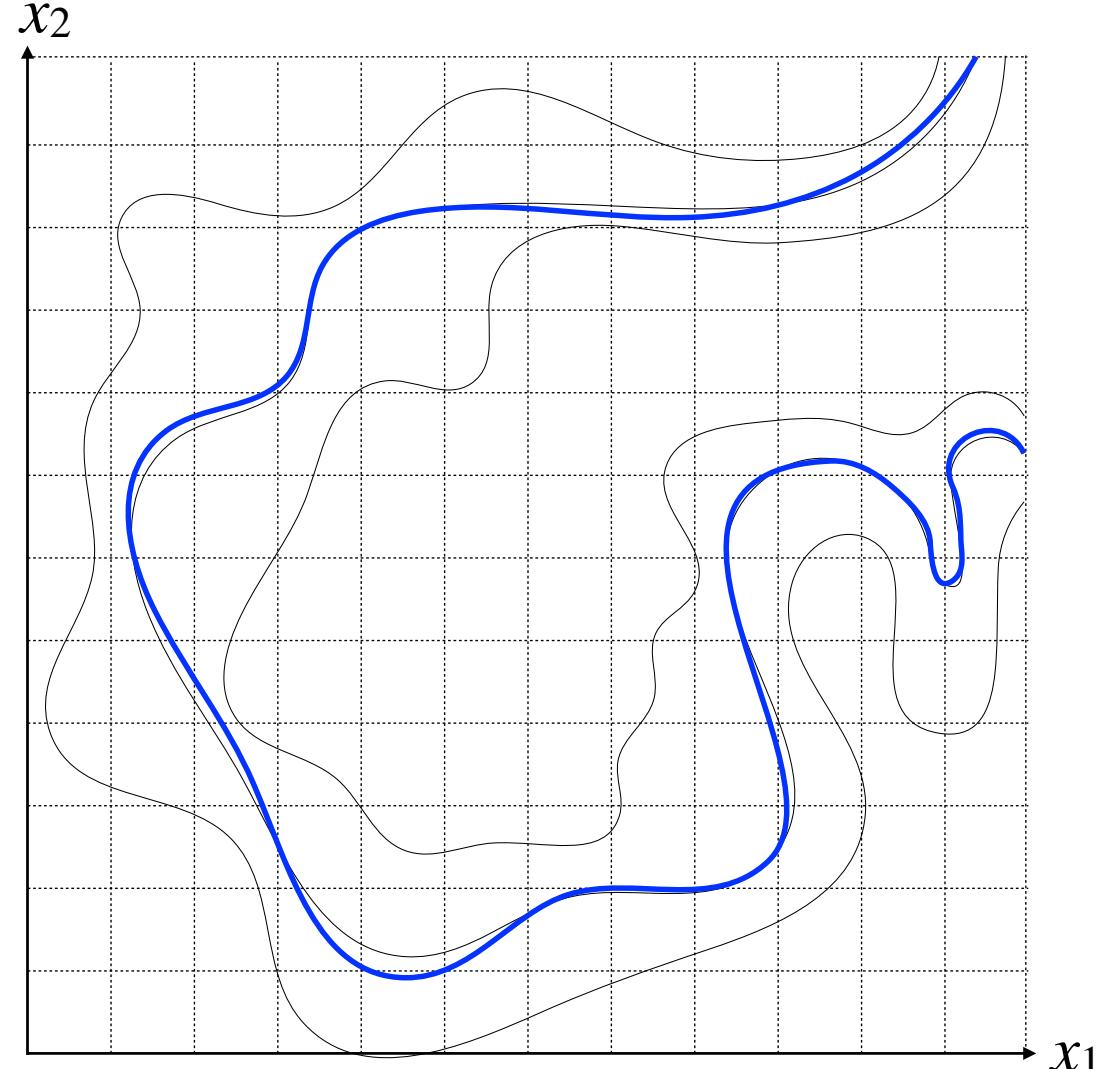
# Cartoon of the Image manifold



# What makes GANs special?



more traditional max-likelihood approach

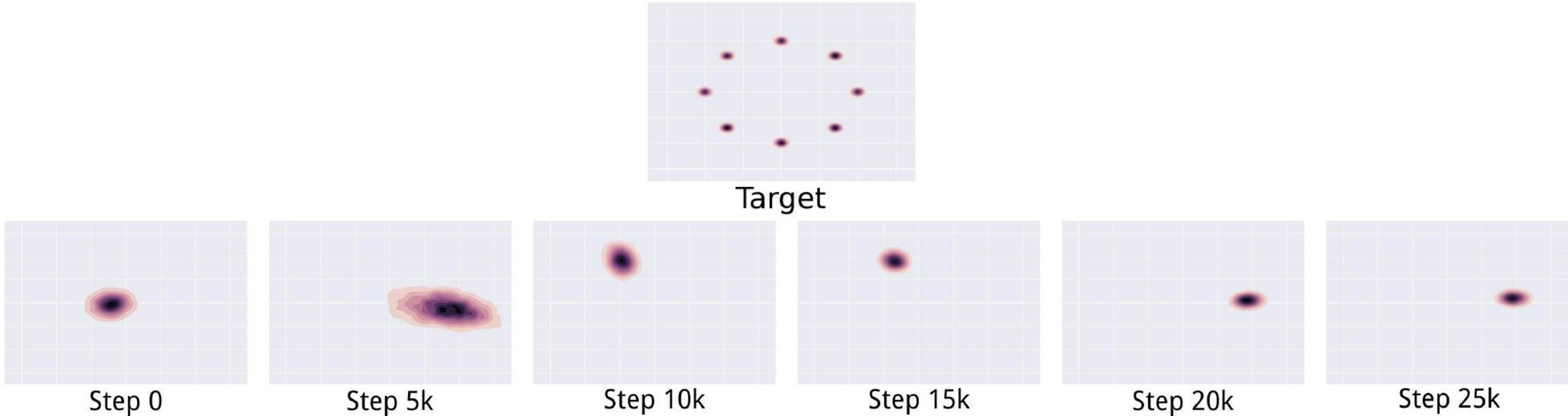


GAN

# GAN Failures: Mode Collapse

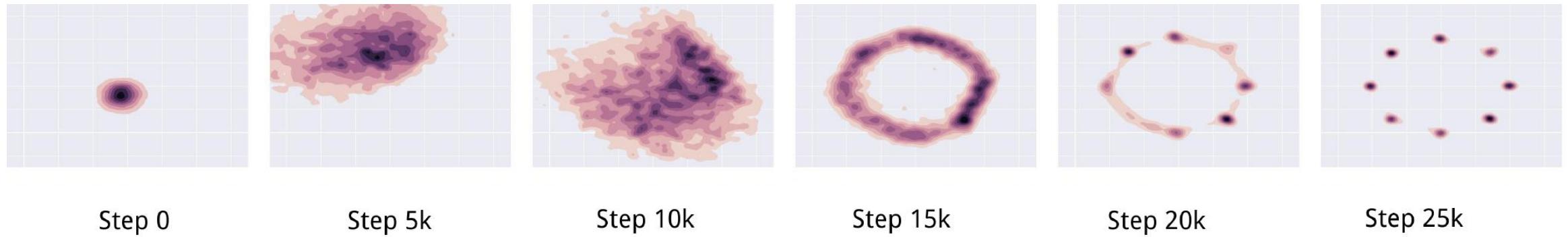
$$\min_G \max_D V(G, D) \neq \max_D \min_G V(G, D)$$

- $D$  in inner loop: convergence to correct distribution
- $G$  in inner loop: place all mass on most likely point



# Mode Collapse: Solutions

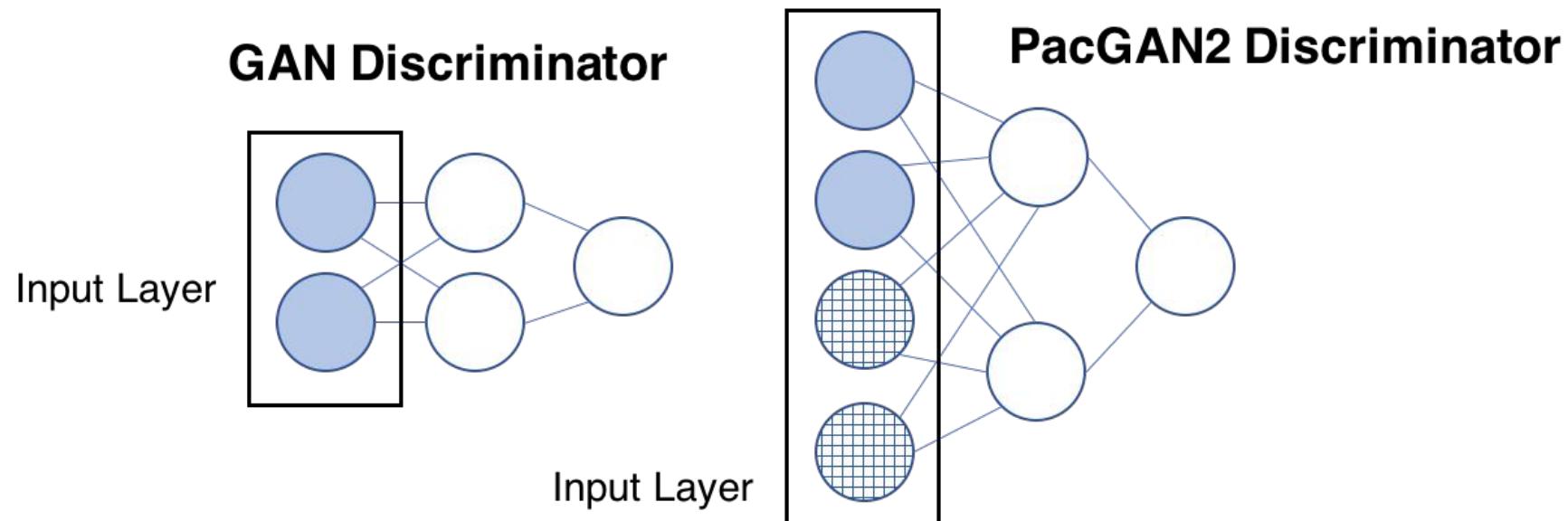
- **Unrolled GANs** (Metz et al 2016): Prevents mode collapse by backproping through a set of ( $k$ ) updates of the discriminator to update generator parameters



- **VEEGAN** (Srivastava et al 2017): Introduce a reconstructor network which is learned both to map the true data distribution  $p(x)$  to a Gaussian and to approximately invert the generator network.

# Mode Collapse: Solutions

- **Minibatch Discrimination** (Salimans et al 2016): Add minibatch features that classify each example by comparing it to other members of the minibatch (Salimans et al 2016)
- **PacGAN**: The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.



# Mode Collapse: Solutions

- PacGAN: The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.

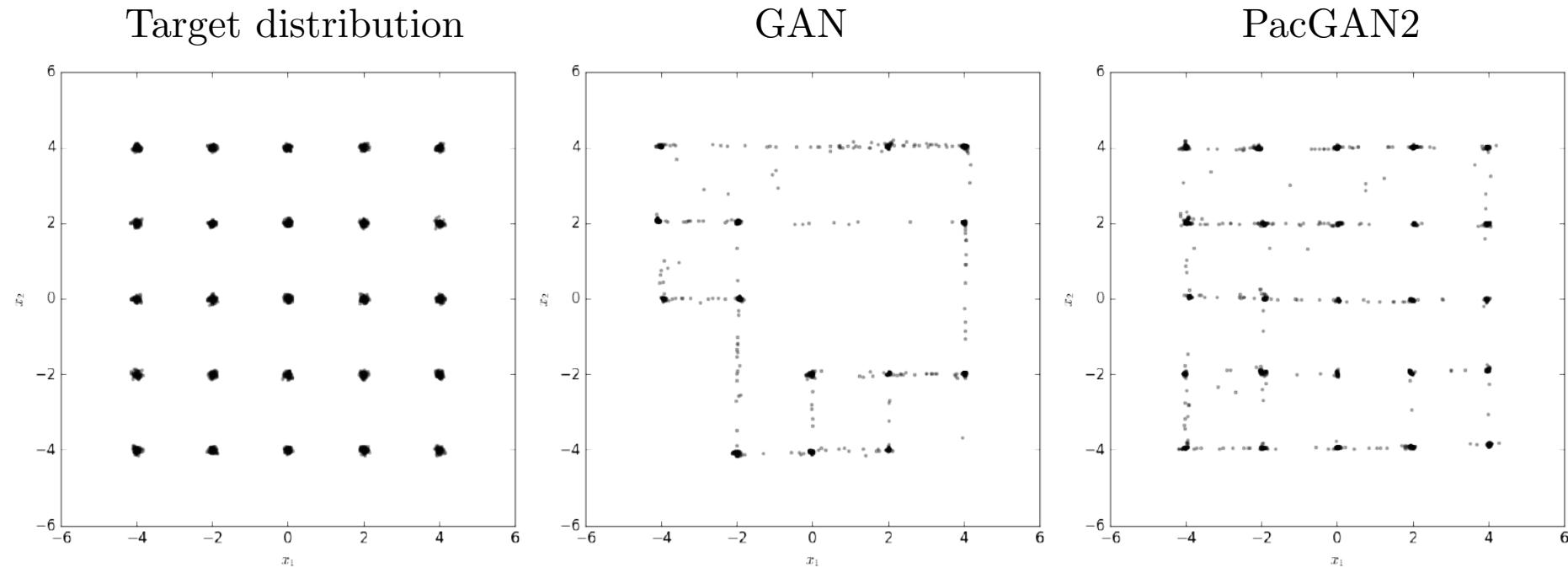


Figure 2: Scatter plot of the 2D samples from the true distribution (left) of 2D-grid and the learned generators using GAN (middle) and PacGAN2 (right). PacGAN2 captures all of the 25 modes.

# GAN Evaluation

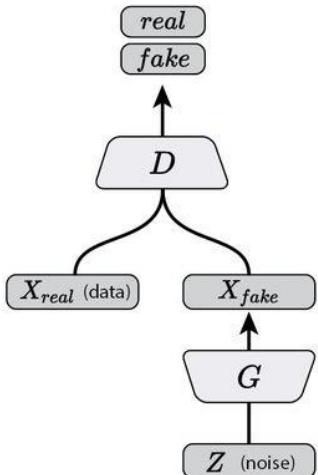
- Quantitatively evaluating GANs is not straightforward:
  - Max Likelihood is a poor indication of sample quality
- Some evaluation metrics
  - **Inception Score (IS):**  
 $y$  = labels given gen. image.  $p(y|x)$  is from classifier - InceptionNet
$$\text{IS}(\mathbb{P}_g) = e^{\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_g} [KL(p_M(y|\mathbf{x}) || p_M(y))]}$$
  - **Fréchet inception distance (FID):** (Currently most popular)  
Estimate mean  $m$  and covariance  $C$  from classifier output - InceptionNet
$$d^2((m, C), (m_w, C_w)) = \|m - m_w\|_2^2 + \text{Tr}(C + C_w - 2(CC_w)^{1/2})$$
  - **Kernel MMD** (Maximum Mean Discrepancy):

$$\text{MMD}(\mathbb{P}_r, \mathbb{P}_g) = \left( \mathbb{E}_{\substack{\mathbf{x}_r, \mathbf{x}'_r \sim \mathbb{P}_r, \\ \mathbf{x}_g, \mathbf{x}'_g \sim \mathbb{P}_g}} \left[ k(\mathbf{x}_r, \mathbf{x}'_r) - 2k(\mathbf{x}_r, \mathbf{x}_g) + k(\mathbf{x}_g, \mathbf{x}'_g) \right] \right)^{\frac{1}{2}}$$

# Subclasses of GANs

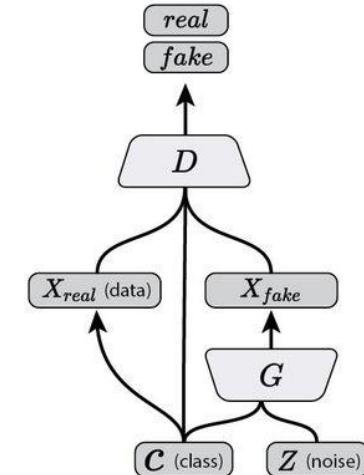
**Vanilla GAN**

**Vanilla GAN**  
(Goodfellow, et al., 2014)

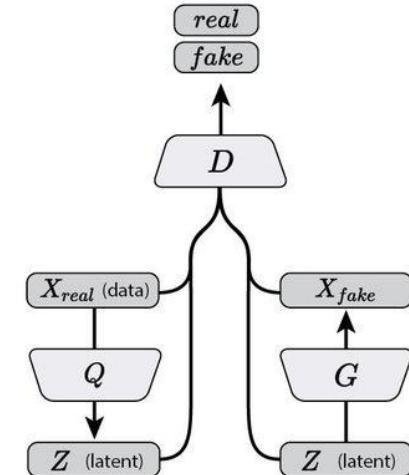


**Discriminator Looks at Latent Variables**

**Conditional GAN**  
(Mirza & Osindero, 2014)

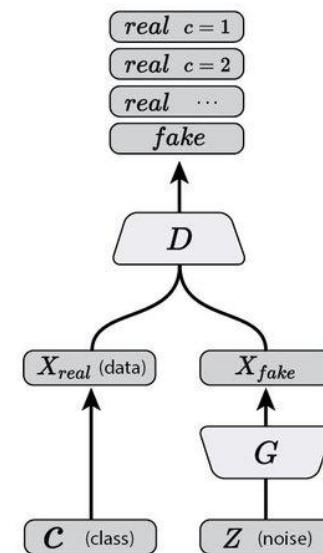


**Bidirectional GAN**  
(Donahue, et al., 2016; Dumoulin, et al., 2016)

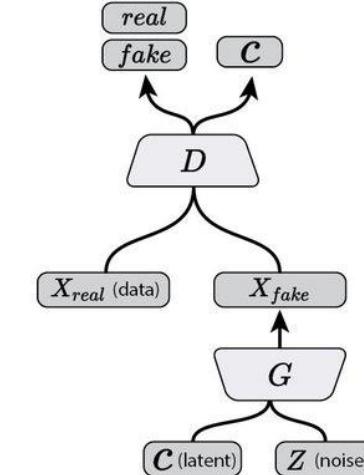


**Discriminator Predicts Latent Variables**

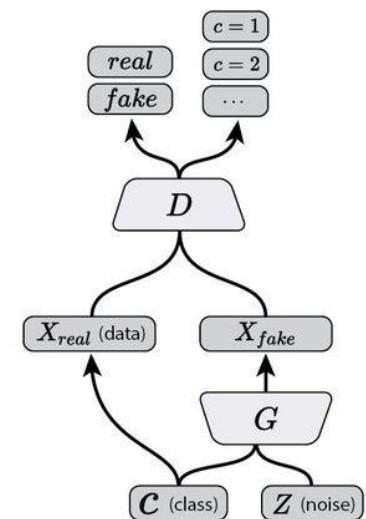
**Semi-Supervised GAN**  
(Odena, 2016; Salimans, et al., 2016)



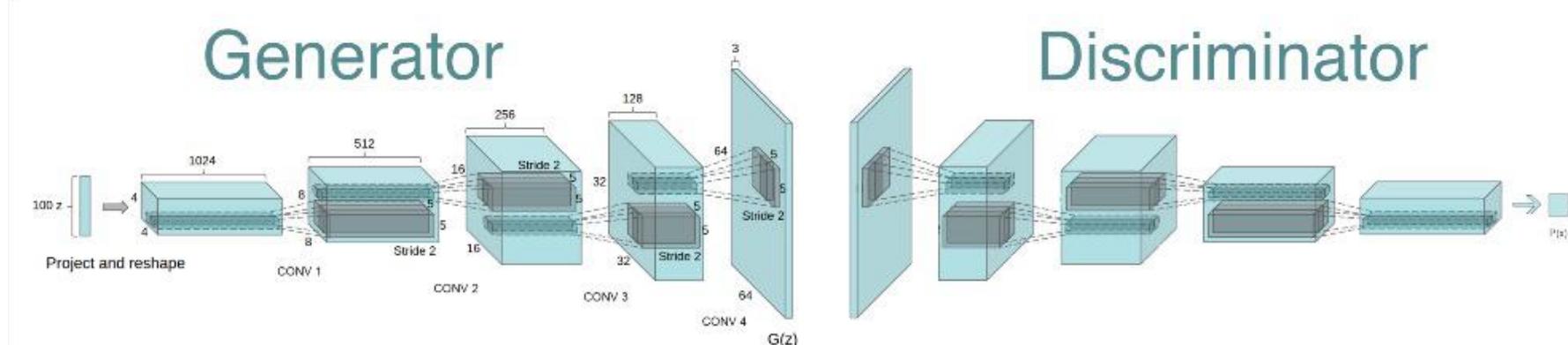
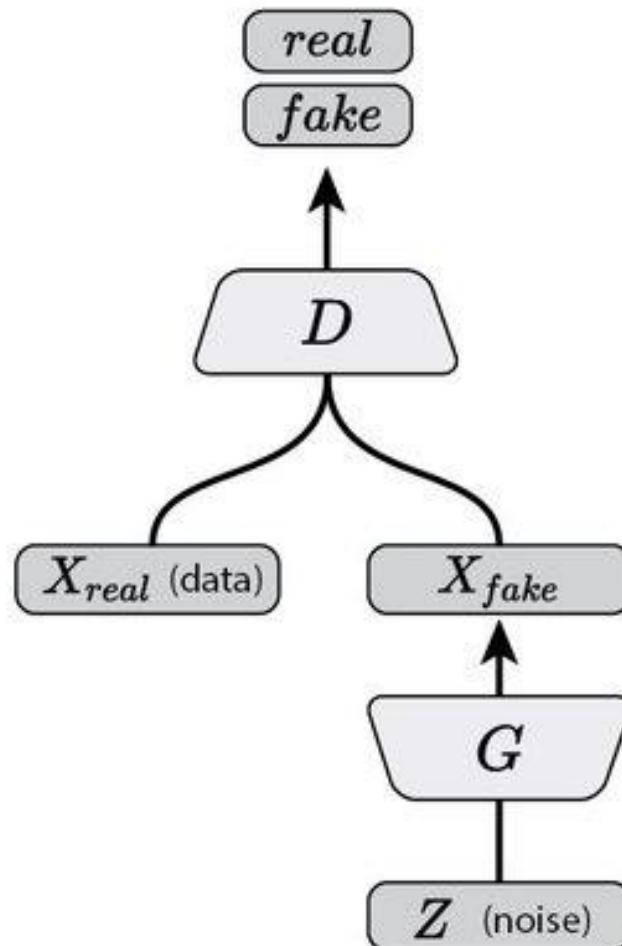
**InfoGAN**  
(Chen, et al., 2016)



