

Convolutional Network for Image Synthesis

Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2025

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”

?

==

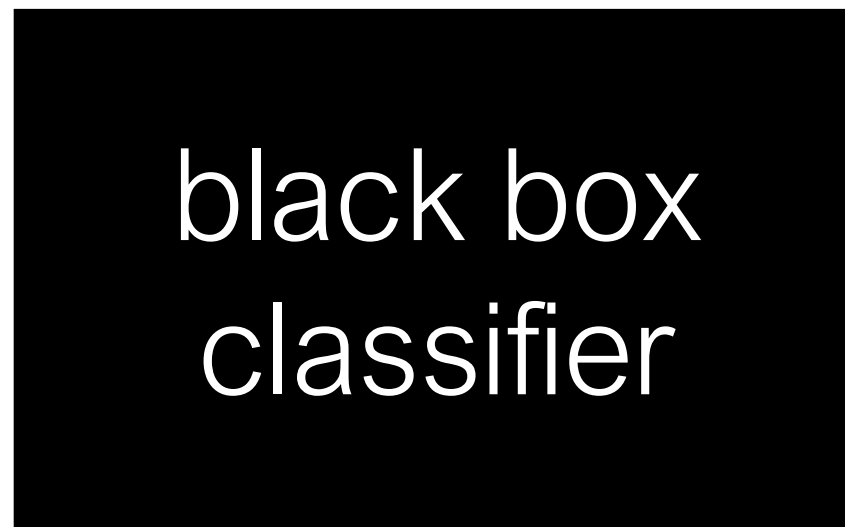


“Penguin”

Machine Learning as data association



image X



“Penguin”

label Y

At test time...



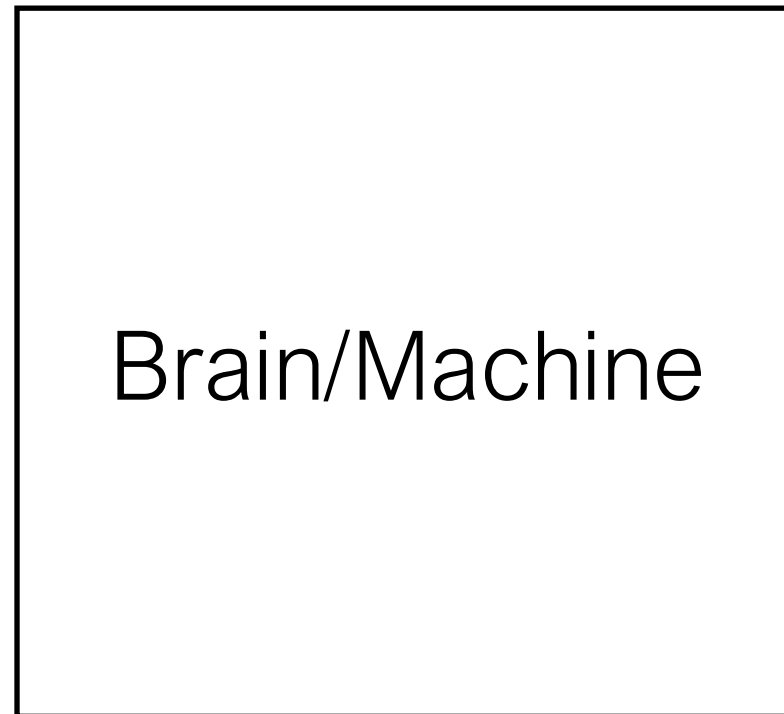
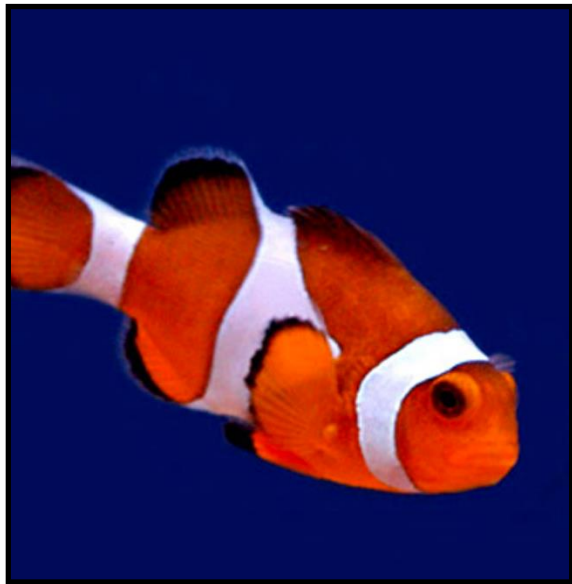
black box
classifier



?

