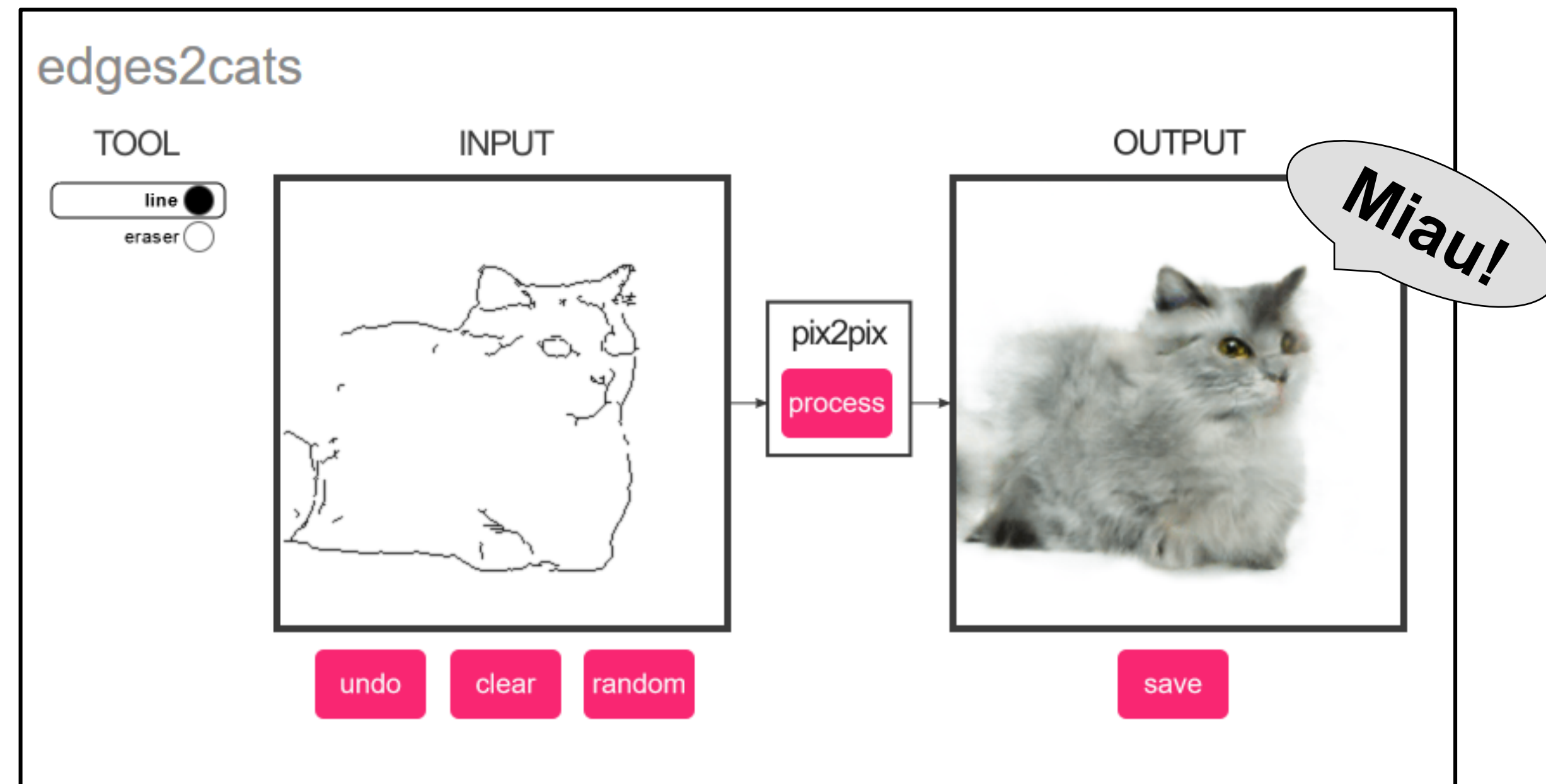


# Conditional Adversarial Networks

(or “mapping from A to B”)



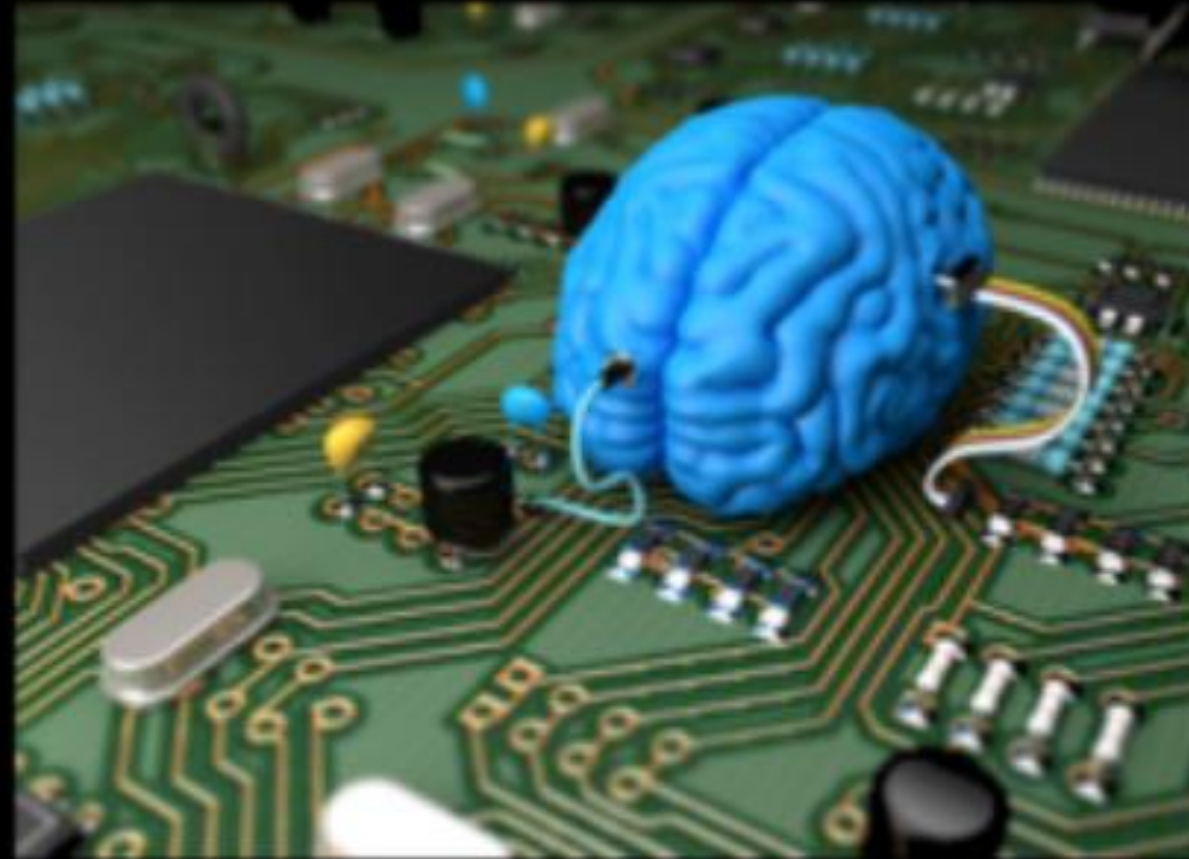
CS448V — Computational Video Manipulation

May 22<sup>th</sup>, 2019

# Deep Learning



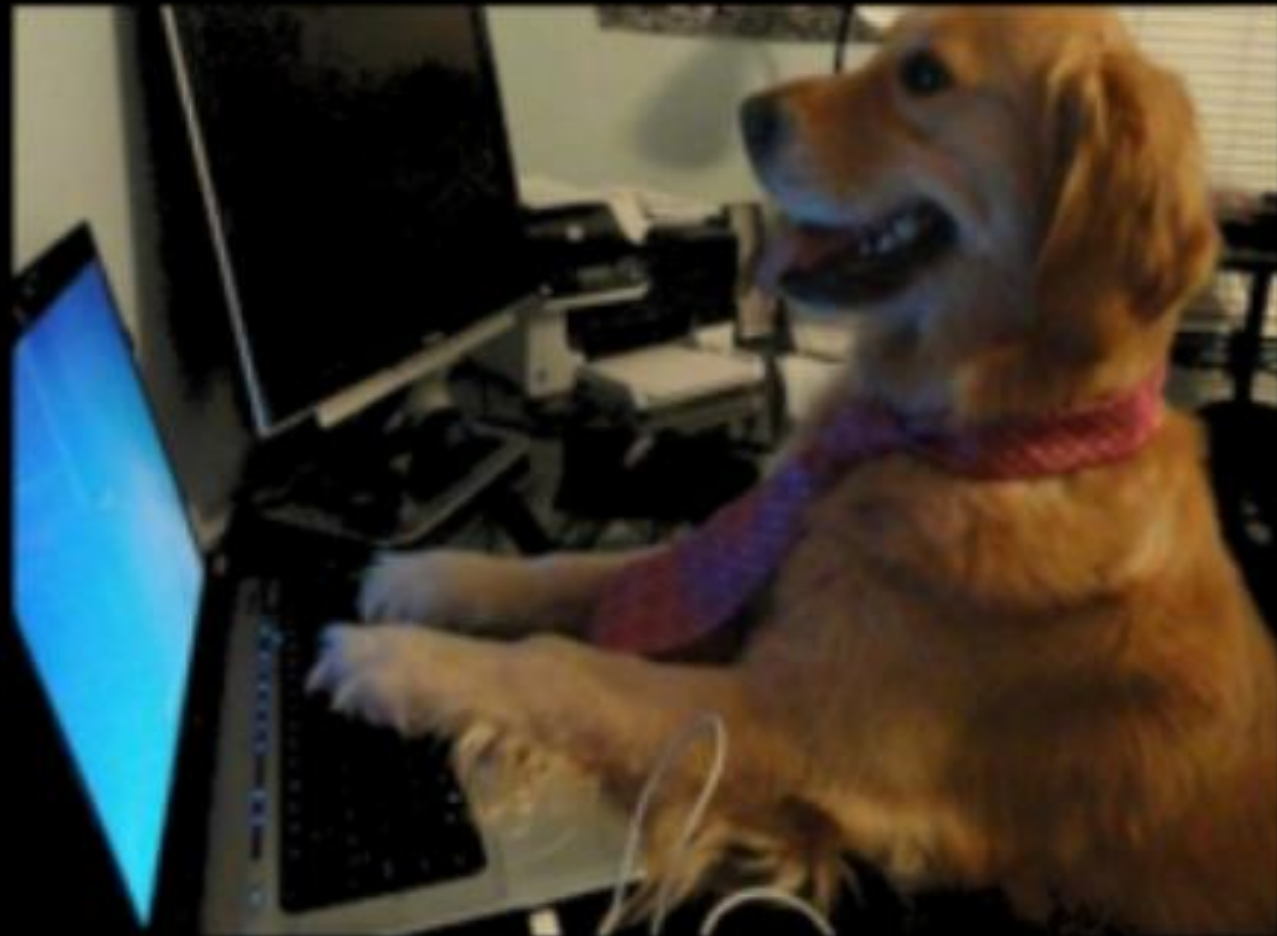
What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
In [1]:  
import keras  
Using TensorFlow backend.
```

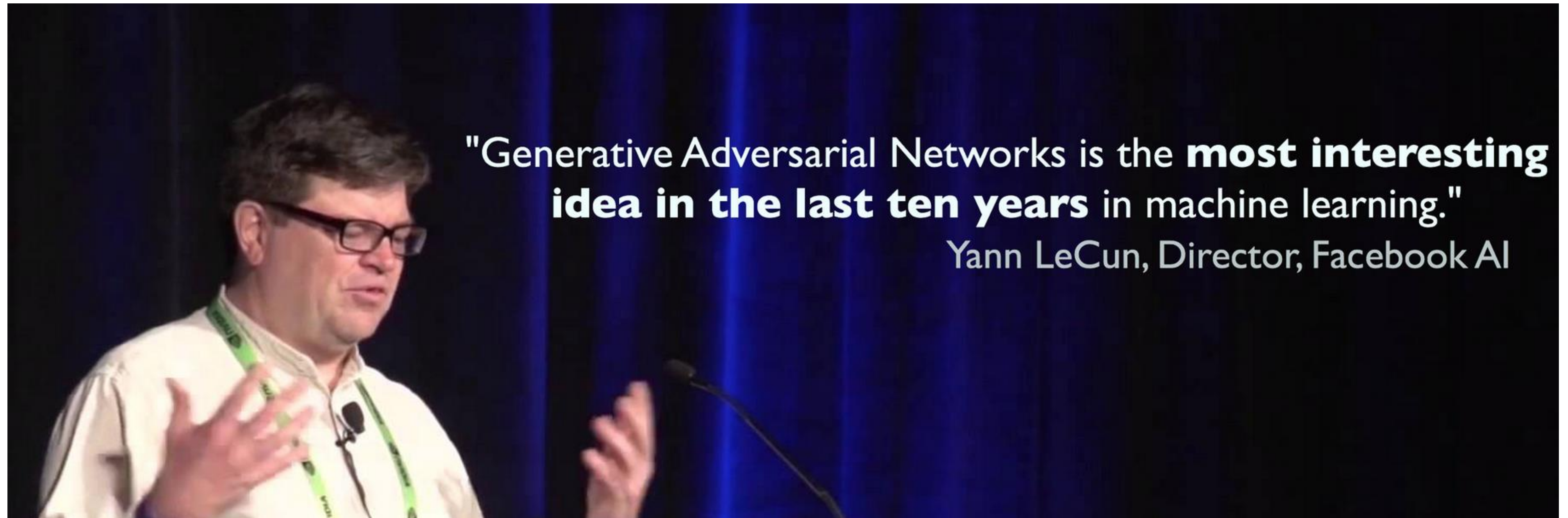
What I actually do

# Why? - Cool! Trendy! - Google Scholar

TITLE	CITED BY	YEAR
<b>Generative adversarial nets</b> I Goodfellow, J Pouget-Abadie, M Mirza, B Xu, D Warde-Farley, S Ozair, ... Advances in neural information processing systems, 2672-2680	8405	2014
<b>Image-to-image translation with conditional adversarial networks</b> P Isola, JY Zhu, T Zhou, AA Efros Proceedings of the IEEE conference on computer vision and pattern ...	<b>Pix2Pix</b> 2137	2017
<b>Unpaired image-to-image translation using cycle-consistent adversarial networks</b> JY Zhu, T Park, P Isola, AA Efros Proceedings of the IEEE international conference on computer vision, 2223-2232	<b>CycleGAN</b> 1722	2017

⋮  
Hundreds of applications  
and follow-up works  
⋮

# Why? - Cool! Trendy! - Google Scholar



▪  
Hundreds of applications  
and follow-up works

⋮  
▪

# Enhancing Transitions

## Neural Rerendering in the Wild

Moustafa Meshry<sup>1</sup>, Dan B Goldman<sup>2</sup>, Sameh Khamis<sup>2</sup>, Hugues Hoppe<sup>2</sup>, Rohit Pandey<sup>2</sup>,  
Noah Snavely<sup>2</sup>, Ricardo Martin-Brualla<sup>2</sup>

<sup>1</sup>University of Maryland, <sup>2</sup>Google Inc.

# Single-Photo Facial Animation

Warp-Guided GANs for Single-Photo Facial Animation

Jiahao Geng Tianjia Shao Youyi Zheng Yanlin Weng Kun Zhou

State Key Lab of CAD&CG, Zhejiang University

ZJU-FaceUnity Joint Lab of Intelligent Graphics

# Text-based Editing

## Adding New Words



Original Video

I love the smell of napalm in the morning.

# Few-Shot Reenactment

## Few-Shot Adversarial Learning of Realistic Neural Talking Head models

Egor Zakharov<sup>1,2</sup> Aliaksandra Shysheya<sup>1,2</sup> Egor Burkov<sup>1,2</sup> Victor Lempitsky<sup>1,2</sup>

<sup>1</sup>Samsung Research

<sup>2</sup>Skolkovo Institute of Science and Technology

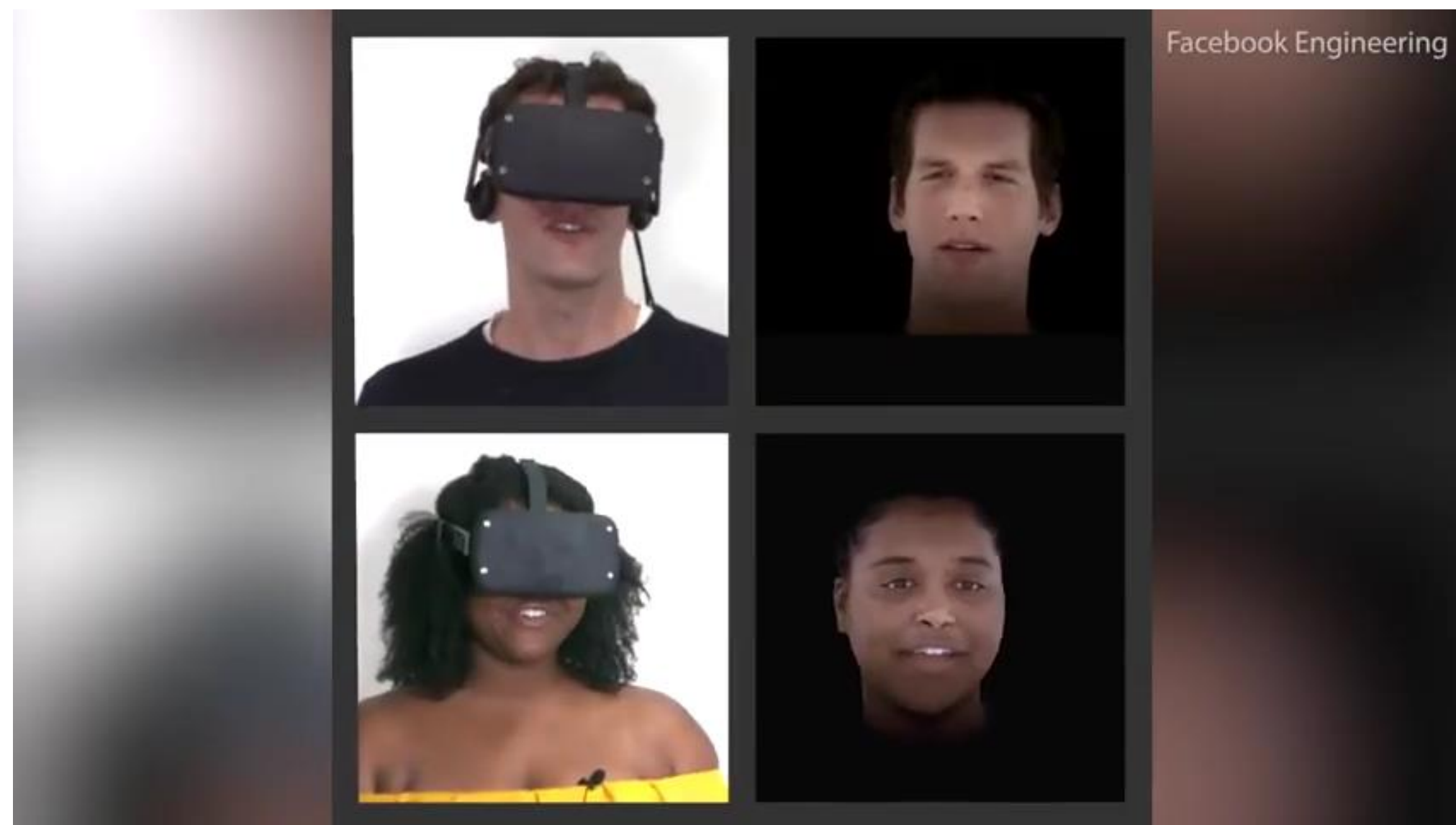
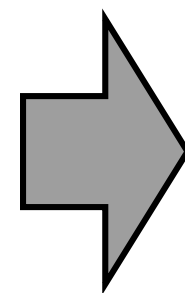
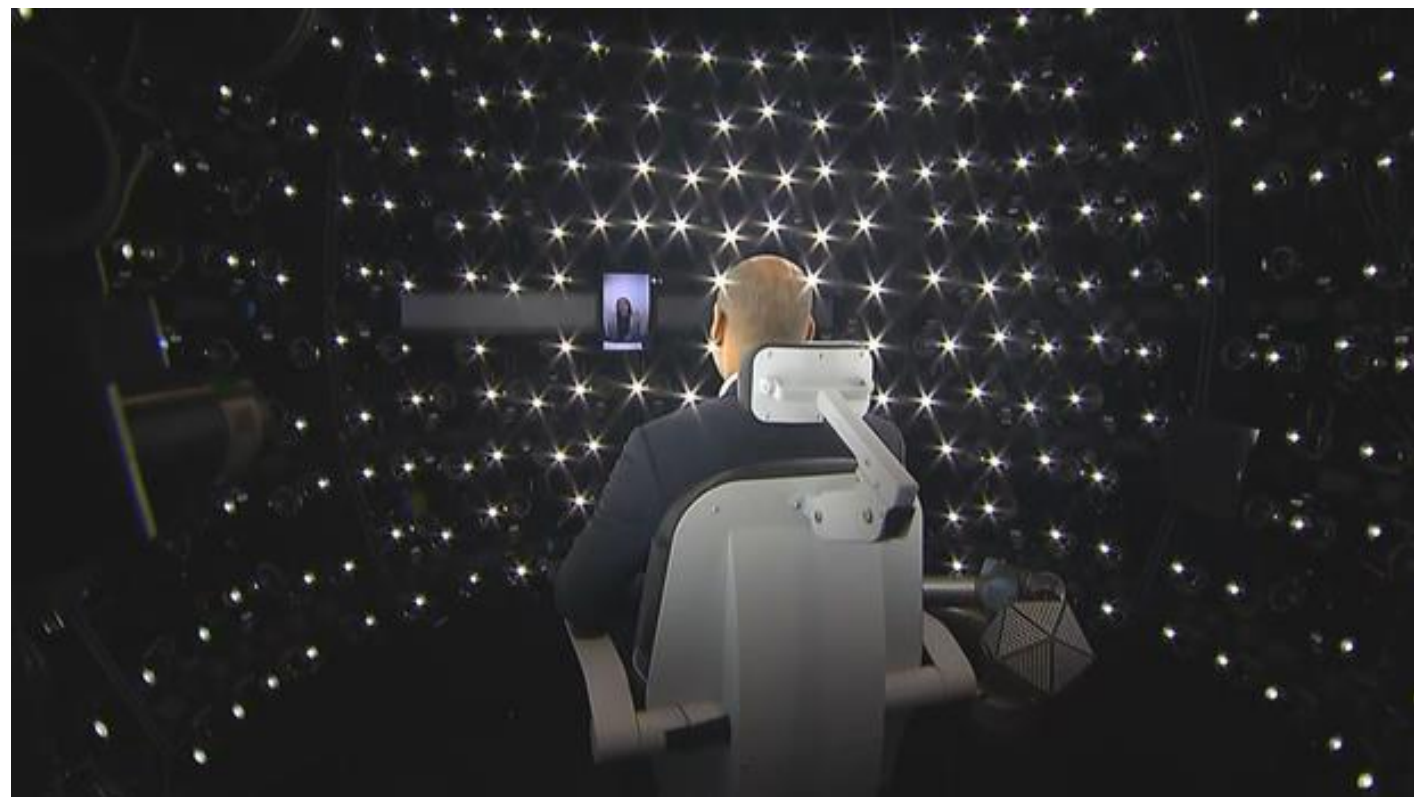
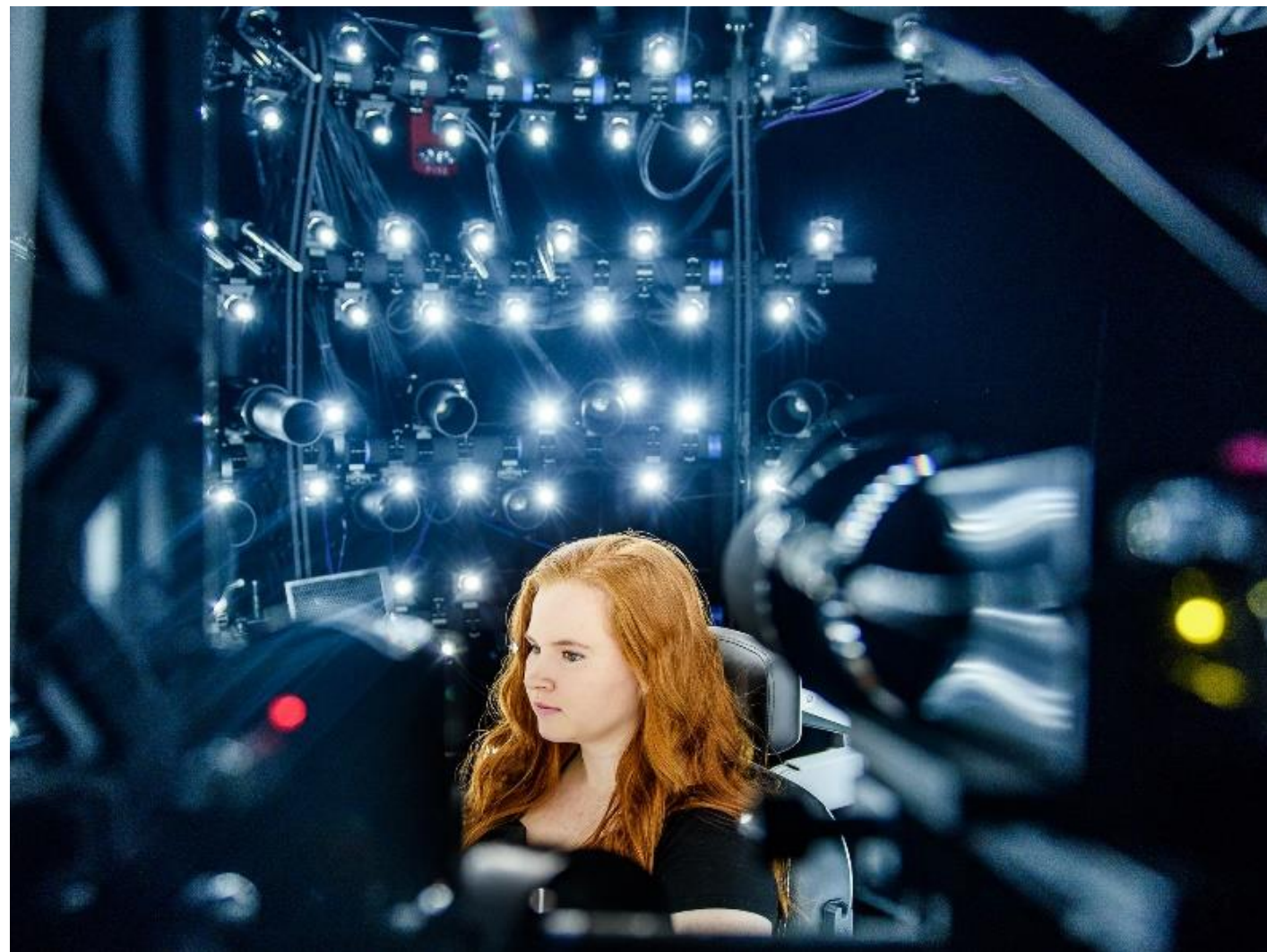


Source

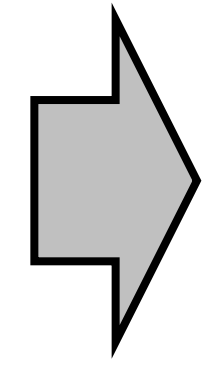
Generated images



# Digital Humans



# Overview



- Convolutional Neural Networks

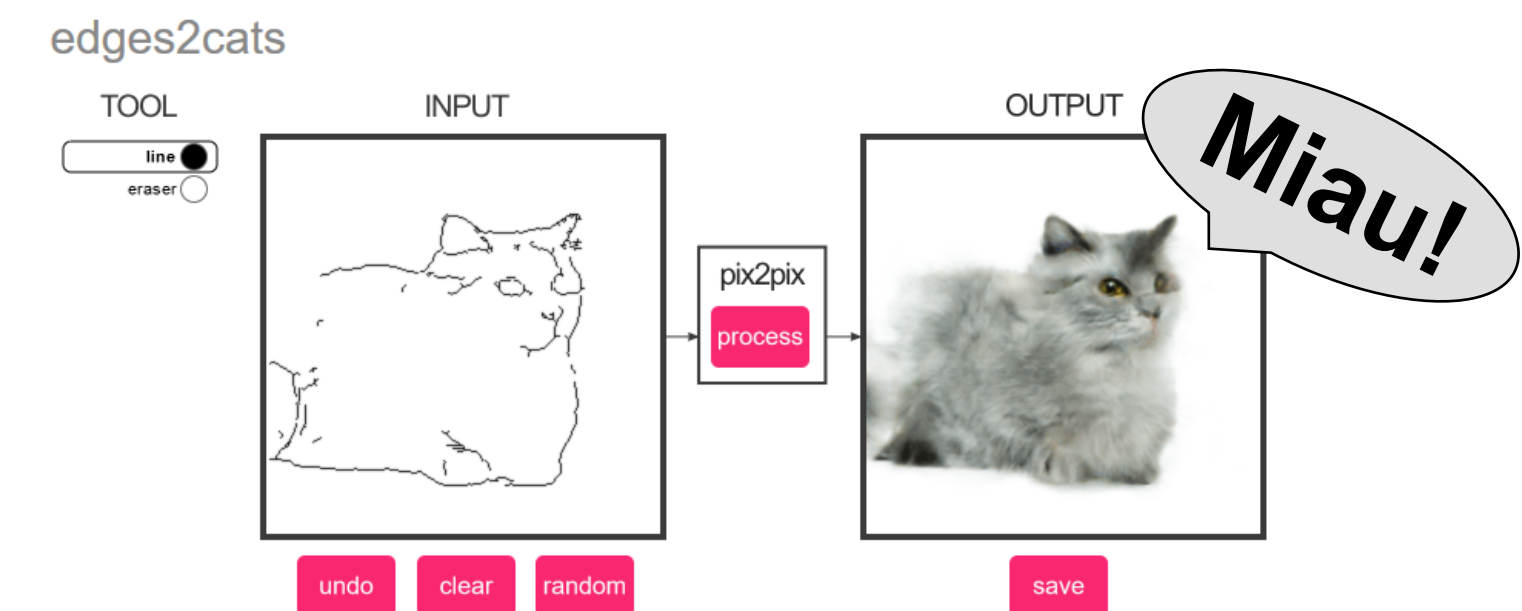
$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

- Generative Modeling



(Brundage et al., 2018)

- Pix2Pix (“mapping from A to B”)



# Convolutional Neural Network

## Components?

- 2D Convolution Layers (Conv2D)
- Subsampling Layers (MaxPool, ...)
- Non-linearity Layers (ReLU, ...)
- Normalization Layers (BatchNorm, ...)
- Upsampling Layers (TransposedConv, ...)
- ...



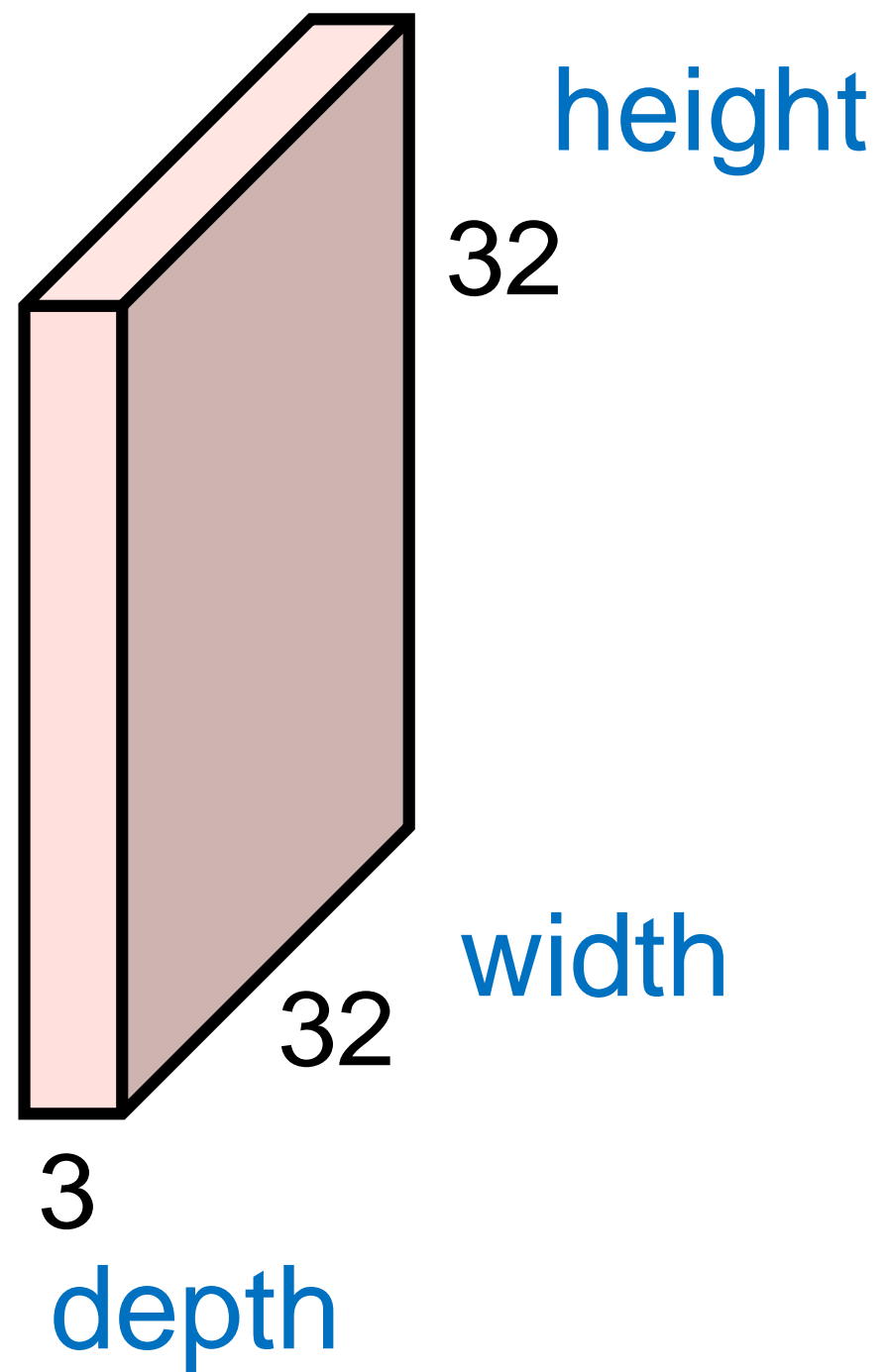
# Convolutional Neural Network

## Components?

- 2D Convolution Layers (Conv2D)
- Subsampling Layers (MaxPool, ...)
- Non-linearity Layers (ReLU, ...)
- Normalization Layers (BatchNorm, ...)
- Upsampling Layers (TransposedConv, ...)
- ...

