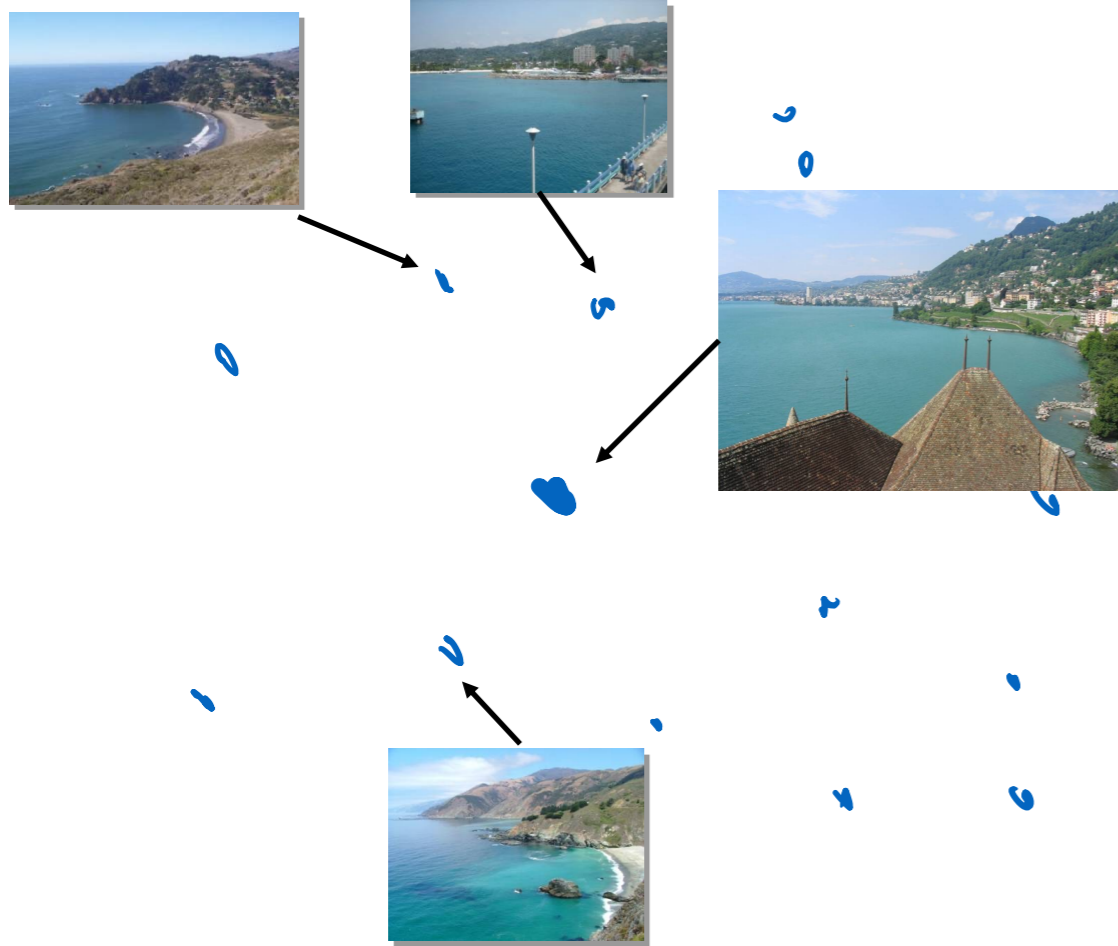# Convolutional Network for Image Synthesis
## Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2022

1

# Review (data-driven graphics)

# Review (data-driven graphics)

Nearest neighbor methods:
1. Stored examples
2. Calculate distance between two examples
3. Voting (label transfer): image blending/averaging

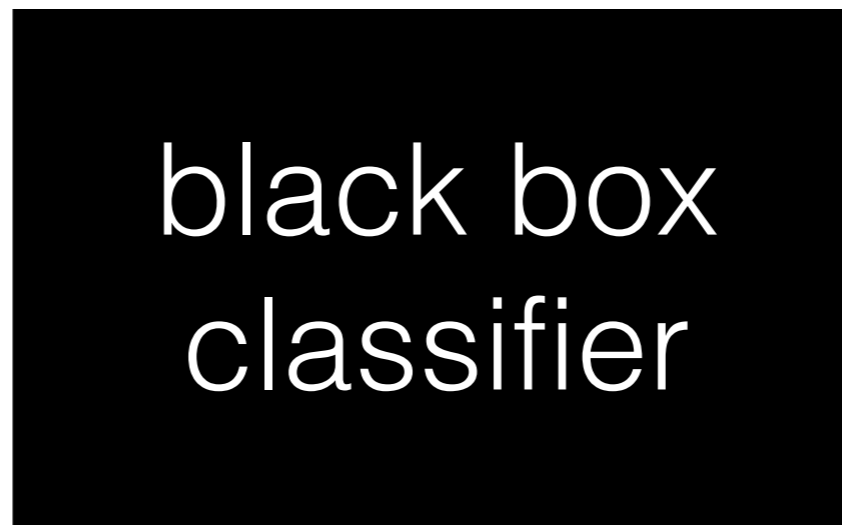# Visual similarity via labels



"Penguin"

?
==

"Penguin"

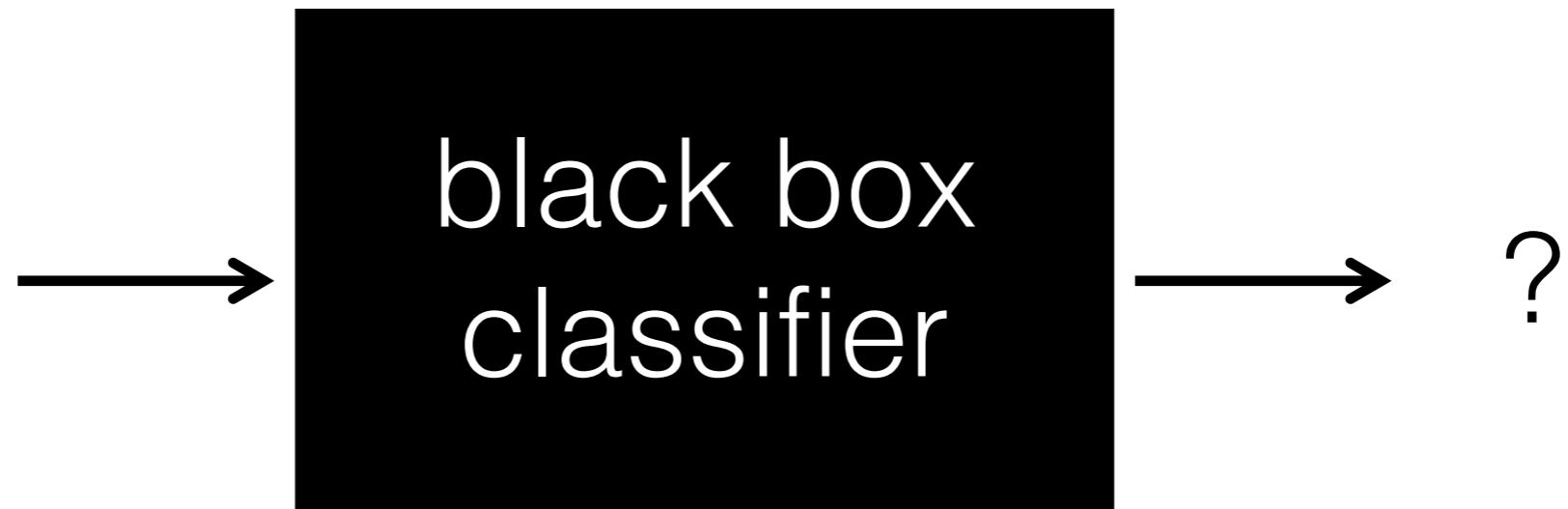# Machine Learning as data association



image $X$        black box classifier        "Penguin"        label $Y$

# At test time...



image *X*

black box classifier

?

Examples from MNIST dataset [LeCun. 1998]

# Warm-up Example: Binary Digit Classification

7 vs. 1

# Learning Approach to Digit Recognition

- **Collect Training Images**
  - Positive:
  - Negative:
- **Training Time**
  - Compute feature vectors for positive and negative example images
  - Train a classifier
- **Test Time**
  - Compute feature vector on new test image:
  - Evaluate classifier

# Let us take an example…



image patch

# Let us take an example…



image patch

Feature vector $\mathbb{R}^{16}$

Note that there are several ways to construct a feature vector. This is one example…

# In feature space, positive and negative examples are just points…

negative examples

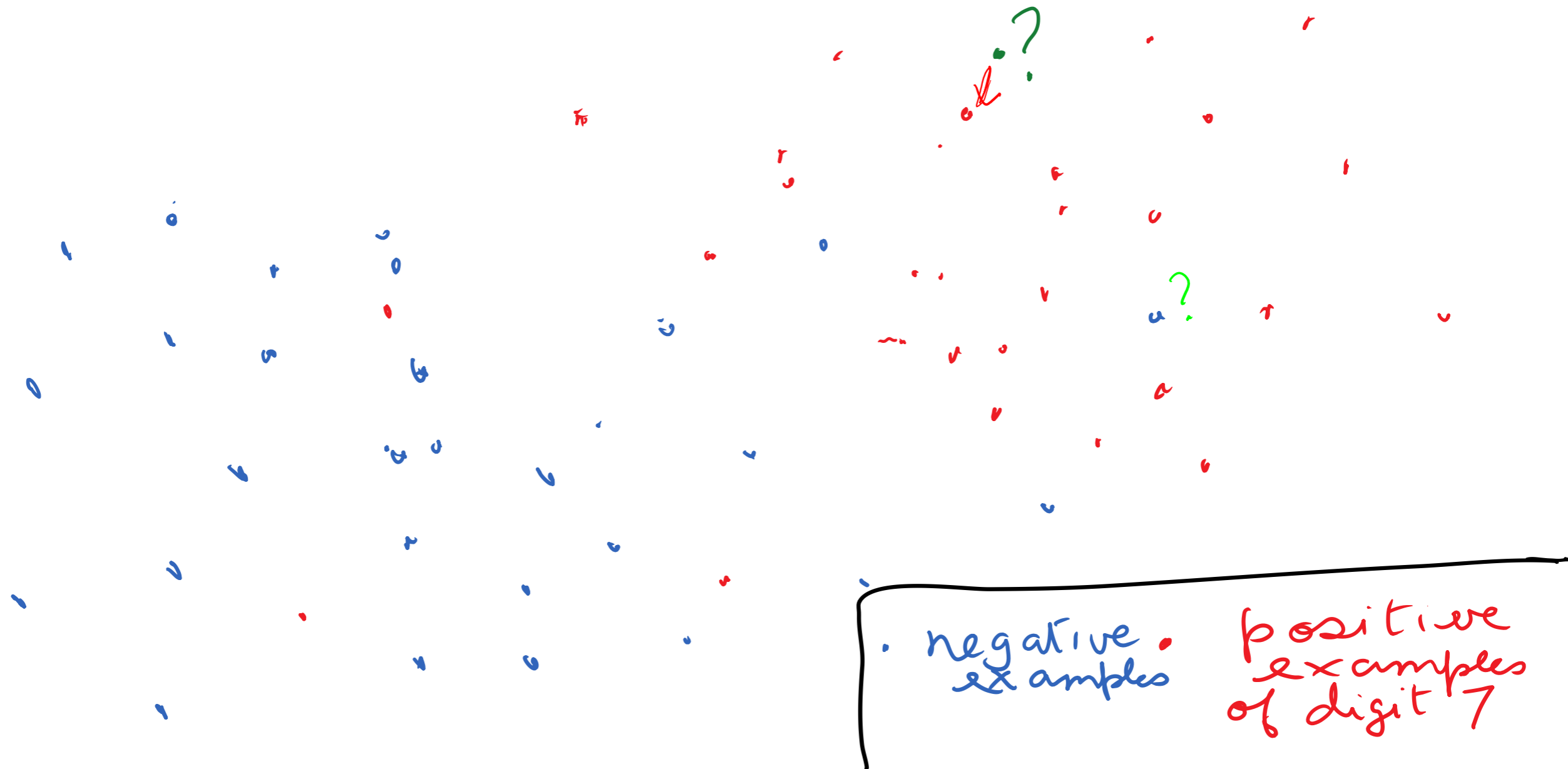positive examples of digit 7

# How do we classify a new point?

?

negative examples

positive examples of digit 7

# Nearest neighbor rule
## "transfer label of nearest example"



negative examples

positive examples of digit 7

# Linear classifier rule



$$\underset{\sim}{w} . \underset{\sim}{x} + b = 0 \text{ is decision boundary}$$

learnt at training time

negative examples

positive examples of digit 7

# Basic idea



Brain/Machine → "clown fish"

# Object recognition



Edges

Texture

Colors

Segments

Parts

"clown fish"

Feature extractors

Classifier

# Object recognition

Learned



Edges

Texture

Colors

Segments

Parts

"clown fish"

Feature extractors

Classifier

# Neural network



Learned

"clown fish"

# Neural network

Learned



"clown fish"

# Deep neural network

Learned



"clown fish"

