# Image Manifolds & Image Synthesis

(including GANs)

By Abe Davis
Some slides from Jin Sun, Phillip Isola

# Announcements

- Take home final May 11-14
- Sample final is online (check Piazza)
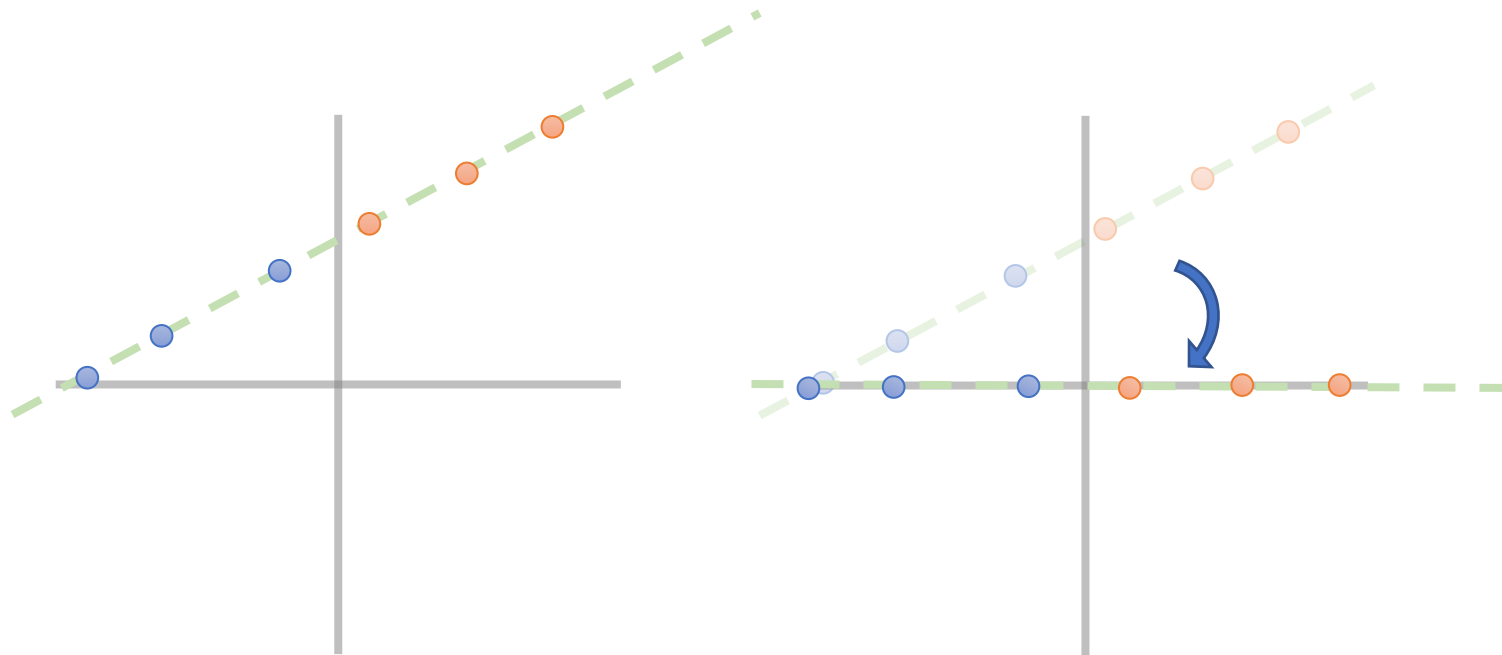- Project 5 deadline extended to Friday May 1

- Course evaluations are open now through May 8
  - We encourage feedback
  - Small amount of extra credit for filling out
    - What you write is still anonymous, instructors only see whether students filled it out
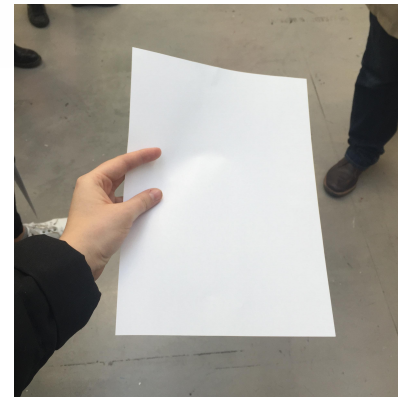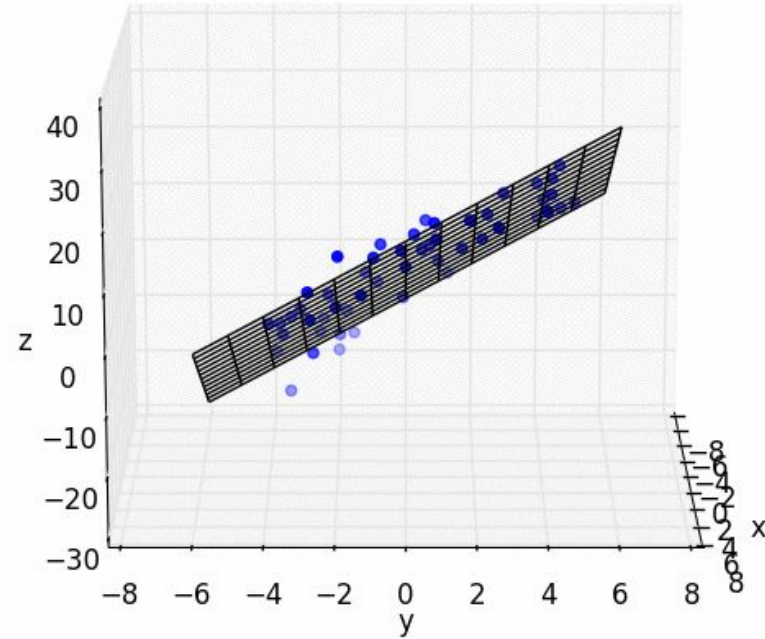
# Dimensionality Reduction

By Abe Davis

# Linear Dimensionality Reduction: 2D->1D

- Consider a bunch of data points in 2D
- Let's say these points only differ along one line
- If so, we can translate and rotate our data so that it is 1D

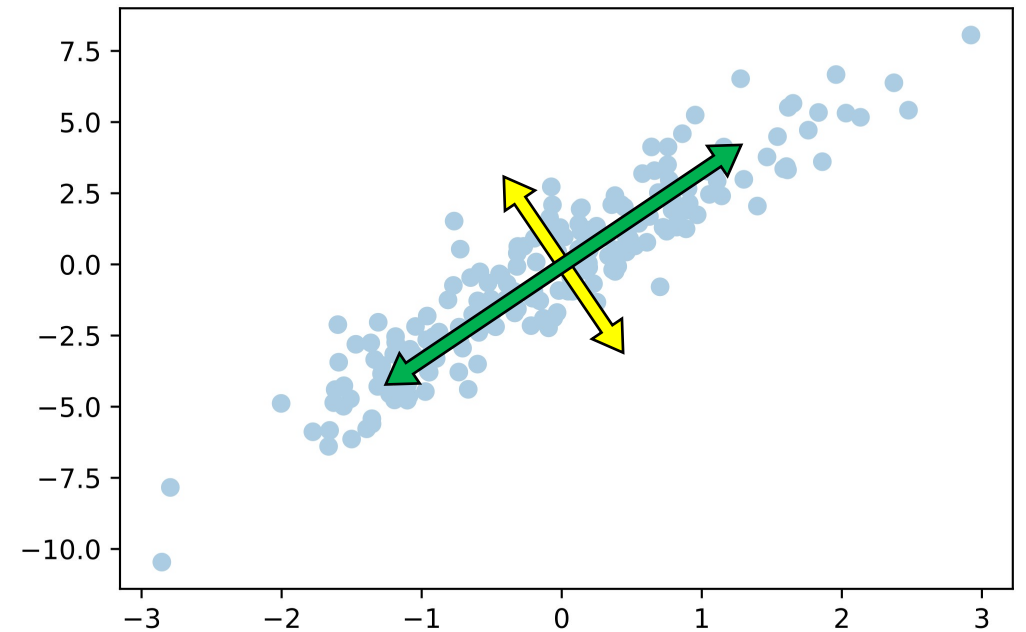# Linear Dimensionality Reduction: 3D->2D

- Similar to 1D case, we can fit a plane to the data, and transform our coordinate system so that plane becomes the x-y plane

- "Plane fitting"

- More generally: look for the 2D subspace that best fits the data, and ignore the remaining dimensions





Think of this as data that sits on a flat sheet of paper, suspended in 3D space. We will come back to this analogy in a couple slides…

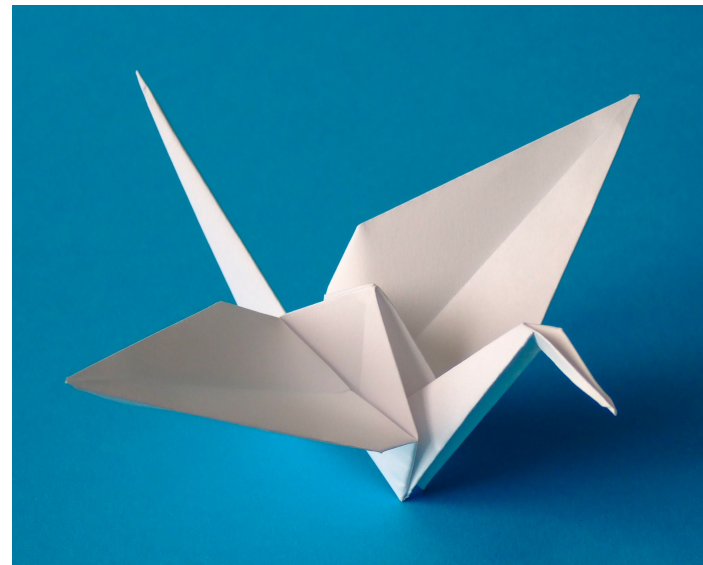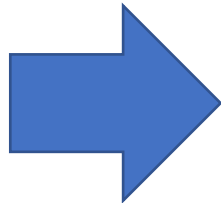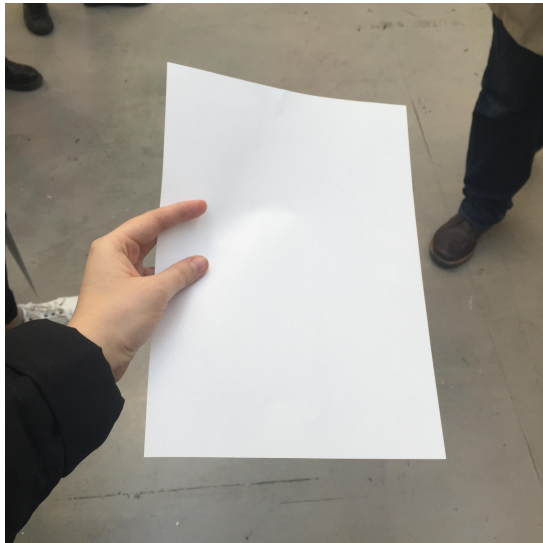# Generalizing Linear Dimensionality Reduction

- ***Principle Component Analysis (PCA)***: find and order orthogonal axes by how much the data varies along each axis.

- The axes we find (ordered by variance of our data) are called ***principle components***.

- Dimensionality reduction can be done by using only the first $k$ principle components



Side Note: principle components are closely related to the eigenvectors of the covariance matrix for our data
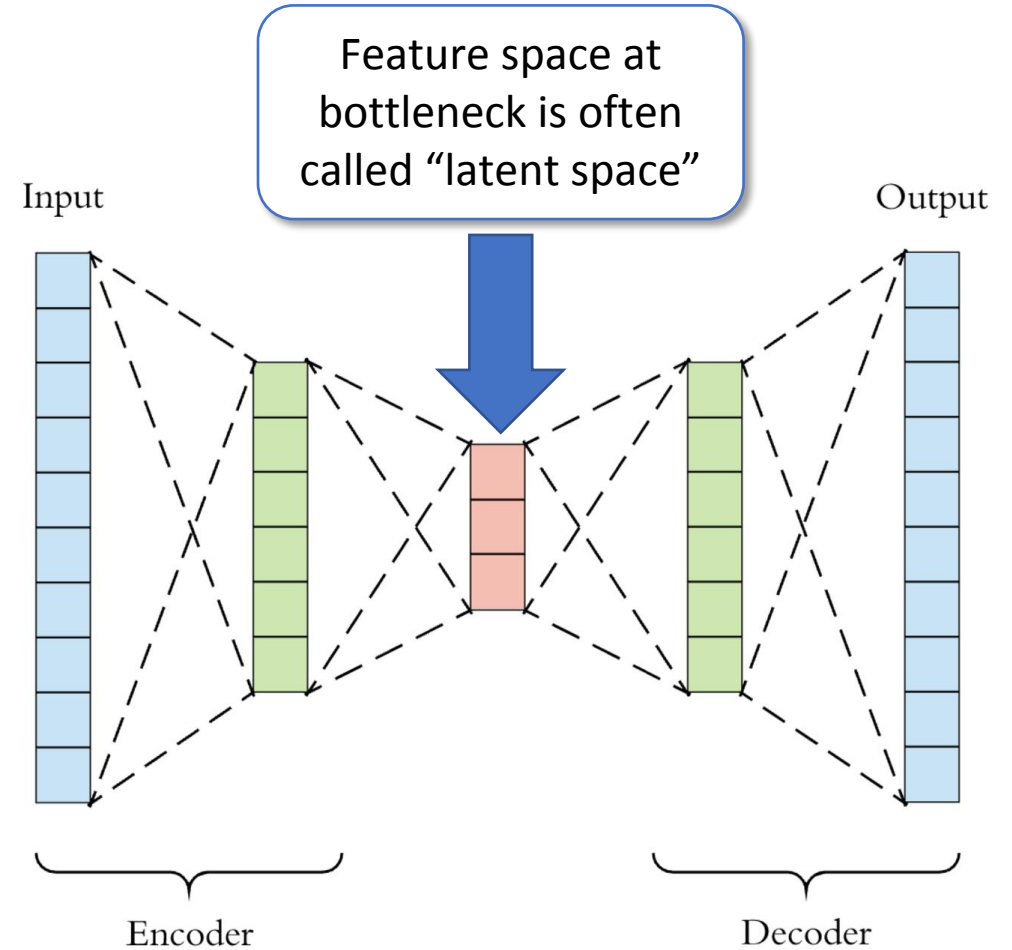
# Manifolds

- Think of a piece of paper as a 2D subspace
- If we bend and fold that paper, it's still locally a 2D subspace…
- A "manifold" is the generalization of this concept to higher dimensions…

# Autoencoders: Dimensionality Reduction for Manifolds

- Learn a non-linear transformation into some lower-dimensional space (encoder)

- Learn a transformation from lower-dimensional space back to original content (decoder)

- Loss function measures the difference between input and output

- **Unsupervised**
  - No labels required!

