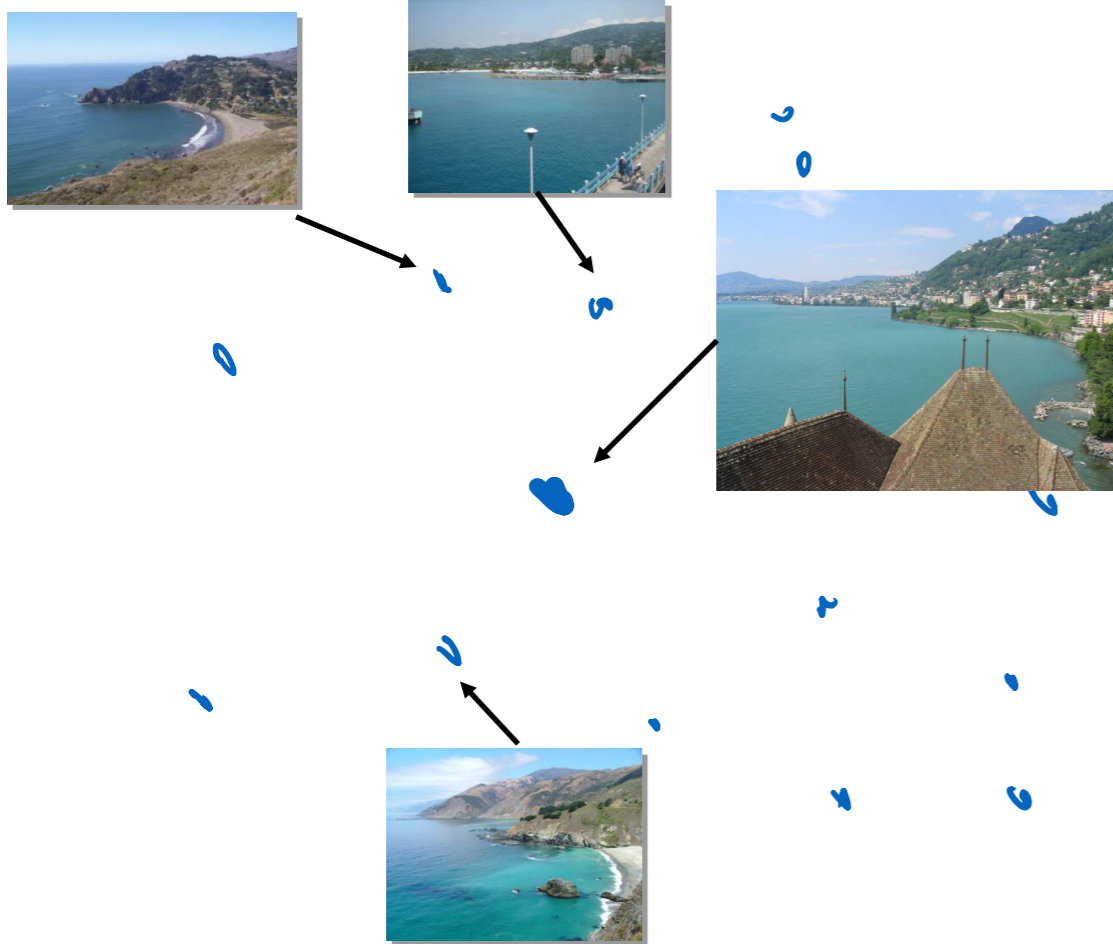


Convolutional Network for Image Synthesis

Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2021

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



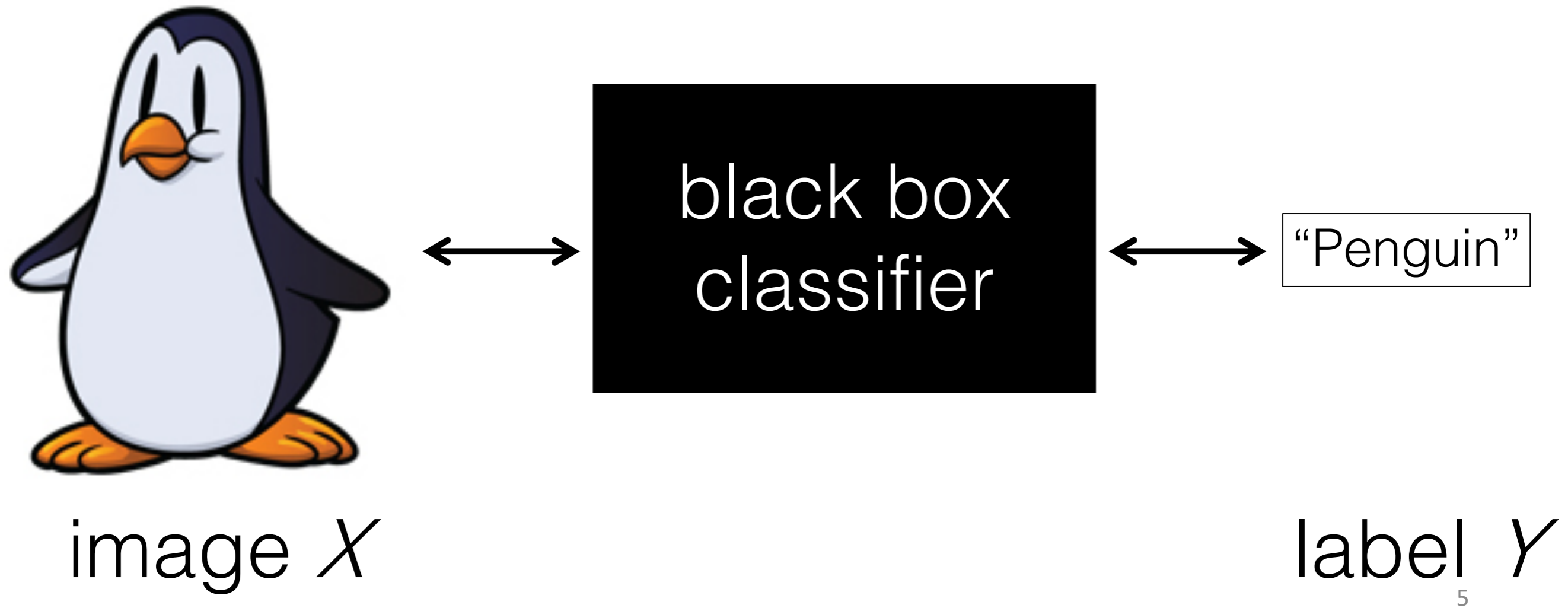
“Penguin”

?
==



“Penguin”

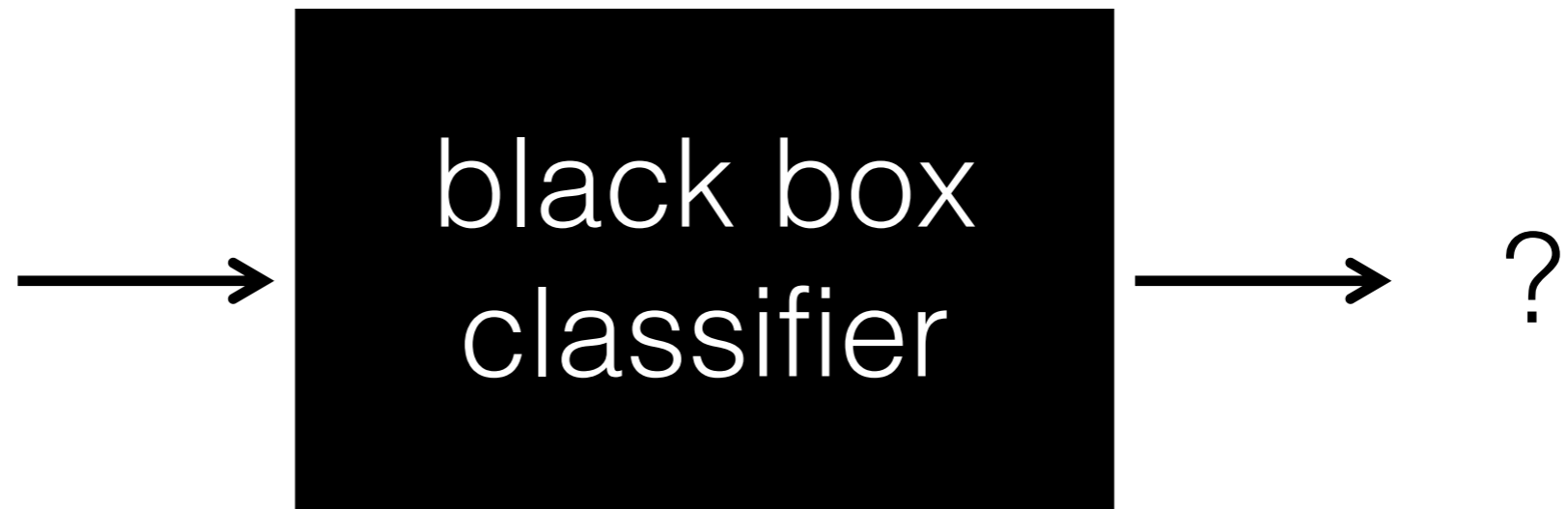
Machine Learning as data association



At test time...



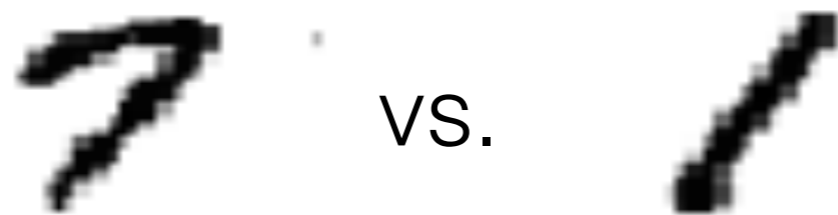
image X



3	6	8	1	7	9	6	6	9	1
6	7	5	7	8	6	3	4	8	5
2	1	7	9	7	1	2	8	4	5
4	8	1	9	0	1	8	8	9	4
7	6	1	8	6	4	1	5	6	0
7	5	9	2	6	5	8	1	9	7
2	2	2	2	2	3	4	4	8	0
0	2	3	8	0	7	3	8	5	7
0	1	4	6	4	6	0	2	4	3
7	1	2	8	7	6	9	8	6	1



Examples from MNIST dataset [LeCun, 1998]

Warm-up Example: Binary Digit Classification



Learning Approach to Digit Recognition

- **Collect Training Images**

- Positive: 
- Negative: 

- **Training Time**

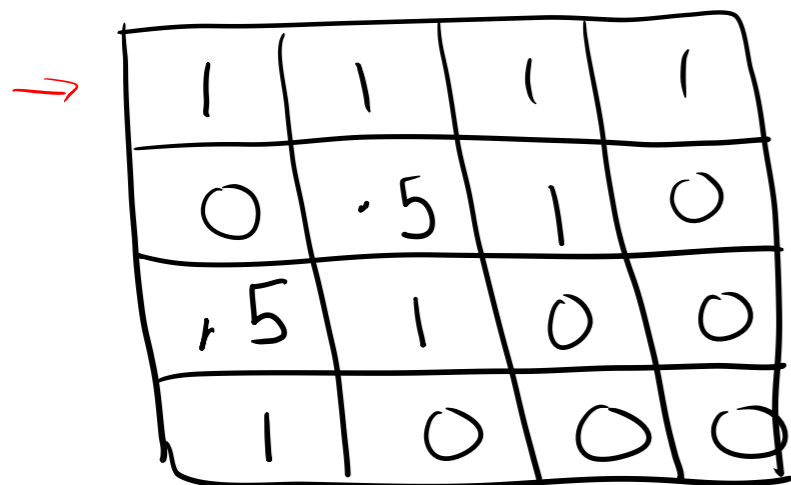
- Compute **feature vectors** for positive and negative example images
- Train a **classifier**

- **Test Time**

- Compute feature vector on new test image:
- Evaluate classifier



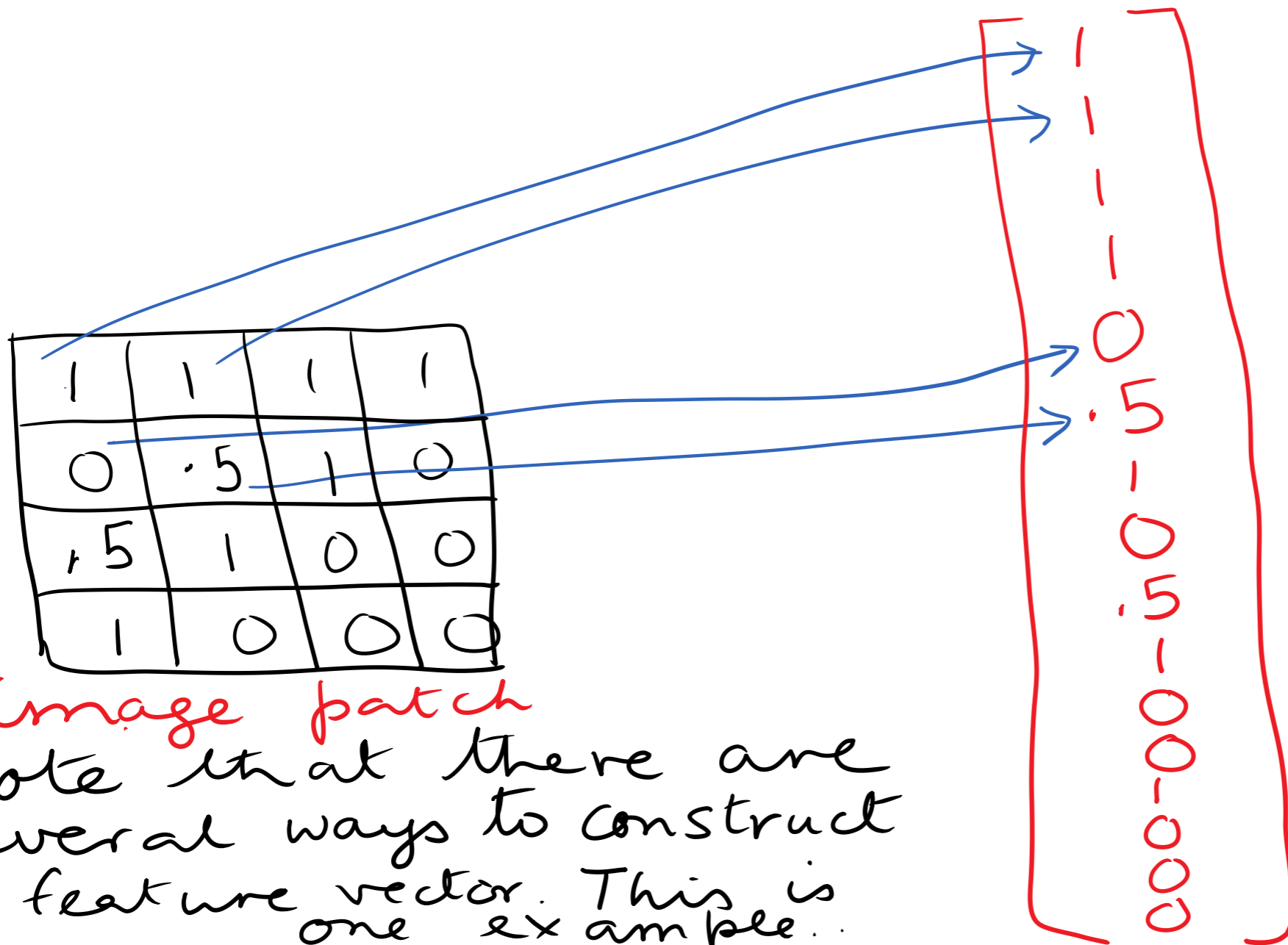
Let us take an example...



1	1	1	1
0	.5	1	0
.5	1	0	0
1	0	0	0

*image
patch*

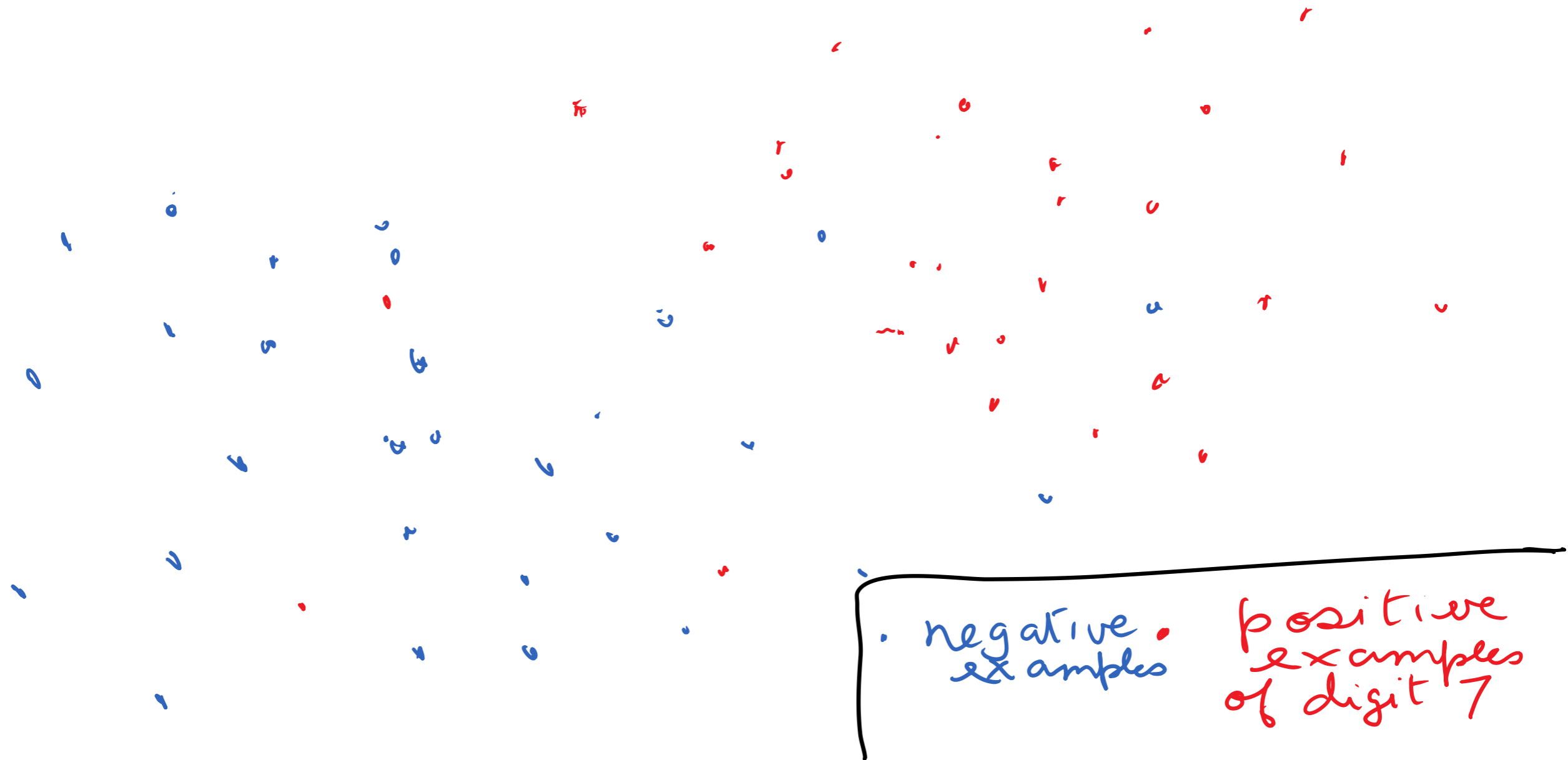
Let us take an example...



Feature
vector
 \mathbb{R}^{16}

image patch
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...



How do we classify a new point?

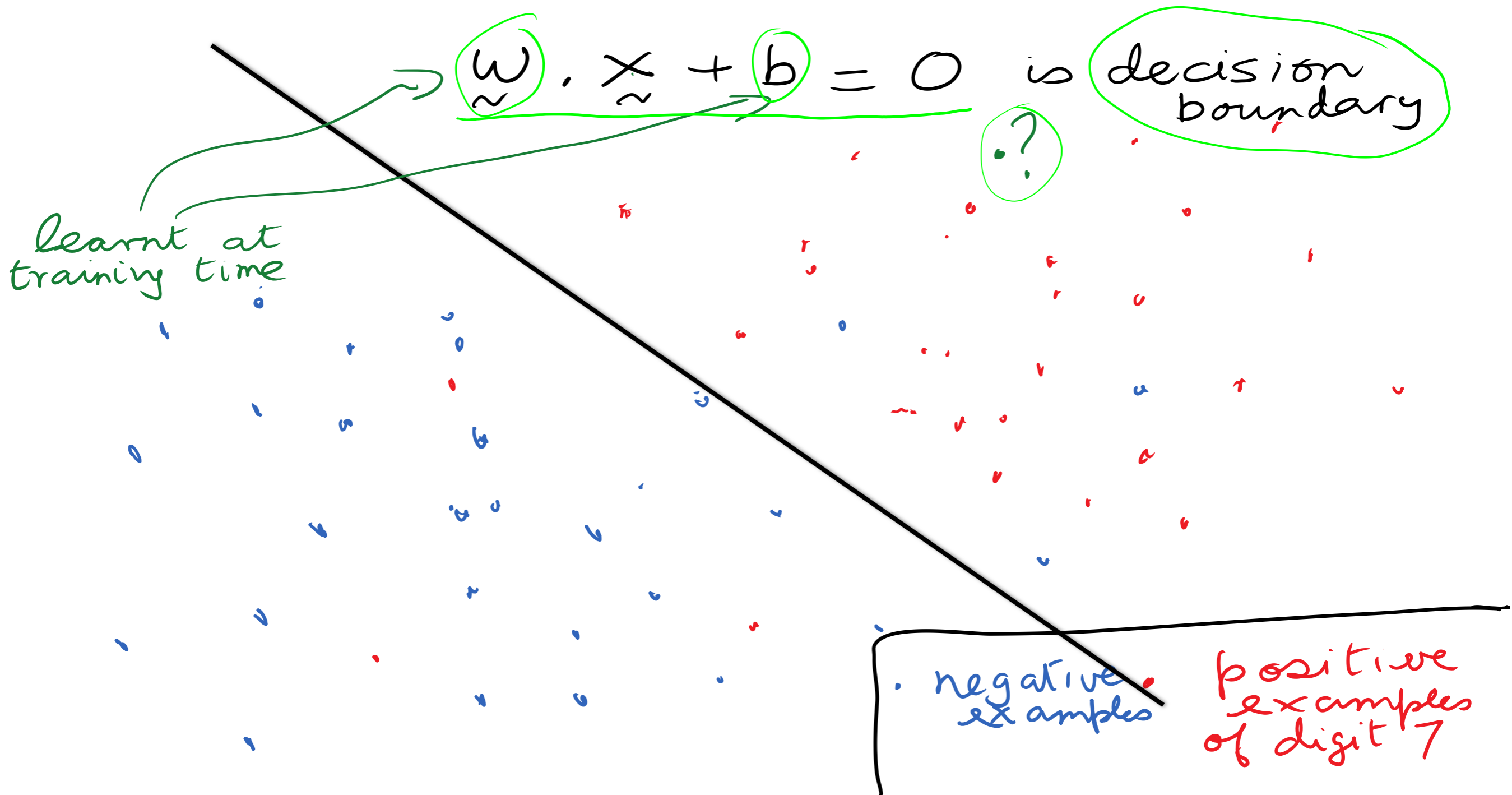


Nearest neighbor rule

“transfer label of nearest example”



