

# Perceptual Loss, GANs (part I)


Jun-Yan Zhu

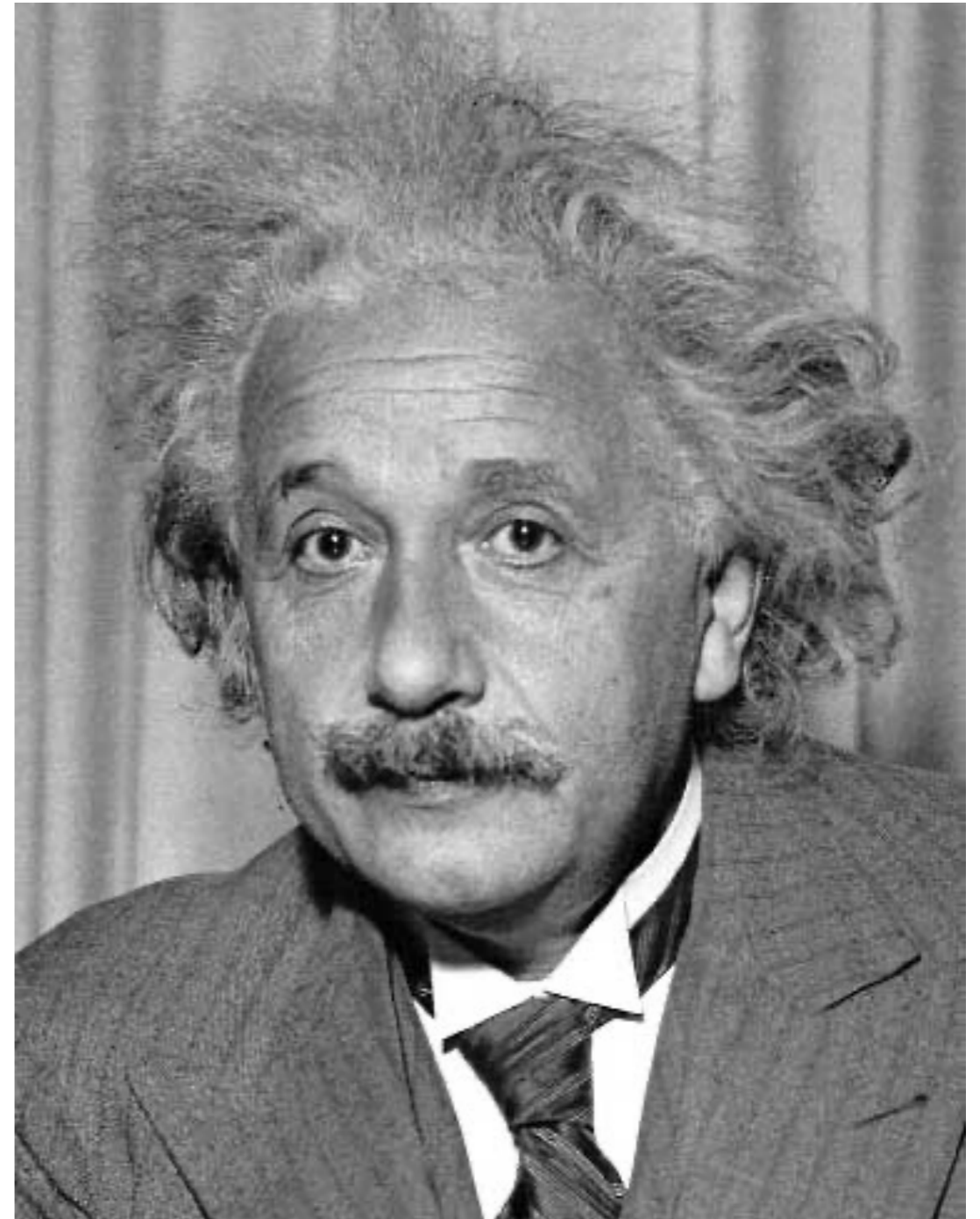
16-726 Learning-based Image Synthesis, Spring 2022

many slides from Alyosha Efros, Phillip Isola, Richard Zhang, James Hays, and Andrea Vedaldi, Jitendra Malik.

# 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

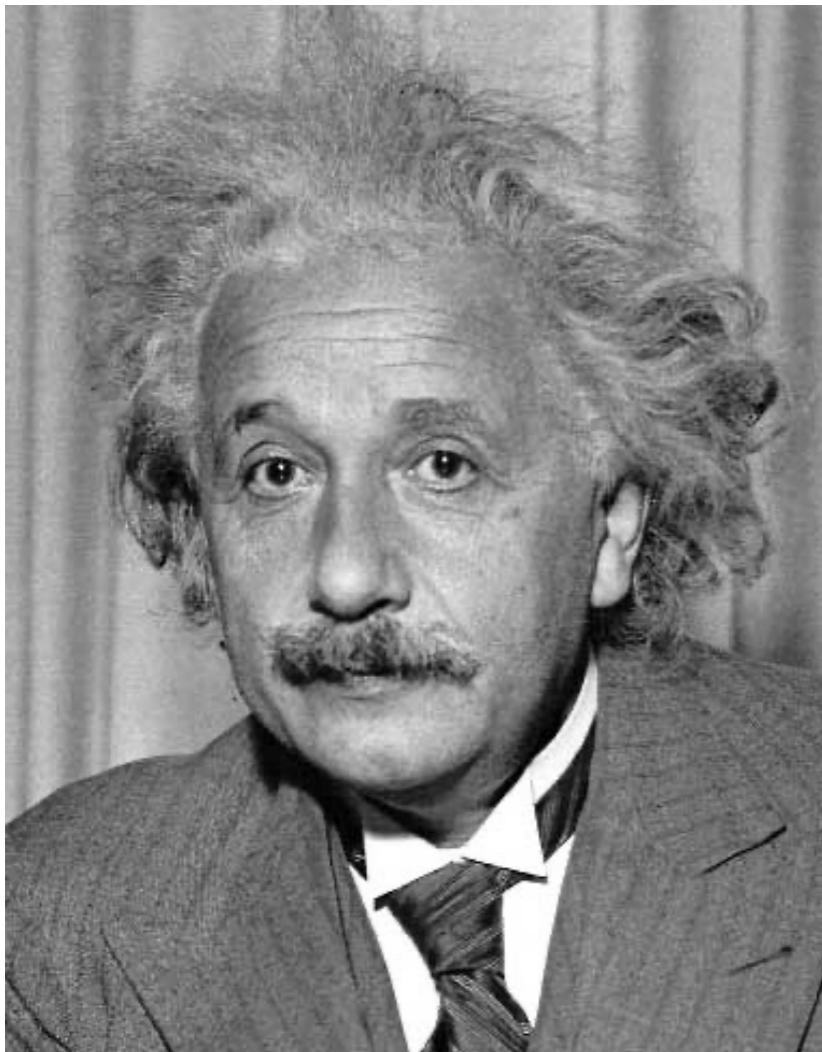


# Matching with filters

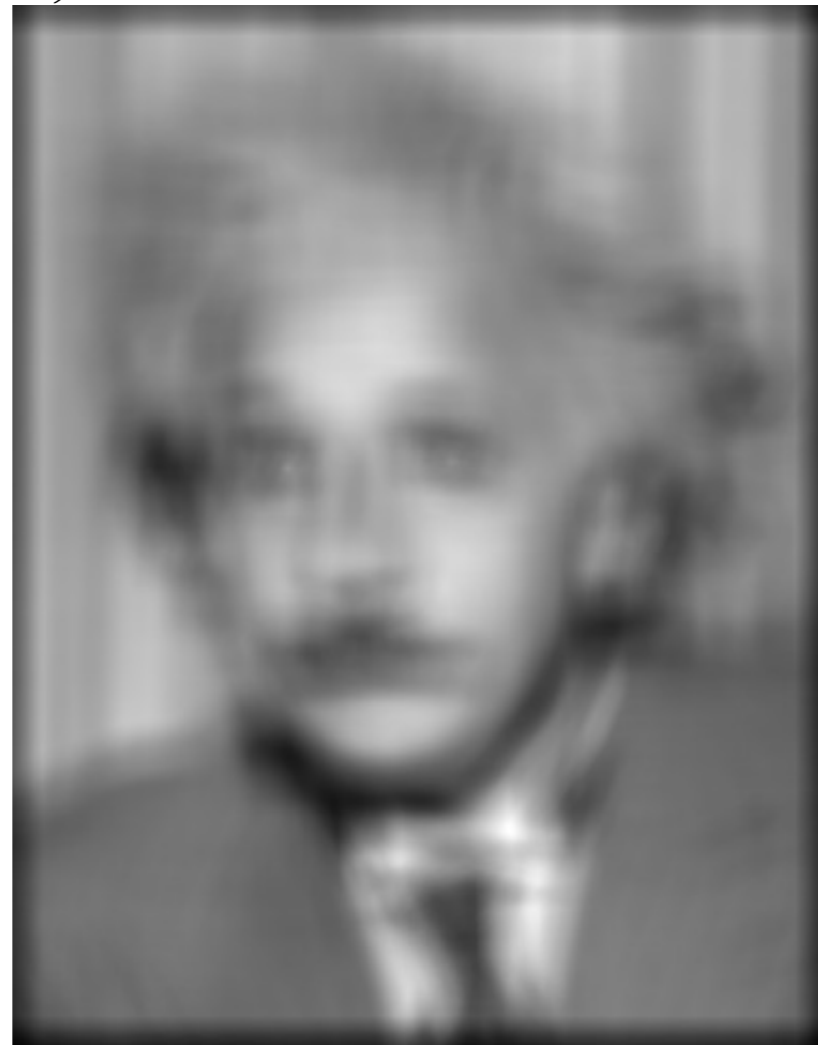
- Goal: find  in image
- Method 0: filter the image with eye patch

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

f = image  
g = filter



Input



Filtered Image

What went wrong?

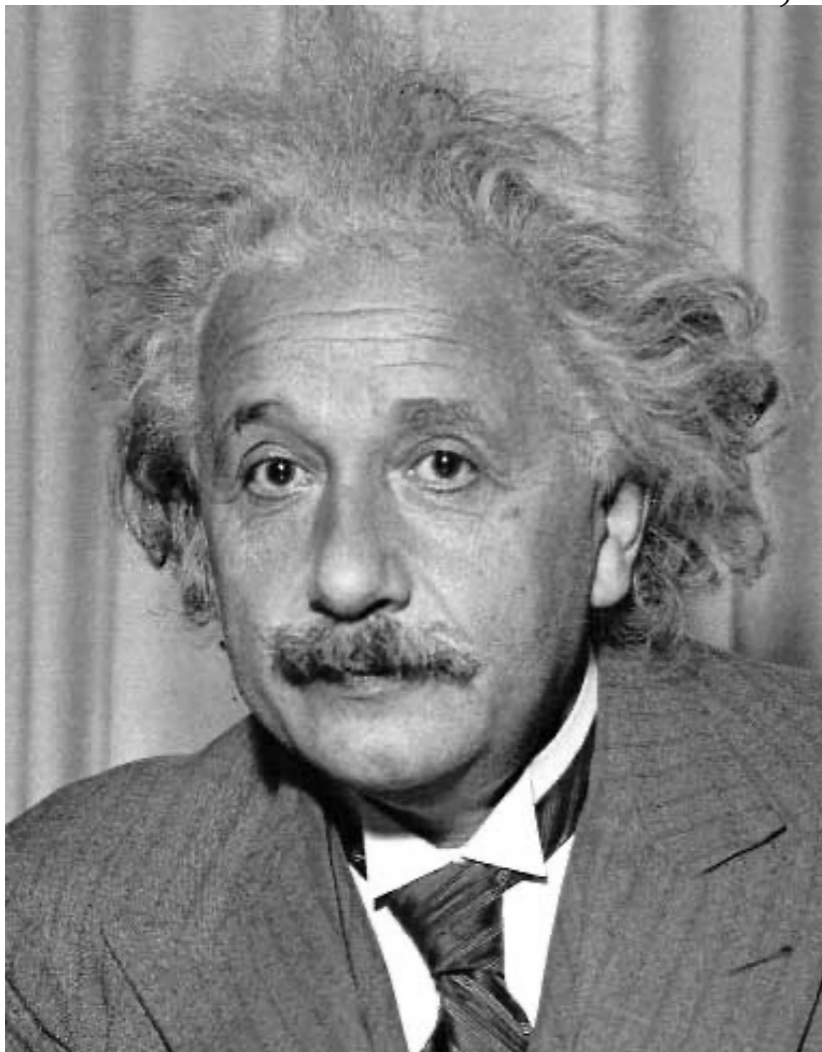
# Matching with filters

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

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

f = image  
g = filter

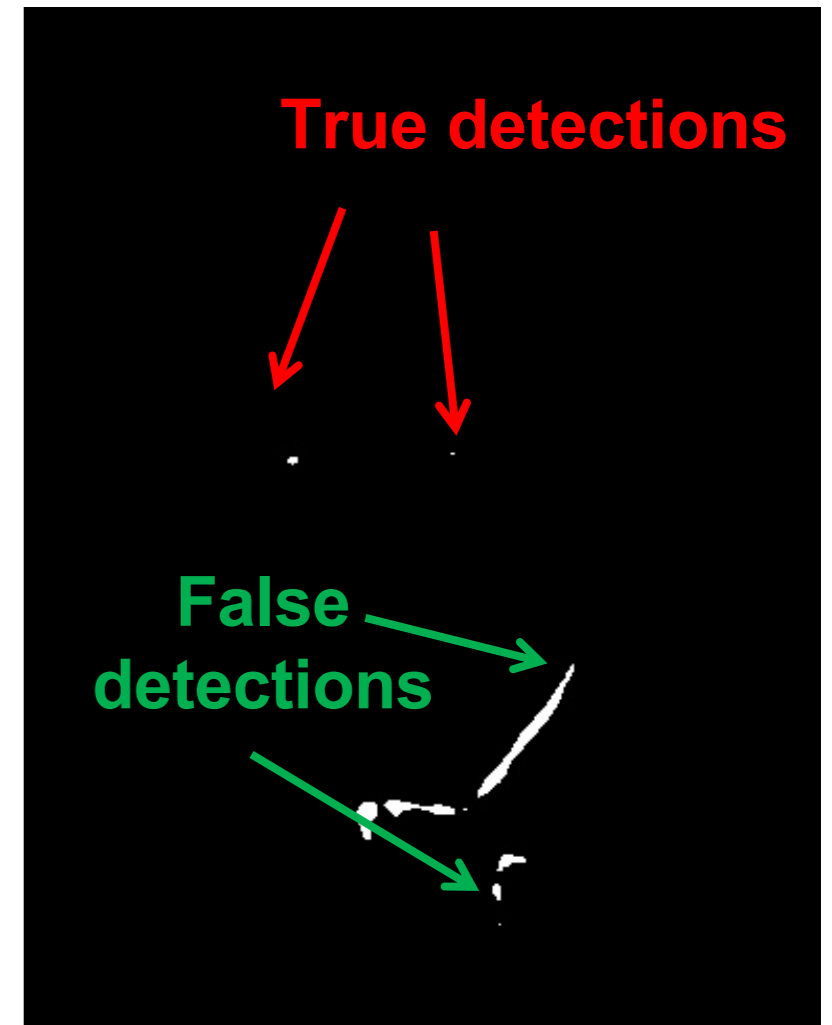
← mean of f



Input



Filtered Image (scaled)



