

Perceptual Loss, GANs (part I)


Jun-Yan Zhu

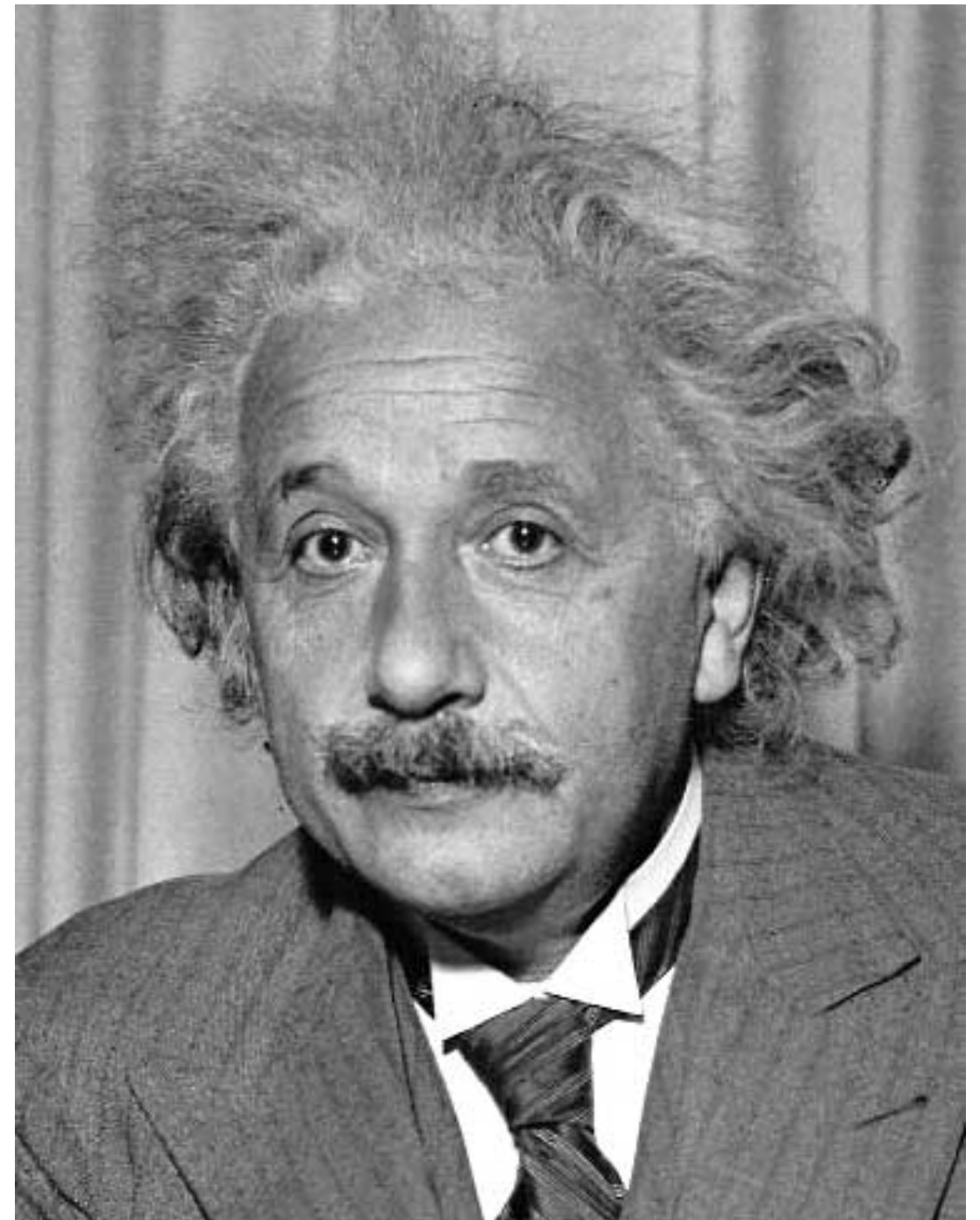
16-726 Learning-based Image Synthesis, Spring 2021

many slides from Phillip Isola, Richard Zhang, Alyosha Efros


HW1 (hints)

Template matching

- Goal: find  in image
- Main challenge: What is a good similarity or distance measure between two patches?
 - Correlation
 - Zero-mean correlation
 - Sum Square Difference
 - Normalized Cross Correlation

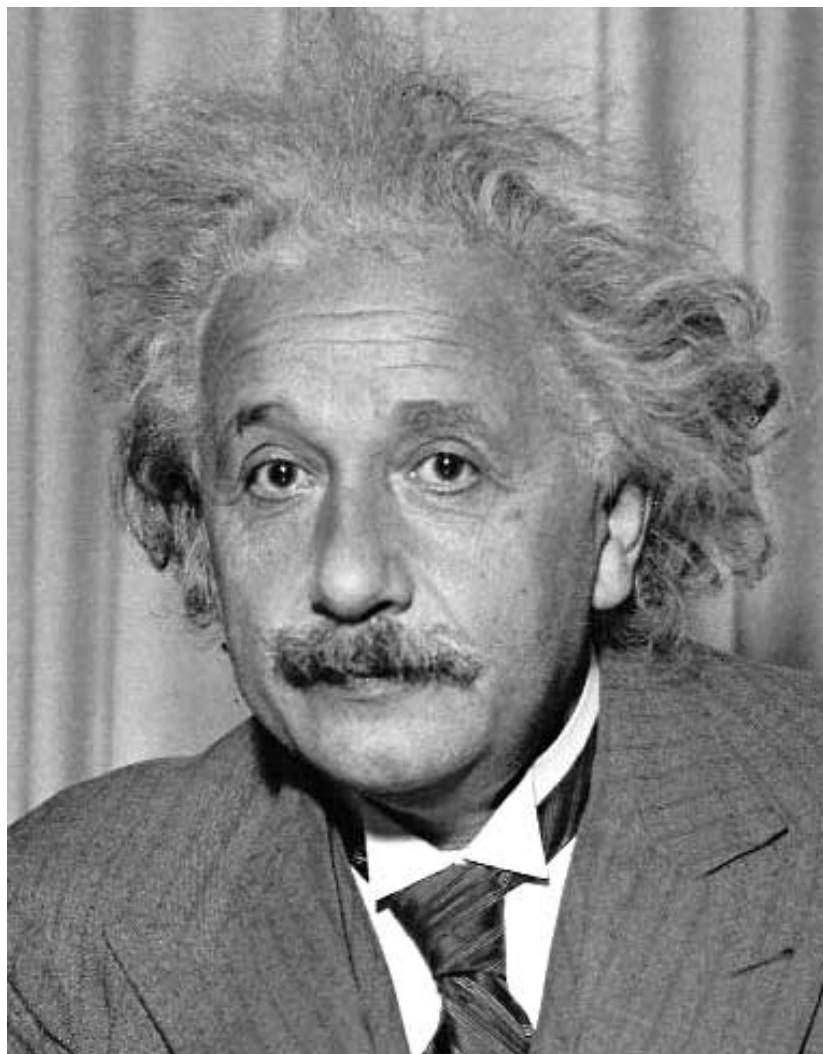


Matching with filters

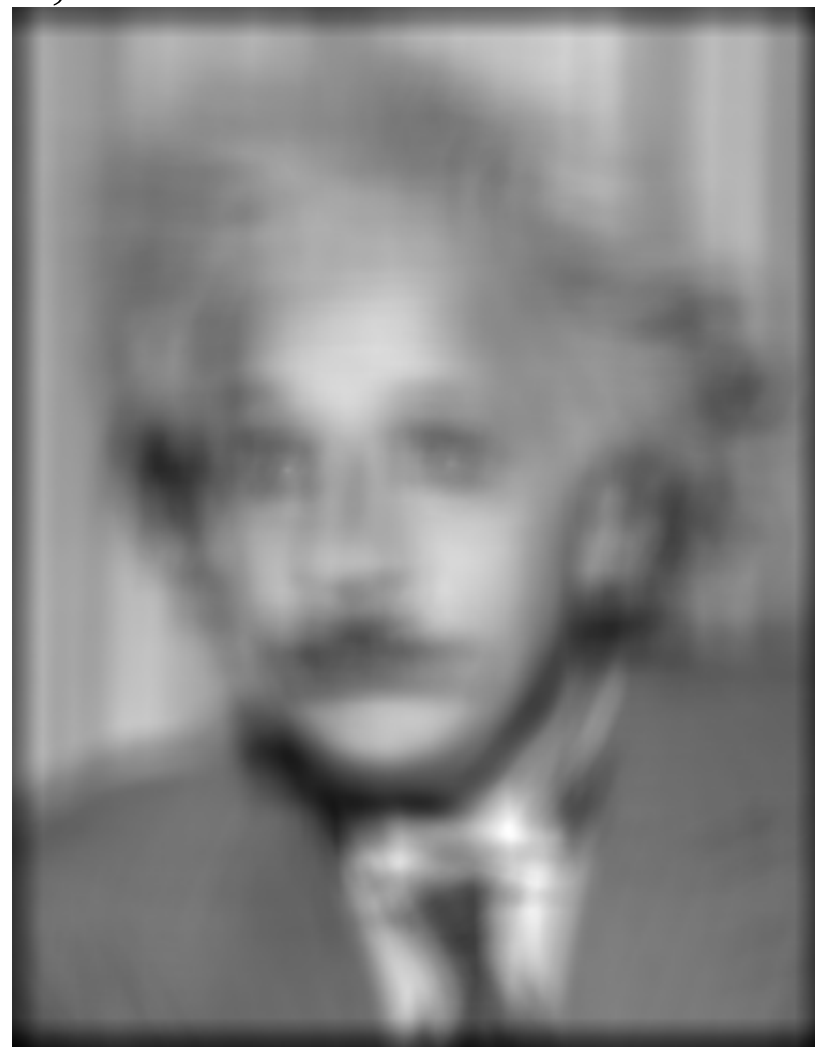
- Goal: find  in image
- Method 0: filter the image with eye patch

$$h[m,n] = \sum_{k,l} g[k,l] f[m+k, n+l]$$

f = image
g = filter




Input



Filtered Image

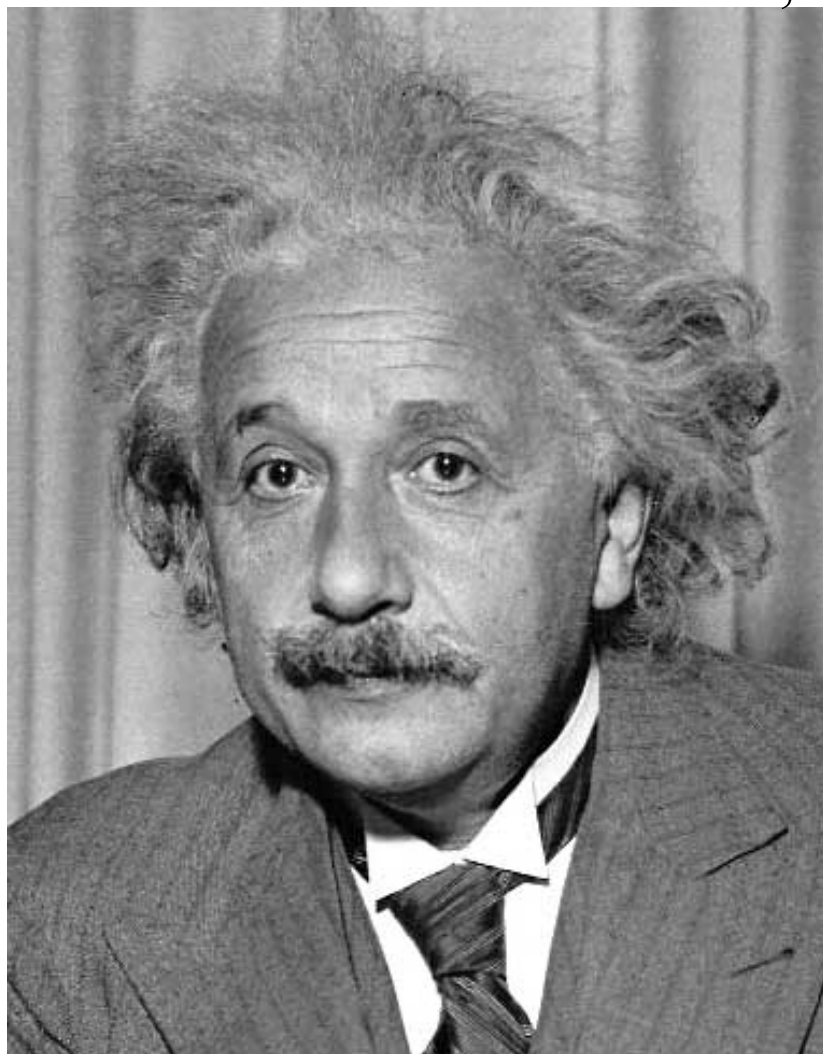
What went wrong?

Matching with filters

- Goal: find  in image
- Method 1: filter the image with zero-mean eye

$$h[m,n] = \sum_{k,l} (f[k,l] - \bar{f}) \underbrace{(g[m+k, n+l])}_{\text{mean of } f}$$

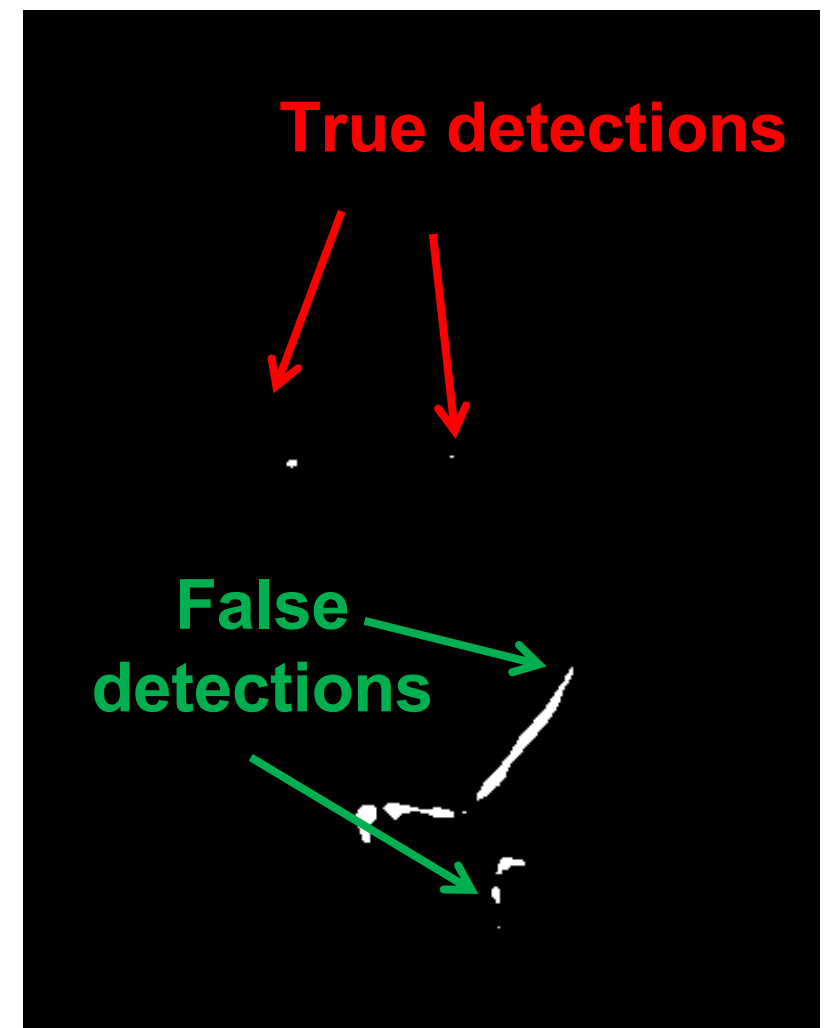
f = image
g = filter



Input




Filtered Image (scaled)



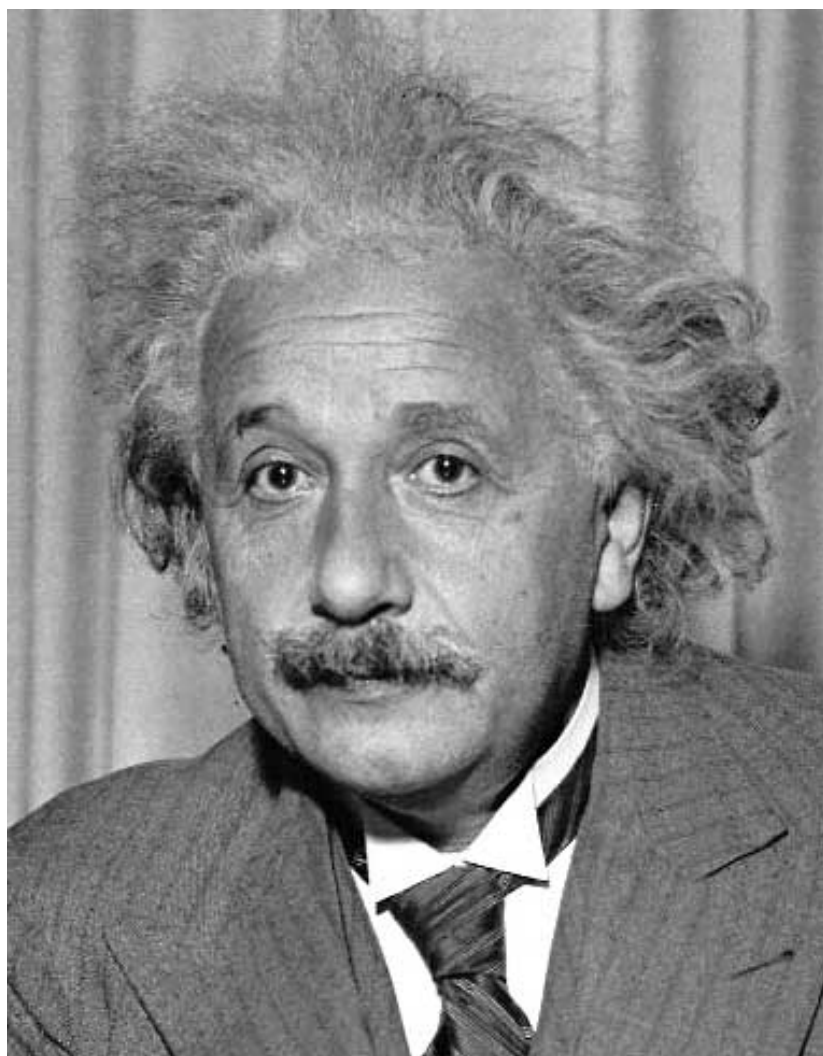
Thresholded Image

Matching with filters

- Goal: find  in image
- Method 2: SSD (Sum Square Difference)

$$h[m,n] = \sum_{k,l} (g[k,l] - f[m+k,n+l])^2$$

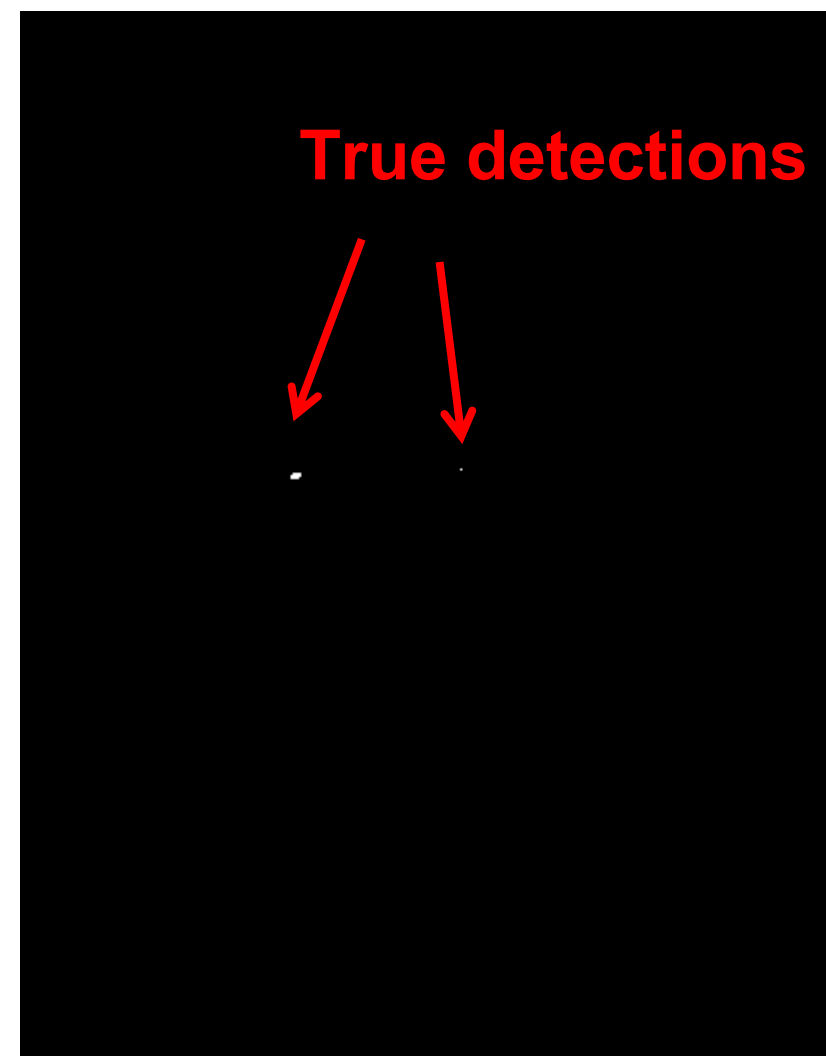
f = image
g = filter



Input



1- sqrt(SSD)



Thresholded Image

Matching with filters

$$h[m,n] = \sum_{k,l} (g[k,l] - f[m+k, n+l])^2$$

f = image
g = filter

- Can SSD be implemented with linear filters?