# Computation in a neural net

Input
representation

Output
representation

# Computation in a neural net

Input
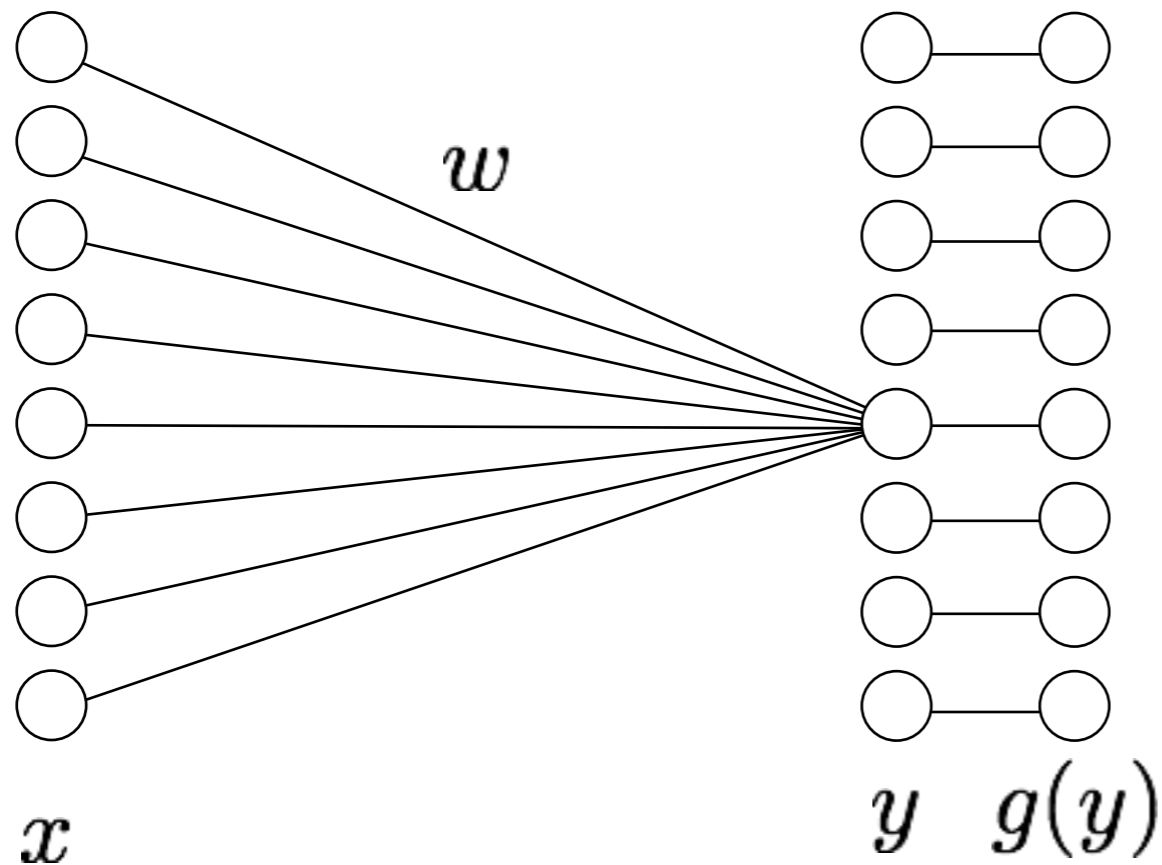representation

Output
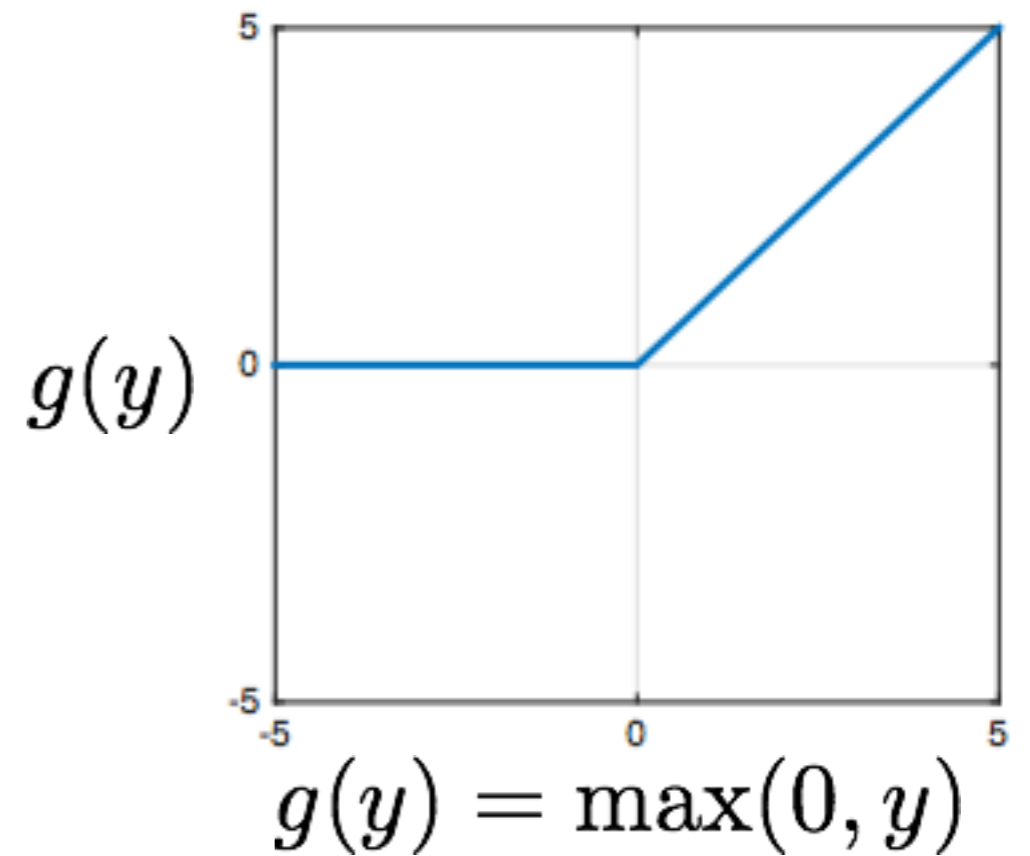representation



$w$

$x$

$y$

$$y_j = \sum_i w_{ij} x_i$$

i: the $i^{th}$ dimension of $x$, j; the $j^{th}$ dimension of $y$

# Computation in a neural net

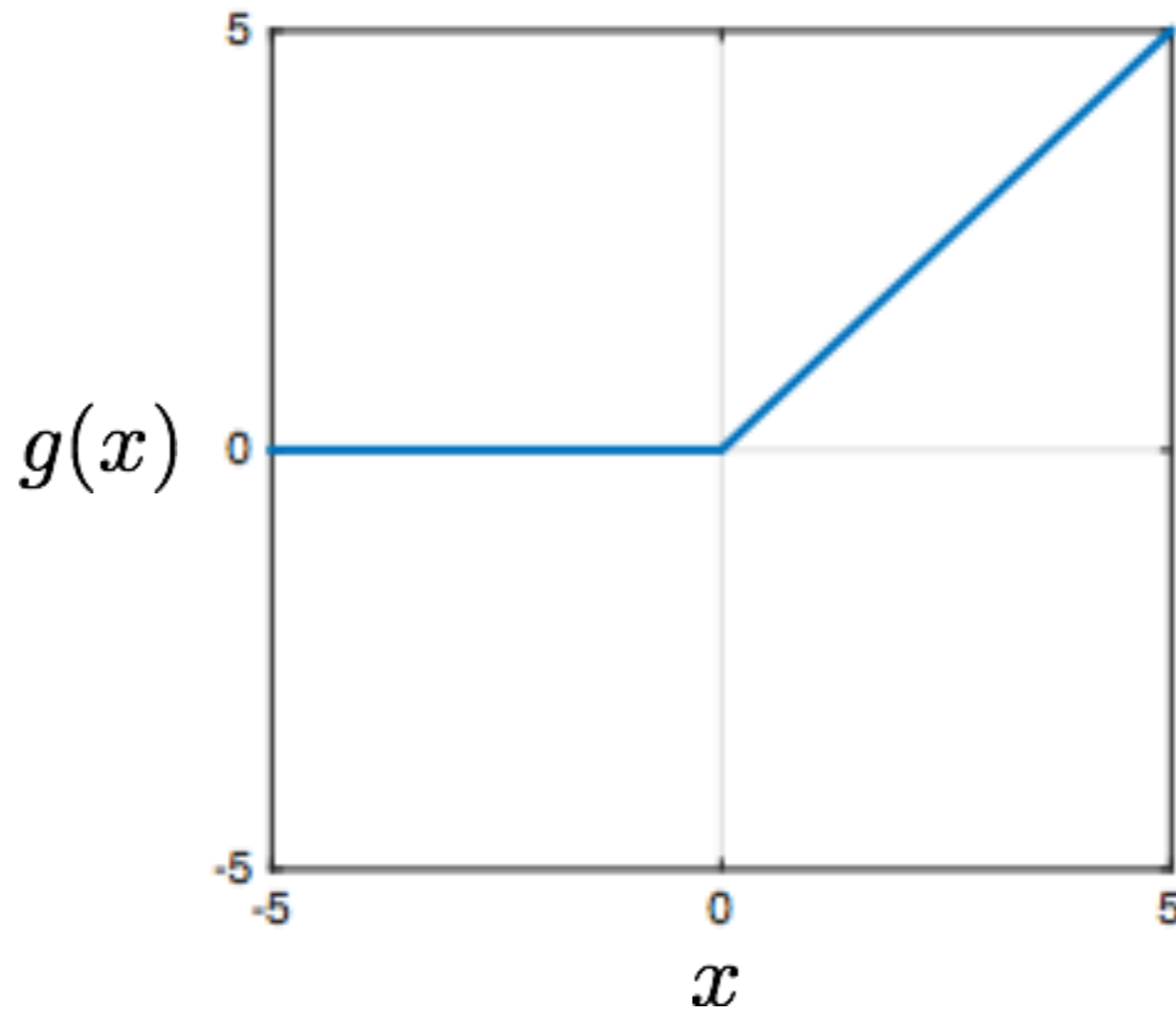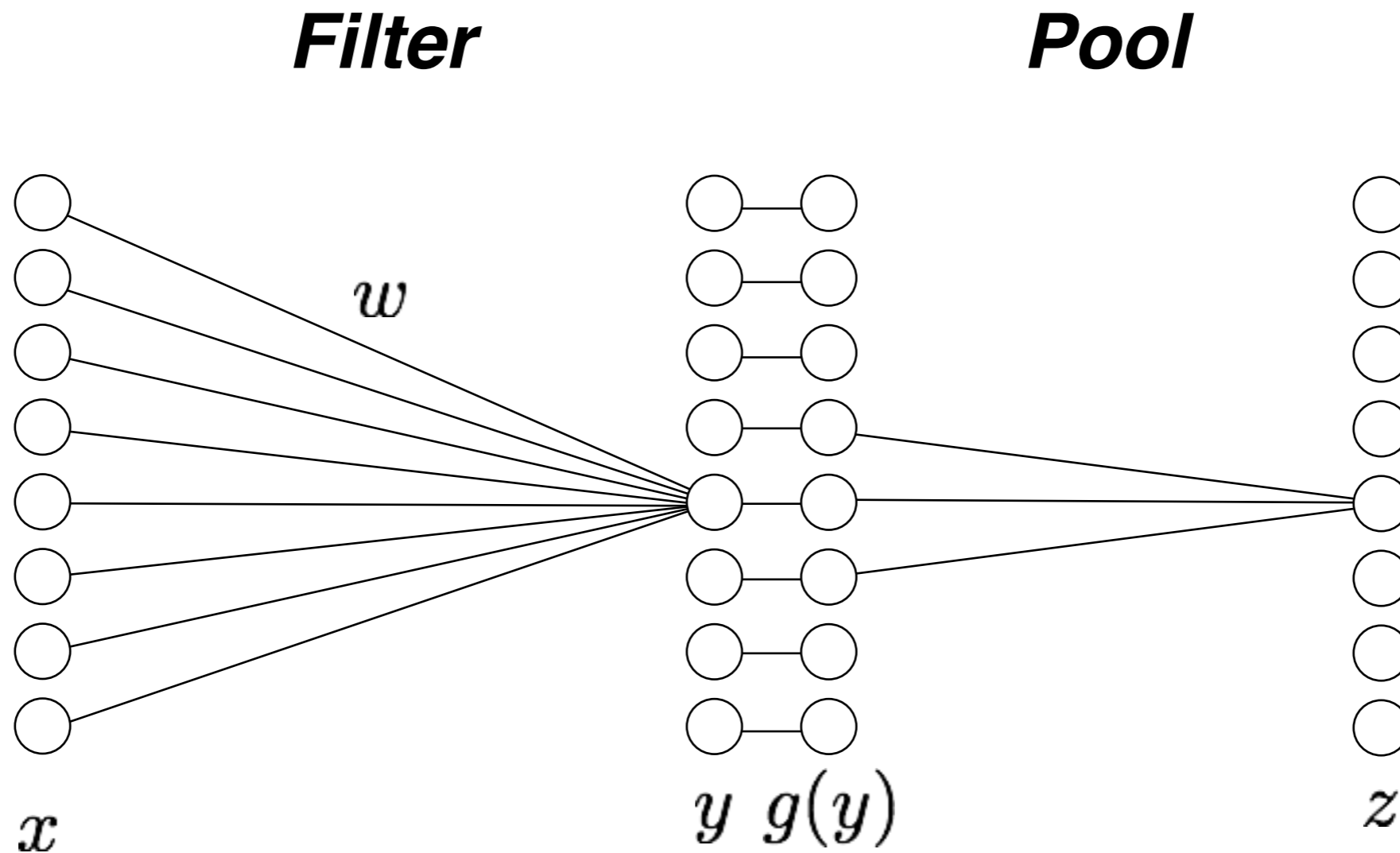Rectified linear unit (ReLU)