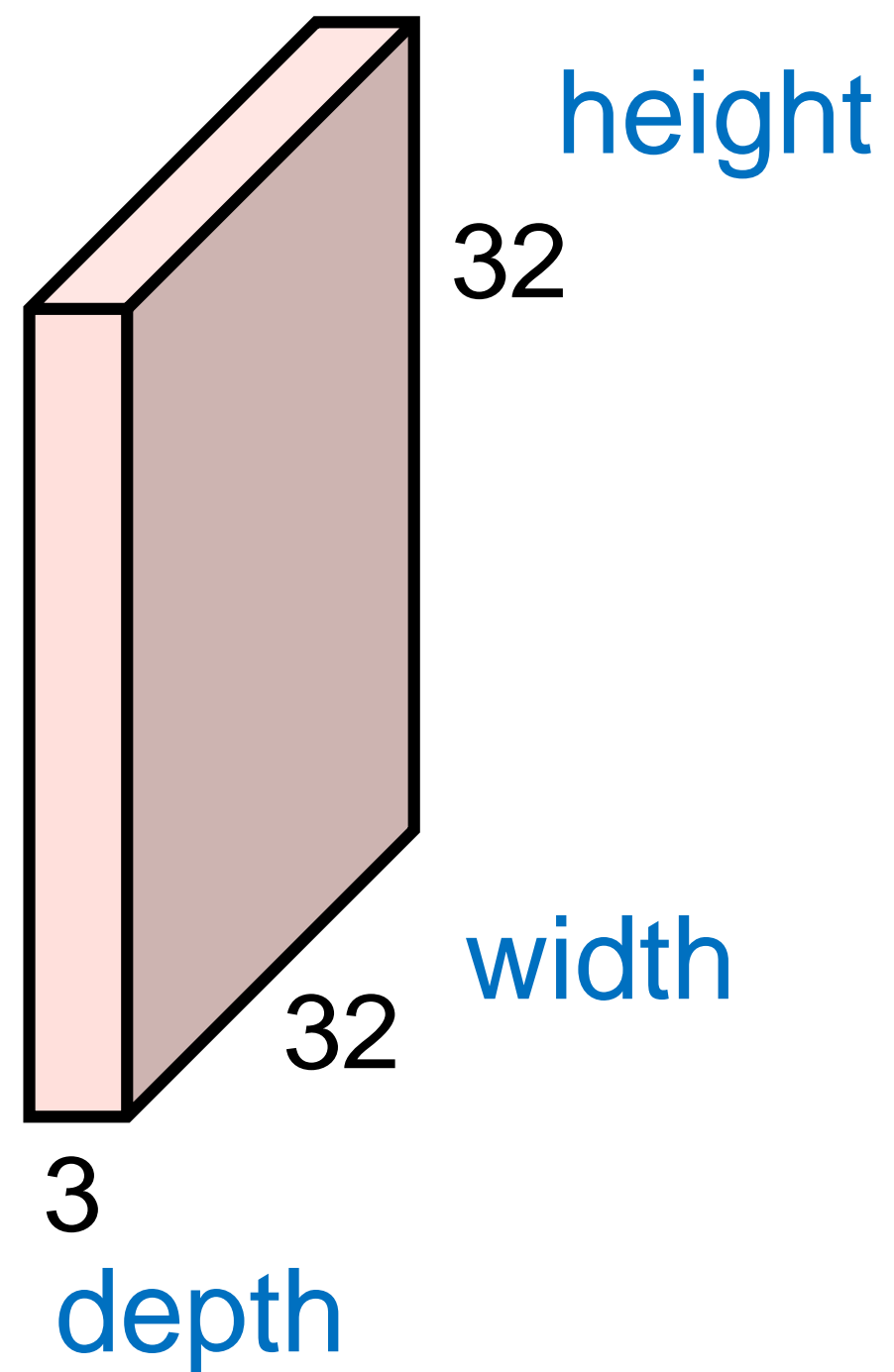
# Convolution

32x32x3 image

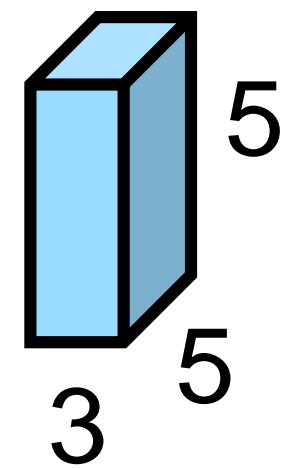


# Convolution

32x32x3 image



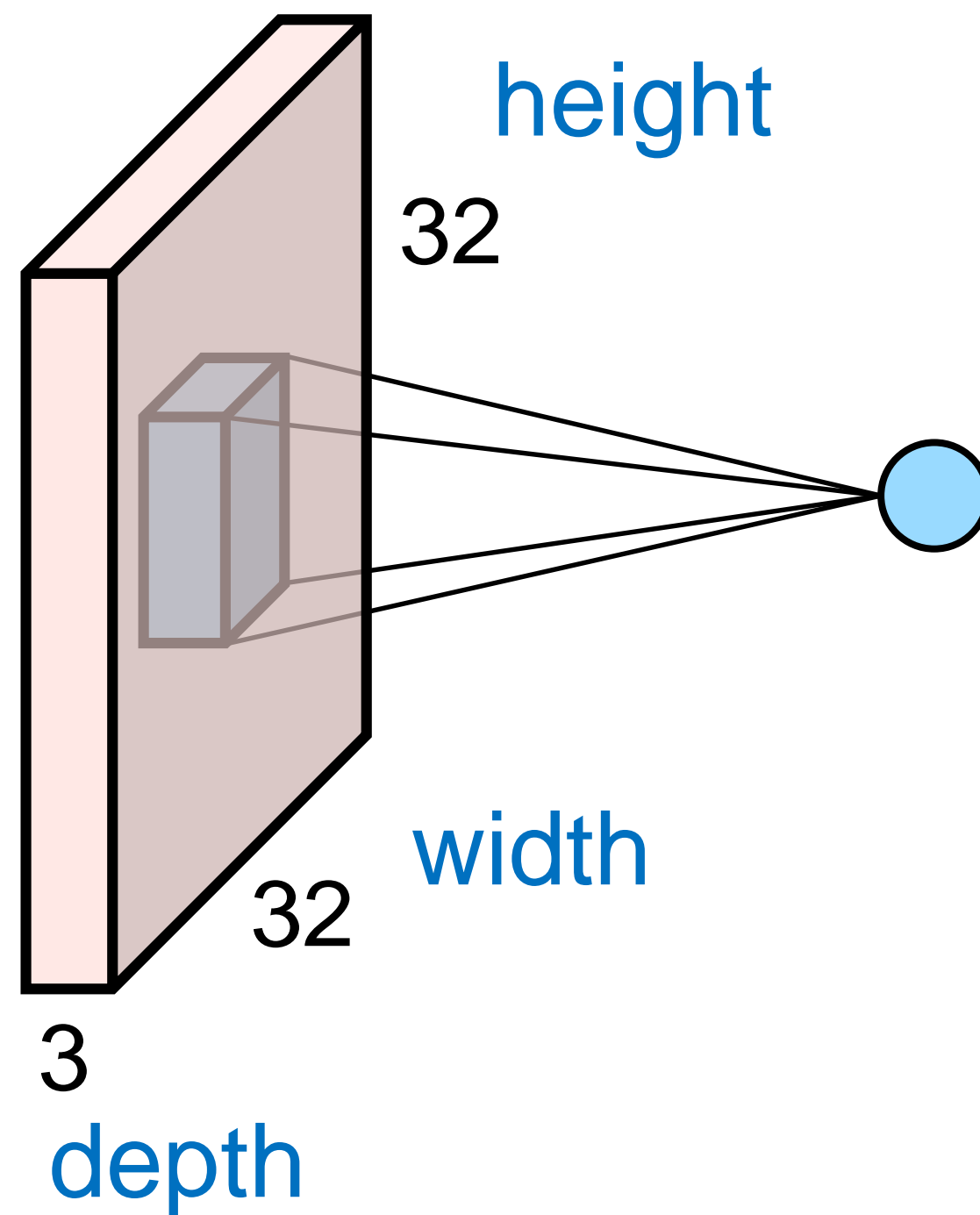
5x5x3 filter



**Convolve** the filter with the image, i.e., “slide over the image spatially, computing dot products”

# Convolution

32x32x3 image



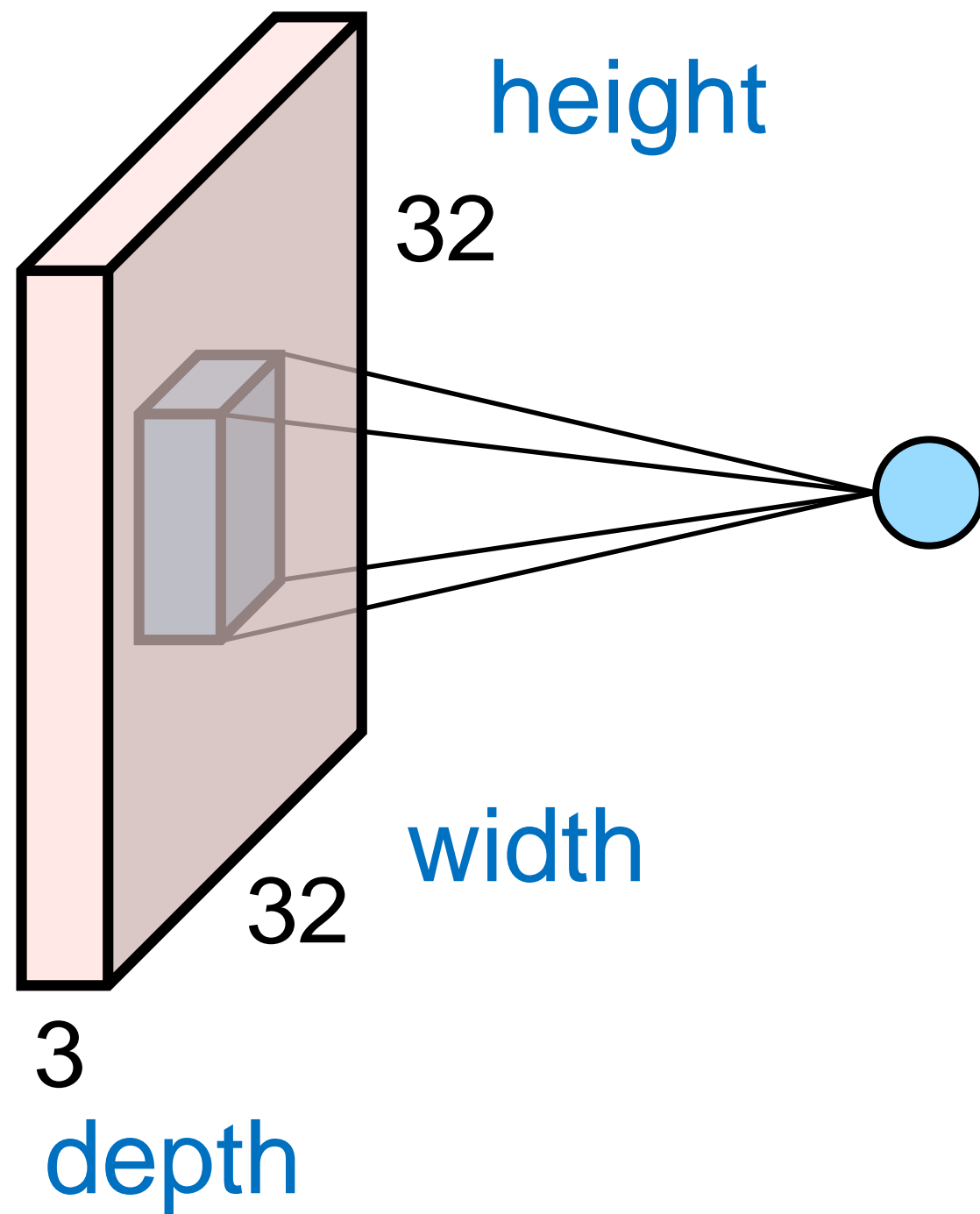
5x5x3 filter

**Result: 1 number**, the result of taking the dot product between the filter and a small 5x5x3 chunk of the image, i.e., 5x5x3 = 70-dimensional dot product + bias

$$w^T x + b$$

# Convolution

32x32x3 image

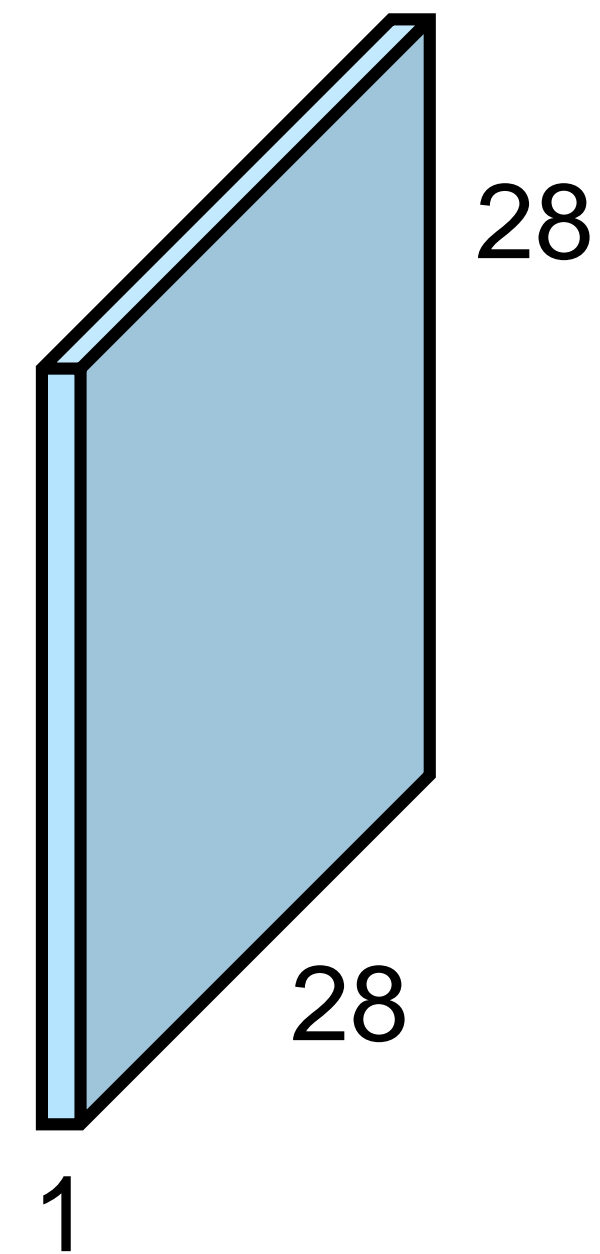


5x5x3 filter



**Convolve (slide) over all spatial locations**

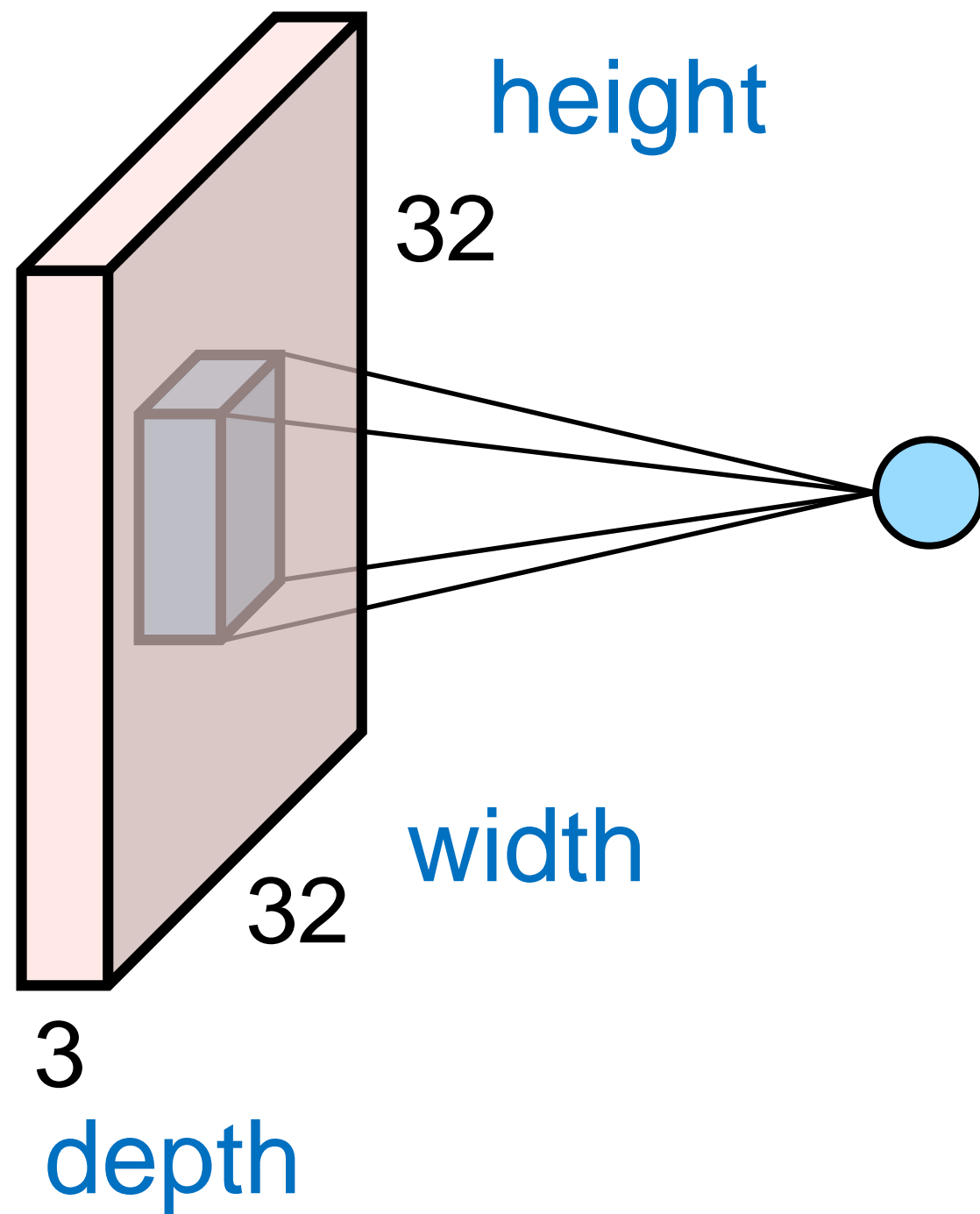
Activation map





# Convolution

32x32x3 image

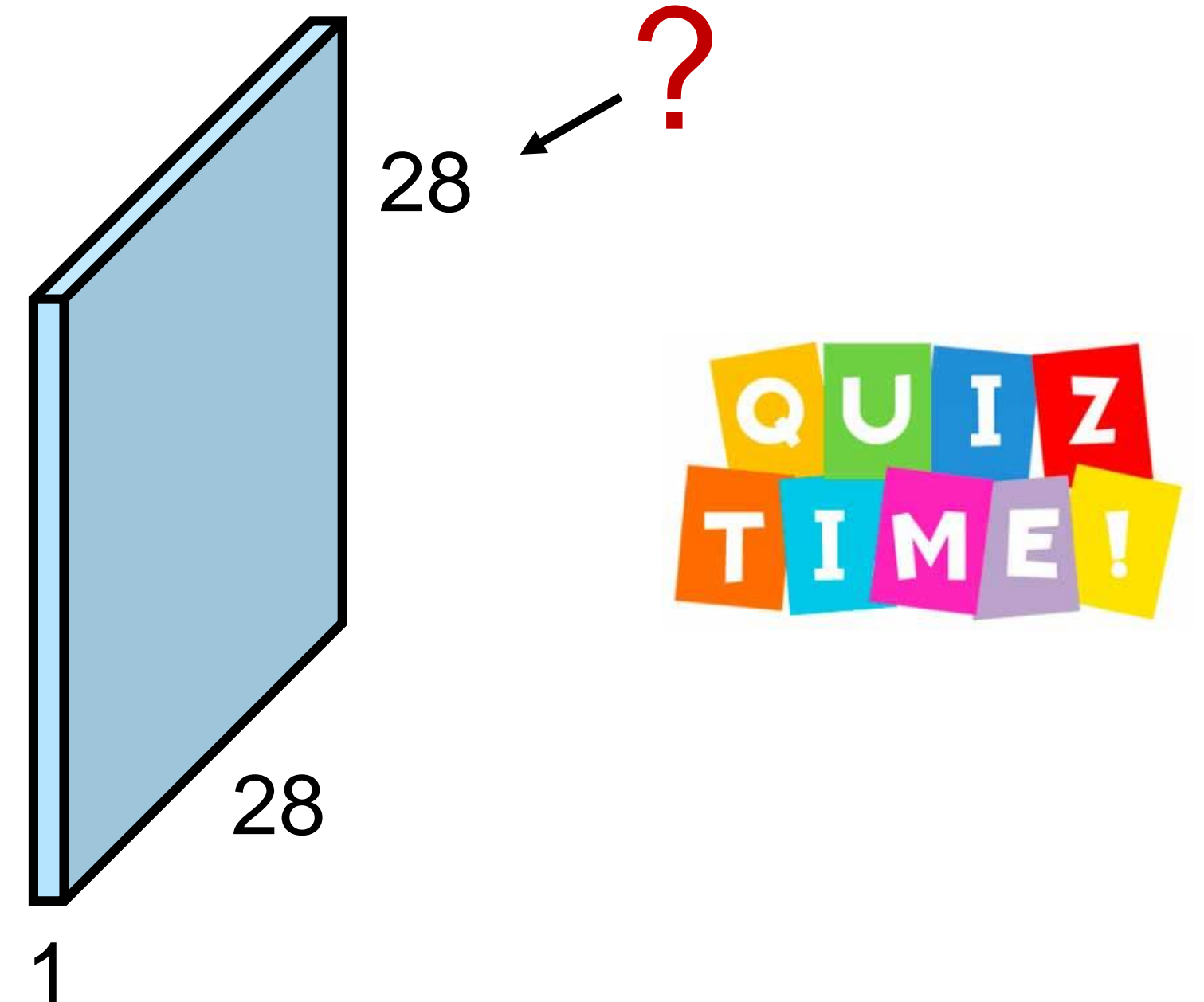


5x5x3 filter



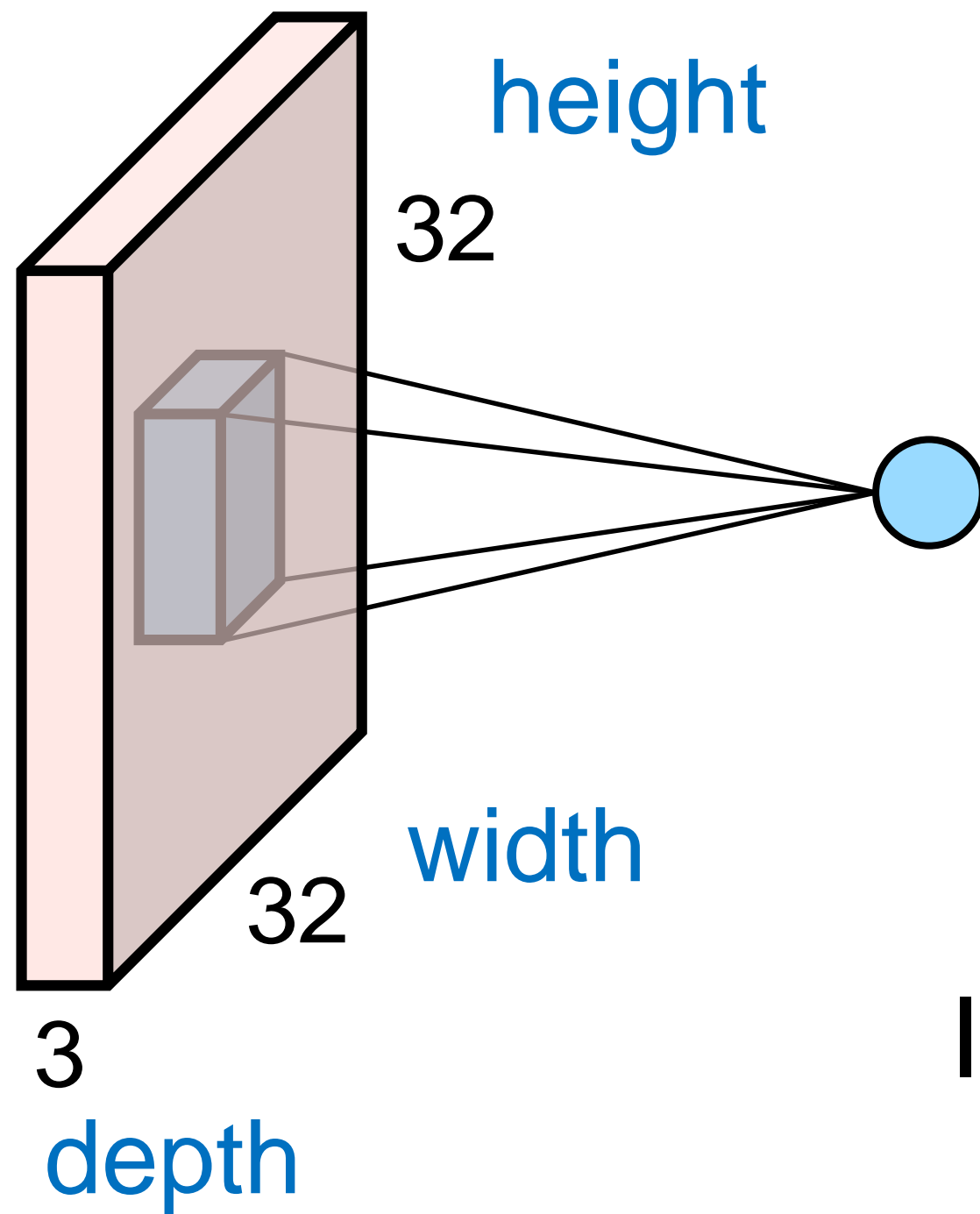
Convolve (slide) over all spatial locations

Activation map



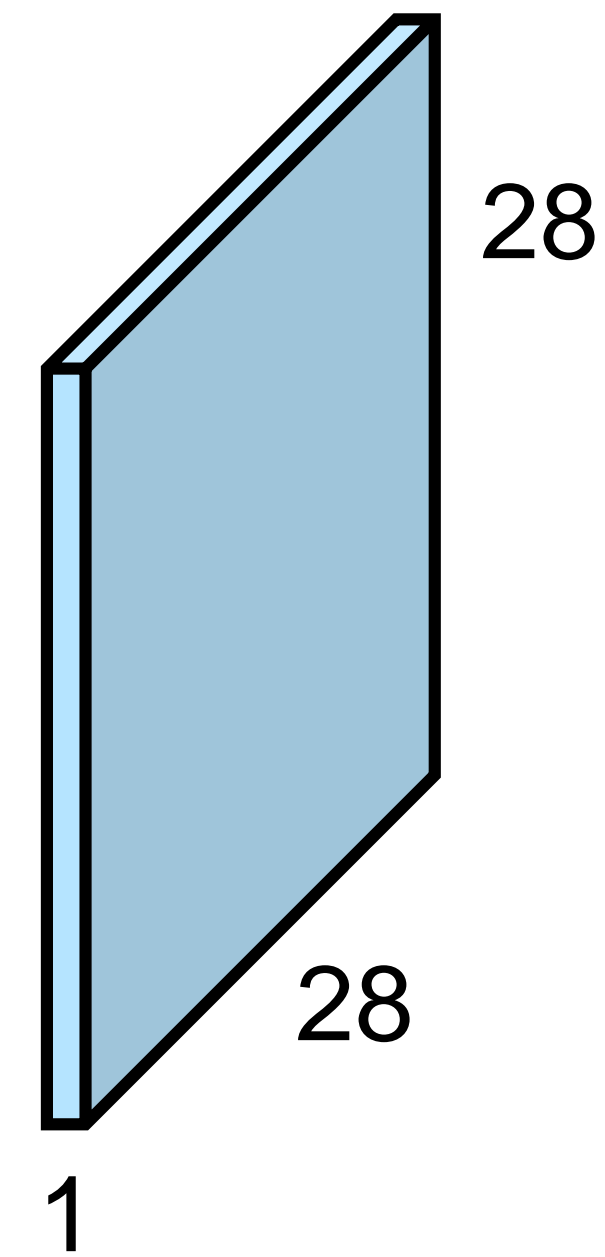
# Convolution

32x32x3 image



Convolve (slide) over all spatial locations

Activation map



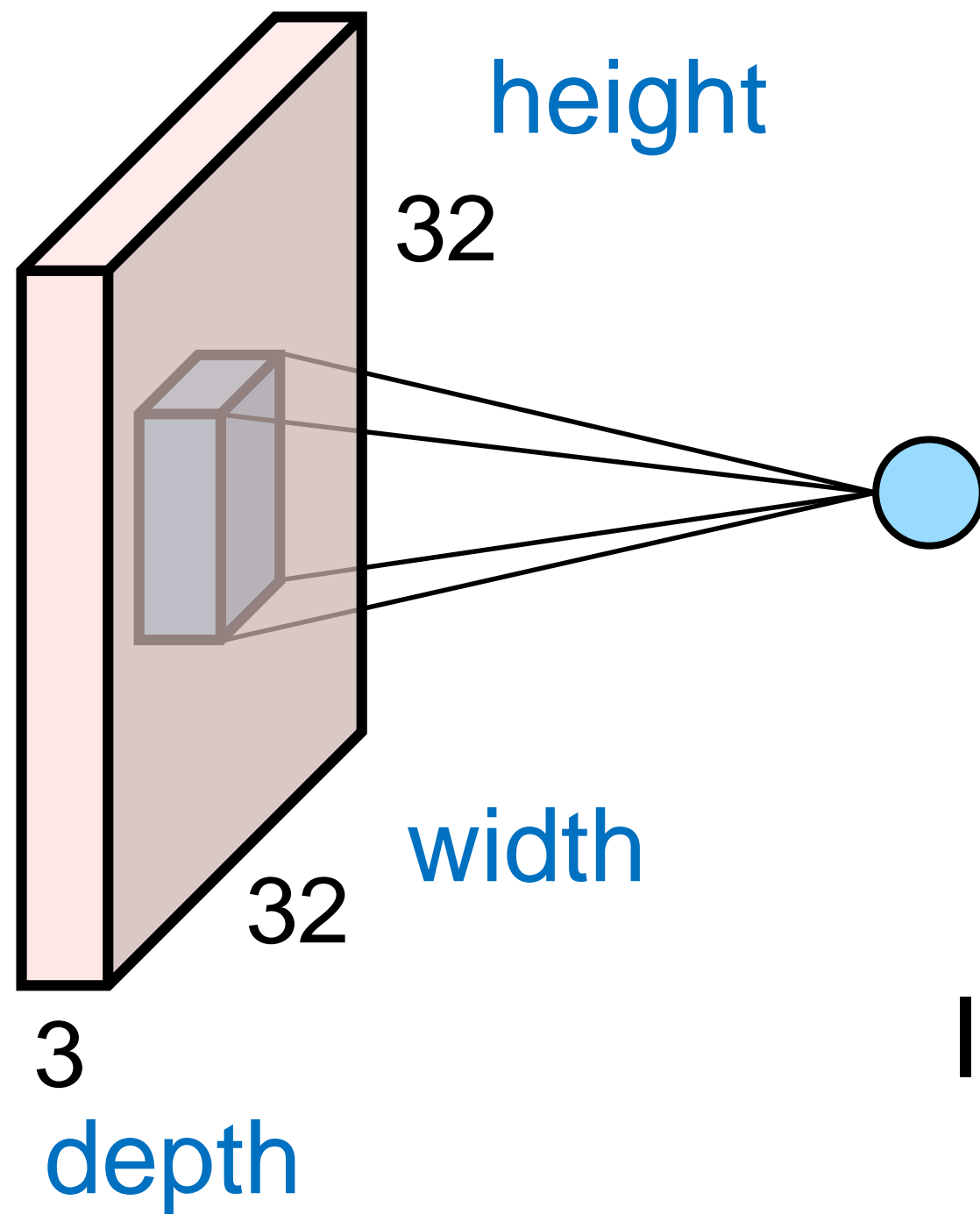
Invariant to?

- Rotation
- Translation
- Scaling



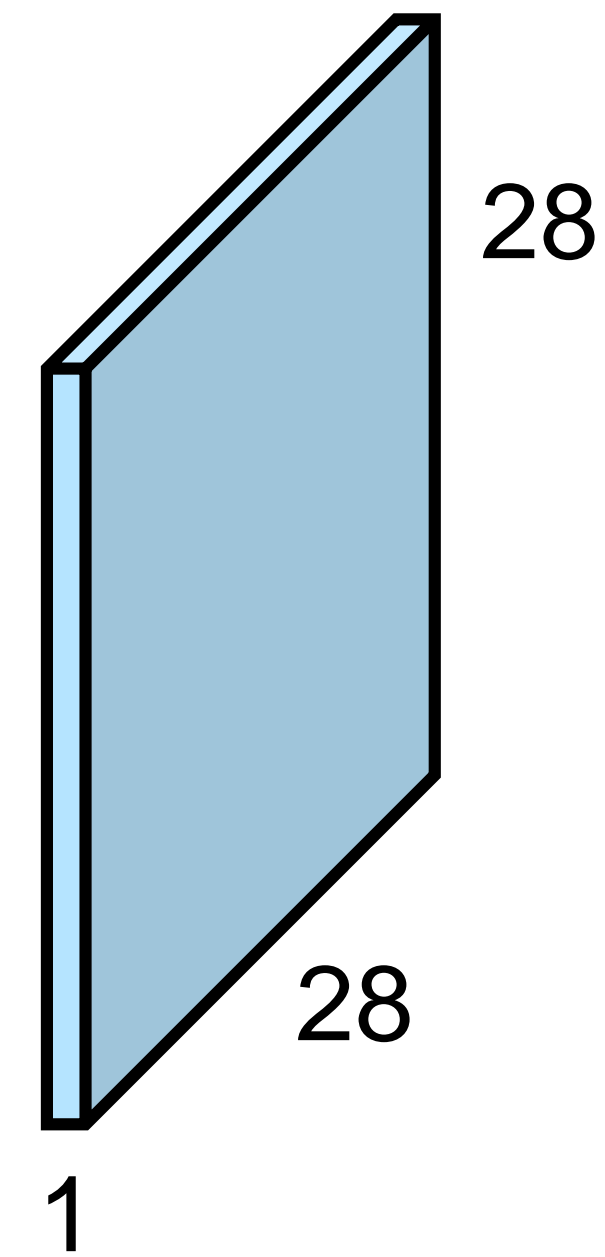
# Convolution

32x32x3 image



Convolve (slide) over all spatial locations

Activation map



Invariant to?