Linear classifier rule





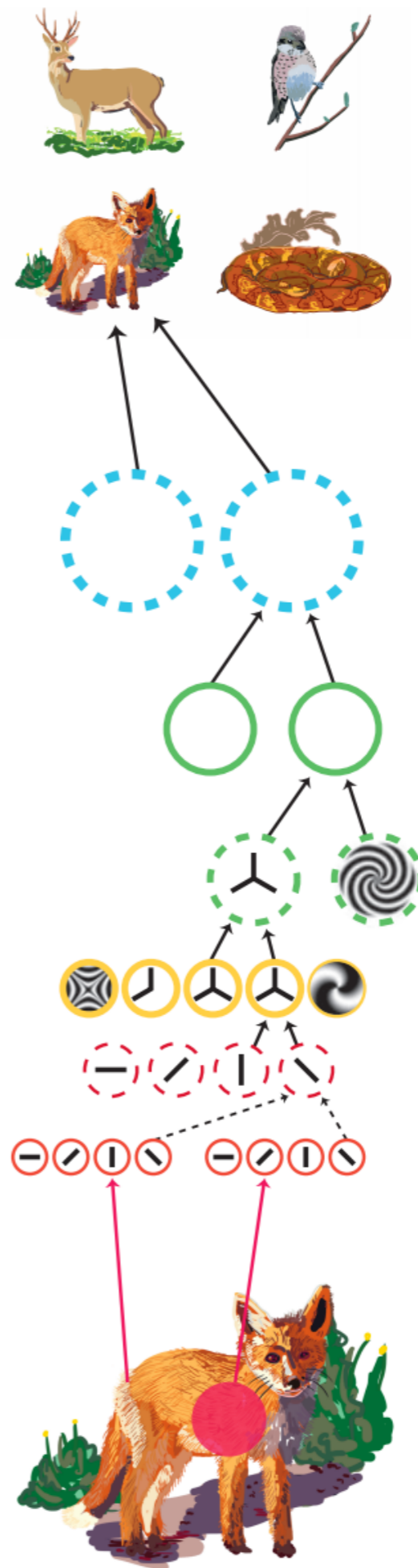
Classification
units

PIT/AIT

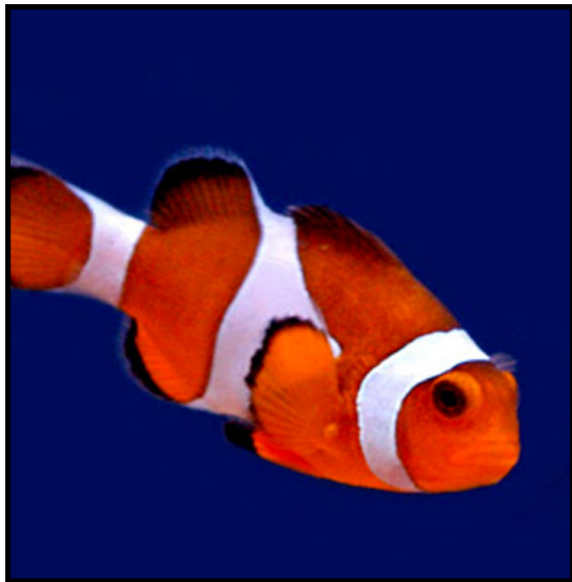
V4/PIT

V2/V4

V1/V2



Basic idea

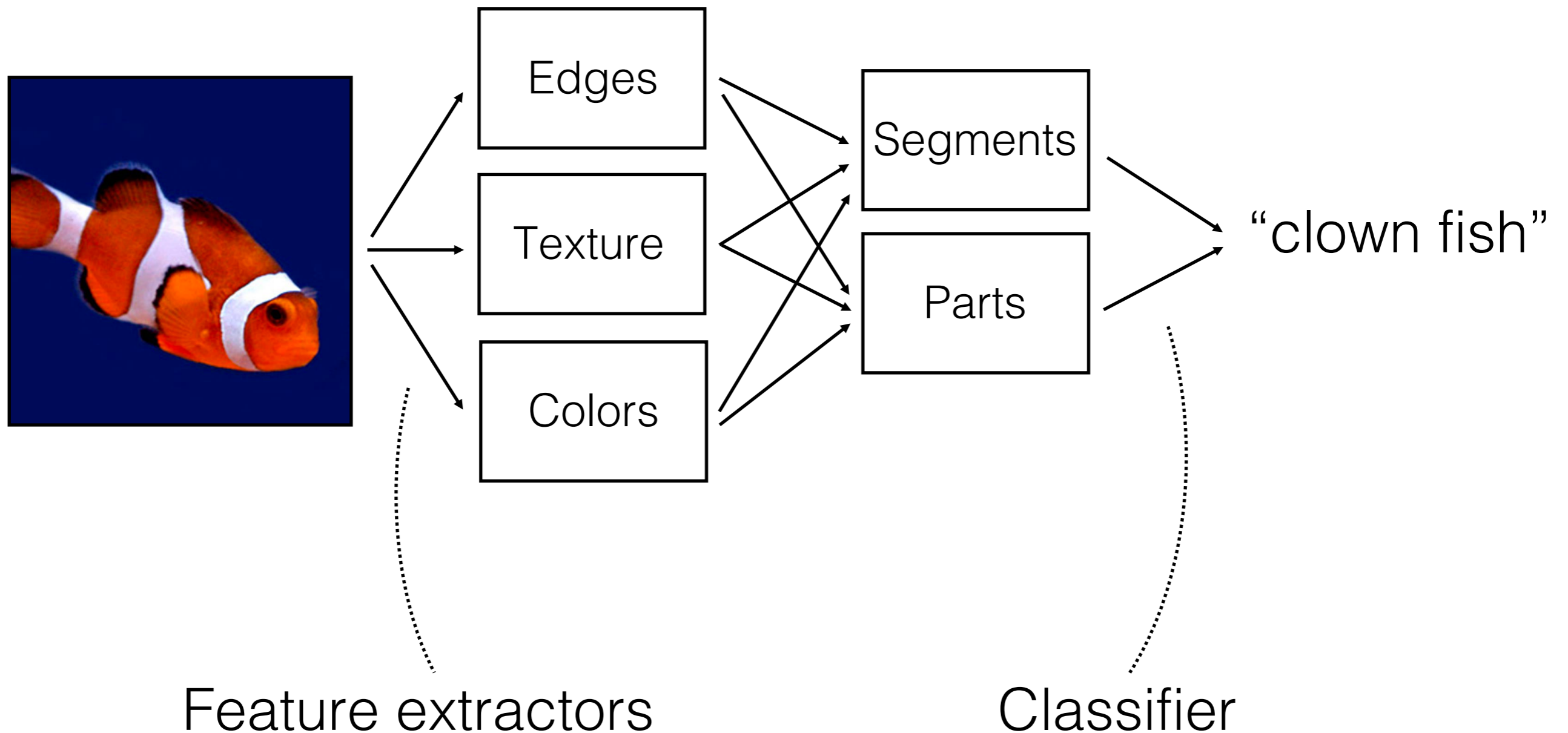


Brain/Machine

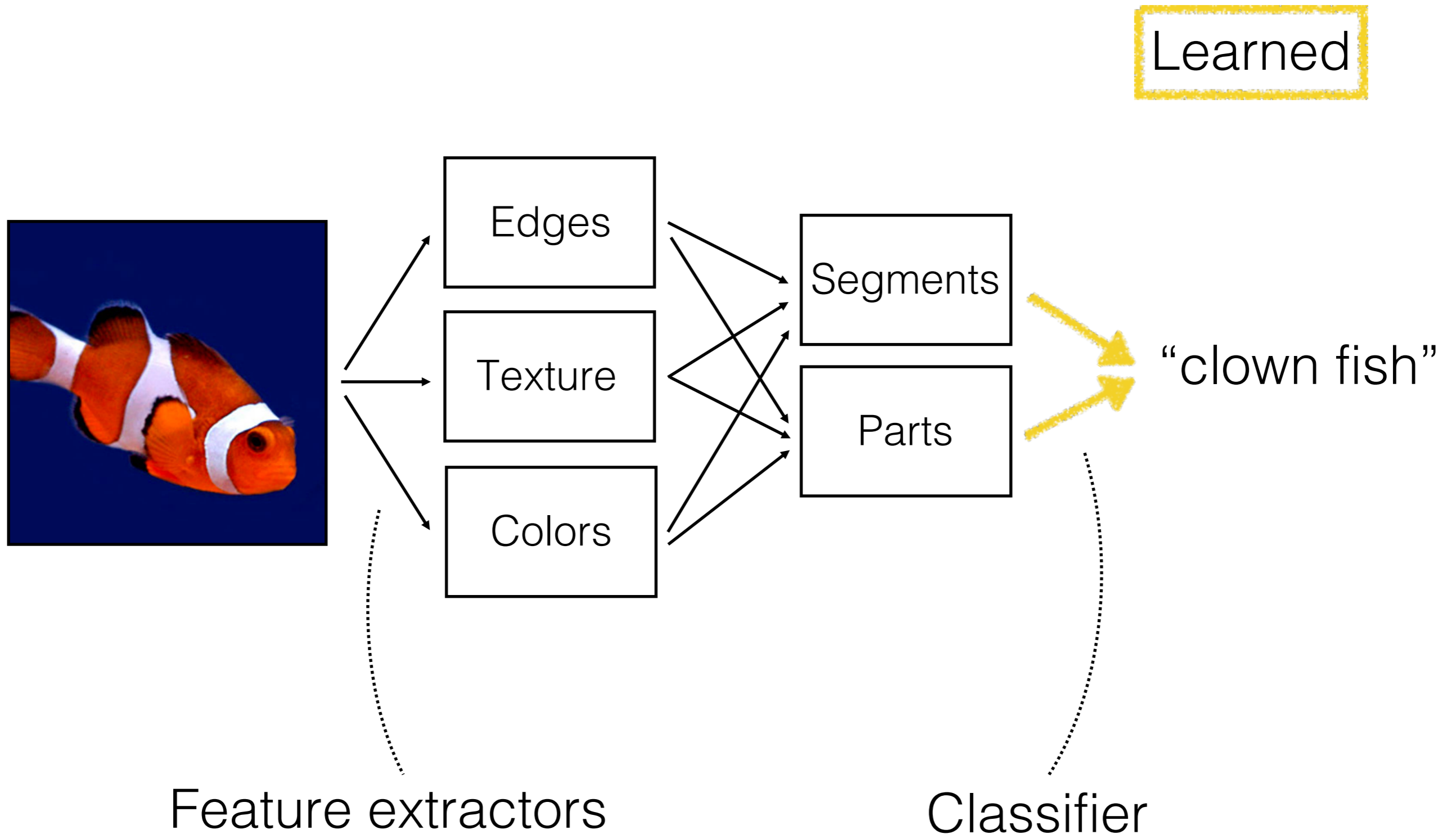


“clown fish”

