

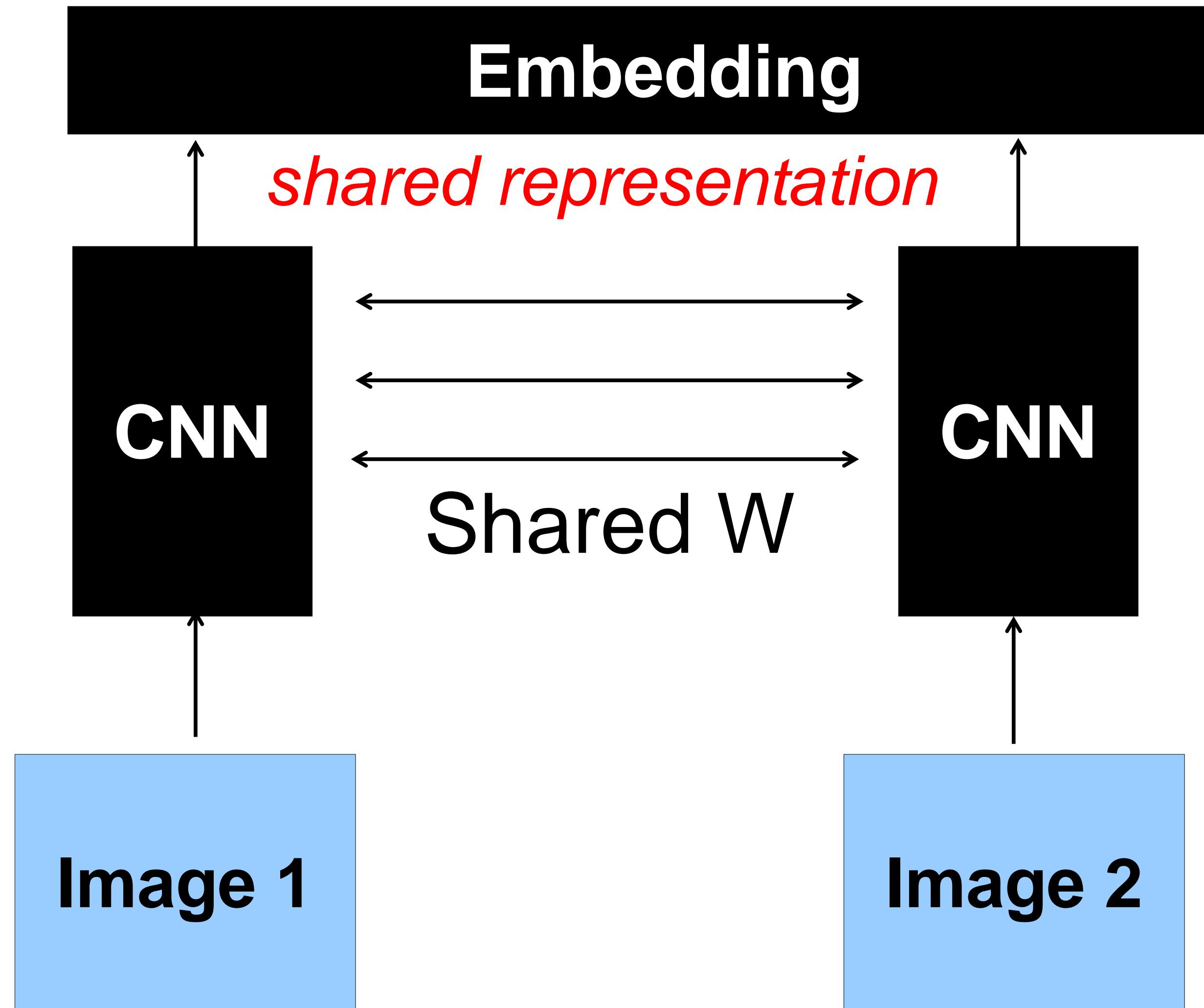
# Fancier Neural Networks, fun applications, and scary cats



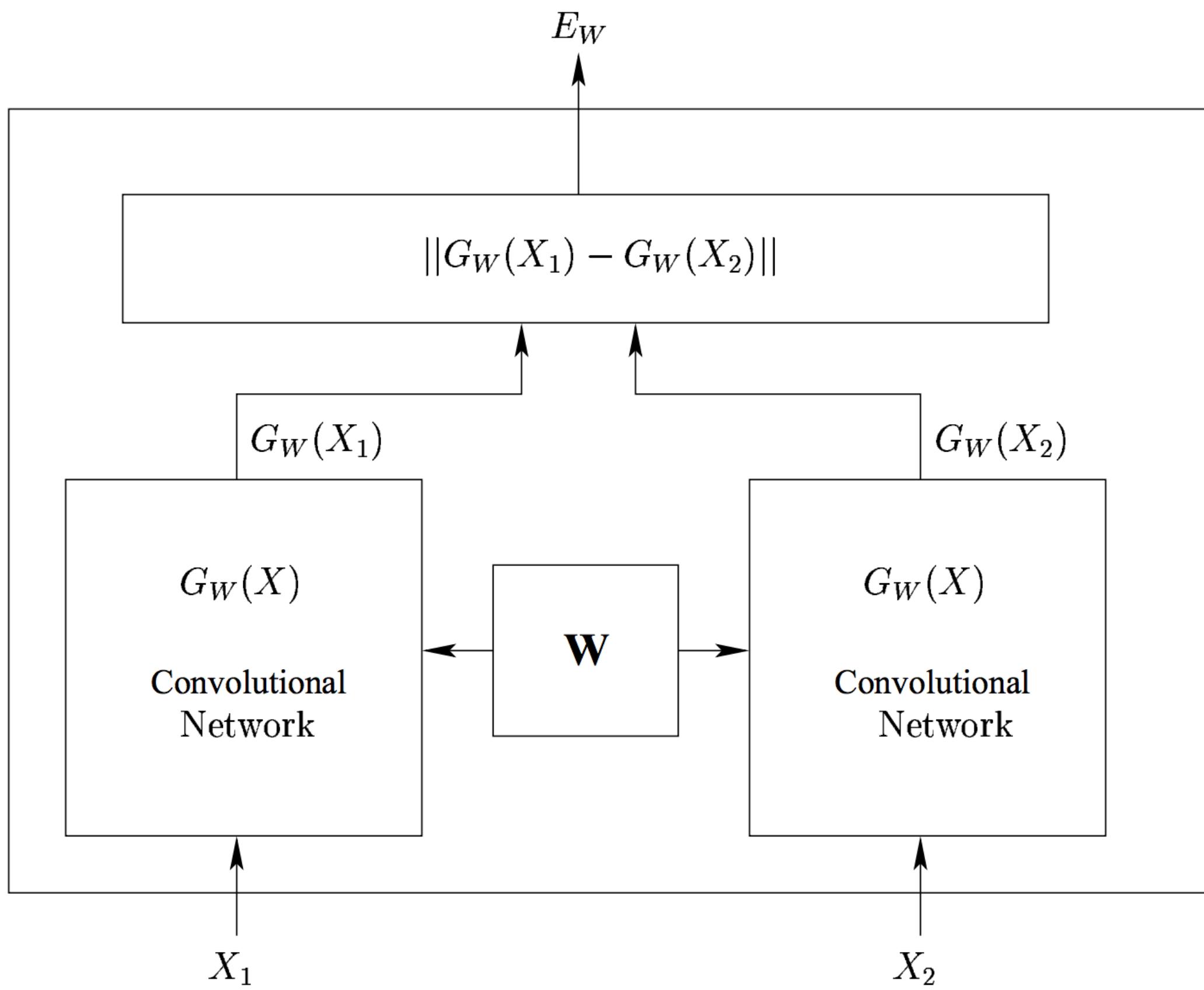
Many slides stolen from Sean Bell,  
Jon Long, Carl Doersch, and Phillip Isola

CS280, Spring 2017

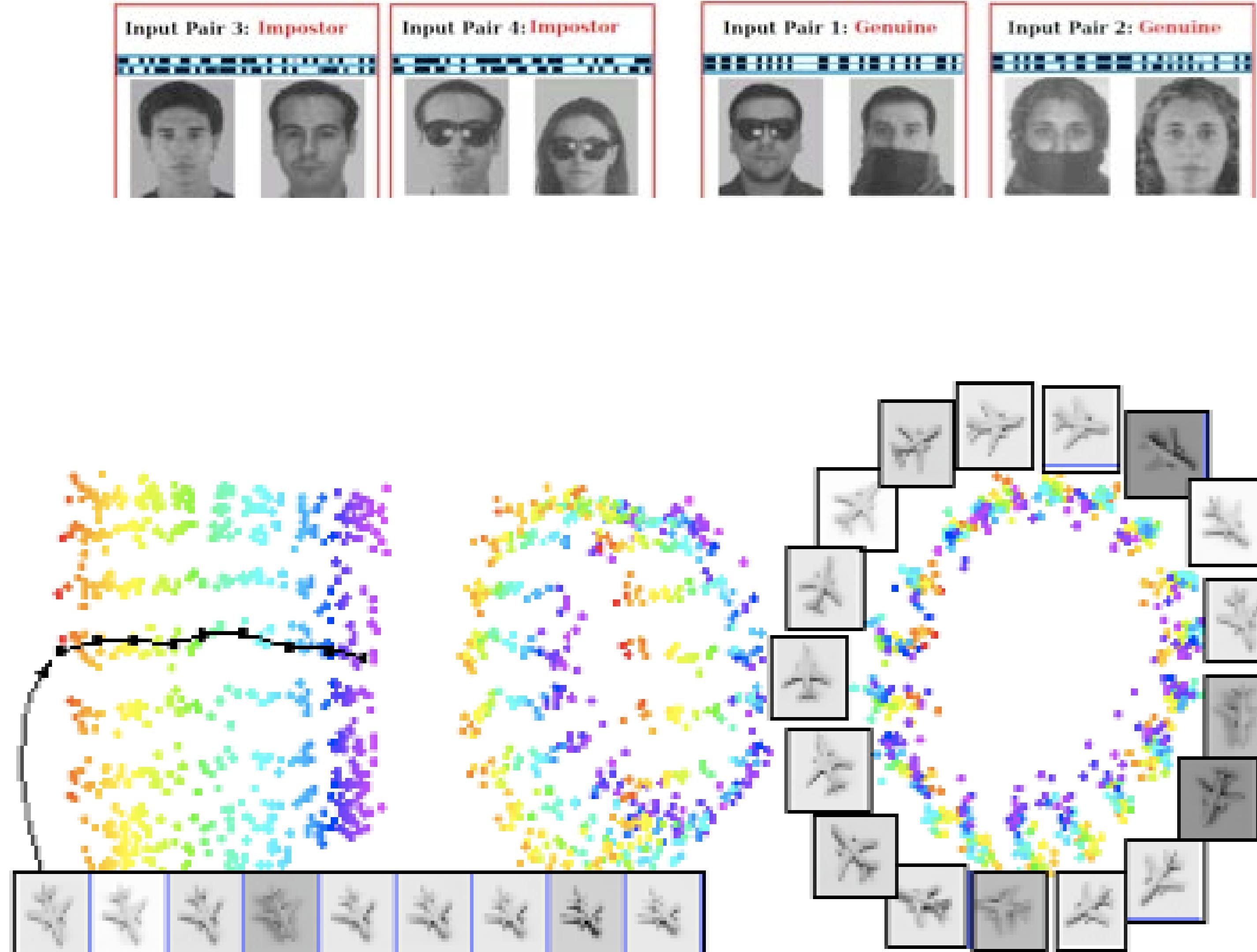
# Fancier Architectures: “Siamese”



# Siamese Network



Siamese Architecture  
[Chopra 2005, Hadsell 2006]



# LEARNING VISUAL SIMILARITY FOR PRODUCT DESIGN WITH CONVOLUTIONAL NEURAL NETWORKS

SEAN BELL AND KAVITA BALA  
CORNELL UNIVERSITY

# THE PROBLEM

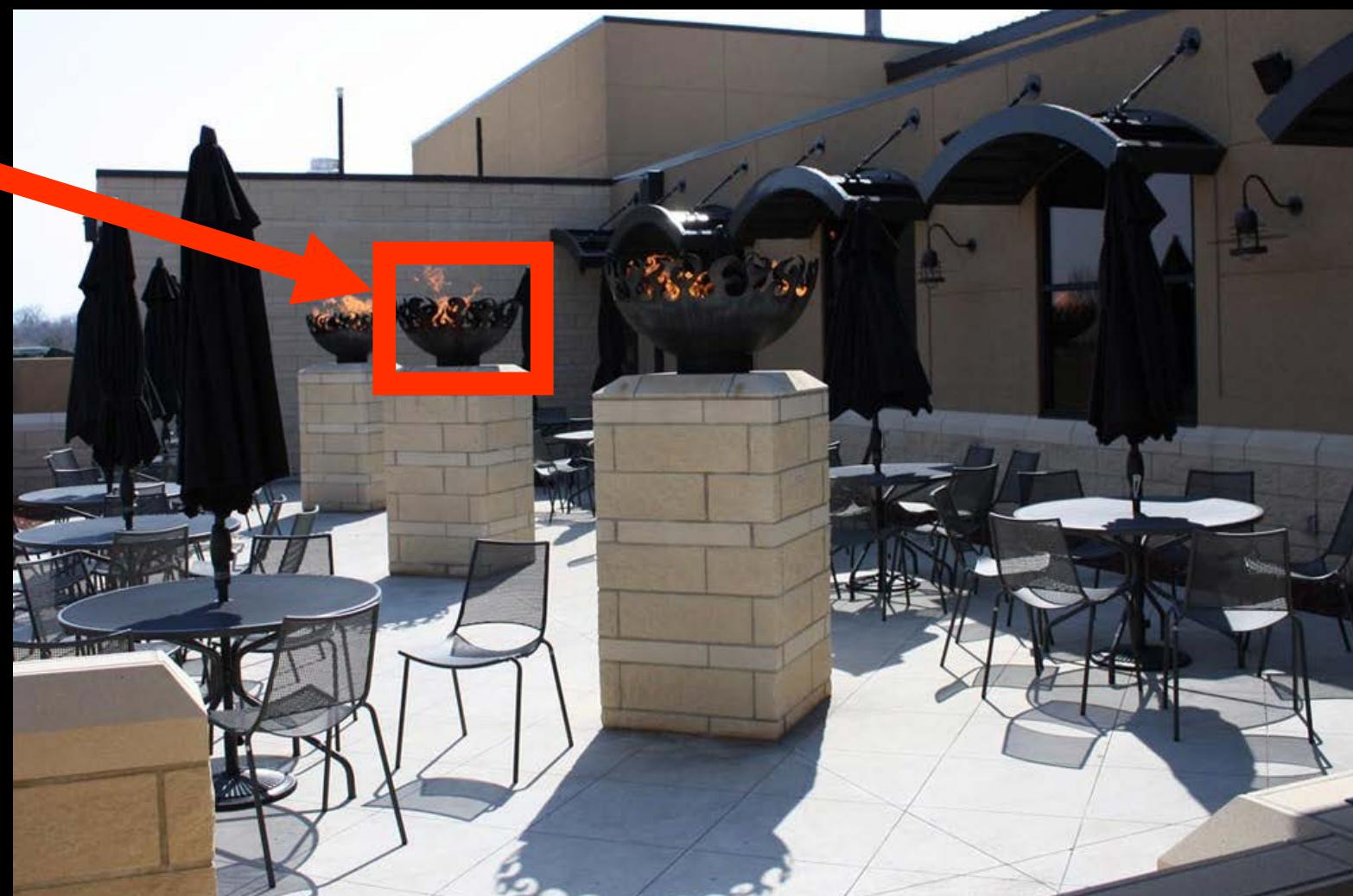
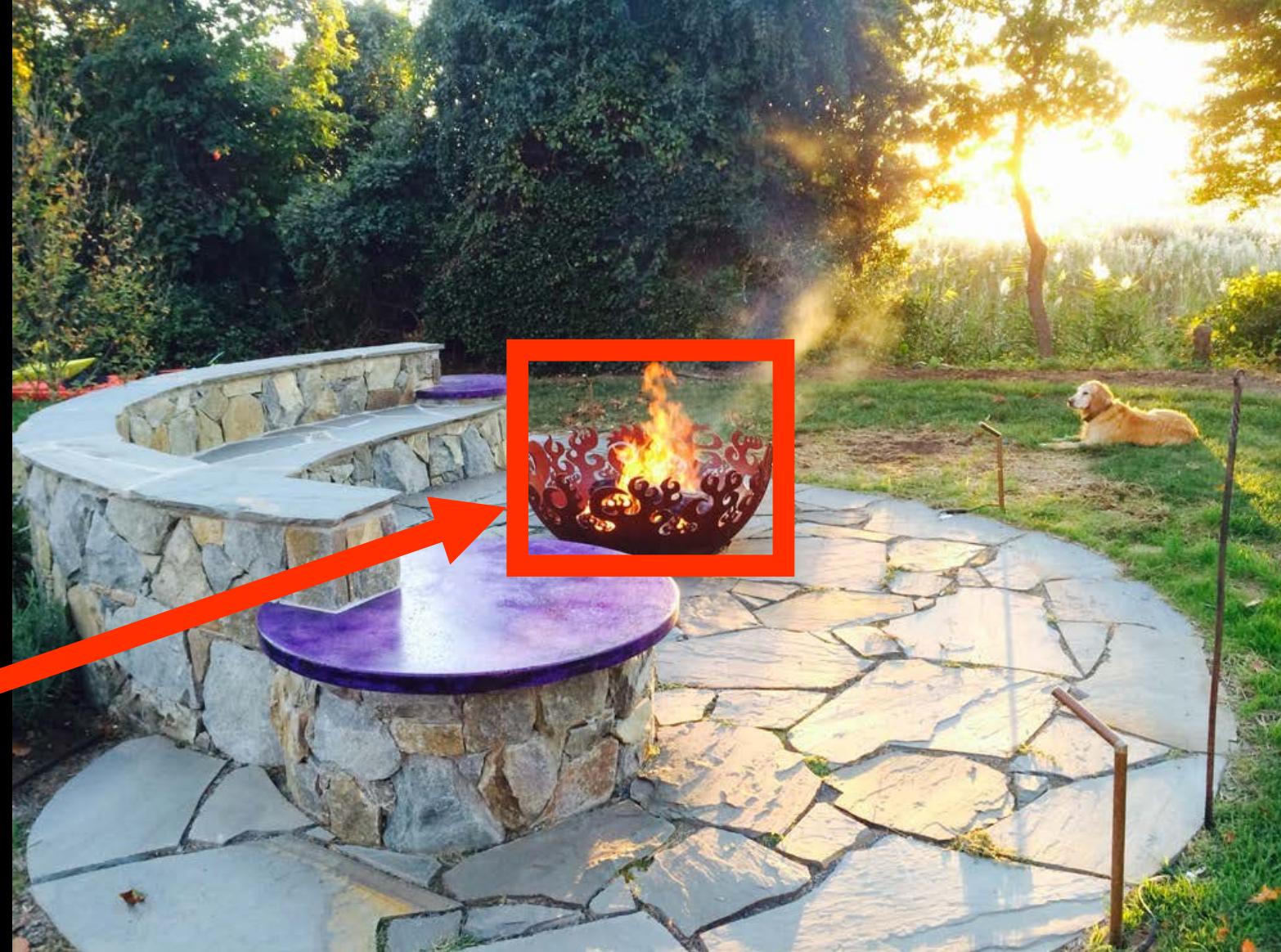
(1) "What is this?"



(2) "Where is it used?"



Name: "Great Bowl O'Fire  
Sculptural Fire Bowl"  
Category: Fire pit  
Sold by: John T. Unger, LLC



# THE PROBLEM

(1) "What is this?"

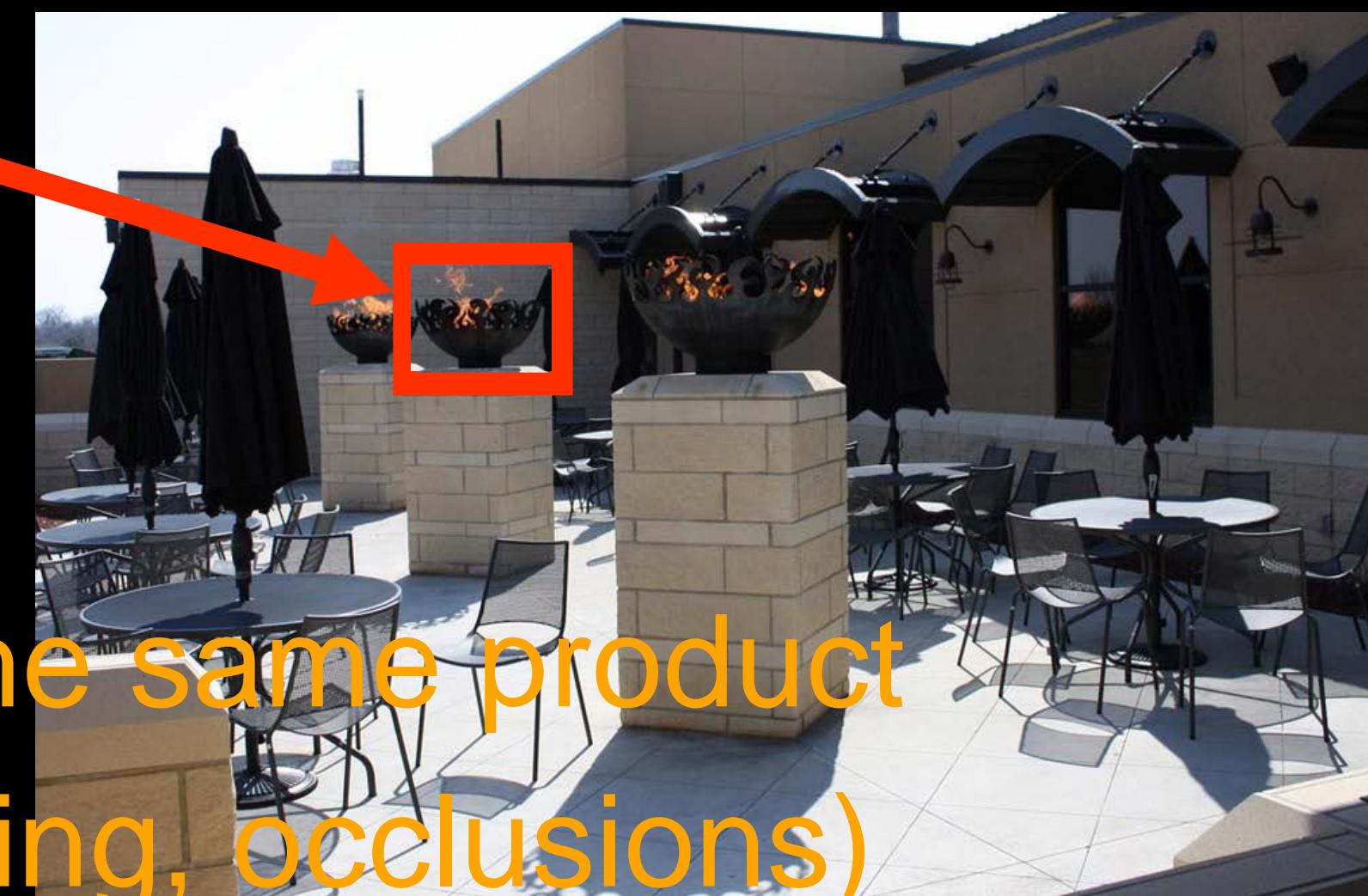
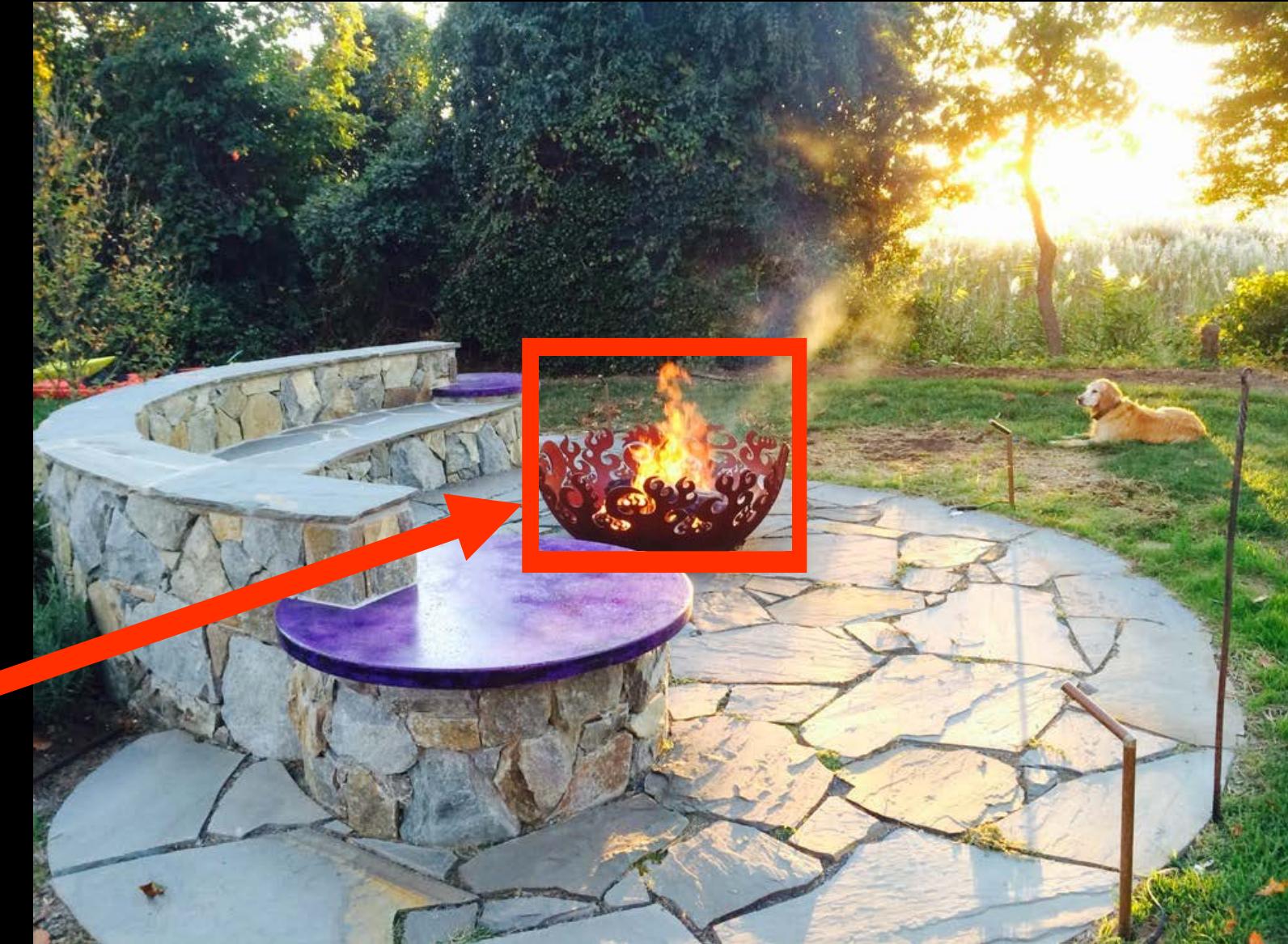


(2) "Where is it used?"



Name: "Great Bowl O'Fire  
Sculptural Fire Bowl"

Challenge: determine whether these are the same product  
Category: Fire pit  
(different resolution, viewpoint, color, lighting, occlusions)  
Sold by: John T. Unger, LLC



# TWO KINDS OF IMAGES

Iconic



In context



(From a product website)

(Cropped from a scene photo)

# PROJECTING INTO A JOINT EMBEDDING

Iconic

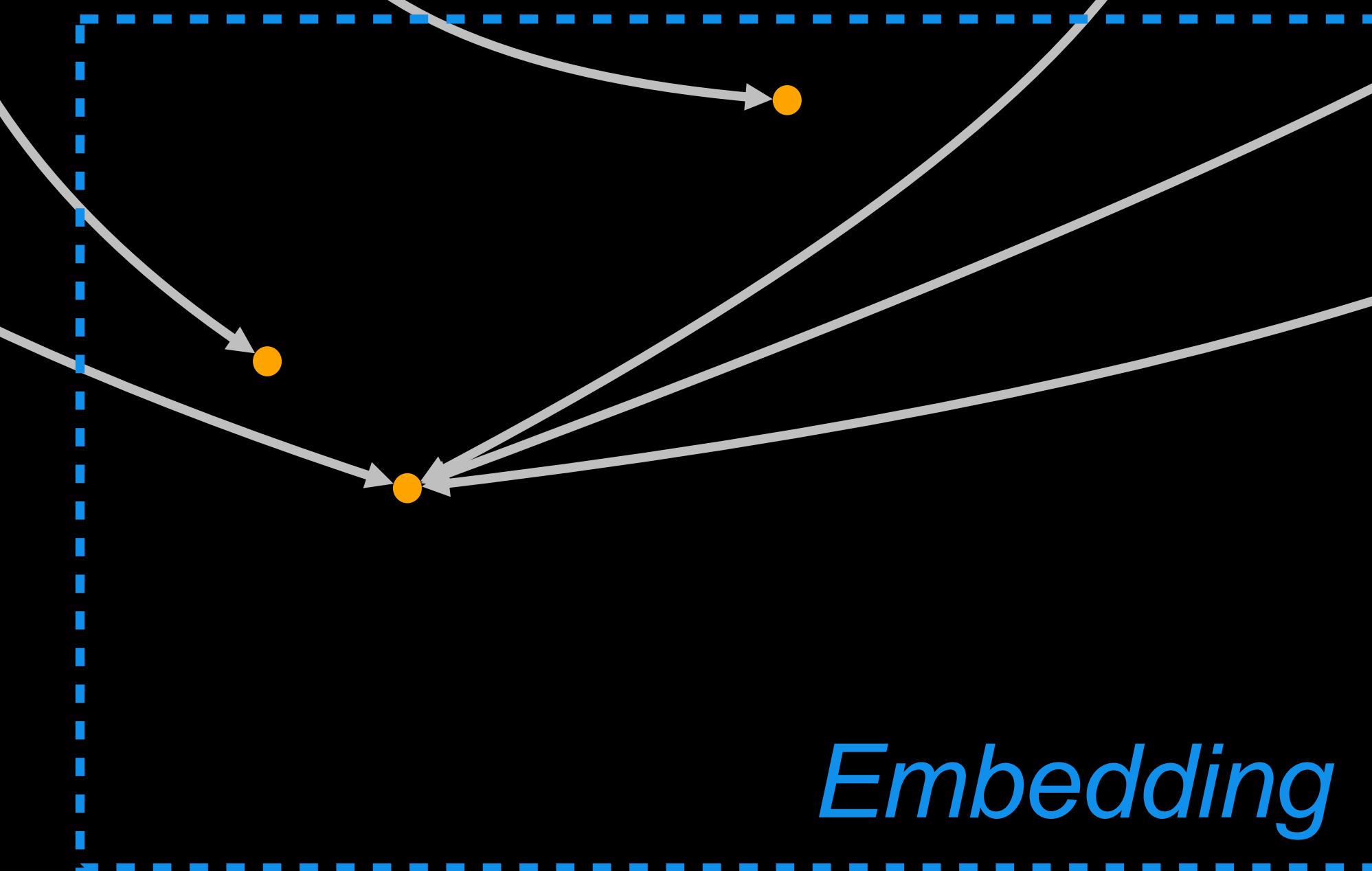


In context



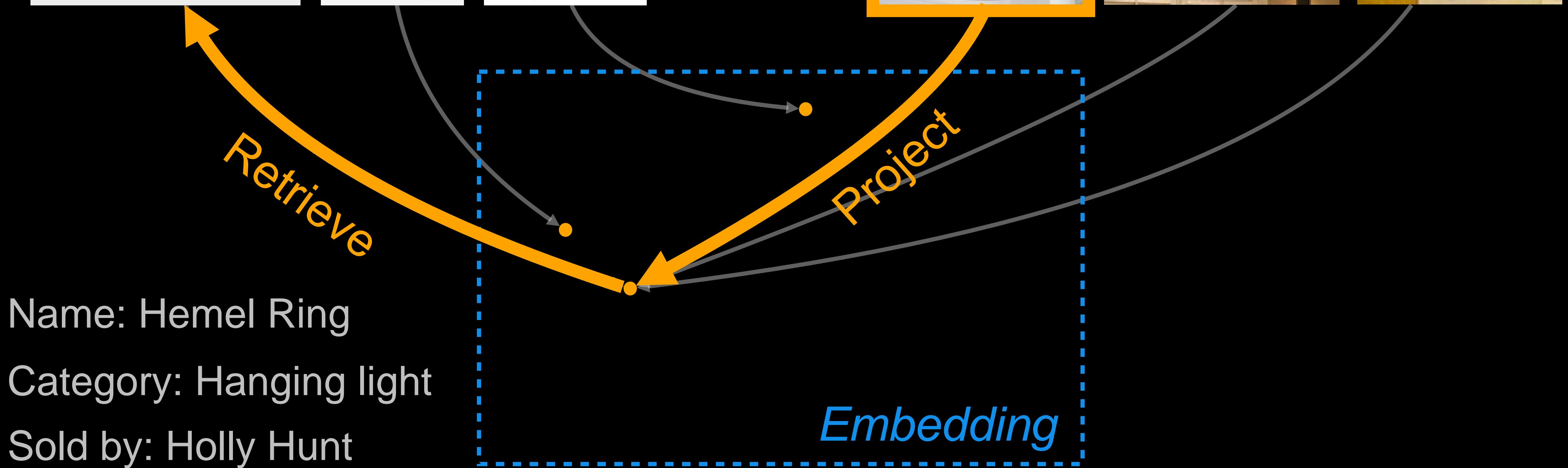
*Project*

*Embedding*



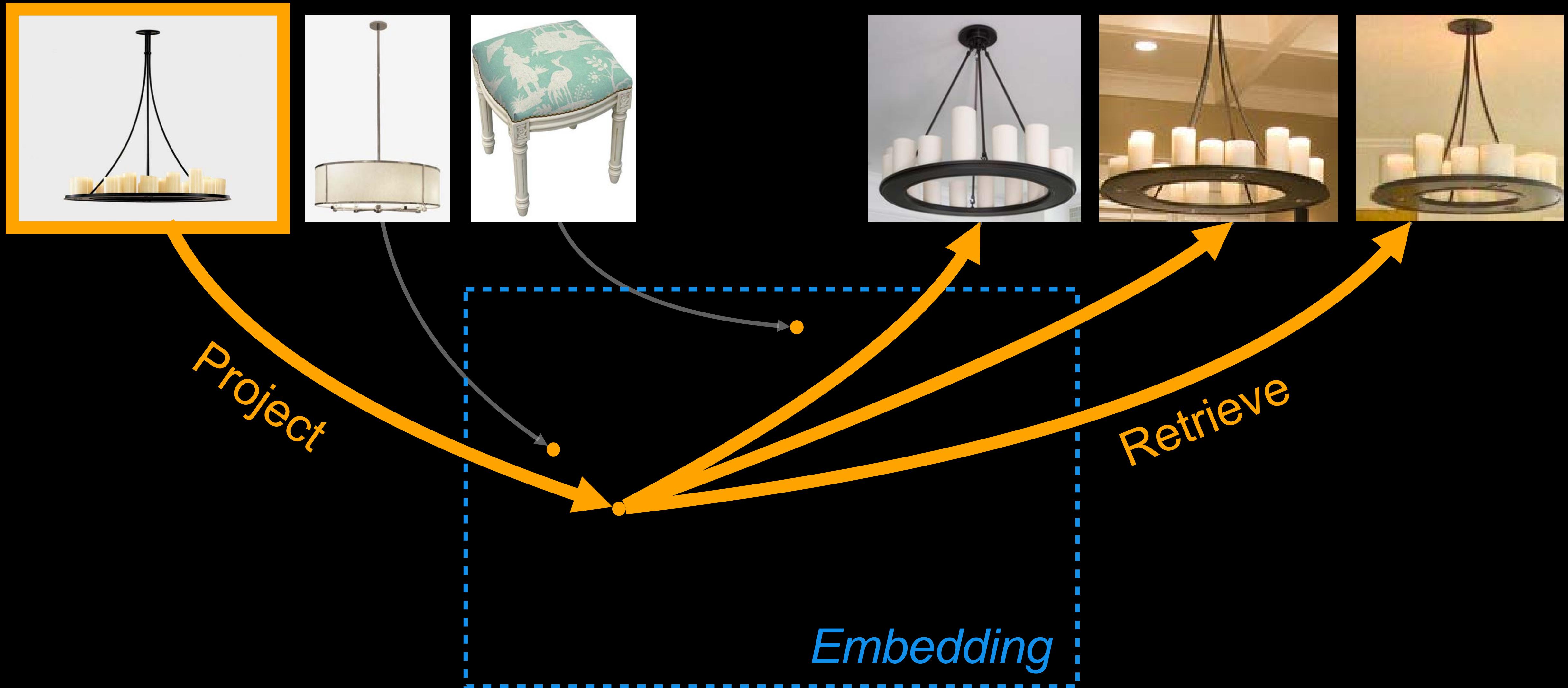
# SEARCH USING THE EMBEDDING

“What is it?”

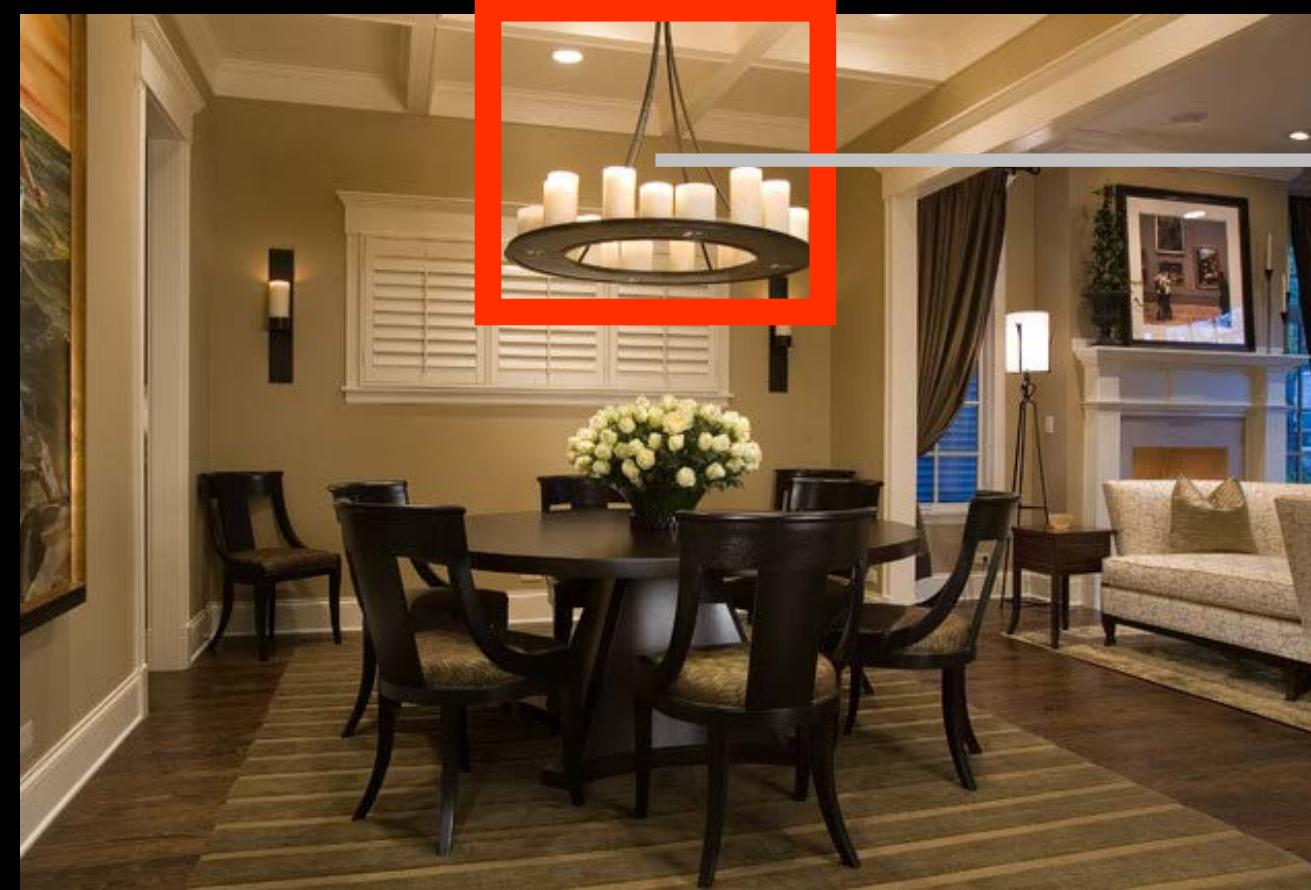


# SEARCH USING THE EMBEDDING

“Where is it used?”



# CONTRASTIVE LOSS: POSITIVE EXAMPLE



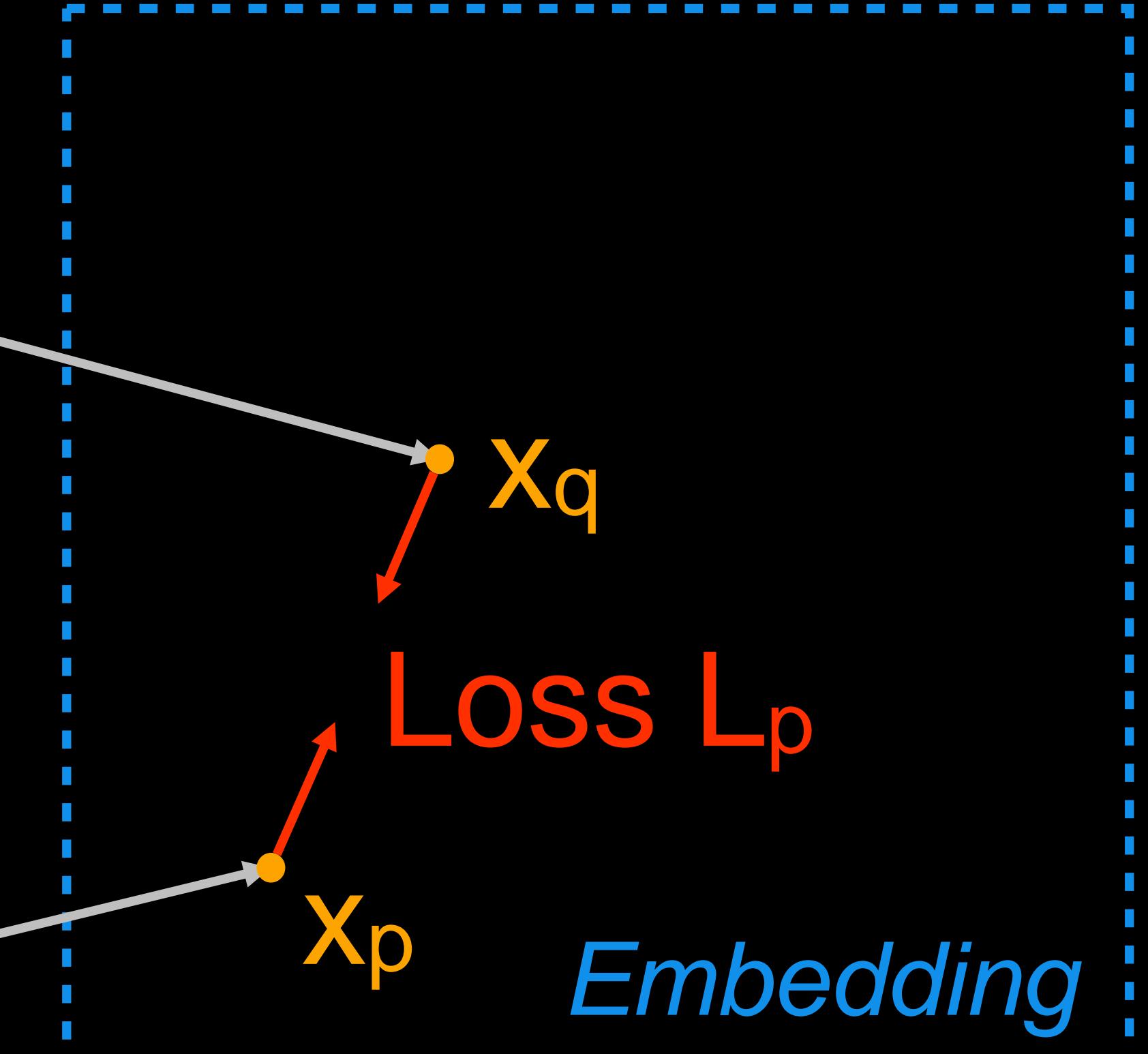
In context



Iconic (same)

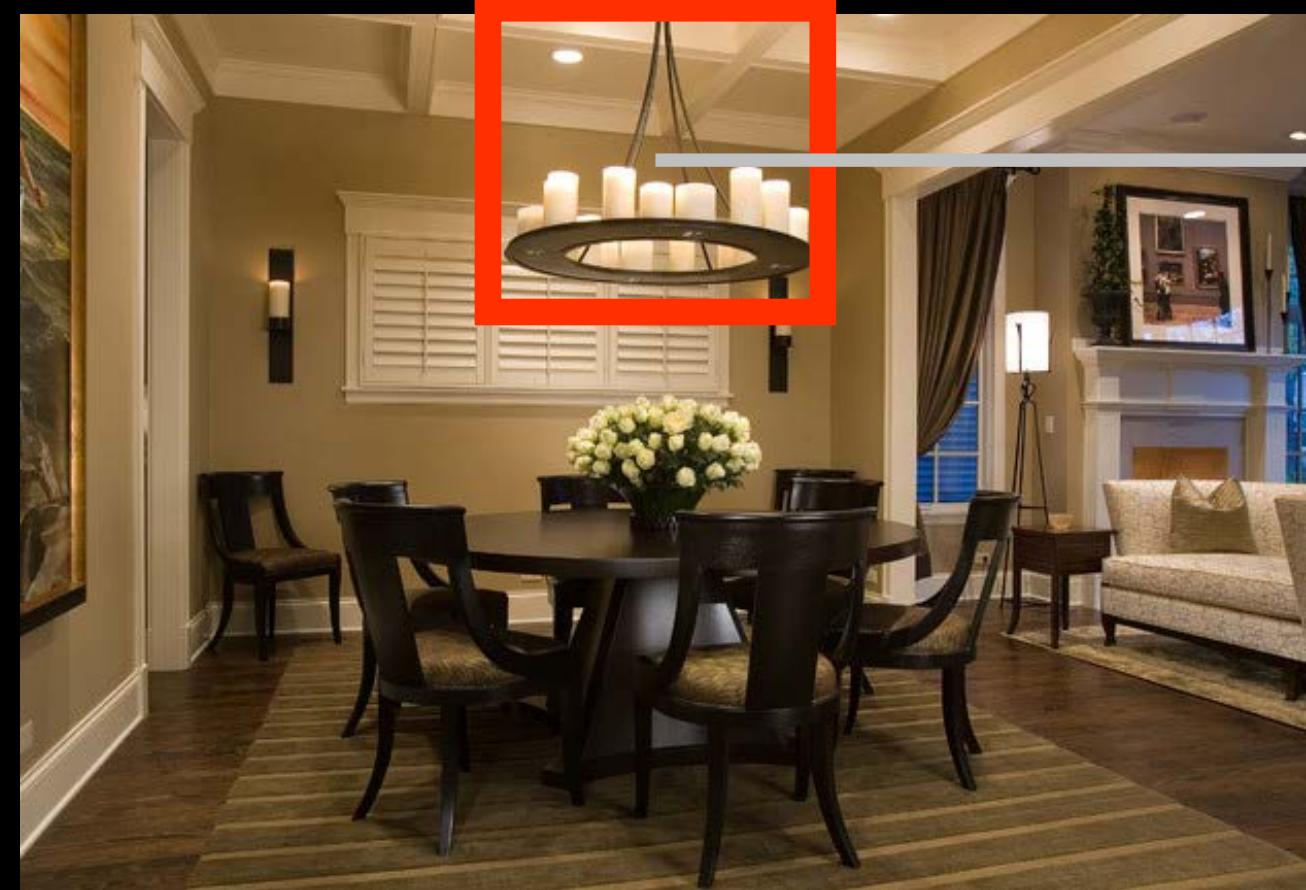


Parameters  $\theta$



$$L_p(x_q, x_p) = \|x_q - x_p\|_2^2$$

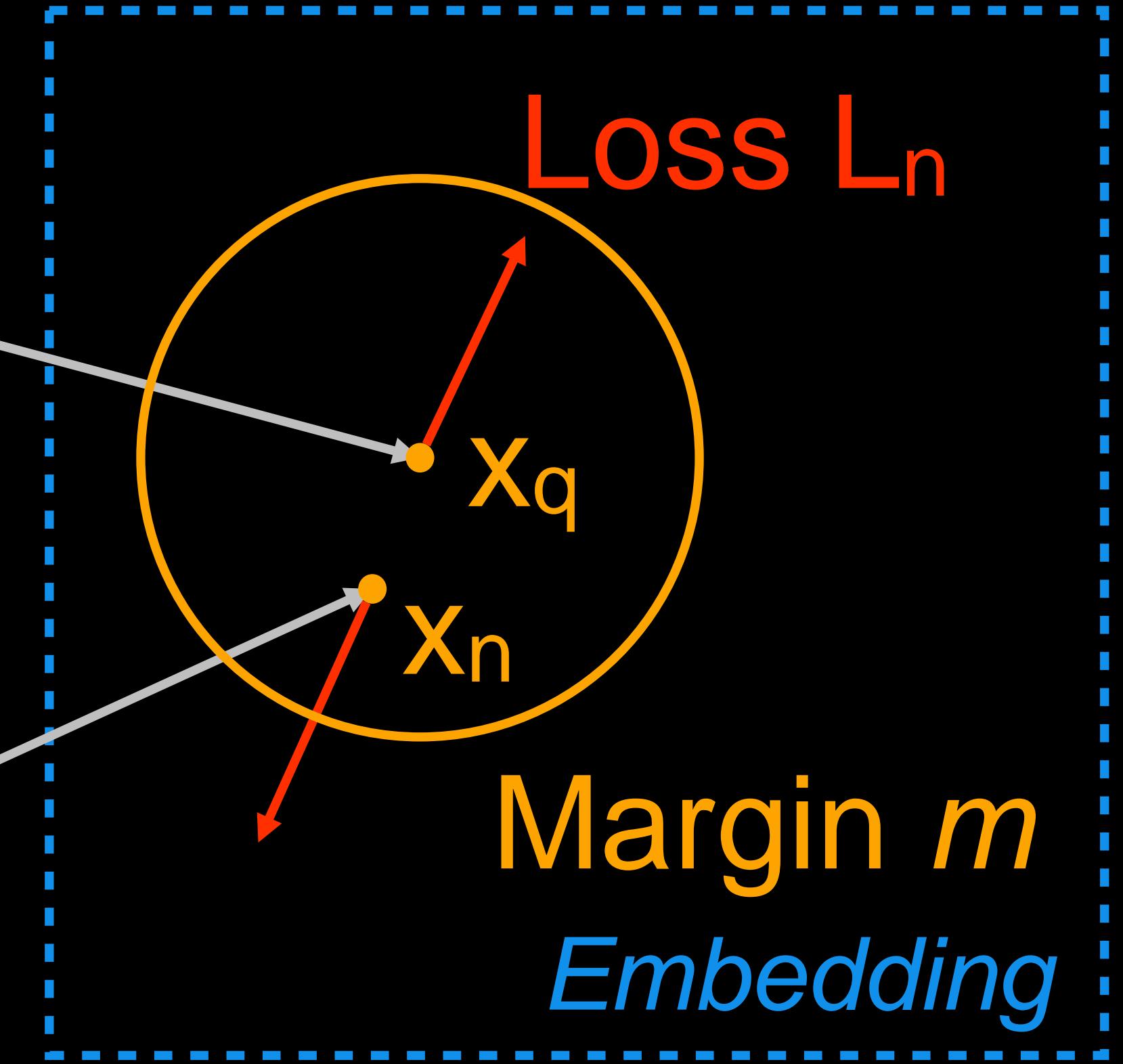
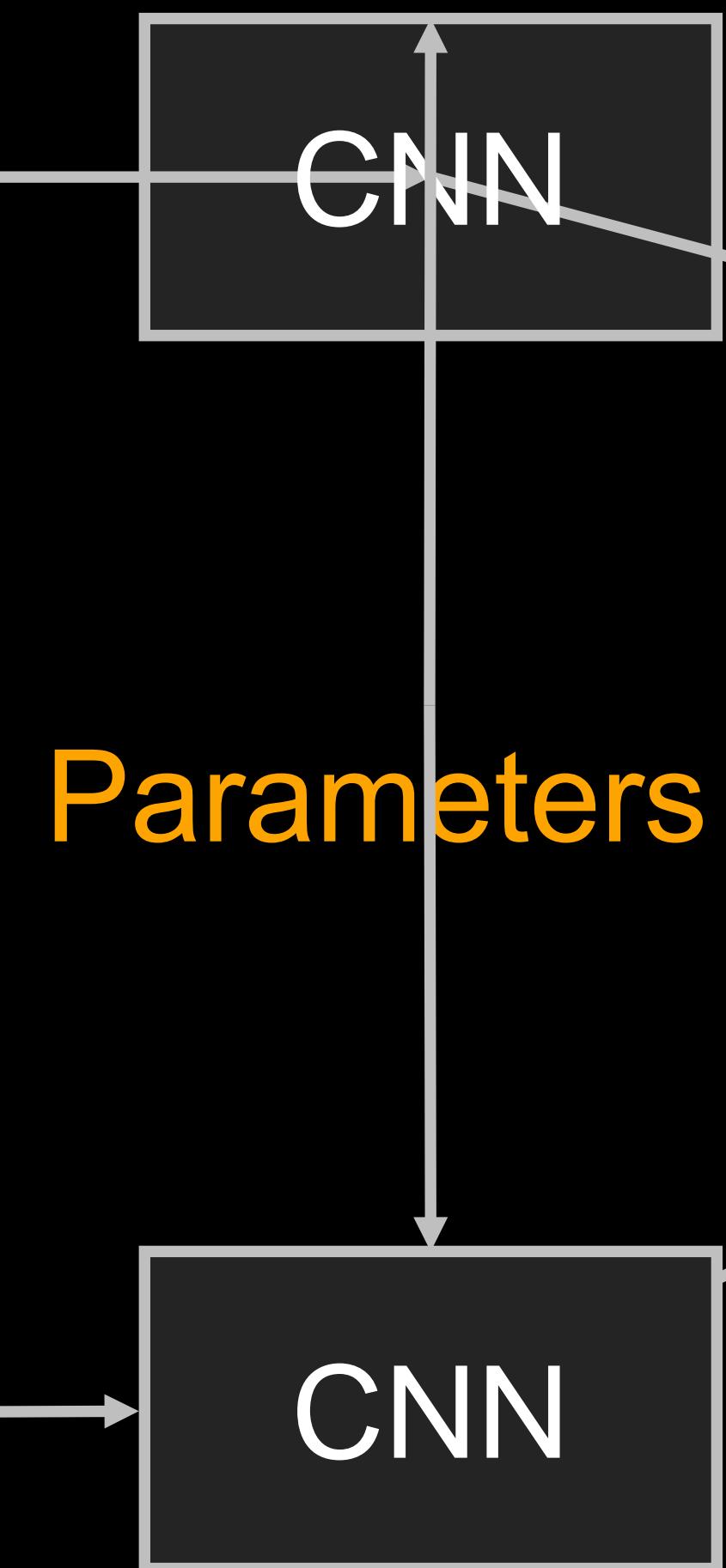
# CONTRASTIVE LOSS: NEGATIVE EXAMPLE