$w$

$g(y)$

$g(y) = \max(0, y)$

$y \quad g(y)$

$x$

# Computation in a neural net

## Rectified linear unit (ReLU)



$$g(x) = \max(0, x)$$

# Computation in a neural net

**Filter**                    **Pool**



$x$                    $y \; g(y)$                    $z$
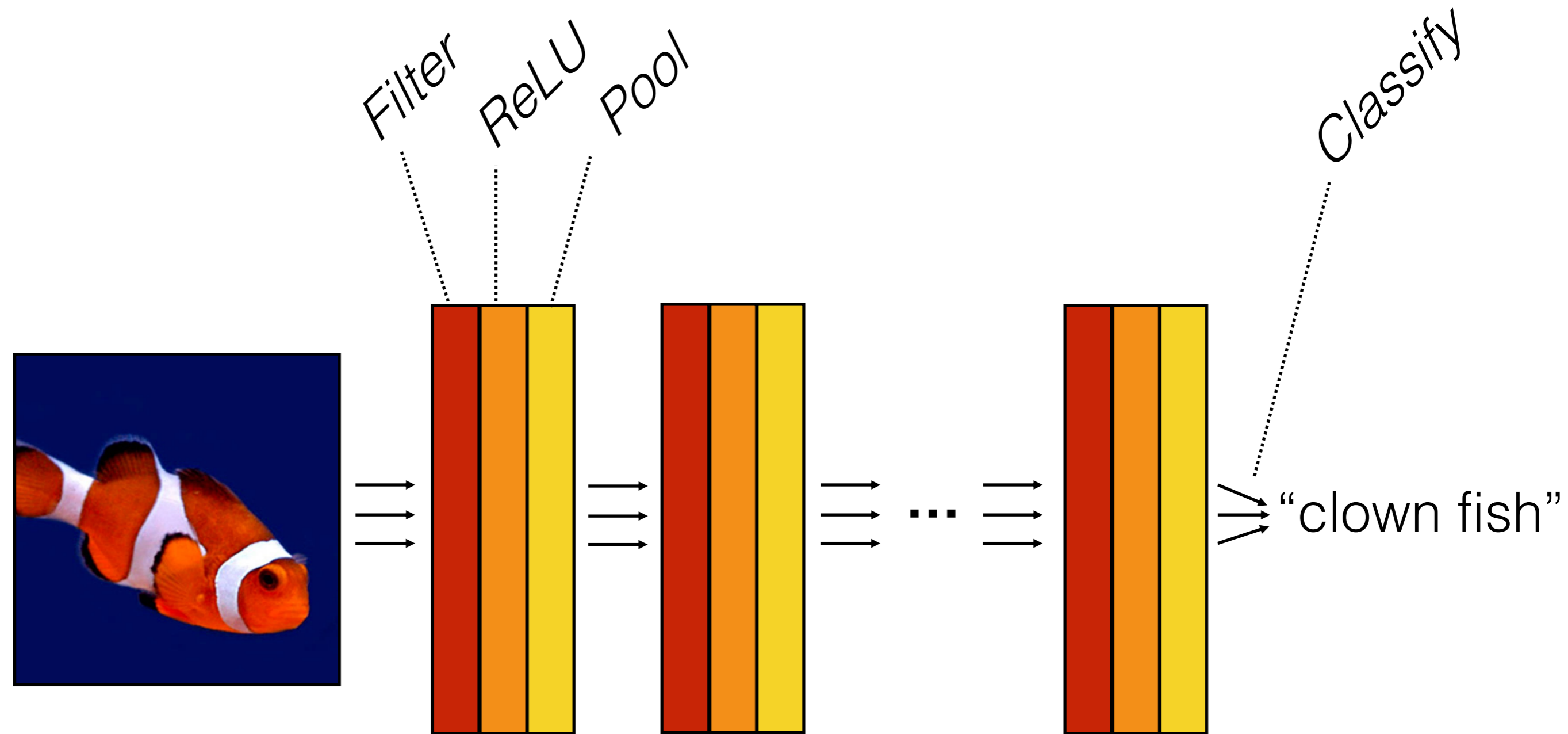
$$y_j = \sum_i w_{ij} x_i \qquad\qquad z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$

i: the $i^{th}$ dimension of $x$, j; the $j^{th}$ dimension of $y$

# Computation in a neural net



Single depth slice

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

max pool with 2x2 filters and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

224x224x64
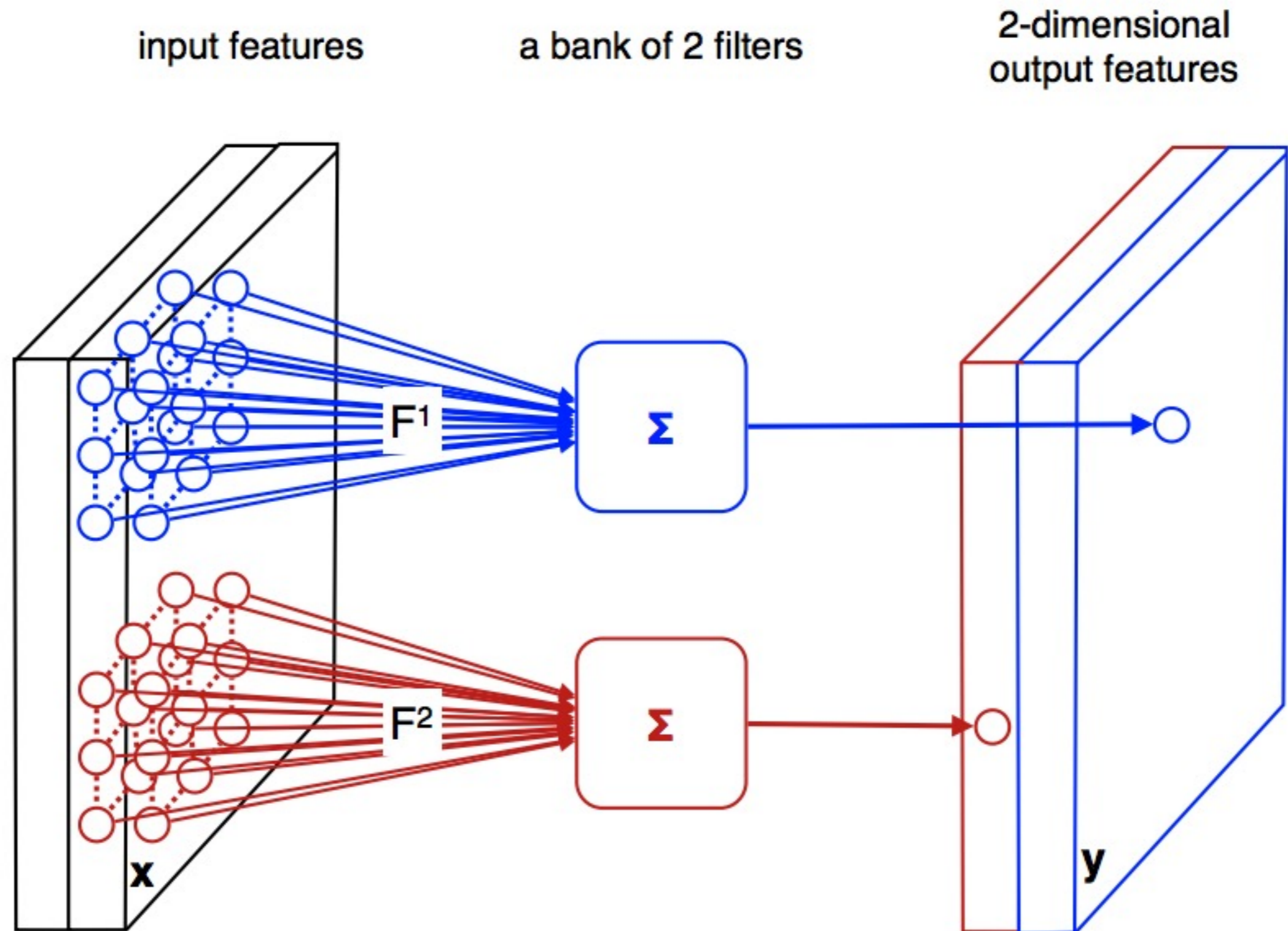
pool

112x112x64

224
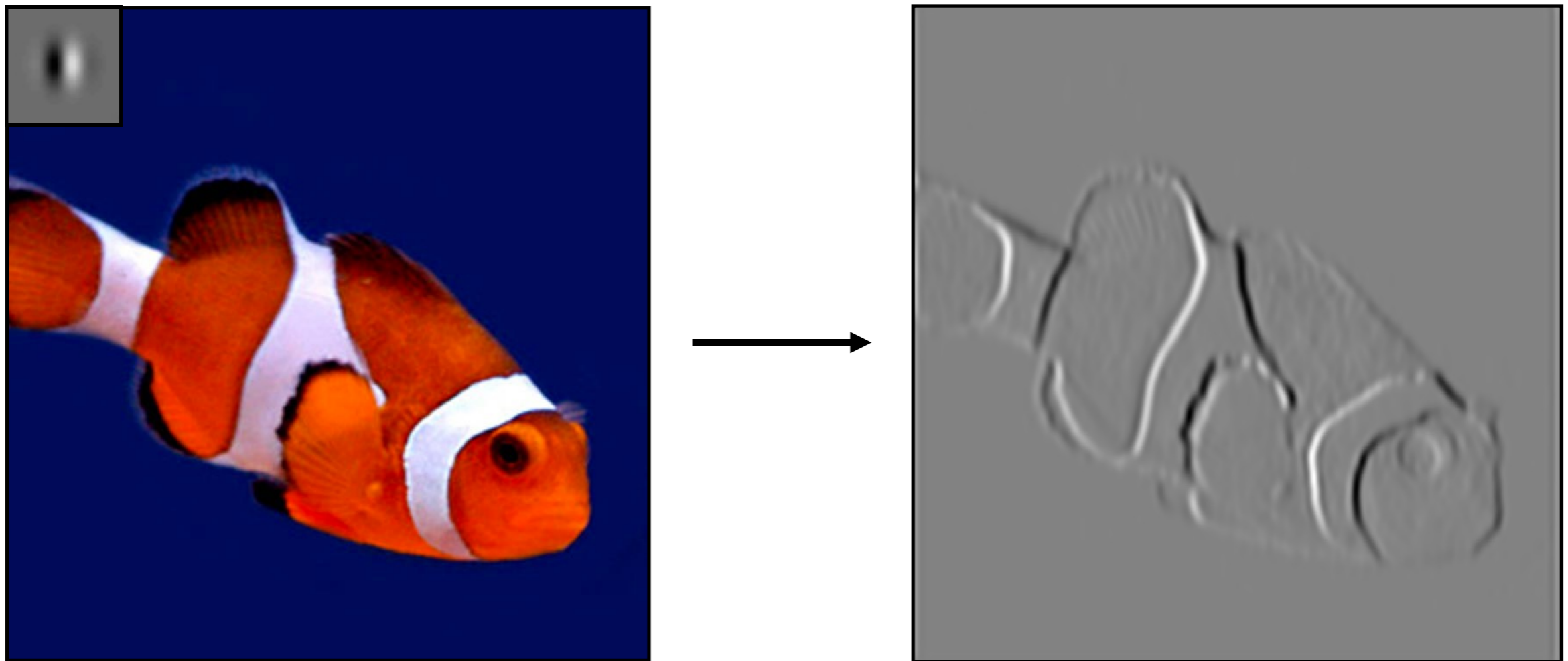
downsampling

112

224

112

# Computation in a neural net



$$f(\mathbf{x}) = f_L(\ldots f_2(f_1(\mathbf{x})))$$

# Convolutional Neural Nets



input features

a bank of 2 filters

2-dimensional output features

$F^1$

$\Sigma$

$F^2$
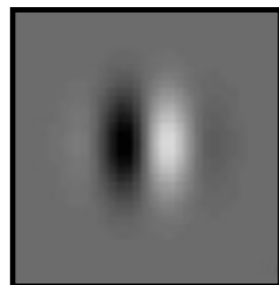
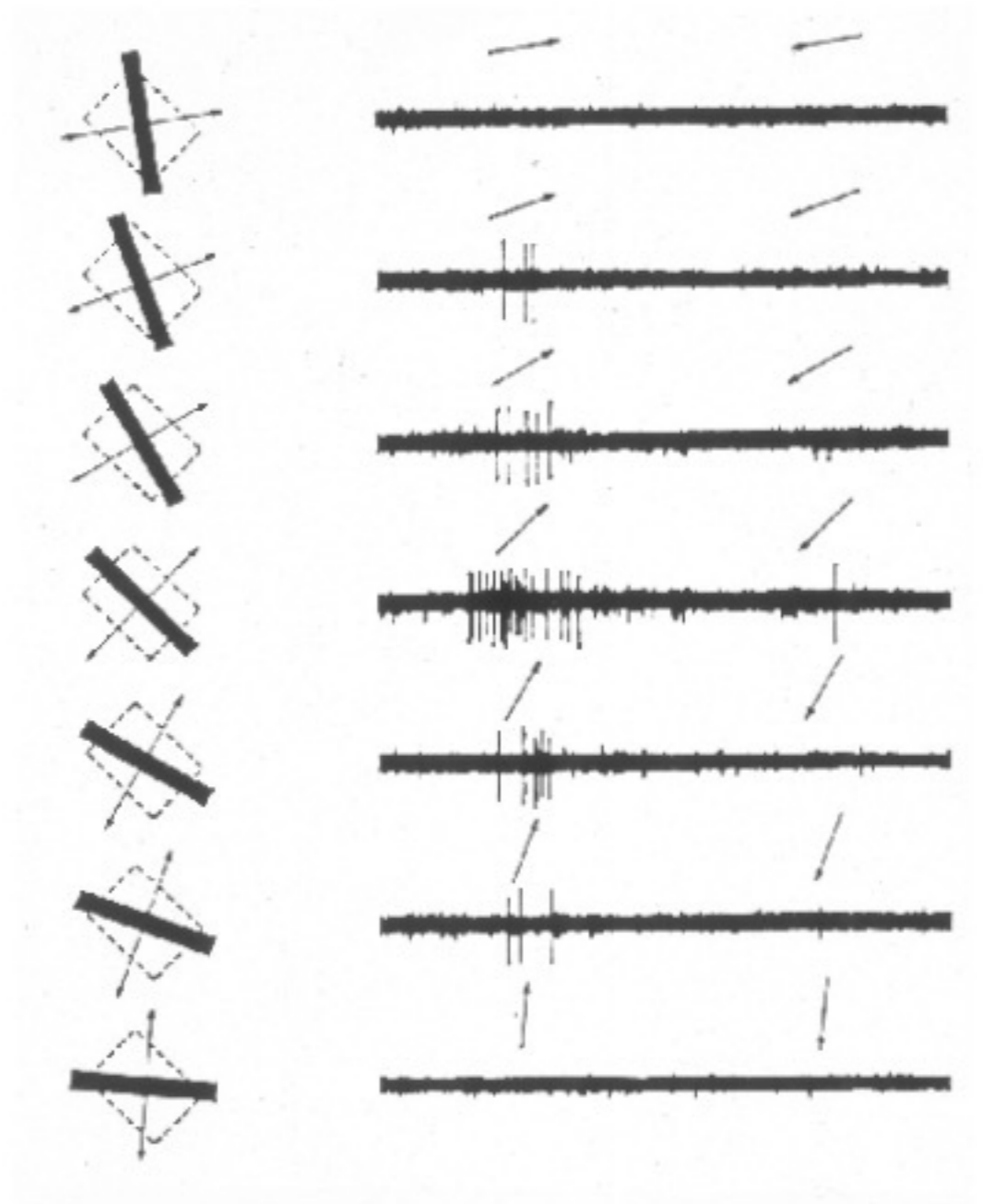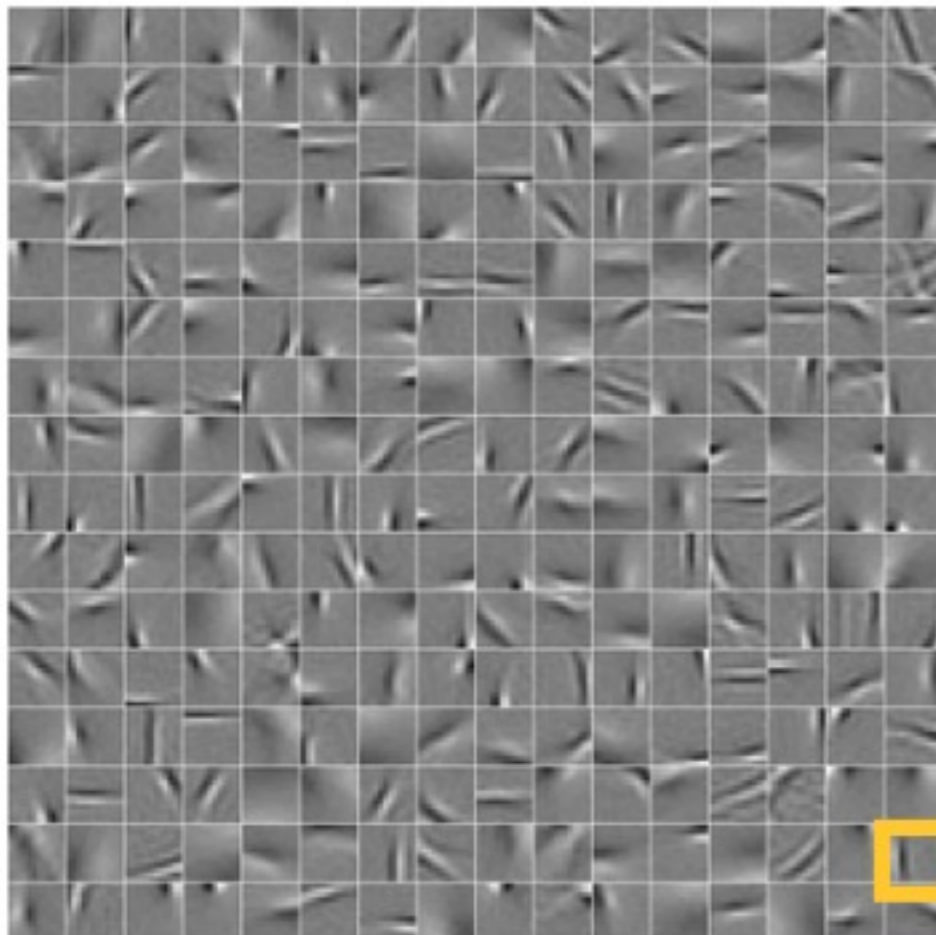$\Sigma$