Rotation



Translation

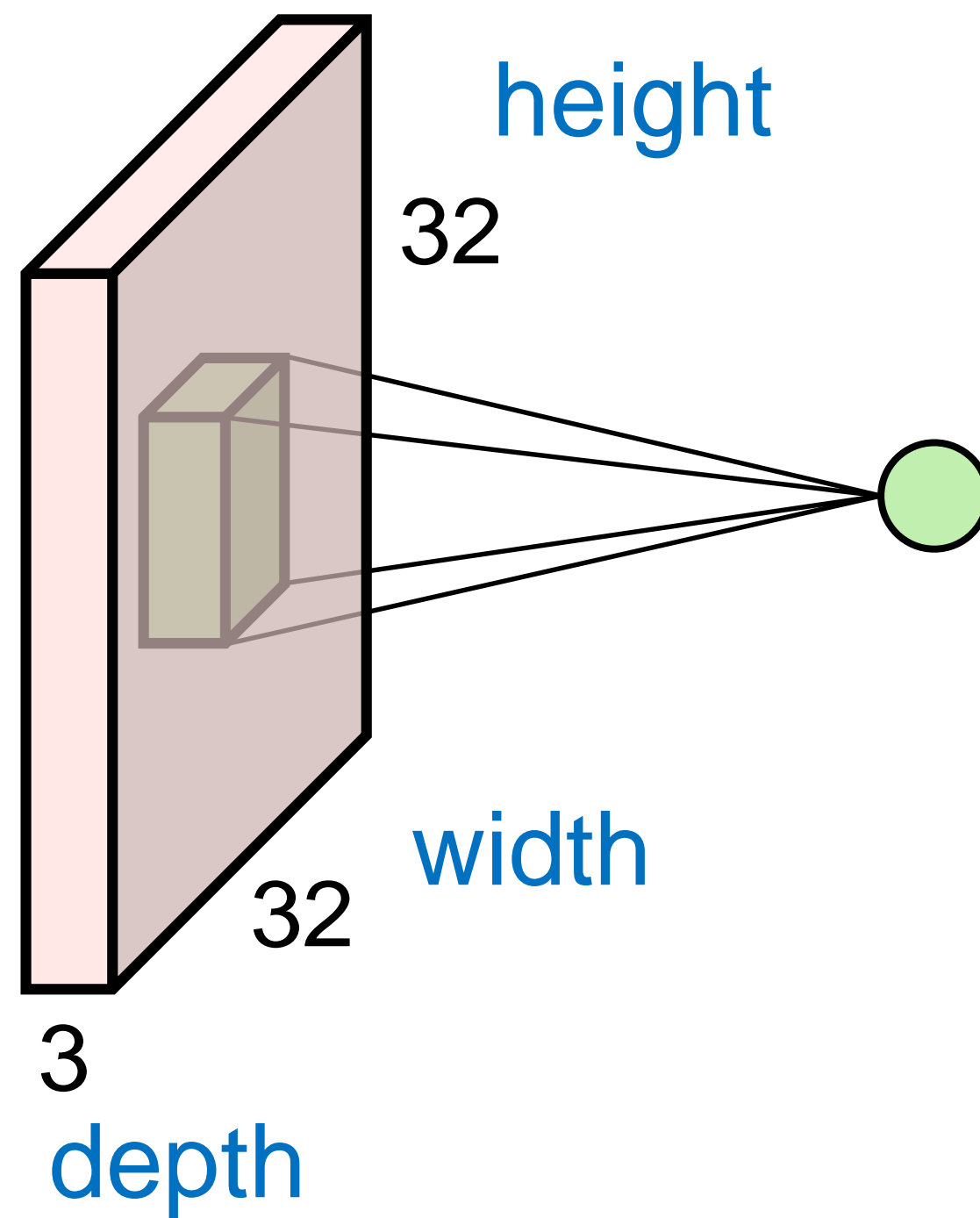


Scaling



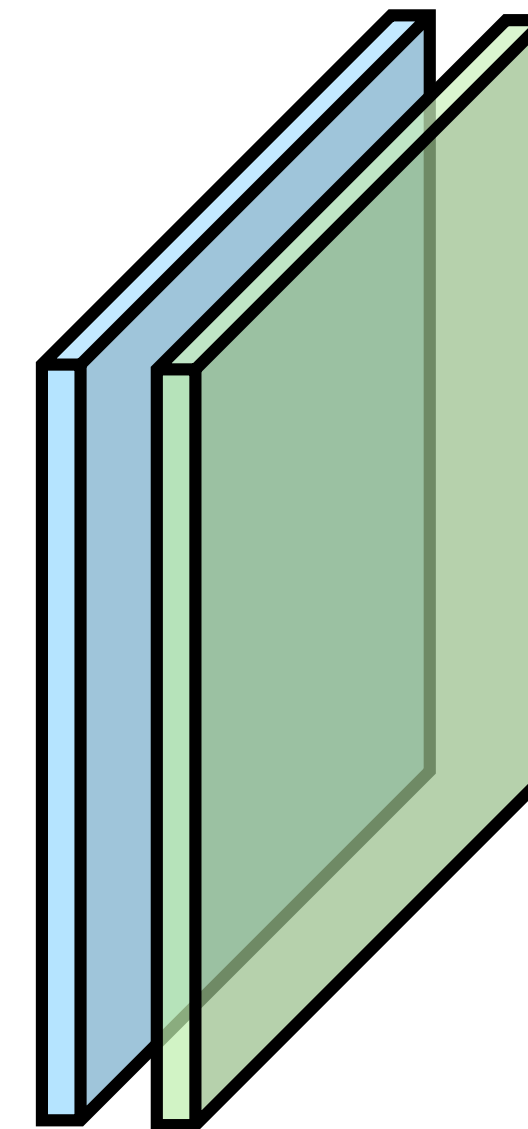
# Convolution

32x32x3 image



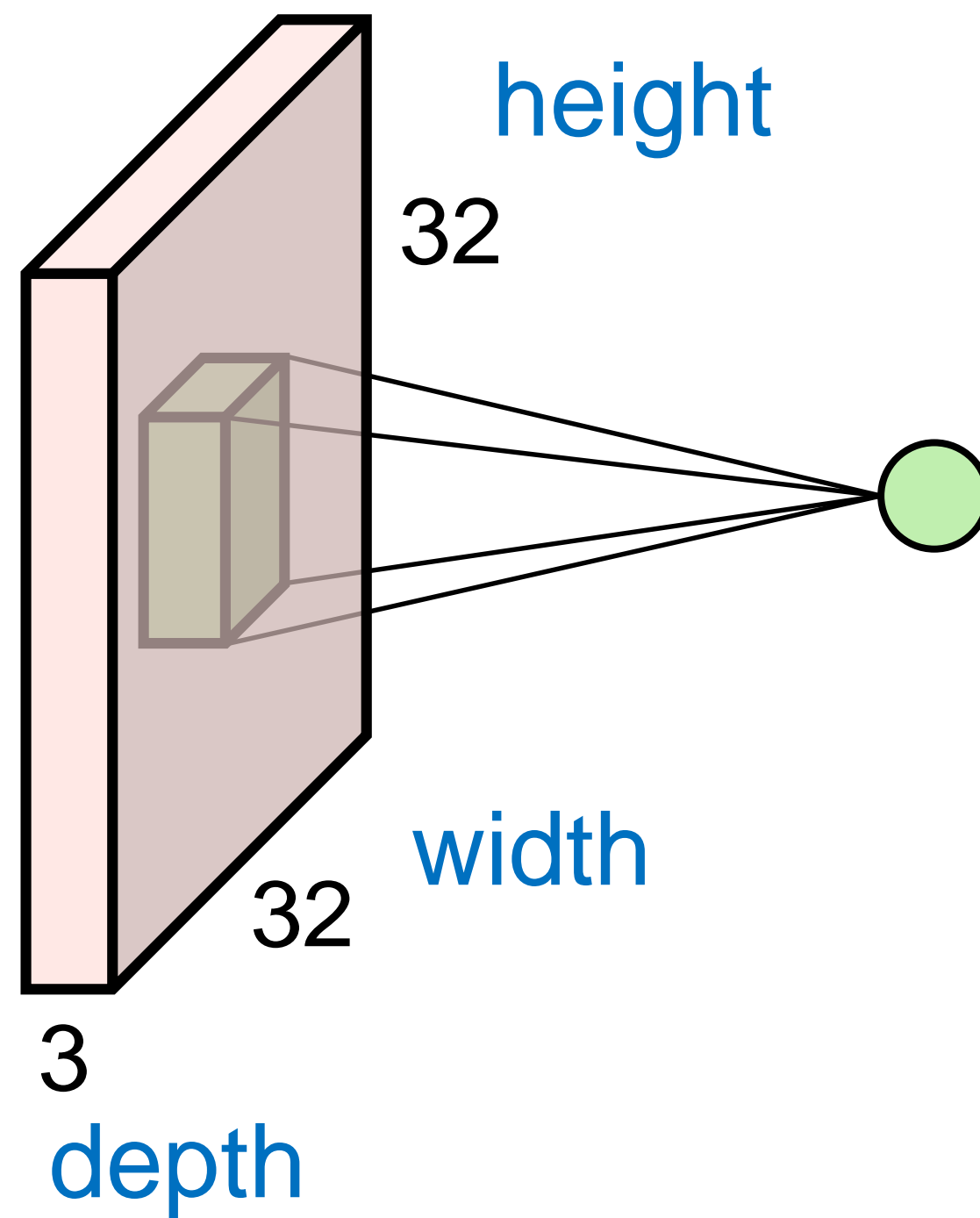
→  
Convolve (slide) over all  
spatial locations

Activation map



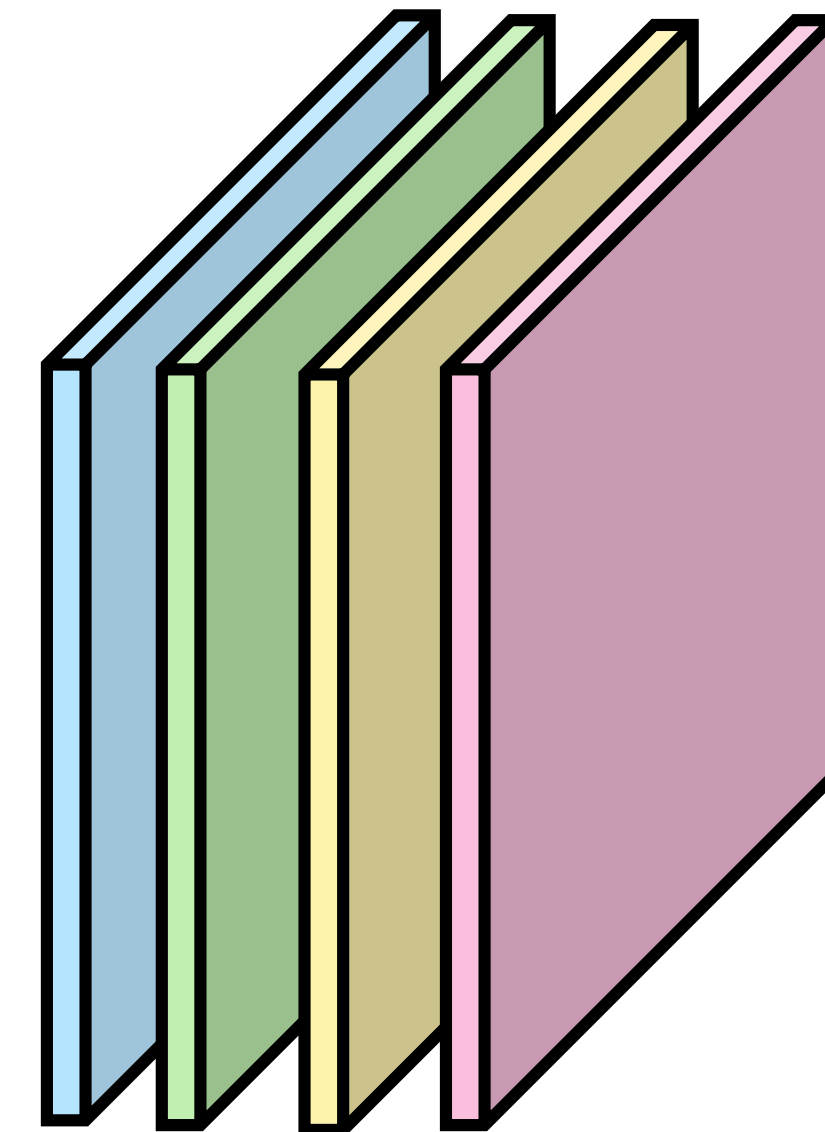
# Convolution Layer

32x32x3 image



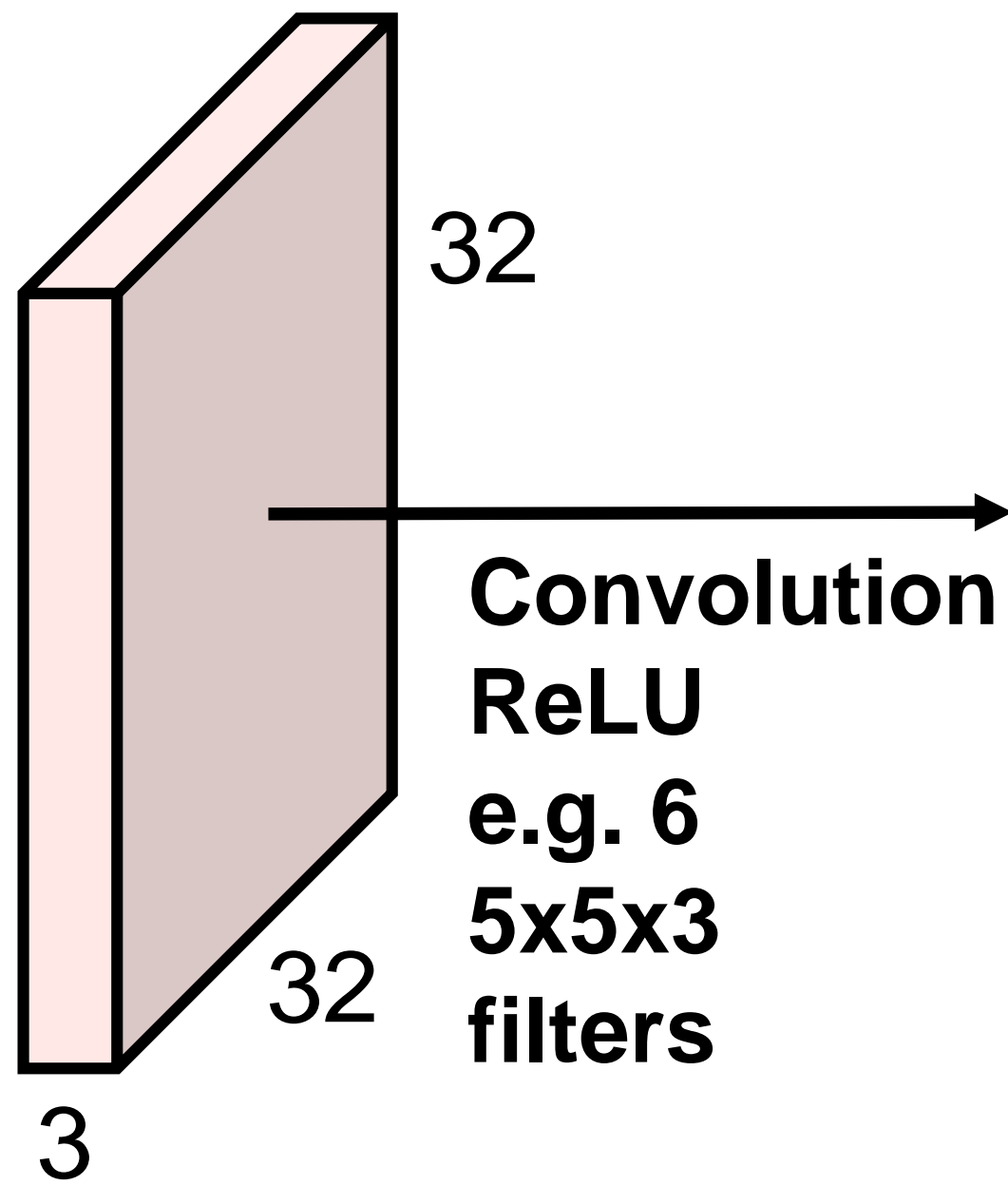
Convolution Layer

Activation tensor



# Convolutional Neural Network

32x32x3 image

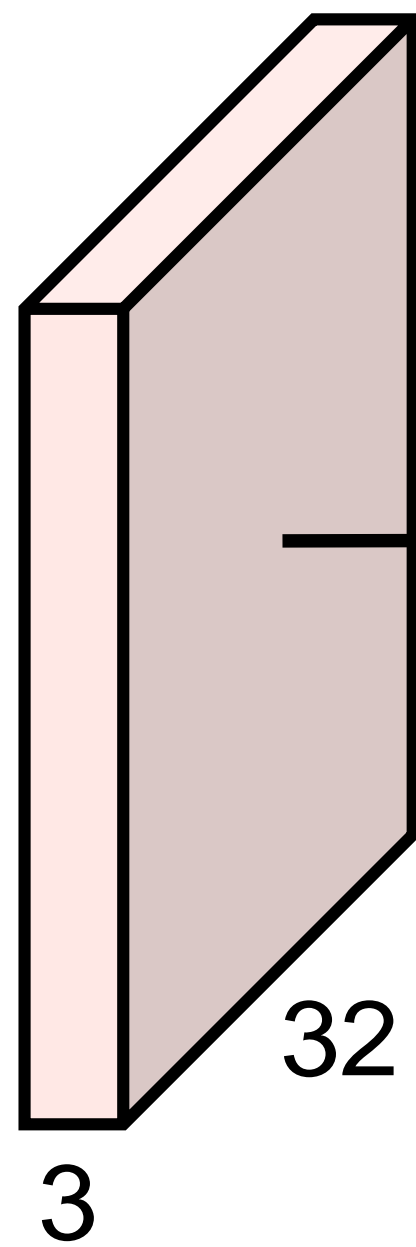


?



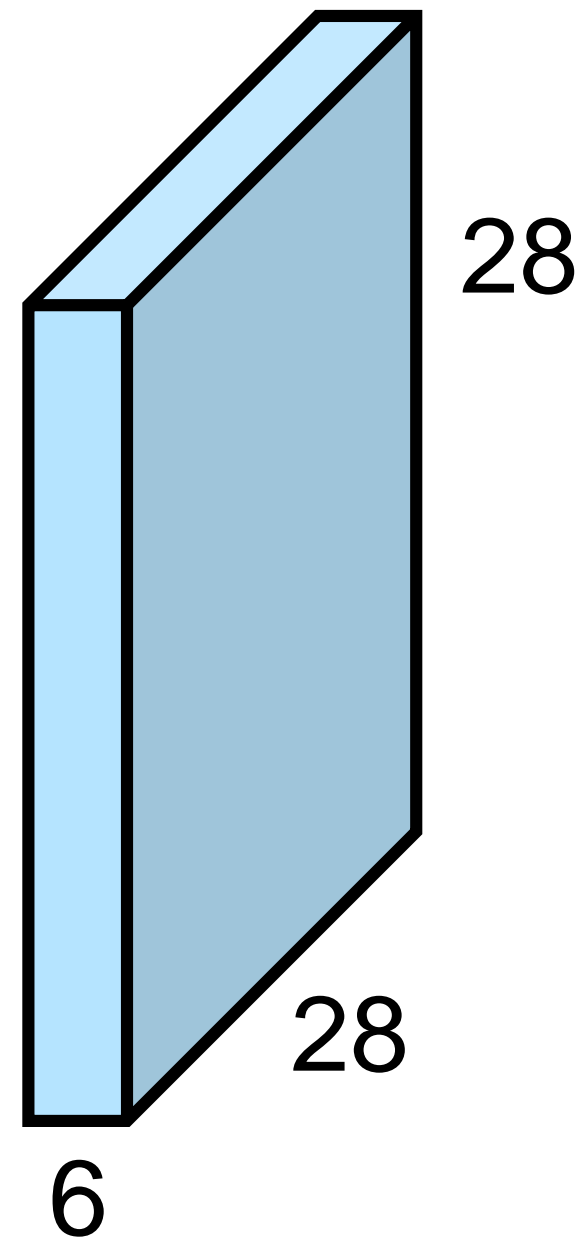
# Convolutional Neural Network

32x32x3 image



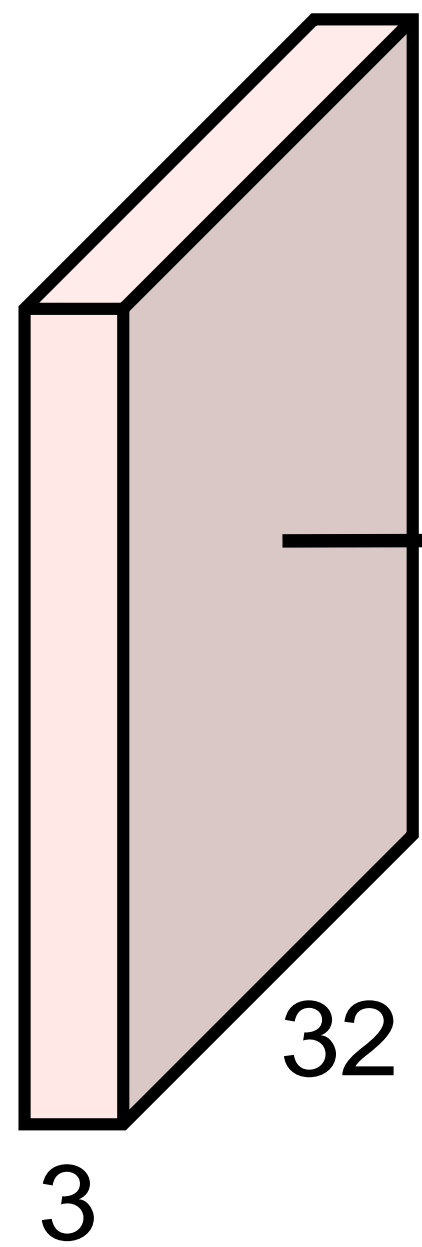
Convolution  
ReLU  
e.g. 6  
5x5x3  
filters

28x28x6 tensor



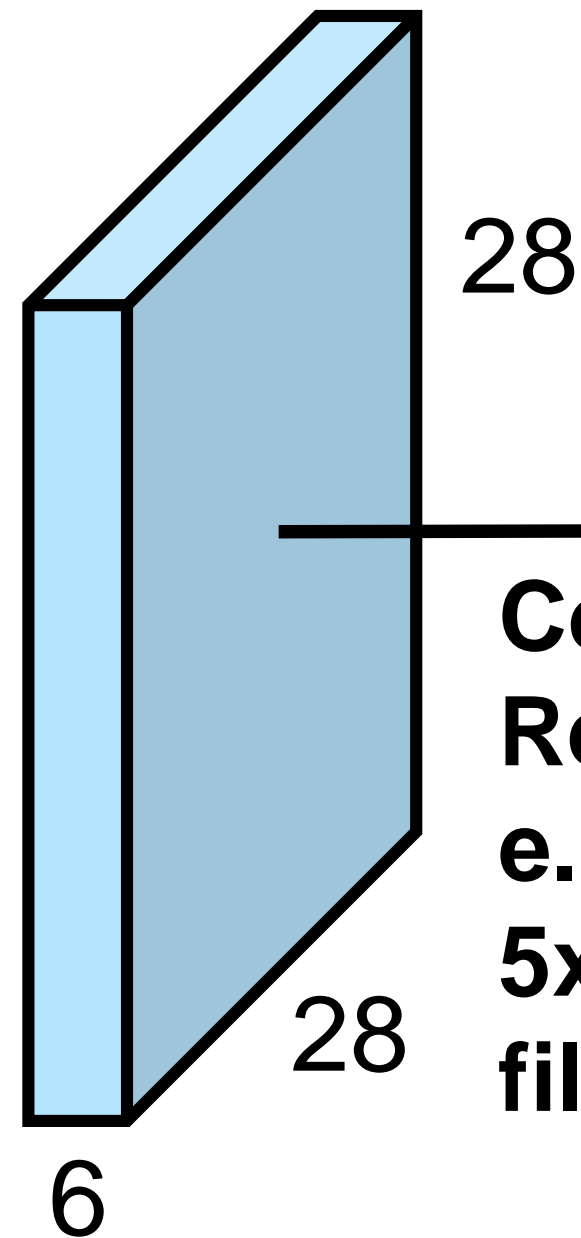
# Convolutional Neural Network

32x32x3 image

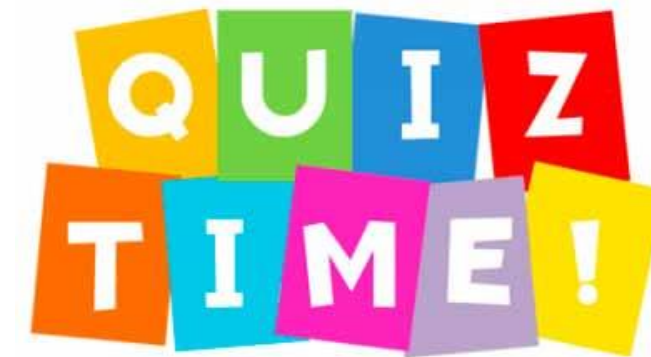


Convolution  
ReLU  
e.g. 6  
5x5x3  
filters

28x28x6 tensor



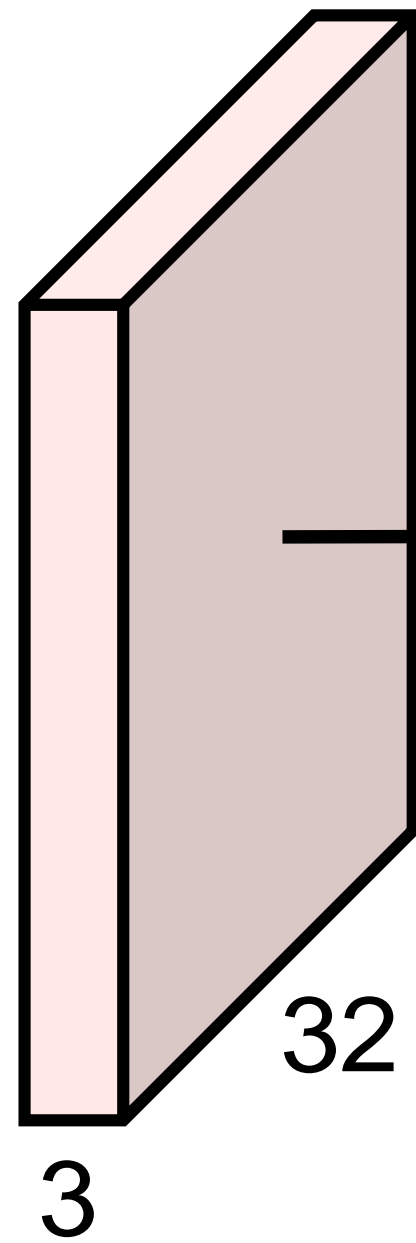
Convolution  
ReLU  
e.g. 10  
5x5x6  
filters





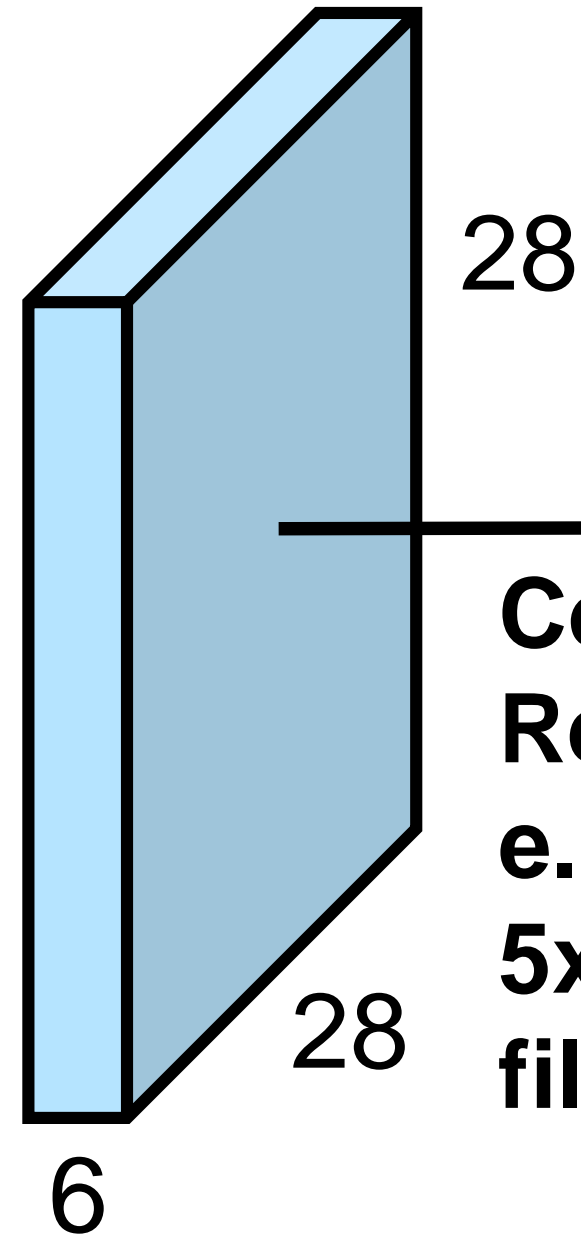
# Convolutional Neural Network

32x32x3 image



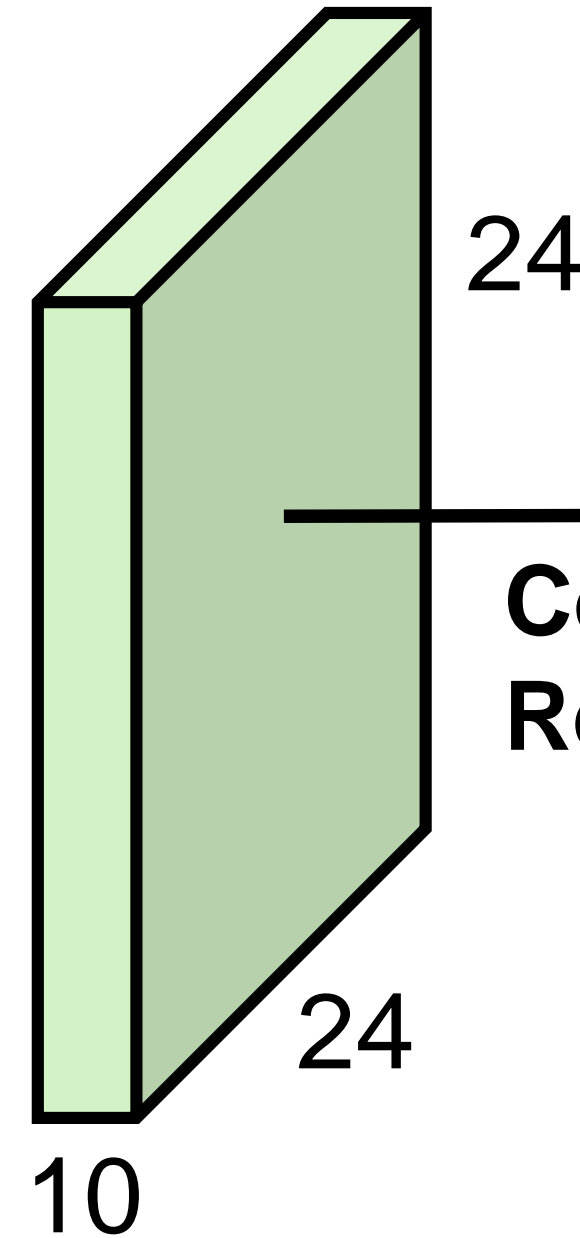
Convolution  
ReLU  
e.g. 6  
5x5x3  
filters

28x28x6 tensor



Convolution  
ReLU  
e.g. 10  
5x5x6  
filters

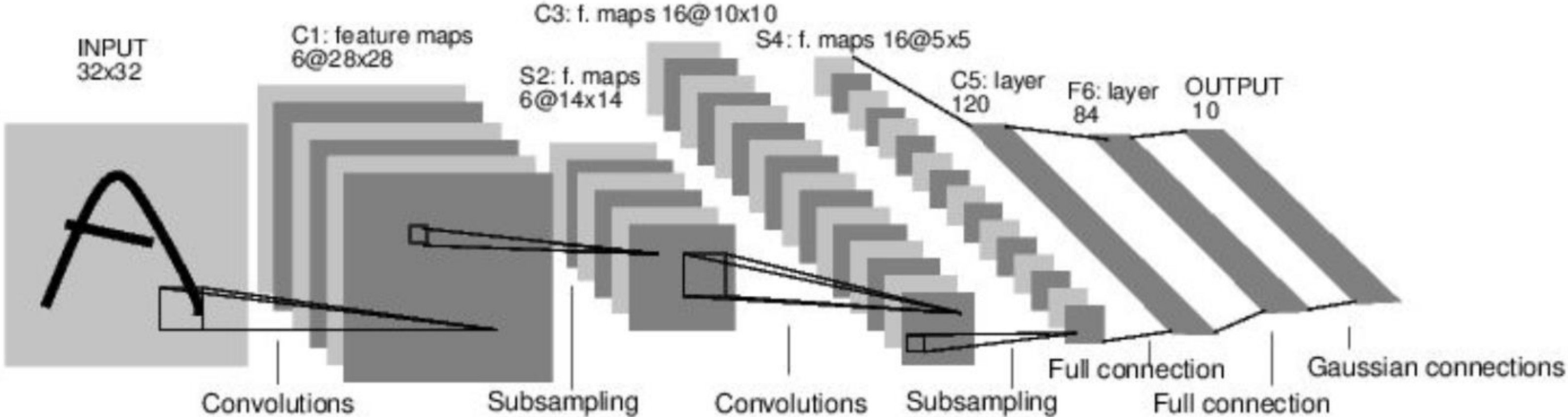
24x24x10 tensor



Convolution  
ReLU

...