Object recognition

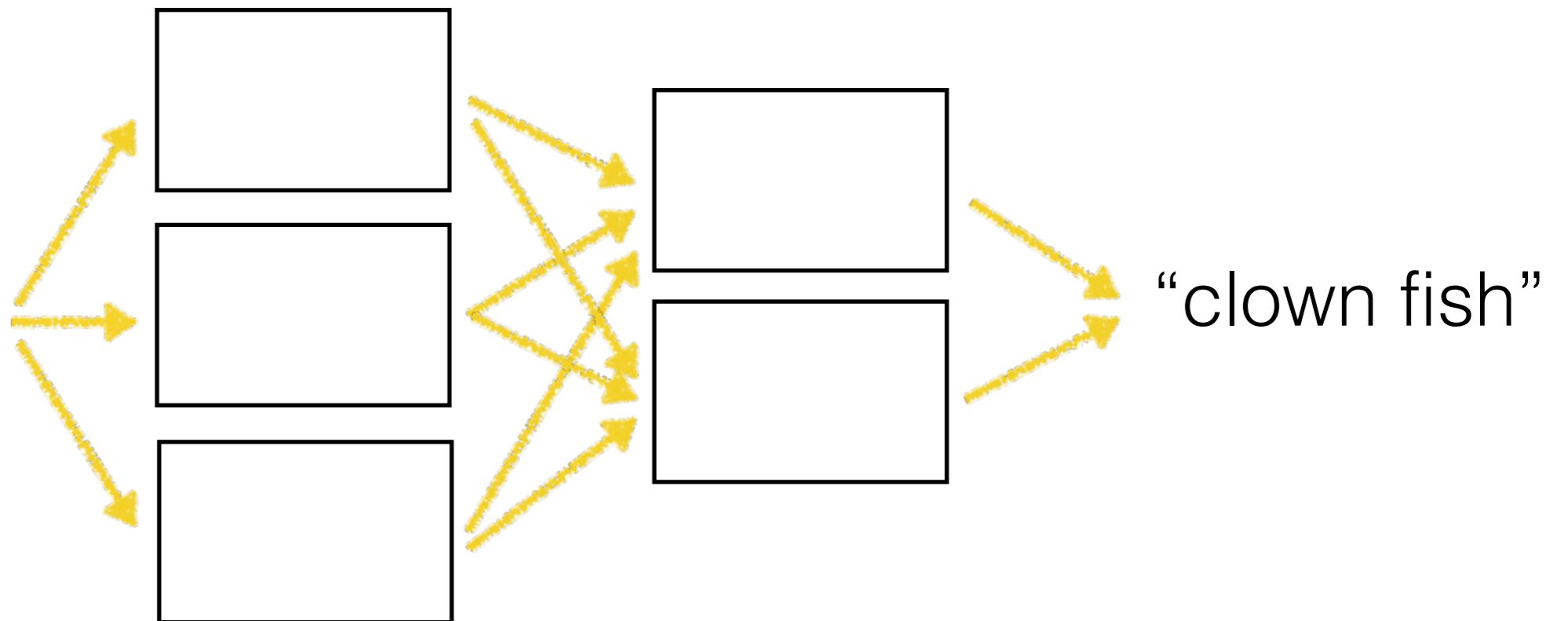


Object recognition



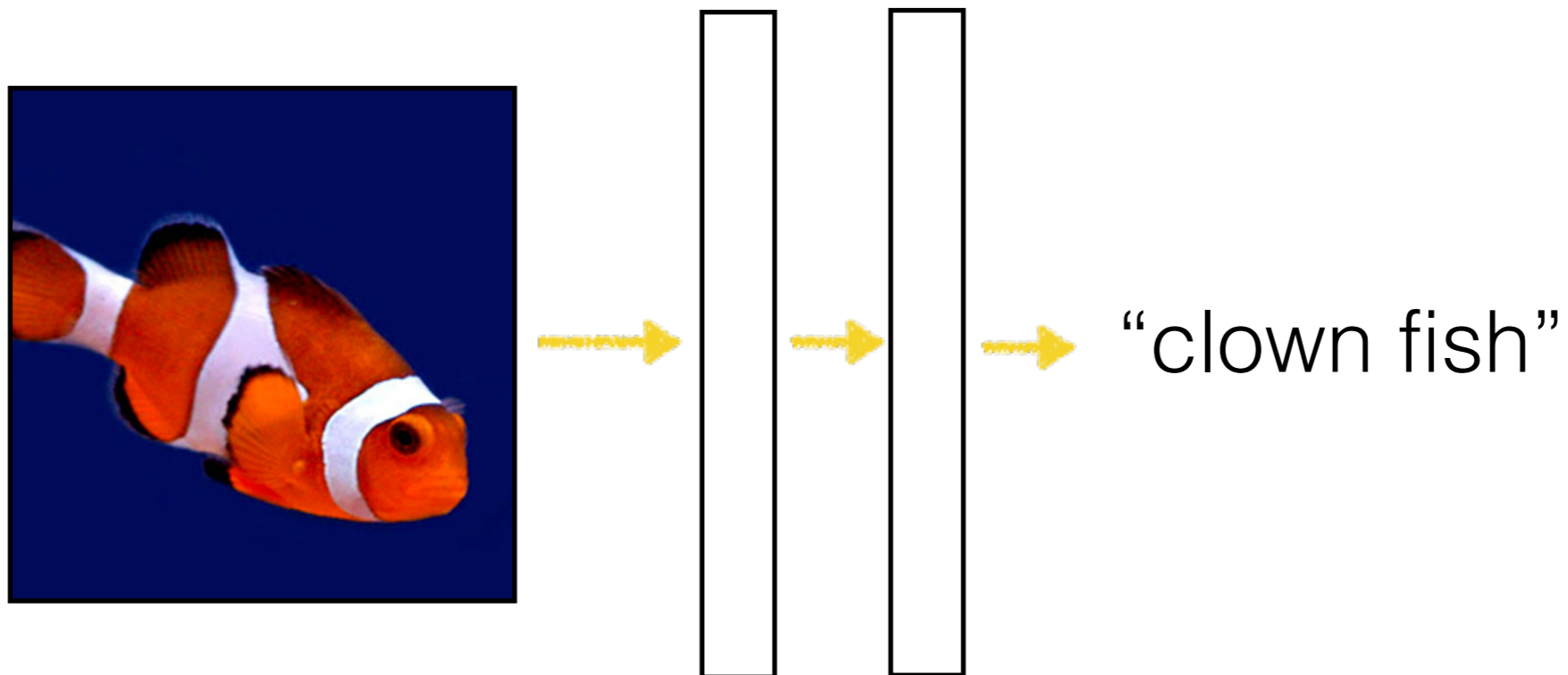
Neural network

Learned



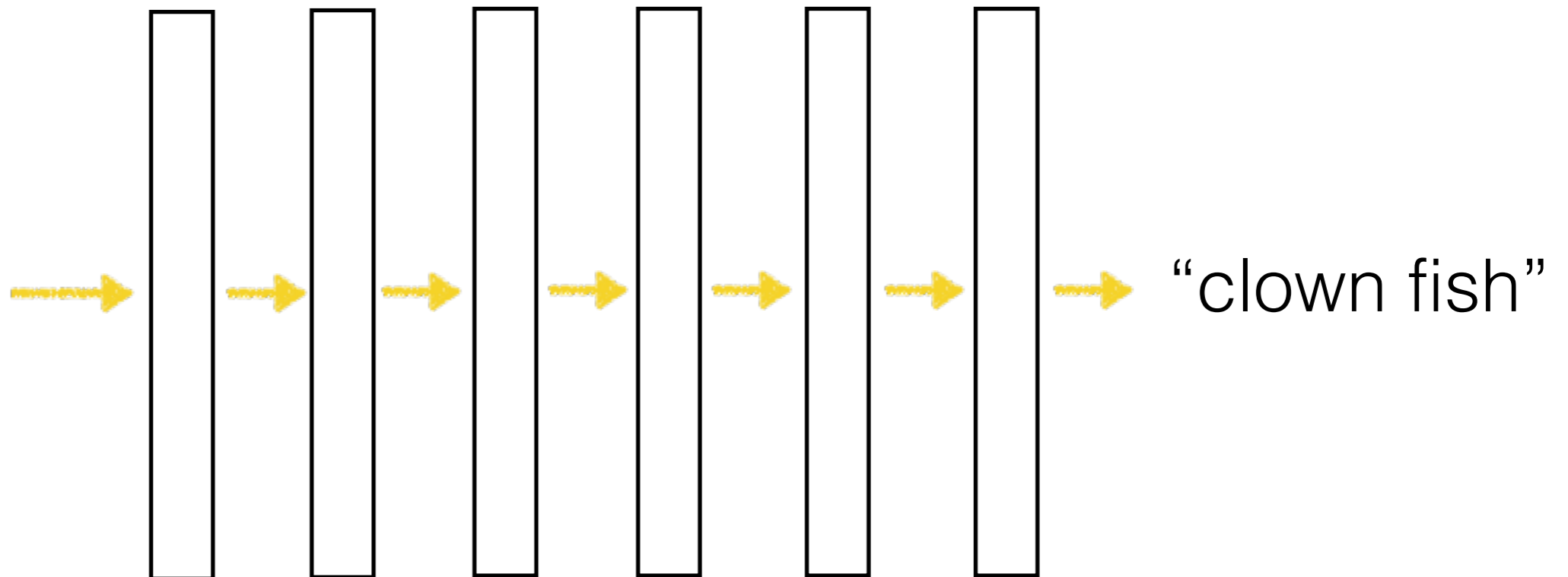
Neural network

Learned

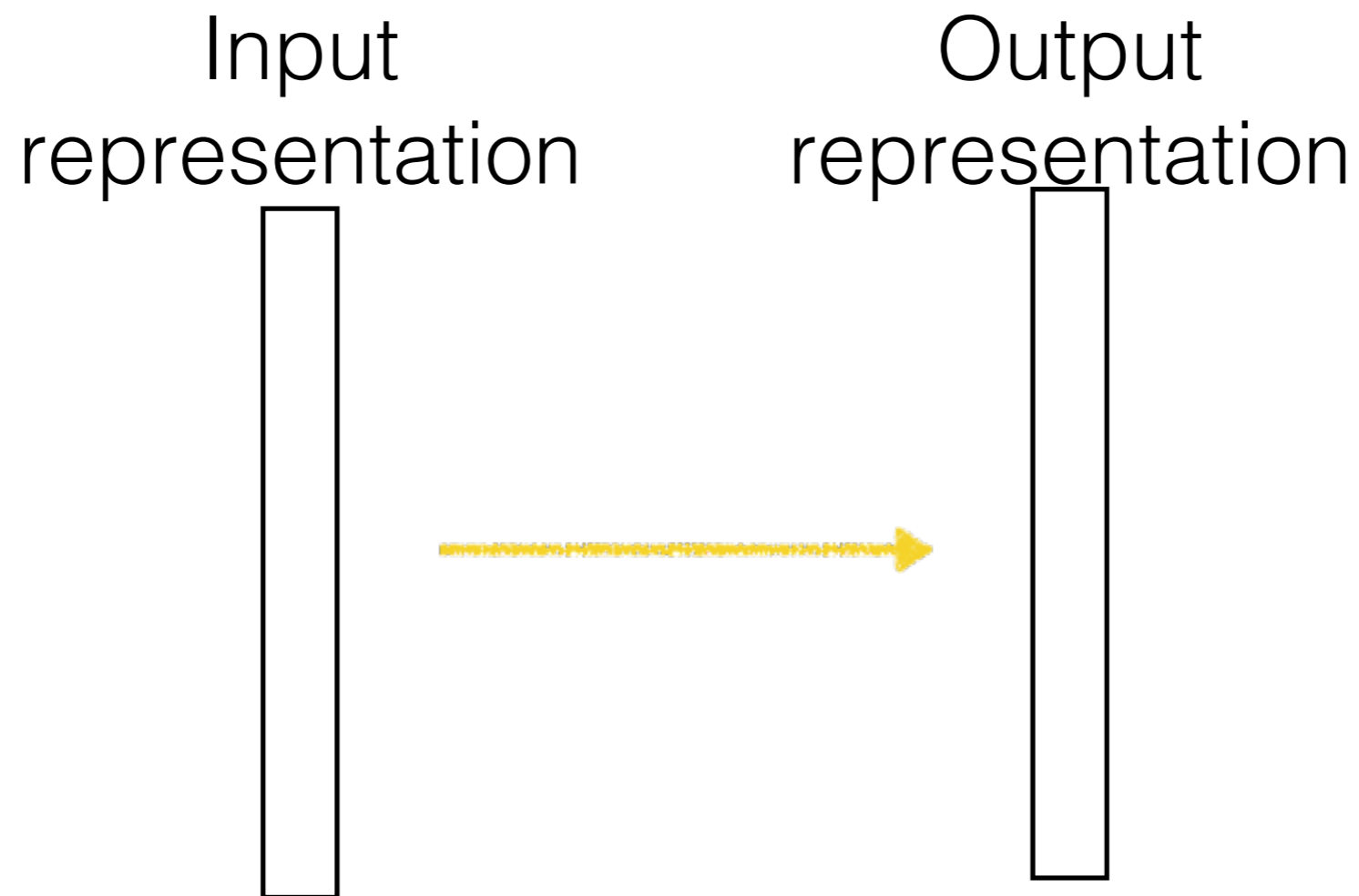


Deep neural network

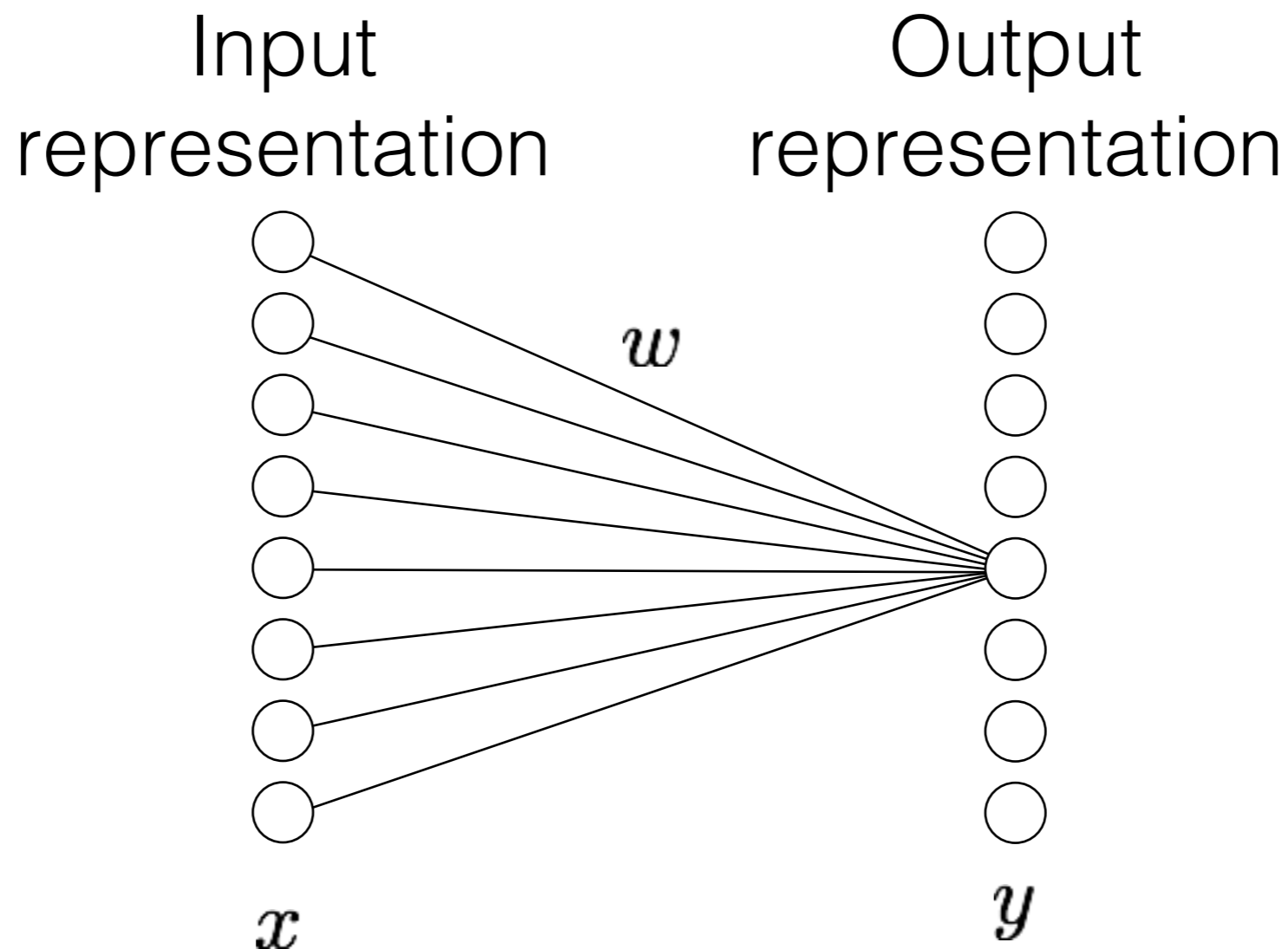
Learned



Computation in a neural net



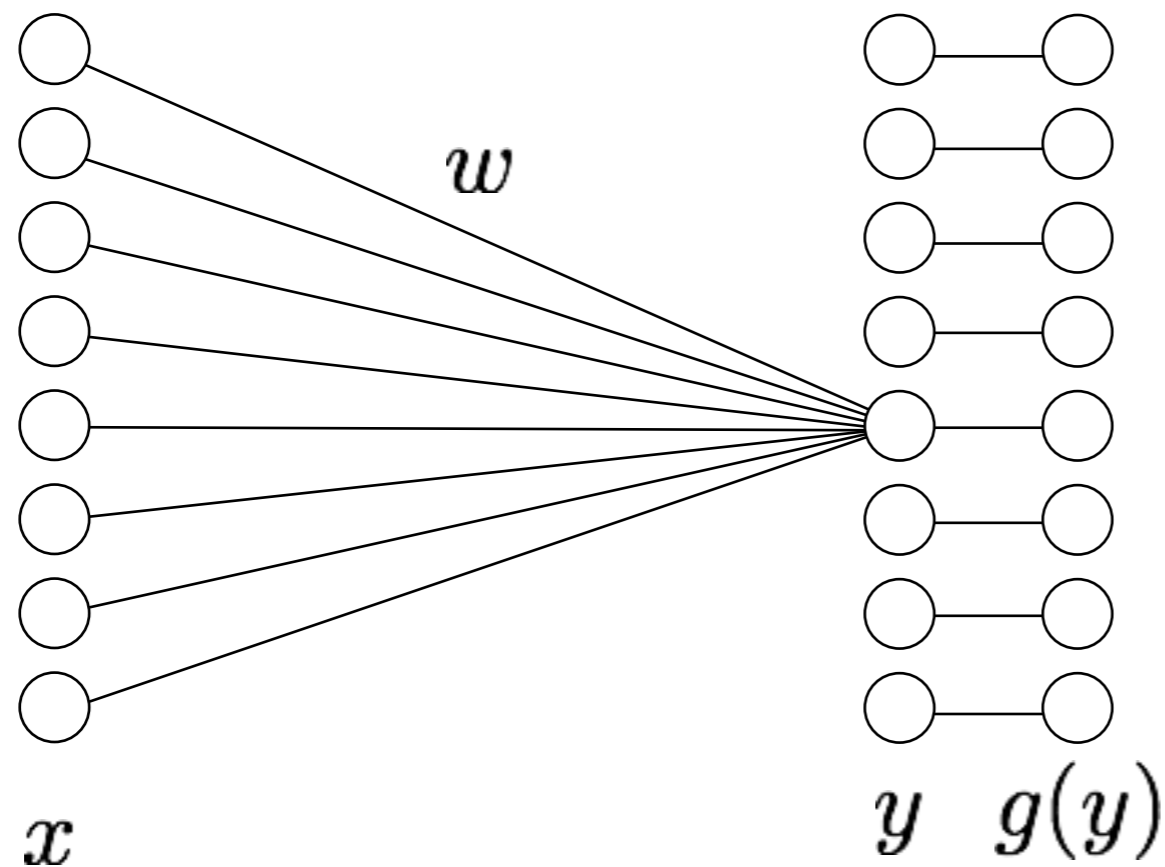
Computation in a neural net



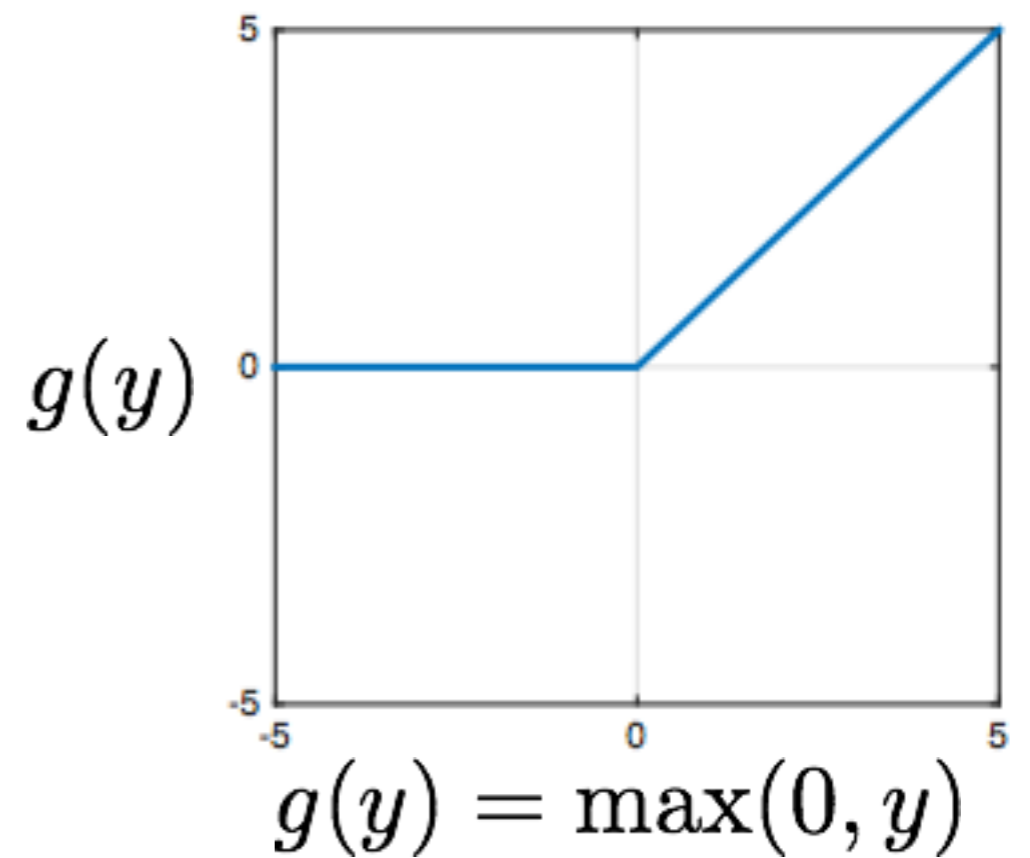
$$y_j = \sum_i w_{ij} x_i$$

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

Computation in a neural net

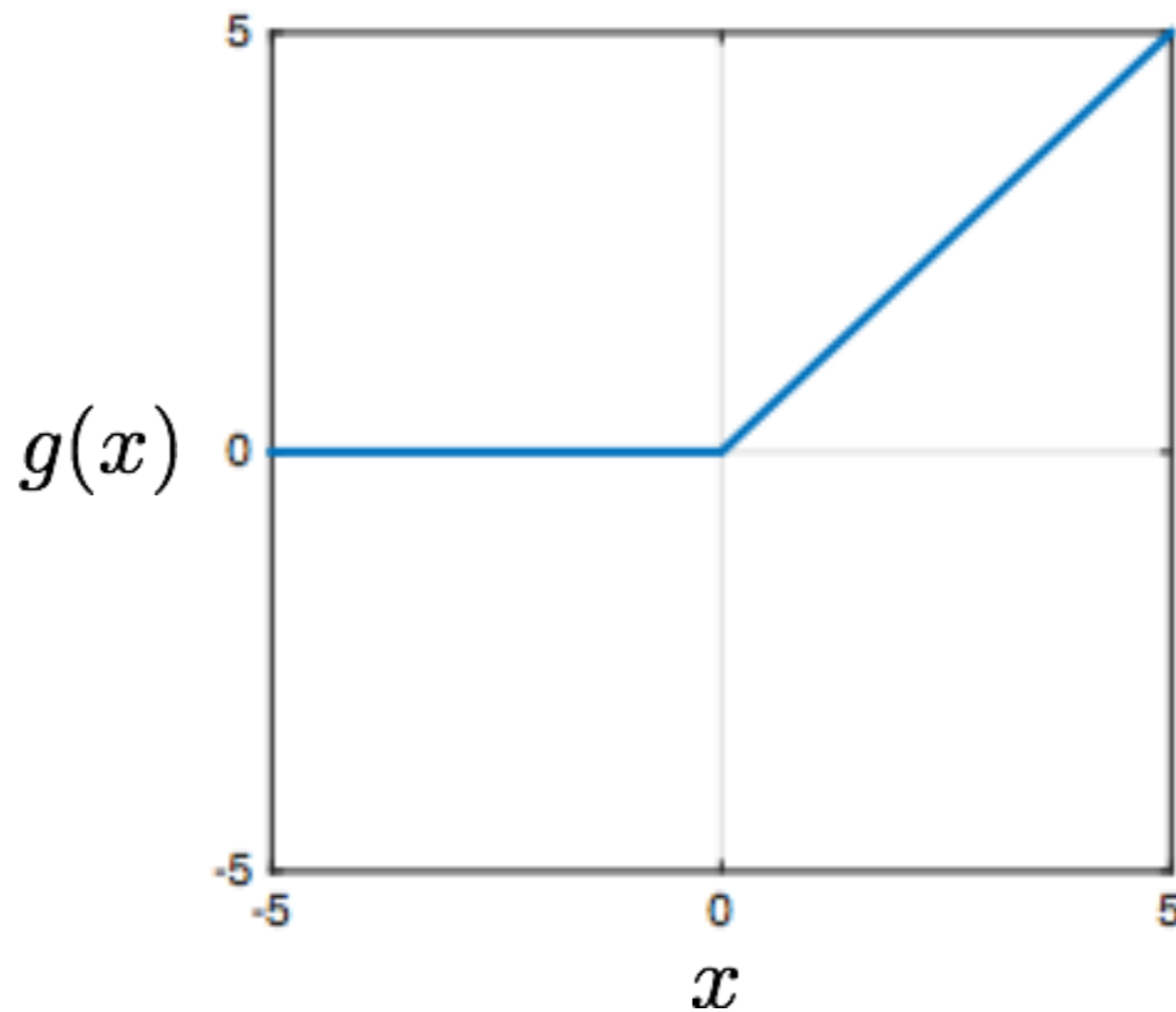


Rectified linear unit (ReLU)



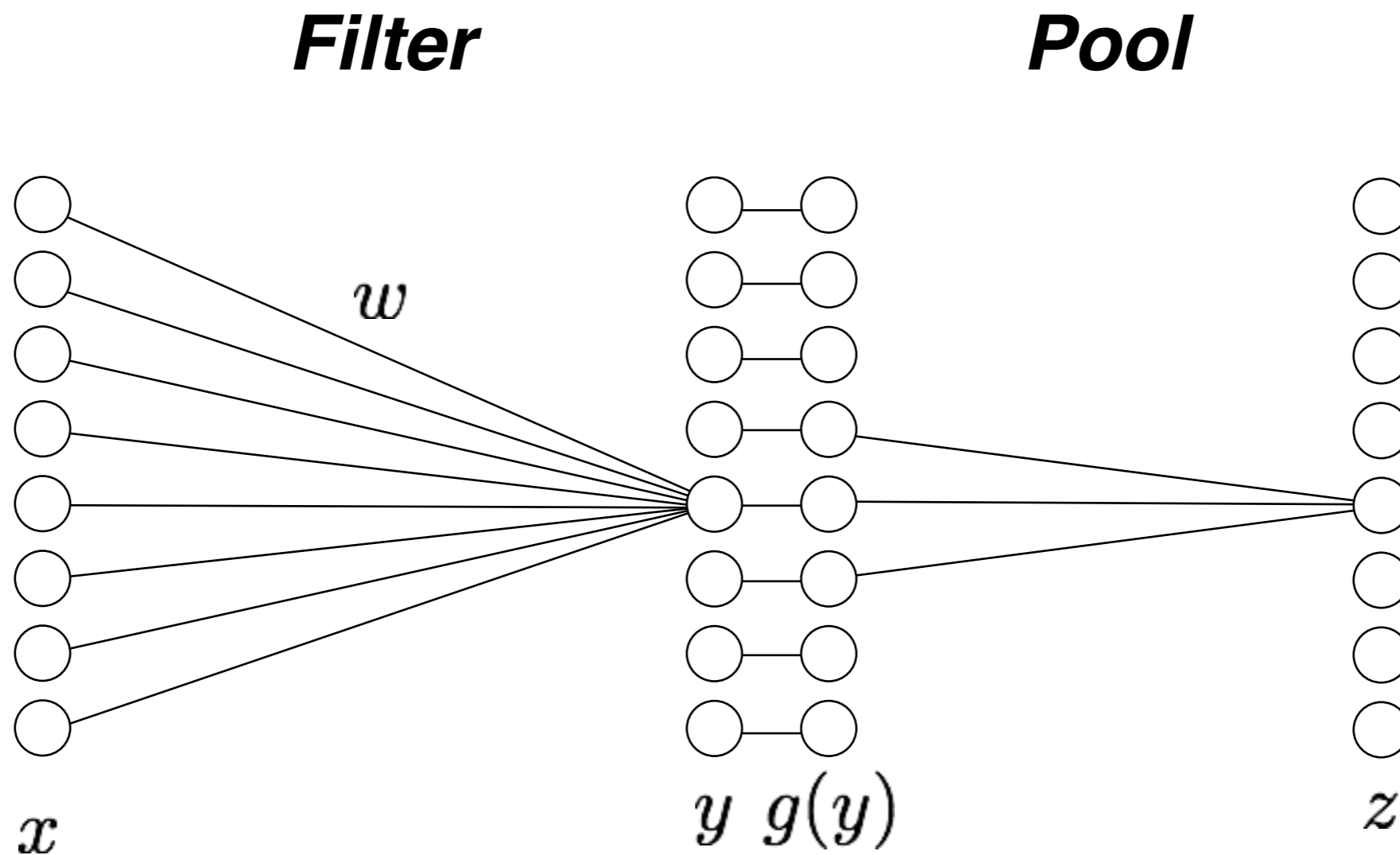
Computation in a neural net

Rectified linear unit (ReLU)



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

Computation in a neural net

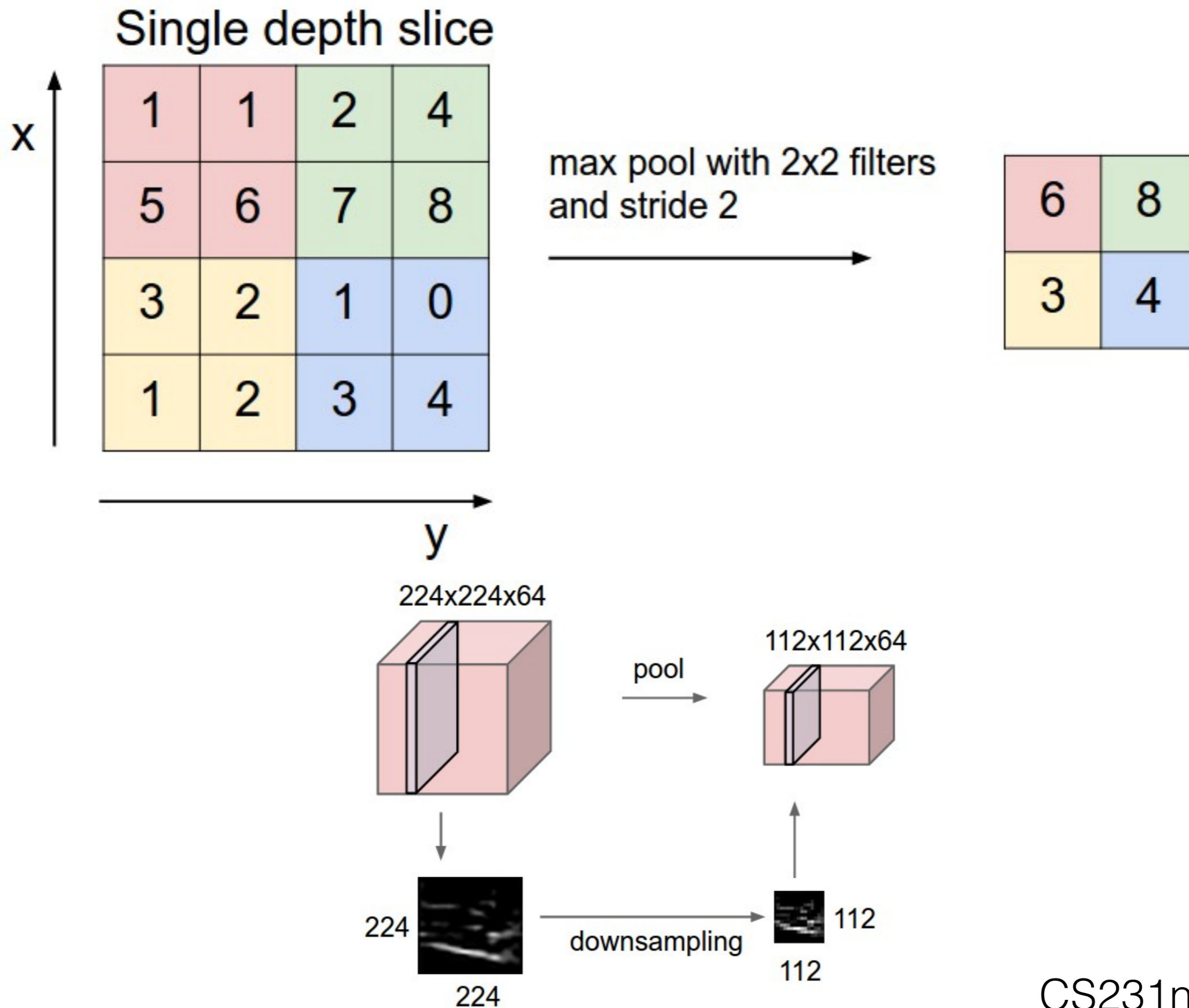


$$y_j = \sum_i w_{ij} x_i$$

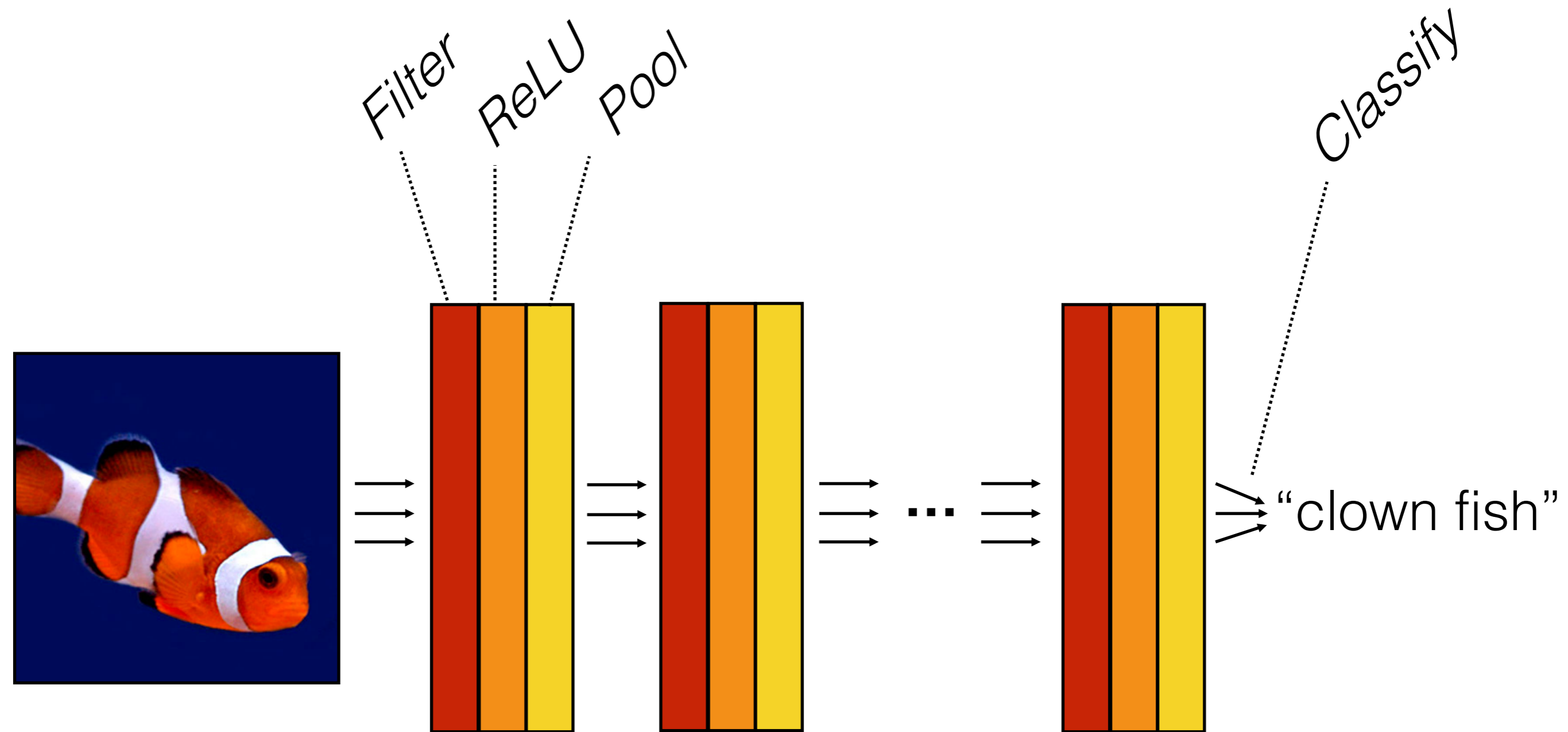
$$z_k = \max_{j \in \mathcal{N}(k)} g(y_j)$$