x

y

Slide from Andrea Vedaldi

# Convolutional Neural Nets

## Convolution



filter

# [Hubel and Wiesel 59]



oriented filter

Slide from Andrea Vedaldi

# Computation in a neural net

Last layer



dolphin

cat

grizzly bear

angel fish

chameleon

**clown fish**

iguana

elephant

argmax

"clown fish"

# Learning with deep nets

Learned

"clown fish"

# Learning with deep nets

 ⟶ "clown fish"

 ⟶ "grizzly bear"

Train network to associate the right label with each image

 ⟶ "chameleon"

# Learning with deep nets

"clown fish"

Learned

**Loss** — $L()$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

$$\underset{\mathbf{w}}{\mathrm{argmin}} \quad L(w_1, \ldots, w_6)$$

# Loss function

## Network output

**Ground truth label**

dolphin

cat

grizzly bear

angel fish

chameleon

**clown fish**

iguana

elephant

...

"clown fish"

Loss → error

36

# Loss function

**Network output**



- dolphin
- cat
- grizzly bear
- angel fish
- chameleon
- **clown fish**
- iguana
- elephant

**Ground truth label**

"clown fish"

Loss ⟶ **small**

# Loss function

**Network output**

**Ground truth label**



dolphin

cat

grizzly bear

angel fish

chameleon

**clown fish**

iguana

elephant

...

"grizzly bear"

Loss ⟶ **large**

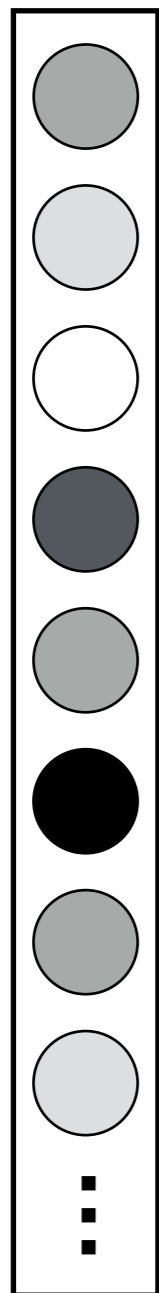# Loss function for classification

Network output     Ground truth label

$\hat{\mathbf{z}}$               $\mathbf{z}$

dolphin

cat

**grizzly bear**
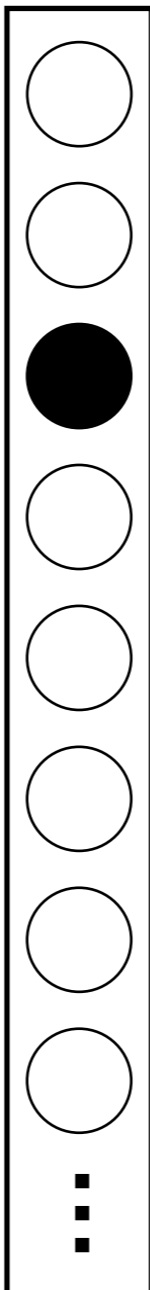
angel fish

chameleon

**clown fish**

iguana

elephant

**Probability of the observed data under the model**

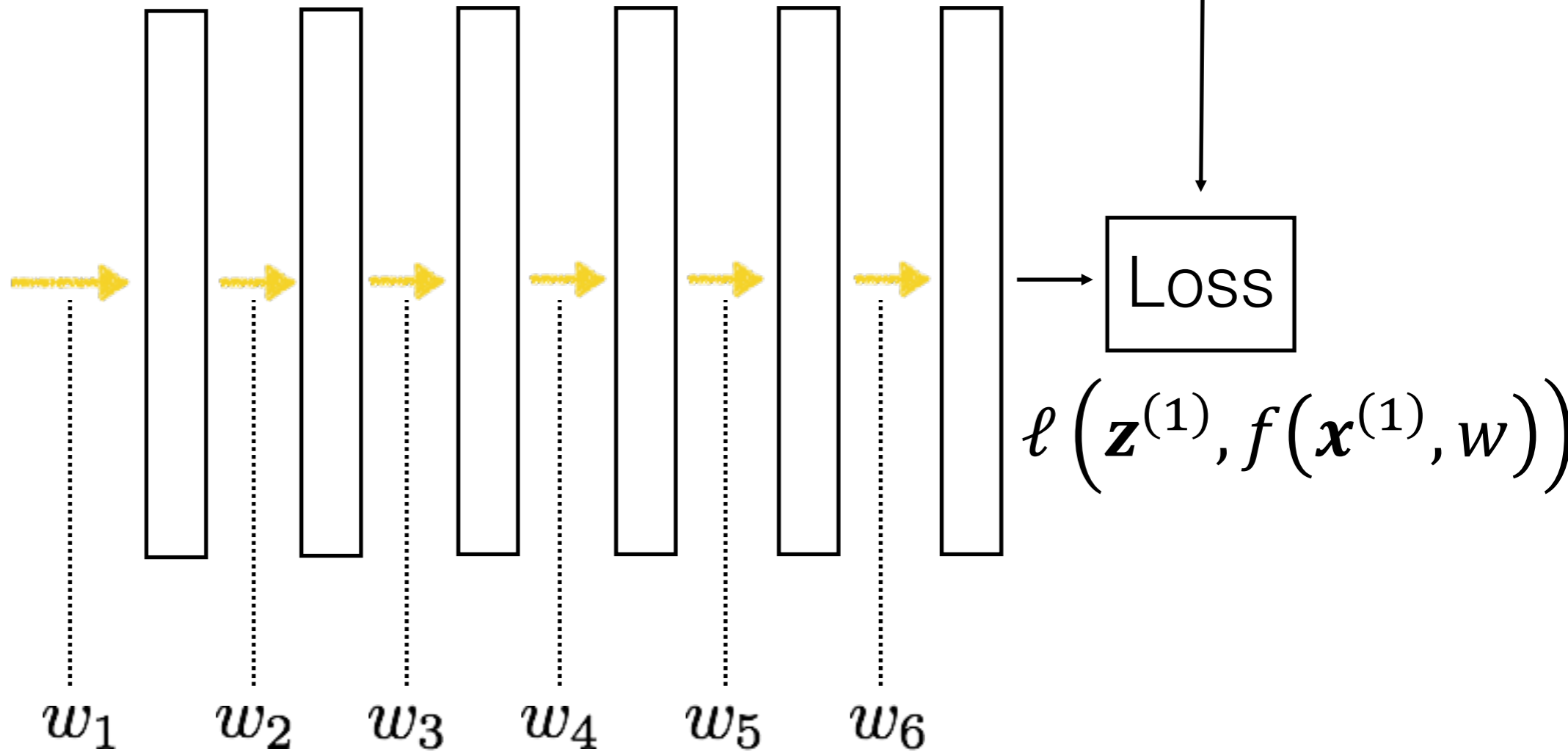$$H(\hat{z}, z) = -\sum_c z_c \log \hat{z}_c$$

*Cross-entropy loss*

$c$ is the $c^{th}$ class in the output

# Learning with deep nets

$\boldsymbol{z}^{(1)}$
"clown fish"

$\boldsymbol{x}^{(1)}$

Learned

Loss

$\ell\left(\boldsymbol{z}^{(1)}, f\big(\boldsymbol{x}^{(1)}, w\big)\right)$

$w_1 \qquad w_2 \qquad w_3 \qquad w_4 \qquad w_5 \qquad w_6$

$\boldsymbol{x}^{(1)}, \boldsymbol{z}^{(1)}$ *is the input and label of the 1st training image*

# Learning with deep nets

$\boldsymbol{z}^{(2)}$
"grizzly bear"

$\boldsymbol{x}^{(2)}$

Learned



Loss

$\ell\left(\boldsymbol{z}^{(2)}, f\left(\boldsymbol{x}^{(2)}, w\right)\right)$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

$\boldsymbol{x}^{(2)}, \boldsymbol{z}^{(2)}$ *is the input and label of the* 2*nd training image*

41

# Learning with deep nets

$z^{(3)}$
"chameleon"

$x^{(3)}$



Learned

Loss

$\ell\left(\mathbf{z}^{(3)}, f(\mathbf{x}^{(3)}, w)\right)$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

$$argmin_w \sum_i \ell\left(z^{(i)}, f(x^{(i)}, w)\right)$$

# Gradient descent

$$argmin_w \sum_i \ell\big(z^{(i)}, f(x^{(i)}, w)\big) = argmin_w L(w)$$

One iteration of gradient descent:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \frac{\partial L(\mathbf{w^t})}{\partial \mathbf{w}}$$

learning rate

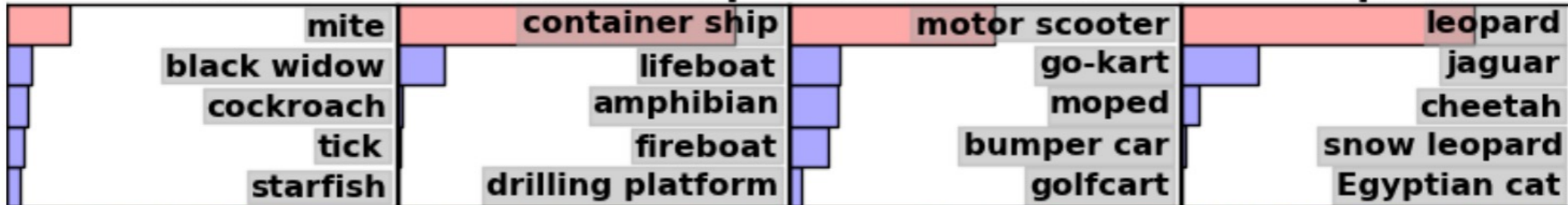# Gradient descent

$$p(c|\mathbf{x})$$



| | | | |
|---|---|---|---|
| **mite** | **container ship** | **motor scooter** | **leopard** |
| mite | container ship | motor scooter | leopard |
| black widow | lifeboat | go-kart | jaguar |
| cockroach | amphibian | moped | cheetah |
| tick | fireboat | bumper car | snow leopard |
| starfish | drilling platform | golfcart | Egyptian cat |
| **grille** | **mushroom** | **cherry** | **Madagascar cat** |
| convertible | agaric | dalmatian | squirrel monkey |
| grille | mushroom | grape | spider monkey |
| pickup | jelly fungus | elderberry | titi |
| beach wagon | gill fungus | ffordshire bullterrier | indri |
| fire engine | dead-man's-fingers | currant | howler monkey |

[Krizhevsky et al. NIPS 2012]

# Computer Vision before 2012



Features     Clustering     Pooling   Classification     Cat
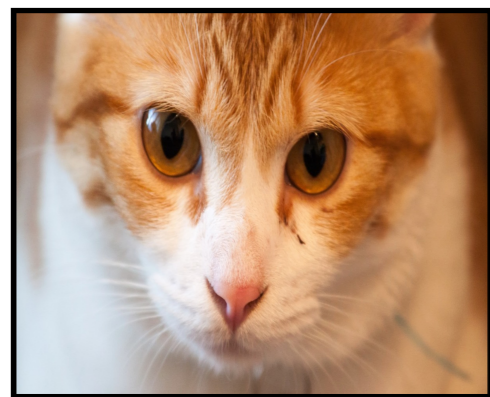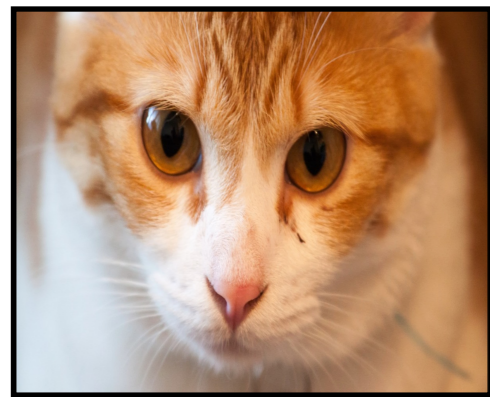
# Computer Vision Now