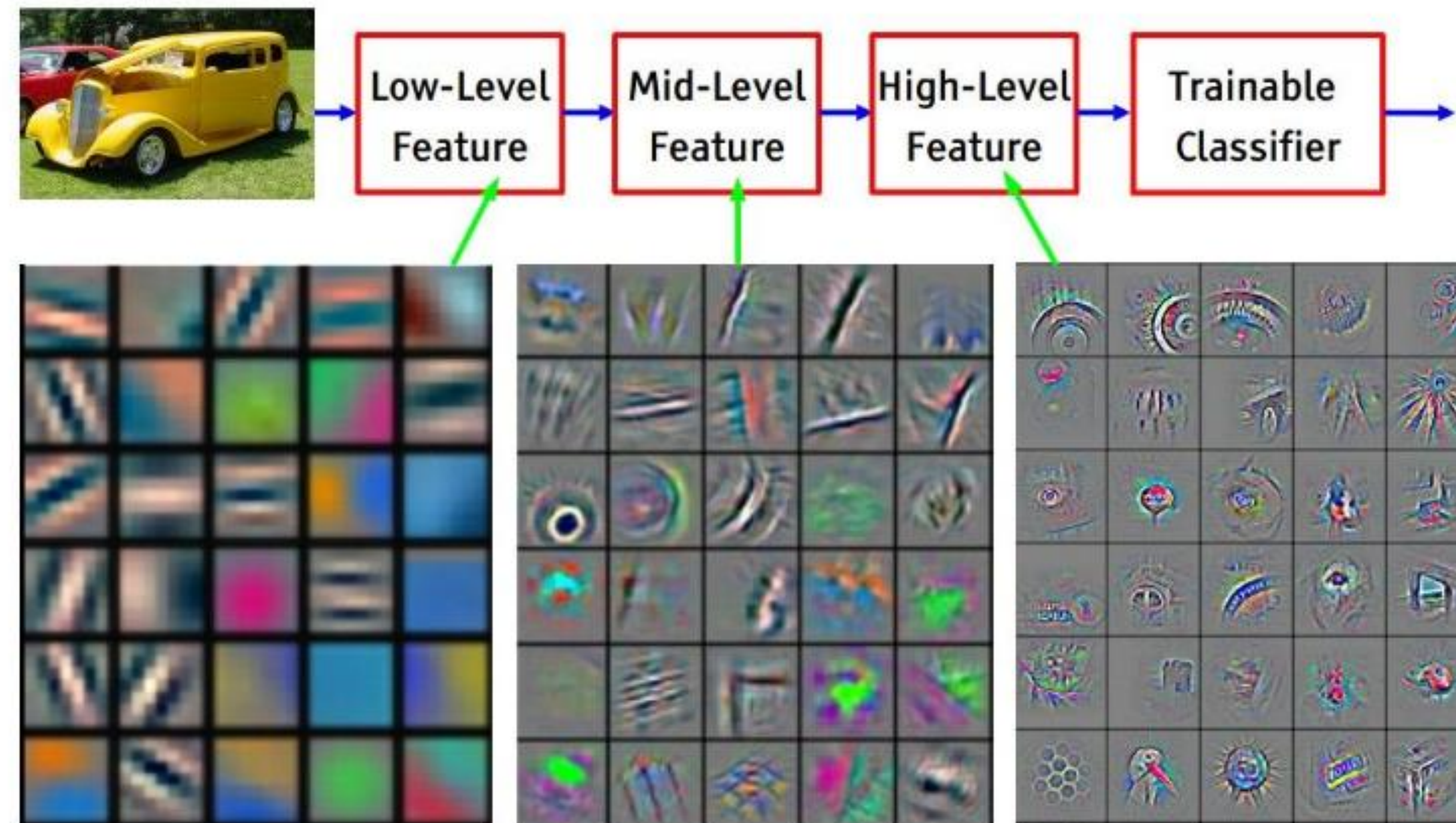
# Convolutional Neural Networks



[LeNet-5, LeCun 1980]

# Feature Hierarchy

[From recent Yann LeCun slides]

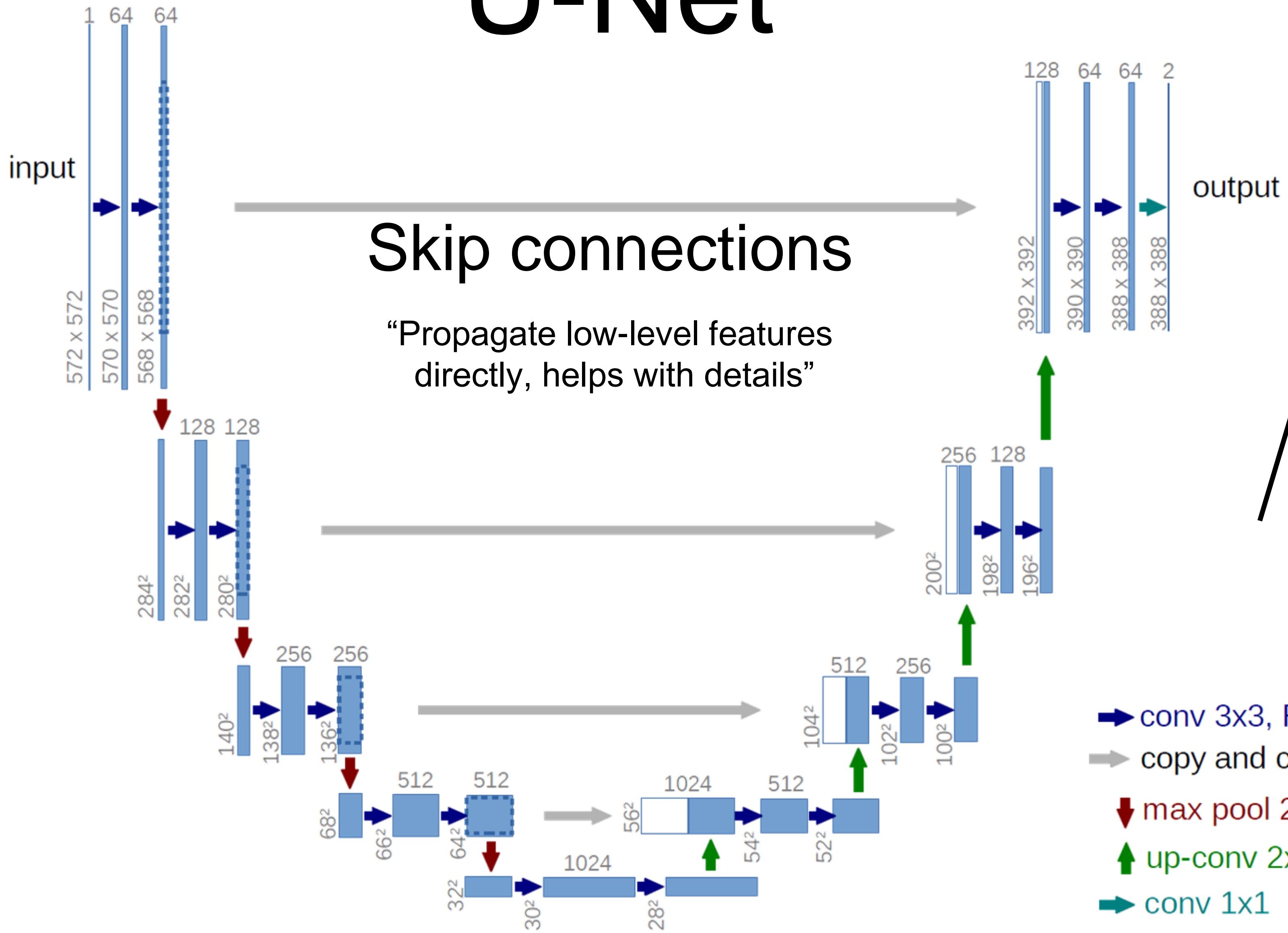


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

**Learn the features from data instead of hand engineering them!  
(If enough data is available)**

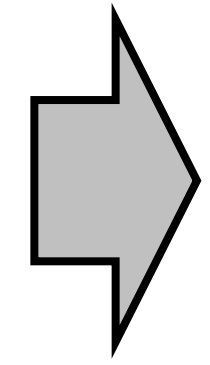
# U-Net

Downsampling



Upsampling

# Overview



- Convolutional Neural Networks

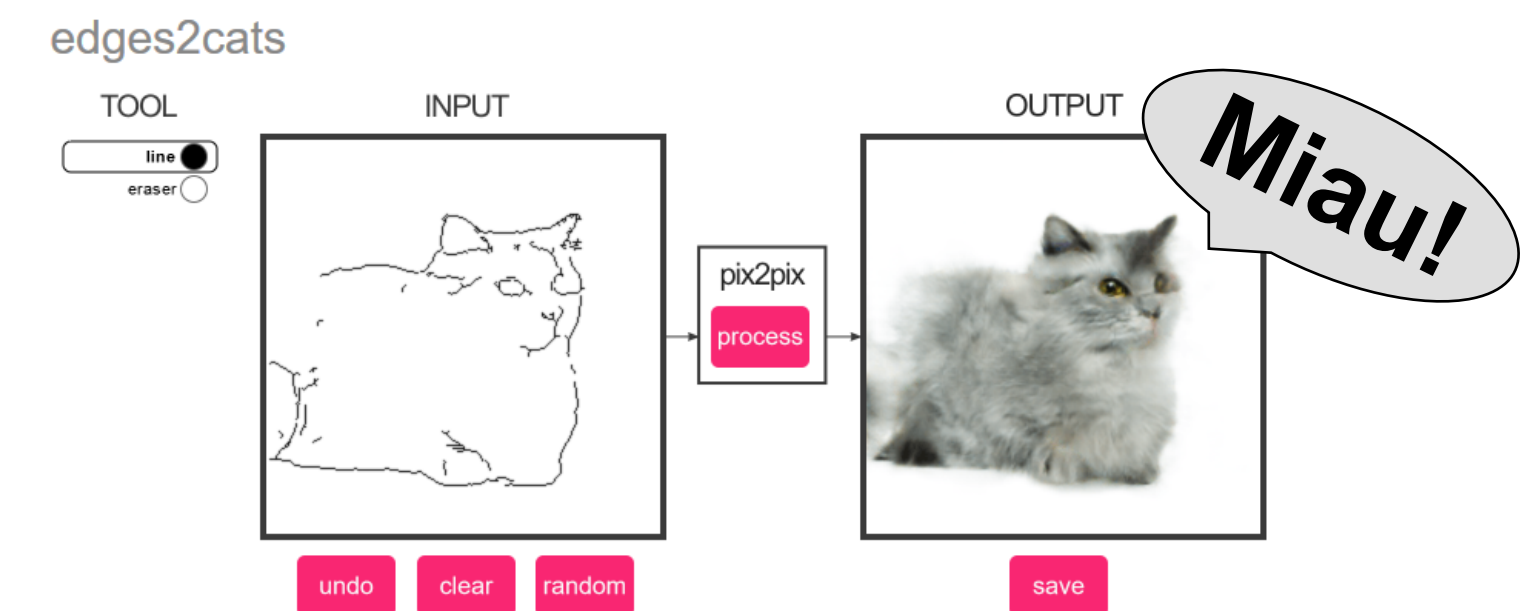
- Generative Modeling

- Pix2Pix

$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$



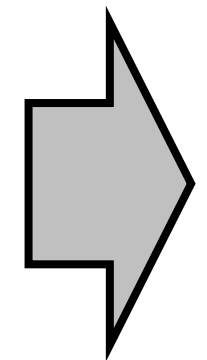
(Brundage et al., 2018)



# Overview

- Convolutional Neural Networks

$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

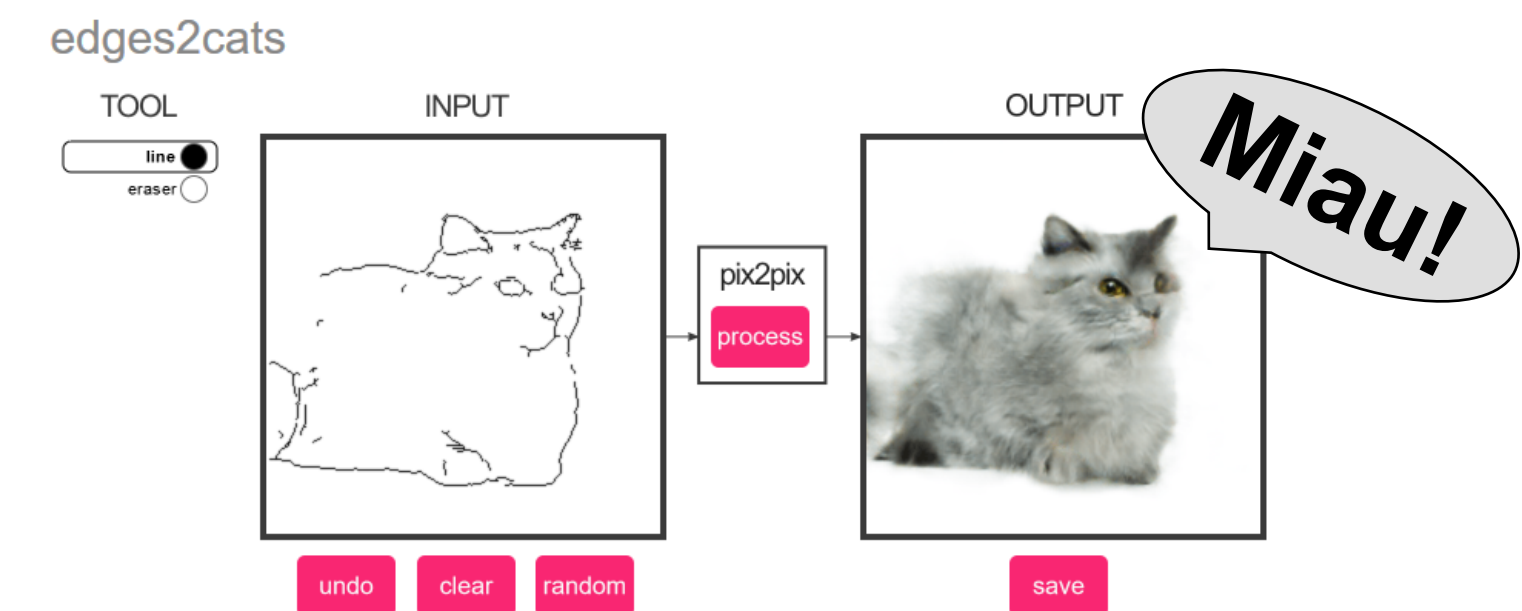


- Generative Modeling

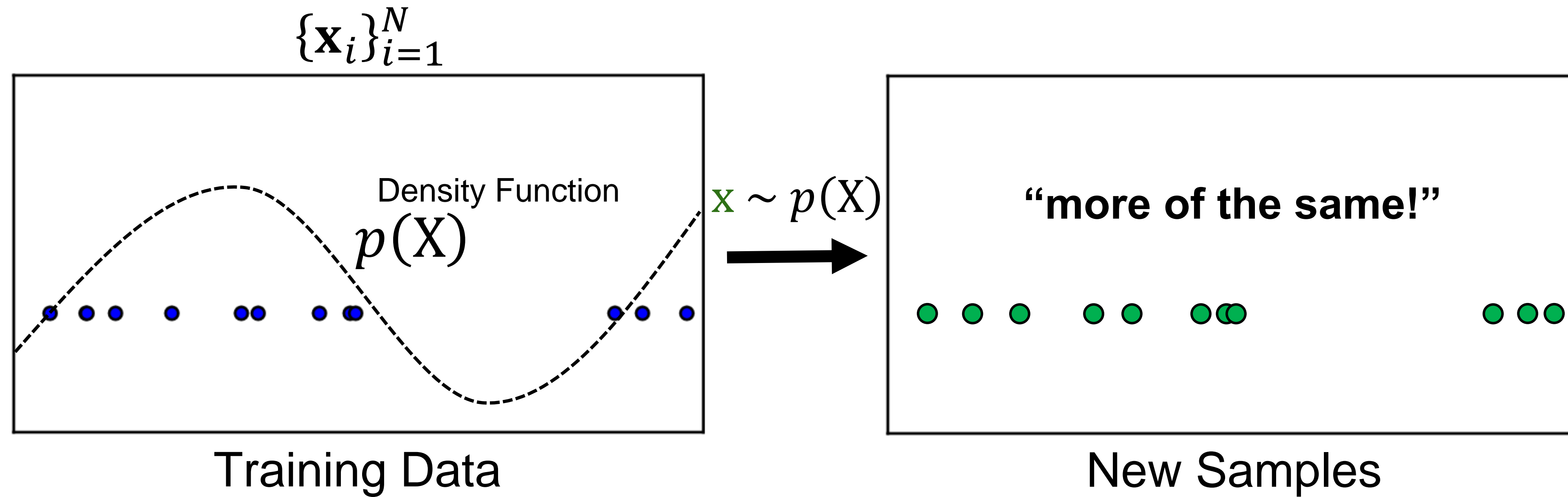


(Brundage et al., 2018)

- Pix2Pix

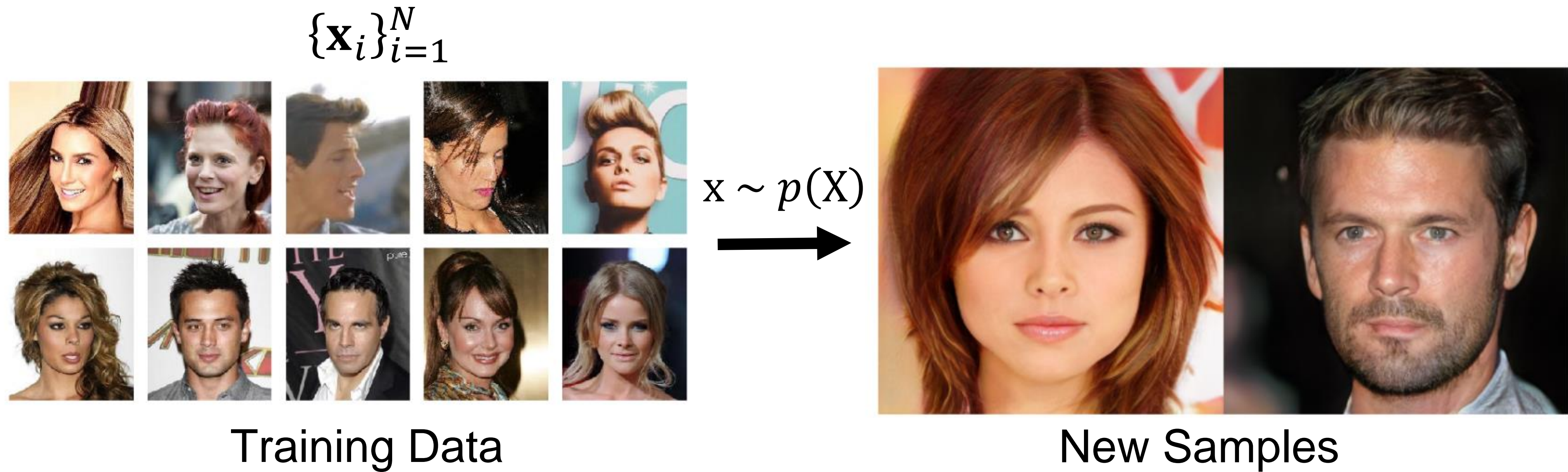


# Generative Modeling



We want to learn  $p(X)$  from data, such that we can “sample from it”!

# Generative 2D Face Modeling



The world needs more celebrities  
... or not ... ?



# 3.5 Years of Progress on Faces



2014



2015



2016



2017

(Brundage et al, 2018)

<https://thispersondoesnotexist.com>



**A Style-Based Generator Architecture for Generative Adversarial Networks**

Tero Karras  
NVIDIA

tkarras@nvidia.com

Samuli Laine  
NVIDIA

slaine@nvidia.com

Timo Aila  
NVIDIA

taila@nvidia.com



2018

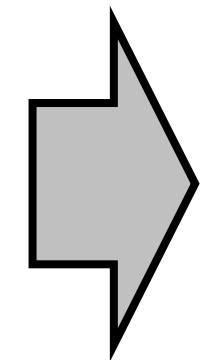
# StyleGAN - Interpolation



# Overview

- Convolutional Neural Networks

$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

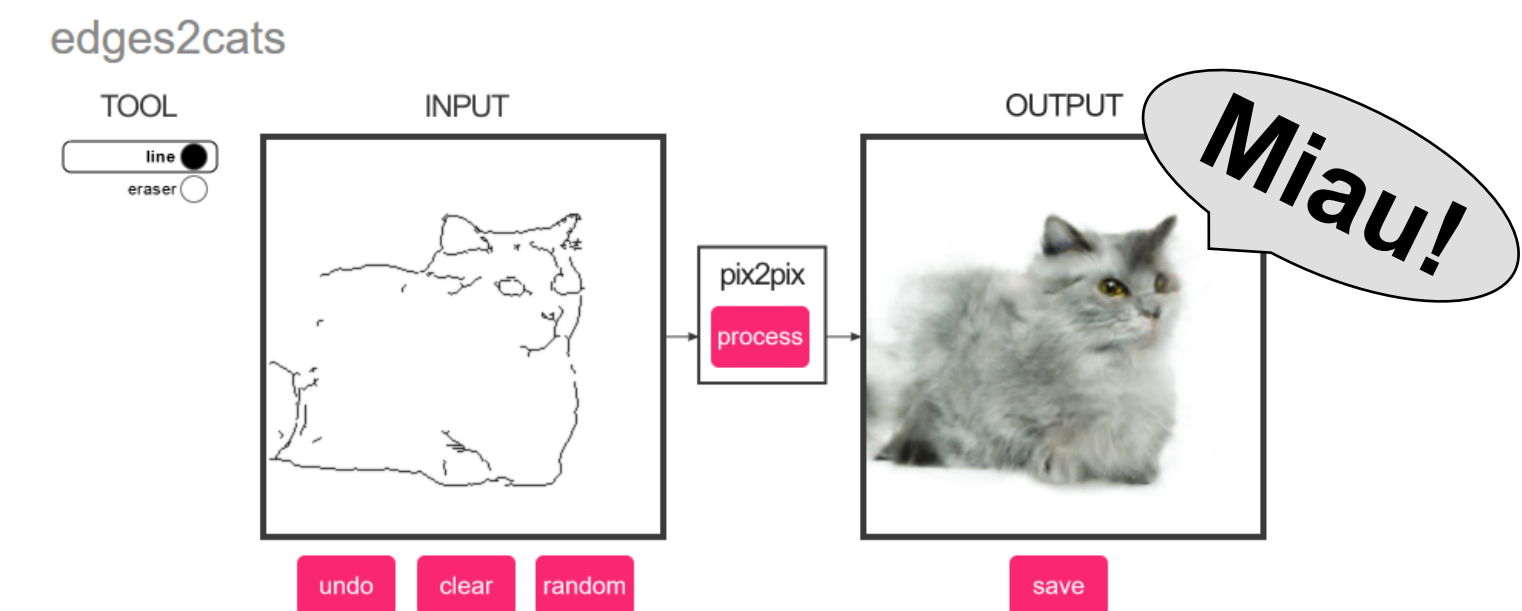


- Generative Modeling



(Brundage et al., 2018)

- Pix2Pix (“mapping from A to B”)



# Overview

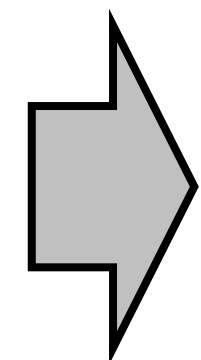
- Convolutional Neural Networks

$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

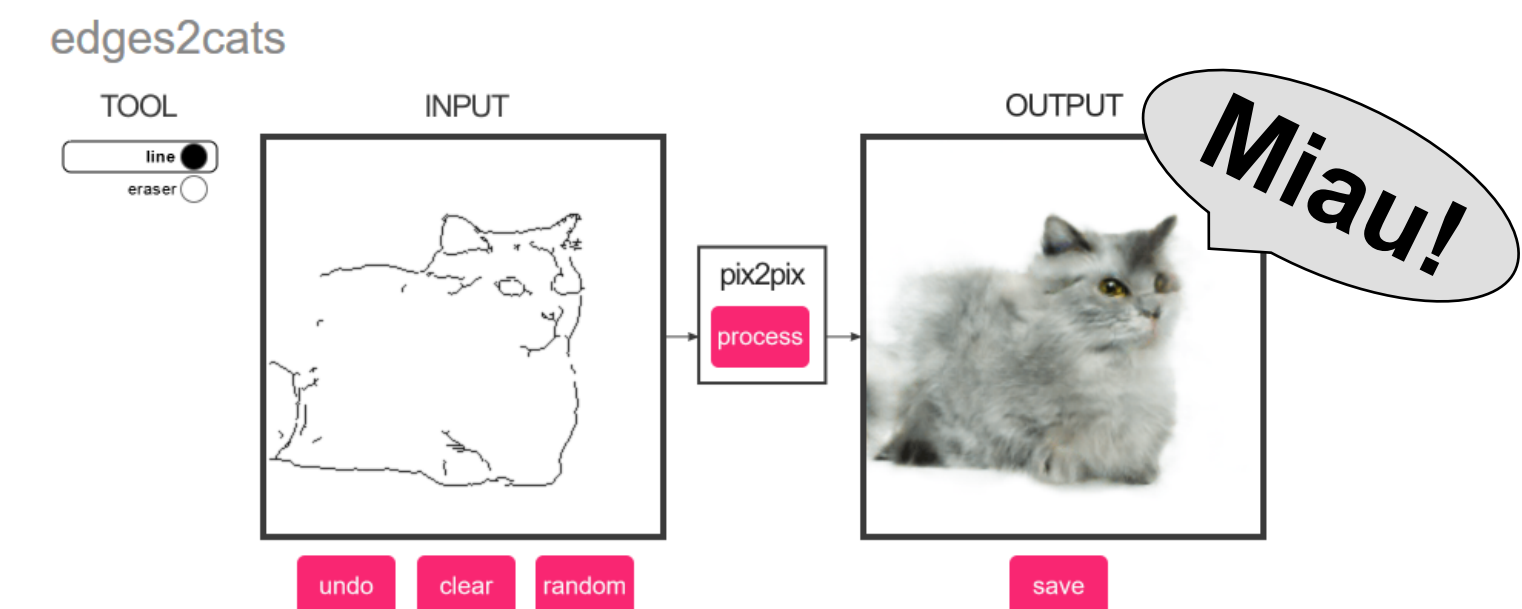
- Generative Modeling



(Brundage et al., 2018)



- Pix2Pix (“mapping from A to B”)



# Image-to-Image Translation with Conditional Adversarial Networks

Phillip Isola

Jun-Yan Zhu

Tinghui Zhou

Alexei A. Efros

Berkeley AI Research (BAIR) Laboratory, UC Berkeley

{isola, junyanz, tinghuiz, efros}@eecs.berkeley.edu

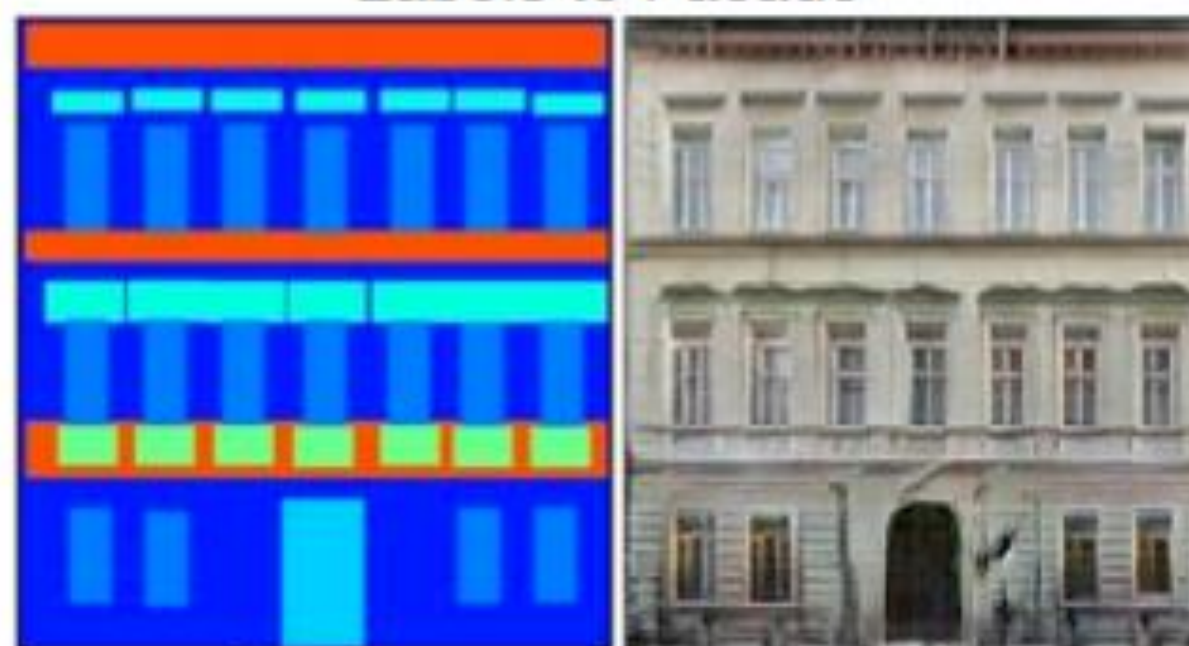
Labels to Street Scene



input

output

Labels to Facade



input

output

BW to Color