Feature space at bottleneck is often called "latent space"

Input

Output

Encoder

Decoder

# Autoencoders: Dimensionality Reduction for Manifolds



- Transformations that reduce dimensionality **cannot be invertible** in general

- An autoencoder tries to learn a transformation that is **invertible for points on some manifold**.
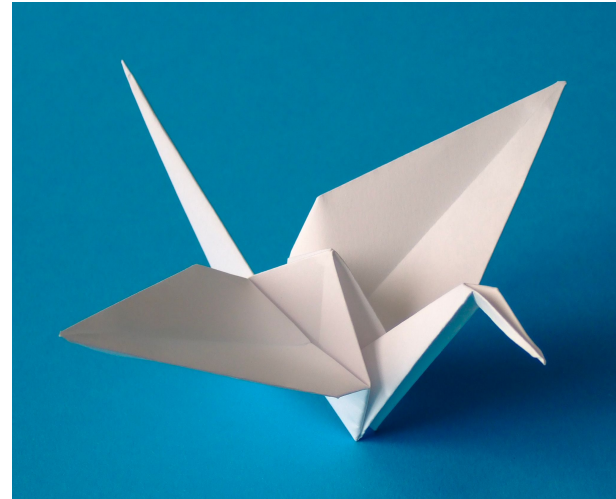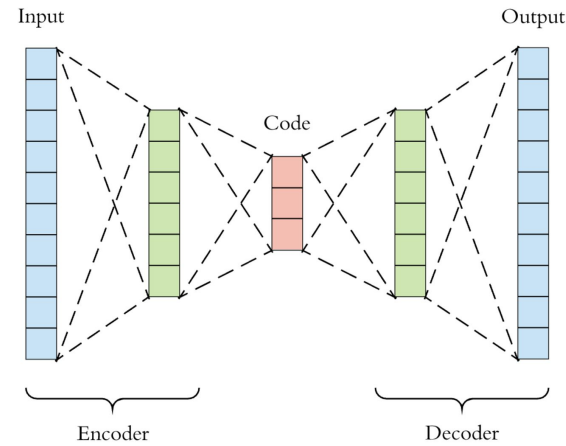
# Image Manifolds

By Abe Davis

# The Space of All Images

- Lets consider the space of all 100x100 images

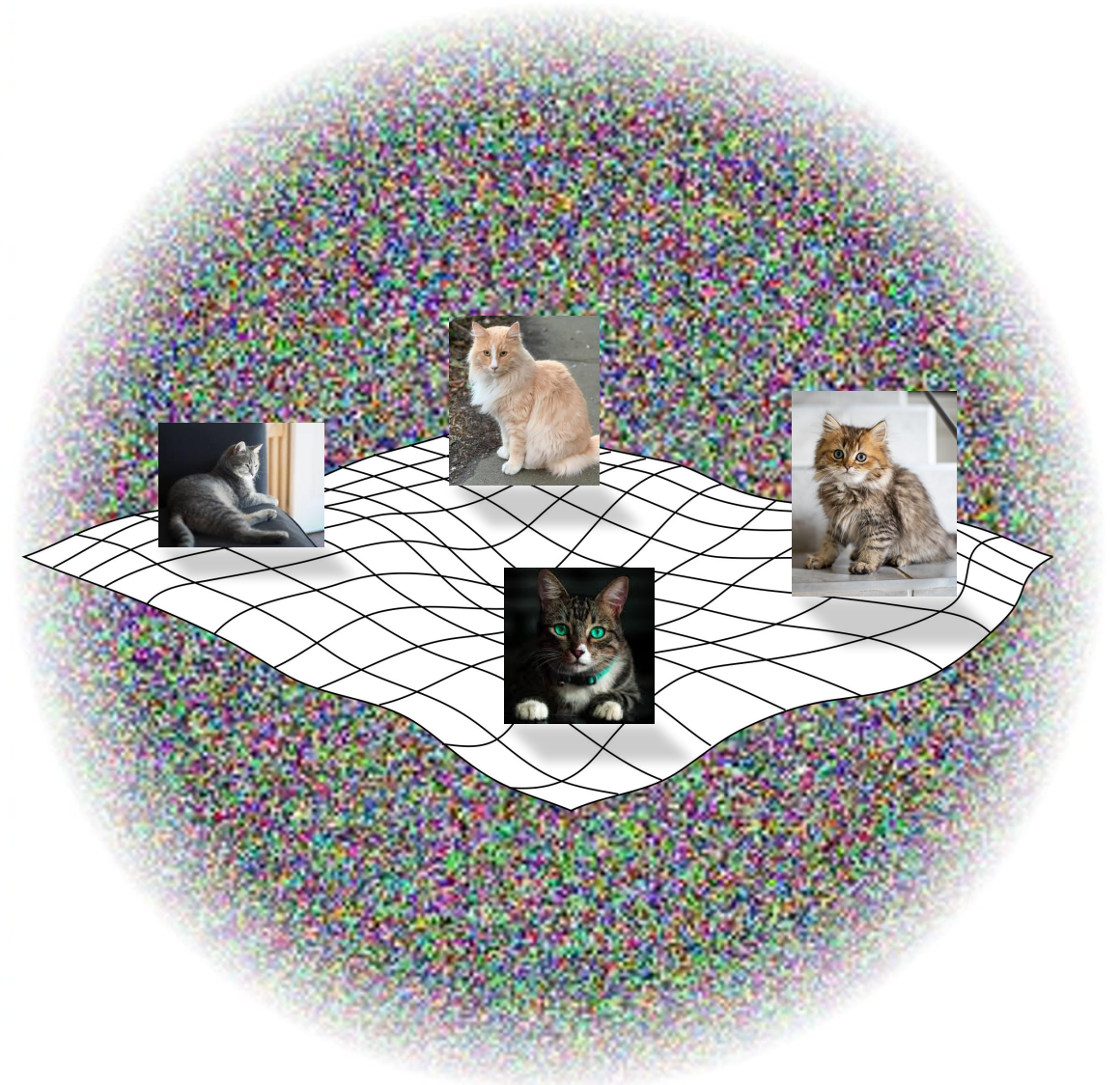- Now lets randomly sample that space…

- Conclusion: Most images are noise

**Question:**
What do we expect a random uniform sample of all images to look like?
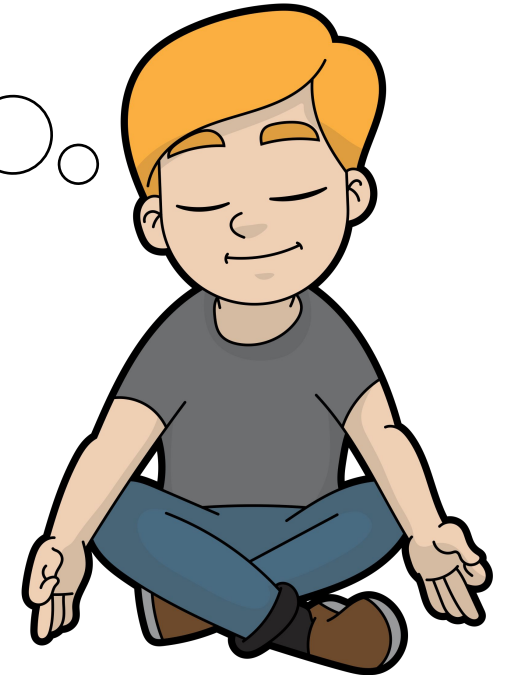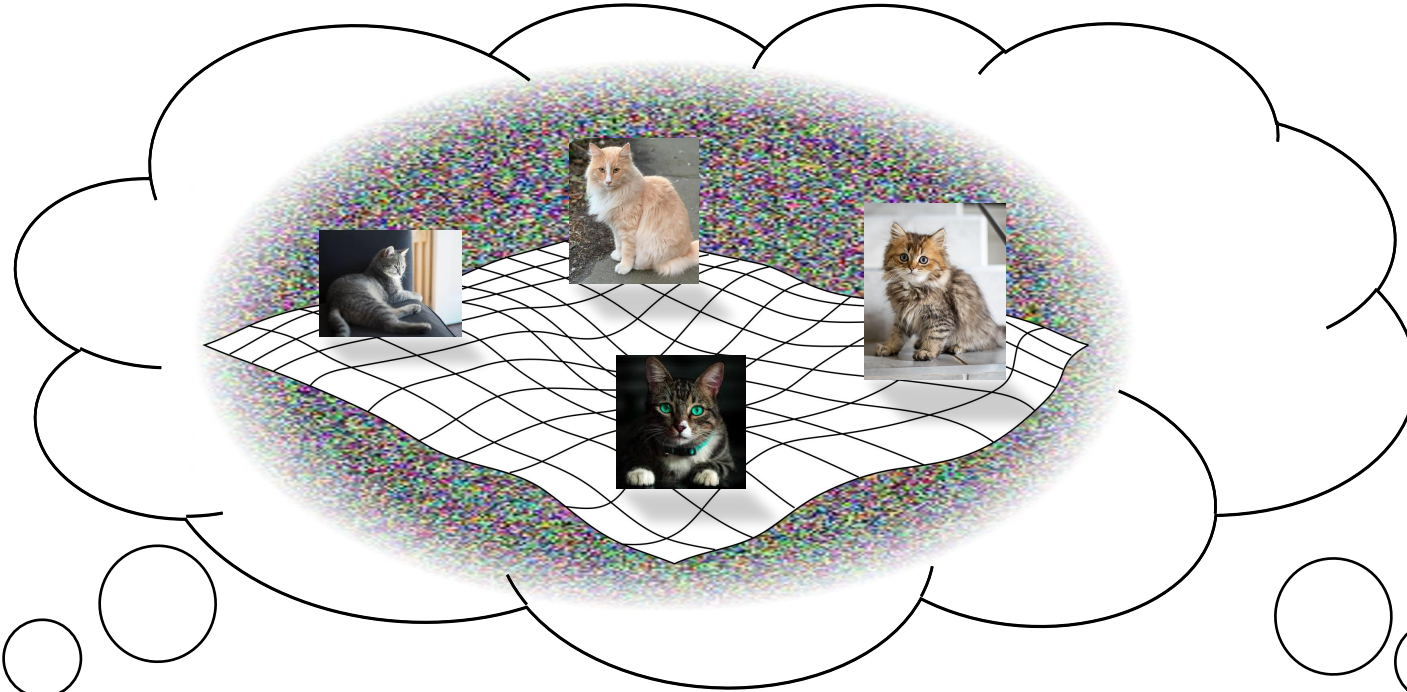
```
pixels = np.random.rand(100,100,3)
```

# Natural Image Manifolds

- Most images are "noise"

- "Meaningful" images tend to form some manifold within the space of all images

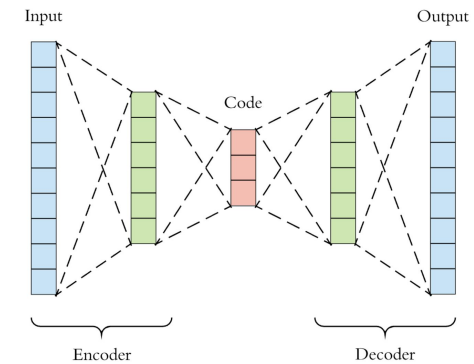- Images of a particular class fall on manifolds within that manifold...



The Space of All Images

# Natural Image Manifolds

# Denoising and the "Null Space" of Autoencoders

- The autoencoder tries to learn a dimensionality reduction that is invertible for our data (data on some manifold)

- Most noise will be in the non-invertible part of image space (off the manifold)

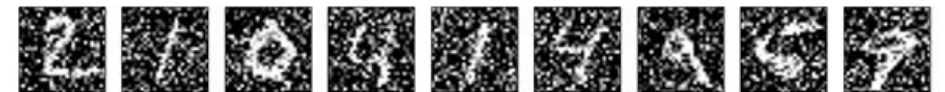- If we feed noisy data in, we will often get denoised data out



Examples from: https://blog.keras.io/building-autoencoders-in-keras.html

# Question:

- Autoencoders are able to compress because data sits on a manifold

- This doesn't mean that every point in the latent space will be on the manifold...

- GANs (covered later in this lecture)will learn a loss function that helps with this...
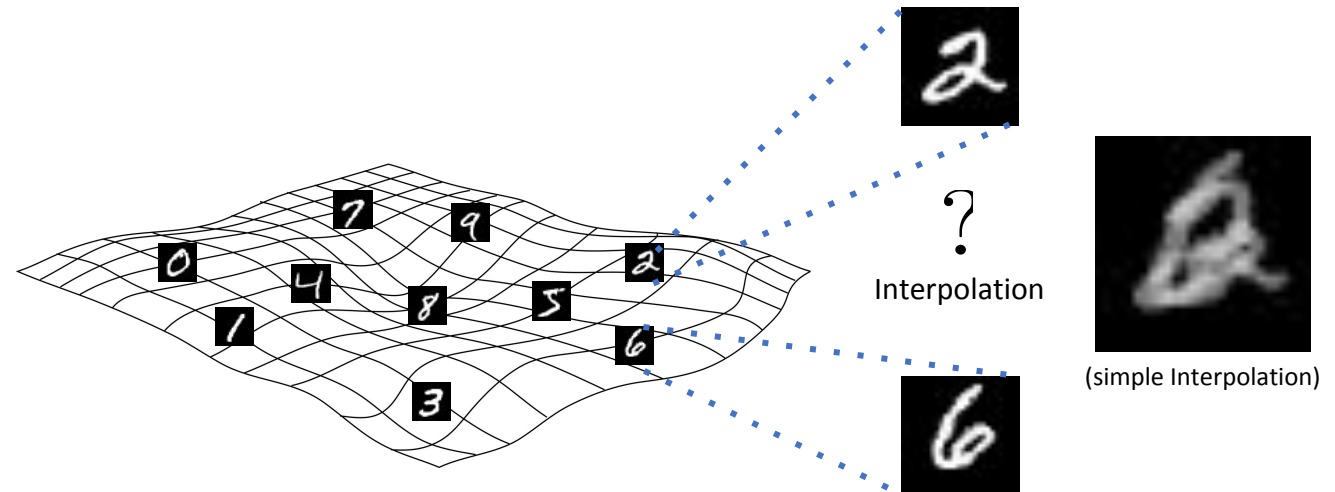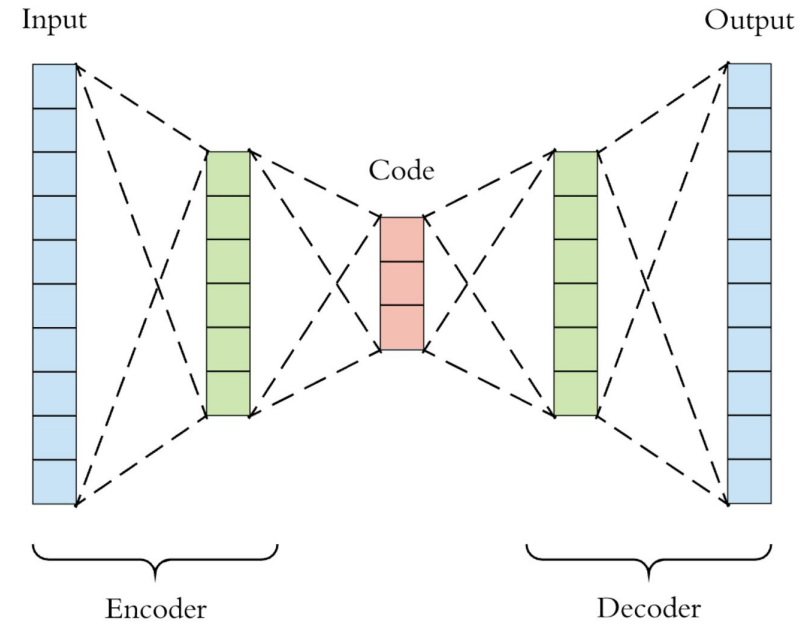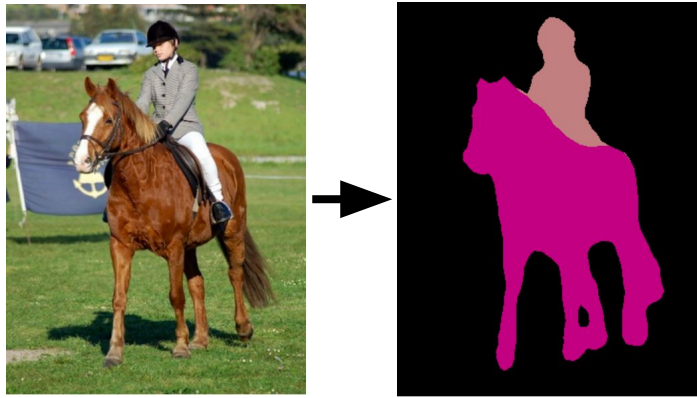
# Image-to-Image Applications

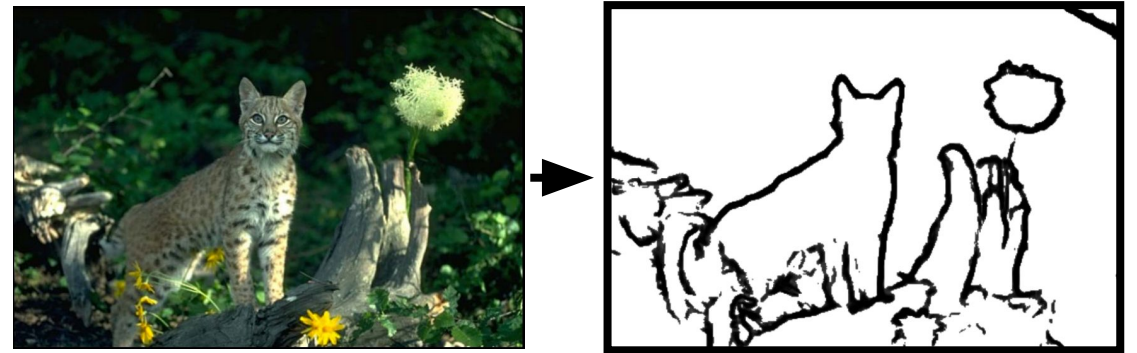Abe Davis, with slides from Jin Sun, Phillip Isola, and Richard Zhang

# Image prediction ("structured prediction")

## Object labeling:



[Long et al. 2015, …]
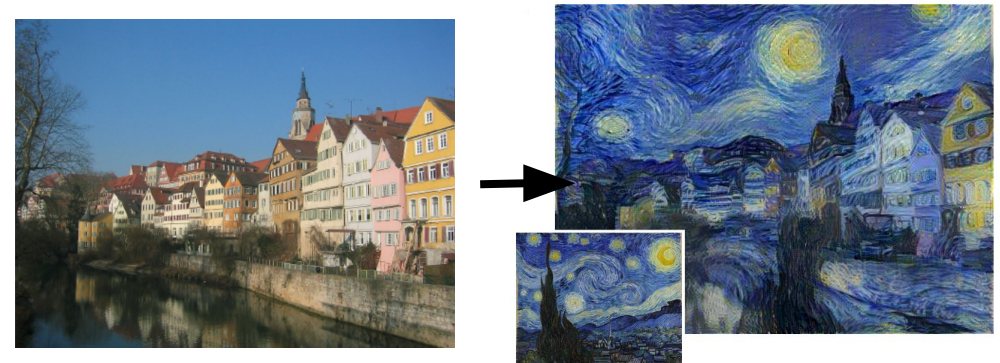
## Edge Detection:



[Xie et al. 2015, …]