Features　Clustering　Pooling　Classification　→ Cat
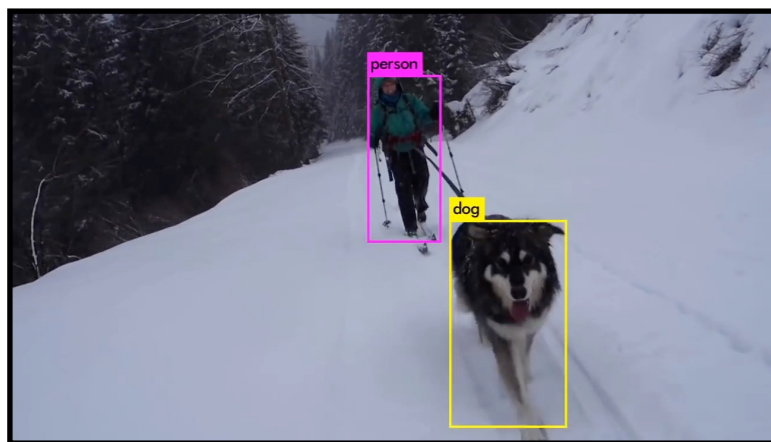
→ Deep Net → Cat

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

# Deep Learning for Computer Vision
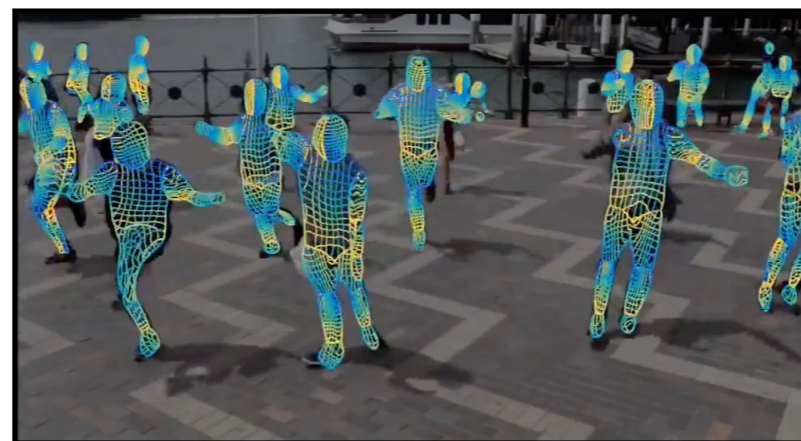


[Deng et al. 2009]



Deep Net

Top 5 **accuracy** on ImageNet benchmark



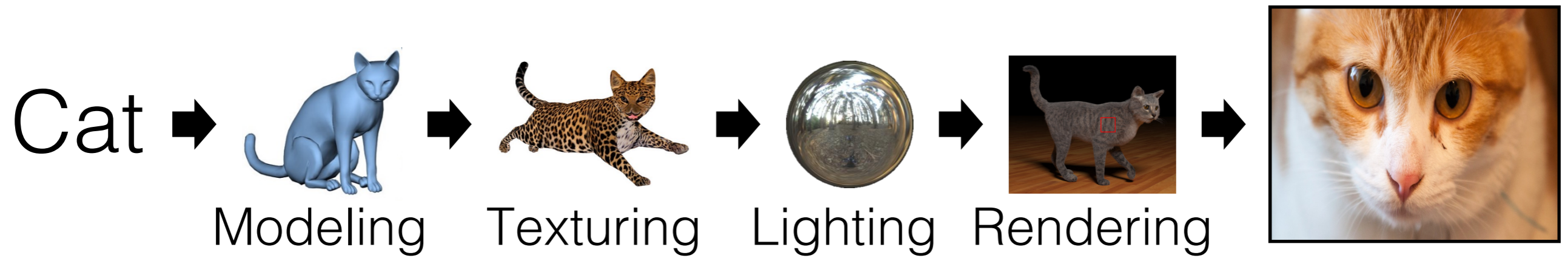[Redmon et al., 2018]

Object detection



[Güler et al., 2018]

Human understanding

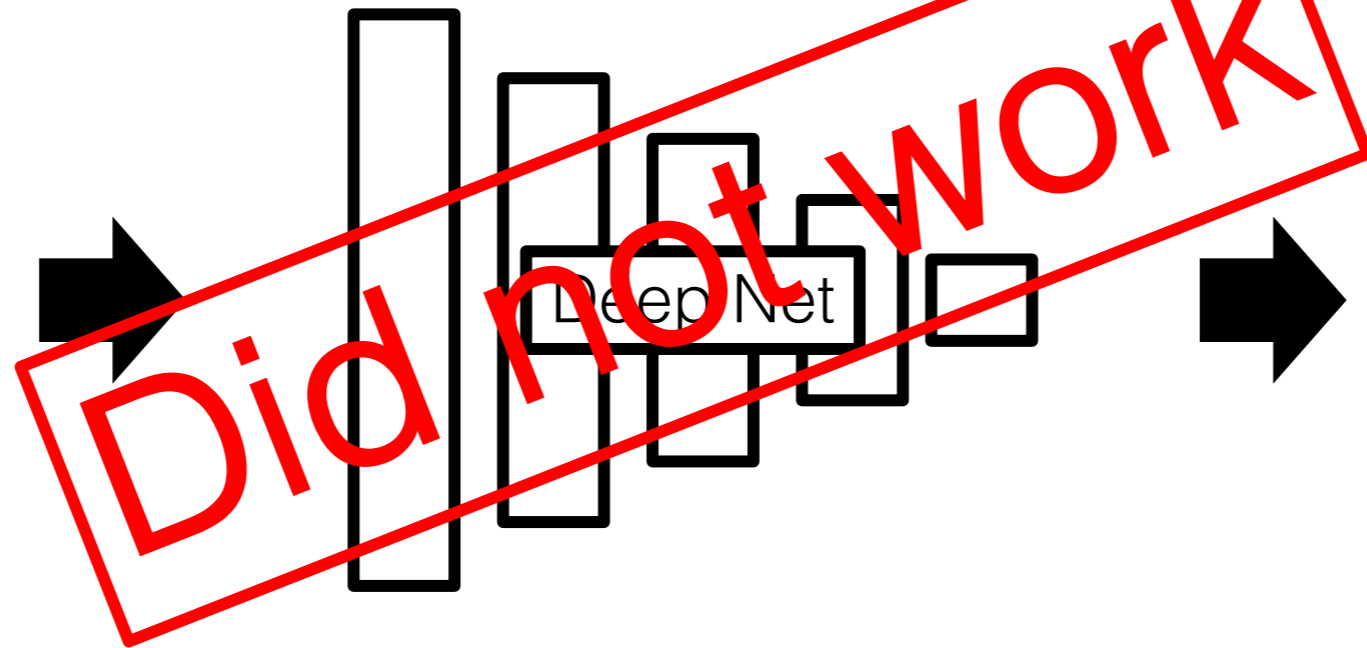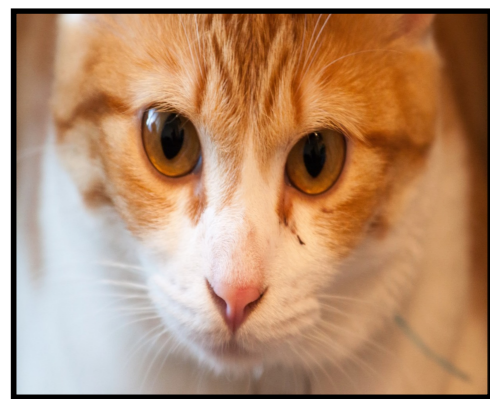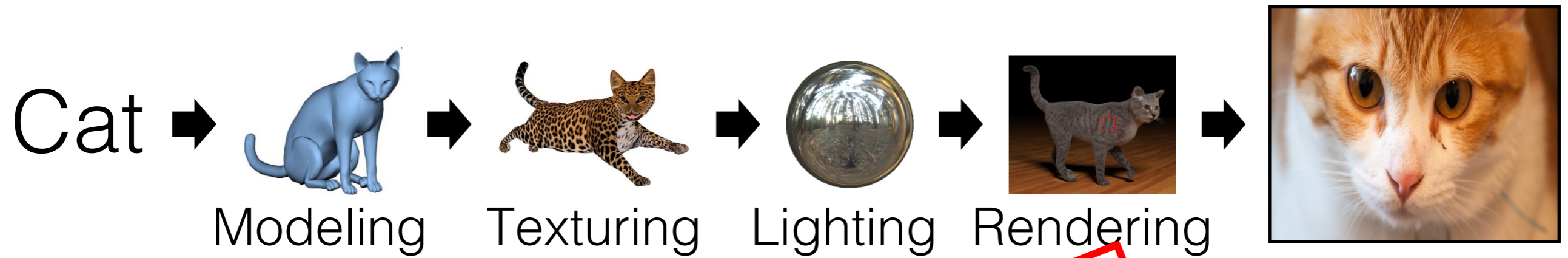

[Zhao et al., 2017]
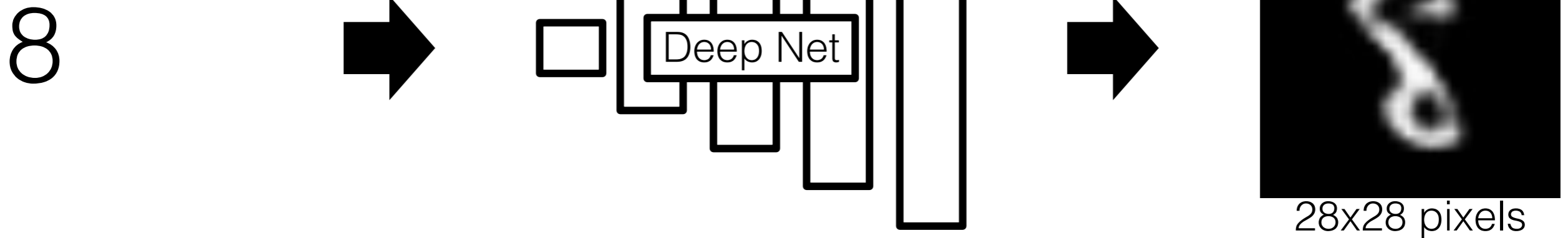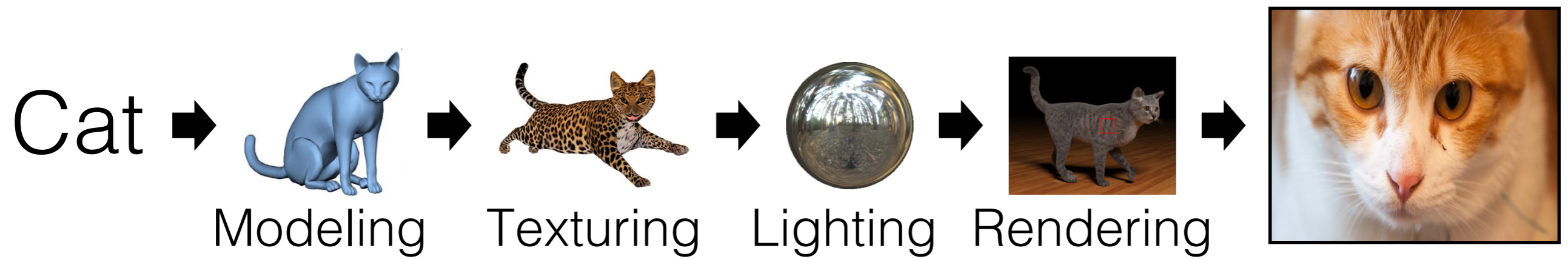
Autonomous driving

# Can Deep Learning **Help** Graphics?

Cat ➡  ➡  ➡  ➡  ➡ 

Modeling　　Texturing　　Lighting　　Rendering

# Can Deep Learning **Help** Graphics?



Cat ➡ Modeling ➡ Texturing ➡ Lighting ➡ Rendering ➡

Did not work

Deep Net ➡ Cat

# Generating images is hard!

Cat → Modeling → Texturing → Lighting → Rendering →

8 → Deep Net → 28x28 pixels
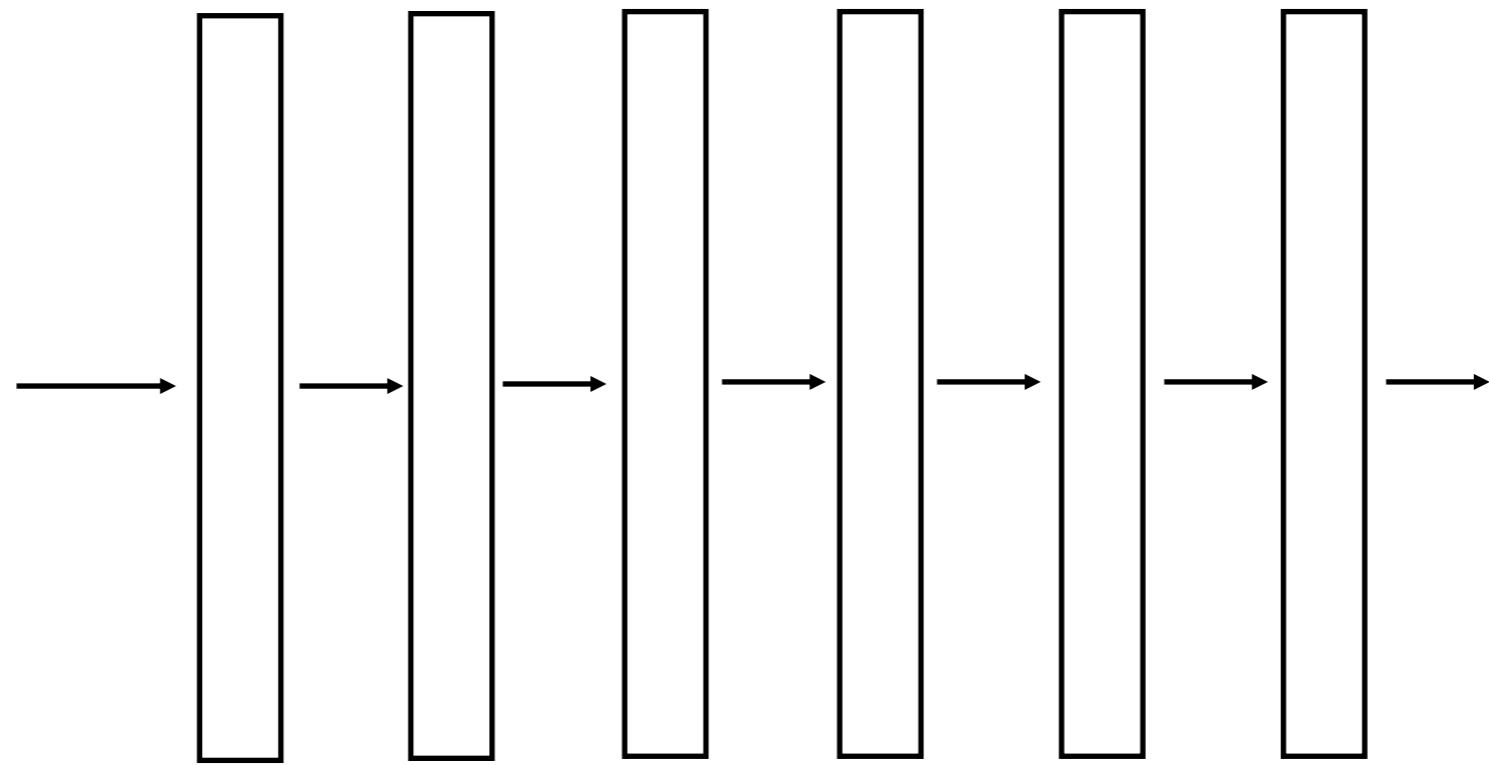
# from Classification to Generation