**Auxiliary Classifier GAN**  
(Odena, et al., 2016)



# Vanilla GAN (Goodfellow et al., 2014)

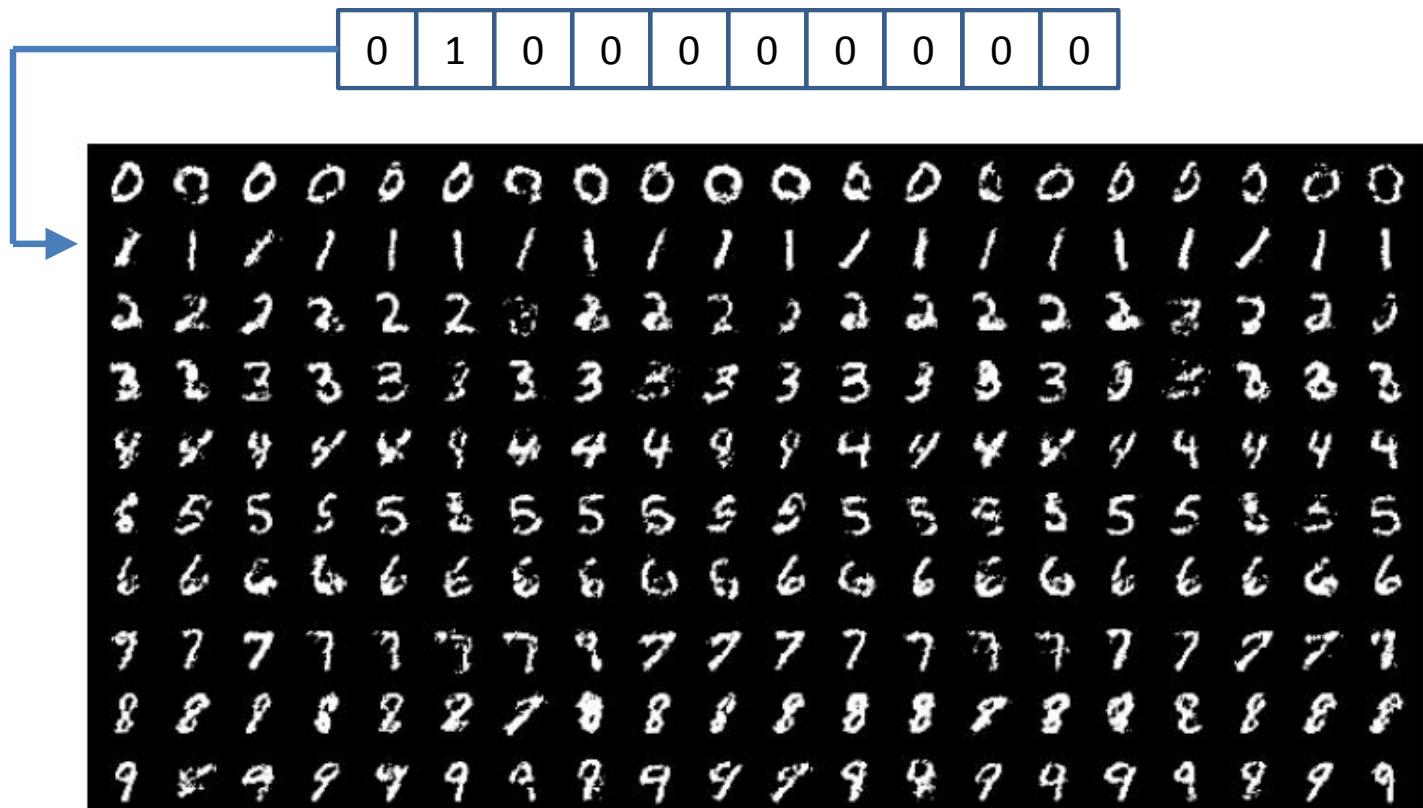
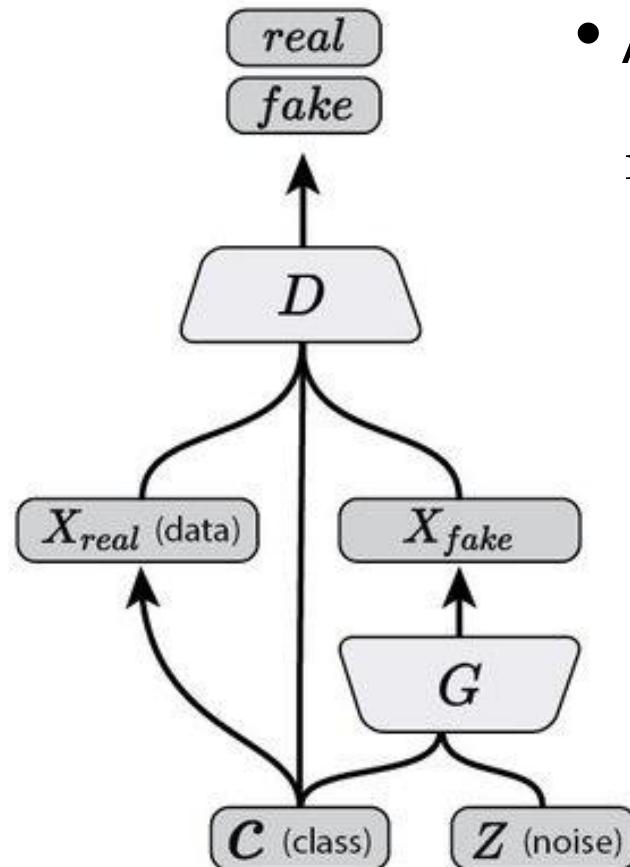


DCGAN (Radford et al., 2015)

# Conditional GAN (Mirza and Osindero, 2014)

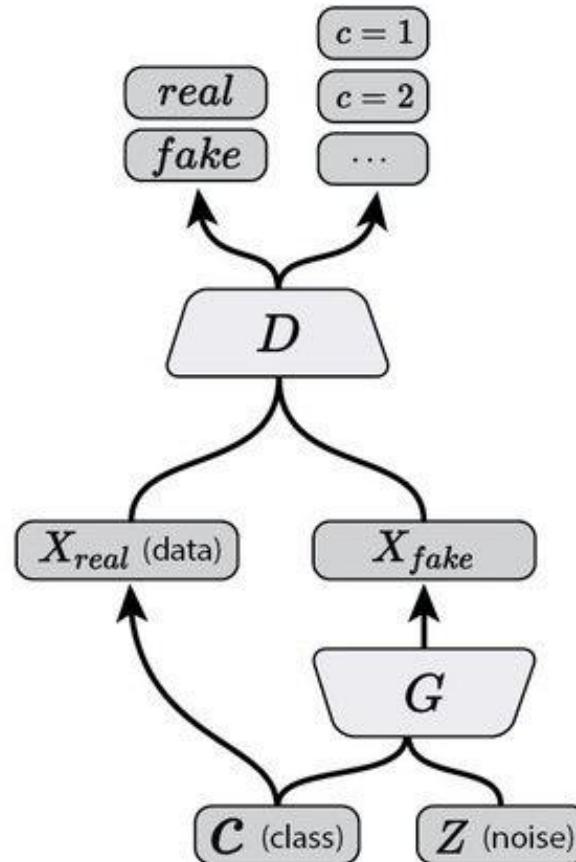
- Add conditional variables  $\mathbf{y}$  into  $G$  and  $D$

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$



# Auxiliary Classifier GAN (Odena et al., 2016)

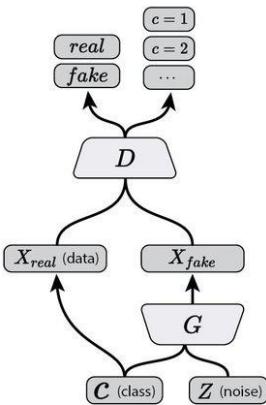
- Every generated sample has a corresponding class label



$$L_S = E[\log P(S = \text{real} \mid X_{real})] + E[\log P(S = \text{fake} \mid X_{fake})]$$
$$L_C = E[\log P(C = c \mid X_{real})] + E[\log P(C = c \mid X_{fake})]$$

- $D$  is trained to maximize  $L_S + L_C$
- $G$  is trained to maximize  $L_C - L_S$
- Learns a representation for  $z$  that is independent of class label

# Auxiliary Classifier GAN (Odena et al., 2016)



128×128 resolution samples from 5 classes taken from an AC-GAN trained on the ImageNet



monarch butterfly



goldfinch



daisy



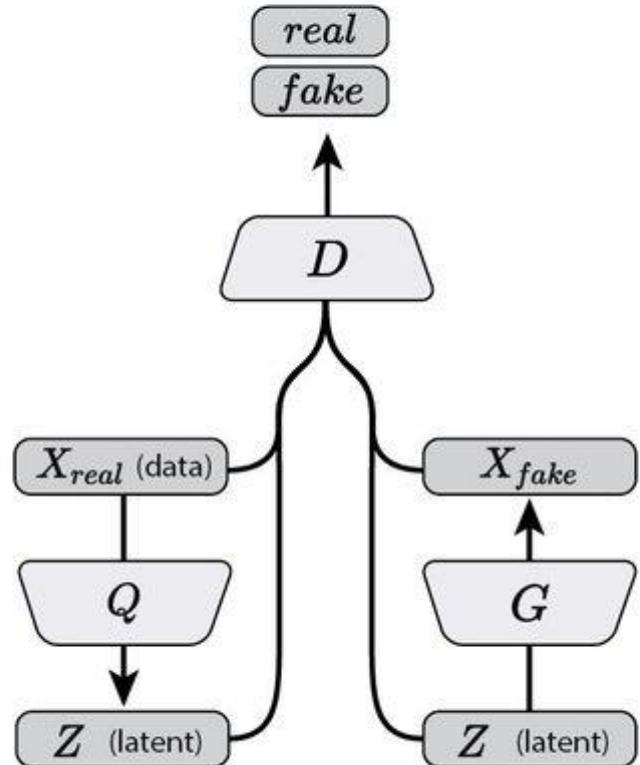
redshank



grey whale

# Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)

- Jointly learns a generator network and an inference network using an adversarial process.



$$\begin{aligned} \min_G \max_D V(D, G) &= \mathbb{E}_{q(\mathbf{x})}[\log(D(\mathbf{x}, G_z(\mathbf{x})))] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - D(G_x(\mathbf{z}), \mathbf{z}))] \\ &= \iint q(\mathbf{x})q(\mathbf{z} \mid \mathbf{x}) \log(D(\mathbf{x}, \mathbf{z})) d\mathbf{x} d\mathbf{z} \\ &+ \iint p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z}) \log(1 - D(\mathbf{x}, \mathbf{z})) d\mathbf{x} d\mathbf{z}. \end{aligned}$$



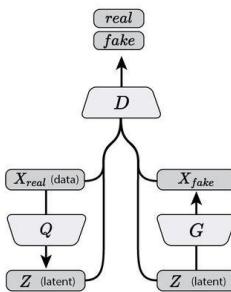
CelebA reconstructions



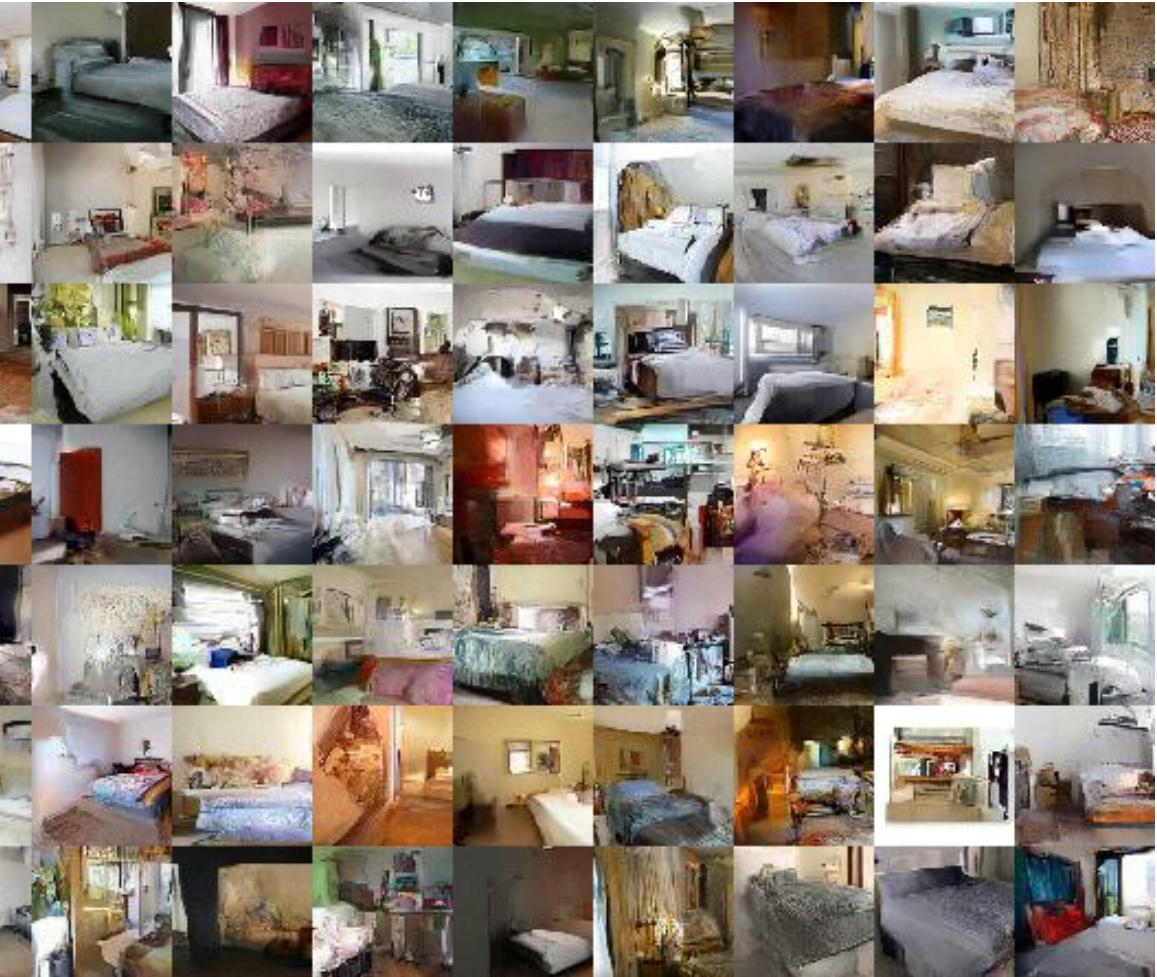
SVNH reconstructions

# Bidirectional GAN

(Donahue et al., 2016;  
Dumoulin et al., 2016)



LSUN bedrooms



Tiny ImageNet

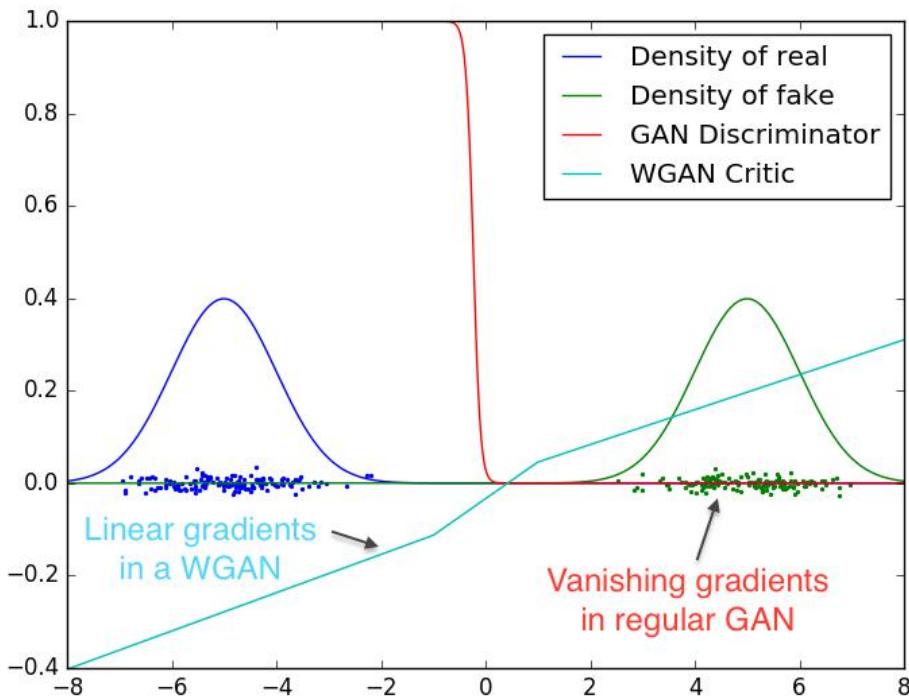


# Wasserstein GAN (Arjovsky et al., 2016)

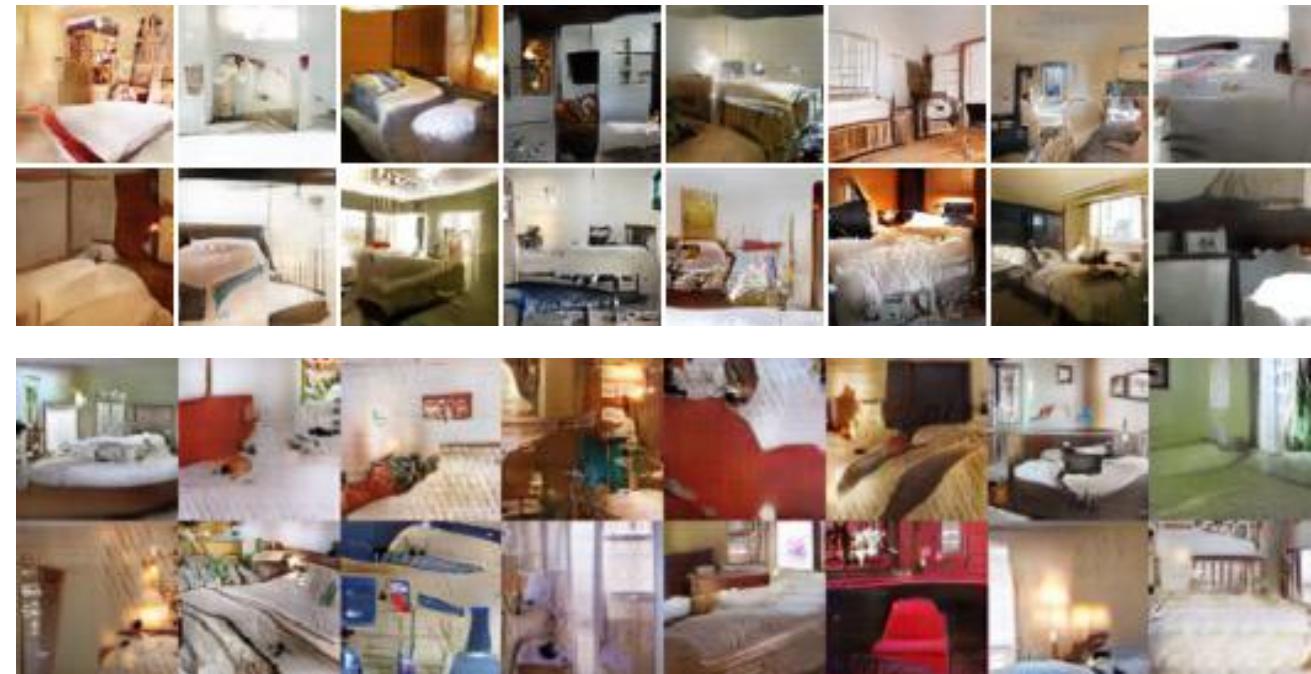
- Objective based on Earth-Mover or Wasserstein distance:

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [D_{\omega}(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [D_{\omega}(G_{\theta}(\mathbf{z}))]$$

- Provides nice gradients over real and fake samples

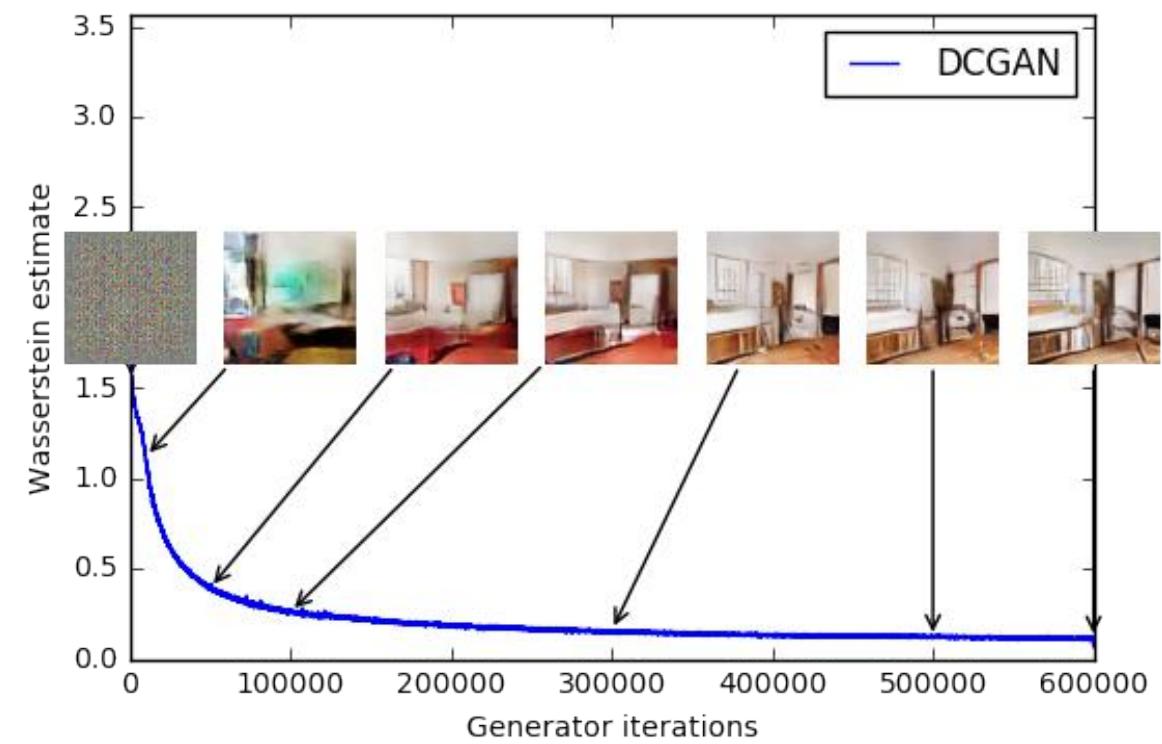
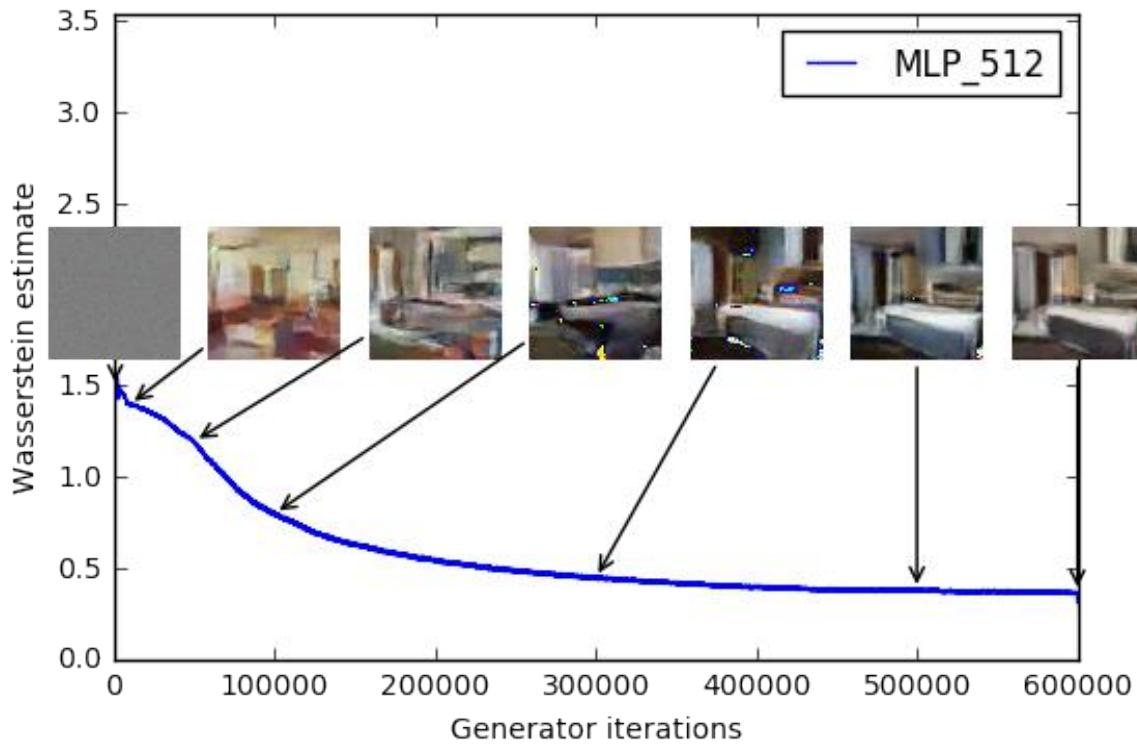


WGAN



# Wasserstein GAN (Arjovsky et al., 2016)

- Wasserstein loss seems to correlate well with image quality.