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

Computation in a neural net

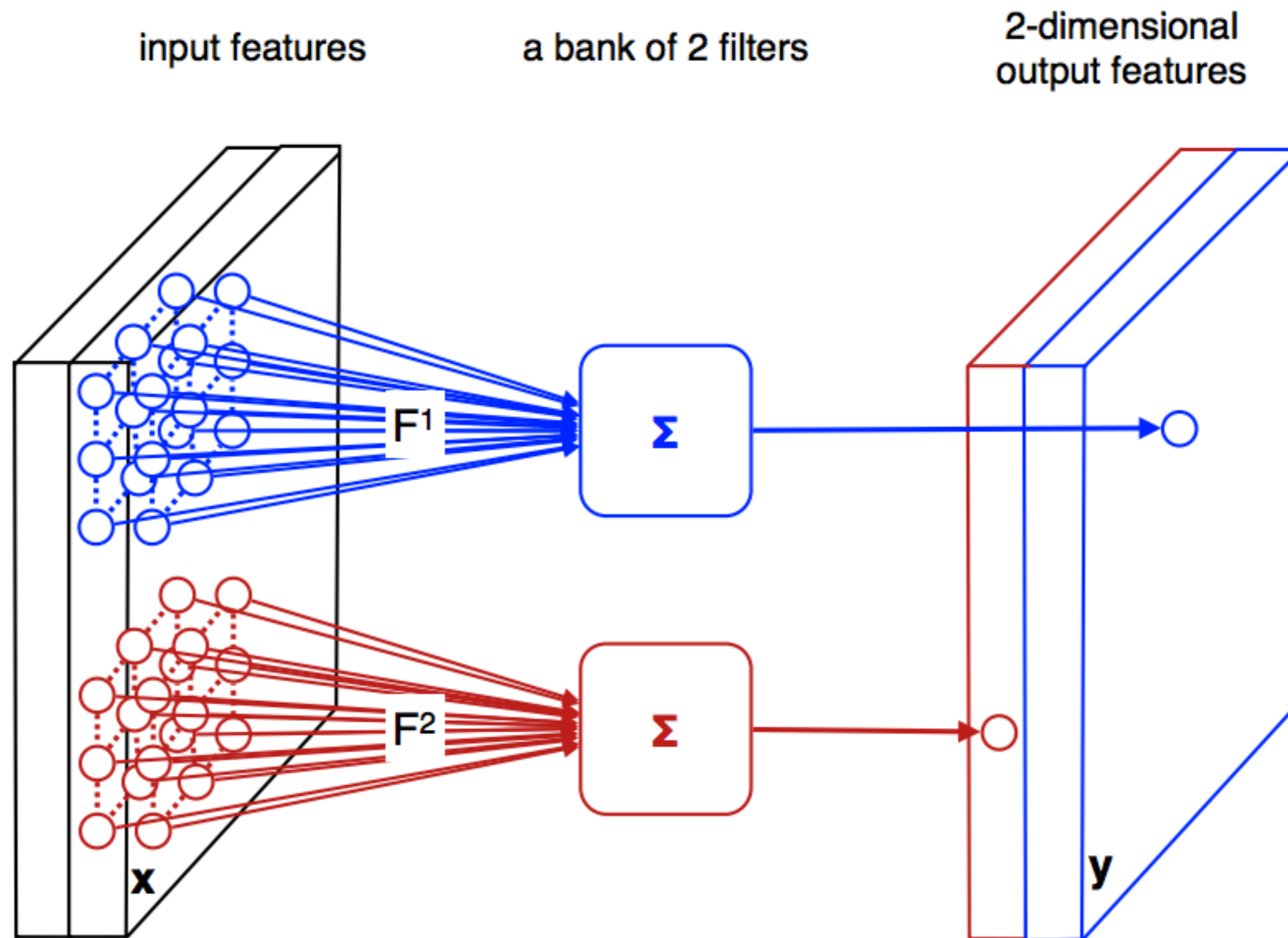


Computation in a neural net



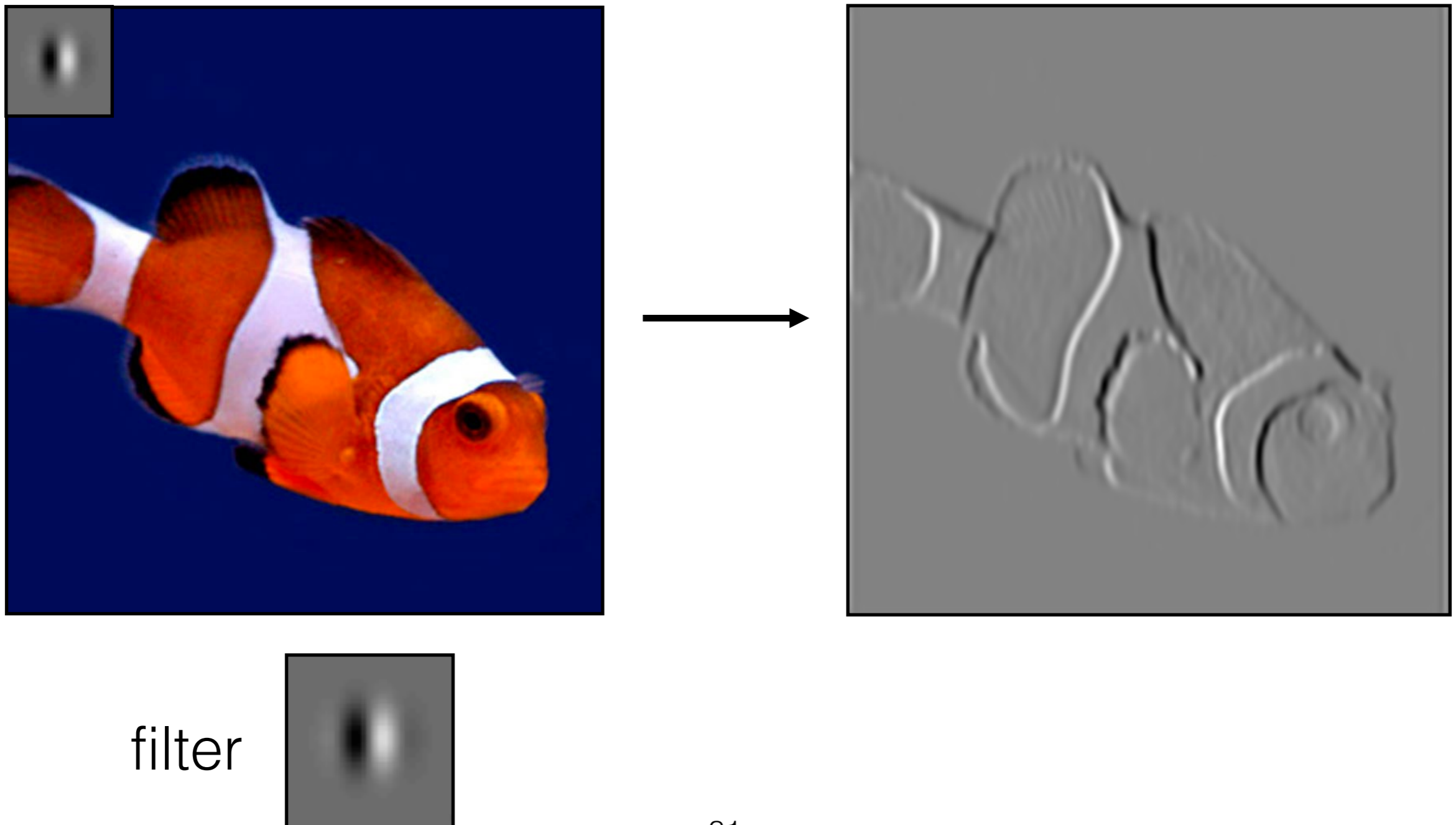
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Convolutional Neural Nets

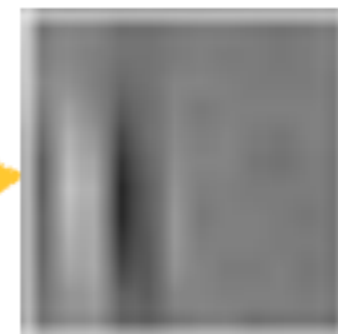
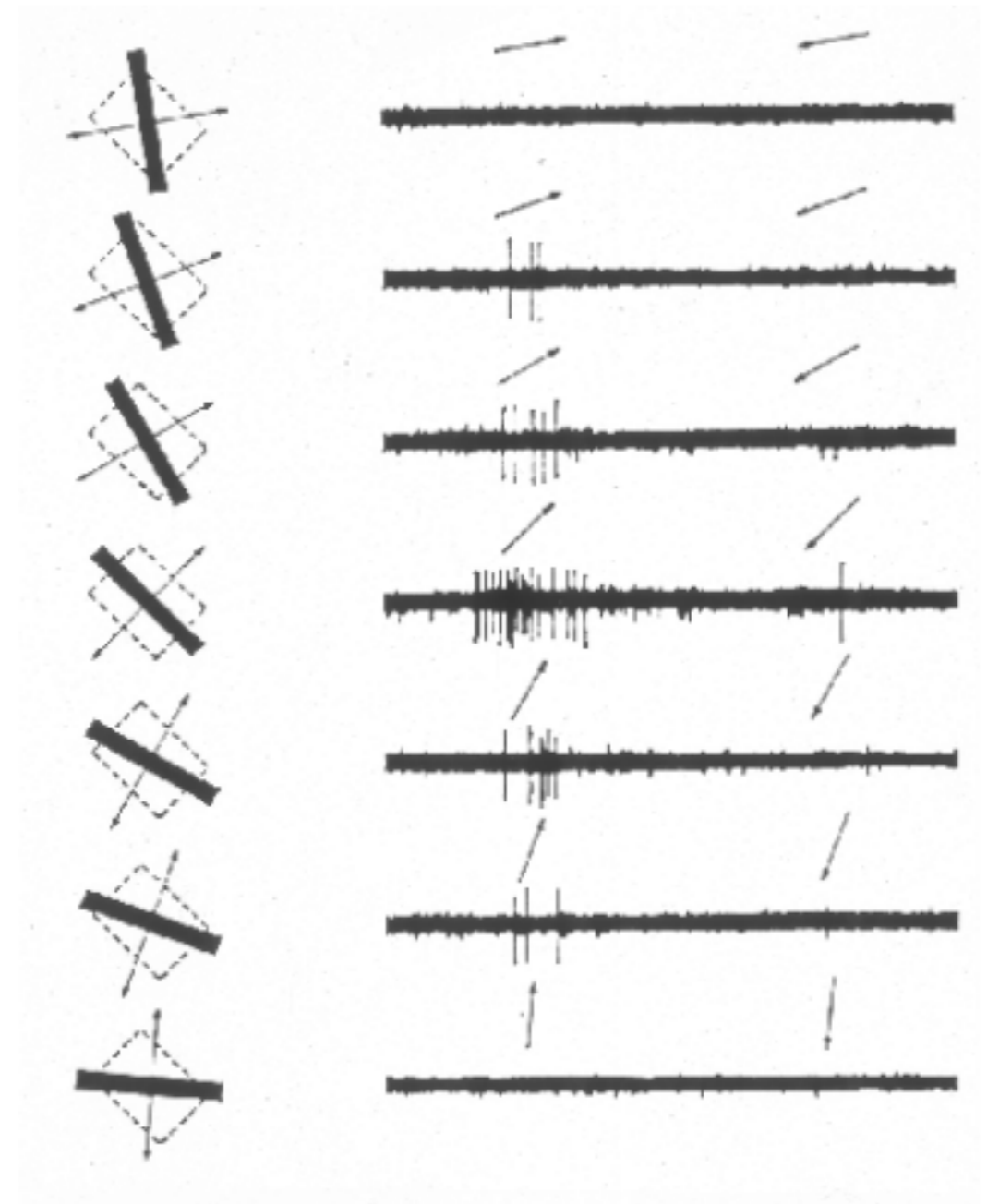
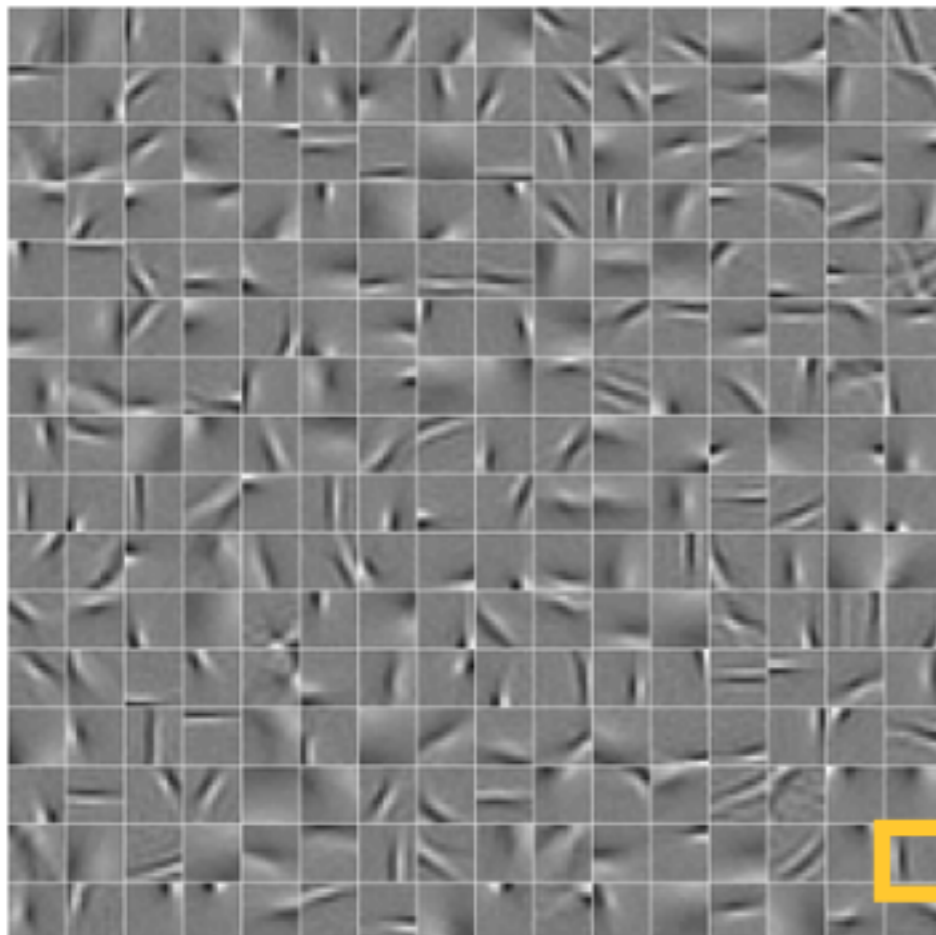
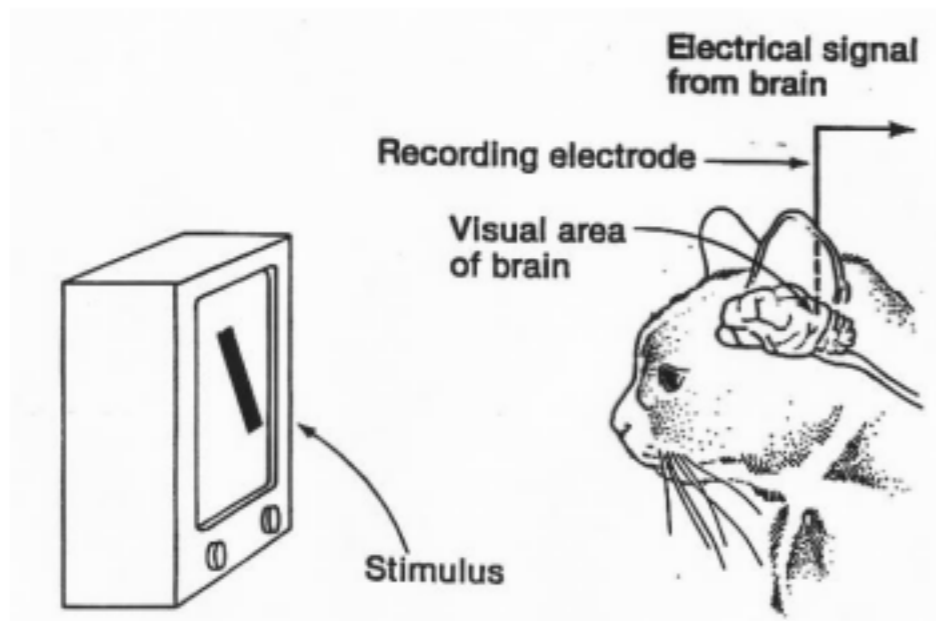


Convolutional Neural Nets

Convolution



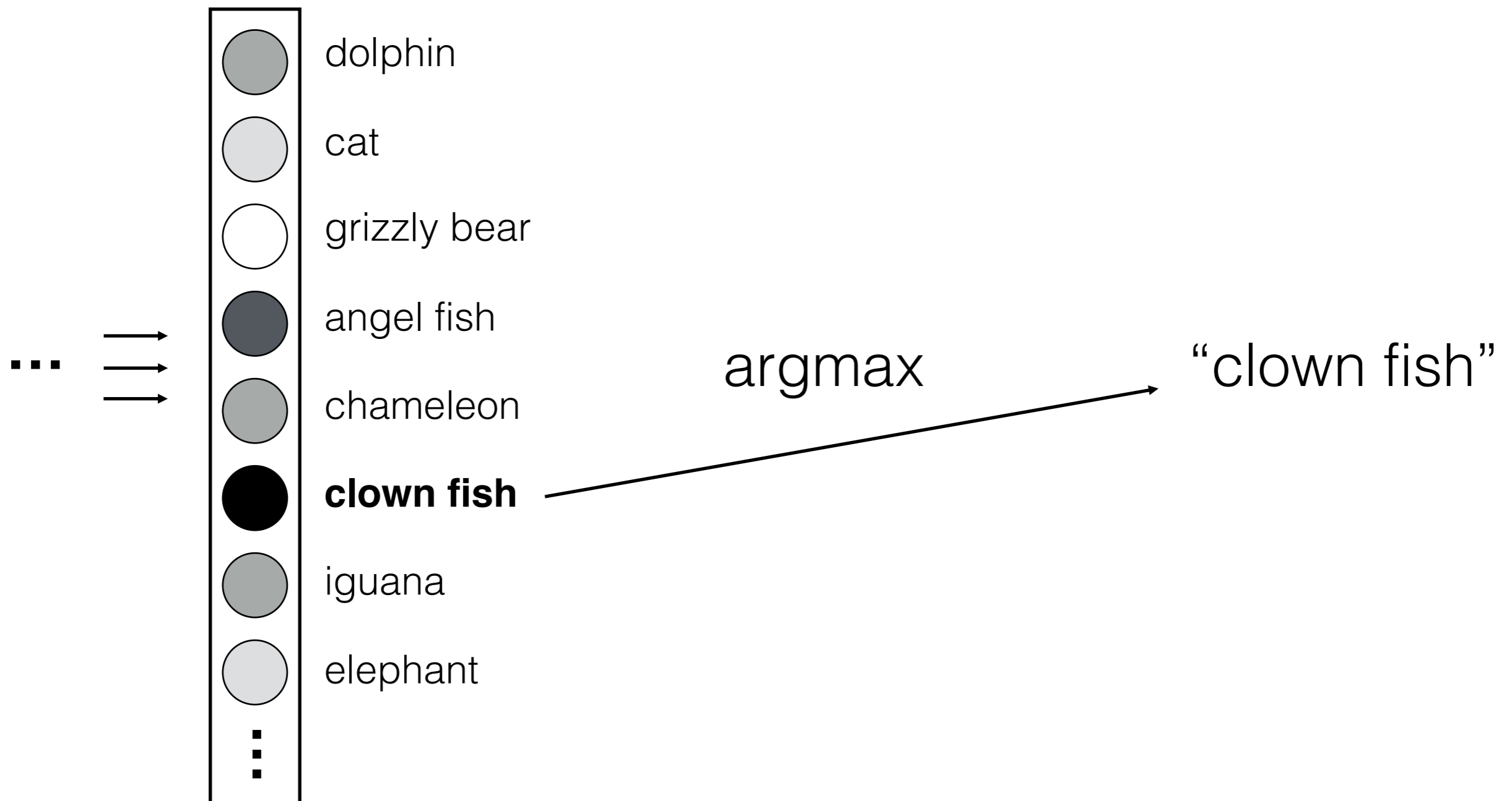
[Hubel and Wiesel 59]



oriented filter

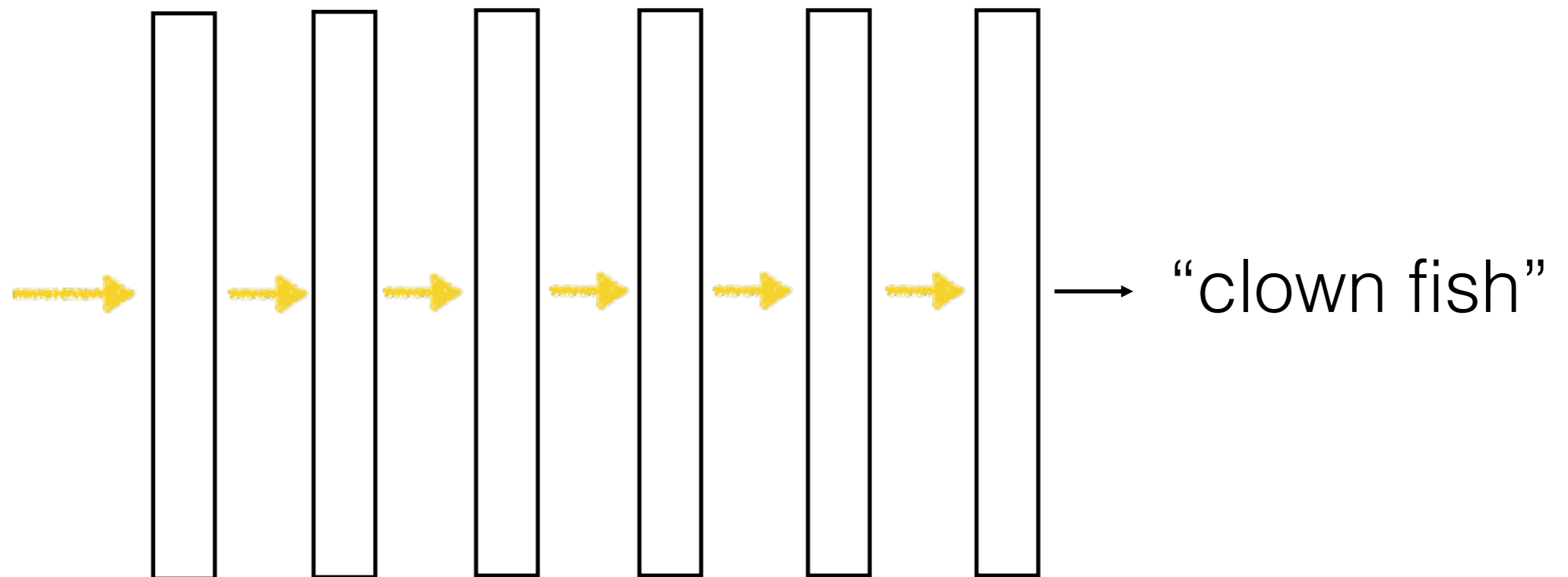
Computation in a neural net

Last layer



Learning with deep nets

Learned



Learning with deep nets



→ “clown fish”



→ “grizzly bear”