"yellow"

Predicting the color value of an output pixel given a patch

# Discriminative Deep Networks



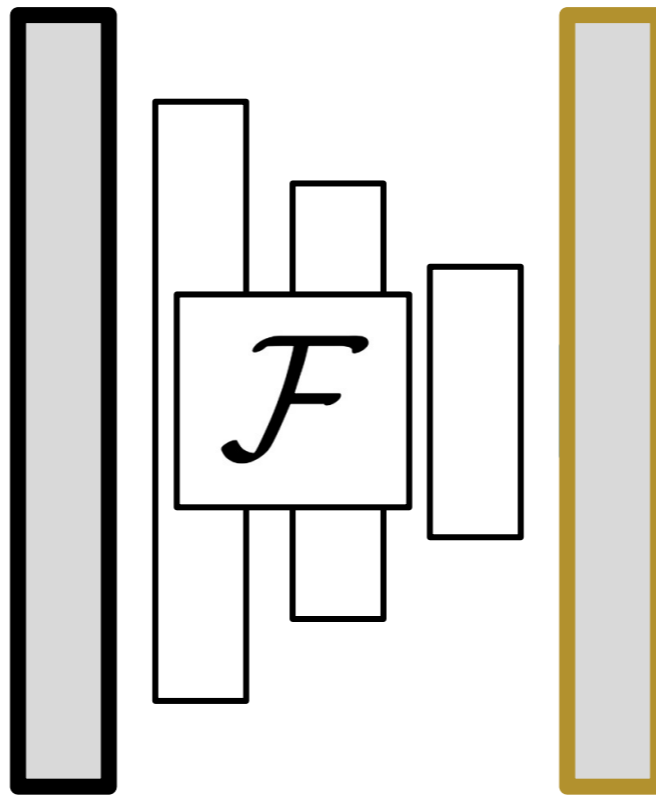$\mathcal{F}$

"Rockfish"