In context



Iconic (different)



$$L_n(x_q, x_n) = \max(0, m^2 - \|x_q - x_n\|_2^2)$$

# CONTRASTIVE LOSS: ALL TOGETHER

$$L(\theta) = \underbrace{\sum_{(x_q, x_p)} L_p(x_q, x_p)}_{\text{Penalty for similar images that are far away}} + \underbrace{\sum_{(x_q, x_n)} L_n(x_q, x_n)}_{\text{Penalty for dissimilar images that are nearby}}$$

$$L_p(x_q, x_p) = \|x_q - x_p\|_2^2$$

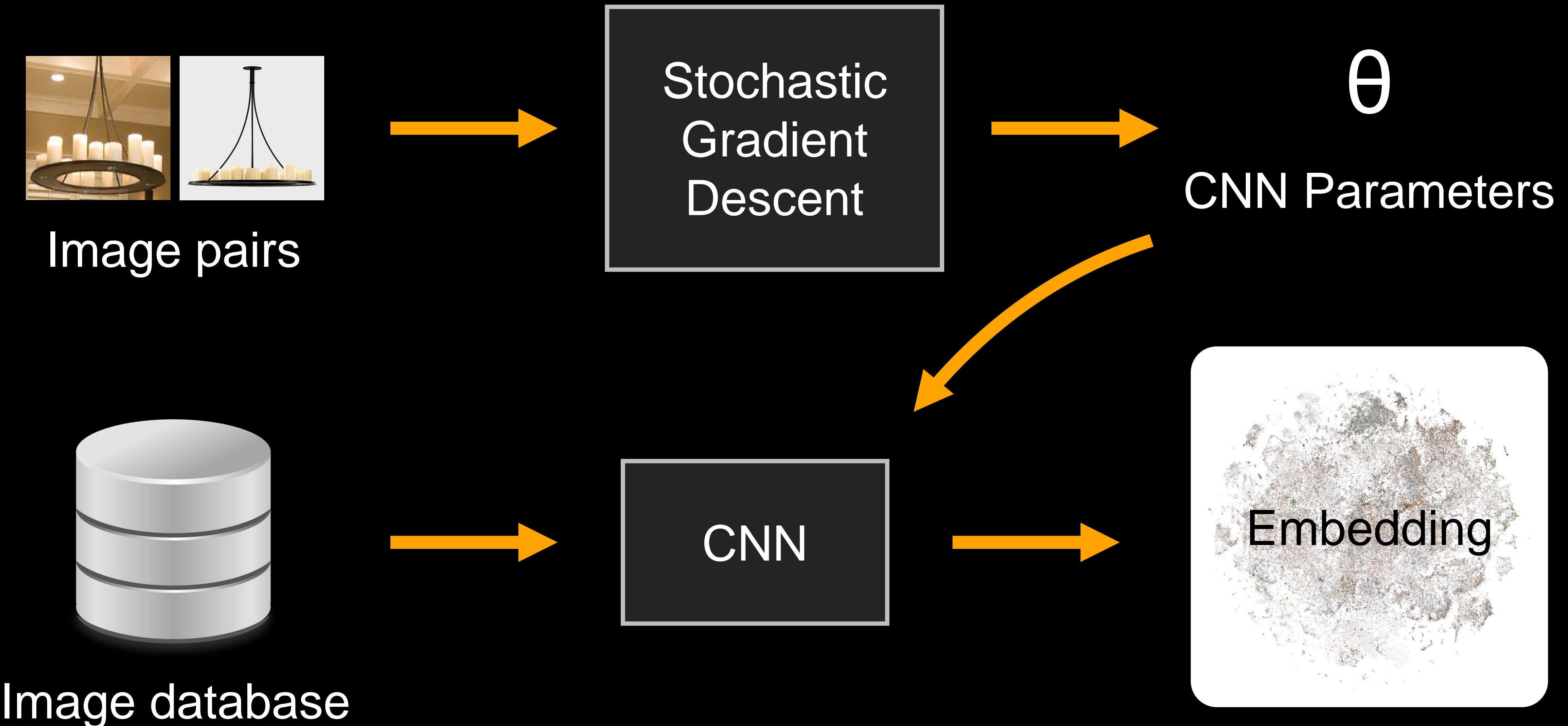
Margin

$$L_n(x_q, x_n) = \max(0, m^2 - \|x_q - x_n\|_2^2)$$

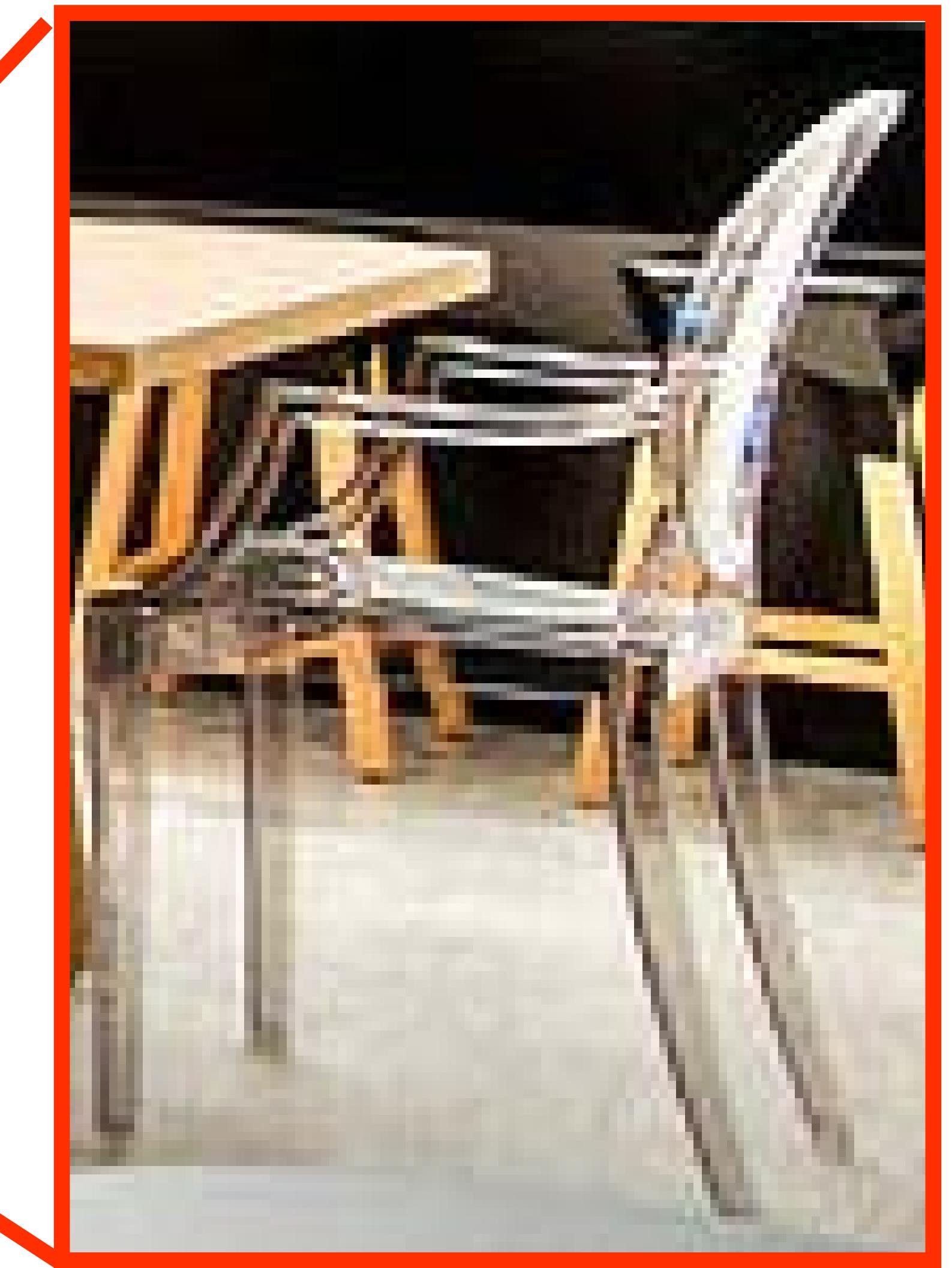
Minimize  $L(\theta)$  with stochastic gradient descent and momentum

[Chopra 2005, Hadsell 2006]

# TRAINING PIPELINE



# RESULTS: “WHAT IS IT?”



In context

# RESULTS: “WHAT IS IT?”

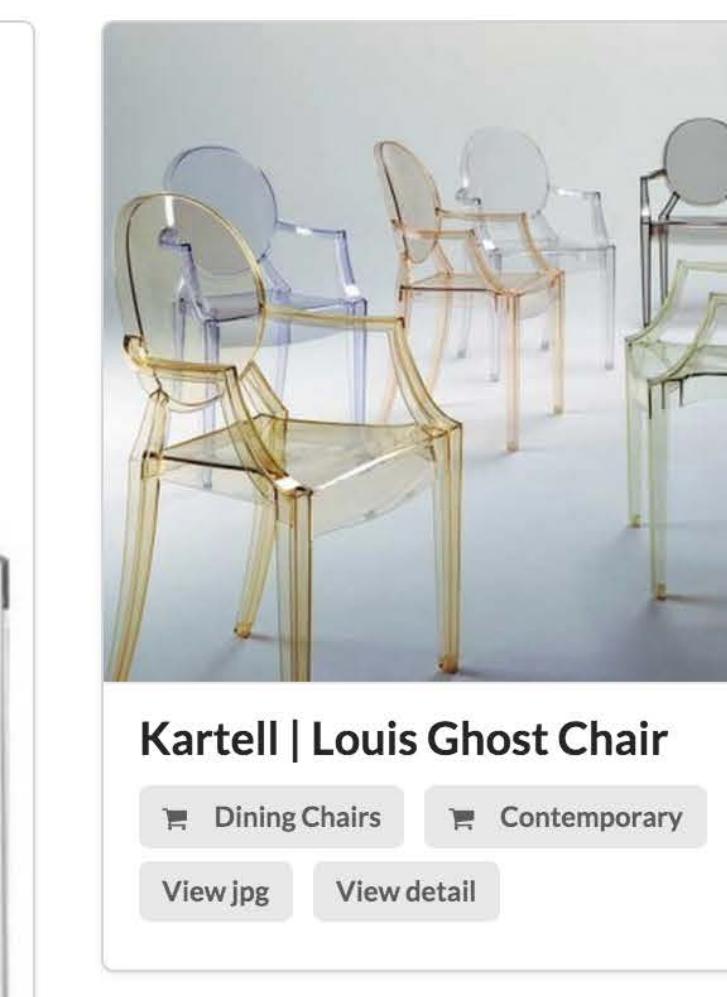
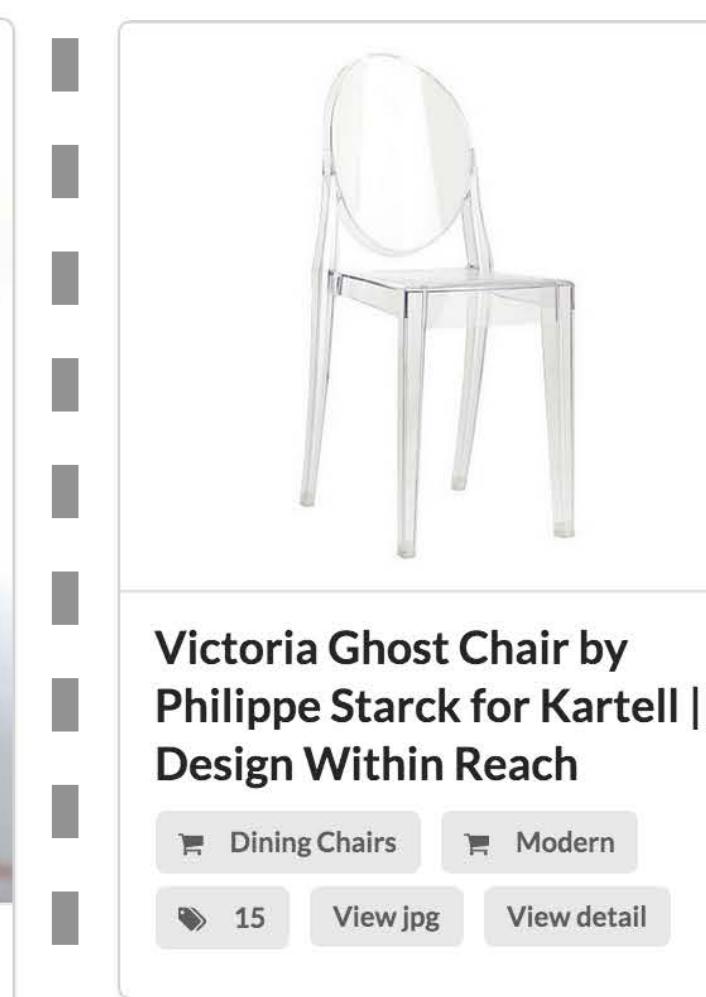
In context



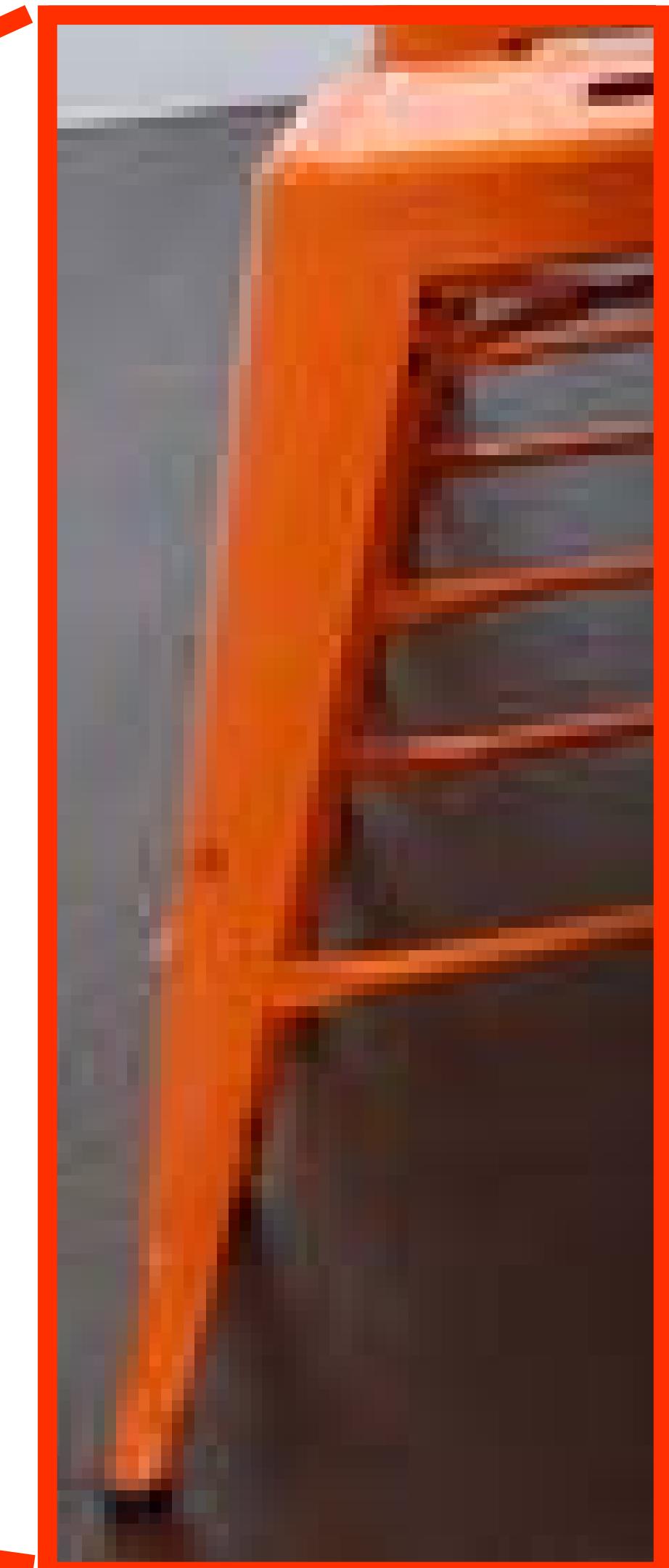
Iconic



Top 4 results:



# RESULTS: “WHAT IS IT?”



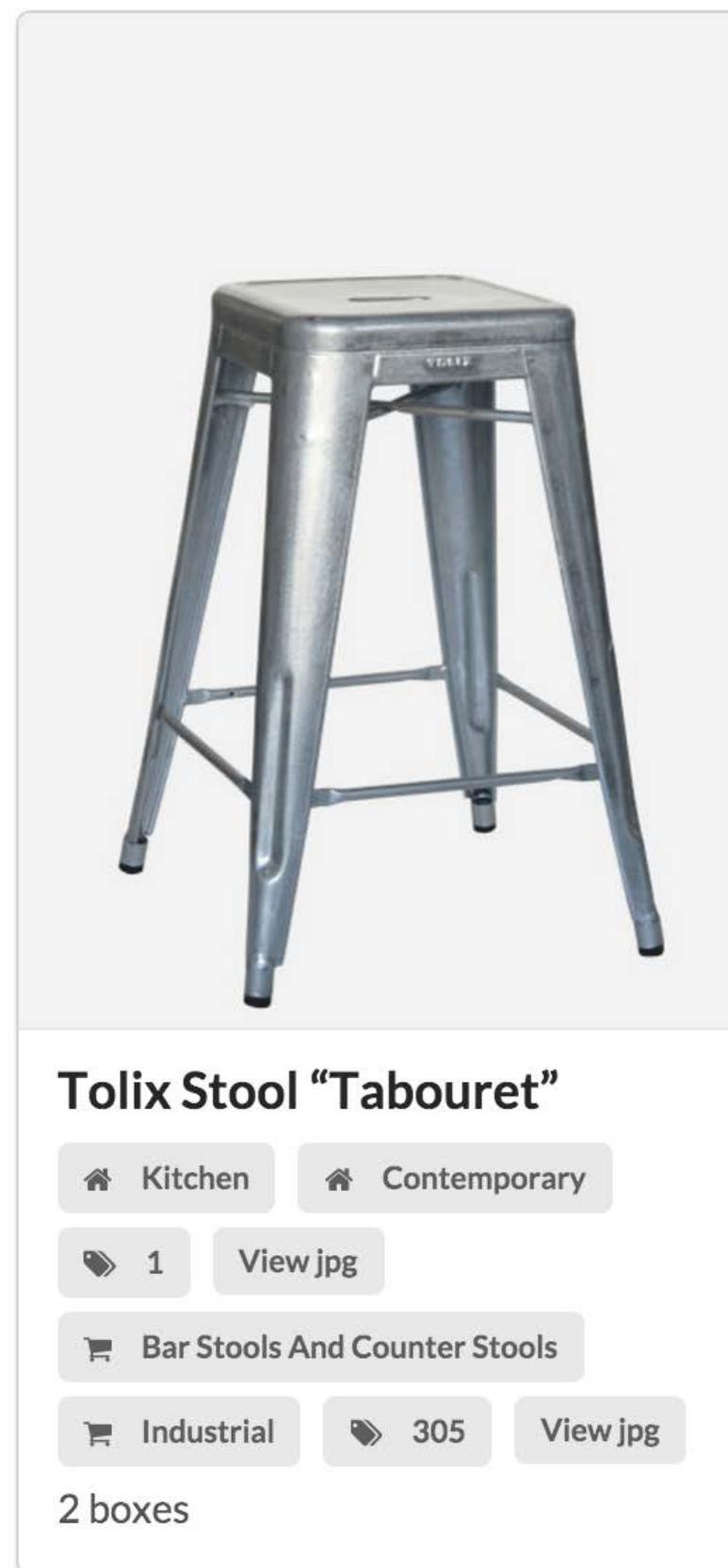
In context

# RESULTS: “WHAT IS IT?”

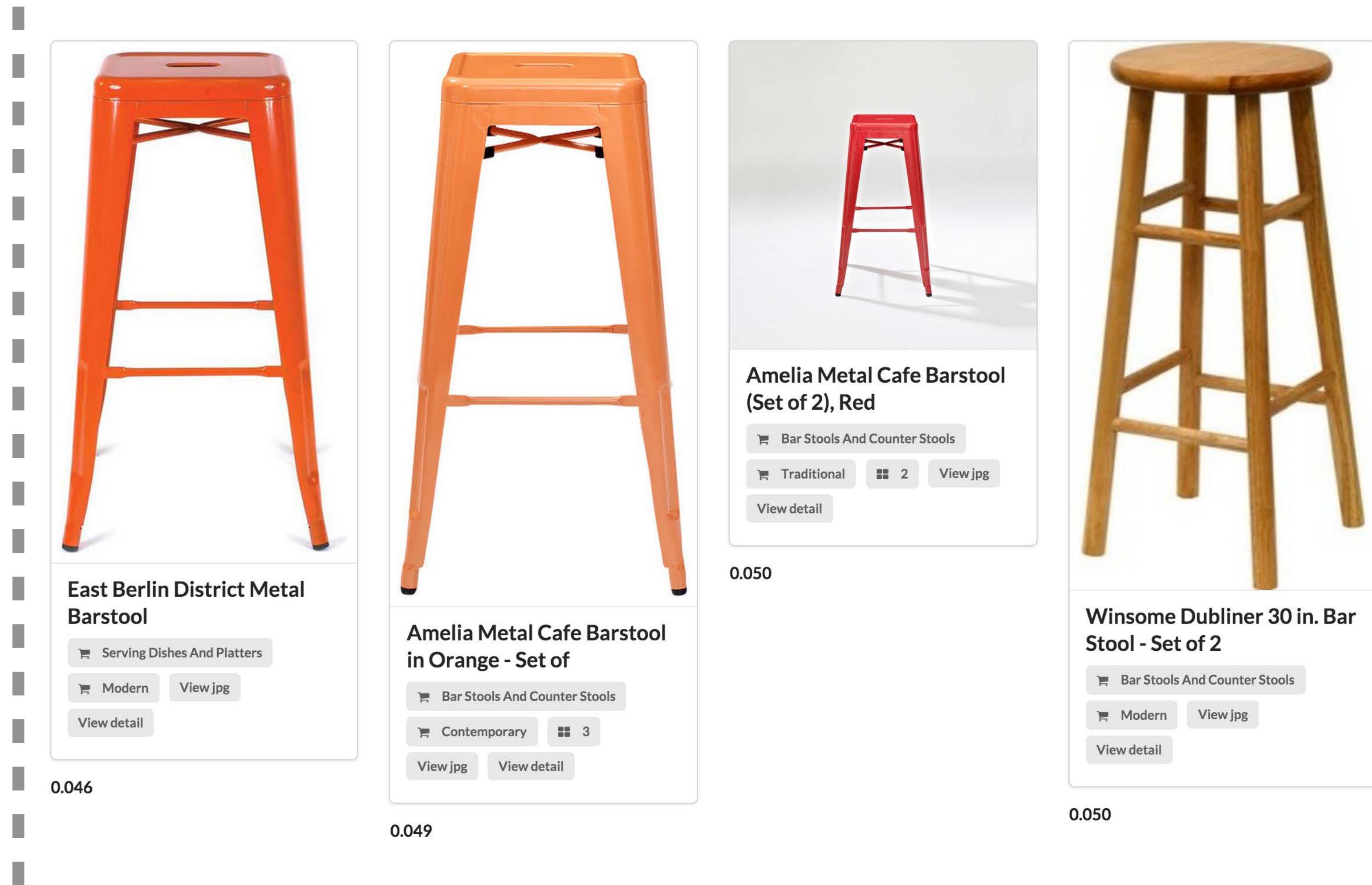
In context



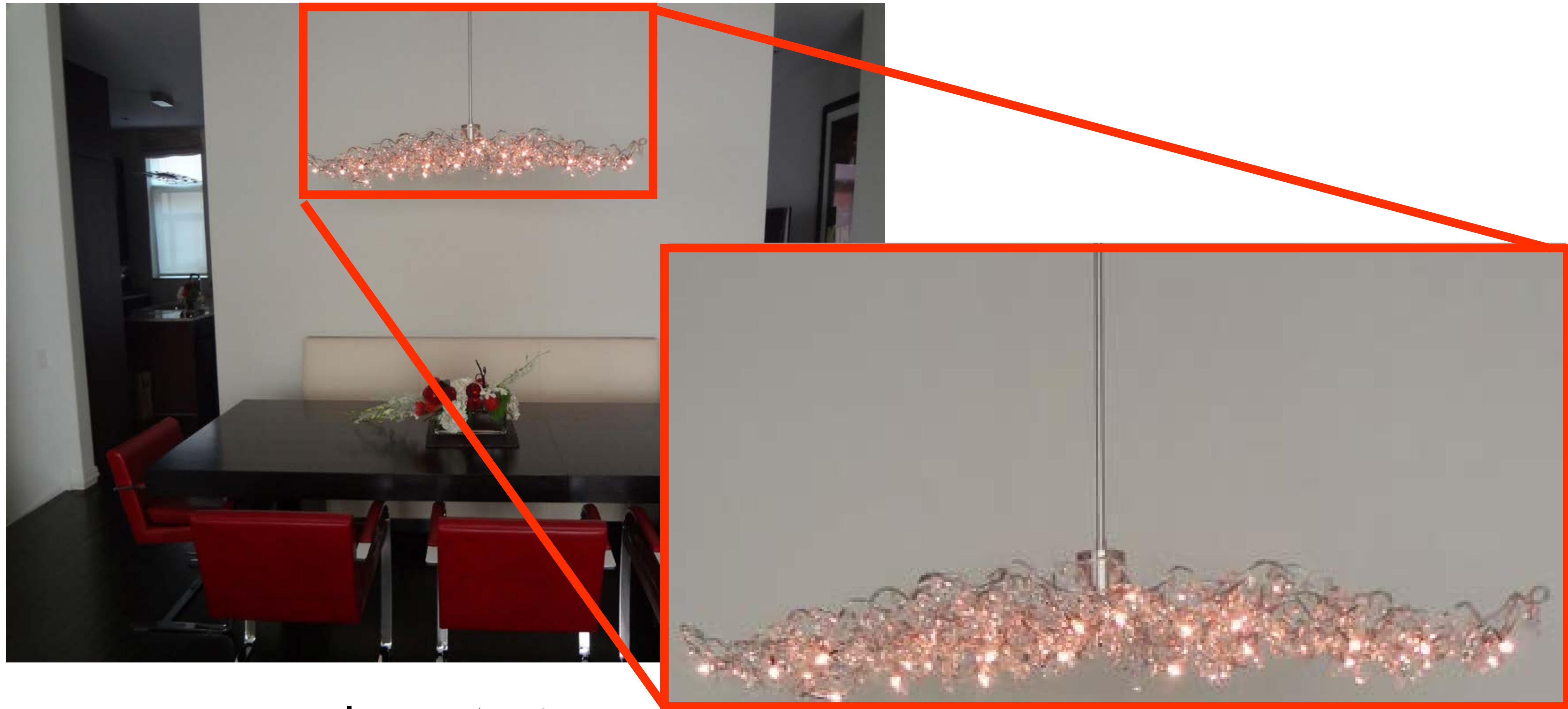
Iconic



Top 4 results:



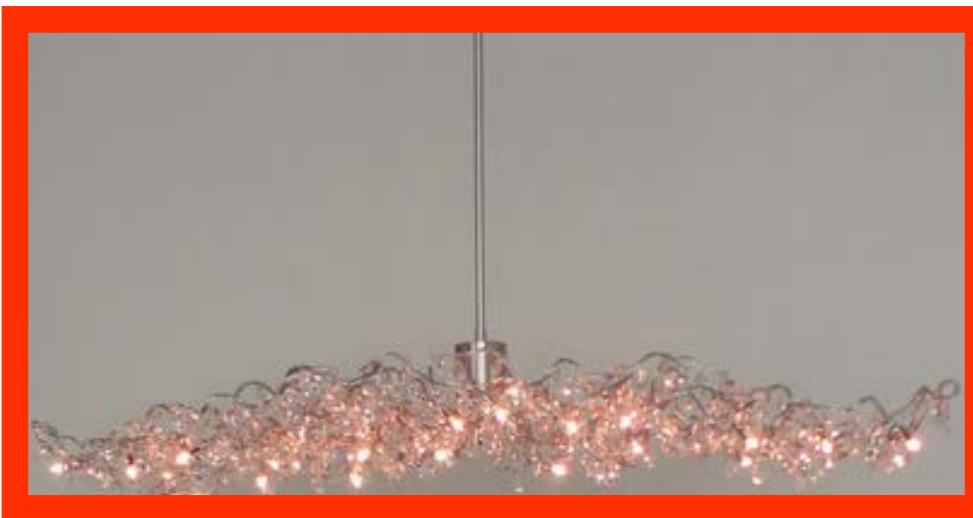
# RESULTS: “WHAT IS IT?”



In context

# RESULTS: “WHAT IS IT?”

In context



Iconic

Tiara Oval Suspension by  
Hanco Loor

Dining Room   Contemporary

1 View jpg

Wall Sconces   Unknown style

1 View jpg

2 boxes

Top 4 results:

HARCO LOOR Tiara  
Chandelier

Chandeliers   Contemporary

1 View jpg

View detail

0.035

Argent N92S Suspension  
Light

Bathroom Vanity Lighting

Modern

View jpg

View detail

0.038

Eurofase 25620-016 Divo 9  
Light Pendant in Nickel  
25620-016

Pendant Lighting

Modern

View jpg

View detail

0.041

Eurofase 25618-013 Divo 6  
Light Pendant in Nickel  
25618-013

Pendant Lighting

Modern

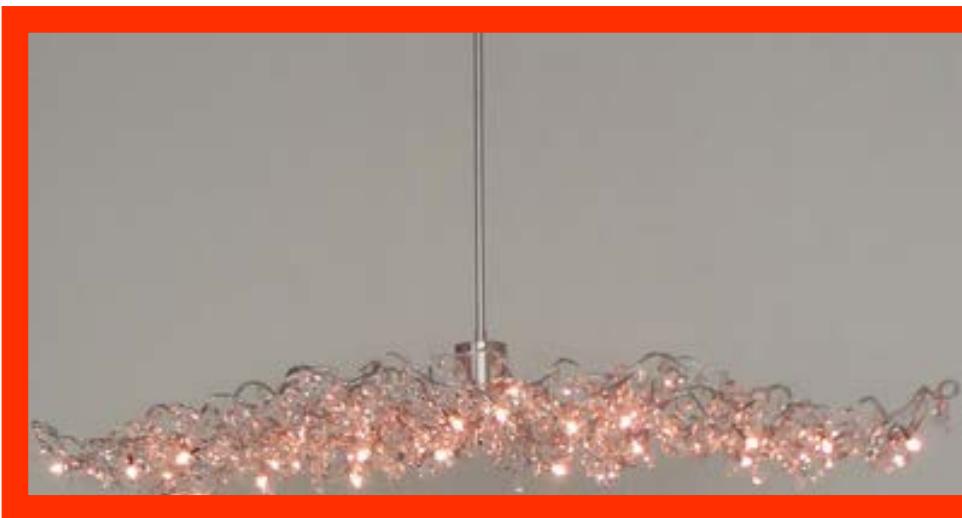
View jpg

View detail

0.047

# COMPARISON: TRAINED ONLY ON CATEGORIES

In context



Iconic

A wireframe or schematic representation of the chandelier, showing its basic structure and light fixtures.

Tiara Oval Suspension by Harco Loor

Dining Room Contemporary

1 View jpg

Wall Sconces Unknown style

1 View jpg

2 boxes

Top 4 results:

A rectangular chandelier with a black frame and hanging crystals.

'Apollon' Black Shaded 6-light Crystal Chandelier

Chandeliers Contemporary

View jpg View detail

1792.880

A chandelier made of deer antlers and lit candles.

Authentic Deer and Elk Antler Banquet Table Chandelier

Chandeliers Unknown style

View jpg View detail

1895.330

A large, oval-shaped chandelier with many hanging crystals.

Harco Loor | Tiara Oval HL 15 Suspension Light

Pendant Lighting Modern

View jpg View detail

1911.520

A hexagonal chandelier with a complex, geometric frame and hanging crystals.

Vermeer Hexagonal Pendant

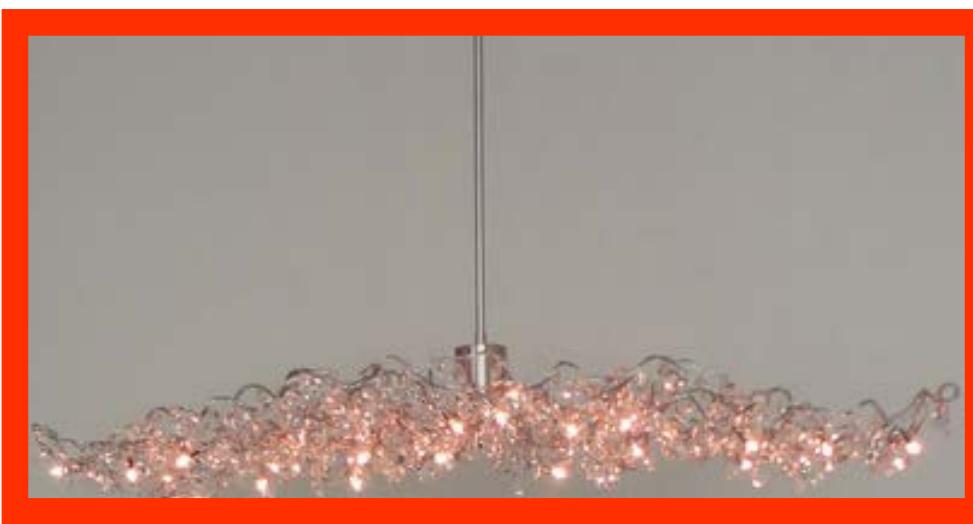
Pendant Lighting Modern

View jpg View detail

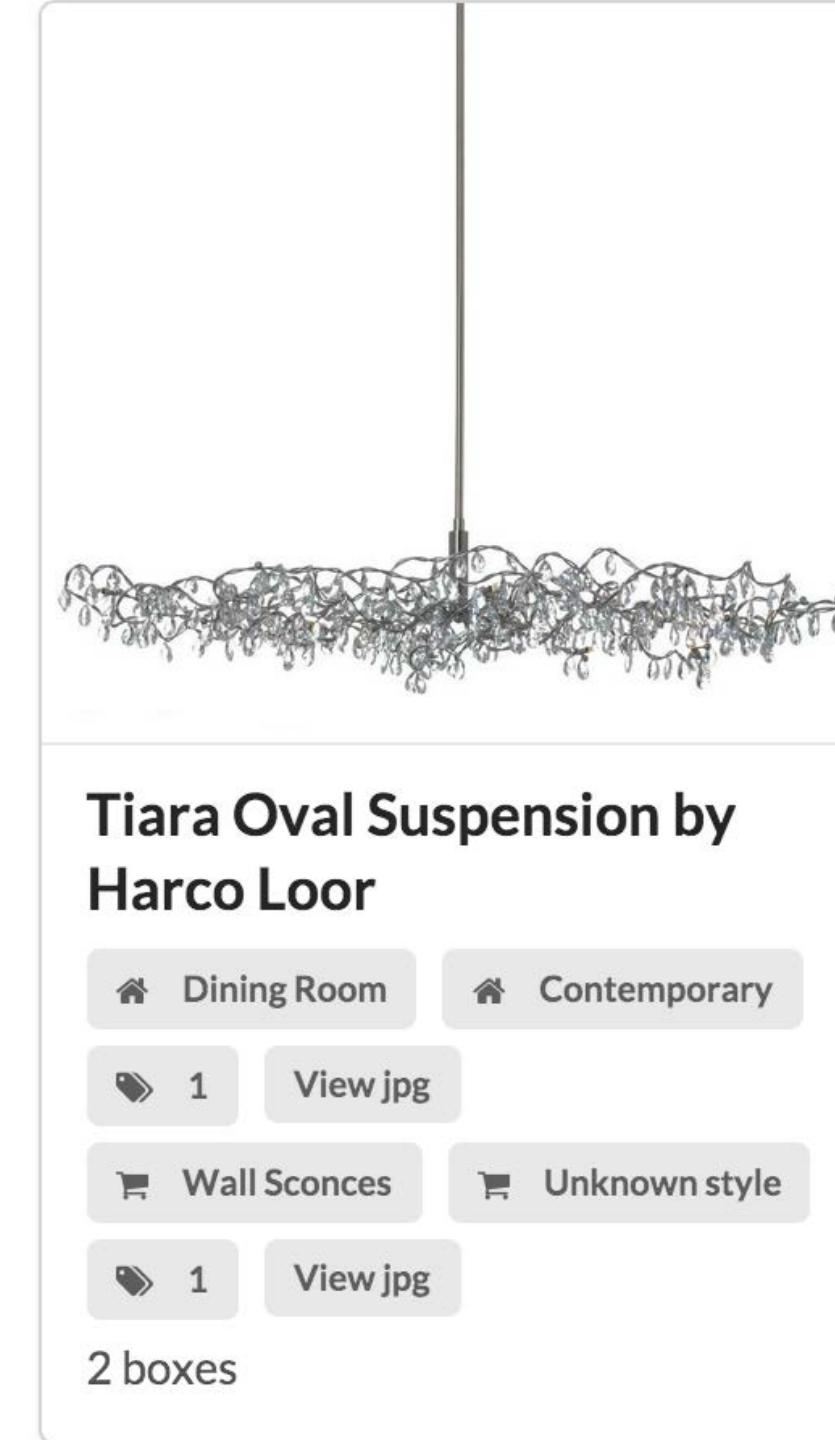
2035.000

# COMPARISON: TRAINED ONLY ON IMAGENET

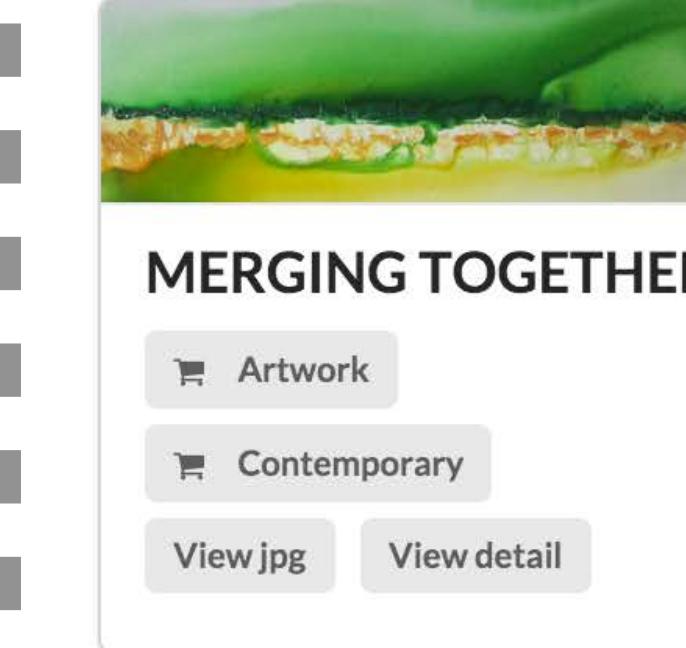
In context



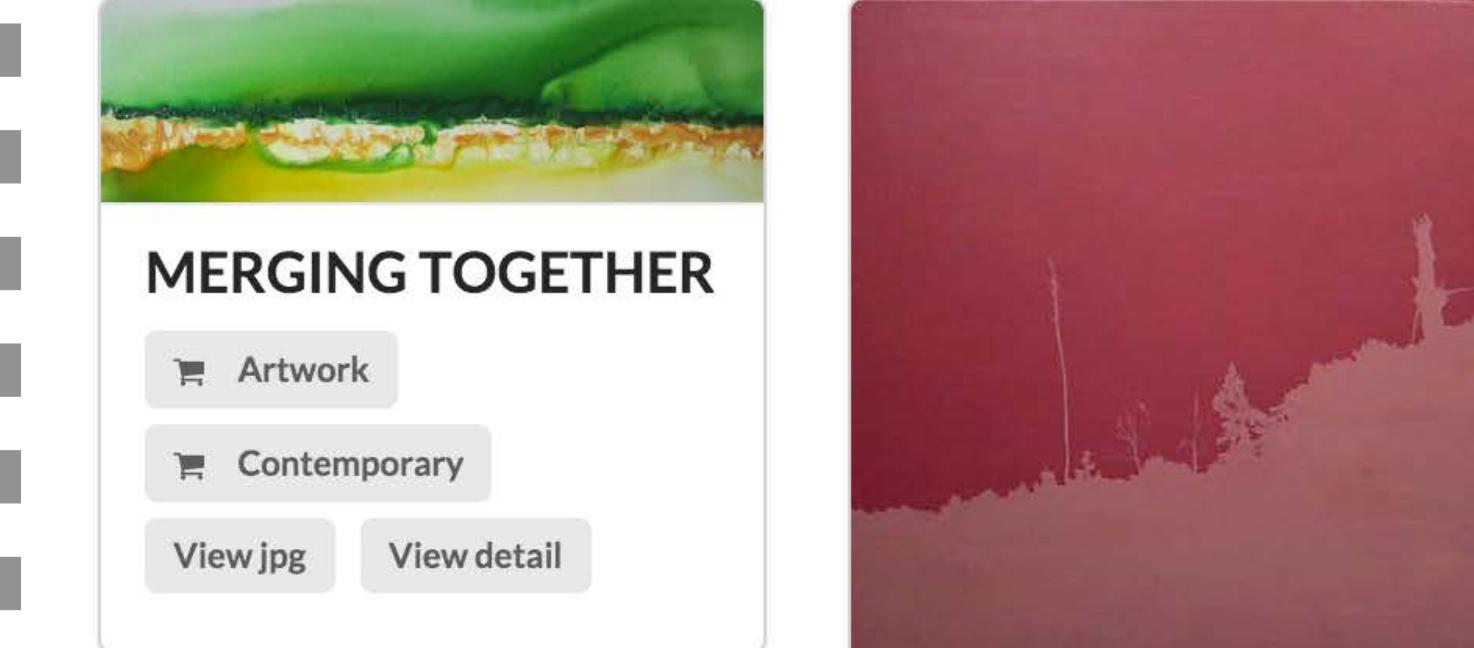
Iconic



Top 4 results:

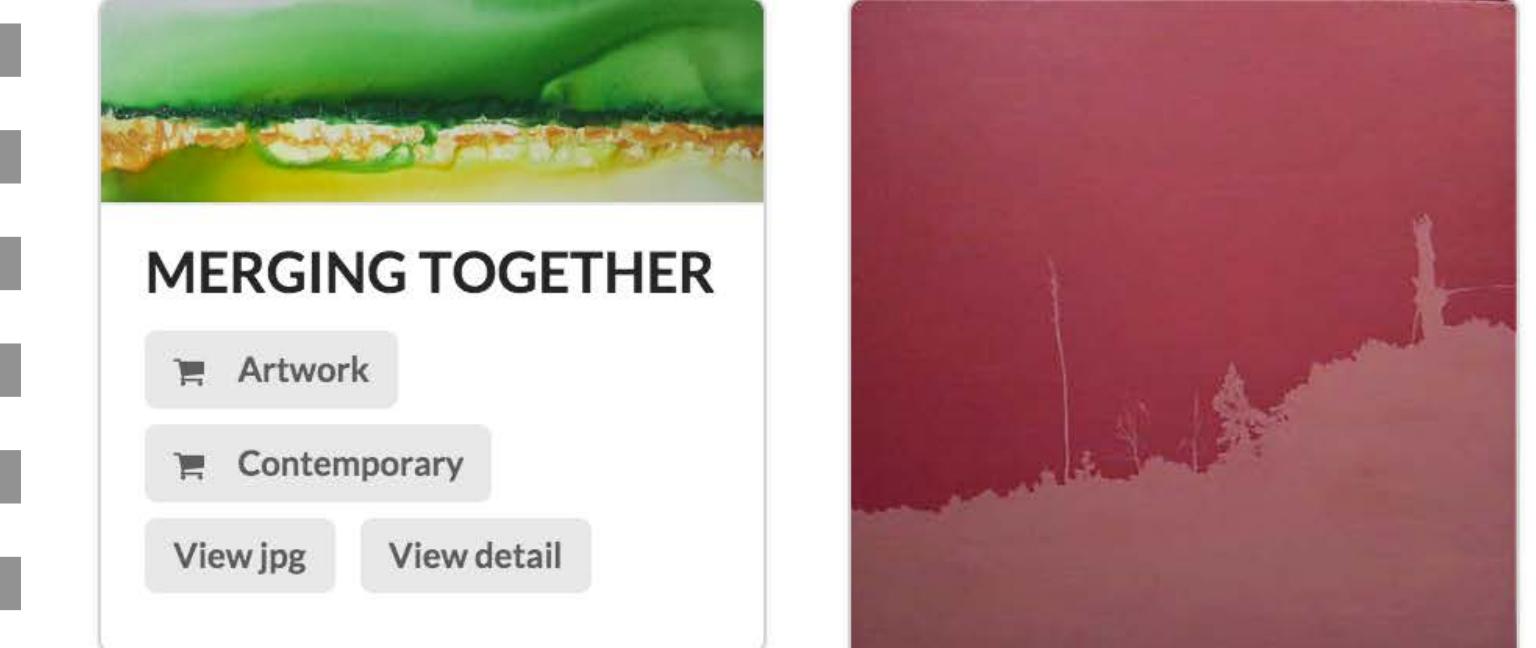


498.493



"Elephant View"  
Artwork

499.135



Terzani Argent N92S  
Chandelier

503.232



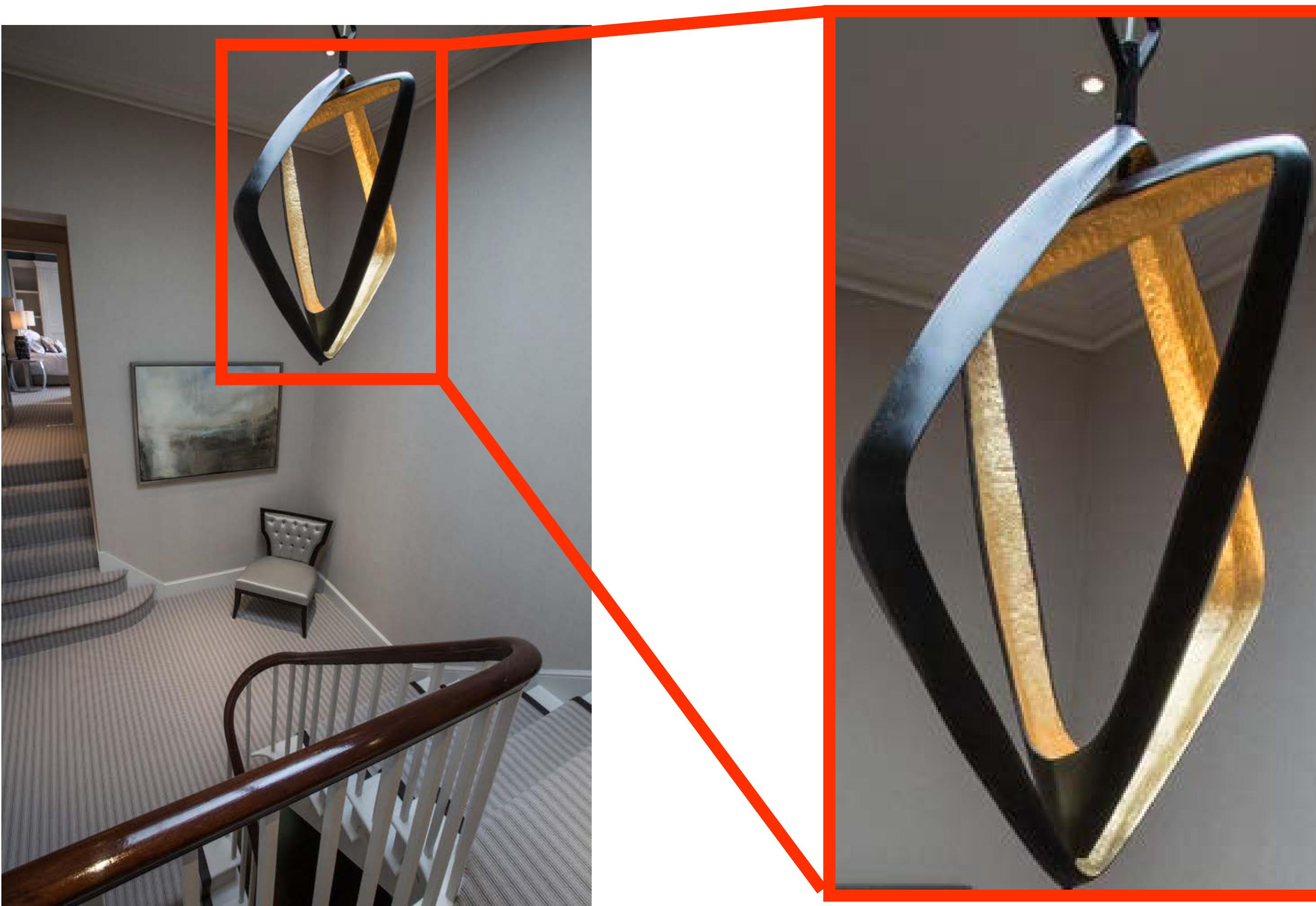
Mensa Hanging Light  
Fixture

Chandeliers  
Contemporary  
View jpg View detail

514.867



# RESULTS: FAILURE CASE



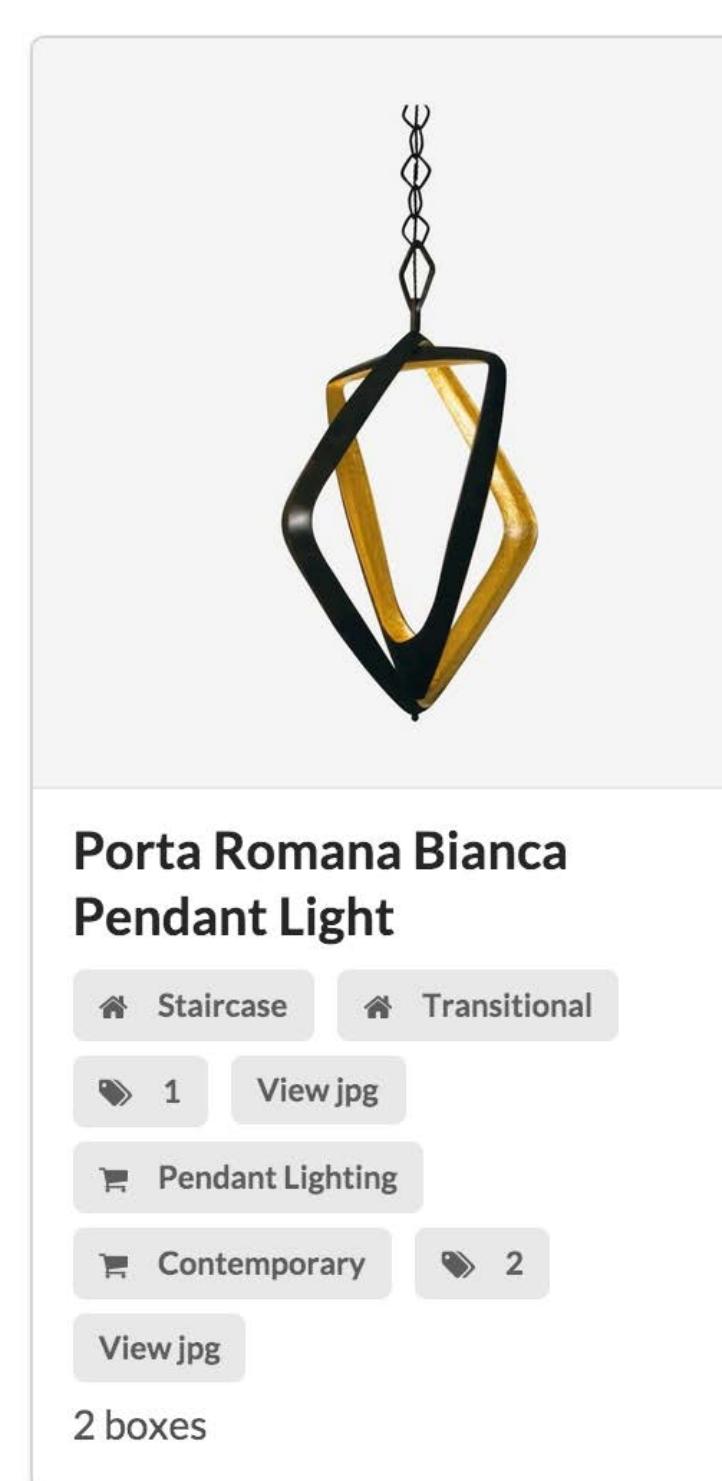
In context

# RESULTS: FAILURE CASE

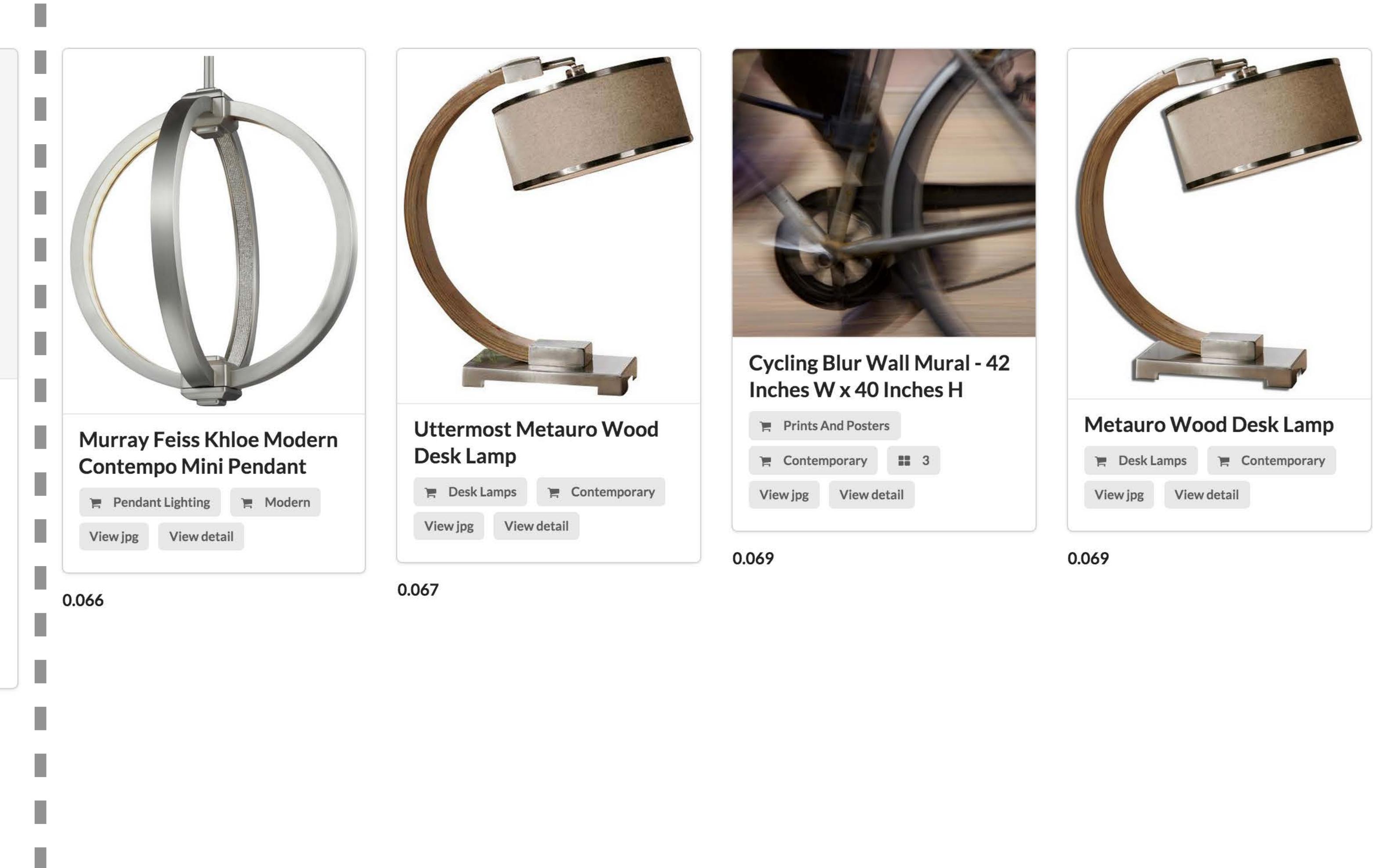
In context



Iconic



Top 4 results:



# RESULTS: “WHERE IS IT USED?”



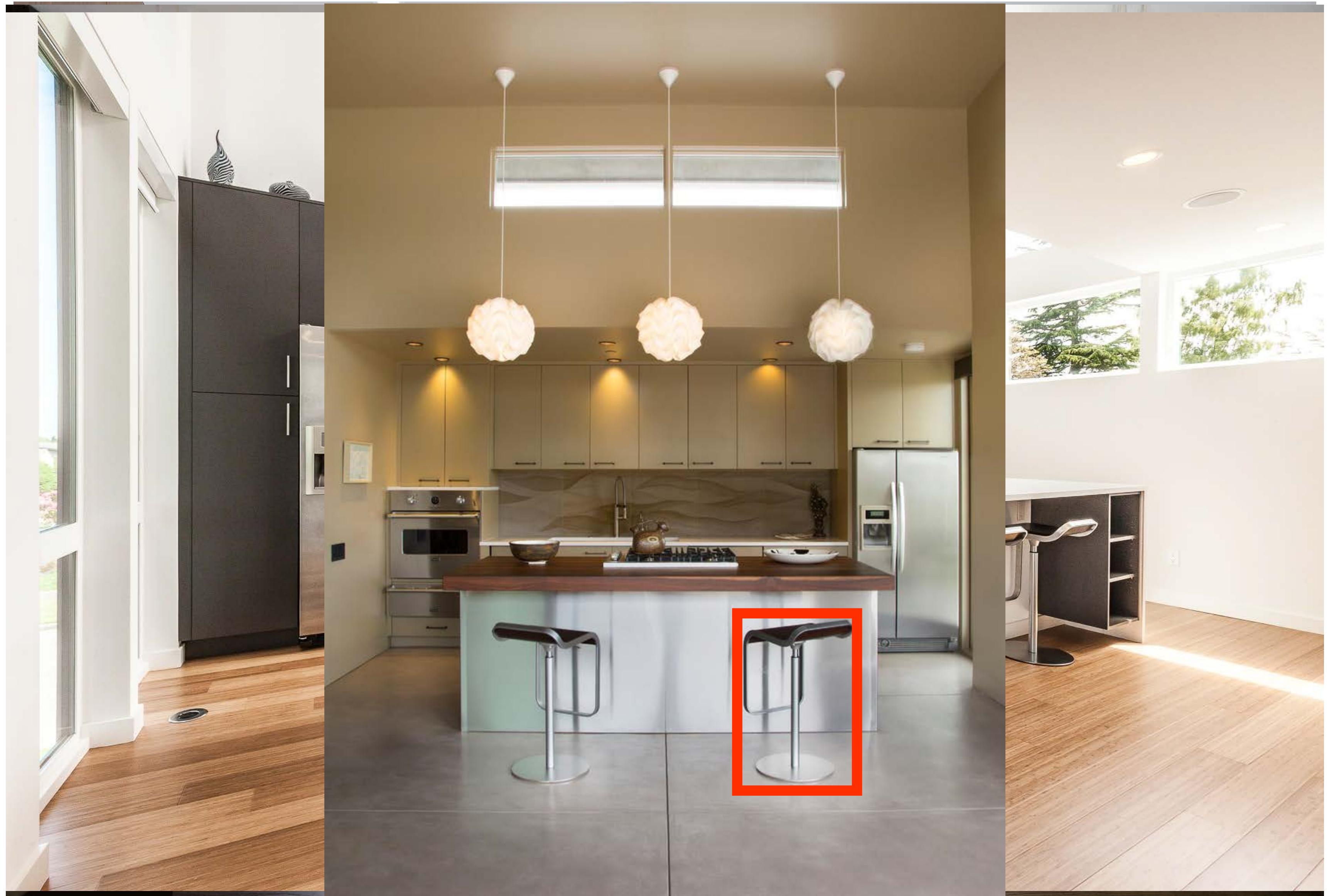
“Maskros  
Pendant  
Lamp”