Matching with filters

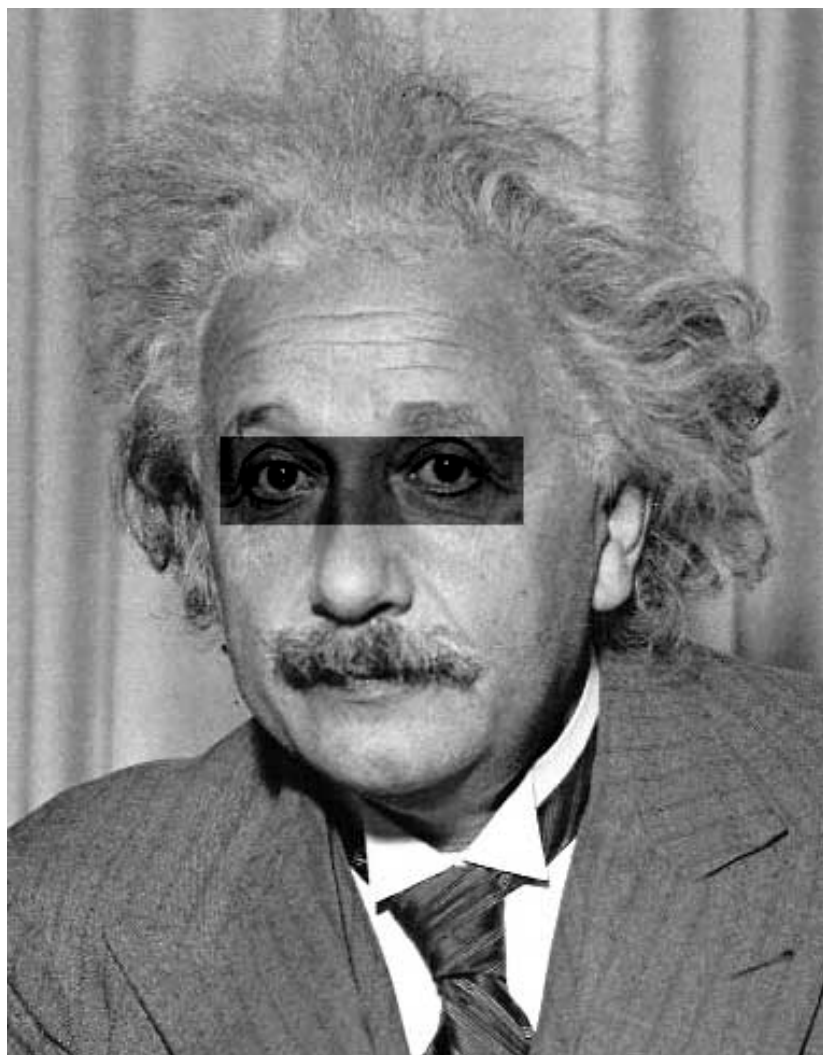
- Goal: find  in image

**What's the potential
downside of SSD?**

- Method 2: SSD (Sum Square Difference)

$$h[m,n] = \sum_{k,l} (g[k,l] - f[m+k,n+l])^2$$

f = image
g = filter



Input



1- sqrt(SSD)

Matching with filters

- Goal: find  in image


- Method 2: Normalized Cross-Correlation f = image
g = filter

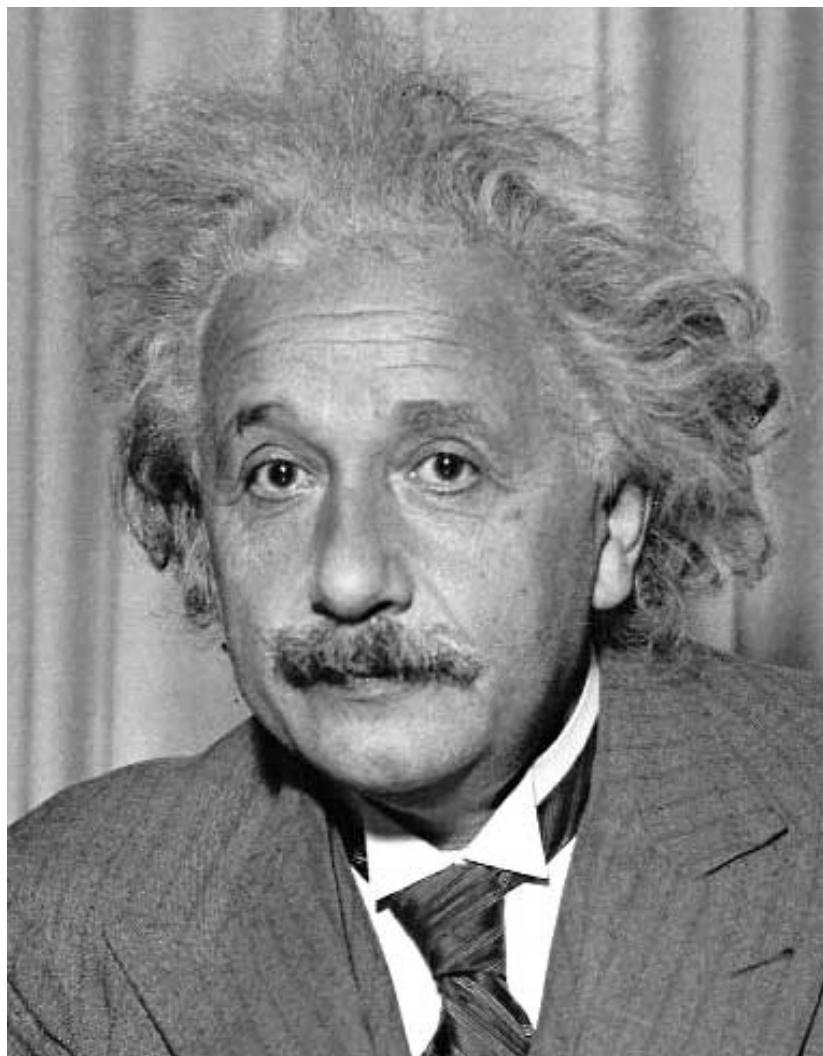
$$h[m, n] = \frac{\sum_{k,l} (g[k, l] - \bar{g})(f[m + k, n + l] - \bar{f}_{m,n})}{\left(\sum_{k,l} (g[k, l] - \bar{g})^2 \sum_{k,l} (f[m + k, n + l] - \bar{f}_{m,n})^2 \right)^{0.5}}$$

mean template mean image patch

↓ ↓

Matching with filters

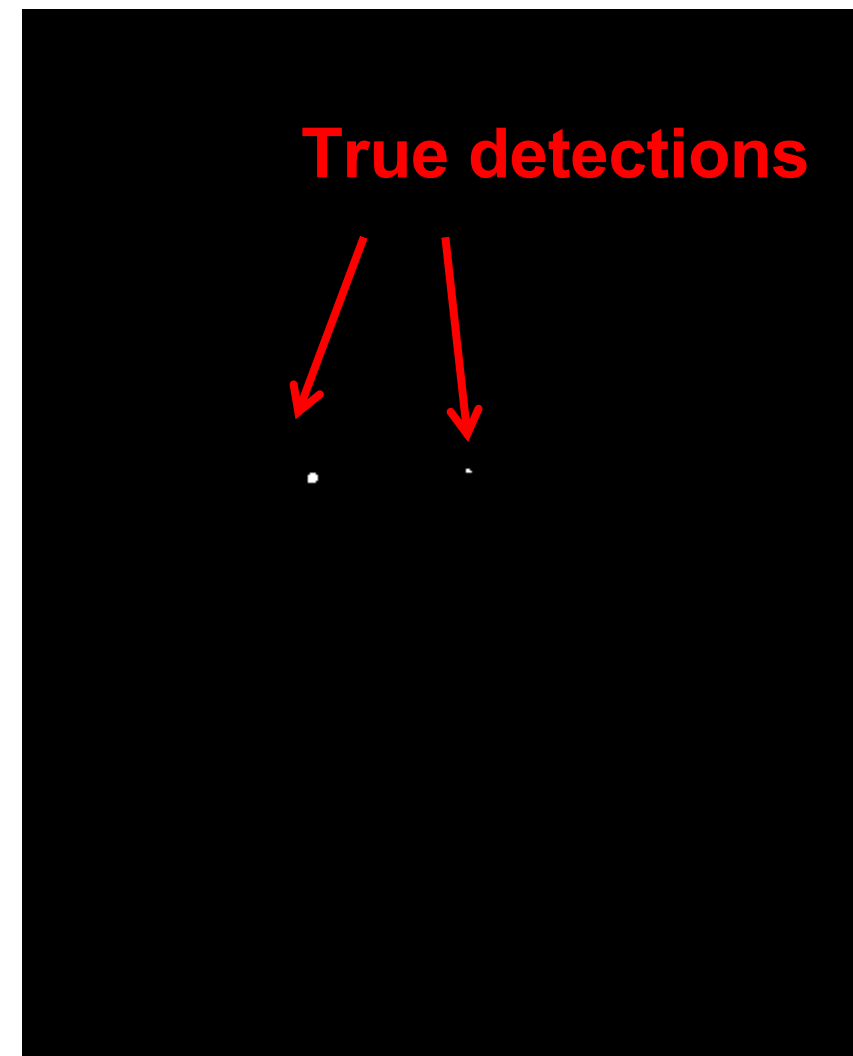
- Goal: find  in image
- Method 2: Normalized Cross-Correlation



Input




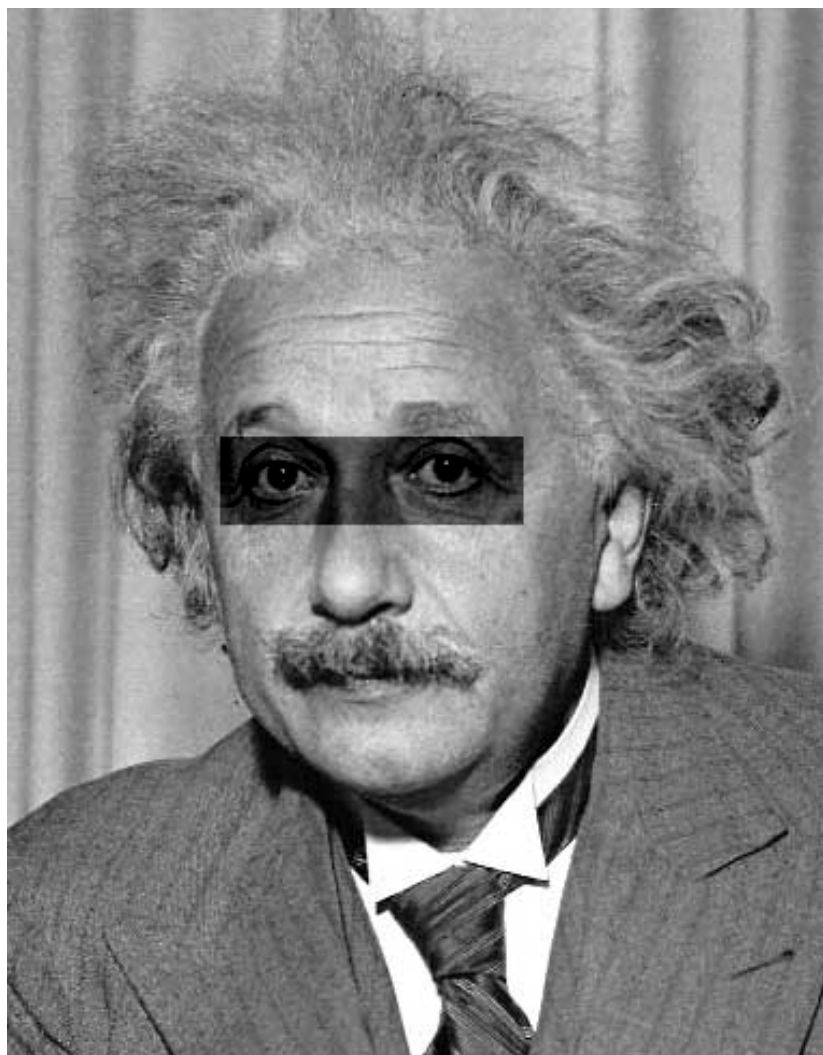
Normalized Cross-Correlation



Thresholded Image

Matching with filters

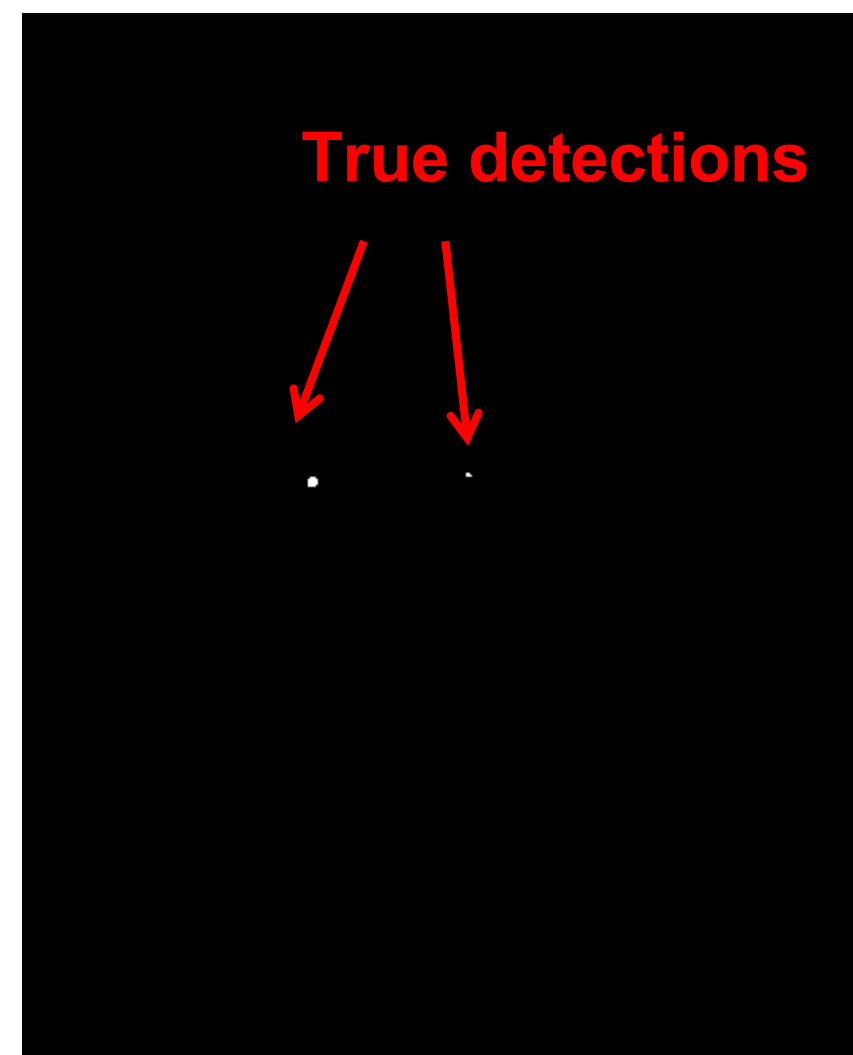
- Goal: find  in image
- Method 2: Normalized Cross-Correlation



Input



Normalized X-Correlation



Thresholded Image

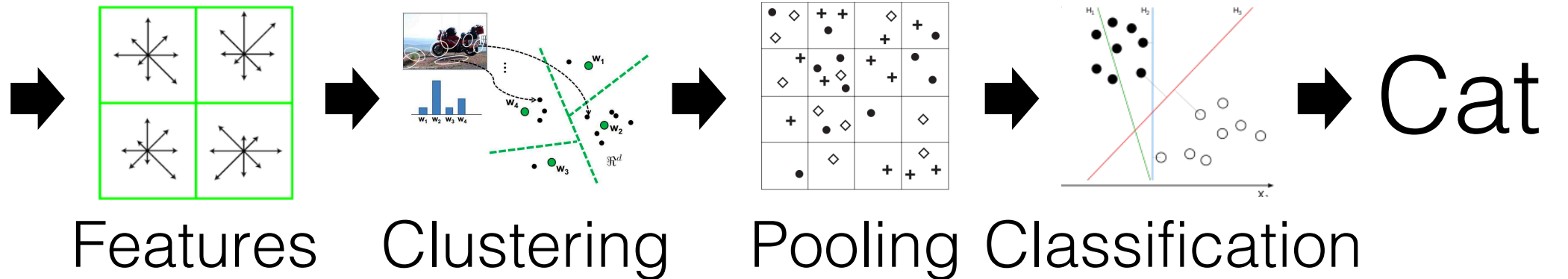
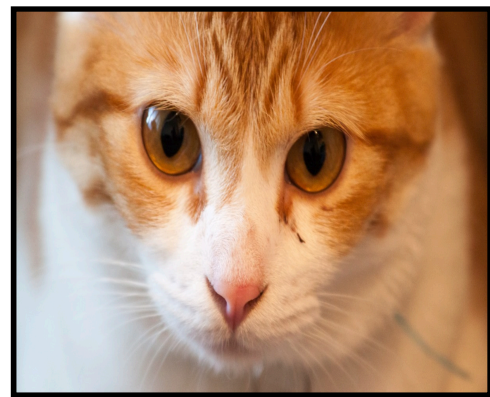
Q: What is the best method to use?