input

output

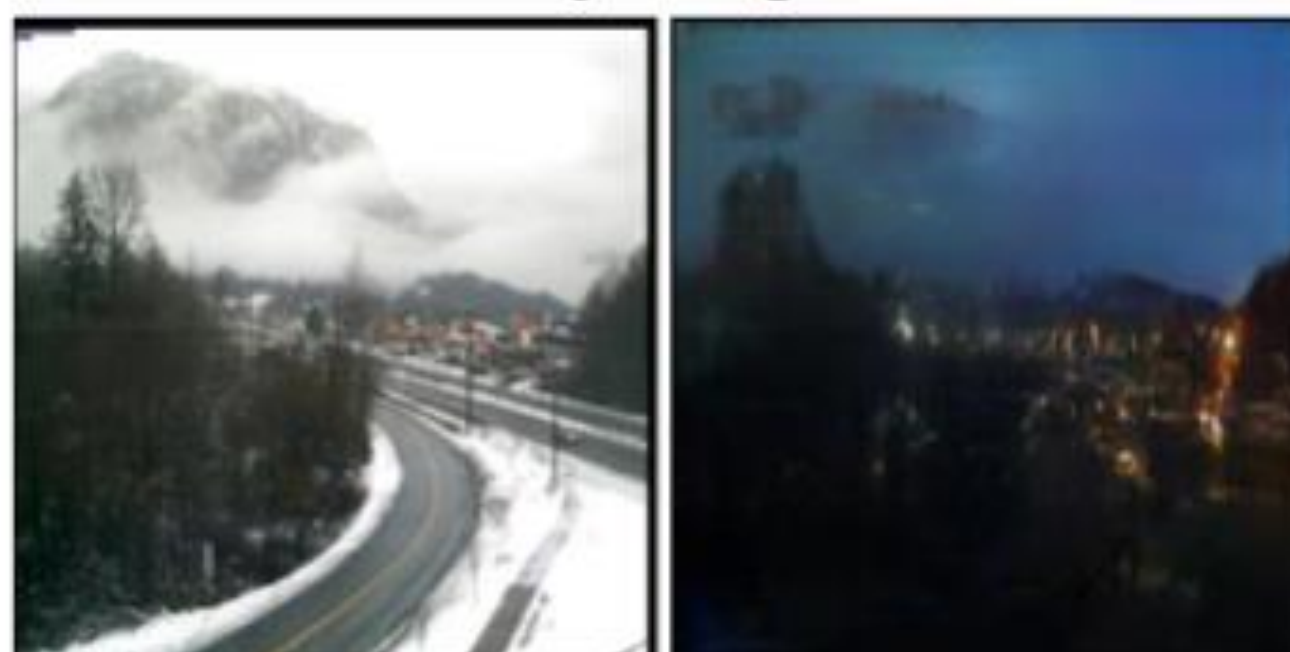
Aerial to Map



input

output

Day to Night



input

output

Edges to Photo

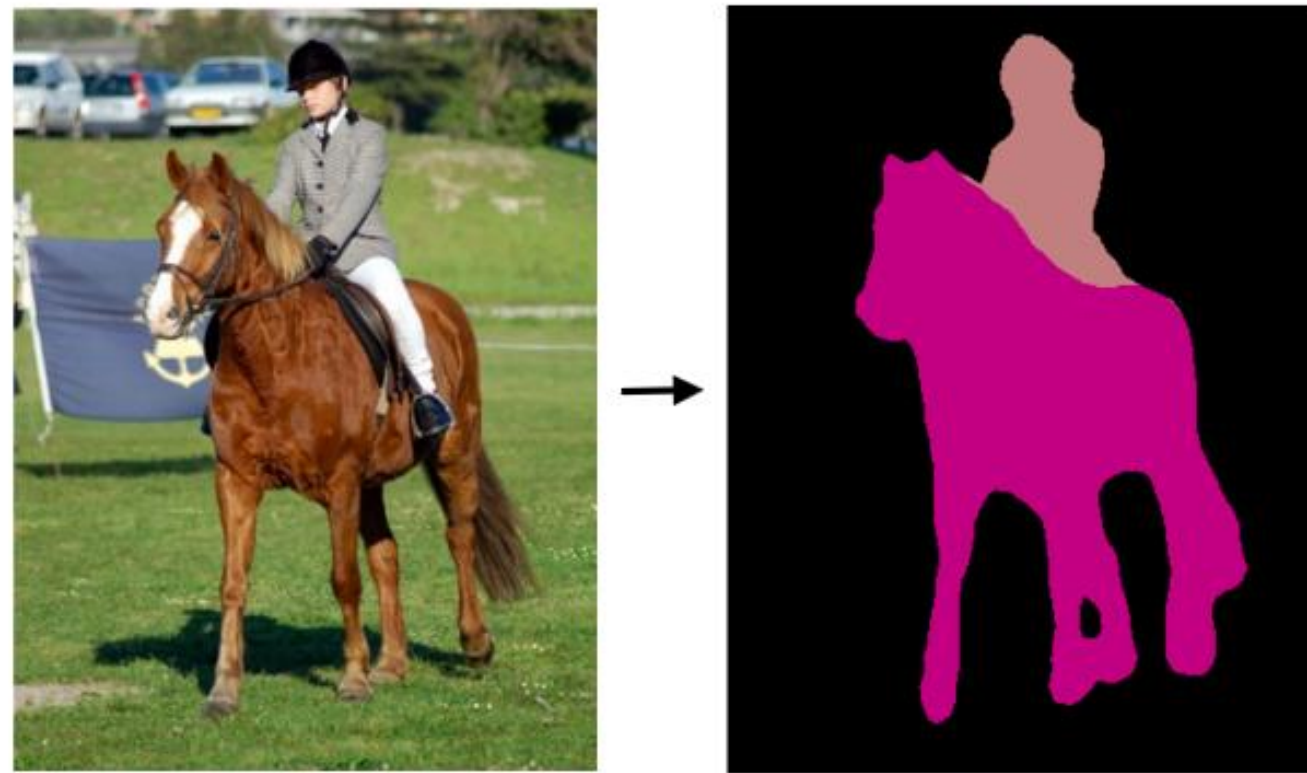


input

output

# Image-to-Image Translation

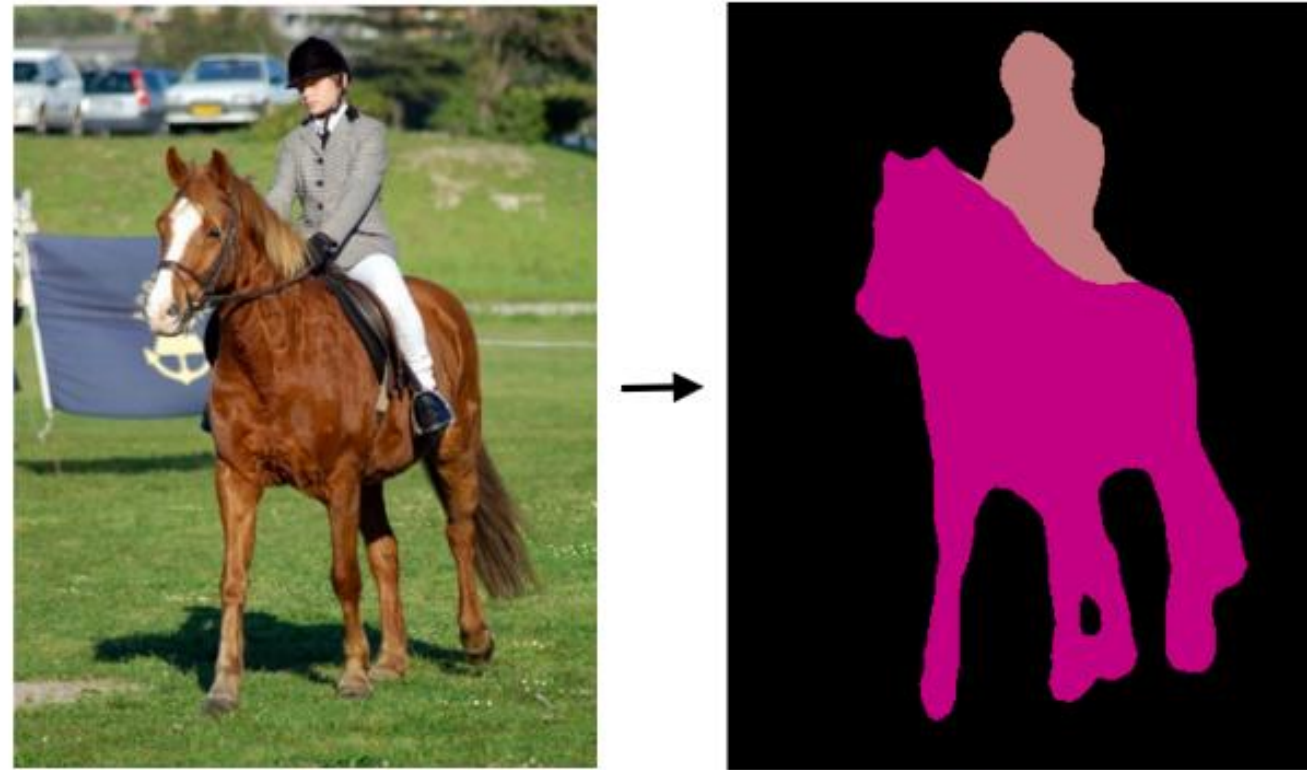
Object labeling



[Long et al. 2015]

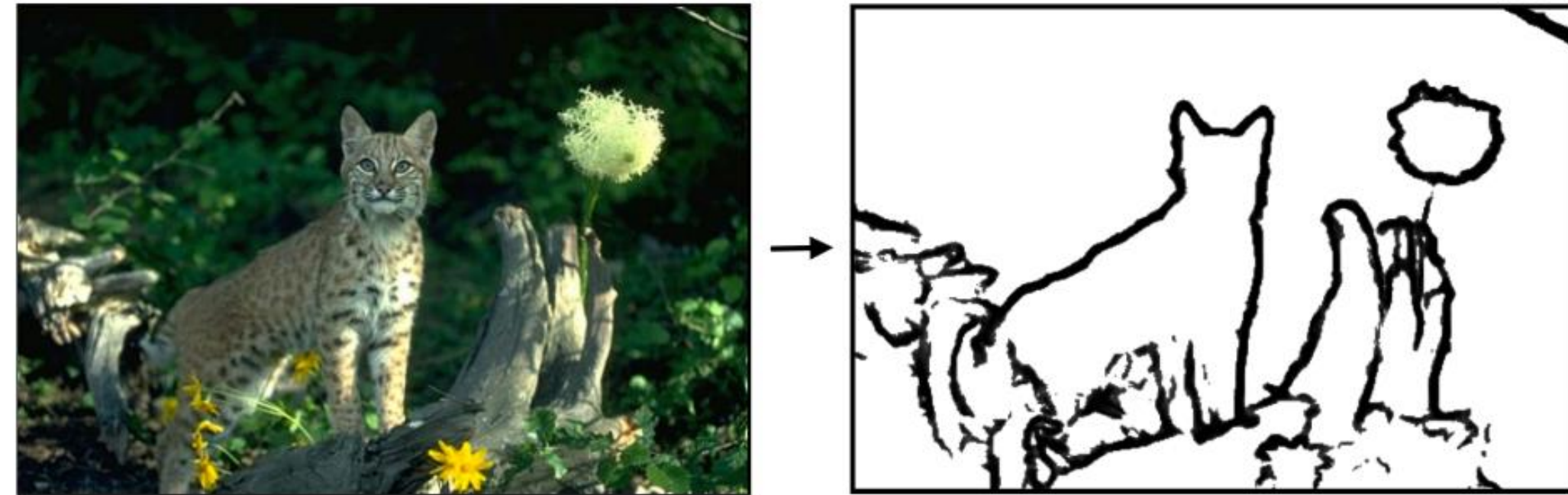
# Image-to-Image Translation

Object labeling



[Long et al. 2015]

Edge Detection

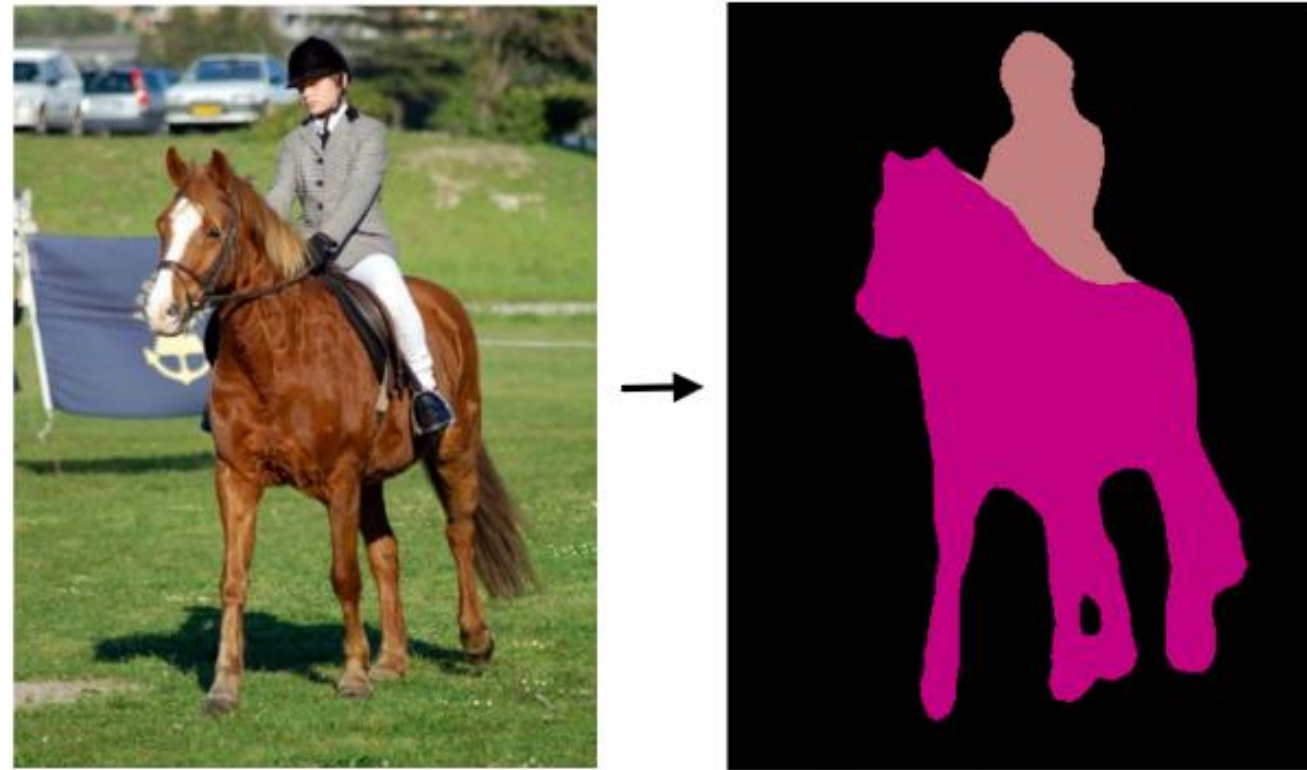


[Xie et al. 2015]



# Image-to-Image Translation

Object labeling



[Long et al. 2015]

Edge Detection



[Xie et al. 2015]

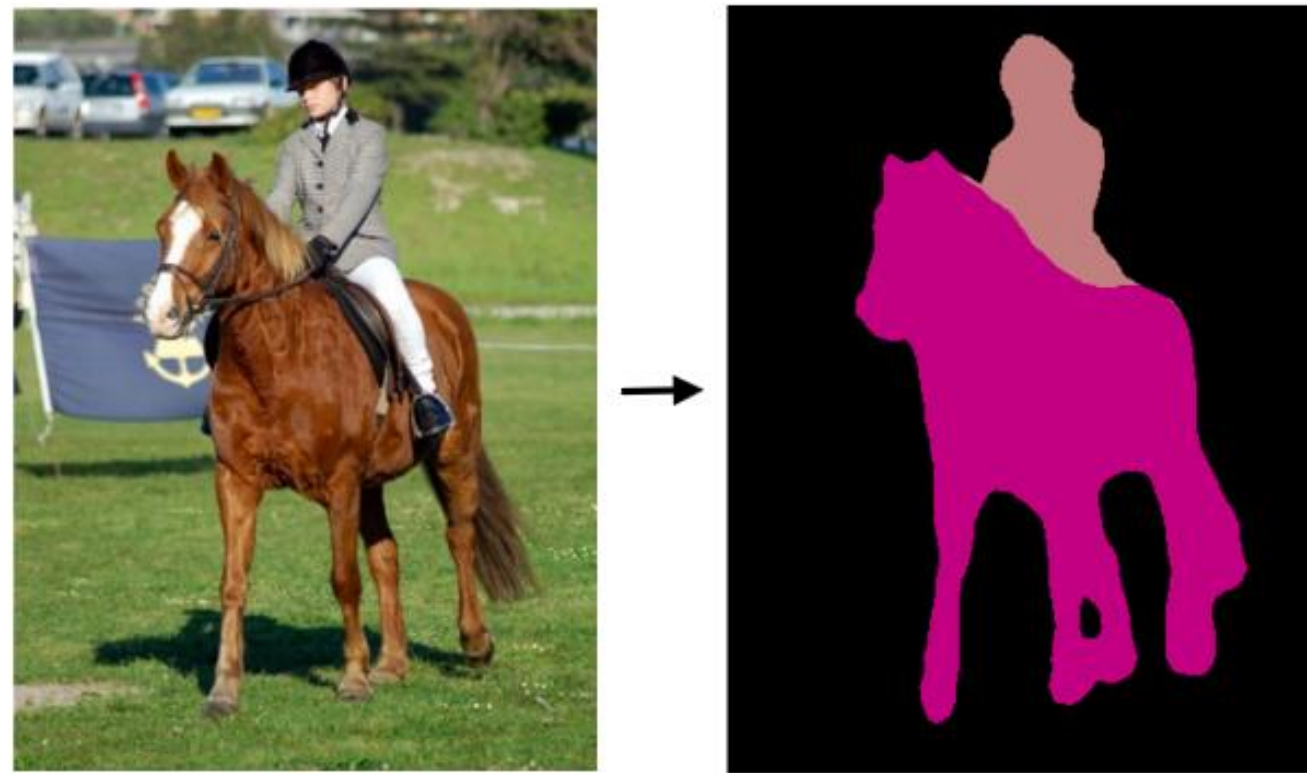
Season change



[Laffont et al. 2014]

# Image-to-Image Translation

Object labeling



[Long et al. 2015]

Edge Detection



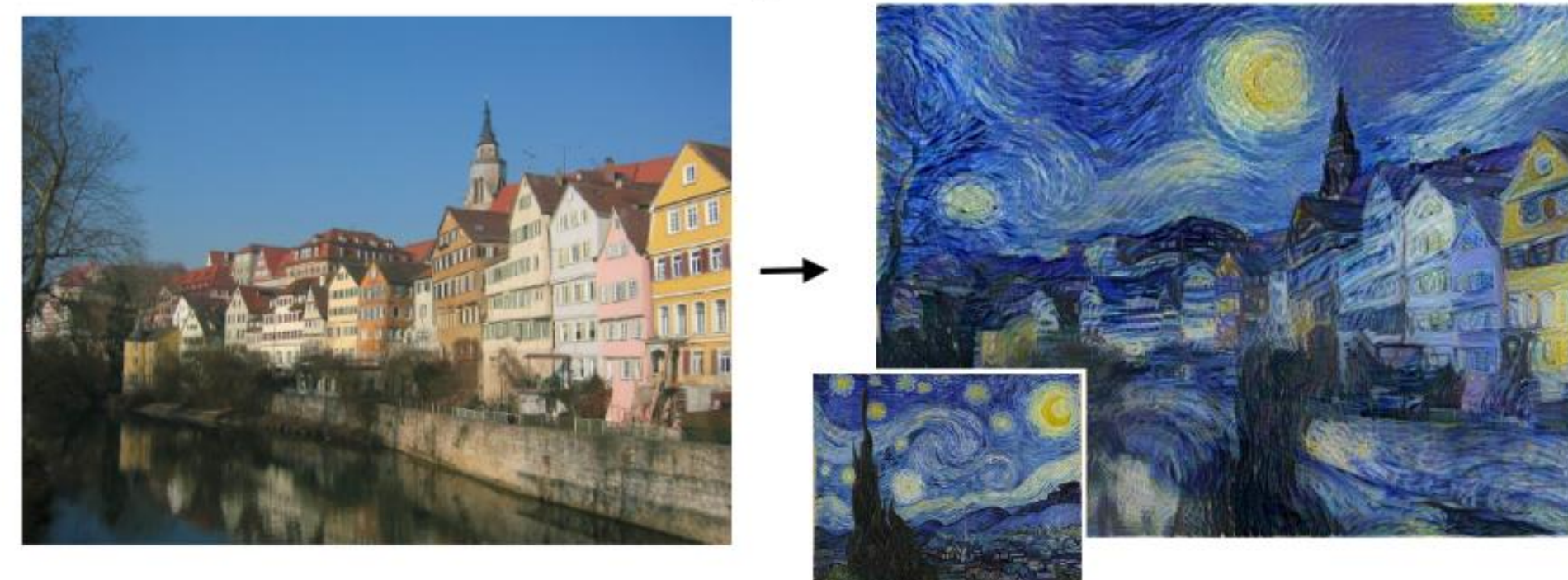
[Xie et al. 2015]

Season change



[Laffont et al. 2014]

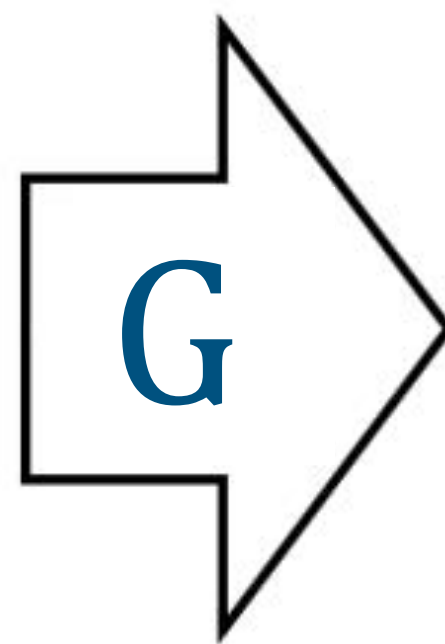
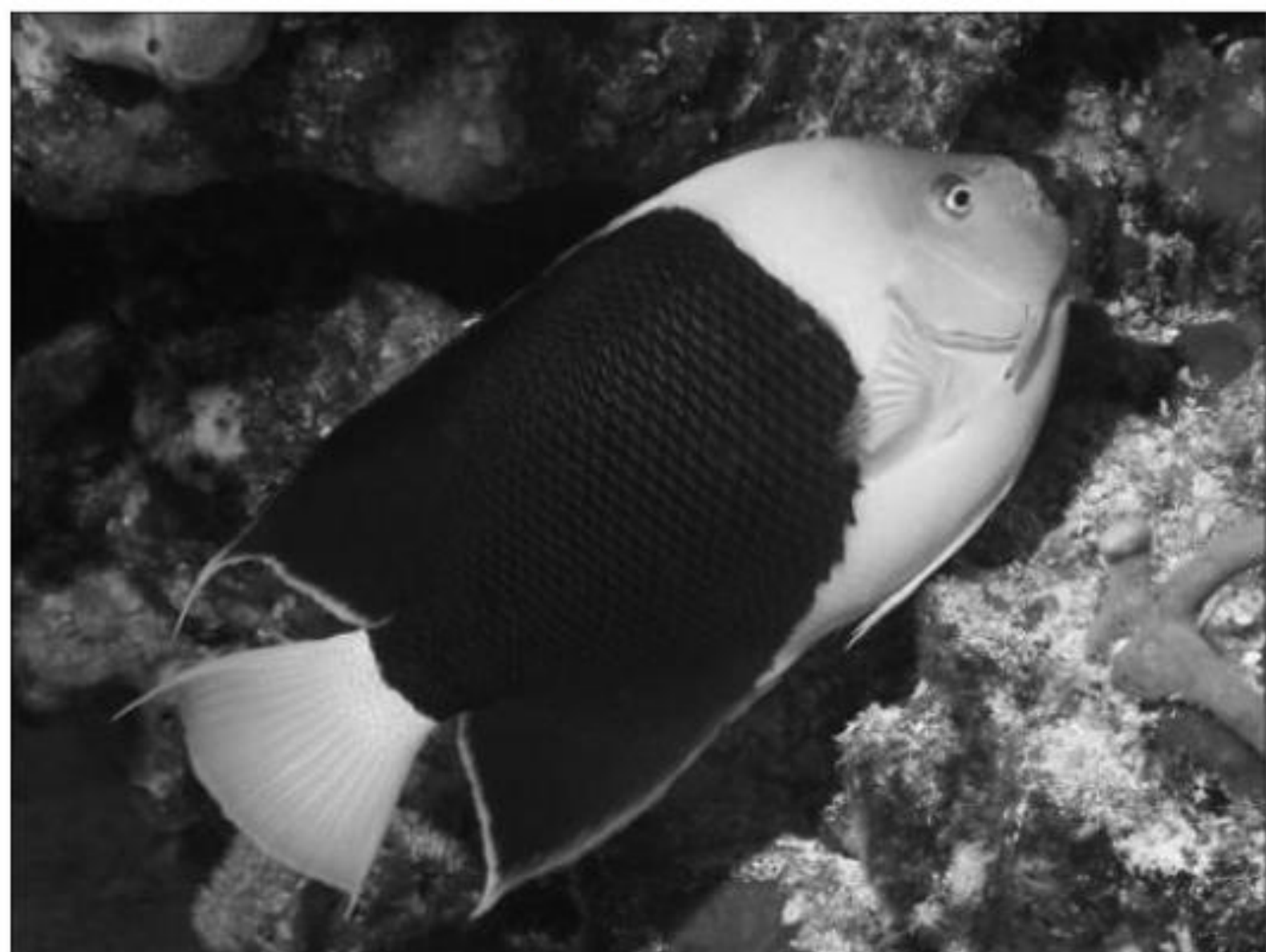
Artistic style transfer



[Gatys et al. 2016]

# Image-to-Image Translation

Input  $\mathbf{x}$

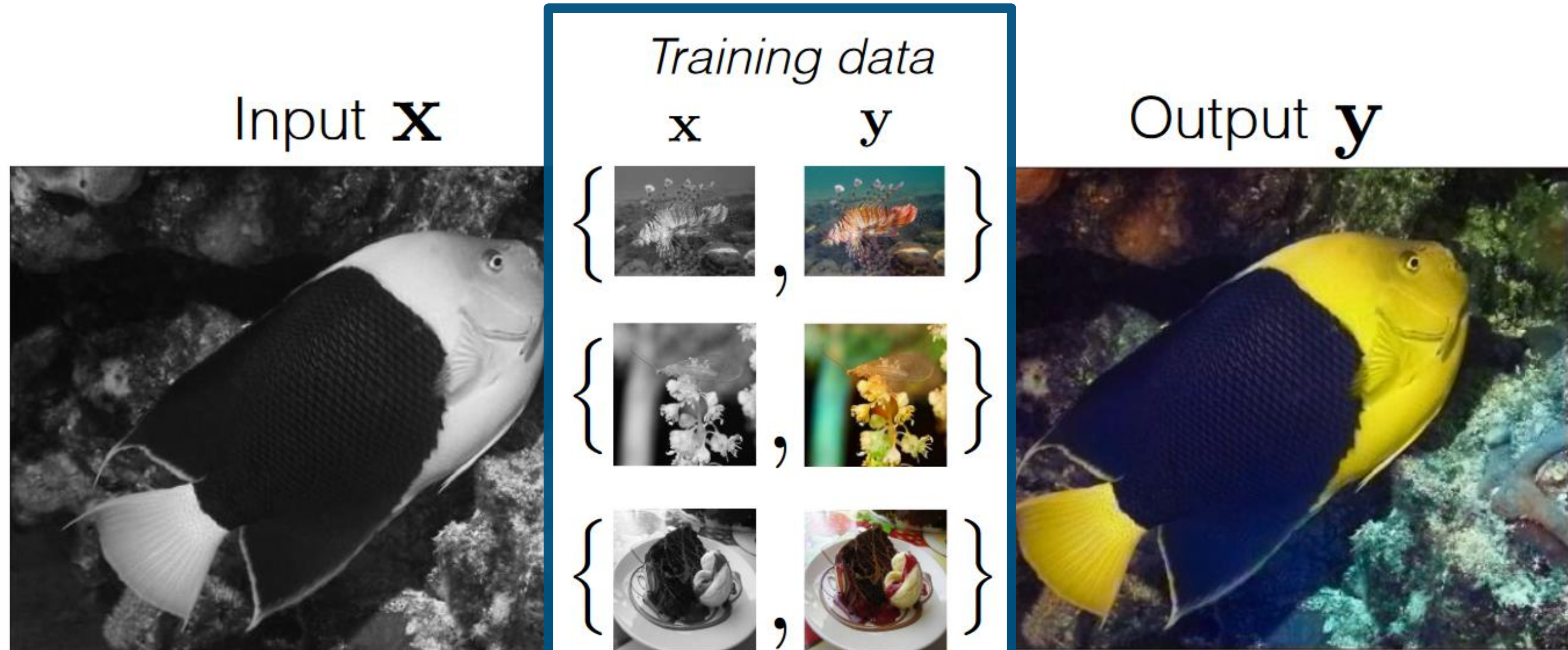


Output  $\mathbf{y}$



$$\underset{G}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\underset{\text{Loss}}{L}(\underset{\text{Neural Network}}{G}(\mathbf{x}), \mathbf{y})]$$

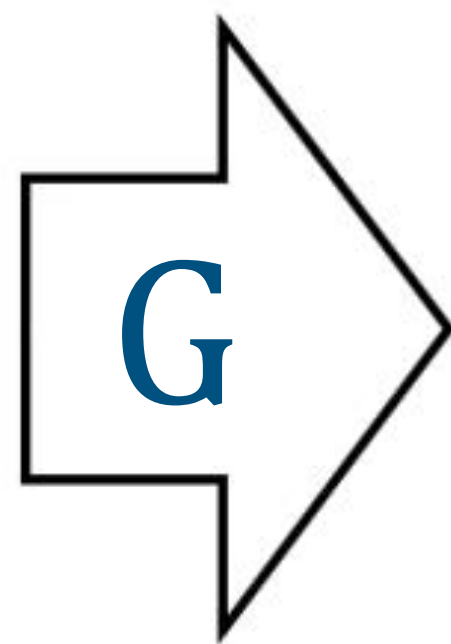
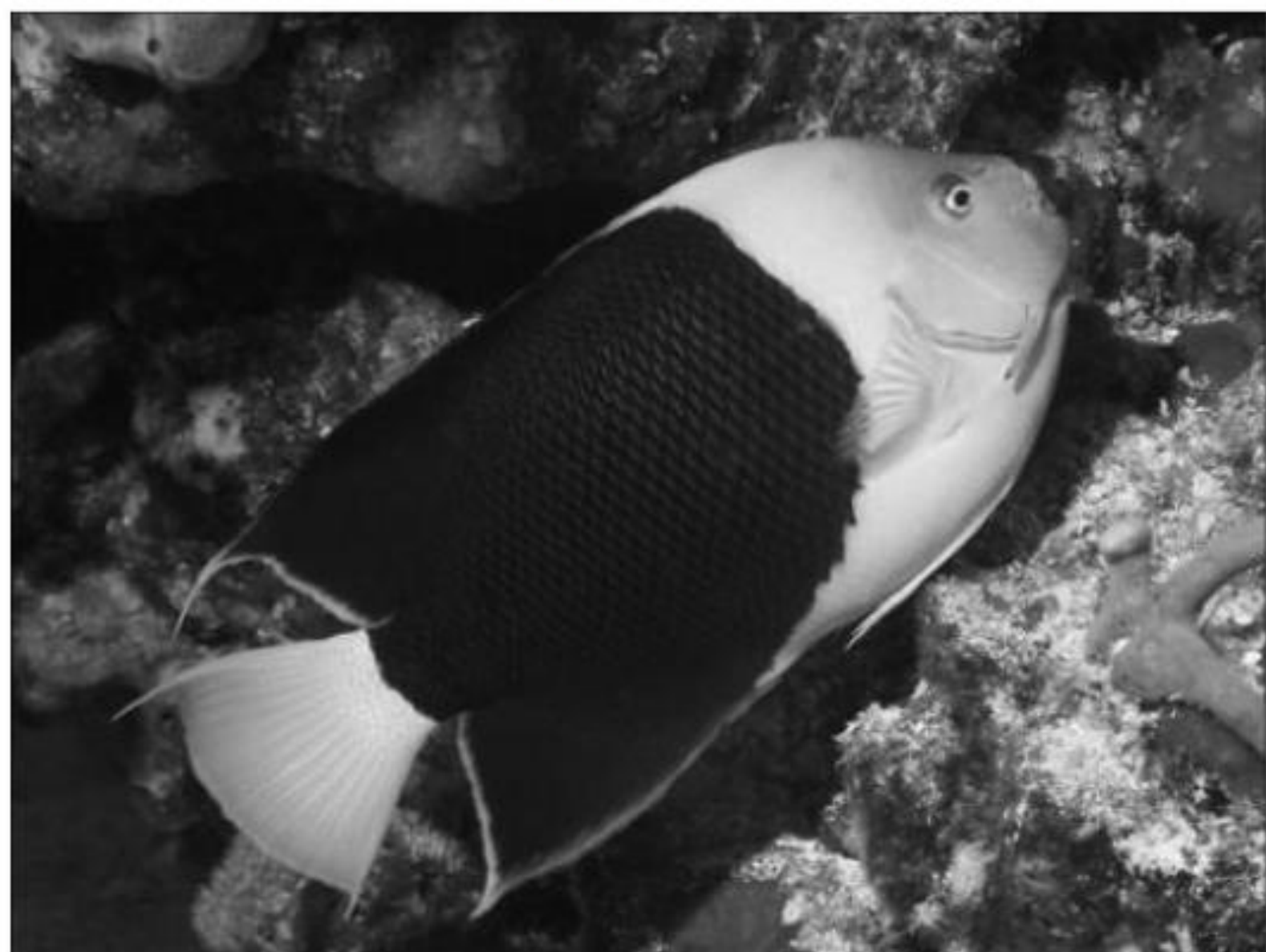
# Image-to-Image Translation



$$\underset{G}{\operatorname{argmin}} \sum_{x,y} [ \underset{\text{Loss}}{\mathcal{L}}(\underset{\text{Neural Network}}{G(x)}, y) ]$$

# Image-to-Image Translation

Input  $\mathbf{x}$



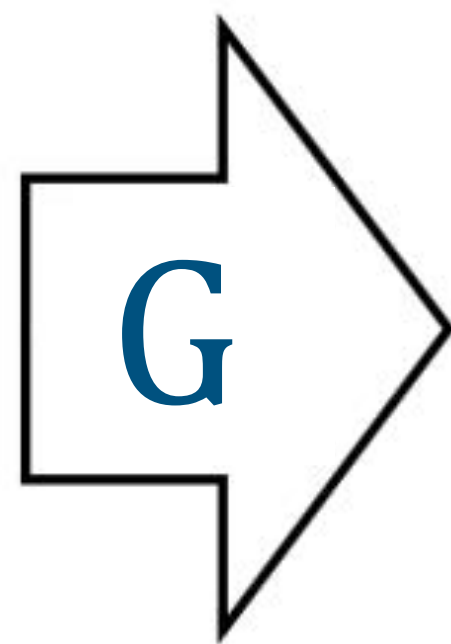
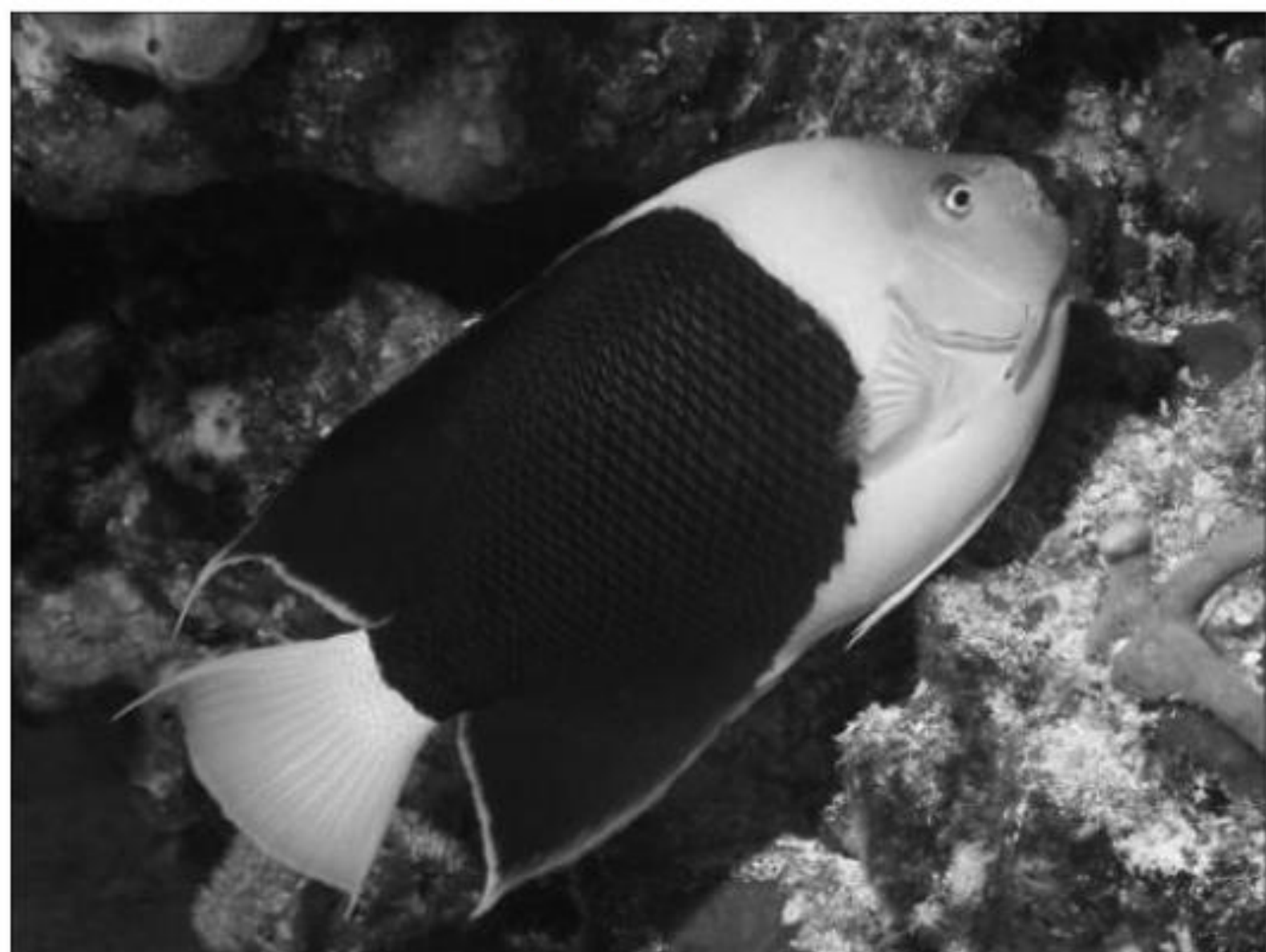
Output  $\mathbf{y}$



$$\underset{G}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\underset{\text{Loss}}{L}(\underset{\text{Neural Network}}{G}(\mathbf{x}), \mathbf{y})]$$

# Image-to-Image Translation

Input  $\mathbf{x}$



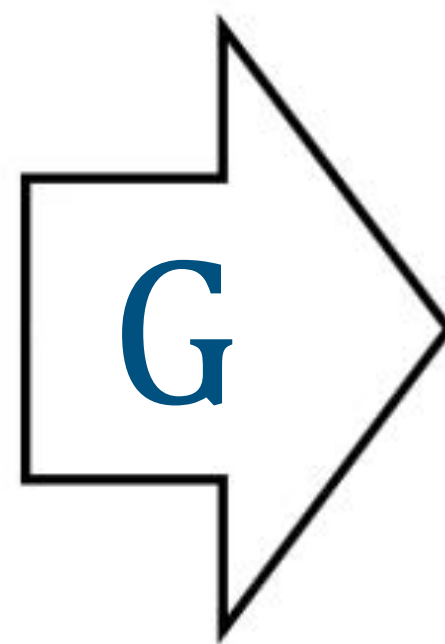
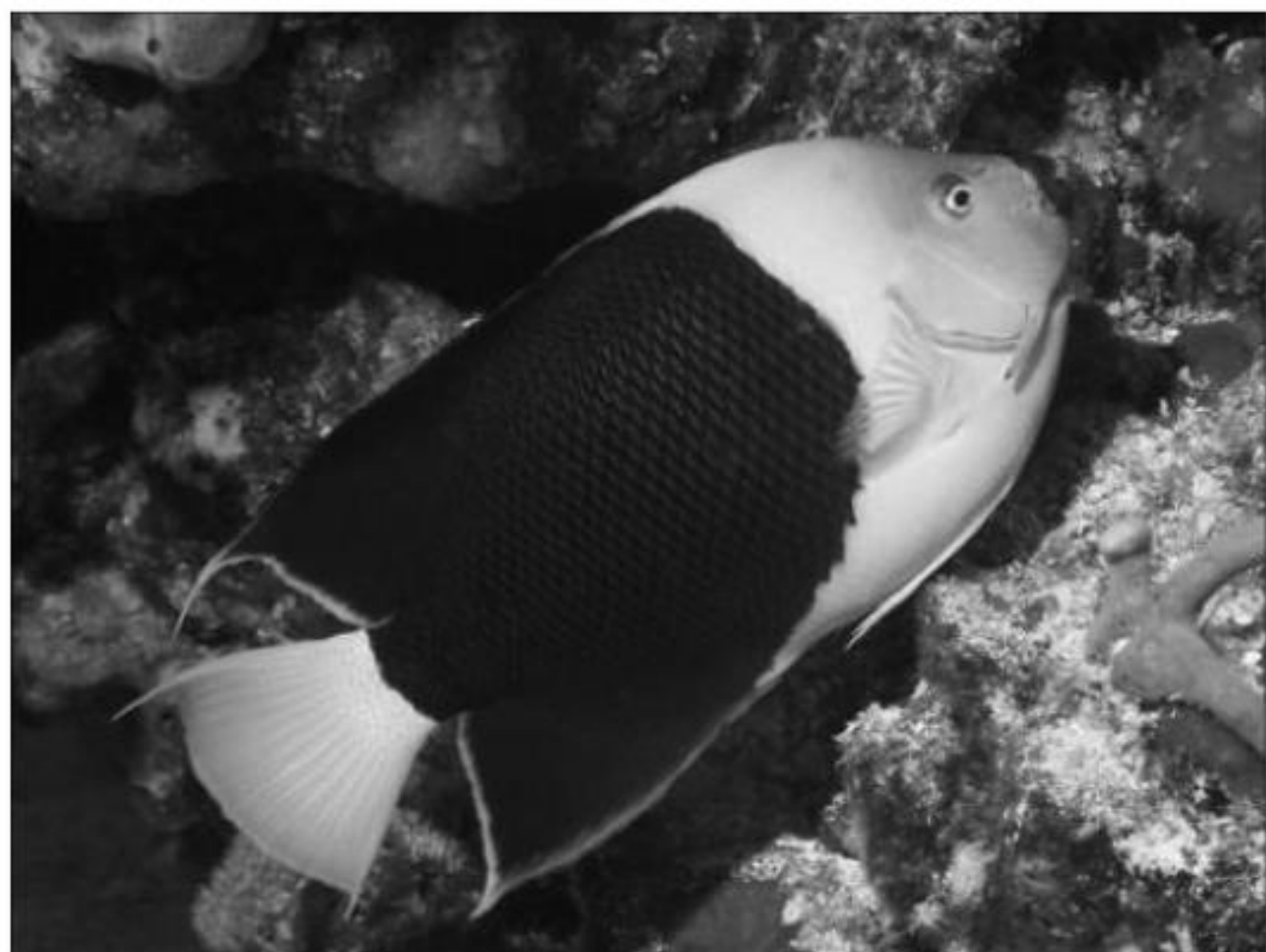
Output  $\mathbf{y}$



$$\underset{G}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\underset{\text{"What should I do?"}}{L}(\underset{\text{Neural Network}}{G}(\mathbf{x}), \mathbf{y})]$$

# Image-to-Image Translation

Input  $\mathbf{x}$



Output  $\mathbf{y}$



$$\underset{G}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\underset{\text{“What should I do?”}}{L}(\underset{\text{“How should I do it?”}}{G}(\mathbf{x}), \mathbf{y})]$$