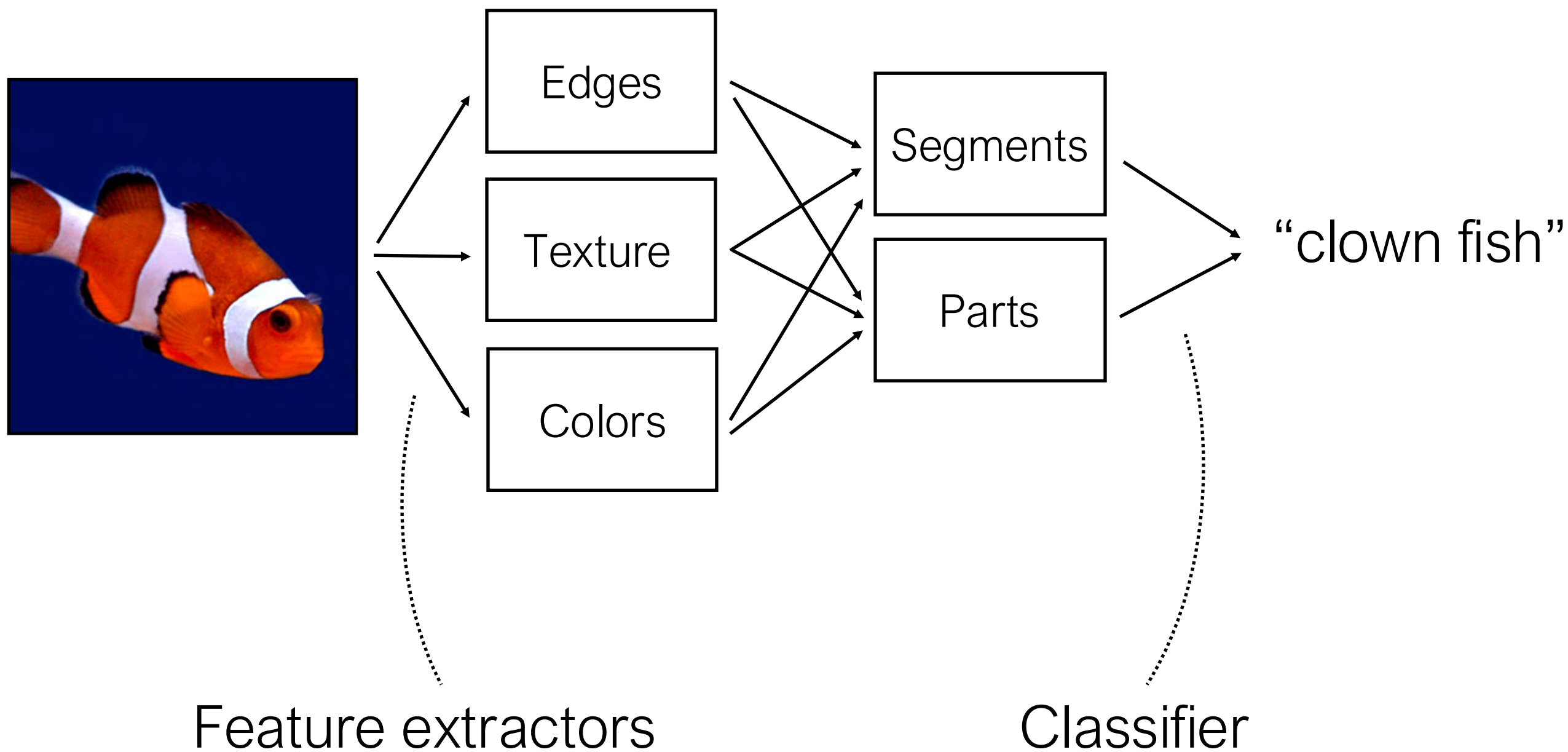
image X

Basic idea



“clown fish”