# RESULTS: “WHERE IS IT USED?”



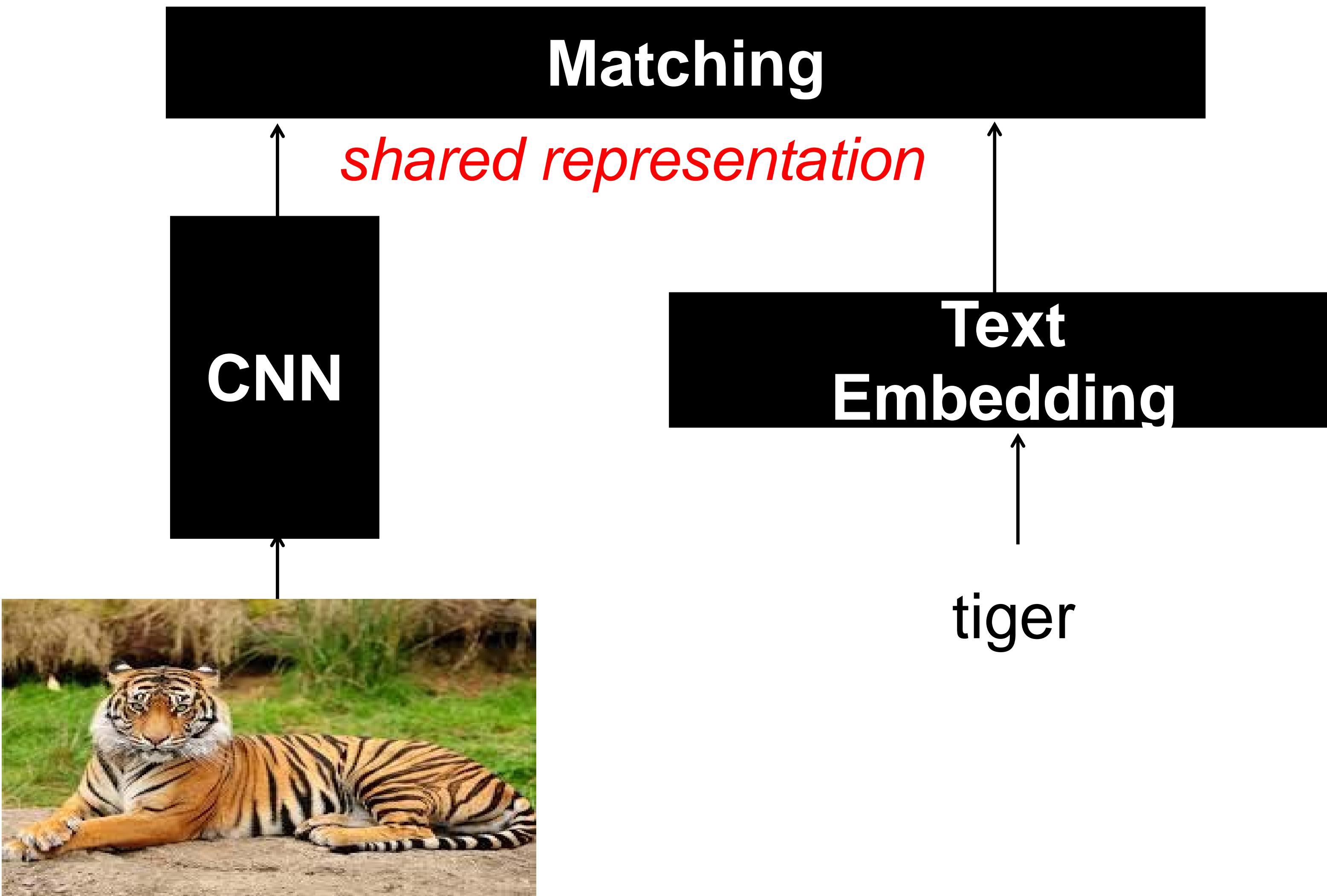
"LEM Piston  
Stool | Design  
Within Reach"



# SEARCHING ACROSS CATEGORIES

Query $I_q$	Top-1 nearest neighbor from different object categories											
	Dining chairs	Armchairs	Rocking chairs	Bar stools	Table lamps	Outdoor lighting	Bookcases	Coffee tables	Side tables	Floor lamps	Rugs	Wallpaper

# Fancier Architectures: Multi-Modal



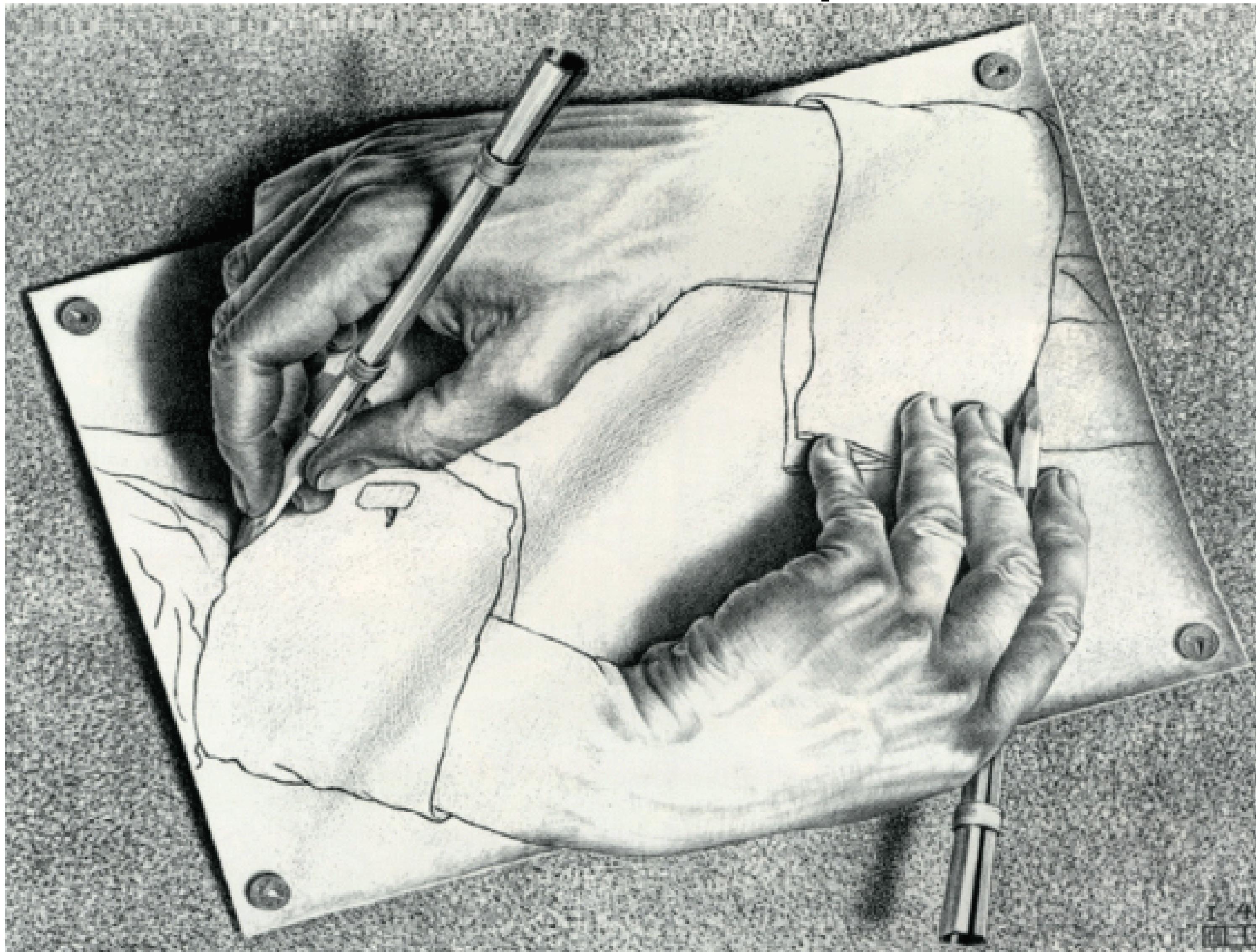
# Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk, but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would

Deep  
Net

# Self-supervised Learning: data as its own supervision



# Context Prediction for Images

?

?

?

?

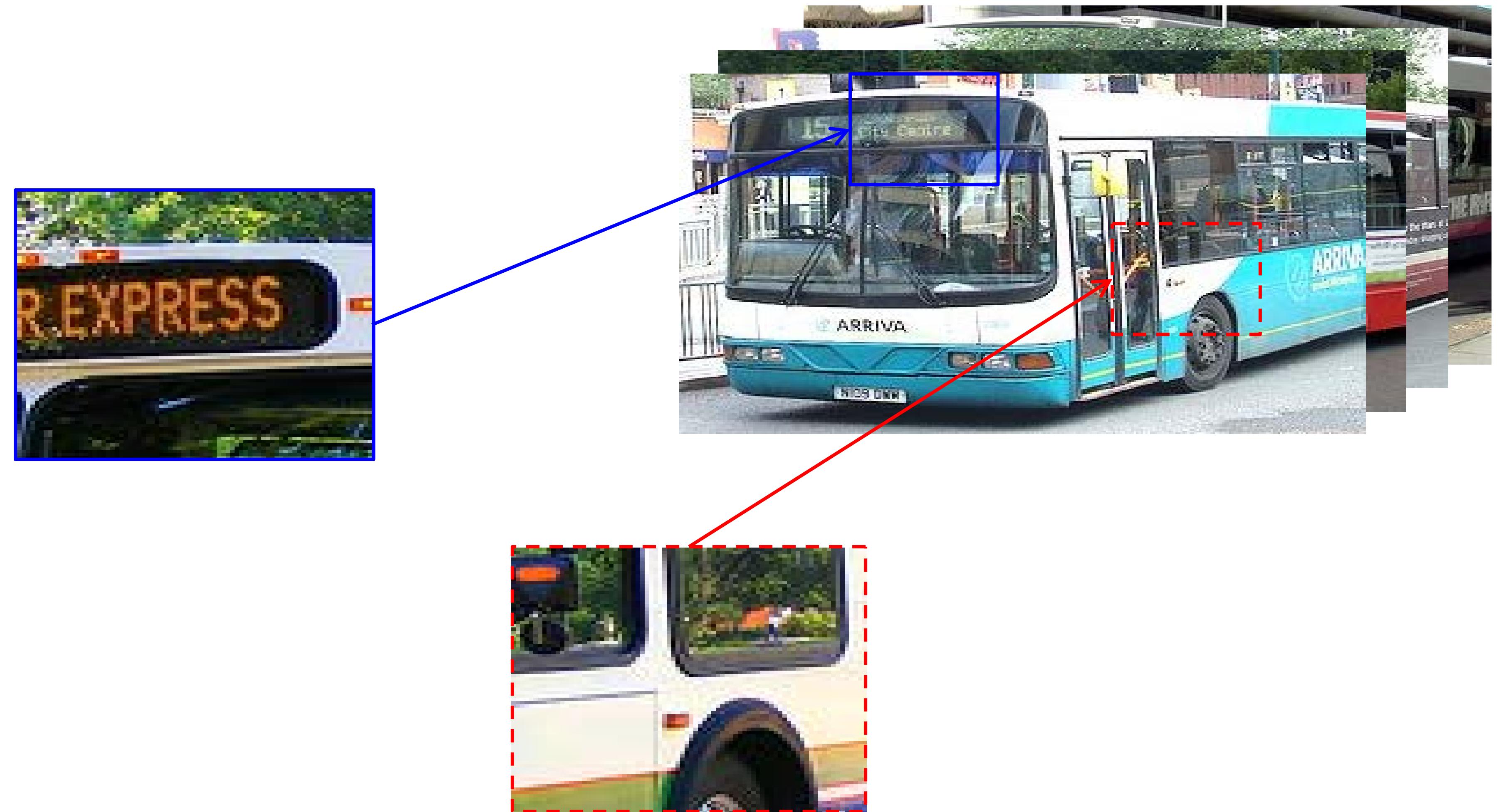


?

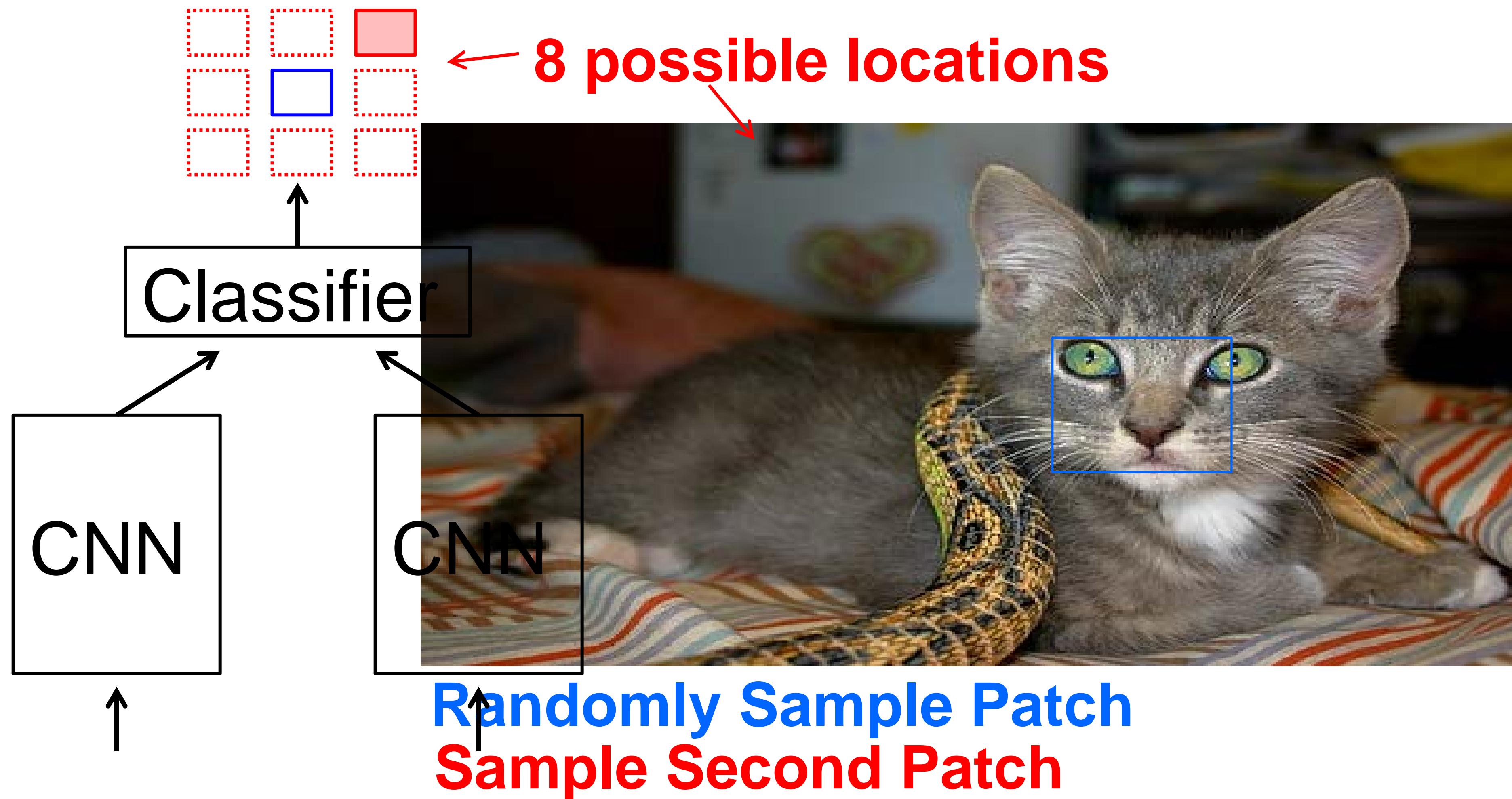
A

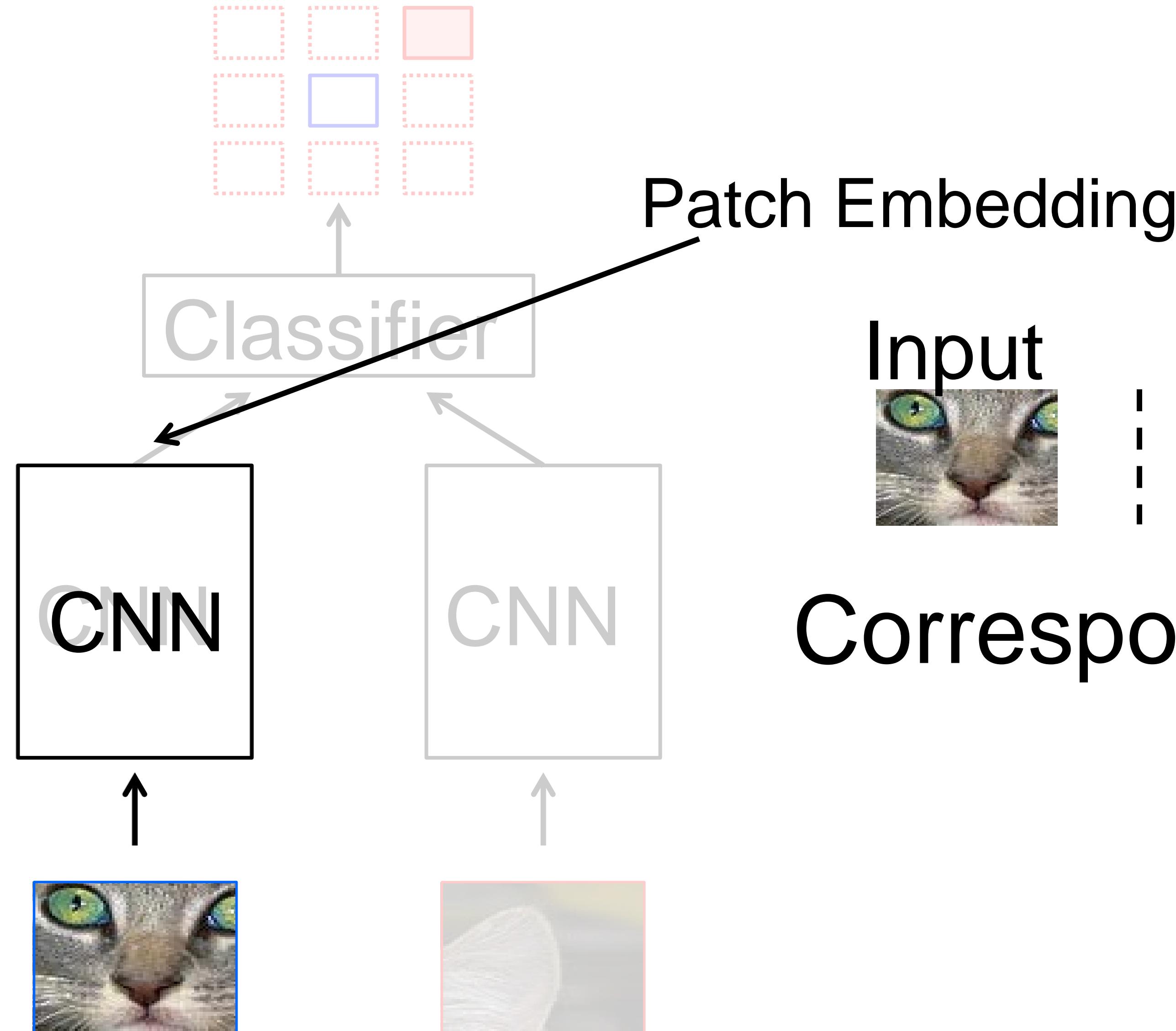
B

# Semantics from a non-semantic task



# Relative Position Task





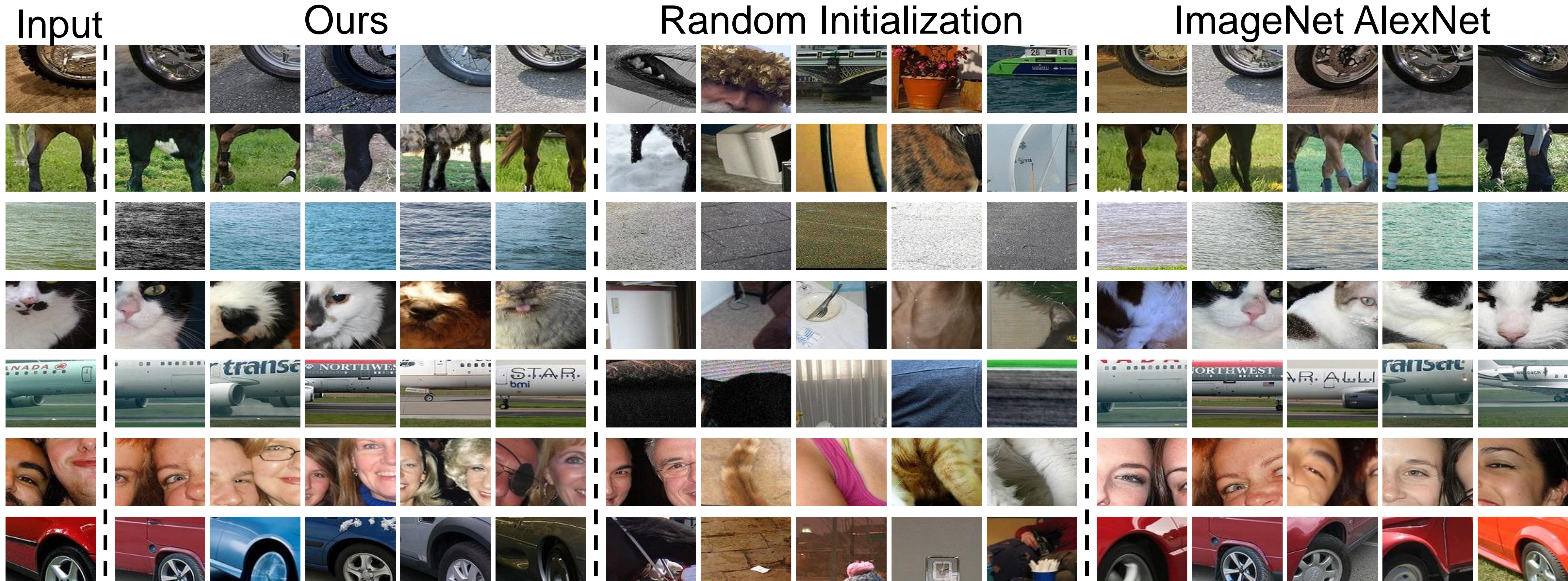
Patch Embedding

Input

Nearest Neighbors

Correspondence ***across*** instances!

# What is learned?



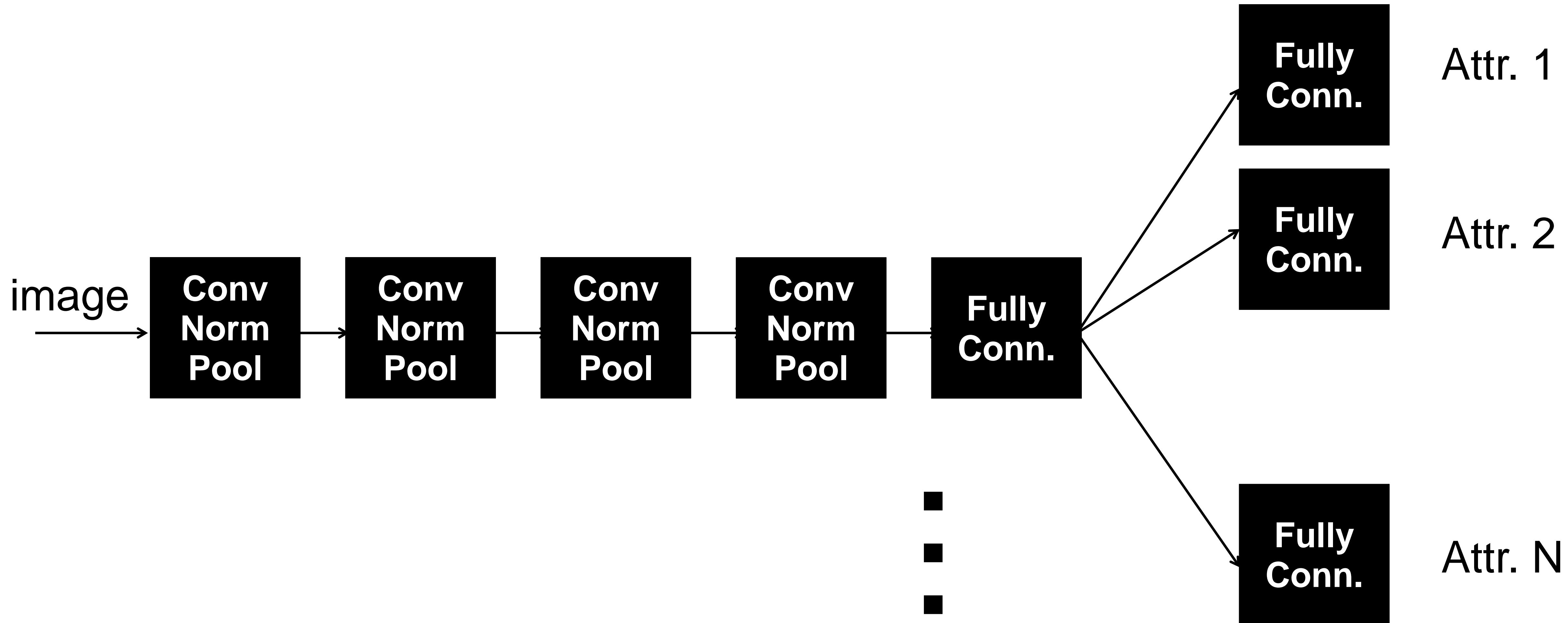
# Still don't capture everything



You don't always need to learn!



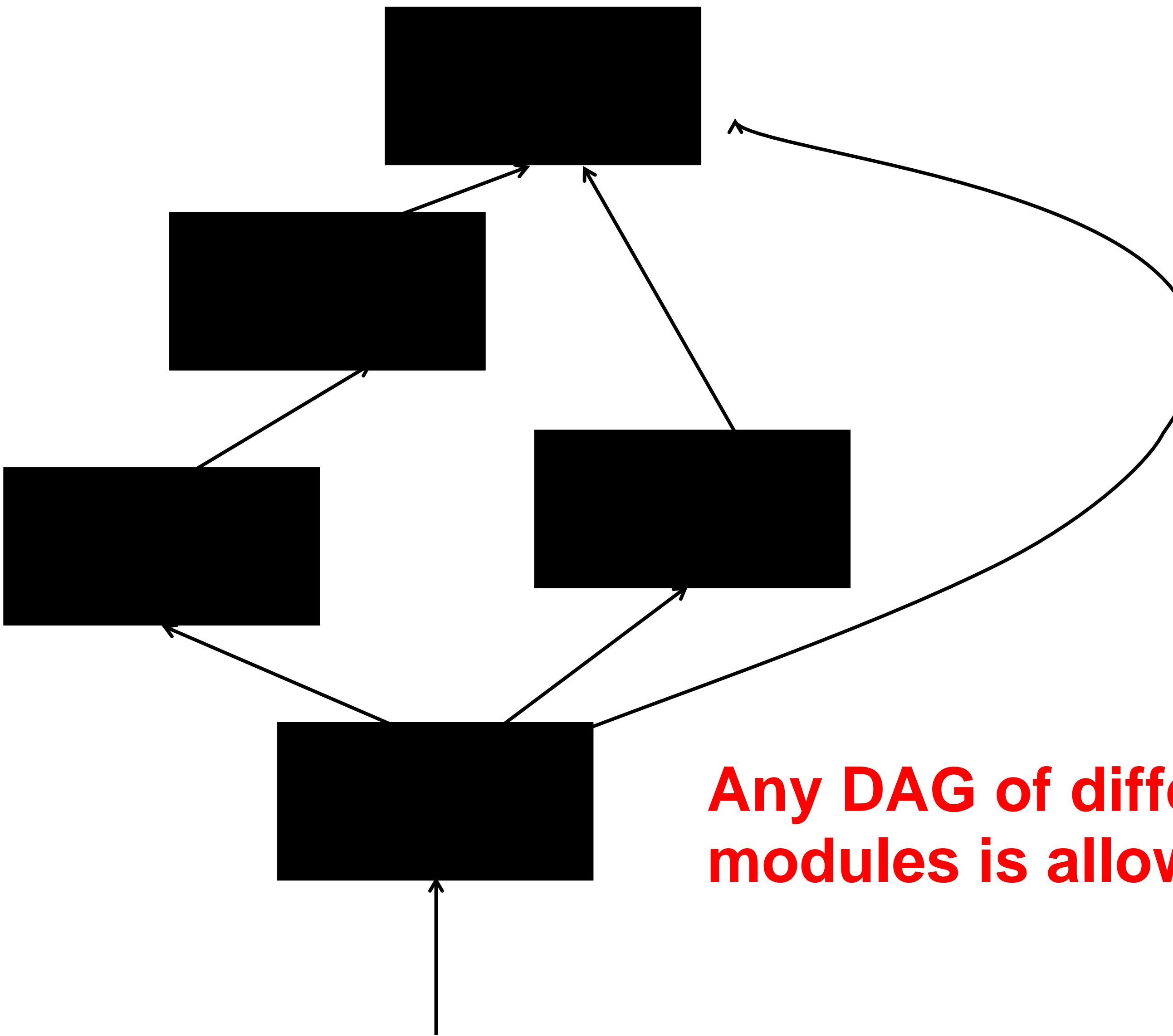
# Fancier Architectures: Multi-Task



# Decompose the image into parts

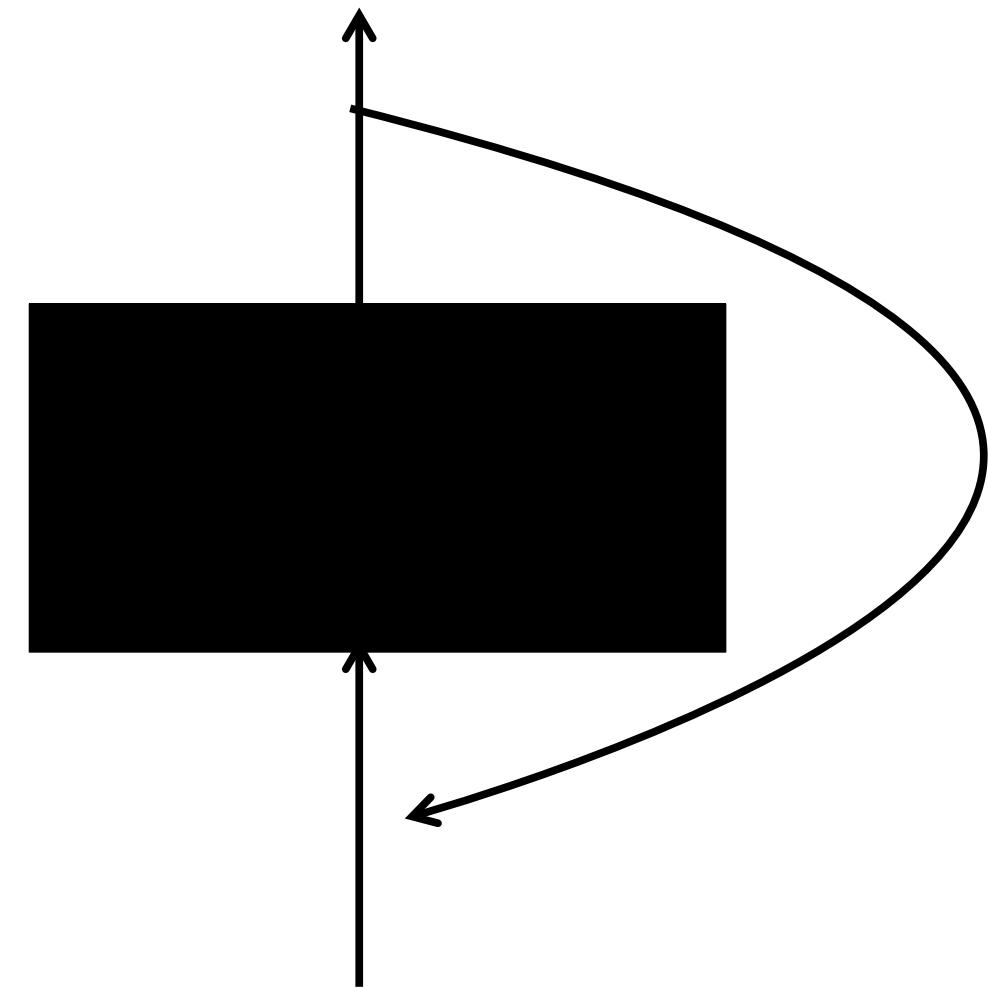


# Fancier Architectures: Generic DAG



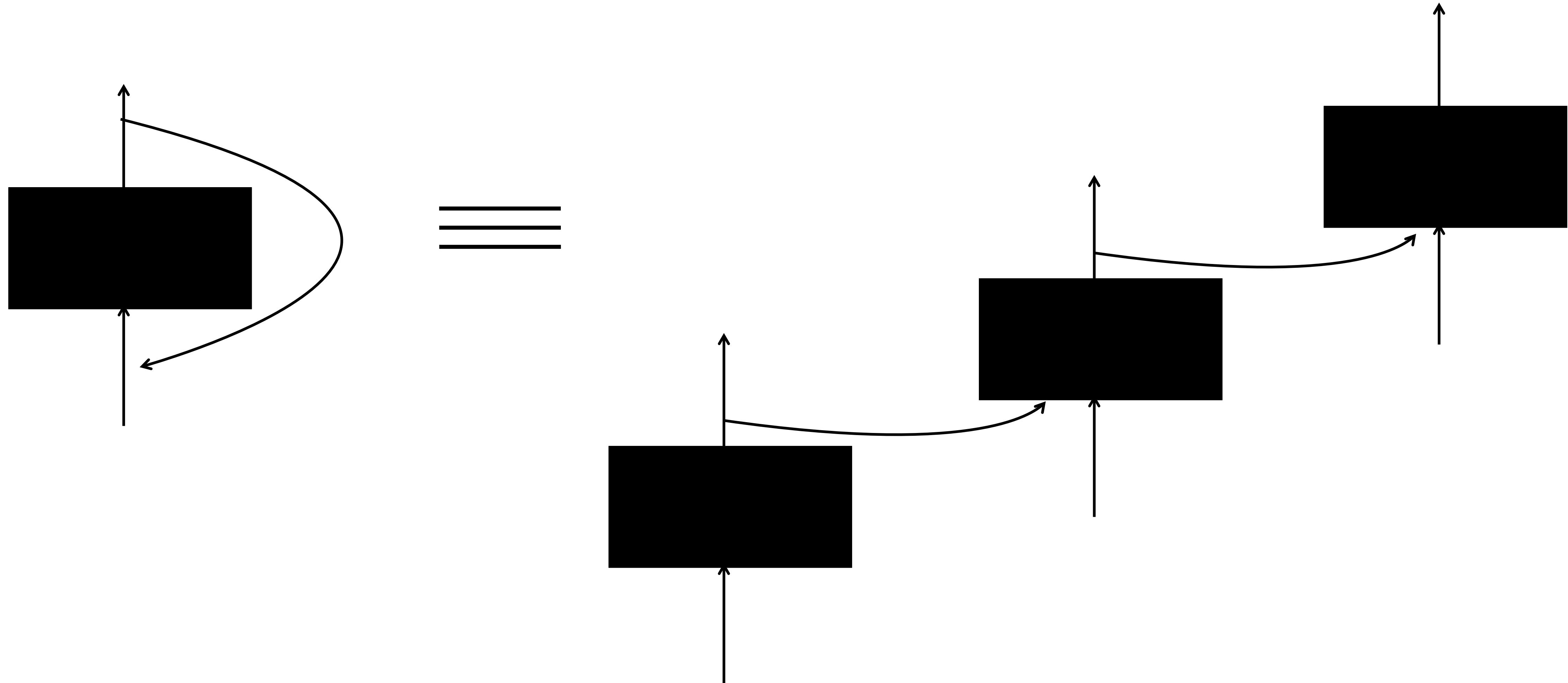
**Any DAG of differentiable modules is allowed!**

# Fancier Architectures: Generic DAG



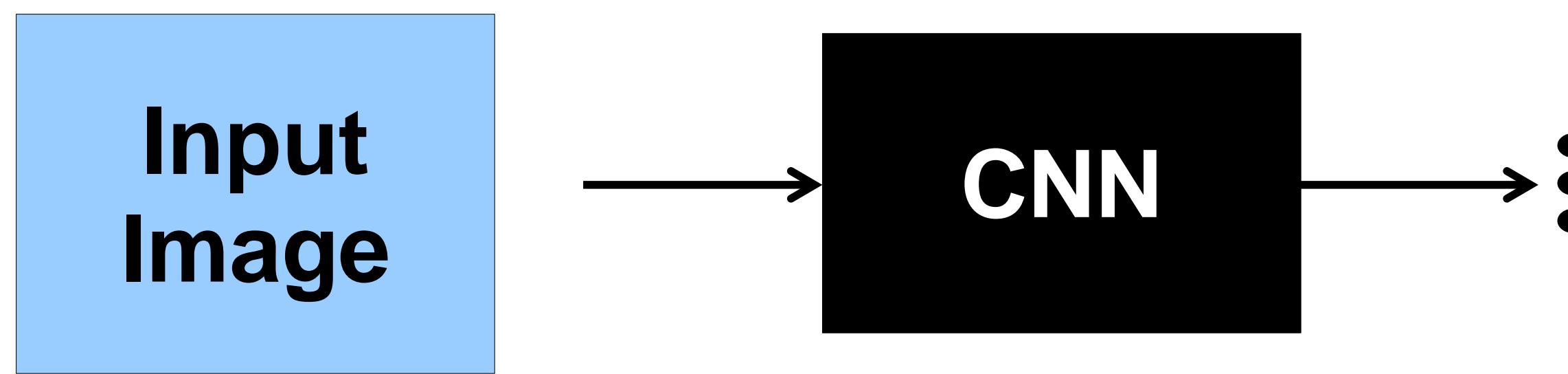
# Fancier Architectures: Recurrent NN

If there are cycles, one needs to un-roll it.



# Fancier Architectures: Fully-Convolutional

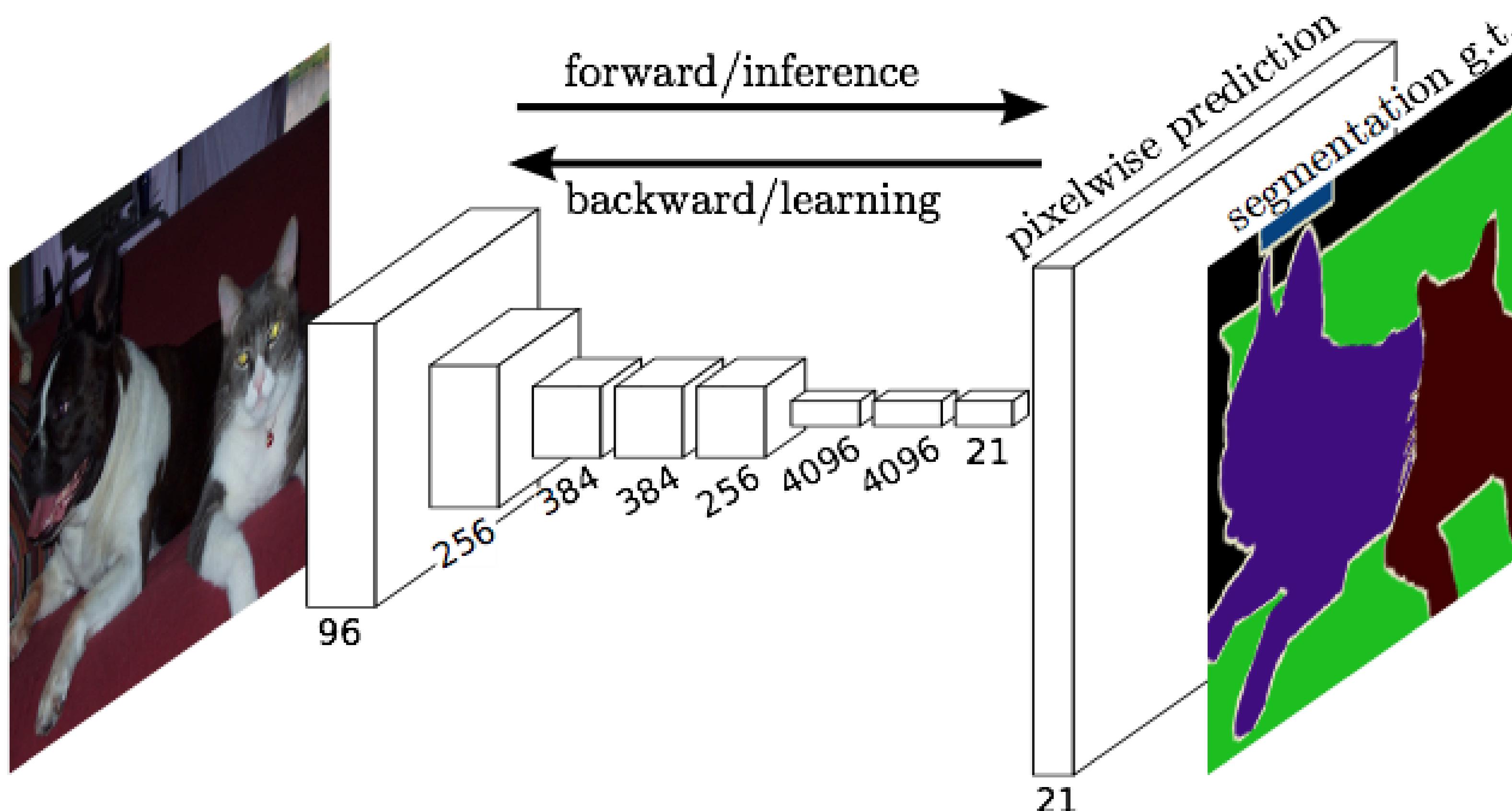
Classification ConvNet (e.g. ImageNet)



Per-pixel Labeling ConvNet

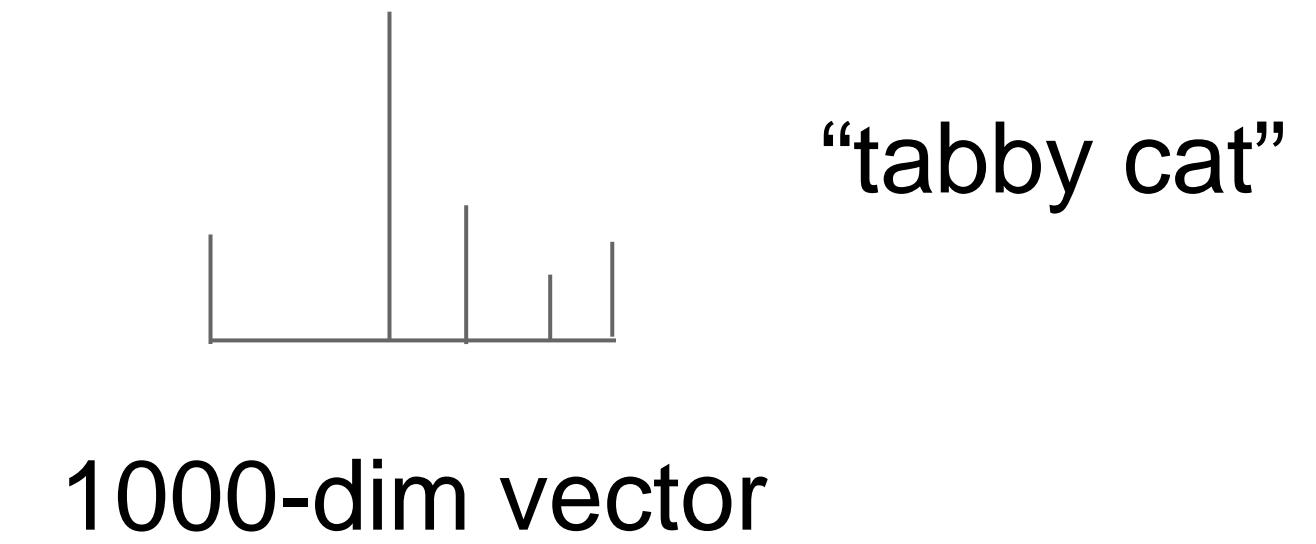
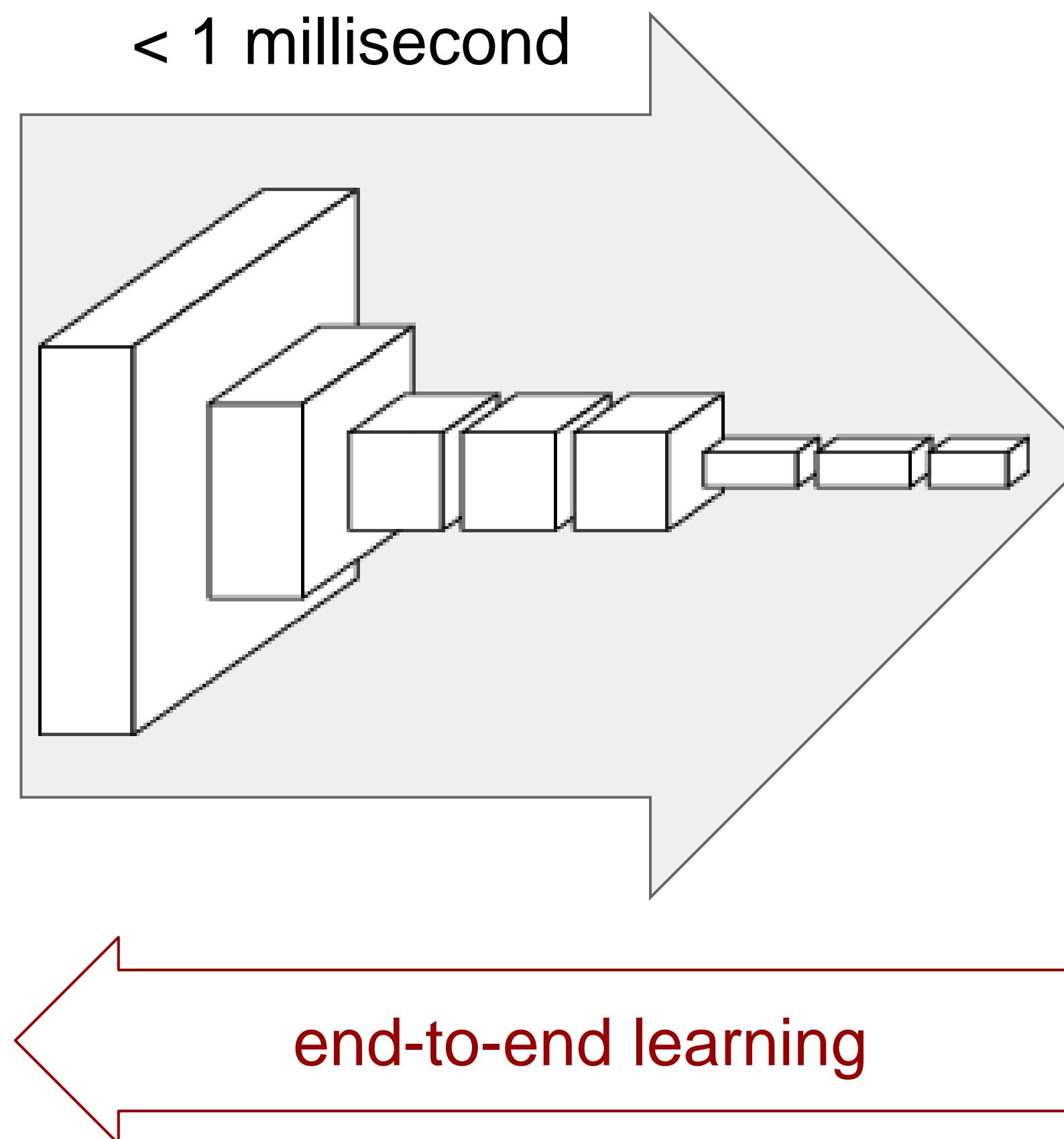


# Fully Convolutional Networks for Semantic Segmentation

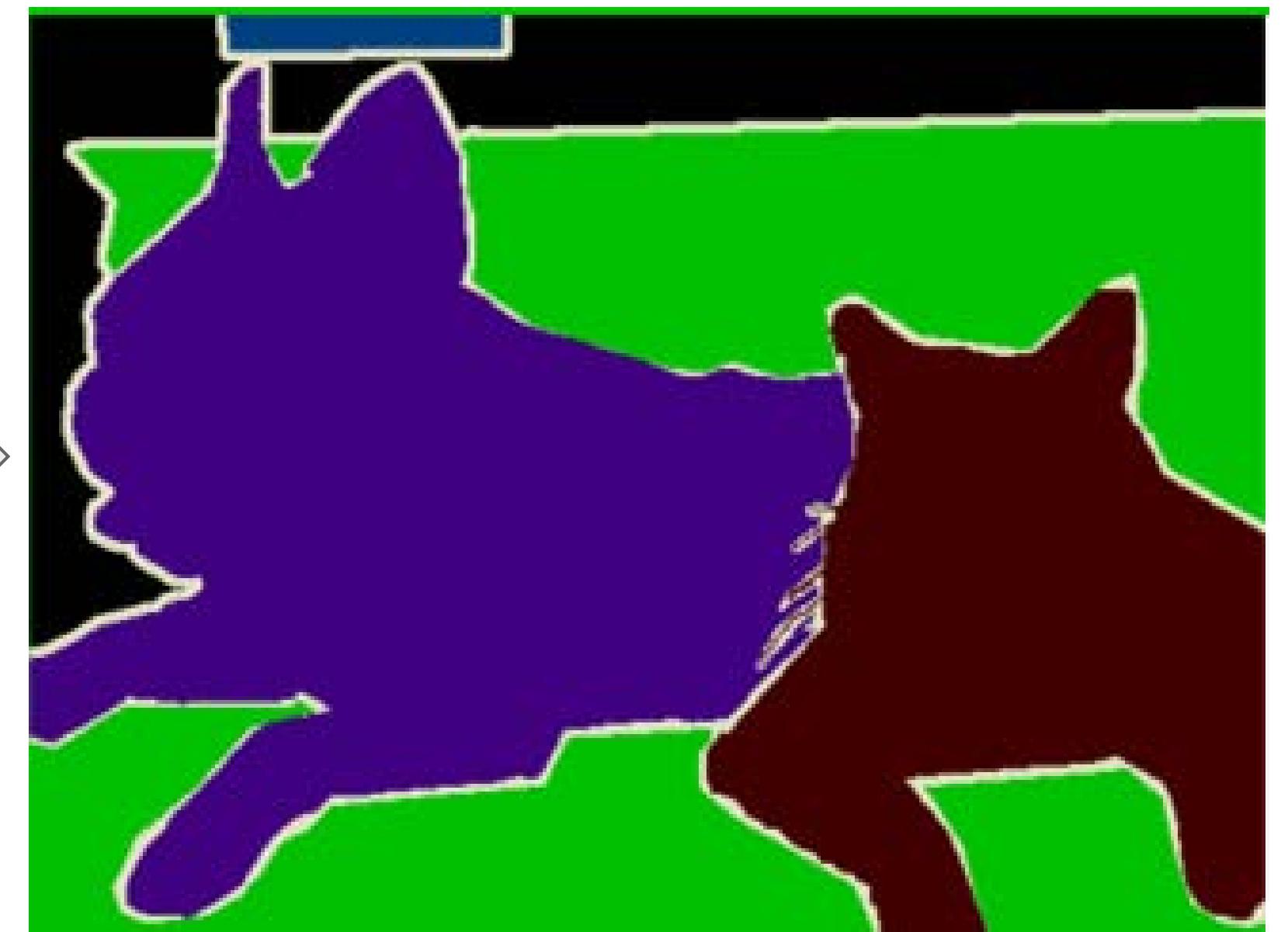
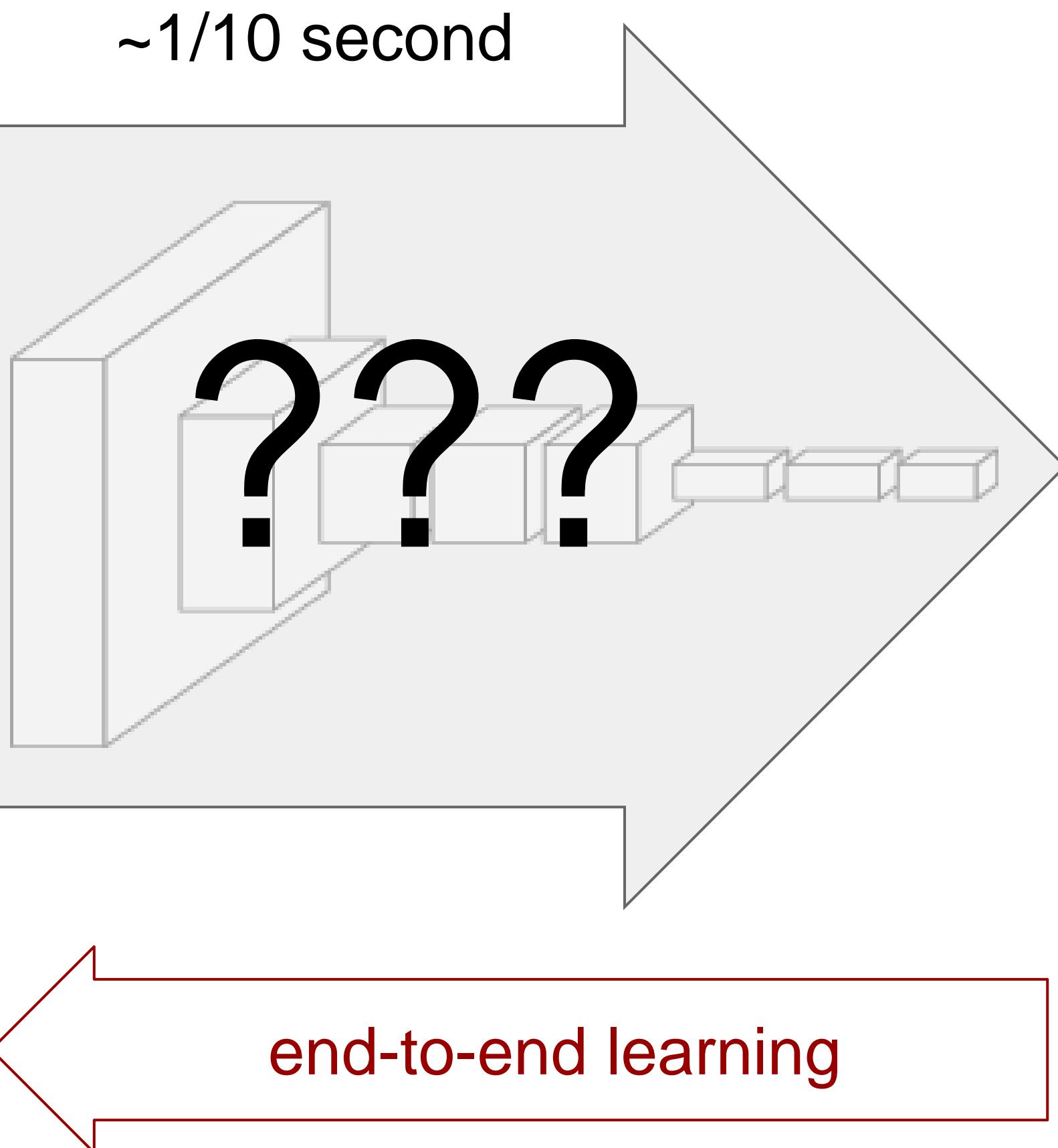


Evan Shelhamer\* Jonathan Long\* Trevor Darrell  
UC Berkeley in CVPR'15, PAMI'16