# Be careful what you wish for!

Input



Output



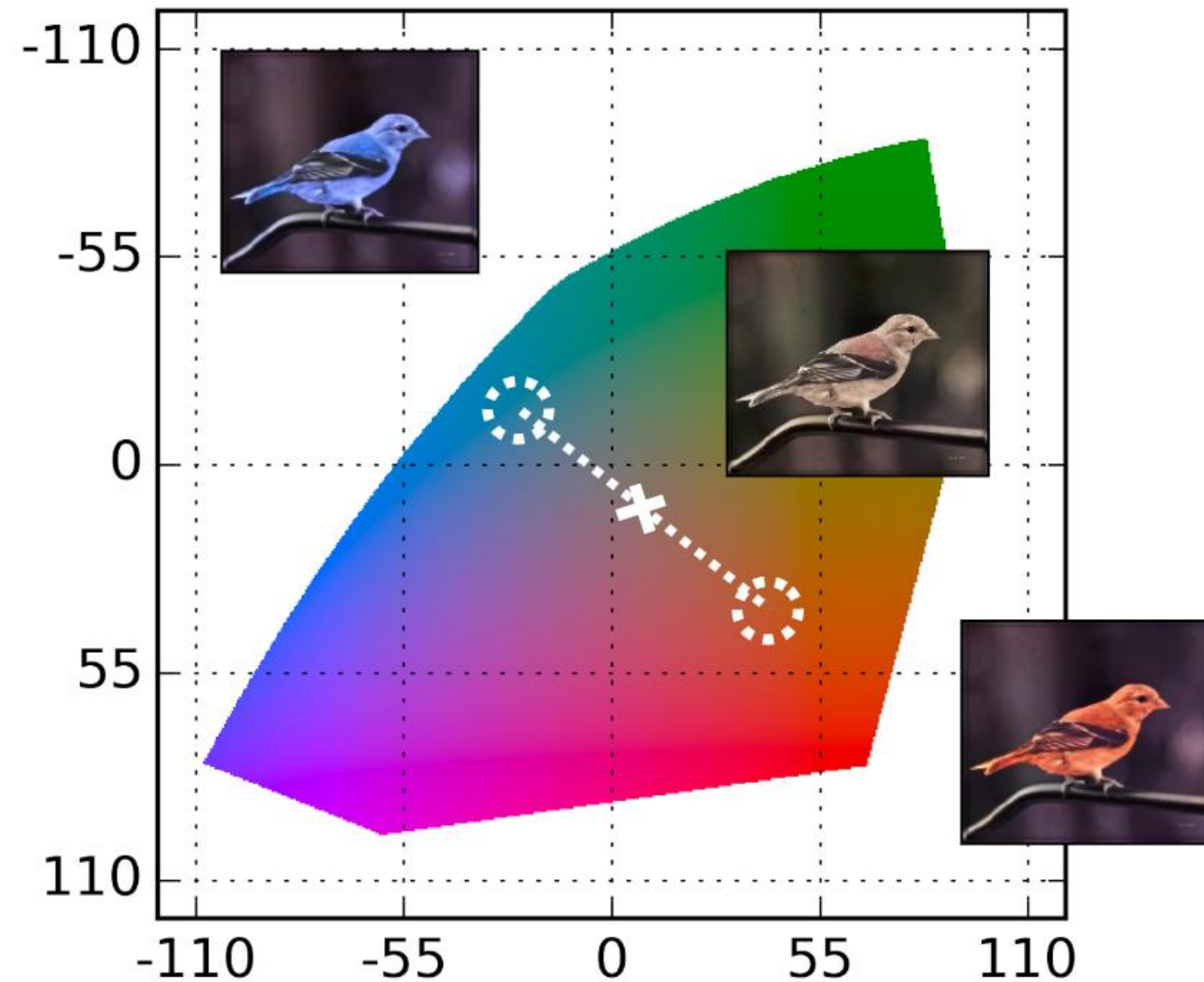
Ground truth



$$L(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$$



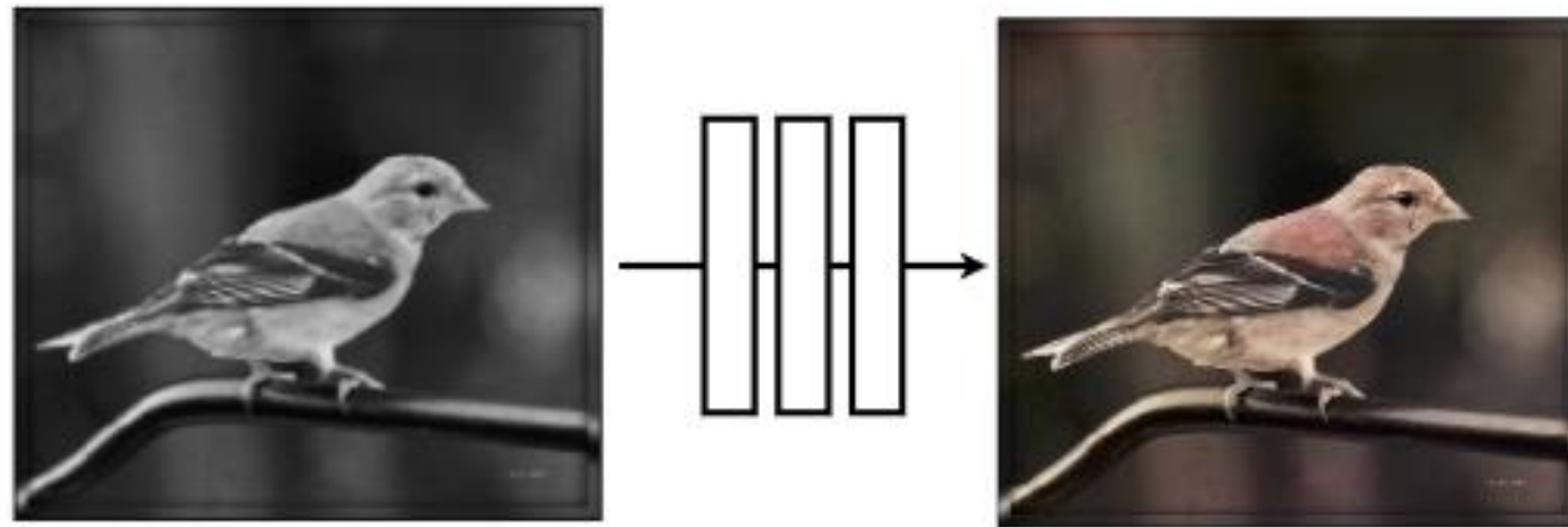
# Degradation to the mean!



$$L(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$$

# Automate Design of the Loss?

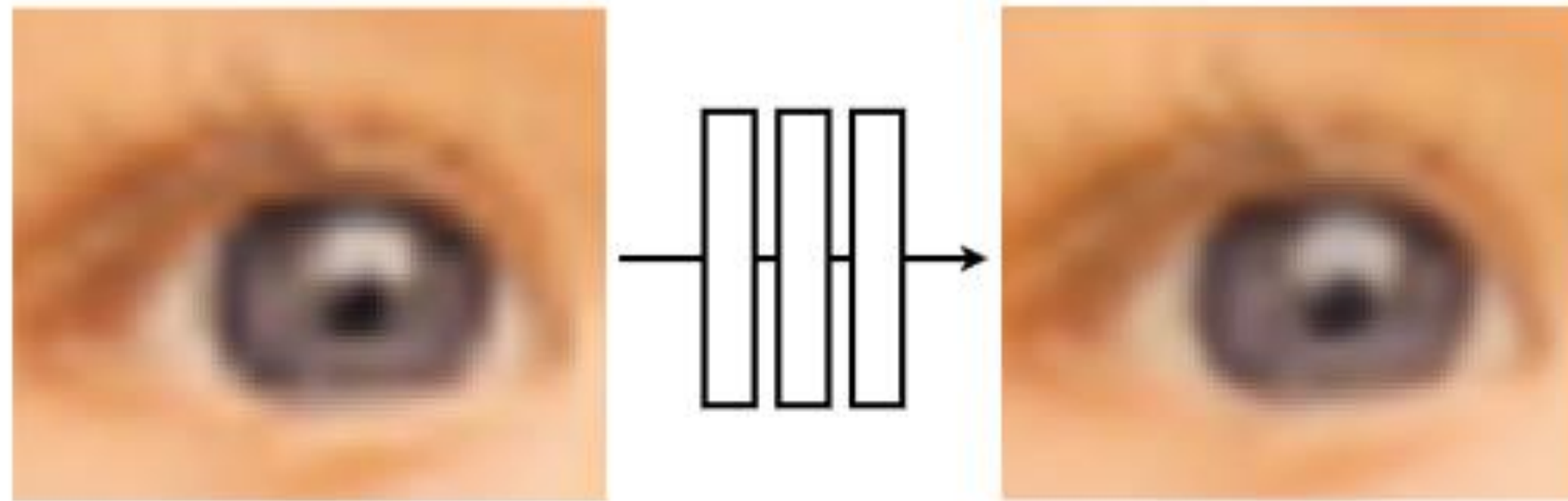
Image colorization



[Zhang, Isola, Efros, ECCV 2016]

L2 regression

Super-resolution

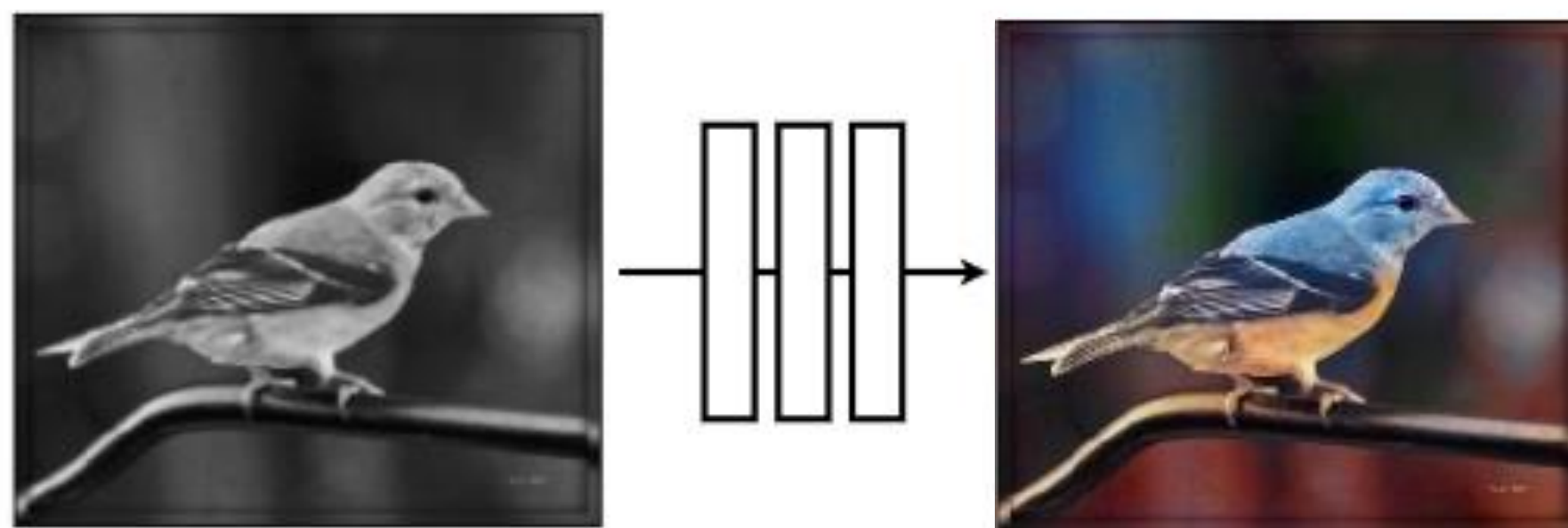


[Johnson, Alahi, Li, ECCV 2016]

L2 regression

# Automate Design of the Loss?

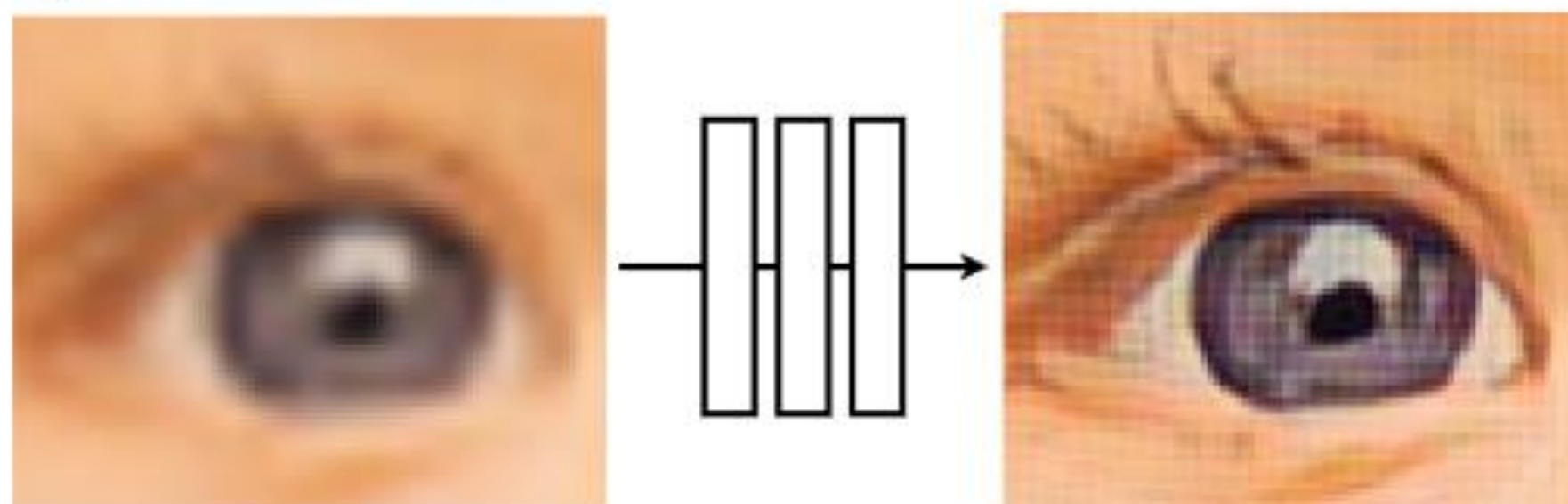
Image colorization



[Zhang, Isola, Efros, ECCV 2016]

Cross entropy objective,  
with colorfulness term

Super-resolution

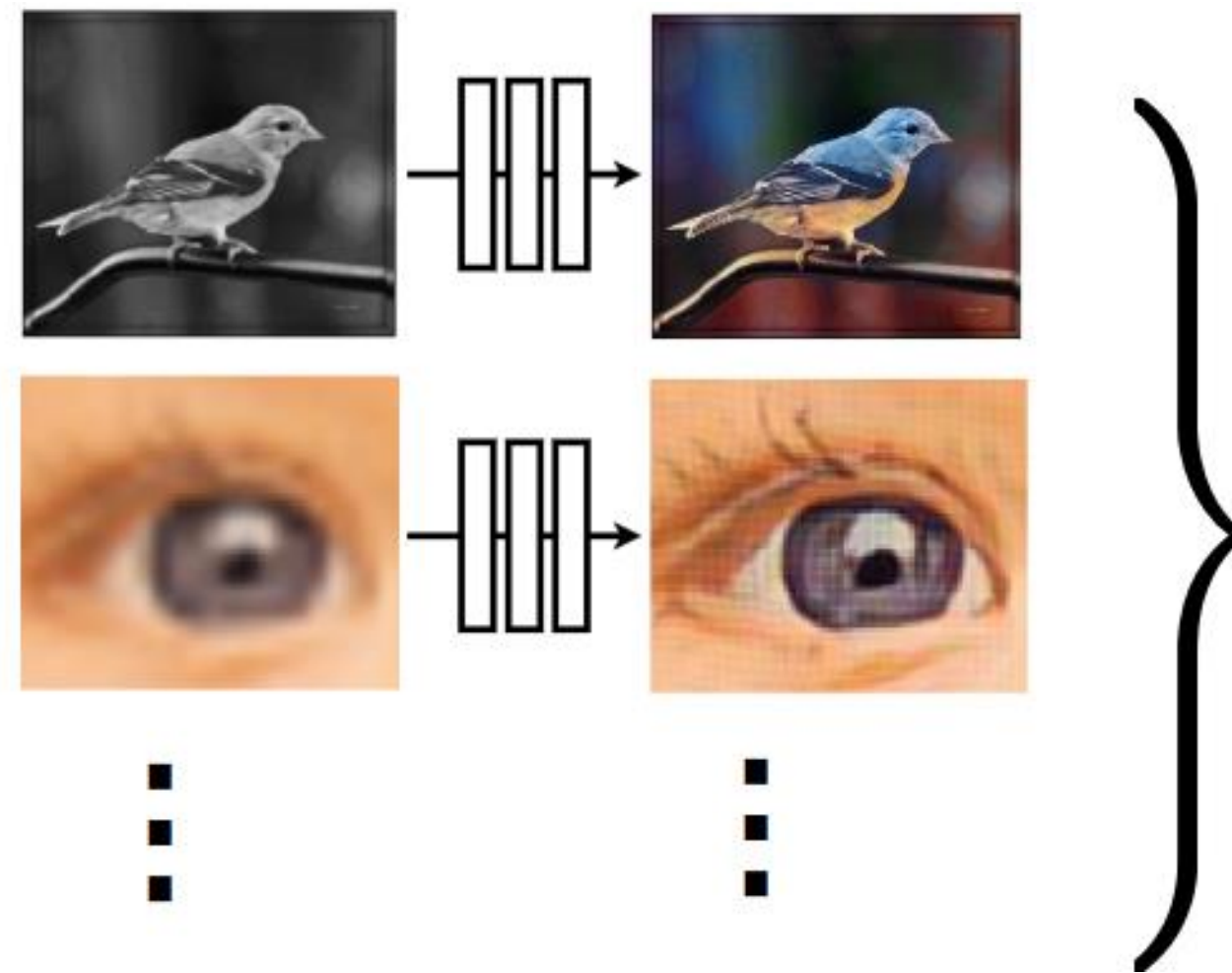


[Johnson, Alahi, Li, ECCV 2016]

Deep feature covariance  
matching objective

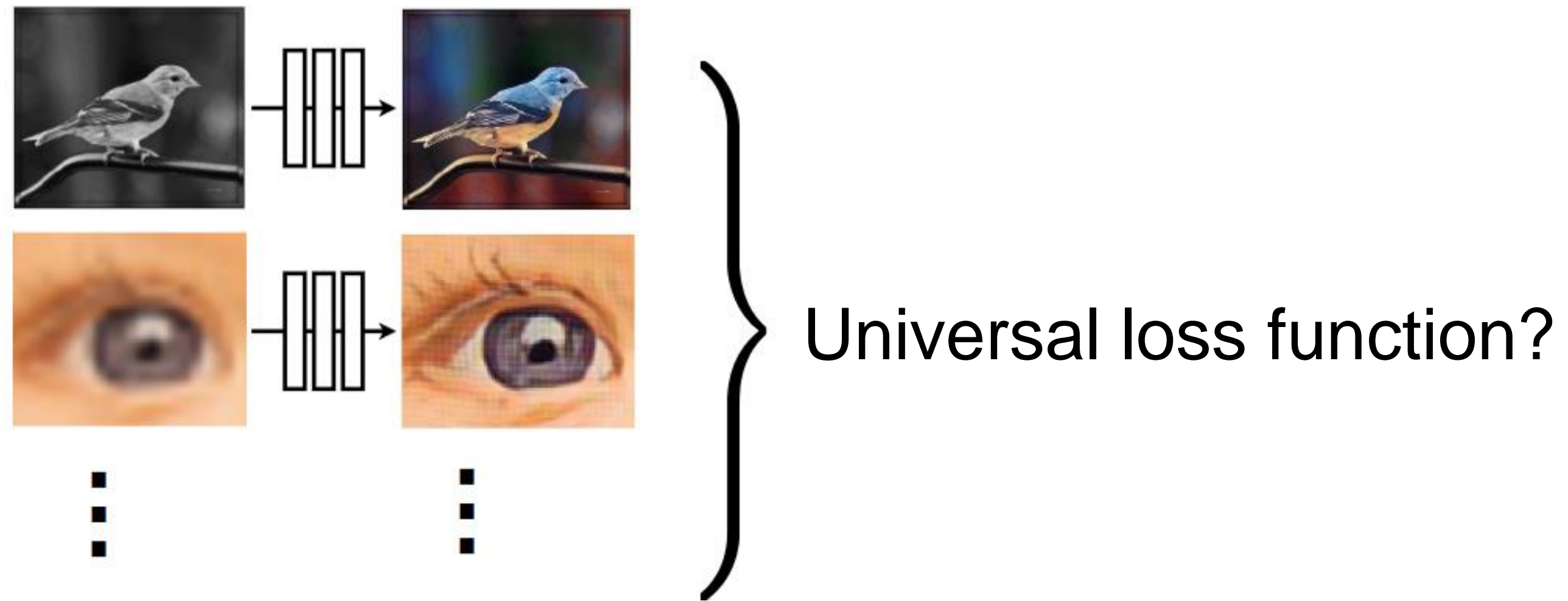
# Automate Design of the Loss?

Deep learning got rid of handcrafted features. Can we also get rid of handcrafting the loss function?



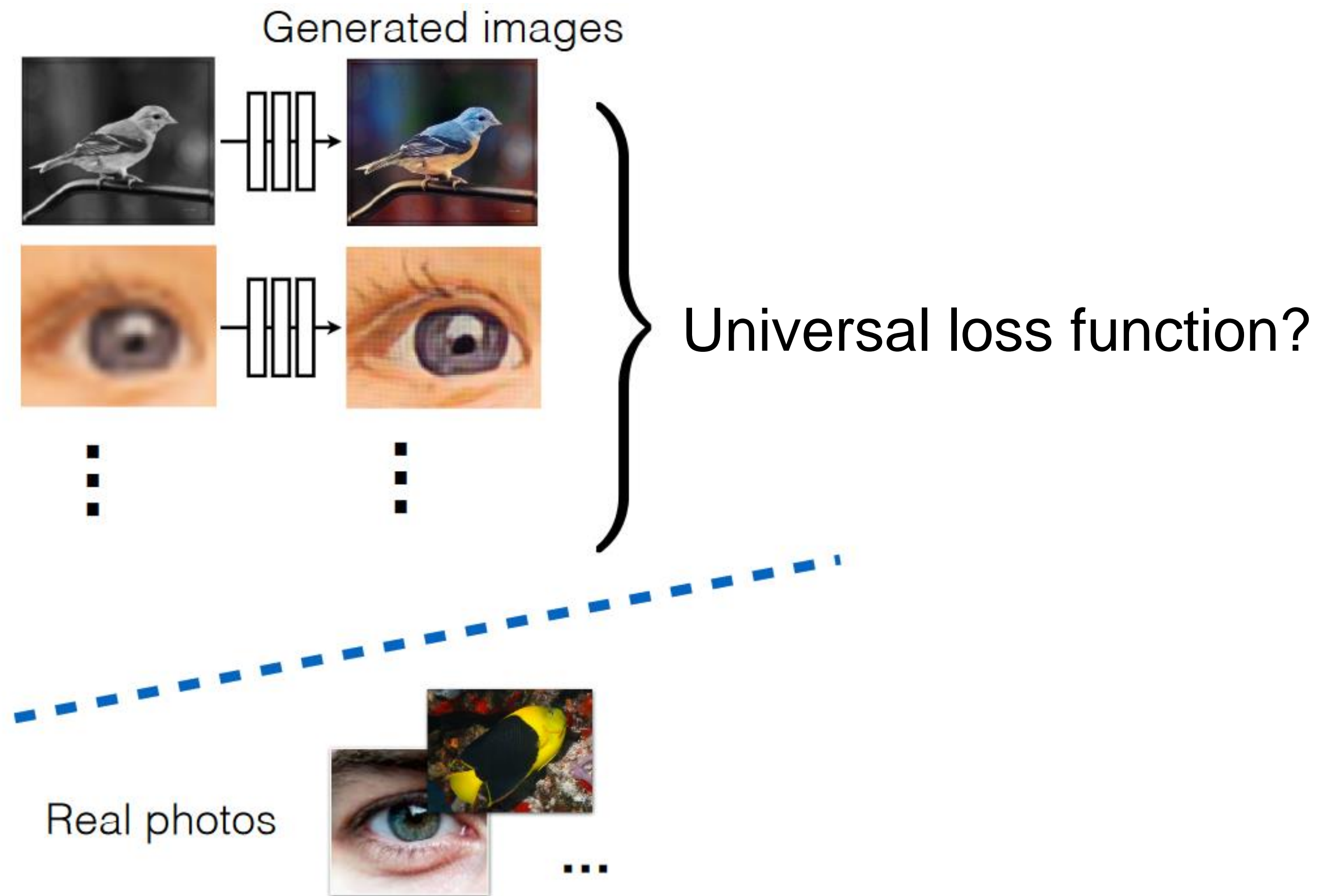
# Automate Design of the Loss?

Deep learning got rid of handcrafted features. Can we also get rid of handcrafting the loss function?

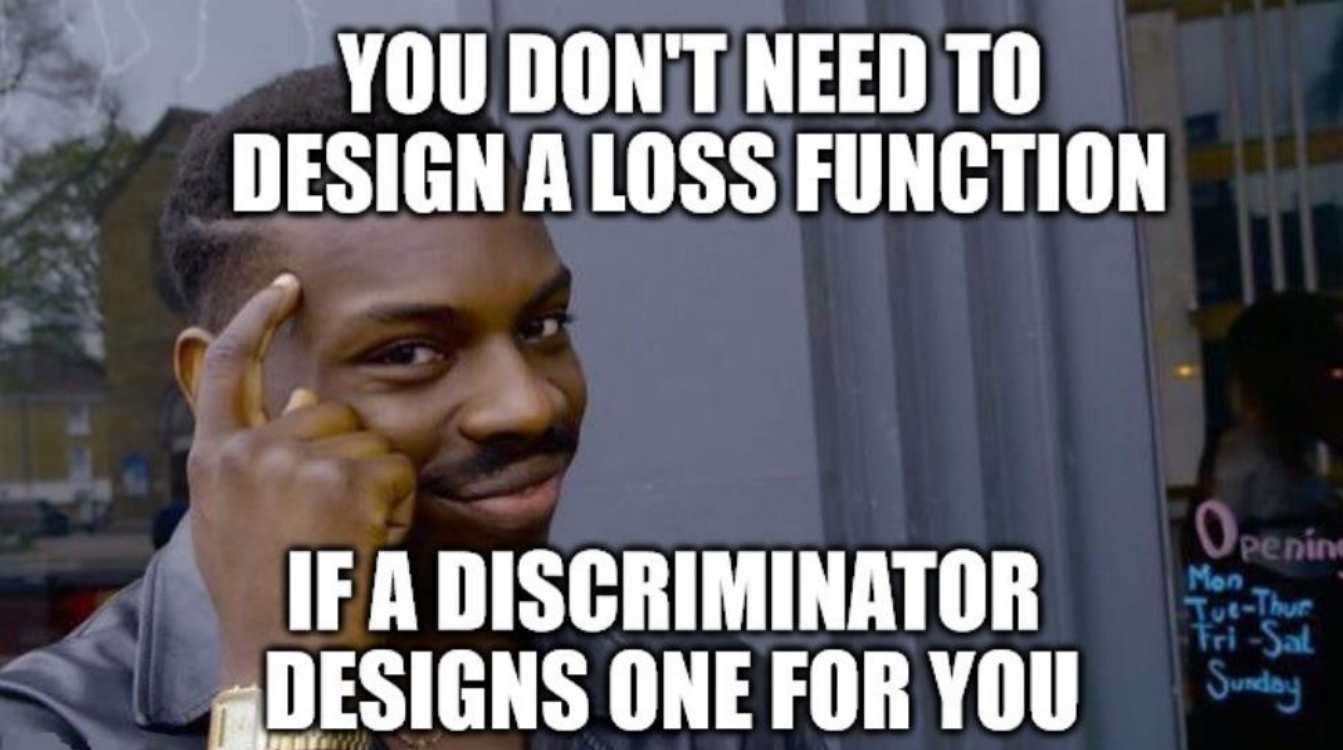
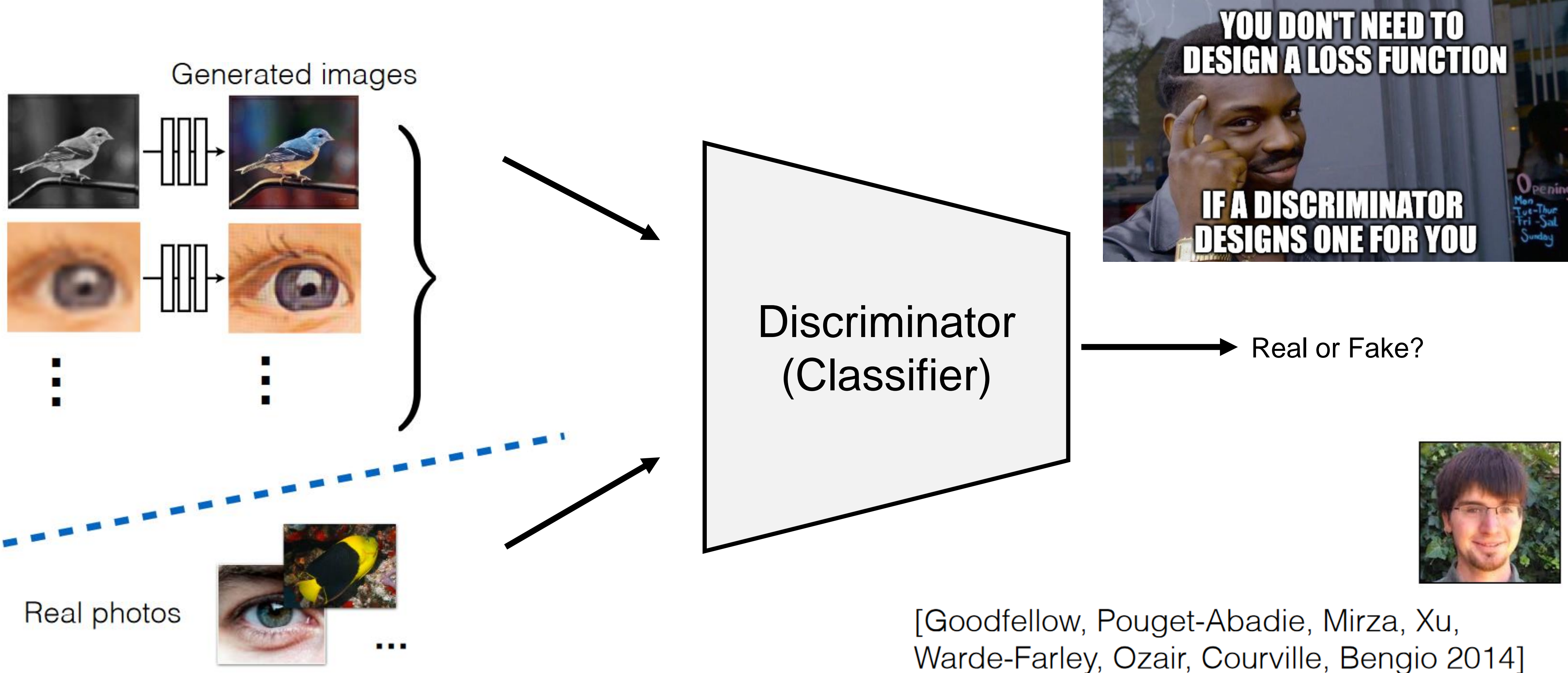


# Automate Design of the Loss?

Deep learning got rid of handcrafted features. Can we also get rid of handcrafting the loss function?