# WGAN with gradient penalty (Gulraani et al., 2017)

$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}}$$

- Faster convergence and higher-quality samples than WGAN with weight clipping
- Train a wide variety of GAN architectures with almost no hyperparameter tuning, including discrete models

Samples from a character-level GAN language model on Google Billion Word

## WGAN with gradient penalty

Busino game camperate spent odea  
In the bankaway of smarling the  
SingersMay , who kill that imvic  
Keray Pents of the same Reagun D  
Manging include a tudancs shat "  
His Zuith Dudget , the Denmbern  
In during the Uitational questio  
Divos from The ' noth ronkies of  
She like Monday , of macunsuer S  
The investor used ty the present  
A papees are cointry congress oo  
A few year inom the group that s  
He said this syenn said they wan  
As a world 1 88 ,for Autouries  
Foand , th Word people car , Il  
High of the upsideader homing pull  
The guipe is worly move dogsfor  
The 1874 incidested he could be  
The allo tooks to security and c

Solice Norkedin pring in since  
This record ( 31. ) UBS ) and Ch  
It was not the annuas were plogr  
This will be us , the ect of DAN  
These leaded as most-worsd p2 a0  
The time I paid0a South Cubry i  
Dour Fraps higs it was these del  
This year out howneed allowed lo  
Kaulna Seto consficates to repor  
A can teal , he was schoon news  
In th 200. Pesish picriers rega  
Konney Panice rimimber the teami  
The new centuct cut Denester of  
The near , had been one injostie  
The incestion to week to shorted  
The company the high product of  
20 - The time of accomplete , wh  
John WVuderenson seqivic spends  
A ceetens in indestdredly the Wat

## Standard GAN objective

ddddddddd ddd ddd ddd ddd ddd ddd  
dd ddd ddd ddd ddd ddd ddd ddd ddd ddd

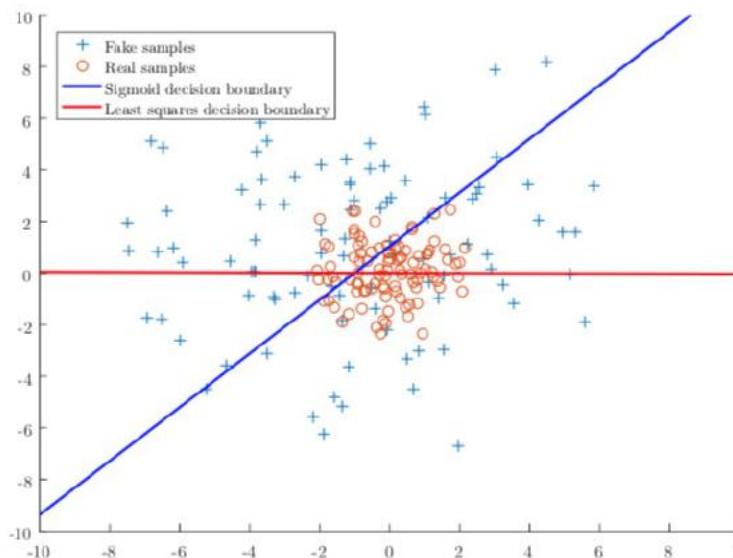
dd ddd ddd ddd ddd ddd ddd ddd ddd ddd  
dd ddd ddd ddd ddd ddd ddd ddd ddd ddd

# Least Squares GAN (LSGAN) (Mao et al., 2017)

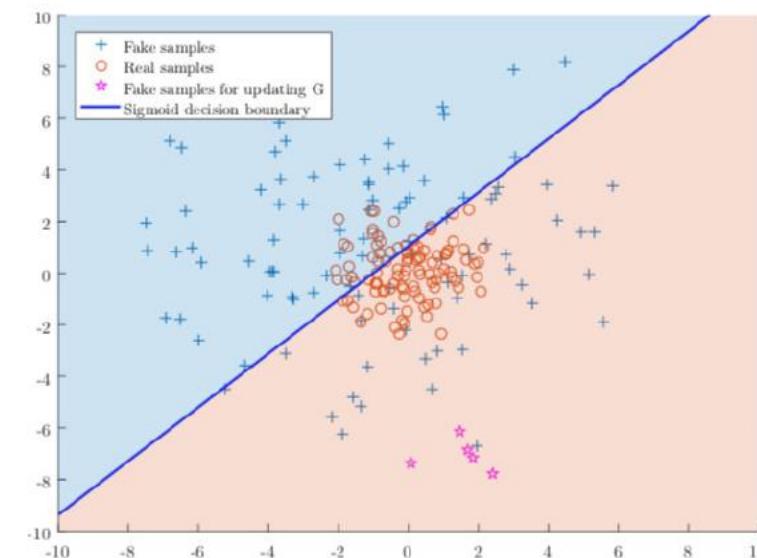
- Use a loss function that provides smooth and non-saturating gradient in discriminator D

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - b)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})) - a)^2]$$

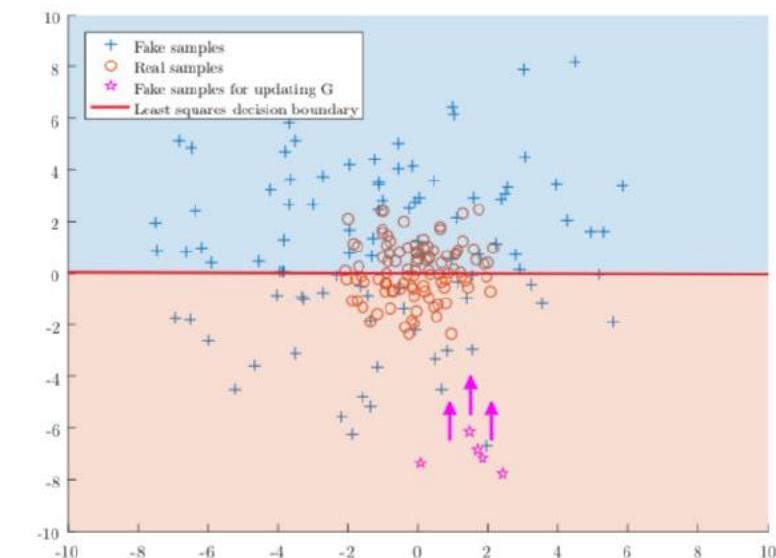
$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z})) - c)^2],$$



Decision boundaries of Sigmoid & Least Squares loss functions



Sigmoid decision boundary

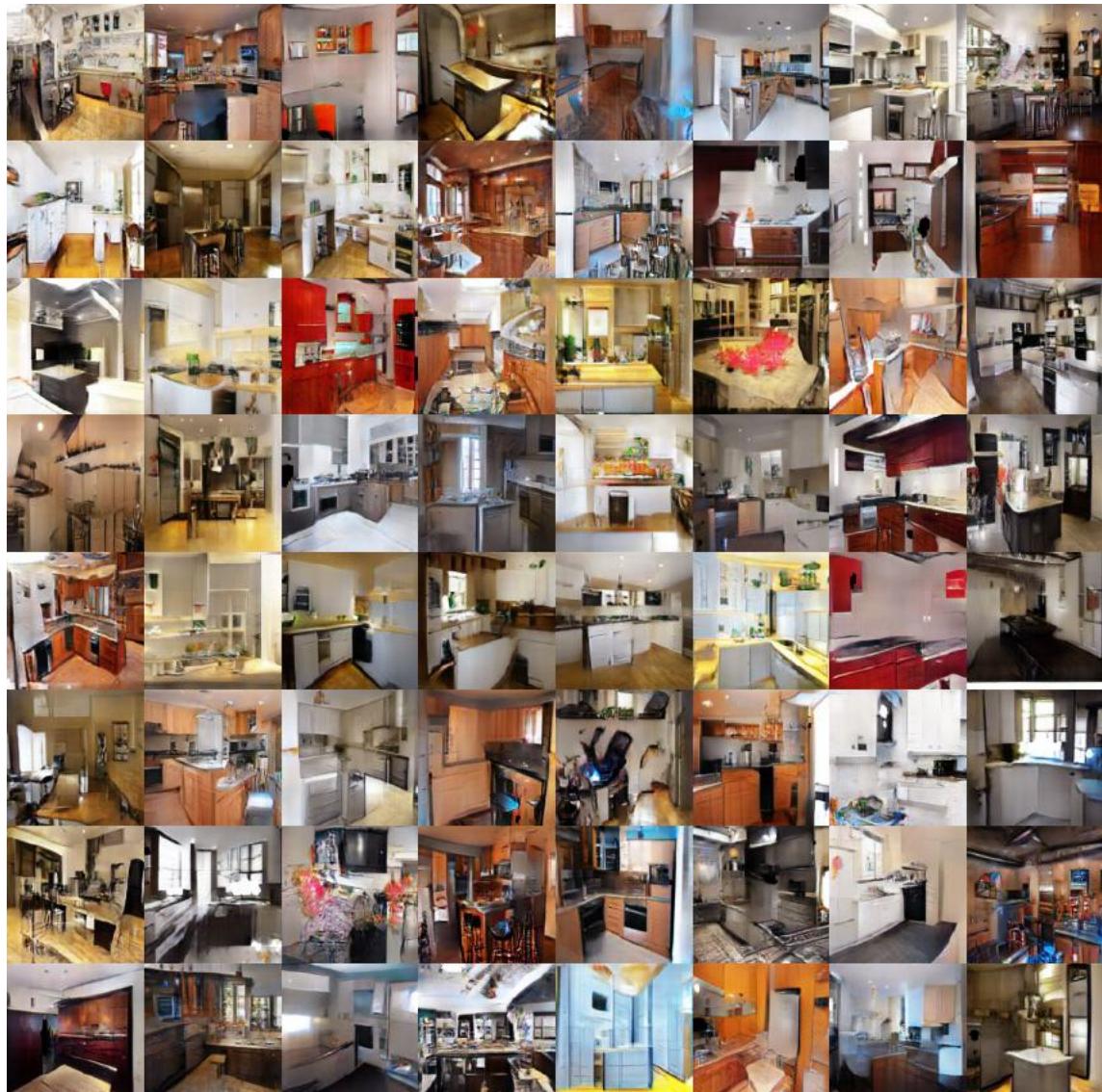


Least Squares decision boundary

# Least Squares GAN (LSGAN) (Mao et al., 2017)



Church



Kitchen

# Boundary Equilibrium GAN (BEGAN)

(Berthelot et al., 2017)

- A loss derived from the Wasserstein distance for training auto-encoder based GANs

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} \\ \eta \in \{1, 2\} \\ v \in \mathbb{R}^{N_x} \end{cases}$$

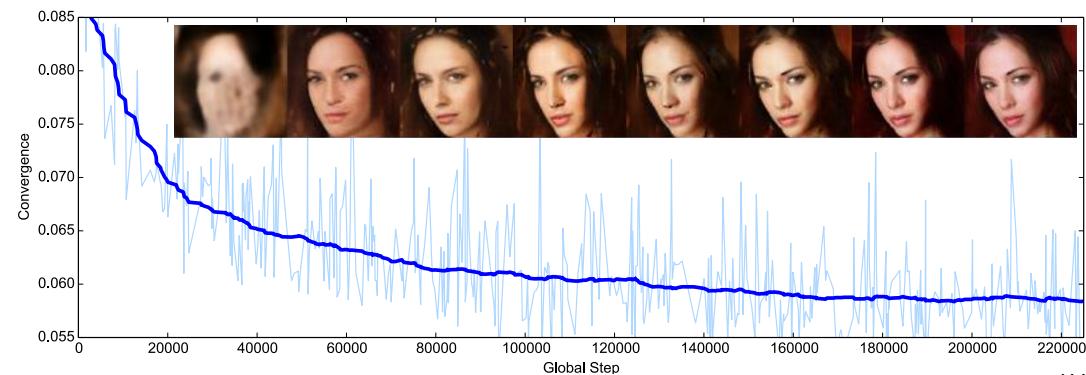
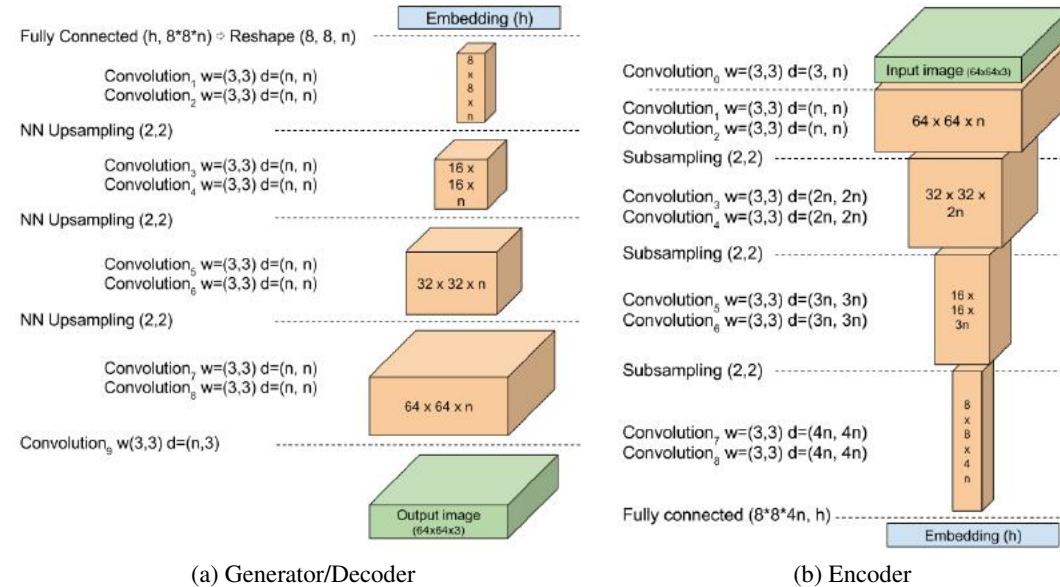
is the autoencoder function.  
is the target norm.  
is a sample of dimension  $N_x$ .

- Wasserstein distance btw. the reconstruction losses of real and generated data
- Convergence measure:  

$$\mathcal{M}_{global} = \mathcal{L}(x) + |\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))|$$
- Objective:

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x) - k_t \cdot \mathcal{L}(G(z_D)) \\ \mathcal{L}_G = \mathcal{L}(G(z_G)) \\ k_{t+1} = k_t + \lambda_k (\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))) \end{cases}$$

for  $\theta_D$   
for  $\theta_G$   
for each training step  $t$



# BEGANs for CelebA

360K celebrity face images  
128x128 with 128 filters

(Berthelot et al., 2017)



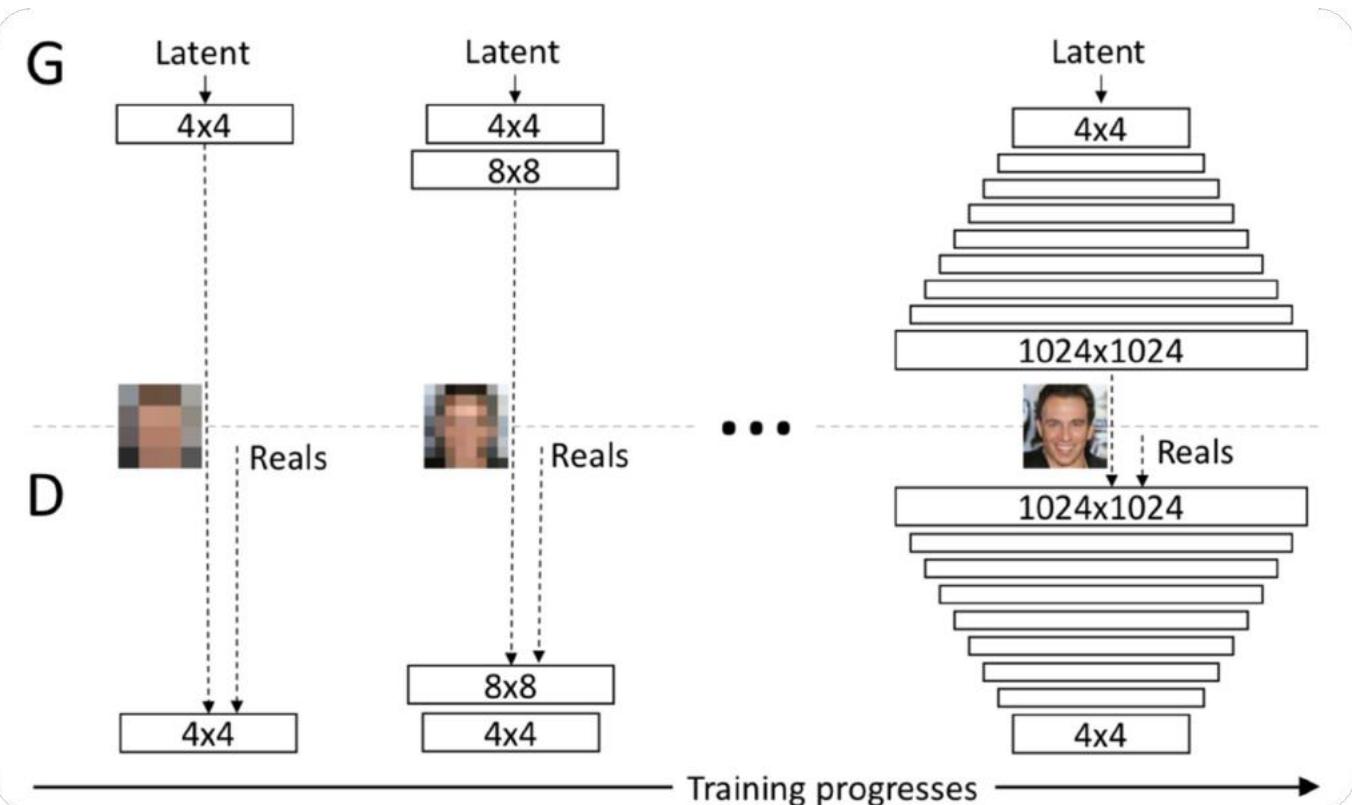
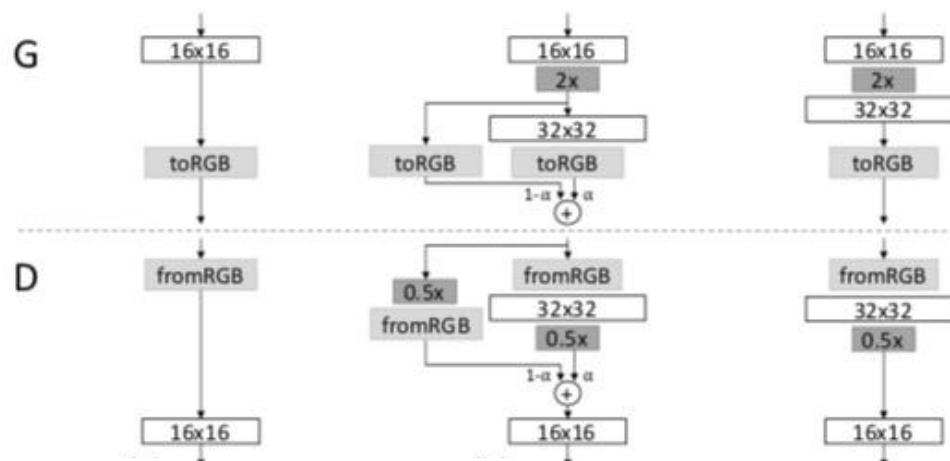
Interpolations in the latent space



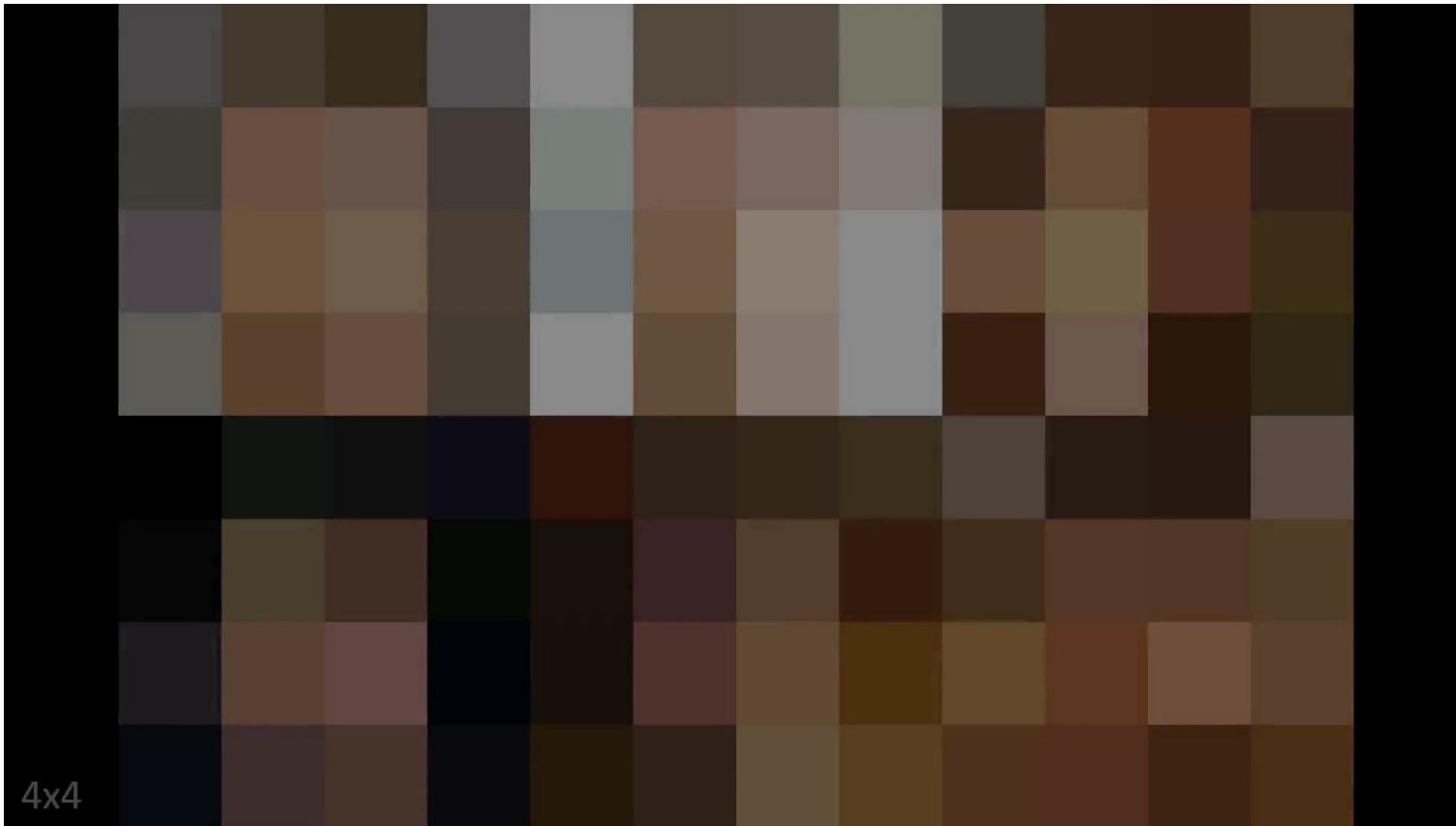
Mirror interpolation example

# Progressive GANs (Karras et al., 2018)

- Progressively generate high-res images
- Multi-step training from low to high resolutions



# Progressive GANs (Karras et al., 2018)



- Training process

# Progressive GANs (Karras et al., 2018)

CelebA-HQ  
random interpolations



# BigGANs (Brock et al., 2019)

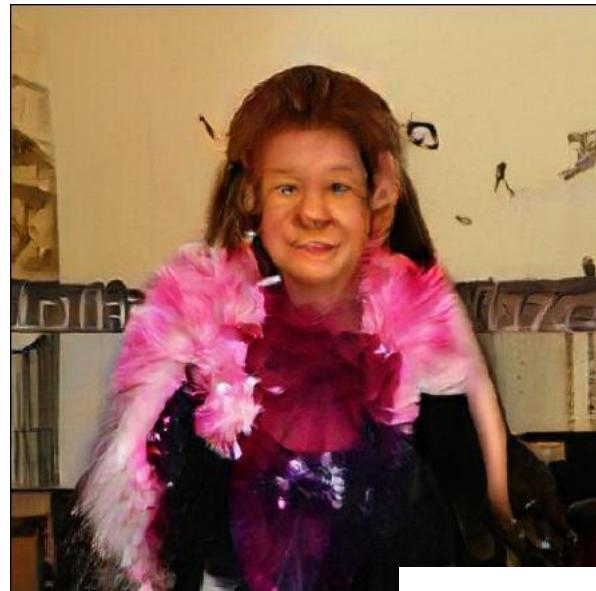
High resolution, class-conditional samples generated by the model



- BigGANs trained with 2-4x as many parameters and 8x the batch size compared to prior art.
- Uses Gaussian truncation to sample z (avoid sampling from the tail of the Gaussian distribution)
- Uses multiple other tricks including multiple regularizations including a Gradient penalty regularization and an Orthogonal Regularization:

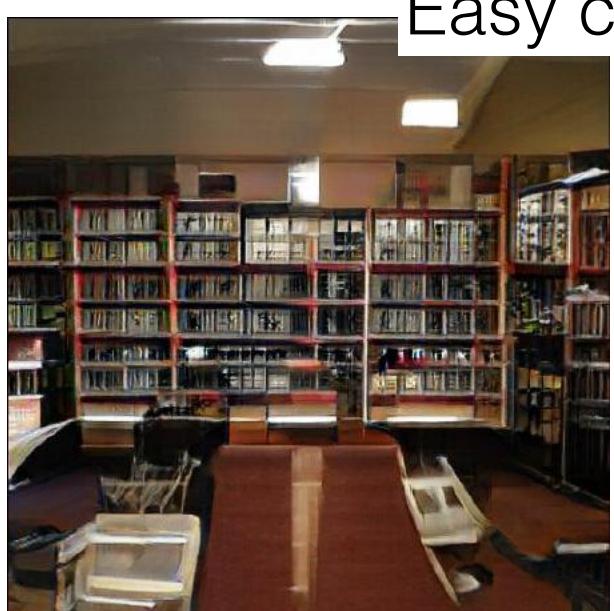
$$R_\beta(W) = \beta \|W^\top W \odot (\mathbf{1} - I)\|_F^2,$$