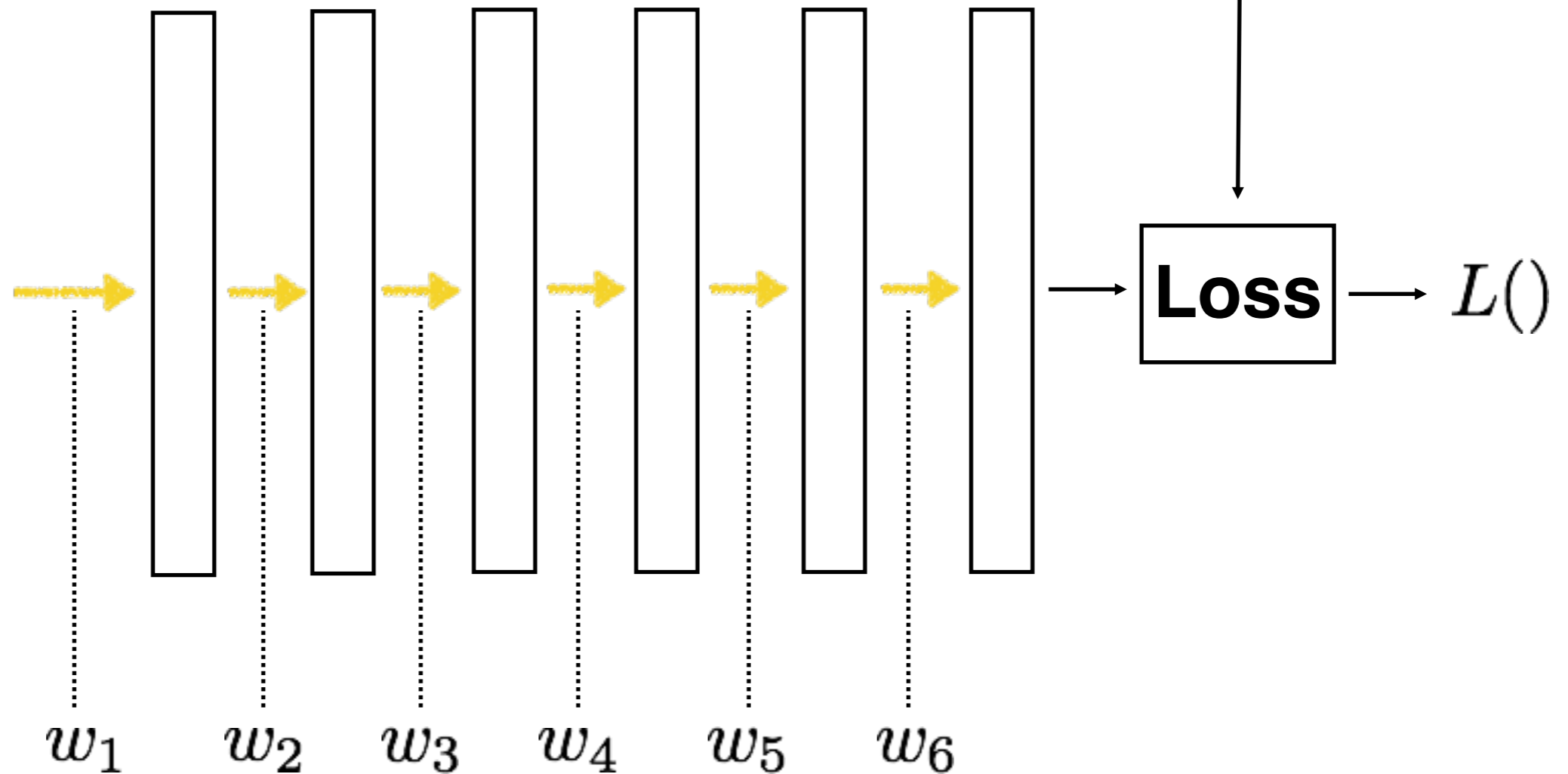
→ “chameleon”

Train network to
associate the right
label with each image

Learning with deep nets

“clown fish”

Learned

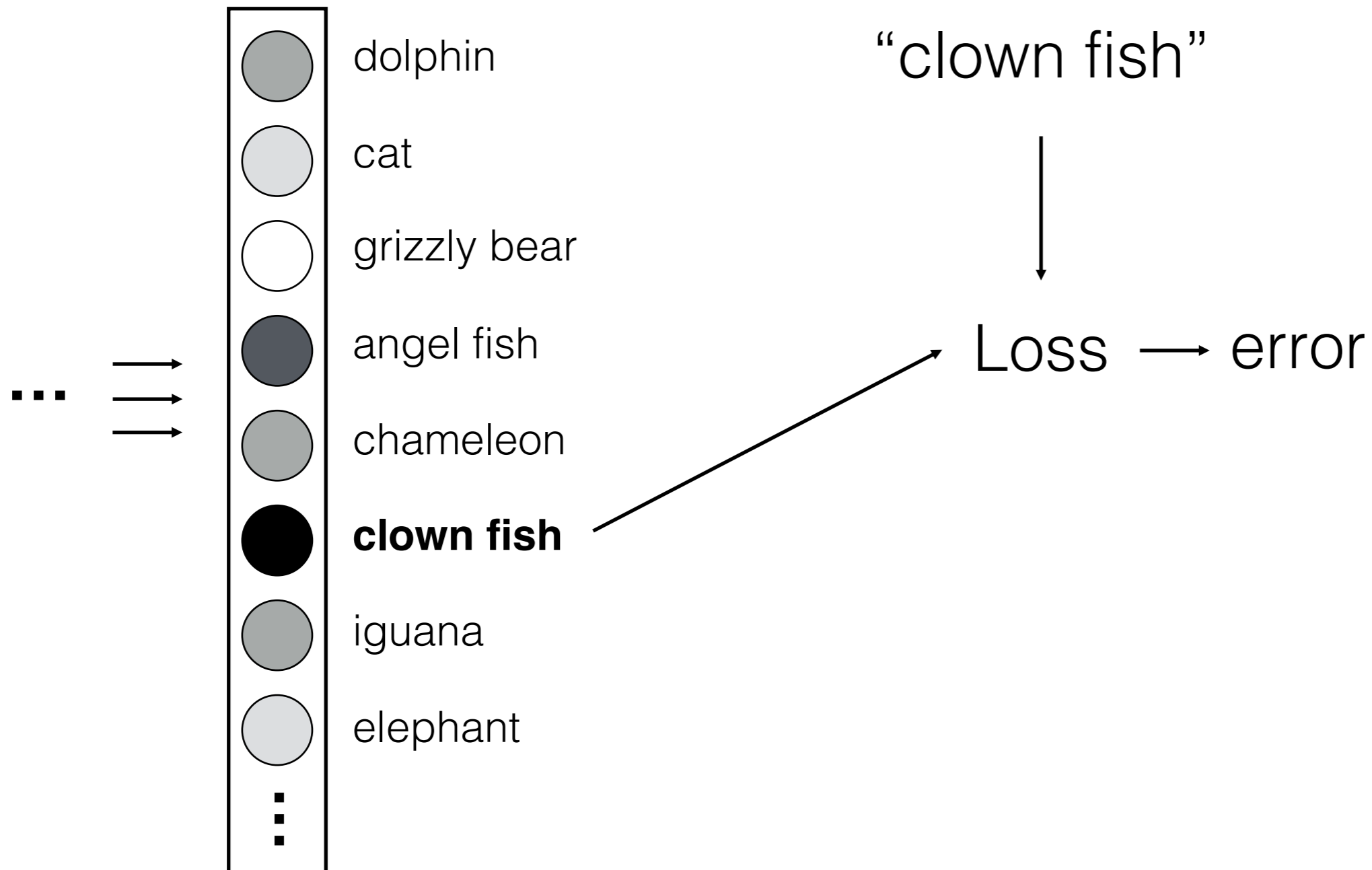


$$\operatorname{argmin}_{\mathbf{w}} L(w_1, \dots, w_6)$$

Loss function

Network output

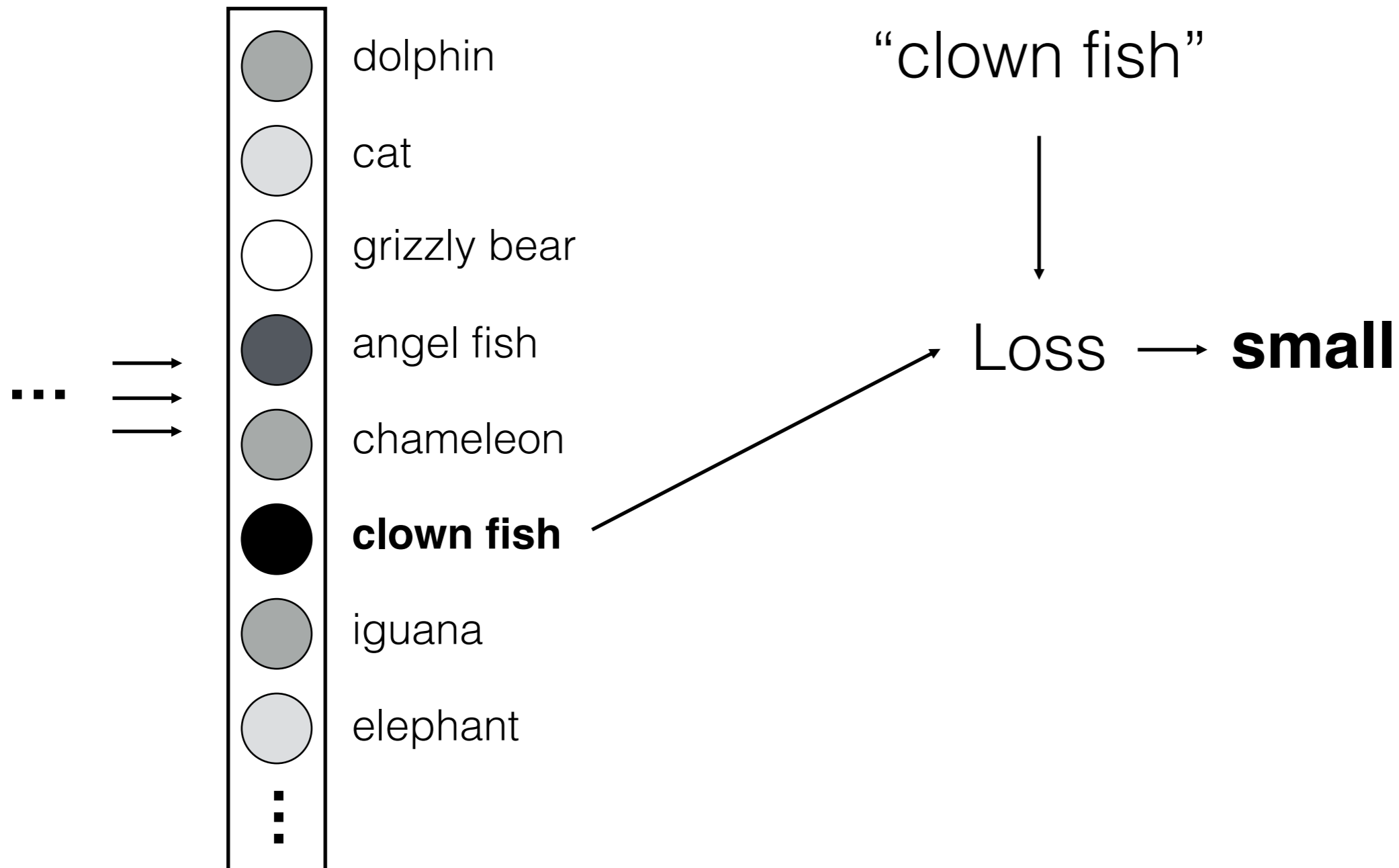
Ground truth label



Loss function

Network output

Ground truth label



Loss function

Network output

Ground truth label



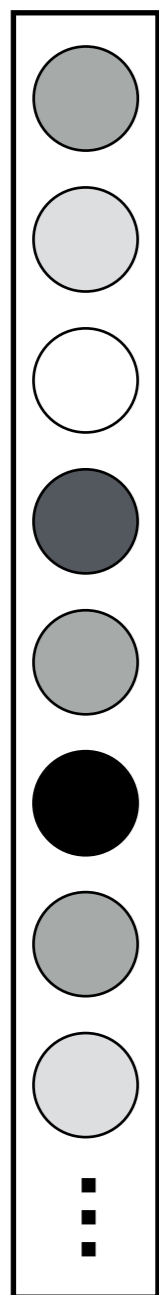
Loss function for classification

Network output

Ground truth label

$\hat{\mathbf{z}}$

\mathbf{z}



dolphin

cat

grizzly bear

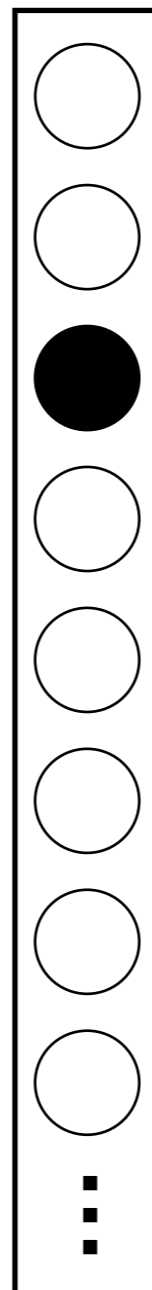
angel fish

chameleon

clown fish

iguana

elephant



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

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

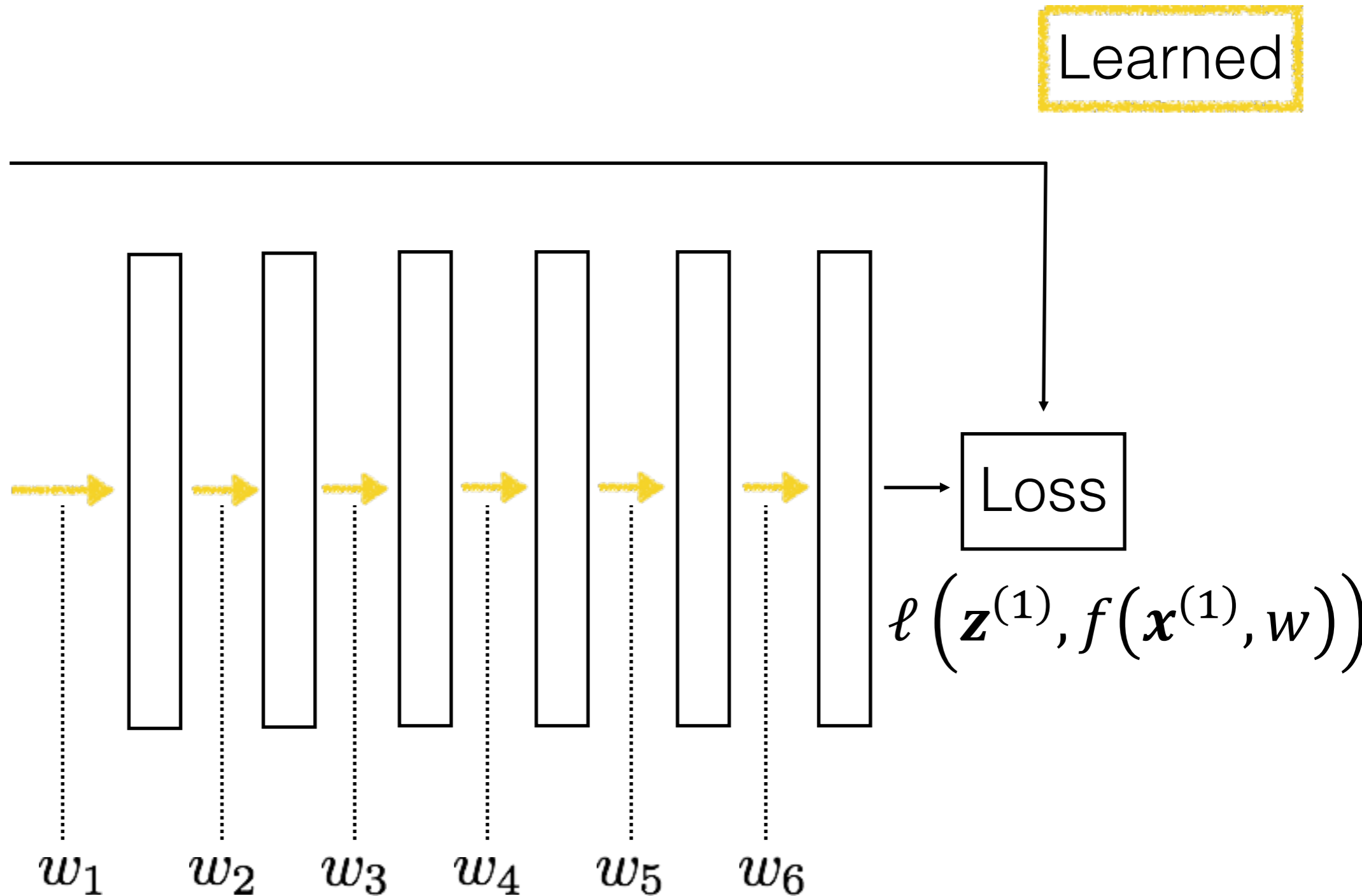
Cross-entropy loss

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

Learning with deep nets

$\mathbf{z}^{(1)}$
“clown fish”

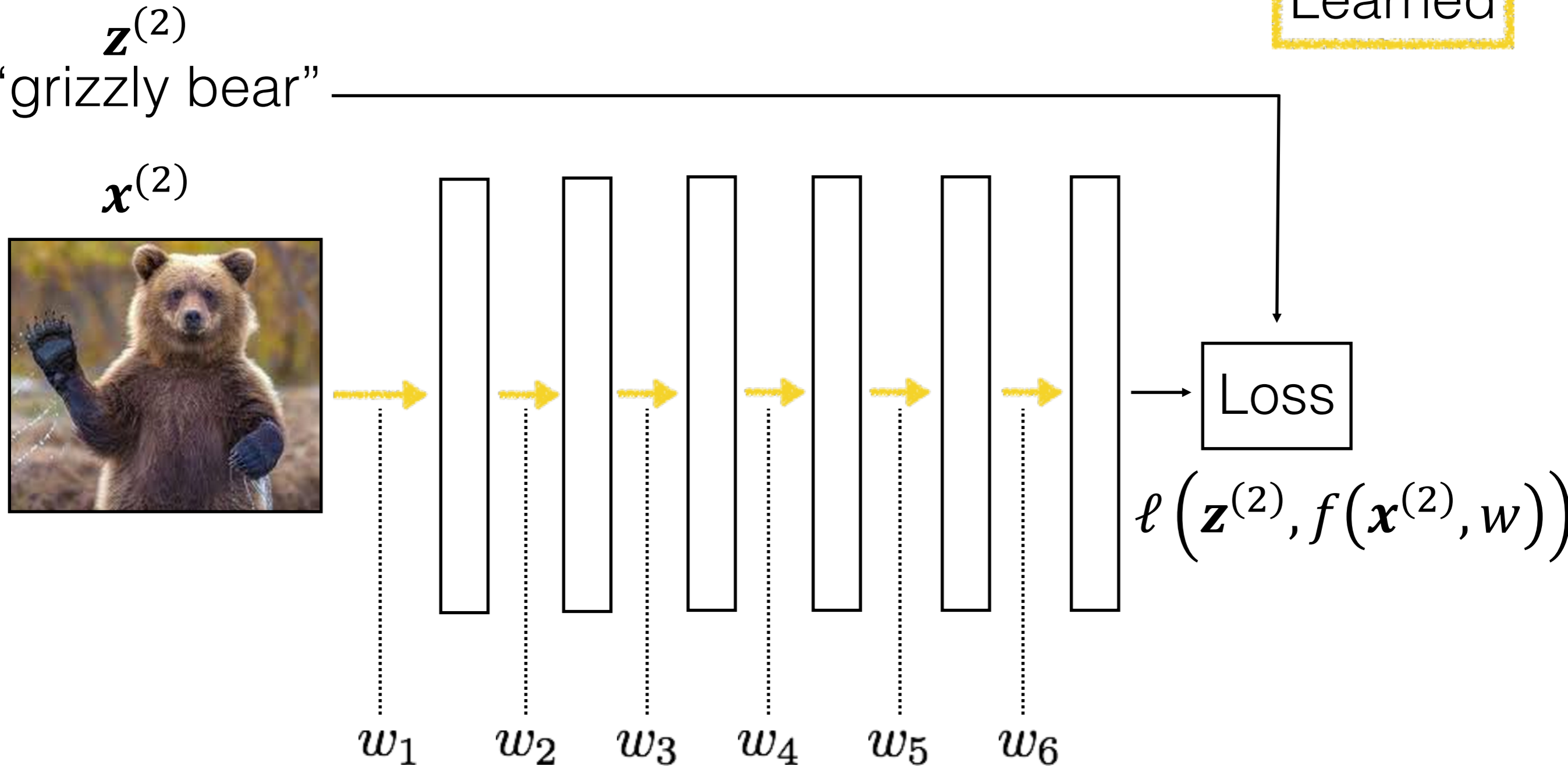
$\mathbf{x}^{(1)}$



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

Learning with deep nets

Learned



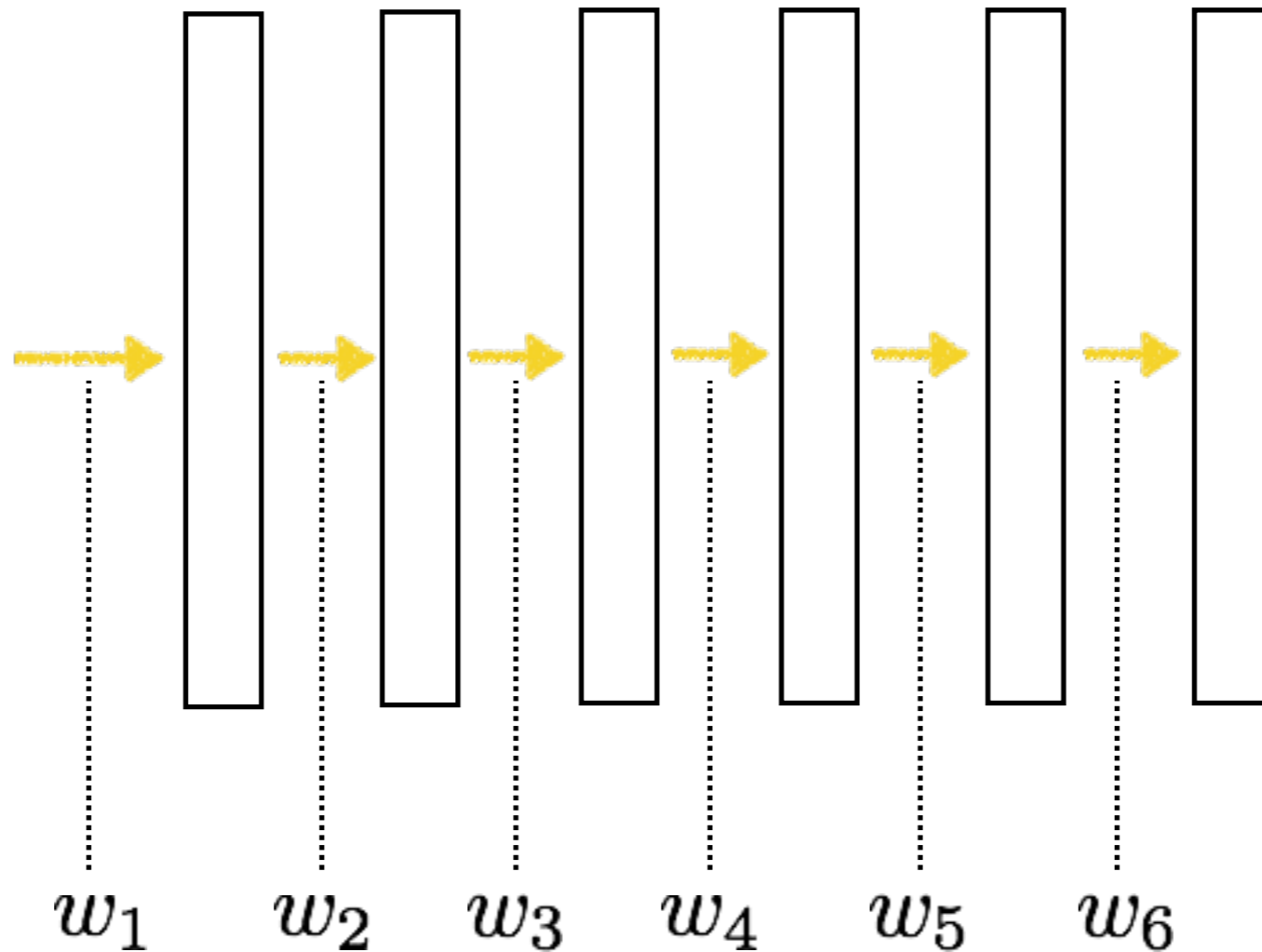
$\mathbf{x}^{(2)}, \mathbf{z}^{(2)}$ is the input and label
of the 2nd training image

Learning with deep nets

Learned

$\mathbf{z}^{(3)}$
“chameleon”