# Discriminator as a Loss Function



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

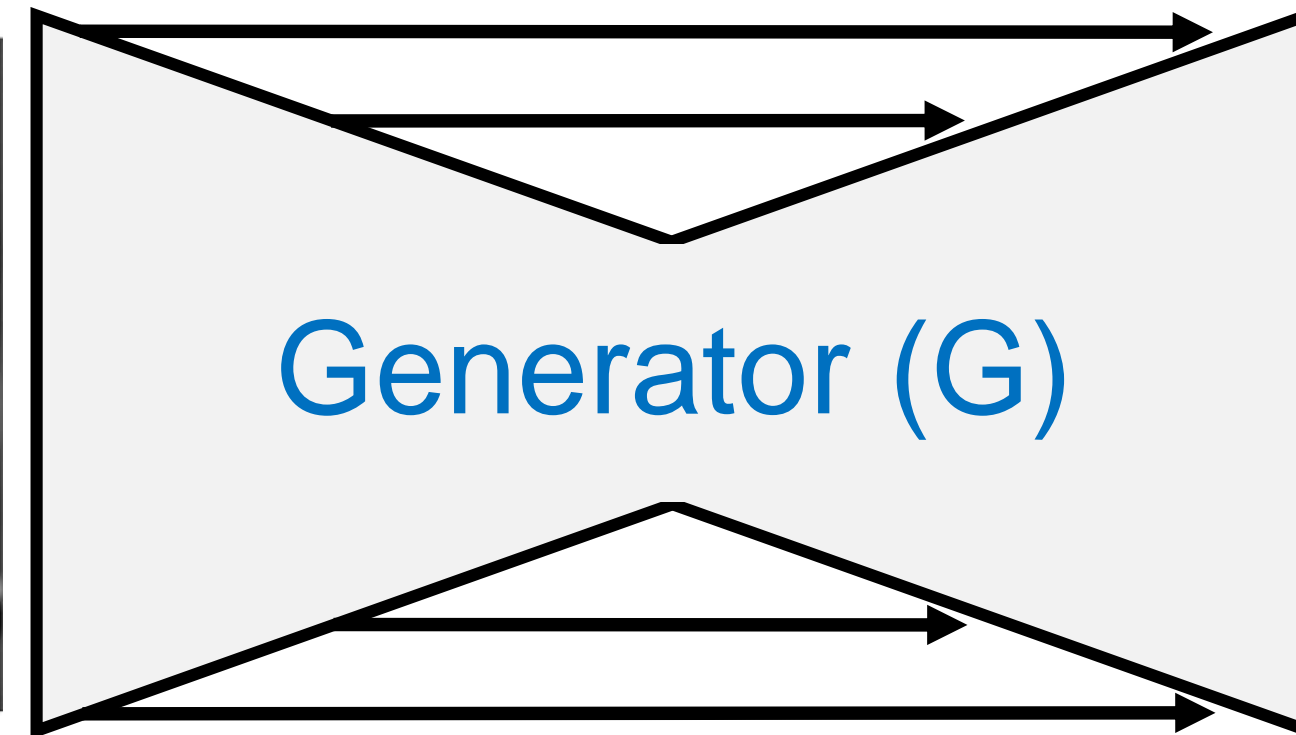
# Conditional GAN





# Conditional GAN

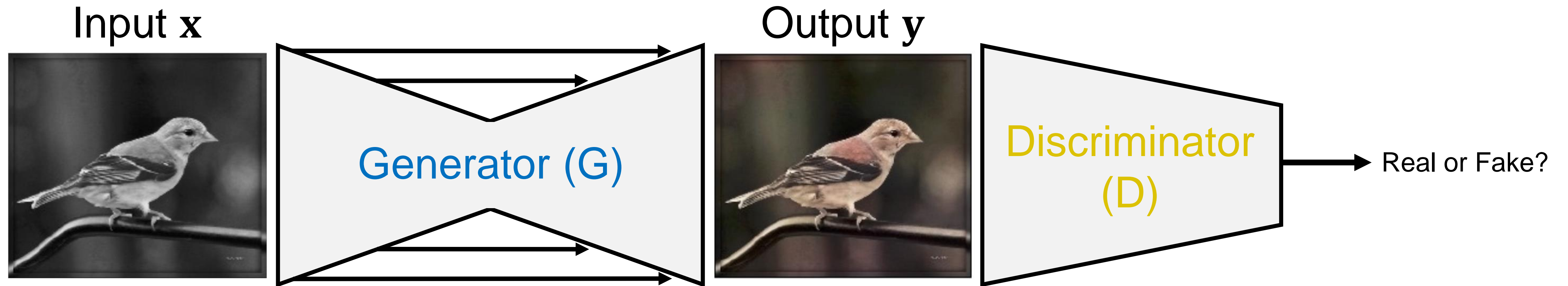
Input  $x$



Output  $y$



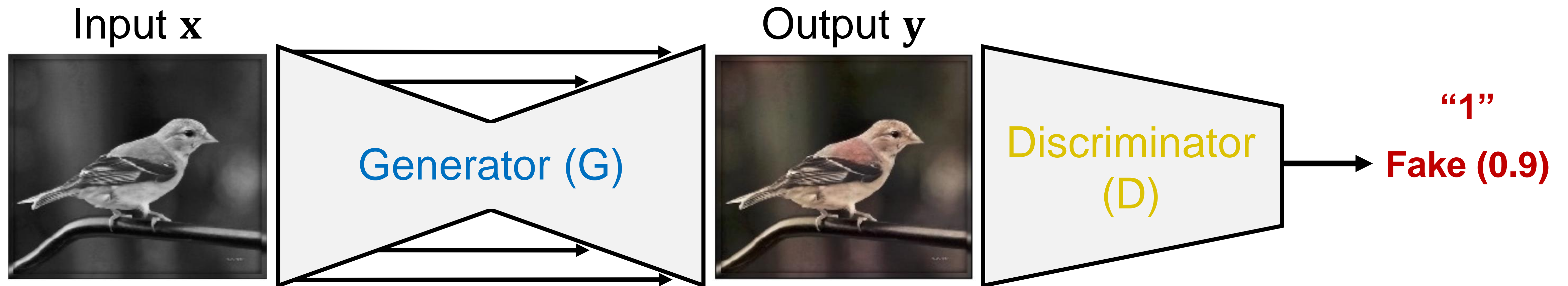
# Conditional GAN



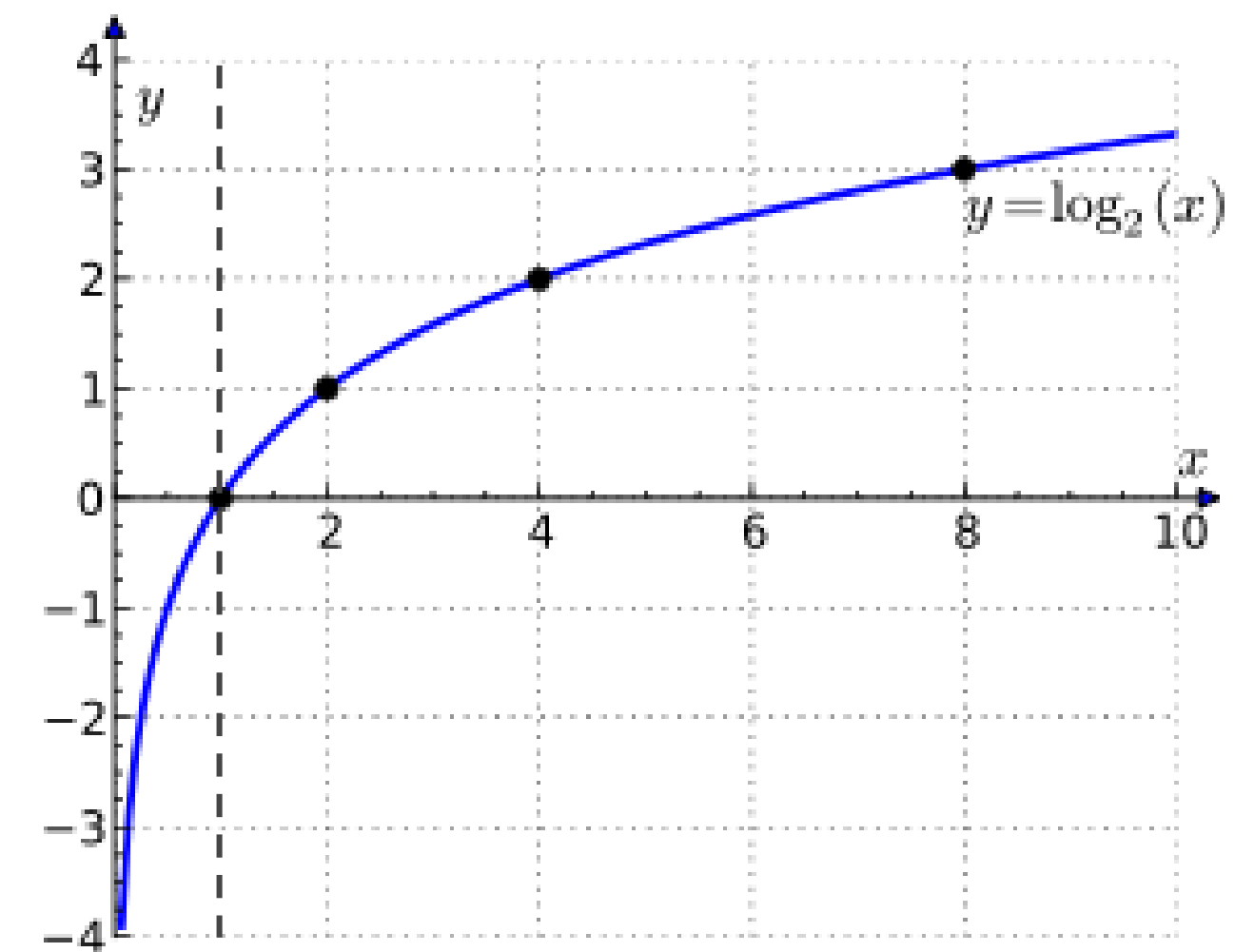
**G** tries to synthesize fake images that fool **D**

**D** tries to tell real from fake

# Conditional GAN (Discriminator)

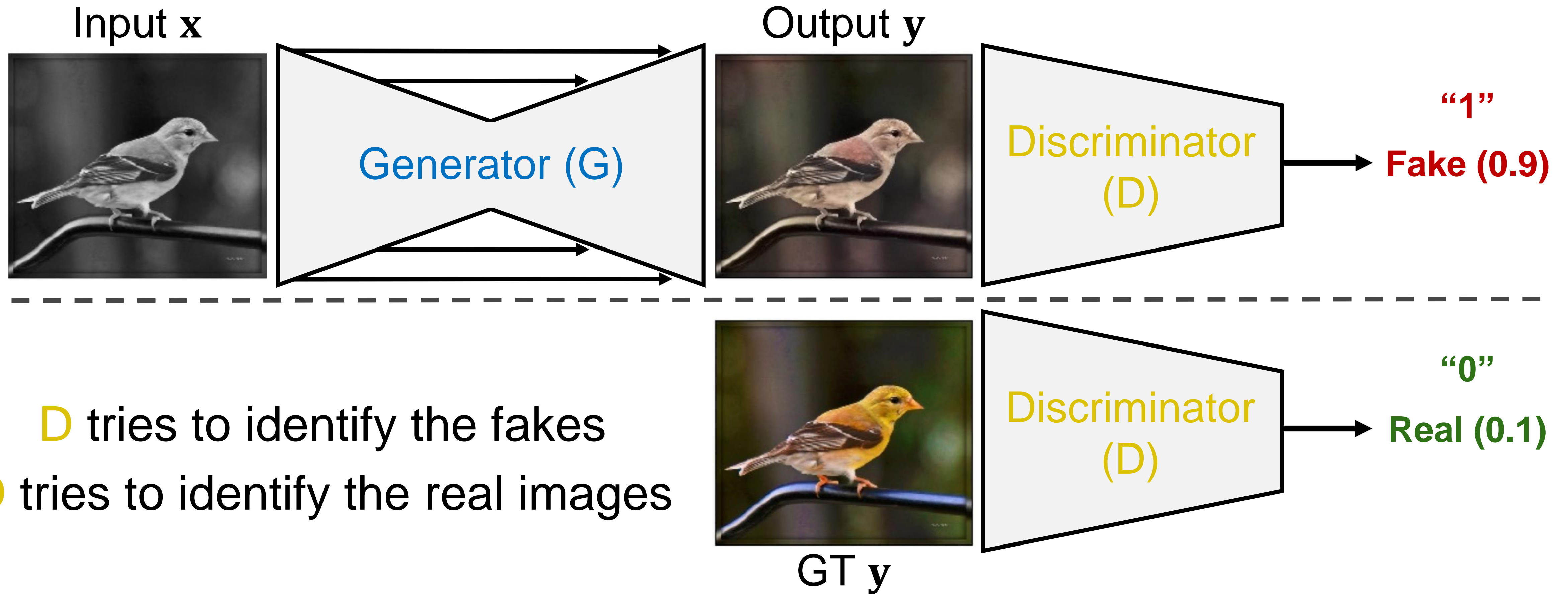


**D** tries to identify the fakes



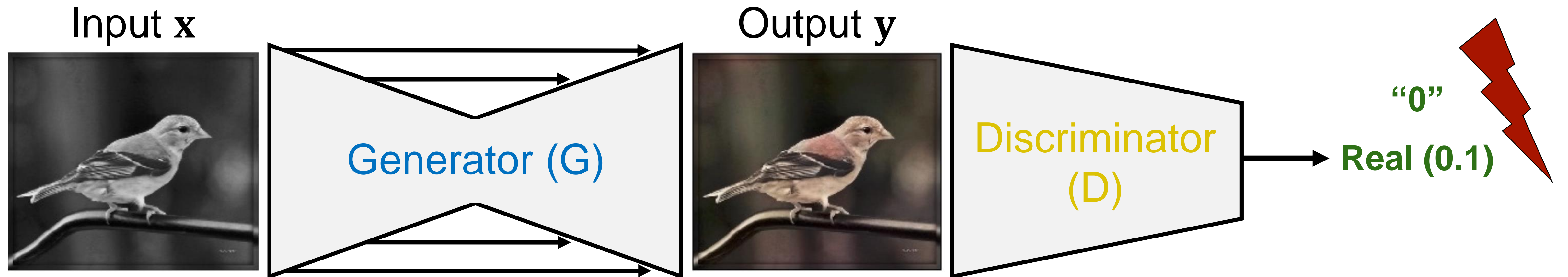
$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \log D(G(\mathbf{x})) + \right]$$

# Conditional GAN (Discriminator)

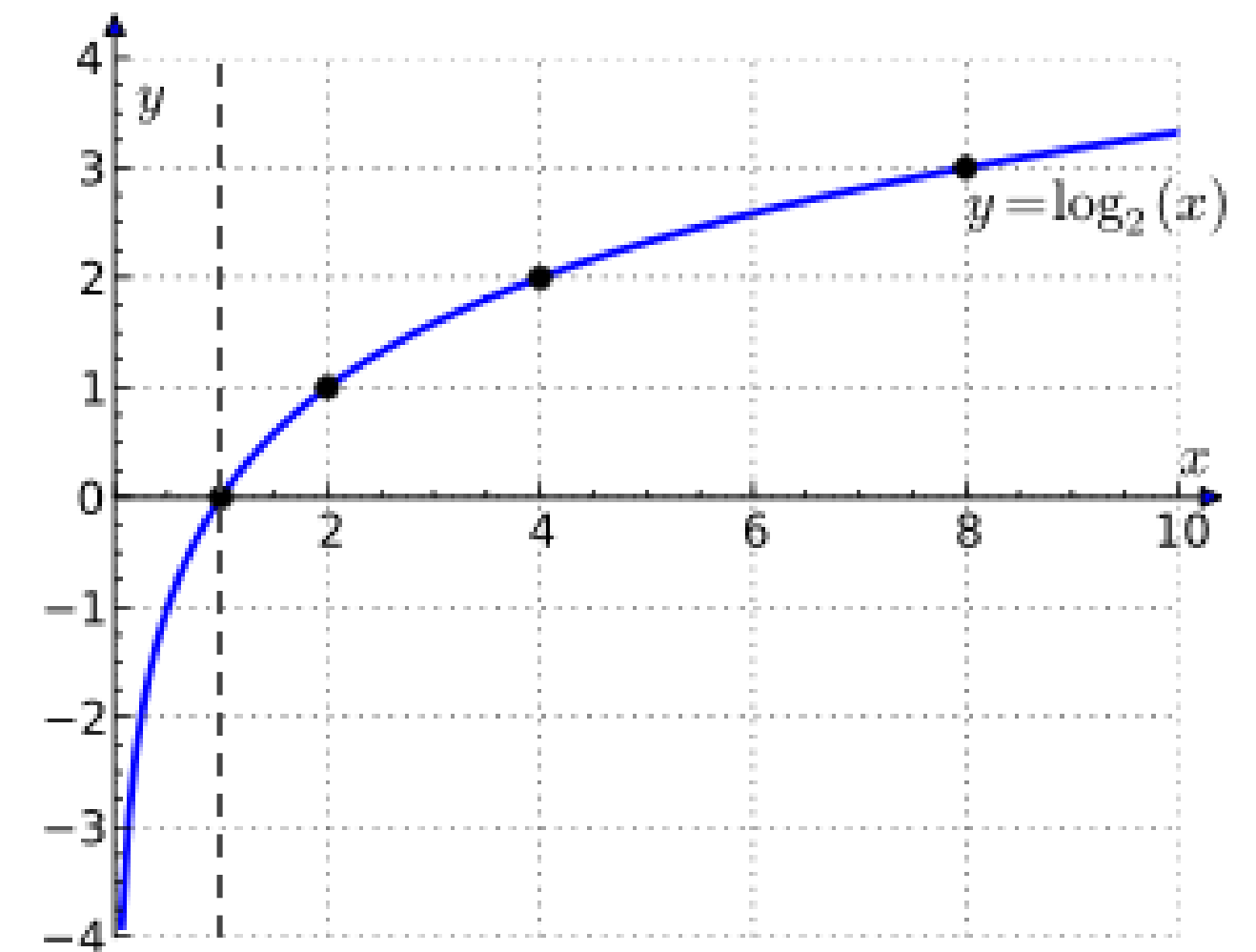


$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \log D(\mathbf{G}(\mathbf{x})) + \log(1 - D(\mathbf{y})) \right]$$

# Conditional GAN (Generator)

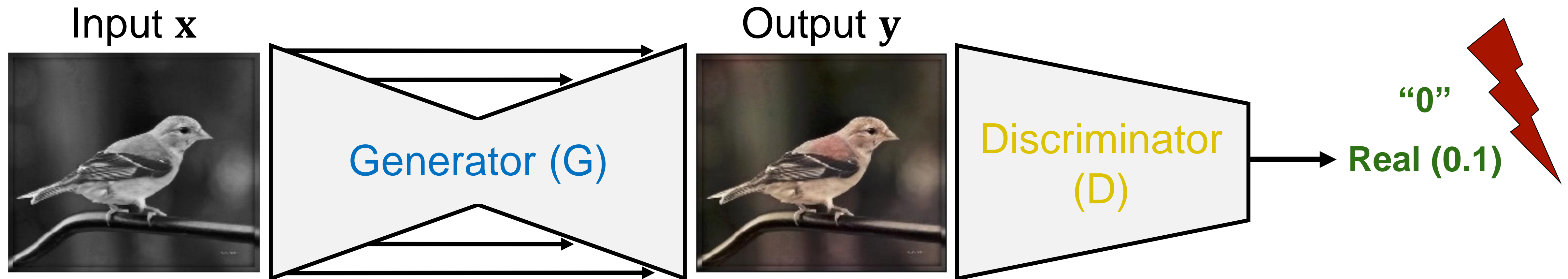


$G$  tries to synthesize fake images that fool  $D$ .

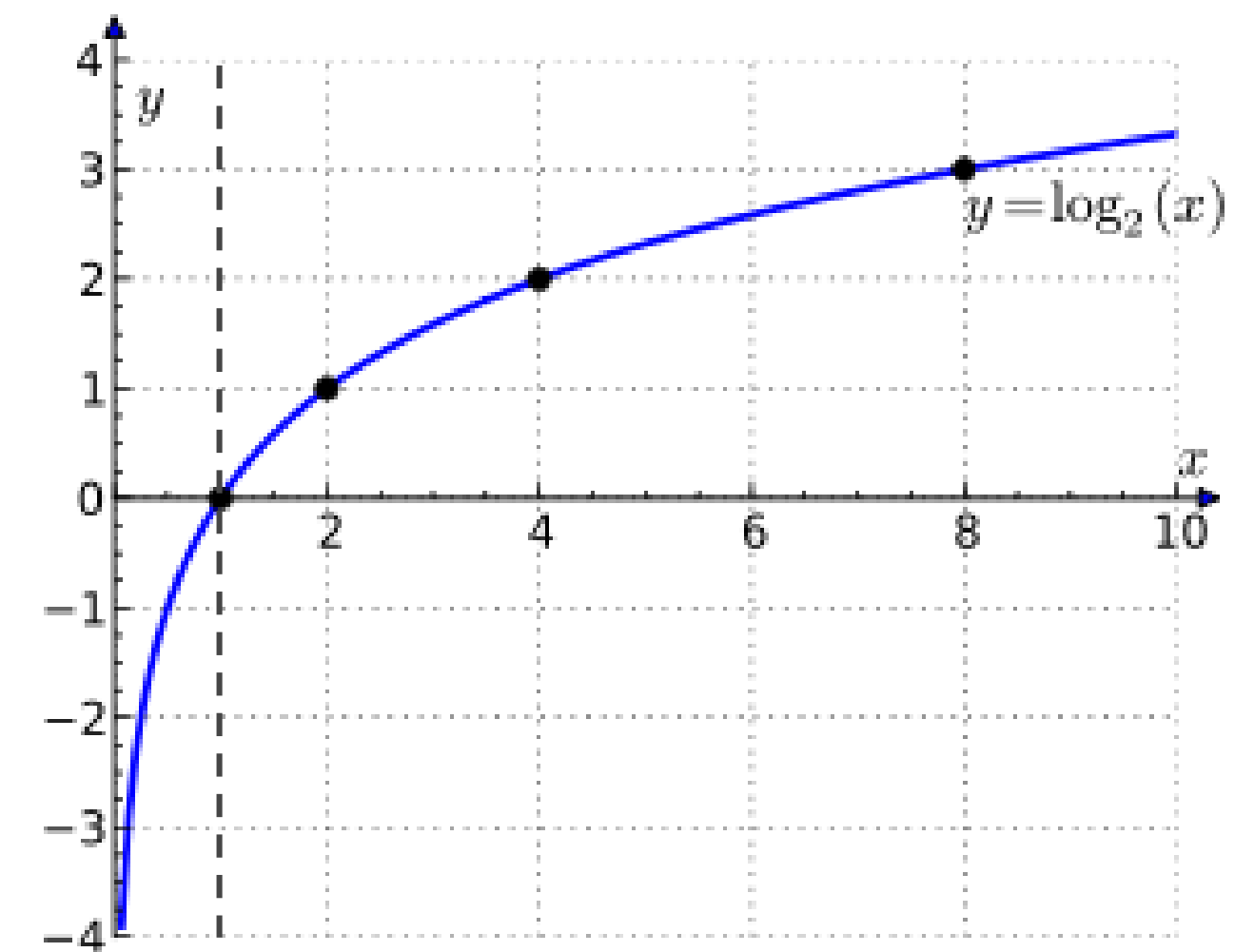


$$\arg \min_G \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

# Conditional GAN (Generator)

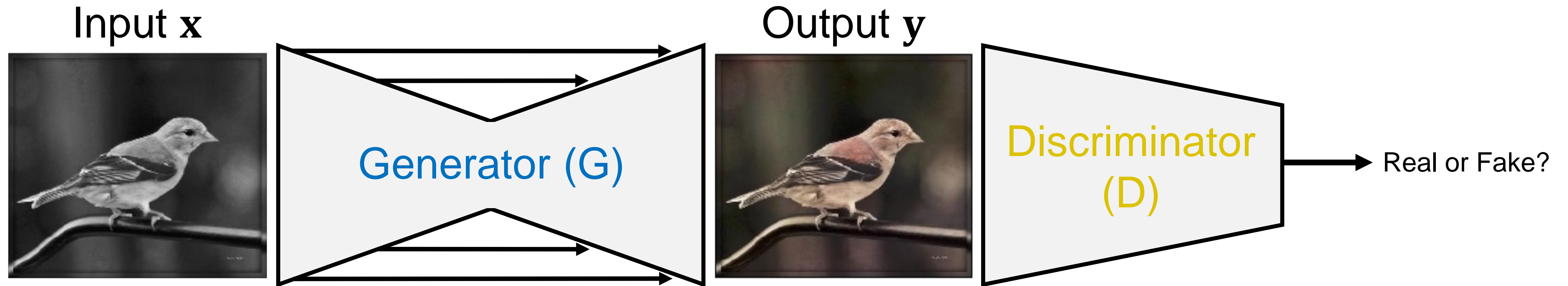


$G$  tries to synthesize fake images that fool  $D$ .



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

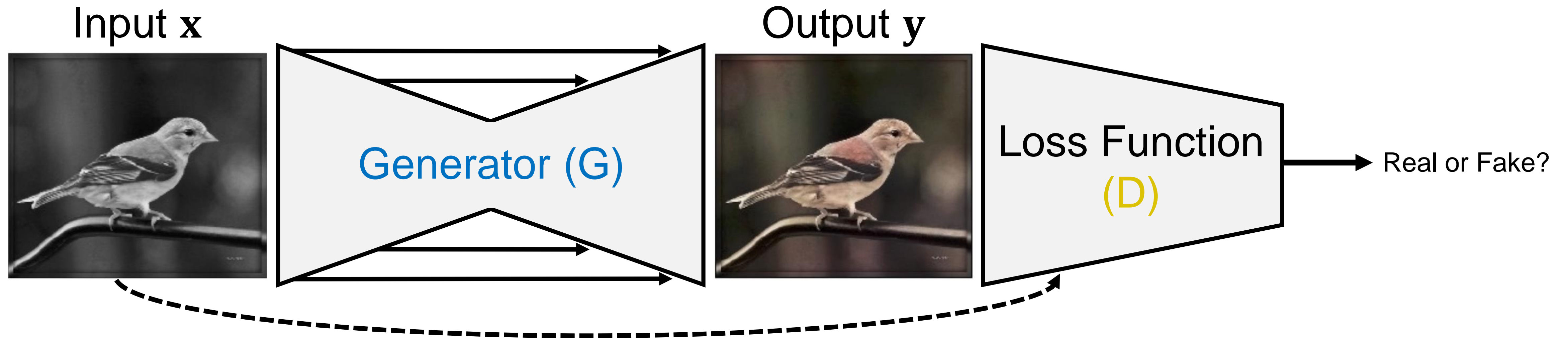
# Conditional GAN



$G$  tries to synthesize fake images that fool the best  $D$ .