# Discriminative Deep Networks



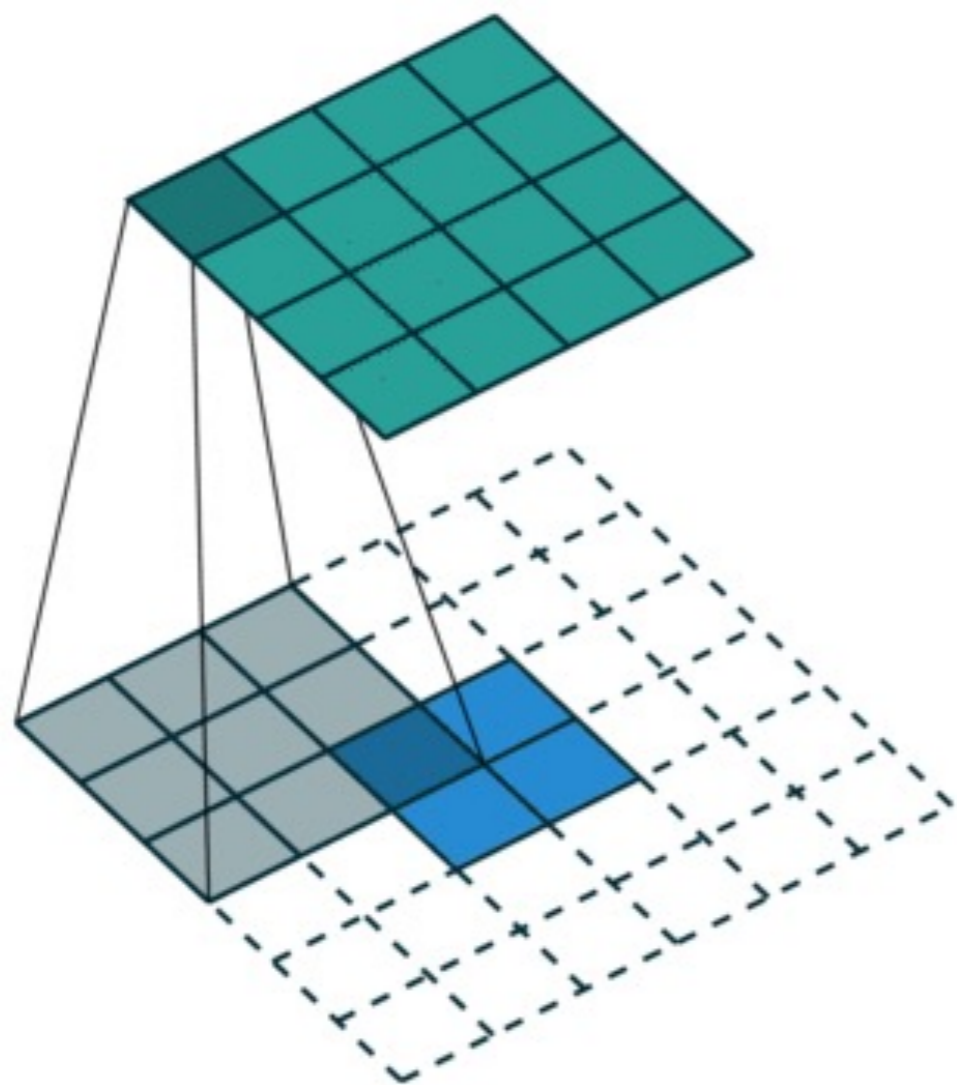Raw, Unlabeled Pixels

55

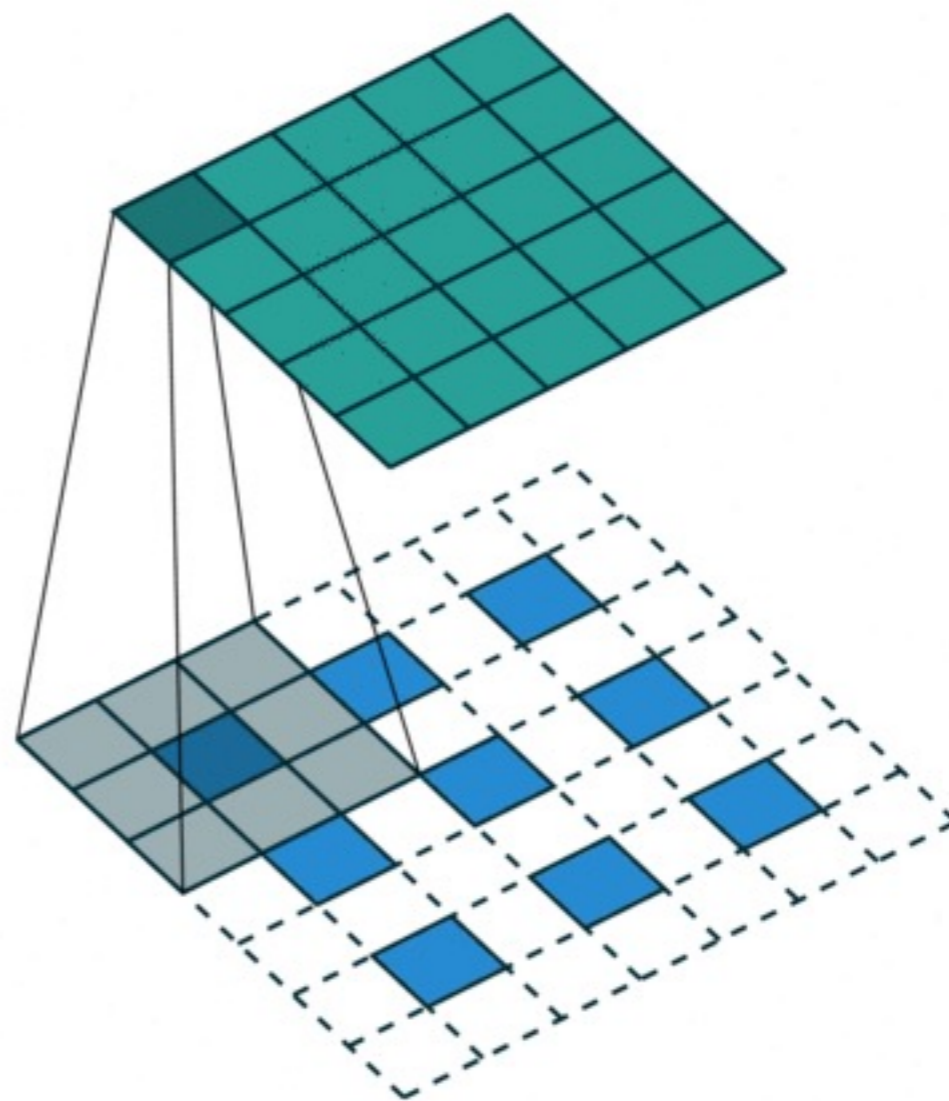# Generative Deep Networks



Raw, Unlabeled Pixels

56

# Better Architectures

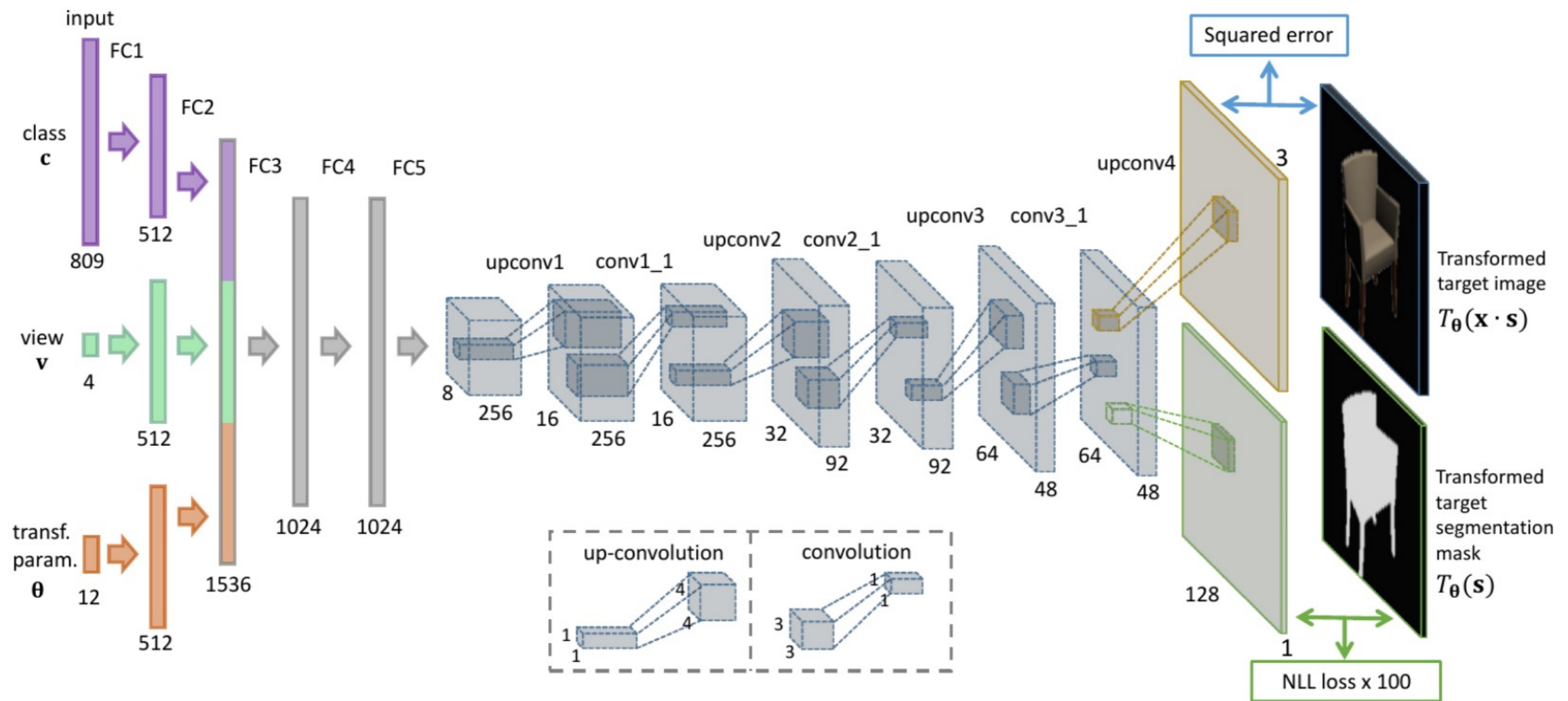# Fractionally-strided Convolution



Regular conv       Fractiaionally-strided conv
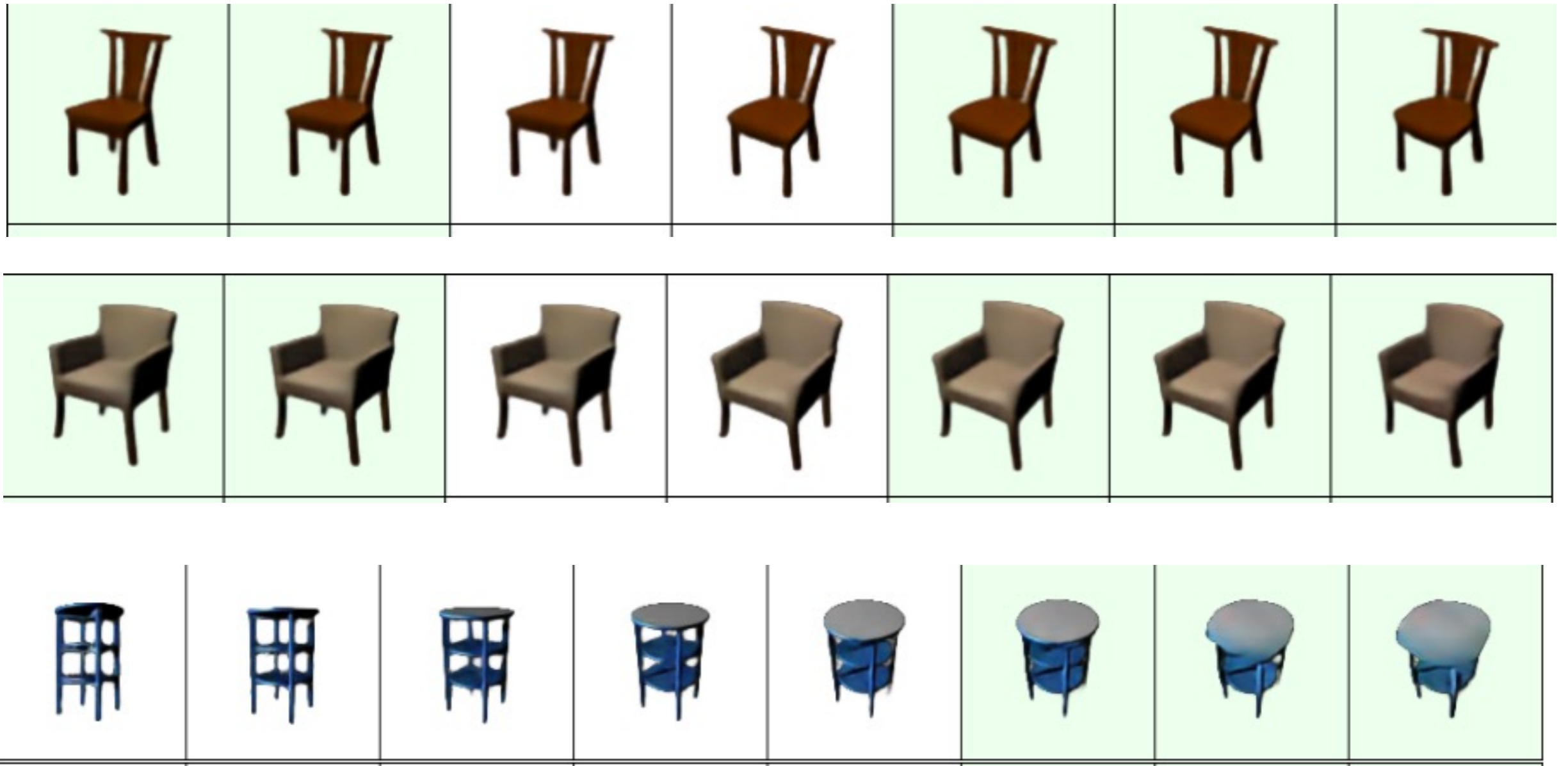
# Generating chairs conditional on chair ID, viewpoint, and transformation parameters



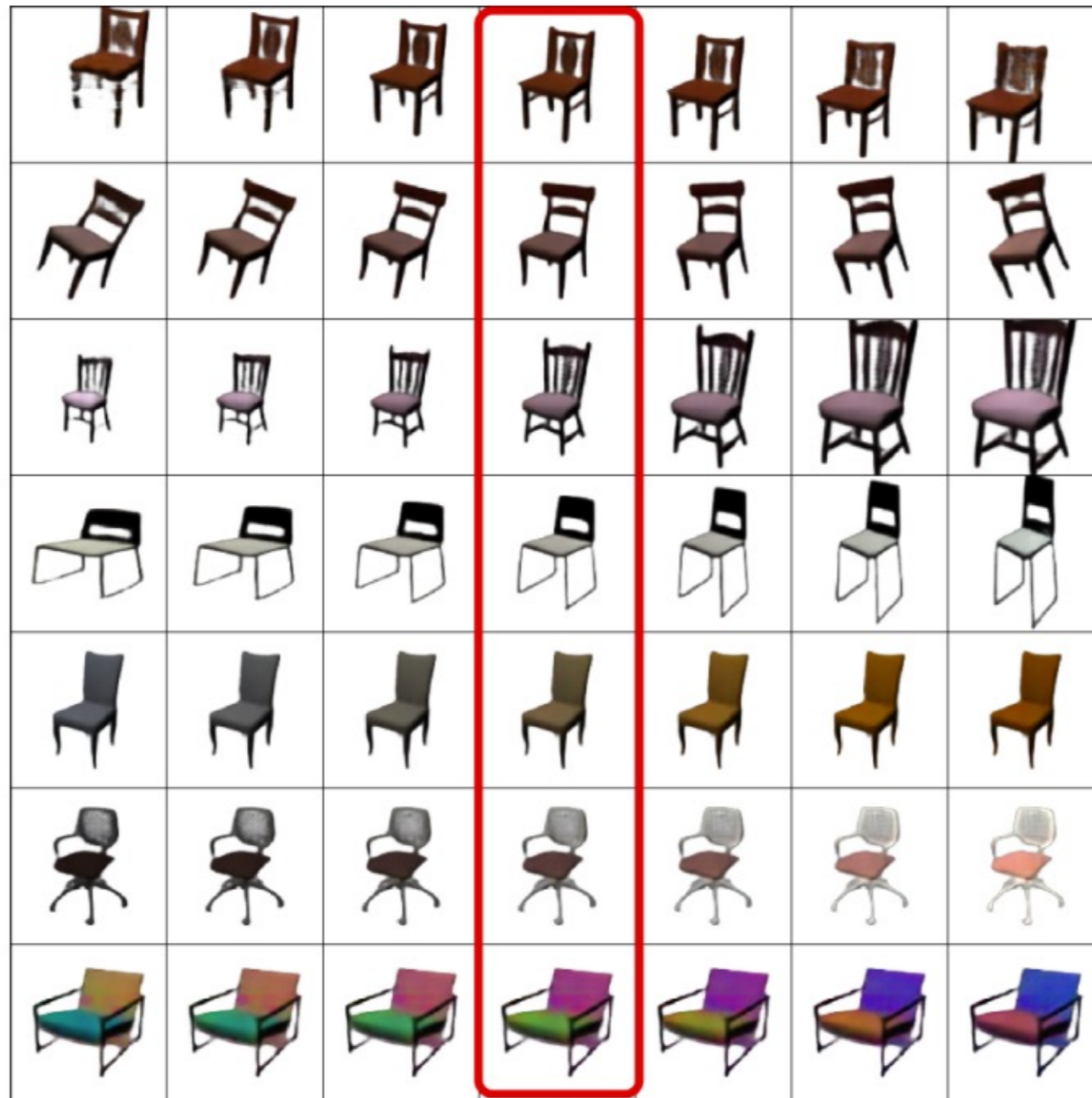Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks PAMI 2017 (CVPR 2015)

# With Varying Viewpoints



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

# With Varying Transformation Parameters

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

# Interpolation between Two Chairs



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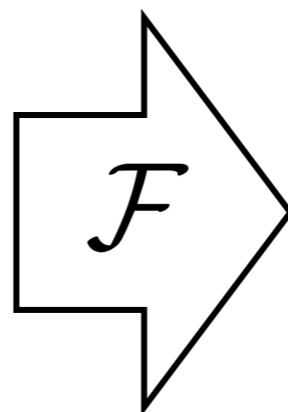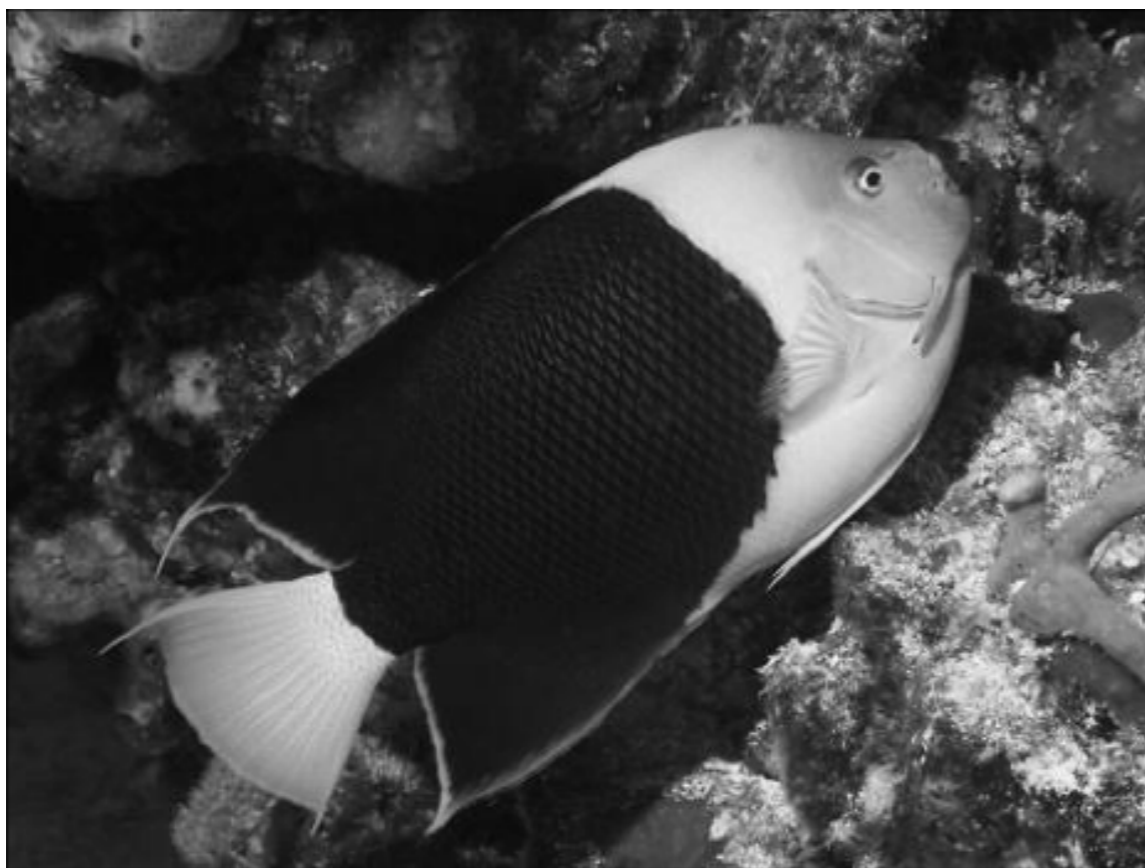PAMI 2017 (CVPR 2015)

# Better Loss Functions

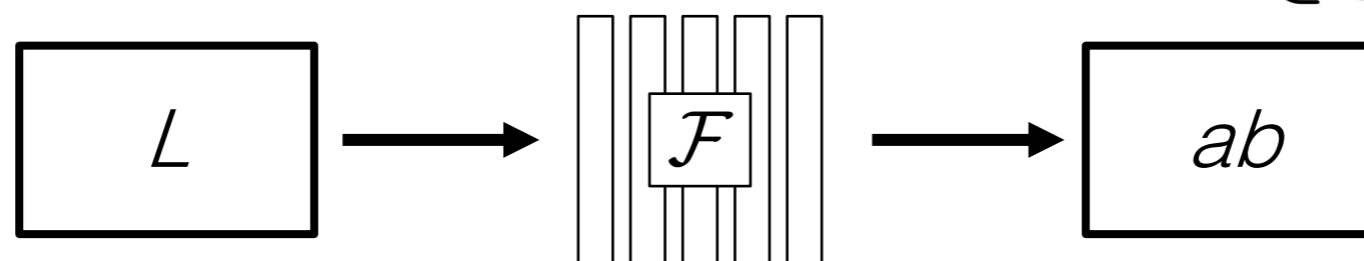Ansel Adams. *Yosemite Valley Bridge.*
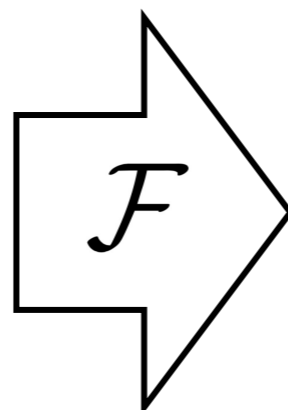
Grayscale image: *L* channel
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: *ab* channels
$$\widehat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

Zhang, Isola, Efros. *Colorful Image Colorization.* In *ECCV*, 2016.

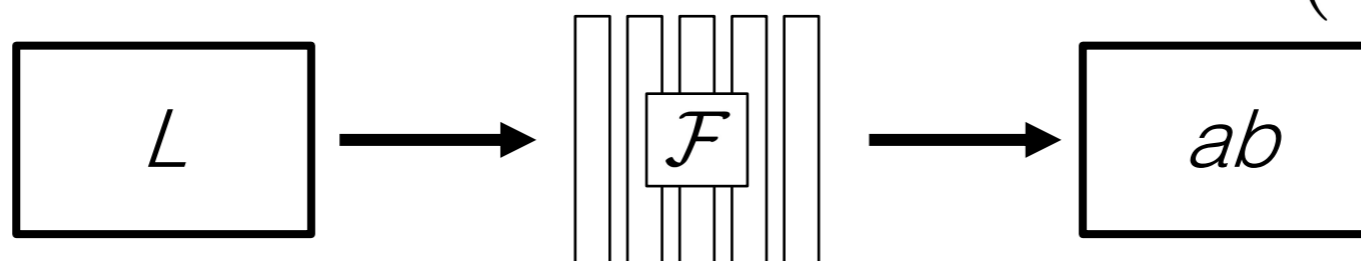Grayscale image: *L* channel

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

Concatenate (*L*,*ab*) channels

$$(\mathbf{X}, \widehat{\mathbf{Y}})$$

$L$ → $\mathcal{F}$ → $ab$

Zhang, Isola, Efros. *Colorful Image Colorization.* In *ECCV*, 2016.