## Text-to-photo:

"this small bird has a pink breast and crown…"
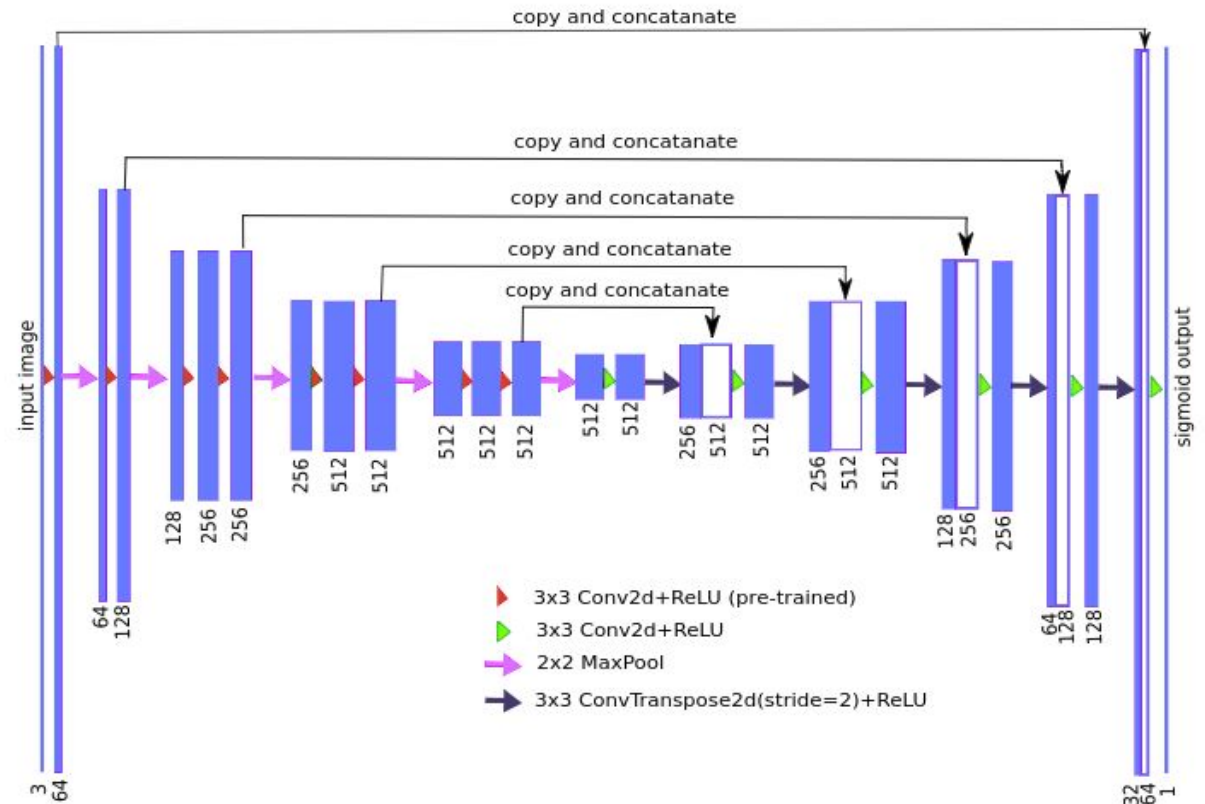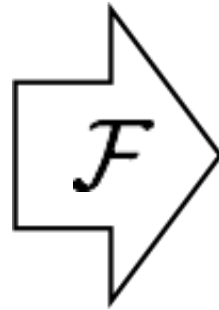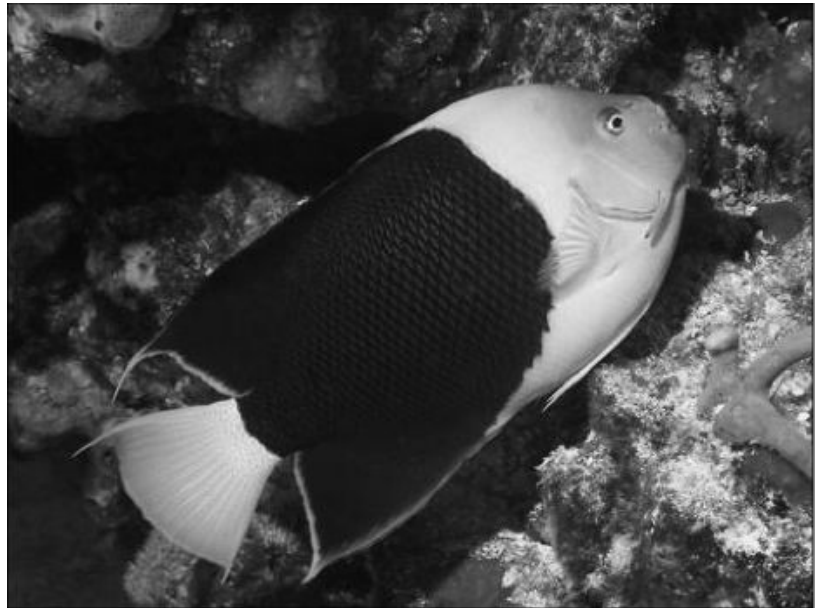


[Reed et al. 2016, …]

## Style transfer:



[Gatys et al. 2016, …]

# U-Net

- A popular network structure to generate same-sized output

- Similar to a convolutional autoencoder, but with "skip connections" that concatenate the output of earlier layers onto later layers

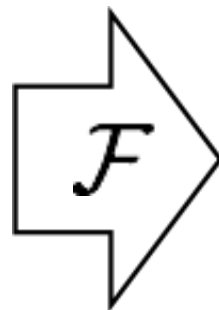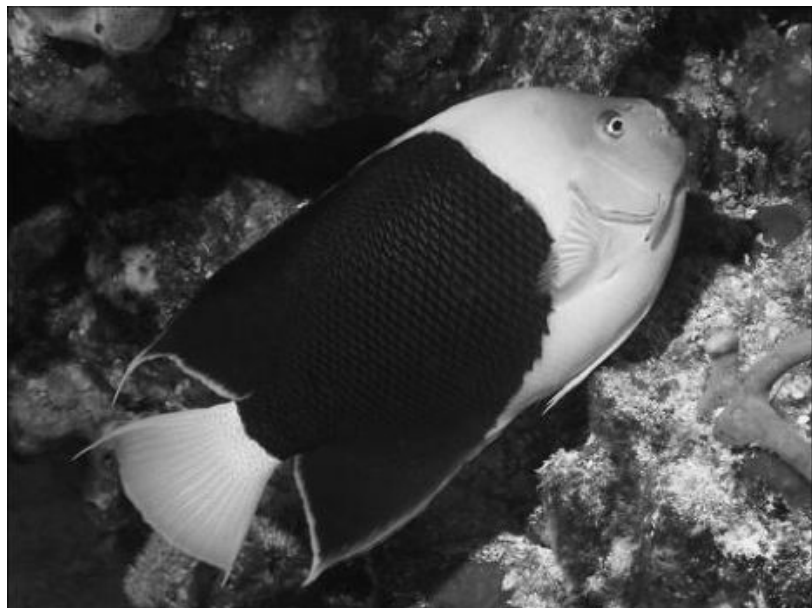- Great for learning transformations from one image to another

x

y

$\mathcal{F}$

Image Colorization

from Jin Sun, Richard Zhang, Phillip Isola
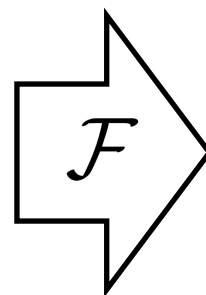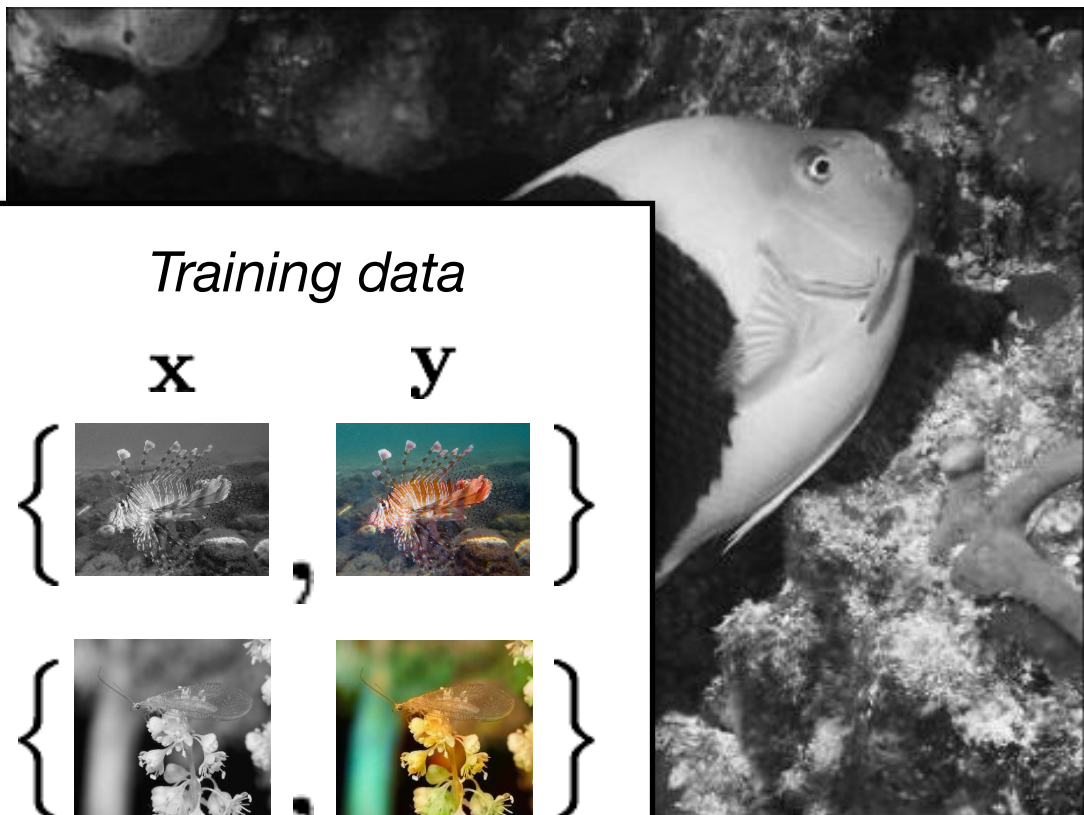
$$\mathbf{x} \qquad\qquad \mathbf{y}$$

$$\arg\min_{\mathcal{F}} \mathbb{E}_{\mathbf{x},\mathbf{y}}[L(\mathcal{F}(\mathbf{x}),\mathbf{y})]$$

"**What** should I do"

"**How** should I do it?"

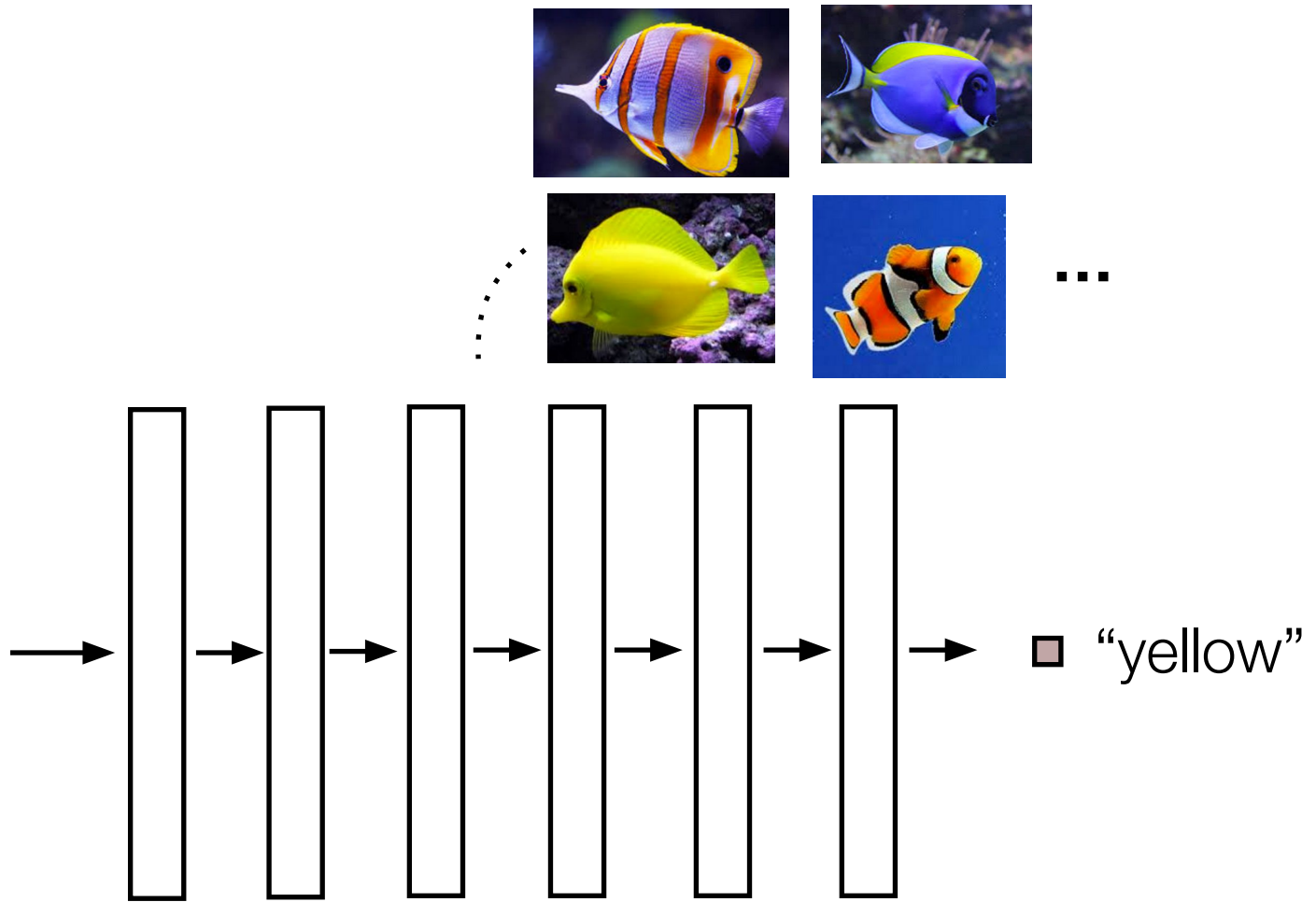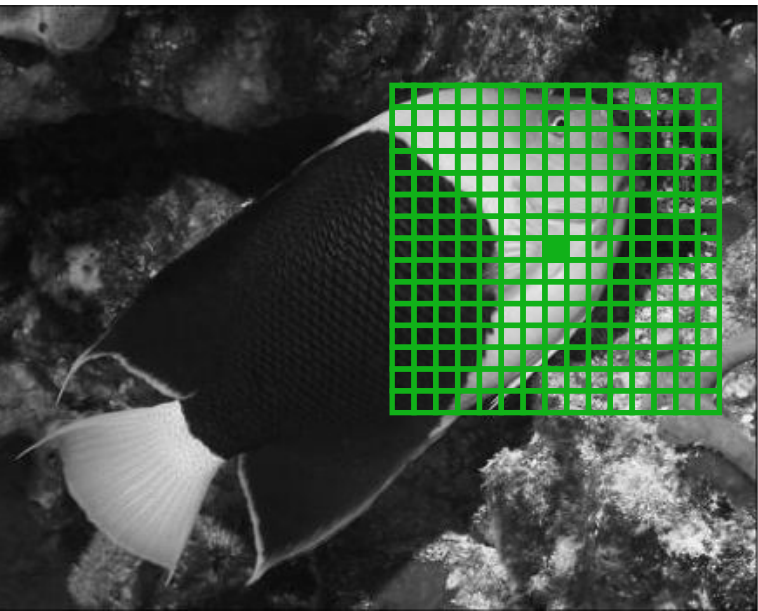from Jin Sun, Richard Zhang, Phillip Isola
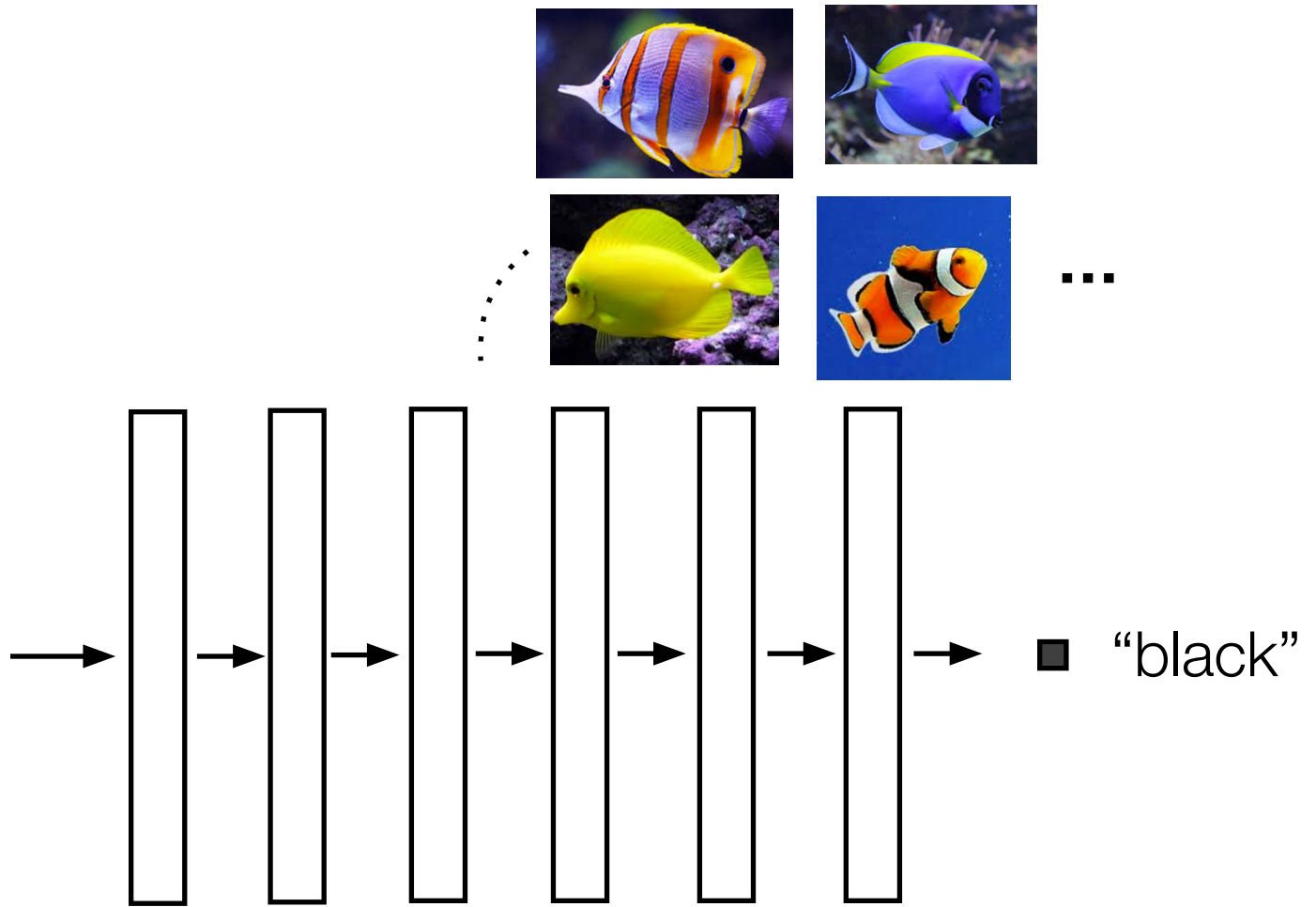
**x**