# BigGANs (Brock et al., 2019)



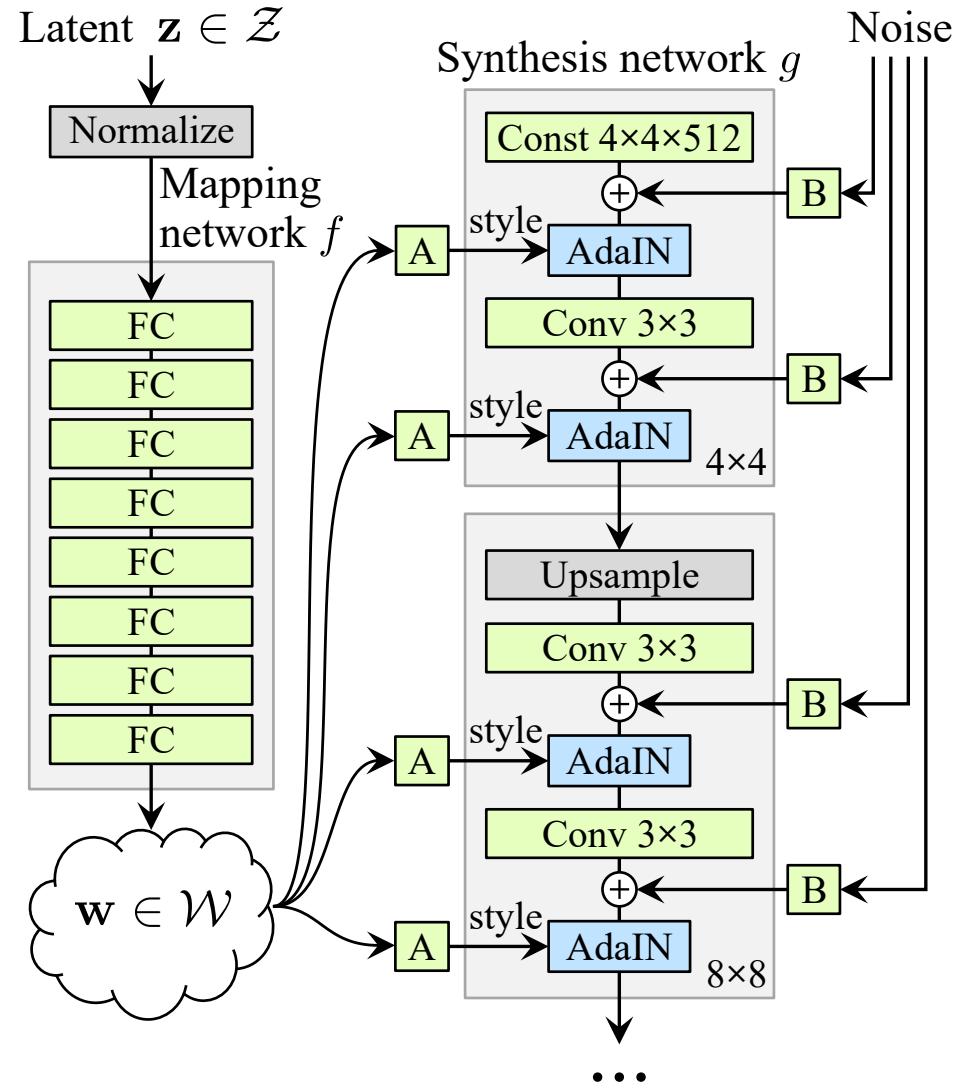
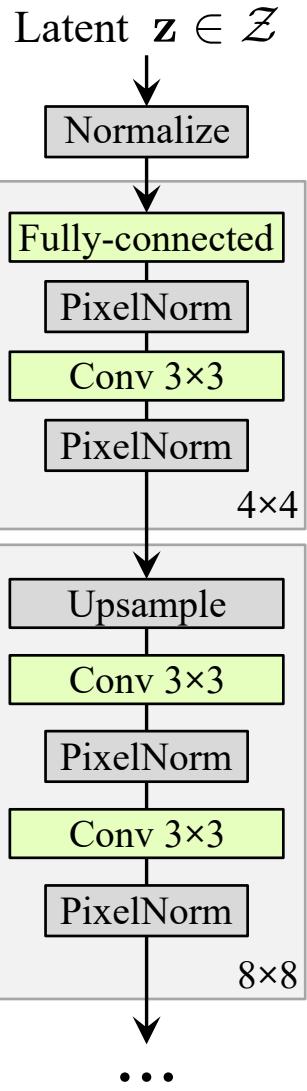
Easy classes

Hard classes



Resolution: 512x512

# StyleGAN (Karras et al., 2019)



$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i},$$



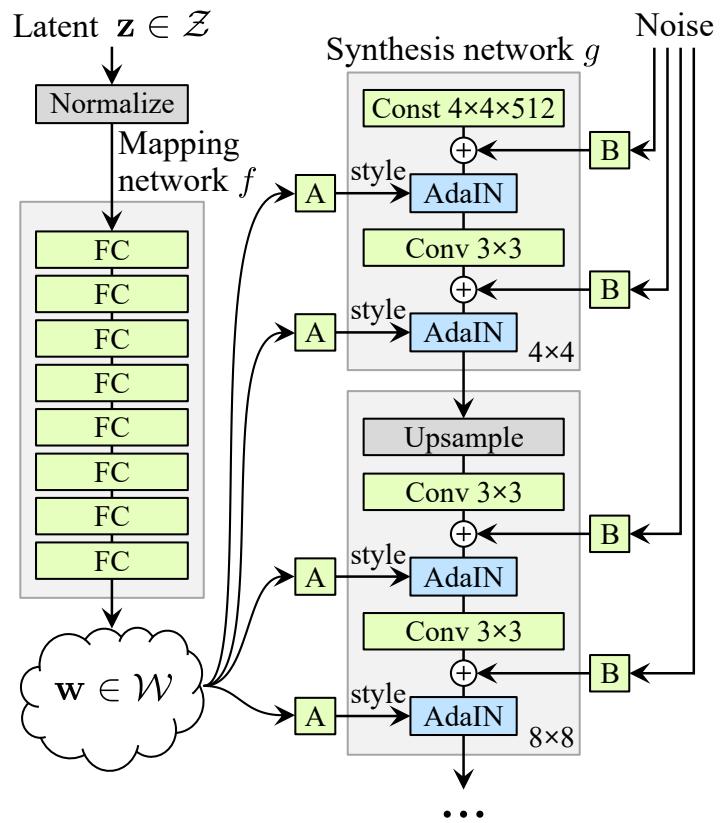
Samples (trained on the FFHQ dataset)

(a) Traditional

(b) Style-based generator

# StyleGAN (Karras et al., 2019)

- Swapping out the destination style for the source style



# Some Applications of GANs

# Semi-supervised Classification

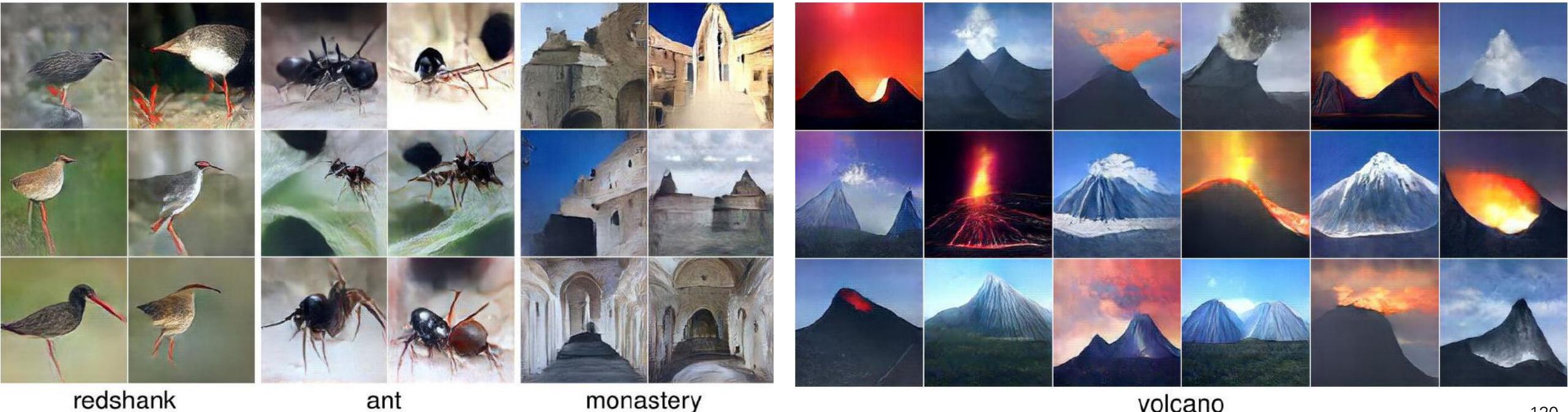
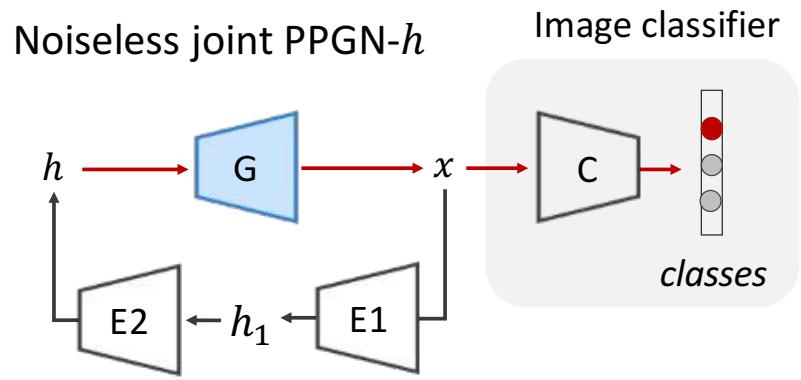
(Salimans et al., 2016;  
Dumoulin et al., 2016)

## SVNH

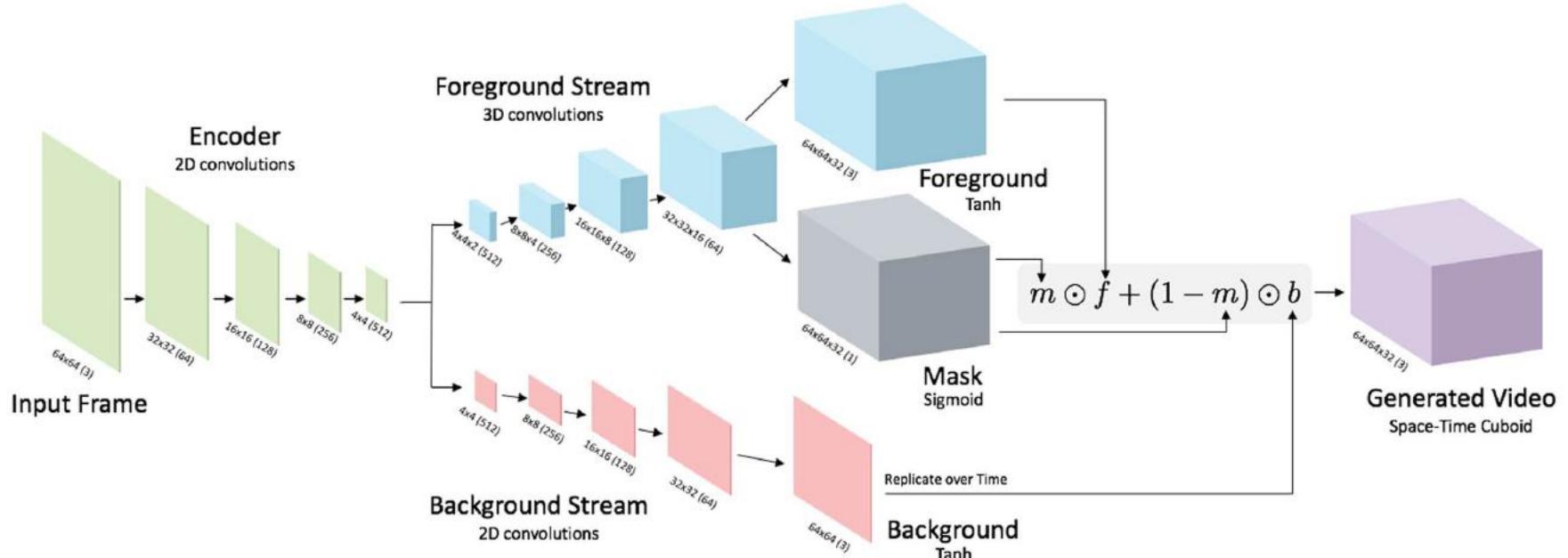
Model	Misclassification rate
VAE (M1 + M2) (Kingma et al., 2014)	36.02
SWWAE with dropout (Zhao et al., 2015)	23.56
DCGAN + L2-SVM (Radford et al., 2015)	22.18
SDGM (Maaløe et al., 2016)	16.61
<b>GAN (feature matching) (Salimans et al., 2016)</b>	<b><math>8.11 \pm 1.3</math></b>
ALI (ours, L2-SVM)	$19.14 \pm 0.50$
<b>ALI (ours, no feature matching)</b>	<b><math>7.42 \pm 0.65</math></b>

# Class-specific Image Generation (Nguyen et al., 2016)

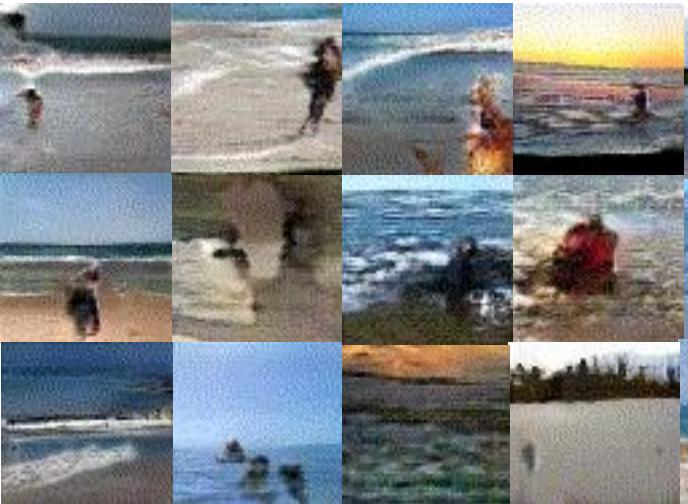
- Generates 227x227 realistic images from all ImageNet classes
- Combines adversarial training, moment matching, denoising autoencoders, and Langevin sampling



# Video Generation (Vondrick et al., 2016)



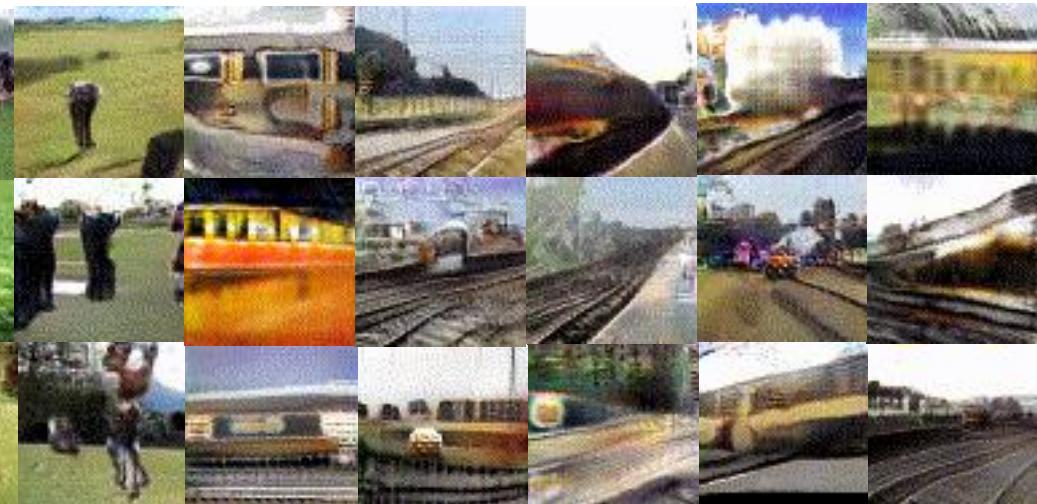
Beach



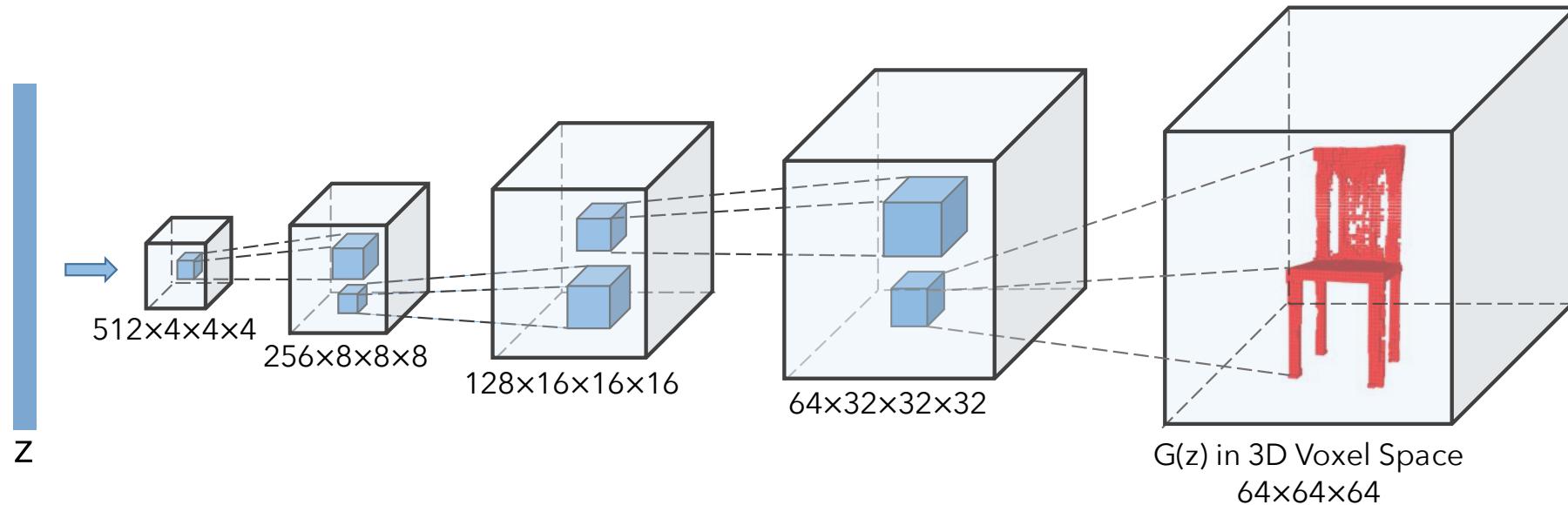
Golf



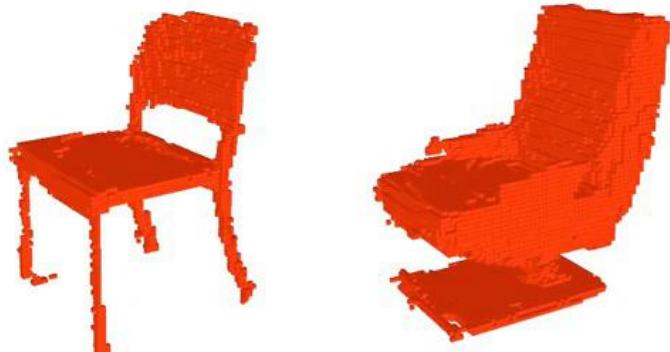
Train Station



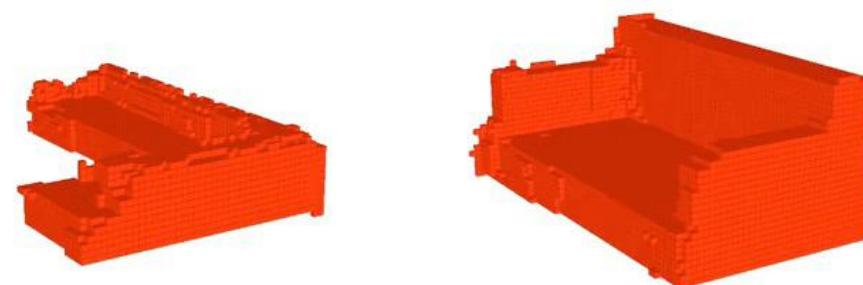
# Generative Shape Modeling (Wu et al., 2016)



Chairs



Sofas



# Text-to-Image Synthesis (Zhang et al., 2016)

The small bird has a red head with feathers that fade from red to gray from head to tail



The petals of this flower are white with a large stigma

A unique yellow flower with no visible pistils protruding from the center

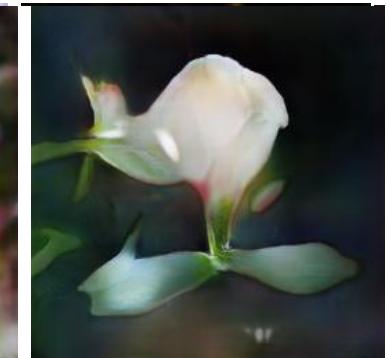
This flower is pink and yellow in color, with petals that are oddly shaped