Object recognition



Object recognition

Learned



Edges

Texture

Colors

Segments

Parts

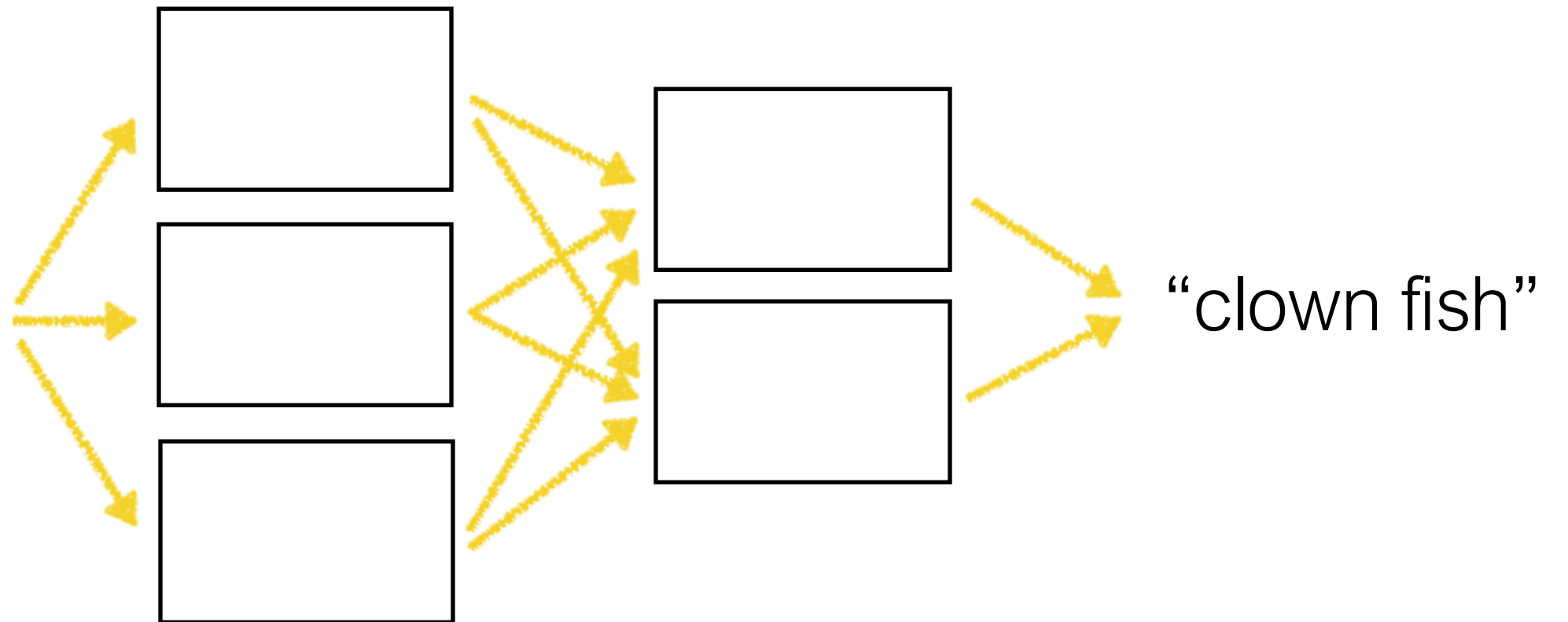
“clown fish”