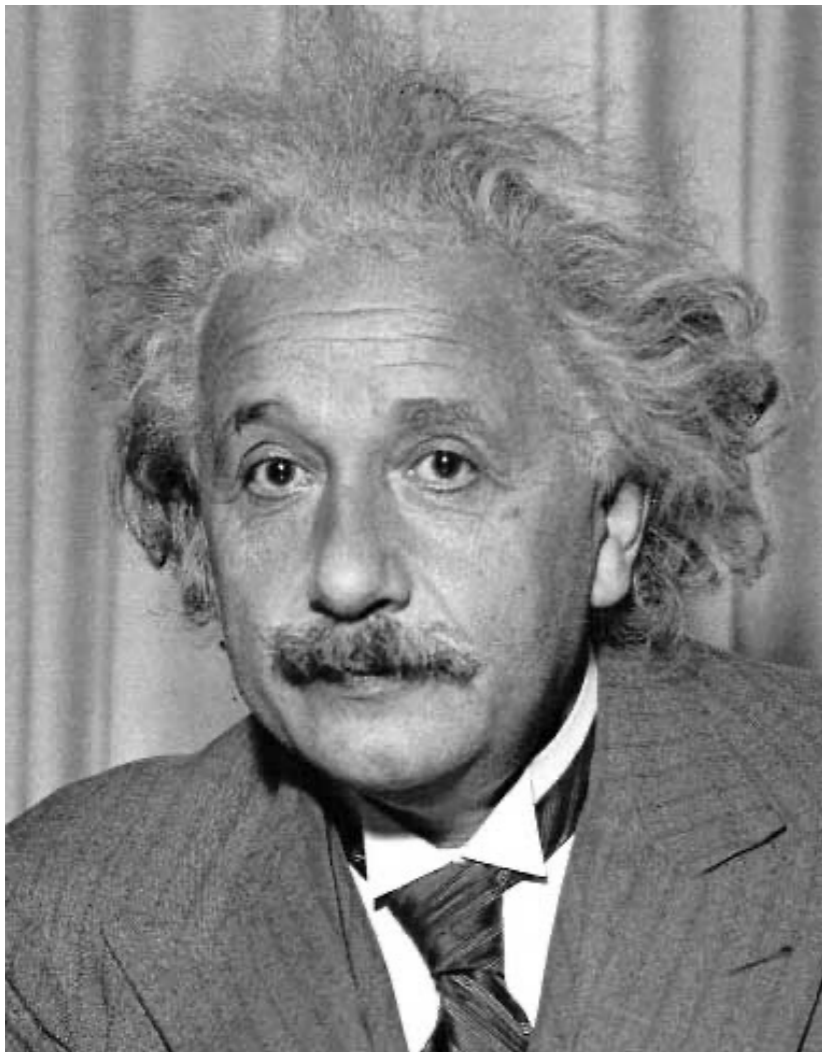
Thresholded Image

# Matching with filters

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

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

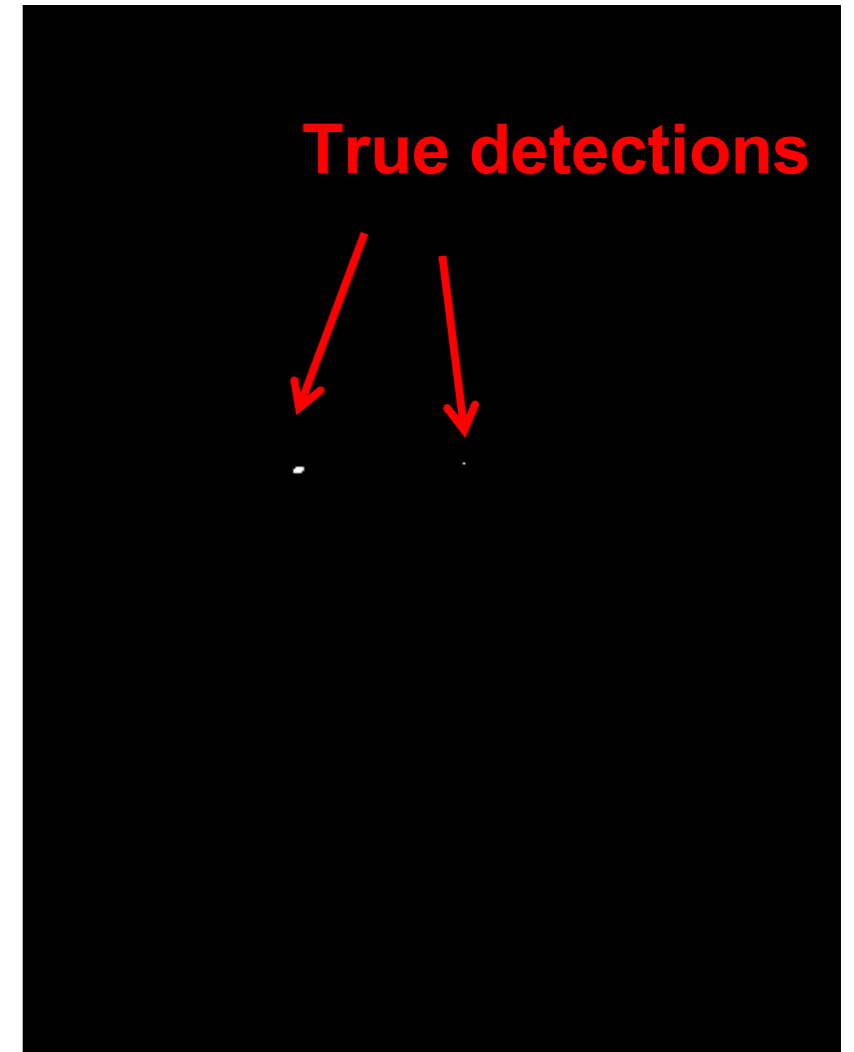
f = image  
g = filter



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

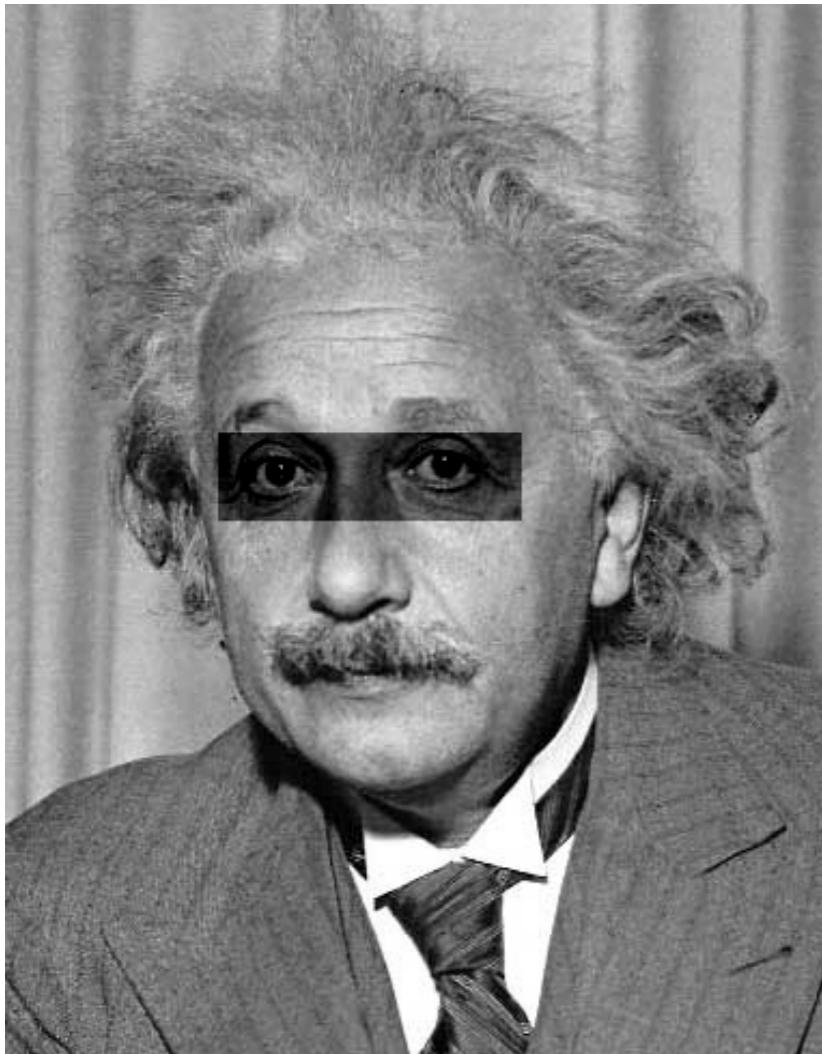
- Goal: find  in image

**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

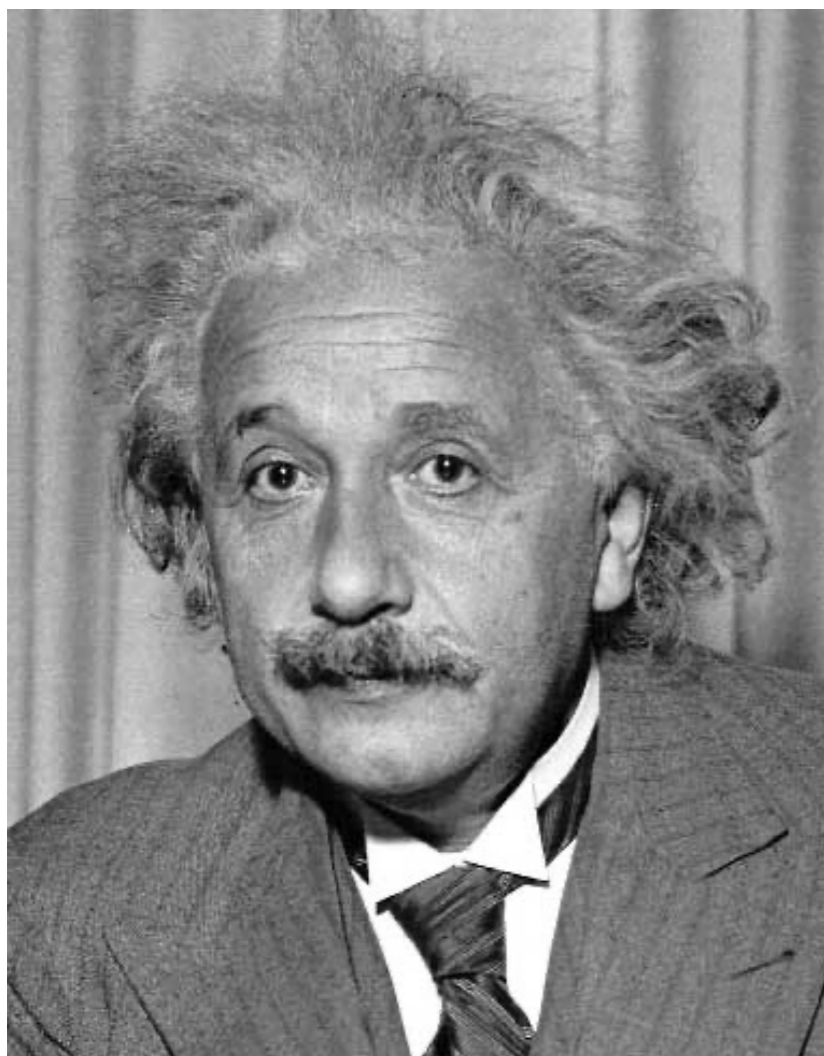
$$h[m,n] = \frac{\sum_{k,l} (g[k,l] - \bar{g})(f[m+k,n+l] - \bar{f}_{m,n})}{\left( \sum_{k,l} (g[k,l] - \bar{g})^2 \sum_{k,l} (f[m+k,n+l] - \bar{f}_{m,n})^2 \right)^{0.5}}$$

mean template mean image patch

↓ ↓

# Matching with filters

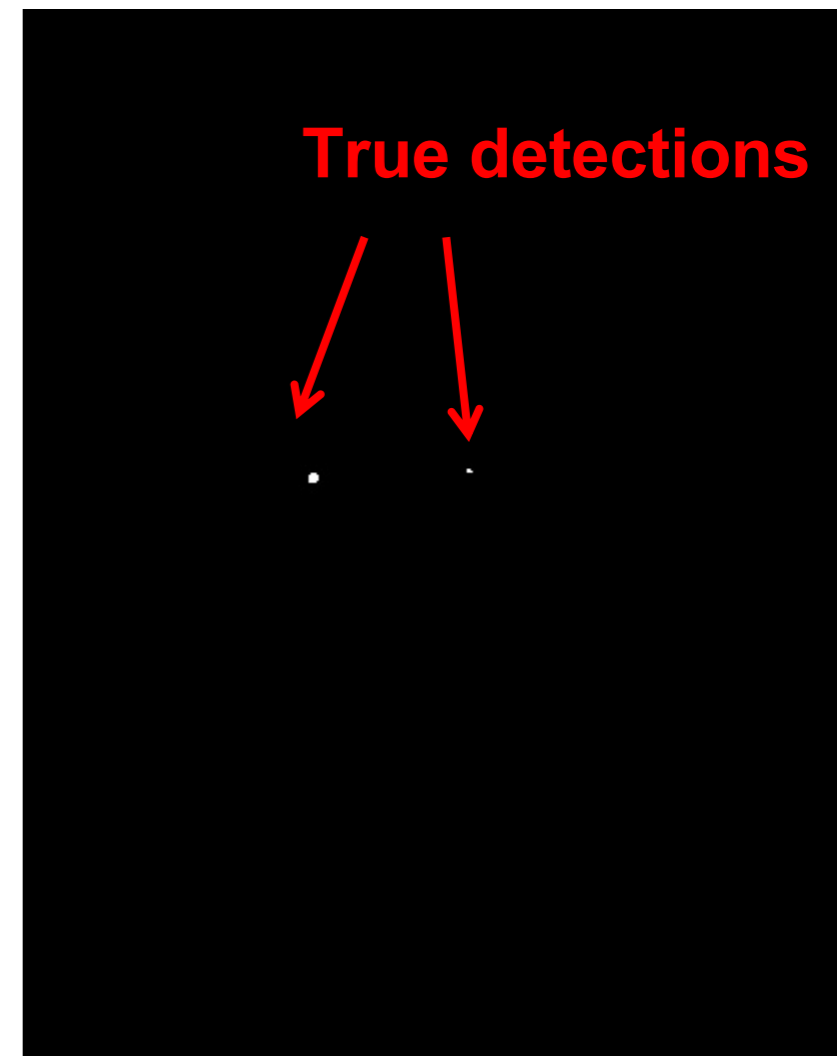
- Goal: find  in image
- Method 2: Normalized Cross-Correlation



Input




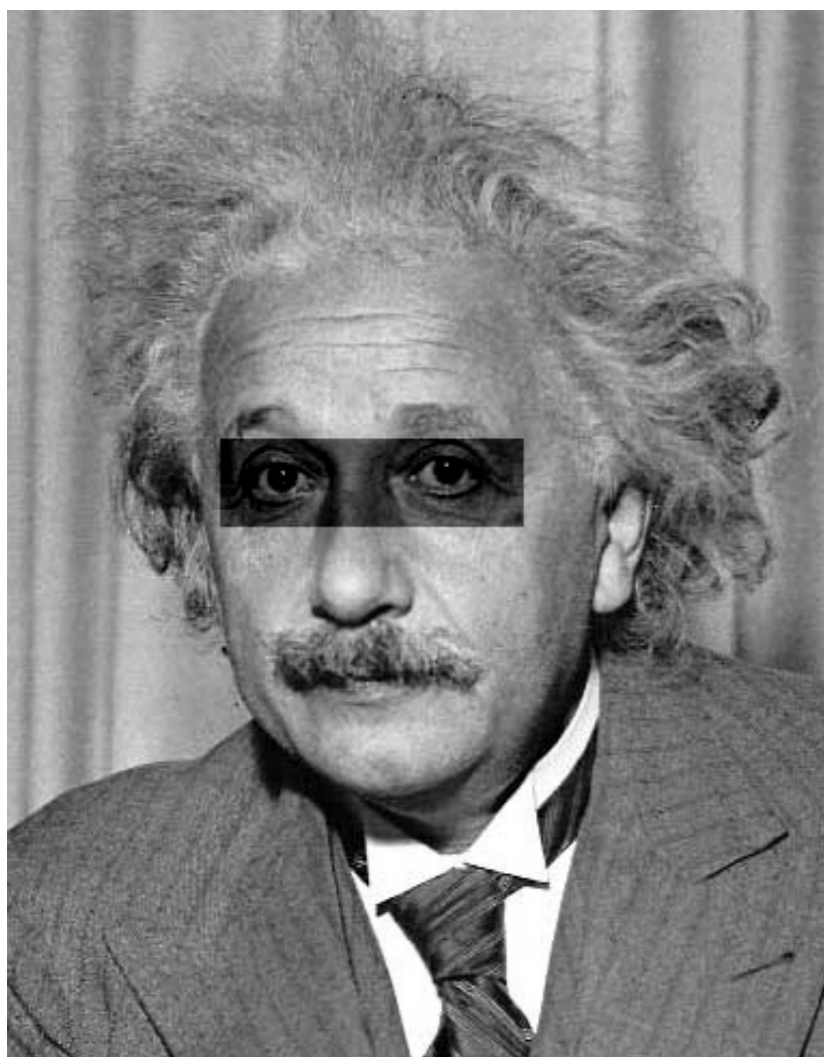
Normalized  $\chi$ -Correlation



Thresholded Image

# Matching with filters

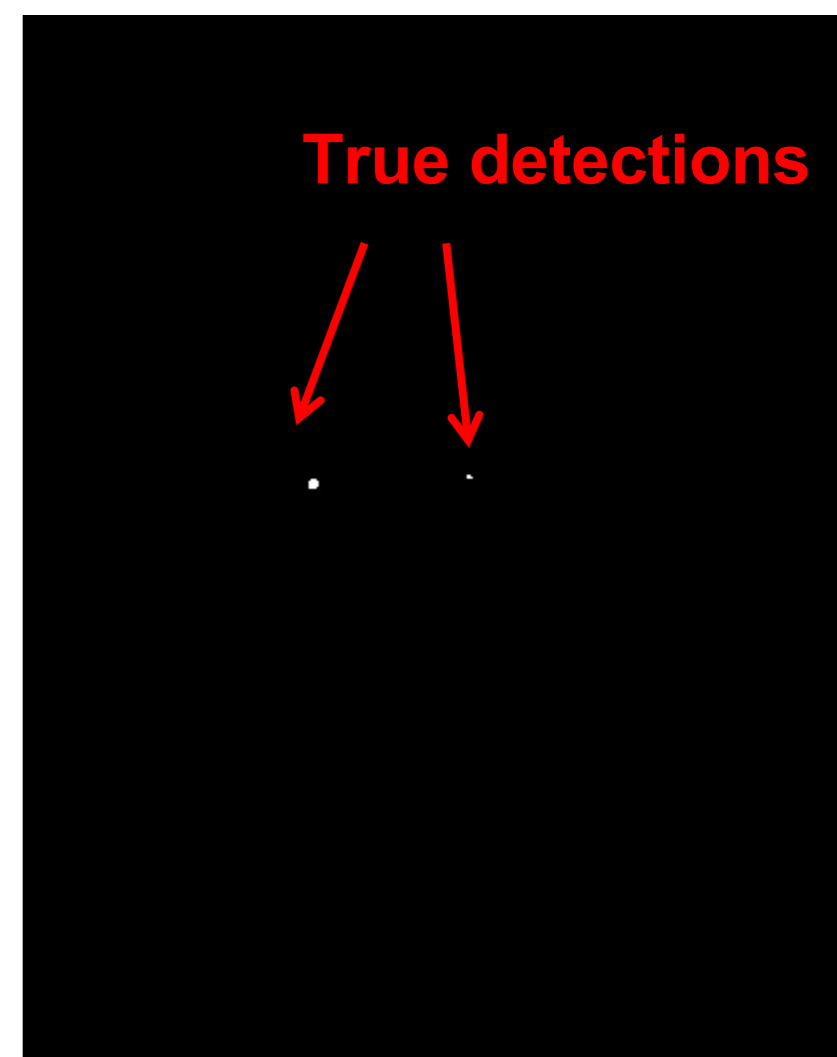
- Goal: find  in image
- Method 2: Normalized Cross-Correlation



Input



Normalized  $\chi$ -Correlation



Thresholded Image

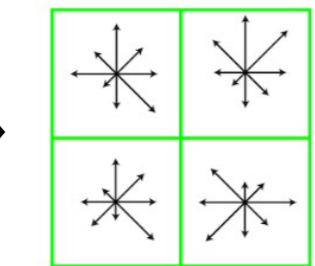
# Q: What is the best method to use?

- 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

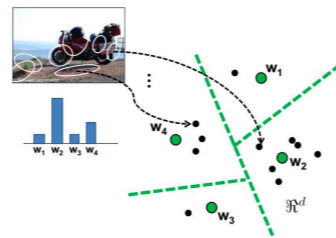
# Review

## (CNN for Image Synthesis)

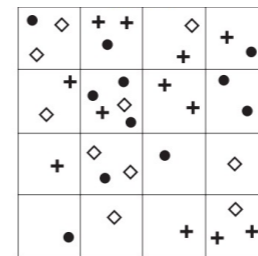
# Computer Vision before 2012



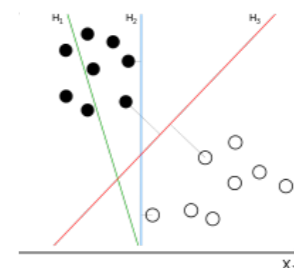
Features



Clustering



Pooling

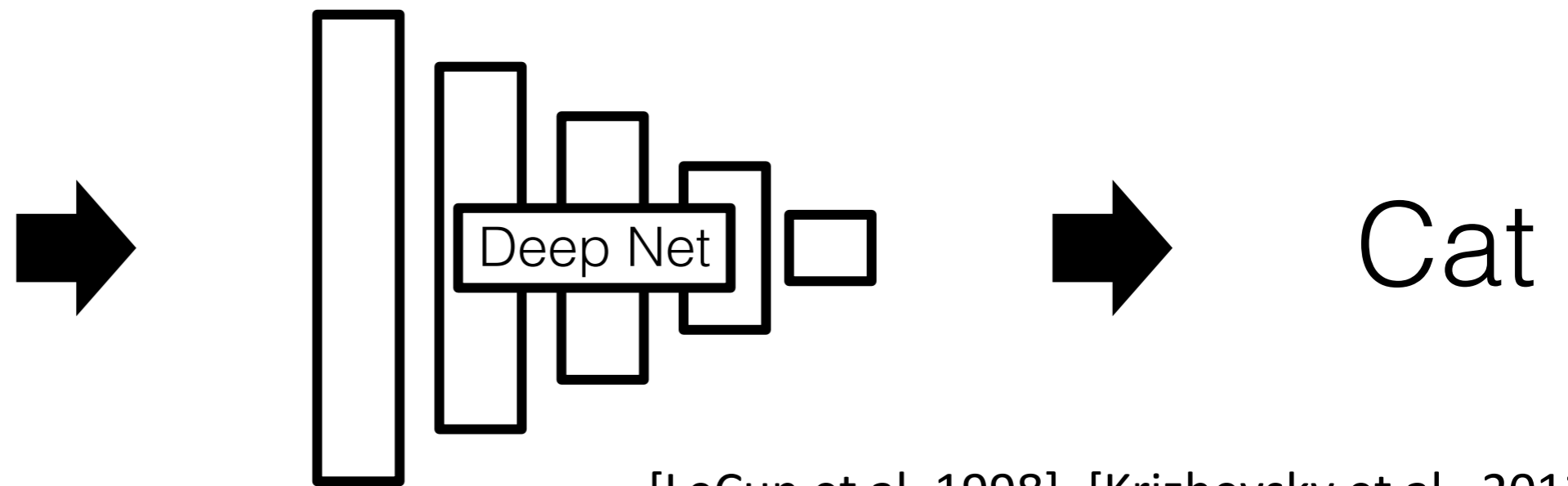
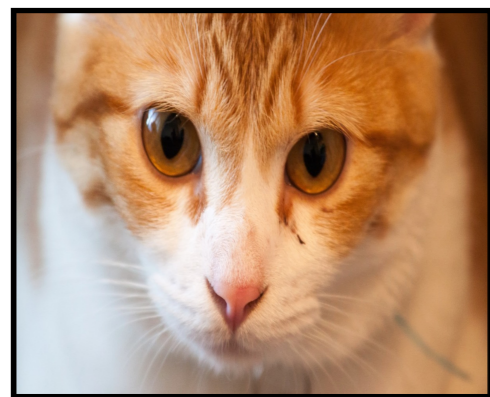
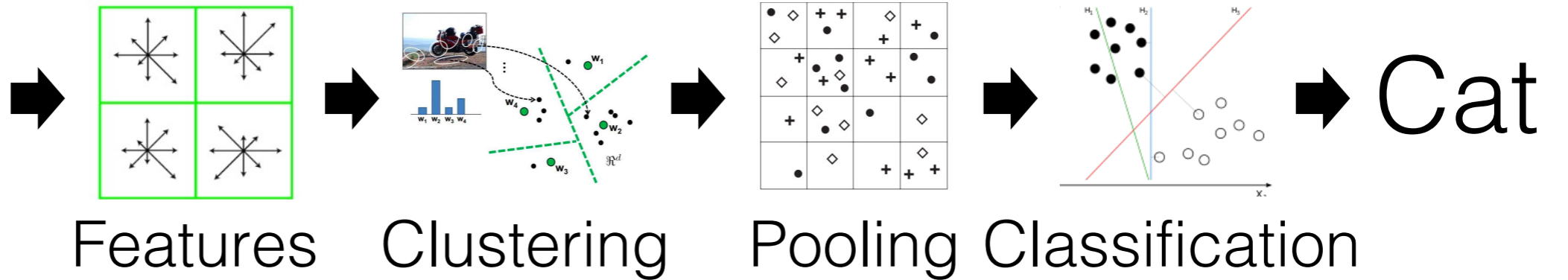
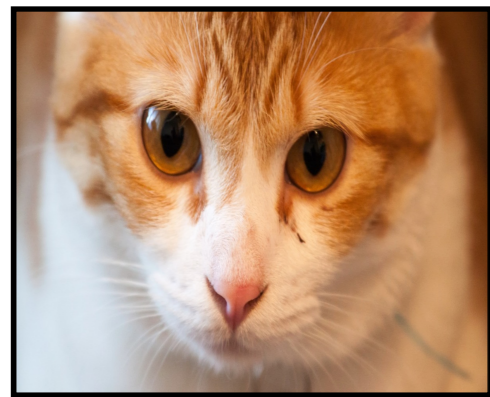