This is a light colored flower with many different petals on a green stem

This flower is yellow and green in color, with petals that are ruffled

The flower have large petals that are pink with yellow on some of the petals

A flower that has white petals with some tones of yellow and green filaments

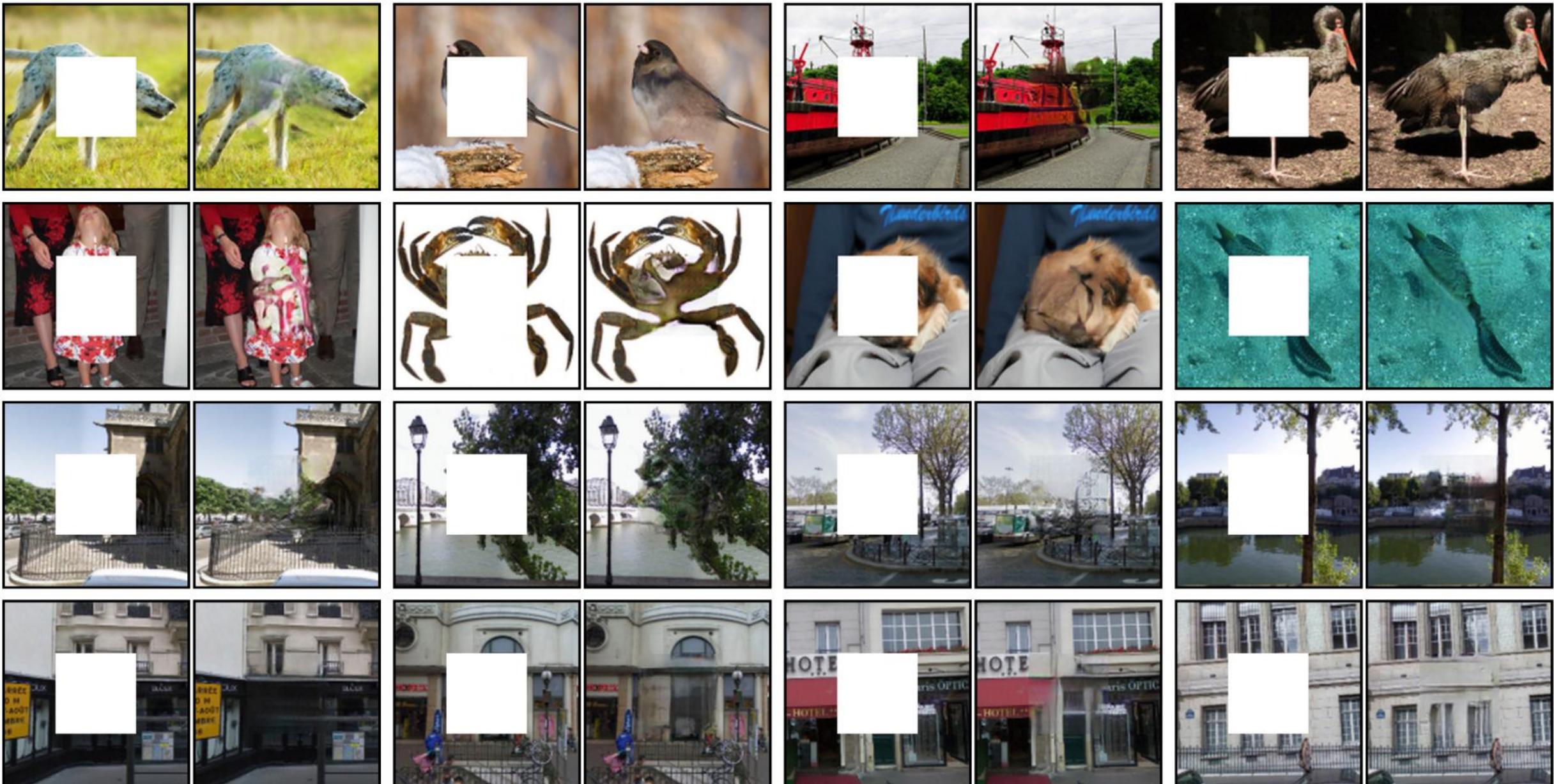


# Single Image Super-Resolution (Ledig et al., 2016)

- Combine content loss with adversarial loss



# Image Inpainting (Pathak et al., 2016)



# Unsupervised Domain Adaptation (Bousmalis et al., 2016)

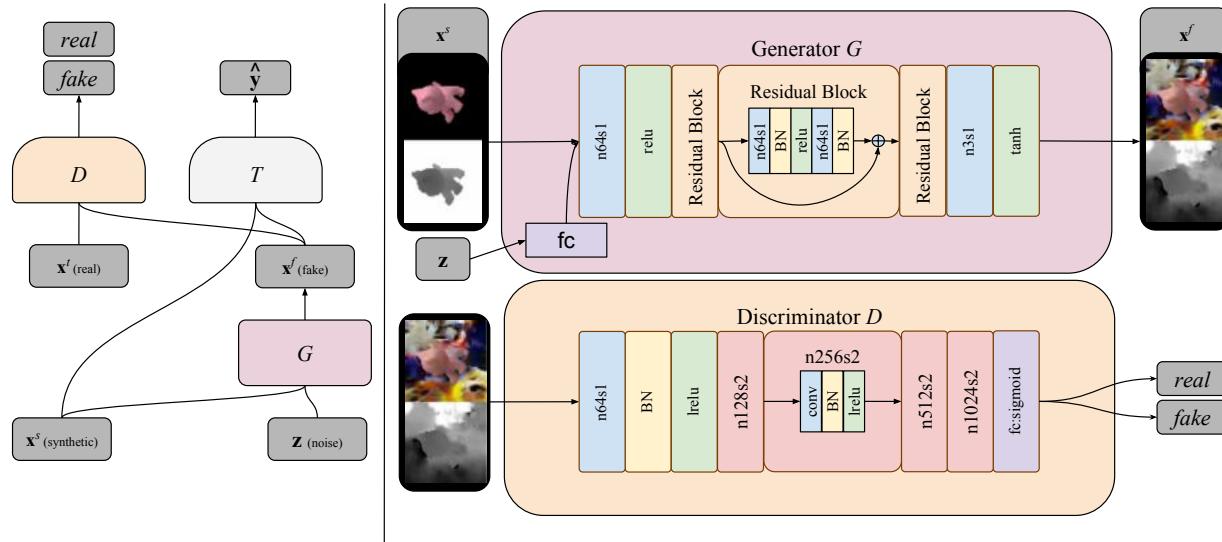


Image examples from the Linemod dataset

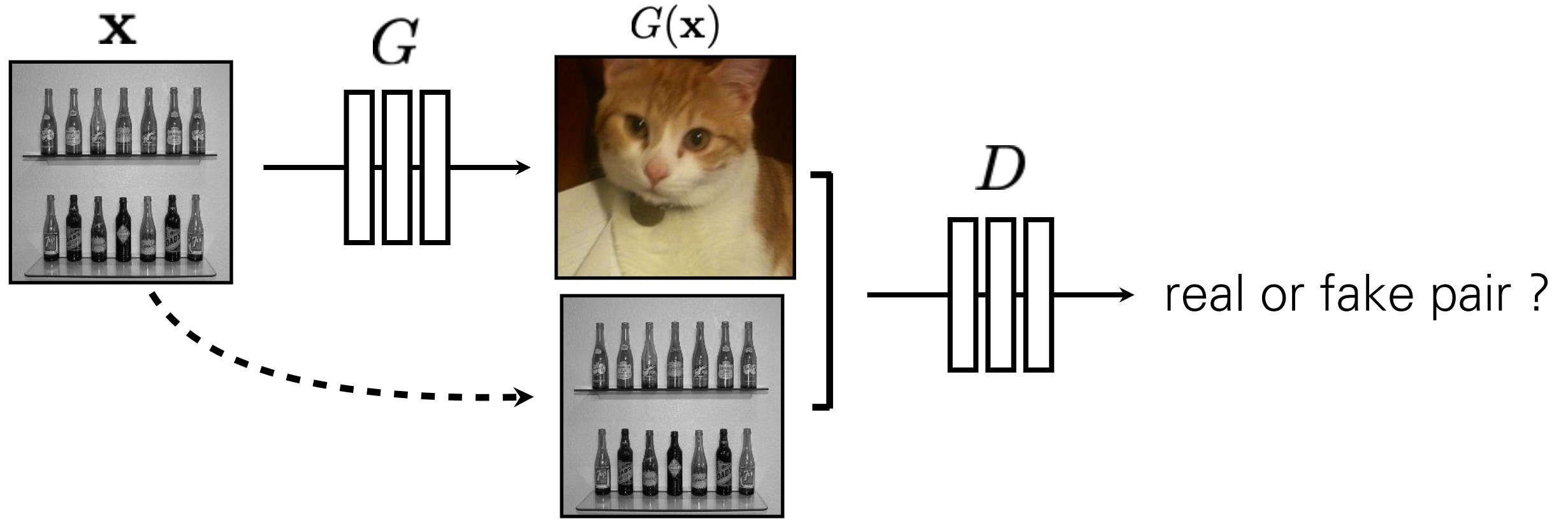


RGDB image samples  
(conditioned on a synthetic image)

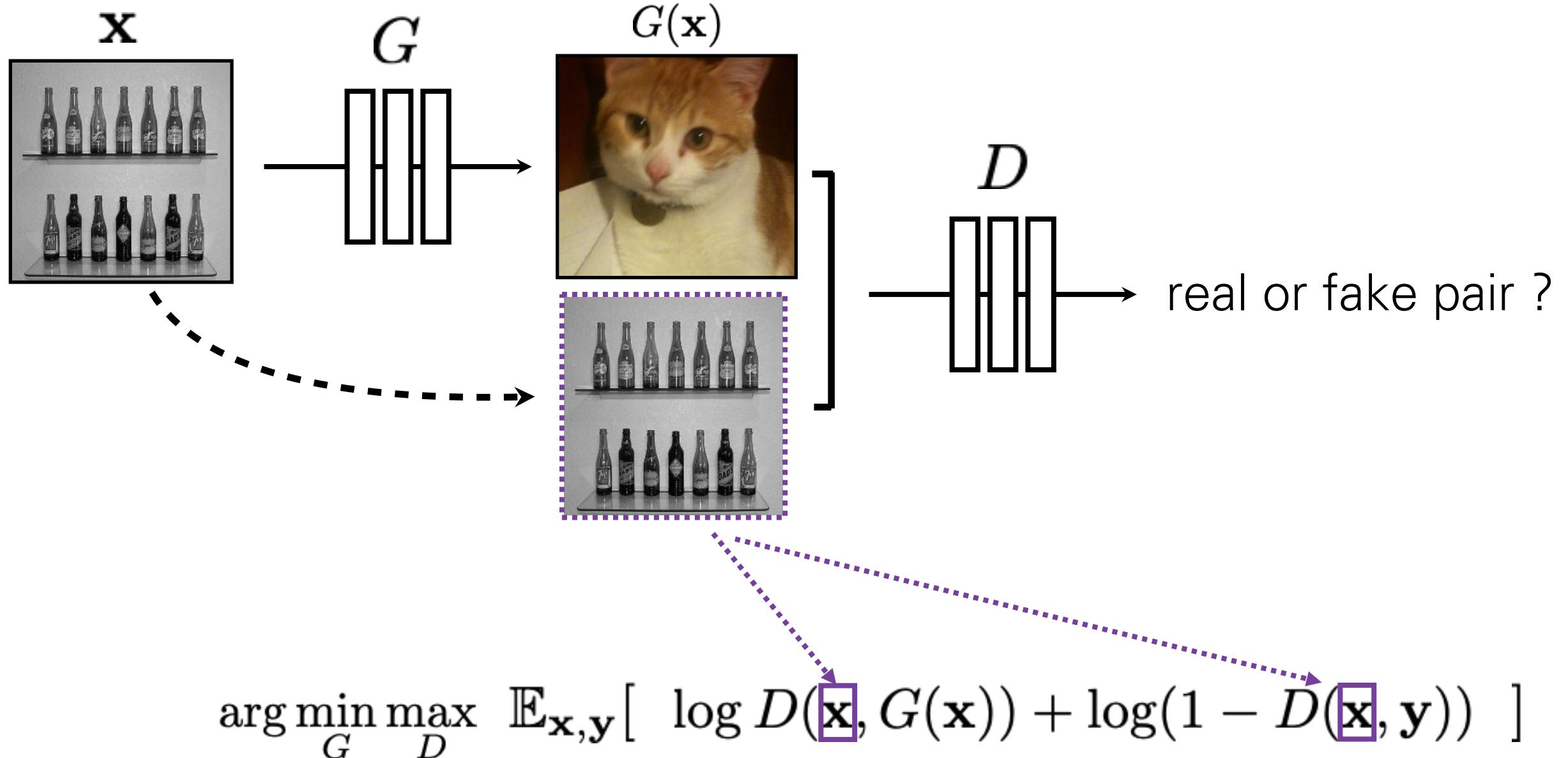
# Image to Image Translation (Pix2Pix)

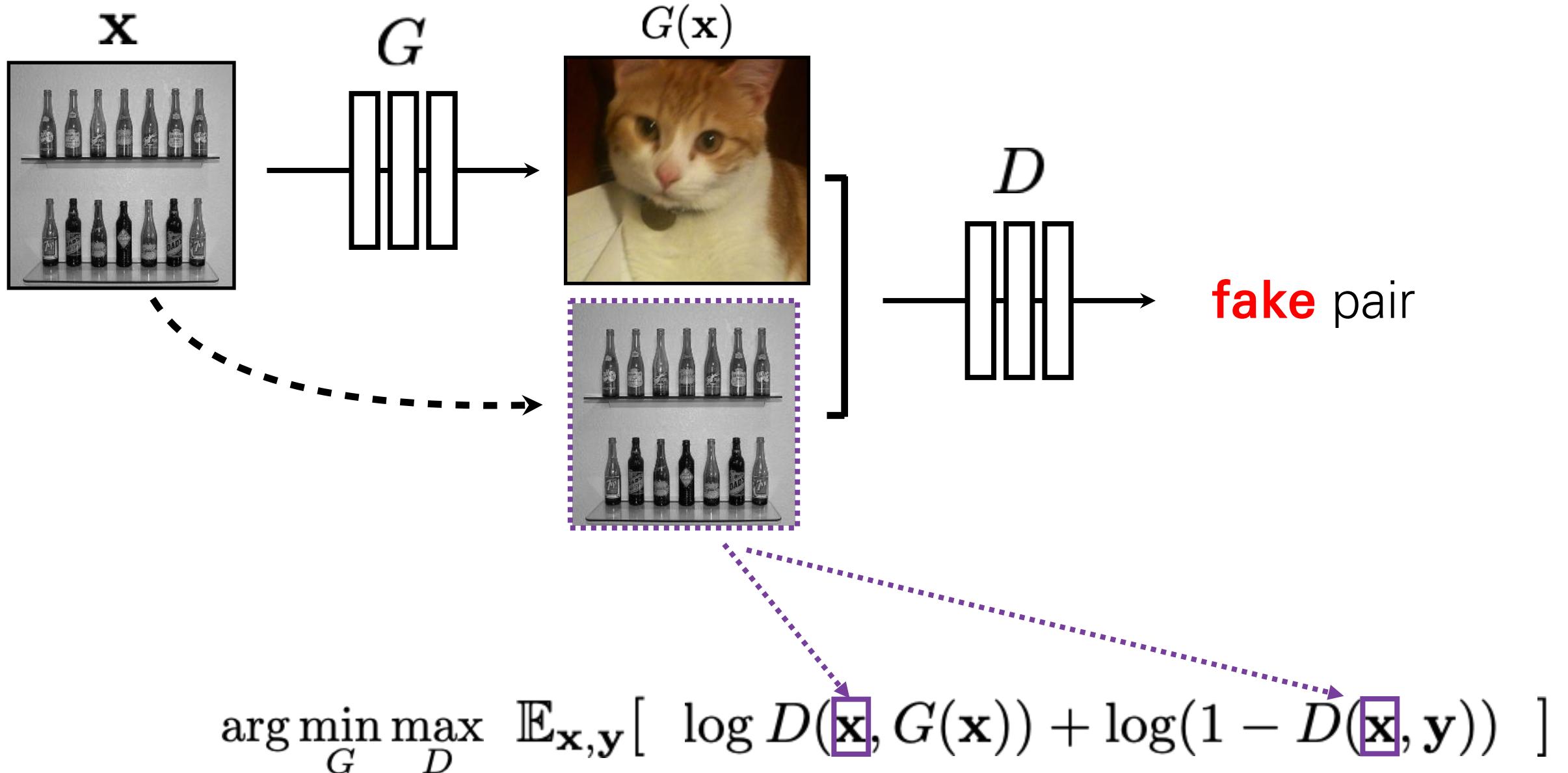


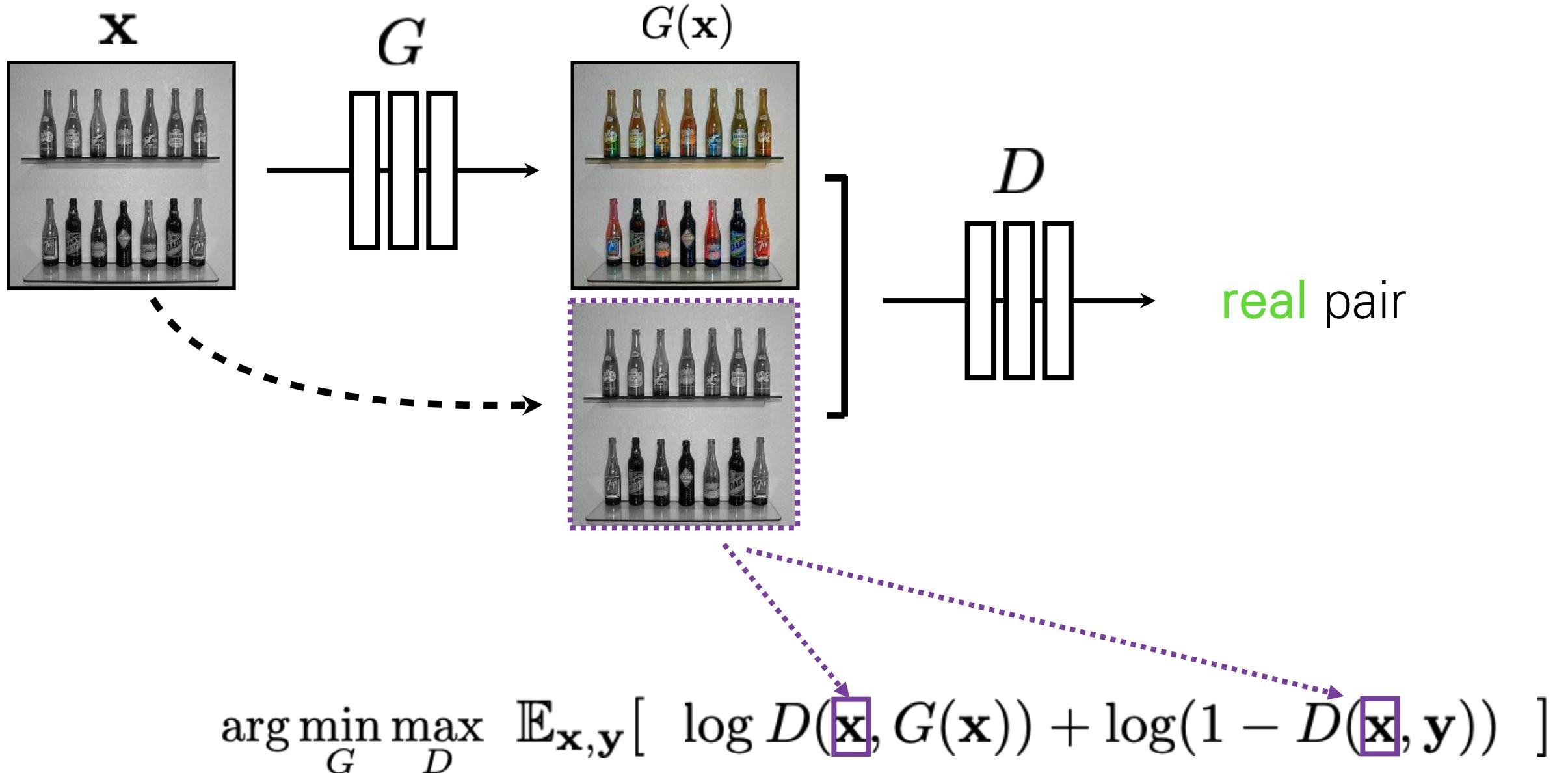
(Isola et al. 2016)

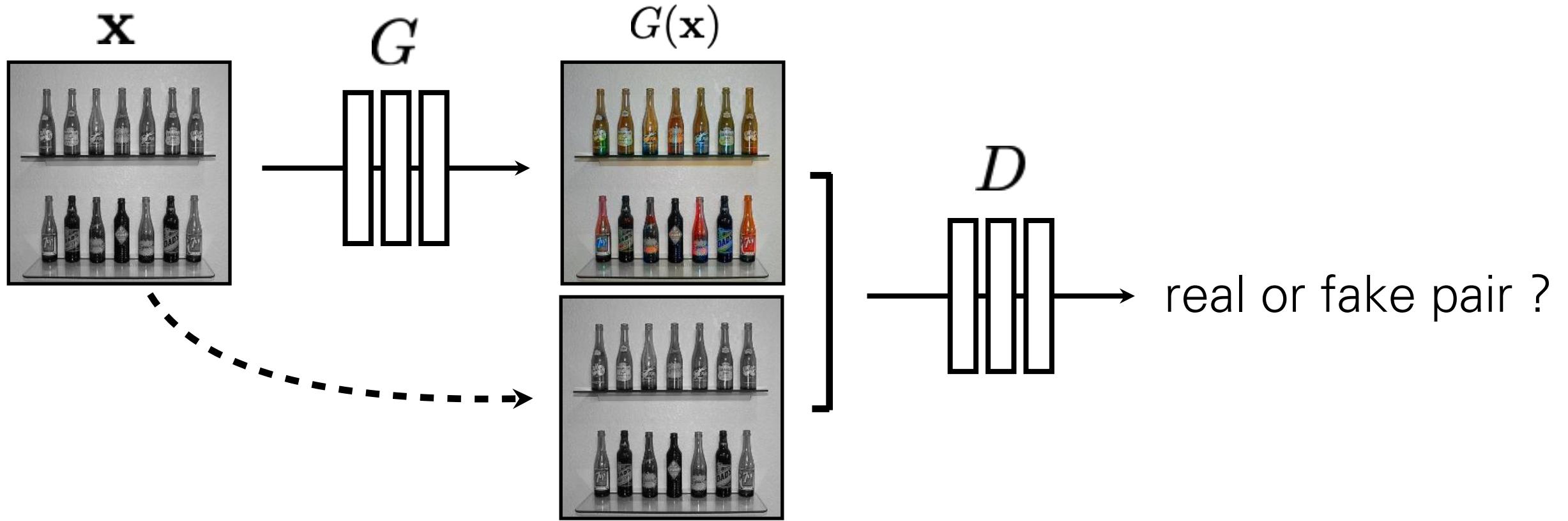


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$









$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$

# BW → Color

Input



Output



Input



Output



Input

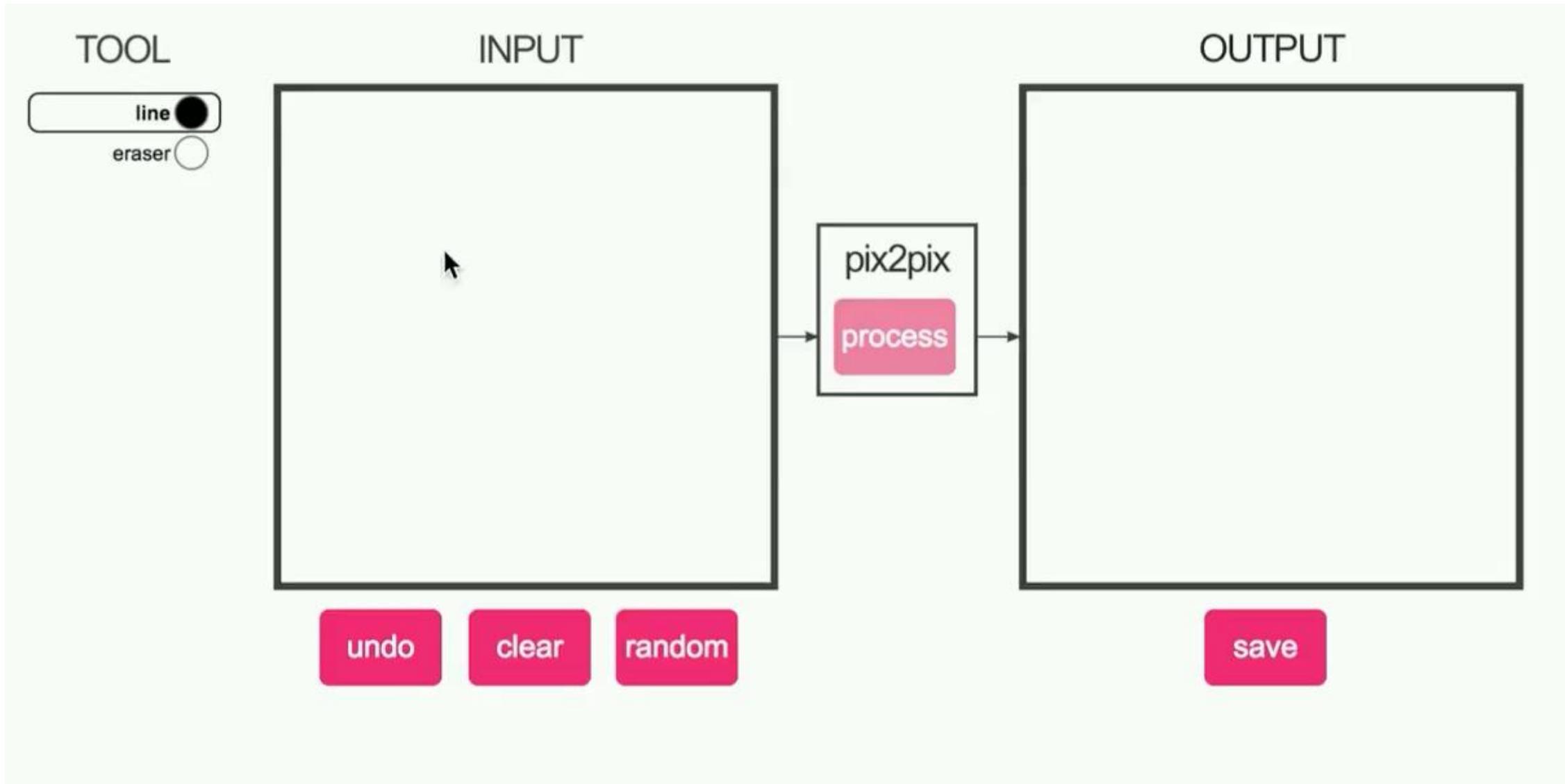


Output

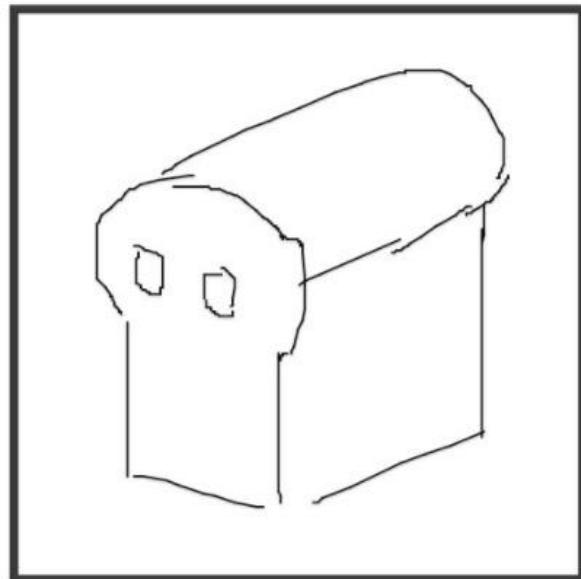


Data from [Russakovsky et al. 2015]