Feature extractors

Classifier

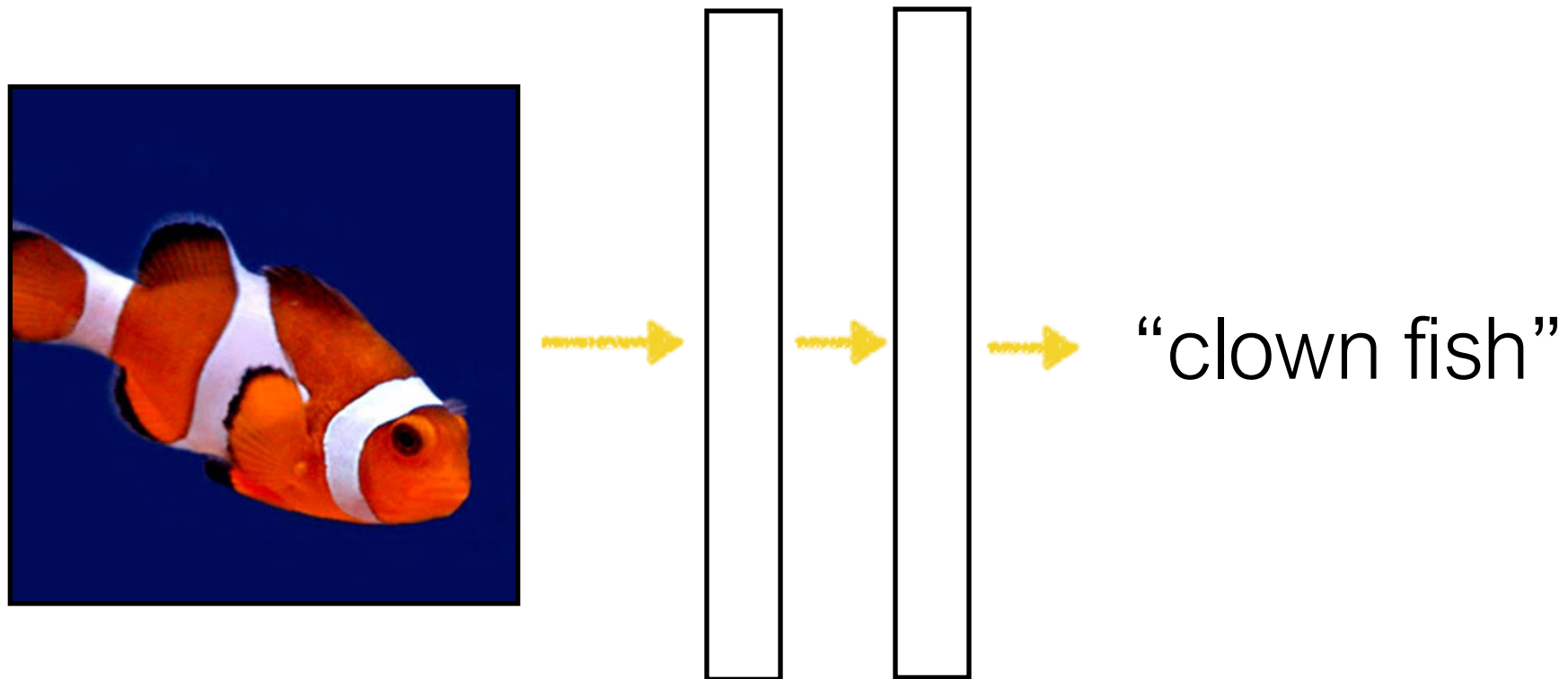
Neural network

Learned