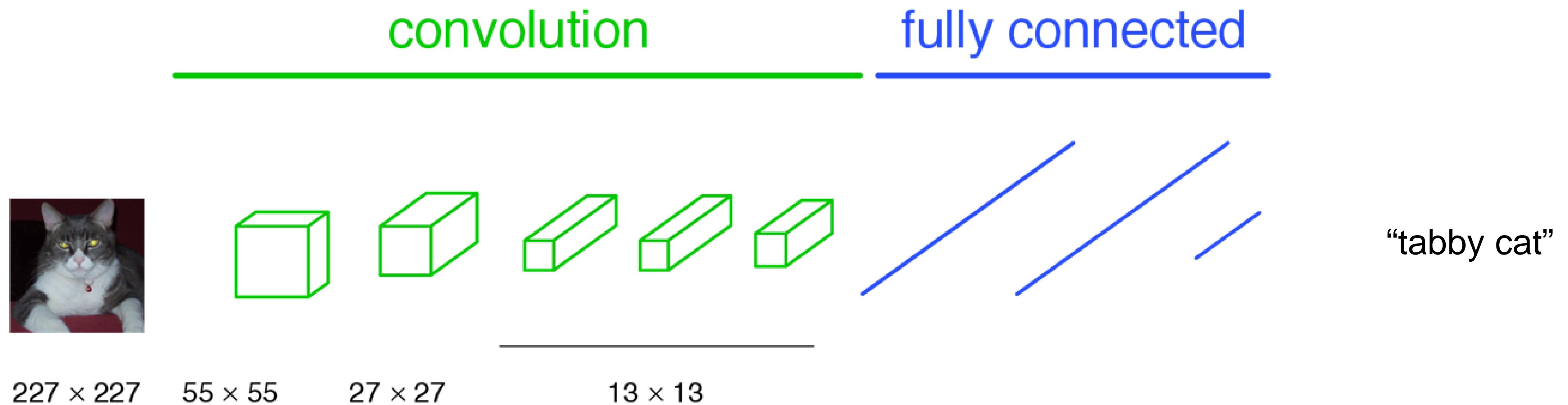
## convnets perform classification



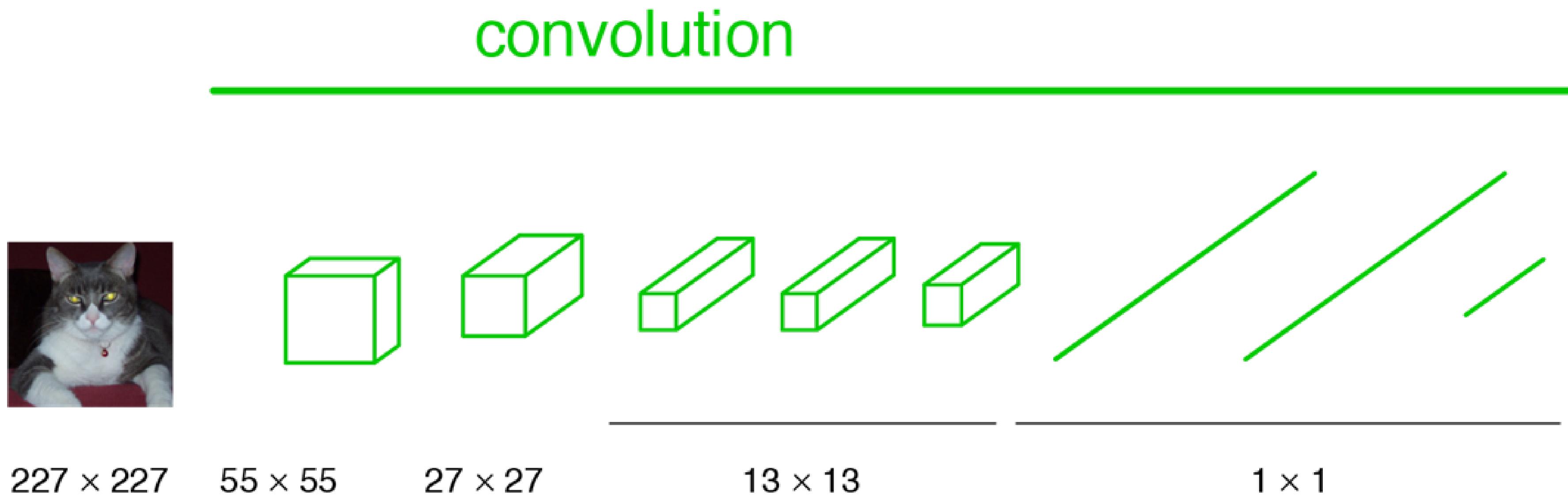
**lots of pixels, little time?**



## a classification network



## becoming fully convolutional



## becoming fully convolutional



convolution

---

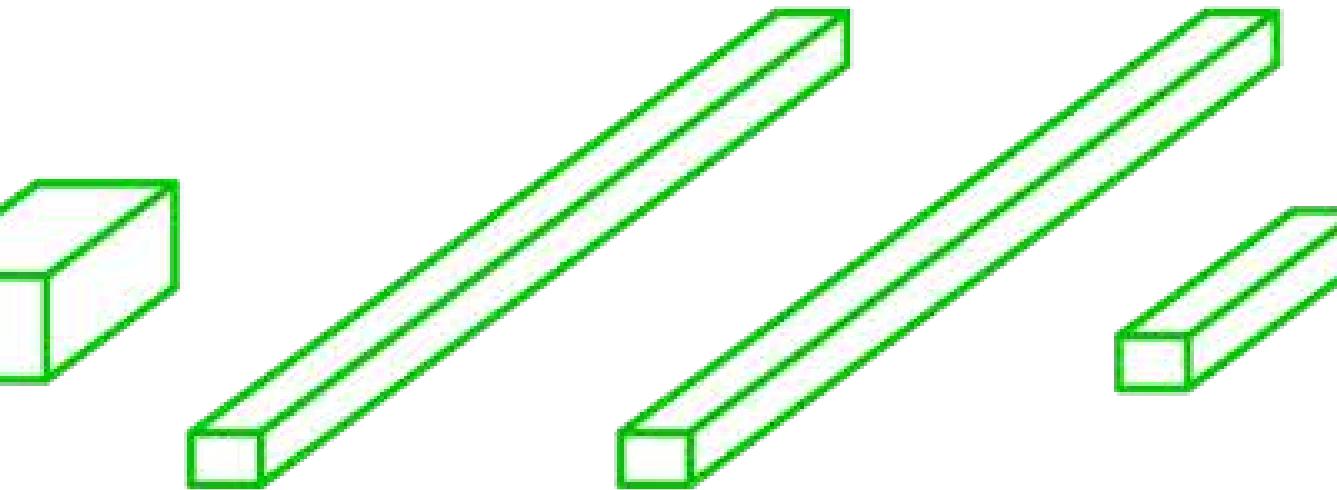
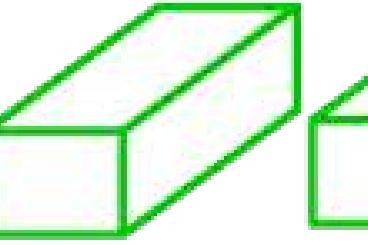
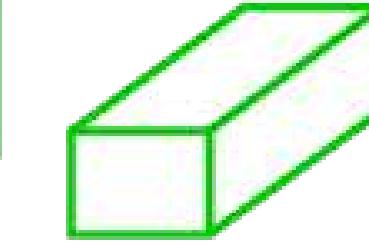
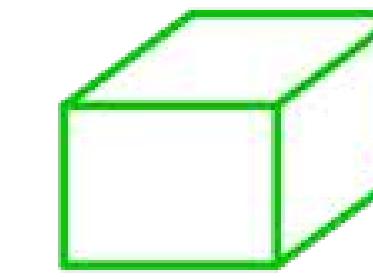
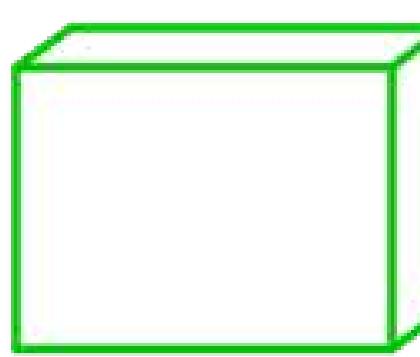
$H \times W$

$H/4 \times W/4$

$H/8 \times W/8$

$H/16 \times W/16$

$H/32 \times W/32$



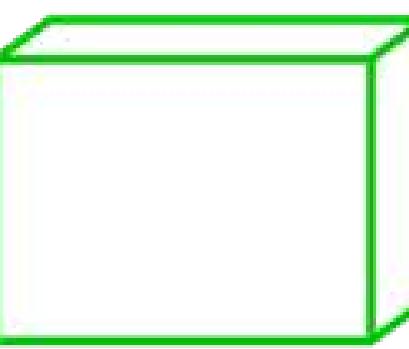
**upsampling output**

convolution

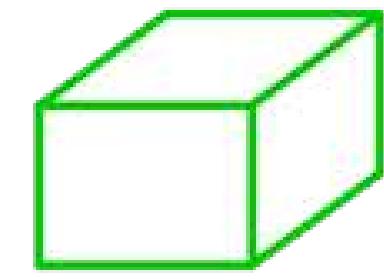
---



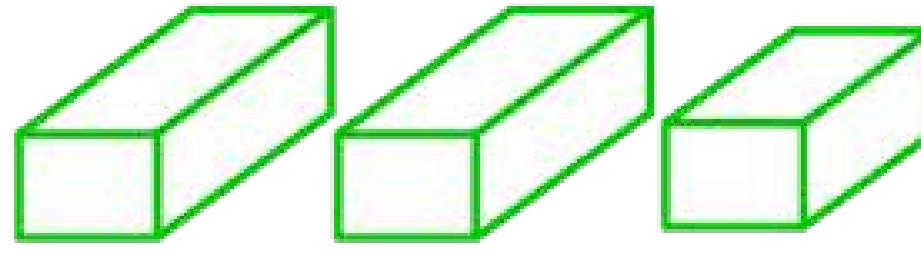
$H \times W$



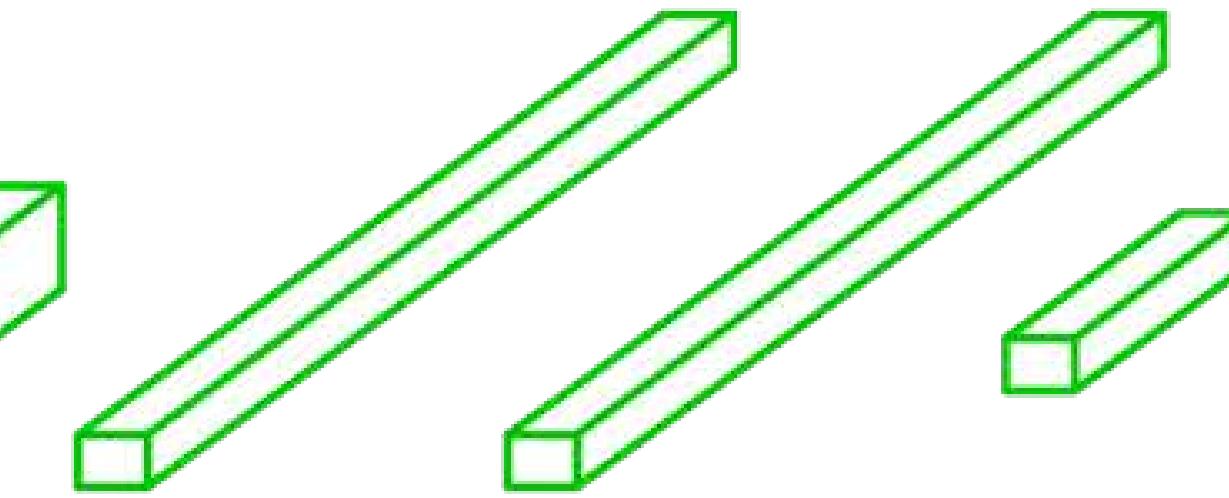
$H/4 \times W/4$



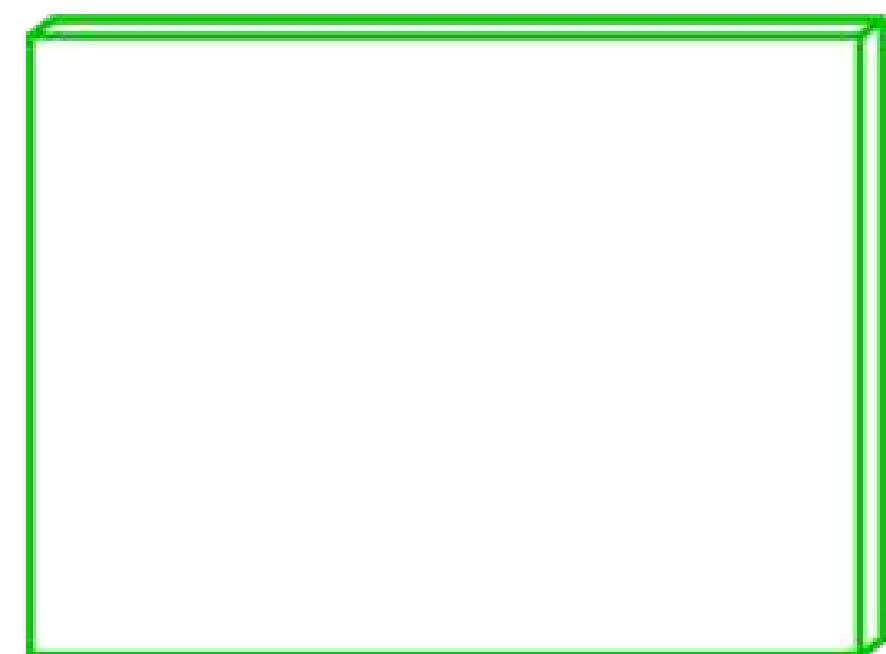
$H/8 \times W/8$



$H/16 \times W/16$

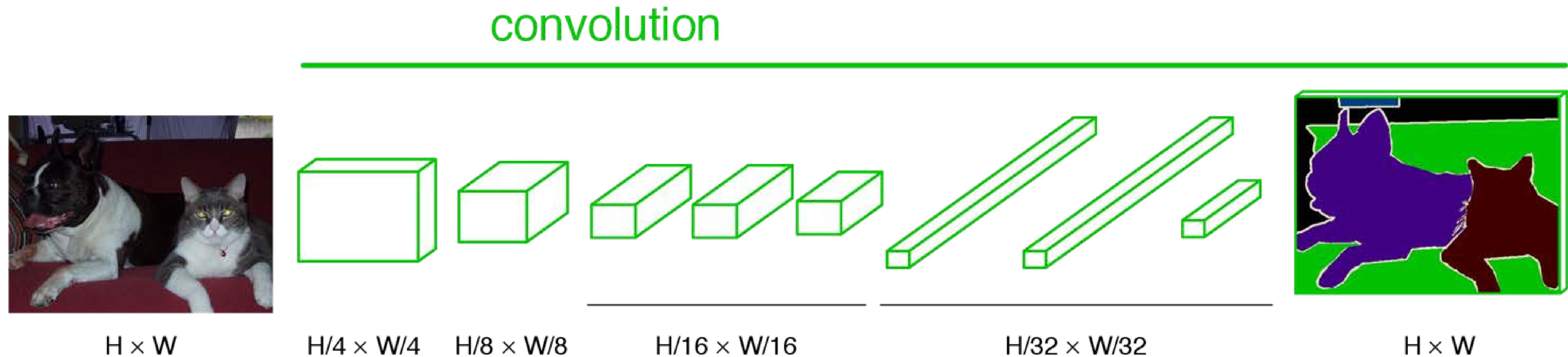


$H/32 \times W/32$

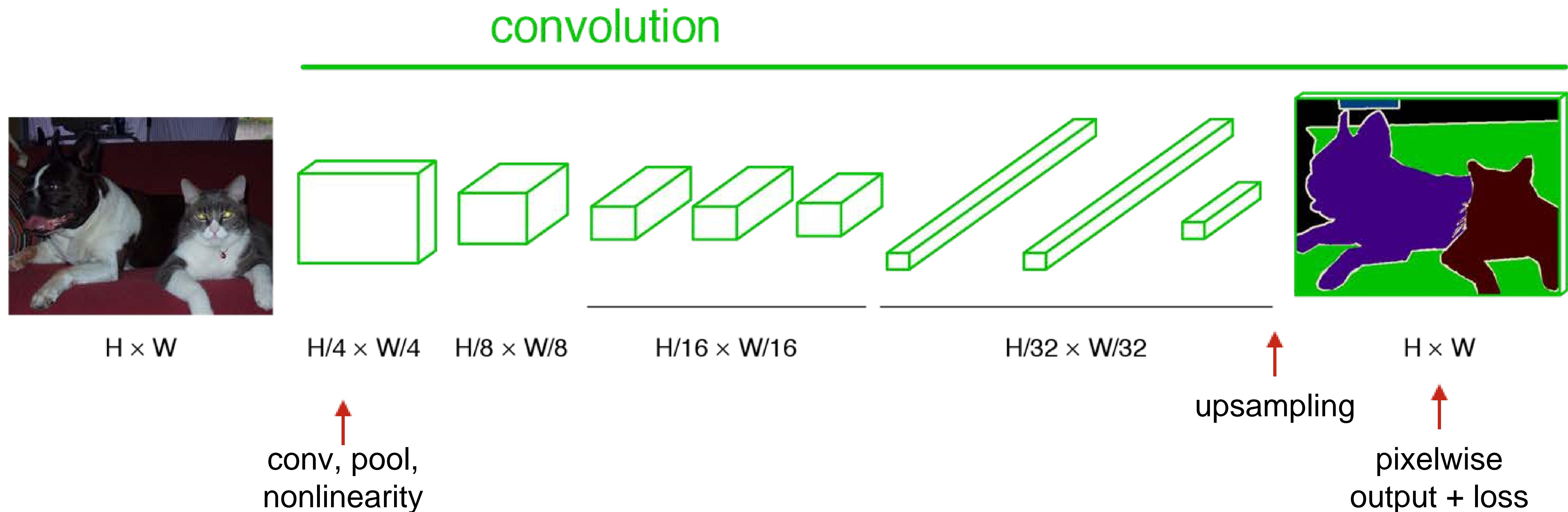


$H \times W$

## end-to-end, pixels-to-pixels network



## end-to-end, pixels-to-pixels network



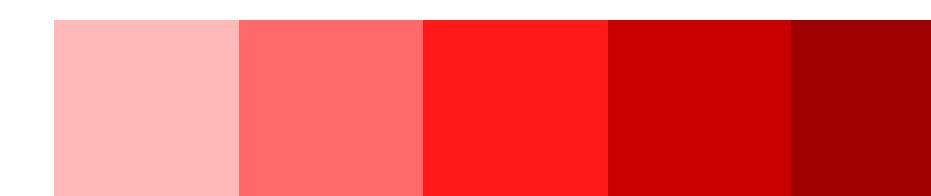
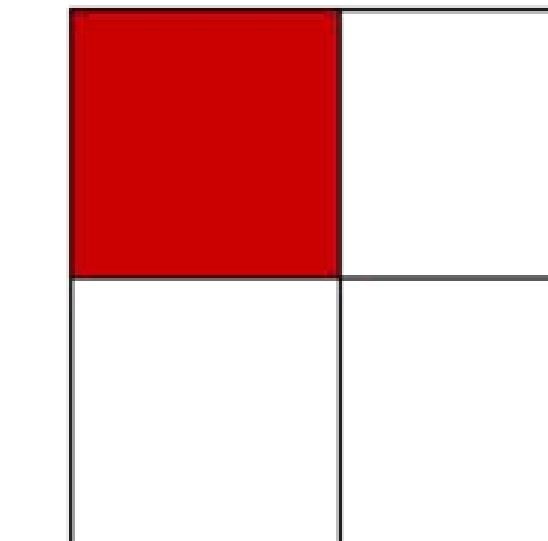
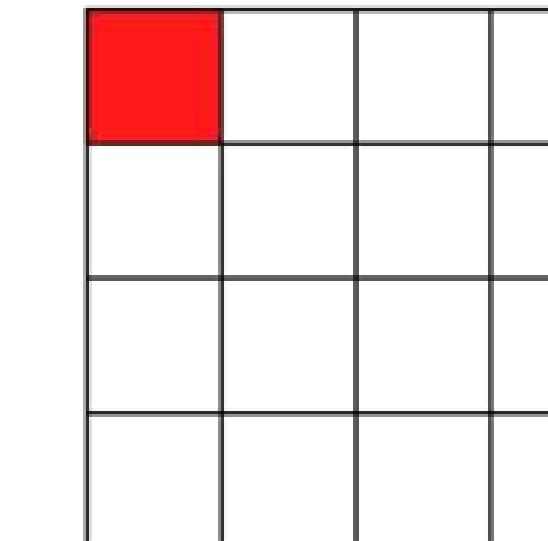
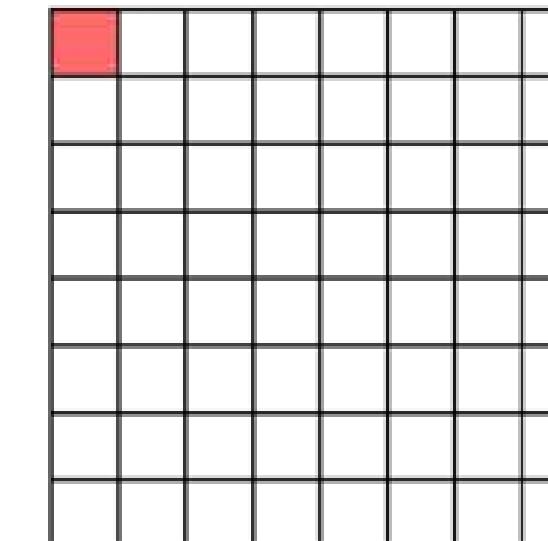
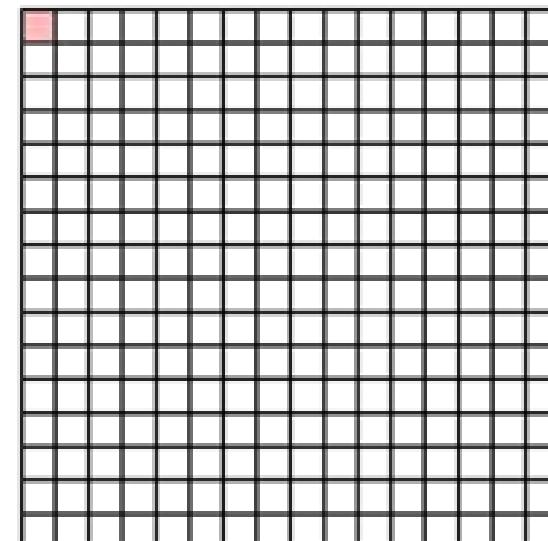
## spectrum of deep features

combine *where* (local, shallow) with *what* (global, deep)

image



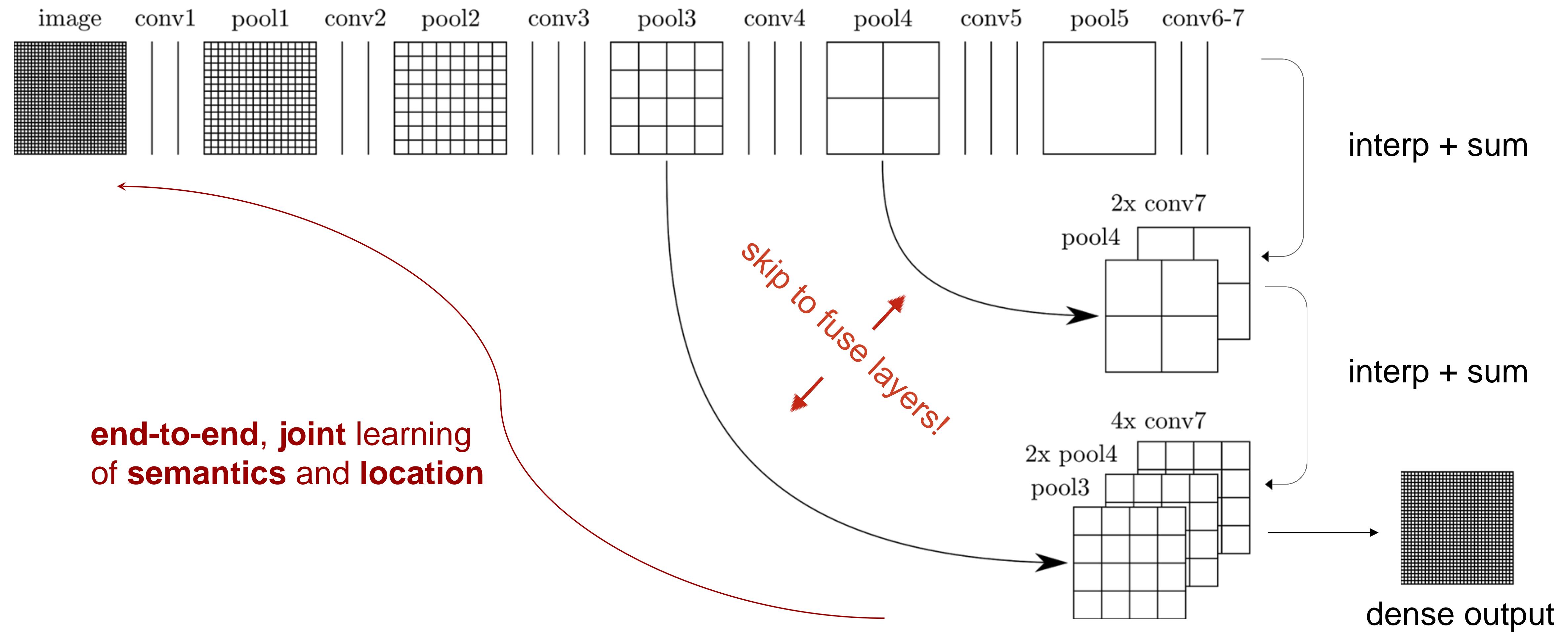
intermediate layers



fuse features into **deep jet**

(cf. Hariharan et al. CVPR15 “hypercolumn”)

## skip layers

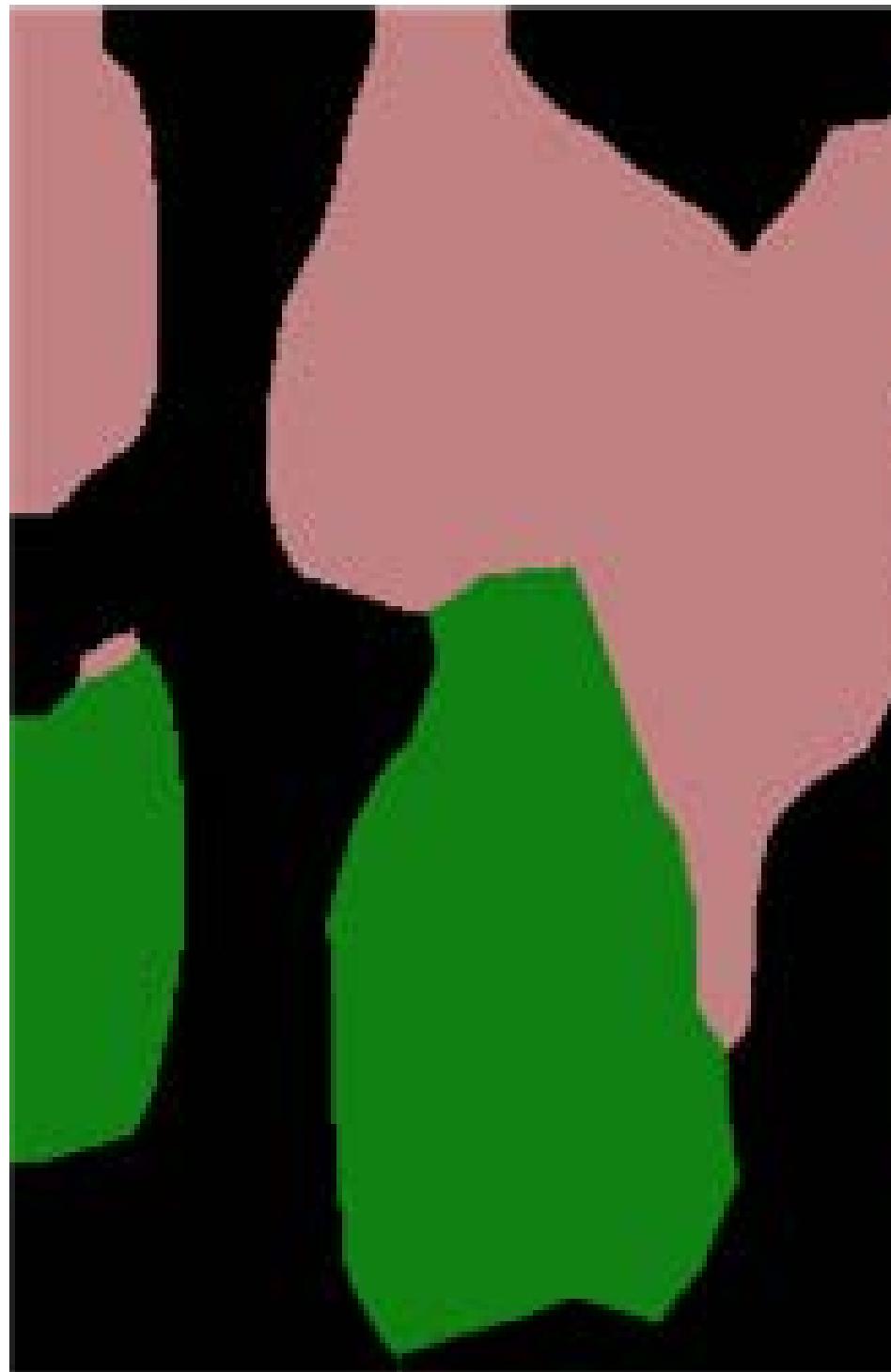


## **skip layer refinement**

input image



stride 32



stride 16



stride 8



ground truth

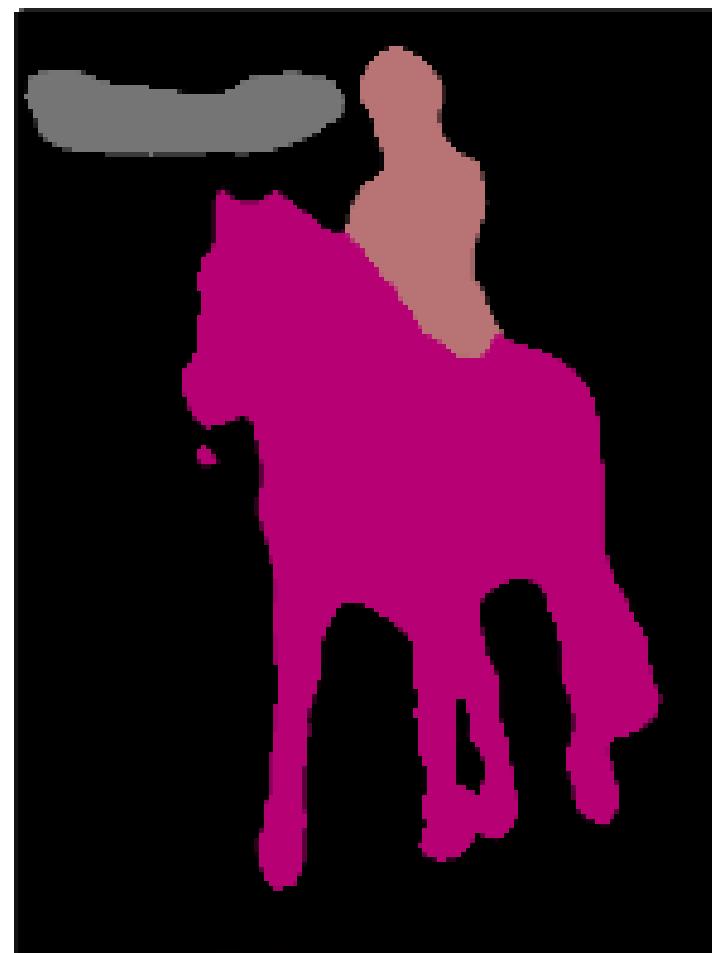


no skips

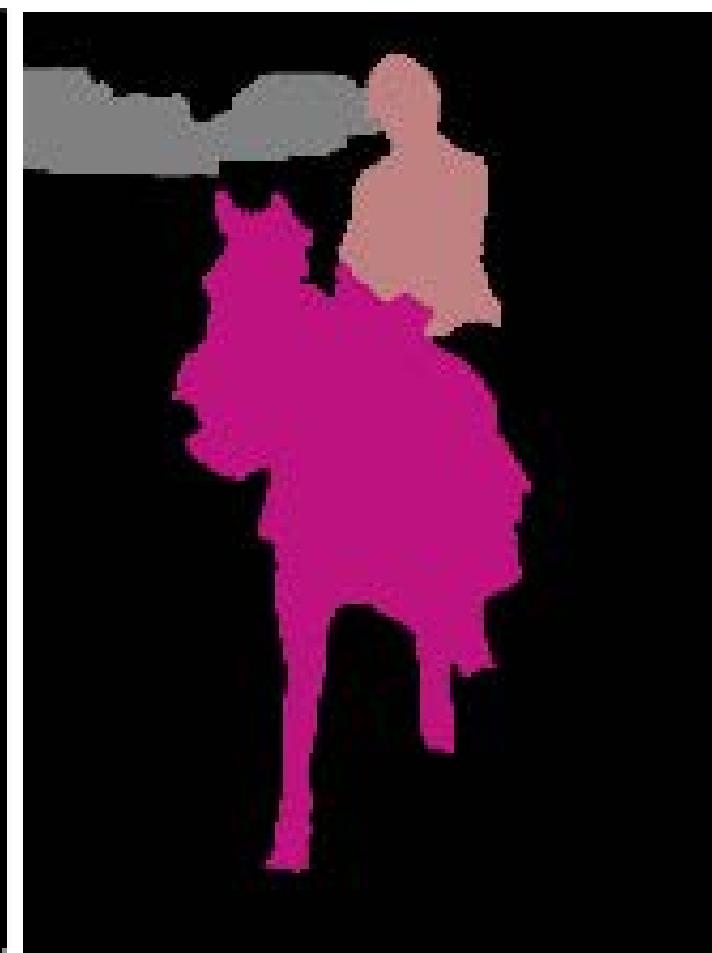
1 skip

2 skips

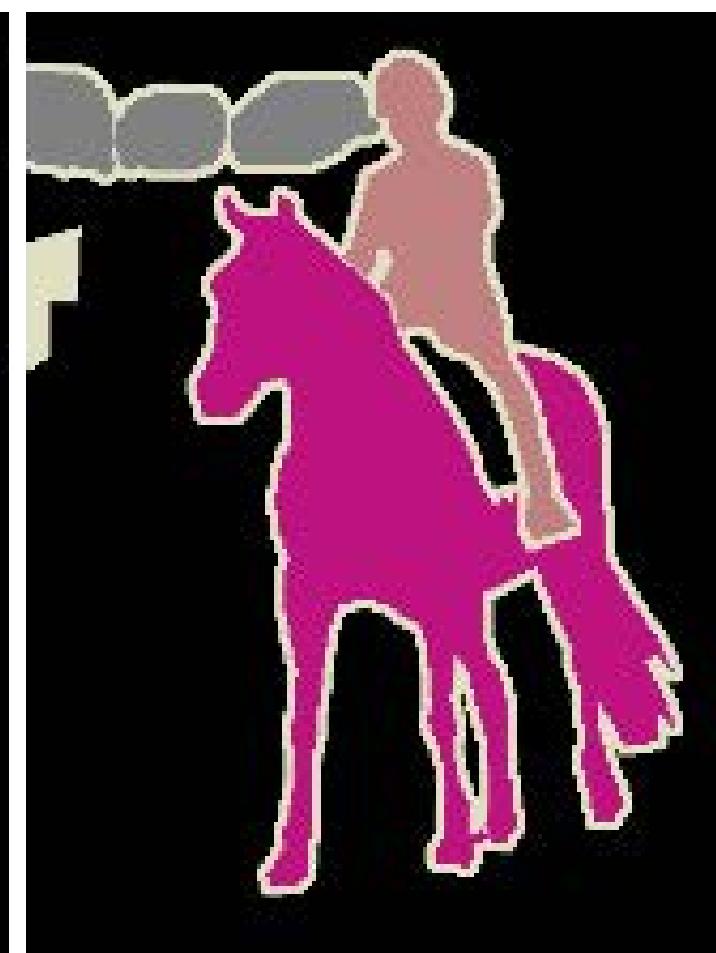
FCN



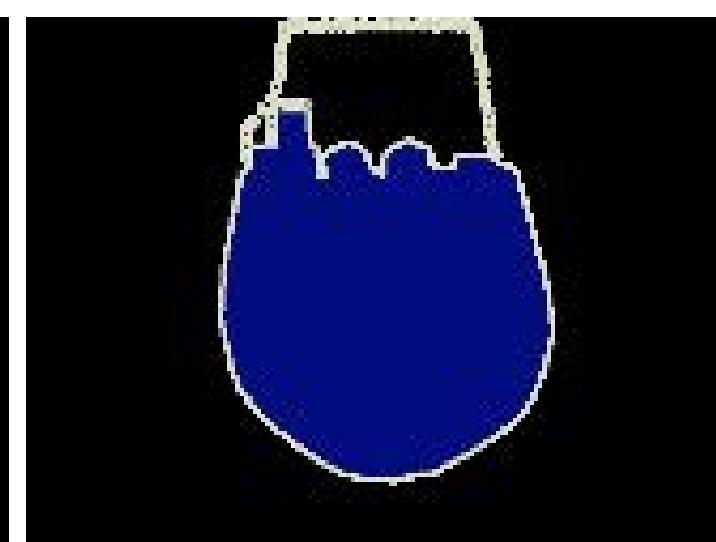
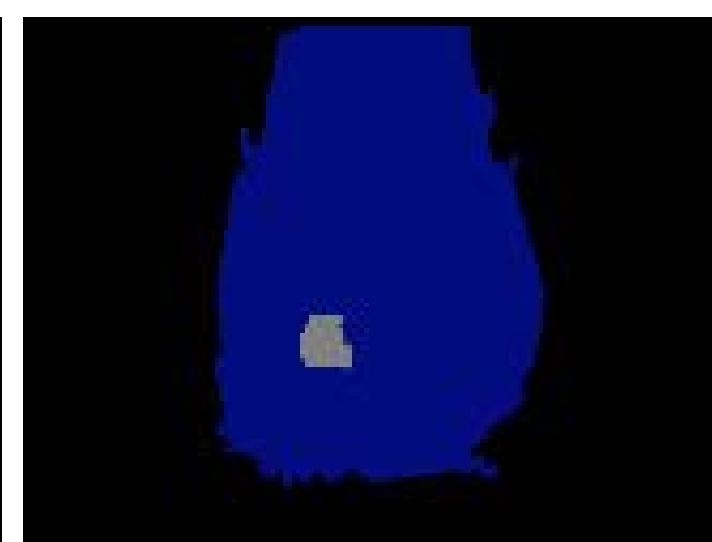
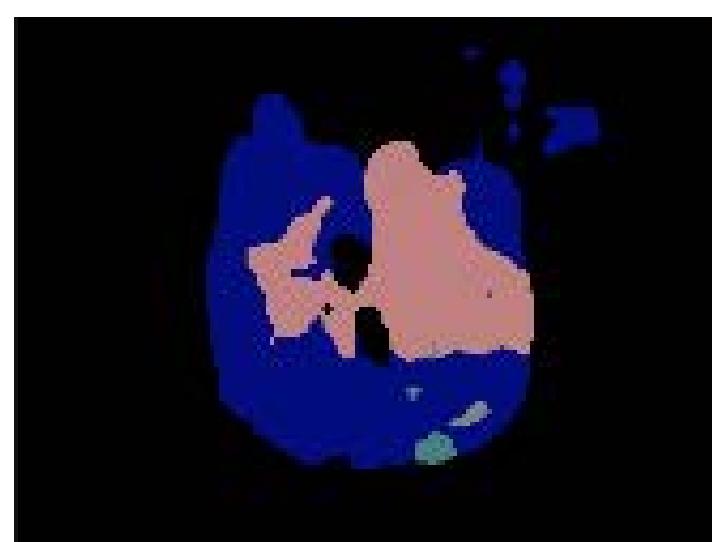
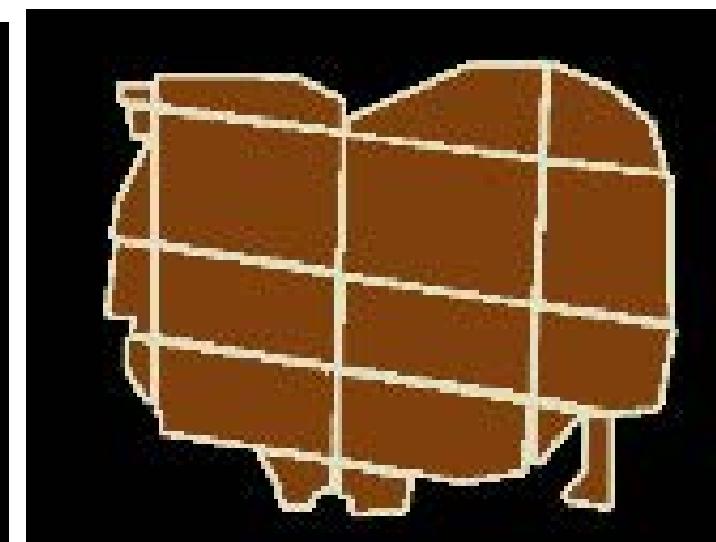
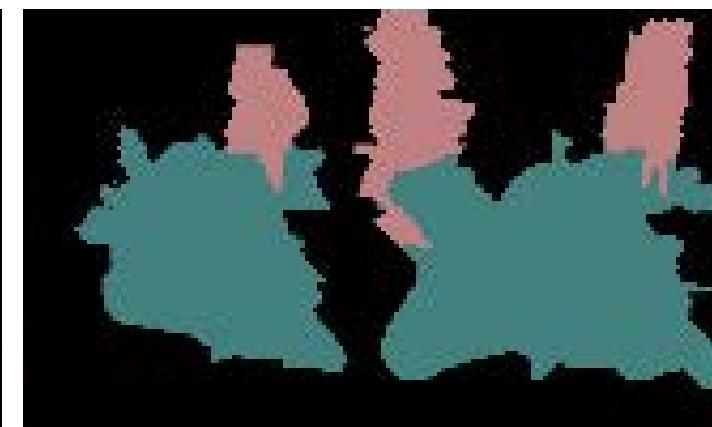
SDS\*



Truth



Input

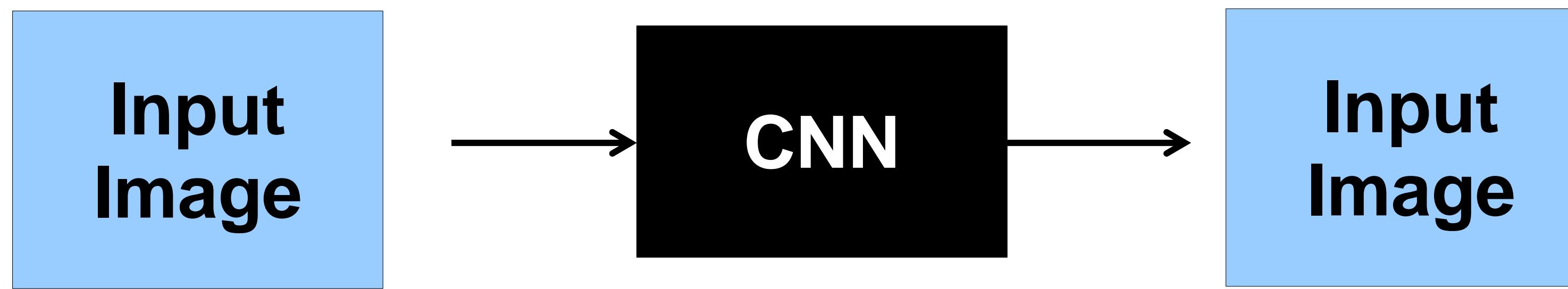


Relative to prior state-of-the-art SDS:

- 30% relative improvement for mean IoU
- 286× faster

\*Simultaneous Detection and Segmentation  
Hariharan et al. ECCV14

# Fancier Architectures: Auto-encoders

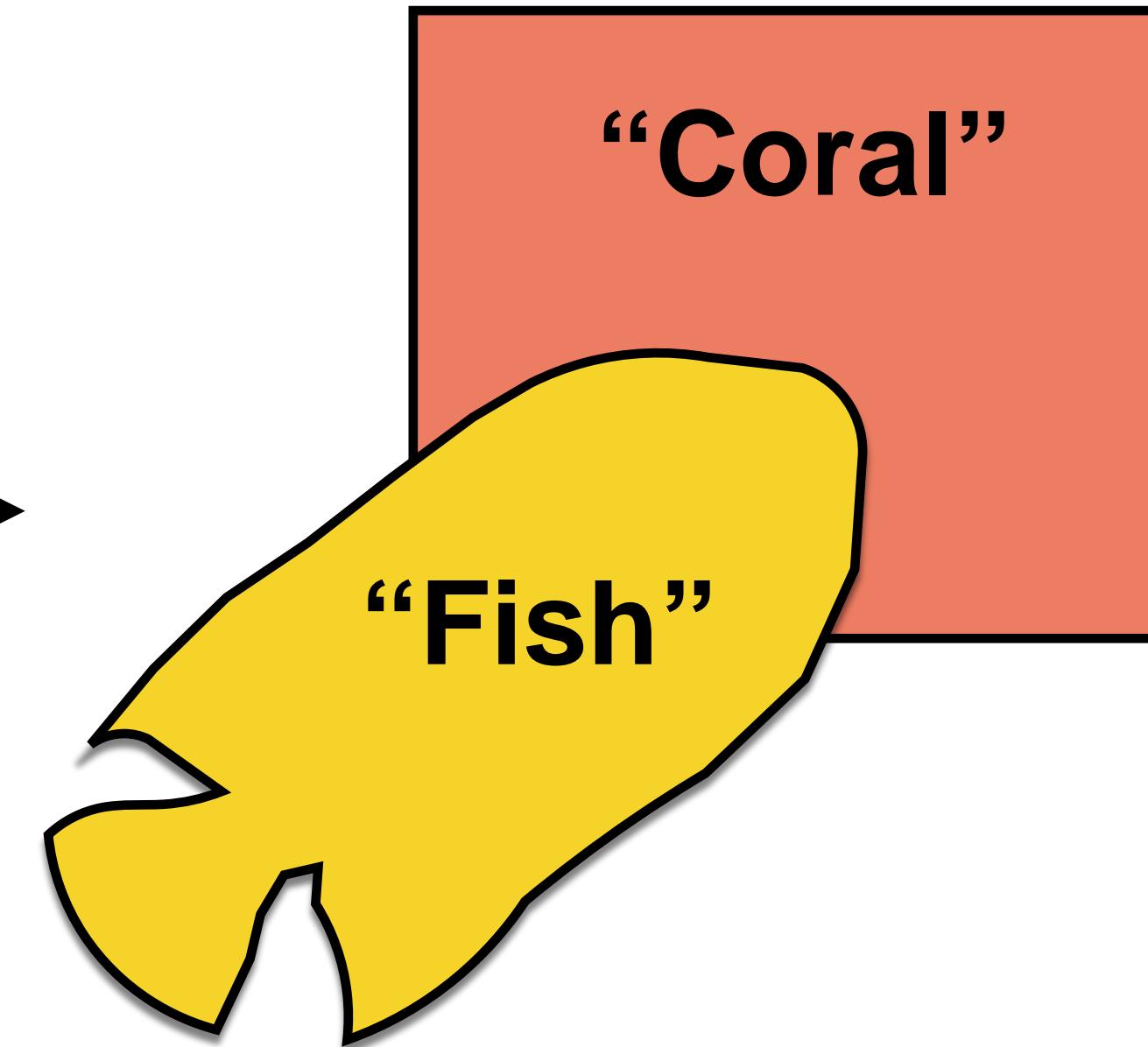
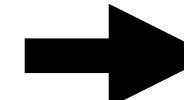


# Representation Learning

**X**

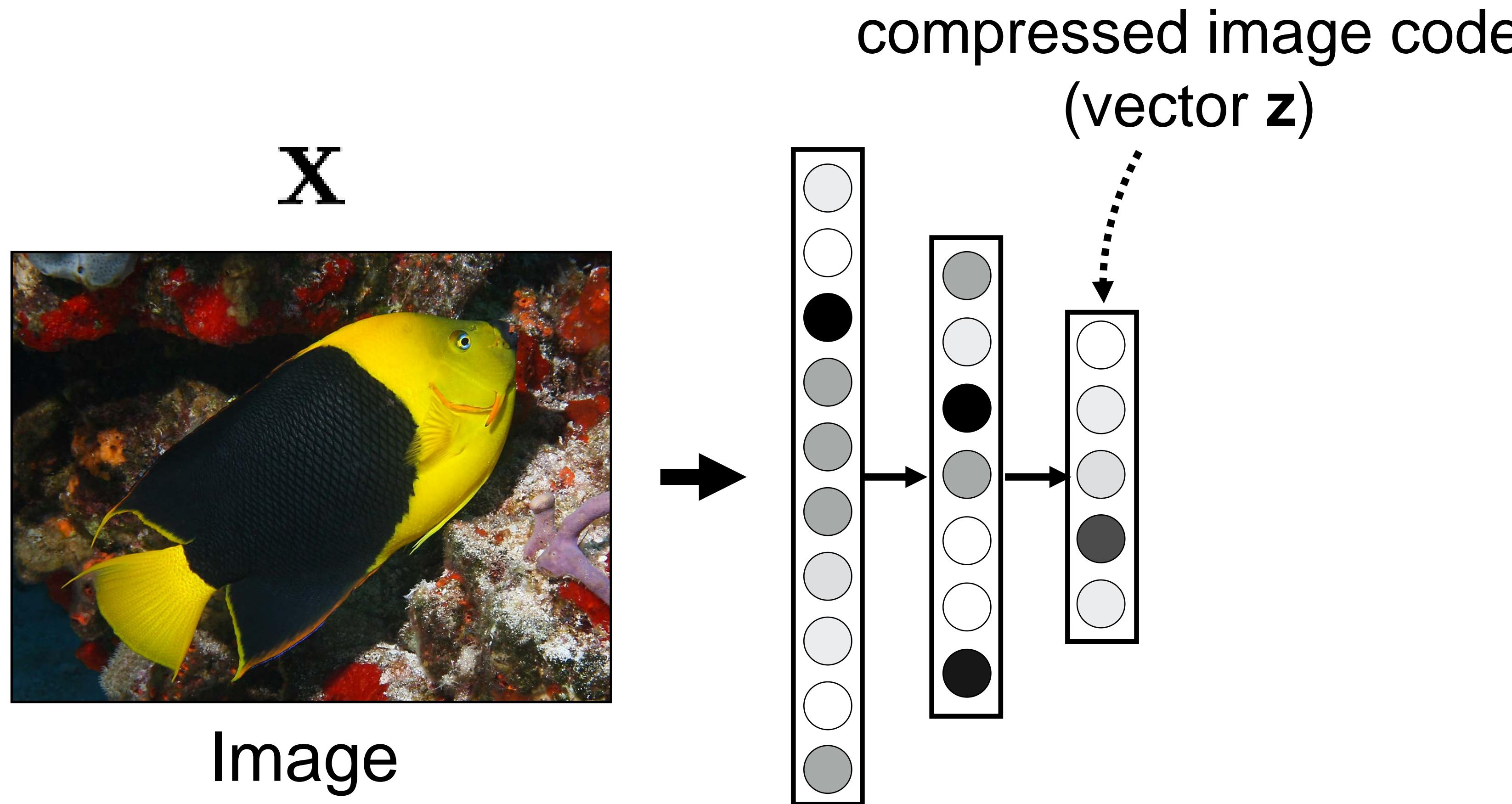


Image

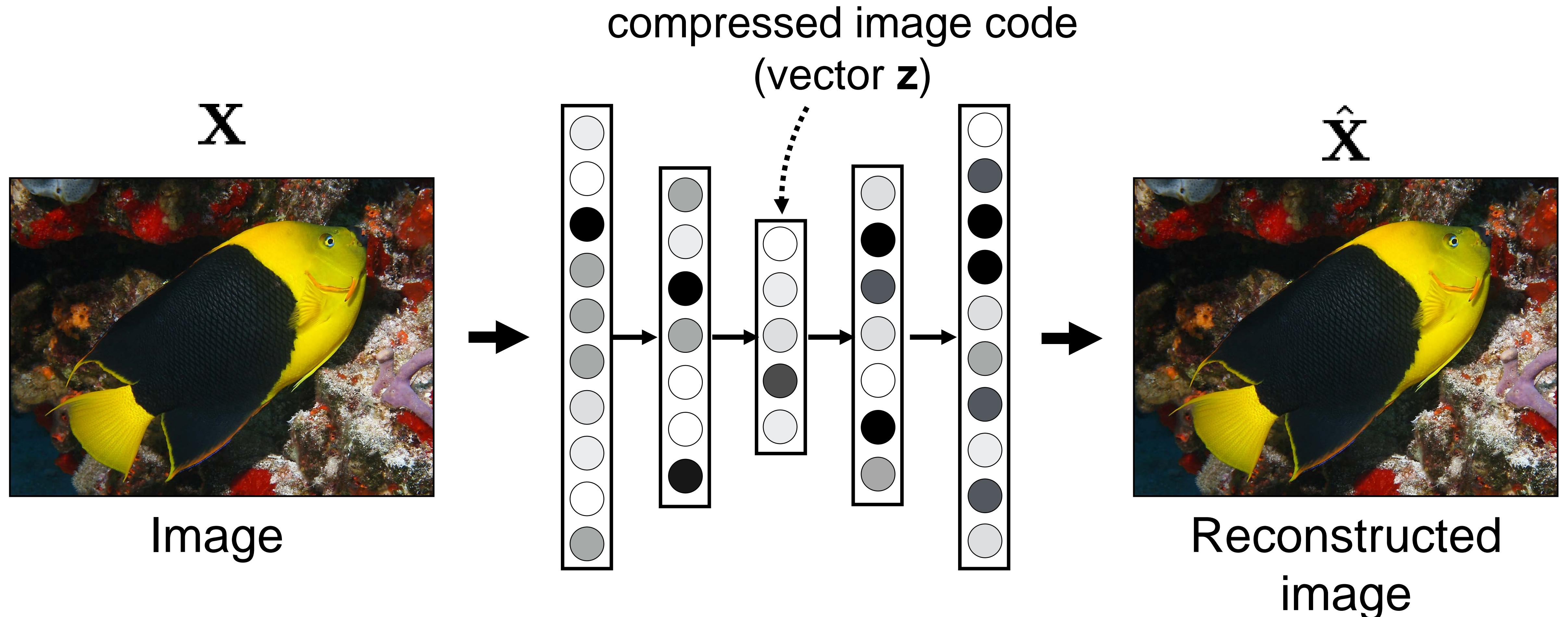


Compact mental  
representation

# Representation Learning



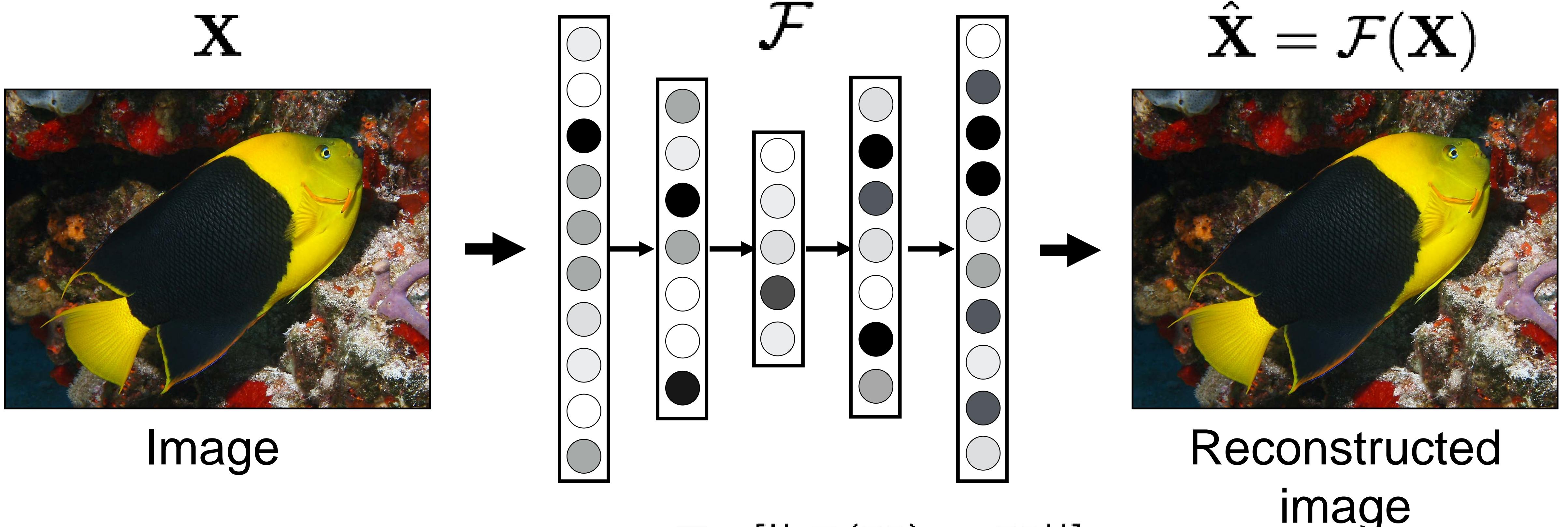
# Representation Learning



“Autoencoder”

[e.g., Hinton & Salakhutdinov, Science 2006]