- Answer: Depends
- Zero-mean filter: fastest but not a great matcher
- SSD: next fastest, sensitive to overall intensity
- Normalized cross-correlation: slowest, invariant to local average intensity and contrast

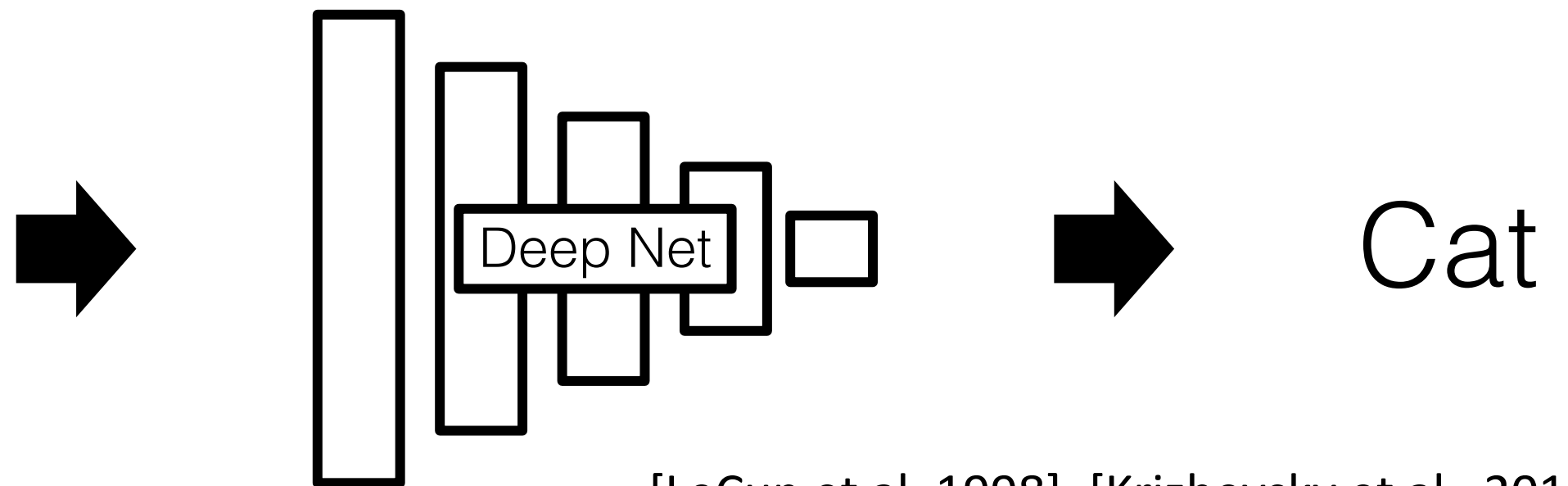
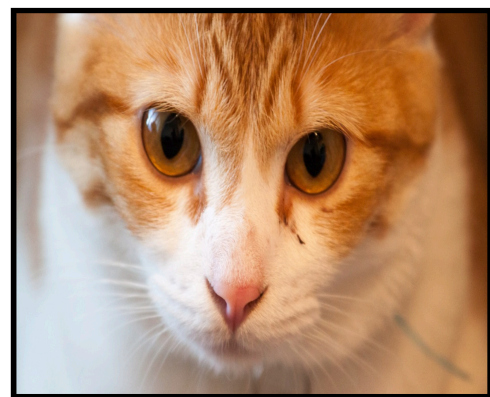
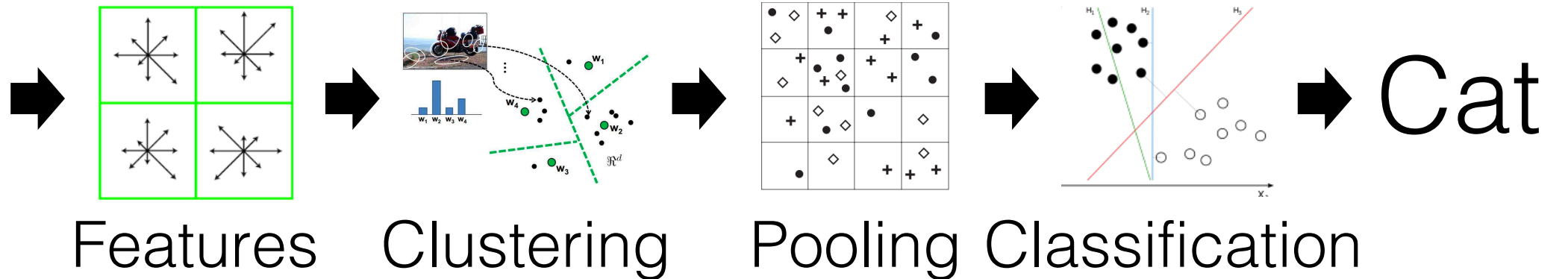
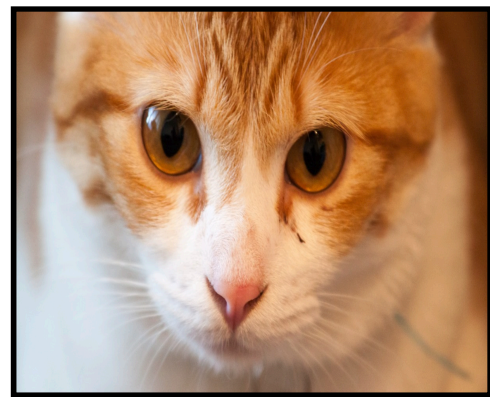
Review

(CNN for Image Synthesis)

Computer Vision before 2012

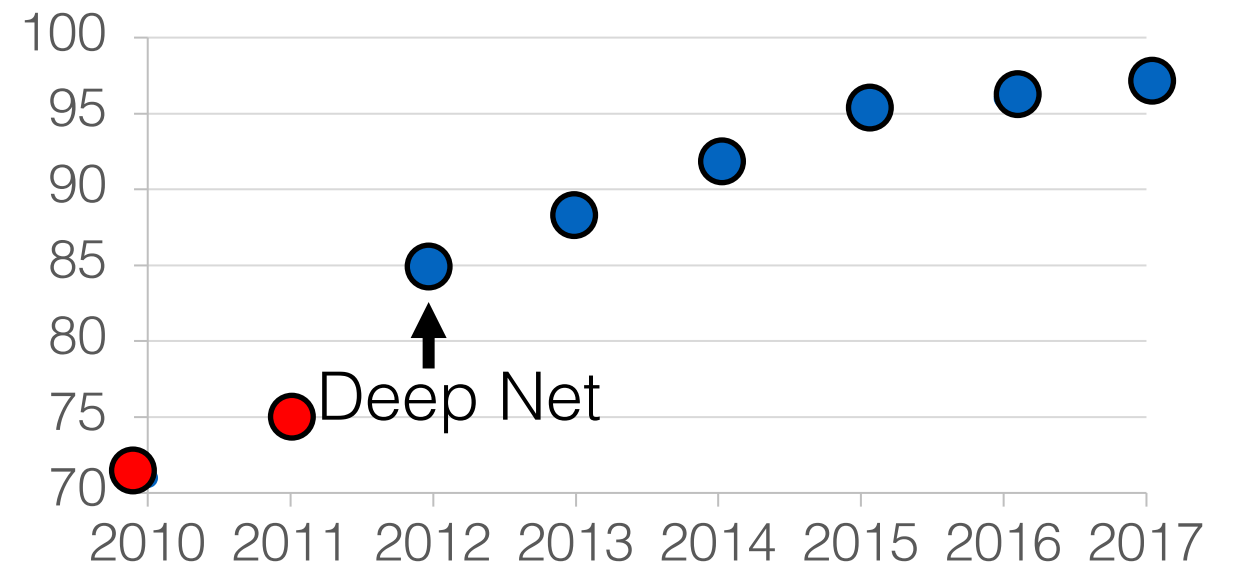


Computer Vision Now

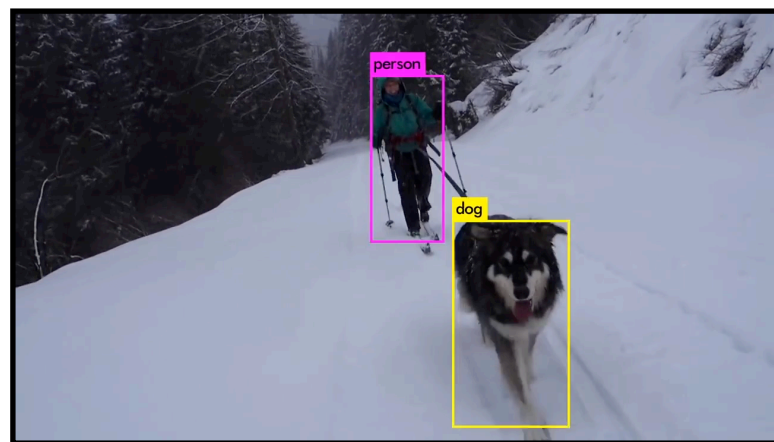


[LeCun et al, 1998], [Krizhevsky et al, 2012]

Deep Learning for Computer Vision

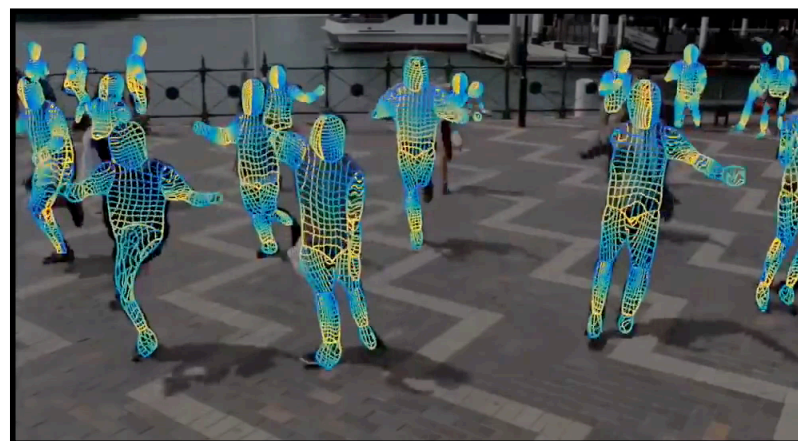


Top 5 **accuracy** on ImageNet benchmark



[Redmon et al., 2018]

Object detection



[Güler et al., 2018]

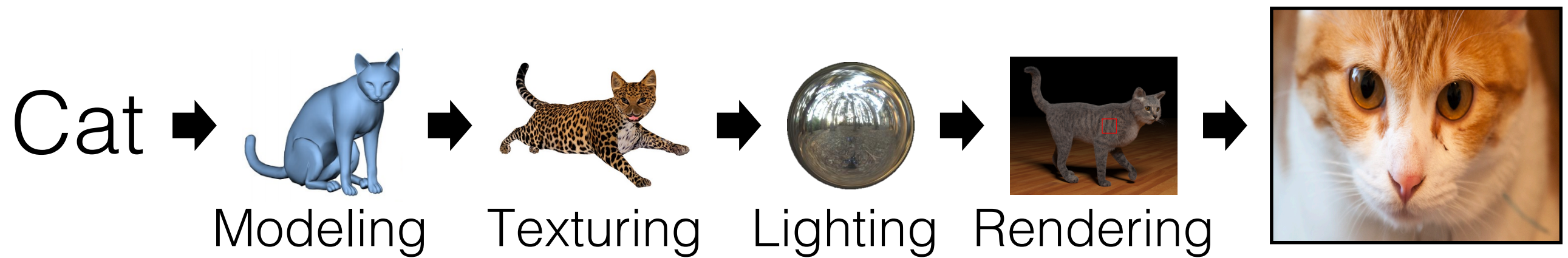
Human understanding



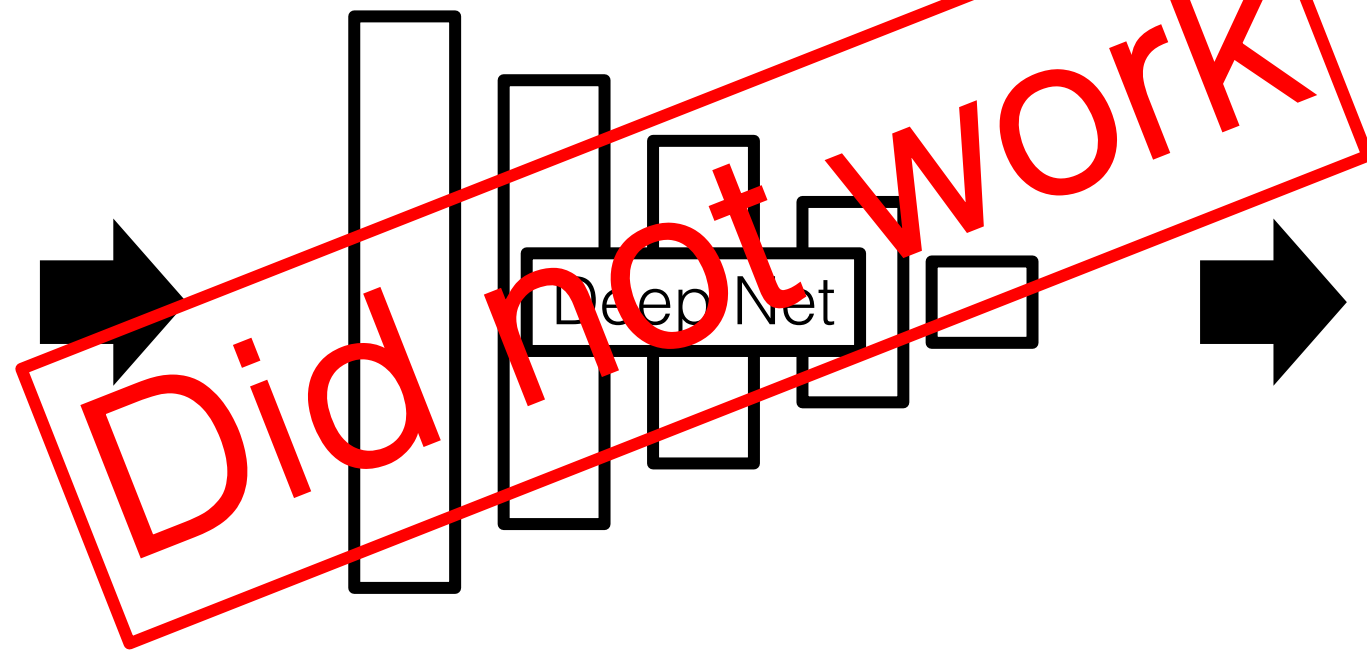
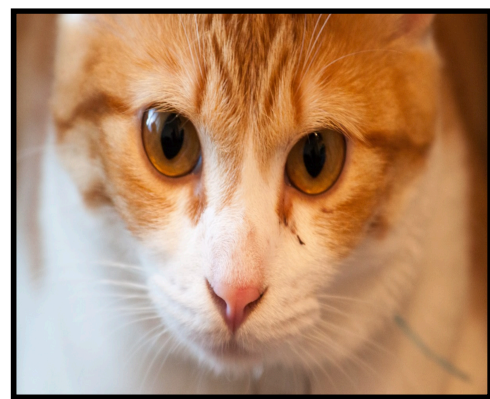
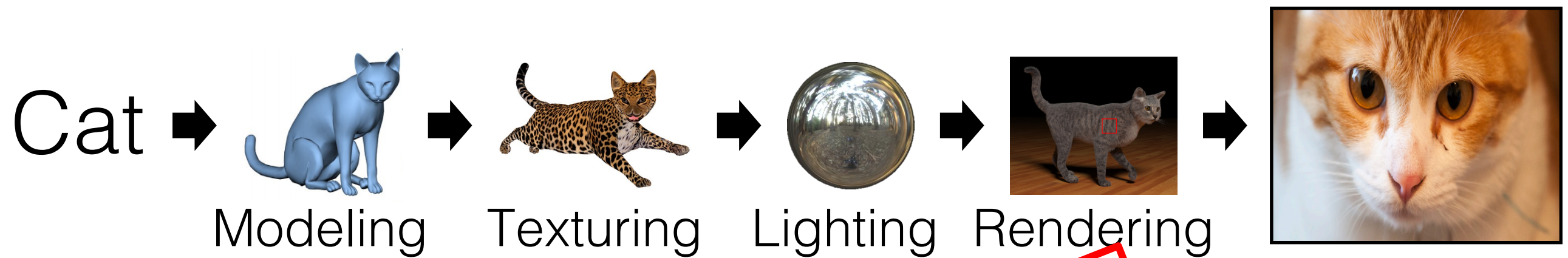
[Zhao et al., 2017]

Autonomous driving

Can Deep Learning Help Graphics?

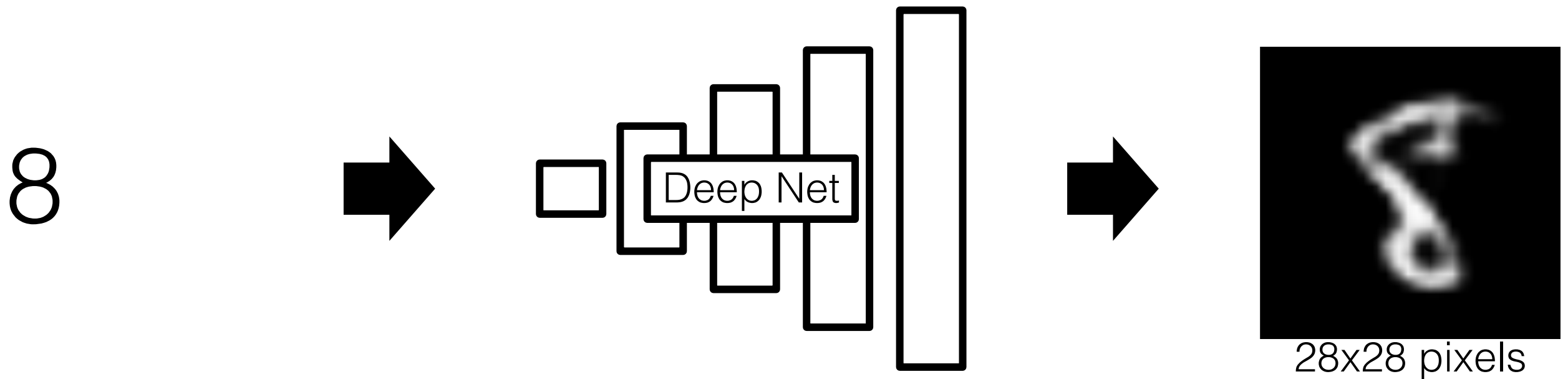
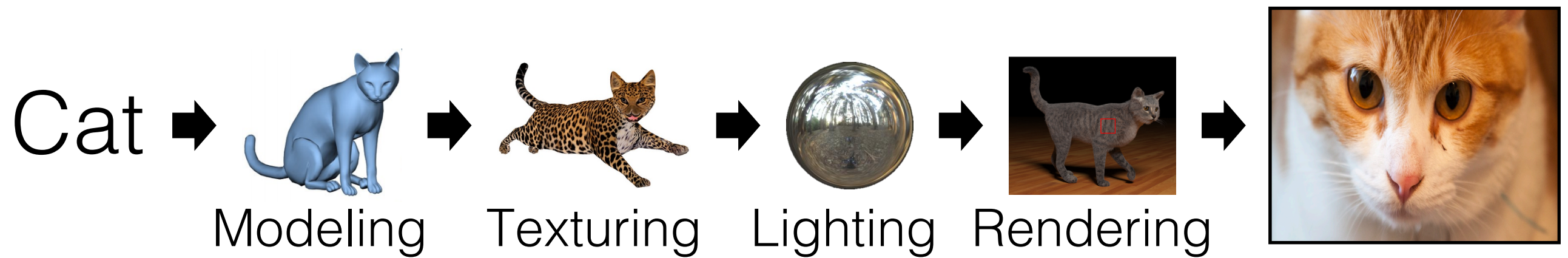


Can Deep Learning Help Graphics?



Cat

Generating images is hard!