# Simple L2 regression doesn't work ☹

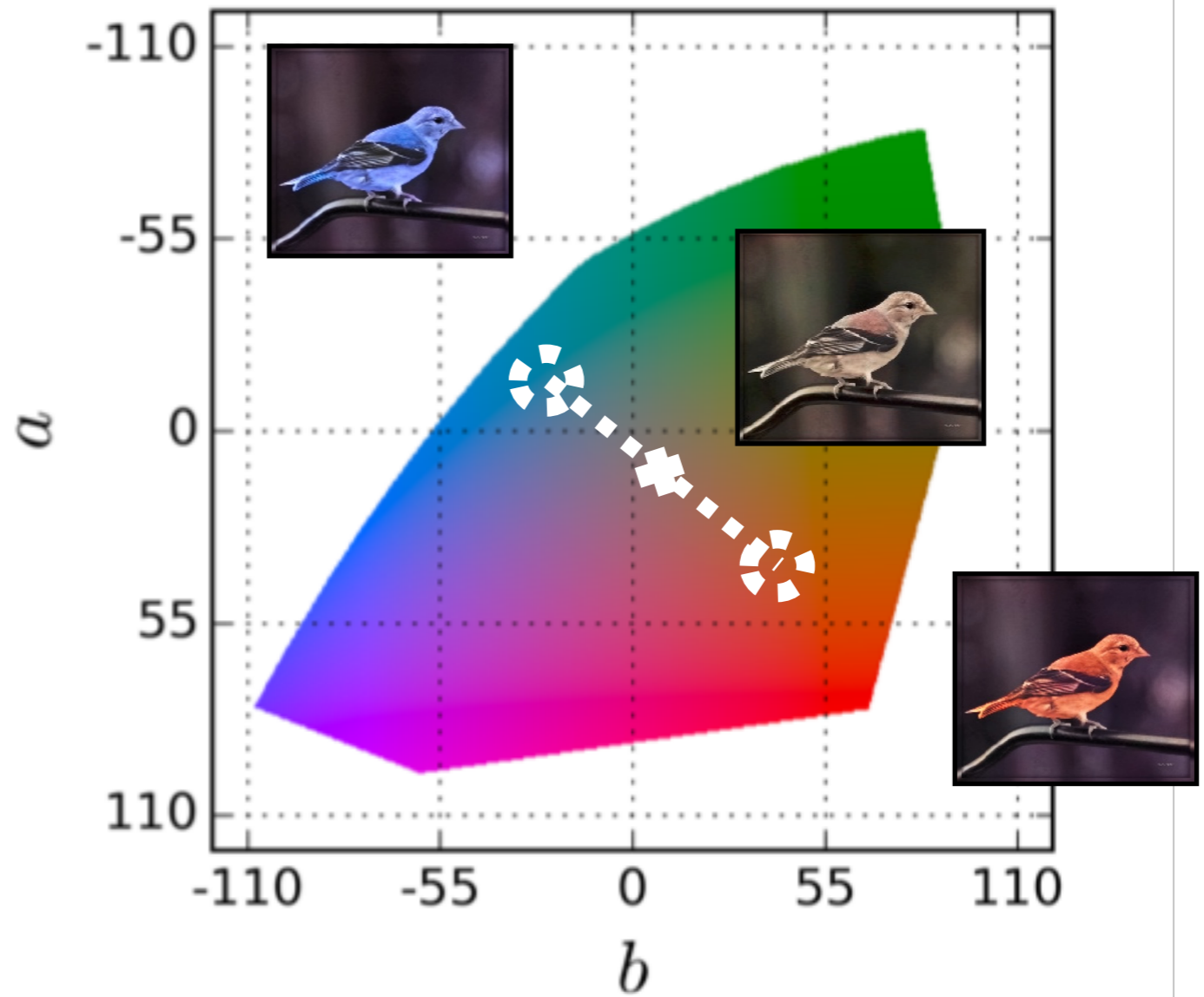Input           Output          Ground truth



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

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

# Better Loss Function

$$\theta^* = \arg\min_{\theta} \ell(\mathcal{F}_\theta(\mathbf{X}), \mathbf{Y})$$
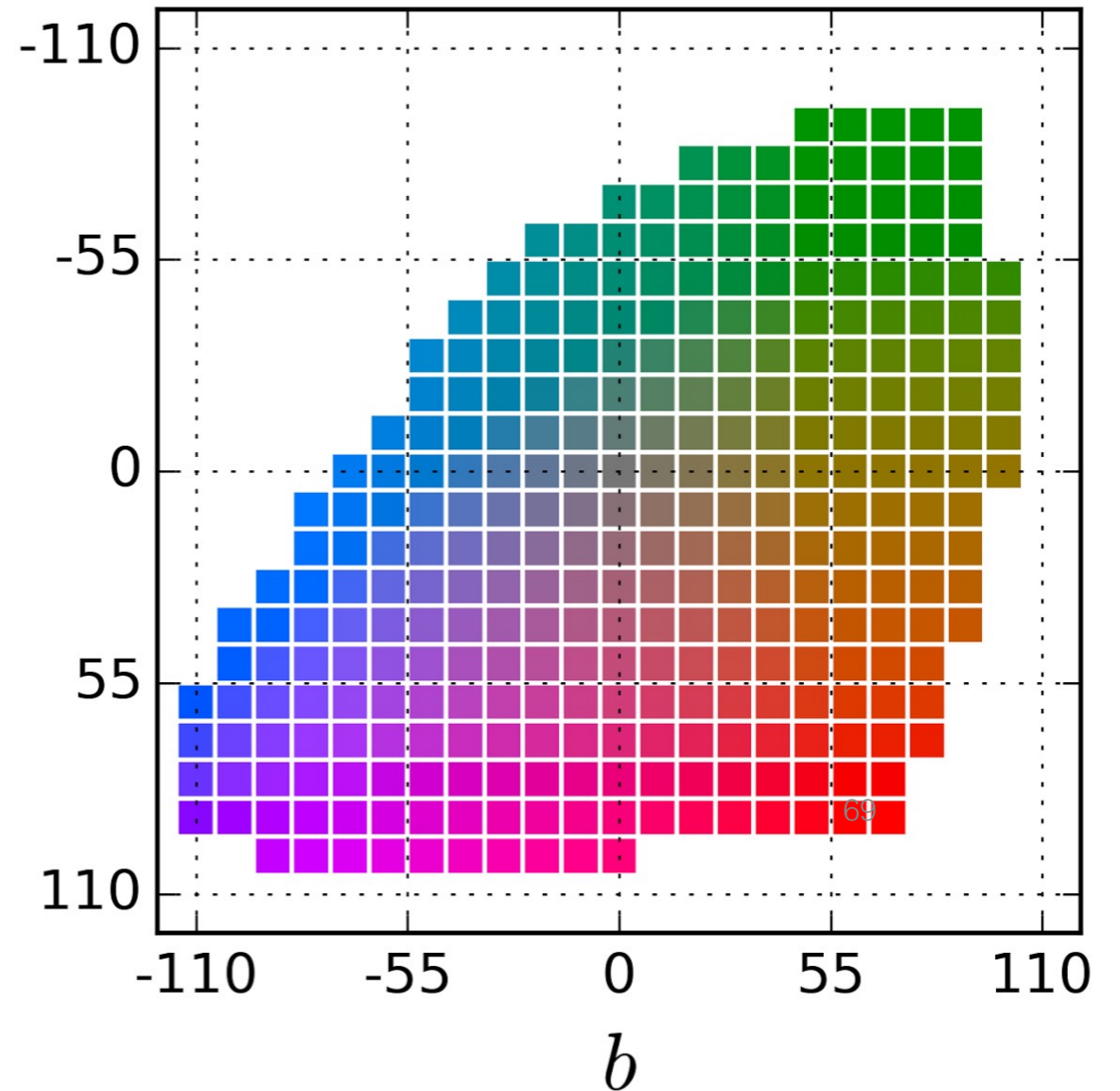
- Regression with L2 loss inadequate

$$L_2(\widehat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \widehat{\mathbf{Y}}_{h,w}\|_2^2$$
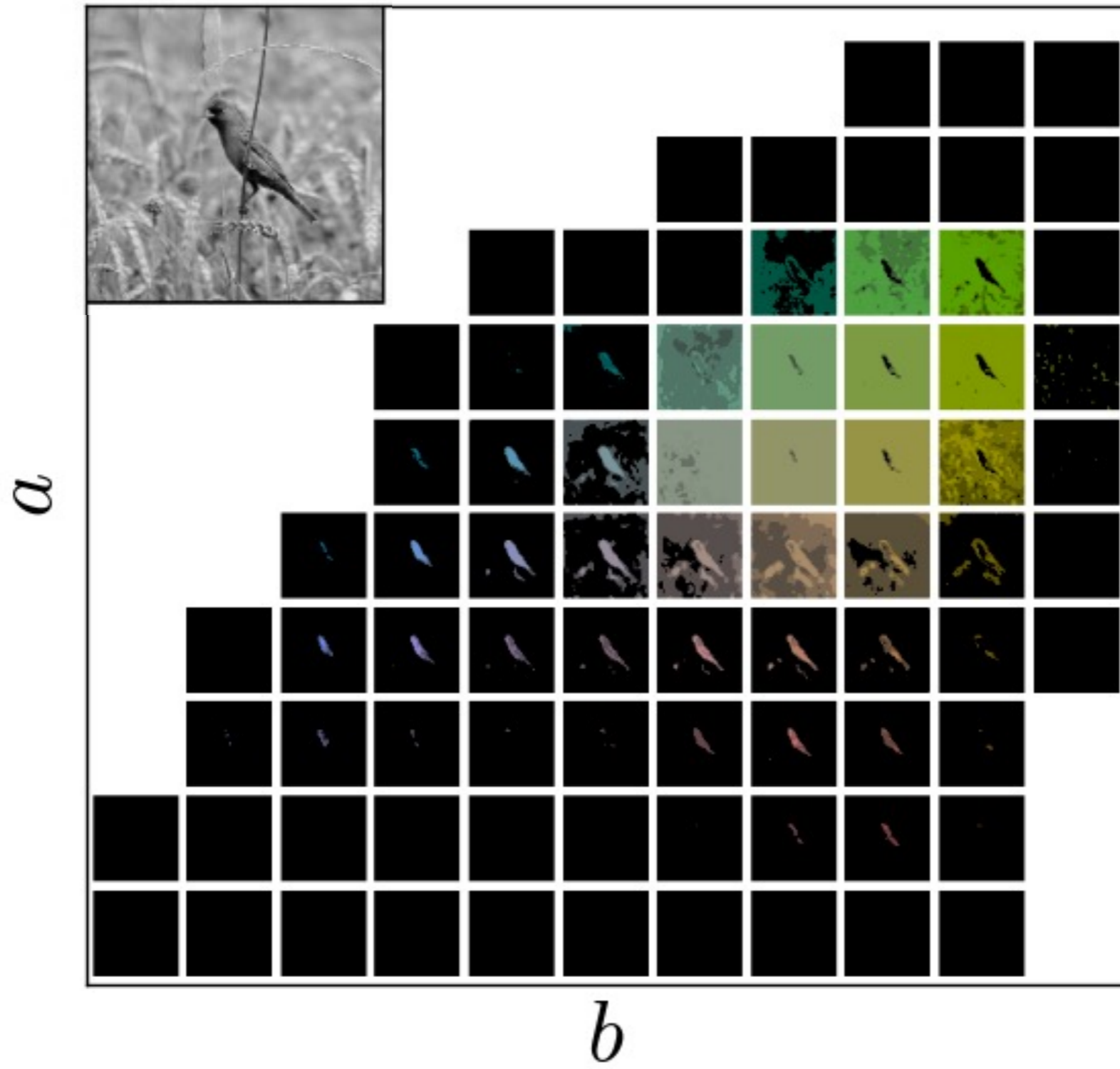
- Use per-pixel multinomial classification

$$L(\widehat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{Z}_{h,w,q} \log(\widehat{\mathbf{Z}}_{h,w,q})$$

**Colors in _ab_ space**
(discrete)

$a$

$b$

# Designing loss functions
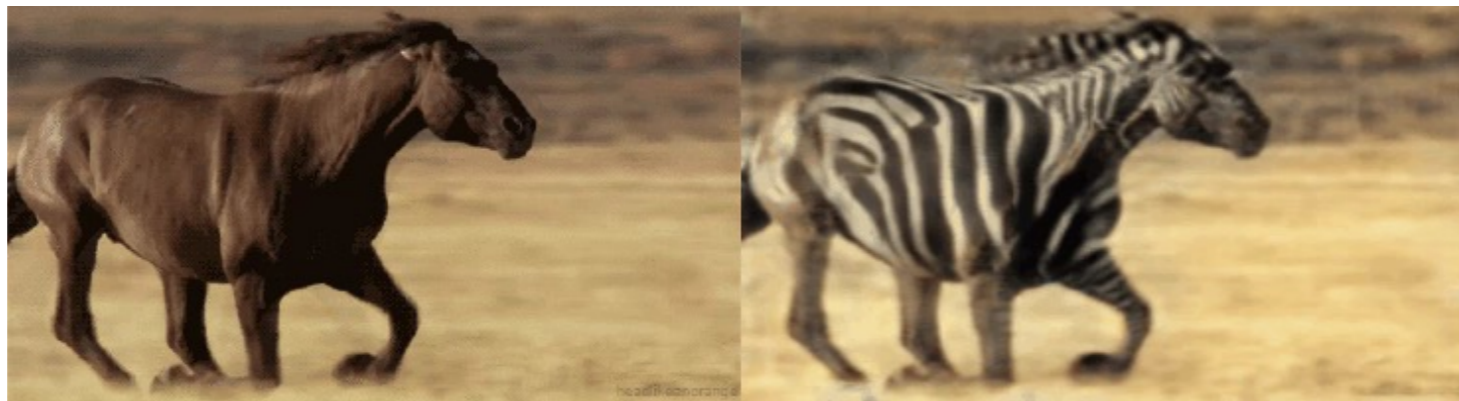
Input  Zhang et al. 2016  Ground truth



Color distribution cross-entropy loss with colorfulness enhancing term.

[Zhang, Isola, Efros, ECCV 2016]

# Thank You!



**16-726, Spring 2022**

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