# Autoencoder



$$\arg \min_{\mathcal{F}} \mathbb{E}_{\mathbf{X}} [\|\mathcal{F}(\mathbf{X}) - \mathbf{X}\|]$$

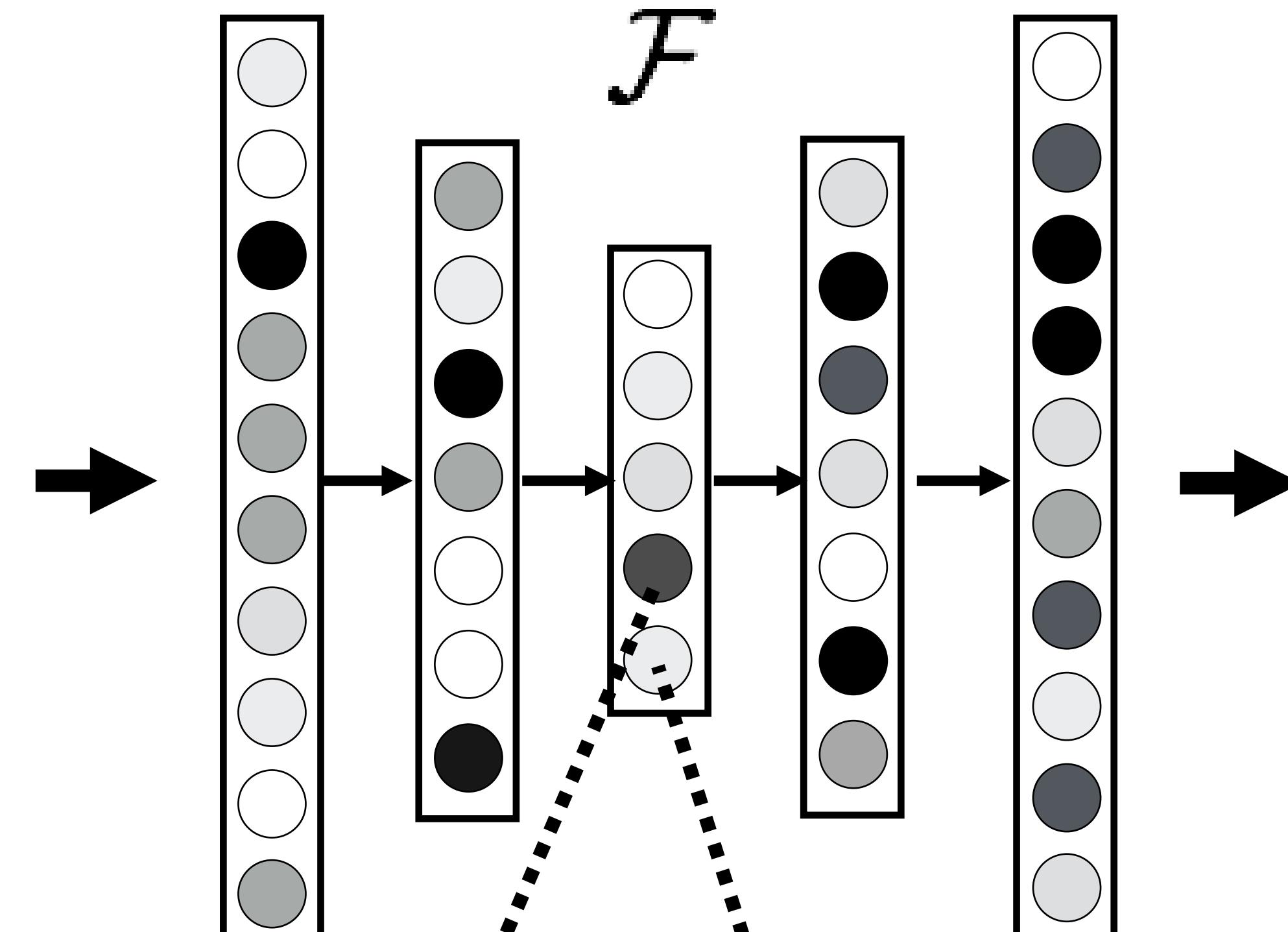
[e.g., Hinton & Salakhutdinov, Science 2006]

$\mathbf{X}$



Image

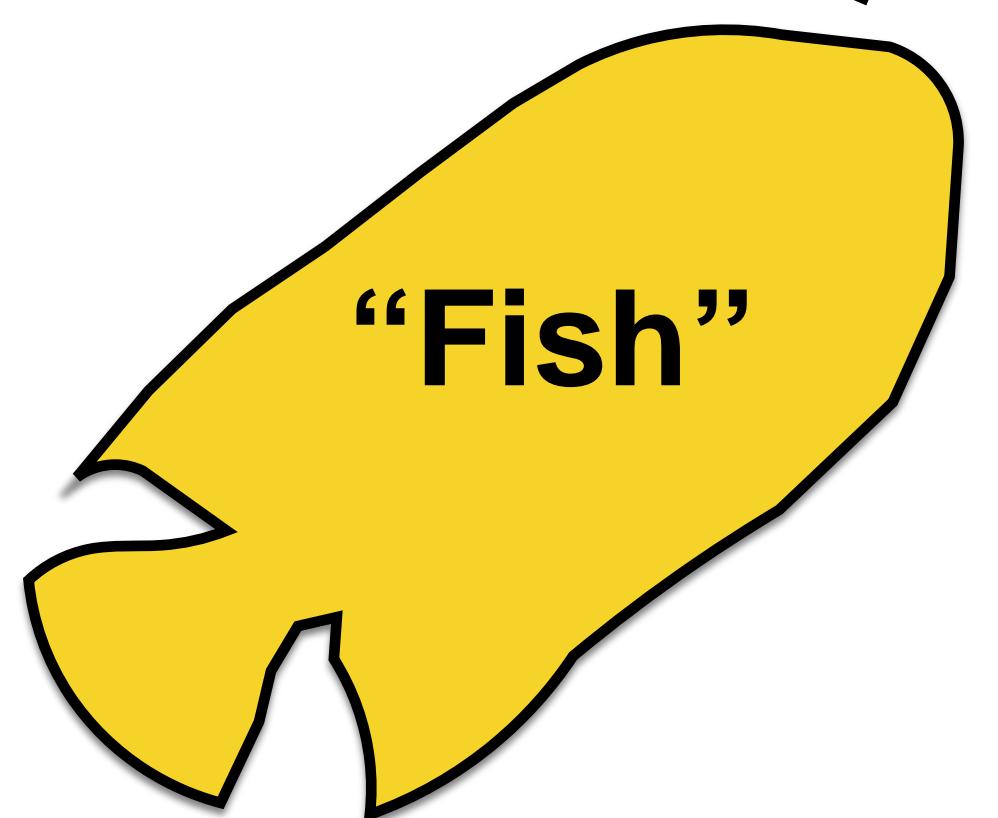
$\mathcal{F}$



$\hat{\mathbf{X}} = \mathcal{F}(\mathbf{X})$

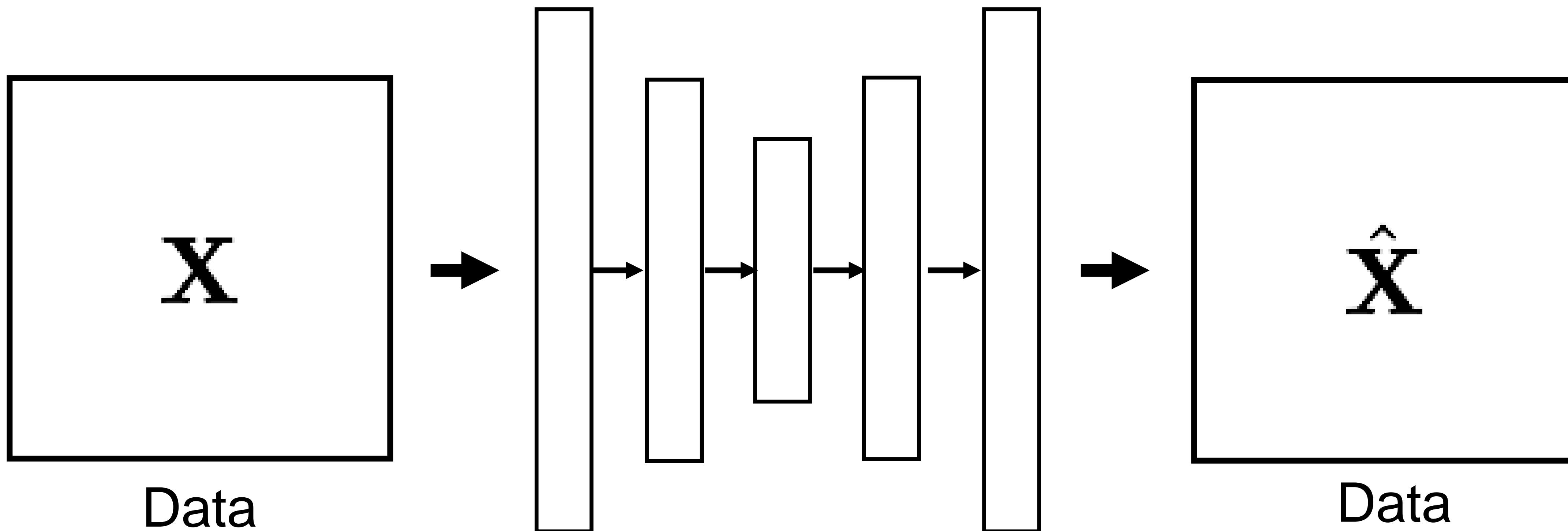


Reconstructed  
image



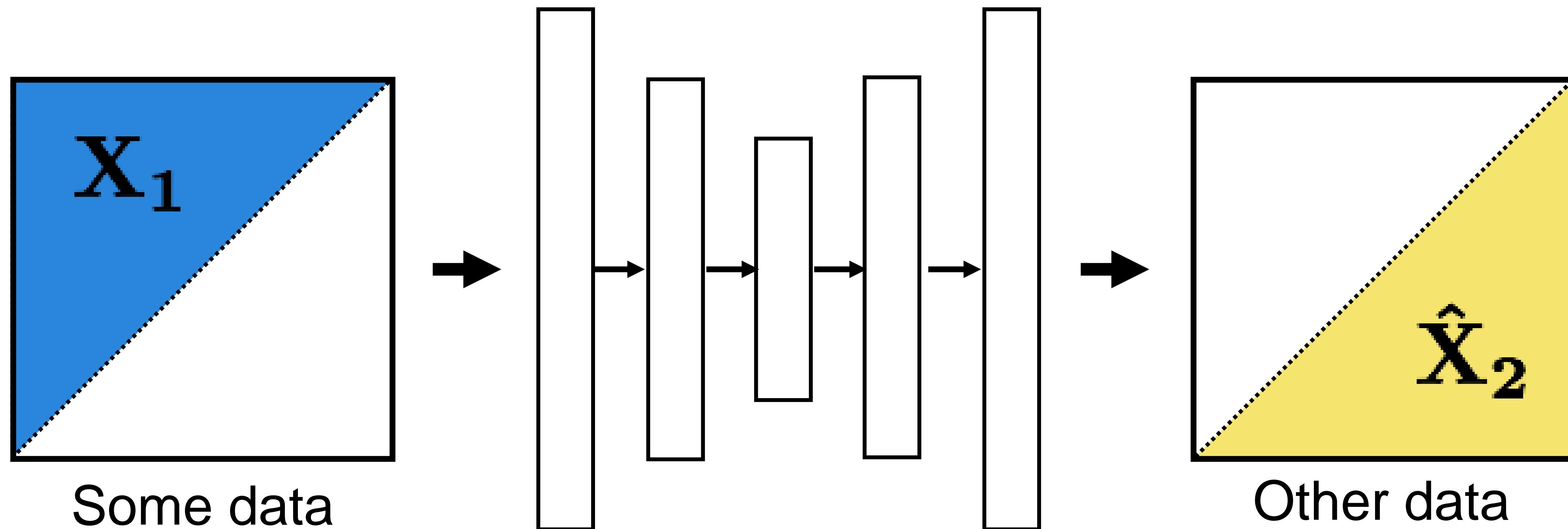
[e.g., Hinton & Salakhutdinov, Science 2006]

# Data compression



[Hinton & Salakhutdinov, Science 2009]

# Data prediction



see also [Vincent et al., 2008]



$$\xrightarrow{\mathcal{F}}$$

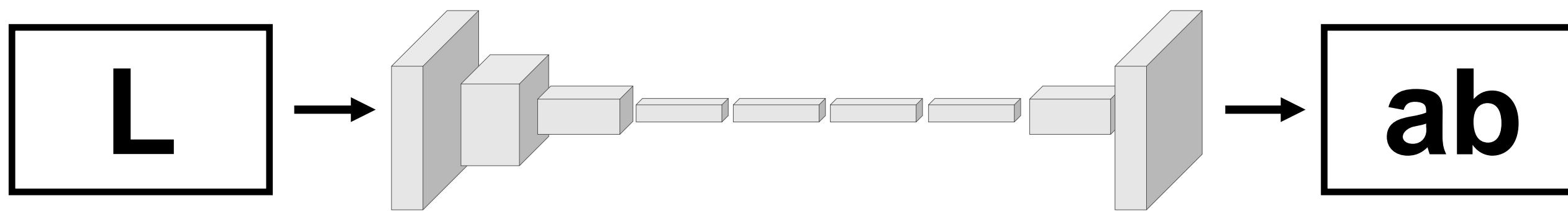


Grayscale image: L channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: ab channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$



[Zhang, Isola, Efros, ECCV 2016]



$$\xrightarrow{\mathcal{F}}$$



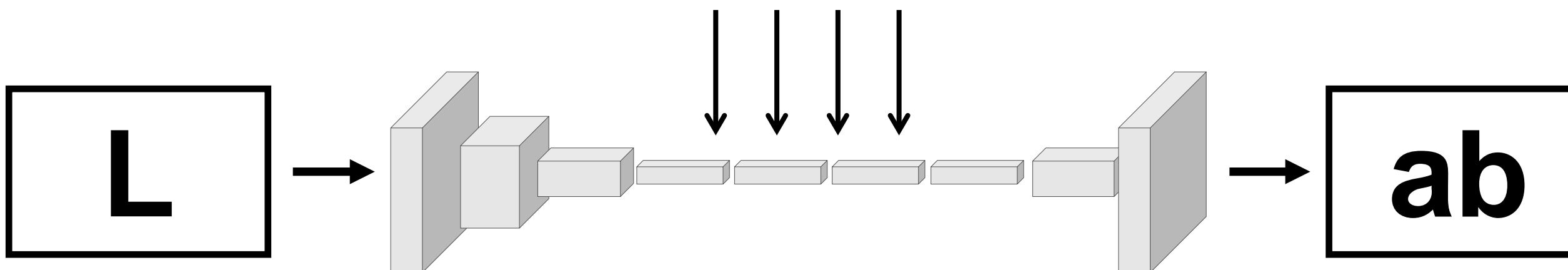
Grayscale image: L cha

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

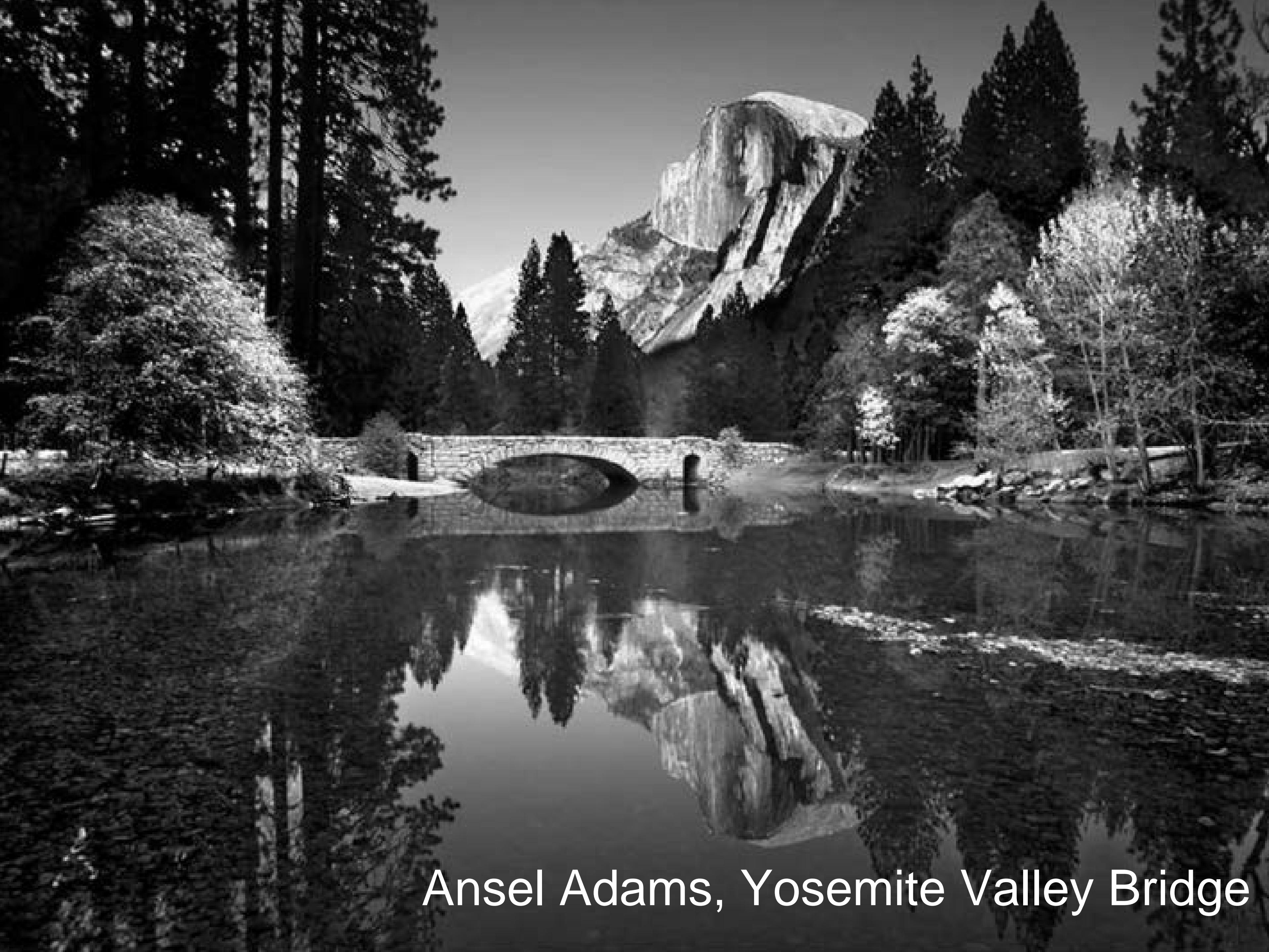
Semantics? Higher-  
level abstraction?

information: ab channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$



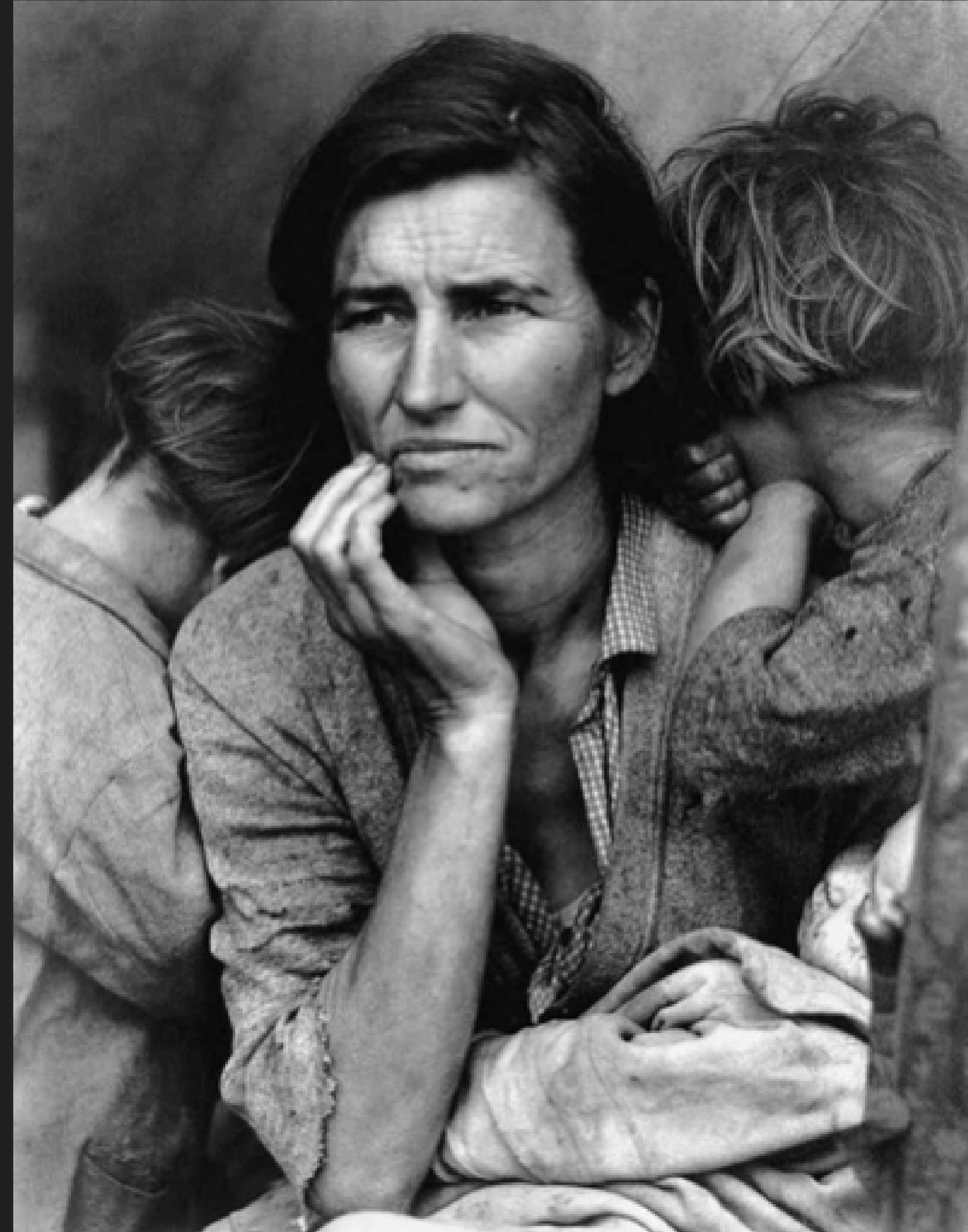
[Zhang, Isola, Efros, ECCV 2016]



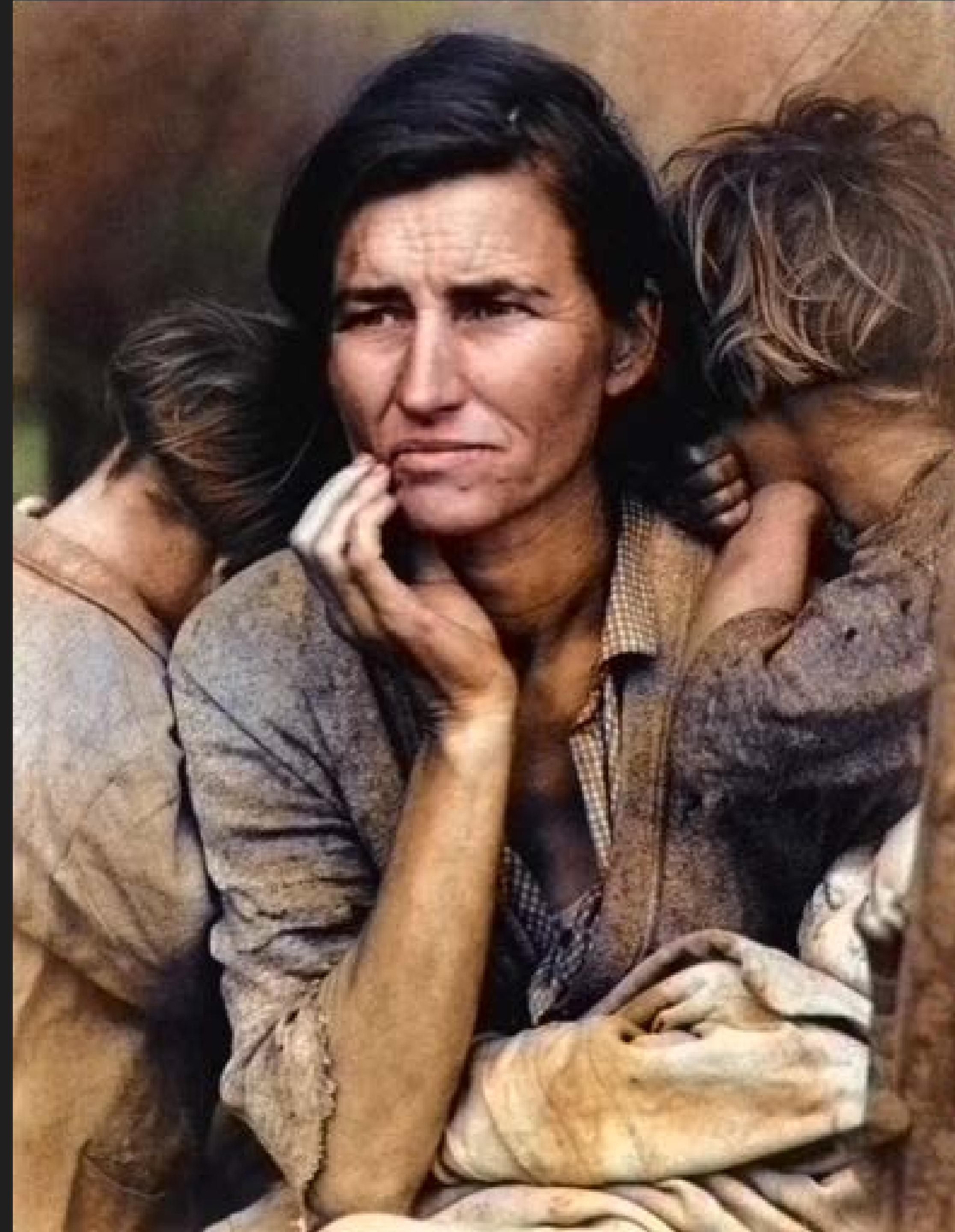
Ansel Adams, Yosemite Valley Bridge



Our result



*Migrant Mother*  
Dorothea Lange  
1936



Our result



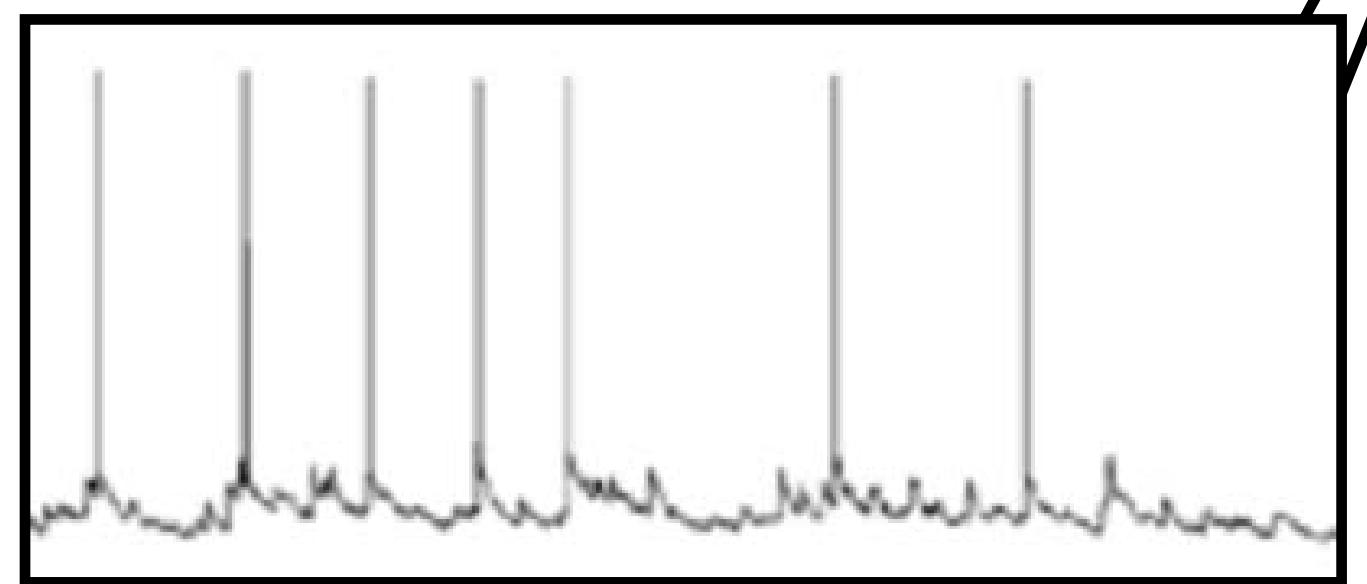
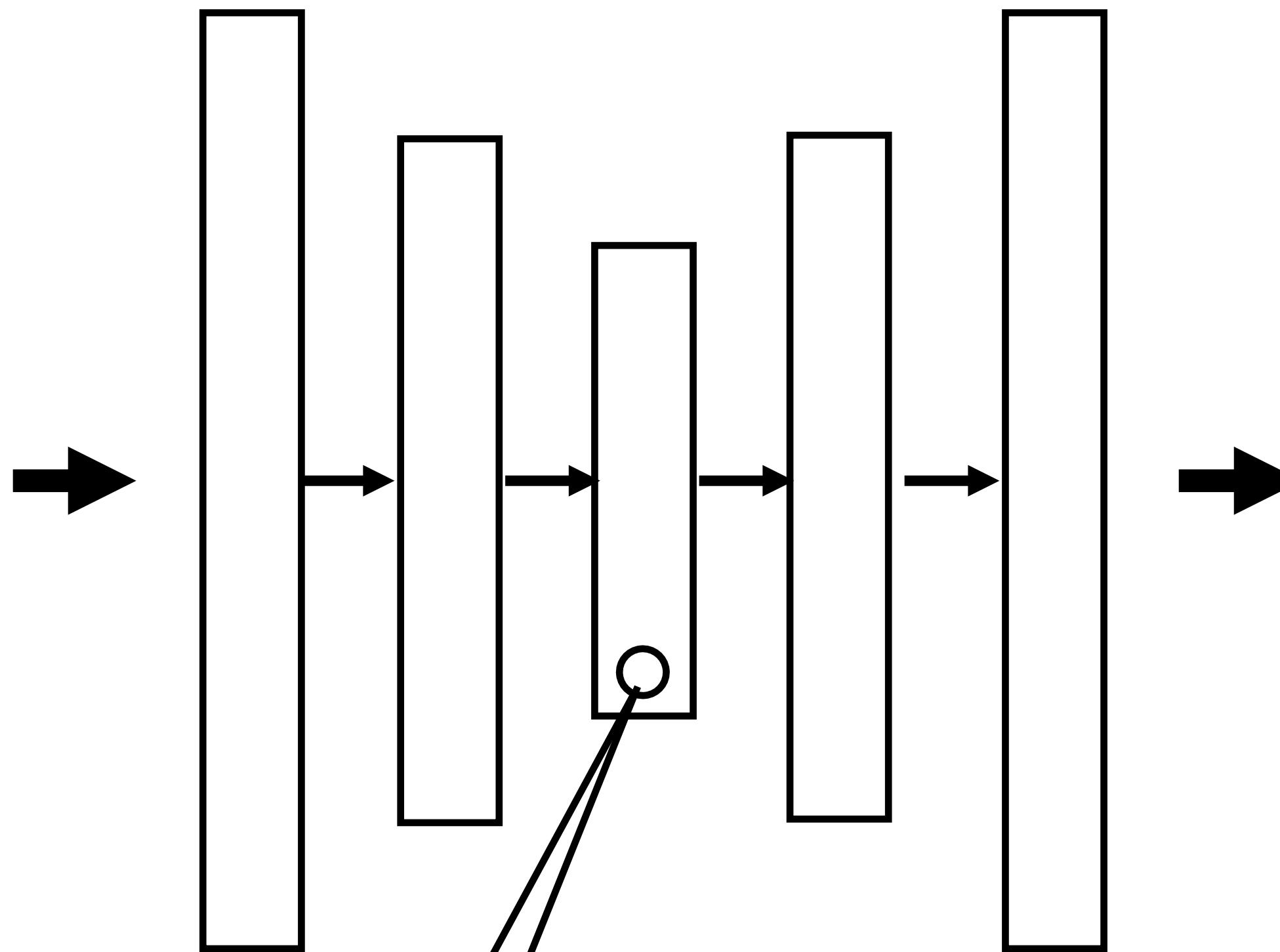
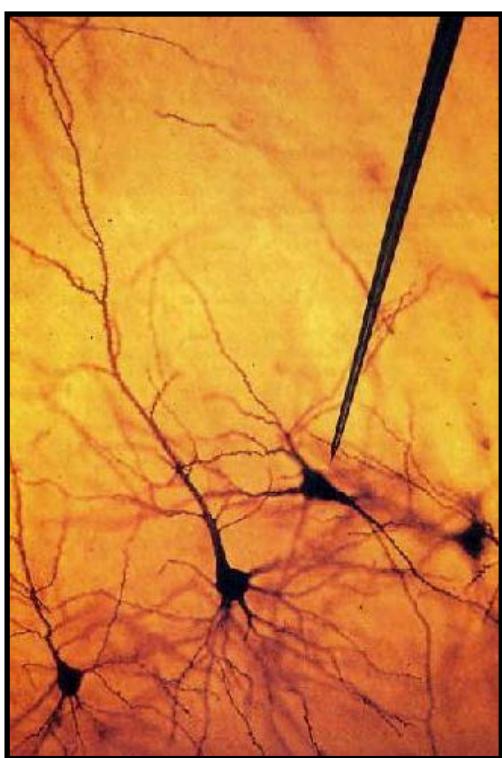
# Instructive failure



# Instructive failure



# Deep Net “Electrophysiology”



[Zeiler & Fergus, ECCV 2014]  
[Zhou et al., ICLR 2015]

# Stimuli that drive selected neurons (conv5 layer)

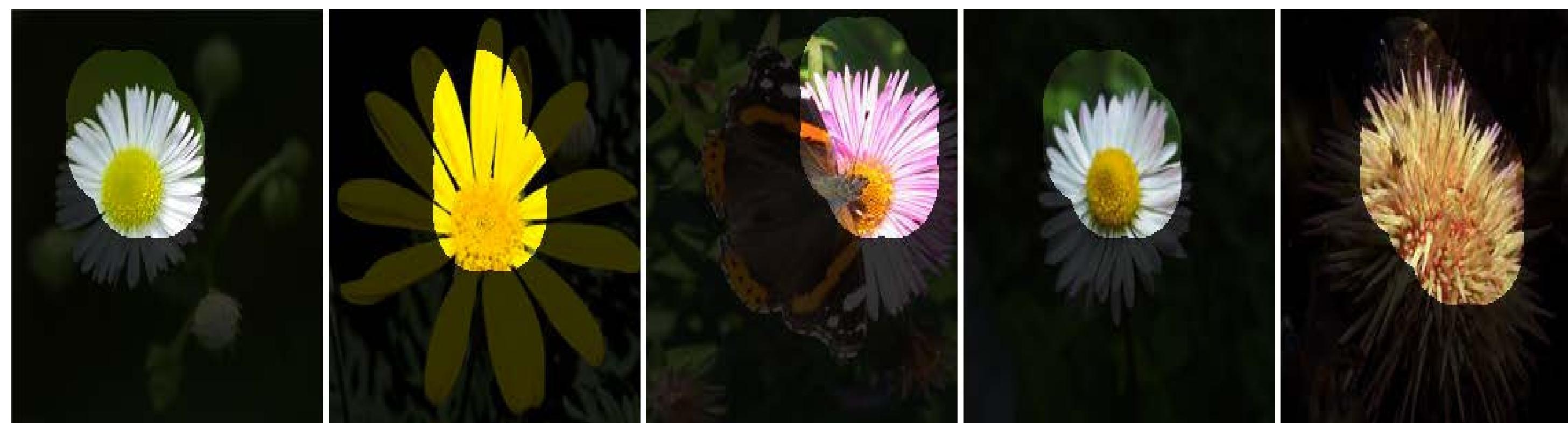
faces



dog  
faces

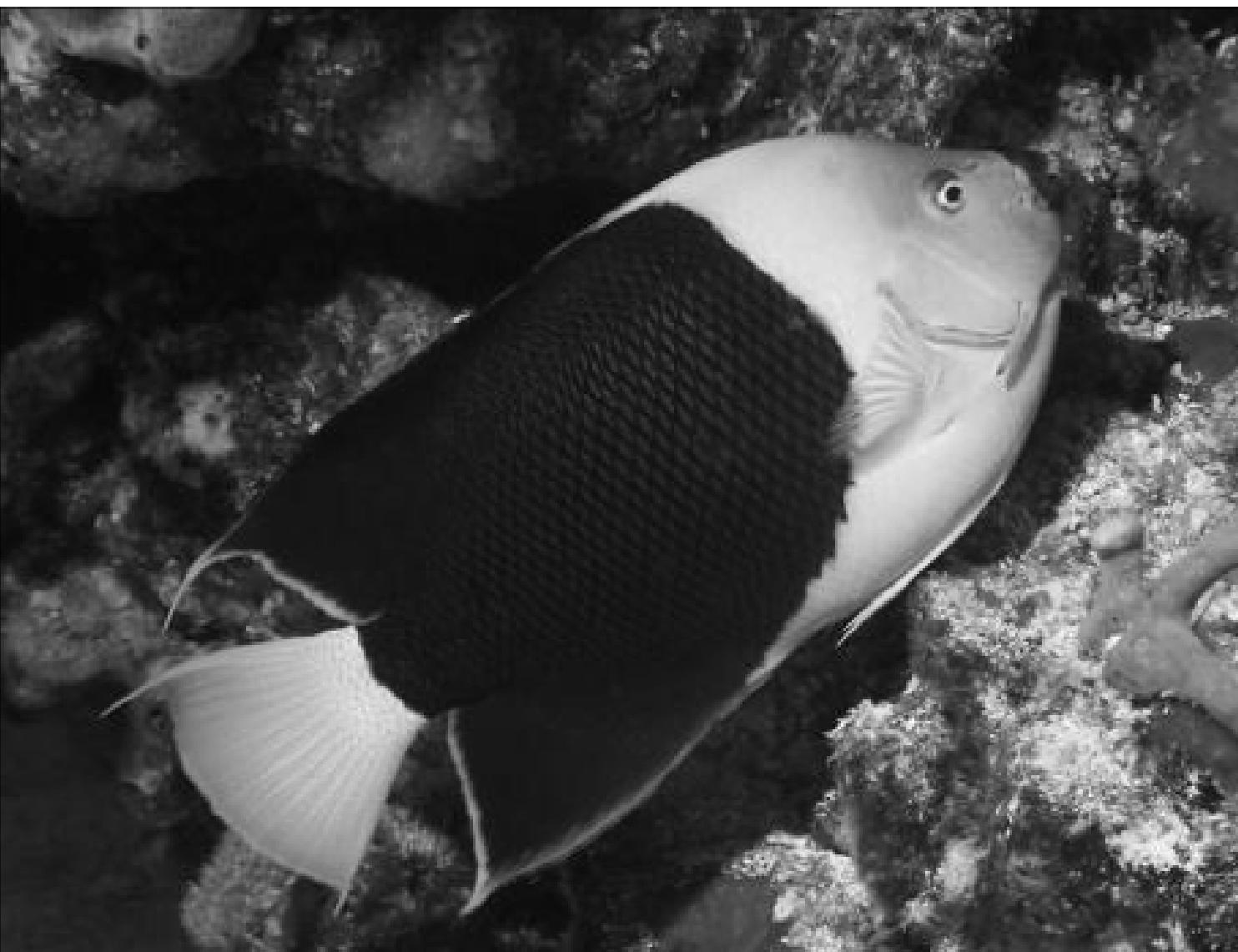


flowers

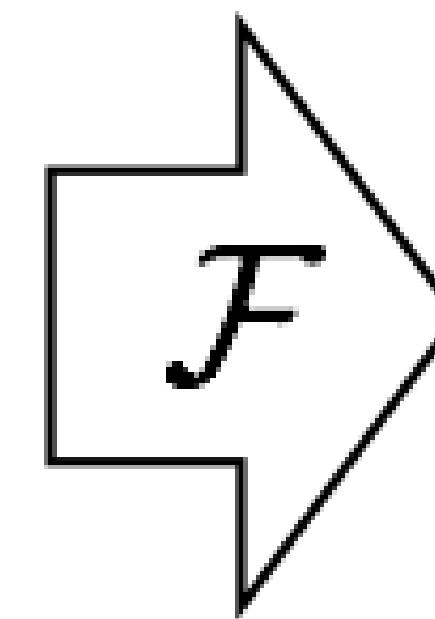


# How to evaluate predictions?

Input



Output



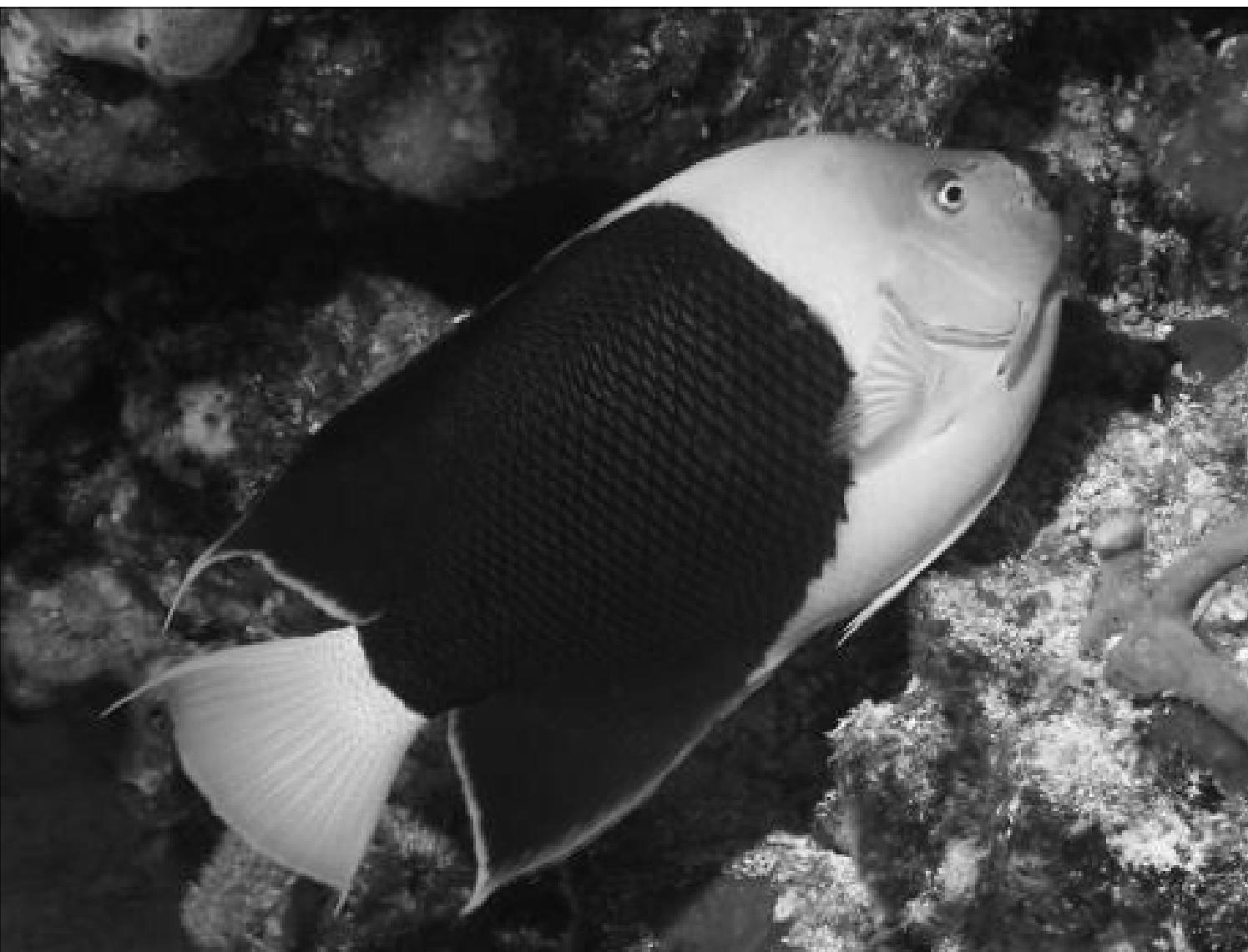
$$\arg \min_{\mathcal{F}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [L(\mathcal{F}(\mathbf{x}), \mathbf{y})]$$

Objective function  
(loss)

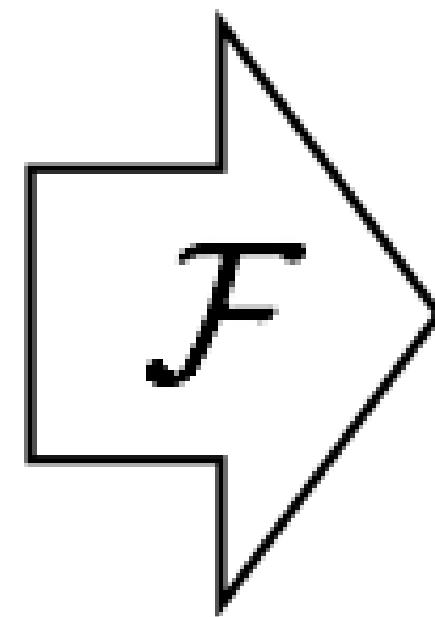
Neural Network

# How to evaluate predictions?

Input



Output



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

**“What should I do”**

**“How should I do it?”**

# Designing objective functions

Input



Output

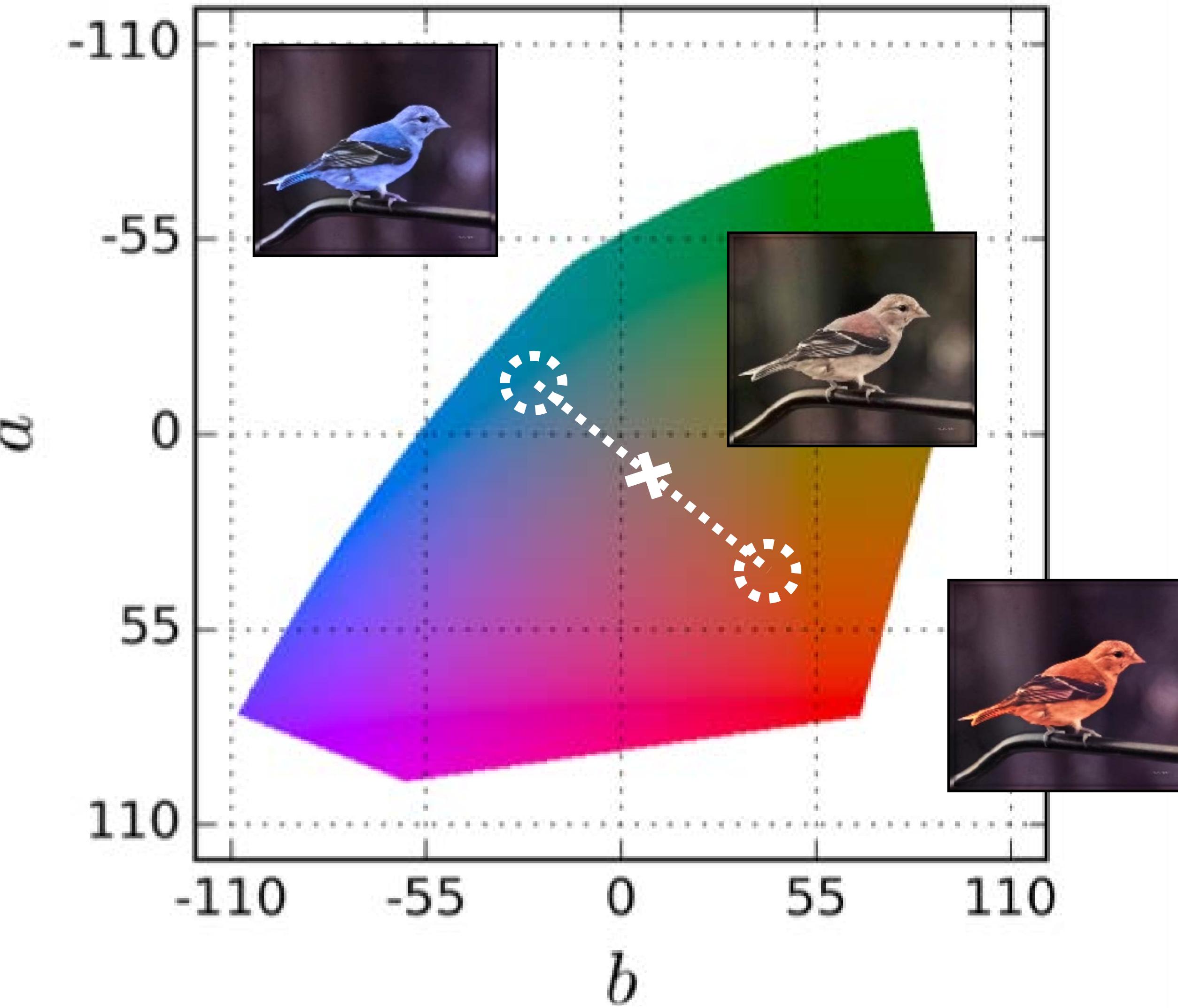


Ground truth



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

[Zhang, Isola, Efros, ECCV 2016]



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

# Designing objective functions

Input



Zhang et al. 2016

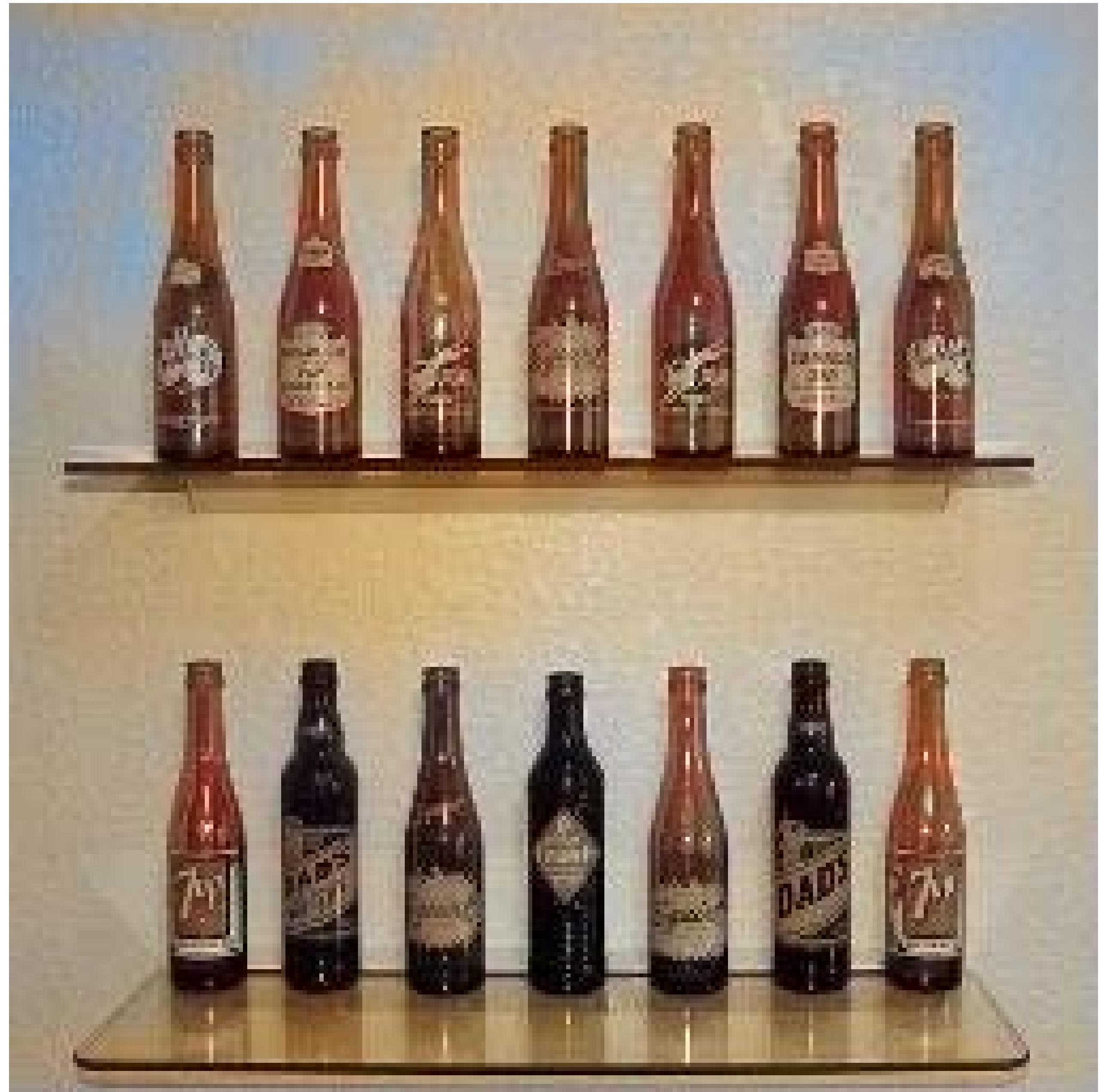


Ground truth

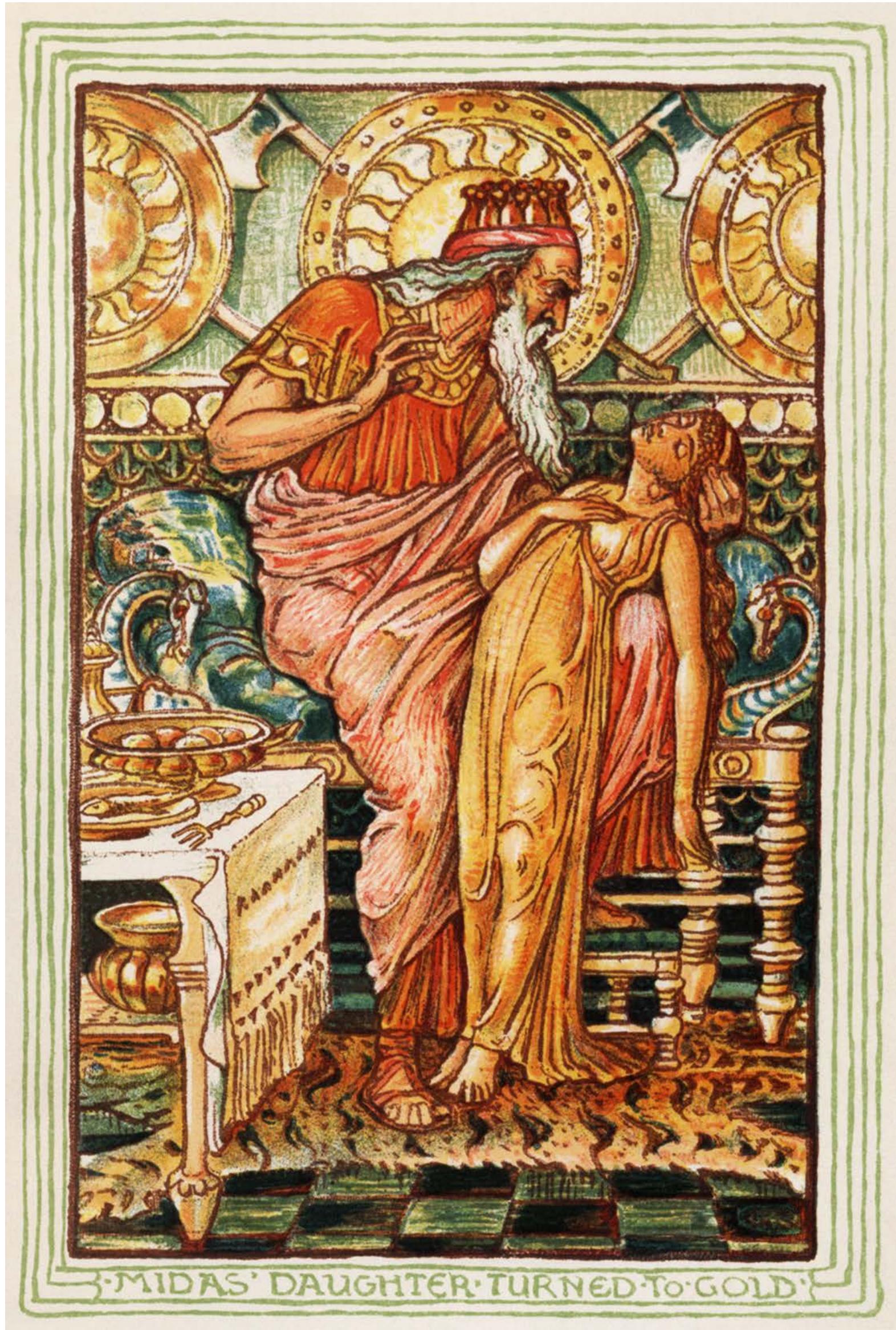


Color distribution cross-entropy loss with colorfulness enhancing term.