Classification



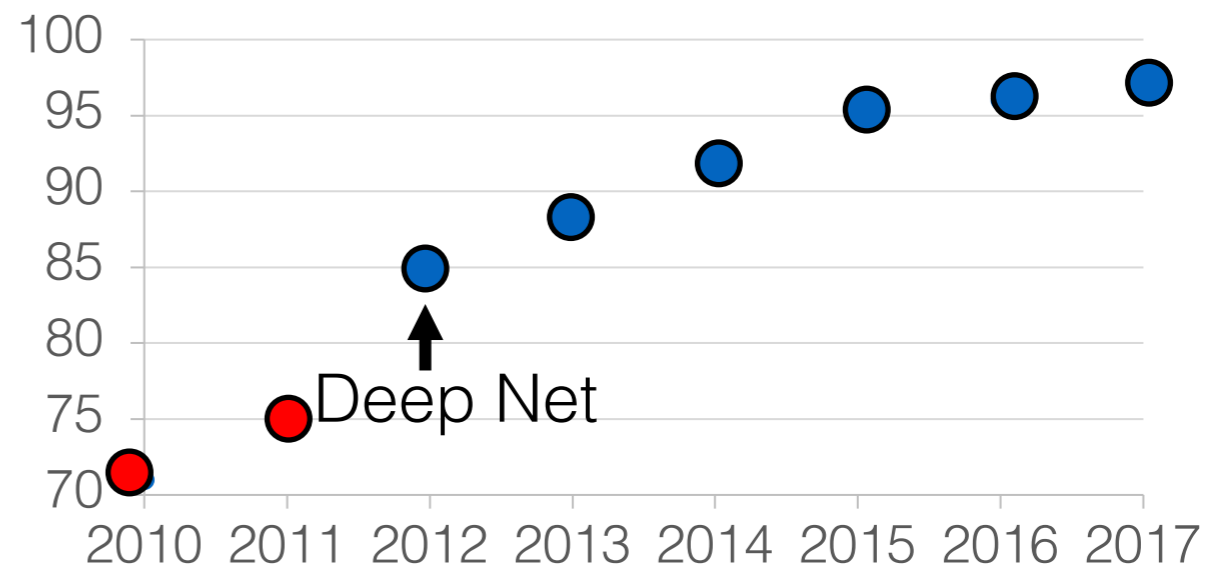
Cat

# Computer Vision Now

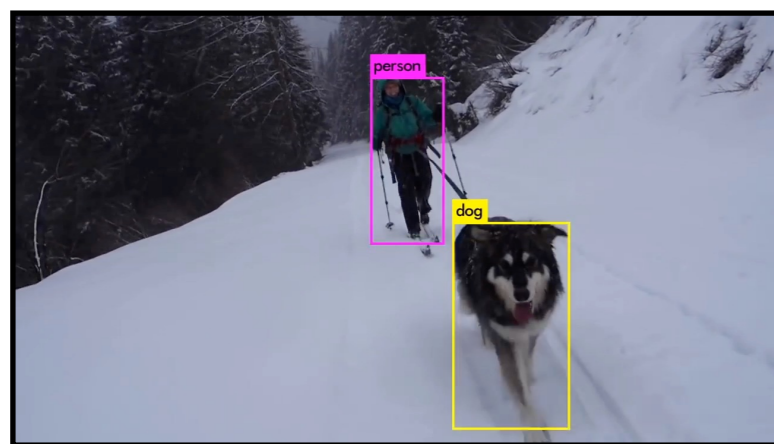


[LeCun et al, 1998], [Krizhevsky et al, 2012]

# Deep Learning for Computer Vision

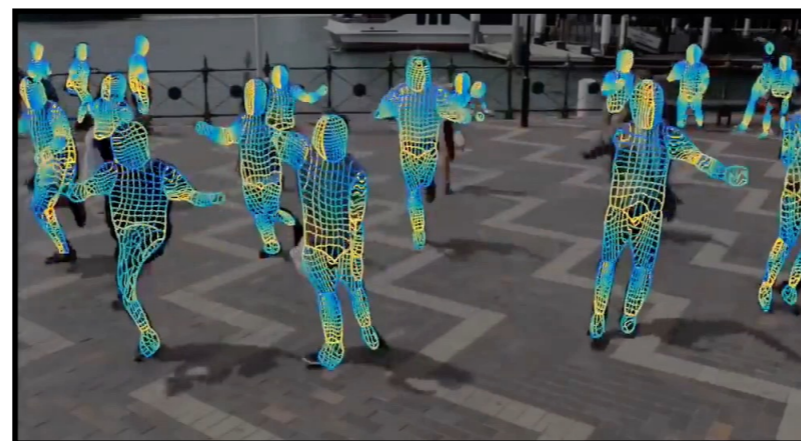


Top 5 accuracy on ImageNet benchmark



[Redmon et al., 2018]

Object detection



[Güler et al., 2018]

Human understanding

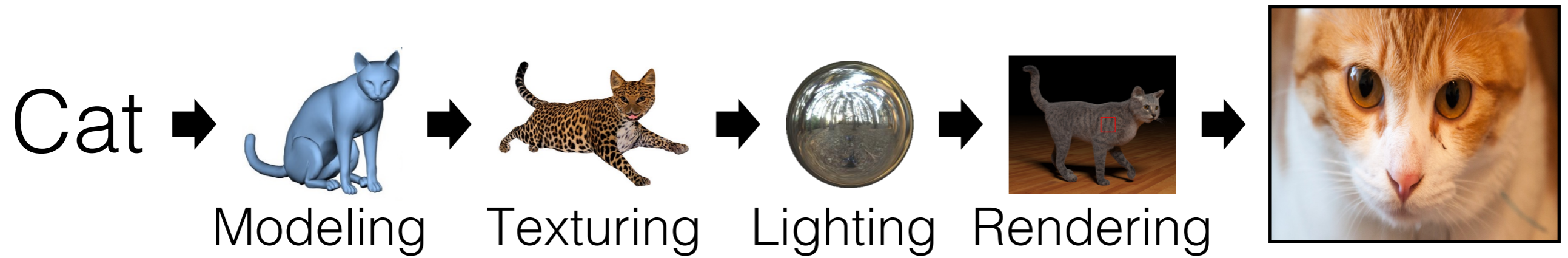


[Zhao et al., 2017]

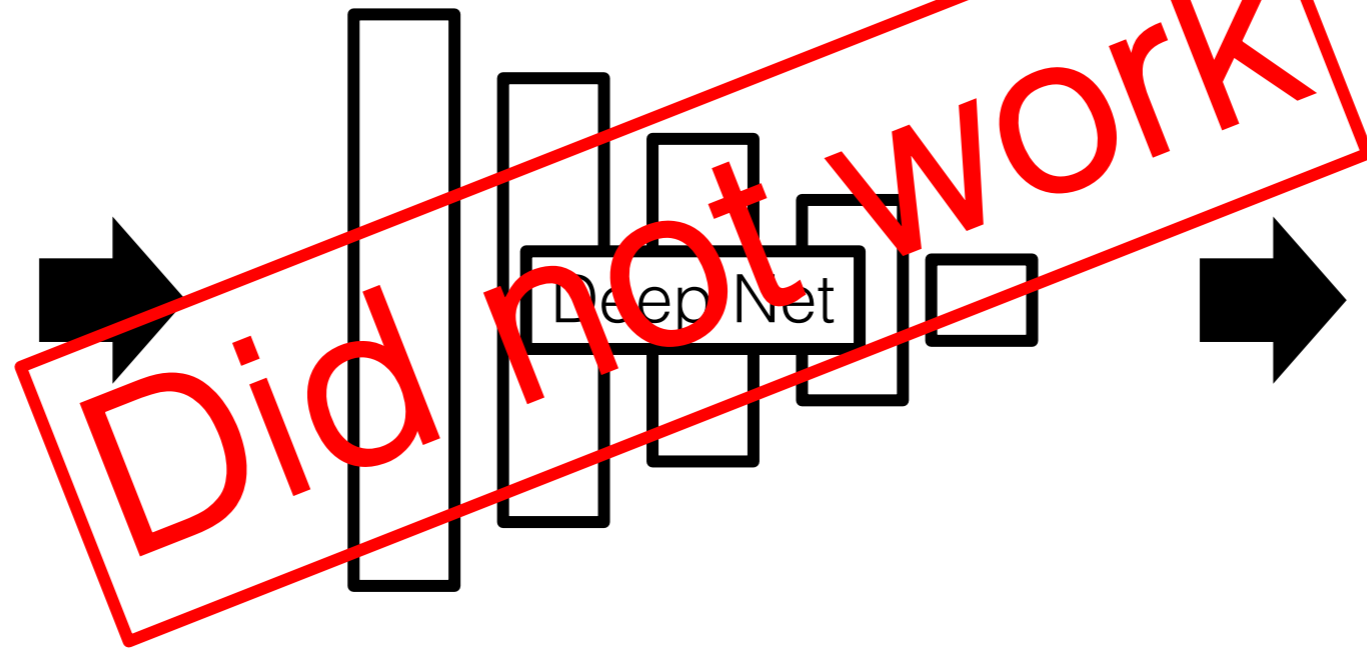
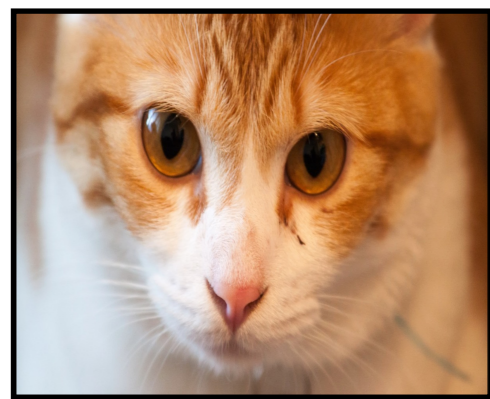
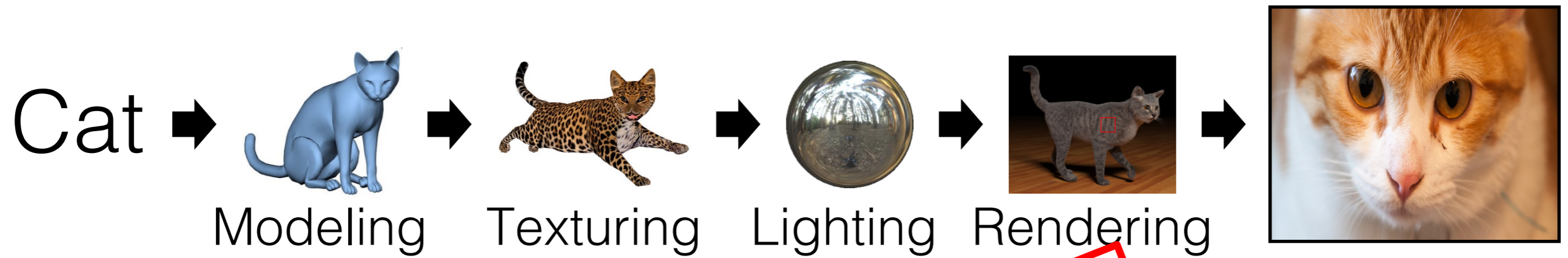
Autonomous driving



# Can Deep Learning Help Graphics?

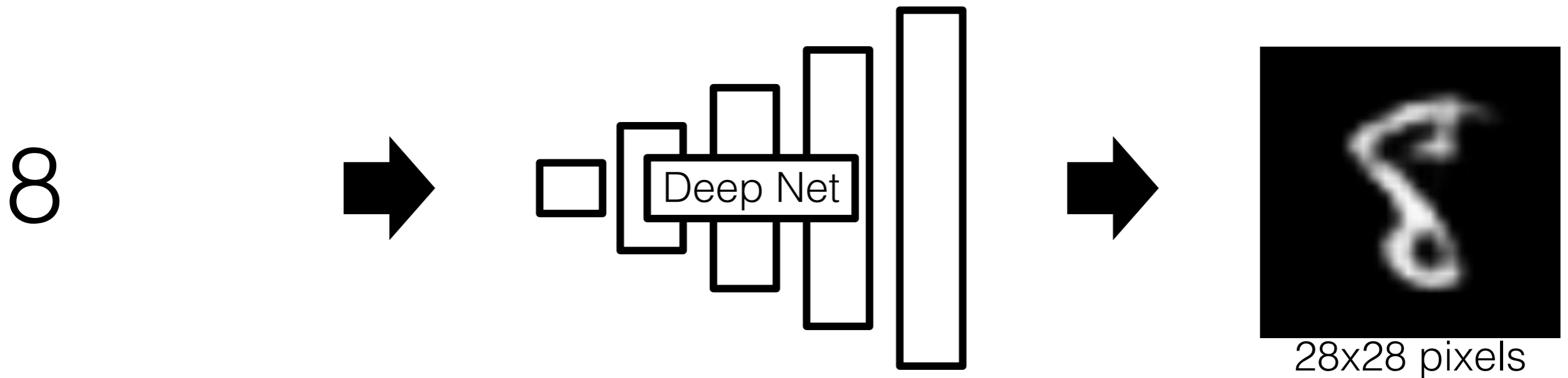
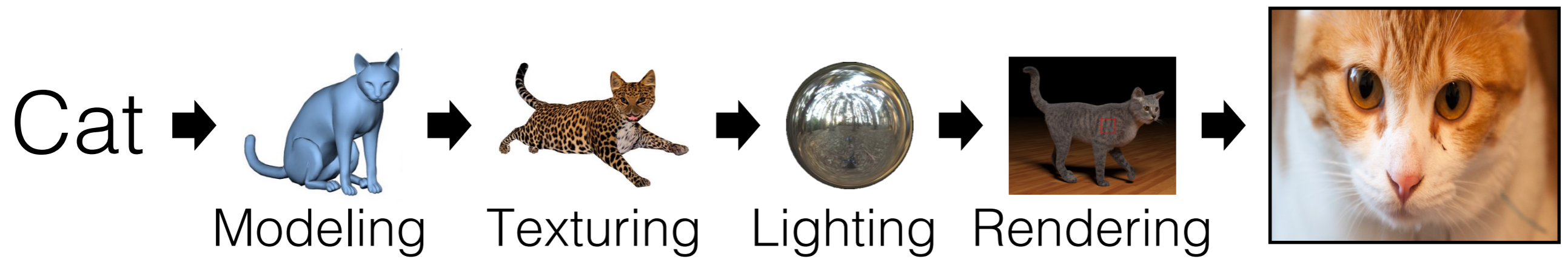