Simple L2 regression doesn't work 😞

Input



Output



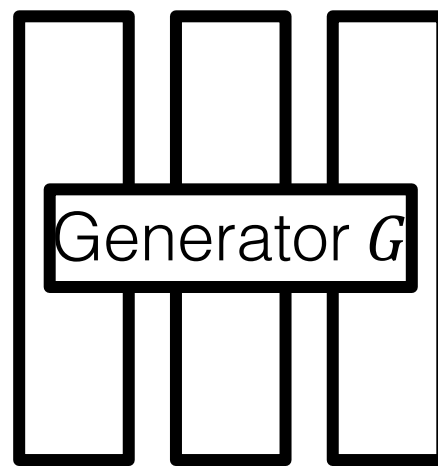
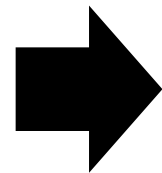
Ground truth



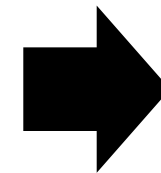
Loss functions for Image Synthesis



Input x



Learnable rendering



Output Image $G(x)$

What is a good objective \mathcal{L} ?

- Capture realism
- Task-agnostic
- Data-dependent

Problem Statement

Loss function

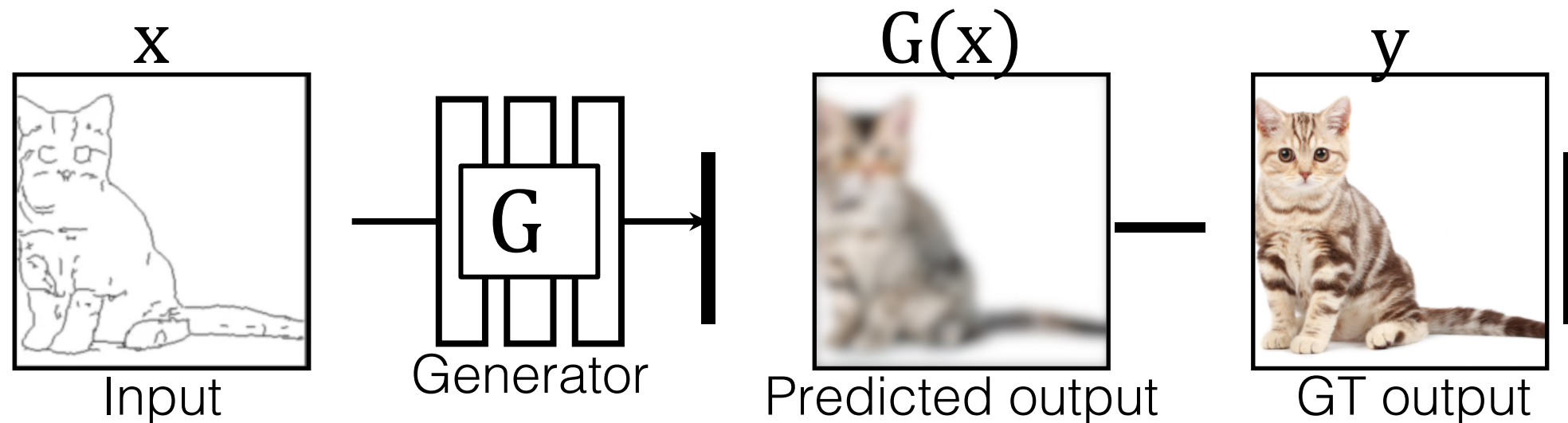
$$\arg \min_G \mathcal{L}(G(x), y)$$

Generator

Input

Output image

Designing Loss Functions

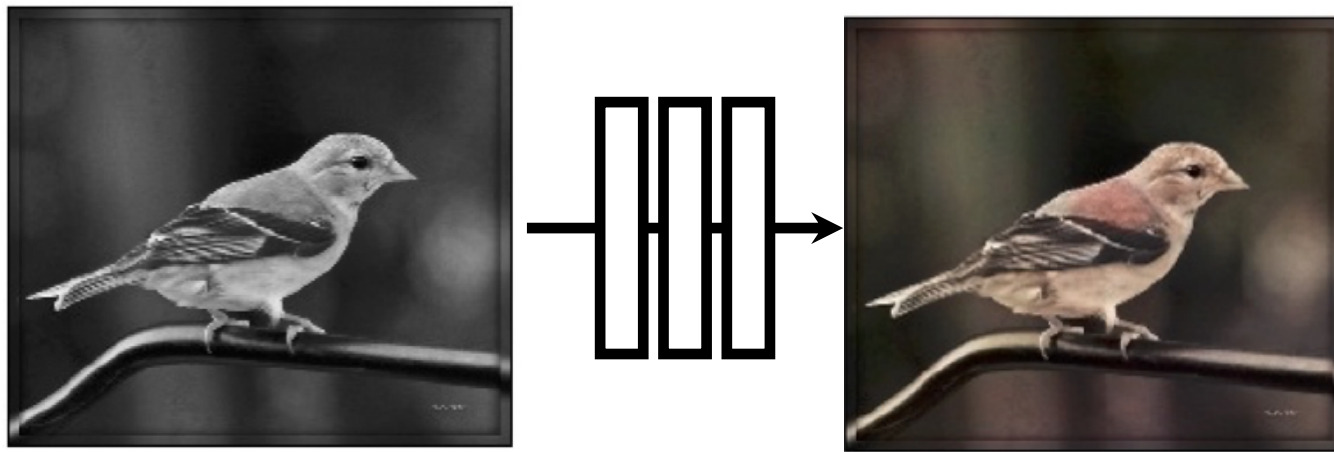


L2 regression

$$\arg \min_G \mathbb{E}_{(x,y)} [||G(x) - y||]$$

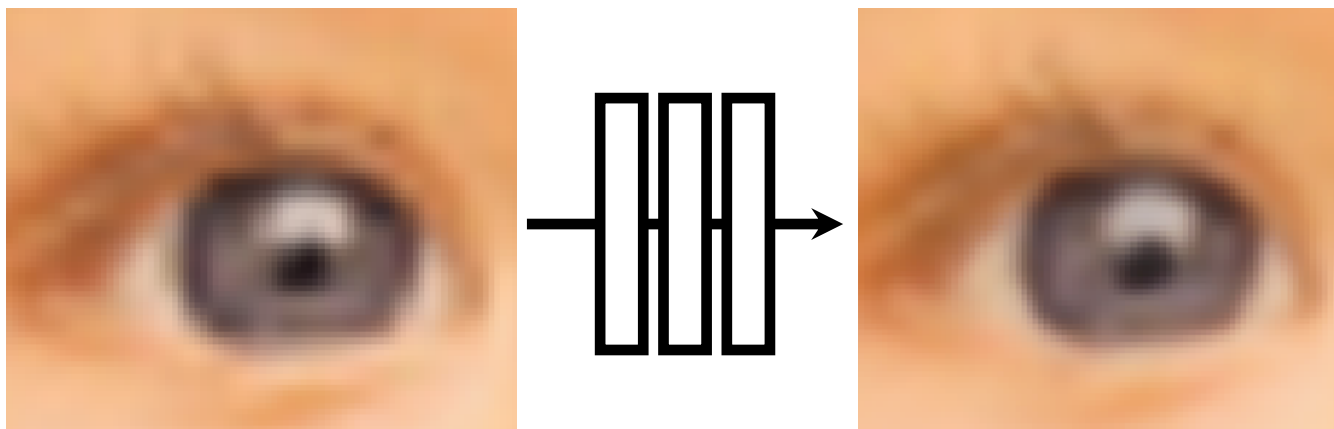
Designing Loss Functions

Image colorization



L2 regression

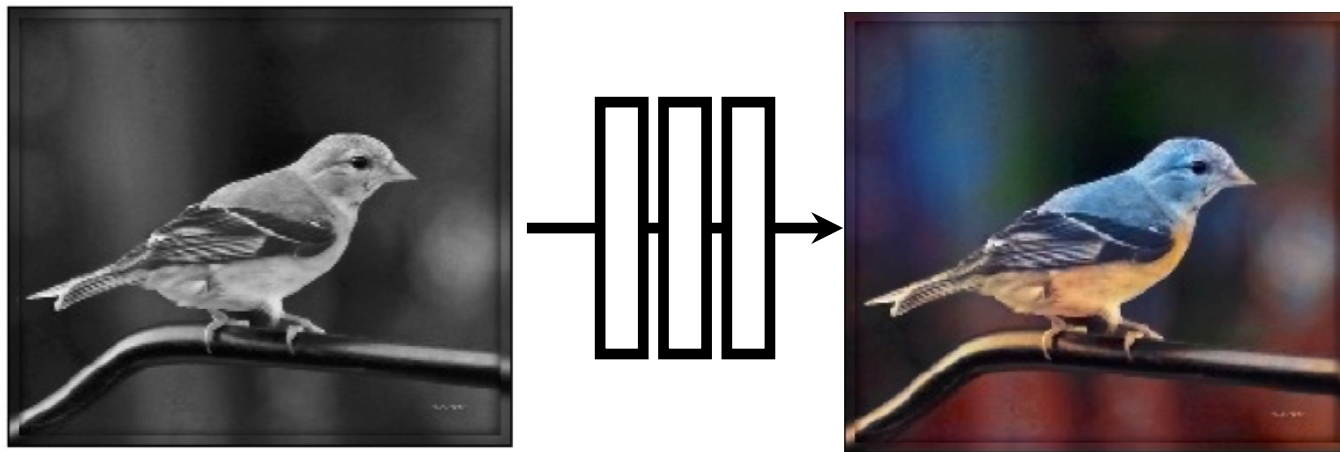
Super-resolution



L2 regression

Designing Loss Functions

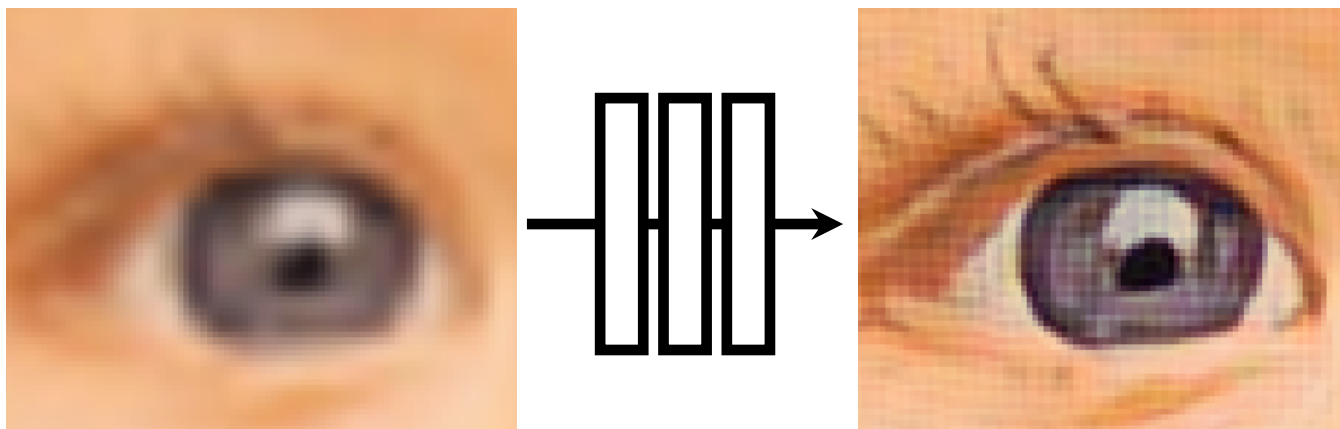
Image colorization



[Zhang et al. 2016]

Classification Loss:
Cross entropy objective,
with colorfulness term

Super-resolution

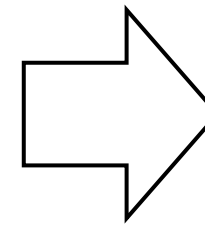


[Gatys et al., 2016], [Johnson et al. 2016]
[Dosovitskiy and Brox. 2016]

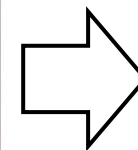
Feature/Perceptual loss
Deep feature covariance
matching objective

“Perceptual Loss”

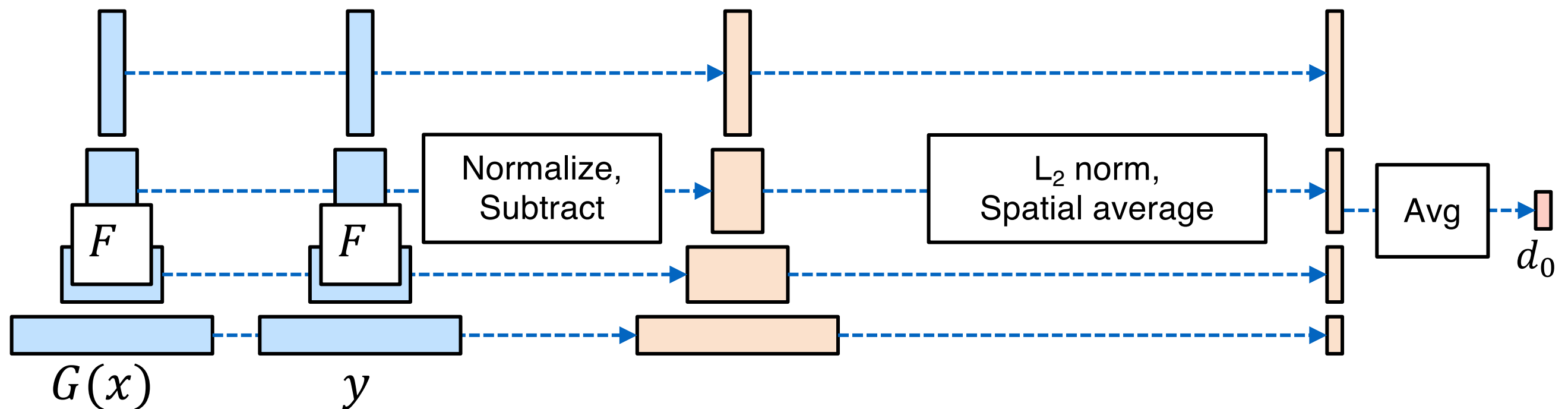
Gatys et al. In CVPR, 2016.
Johnson et al. In ECCV, 2016.
Dosovitskiy and Brox. In NIPS, 2016.



Chen and Koltun. In ICCV, 2017.



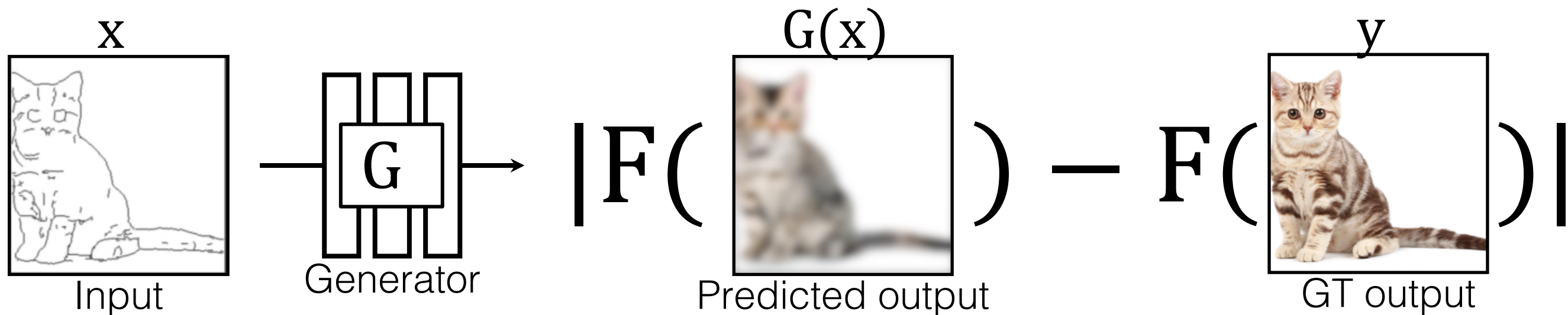
CNNs as a Perceptual Metric



(1) How well do “perceptual losses” describe perception?

c.f. Gatys et al. CVPR 2016. Johnson et al. ECCV 2016. Dosovitskiy and Brox. NIPS 2016.

CNNs as a Perceptual Metric



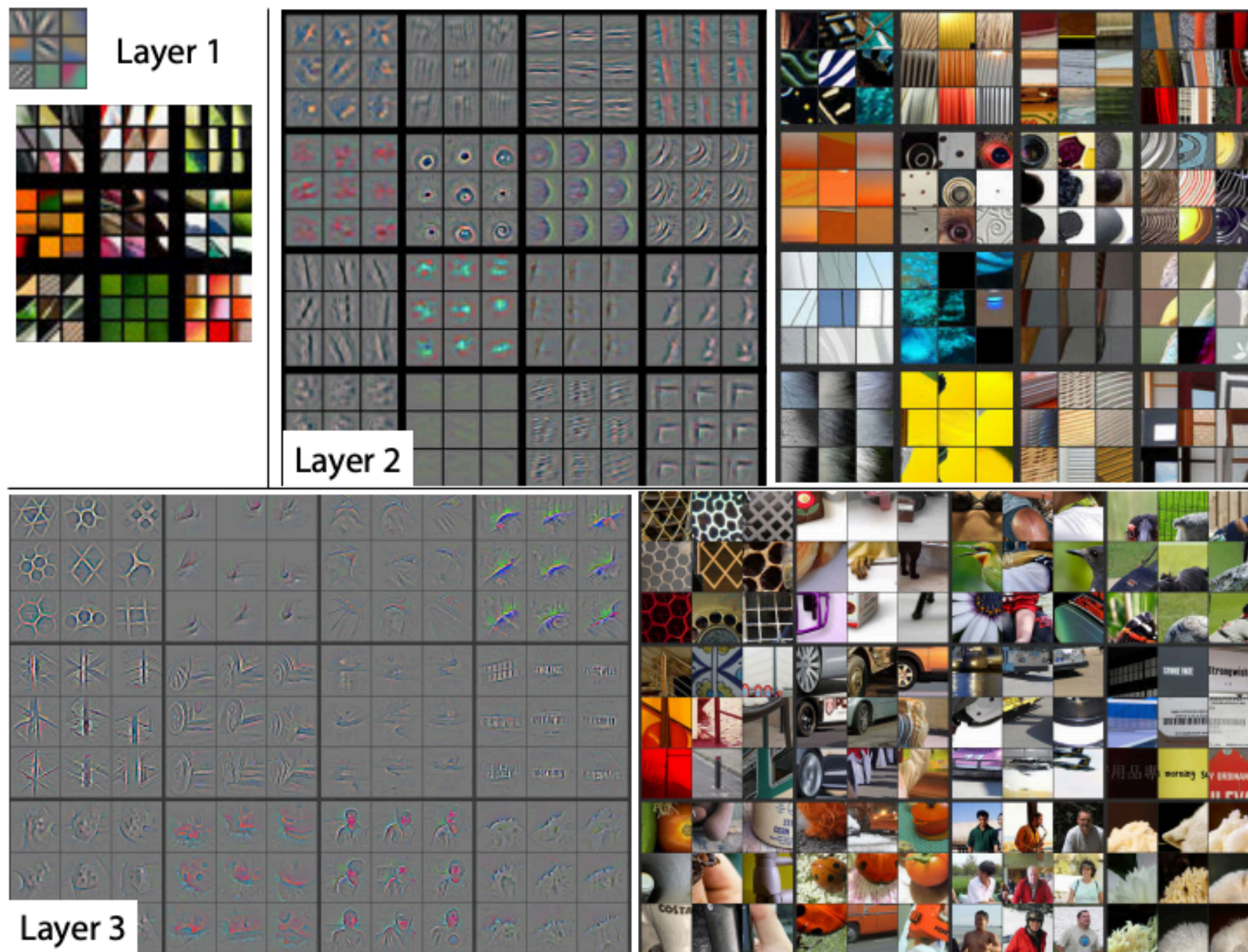
F is a deep network (e.g., ImageNet classifier)