# #edges2cats [Chris Hesse]



INPUT



OUTPUT

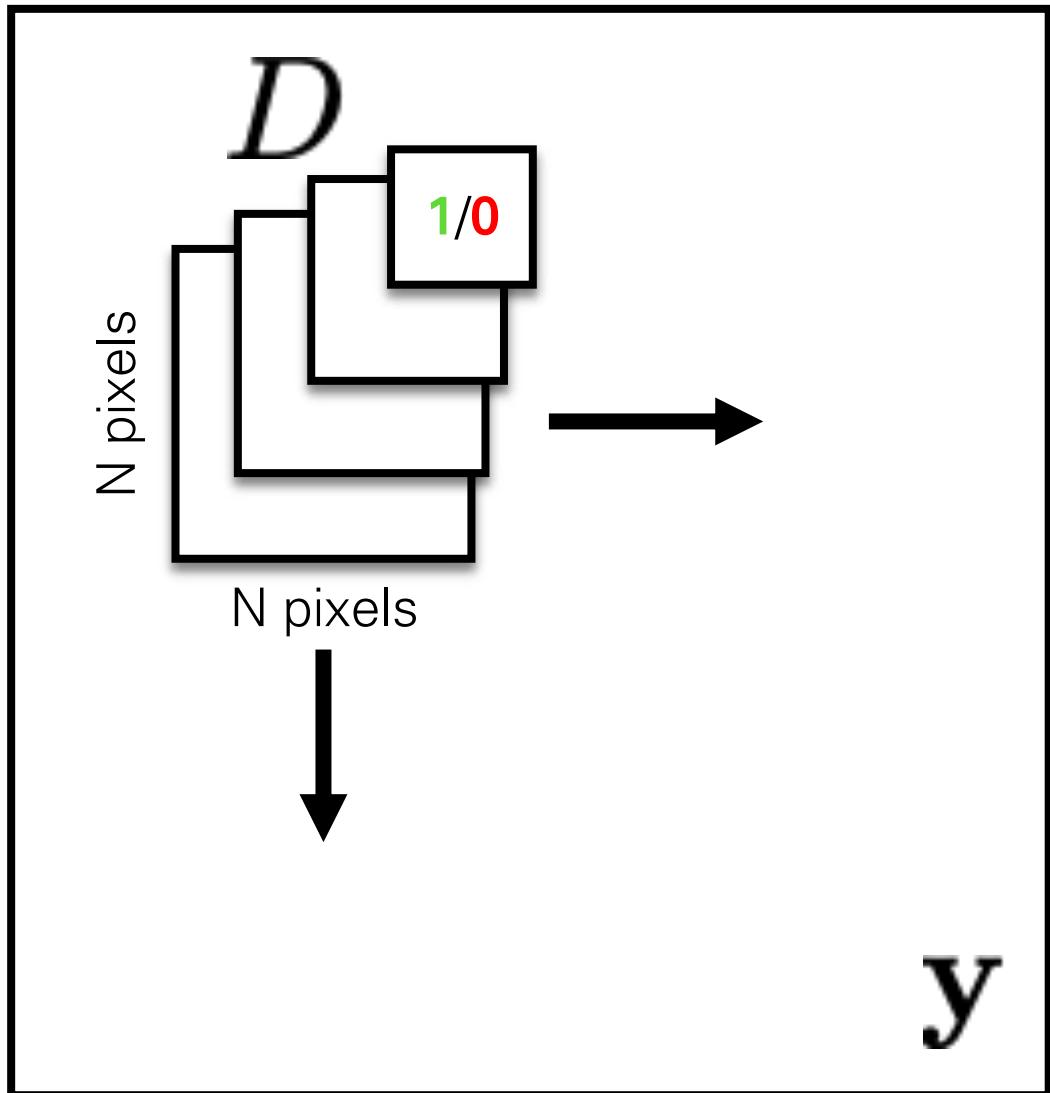


Ivy Tasi @ivymyt



Vitaly Vidmirov @vvid

# Shrinking the capacity: Patch Discriminator



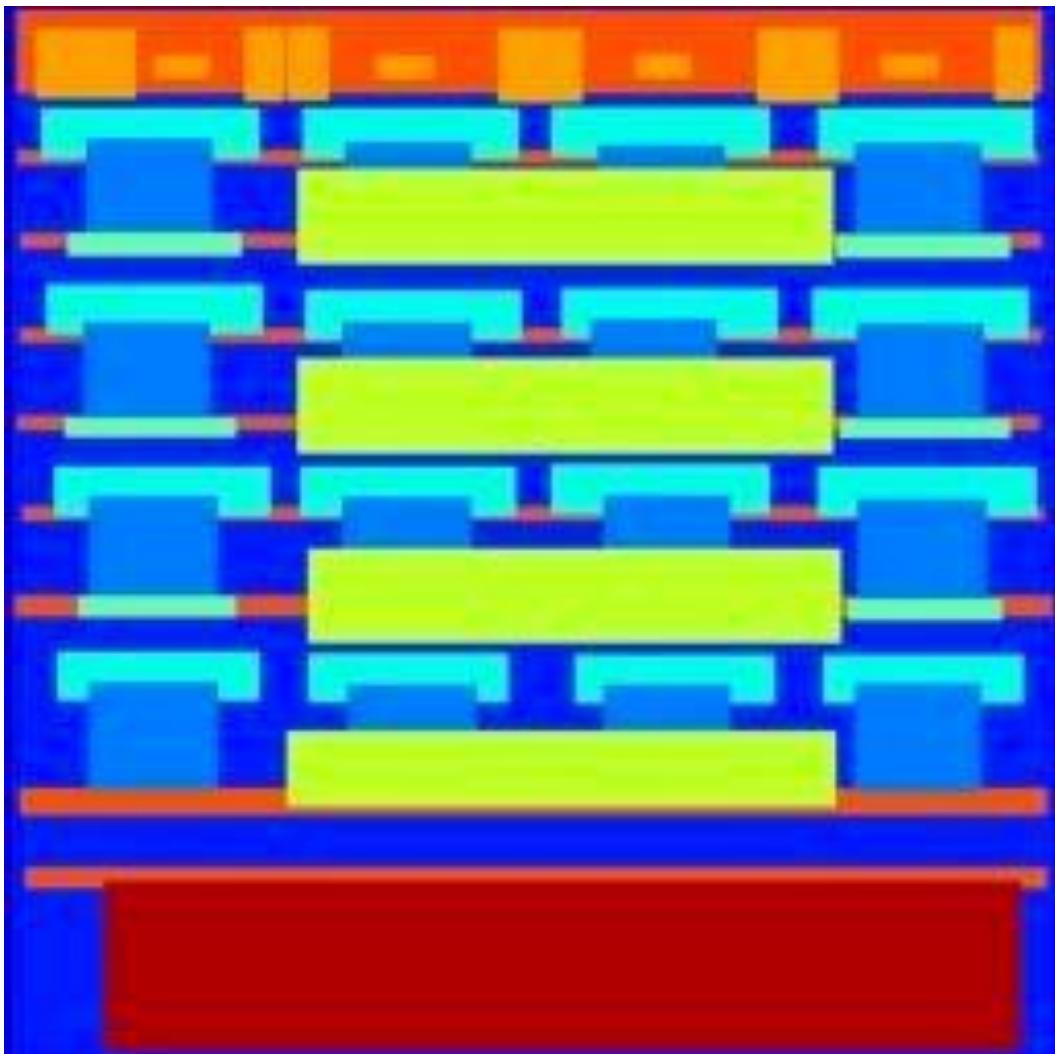
Rather than penalizing if output image looks fake, penalize if each overlapping patch in output looks fake

- Faster, fewer parameters
- More supervised observations
- Applies to arbitrarily large images

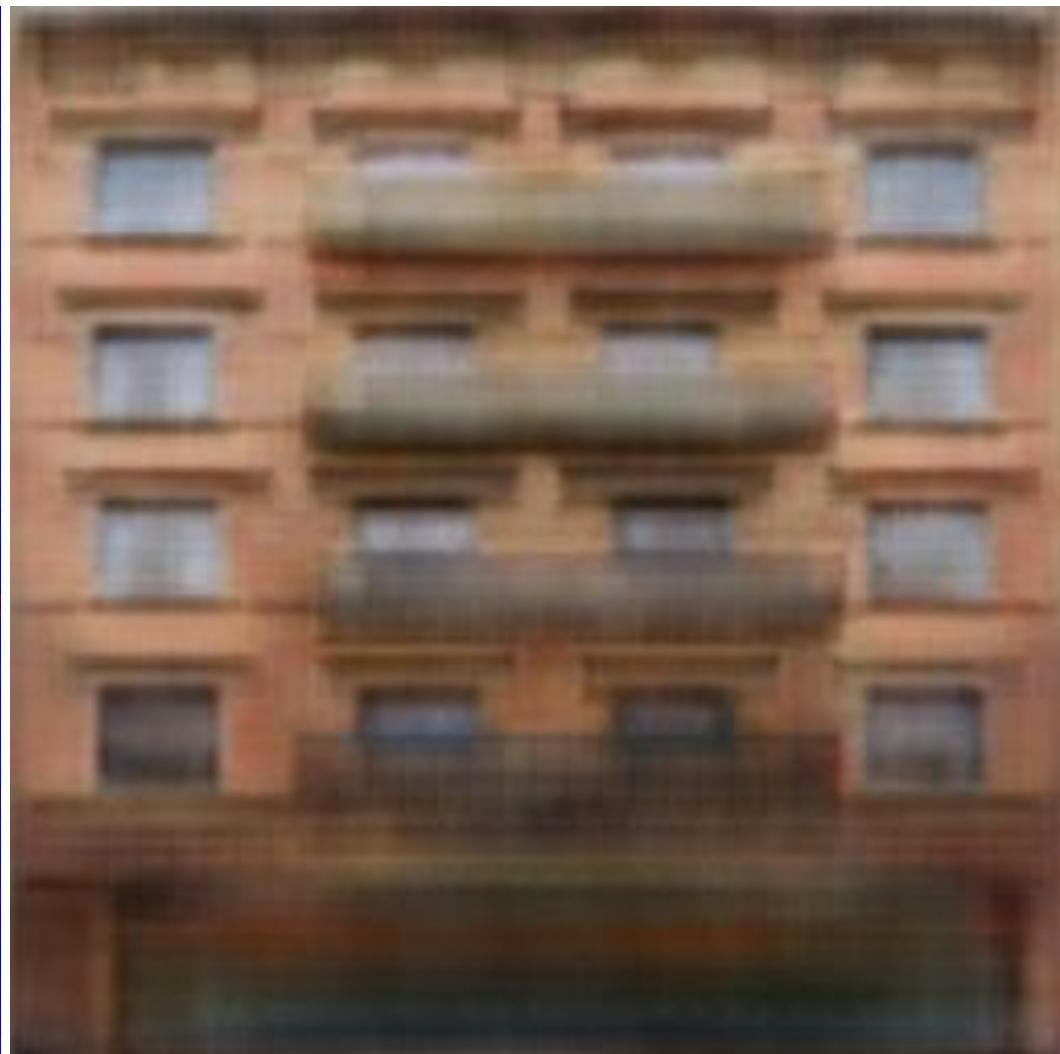
[Li & Wand 2016]  
[Shrivastava et al. 2017]  
[Isola et al. 2017]

# Labels → Facades

Input



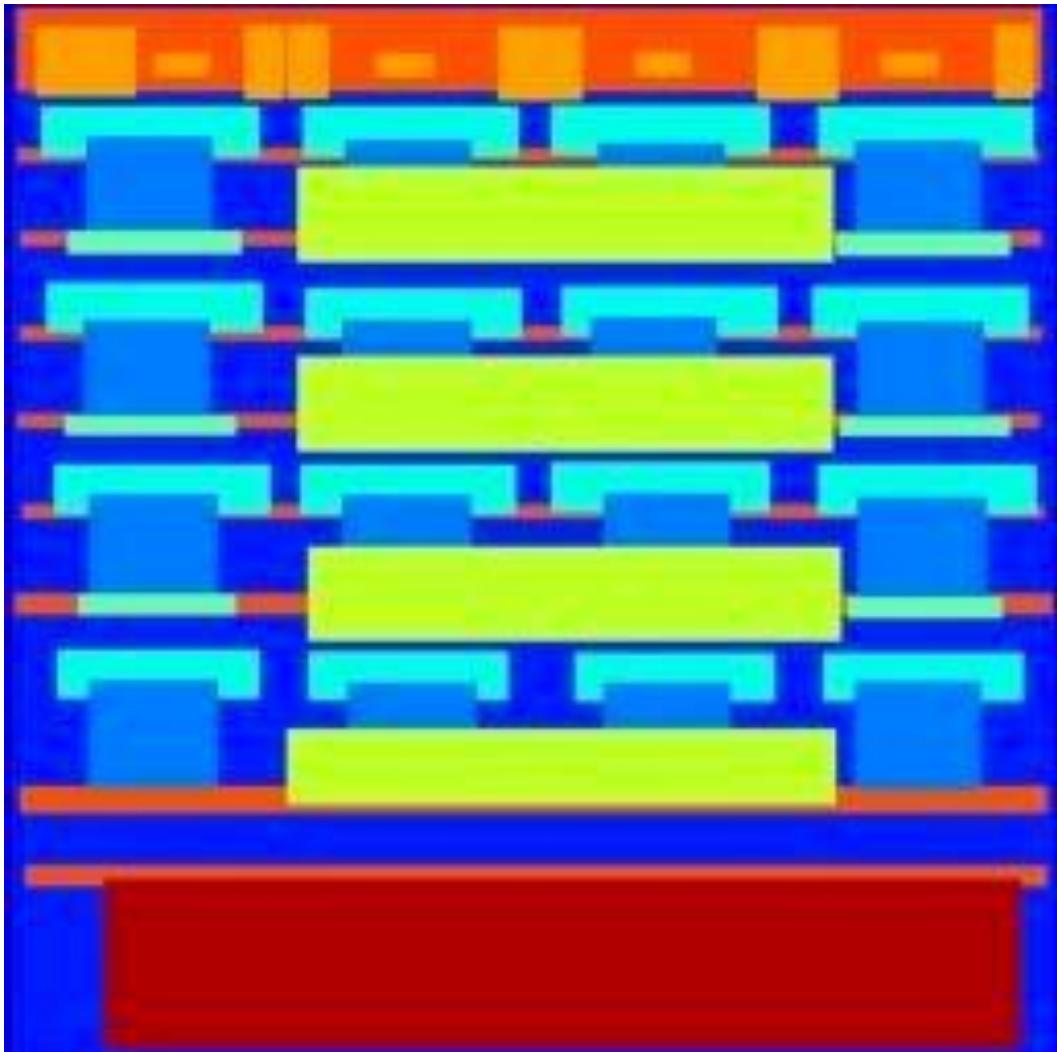
1x1 Discriminator



Data from [Tylecek, 2013]

# Labels → Facades

Input



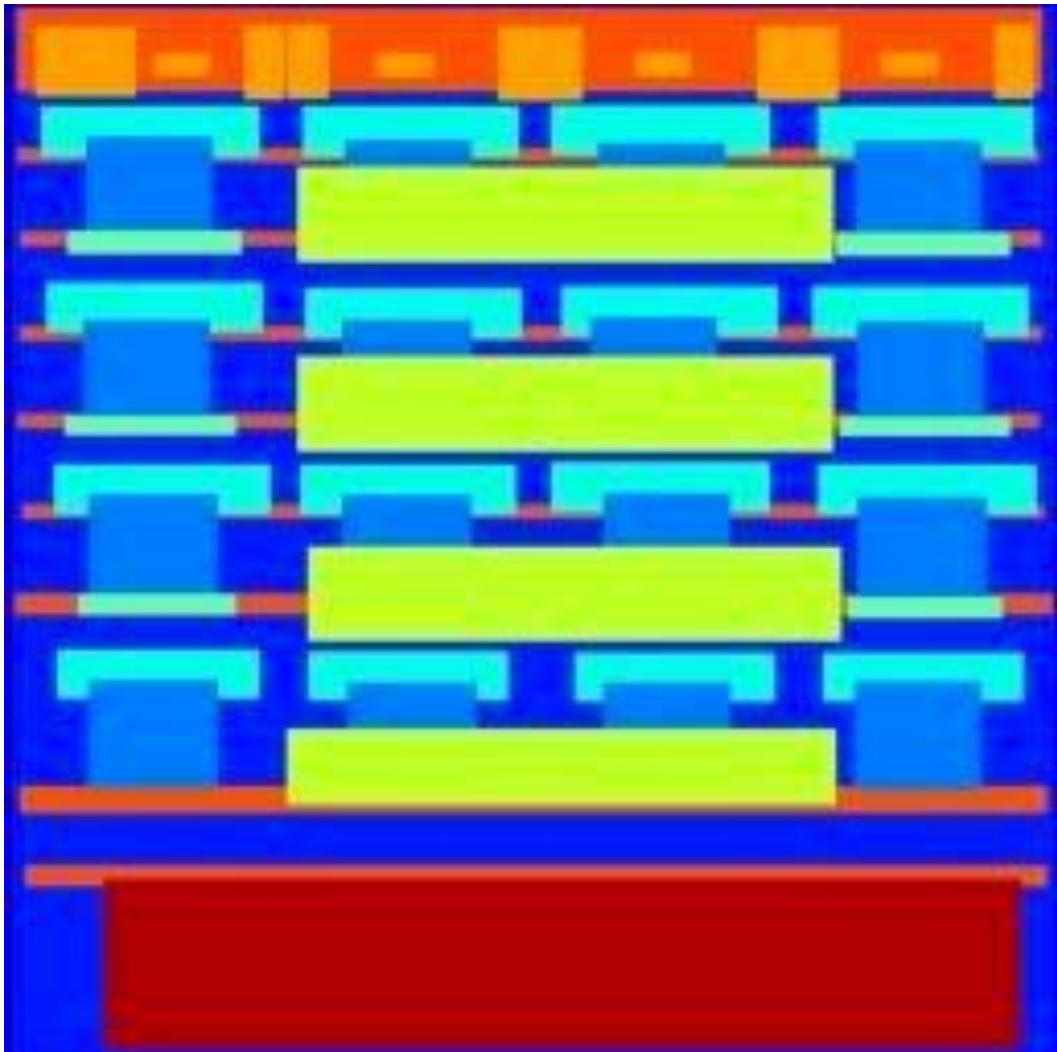
16x16 Discriminator



Data from [Tylecek, 2013]

# Labels → Facades

Input



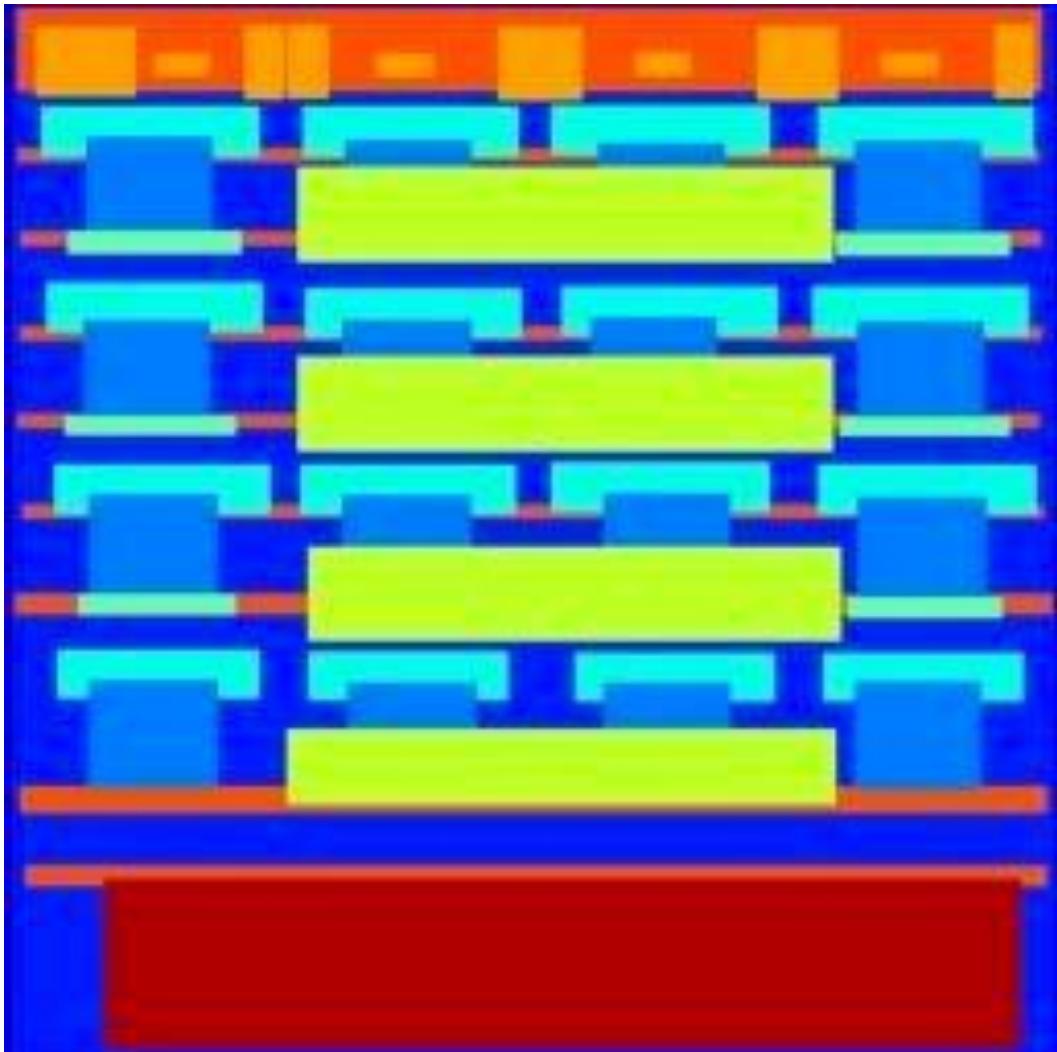
70x70 Discriminator



Data from [Tylecek, 2013]

# Labels → Facades

Input

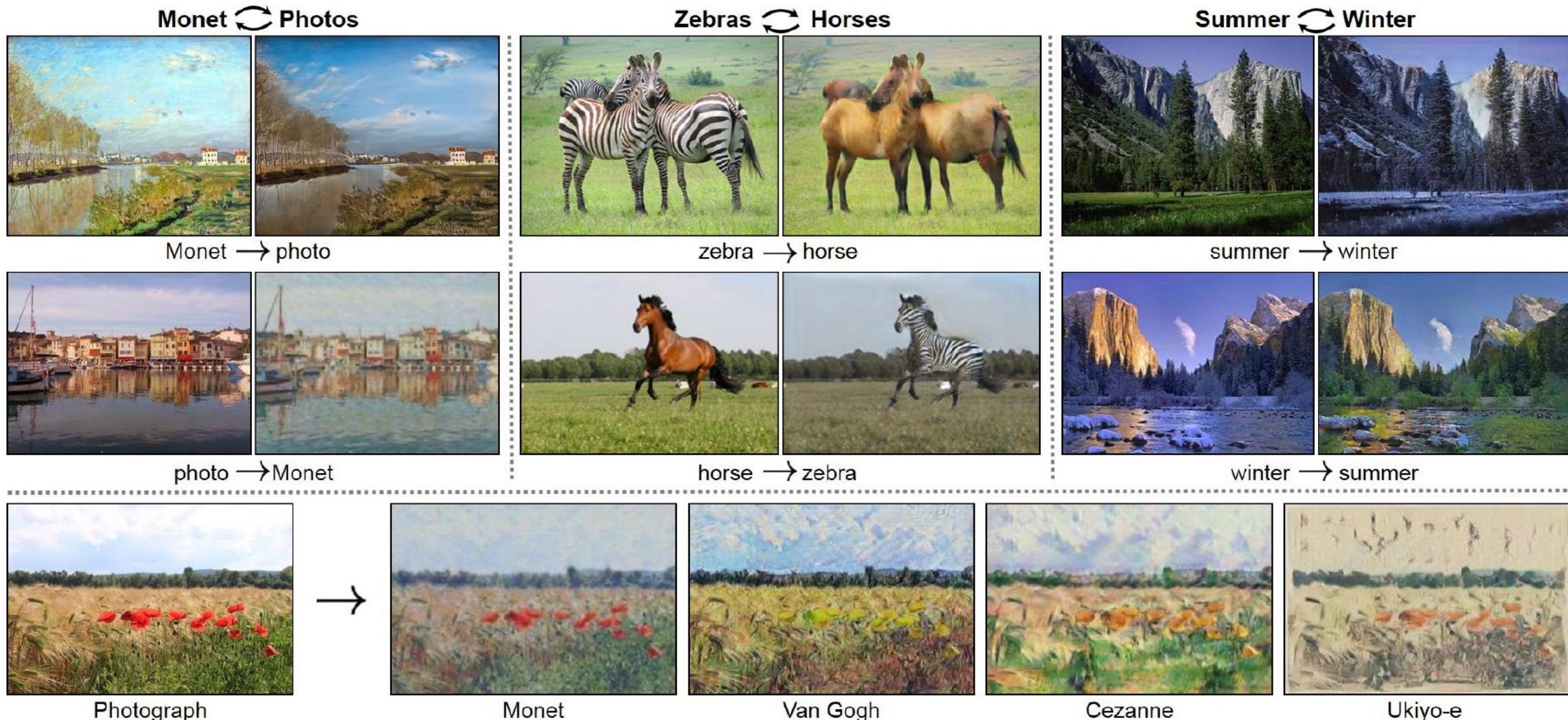


Full image Discriminator



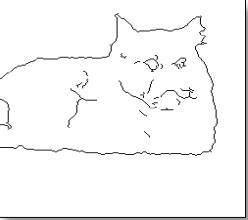
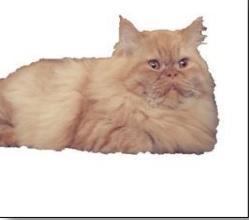
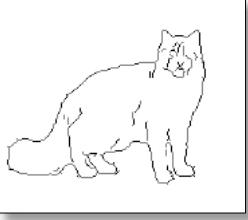
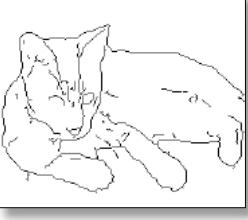
Data from [Tylecek, 2013]

# Pix2Pix w/o input-output pairs



(Zhu et al. 2017)