[Zhang, Isola, Efros, ECCV 2016]



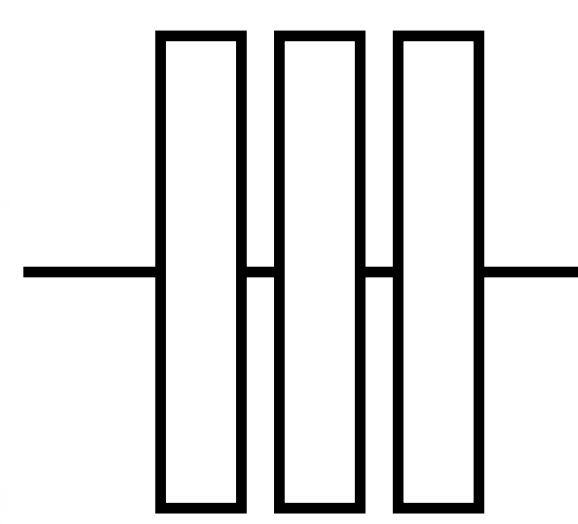
# Designing objective functions



Be careful what you wish for!

# Designing objective functions

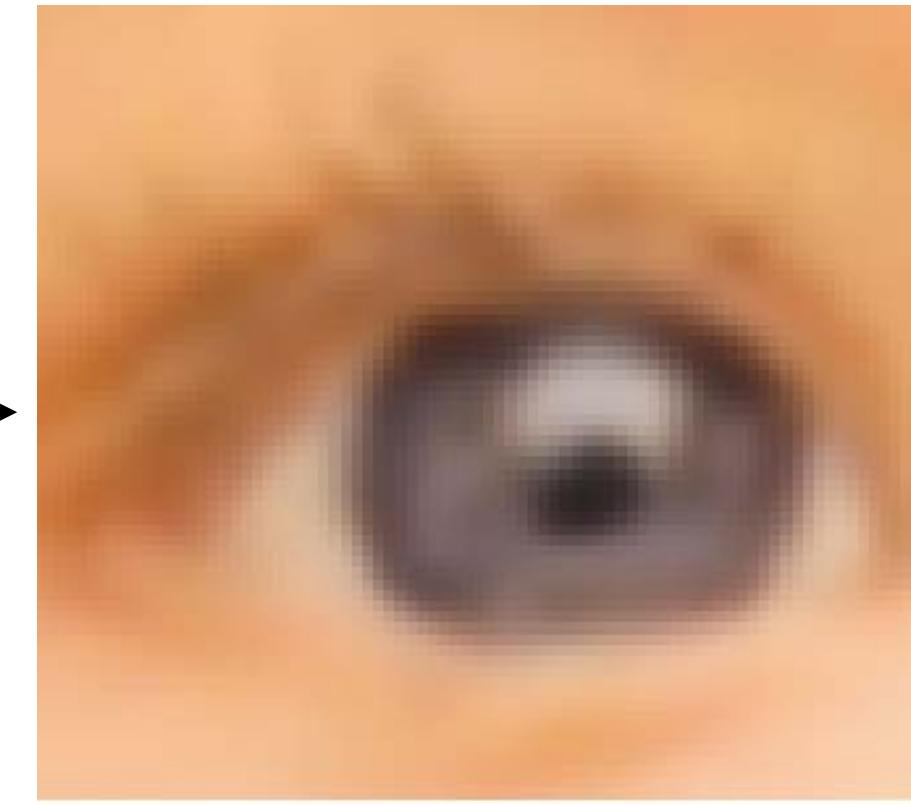
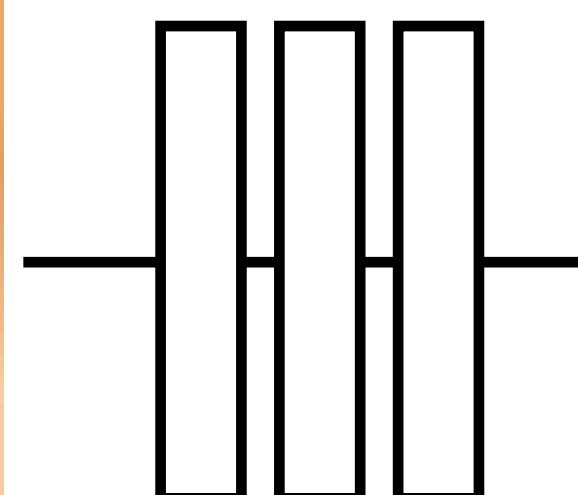
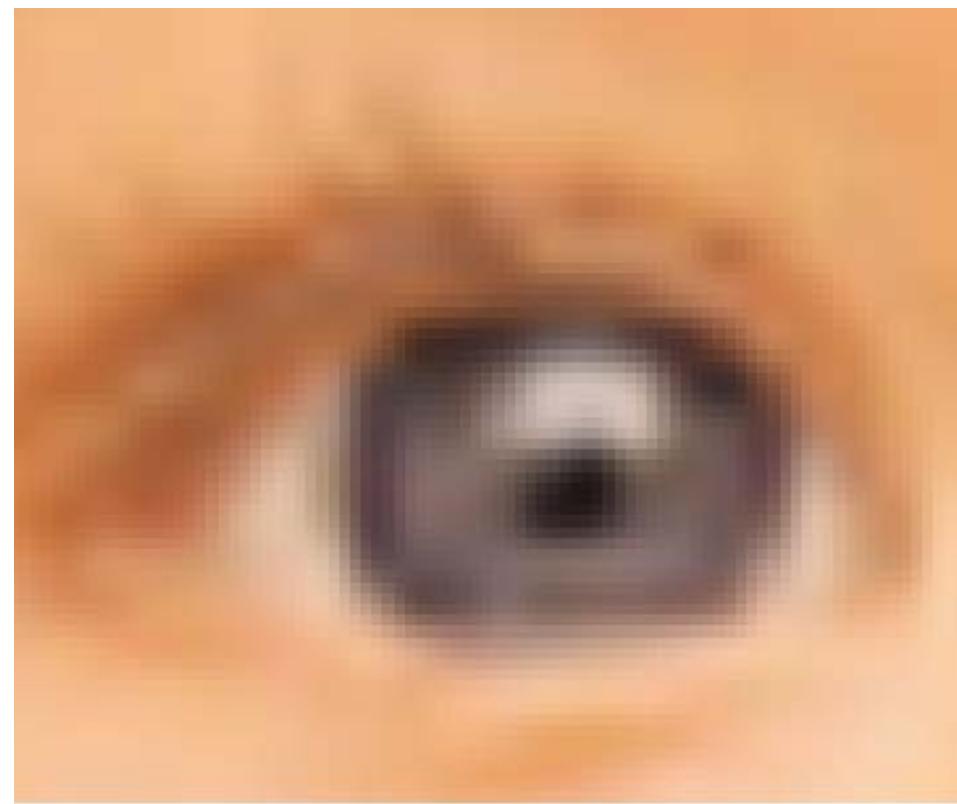
Image colorization



L2 regression

[Zhang, Isola, Efros, ECCV 2016]

Super-resolution

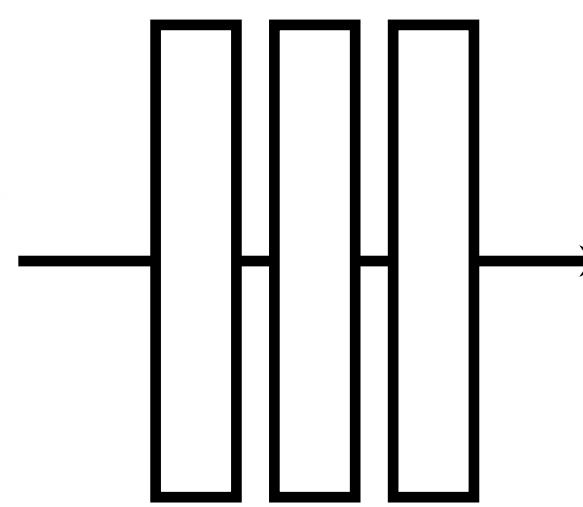


L2 regression

[Johnson, Alahi, Li, ECCV 2016]

# Designing objective functions

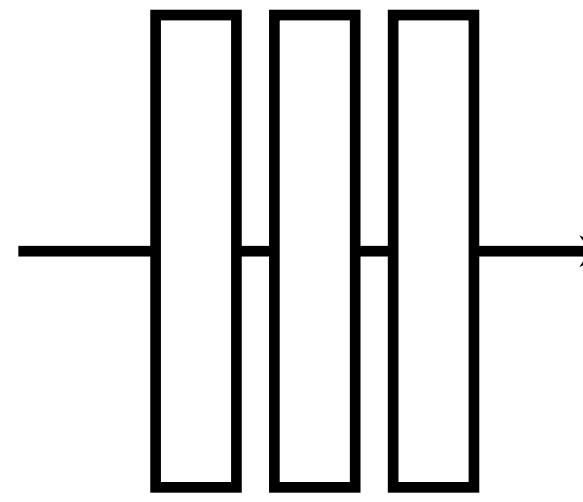
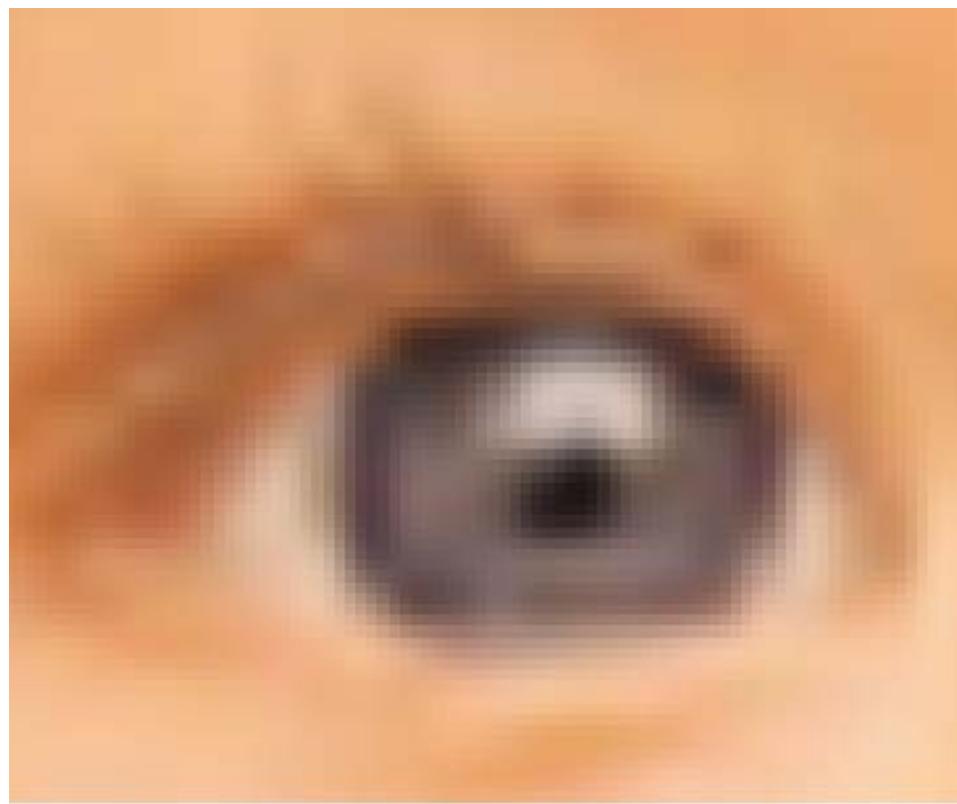
Image colorization



Cross entropy objective,  
with colorfulness term

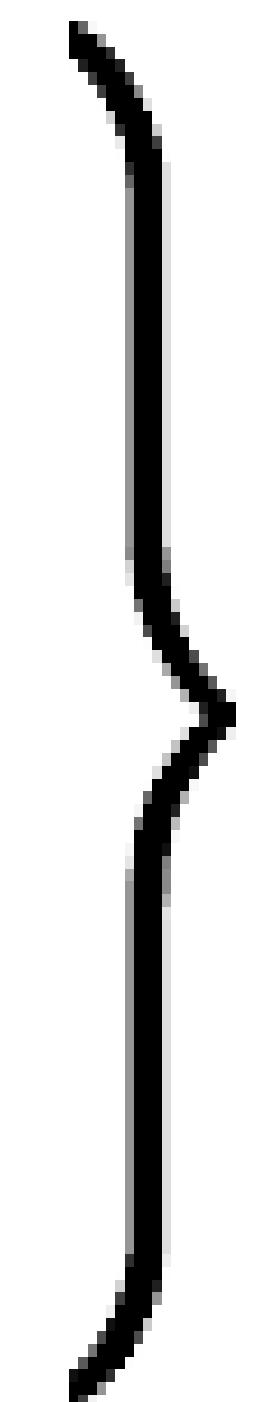
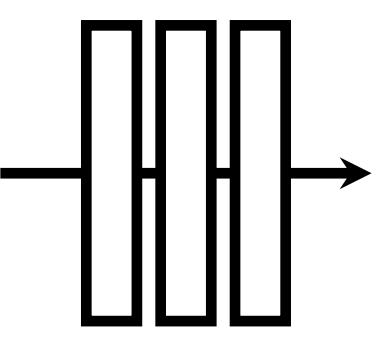
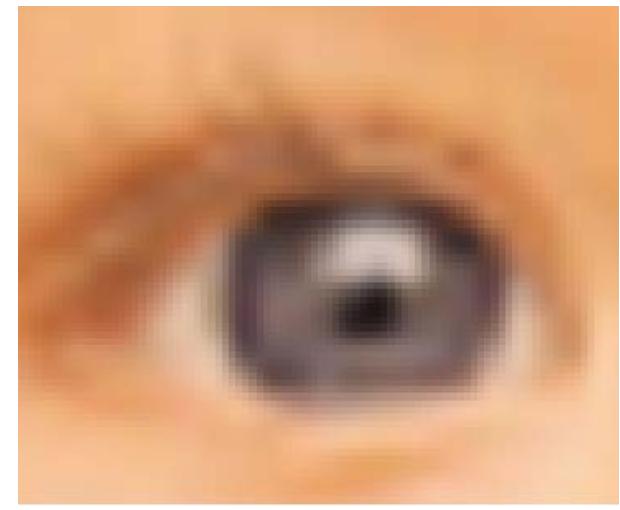
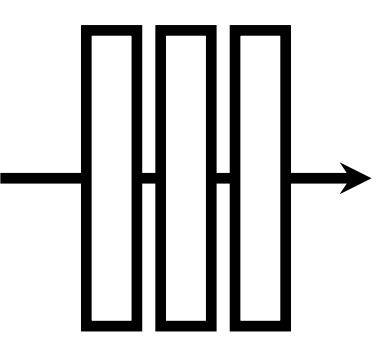
[Zhang, Isola, Efros, ECCV 2016]

Super-resolution



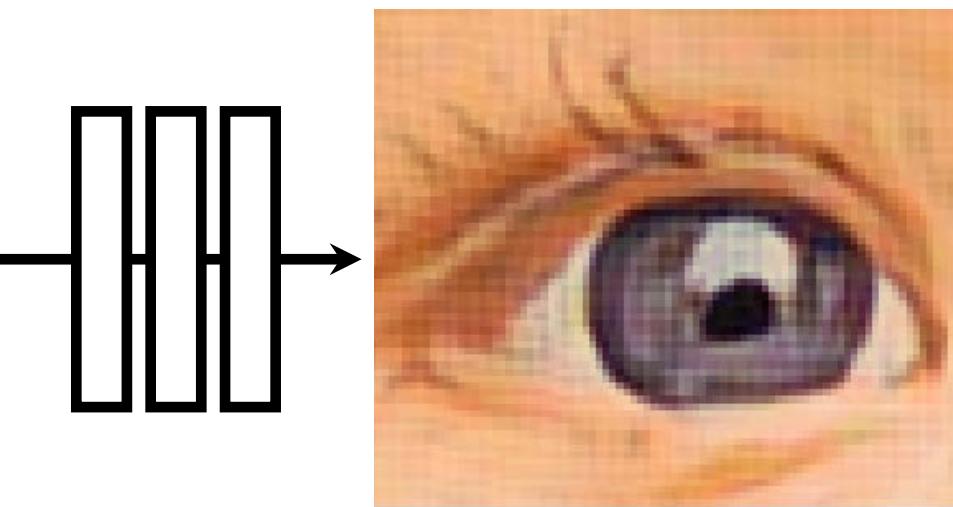
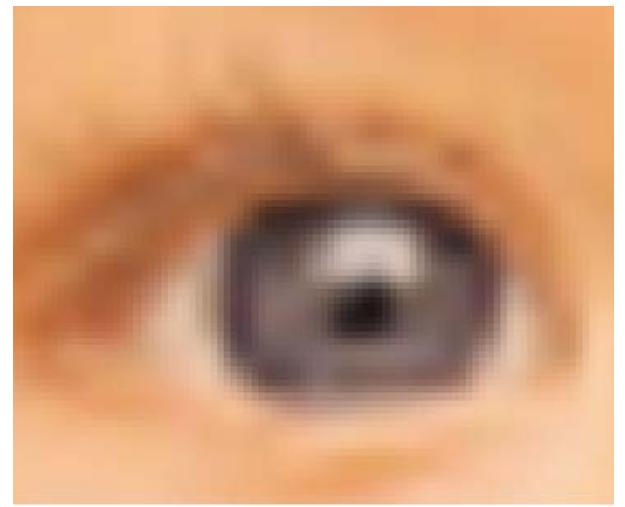
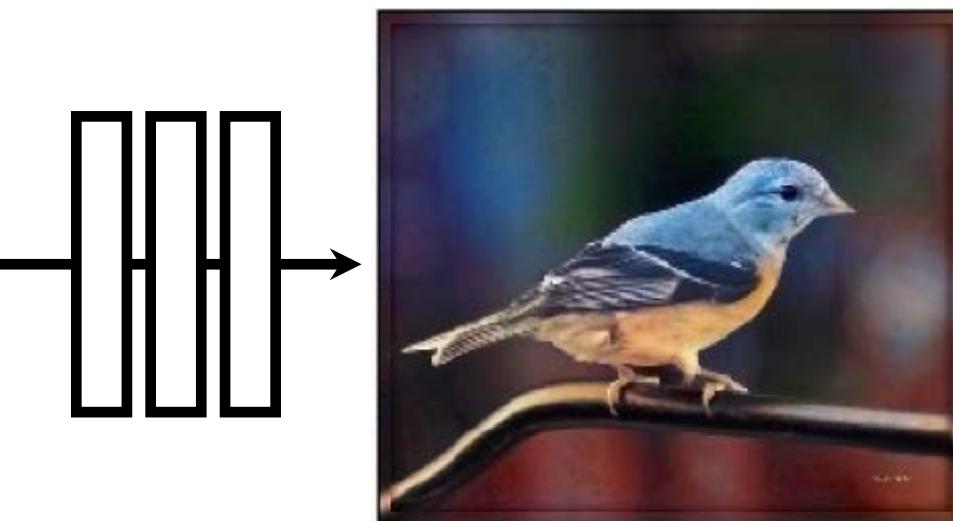
Deep feature covariance  
matching objective

[Johnson, Alahi, Li, ECCV 2016]



Universal loss?

Generated images



:

:

:

# “Generative Adversarial Network” (GANs)

Real photos



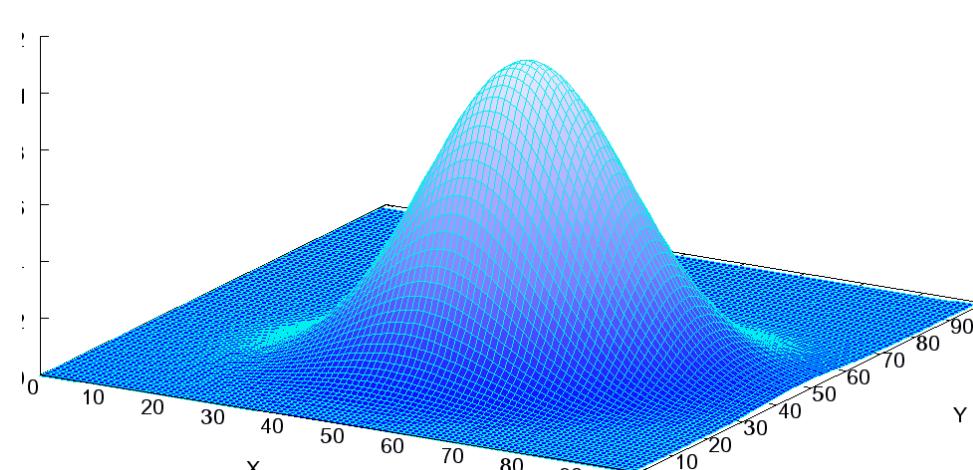
...

Generated  
vs Real  
(classifier)



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

# GANs



**Z**



[Goodfellow et al., 2014]

# Conditional GANs

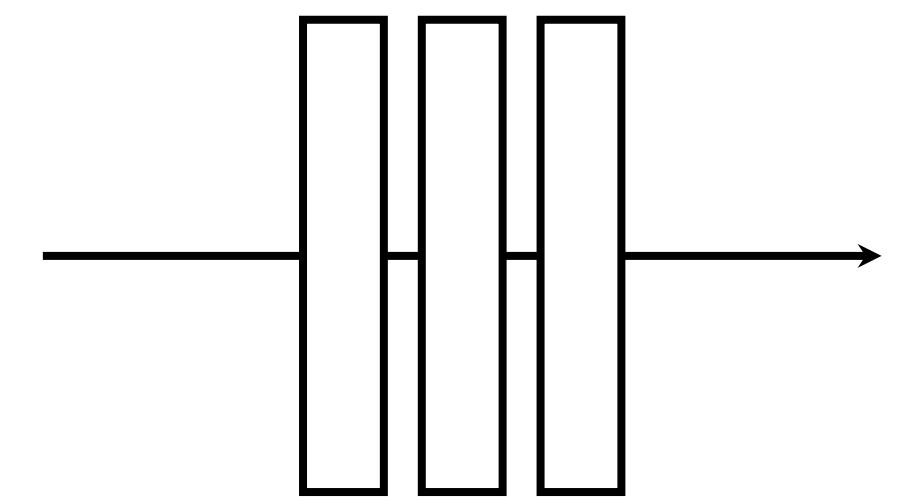


[Goodfellow et al., 2014]  
[Isola et al., 2016]

**x**

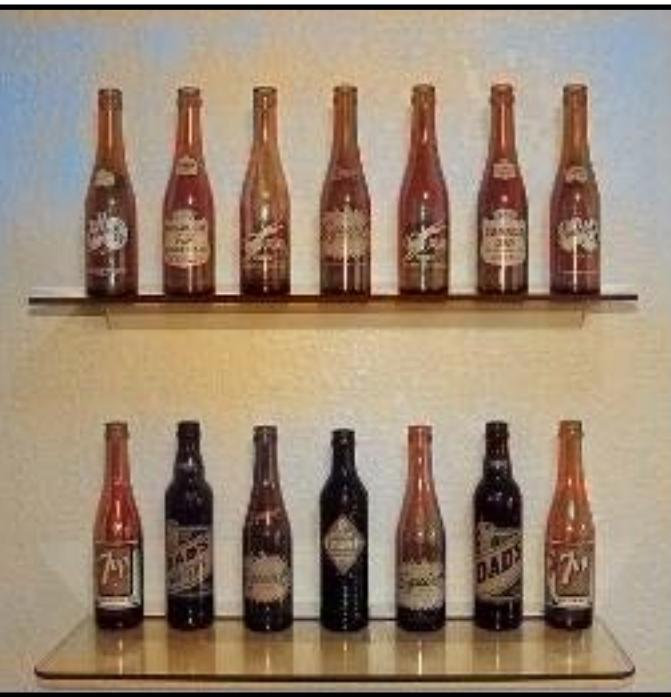


**G**



**Generator**

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

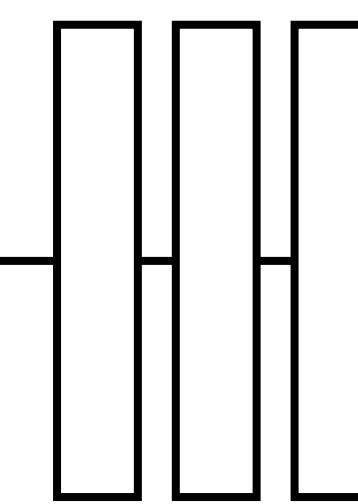


[Goodfellow et al., 2014]

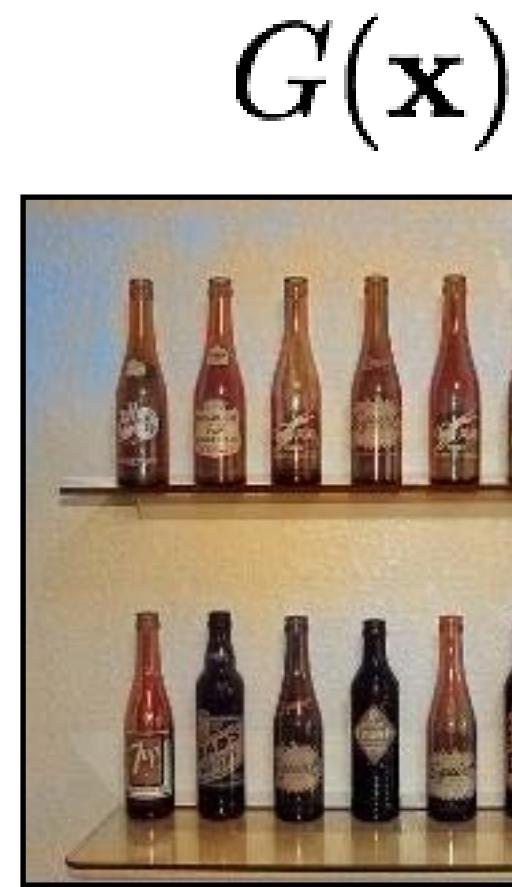


**x**

**G**

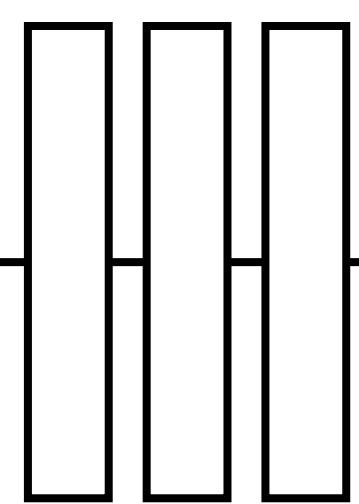


Generator



**G(x)**

**D**

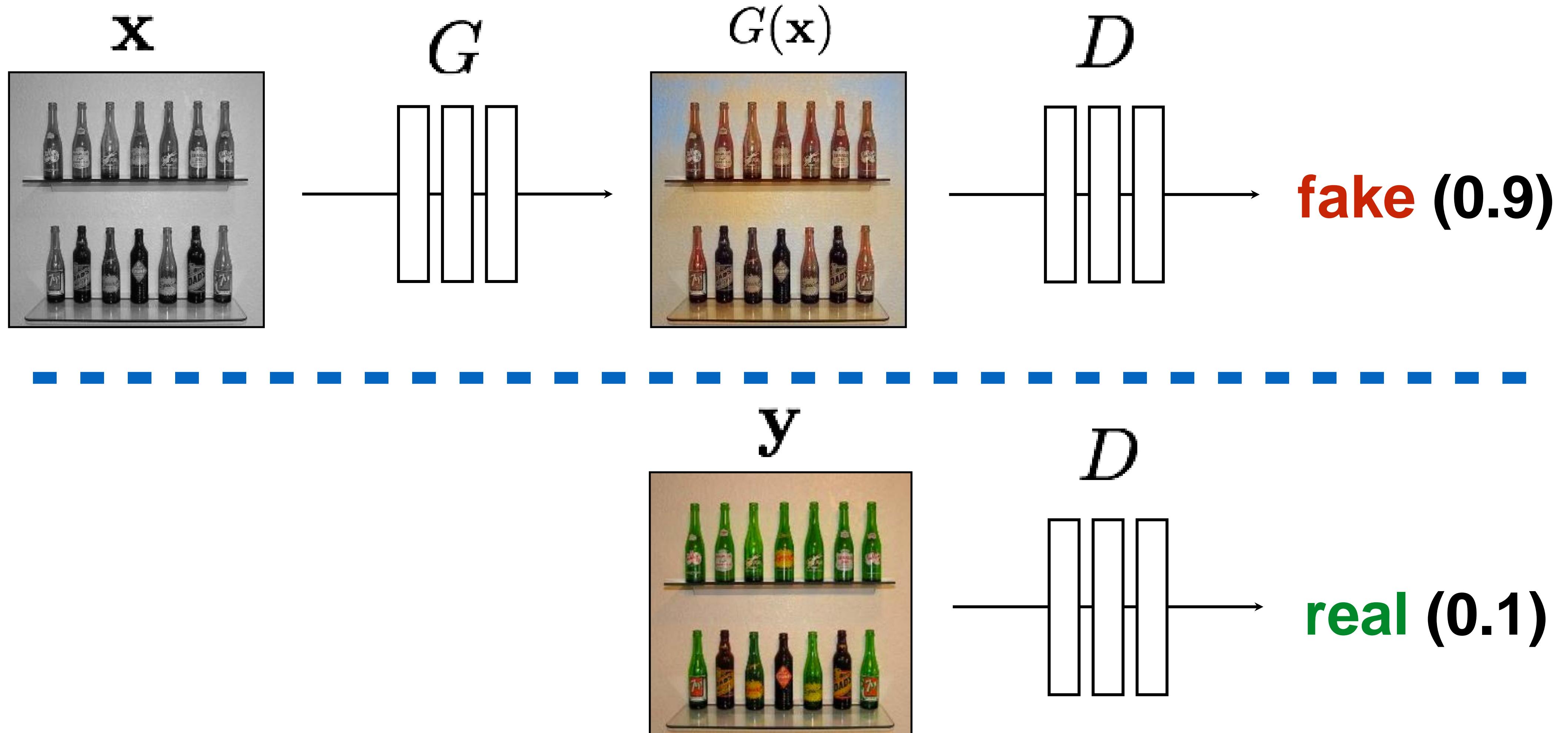


Discriminator

real or  
fake?

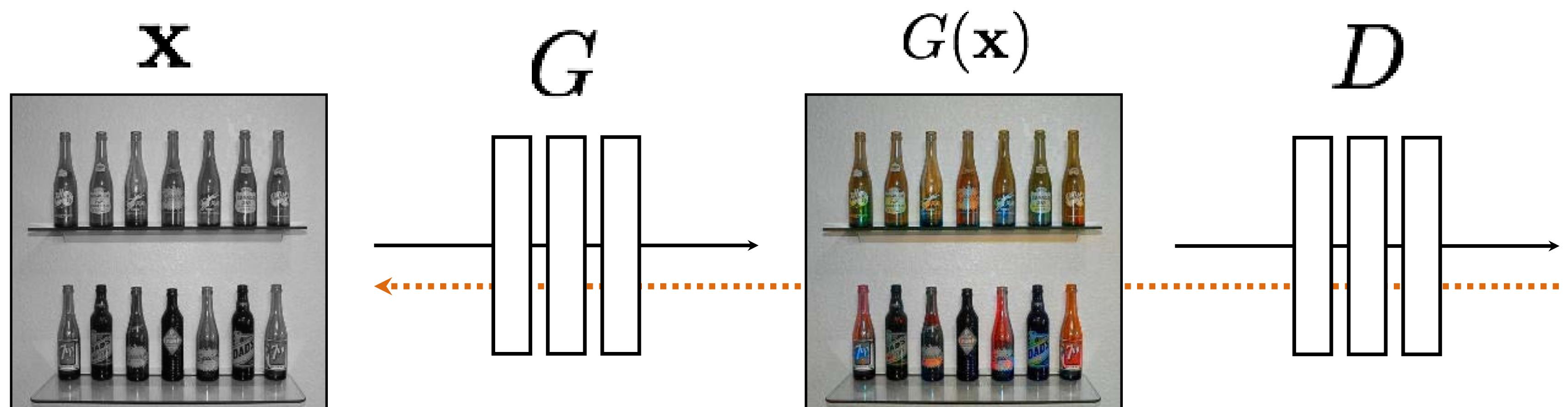
**G tries to synthesize fake images that fool D**

**D tries to identify the fakes**



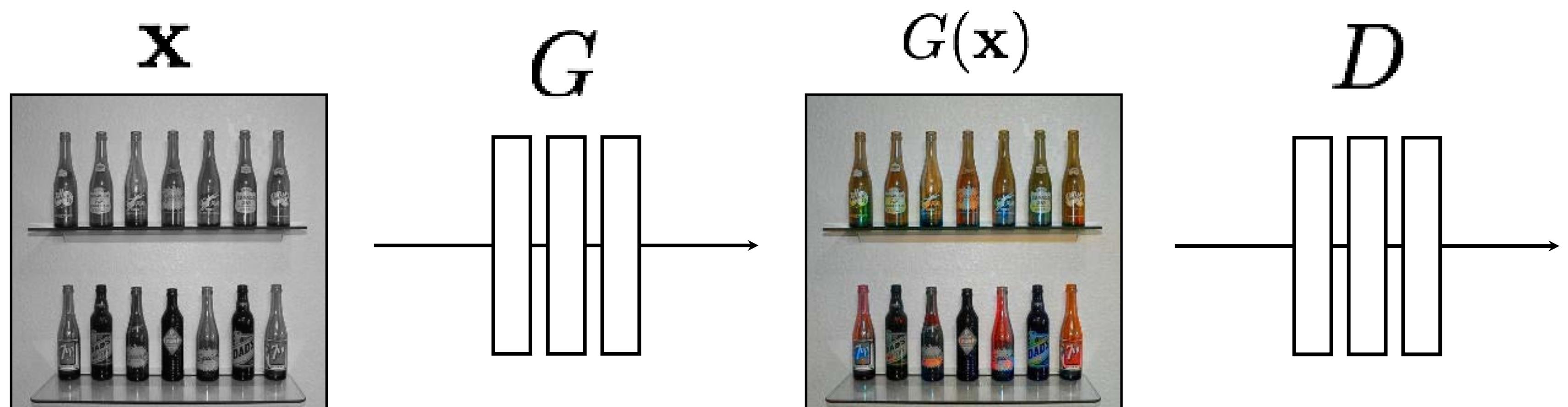
$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \boxed{\log D(G(\mathbf{x}))} + \boxed{\log(1 - D(\mathbf{y}))} ]$$

[Goodfellow et al., 2014]



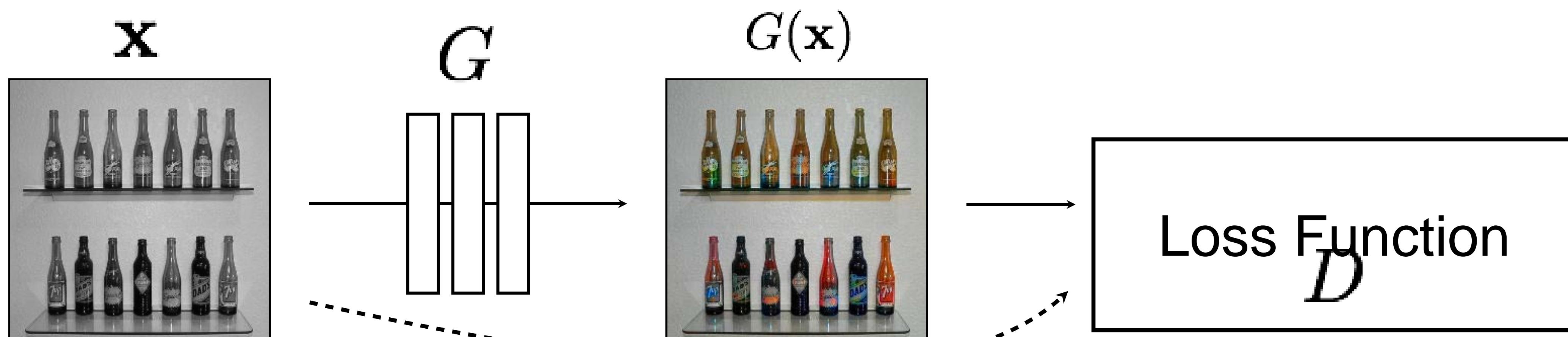
**G tries to synthesize fake images that *fool*  
D:**

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



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

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

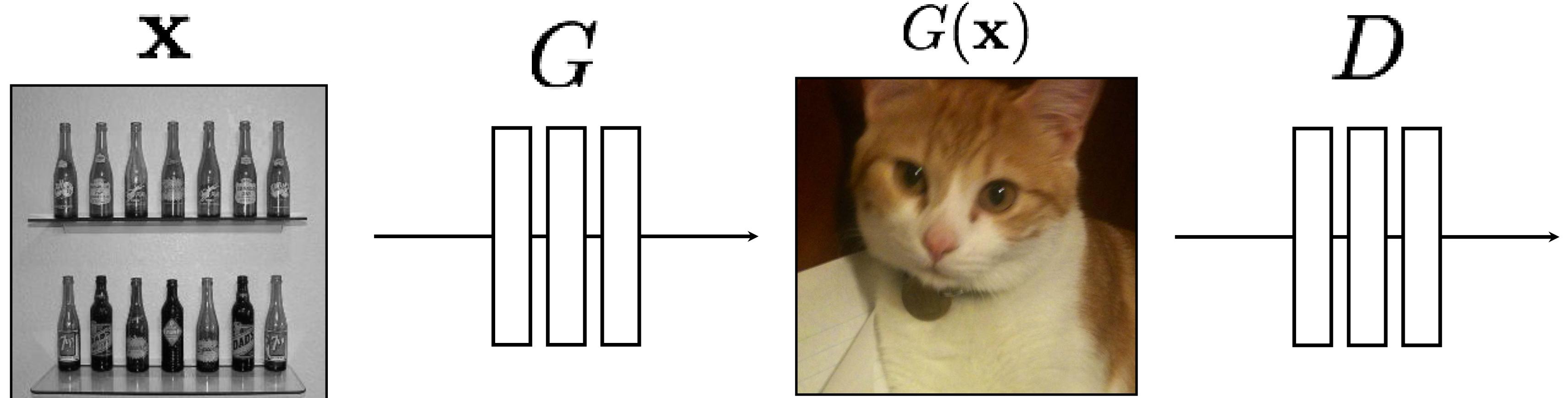


**G's perspective:** **D** is a loss function.

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

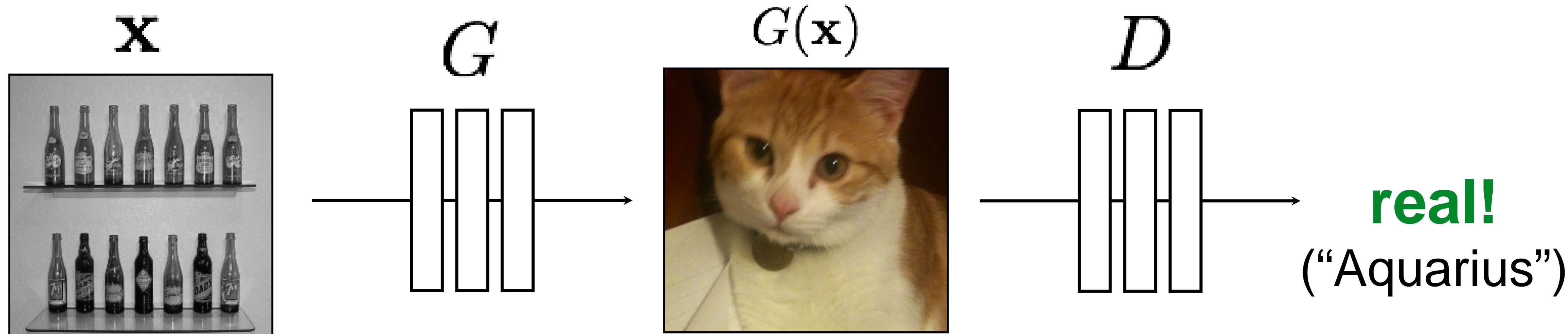
[Goodfellow et al., 2014]

[Isola et al., 2016]



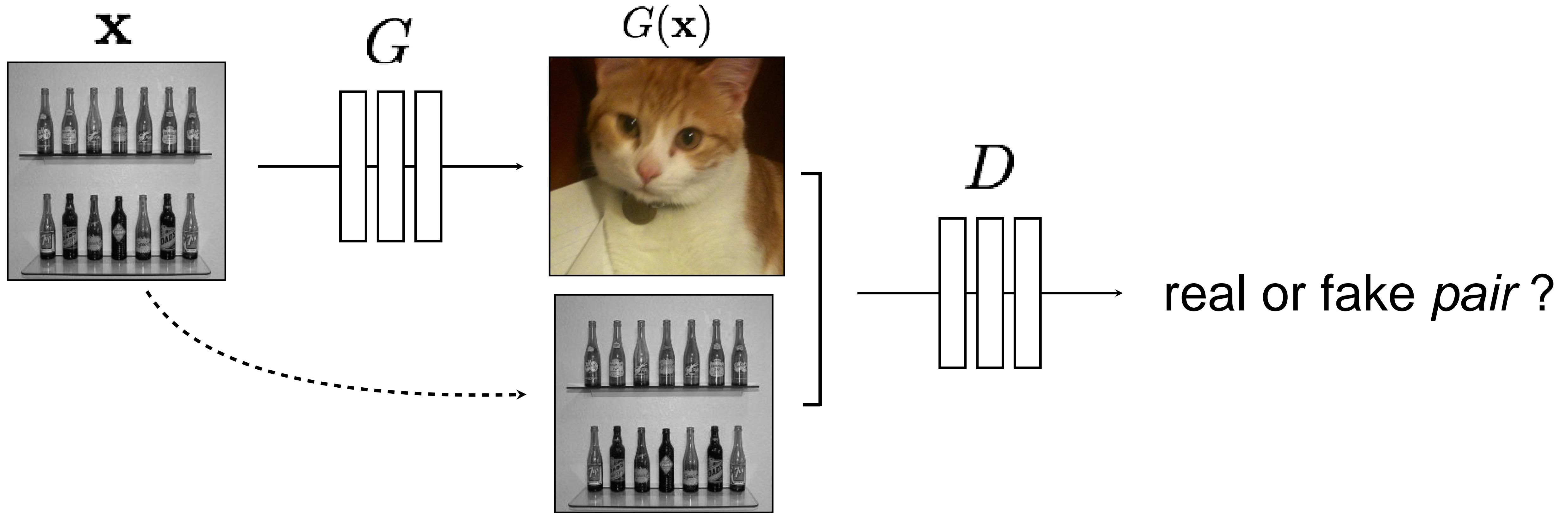
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

[Goodfellow et al., 2014]



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

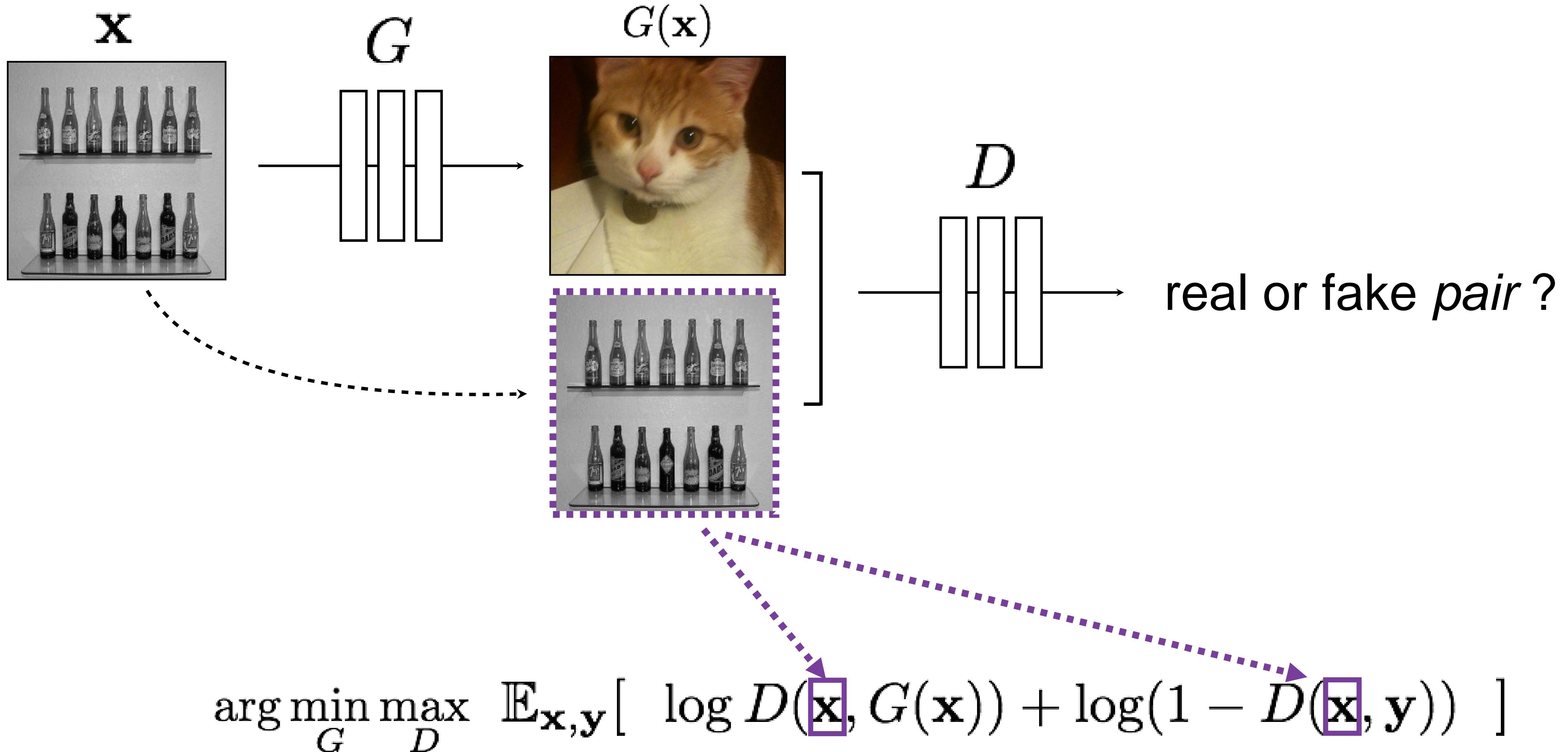
[Goodfellow et al., 2014]



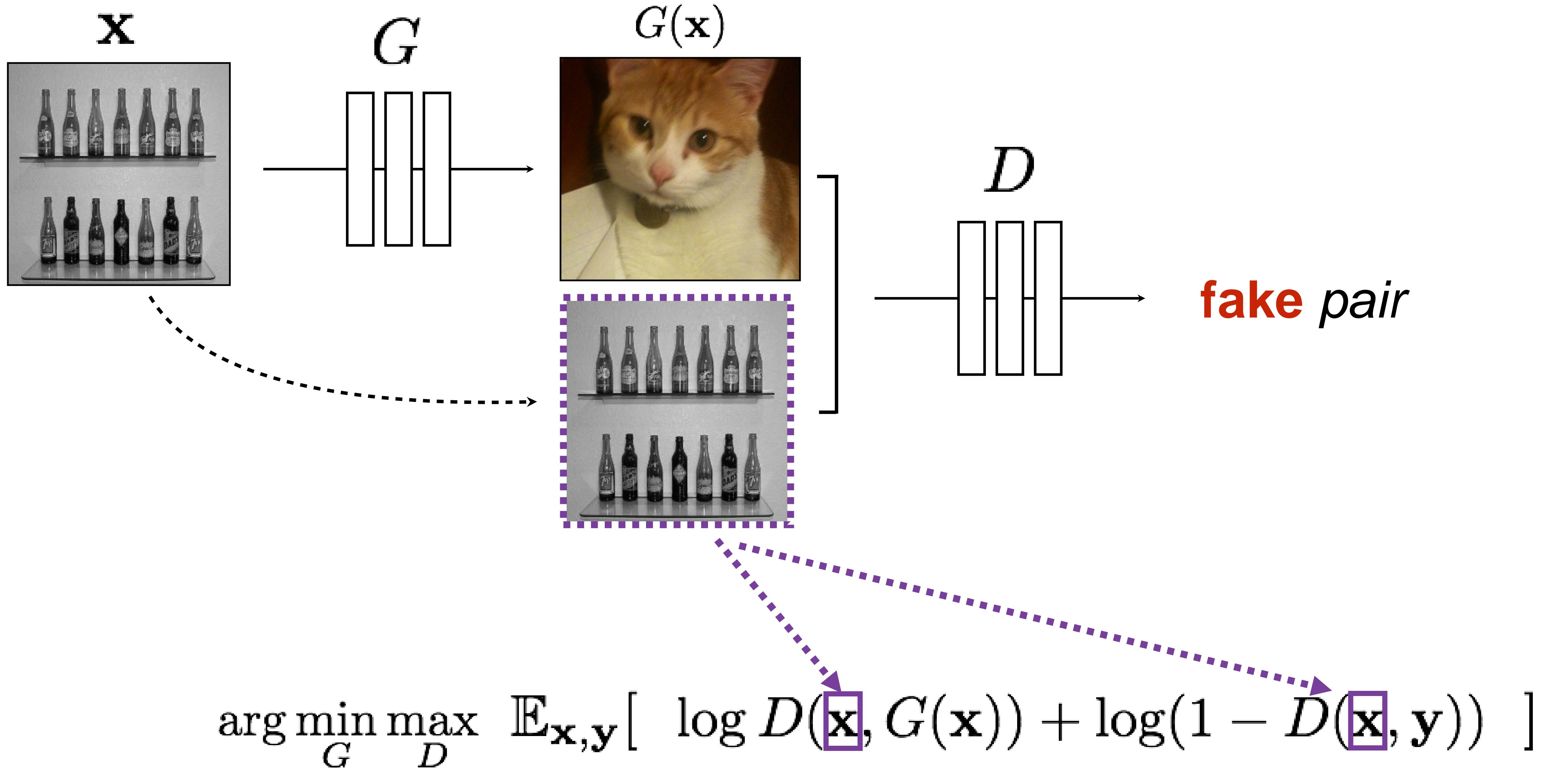
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

[Goodfellow et al., 2014]

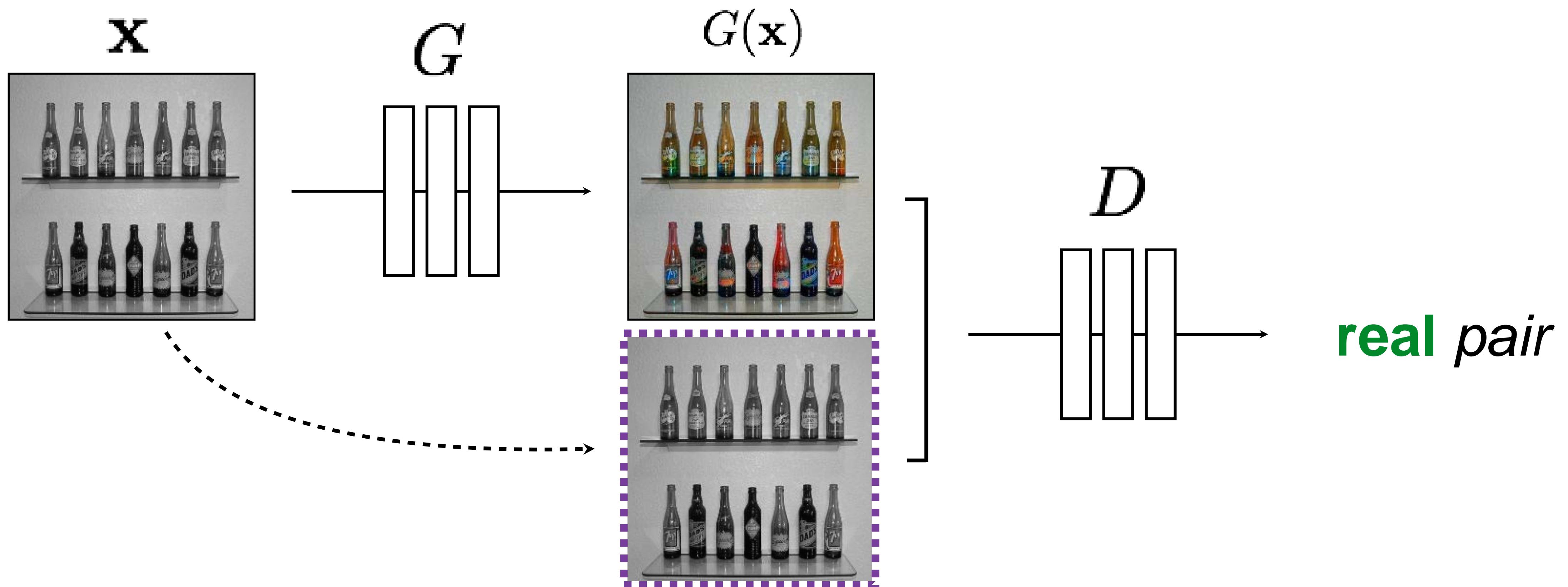
[Isola et al., 2016]



[Goodfellow et al., 2014]  
 [Isola et al., 2016]

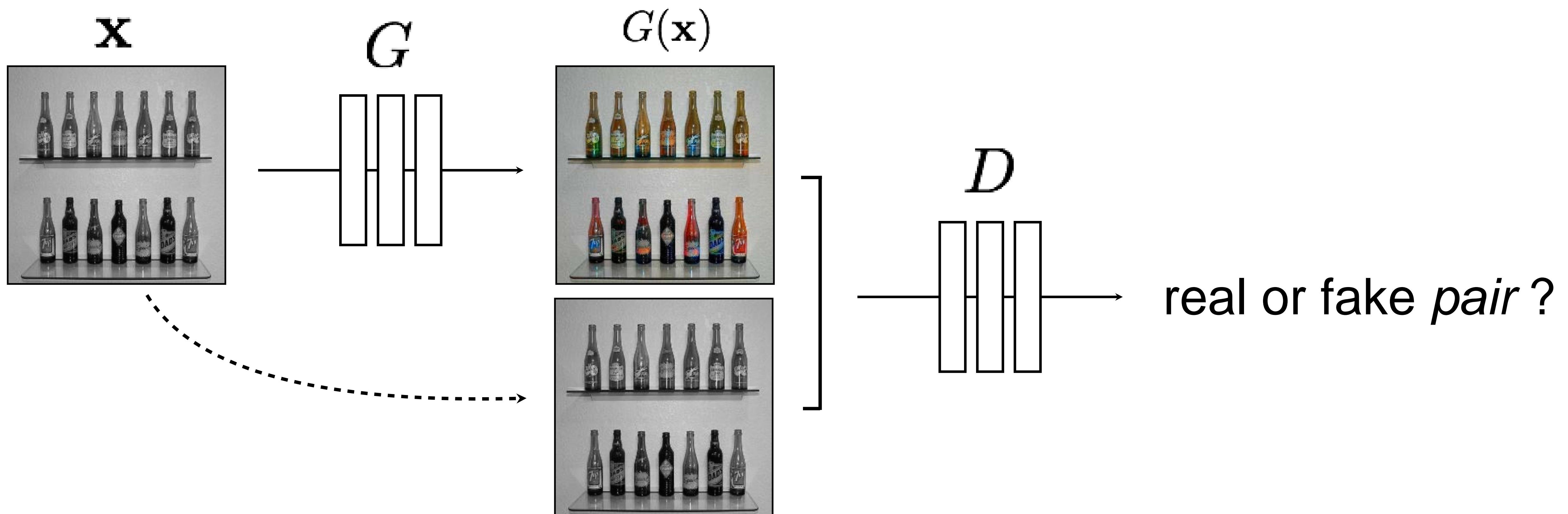


[Goodfellow et al., 2014]  
 [Isola et al., 2016]



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

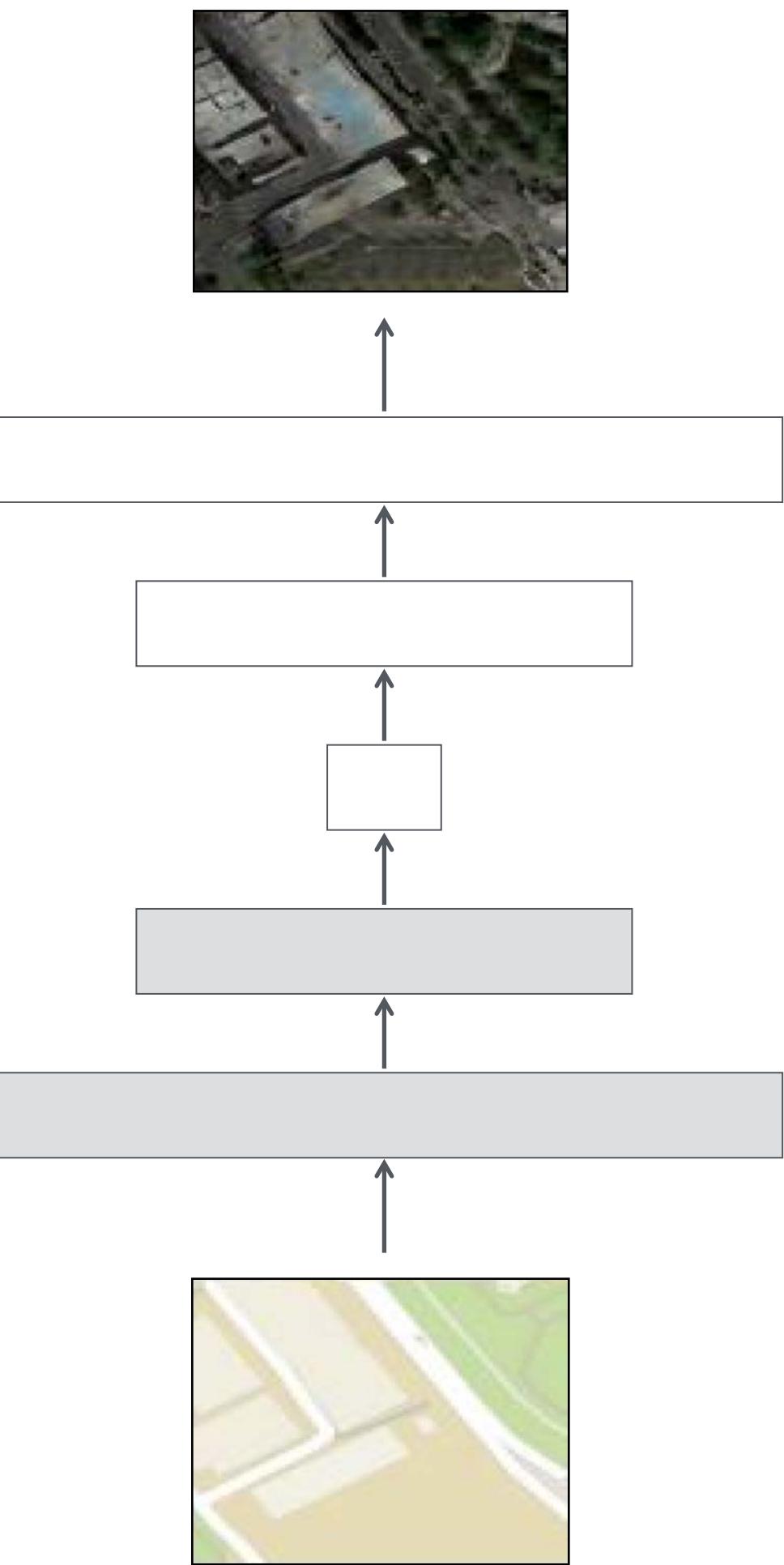
[Goodfellow et al., 2014]  
 [Isola et al., 2016]