# Can Deep Learning Help Graphics?



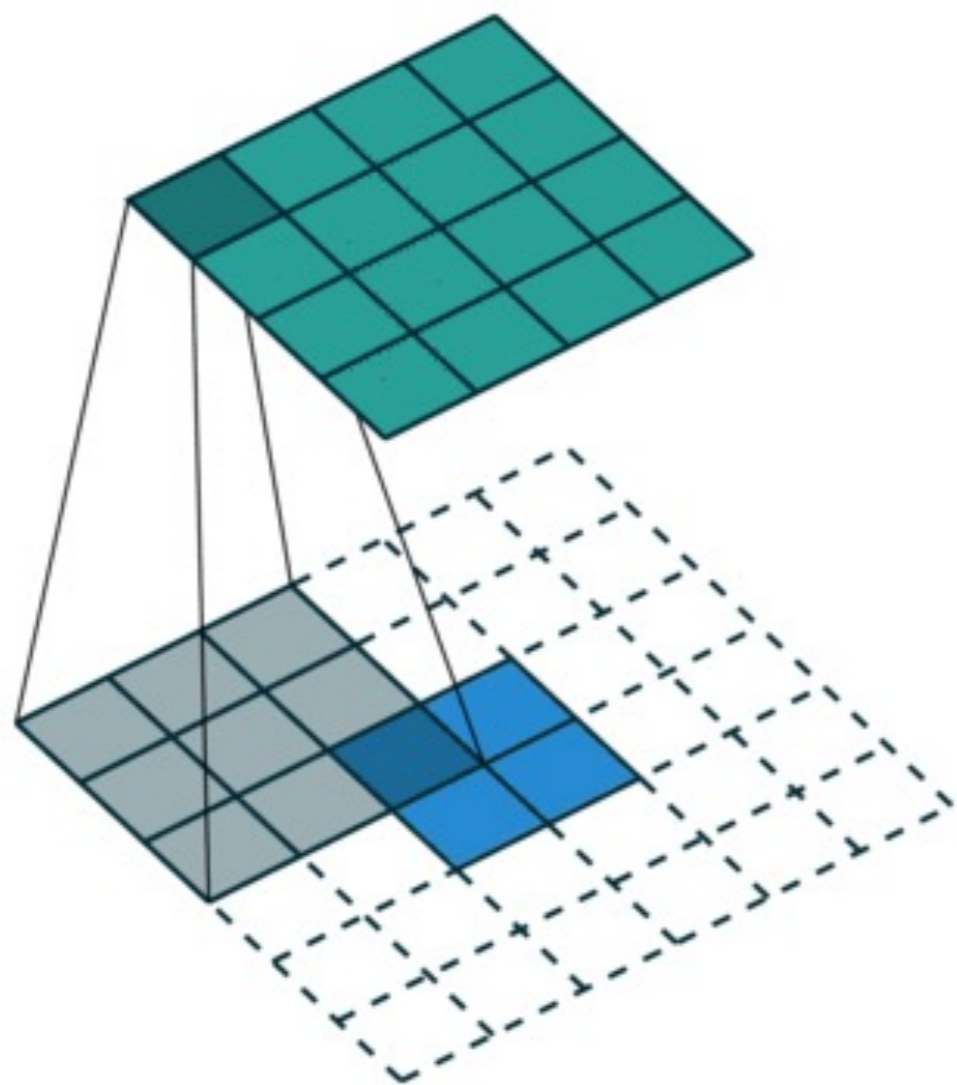
Cat

# Generating images is hard!

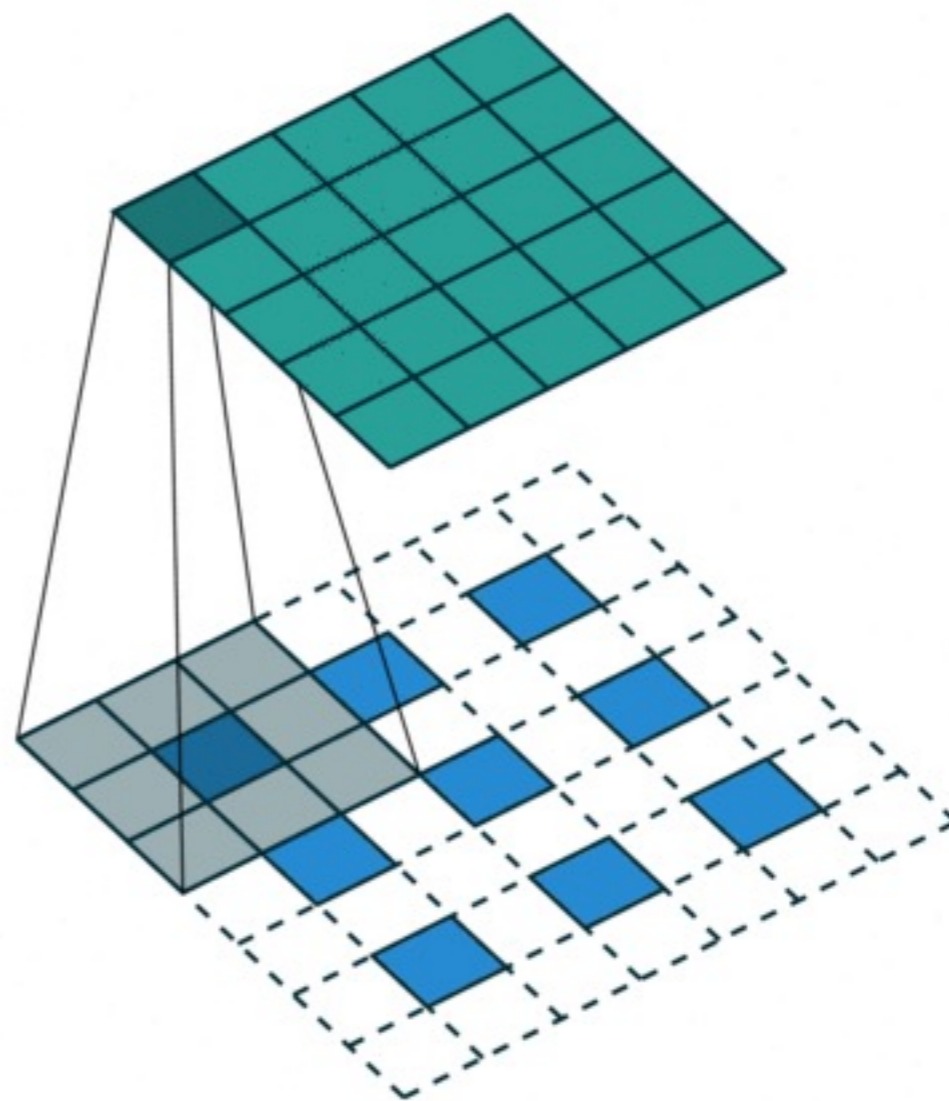


# Better Architectures

# Fractionally-strided Convolution

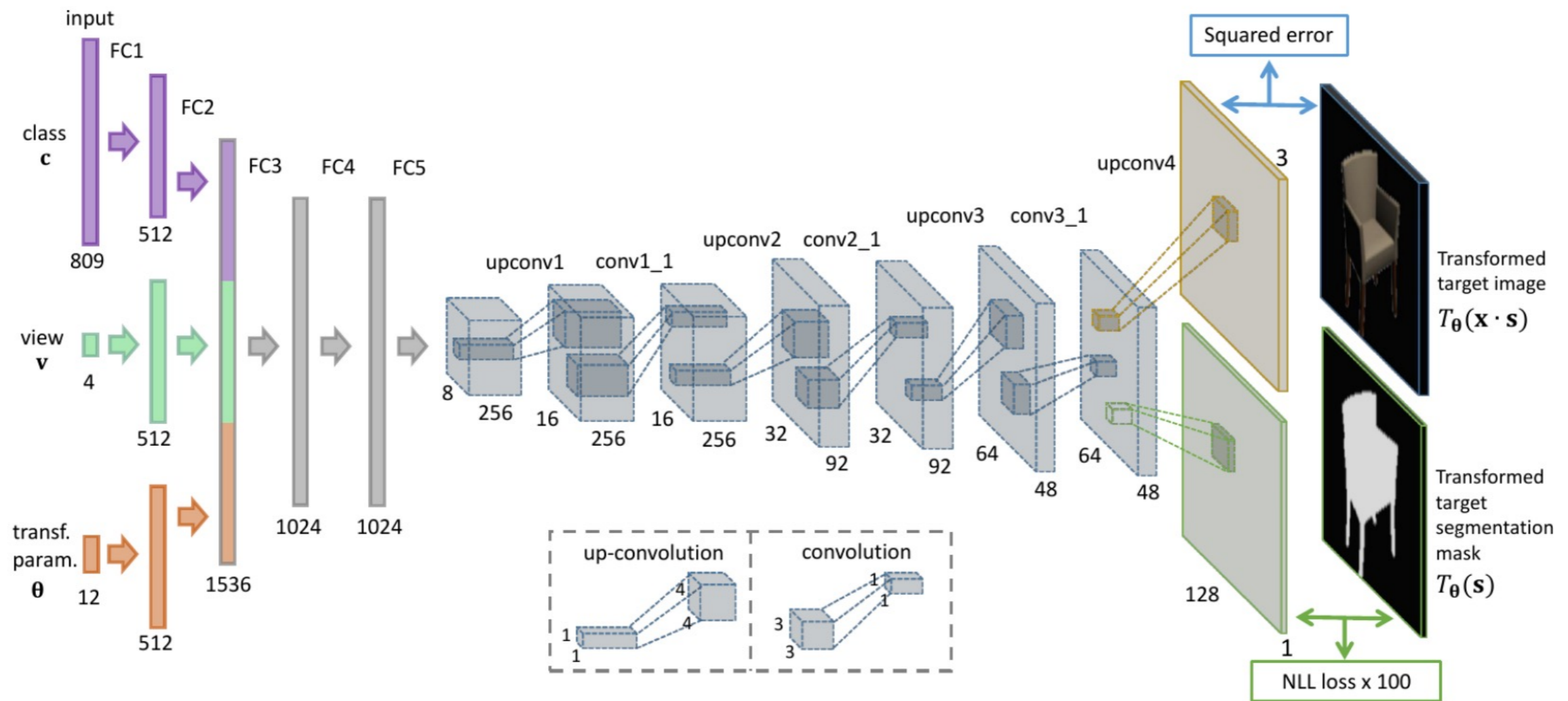


Regular conv



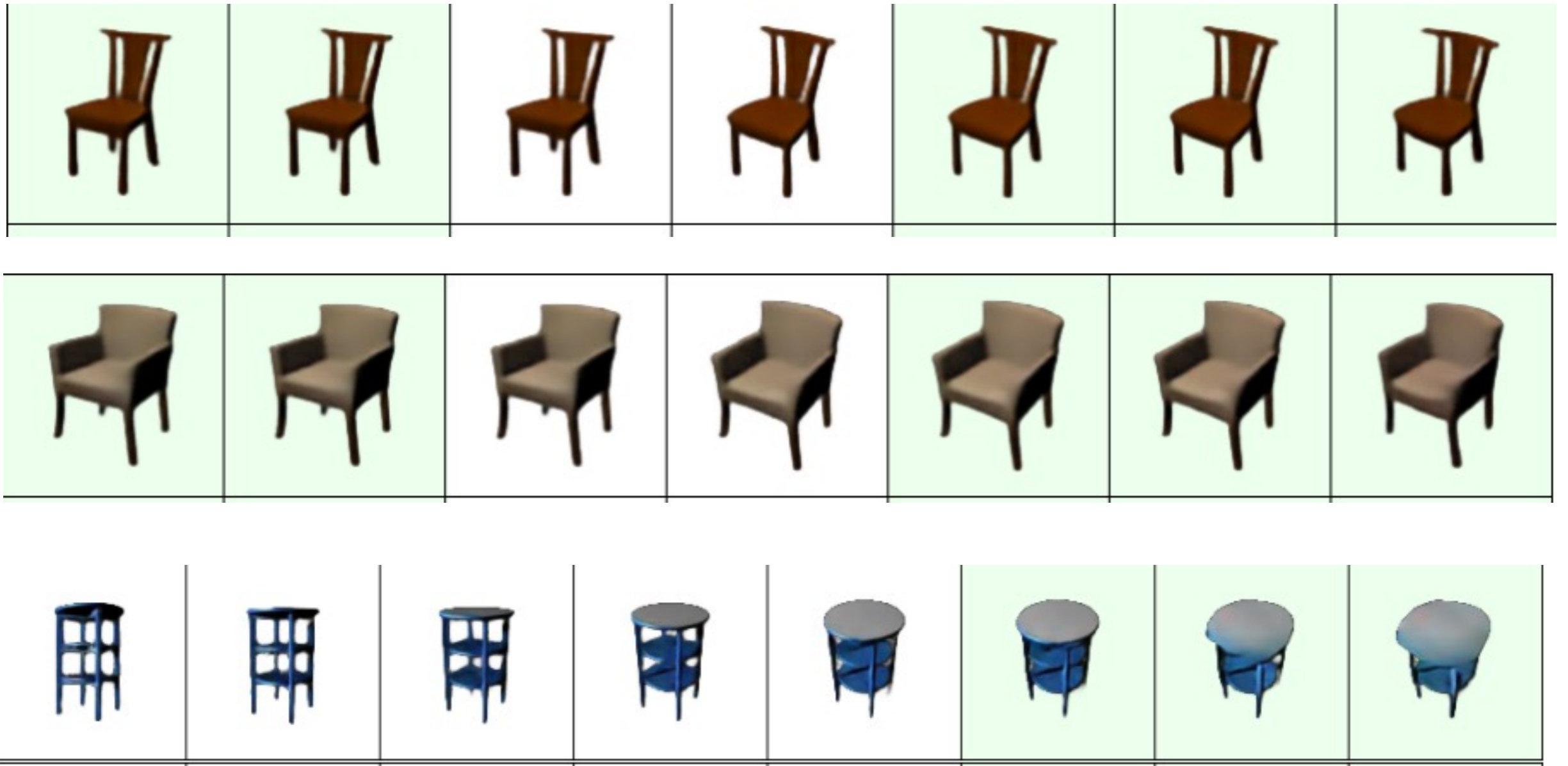
Fractionally-strided conv

# Generating chairs conditional on chair ID, viewpoint, and transformation parameters



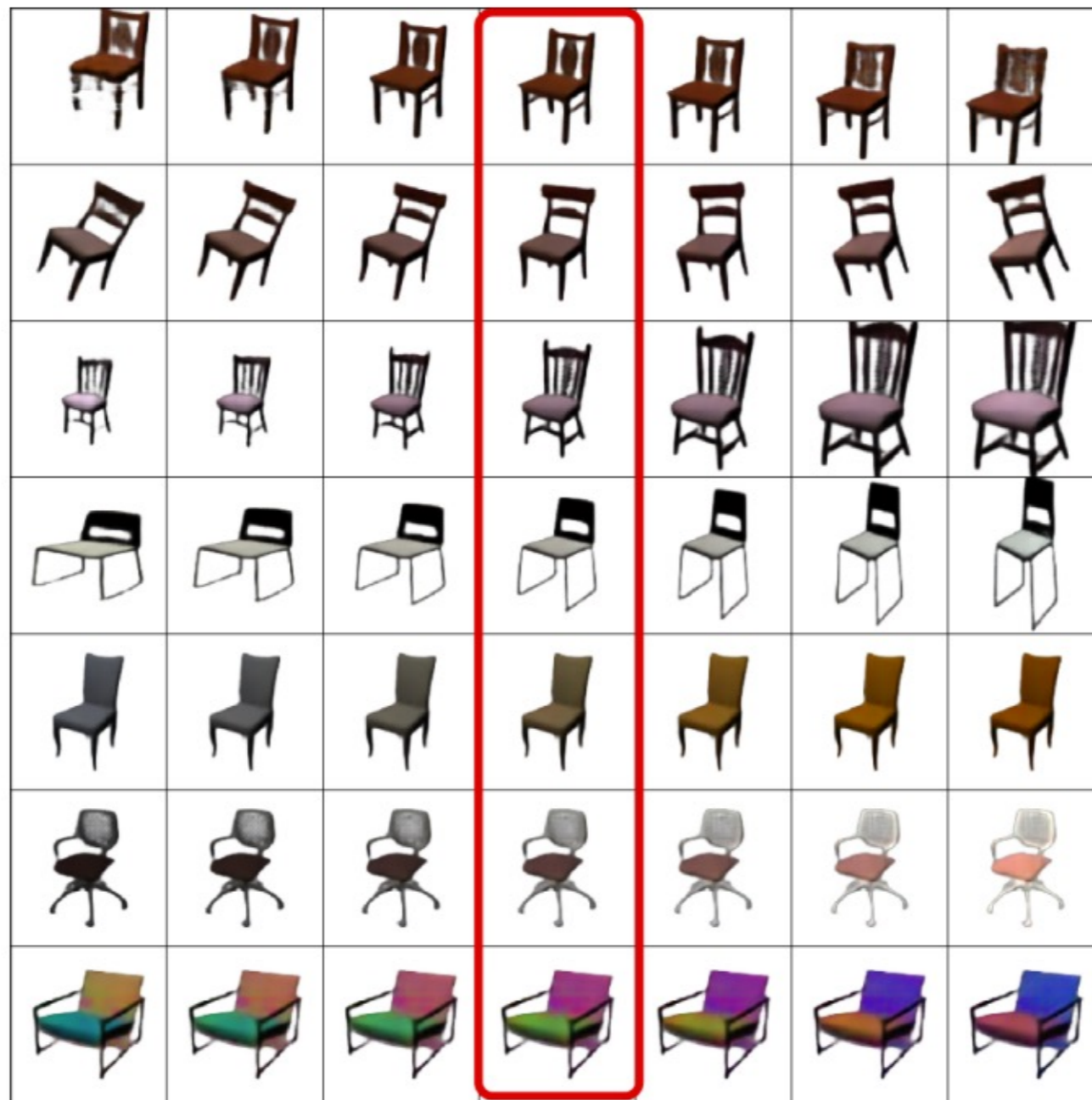
Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks  
PAMI 2017 (CVPR 2015)

# With Varying Viewpoints



Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks  
PAMI 2017 (CVPR 2015)

# With Varying Transformation Parameters





# Interpolation between Two Chairs



Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks  
PAMI 2017 (CVPR 2015)

# Better Loss Functions

# Simple L2 regression doesn't work 😞

Input



Output



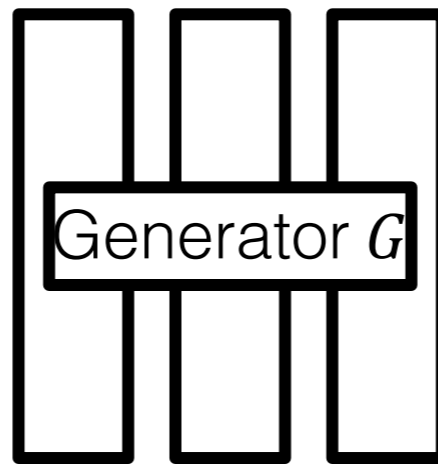
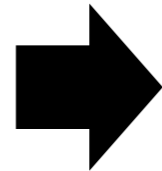
Ground truth



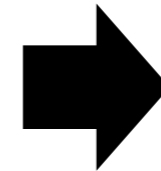
# Loss functions for Image Synthesis



Input  $x$



Learnable rendering



Output Image  $G(x)$

What is a good objective  $\mathcal{L}$ ?

- Capture realism
- Calculate image distance
- Adapt to new tasks/data.

Problem Statement

Loss function

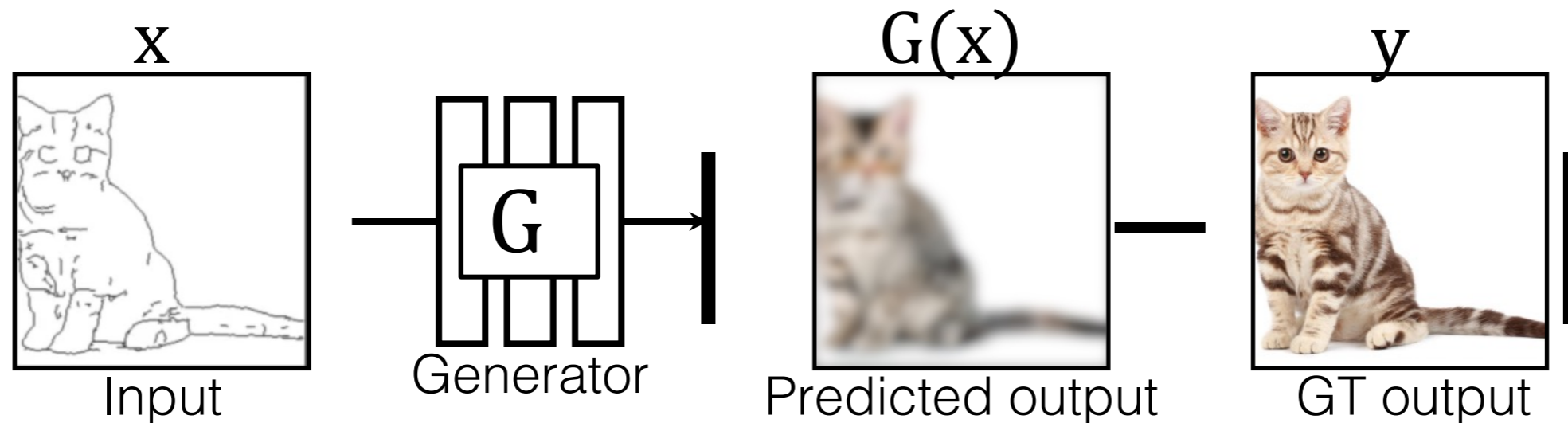
$$\arg \min_G \mathcal{L}(G(x), y)$$

Generator

Input

Output image

# Designing Loss Functions

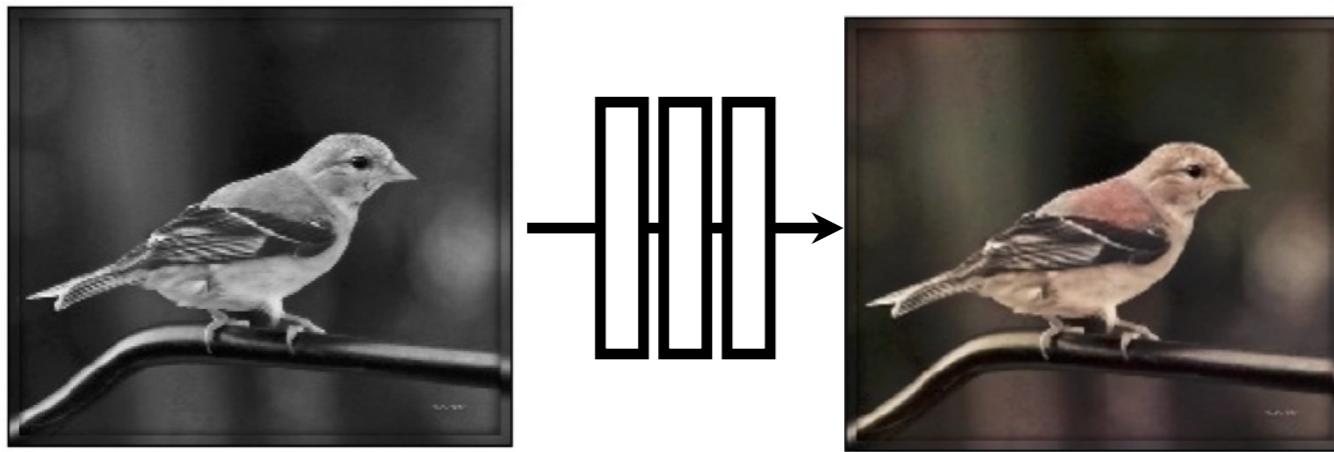


L2 regression

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

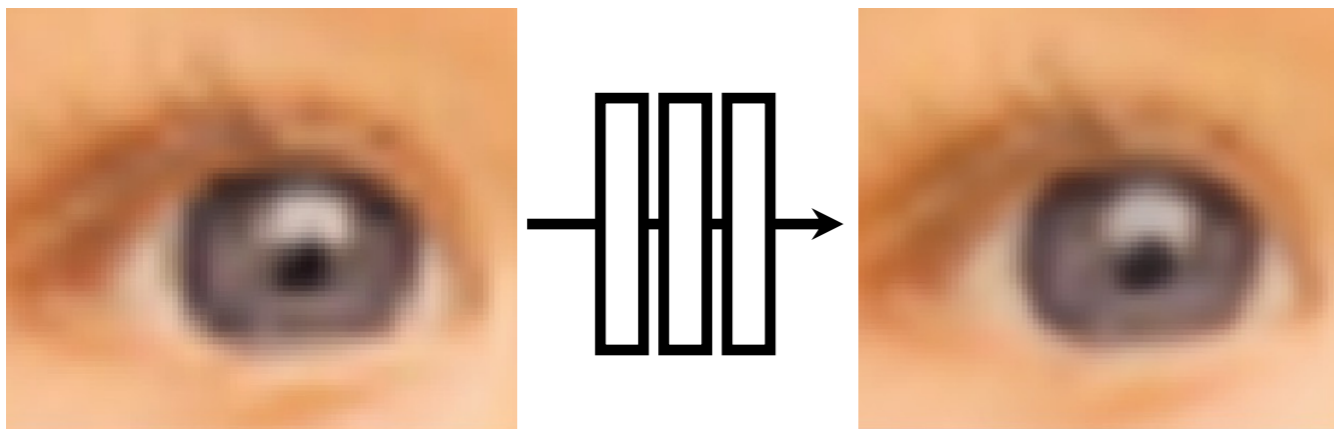
# Designing Loss Functions

Image colorization



L2 regression

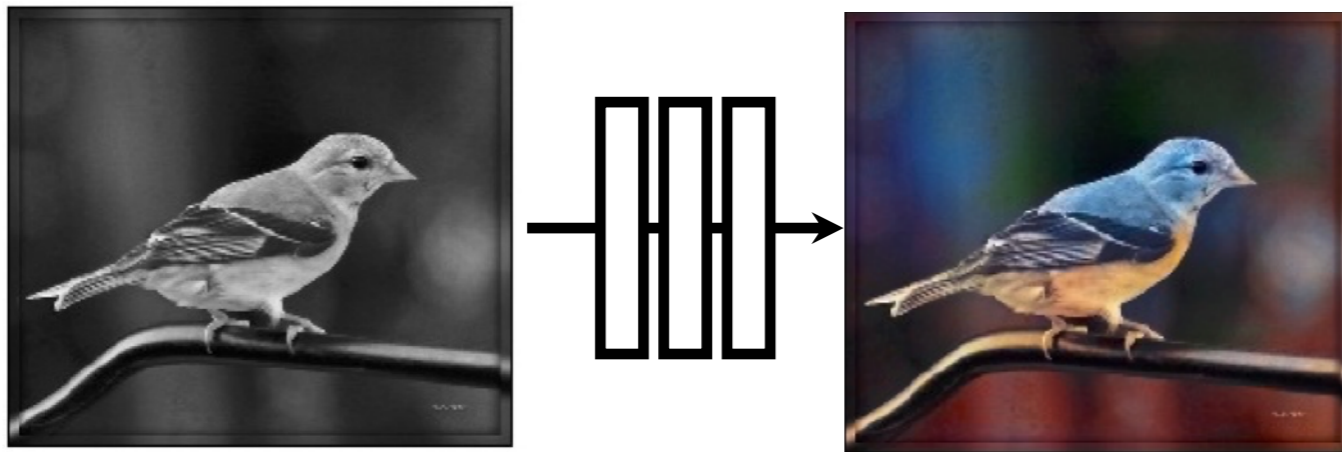
Super-resolution



L2 regression

# Designing Loss Functions

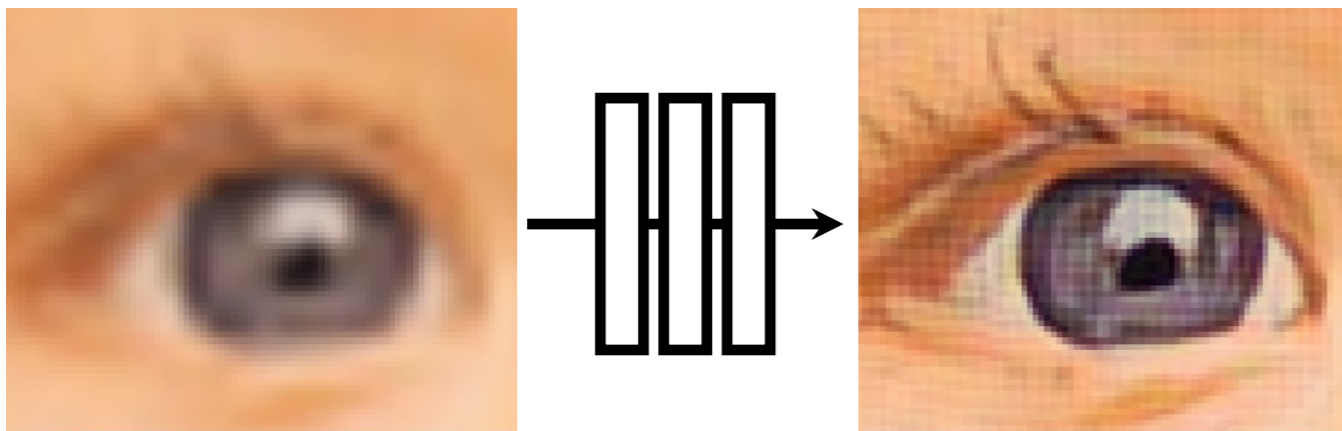
Image colorization



[Zhang et al. 2016]