# Paired data

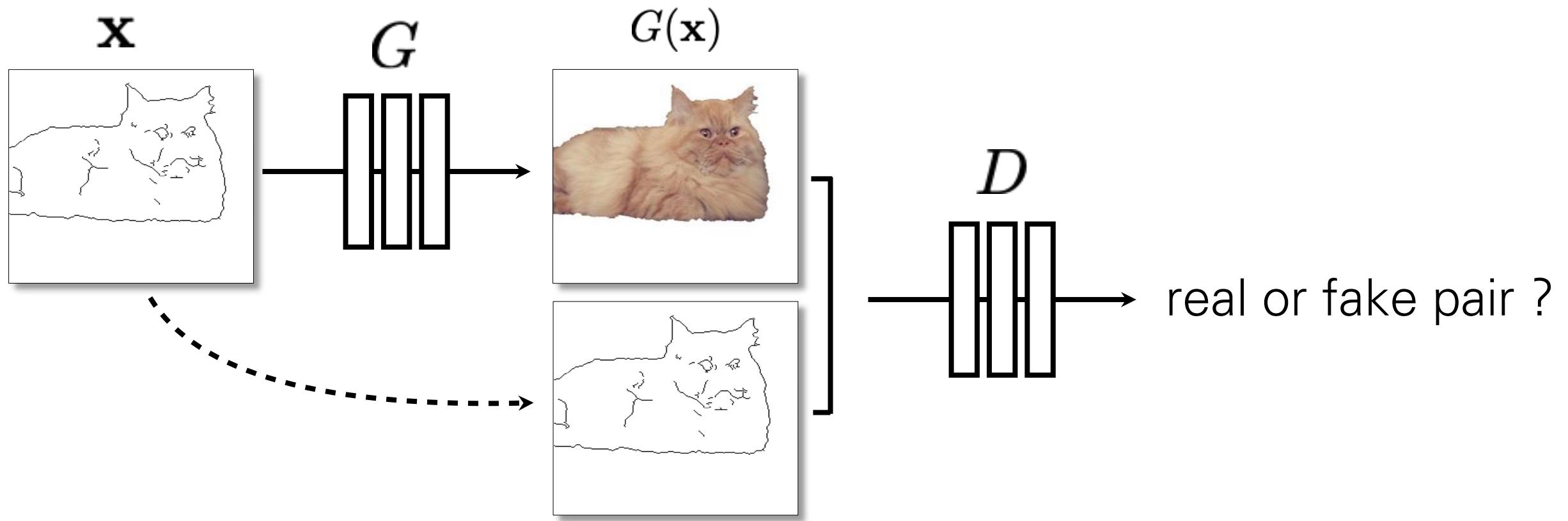
$x_i$	$y_i$
{  , }	{  }
{  , }	{  }
{  , }	{  }
⋮	

## Paired data

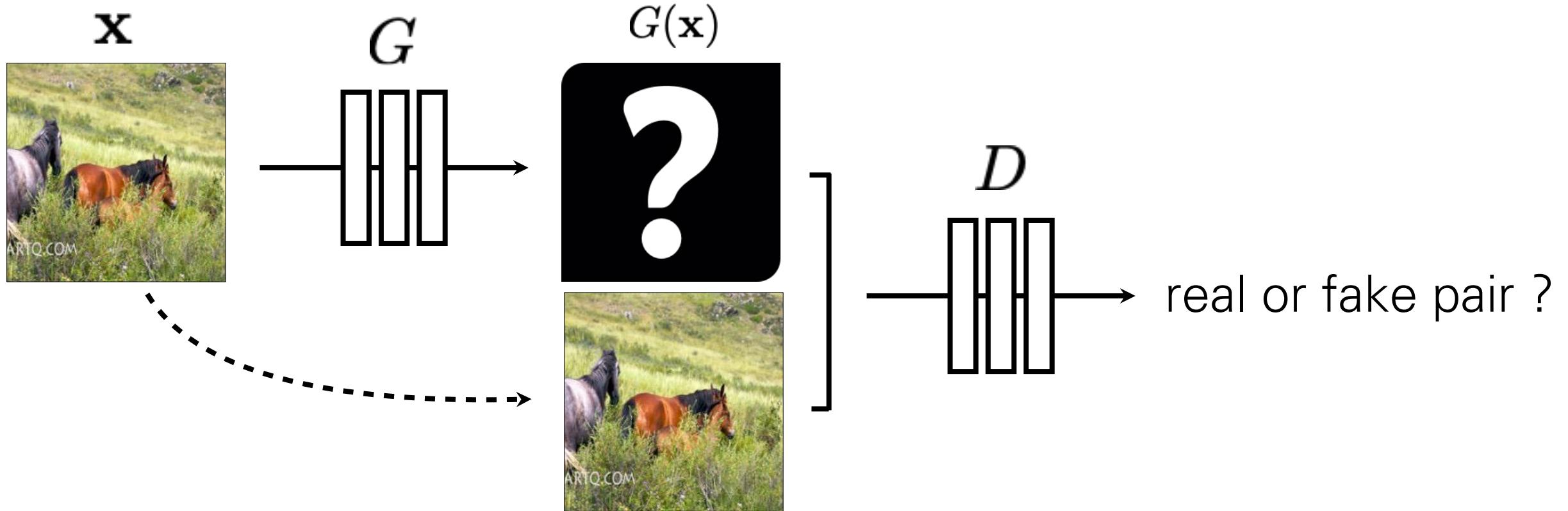
$x_i$	$y_i$	
		}
		}
		}
⋮		

## Unpaired data

$X$	$Y$
⋮	⋮

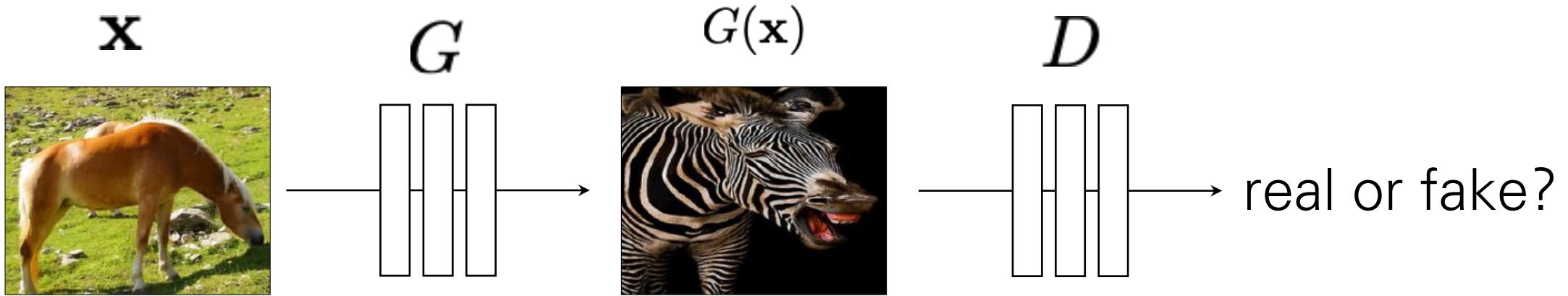


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$



$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$

No input-output pairs!

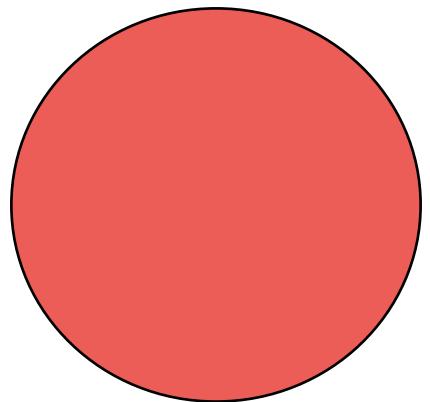


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

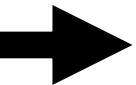
Usually loss functions check if output matches a target instance

GAN loss checks if output is part of an admissible set

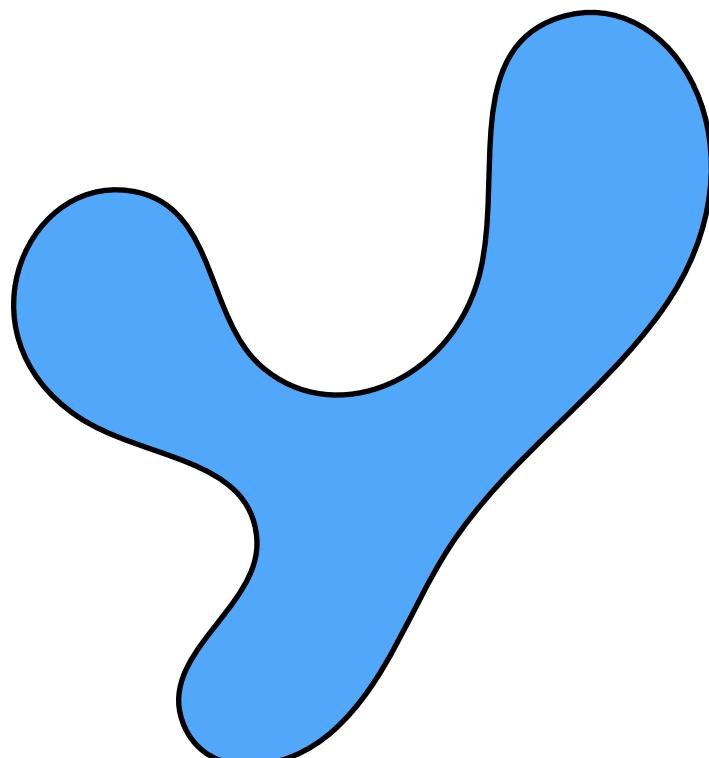
Gaussian



**z**

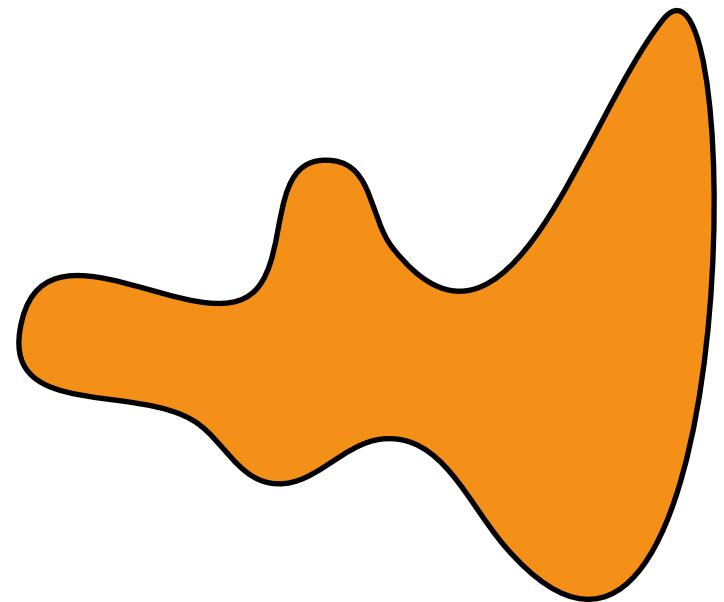


Target distribution

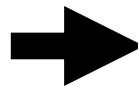


**Y**

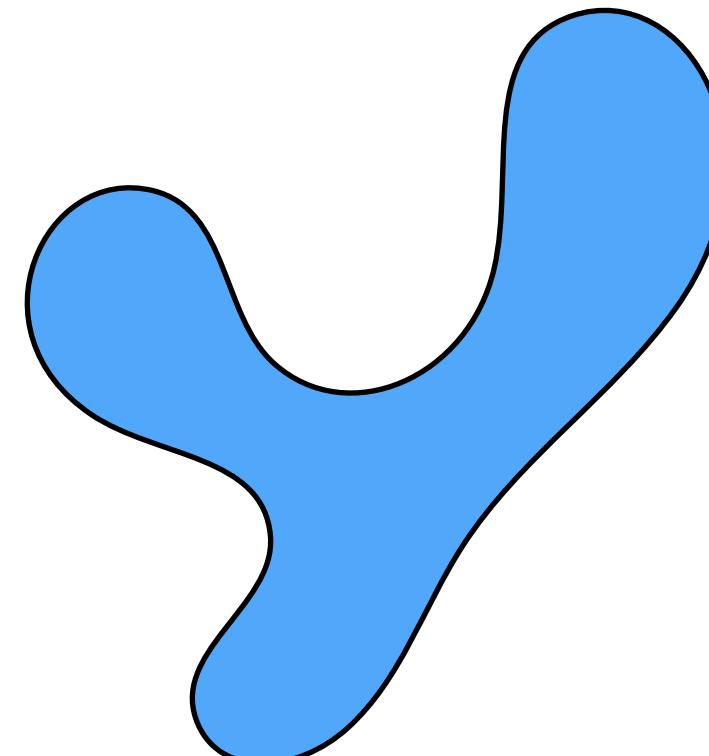
Horses



**X**



Zebras

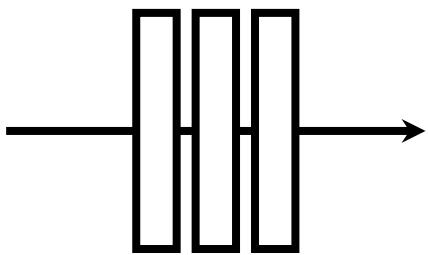


**Y**

**x**



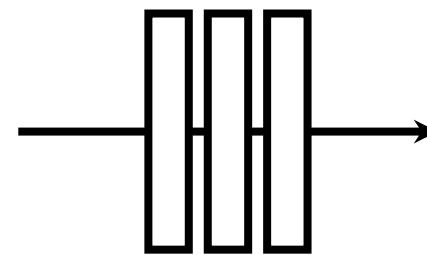
**G**



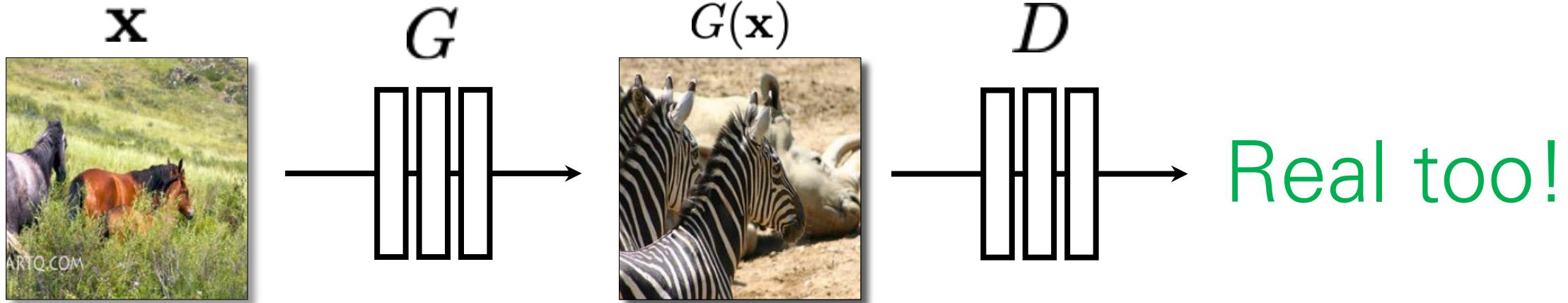
**G(x)**



**D**

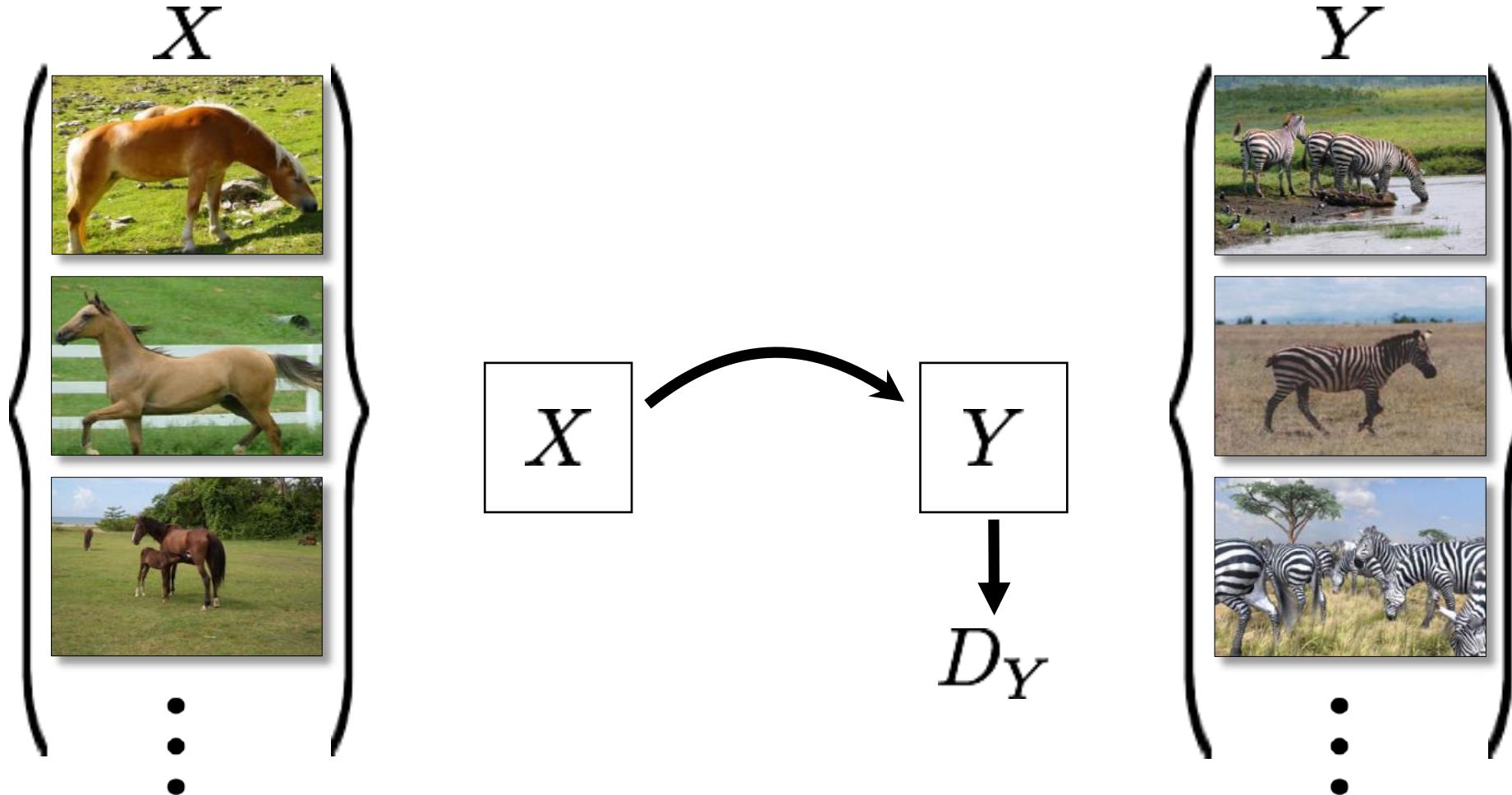


**Real!**



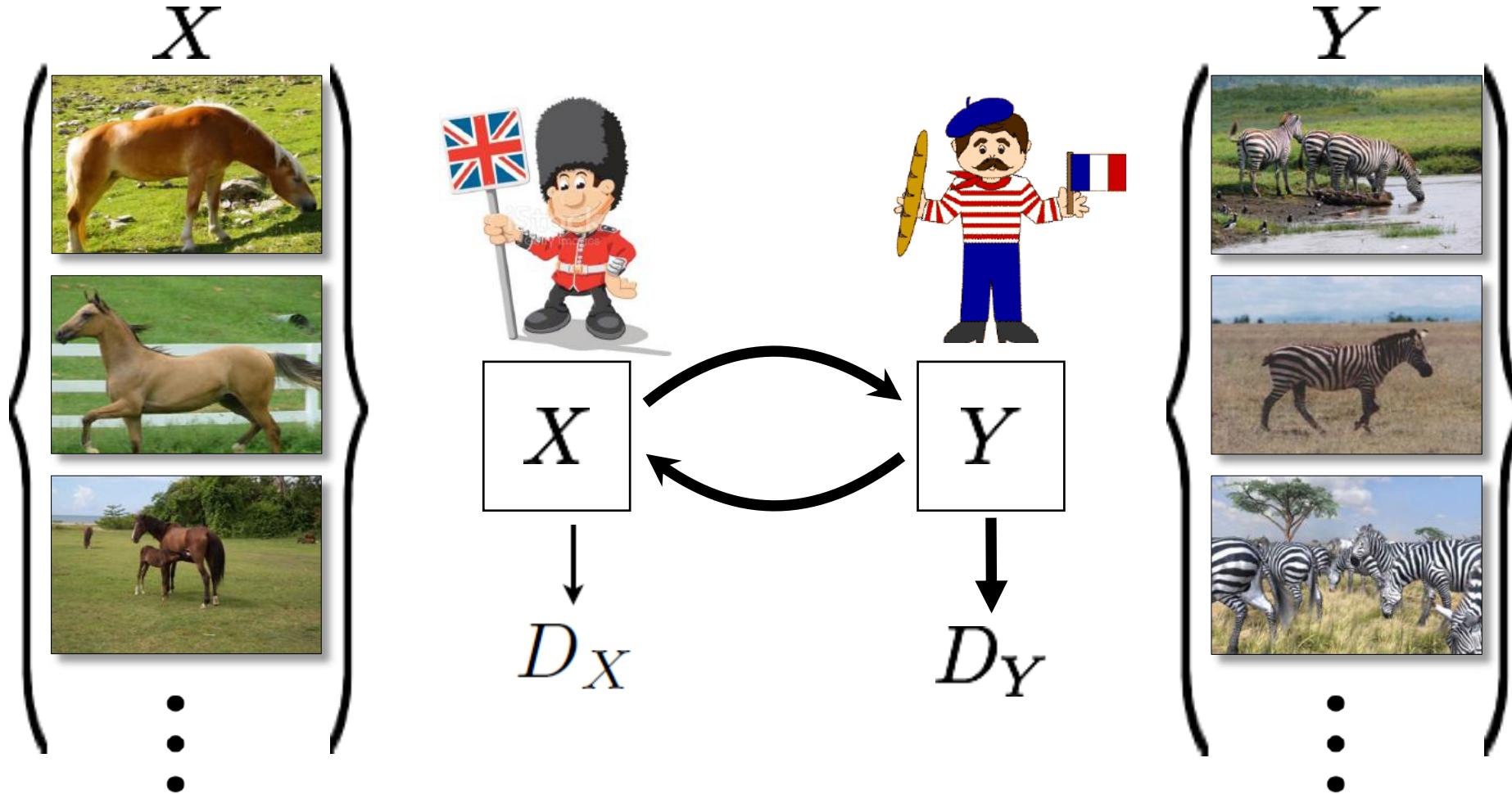
Nothing to force output to correspond to input

# Cycle-Consistent Adversarial Networks

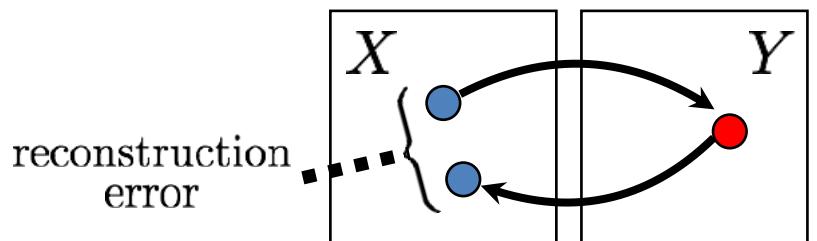
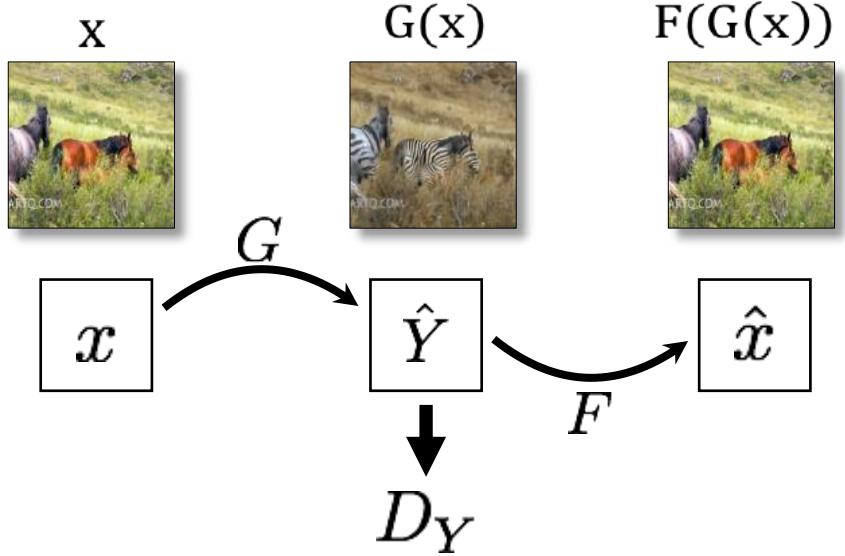