**y**

*Training data*

**x** **y**

L channel
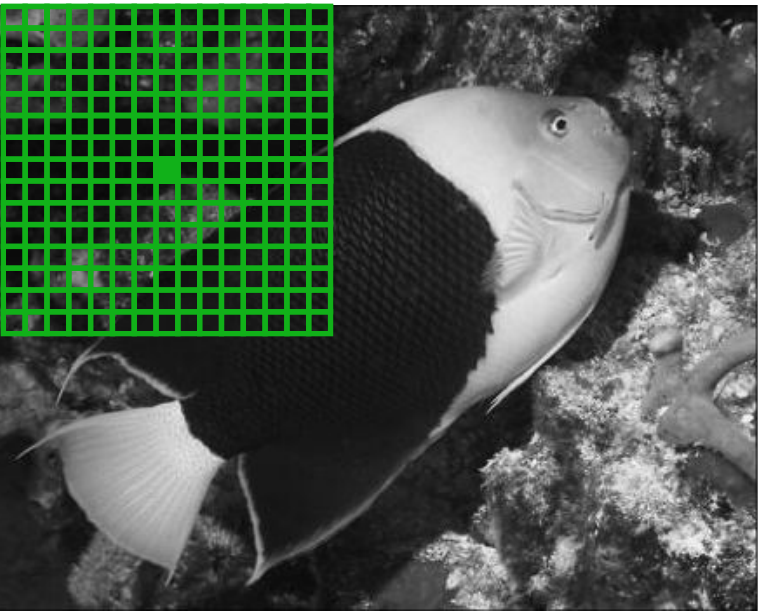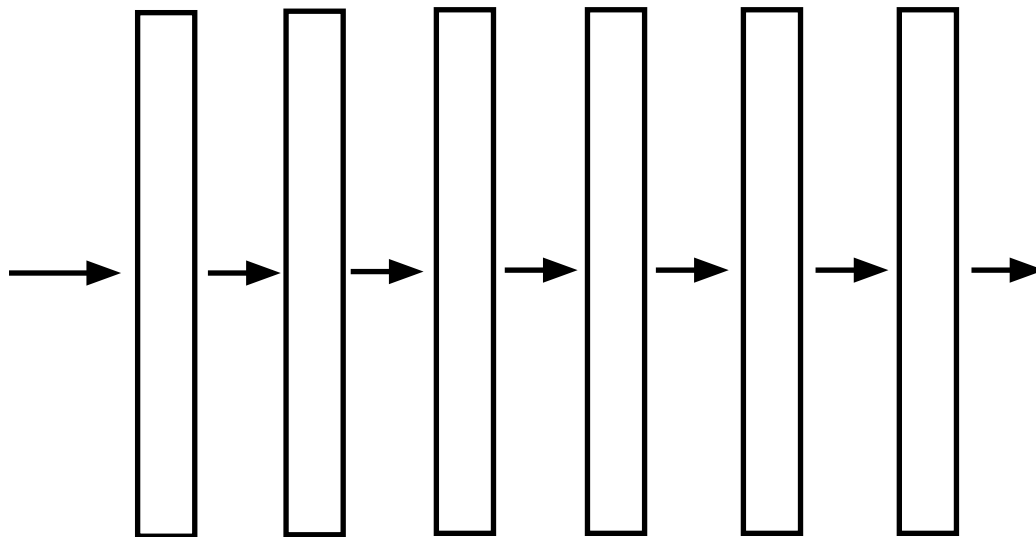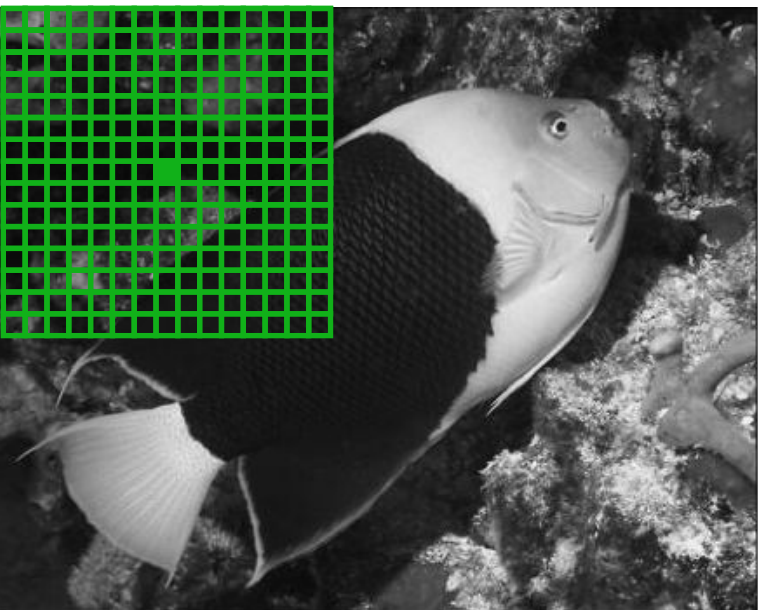
Color information: ab channels

$$\arg\min_{\mathcal{F}} \mathbb{E}_{\mathbf{x},\mathbf{y}}[L(\mathcal{F}(\mathbf{x}),\mathbf{y})]$$

Objective function
(loss)

Neural Network

from Jin Sun, Richard Zhang, Phillip Isola

"yellow"

from Jin Sun, Richard Zhang, Phillip Isola

"black"

from Jin Sun, Richard Zhang, Phillip Isola

from Jin Sun, Richard Zhang, Phillip Isola

# Basic loss functions

Prediction: $\hat{\mathbf{y}} = \mathcal{F}(\mathbf{x})$     Truth: $\mathbf{y}$

Classification (cross-entropy):

$$L(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_i \hat{\mathbf{y}}_i \log \mathbf{y}_i \quad \longleftarrow$$

How many extra bits it takes to correct the predictions

Least-squares regression:

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \left\| \hat{\mathbf{y}} - \mathbf{y} \right\|_2 \quad \longleftarrow$$

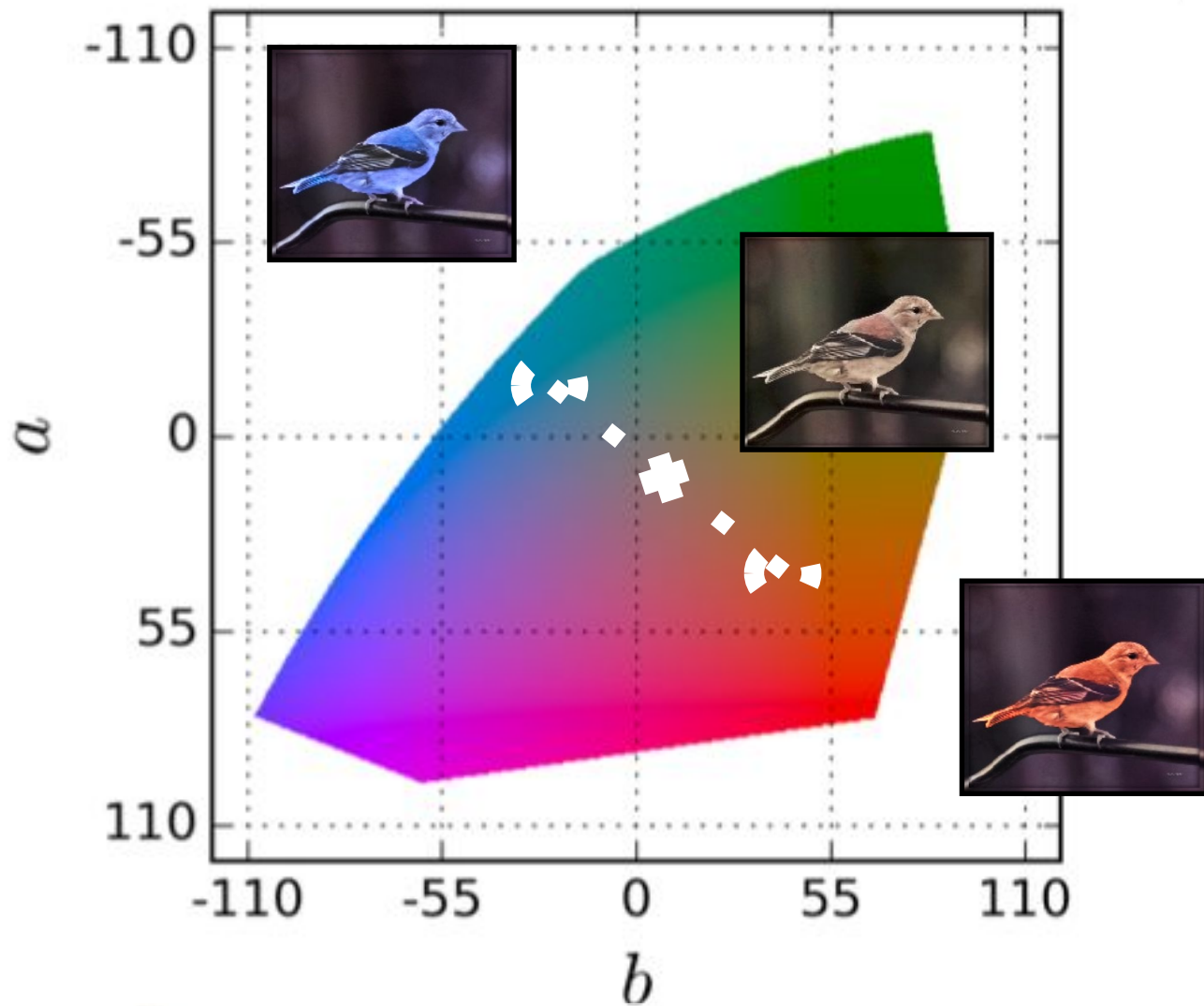How far off we are in Euclidean distance

# Designing loss functions

Input          Output          Ground truth



$$L_2(\widehat{Y}, Y) = \frac{1}{2} \sum_{h,w} \|Y_{h,w} - \widehat{Y}_{h,w}\|_2^2$$

$$\mathrm{L}_2(\widehat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \widehat{\mathbf{Y}}_{h,w}\|_2^2$$

# Designing loss functions

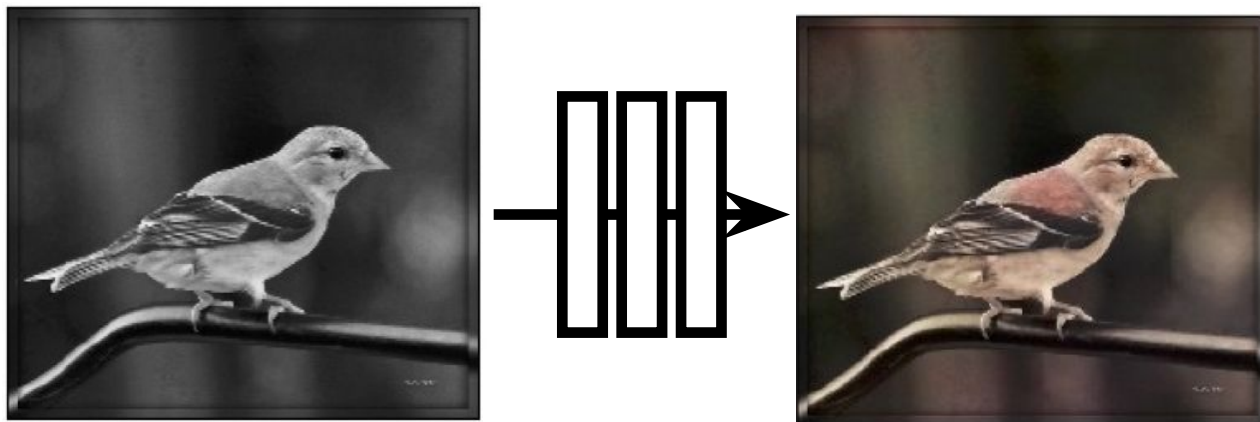Input                  Zhang et al. 2016                  Ground truth



Color distribution cross-entropy loss with colorfulness enhancing term.

[Zhang, Isola, Efros, ECCV 2016]
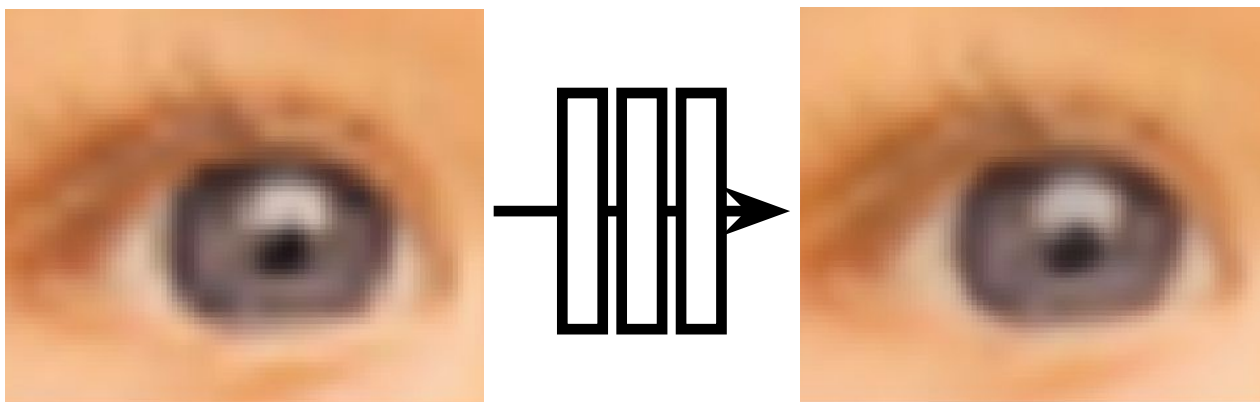
# Designing loss functions

## Image colorization



[Zhang, Isola, Efros, ECCV 2016]
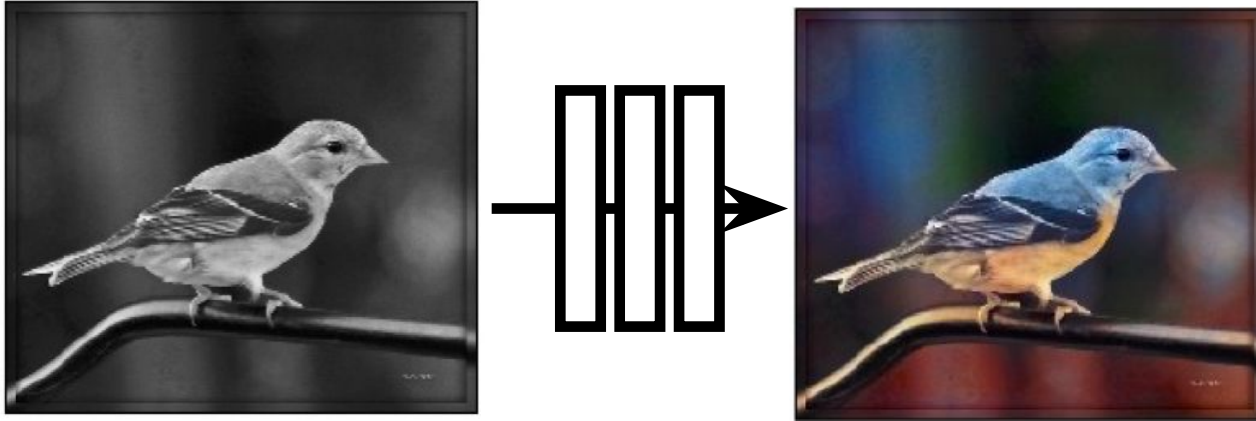
L2 regression