Classification Loss:  
Cross entropy objective,  
with colorfulness term

Super-resolution

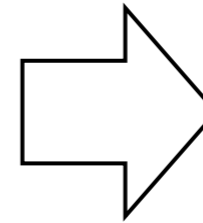


[Gatys et al., 2016], [Johnson et al. 2016]  
[Dosovitskiy and Brox. 2016]

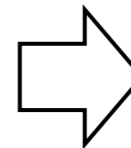
Feature/Perceptual loss  
Deep feature matching  
objective

# “Perceptual Loss”

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

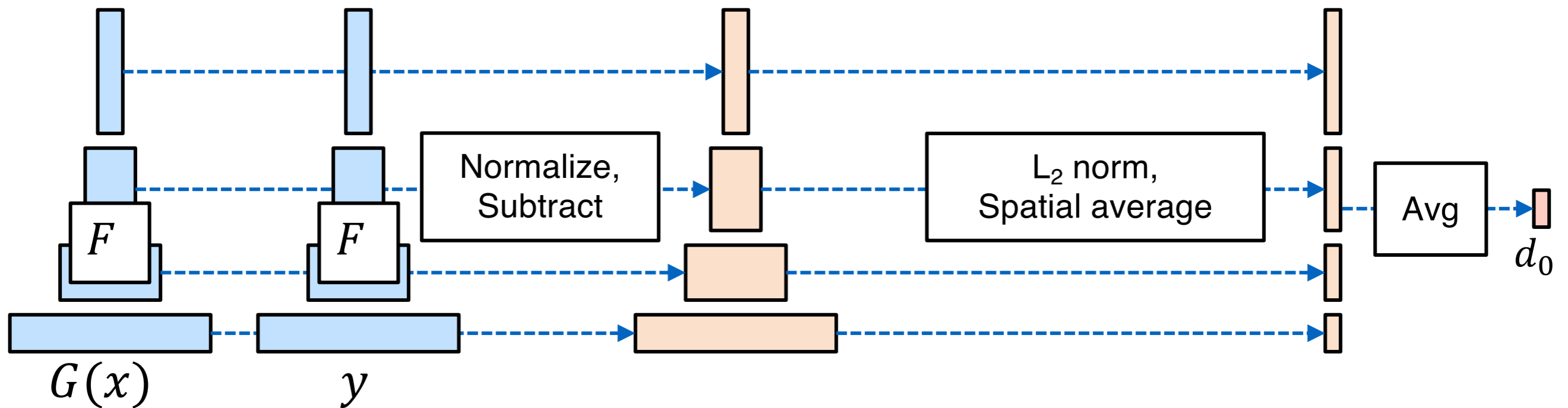


Chen and Koltun. In ICCV, 2017.





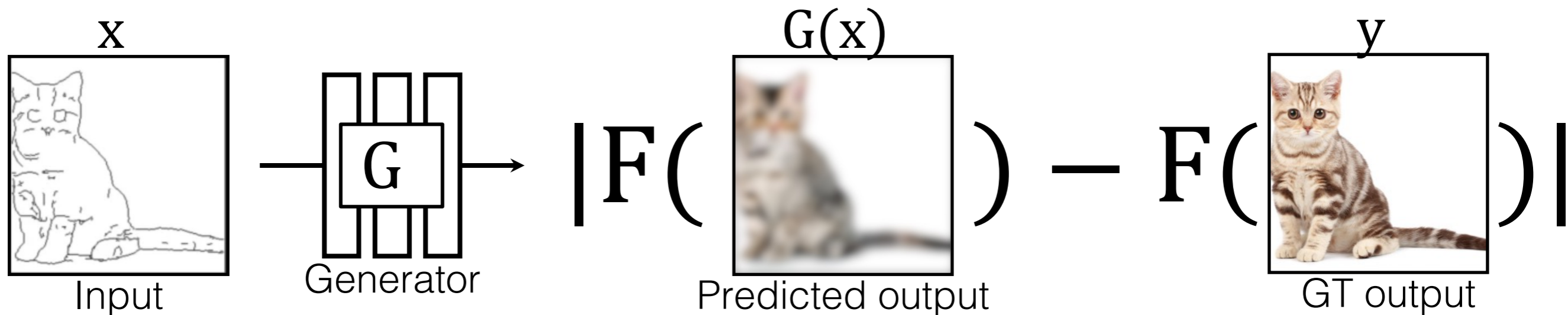
# CNNs as a Perceptual Metric



(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



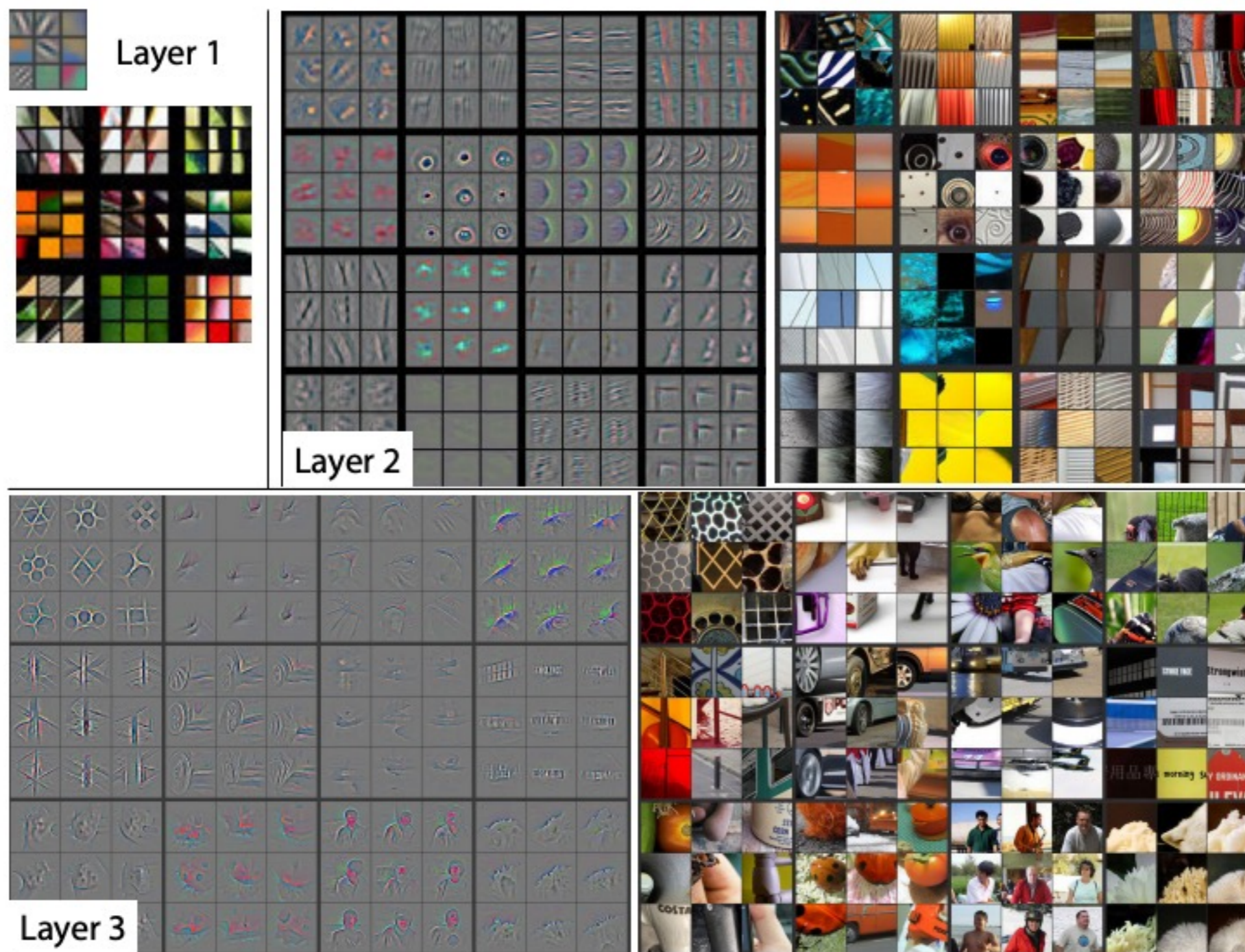
$F$  is a deep network (e.g., ImageNet classifier)

## Perceptual Loss

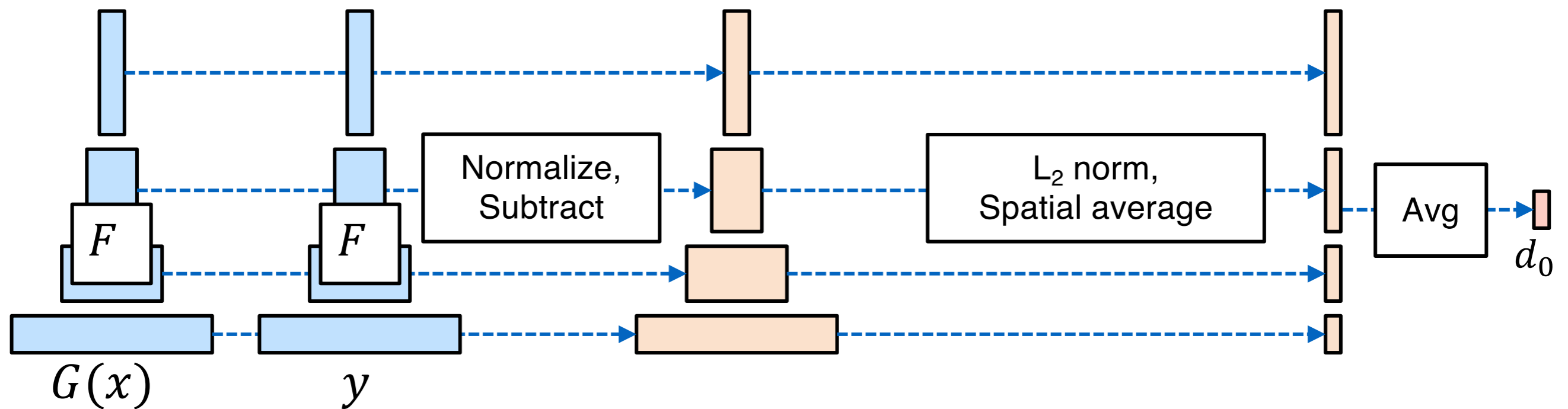
$$\arg \min_G \mathbb{E}_{(x,y)} \sum_{i=1}^N \overset{\text{weight}}{\lambda_i} \frac{1}{M_i} \left\| F^{(i)}(G(x)) - F^{(i)}(y) \right\|_2^2$$

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

# What has a CNN Learned?



# CNNs as a Perceptual Metric



Perceptual Loss

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

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

# How Different are these Patches?



Zhang, Isola, Efros, Shechtman, Wang.

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

Slide credit: Richard Zhang

Which patch is more similar to the middle?



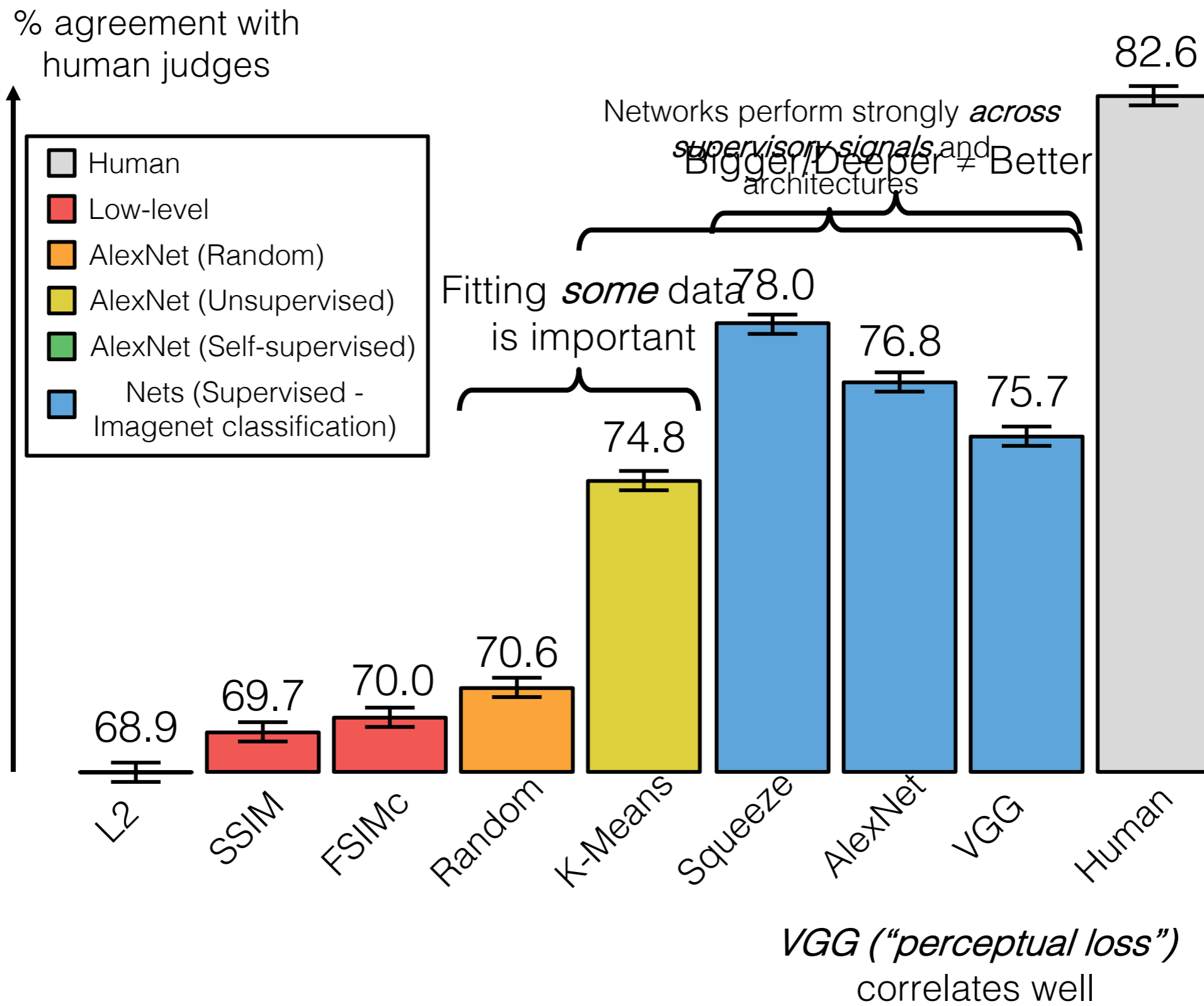
< Type 1 >



Humans  
L2/PSNR  
SSIM/FSIMc  
*Deep Networks?*

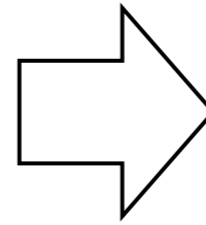


< Type 2 >

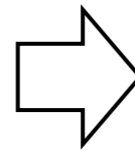


# “Perceptual Loss”

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

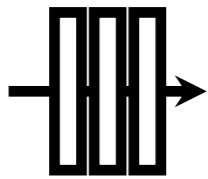
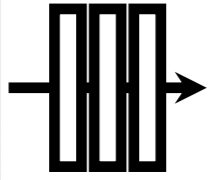


Chen and Koltun. In ICCV, 2017.





Generated images



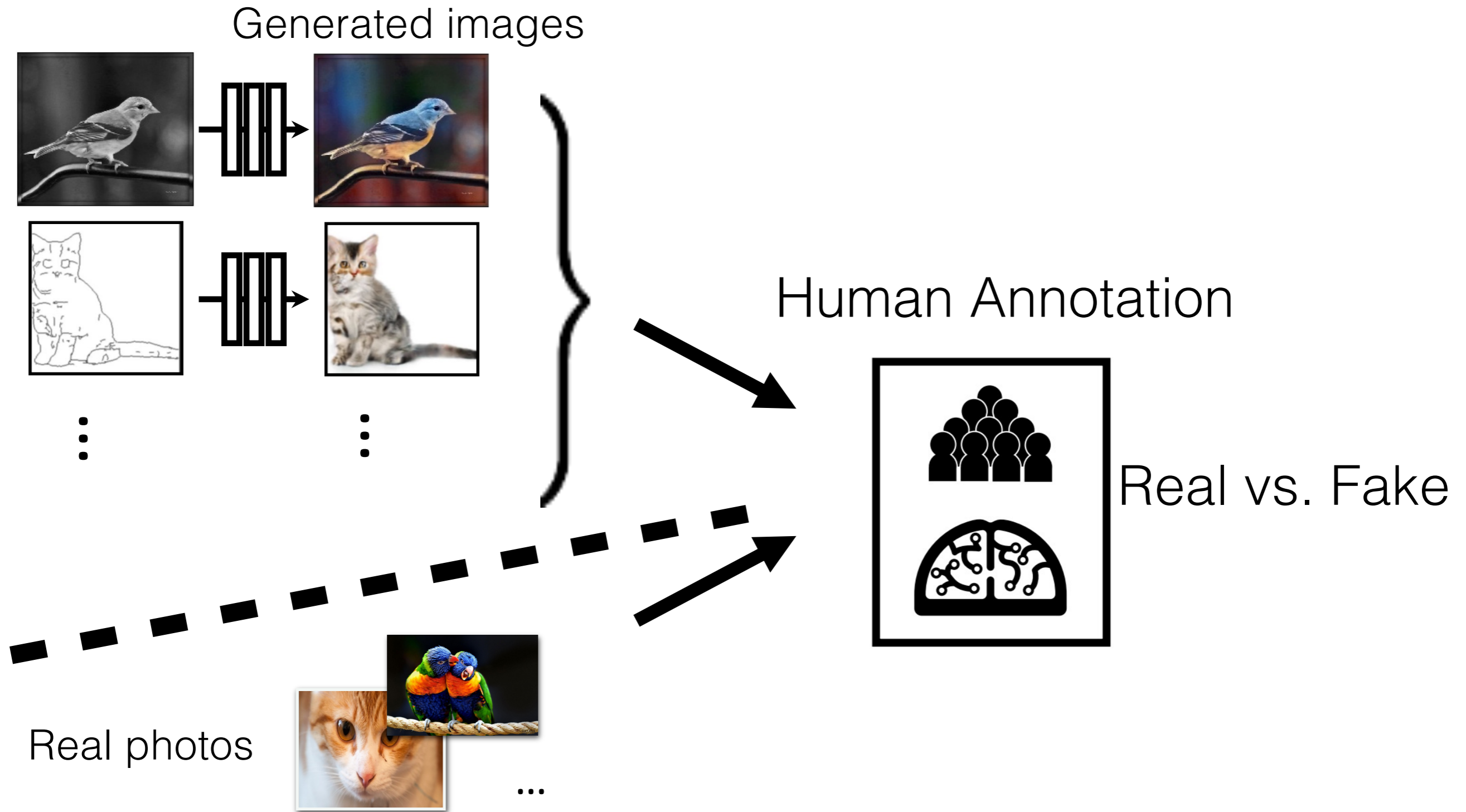
⋮

⋮

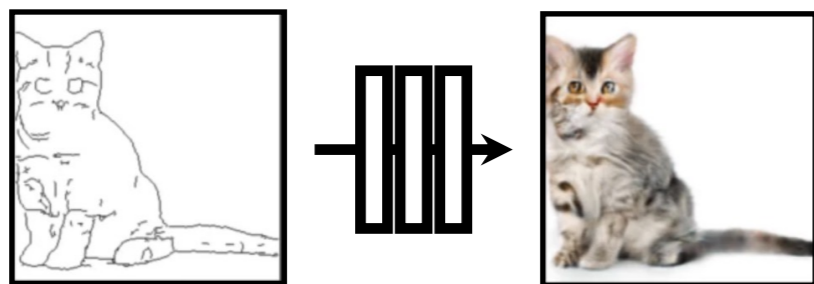
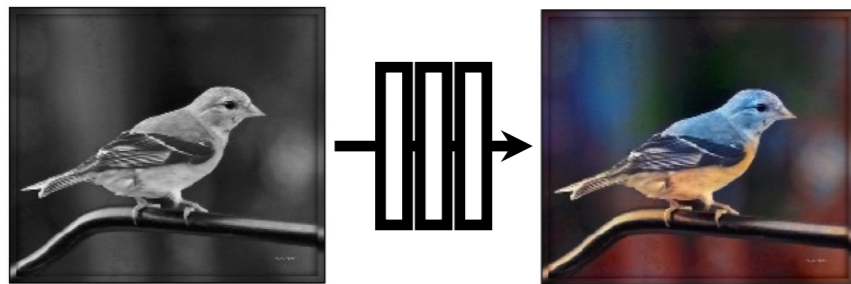


Universal loss?