$\mathbf{x}^{(3)}$



Loss

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

$$\operatorname{argmin}_w \sum_i \ell(\mathbf{z}^{(i)}, f(\mathbf{x}^{(i)}, w))$$

Gradient descent

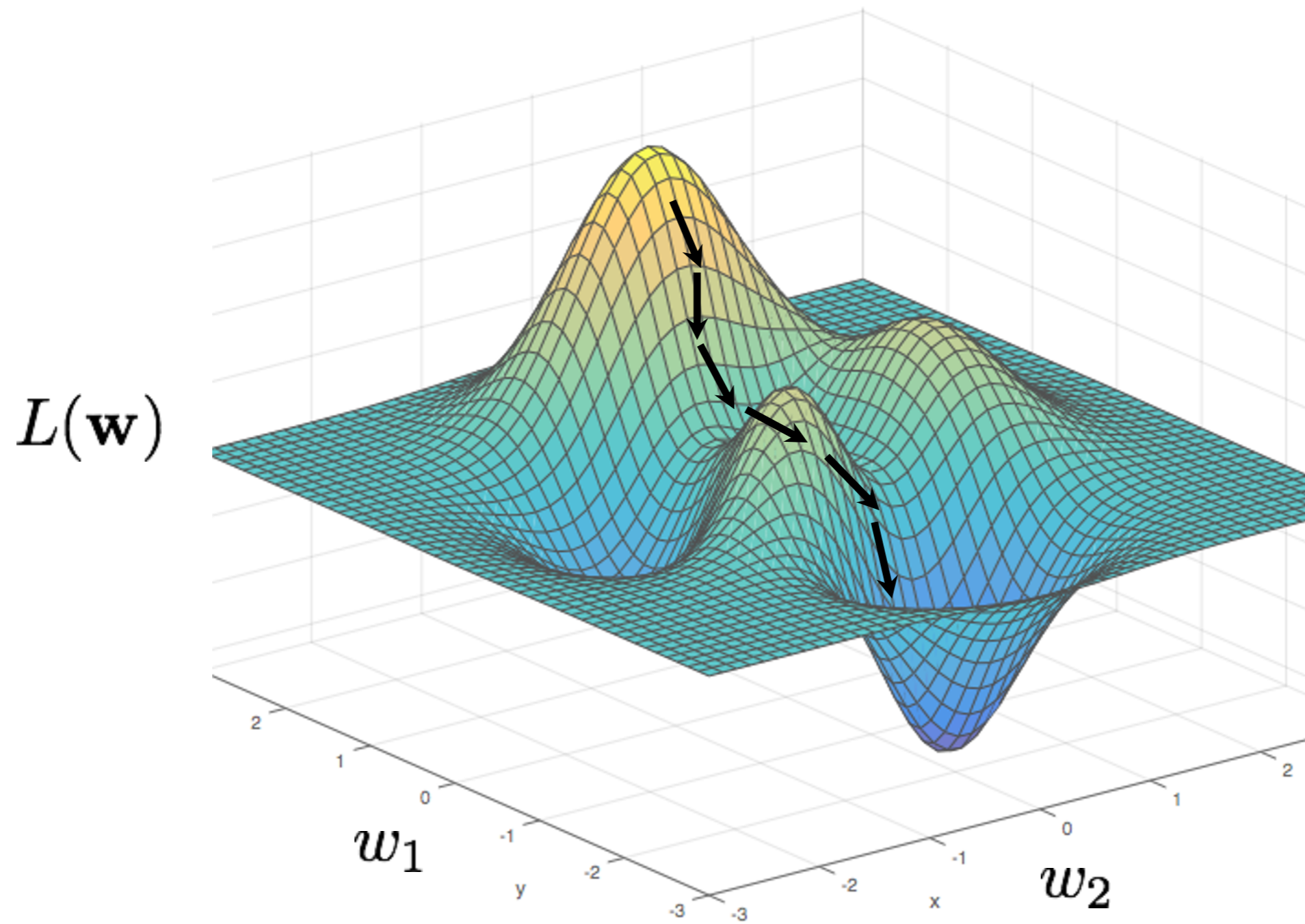
$$\operatorname{argmin}_w \sum_i \ell(z^{(i)}, f(x^{(i)}, w)) = \operatorname{argmin}_w L(w)$$

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

	<p>mite</p> <p>black widow</p> <p>cockroach</p> <p>tick</p> <p>starfish</p>		<p>container ship</p> <p>lifeboat</p> <p>amphibian</p> <p>fireboat</p> <p>drilling platform</p>		<p>motor scooter</p> <p>go-kart</p> <p>moped</p> <p>bumper car</p> <p>golfcart</p>		<p>leopard</p> <p>jaguar</p> <p>cheetah</p> <p>snow leopard</p> <p>Egyptian cat</p>
--	------------------------------------------------------------------------------------	--	--------------------------------------------------------------------------------------------------------	--	-------------------------------------------------------------------------------------------	--	--------------------------------------------------------------------------------------------



grille

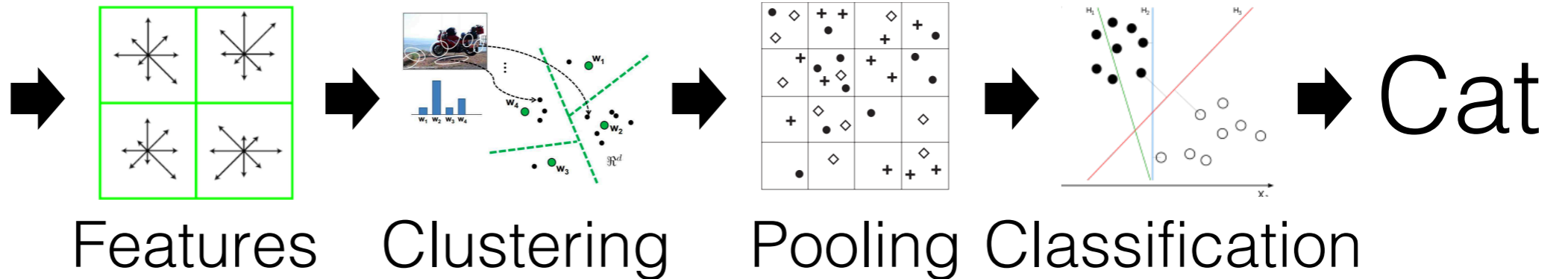
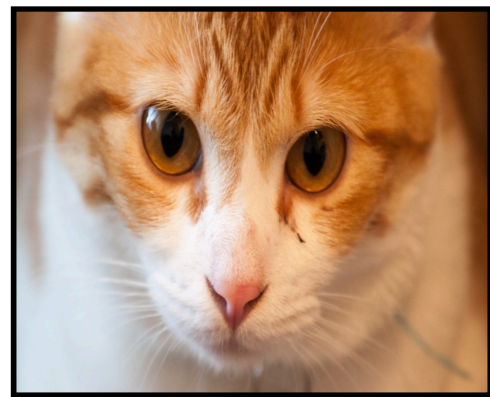
mushroom

cherry

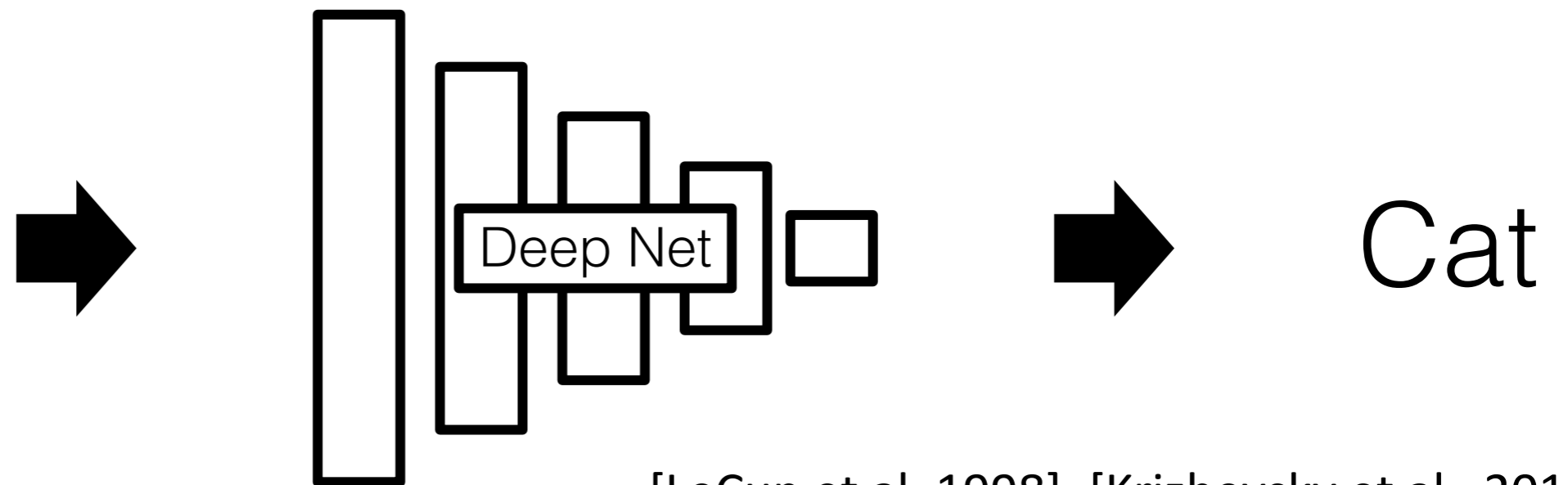
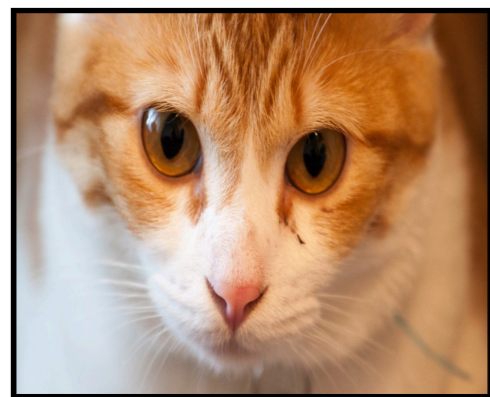
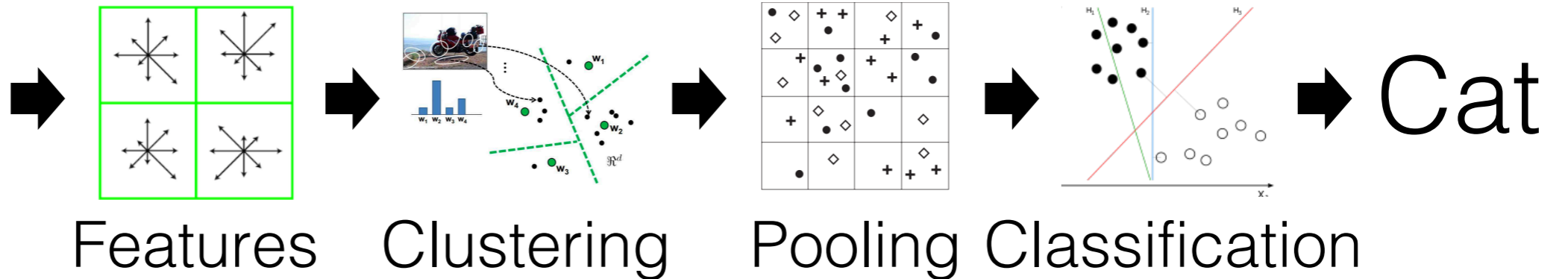
Madagascar cat

	<p>convertible</p> <p>grille</p> <p>pickup</p> <p>beach wagon</p> <p>fire engine</p>		<p>agaric</p> <p>mushroom</p> <p>jelly fungus</p> <p>gill fungus</p> <p>dead-man's-fingers</p>		<p>dalmatian</p> <p>grape</p> <p>elderberry</p> <p>ffordshire bullterrier</p> <p>currant</p>		<p>squirrel monkey</p> <p>spider monkey</p> <p>titi</p> <p>indri</p> <p>howler monkey</p>
--	---------------------------------------------------------------------------------------------	--	-------------------------------------------------------------------------------------------------------	--	-----------------------------------------------------------------------------------------------------	--	--------------------------------------------------------------------------------------------------

Computer Vision before 2012

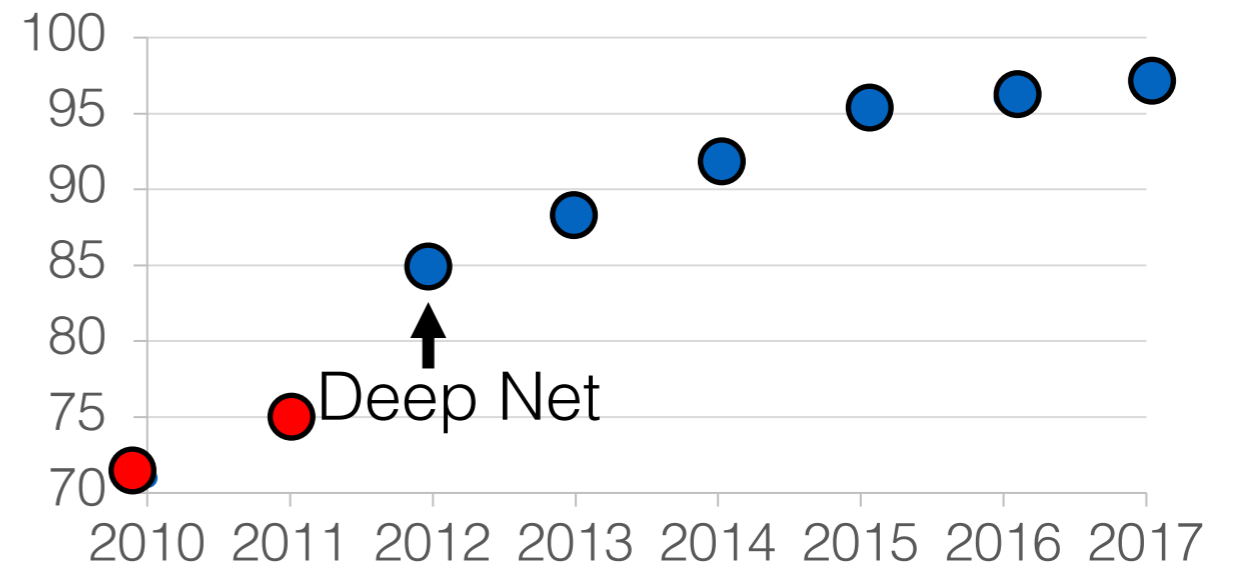


Computer Vision Now

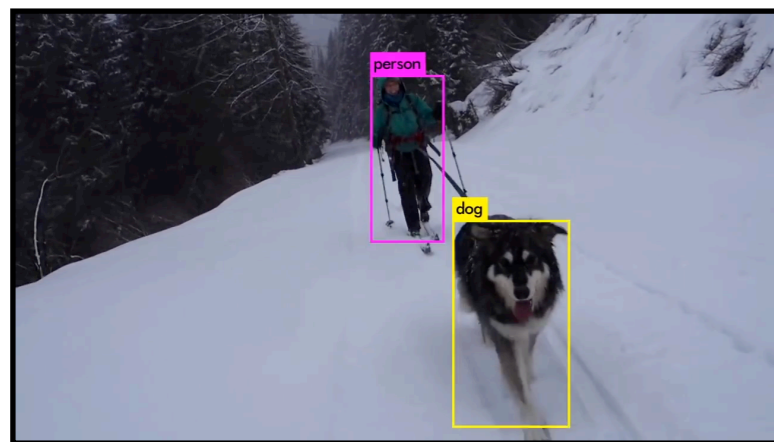


[LeCun et al, 1998], [Krizhevsky et al, 2012]

Deep Learning for Computer Vision

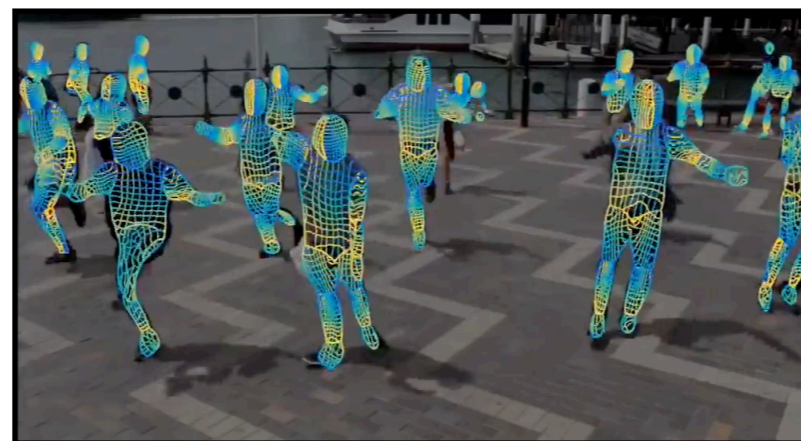


Top 5 **accuracy** on ImageNet benchmark



[Redmon et al., 2018]

Object detection



[Güler et al., 2018]

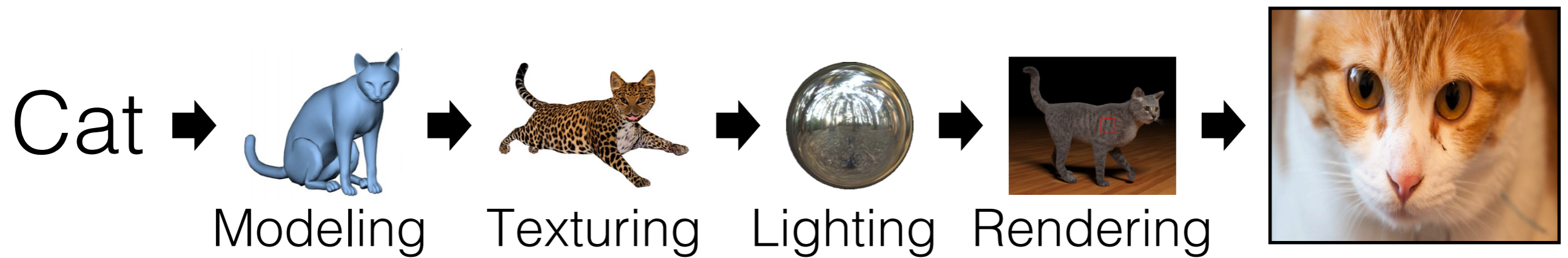
Human understanding



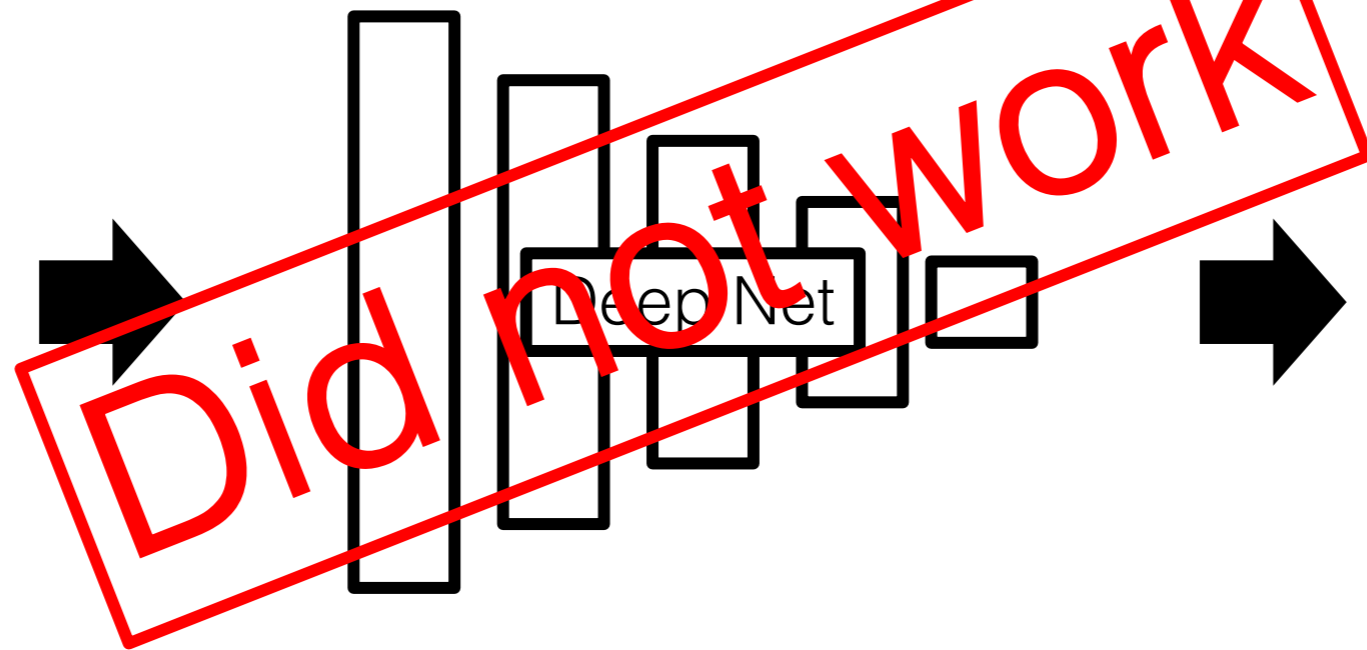
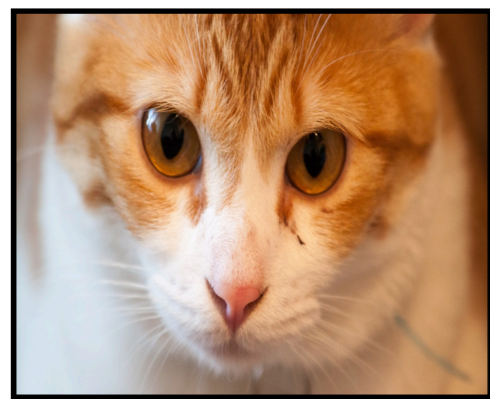
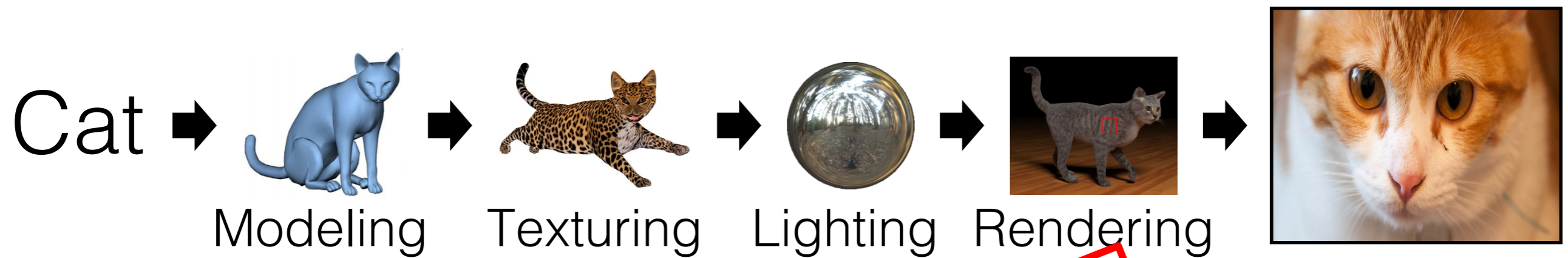
[Zhao et al., 2017]

Autonomous driving

Can Deep Learning Help Graphics?