## Super-resolution



[Johnson, Alahi, Li, ECCV 2016]

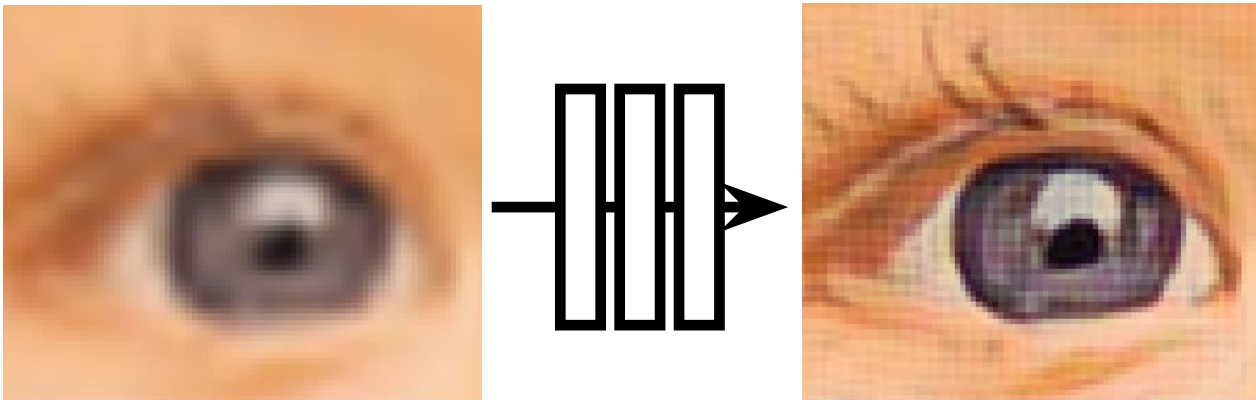L2 regression

# Designing loss functions

Image colorization



[Zhang, Isola, Efros, ECCV 2016]

Cross entropy objective,
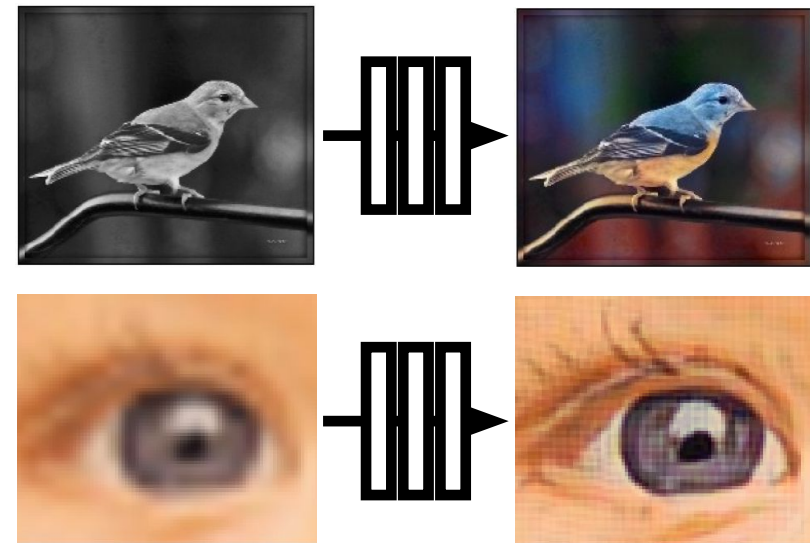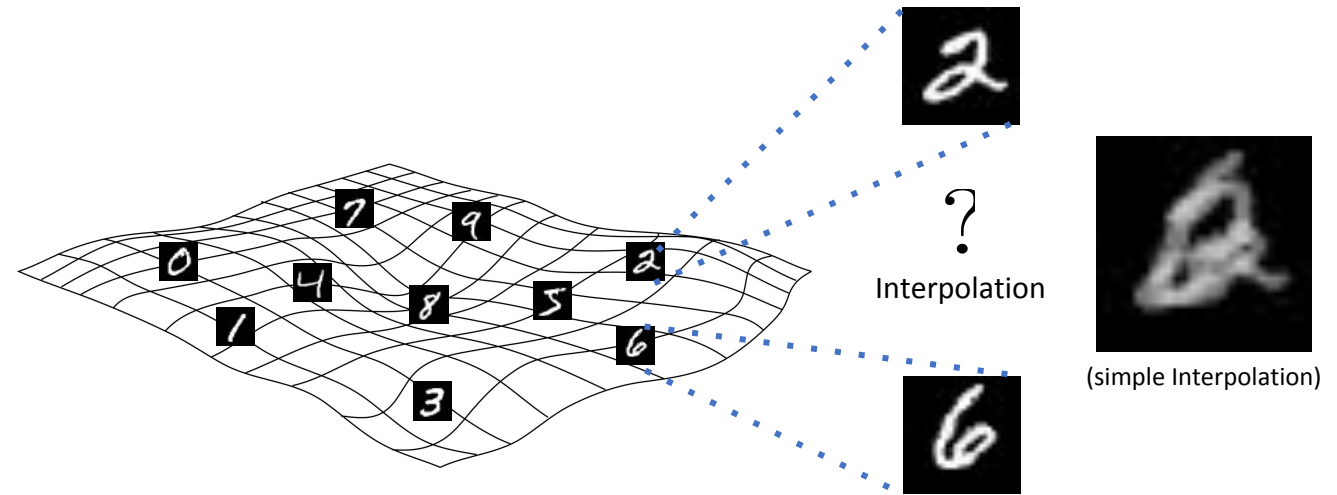with colorfulness term

Super-resolution



[Johnson, Alahi, Li, ECCV 2016]

Deep feature covariance
matching objective

# A Better Loss Function: Sticking to the Manifold

- How do we design a loss function that penalizes images that aren't on the image manifold?

- Key insight: we will *learn* our loss function by training a network to discriminate between images that are on the manifold and images that aren't
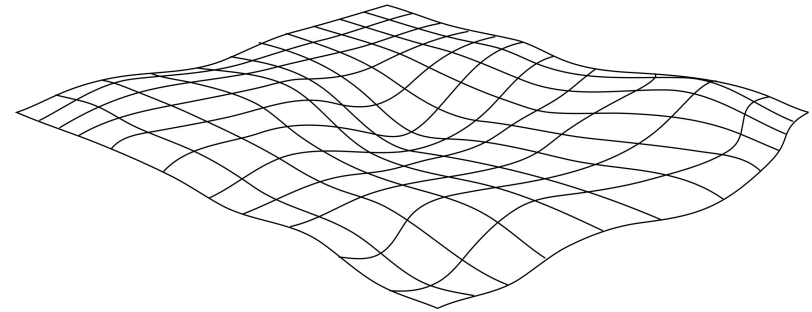
?
Interpolation

(simple Interpolation)

# Part 3: Generative Adversarial Networks (GANs)

Abe Davis, with slides from Jin Sun and Phillip Isola
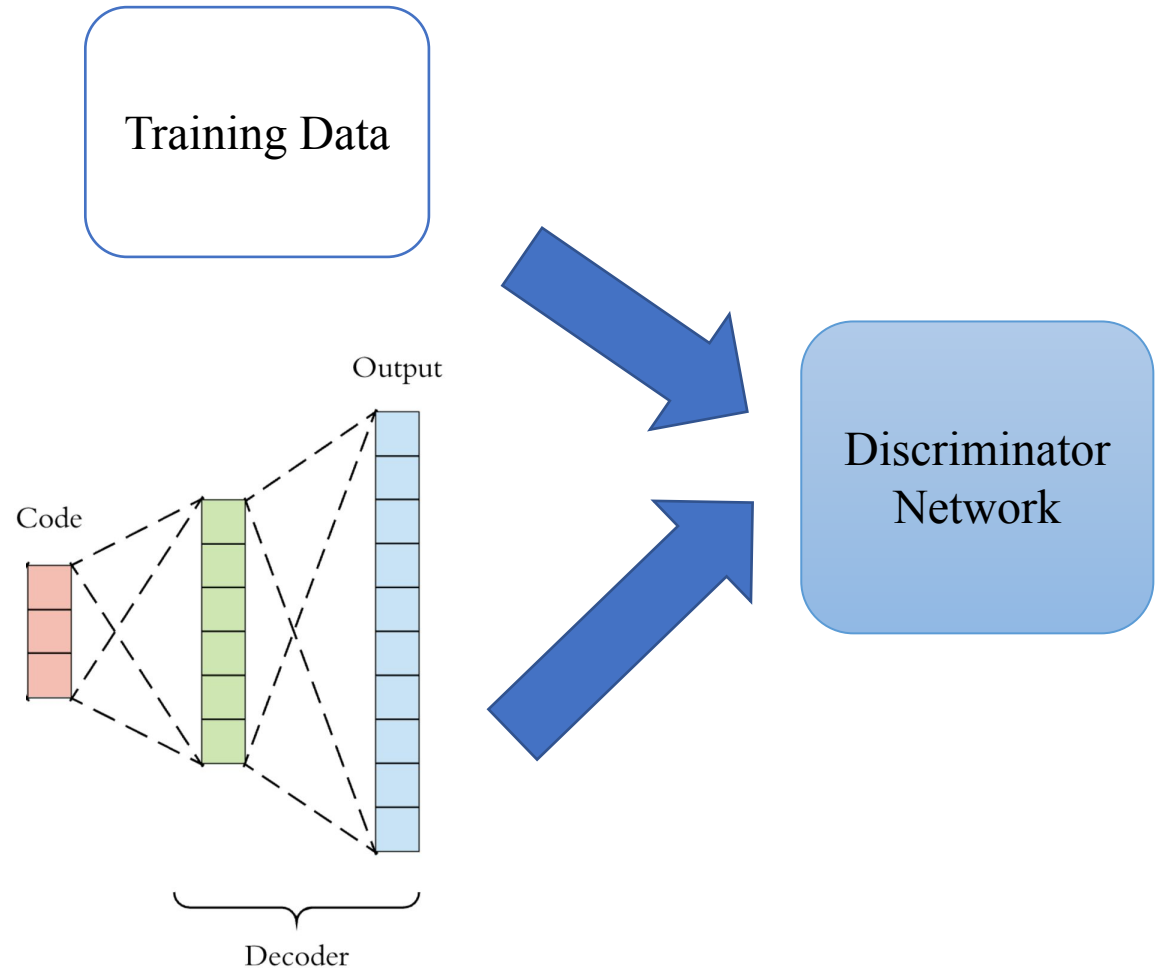
# Generative Adversarial Networks (GANs)

- Basic idea: Learn a mapping from some latent space to images on a particular manifold

- Example of a ***Generative Model:***
  - We can think of classification as a way to compute some $P(x)$ that tells us the probability that image $x$ is a member of a class.
  - Rather than simply evaluating this distribution, a generative model tries to learn a way to sample from it
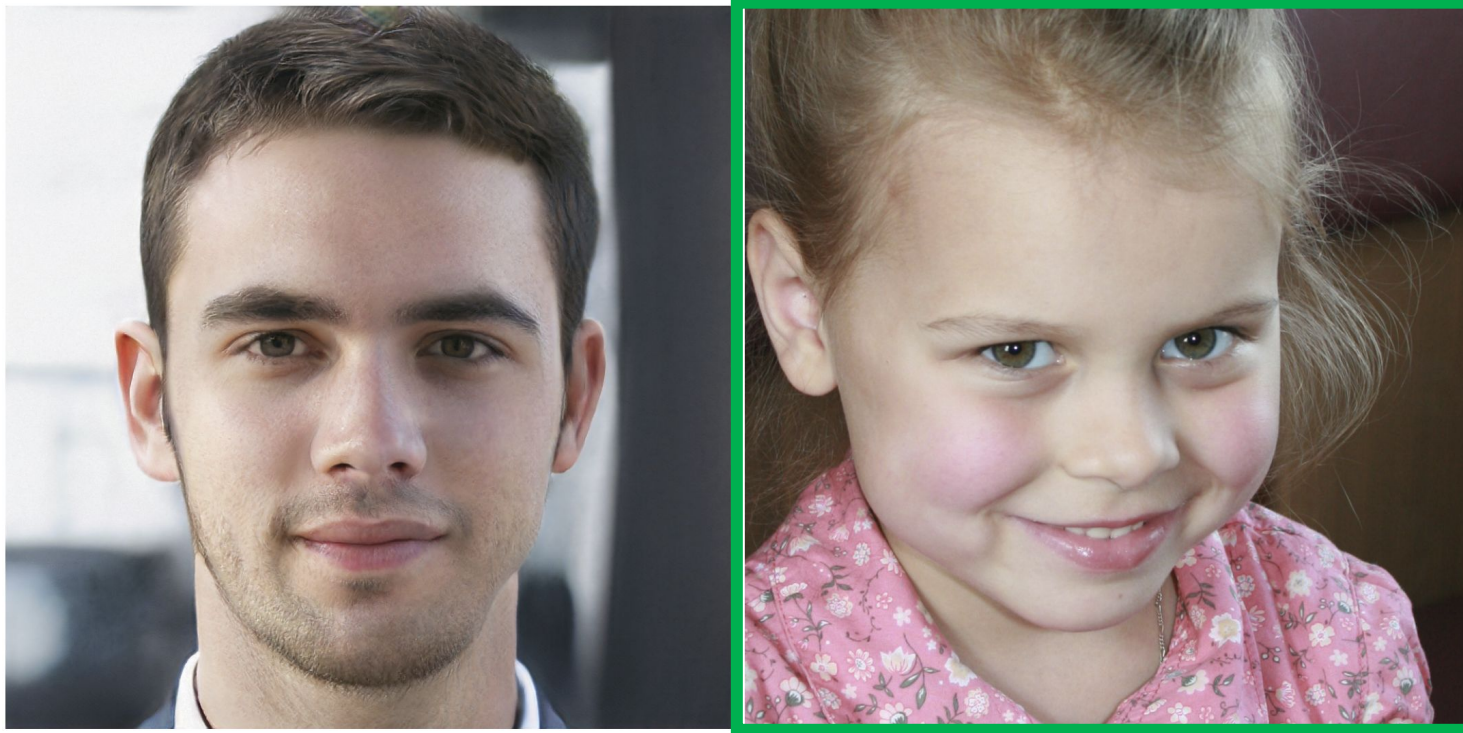
# Generative <u>Adversarial</u> Networks (GANs)

- Generator network has similar structure to the decoder of our autoencoder
  - Maps from some latent space to images

- We train it in an adversarial manner against a discriminator network
  - Generator tries to create output that is indistinguishable from training data
  - Discriminator tries to distinguish between generator output and training data

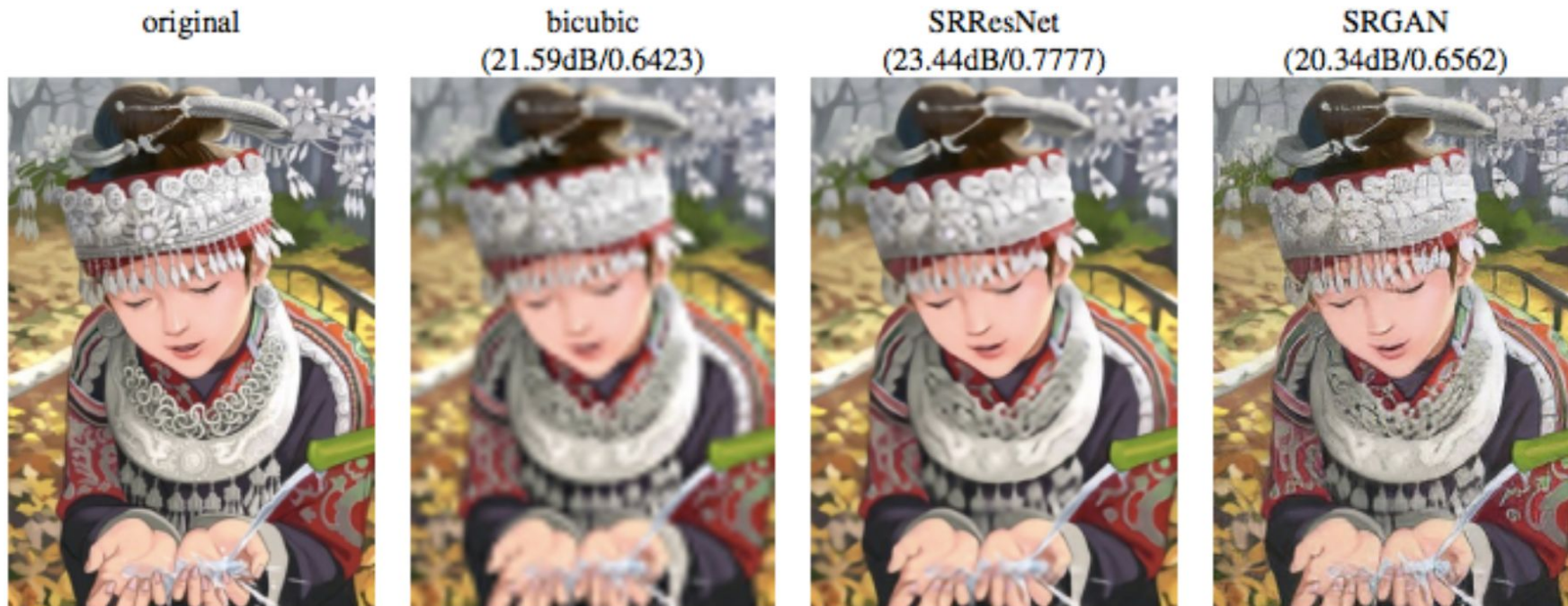# Example: Randomly Sampling the Space of Face Images