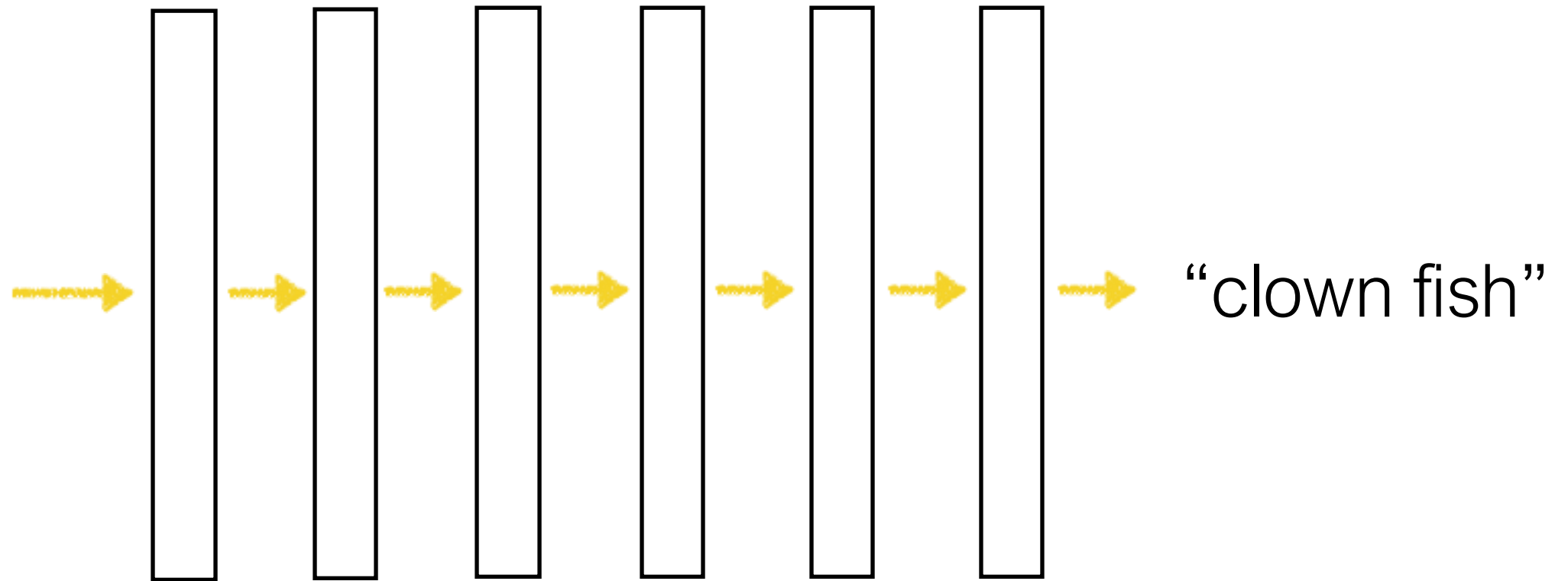
Neural network

Learned

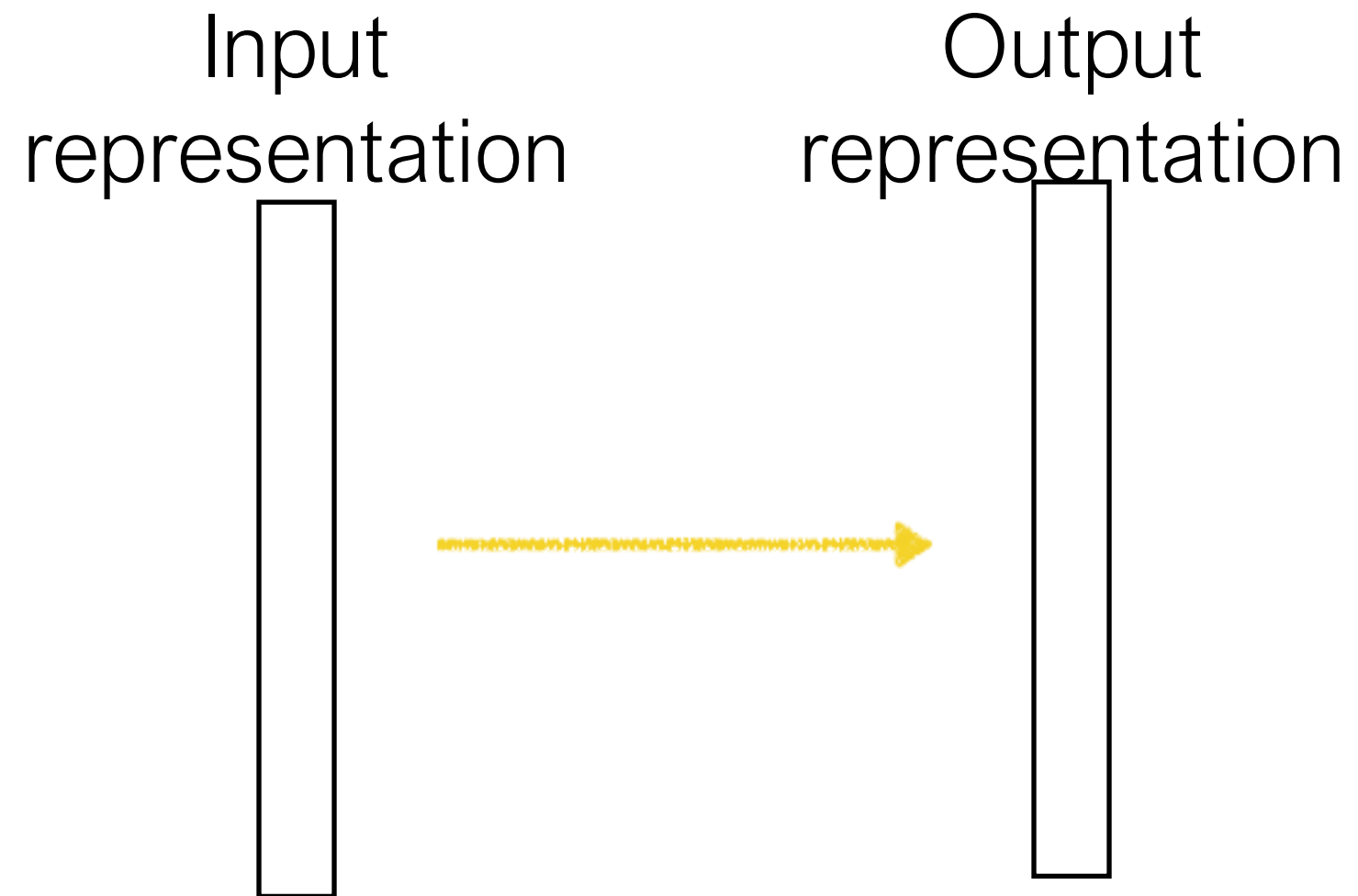