[Zhu et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

# Cycle-Consistent Adversarial Networks

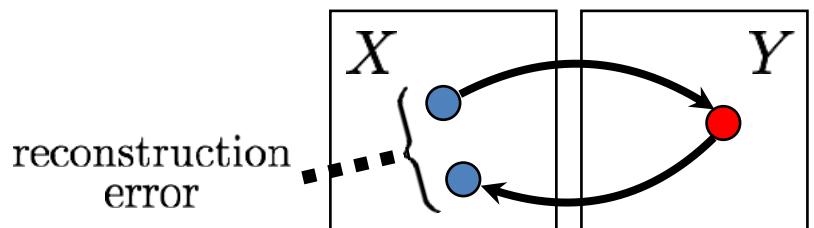
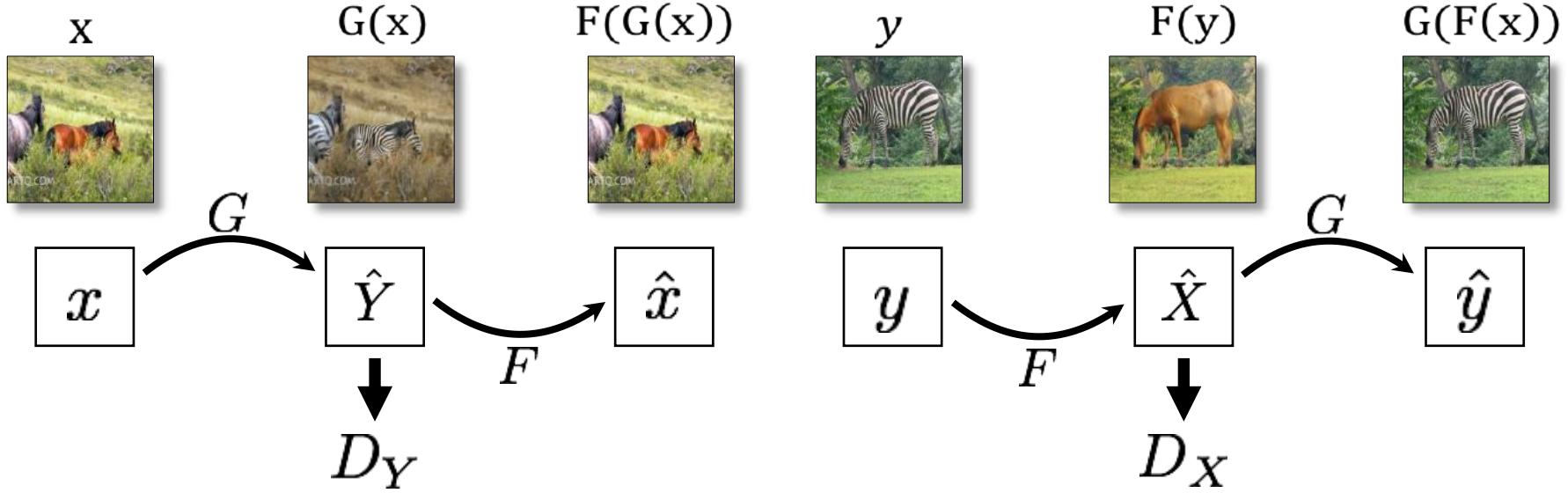


# Cycle Consistency Loss

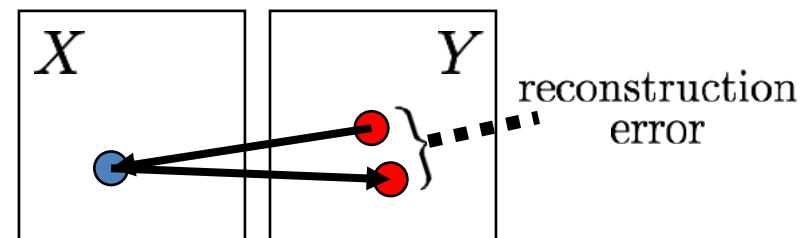


$$\|F(G(x)) - x\|_1$$

# Cycle Consistency Loss



$$\|F(G(x)) - x\|_1$$



$$\|G(F(y)) - y\|_1$$





# Collection Style Transfer



Photograph  
@ Alexei Efros



Monet



Van Gogh



Cezanne



Ukiyo-e

Input



Monet



Van Gogh



Cezanne



Ukiyo-e



# Monet's paintings → photos



# Monet's paintings → photos



# Failure case

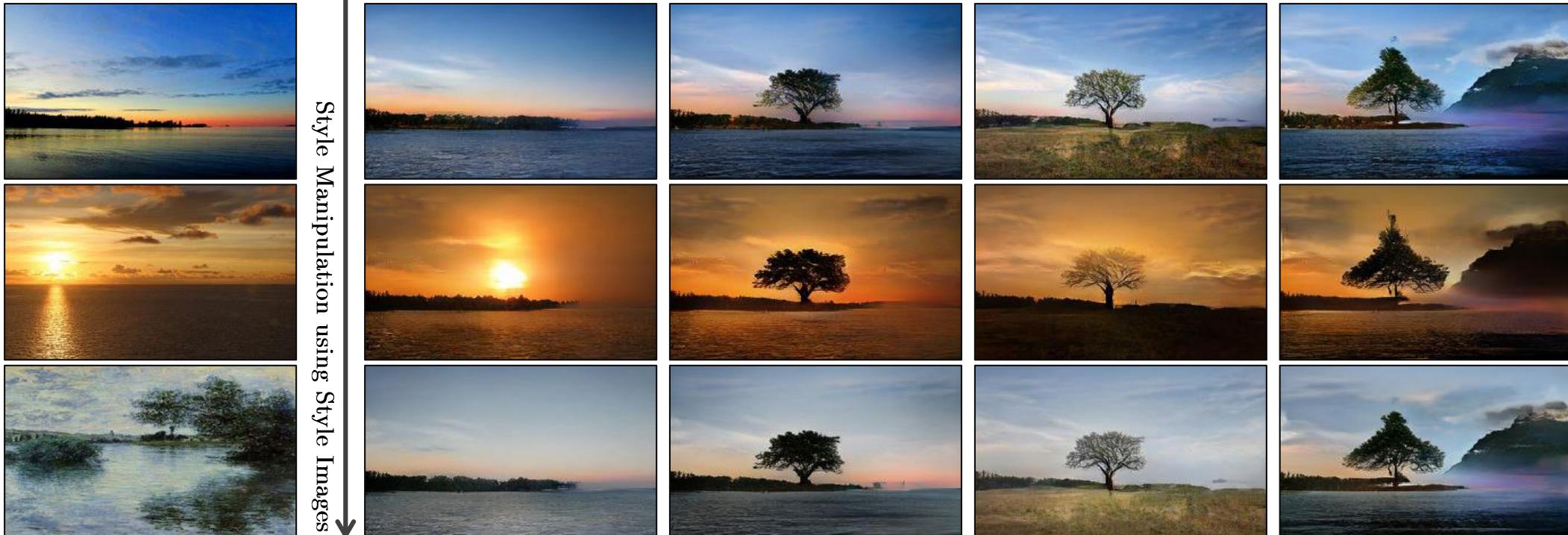


# Failure case



# Semantic Image Synthesis (SPADE) (Park et al., 2019)

- Image generation conditioned on semantic layouts





MIT CSAIL  
@MIT\_CSAIL

Neural networks are now "hallucinating." This framework lets you change an image to appear to be in a different season, weather condition or time of day: [bit.ly/2PeVcVI](https://bit.ly/2PeVcVI) v/@Hacettepe1967 & @UvA\_Amsterdam



## Manipulating Attributes of Natural Scenes via Hallucination.

Levent Karacan, Zeynep Akata, Aykut Erdem & Erkut Erdem.

ACM Trans. on Graphics, Vol. 39, Issue 1, Article 7, February 2020.



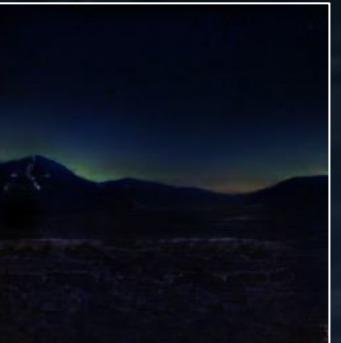


snow



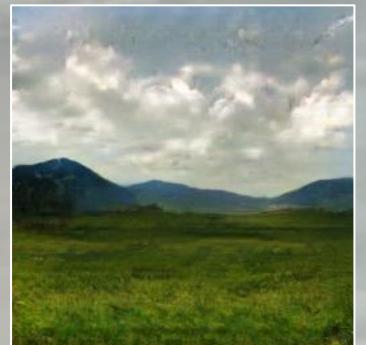


night

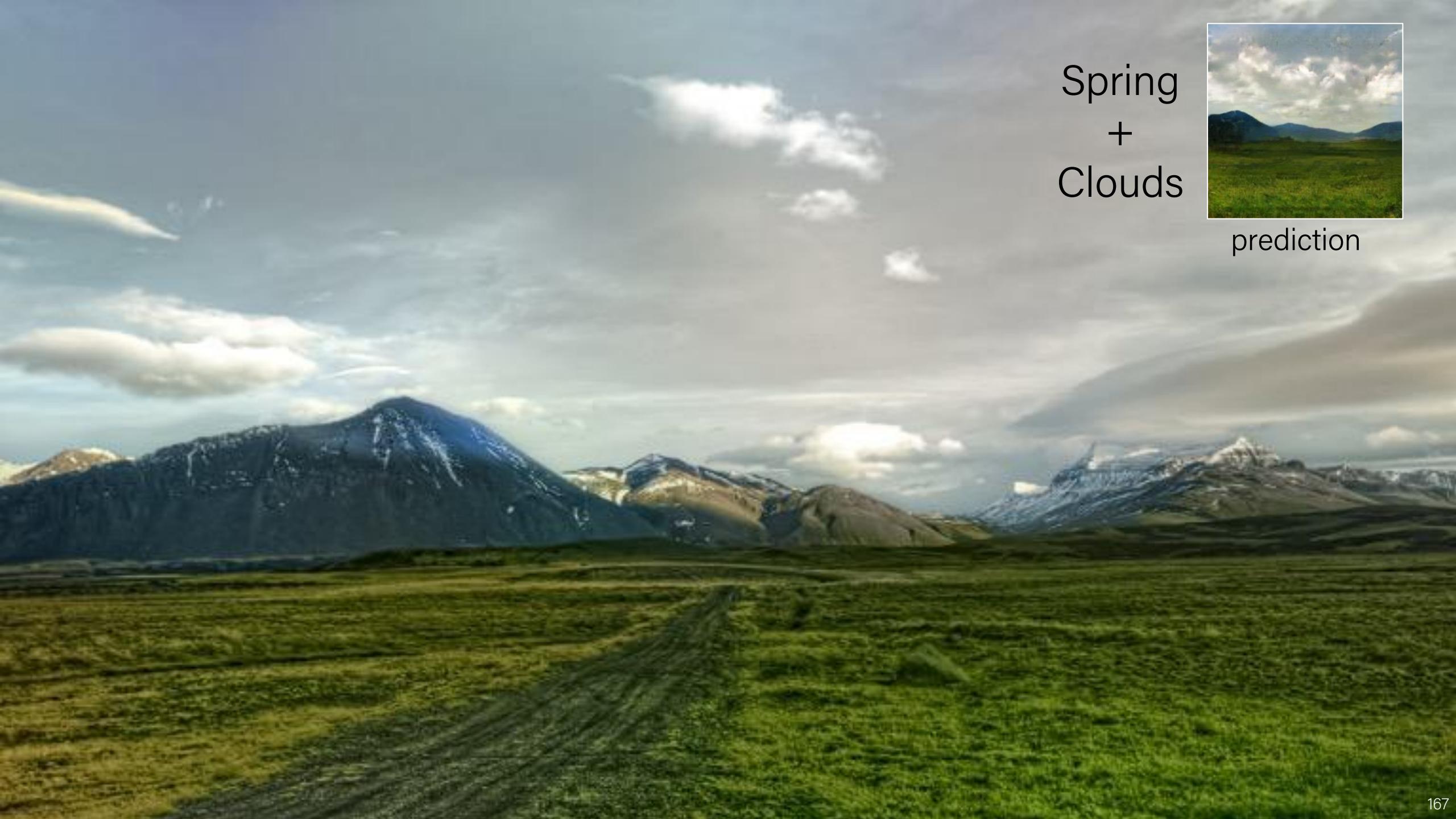


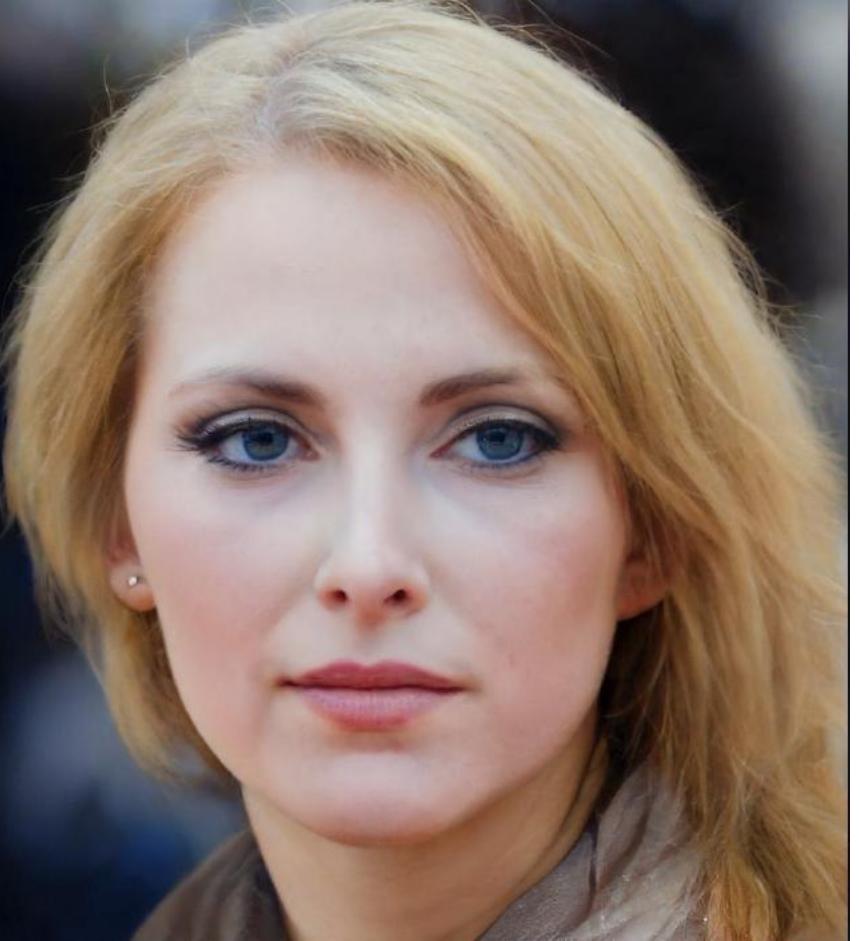
prediction

Spring  
+  
Clouds

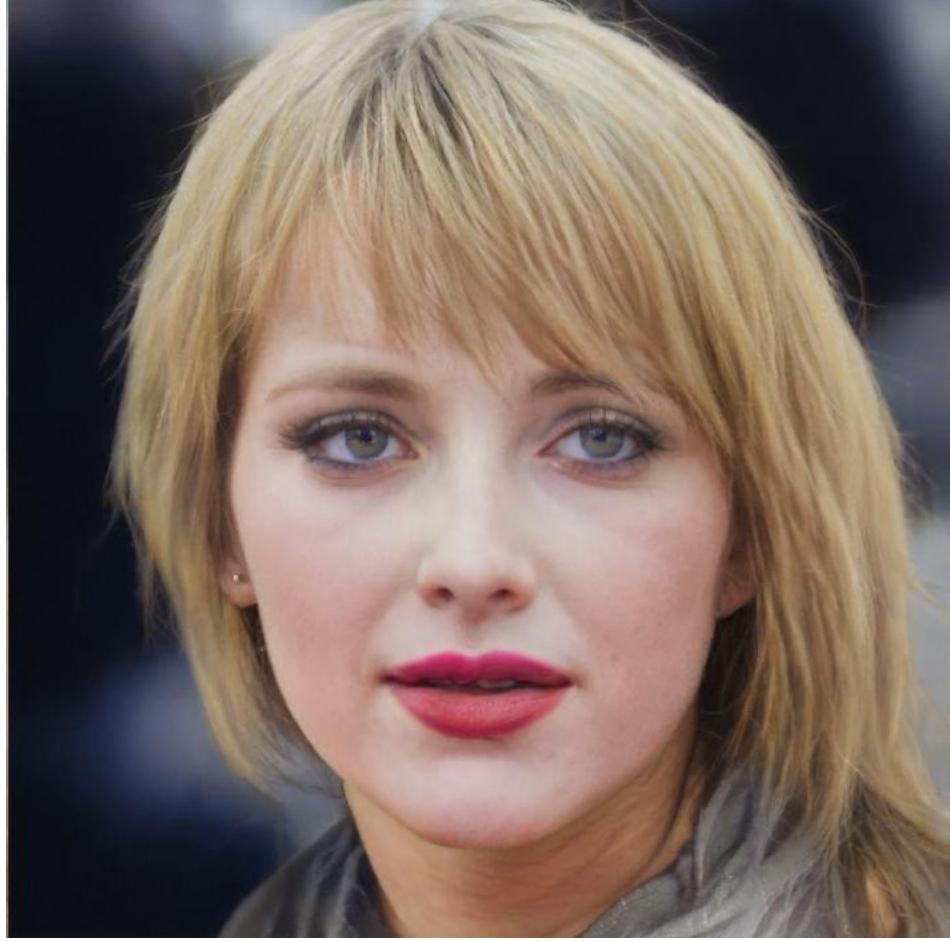


prediction





A young woman  
with bangs  
wearing lipstick

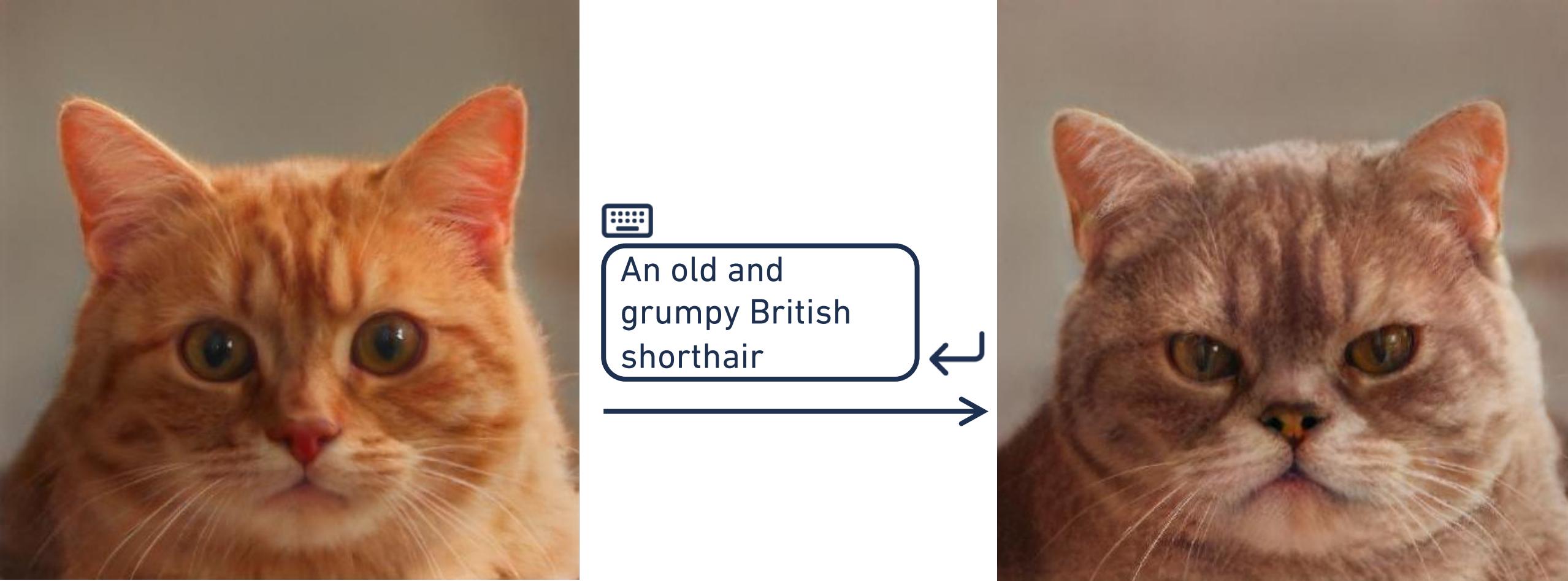


Adobe Research

## CLIP-Guided StyleGAN Inversion for Text-Driven Real Image Editing.

Canberk Baykal, Abdul Basit Anees, Duygu Ceylan,  
Aykut Erdem, Erkut Erdem, Deniz Yuret  
ACM Transactions on Graphics, 2023





Adobe Research

## CLIP-Guided StyleGAN Inversion for Text-Driven Real Image Editing.

Canberk Baykal, Abdul Basit Anees, Duygu Ceylan,  
Aykut Erdem, Erkut Erdem, Deniz Yuret  
ACM Transactions on Graphics, 2023





green jacket

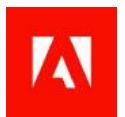
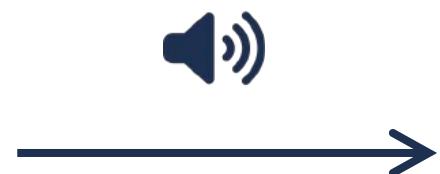
Sleeveless blue blouse

black short



## VidStyleODE: Disentangled Video Editing via StyleGAN and NeuralODE.

Moayed Haji Ali, Andrew Bond, Tolga Birdal, Duygu Ceylan, Levent Karacan, Erkut Erdem,  
Aykut Erdem. ICCV 2023

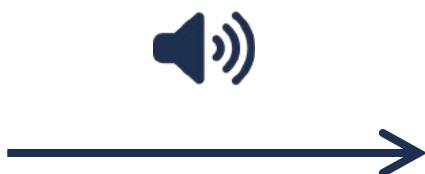


Adobe Research

## Audio-based Image Editing and Generation using Latent Diffusion Models

Burak Can Biner, Farrin Marouf Sofian,  
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,  
Aykut Erdem. In progress



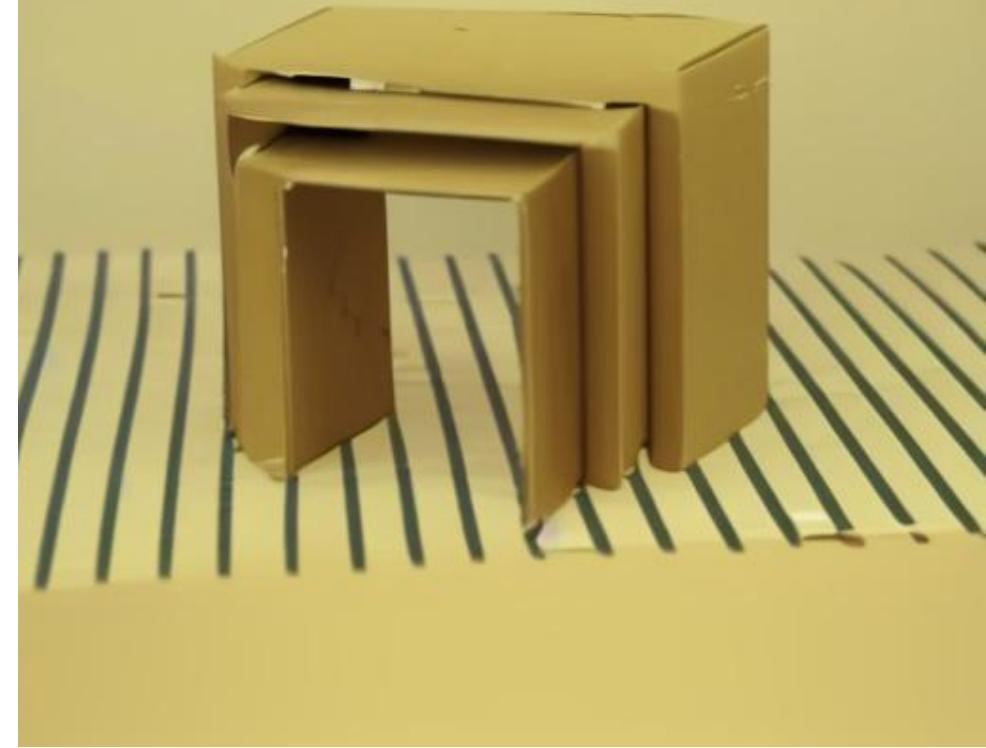
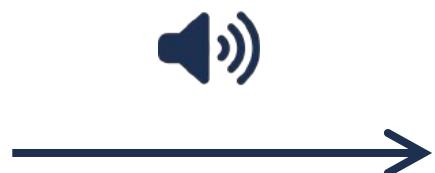


Adobe Research

## Audio-based Image Editing and Generation using Latent Diffusion Models

Burak Can Biner, Farrin Marouf Sofian,  
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,  
Aykut Erdem. In progress





Adobe Research

## Audio-based Image Editing and Generation using Latent Diffusion Models

Burak Can Biner, Farrin Marouf Sofian,  
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,  
Aykut Erdem. In progress



**Next lecture:**  
**Autoregressive and Flow Models**