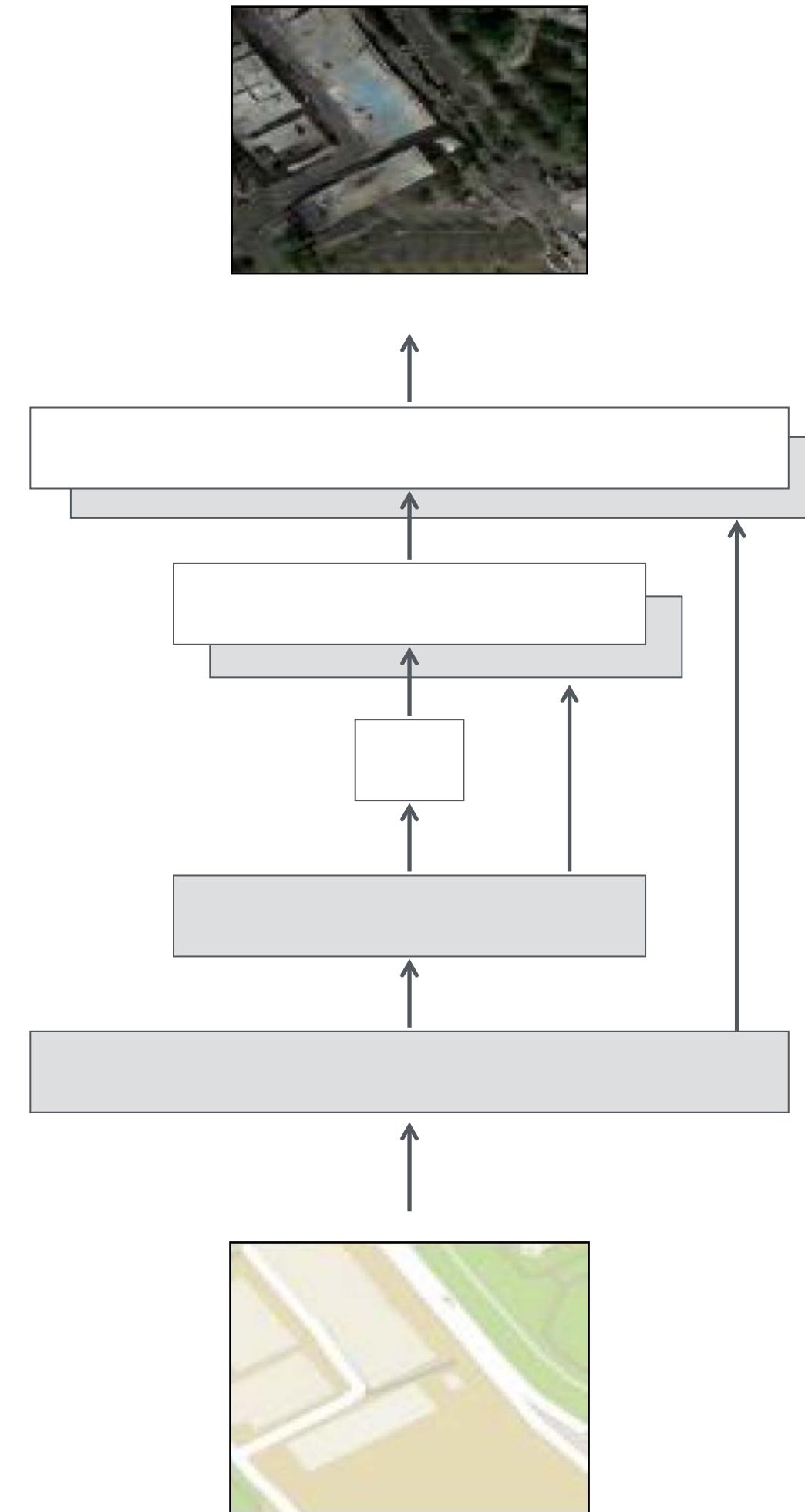
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$

[Goodfellow et al., 2014]  
 [Isola et al., 2016]

# Generator architecture



Encoder-decoder



U-Net

# BW → Color

Input



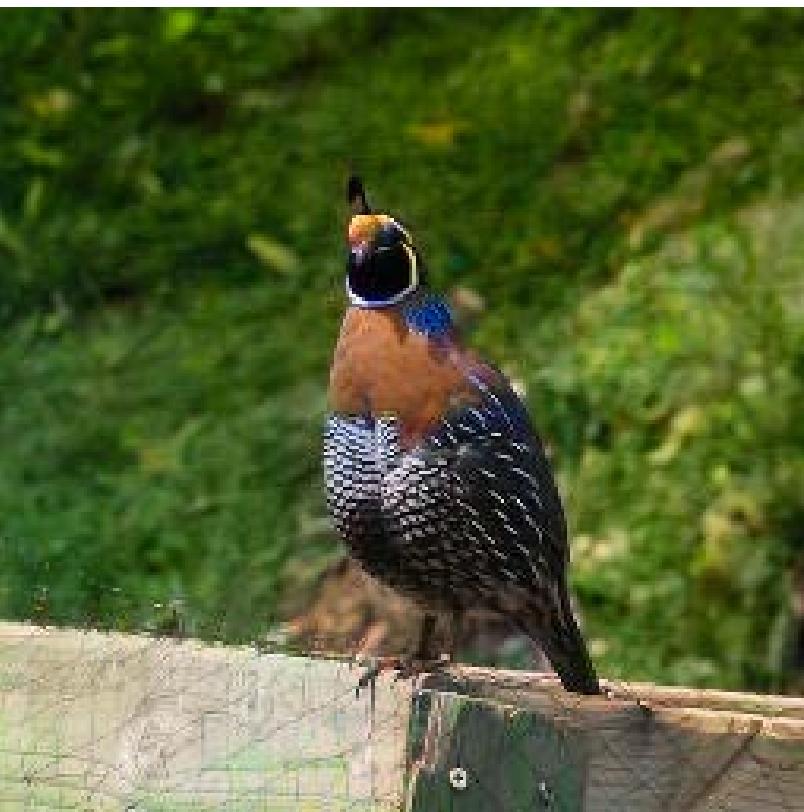
Output



Input



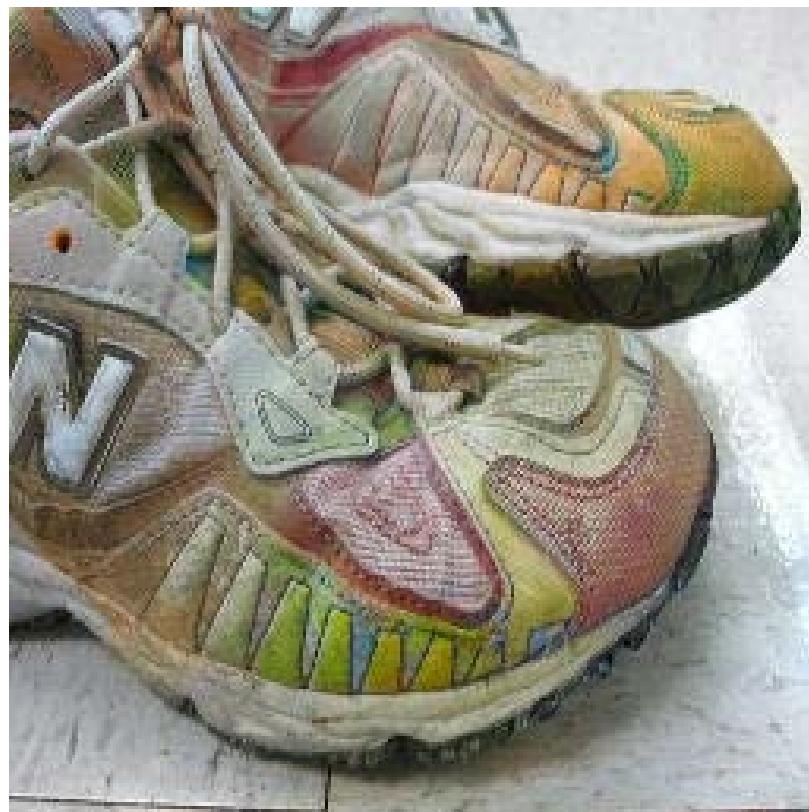
Output



Input

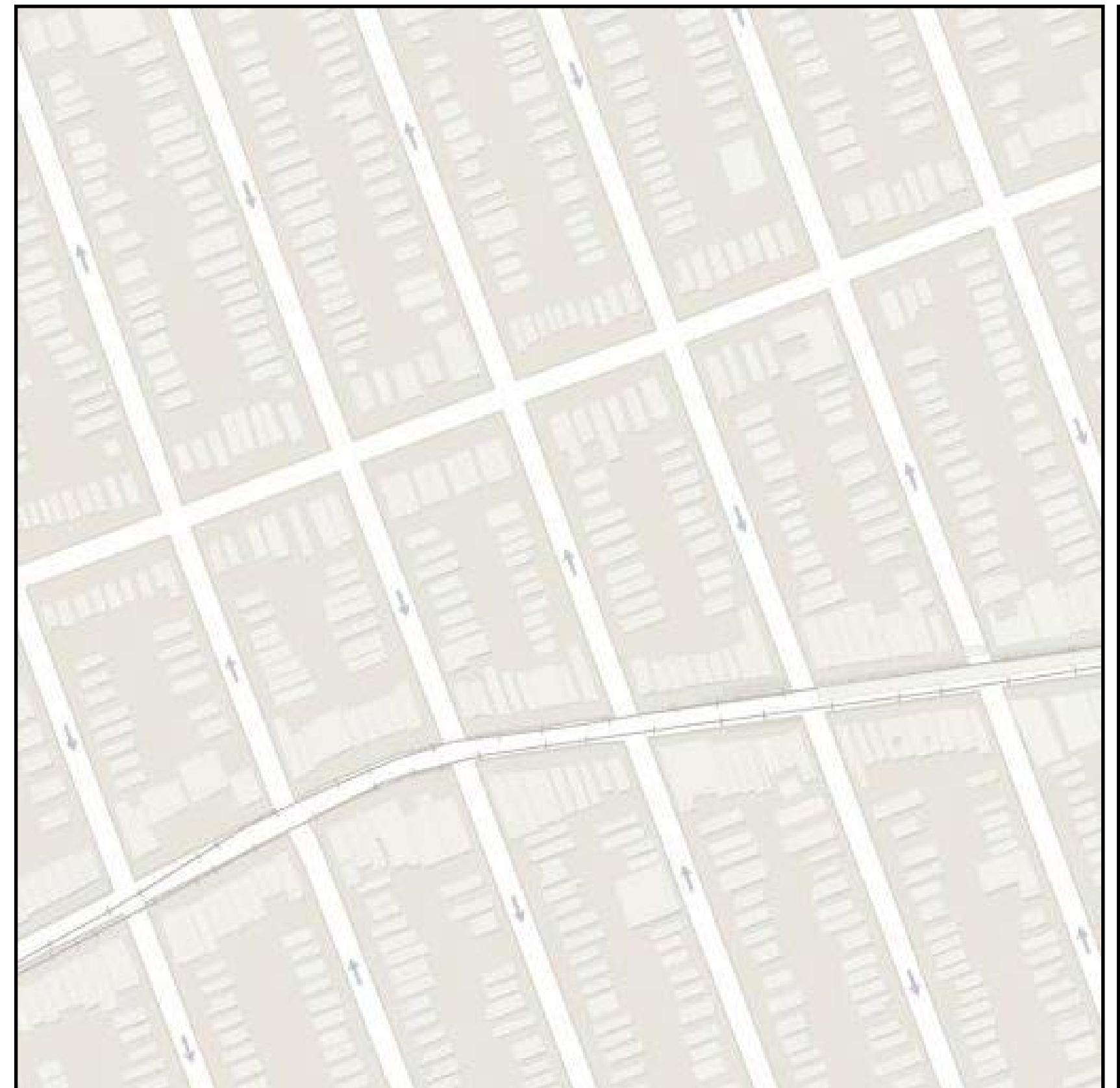


Output

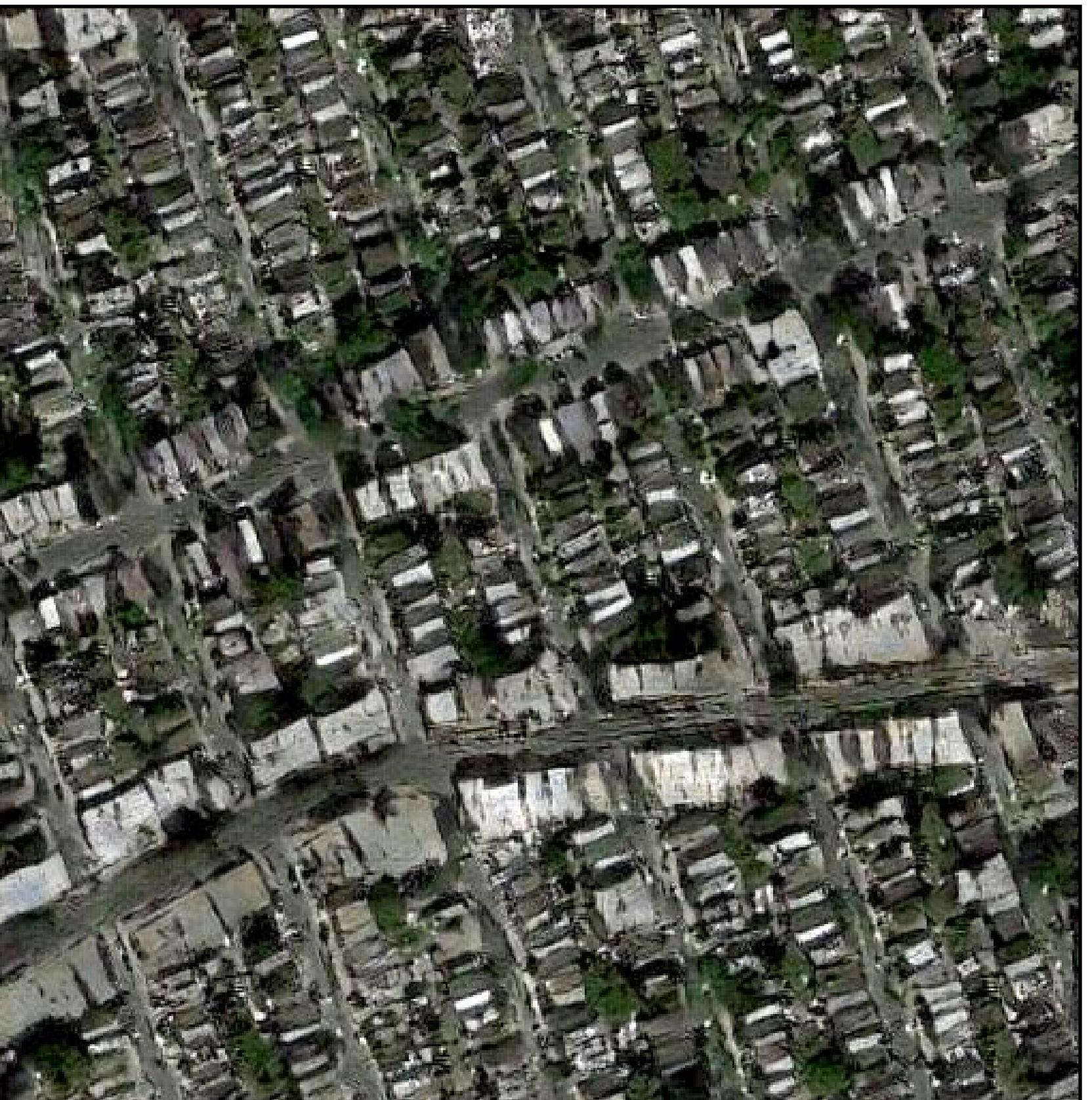


Data from [Russakovsky et al. 2015]

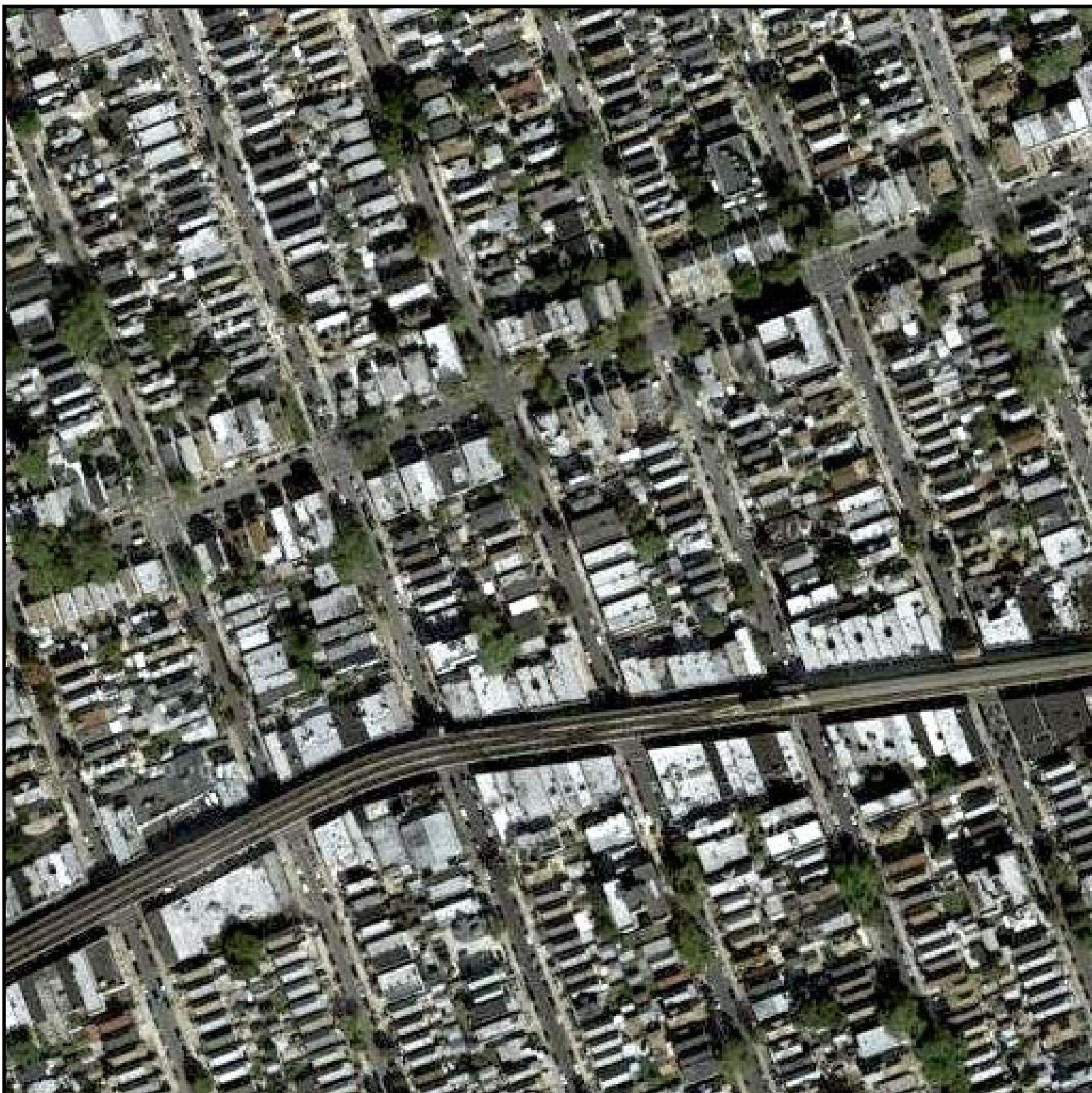
Input



Output



Groundtruth



Data from  
[\[maps.google.com\]](https://maps.google.com)

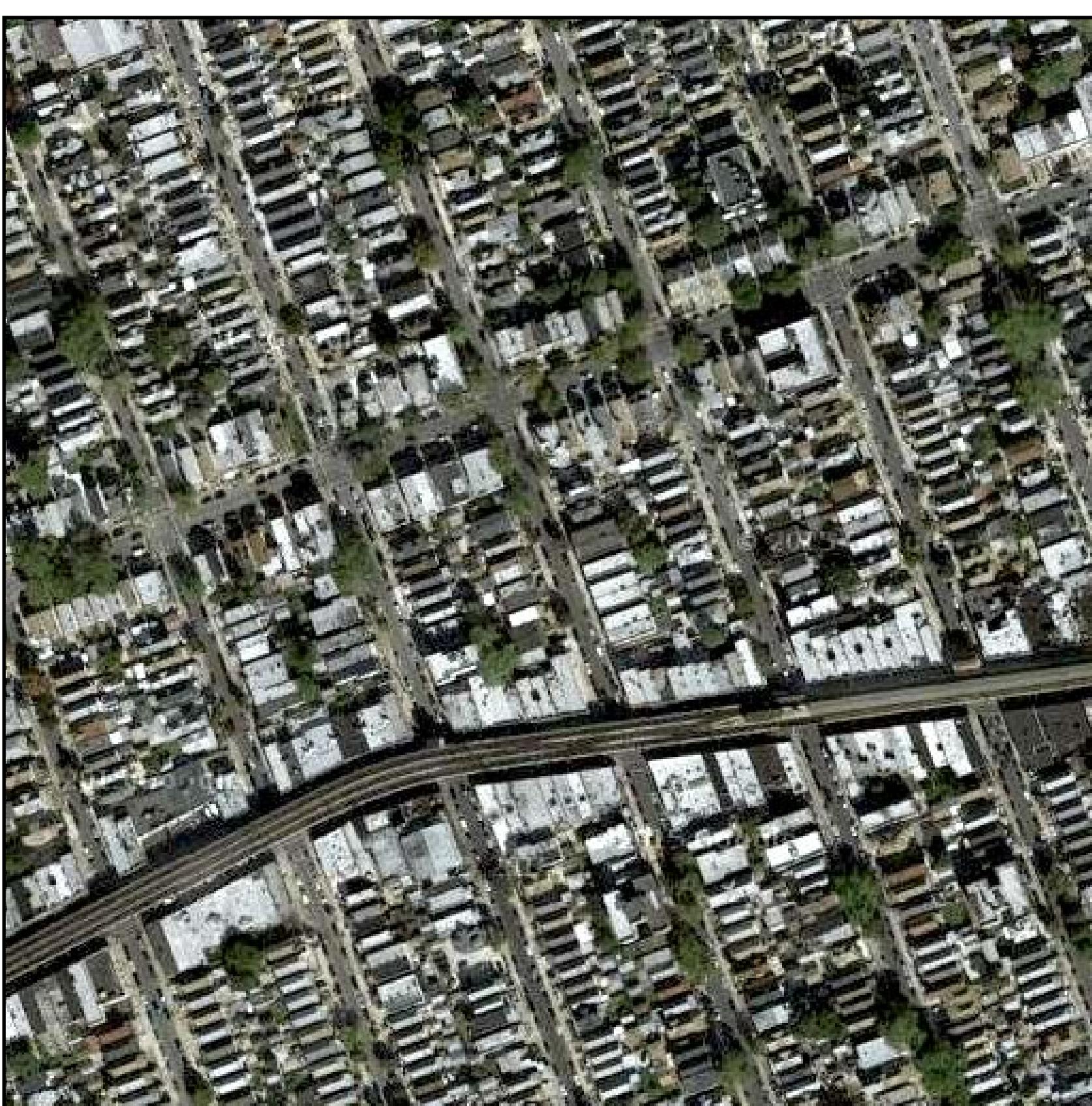


[Isola et al, 2016]

Input



Output

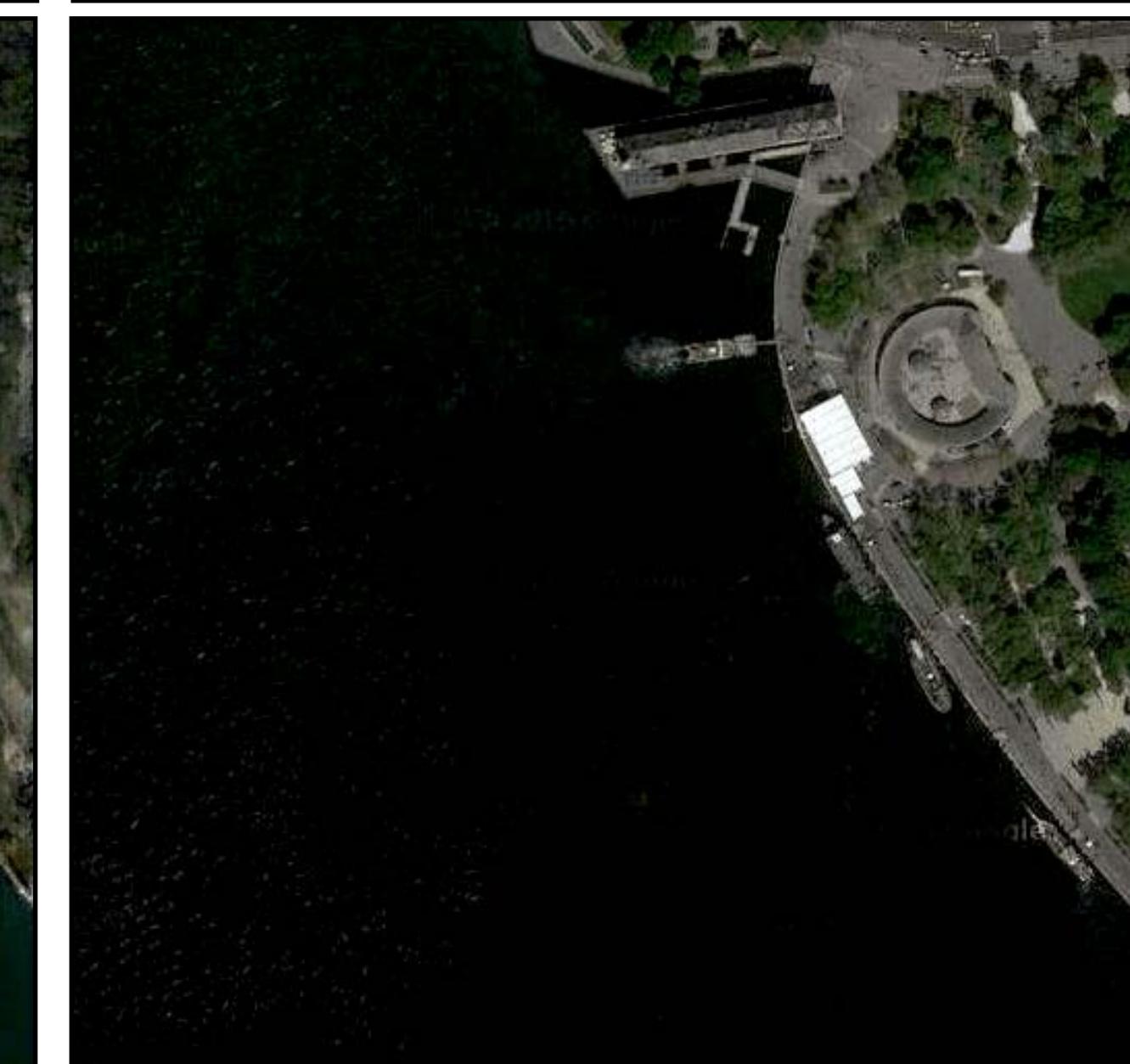
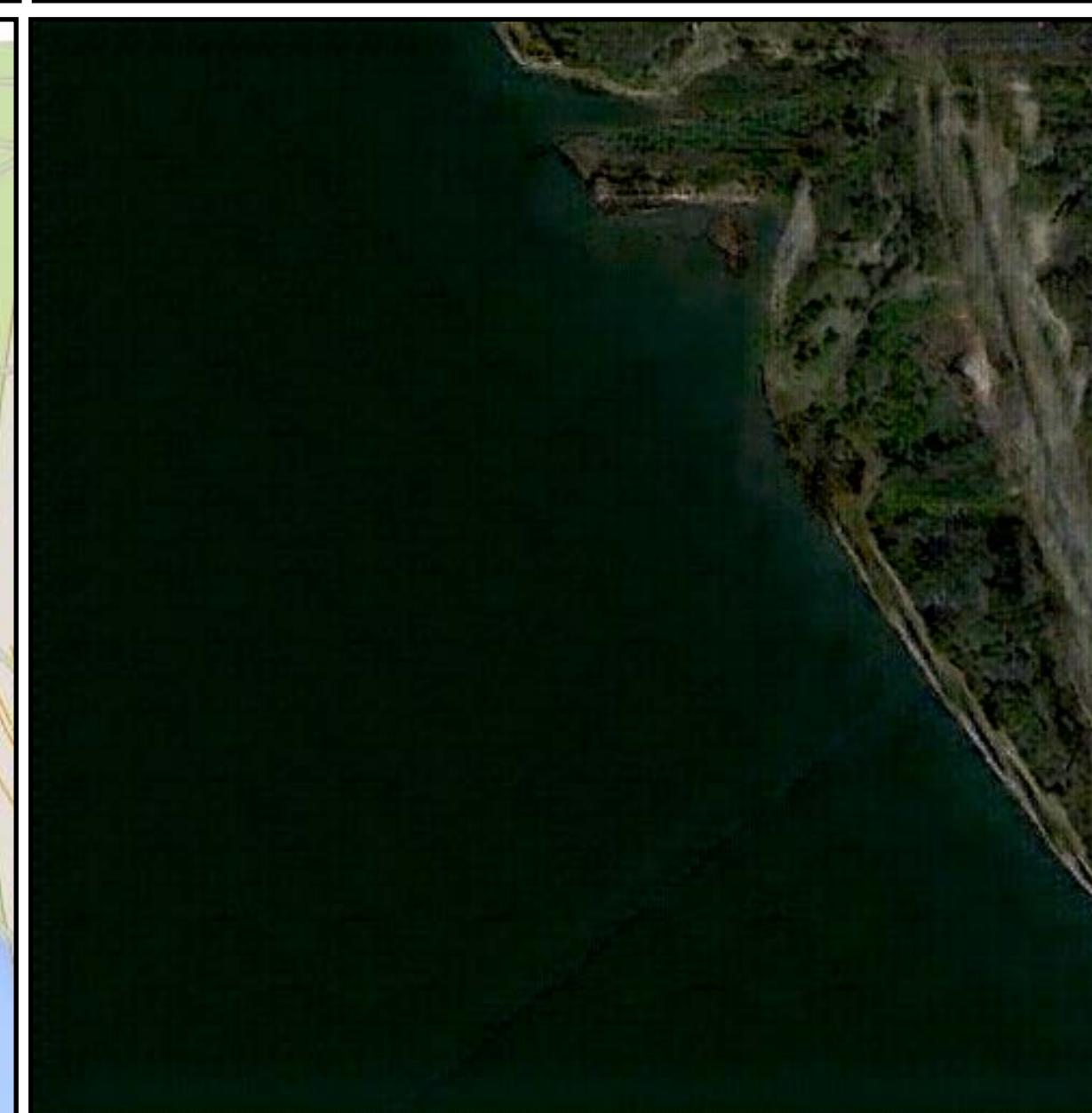
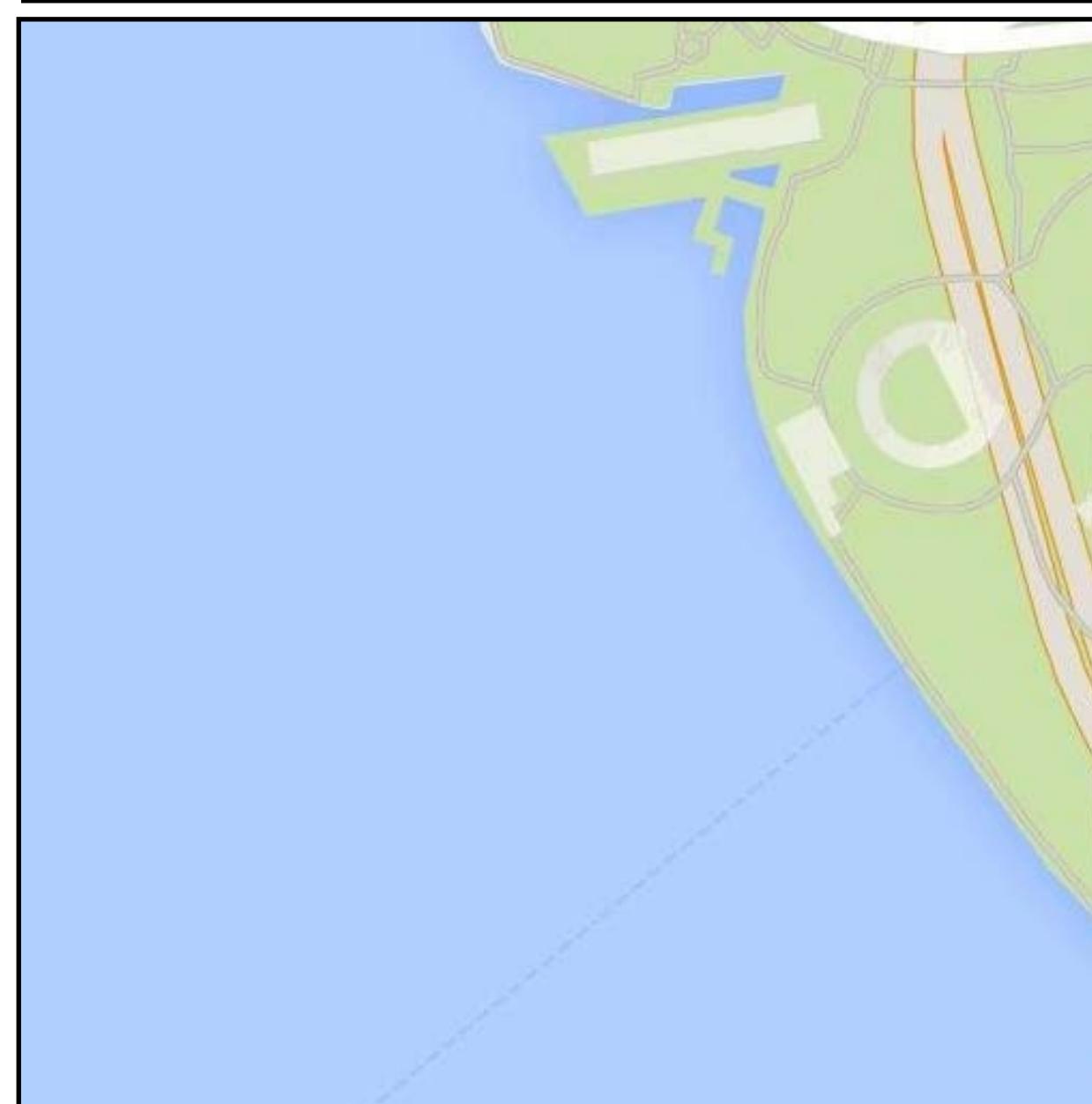
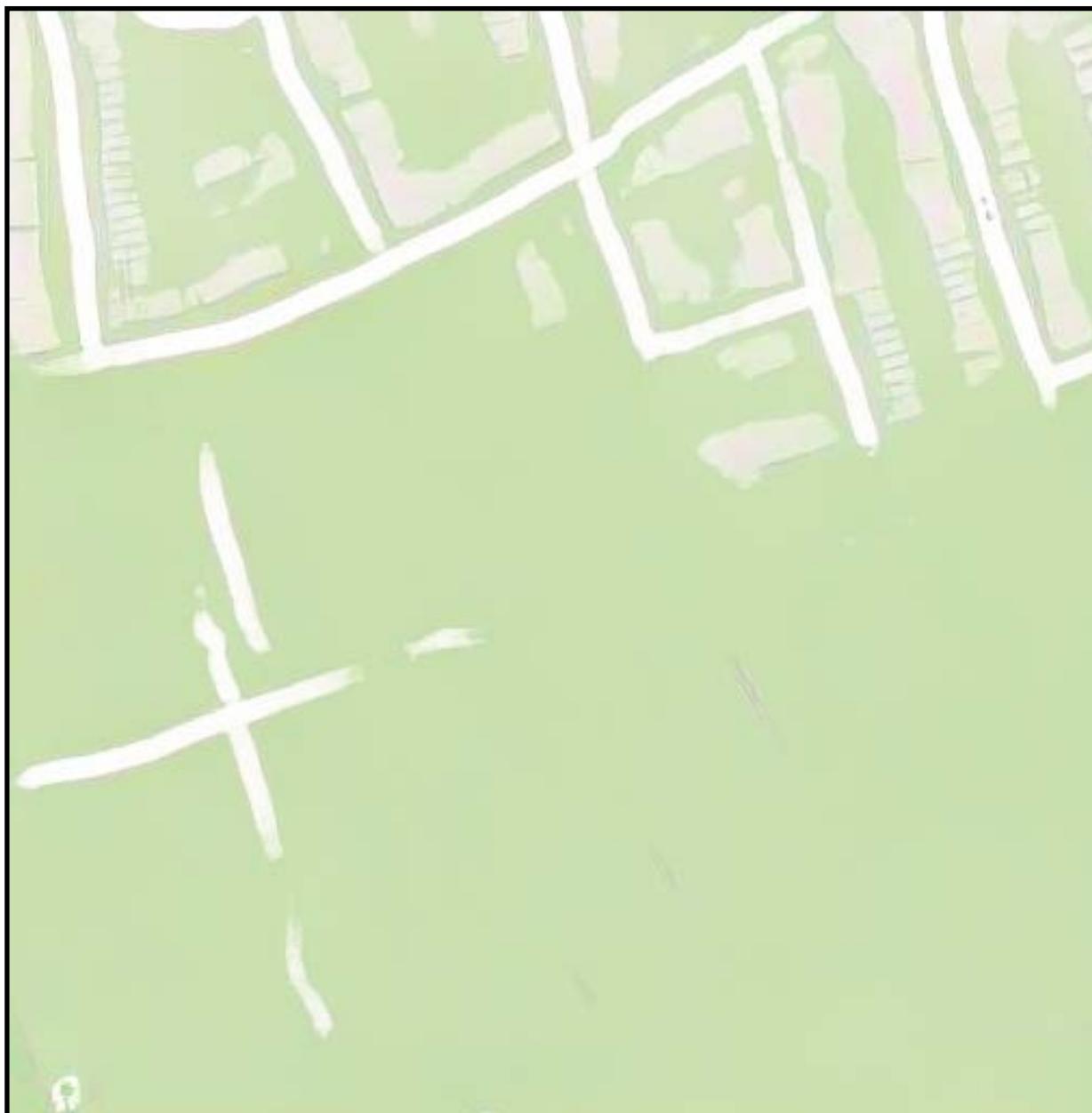
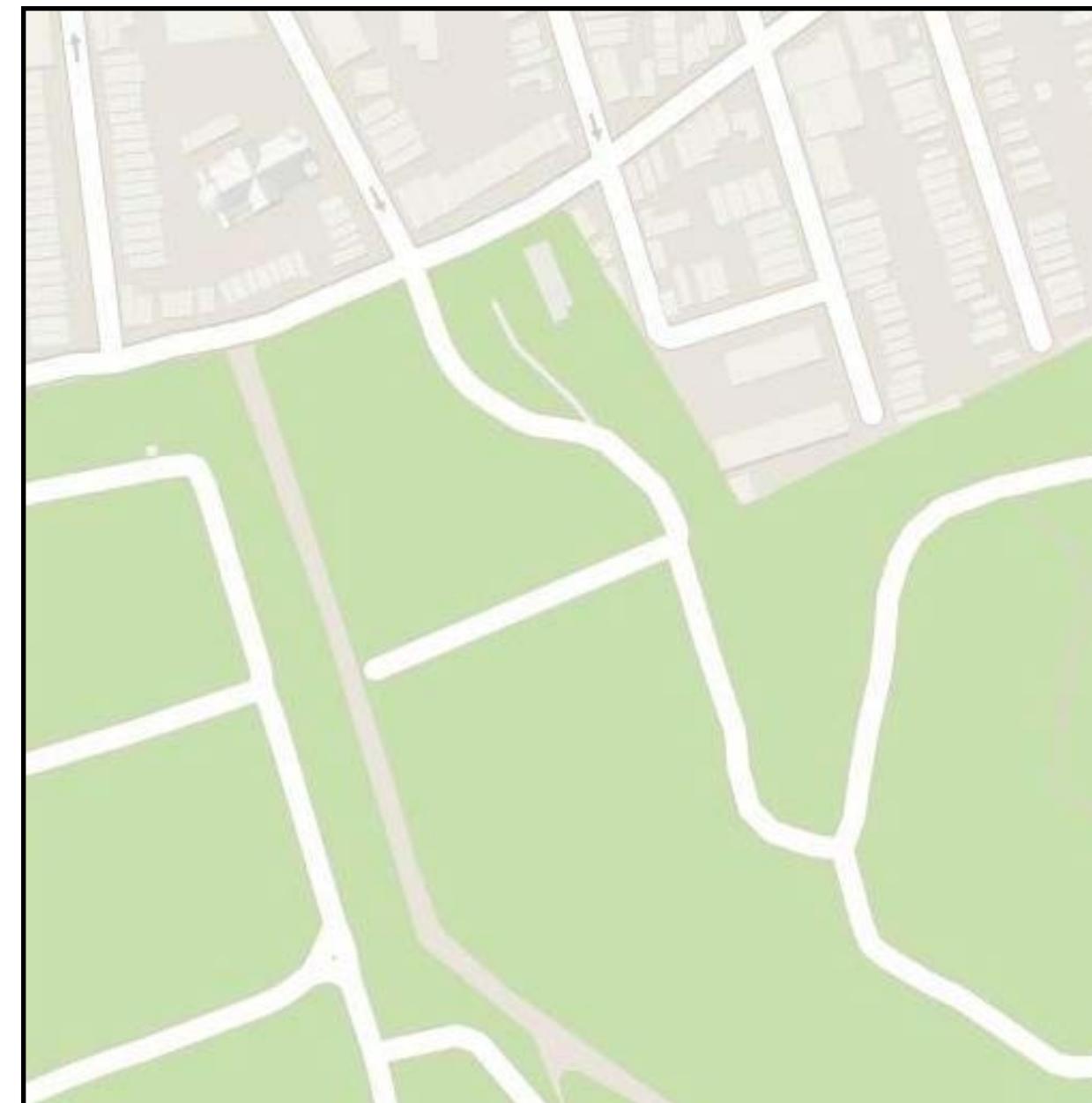


Groundtruth

Data from [[maps.google.com](https://maps.google.com)]

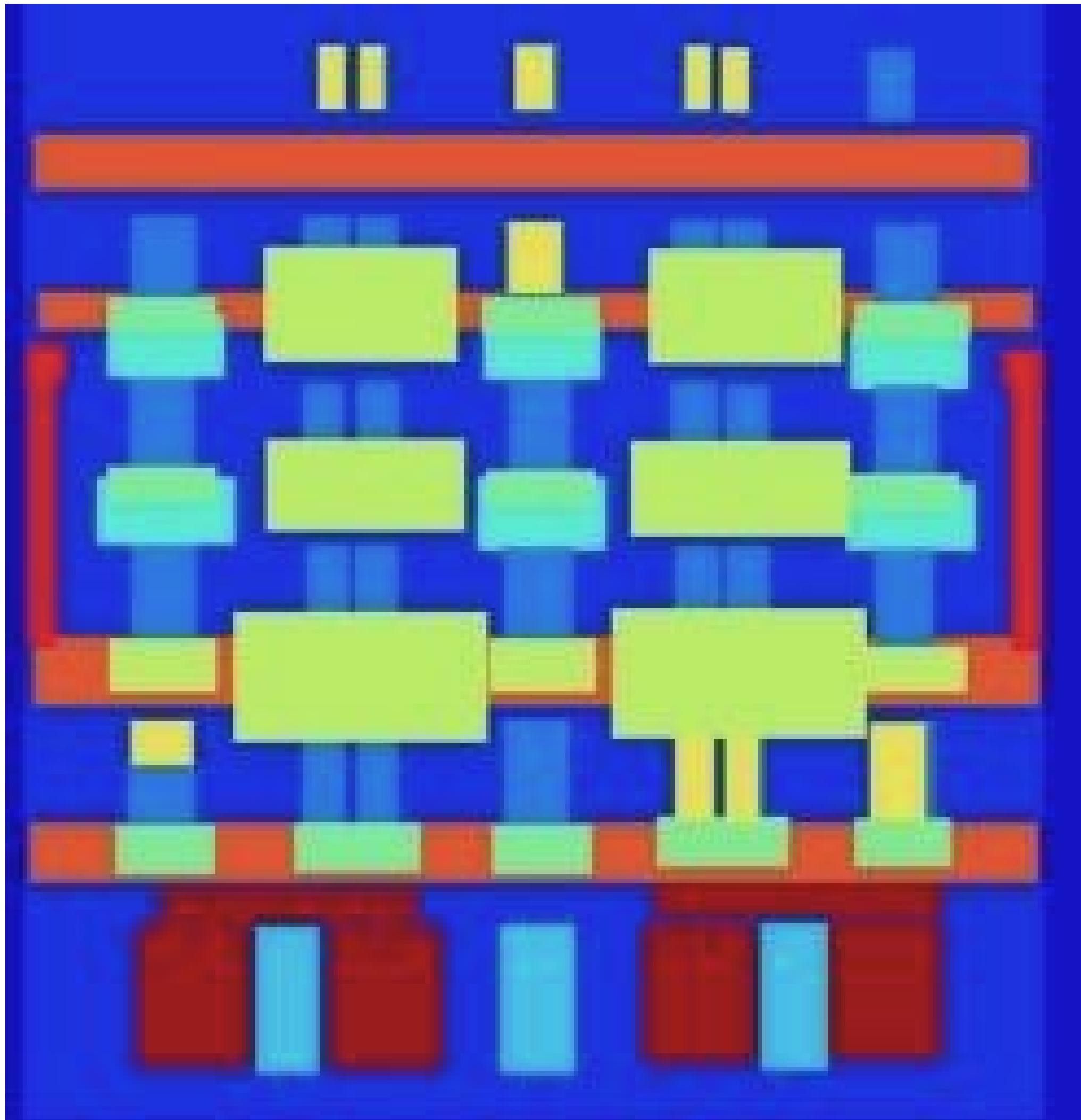
Map → Aerial Photo

Aerial Photo → Map

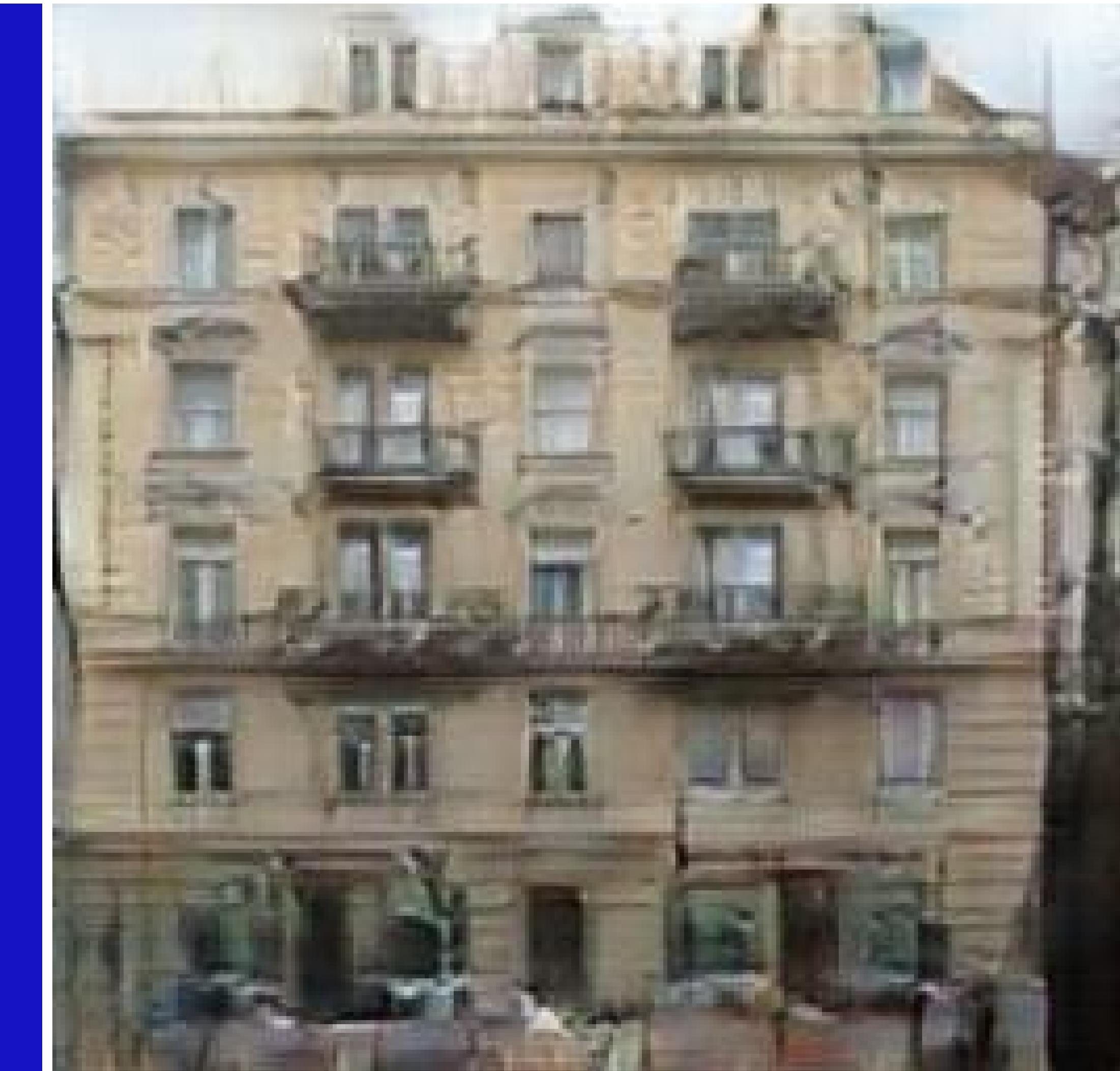


# Labels → Facades

Input



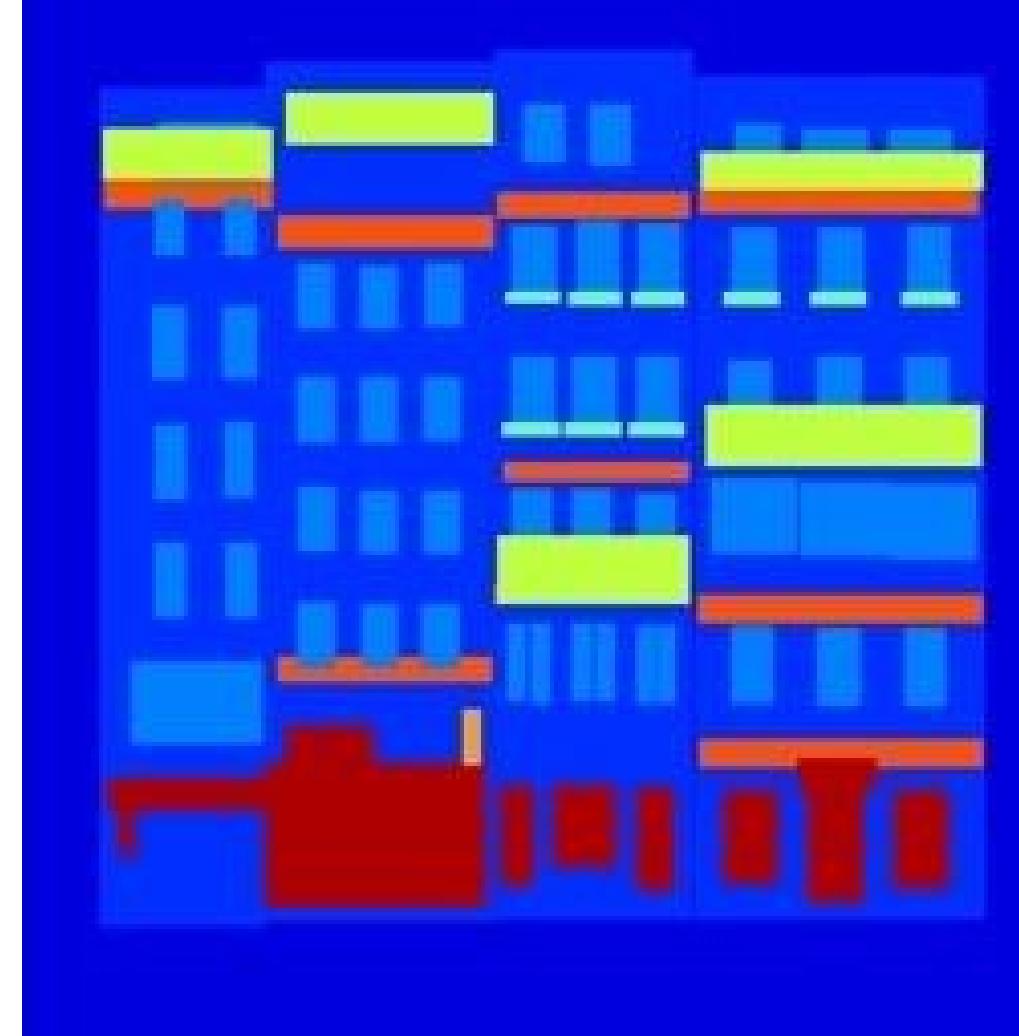
Output



Data from [Tylecek, 2013]