Deep neural network

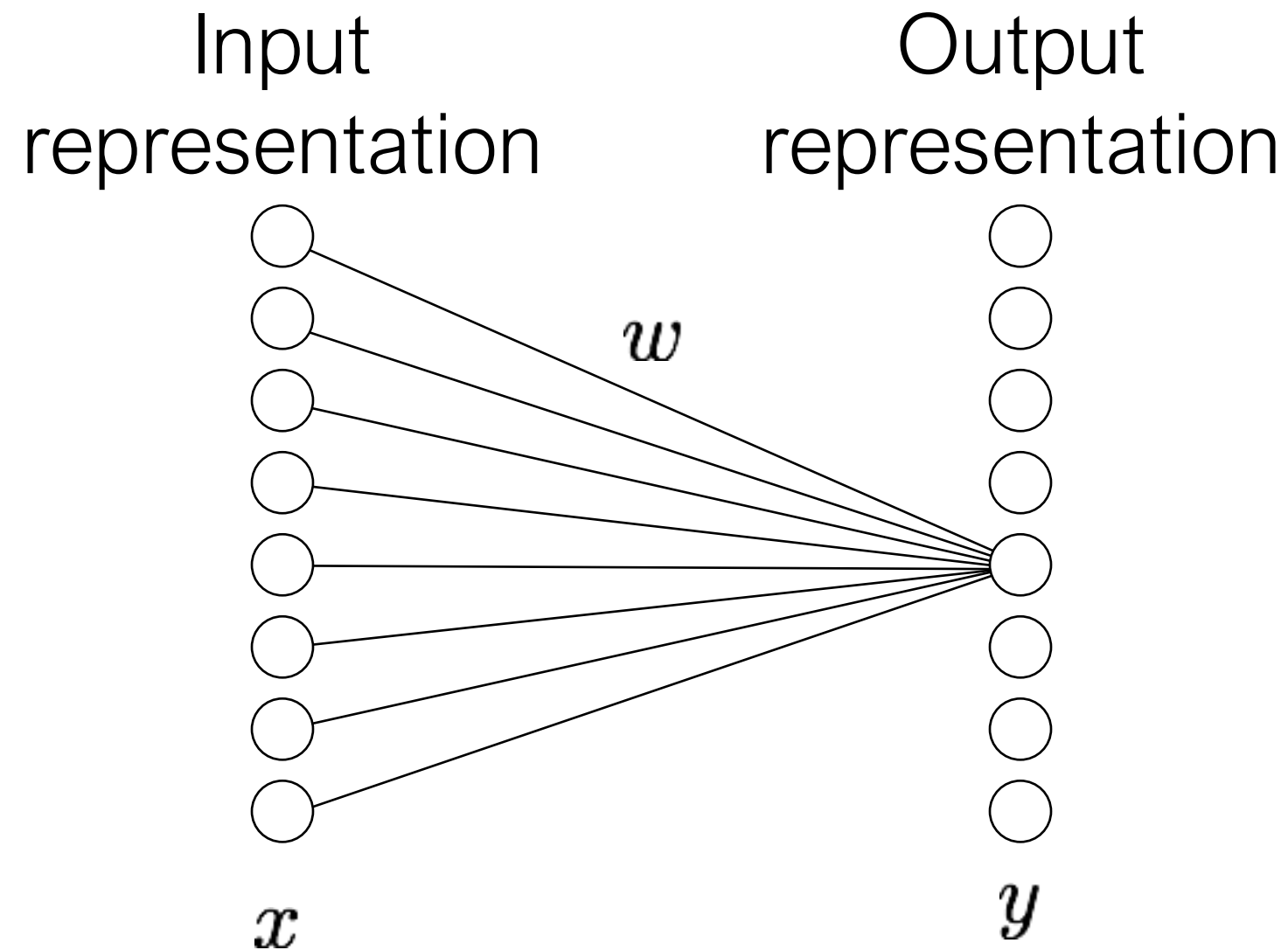
Learned



Computation in a neural net