# Learning with Human Perception



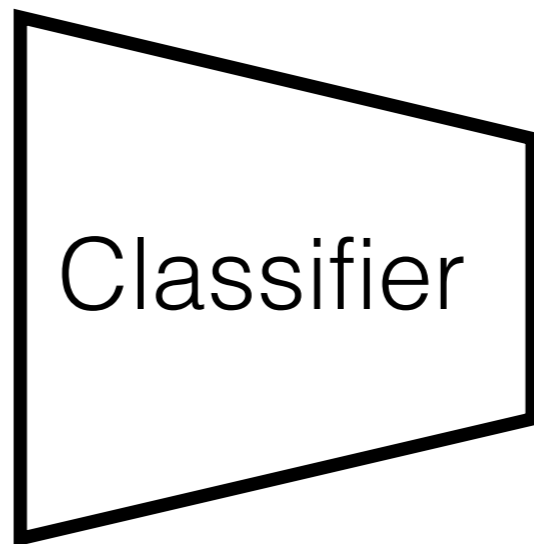
Generated images



⋮

⋮

# Generative Adversarial Network (GANs)



Classifier Real vs. Fake



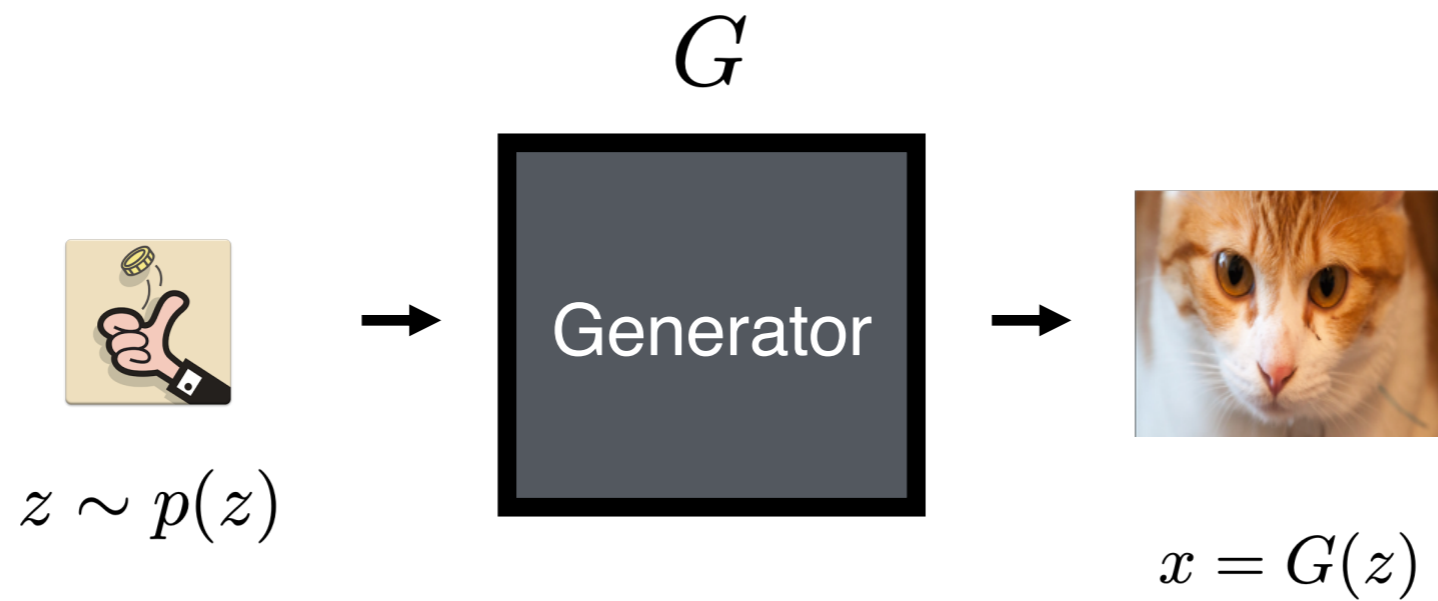
Real photos



...

[Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, Bengio 2014]

# Image synthesis from “noise”



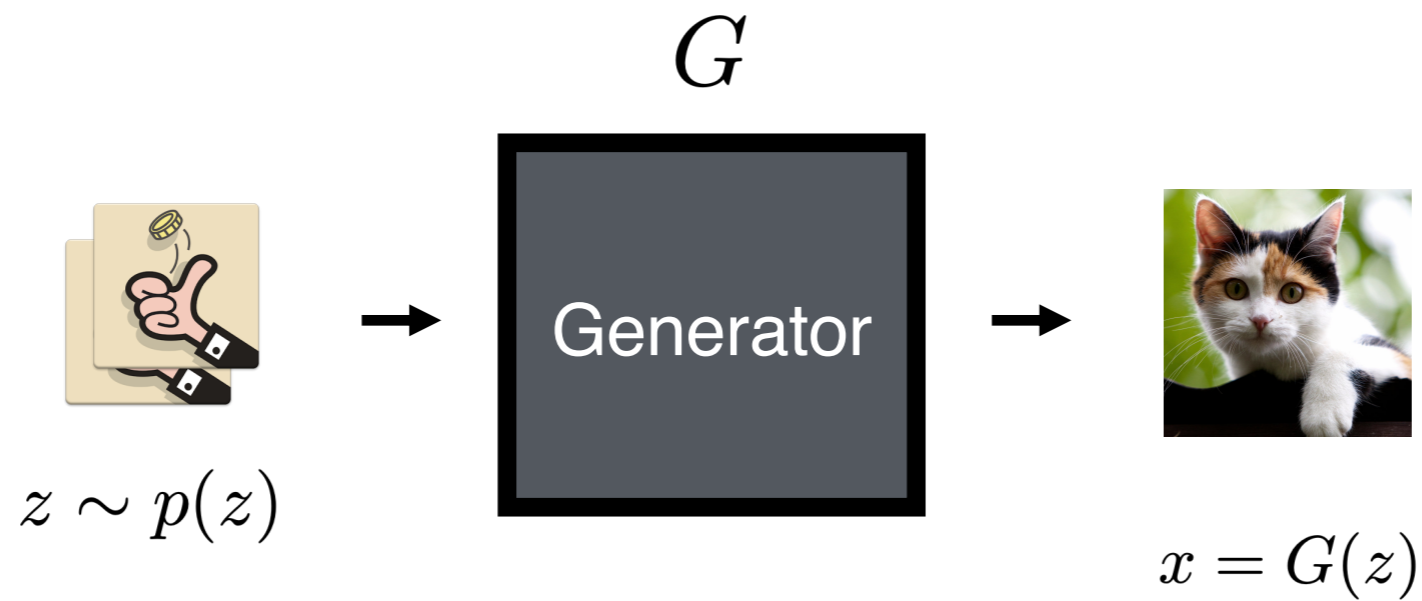
Sampler

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

$$z \sim p(z)$$

$$x = G(z)$$

# Image synthesis from “noise”



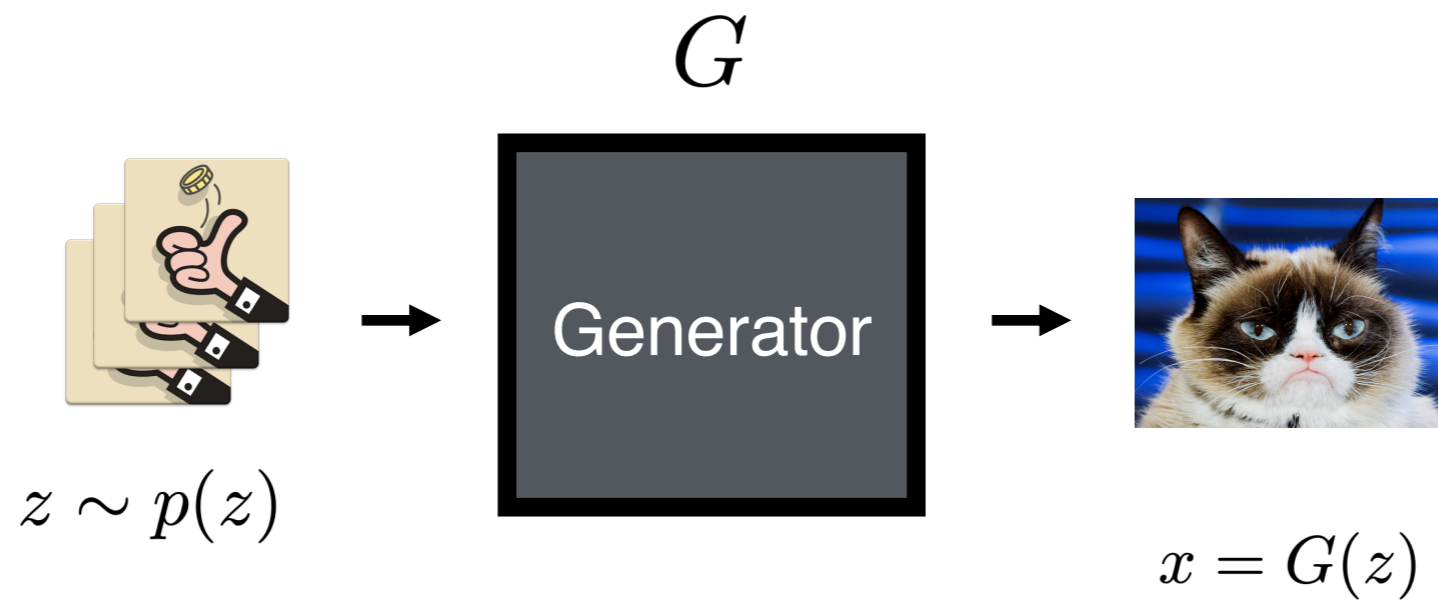
Sampler

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

$$z \sim p(z)$$

$$x = G(z)$$

# Image synthesis from “noise”

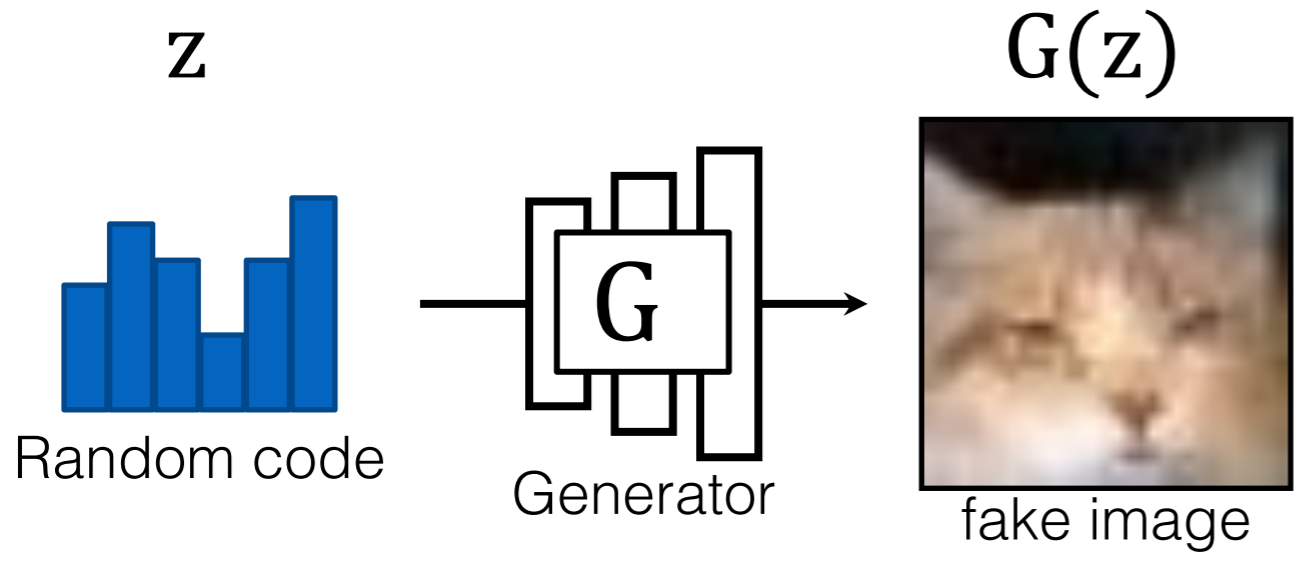


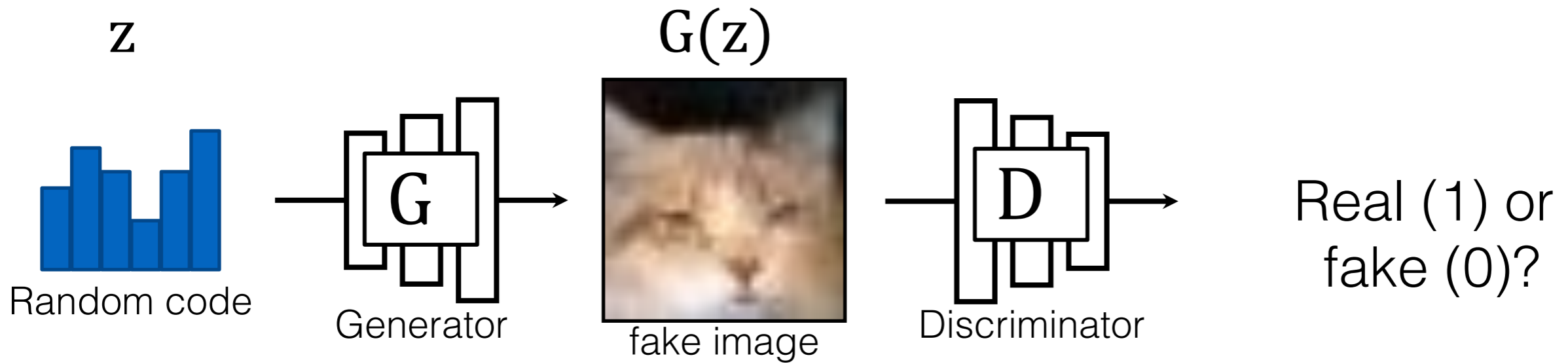
Sampler

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

$$z \sim p(z)$$

$$x = G(z)$$

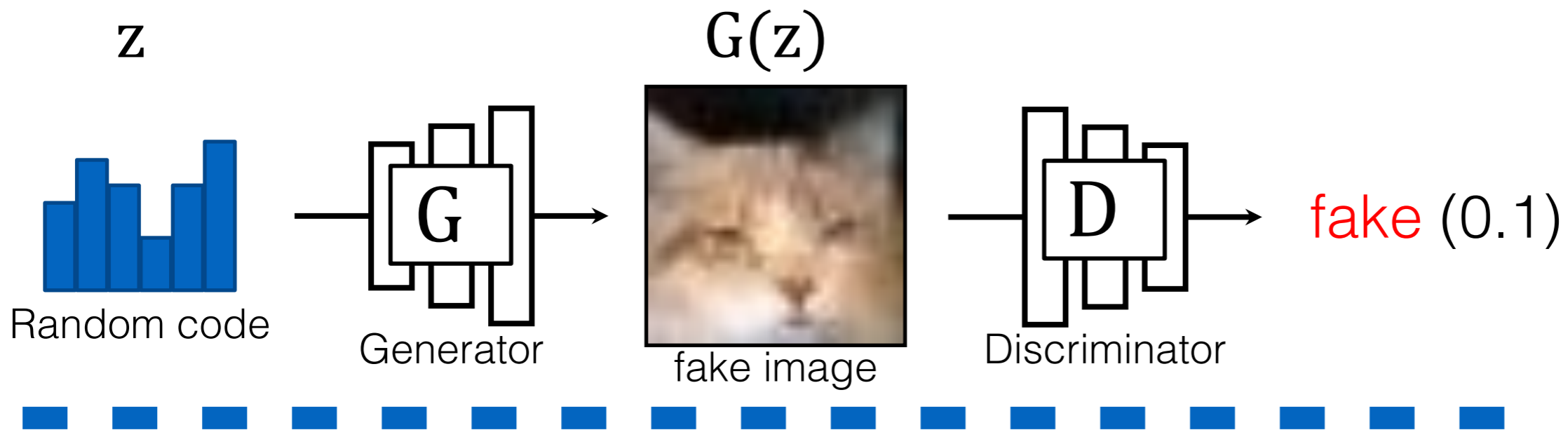




A two-player game:

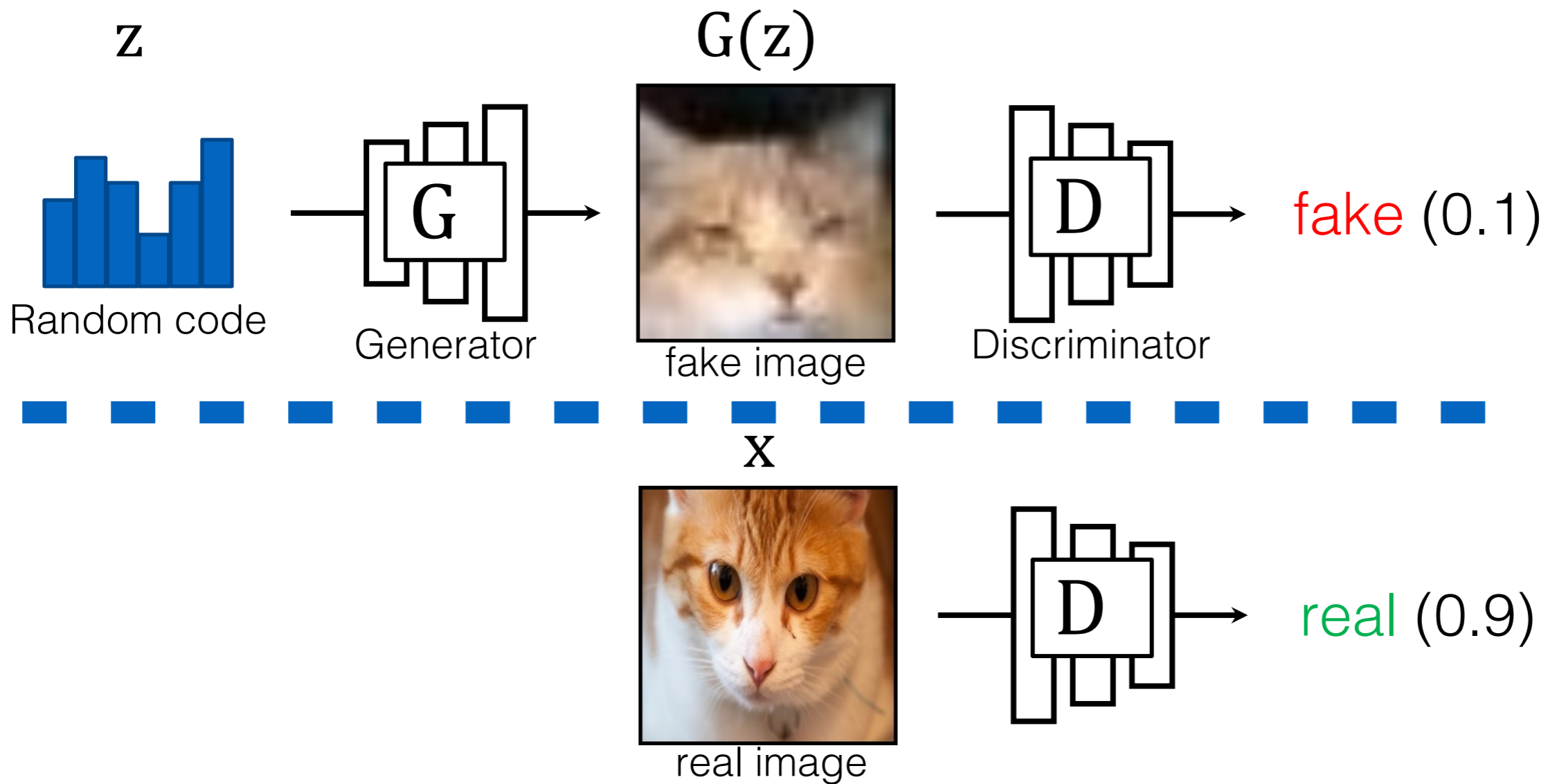
- $G$  tries to generate fake images that can fool  $D$ .
- $D$  tries to detect fake images.