# Labels → Facades

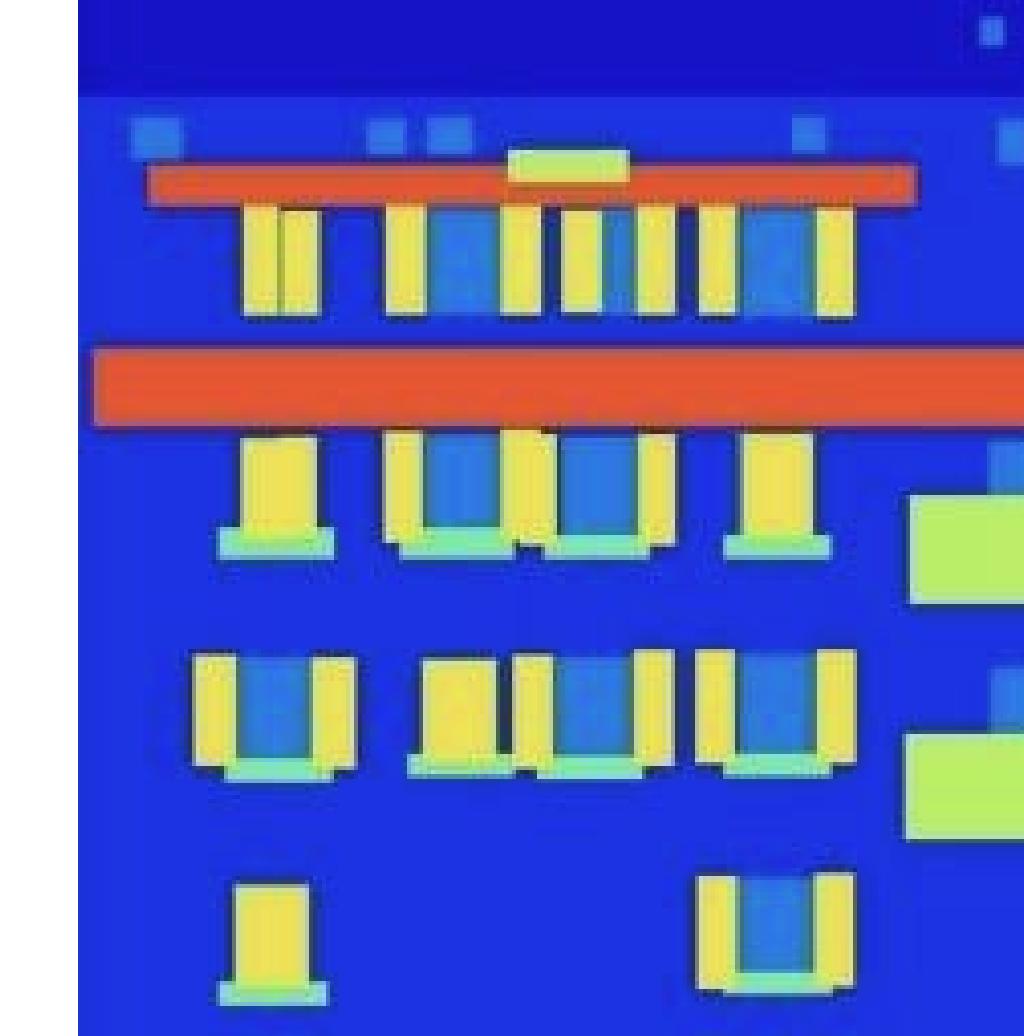
Input



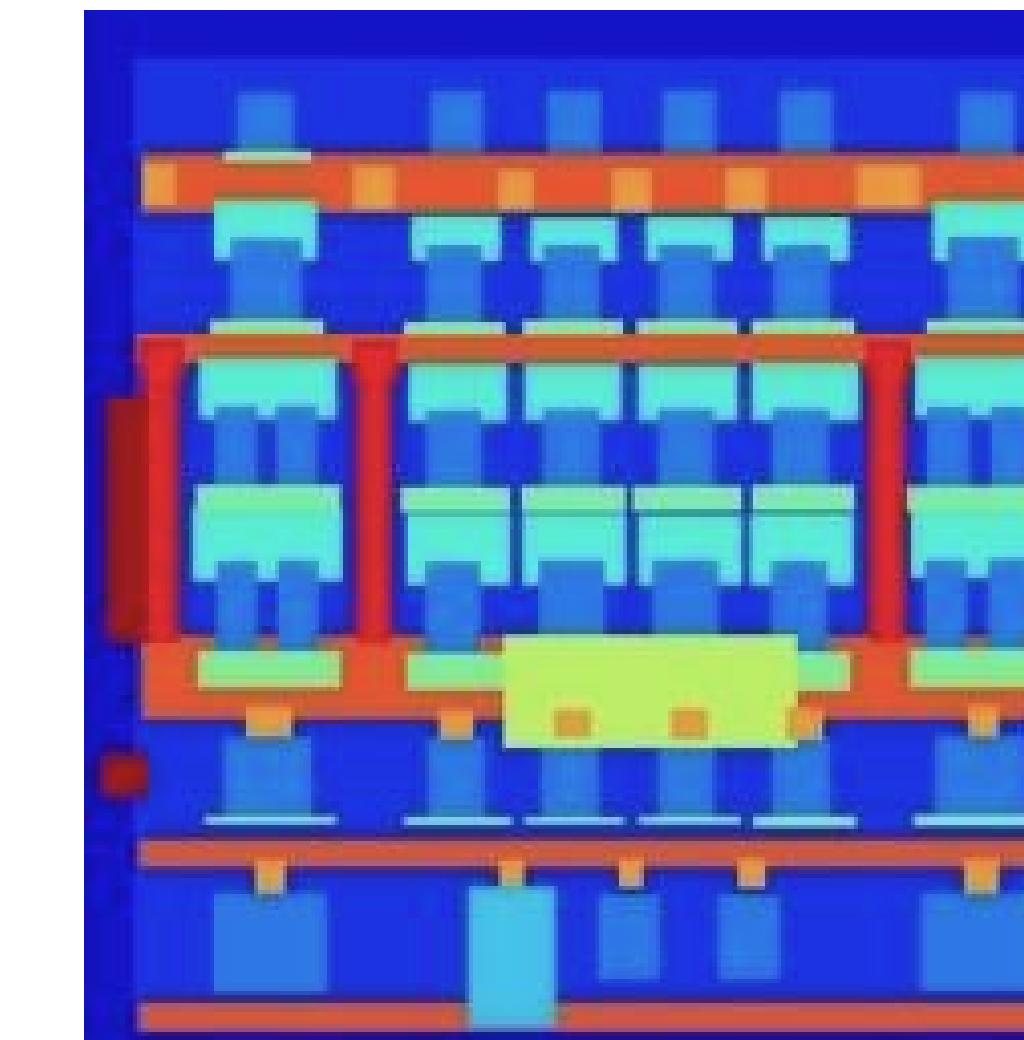
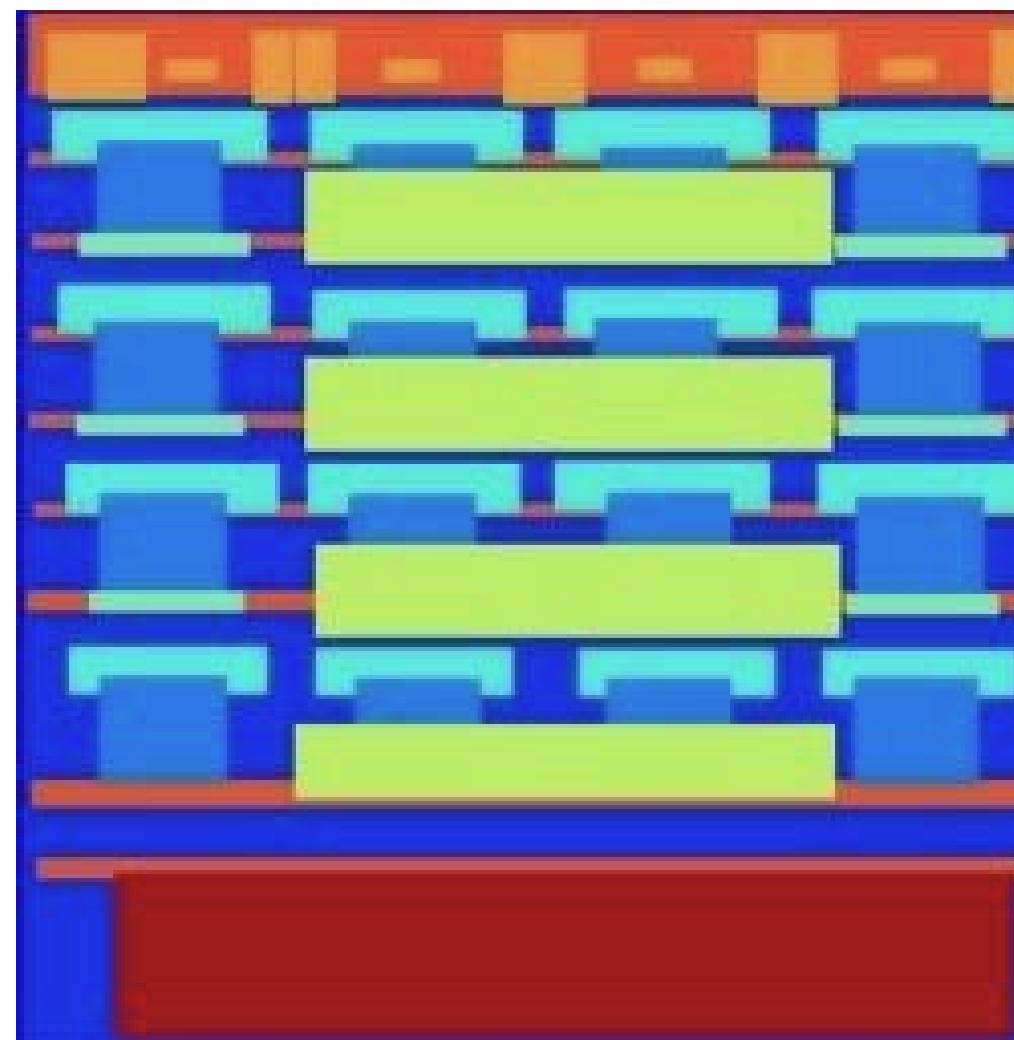
Output



Input



Output



Data from [Tylecek, 2013]

# Day → Night

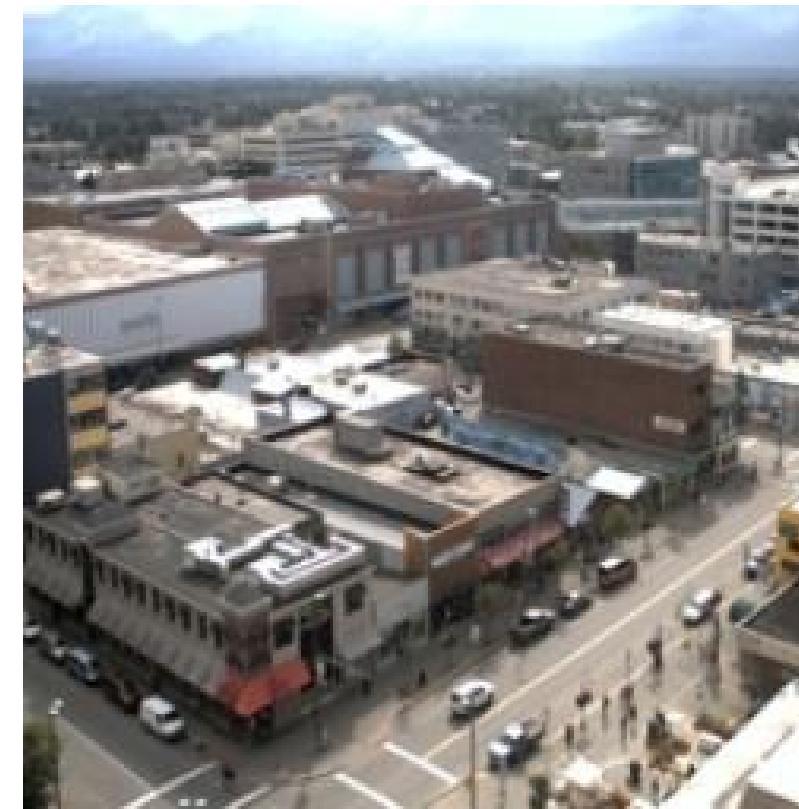
Input



Output



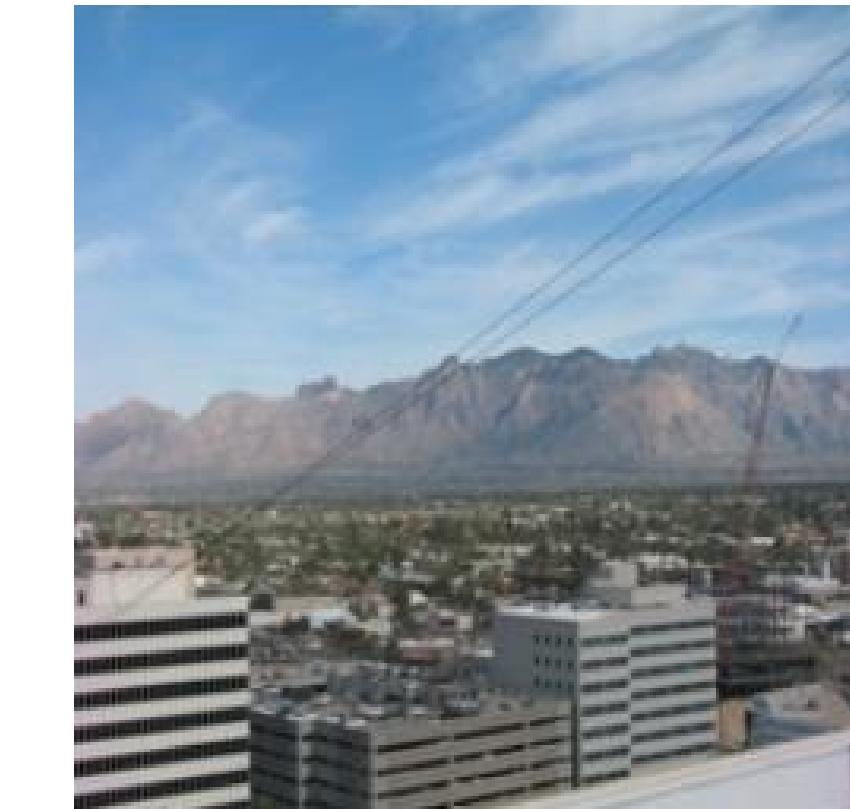
Input



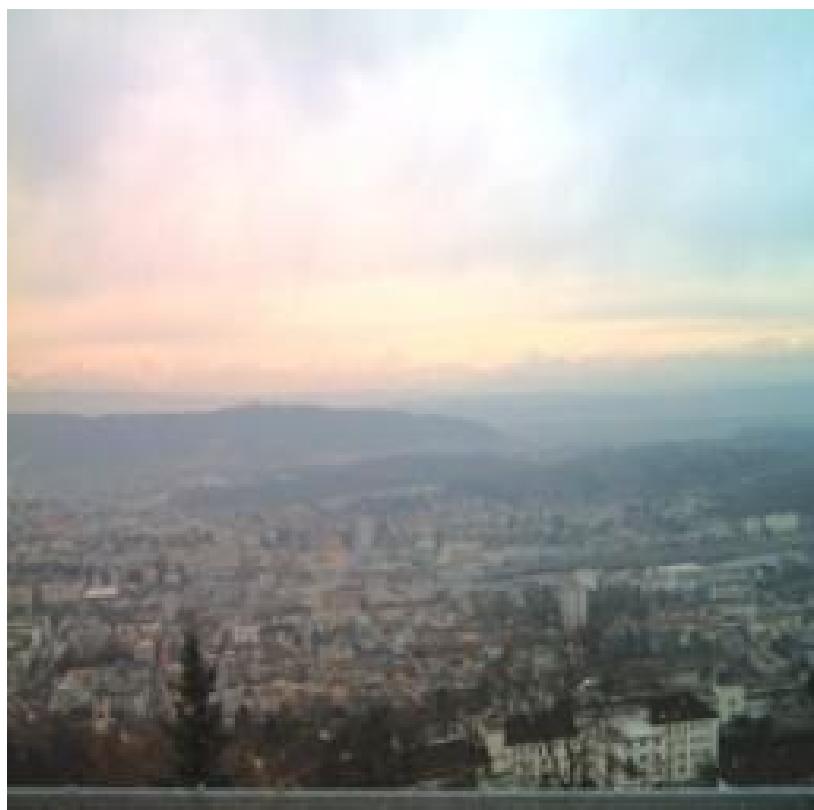
Output



Input



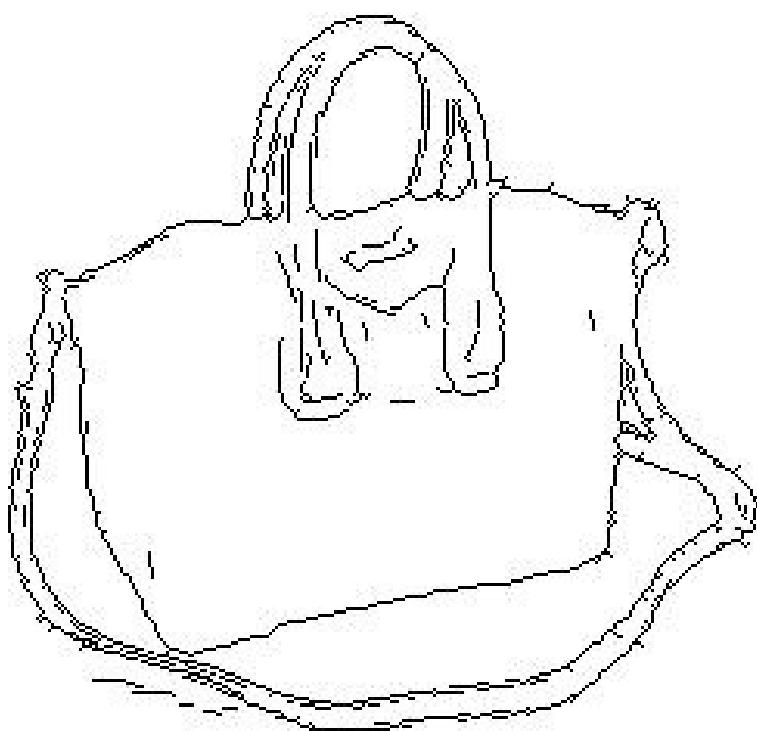
Output



Data from [Laffont et al., 2014]

# Edges → Images

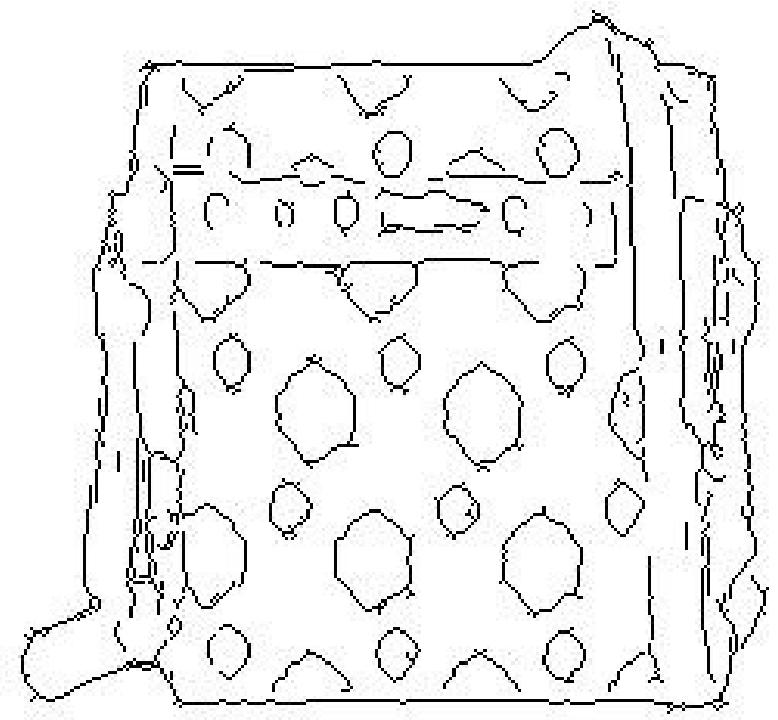
Input



Output



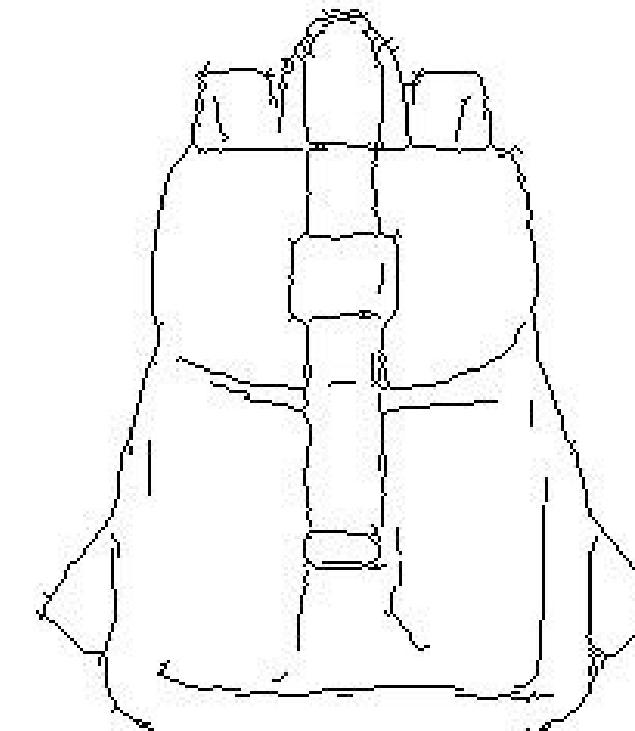
Input



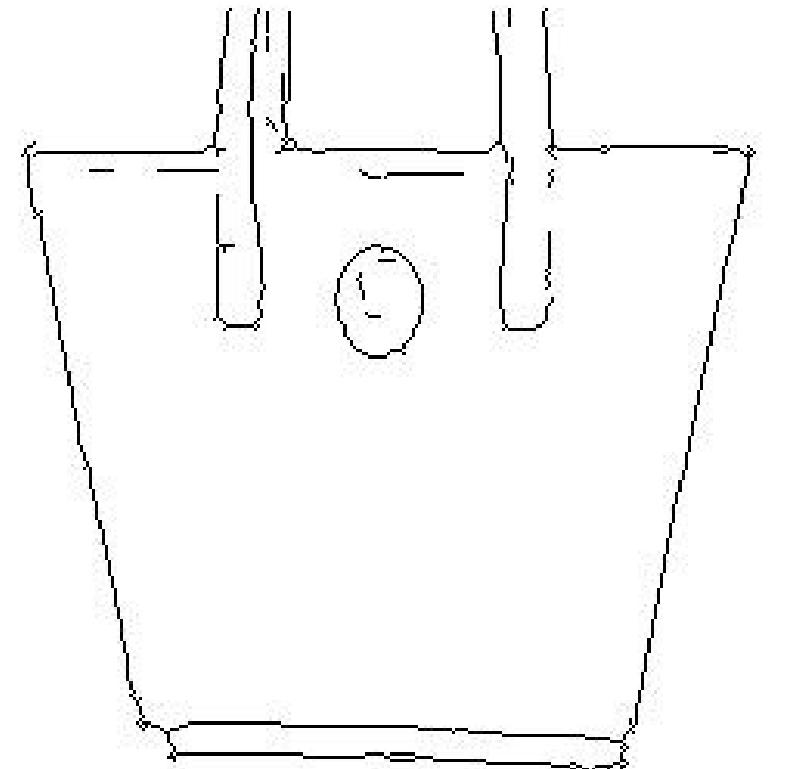
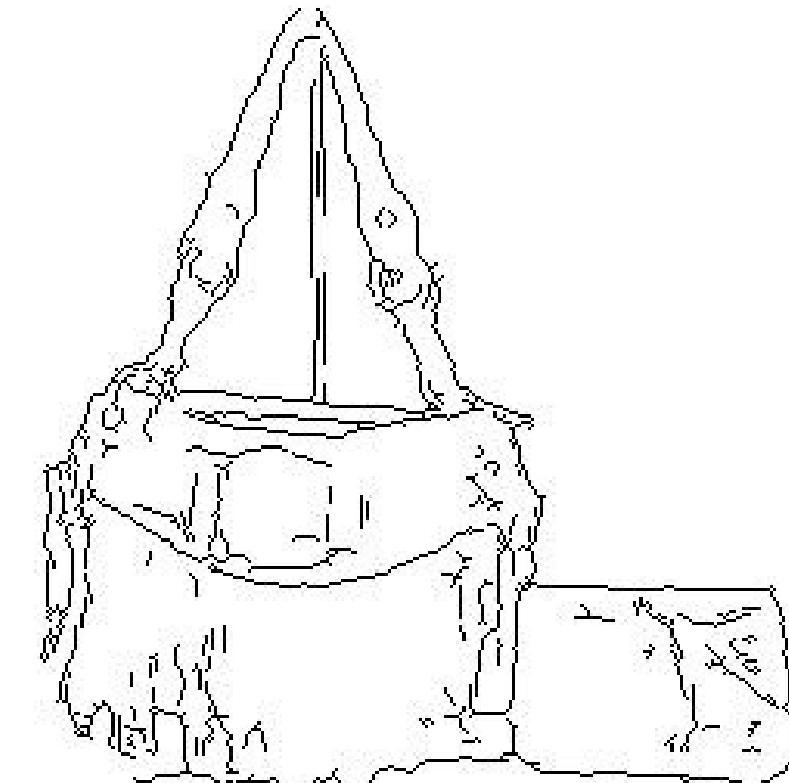
Output



Input



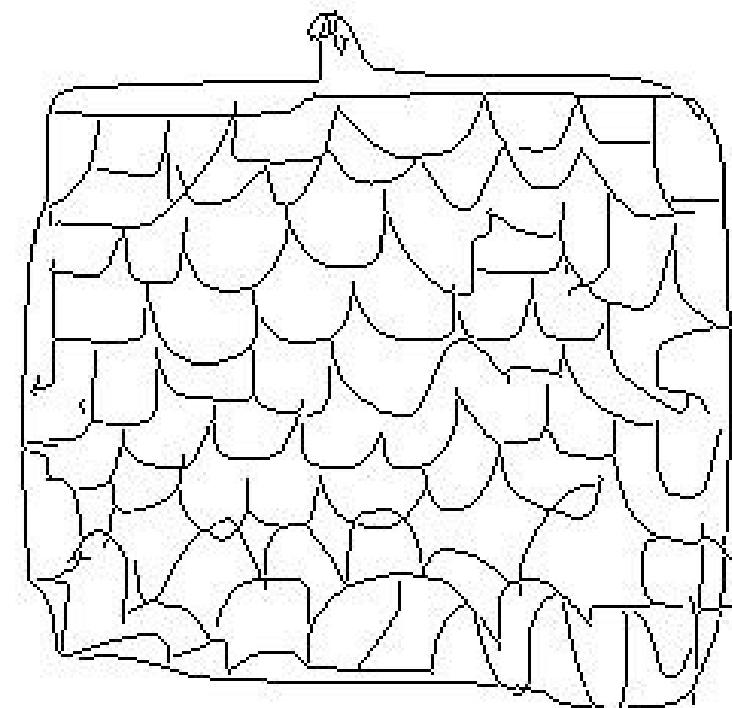
Output



Edges from [Xie & Tu, 2015]

# *Sketches → Images*

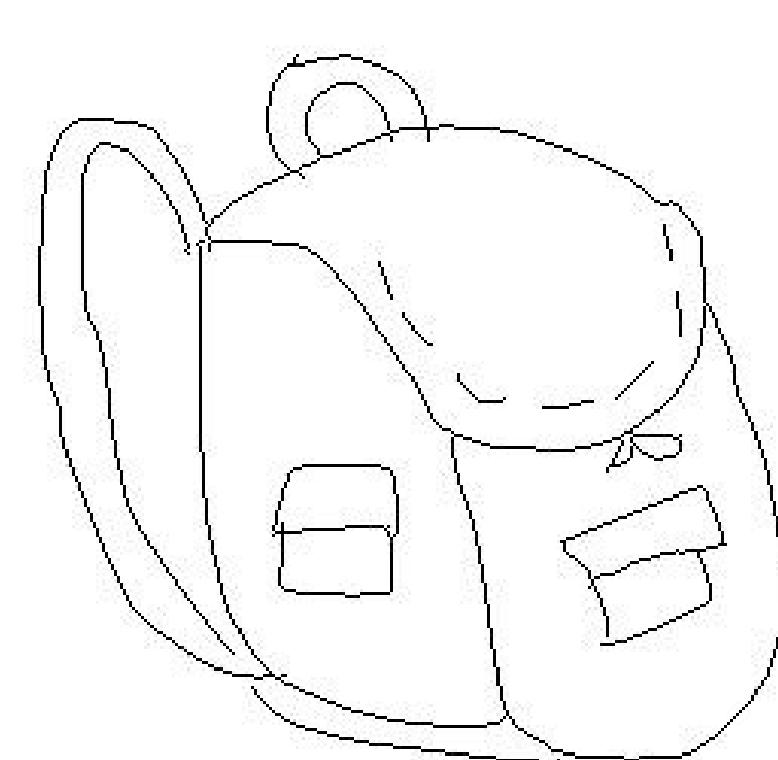
Input



Output



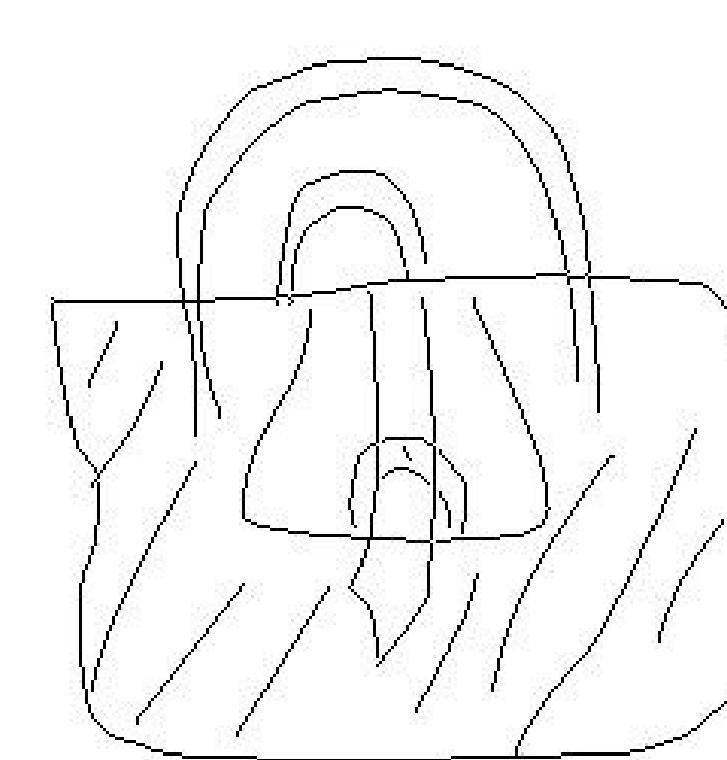
Input



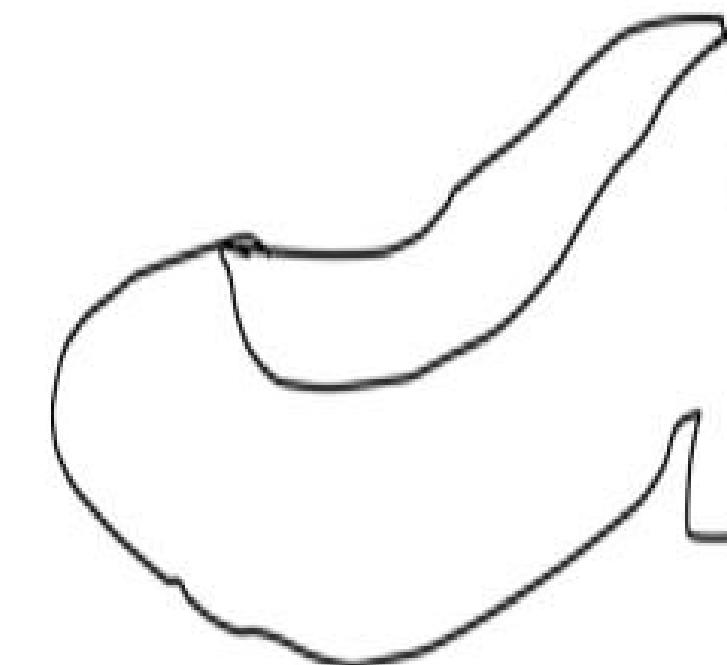
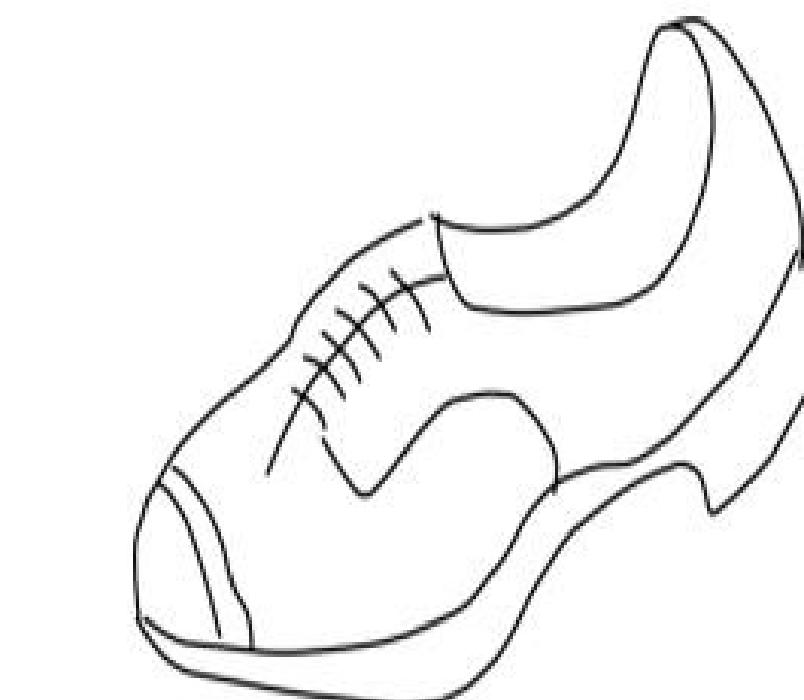
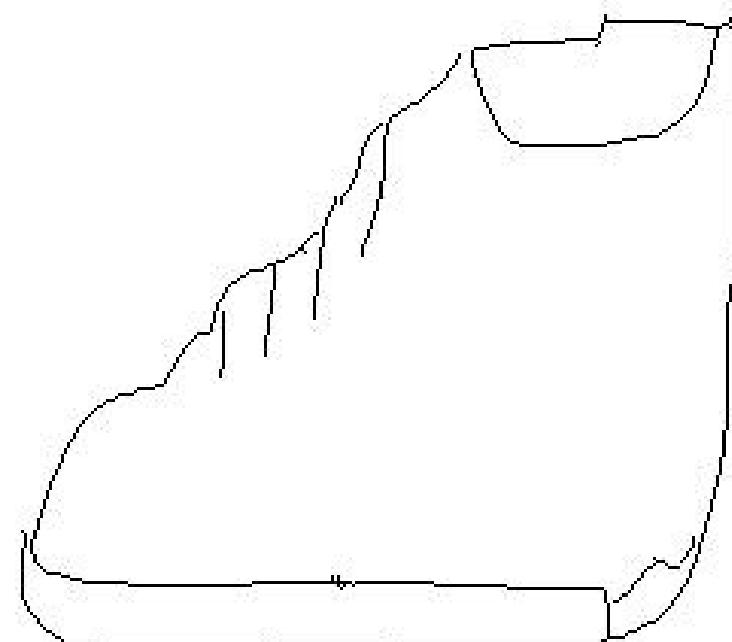
Output



Input



Output



Trained on Edges → Images

Data from [Eitz, Hays, Alexa, 2012]

# Failure cases

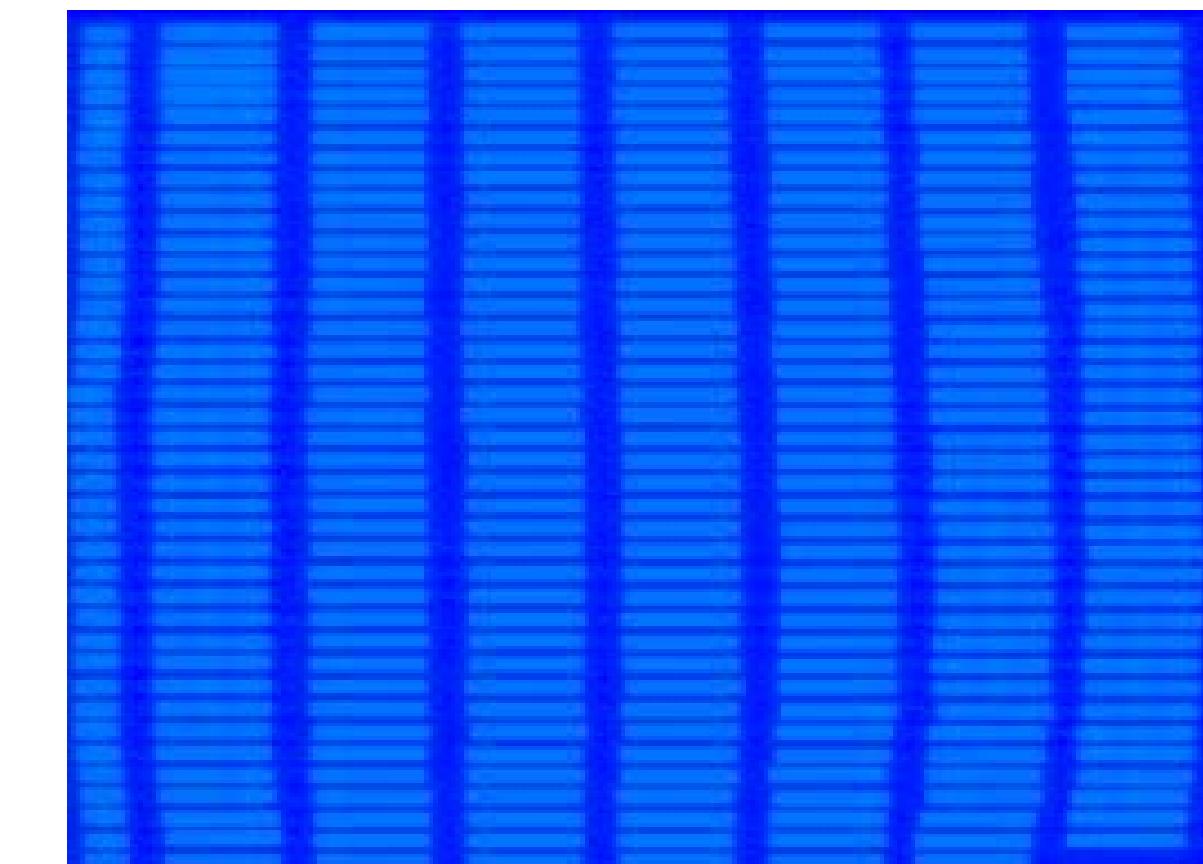
Input



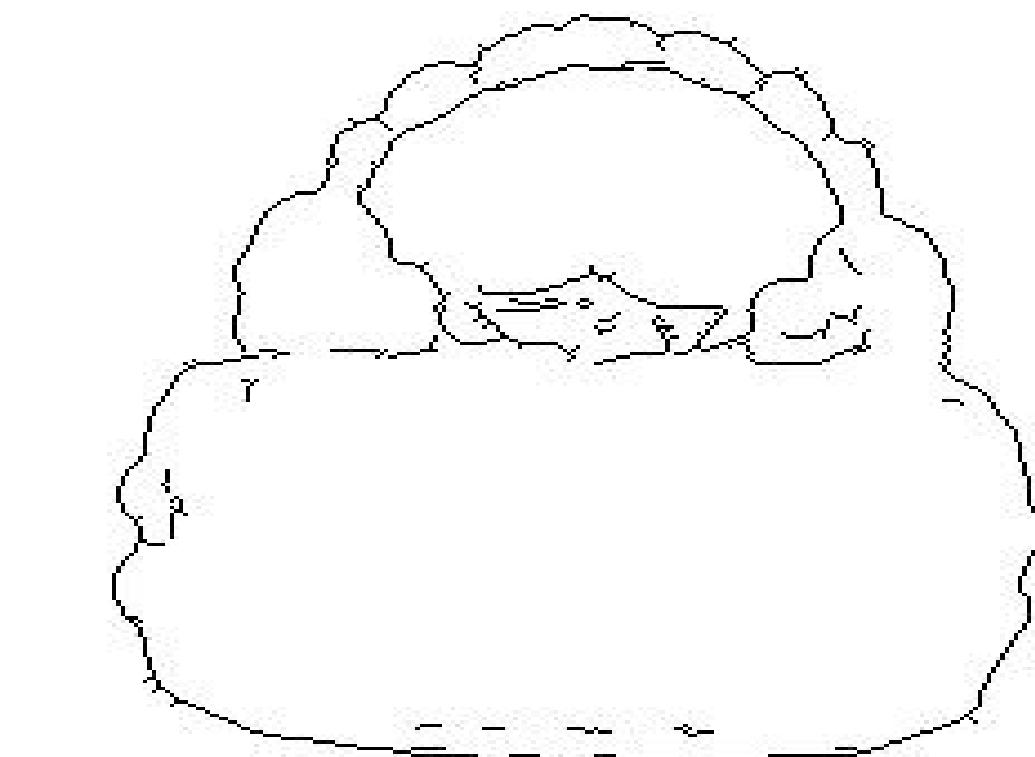
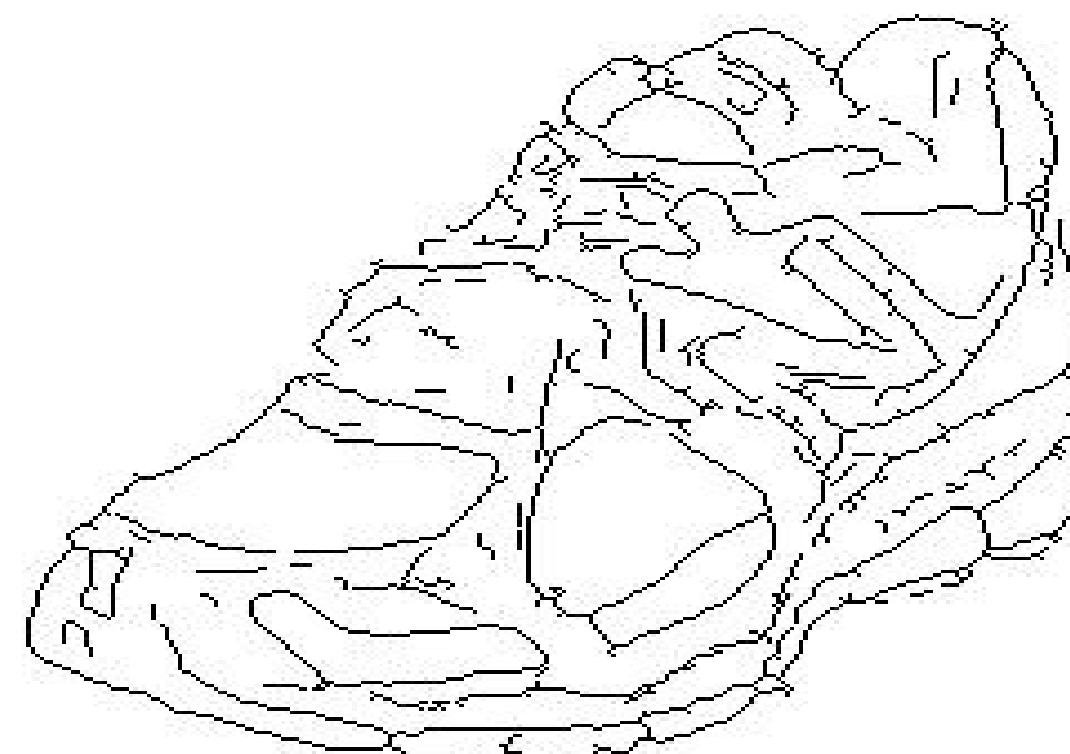
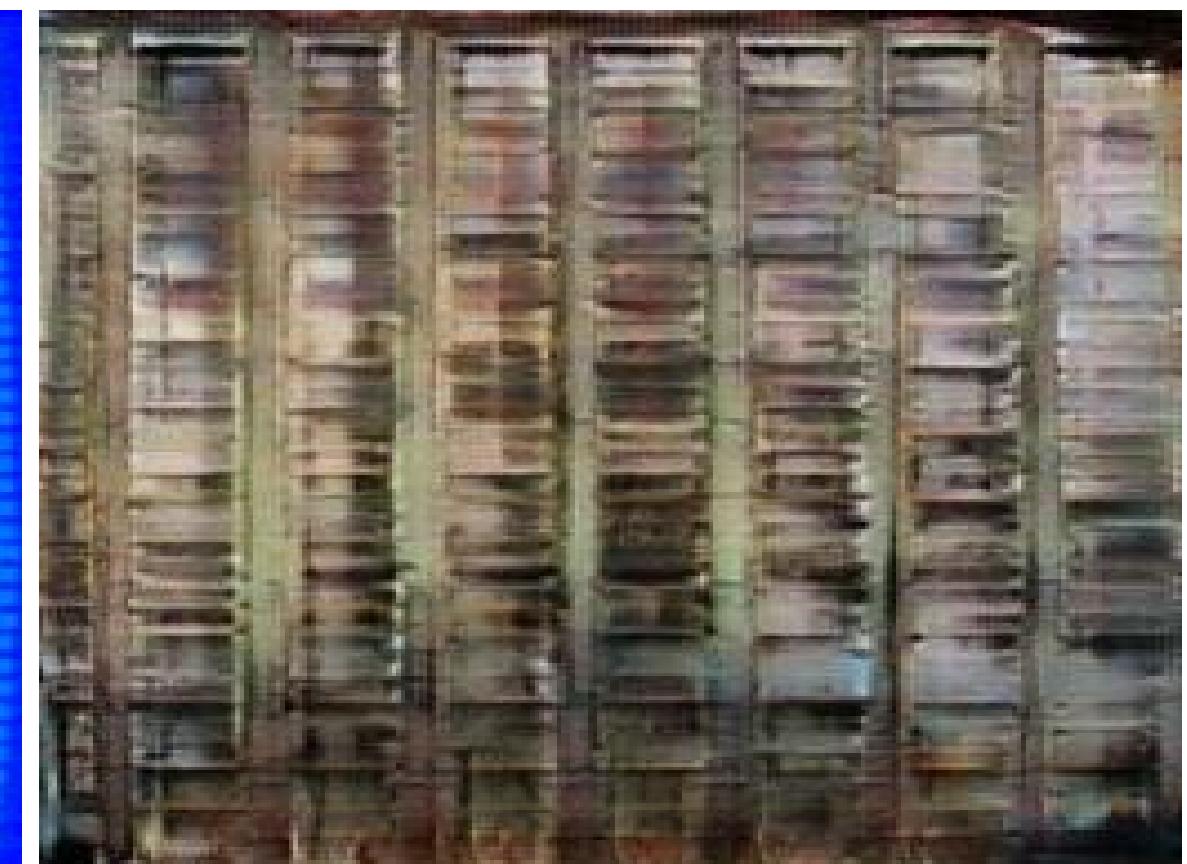
Output



Input



Output



# Code online: <https://github.com/phillipi/pix2pix>

phillipi / pix2pix

Unwatch 80 Star 1,638 Fork 193

Code Issues 8 Pull requests 1 Projects 0 Wiki Pulse Graphs Settings

Image-to-image translation using conditional adversarial nets <https://phillipi.github.io/pix2pix/> Edit

New Add topics

73 commits 1 branch 0 releases 5 contributors

Branch: master New pull request Create new file Upload files Find file Clone or download

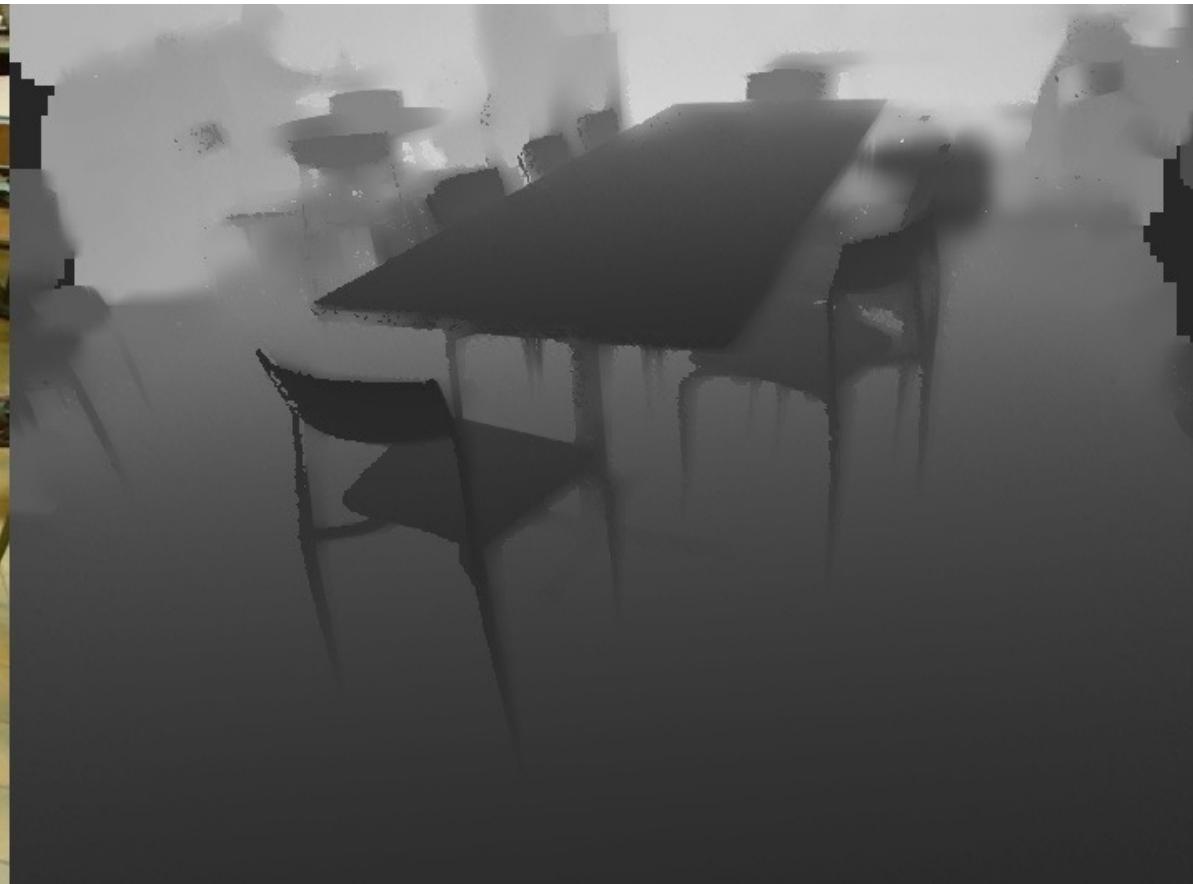
junyanz committed on GitHub add day2night model Latest commit c2bcba4 4 days ago

data	Add preprocess mode for inpainting	a month ago
datasets	add more datasets	2 months ago
imgs	smaller example image	3 months ago
models	add pre-trained models	2 months ago
scripts	update edge scripts	2 months ago
util	Merge pull request #24 from brannondorsey/plot-loss	12 days ago
.gitignore	add pre-trained models	2 months ago
LICENSE	adding license	3 months ago
README.md	add day2night model	4 days ago
models.lua	Added a generator net for 128x128 images	a month ago
test.lua	fix cudnn test error	2 months ago
train.lua	Merge pull request #24 from brannondorsey/plot-loss	12 days ago

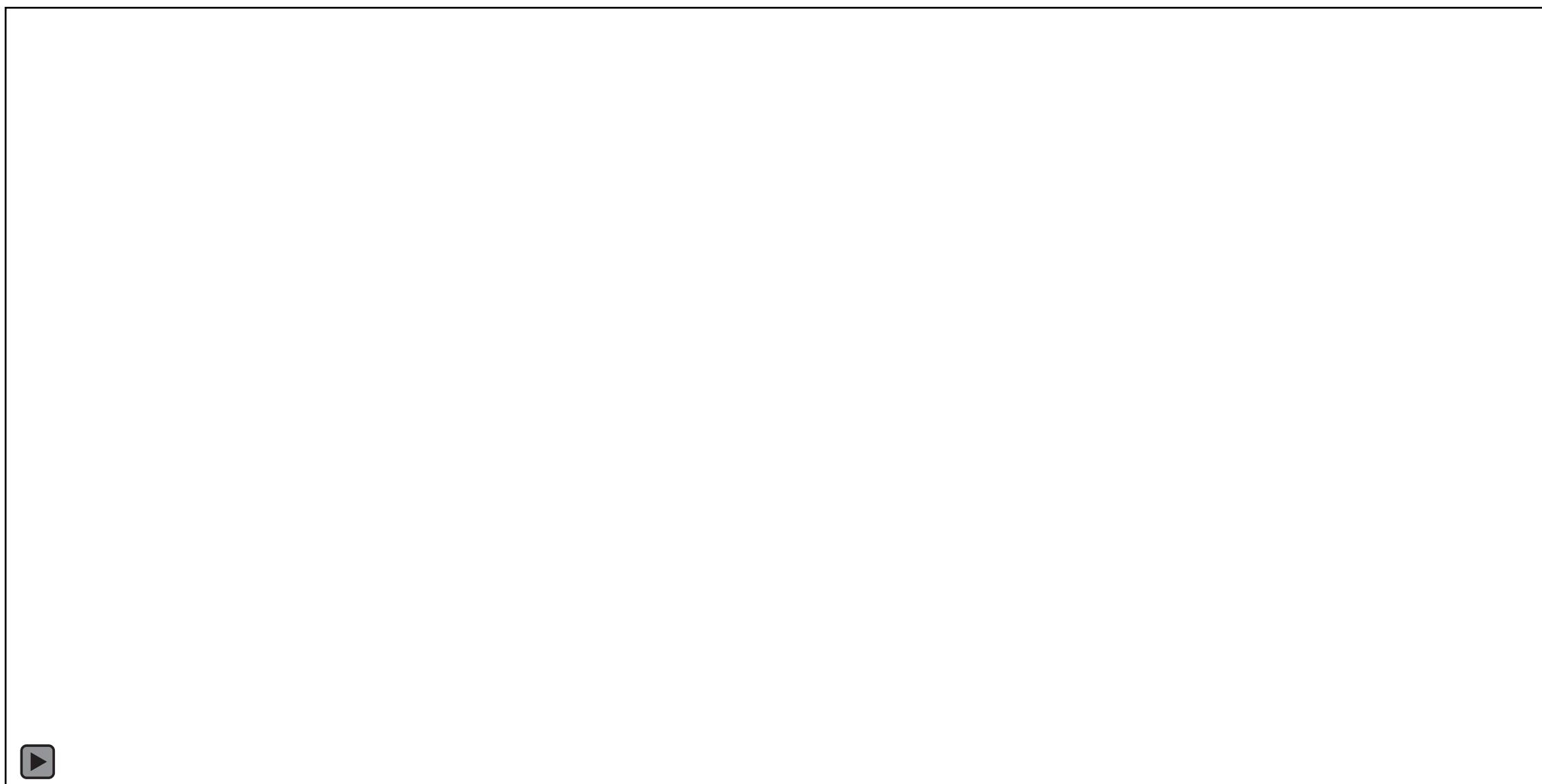
# Twitter-driven research: #pix2pix



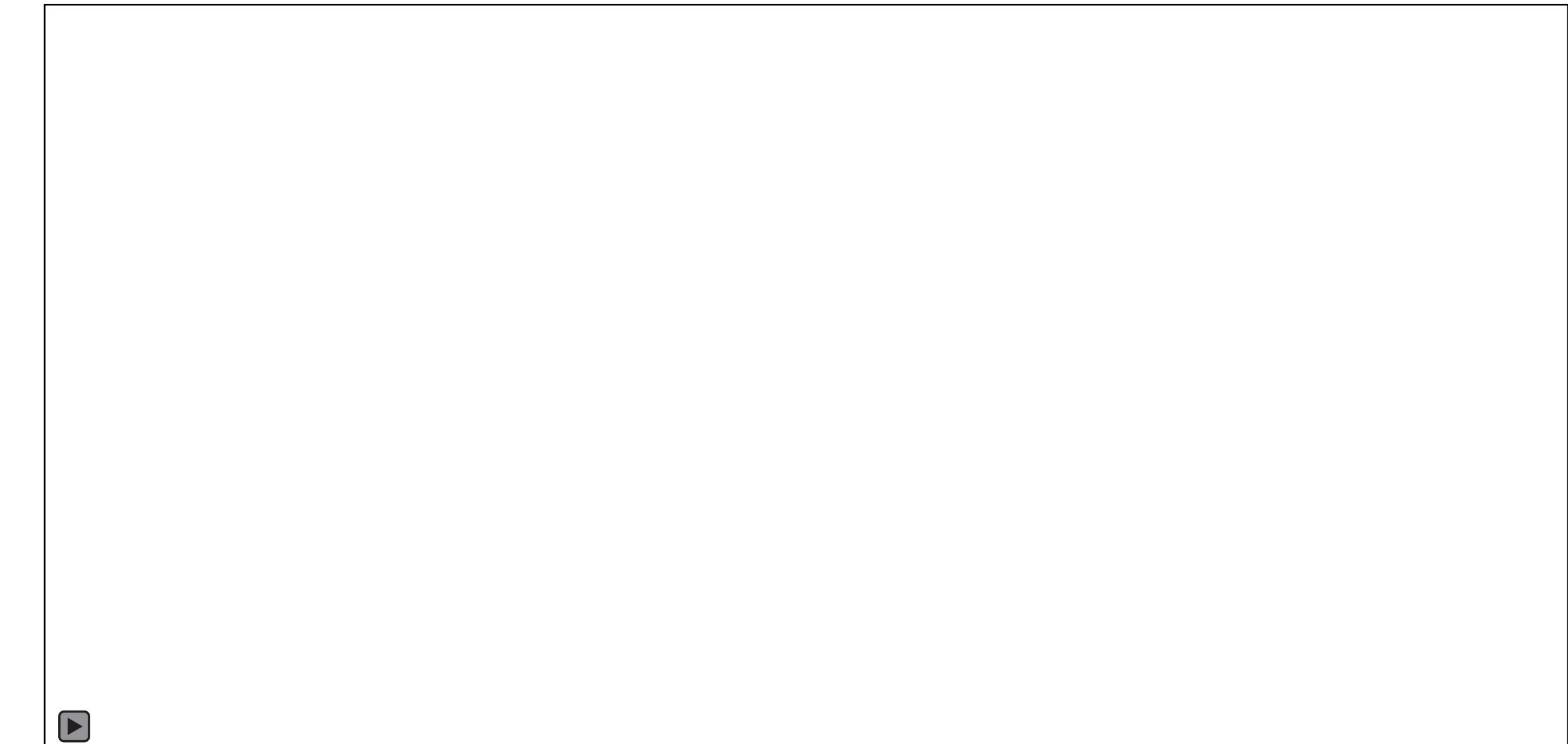
Kaihu Chen @kaihuchen



Mario Klingemann @quasimondo



Brannon Dorsey @brannondorsey



Bertrand Gondouin @bgondouin

# Online Demo!



Pix2pix [Isola et al, 2016]