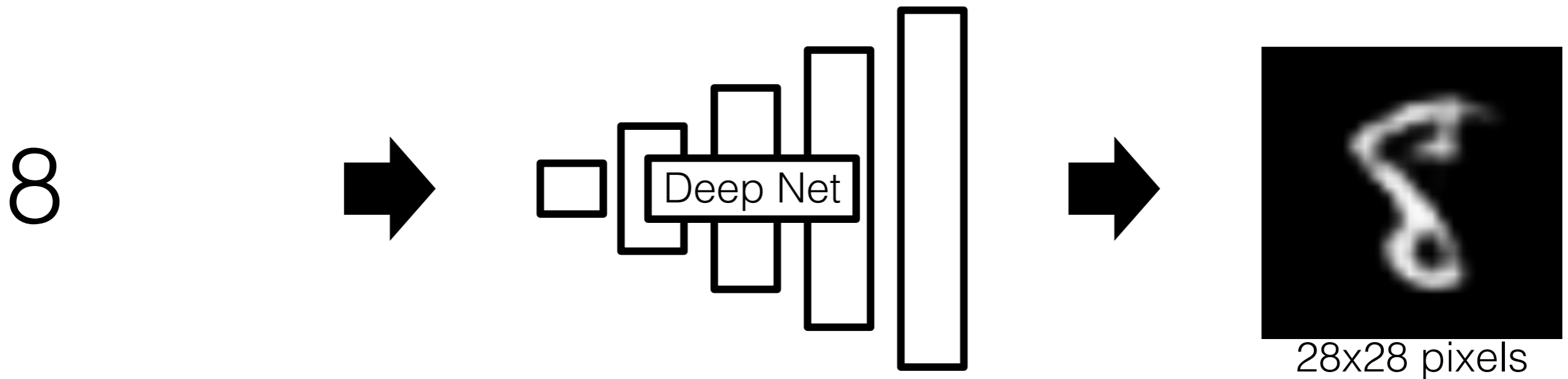


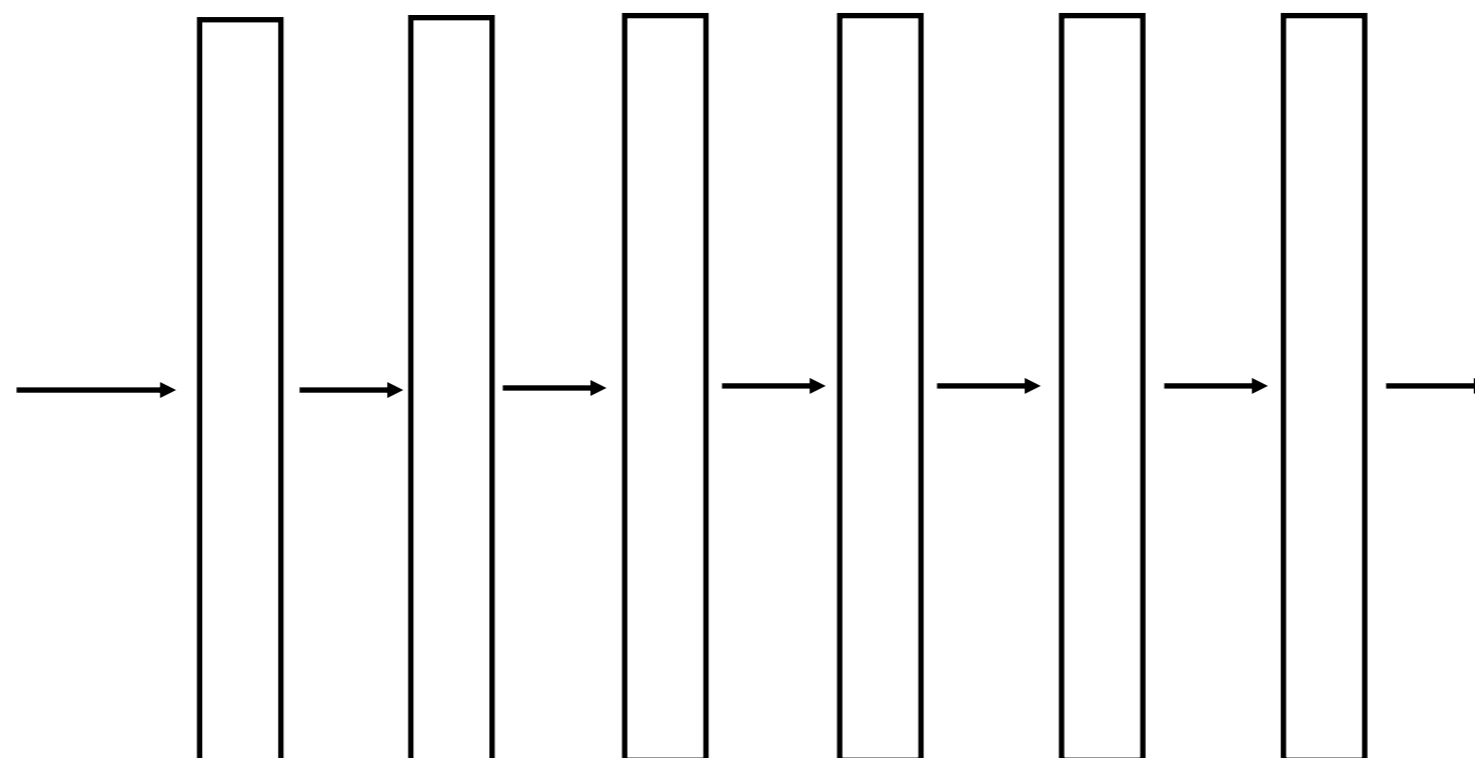
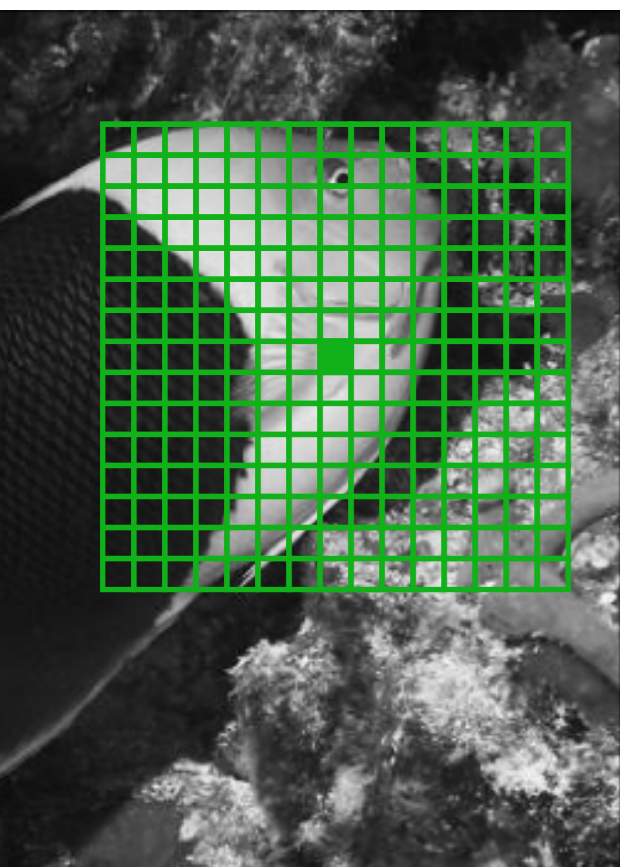
Can Deep Learning Help Graphics?



Cat

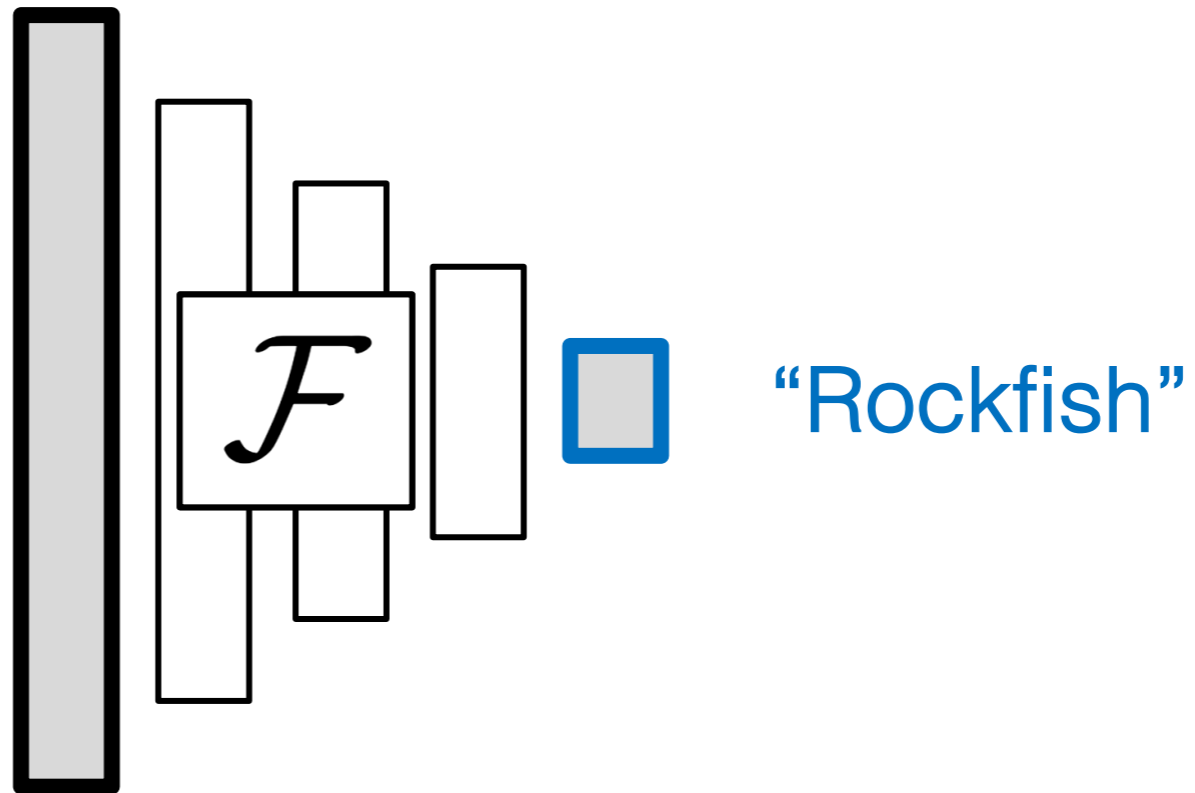
Generating images is hard!



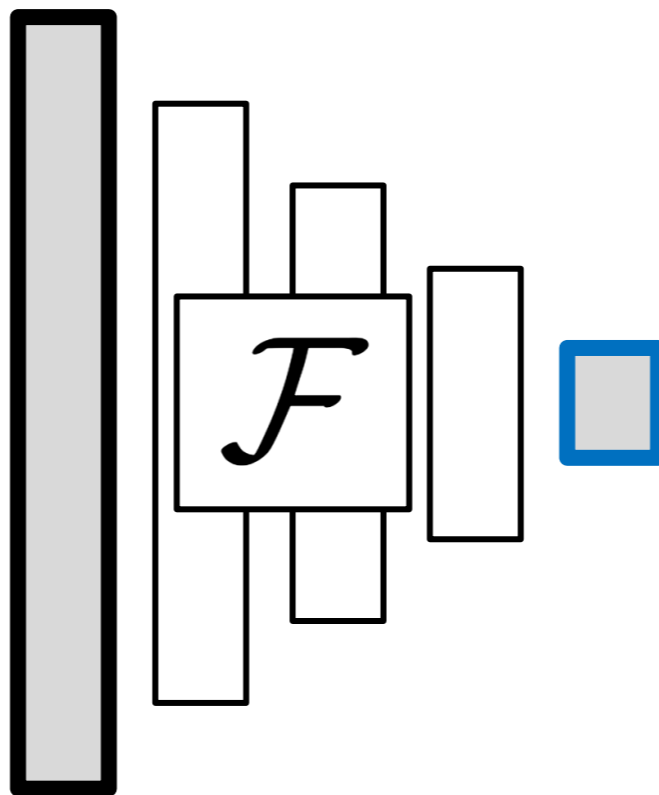


■ “yellow”