Perceptual Loss

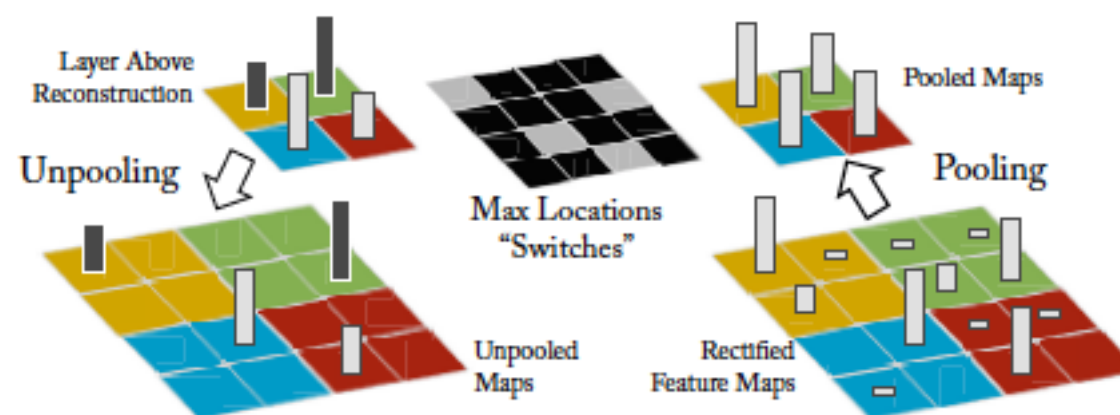
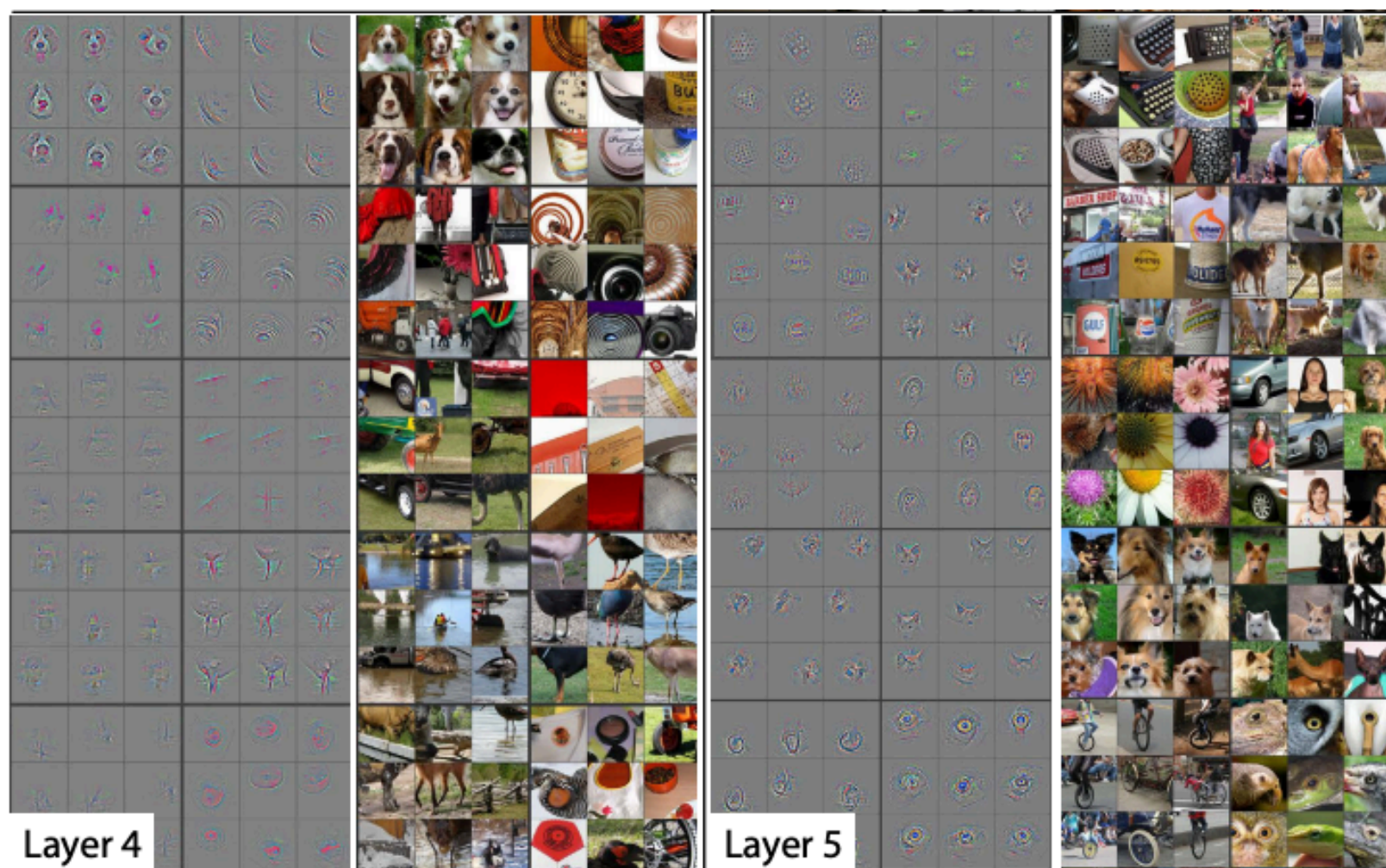
$$\arg \min_G \mathbb{E}_{(x,y)} \sum_{i=1}^N \overset{\text{weight}}{\lambda_i} \frac{1}{M_i} \left\| \overset{\text{(i)-th layer}}{F^{(i)}}(G(x)) - F^{(i)}(y) \right\|_2^2$$

The number of elements in the (i)-th layer

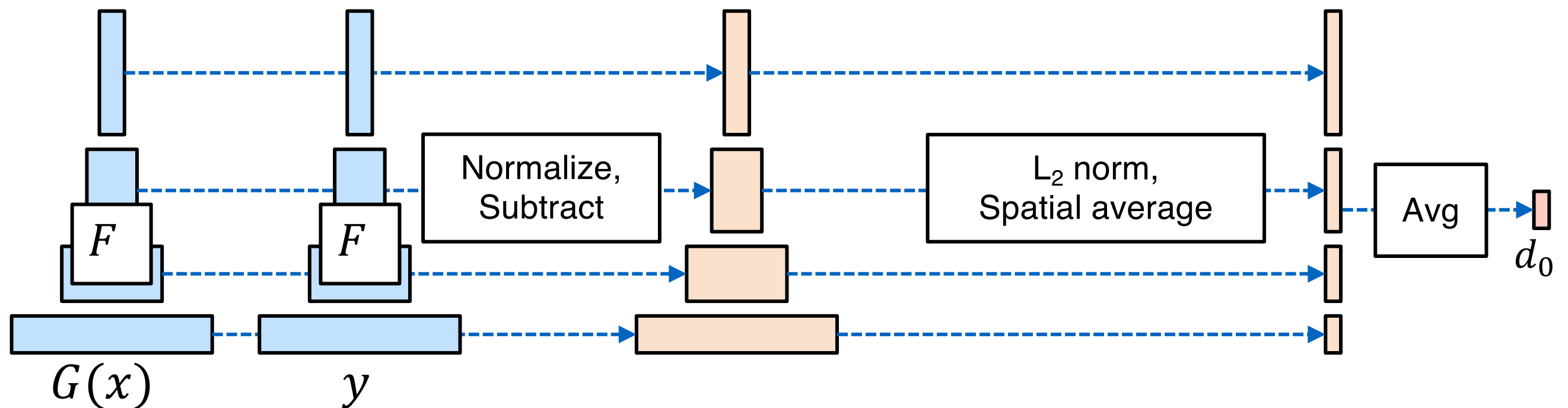
What has a CNN Learned?



What has a CNN Learned?



CNNs as a Perceptual Metric



Perceptual Loss

$$\arg \min_G \mathbb{E}_{(x,y)} \sum_{i=1}^N \overset{\text{weight}}{\lambda_i} \frac{1}{\overset{\text{(i)-th layer}}{M_i}} \left\| F^{(i)}(G(x)) - F^{(i)}(y) \right\|_2^2$$

The number of elements in the (i)-th layer

How Different are these Patches?



Zhang, Isola, Efros, Shechtman, Wang.

The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*, 2018.

Which patch is more similar to the middle?



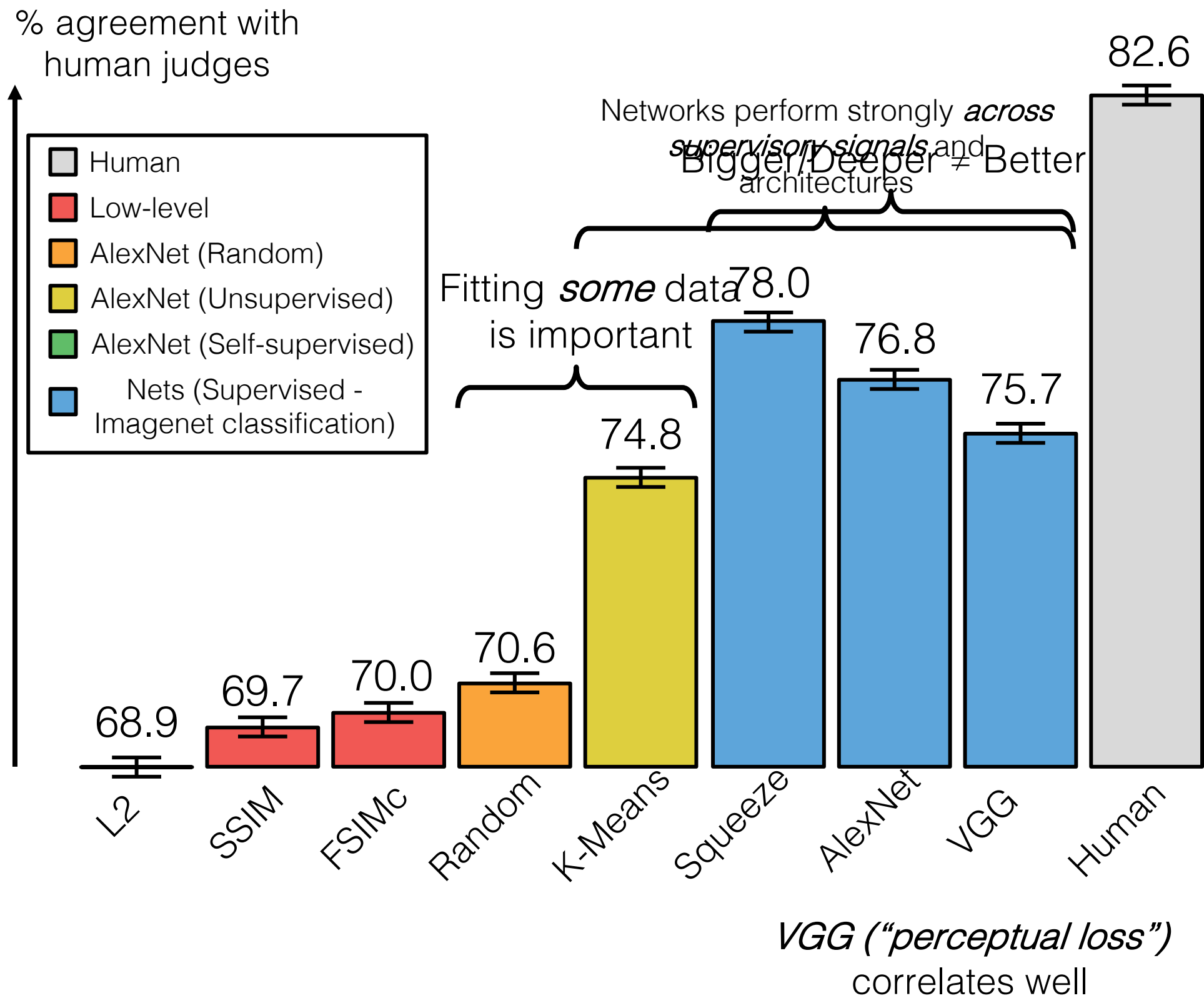
< Type 1 >



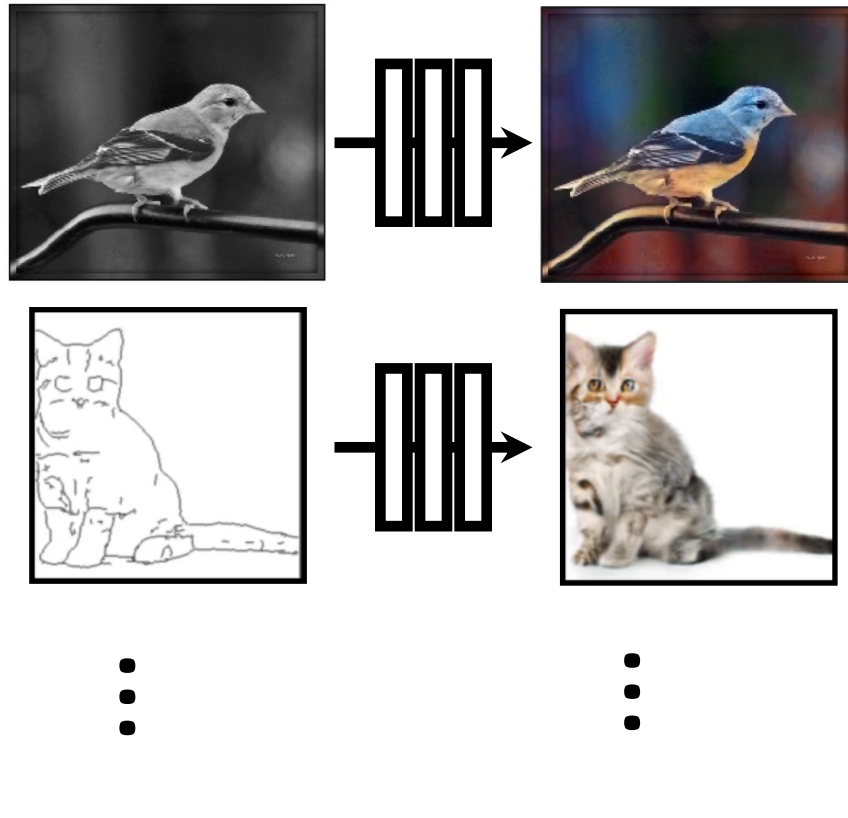
Humans
L2/PSNR
SSIM/FSIMc
Deep Networks?



< Type 2 >

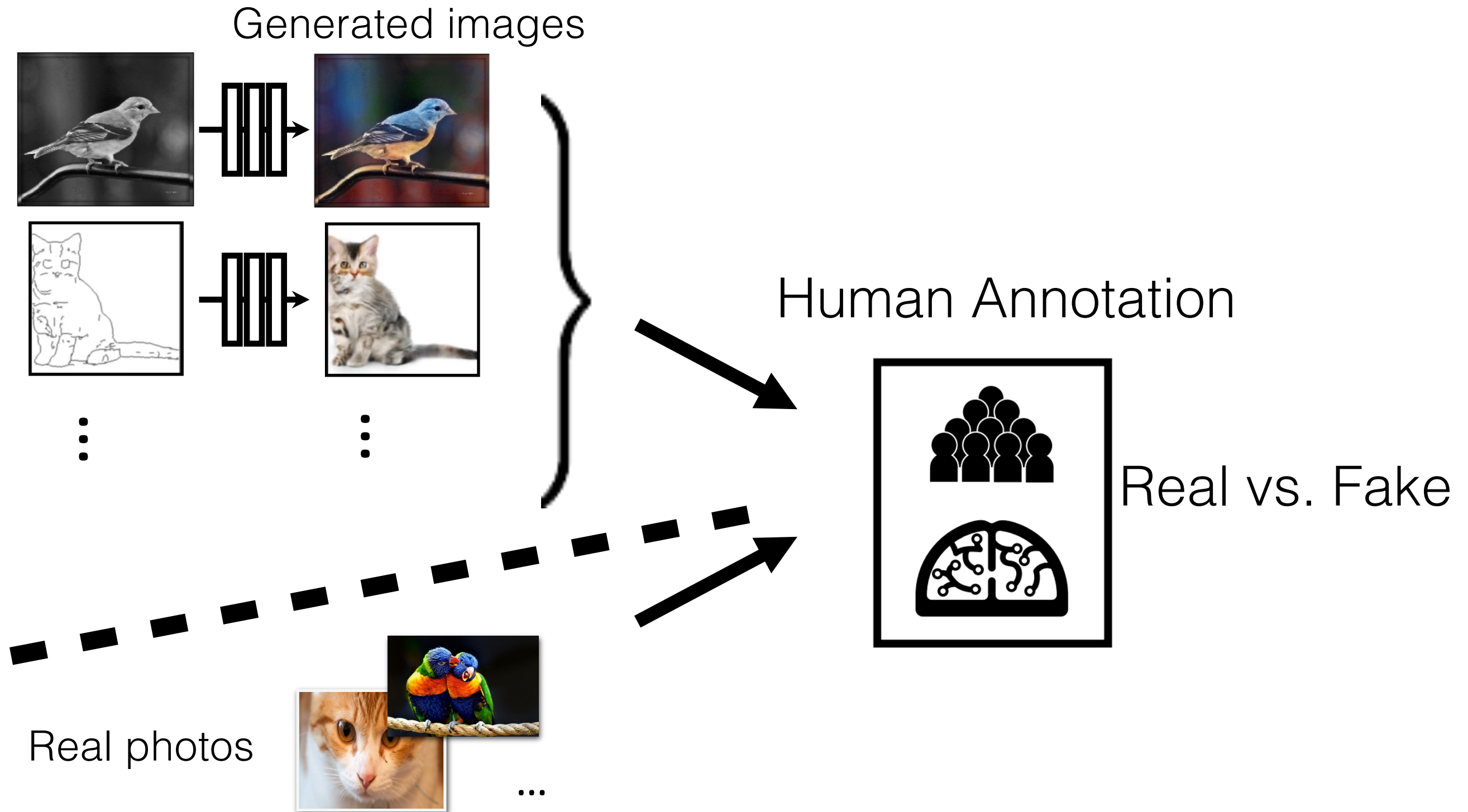


Generated images

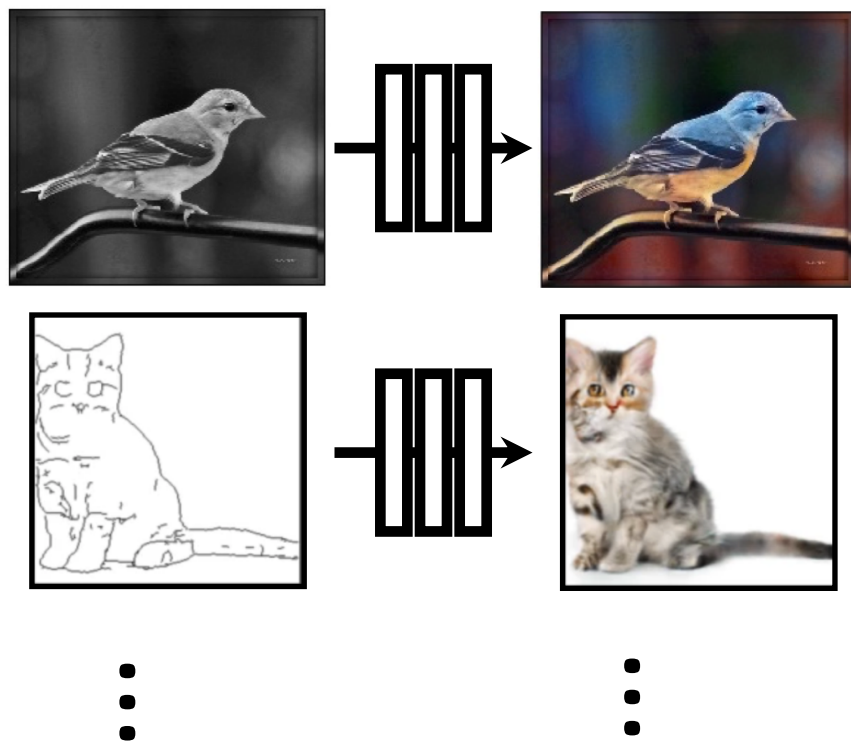


Universal loss?