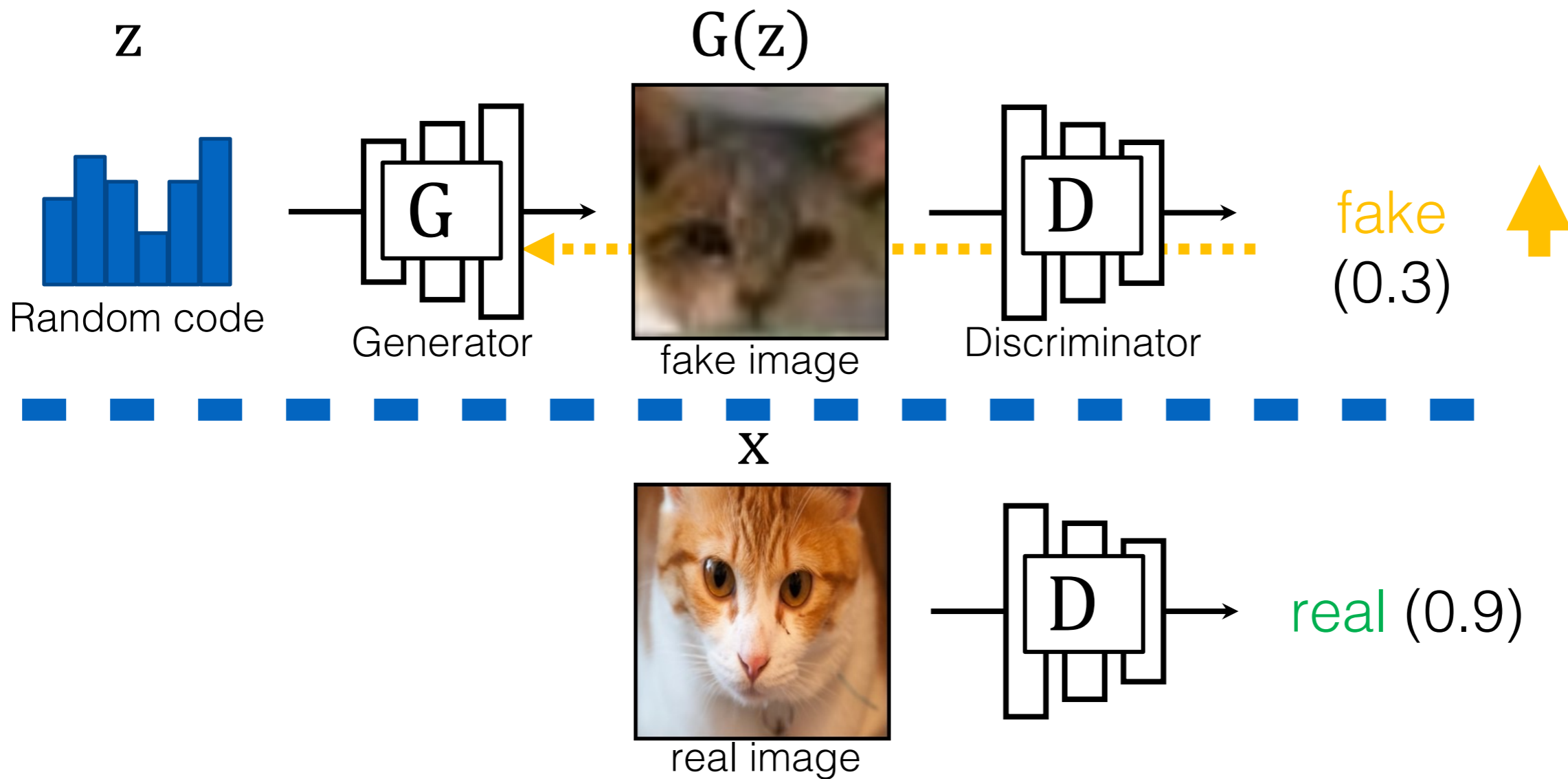
Learning objective (GANs)

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



Learning objective (GANs)

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



Learning objective (GANs)

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

# GANs Training Breakdown

- From the discriminator  $D$ 's perspective:
  - binary classification: real vs. fake.
  - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_D \mathbb{E}[\log(1 - D(\text{cat}))] + \mathbb{E}[\log D(\text{cat})]$$

# GANs Training Breakdown

- From the discriminator  $D$ 's perspective:
  - binary classification: real vs. fake.
  - Nothing special: similar to 1 vs. 7 or cat vs. dog

$$\max_D \mathbb{E}[\log(1 - D(\text{dog}))] + \mathbb{E}[\log D(\text{cat})]$$

- From the generator  $G$ 's perspective:
  - Optimizing a loss that depends on a classifier  $D$
  - We have done it before (Perceptual Loss)

$$\min_G \mathbb{E}_z[\mathcal{L}_D(G(z))]$$

GAN loss for  $G$

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

Perceptual Loss for  $G$

# GANs Training Breakdown



**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.

$p_g = p_{data}$  is the unique global minimizer of the GAN objective.

Proof                      Optimal discriminator given fixed G

$$\begin{aligned}
 C(G) &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))] \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]
 \end{aligned}$$

$$C(G) = -\log(4) + KL \left( p_{data} \left\| \frac{p_{data} + p_g}{2} \right. \right) + KL \left( p_g \left\| \frac{p_{data} + p_g}{2} \right. \right)$$

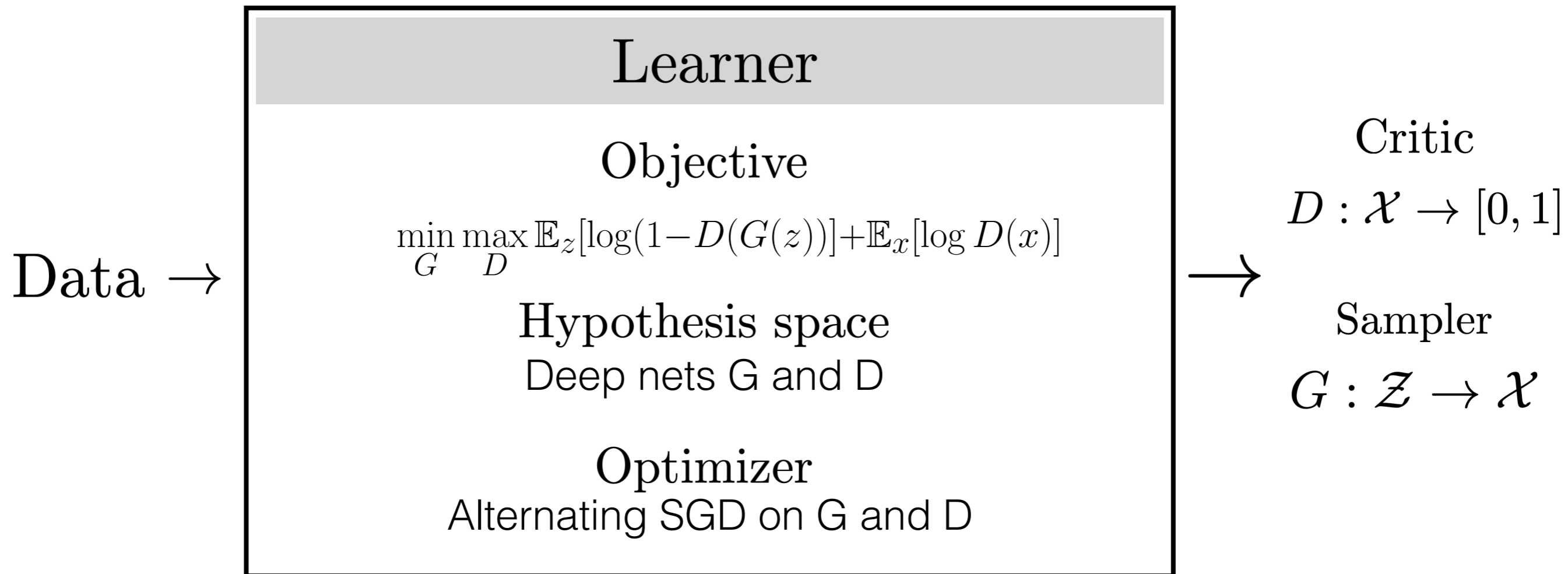
$$C(G) = -\log(4) + 2 \cdot \underbrace{JSD(p_{data} \| p_g)}$$

$$\geq 0, \quad 0 \iff p_g = p_{data} \quad \square$$

KLD (Kullback–Leibler divergence):  $\mathcal{KL}(p \| q) = \int p(x) \log \frac{p(x)}{q(x)} dx$

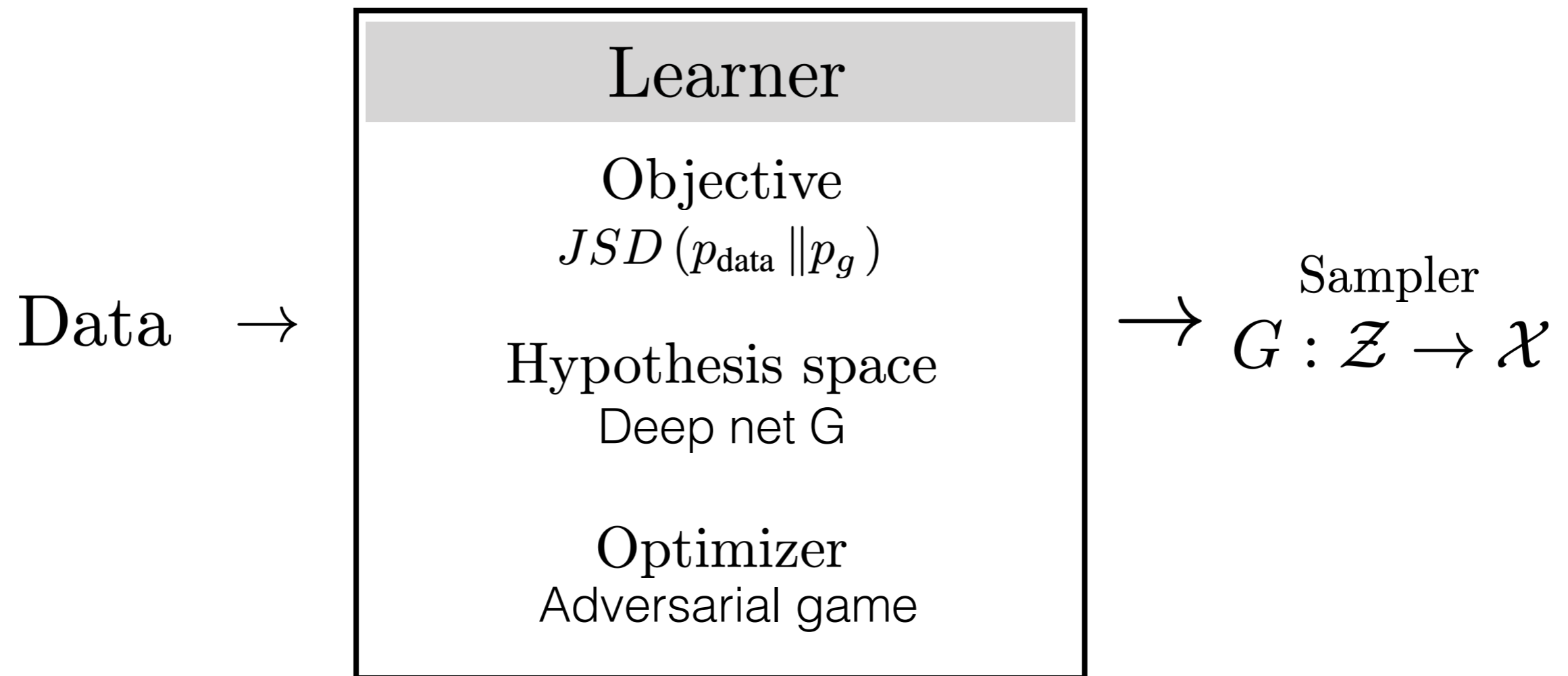
JSD (Jensen–Shannon divergence):  $\mathcal{JSD}(p \| q) = \frac{1}{2} \mathcal{KL}(p \| \frac{p+q}{2}) + \frac{1}{2} \mathcal{KL}(q \| \frac{p+q}{2})$

# Generative Adversarial Network





# Generative Adversarial Network



# What has driven GAN progress?



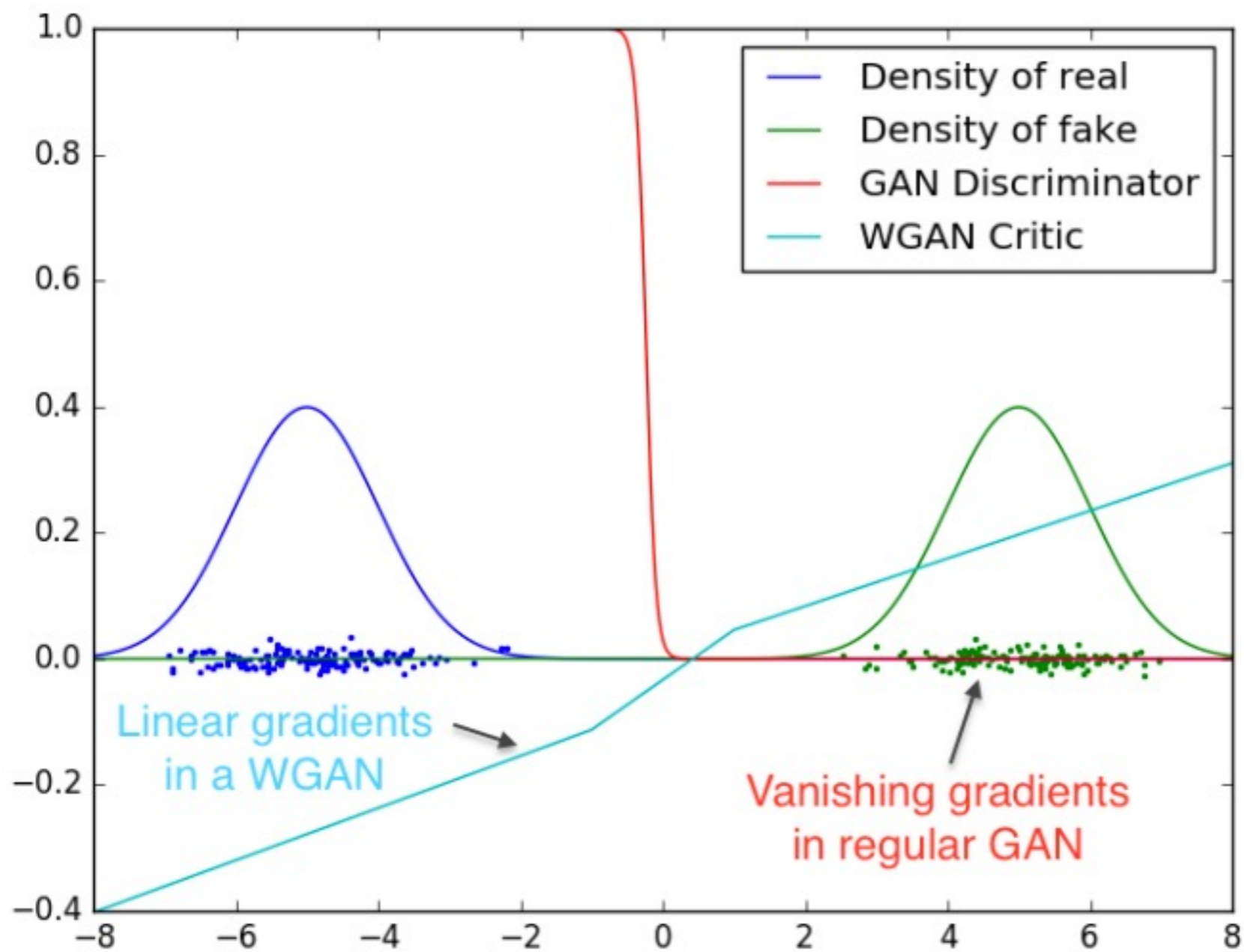
**Ian Goodfellow** @goodfellow\_ian · Jan 14

4.5 years of **GAN progress** on face generation. [arxiv.org/abs/1406.2661](https://arxiv.org/abs/1406.2661)

[arxiv.org/abs/1511.06434](https://arxiv.org/abs/1511.06434) [arxiv.org/abs/1606.07536](https://arxiv.org/abs/1606.07536) [arxiv.org/abs/1710.10196](https://arxiv.org/abs/1710.10196)

[arxiv.org/abs/1812.04948](https://arxiv.org/abs/1812.04948)