Discriminative Deep Networks

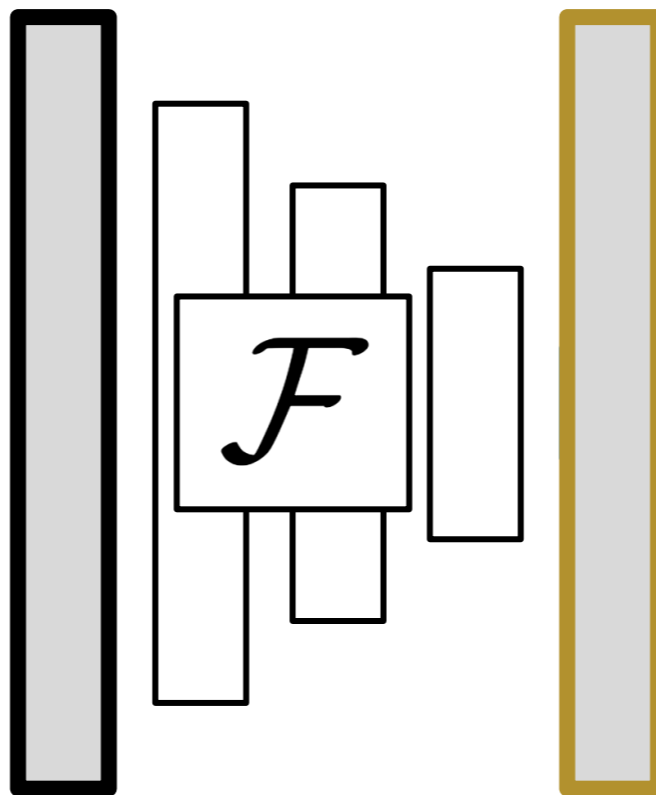


Discriminative Deep Networks



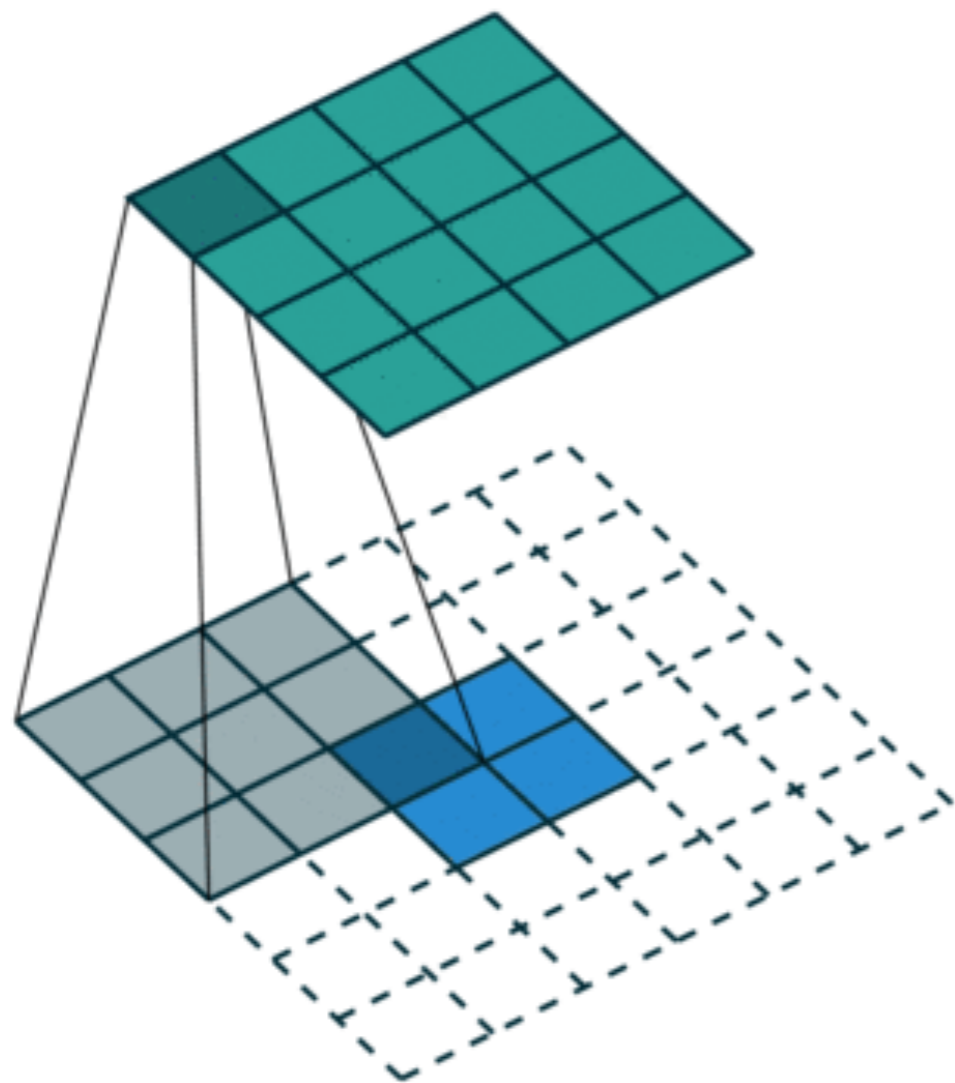
Raw, Unlabeled
Pixels

Generative Deep Networks

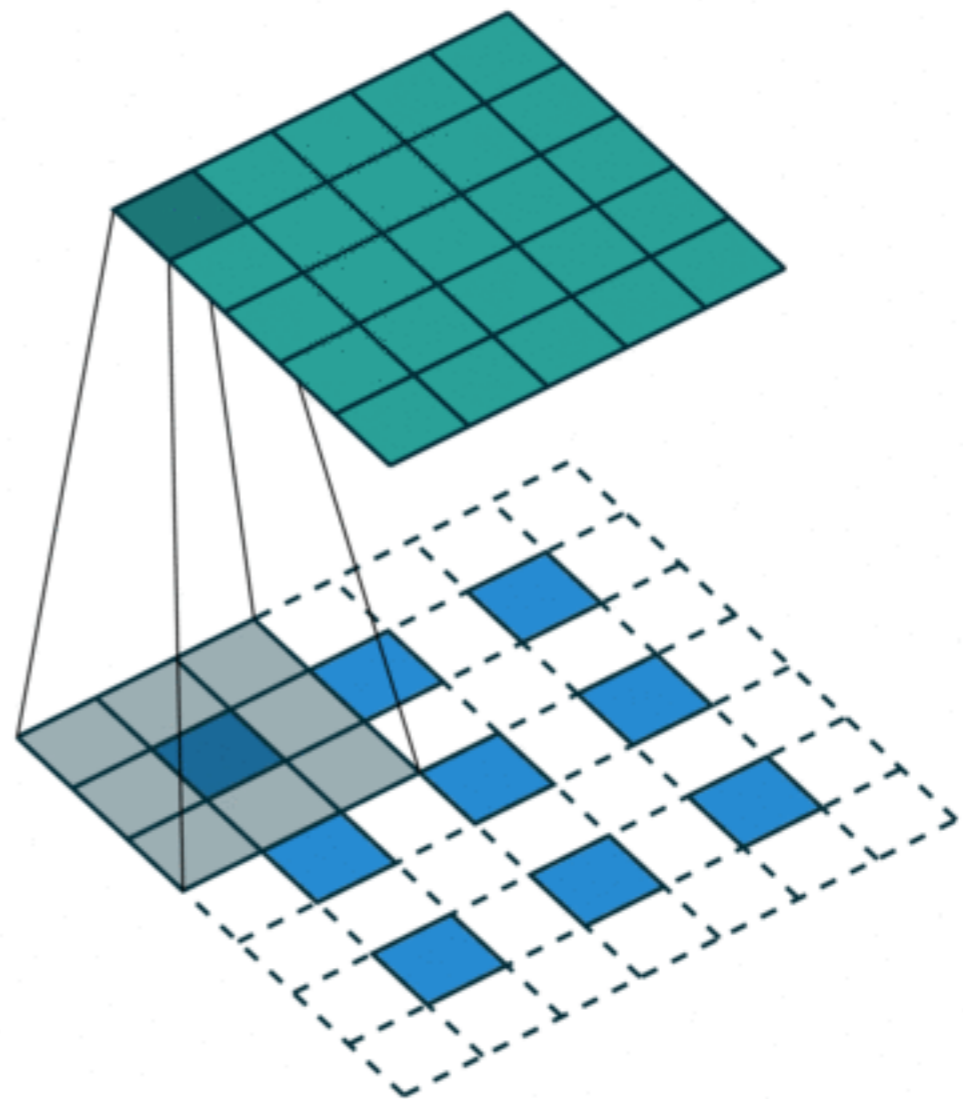


Raw, Unlabeled
Pixels

Fractionally-strided Convolution

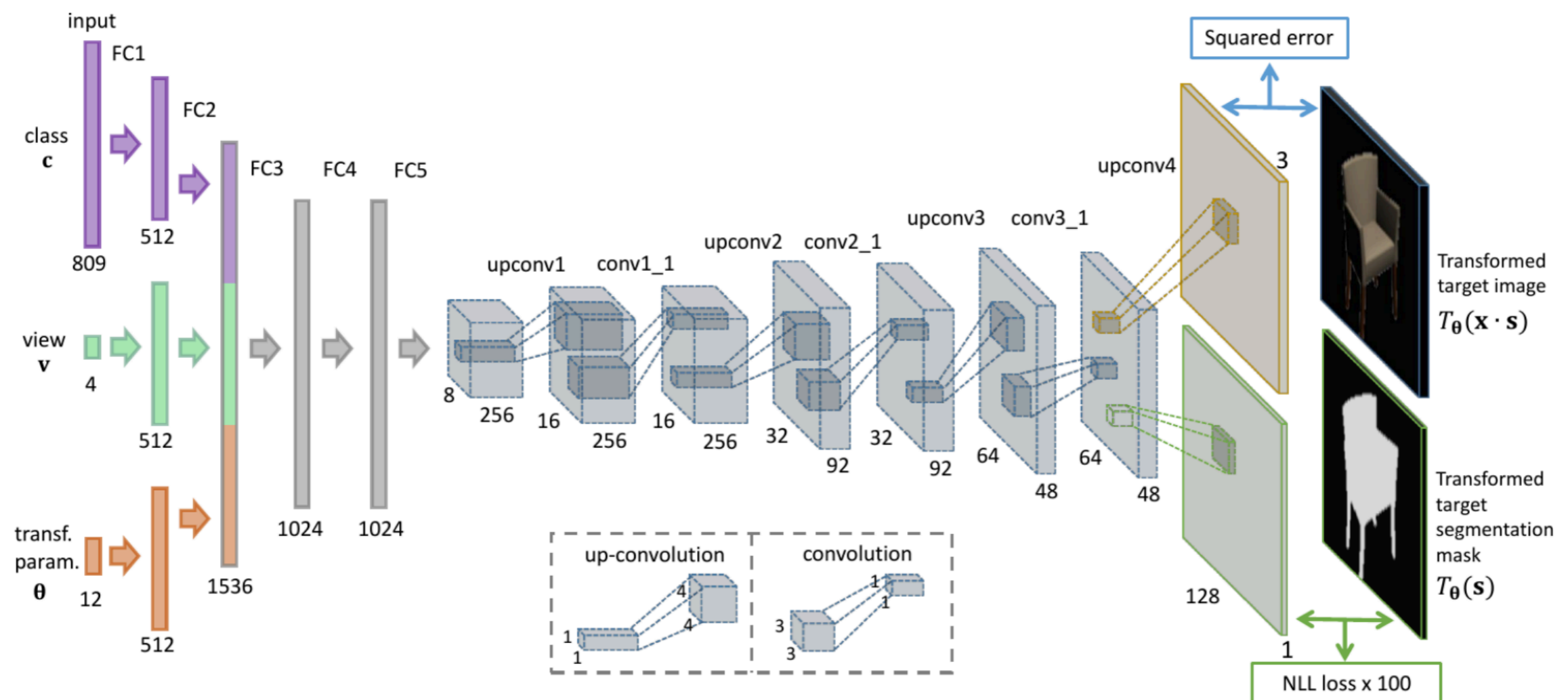


Regular conv



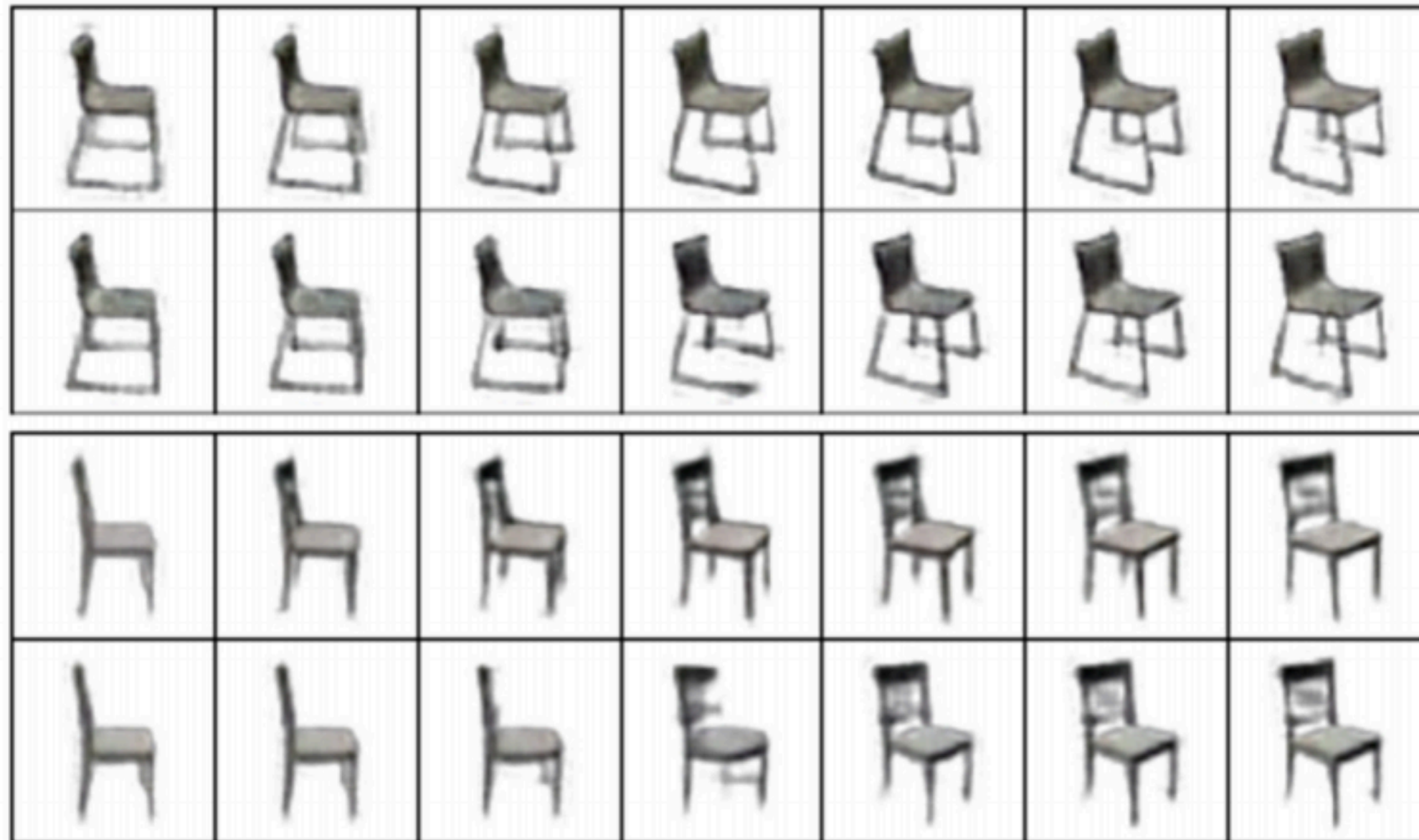
Fractionally-strided conv

Generating chairs given ID, viewpoint, and transformation parameters



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

Generating chairs given ID, viewpoint, and transformation parameters

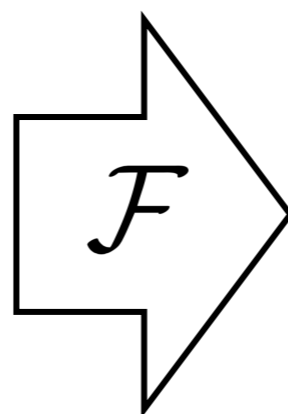
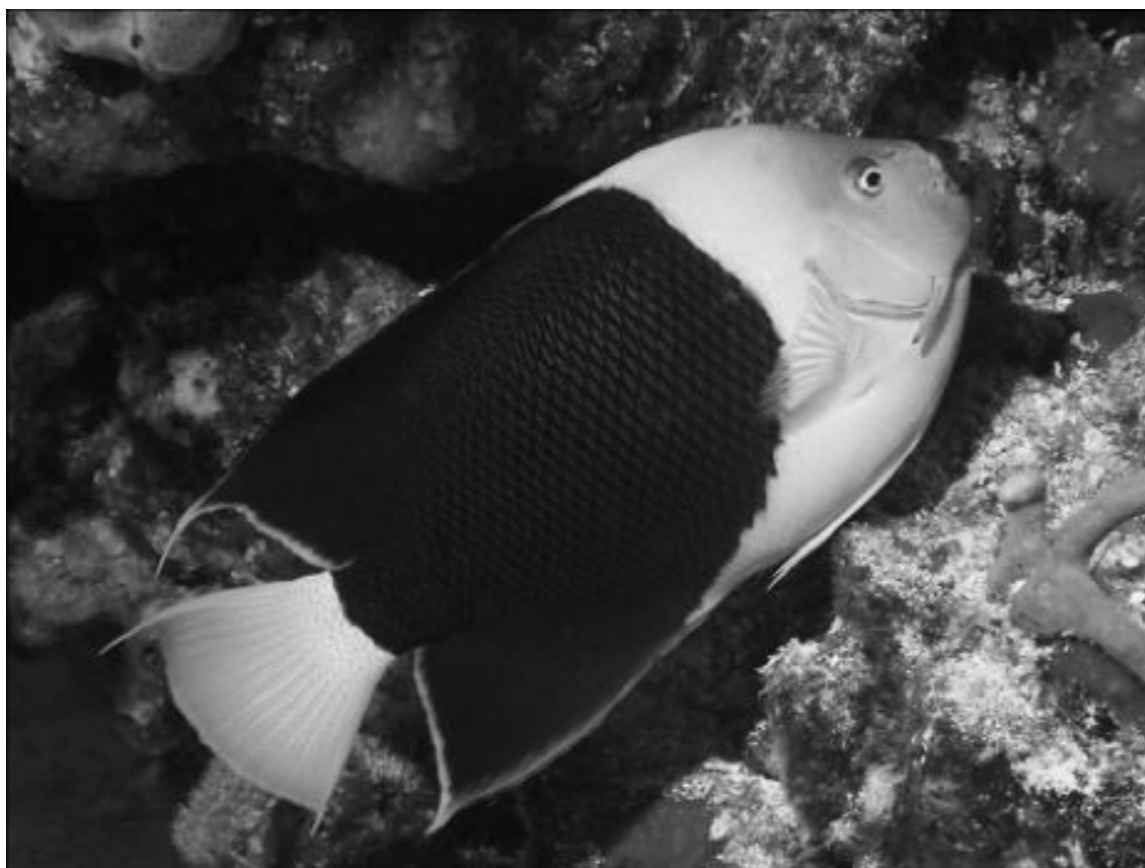


Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks
PAMI 2017 (CVPR 2015)



60

Ansel Adams. *Yosemite Valley Bridge.*

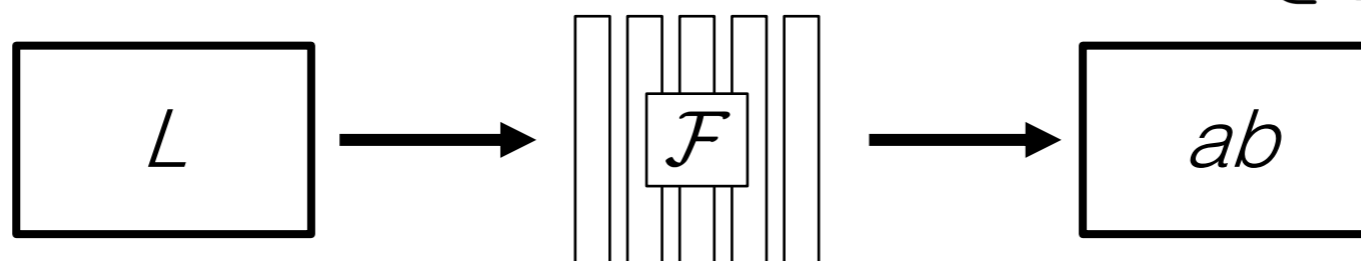


Grayscale image: L channel

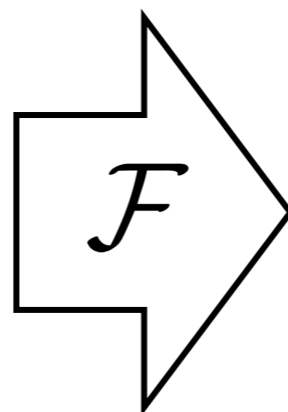
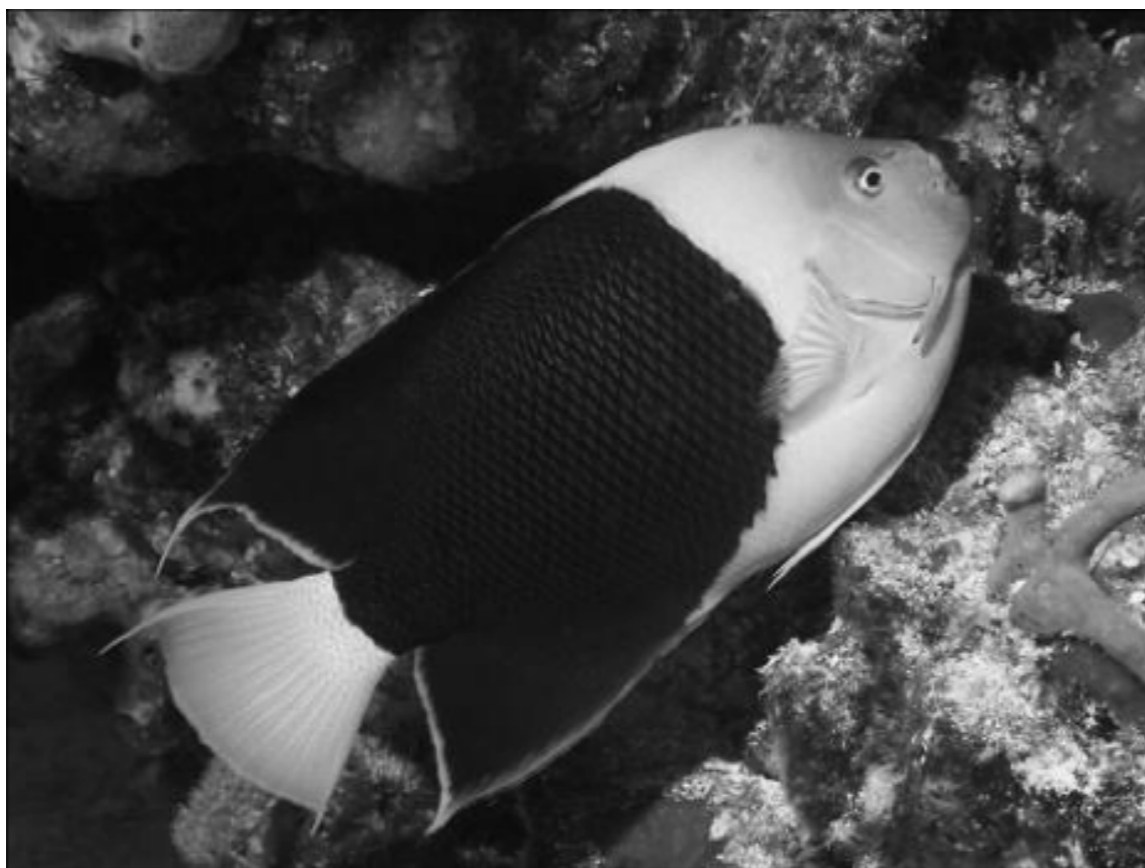
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: ab channels

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



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

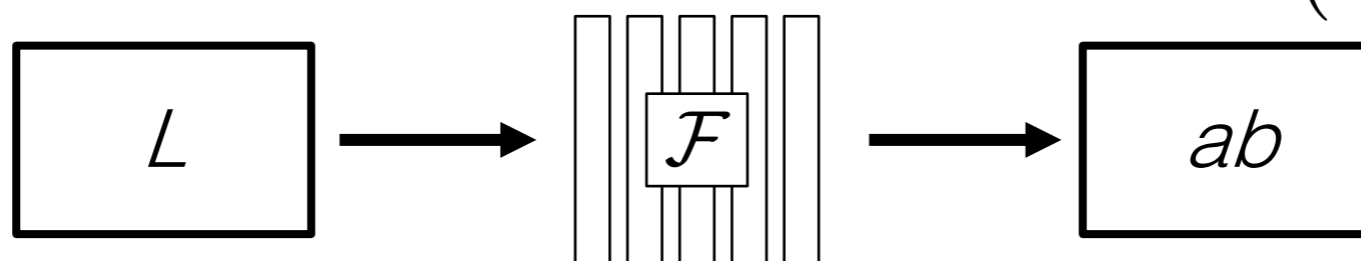


Grayscale image: L channel

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

Concatenate (L, ab) channels

$$(\mathbf{X}, \hat{\mathbf{Y}})^{62}$$



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

Simple L2 regression doesn't work 😞

Input



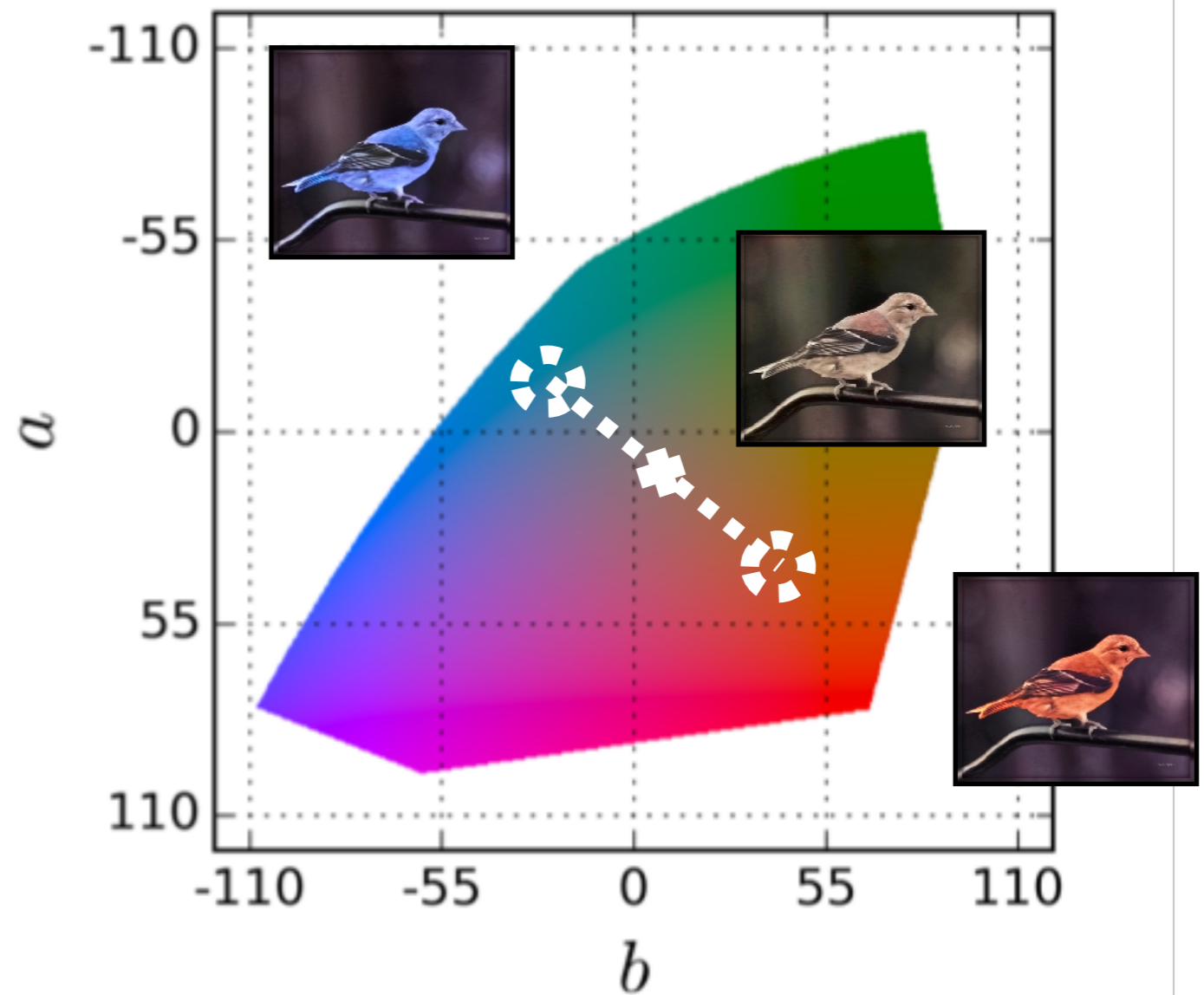
Output



Ground truth



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



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

Better Loss Function

$$\theta^* = \arg \min_{\theta} \ell(\mathcal{F}_{\theta}(\mathbf{X}), \mathbf{Y})$$

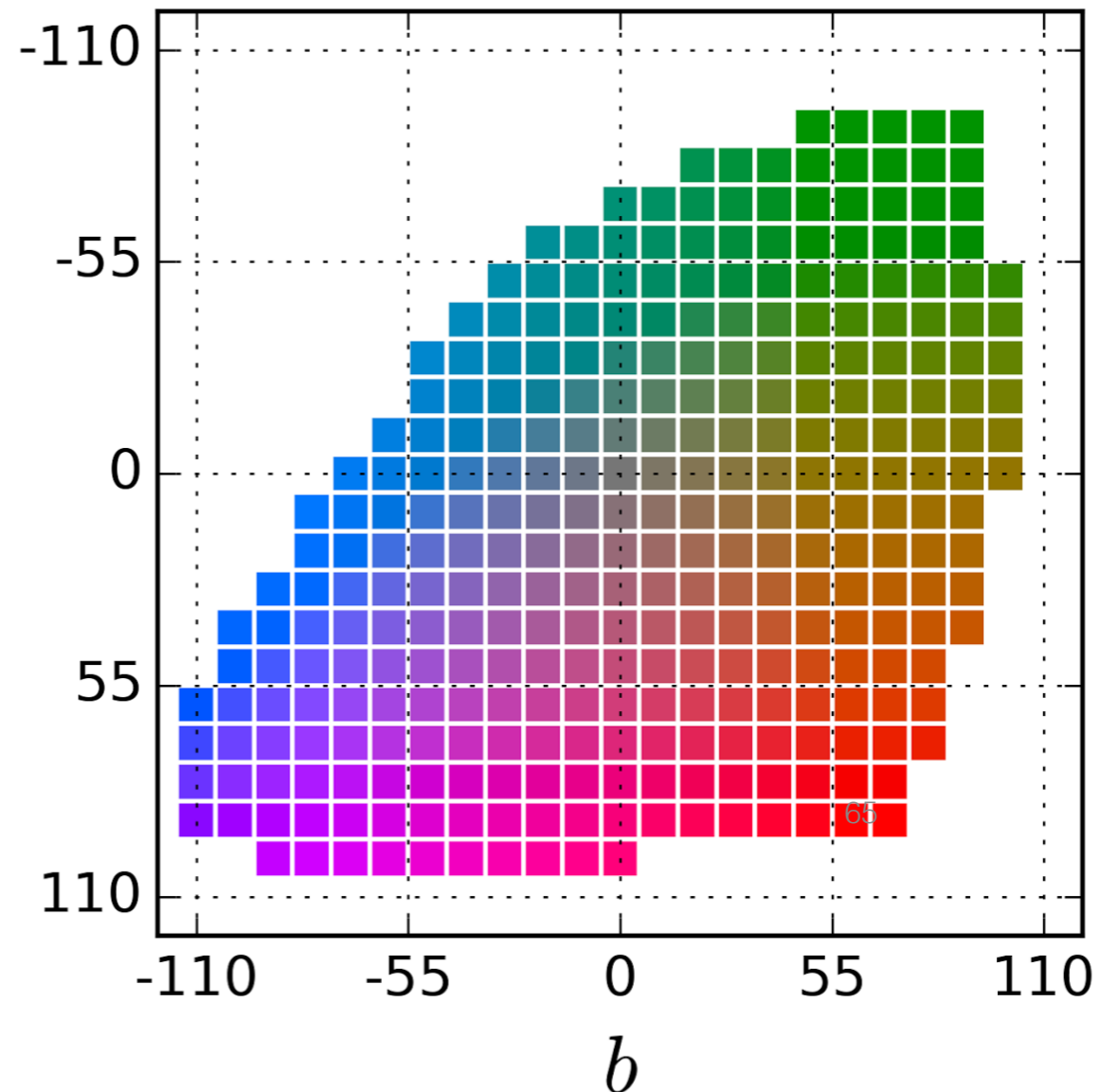
- Regression with L2 loss inadequate

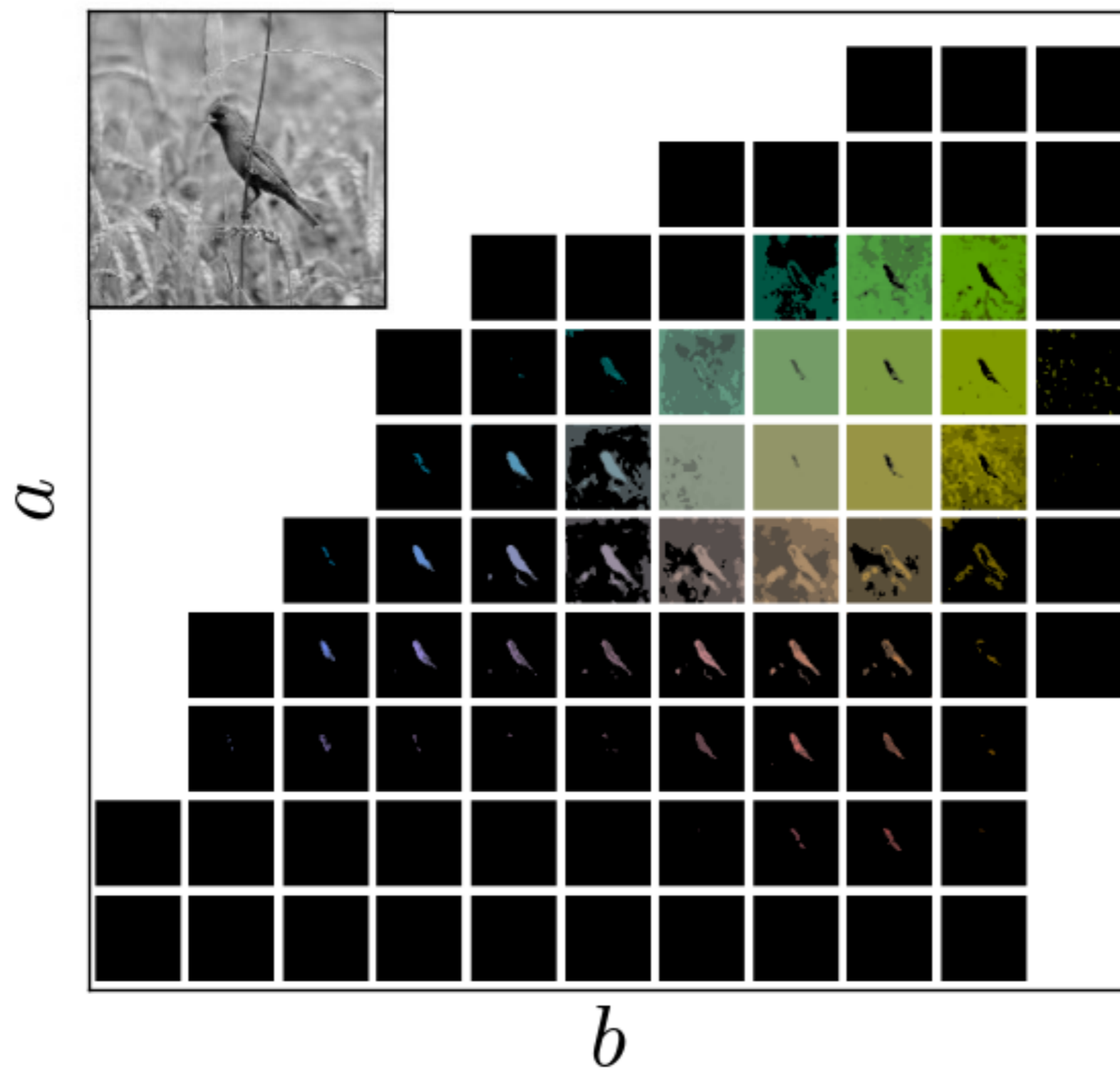
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

- Use per-pixel multinomial classification^a

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

Colors in *ab* space
(discrete)





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 2021

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