(Using Generative Adversarial Networks (GANs)
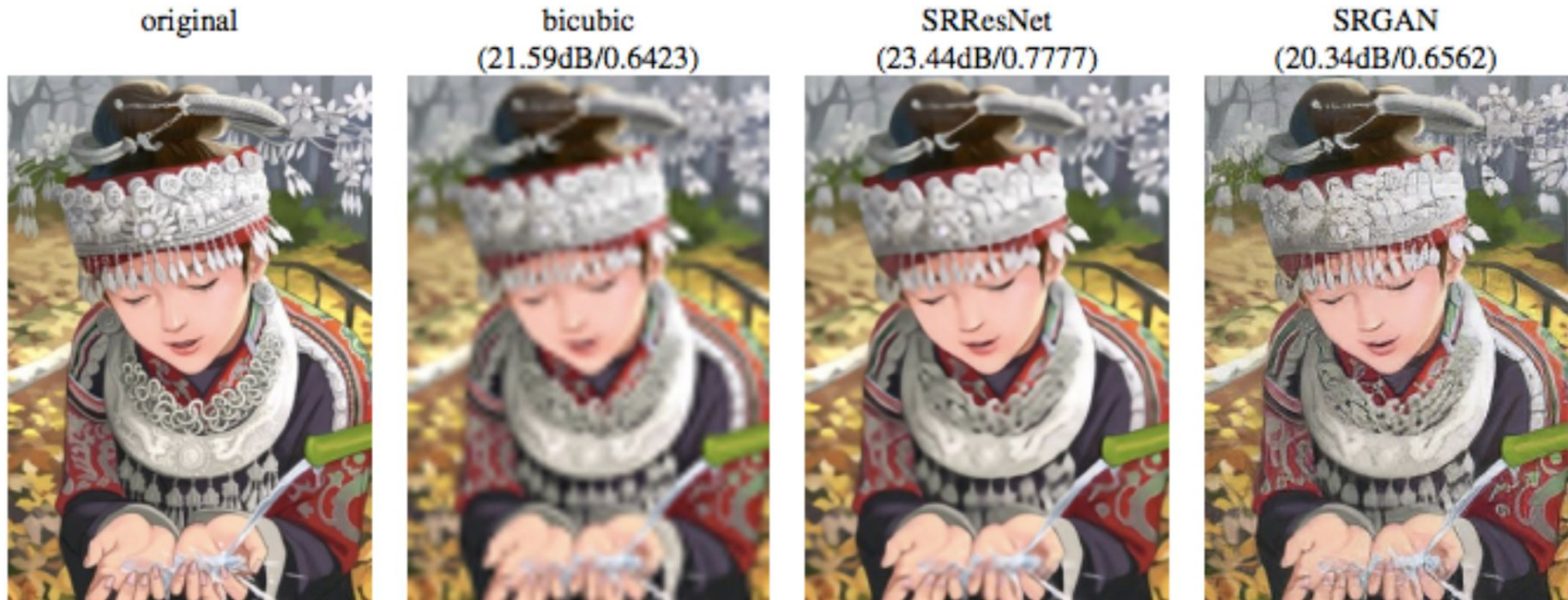


Which face is real?

# Conditional GANs

- Generate samples from a conditional distribution
- Example: generate high-resolution image conditioned on low resolution input



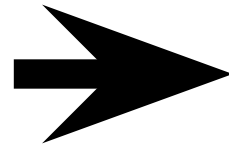| original | bicubic (21.59dB/0.6423) | SRResNet (23.44dB/0.7777) | SRGAN (20.34dB/0.6562) |

[Ledig et al 2016]

# Example: Single Image Super-Resolution

- Generate natural image, conditioned on a lower–resolution version of the image
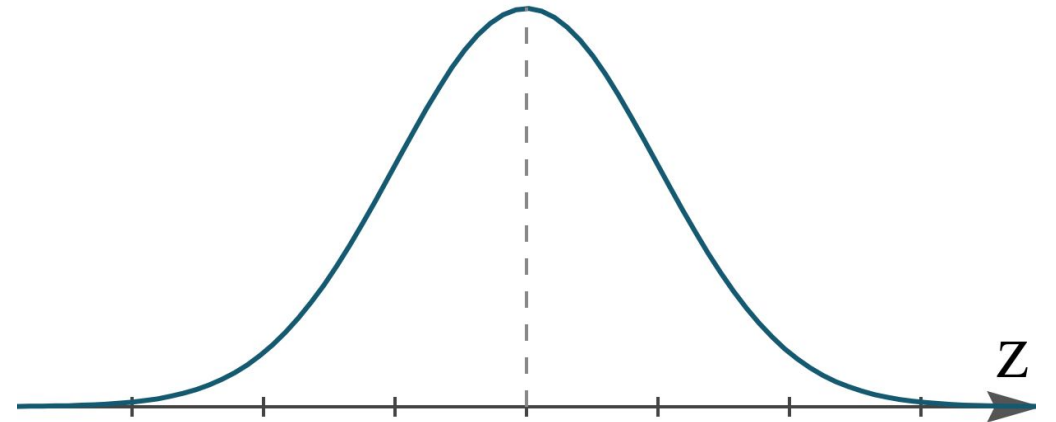


[Ledig et al 2016]

# Conditional GANs
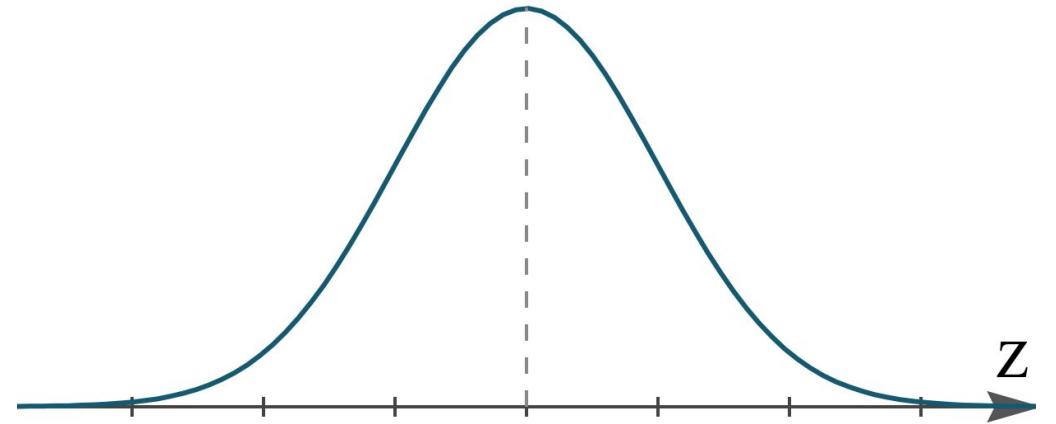


[Goodfellow et al., 2014]

[Isola et al., 2017]

# Generative Models: Generate Samples from a Distribution

- We can look at classification as a way to compute some P(x) that tells us the probability that image x is a member of a class.

- Rather than simply evaluating this distribution, is there some way for us to generate samples from it?

z

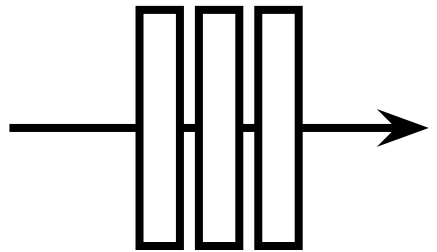# Basic Idea Part 2: Generate Samples from a *Conditional* Distribution

- Can we generate samples from our distribution *conditioned on some input*?

- In other words, can we generate samples from the conditional distribution $P(x|c)?$
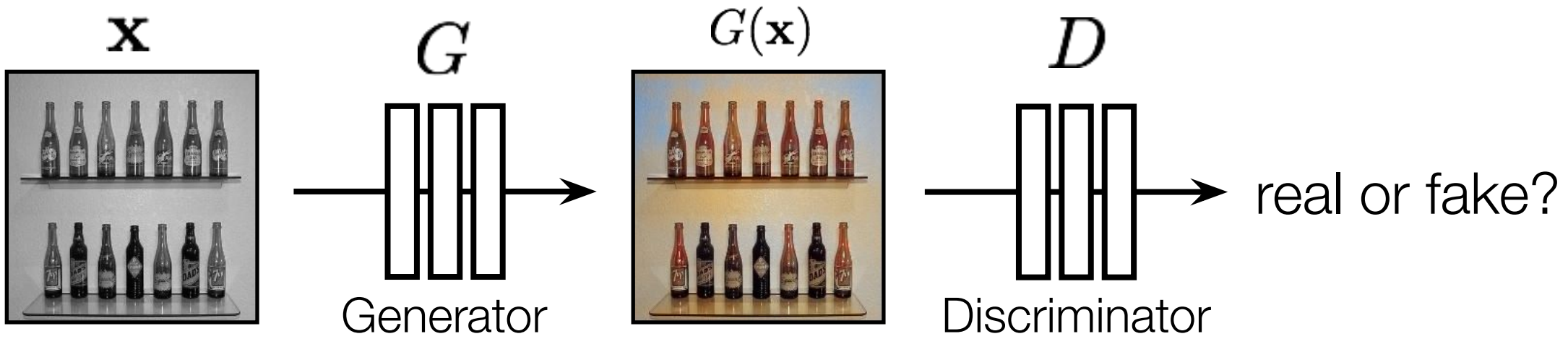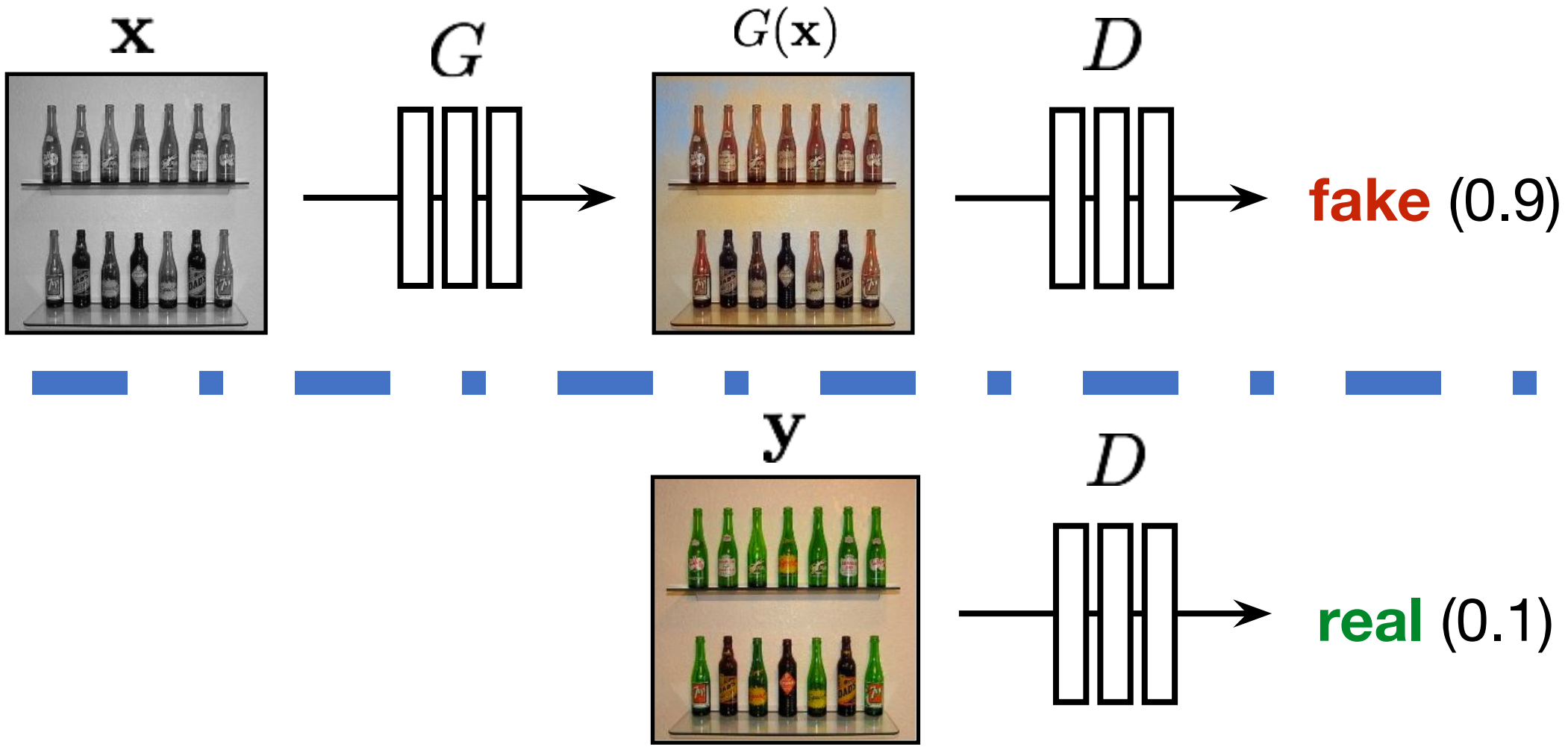
$\mathbf{x}$     $G$     $G(\mathbf{x})$

Generator

[Goodfellow et al., 2014]

$\mathbf{x}$  $G$  $G(\mathbf{x})$  $D$

Generator  →  real or fake?  Discriminator

**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes

[Goodfellow et al., 2014]

$\mathbf{x}$    $G$    $G(\mathbf{x})$    $D$    **fake** (0.9)

$\mathbf{y}$    $D$    **real** (0.1)

(Identify generated images as fake)      (Identify training images as real)

$$\arg\max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}\left[\; \boxed{\log D(G(\mathbf{x}))} \; + \; \boxed{\log(1 - D(\mathbf{y}))} \;\right]$$

[Goodfellow et al., 2014]

**G** tries to synthesize fake images that *fool* **D**:

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

[Goodfellow et al., 2014]

**G** tries to synthesize fake images that *fool* the *best* **D**:

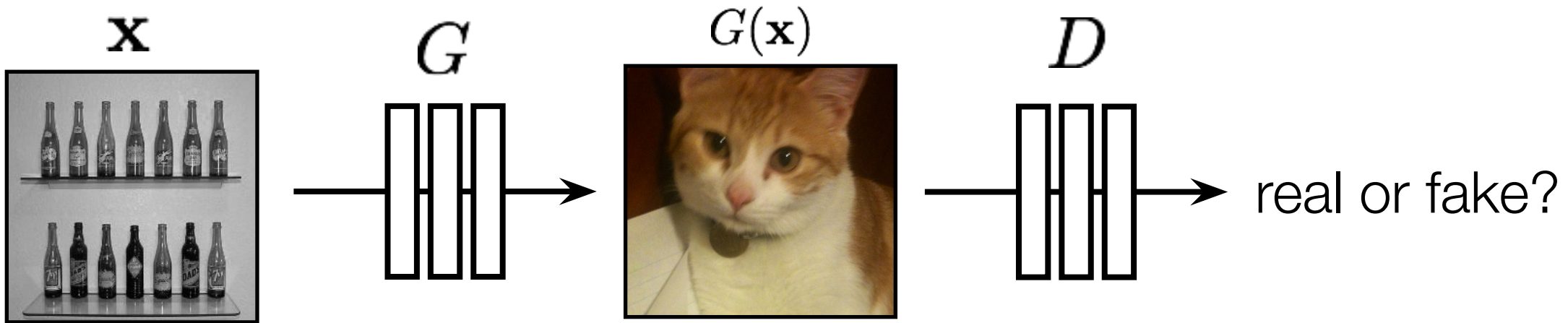$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x},\mathbf{y}}[ \ \log D(G(\mathbf{x})) \ + \ \log(1 - D(\mathbf{y})) \ ]$$

[Goodfellow et al., 2014]

**G**'s perspective: **D** is a loss function.

Rather than being hand-designed, it is *learned*.