Learning with Human Perception



Generated images



Generative Adversarial Network (GANs)

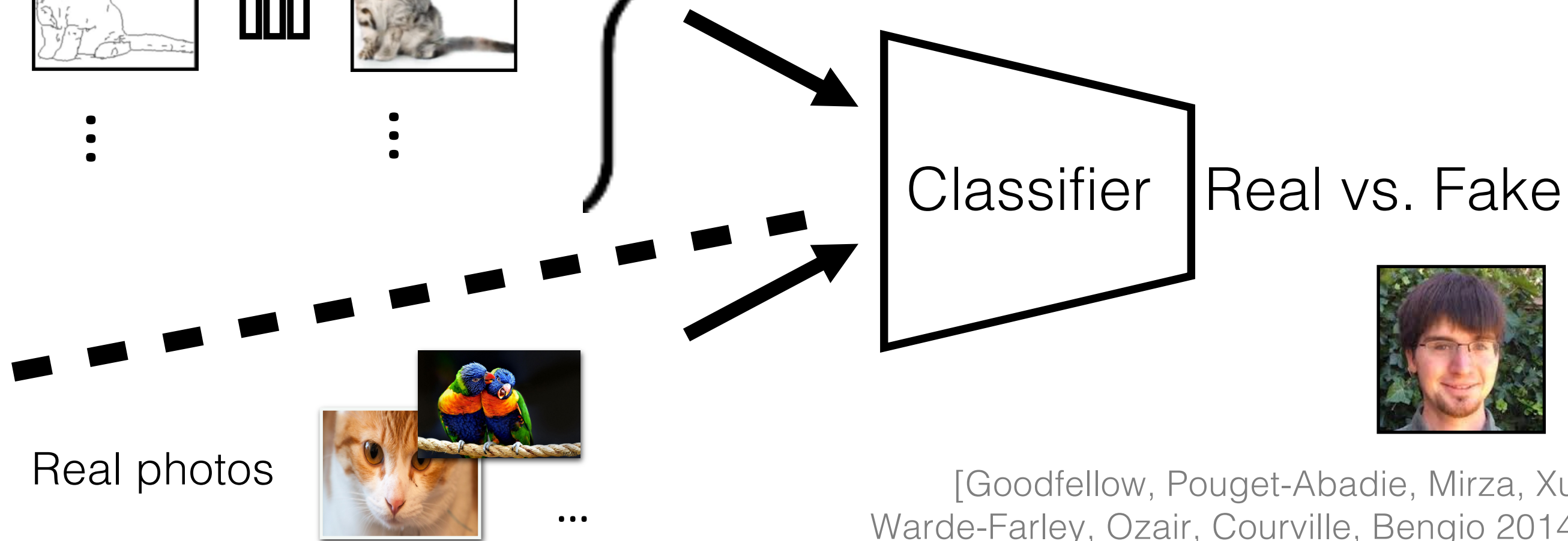
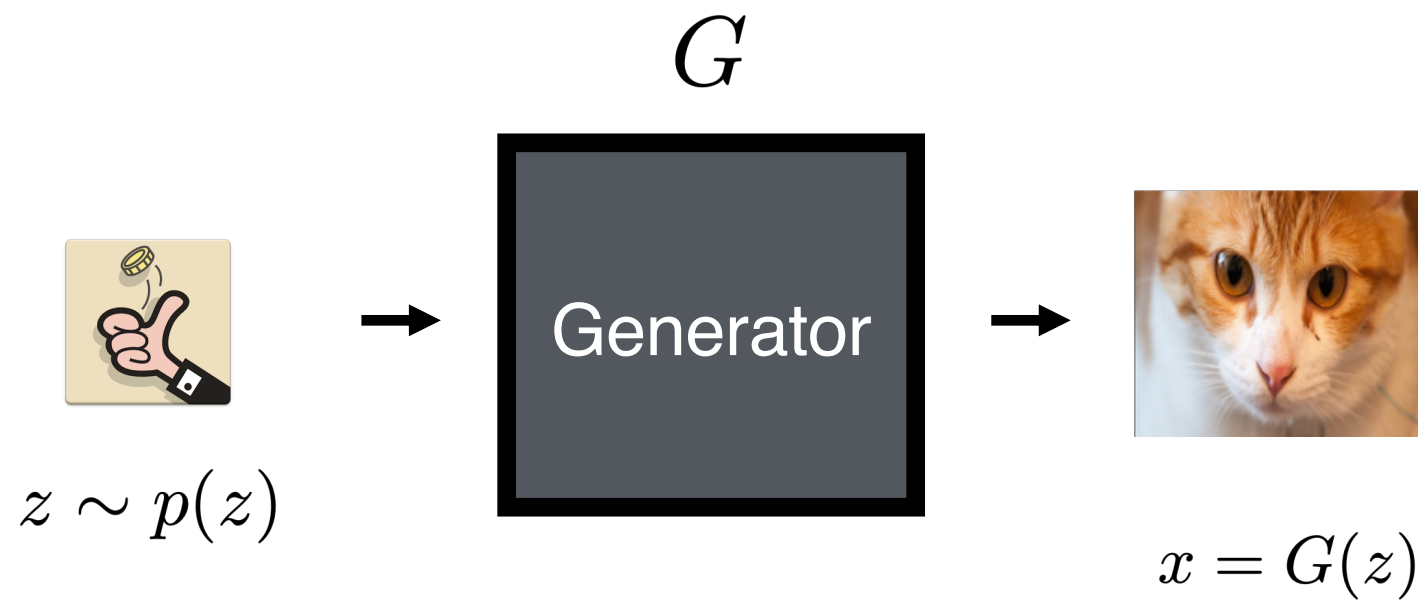


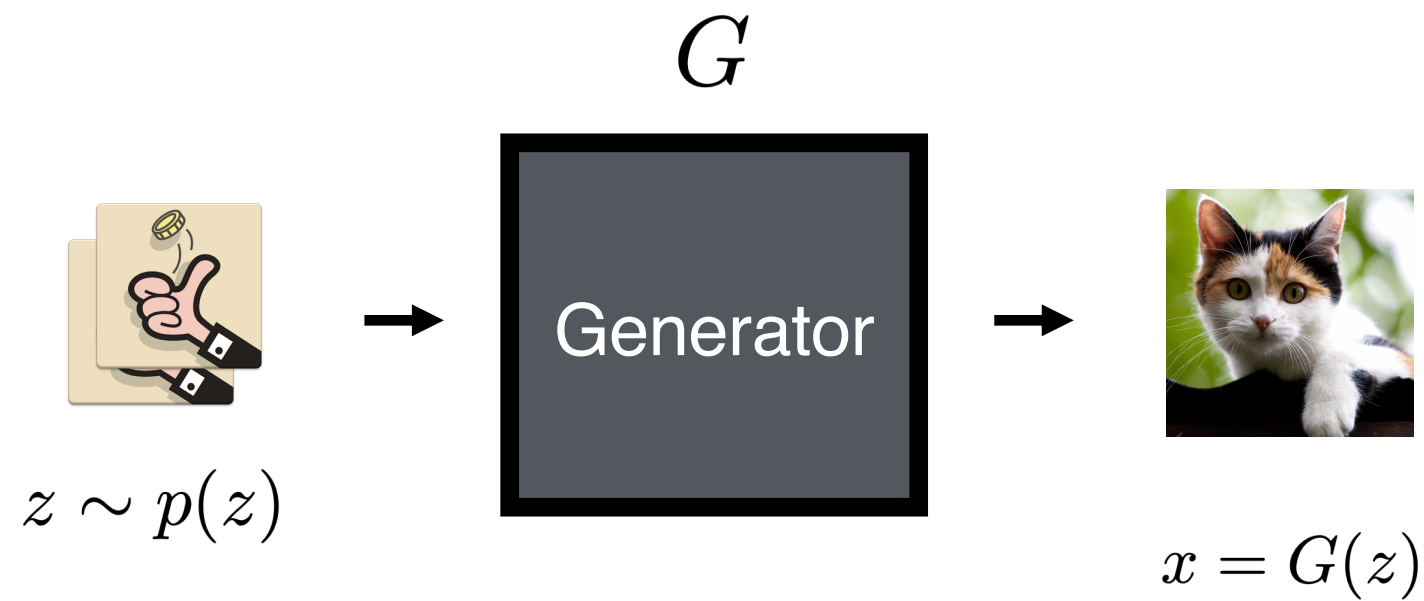
Image synthesis from “noise”



Sampler

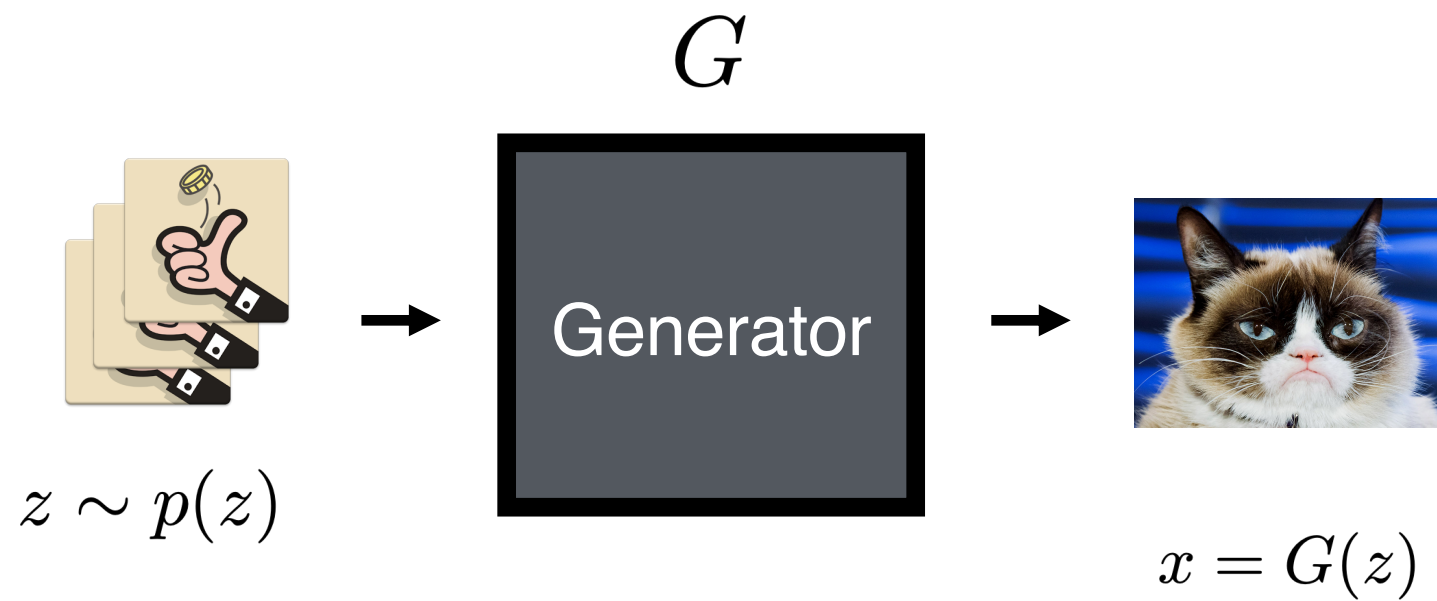
$$G : \mathcal{Z} \rightarrow \mathcal{X}$$
$$z \sim p(z)$$
$$x = G(z)$$

Image synthesis from “noise”



Sampler
 $G : \mathcal{Z} \rightarrow \mathcal{X}$
 $z \sim p(z)$
 $x = G(z)$

Image synthesis from “noise”



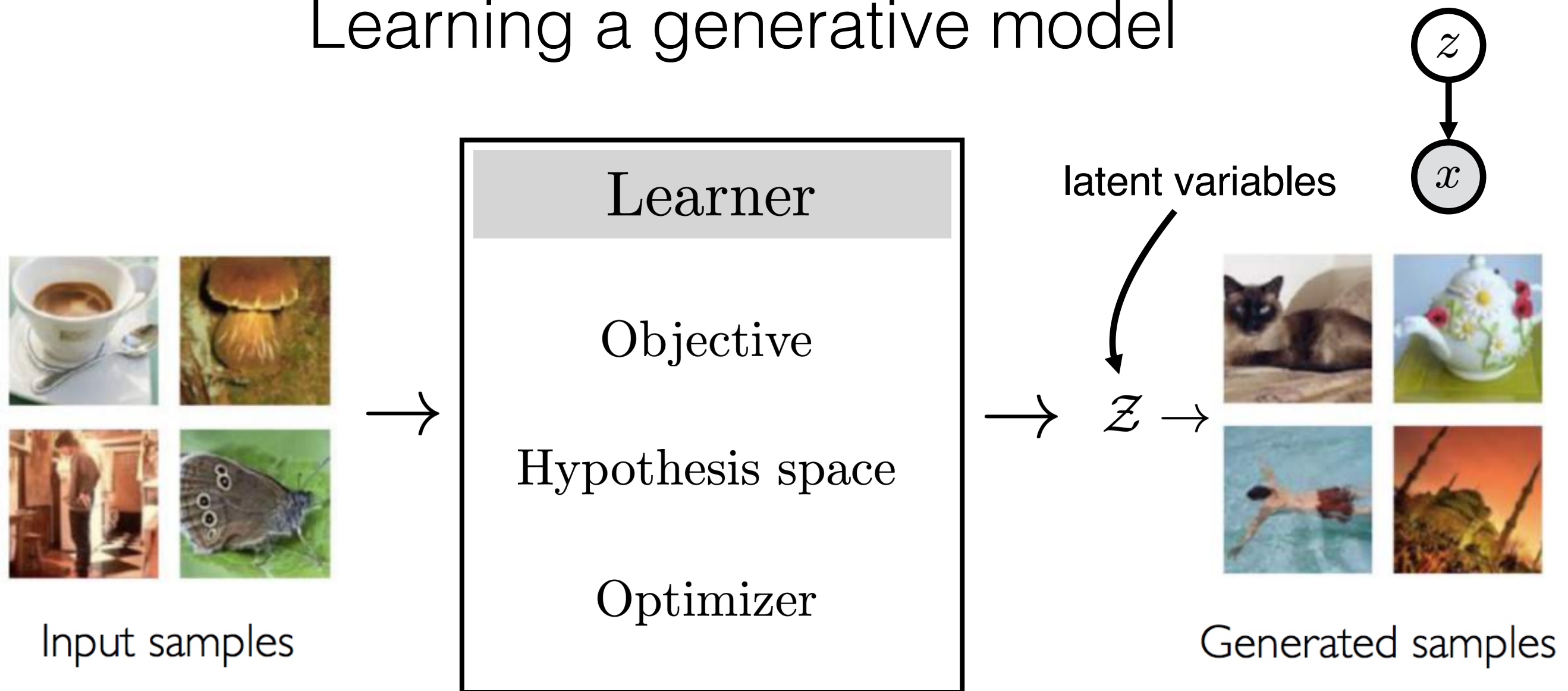
Sampler

$$G : \mathcal{Z} \rightarrow \mathcal{X}$$

$$z \sim p(z)$$

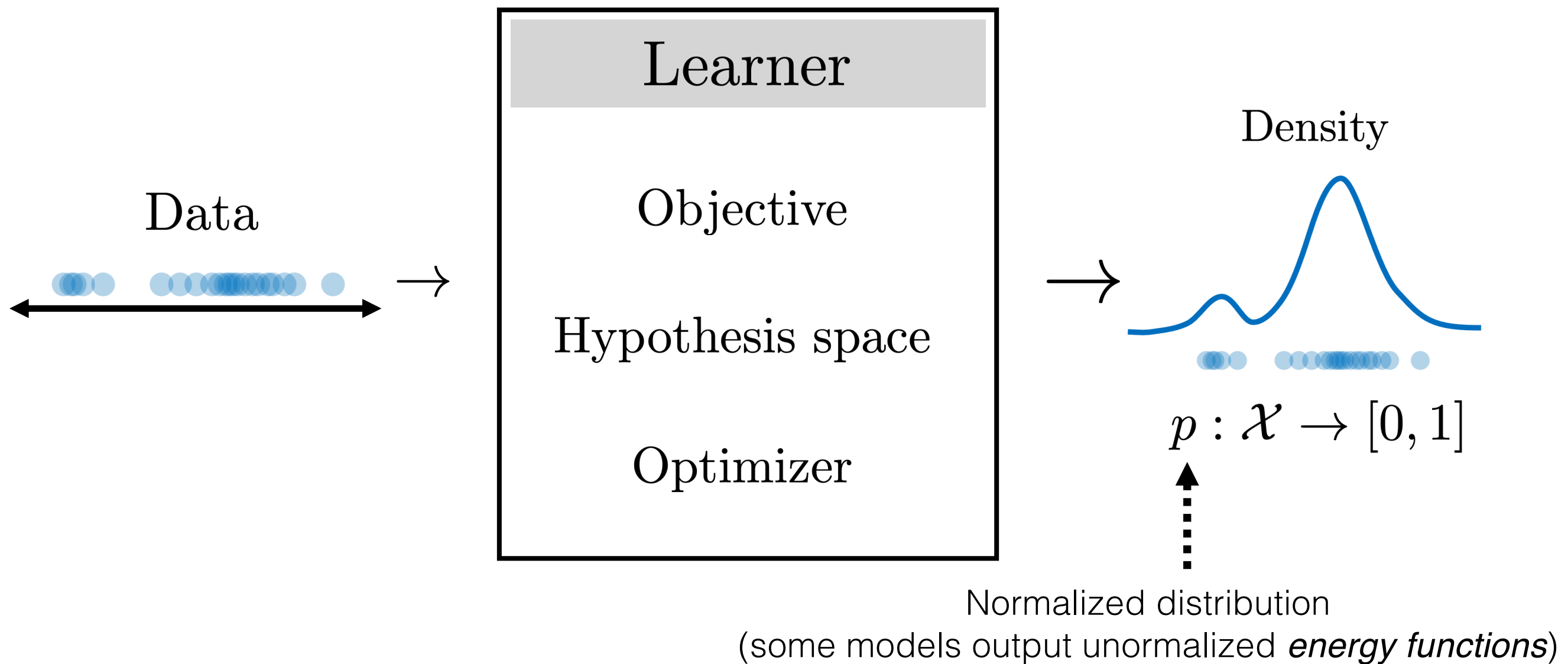
$$x = G(z)$$

Learning a generative model



[figs modified from: http://introtodeeplearning.com/materials/2019_6S191_L4.pdf]

Learning a density model



[figs modified from: http://introtodeeplearning.com/materials/2019_6S191_L4.pdf]