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

# Conditional GAN

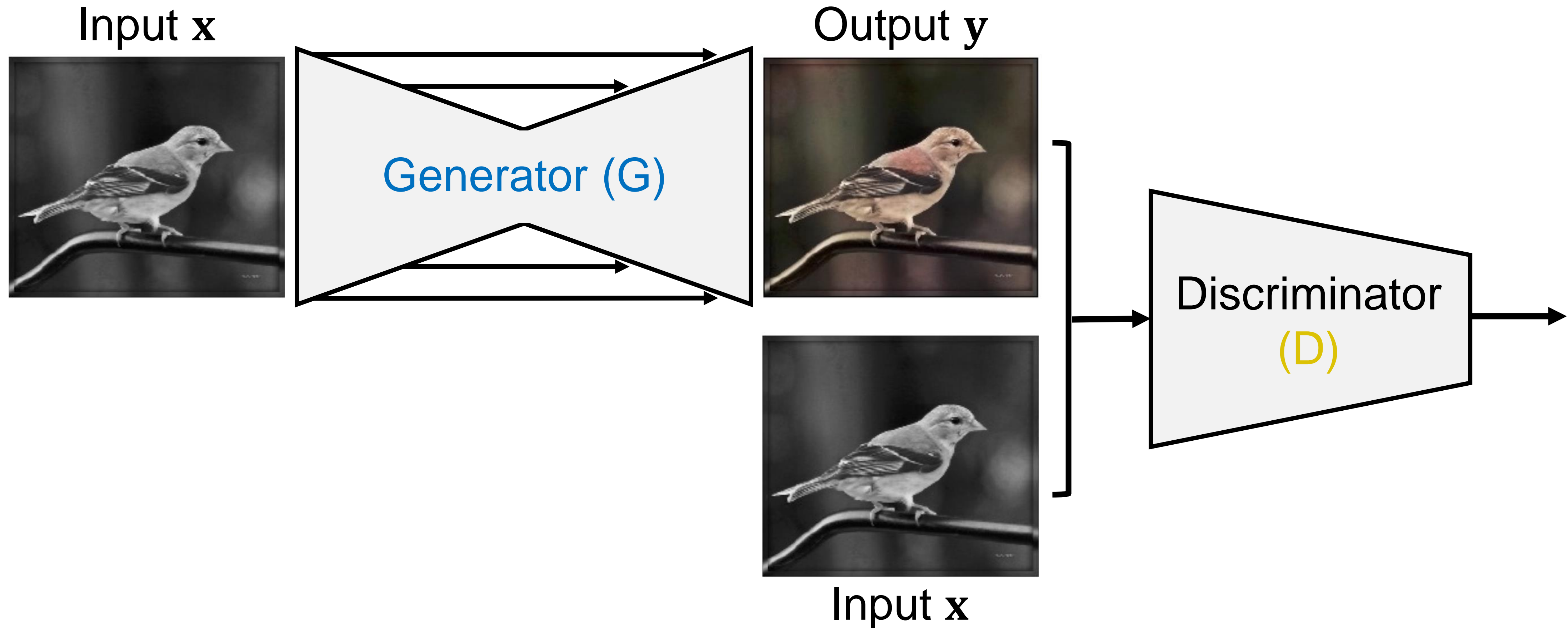


**G's** perspective: **D** is a loss function

Rather than being hand-designed, it is *learned jointly!*

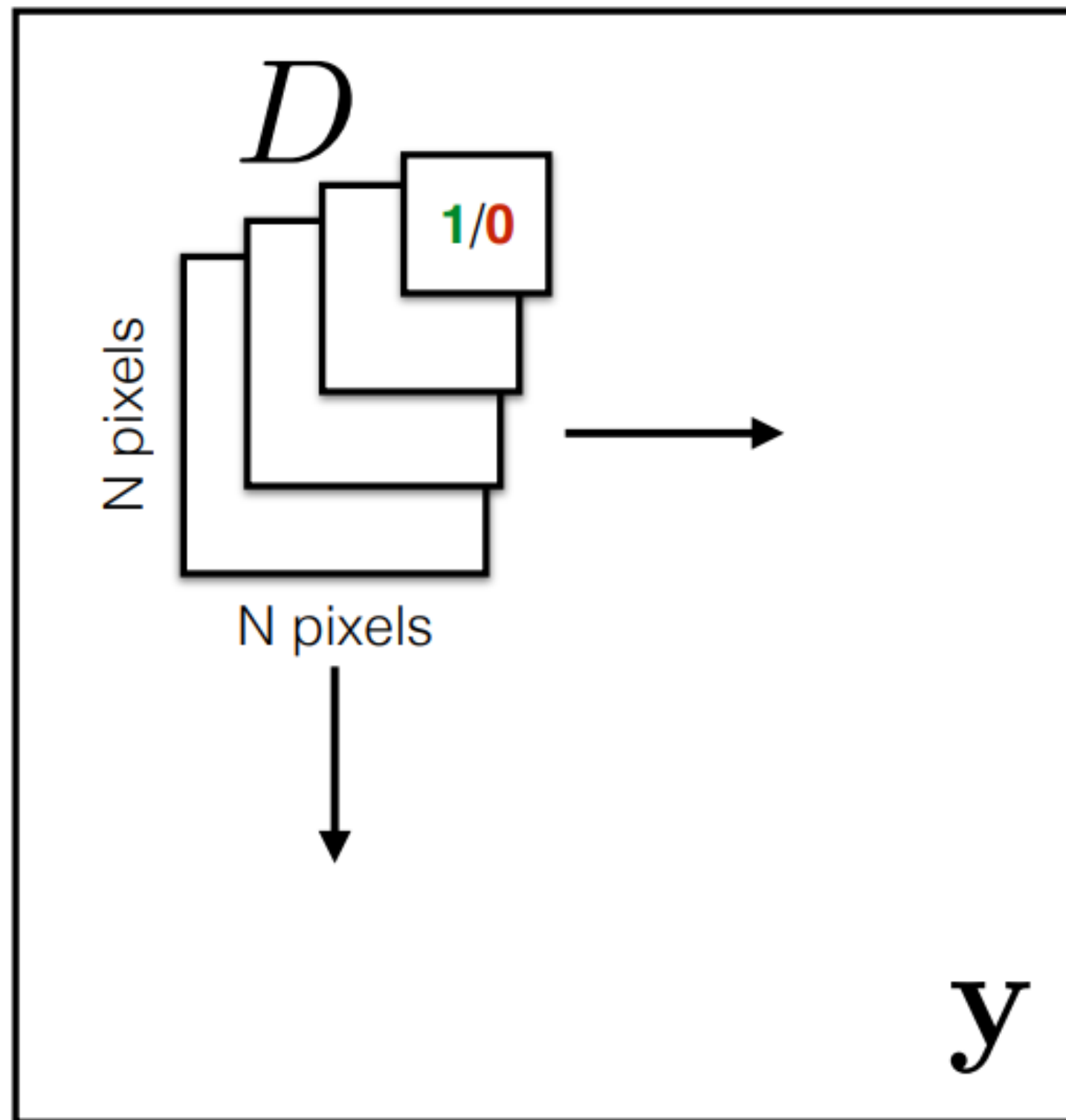


# Conditional Discriminator



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

# Patch Discriminator

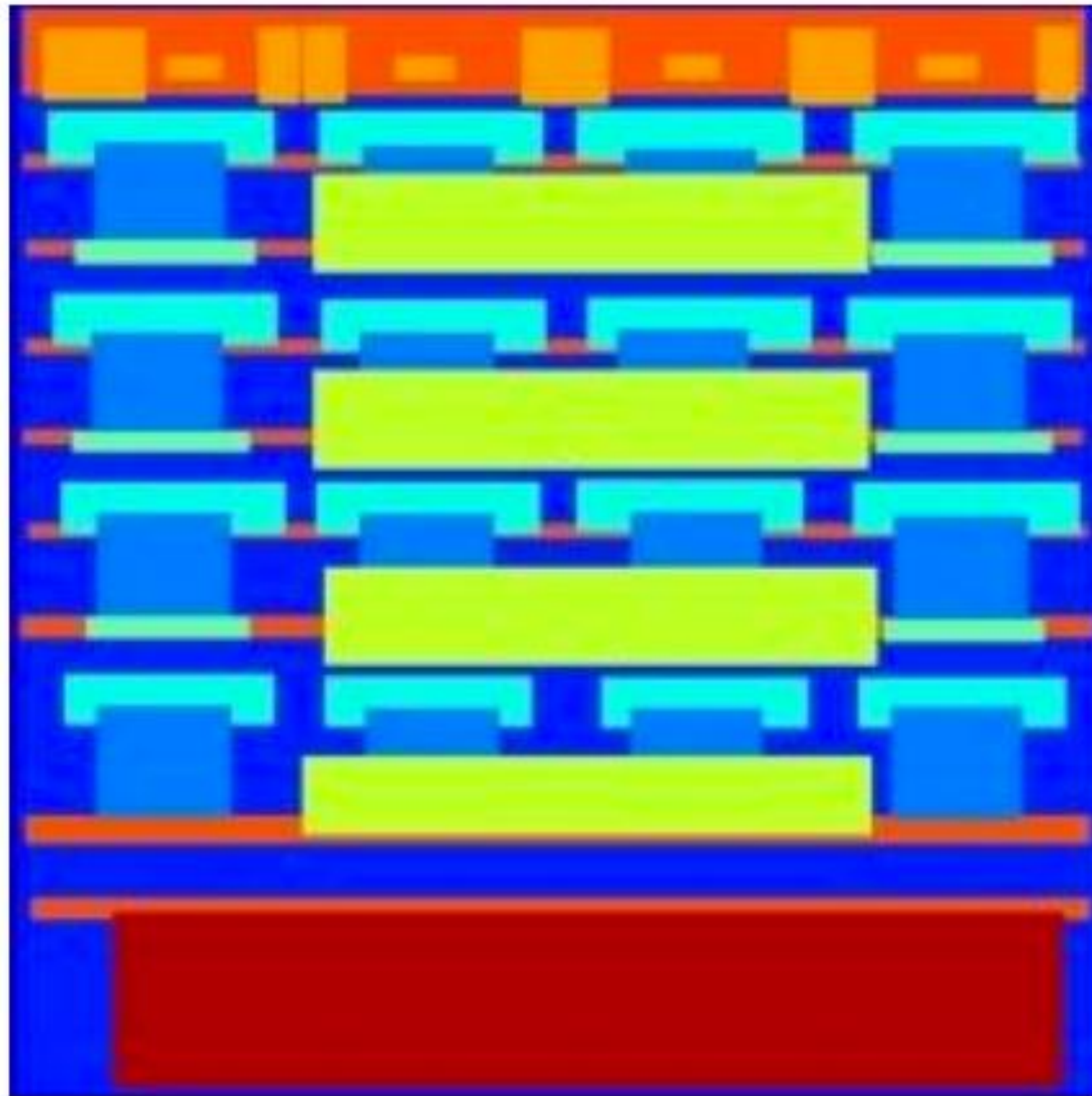


“Rather than penalizing if the output image looks fake, penalize if each overlapping patch in the output looks fake”

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

# 1x1 Pixel Discriminator

Input

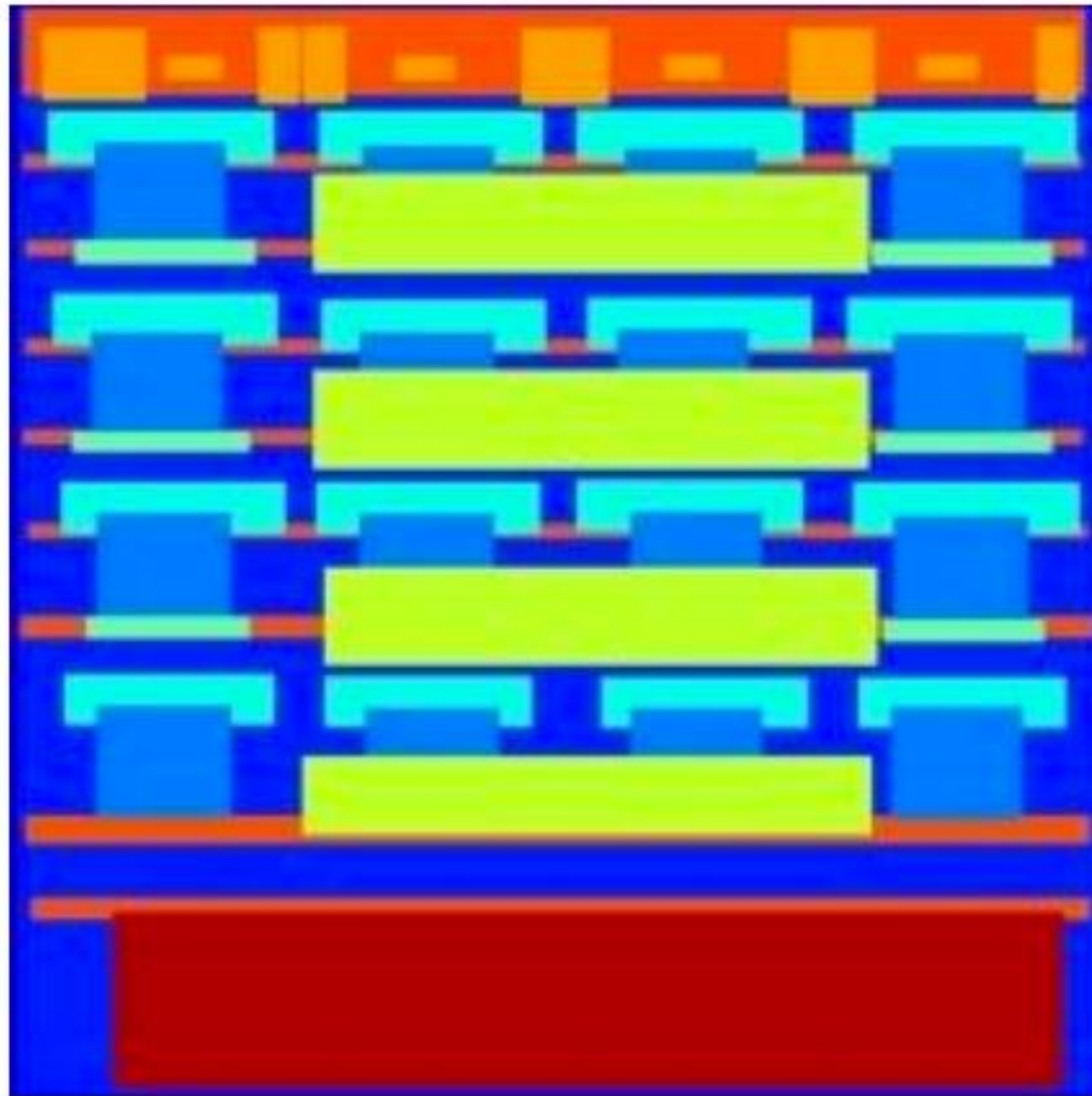


1x1 Discriminator



# Image Discriminator

Input

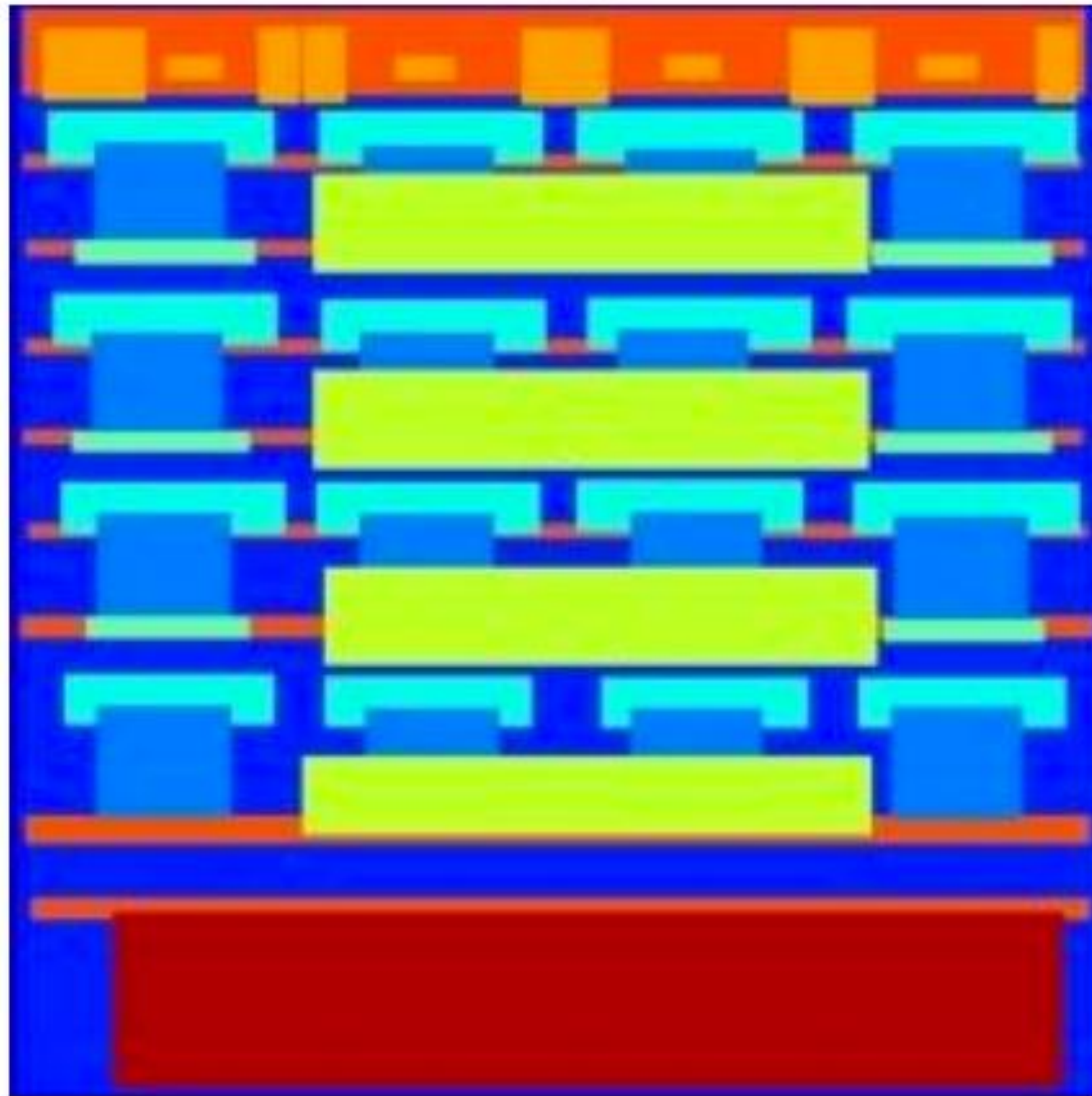


1x1 Discriminator



# 70x70 Patch Discriminator

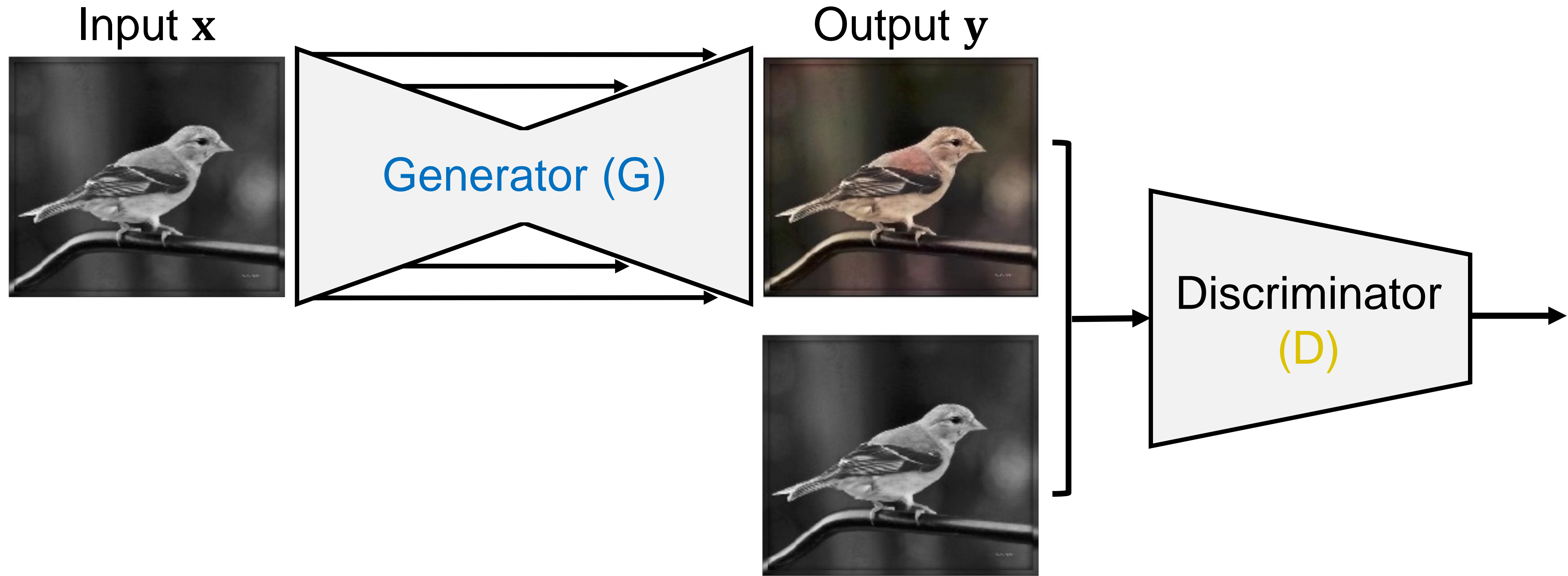
Input



1x1 Discriminator

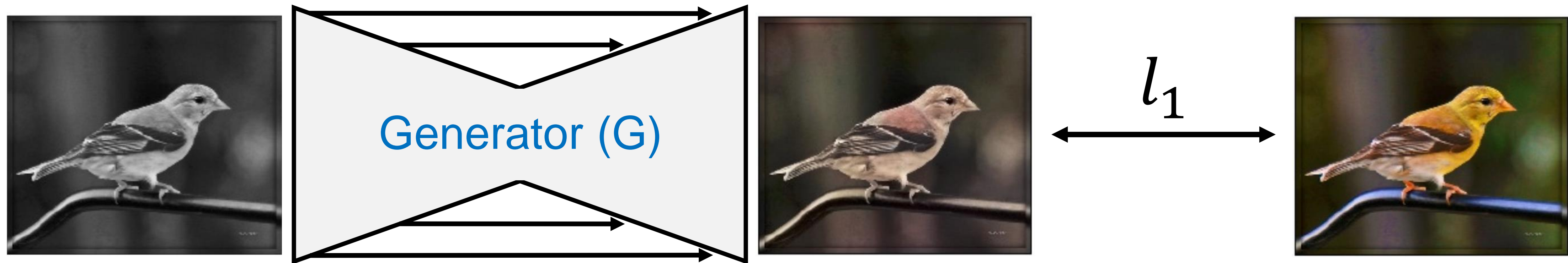


# Conditional Discriminator



$$L_{cGAN}(\mathbf{G}, \mathbf{D}) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log \mathbf{D}(\mathbf{x}, \mathbf{G}(\mathbf{x})) + \log(1 - \mathbf{D}(\mathbf{x}, \mathbf{y})) ]$$

# Reconstruction Loss



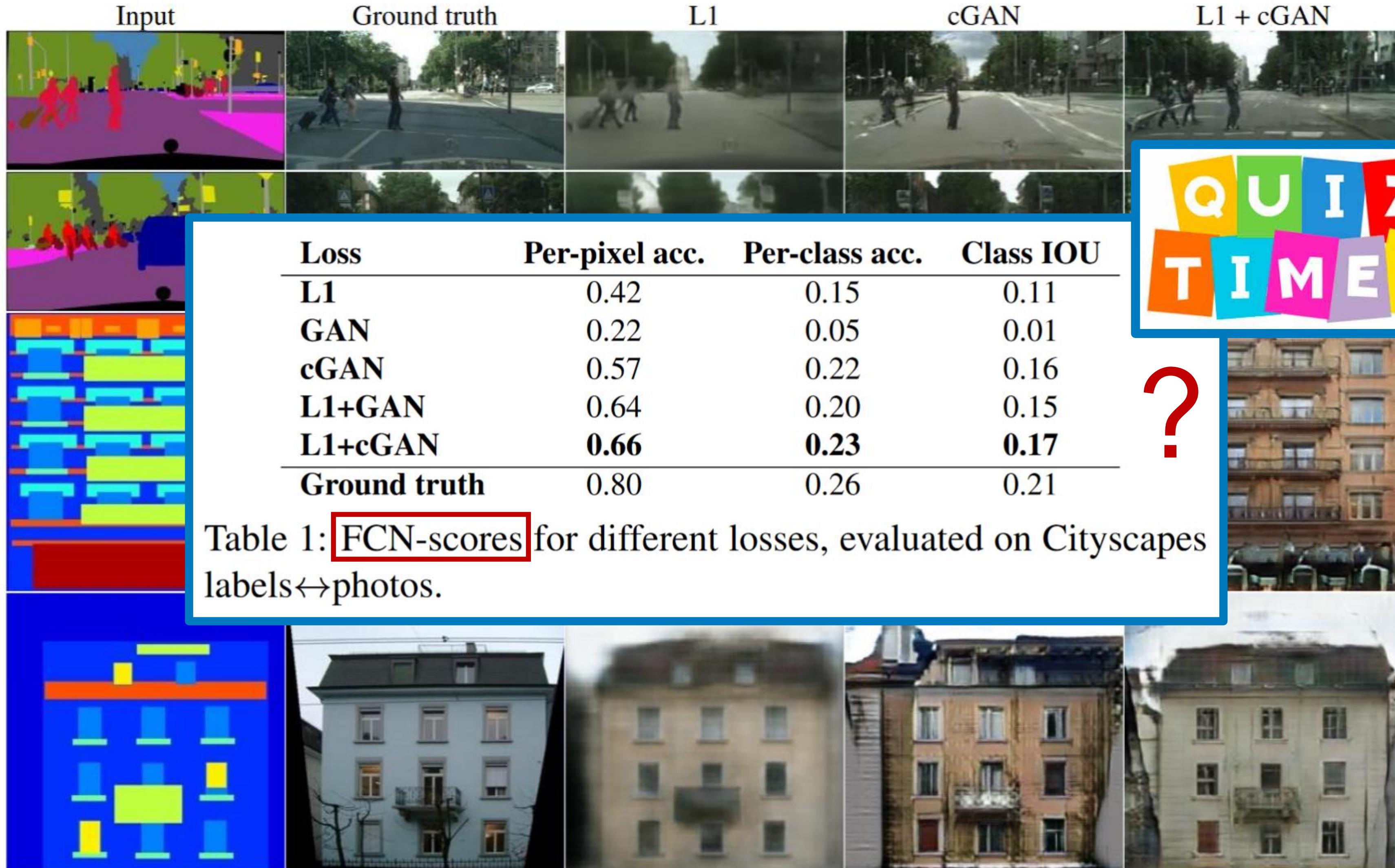
$$L_{l_1}(\mathbf{G}) = \mathbb{E}_{\mathbf{x}, \mathbf{y}} \|\mathbf{G}(\mathbf{x}) - \mathbf{y}\|_1$$

“Stable training + fast convergence”

$$G^* = \arg \min_{\mathbf{G}} \max_{\mathbf{D}} L_{cGAN}(\mathbf{G}, \mathbf{D}) + \lambda L_{l_1}(\mathbf{G})$$

↑  
100

# Ablation Study





# Ablation Study



# Results on the Test Split



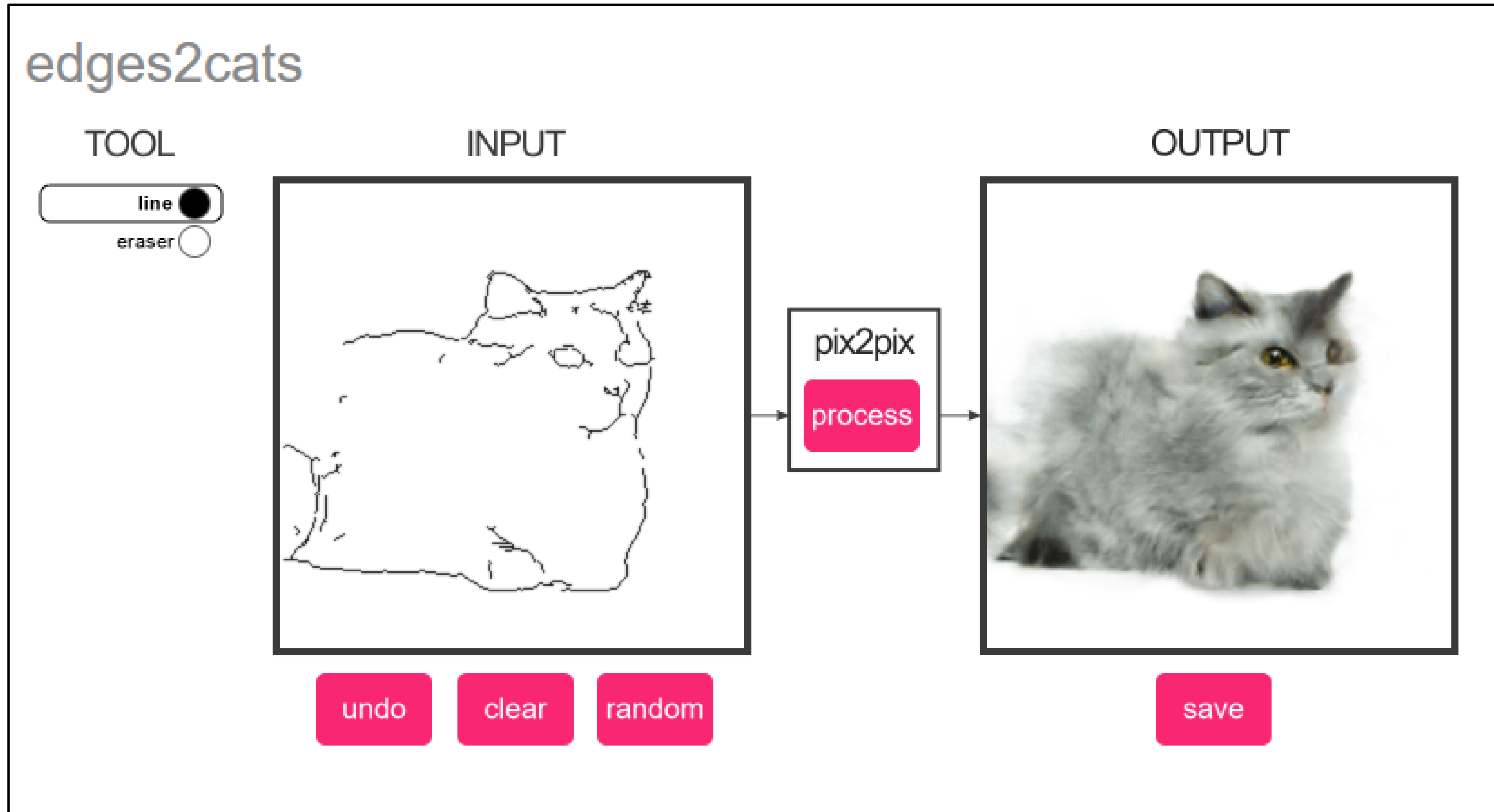
# Results for Hand Drawings



QUIZ  
TIME!



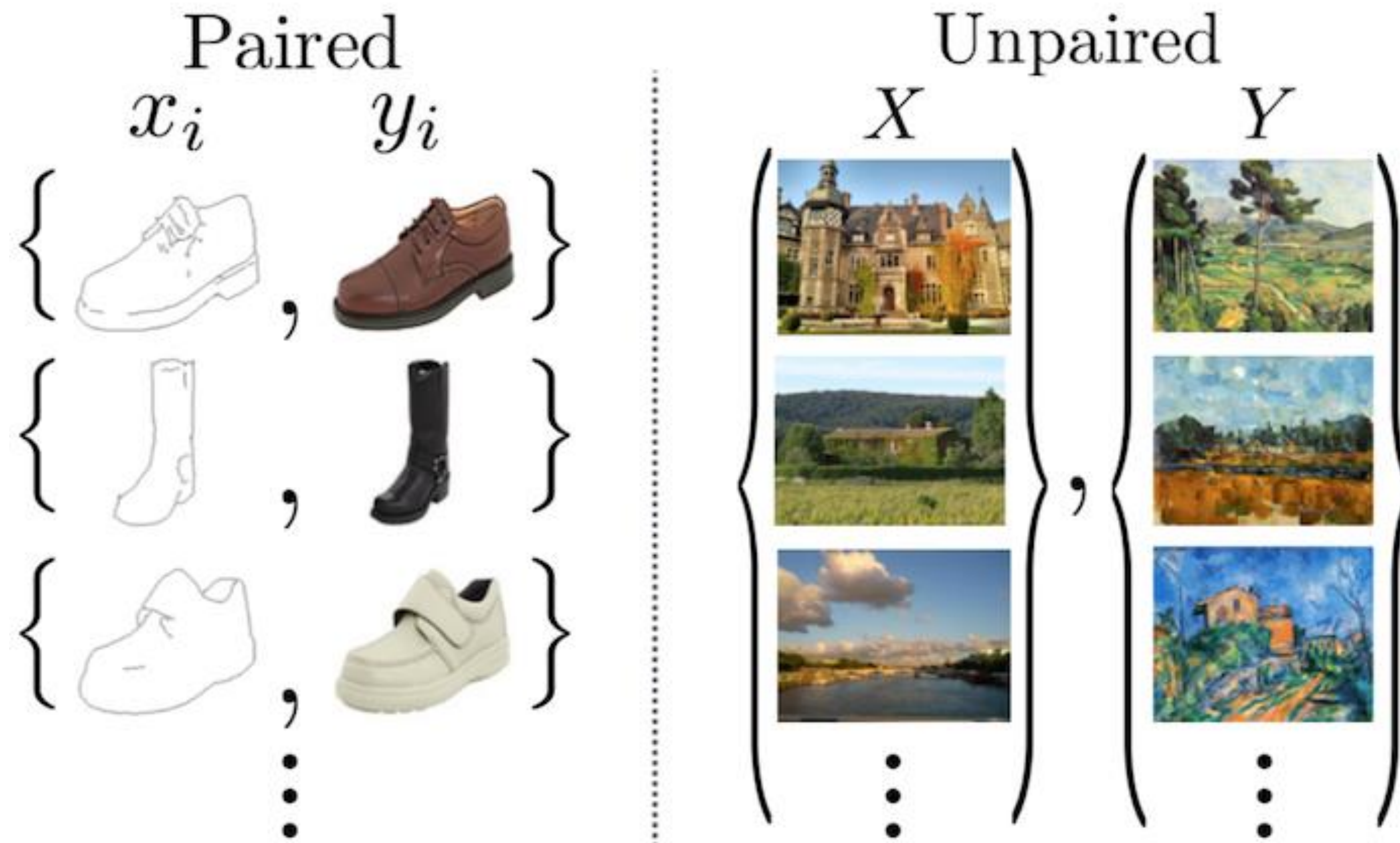
# Demo: Pix2Pix



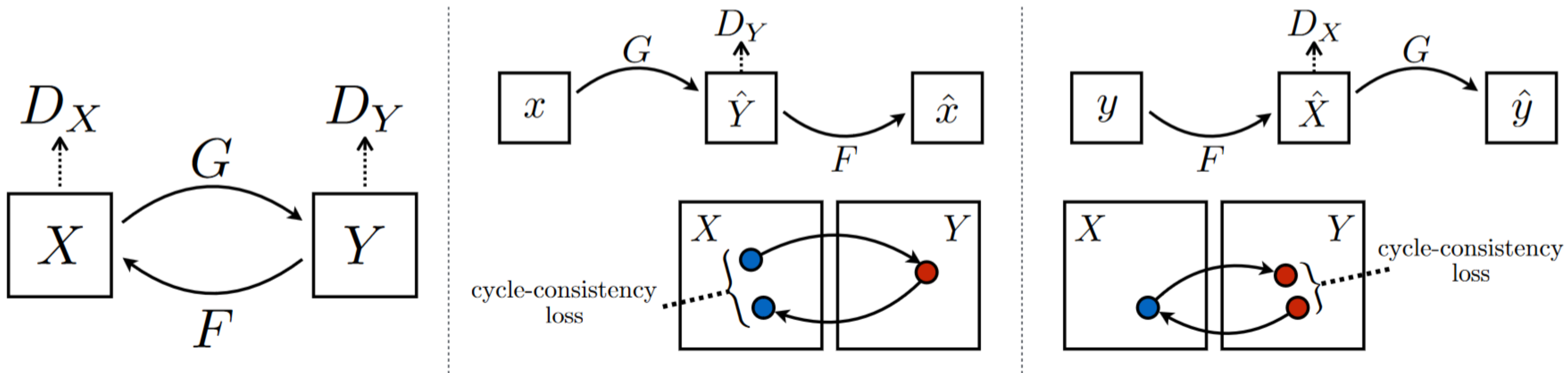
# Limitations

1. Paired data is required

# CycleGAN



# Cycle Consistency



# CycleGAN





# Recycle-GAN

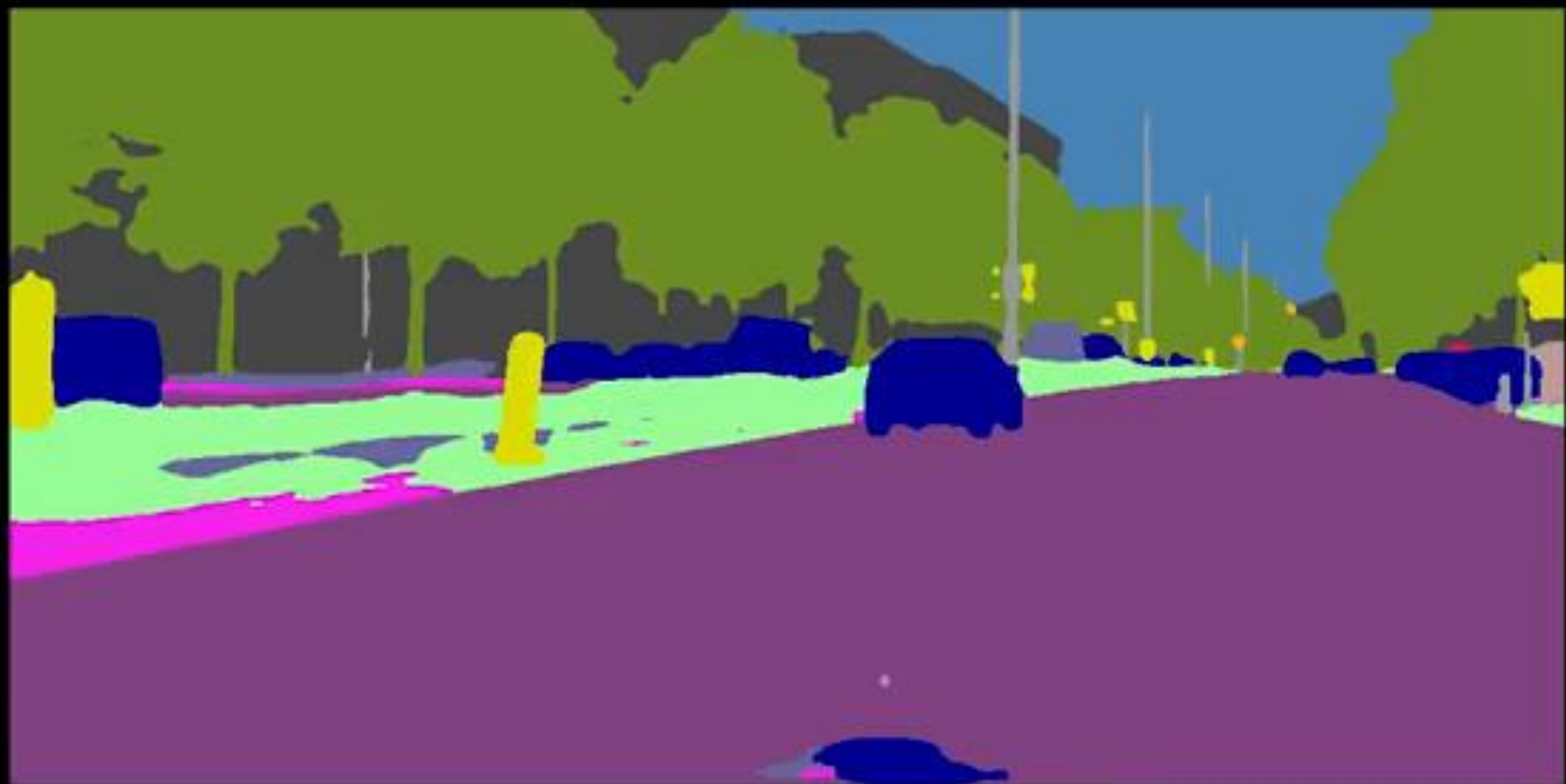


# Limitations

1. Paired data is required

# Limitations

1. Paired data is required
2. Temporally instable if applied per-frame to a video sequence



Labels



pix2pixHD



COVST



Ours

# Video-to-Video Synthesis

Ting-Chun Wang<sup>1</sup>, Ming-Yu Liu<sup>1</sup>, Jun-Yan Zhu<sup>2</sup>, Guilin Liu<sup>1</sup>,  
Andrew Tao<sup>1</sup>, Jan Kautz<sup>1</sup>, Bryan Catanzaro<sup>1</sup>

<sup>1</sup>NVIDIA Corporation <sup>2</sup>MIT

# Limitations

1. Paired data is required
2. Temporally instable if applied per-frame to a video sequence

# Limitations

1. Paired data is required
2. Temporally instable if applied per-frame to a video sequence
3. Does not generalize to 3D transformations

# DeepVoxels

Tatarchenko et al. [2015]



Pix2Pix [Isola et al. 2017]



Ground Truth



Worrall et al. [2017]



DeepVoxels





# Summary

- Convolutional Neural Networks

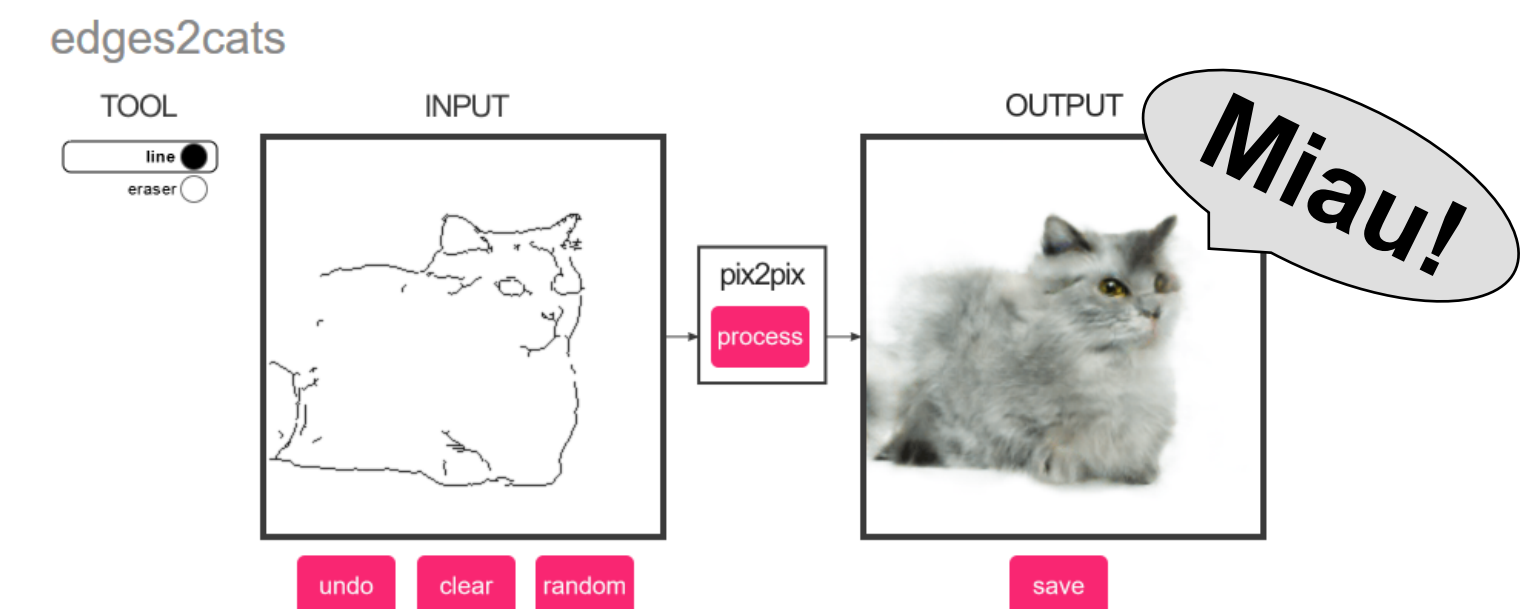
$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$

- Generative Modeling



(Brundage et al., 2018)

- Pix2Pix (“mapping from A to B”)



# References

- CVPR GAN Tutorial
  - <https://sites.google.com/view/cvpr2018tutorialongans>
- CS231n
  - [http://cs231n.stanford.edu/slides/2016/winter1516\\_lecture7.pdf](http://cs231n.stanford.edu/slides/2016/winter1516_lecture7.pdf)
- ...