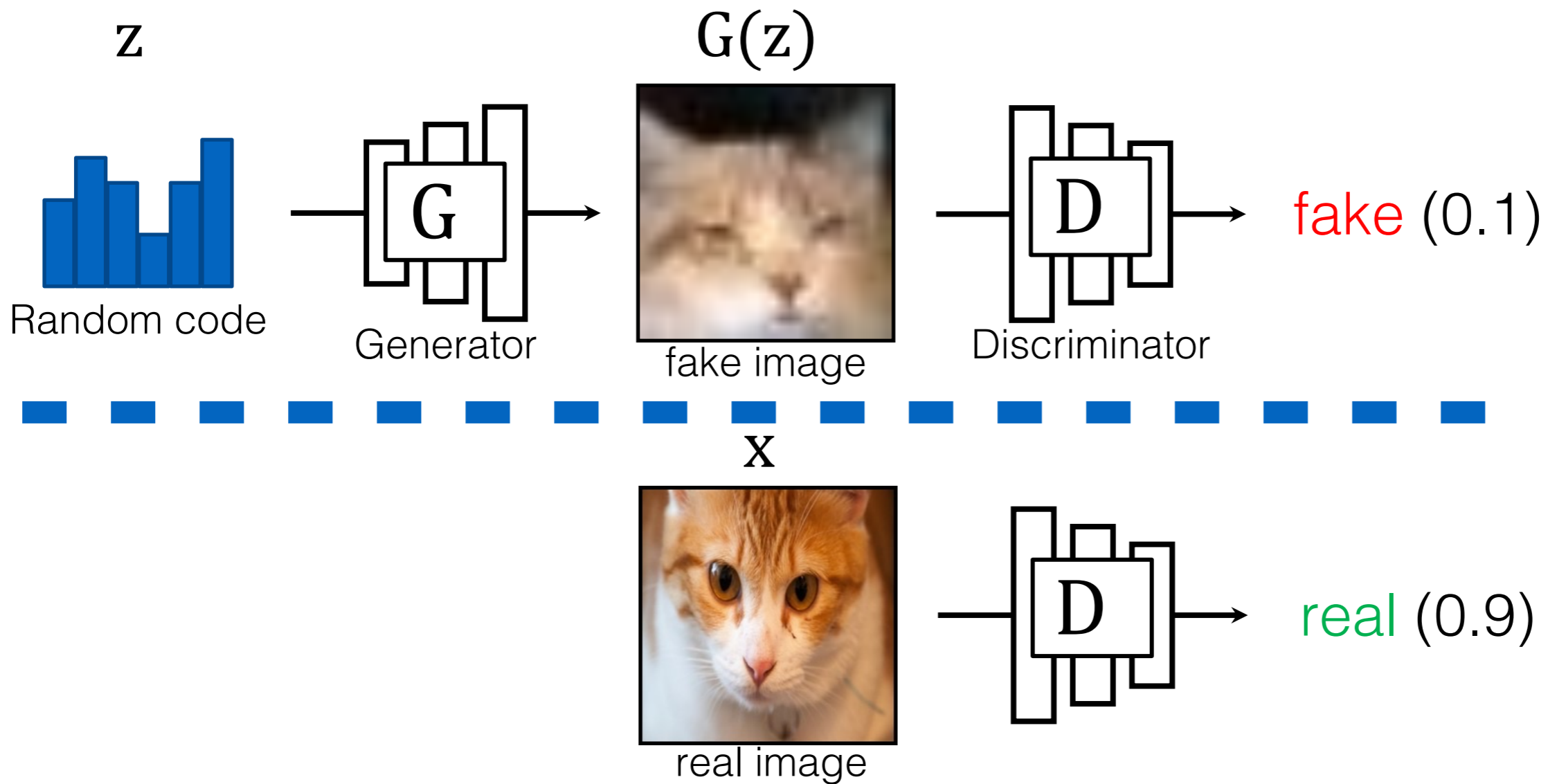




$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$$

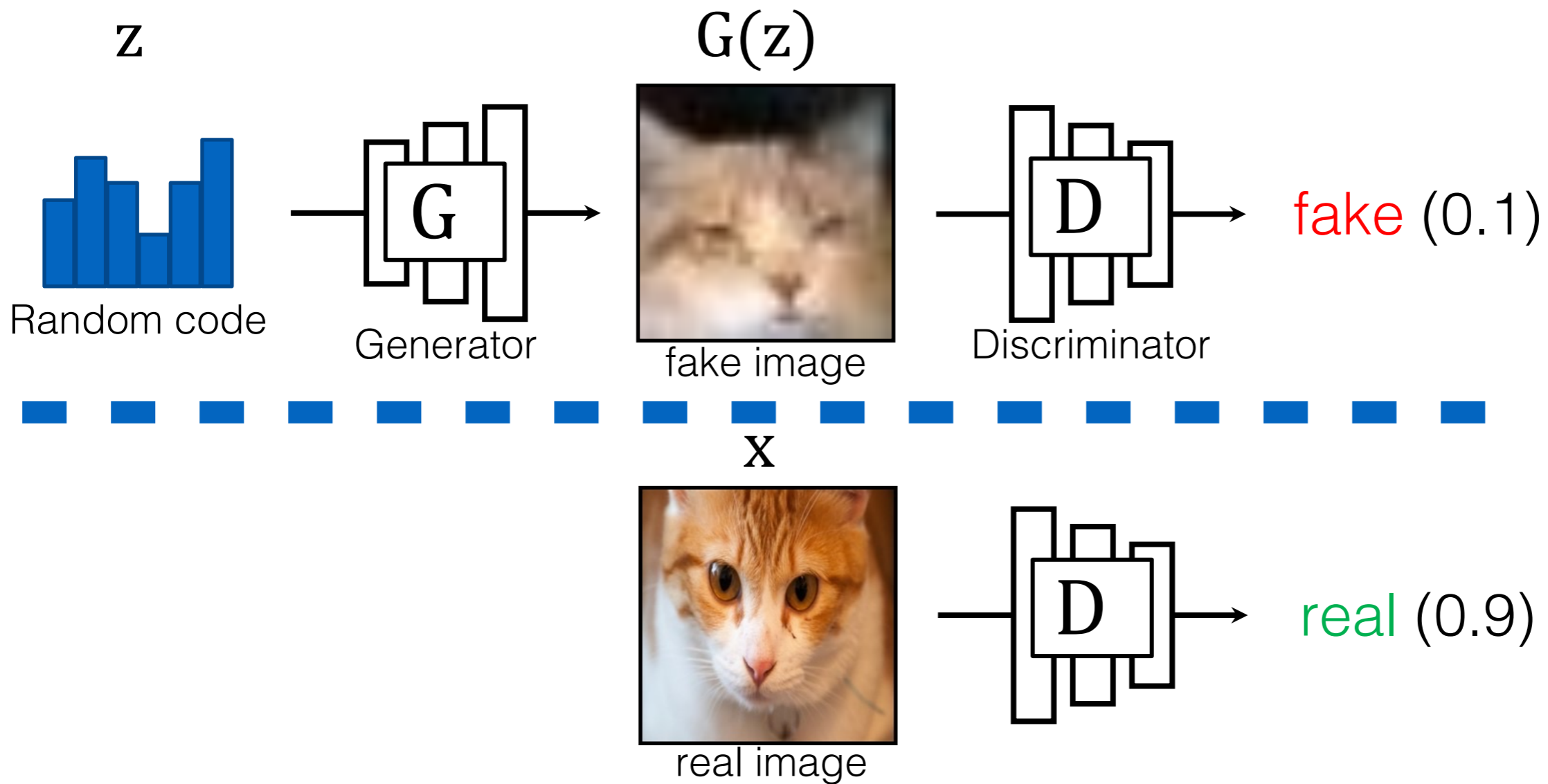
$$\log D(G(\mathbf{z})) \rightarrow -\infty$$

from [Arjovsky, Chintala, Bottou, 2017]



Learning objective (GANs)

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

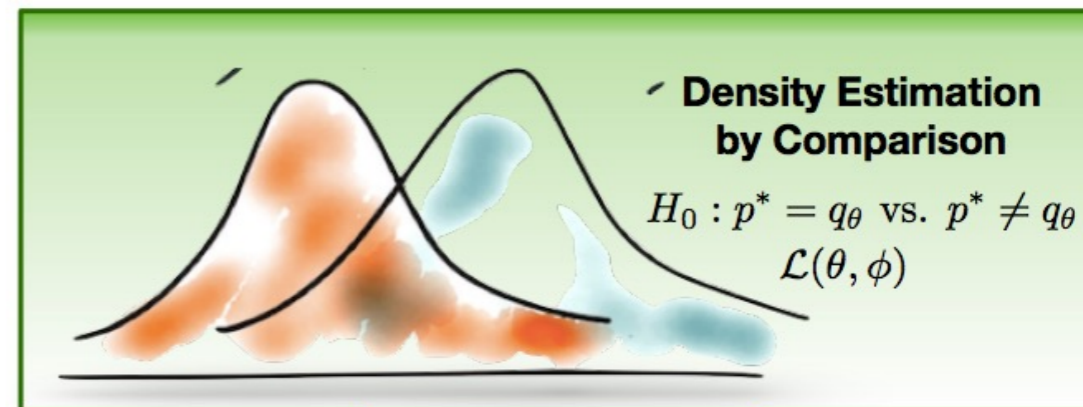


Learning objective (GANs variants)

$$\min_G \max_{f_1, f_2} \mathbb{E}_z [f_1(G(z))] + \mathbb{E}_x [f_2(x)]$$

EBGAN, WGAN, LSGAN, etc

# Other divergences?



from [Mohamed & Lakshminarayanan 2017]

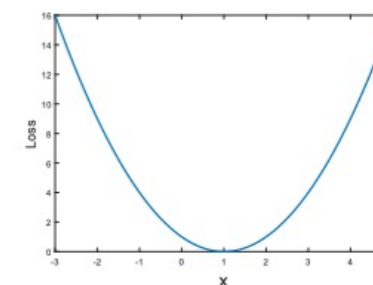
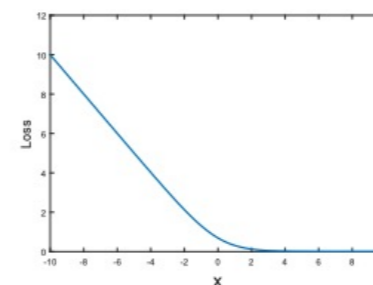
$$\min_G \max_{f_1, f_2} \mathbb{E}_z [f_1(G(z))] + \mathbb{E}_x [f_2(x)] \quad \begin{array}{l} \text{Convenient choice} \\ f_1 = -f \\ f_2 = f \end{array}$$

Different choices of  $f_1$  and  $f_2$  correspond to different divergence measures:

- Original GAN  $\rightarrow$  JSD
- Least-squares GAN  $\rightarrow$  Pearson chi-squared divergence

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

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



# Other divergences?

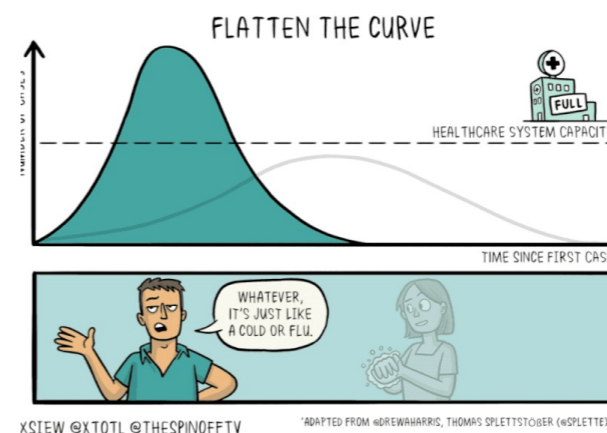
$$KL(p_{\text{data}} || p_{\theta}) \longleftarrow \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

$$KL(p_{\theta} || p_{\text{data}}) \longleftarrow \text{Reverse KL — mode seeking, intractable}$$

$$JS(p_{\text{data}}, p_{\theta}) \longleftarrow \text{Jensen-Shannon, original GAN}$$

$$W(p_{\text{data}}, p_{\theta}) = \inf_{\gamma \in \Pi(p_{\text{data}}, p_{\theta})} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] \longleftarrow \text{Wasserstein}$$

Earth-Mover (EM) distance  
/ Wasserstein distance



# Maximum log likelihood, KL, and JSD

KLD (Kullback–Leibler divergence):  $\mathcal{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$

JSD (Jensen–Shannon divergence):  $\mathcal{JSD}(p || q) = \frac{1}{2}\mathcal{KL}(p || \frac{p+q}{2}) + \frac{1}{2}\mathcal{KL}(q || \frac{p+q}{2})$

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

$$\mathcal{KL}(p_{\text{data}}(x) || p_{\theta}(x)) = \int_x p_{\text{data}}(x) \log \frac{p_{\text{data}}(x)}{p_{\theta}(x)} dx$$

$$= \int_x p_{\text{data}}(x) \log p_{\text{data}}(x) dx - \int_x p_{\text{data}}(x) \log p_{\theta}(x) dx$$

↑  
Constant

(independent of  $\theta$ )

↑  
Maximize log likelihood = minimize KLD

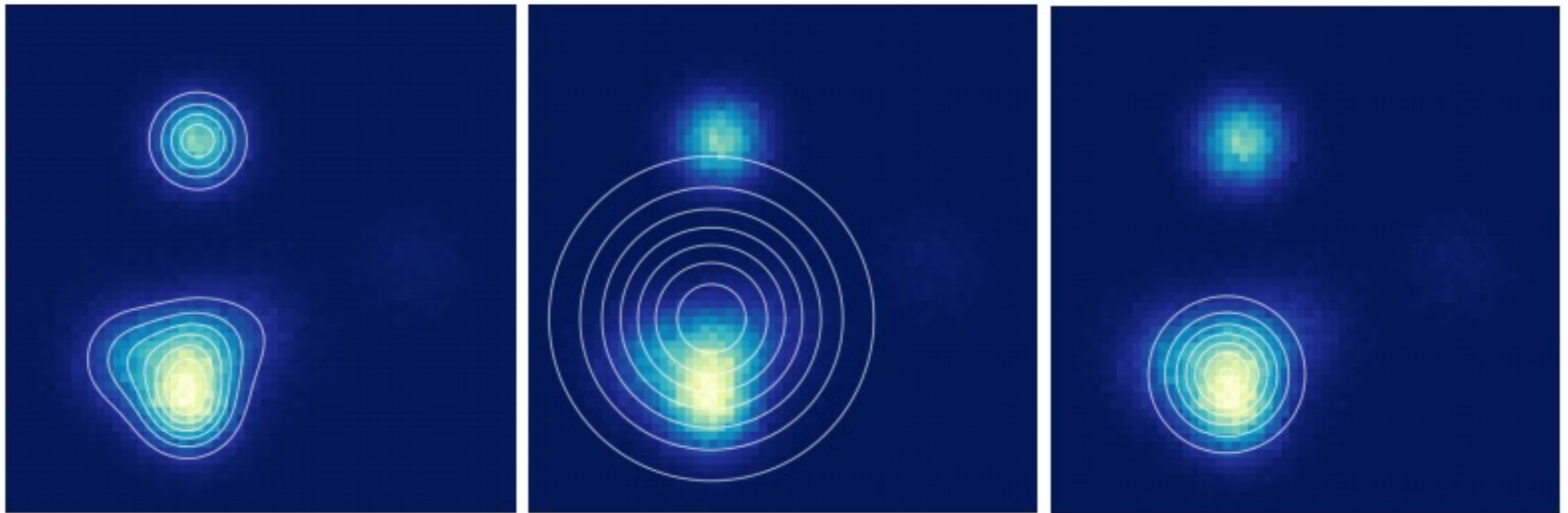


# Maximum log likelihood/KL vs. JSD

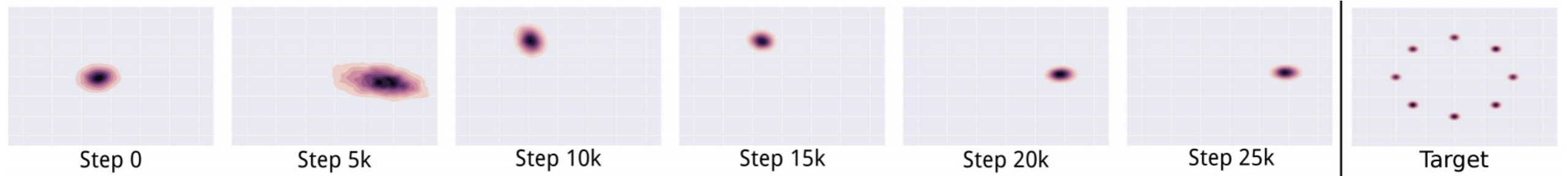
Data

Max likelihood / KL

Jensen-Shannon Divergence



[Theis et al. 2016]



# Other divergences?

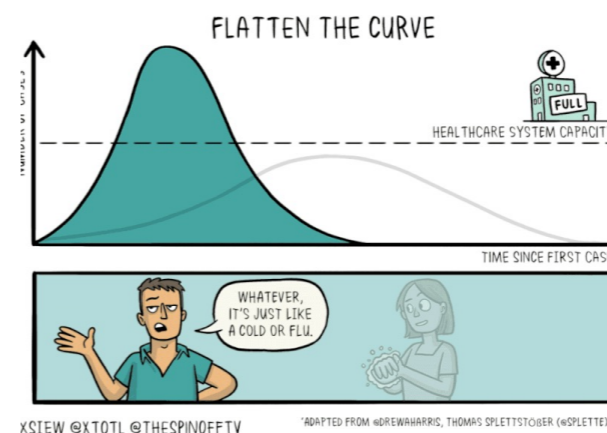
$$KL(p_{\text{data}} || p_{\theta}) \longleftarrow \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

$$KL(p_{\theta} || p_{\text{data}}) \longleftarrow \text{Reverse KL — mode seeking, intractable}$$

$$JS(p_{\text{data}}, p_{\theta}) \longleftarrow \text{Jensen-Shannon, original GAN}$$

$$W(p_{\text{data}}, p_{\theta}) = \inf_{\gamma \in \Pi(p_{\text{data}}, p_{\theta})} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] \longleftarrow \text{Wasserstein}$$

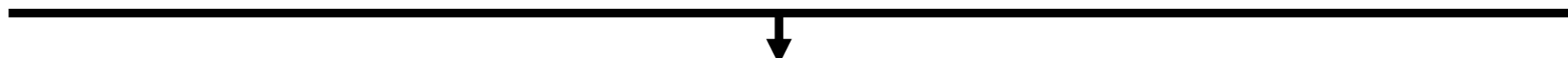
Earth-Mover (EM) distance  
/ Wasserstein distance



# Wasserstein GAN

[Arjovsky, Chintala, Bottou 2017]

$$\arg \min_G \max_{\|f\|_L \leq 1} \mathbb{E}_{\mathbf{z}, \mathbf{x}} [ -f(G(\mathbf{z})) + f(\mathbf{x}) ]$$



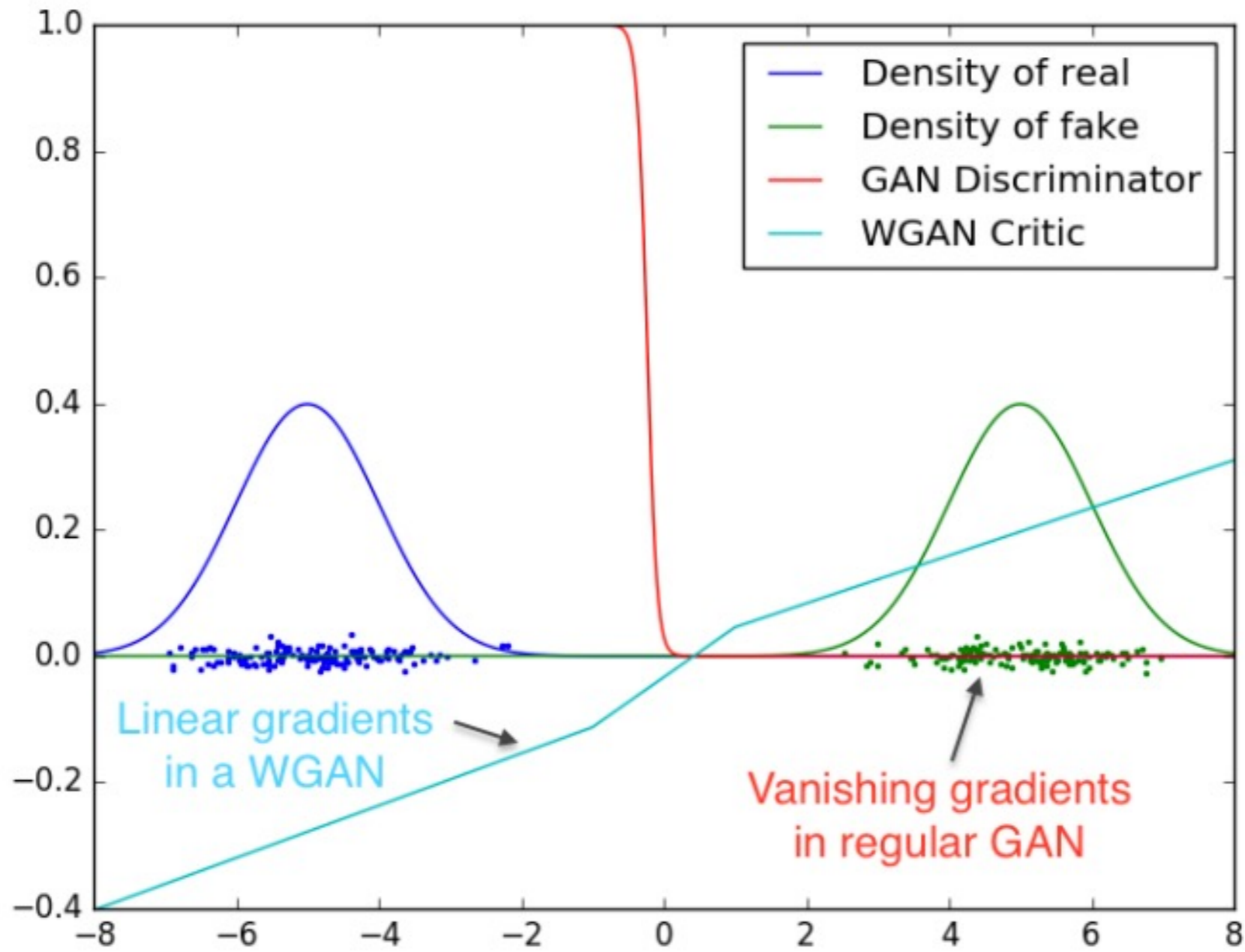
Lipschitz continuity  
 $|f(x) - f(y)| \leq |x - y|$

$$W(p_{\text{data}}, p_{\theta}) = \inf_{\gamma \in \Pi(p_{\text{data}}, p_{\theta})} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

wGAN GP [Gulrajani et al., 2018]:

$$\arg \min_G \max_f \mathbb{E}_{\mathbf{z}, \mathbf{x}} [ -f(G(\mathbf{z})) + f(\mathbf{x}) ] + \lambda \mathbb{E}_{\hat{\mathbf{x}} \sim P_{\hat{\mathbf{x}}}} [ (\|\nabla_{\hat{\mathbf{x}}} f(\hat{\mathbf{x}})\|_2 - 1)^2 ]$$

Gradient penalty (GP)



from [Arjovsky, Chintala, Bottou, 2017]

To be continued...

# Thank You!



16-726, Spring 2022

<https://learning-image-synthesis.github.io/sp22/>