[Goodfellow et al., 2014]
[Isola et al., 2017]

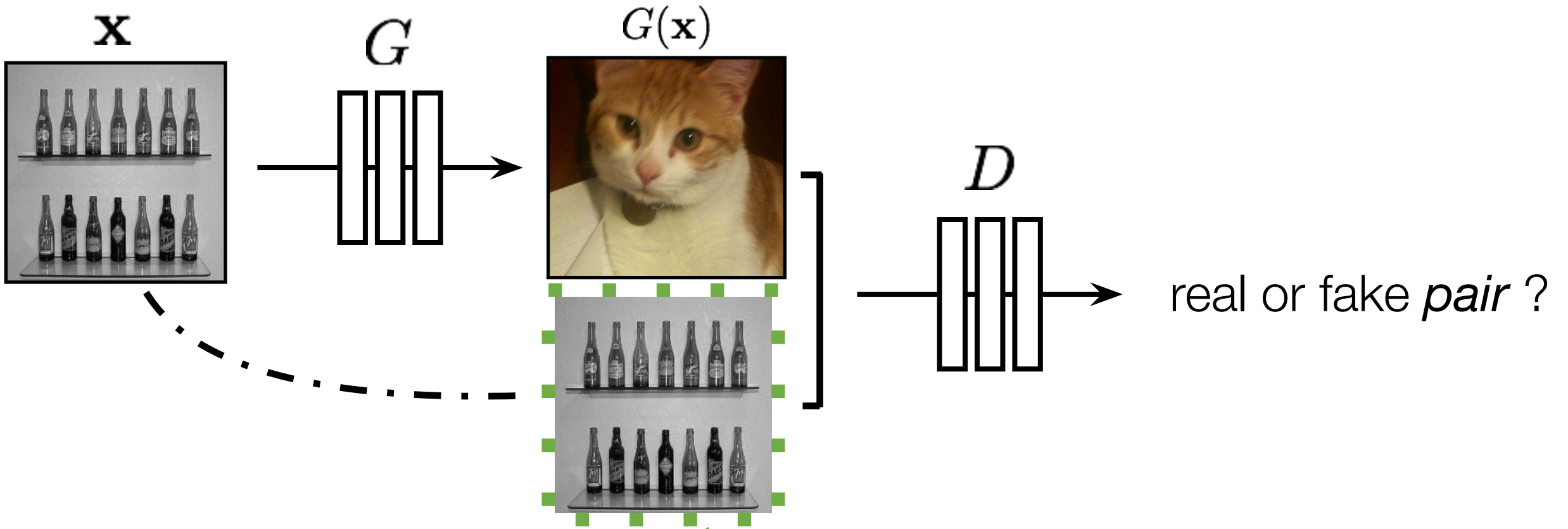$$\arg \min_G \max_D \ \mathbb{E}_{\mathbf{x},\mathbf{y}}[\ \log D(G(\mathbf{x})) \ + \ \log(1 - D(\mathbf{y}))\ ]$$
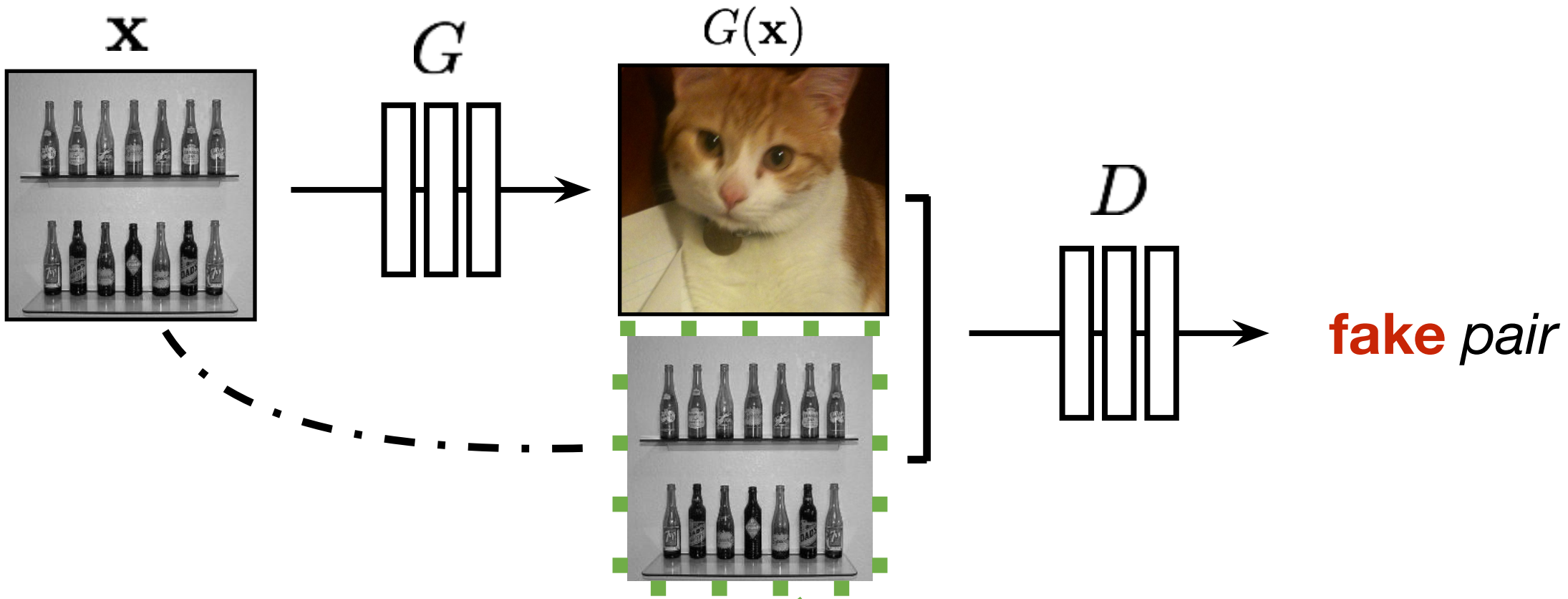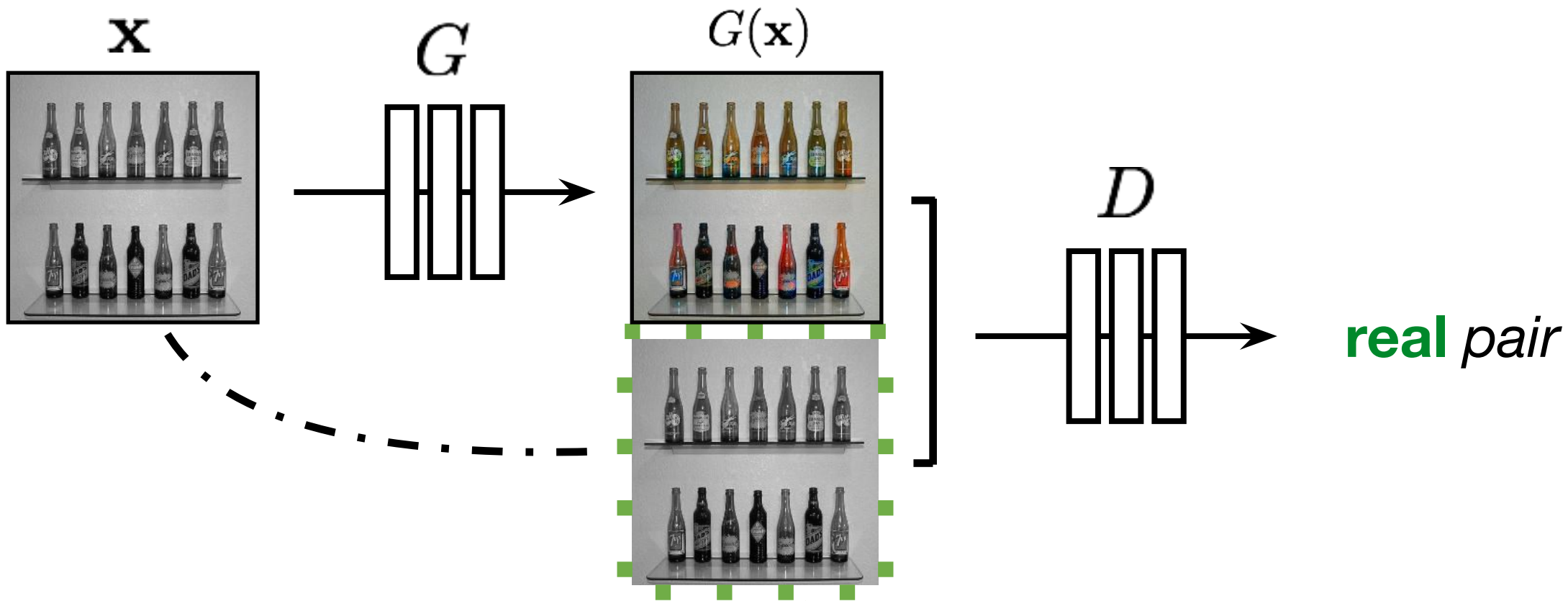
[Goodfellow et al., 2014]

$\mathbf{x}$  $G$  $G(\mathbf{x})$  $D$  **real!** ("Aquarius")

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

[Goodfellow et al., 2014]

$$\arg \min_G \max_D \ \mathbb{E}_{\mathbf{x},\mathbf{y}}[\ \log D(G(\mathbf{x})) \ + \ \log(1 - D(\mathbf{y}))\ ]$$
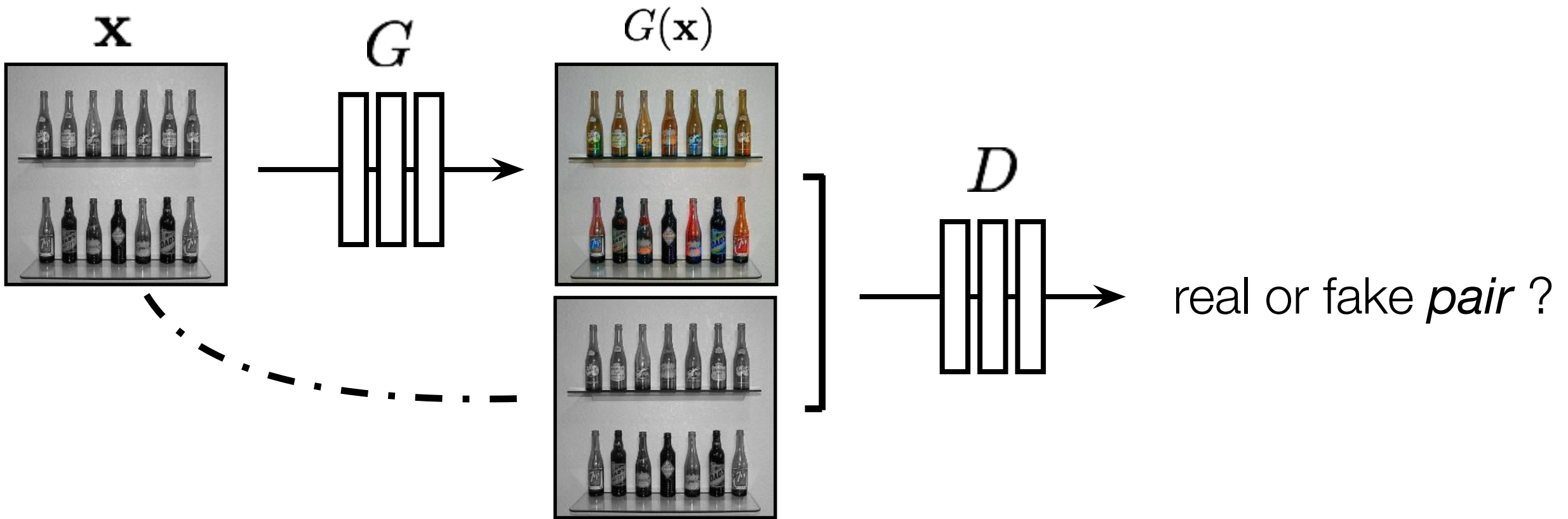
[Goodfellow et al., 2014]
[Isola et al., 2017]

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

[Goodfellow et al., 2014]

[Isola et al., 2017]

$$\arg\min_G \max_D \ \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[\ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))\ \big]$$
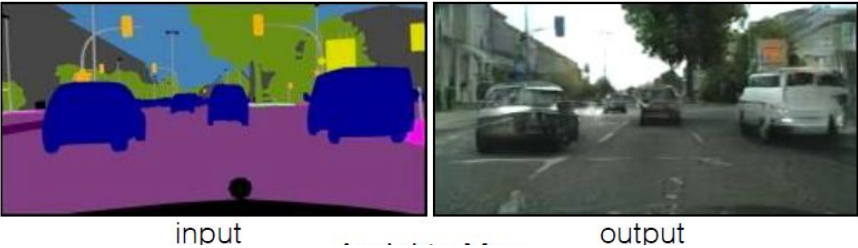
[Goodfellow et al., 2014]
[Isola et al., 2017]

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

[Goodfellow et al., 2014]

[Isola et al., 2017]

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

[Goodfellow et al., 2014]

[Isola et al., 2017]

# More Examples of Image-to-Image Translation with GANs

- We have pairs of corresponding training images
- Conditioned on one of the images, sample from the distribution of likely corresponding images

**Edges to Image**

**Segmentation to Street Image**



input     output

**Aerial Photo To Map**



input     output

Input    Ground truth    Output

# BW →
# Color



Input    Output         Input    Output         Input    Output

Data from [Russakovsky et al. 2015]

Input                           Output                          Groundtruth



Data from
[maps.google.com]

# Labels → Street Views

Input labels



Synthesized image



Data from [Wang et al, 2018]
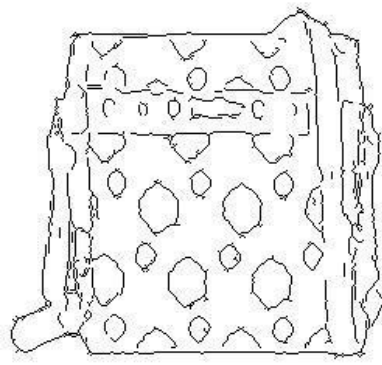
# Day → Night

Input    Output          Input    Output          Input    Output
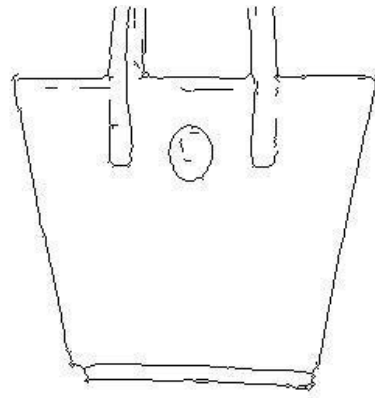
# Edges → Images

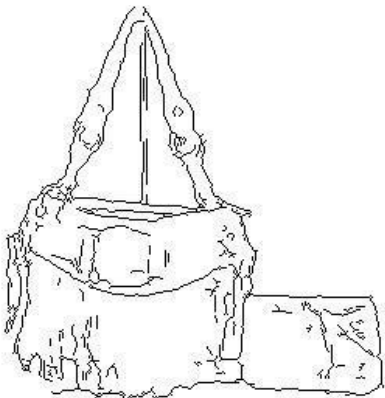Input     Output          Input     Output          Input     Output



Edges from [Xie & Tu, 2015]