Case study #1: Fitting a Gaussian to data

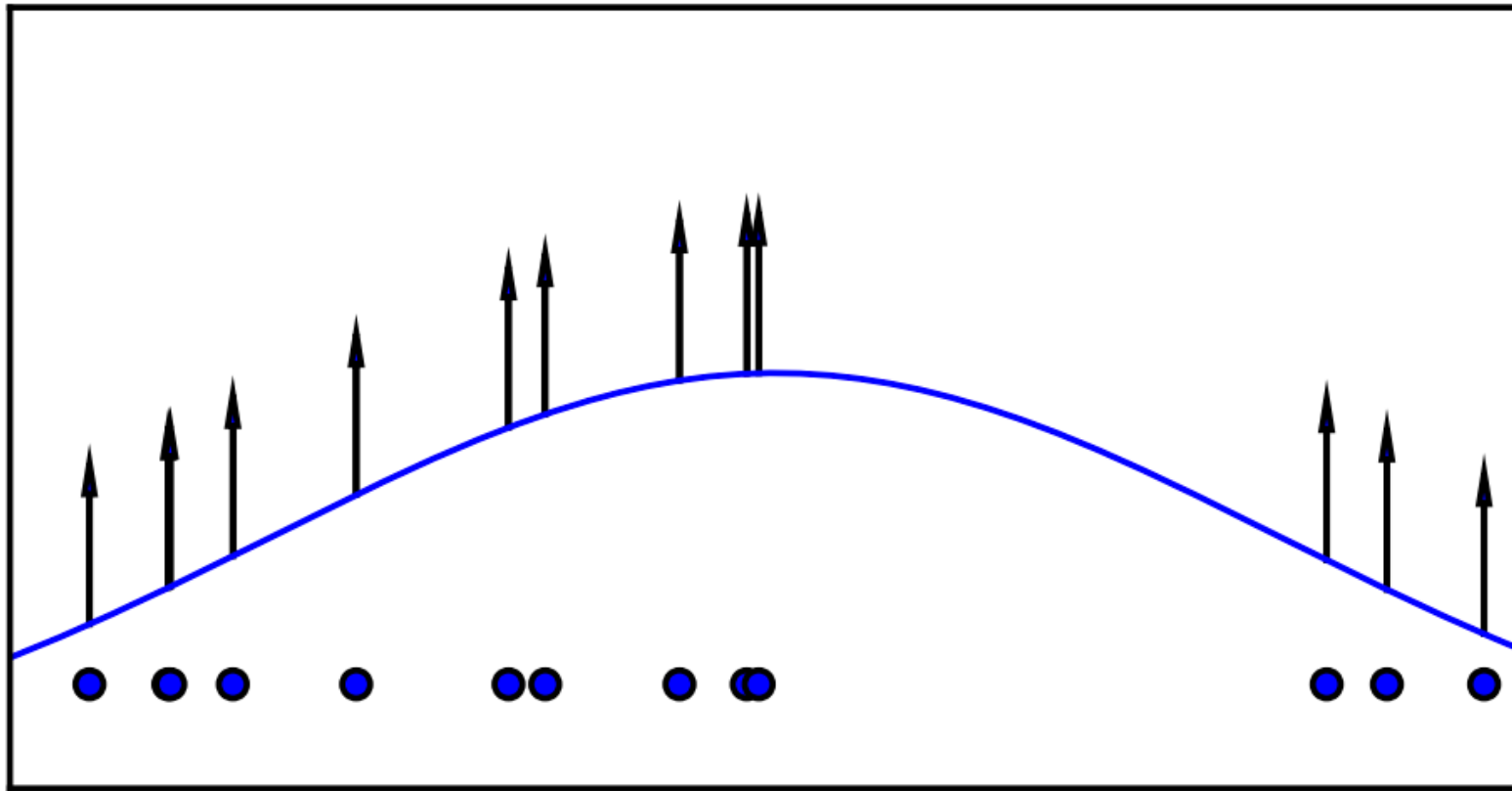


fig from [Goodfellow, 2016]

Max likelihood objective

$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

Considering only Gaussian fits

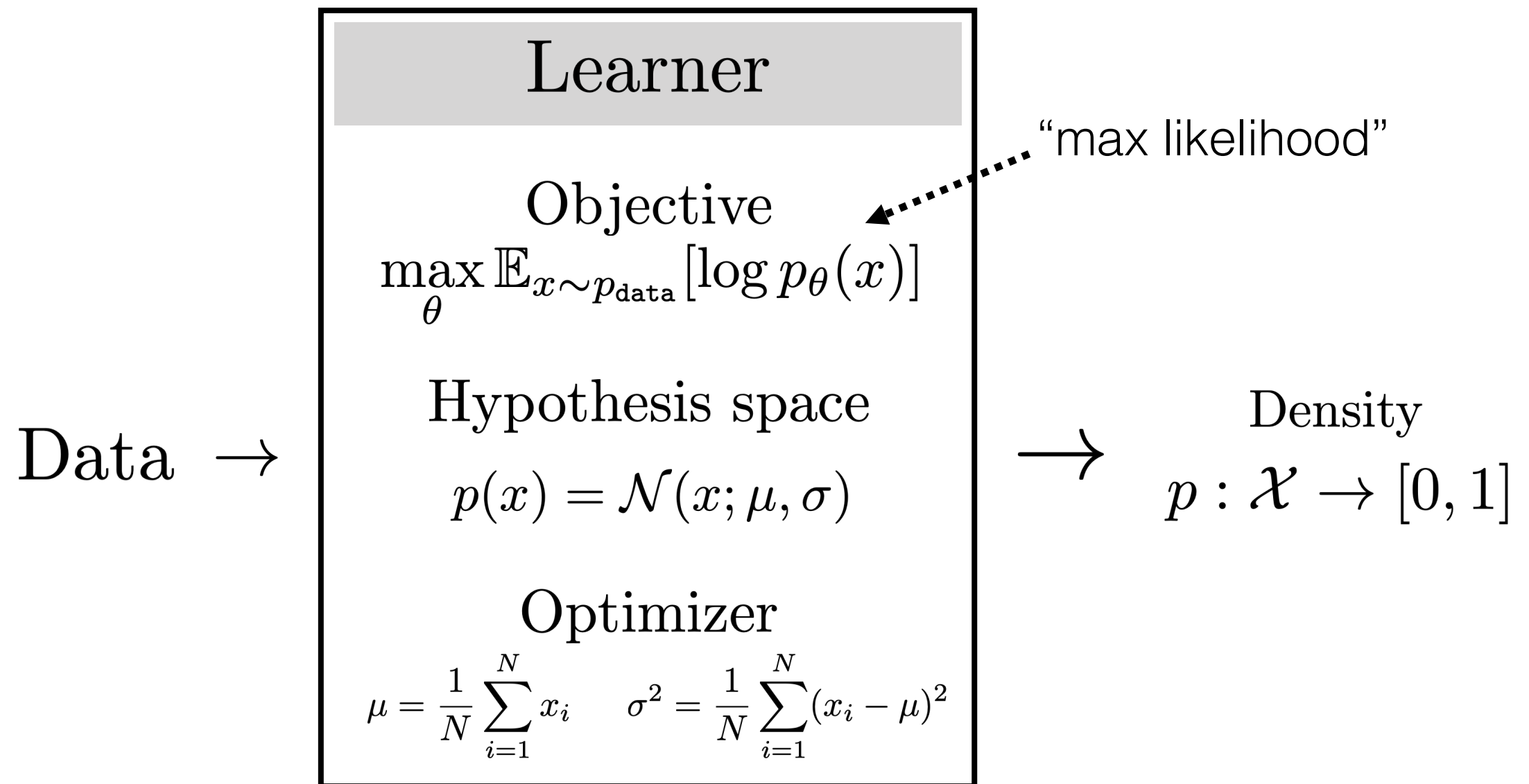
$$p_{\theta}(x) = \mathcal{N}(x; \mu, \sigma)$$

$$\theta = [\mu, \sigma]$$

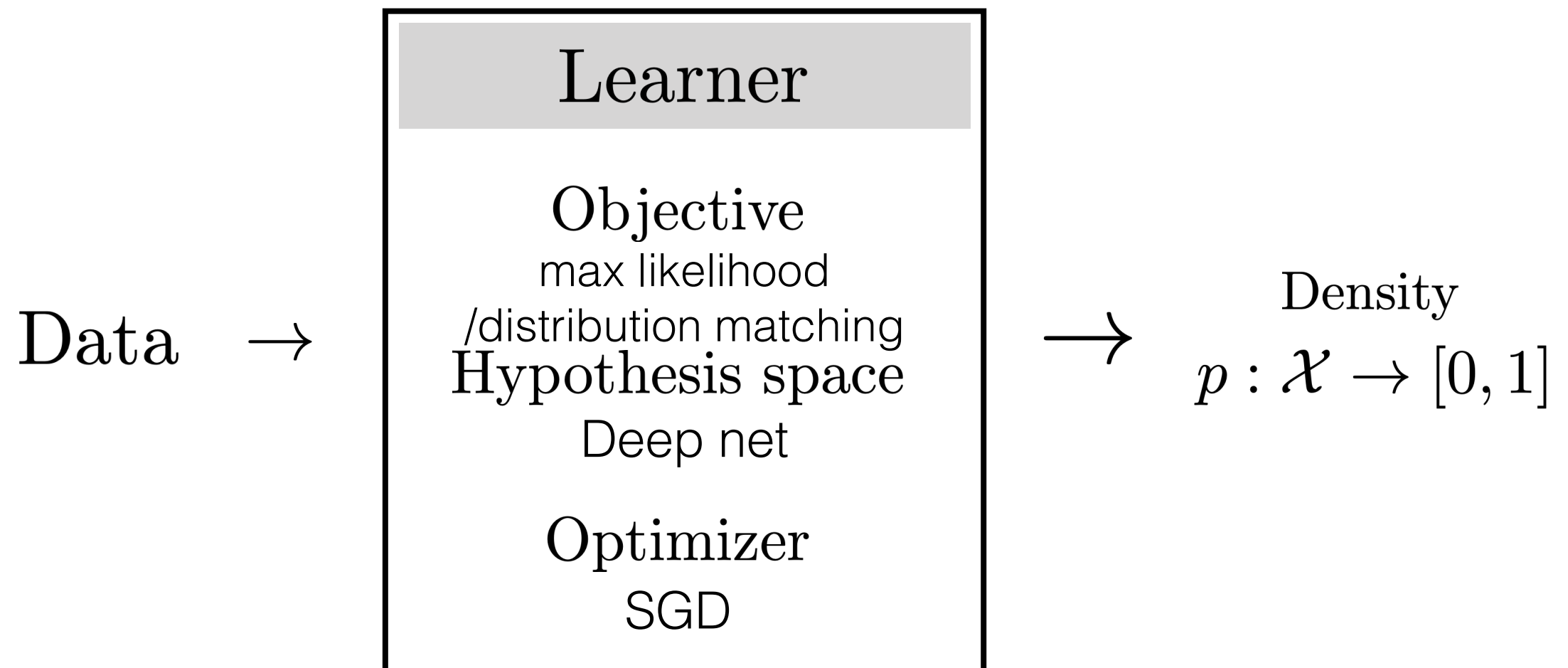
Closed form optimum:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

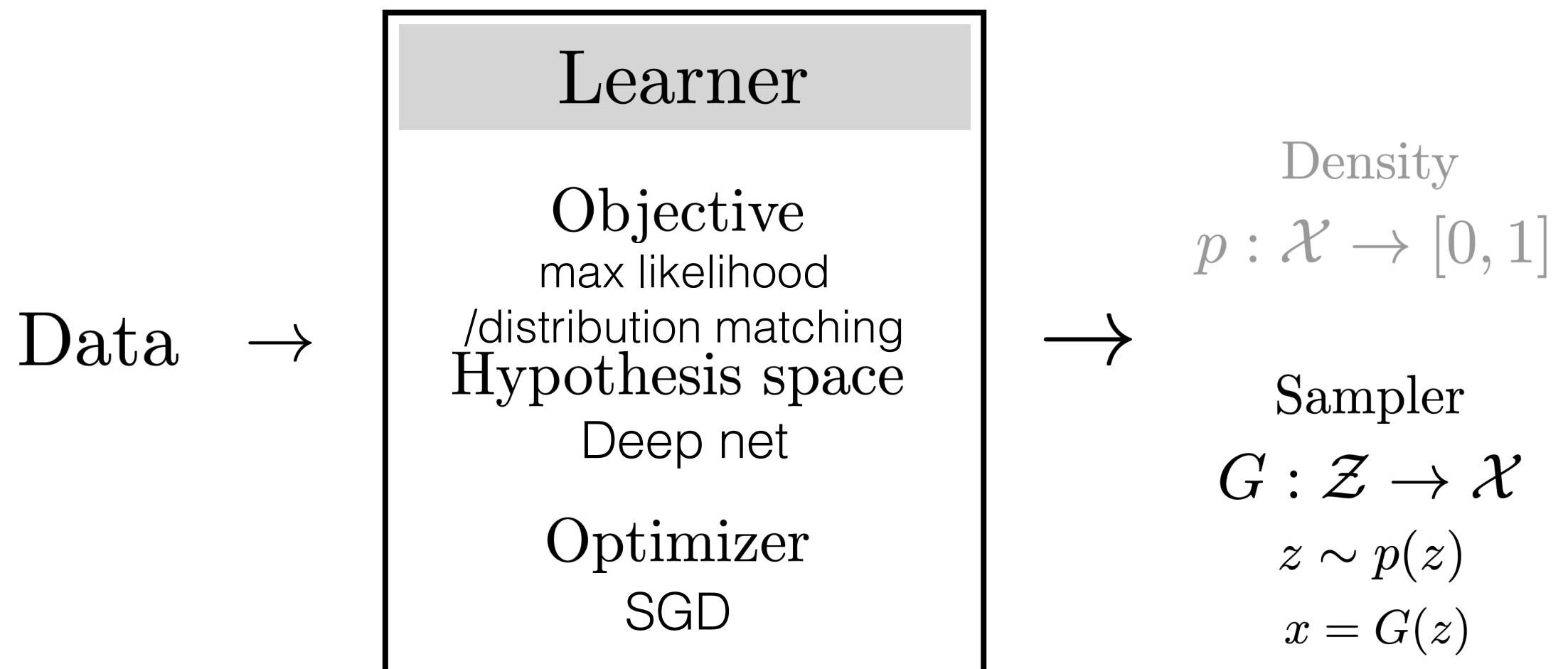
Case study #1: Fitting a Gaussian to data



Case study #2: learning a deep generative model



Case study #2: learning a deep generative model

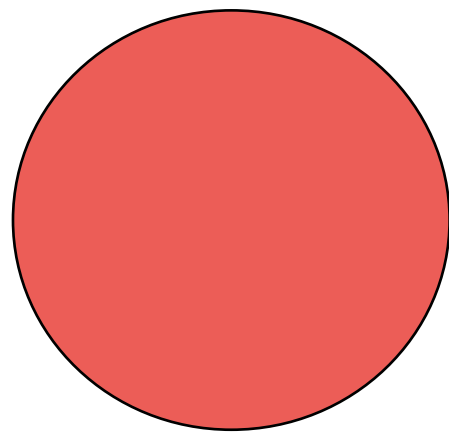


Models that provide a sampler but no density are called **implicit generative models**

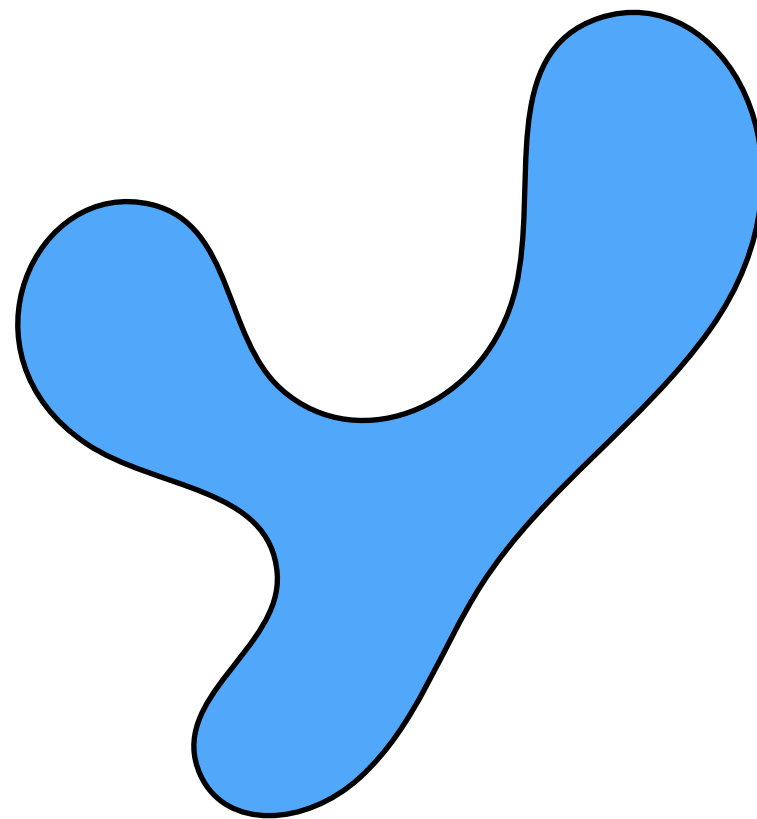
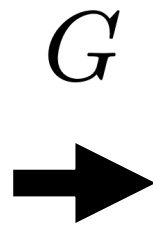
Deep generative models are distribution transformers