# Demo



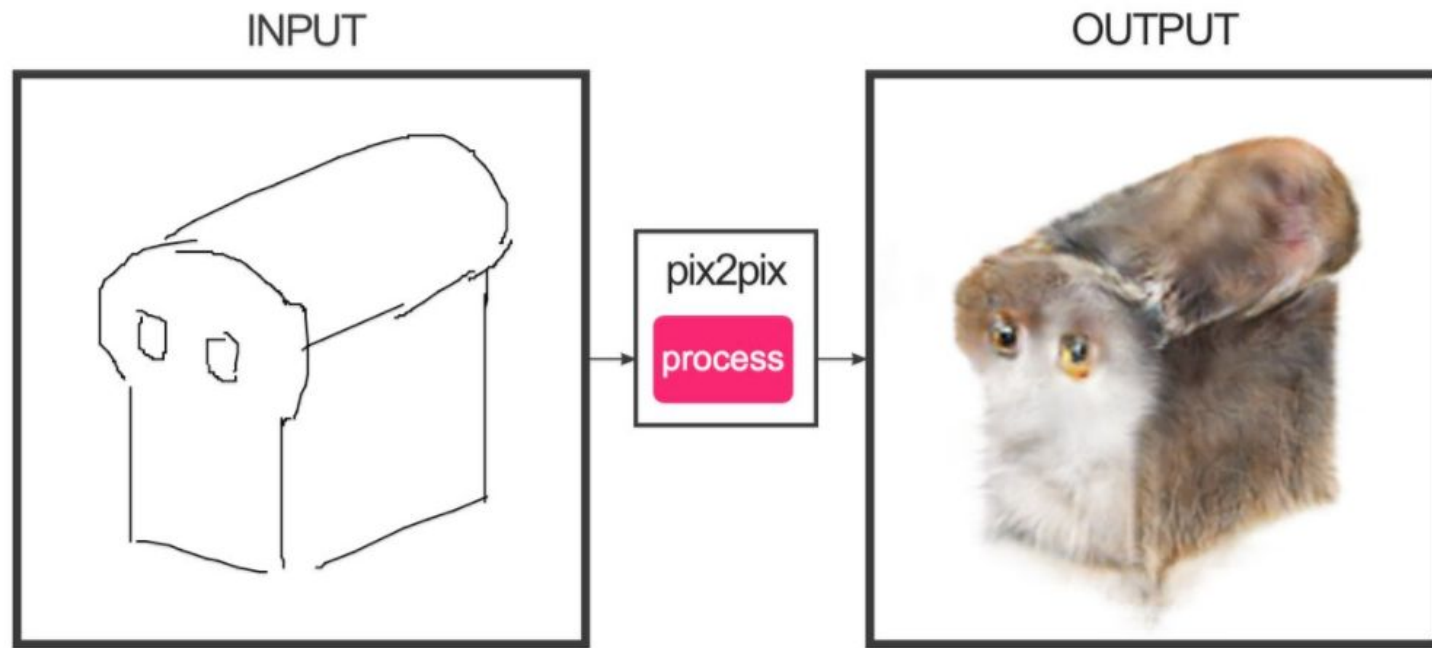https://affinelayer.com/pixsrv/
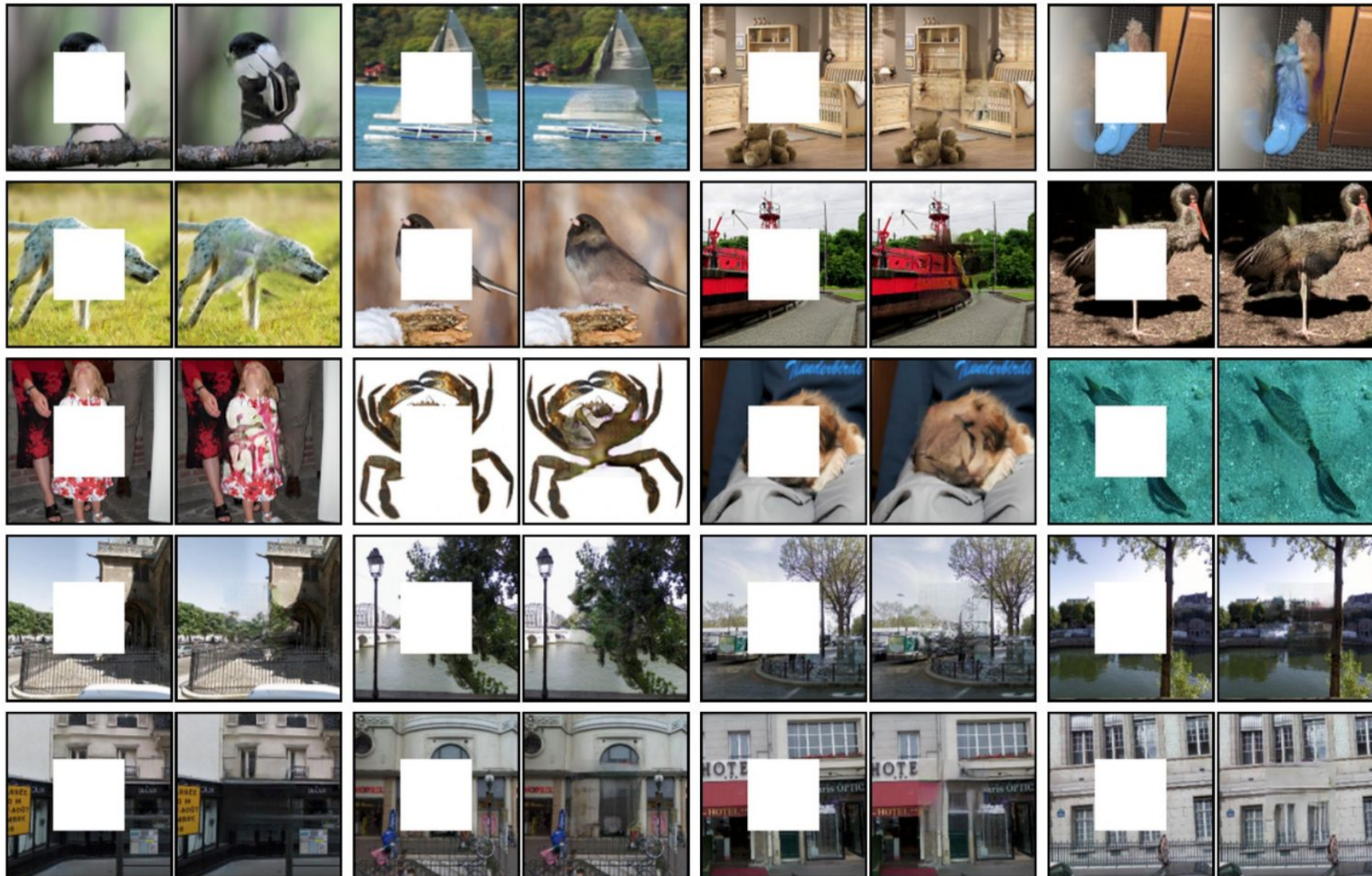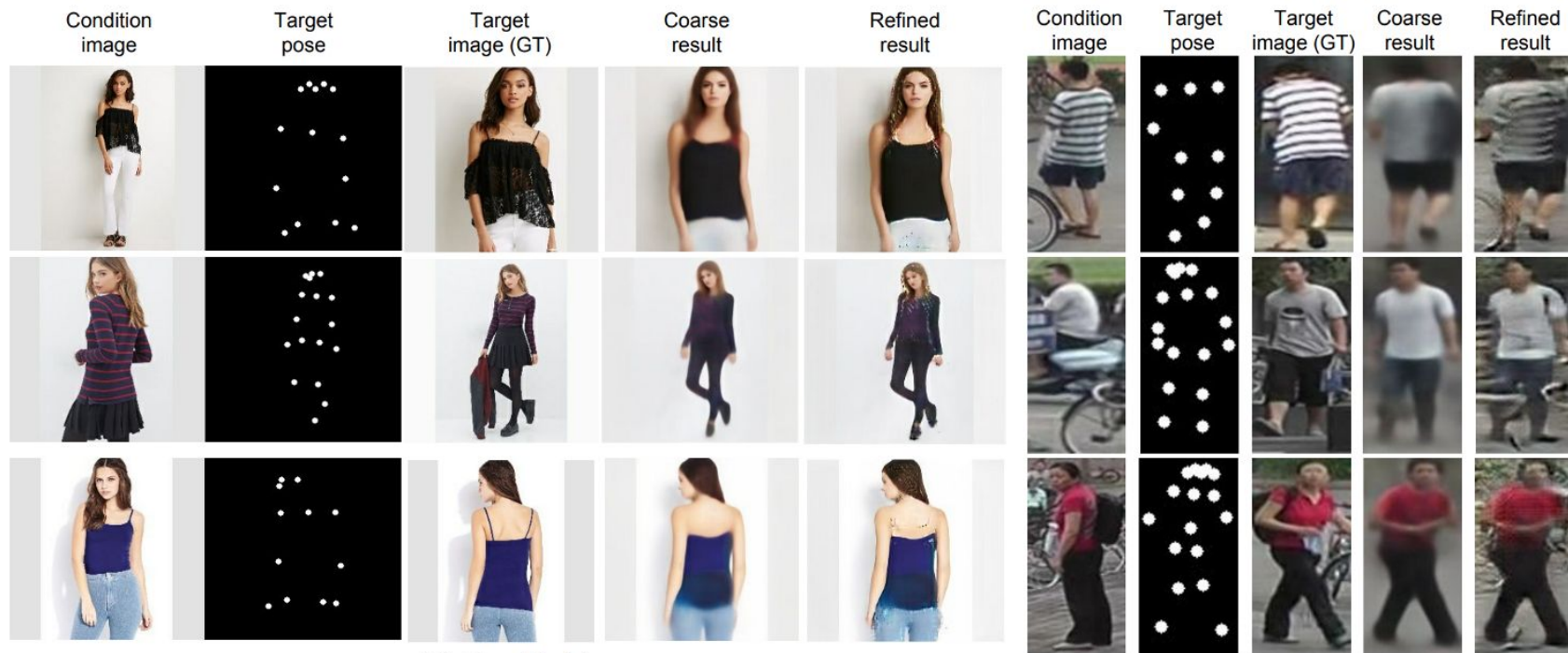
INPUT

OUTPUT

pix2pix

process

Ivy Tasi @ivymyt

Vitaly Vidmirov @vvid

# Image Inpainting

# Pose-guided Generation



(a) DeepFashion

(b) Market-1501

(c) Generating from a sequence of poses

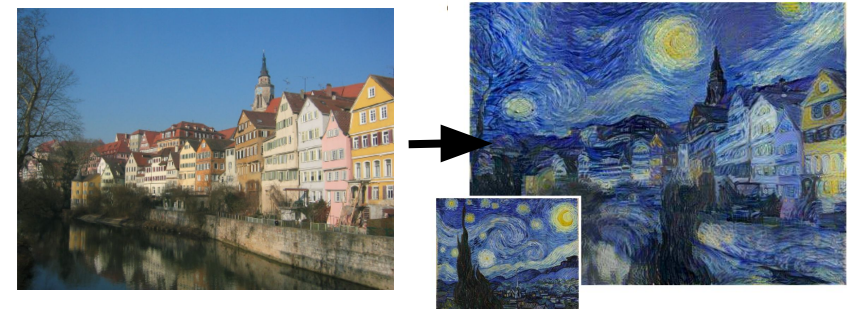Data from [Ma et al., 2018]

# Challenges —> Solutions

- Output is high-dimensional, structured object
  - Approach: Use a deep net, D, to analyze output!

- Uncertainty in mapping; many plausible outputs
  - Approach: D only cares about "plausibility", doesn't hedge

- Lack of supervised training data
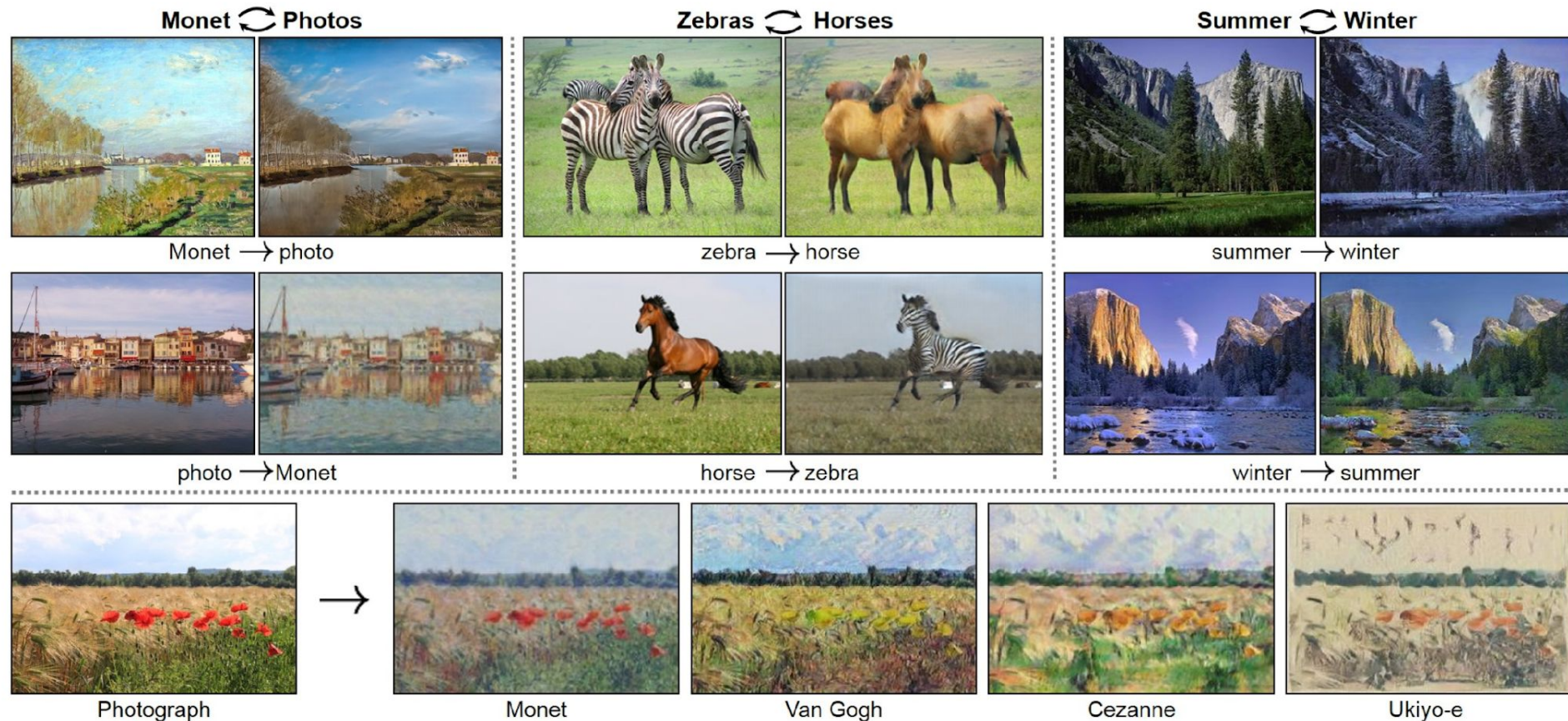  - Approach: ?



"this small bird has a pink breast and crown…"

# Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks

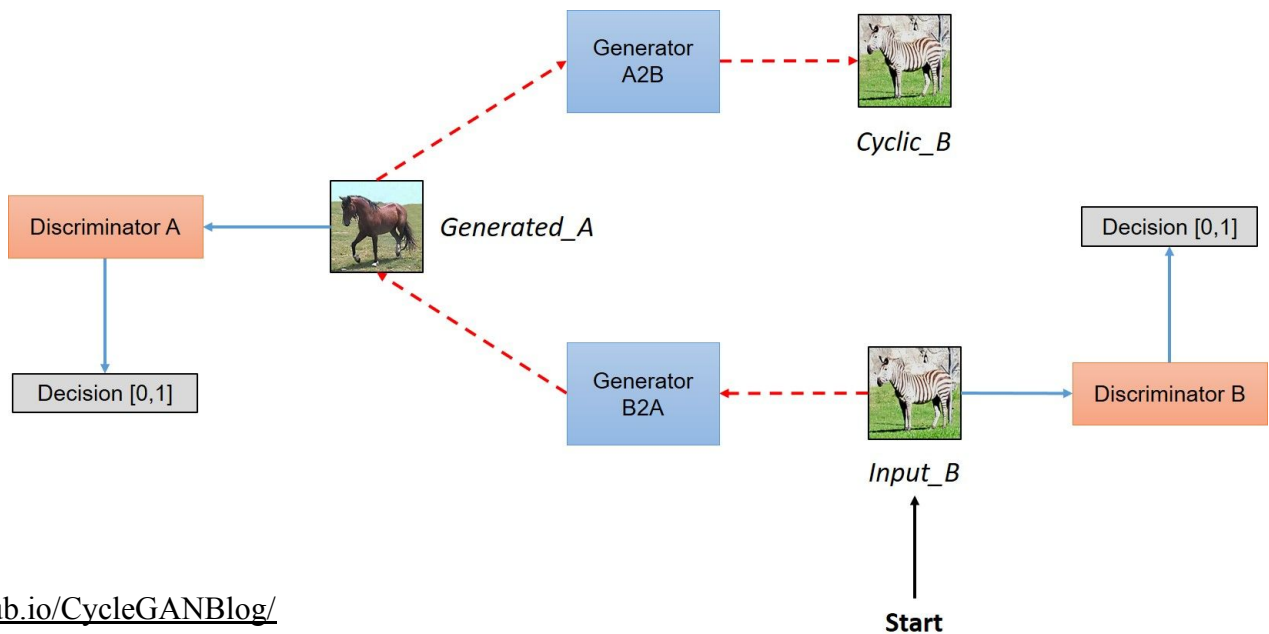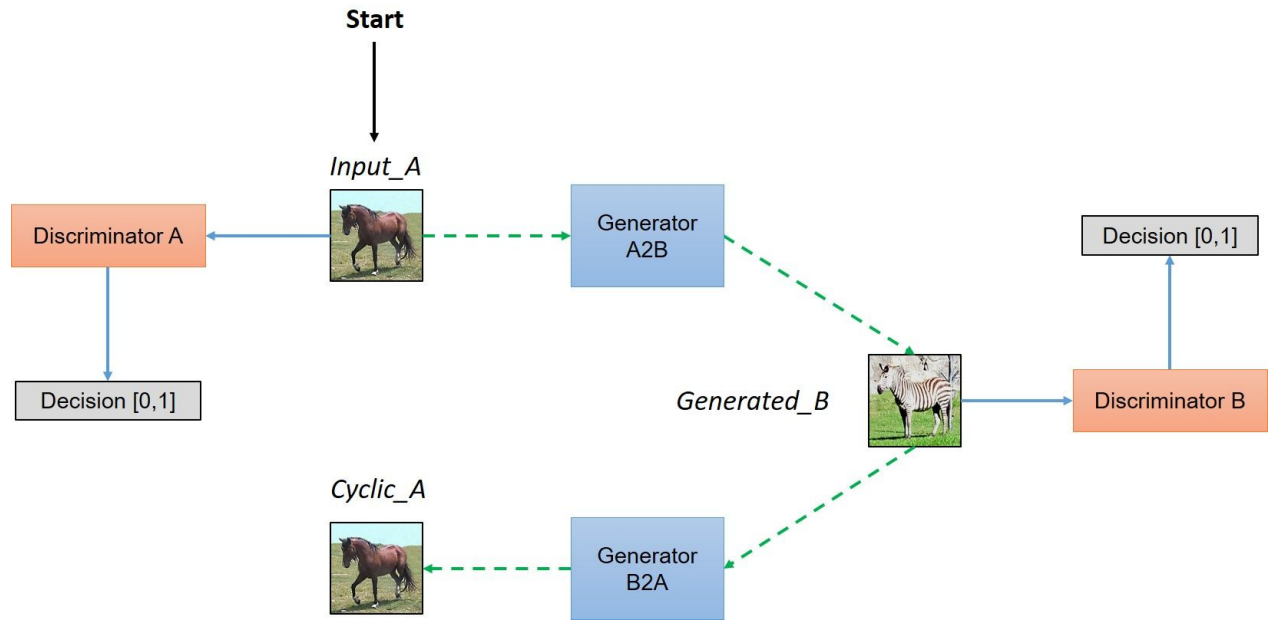**Jun-Yan Zhu*    Taesung Park*    Phillip Isola    Alexei A. Efros**

**UC Berkeley**

In ICCV 2017

[Paper] [Code (Torch)] [Code (PyTorch)]



https://junyanz.github.io/CycleGAN/

# StyleGAN



https://github.com/NVlabs/stylegan

# Questions?