Prior distribution

Target distribution

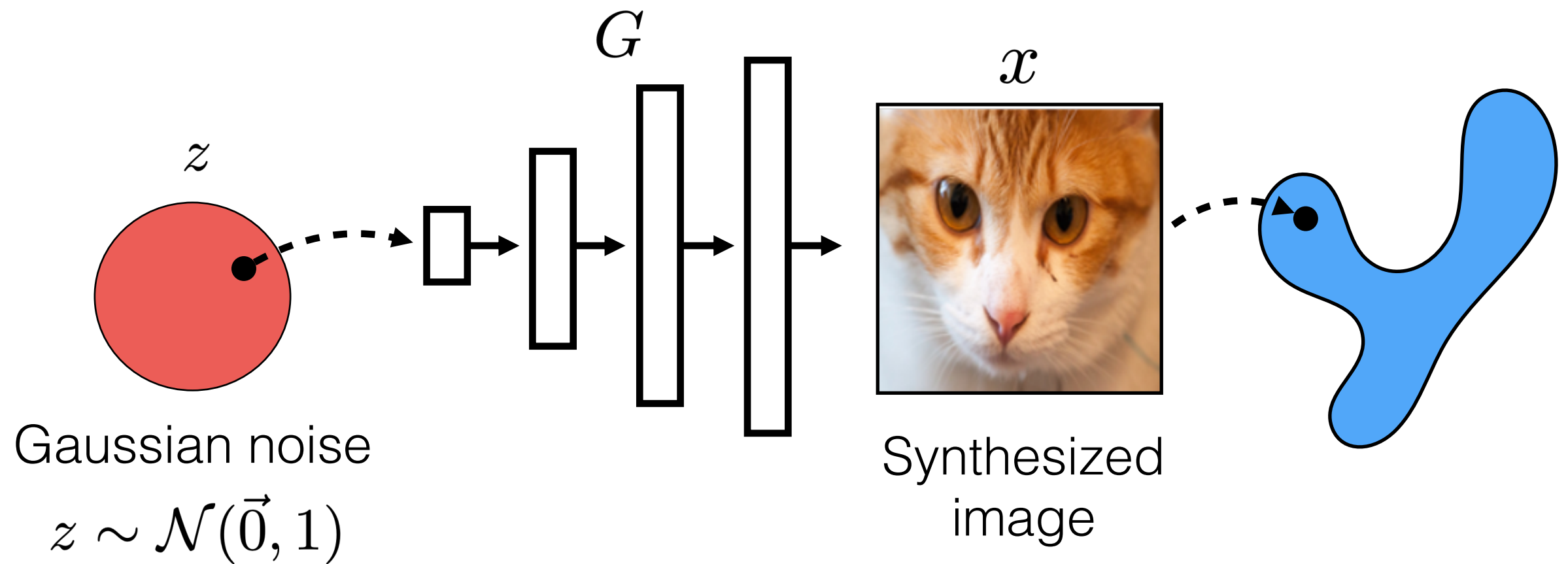


$p(z)$

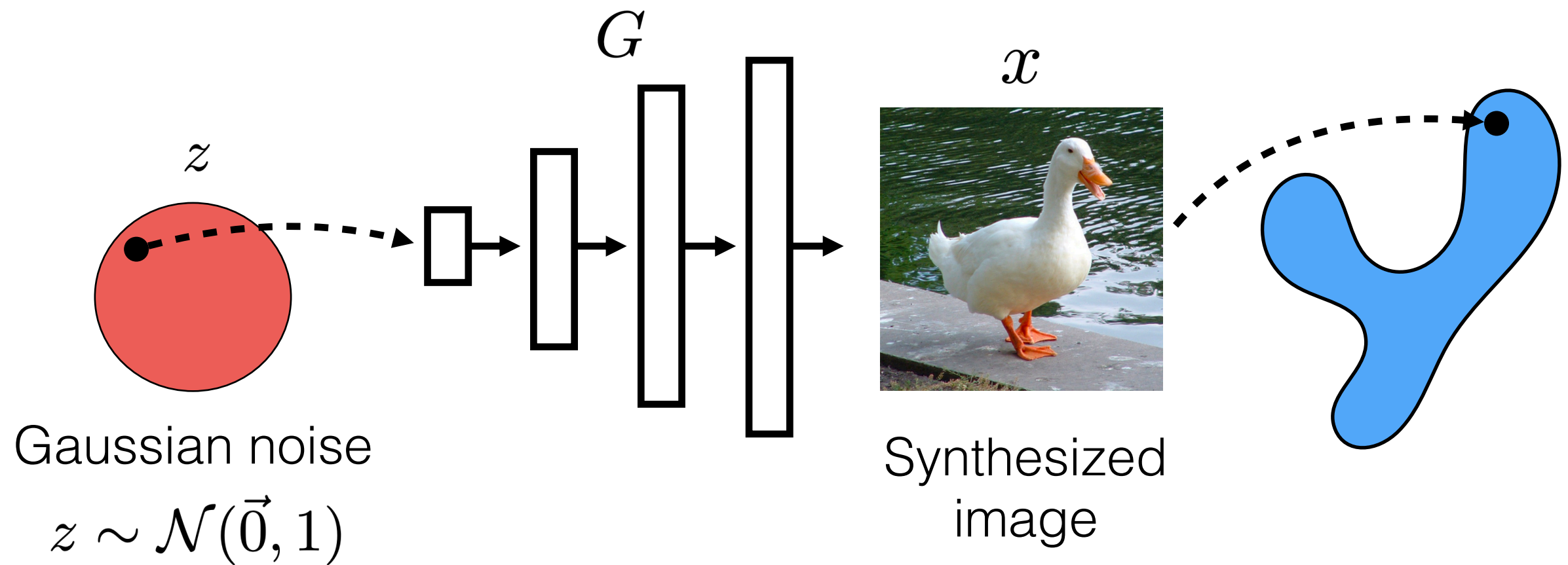


$p(x)$

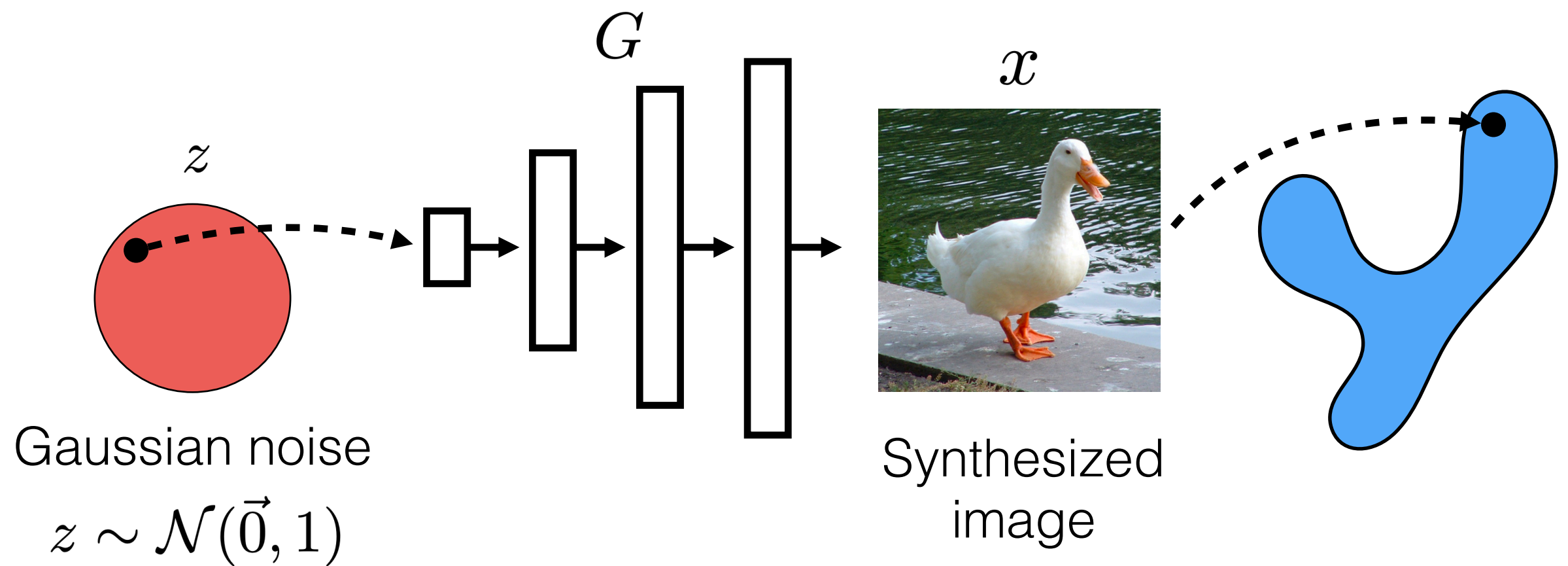
Deep generative models are distribution transformers

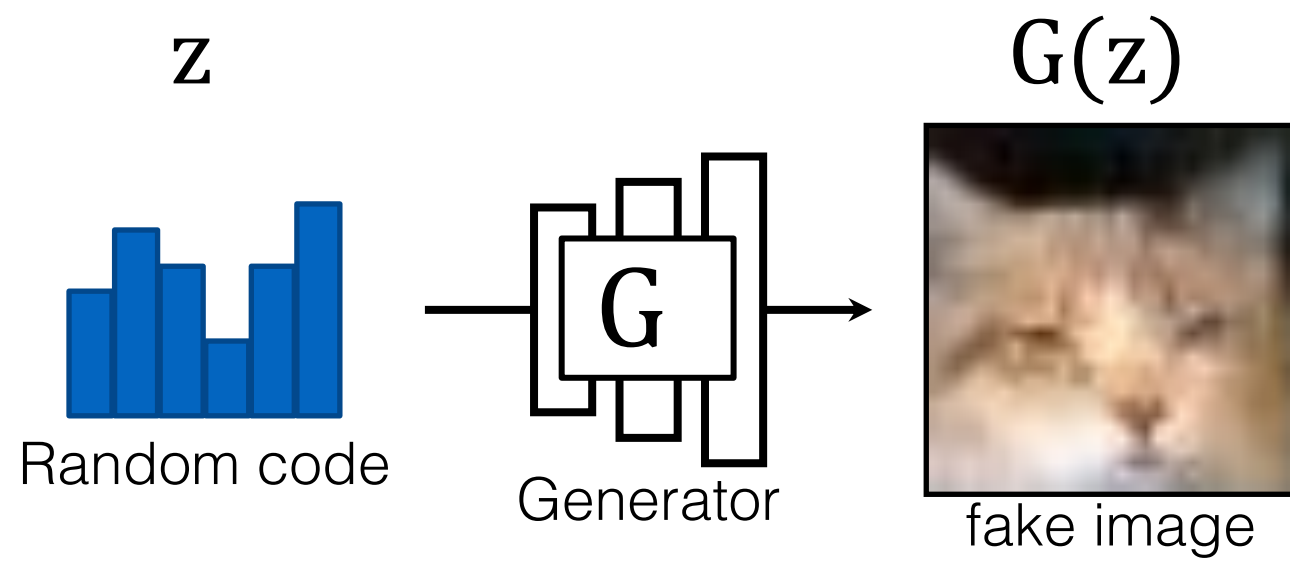


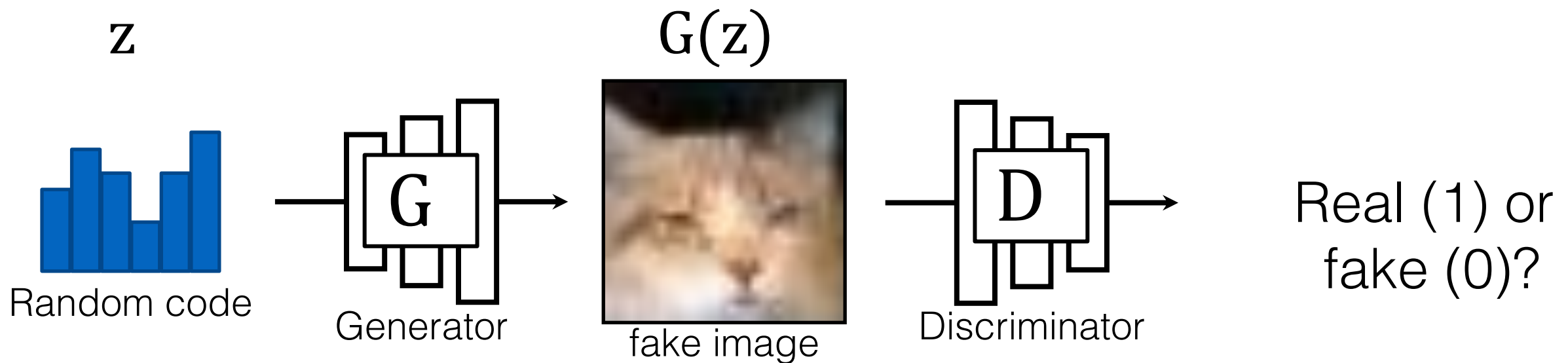
Deep generative models are distribution transformers



Generative Adversarial Networks (GANs)

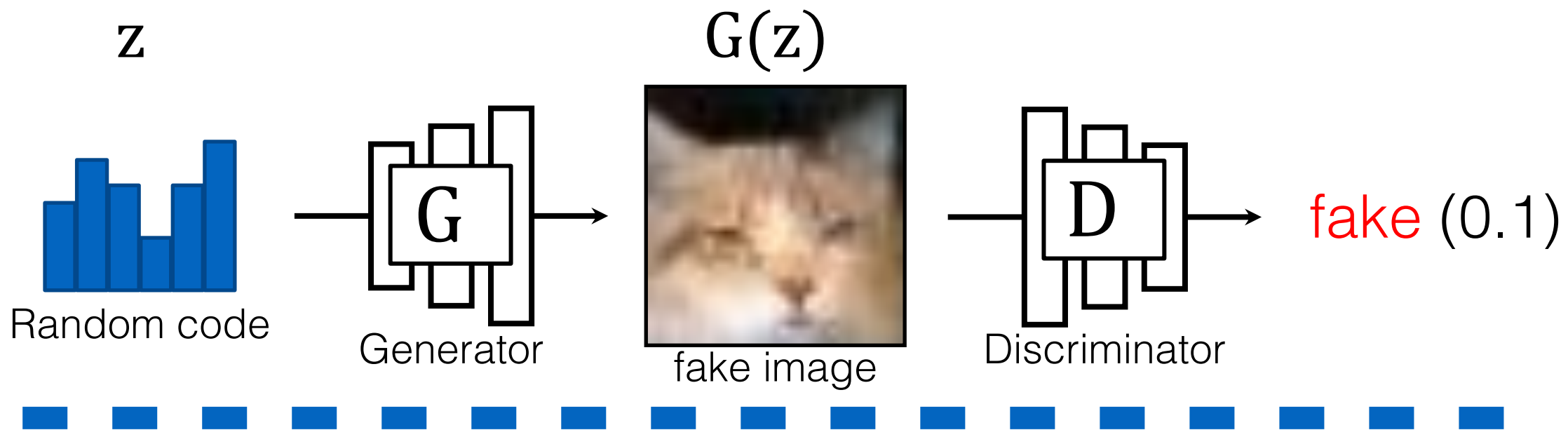






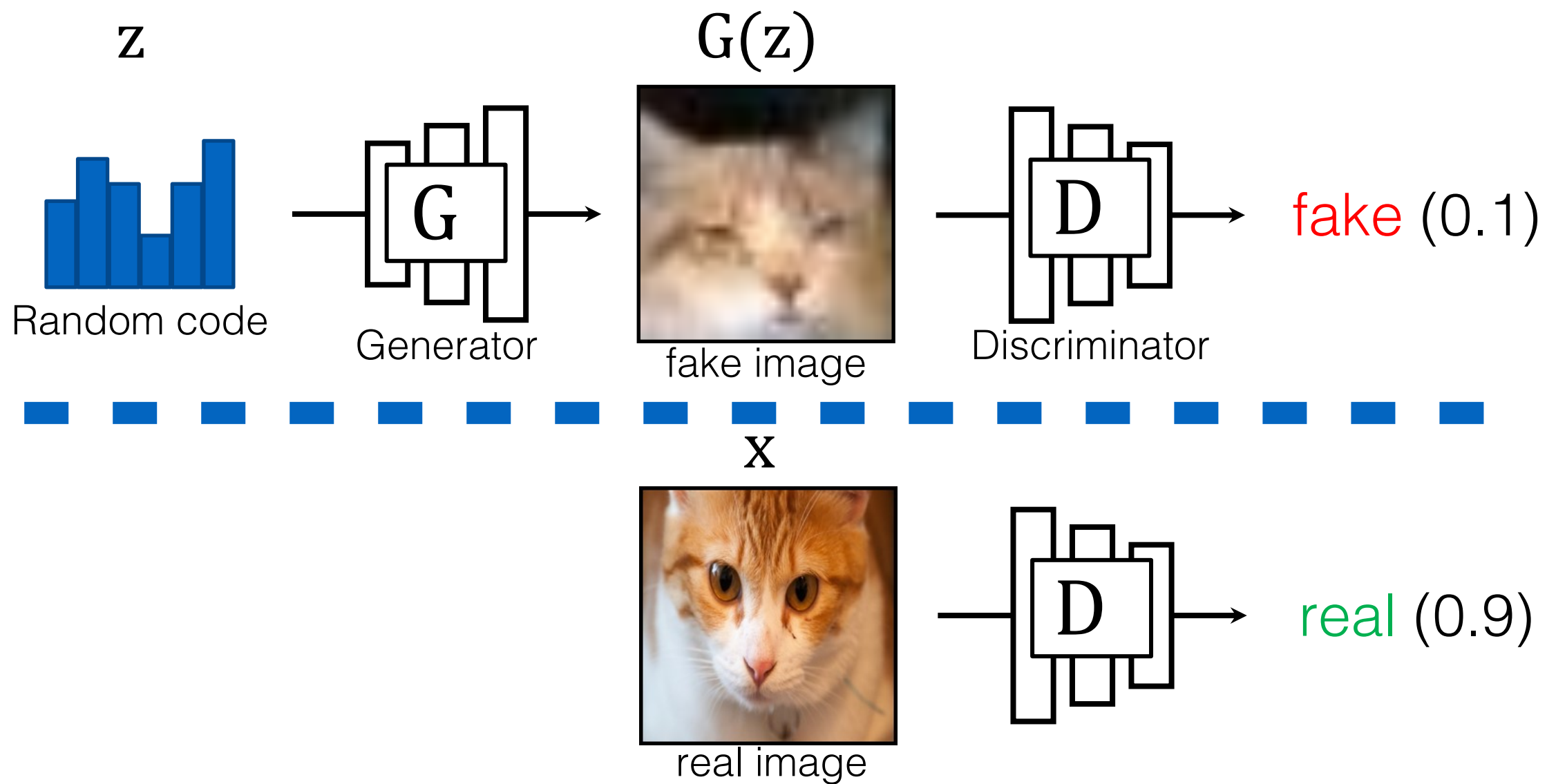
A two-player game:

- G tries to generate fake images that can fool D .
- D tries to detect fake images.



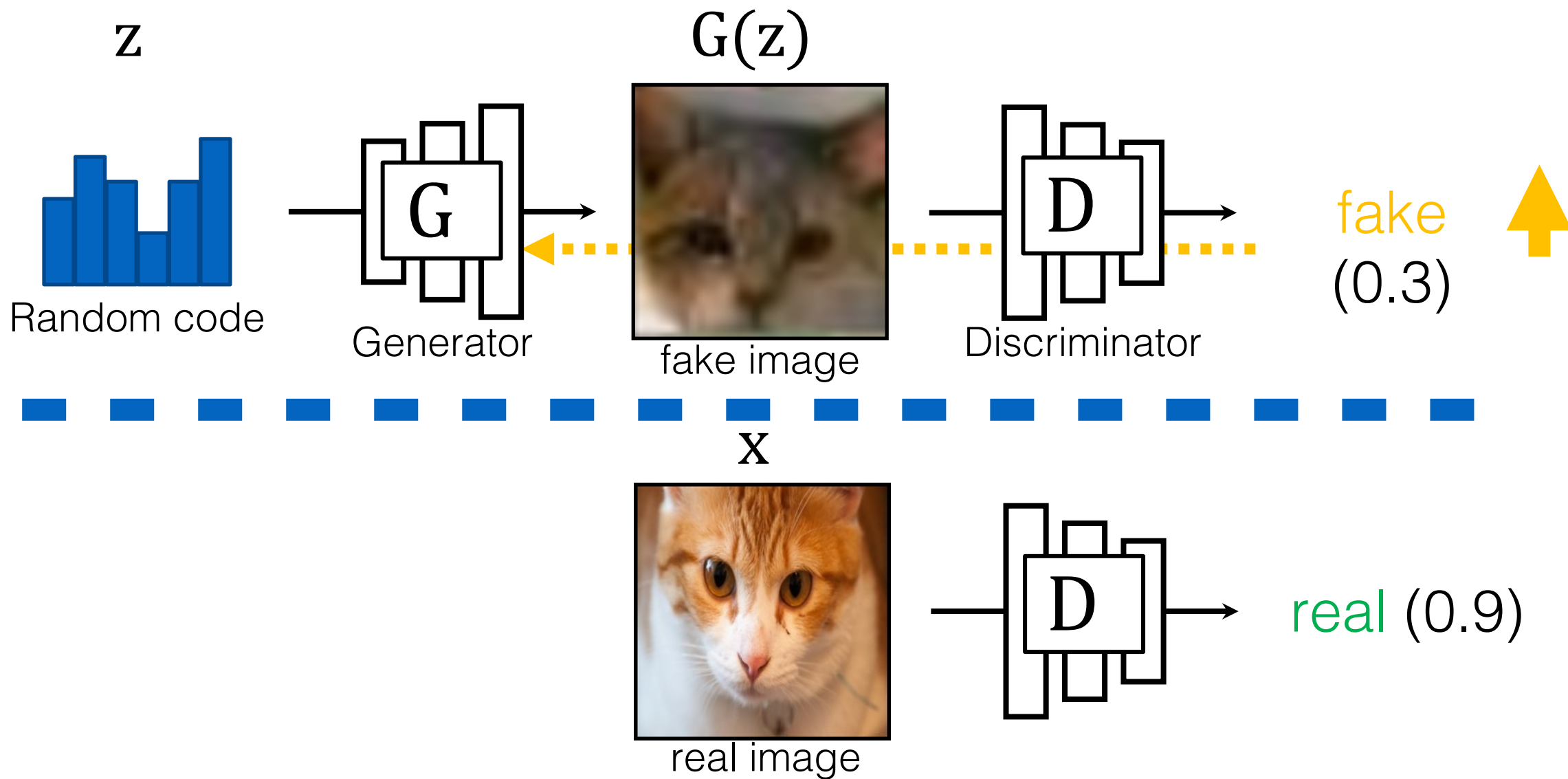
Learning objective (GANs)

$$\min_G \max_D \mathbb{E}[\log(1 - D(G(z)))]$$



Learning objective (GANs)

$$\min_G \max_D \mathbb{E} [\log(1 - D(G(z))) + \log D(x)]$$



Learning objective (GANs)

$$\min_G \max_D \mathbb{E}[\log(1 - D(G(z))) + \log D(x)]$$

GANs Training



G tries to synthesize fake images that fool D

D tries to identify the fakes

- Training: iterate between training D and G with backprop.
- Global optimum when G reproduces data distribution.

Thank You!



16-726, Spring 2021

<https://learning-image-synthesis.github.io/>