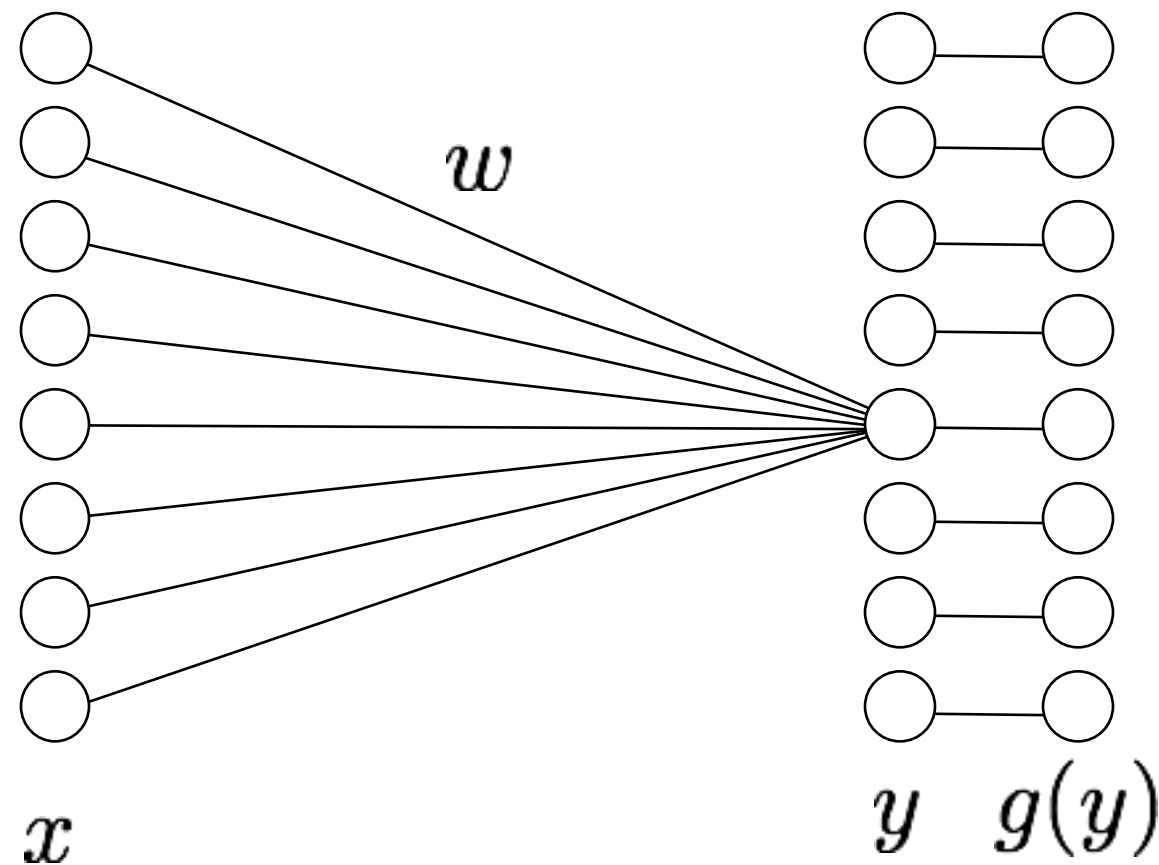
Computation in a neural net



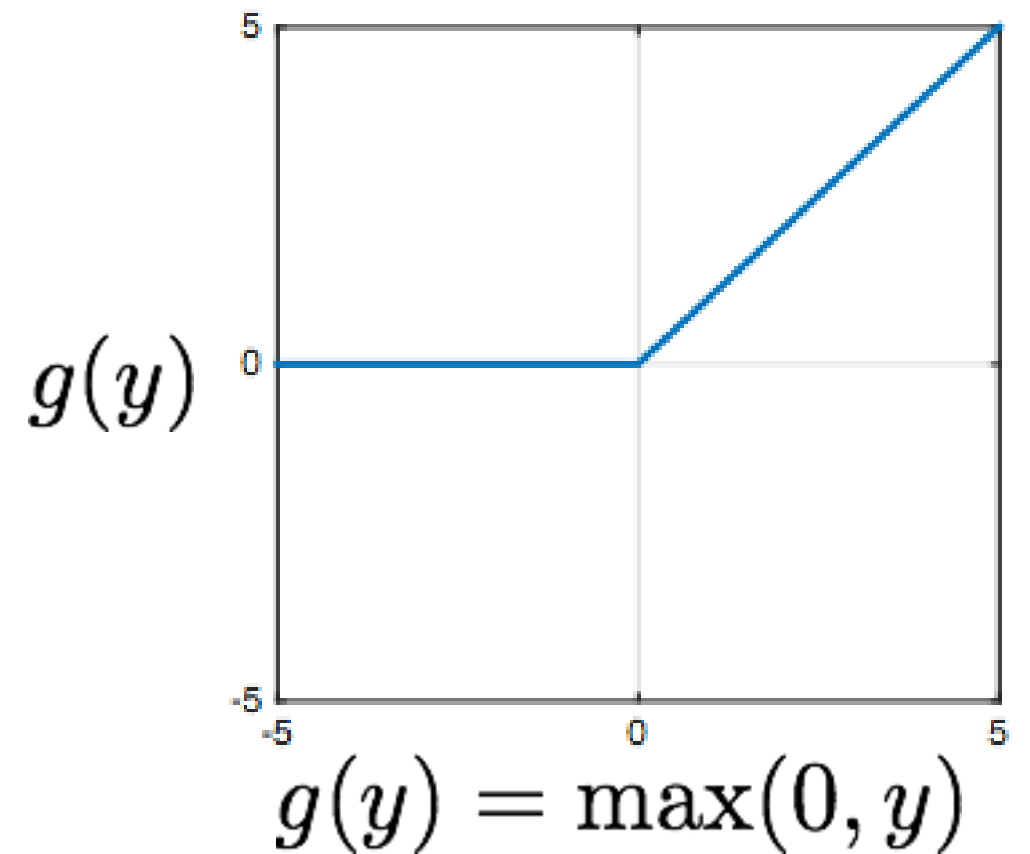
$$y_j = \sum_i w_{ij} x_i$$

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

Computation in a neural net

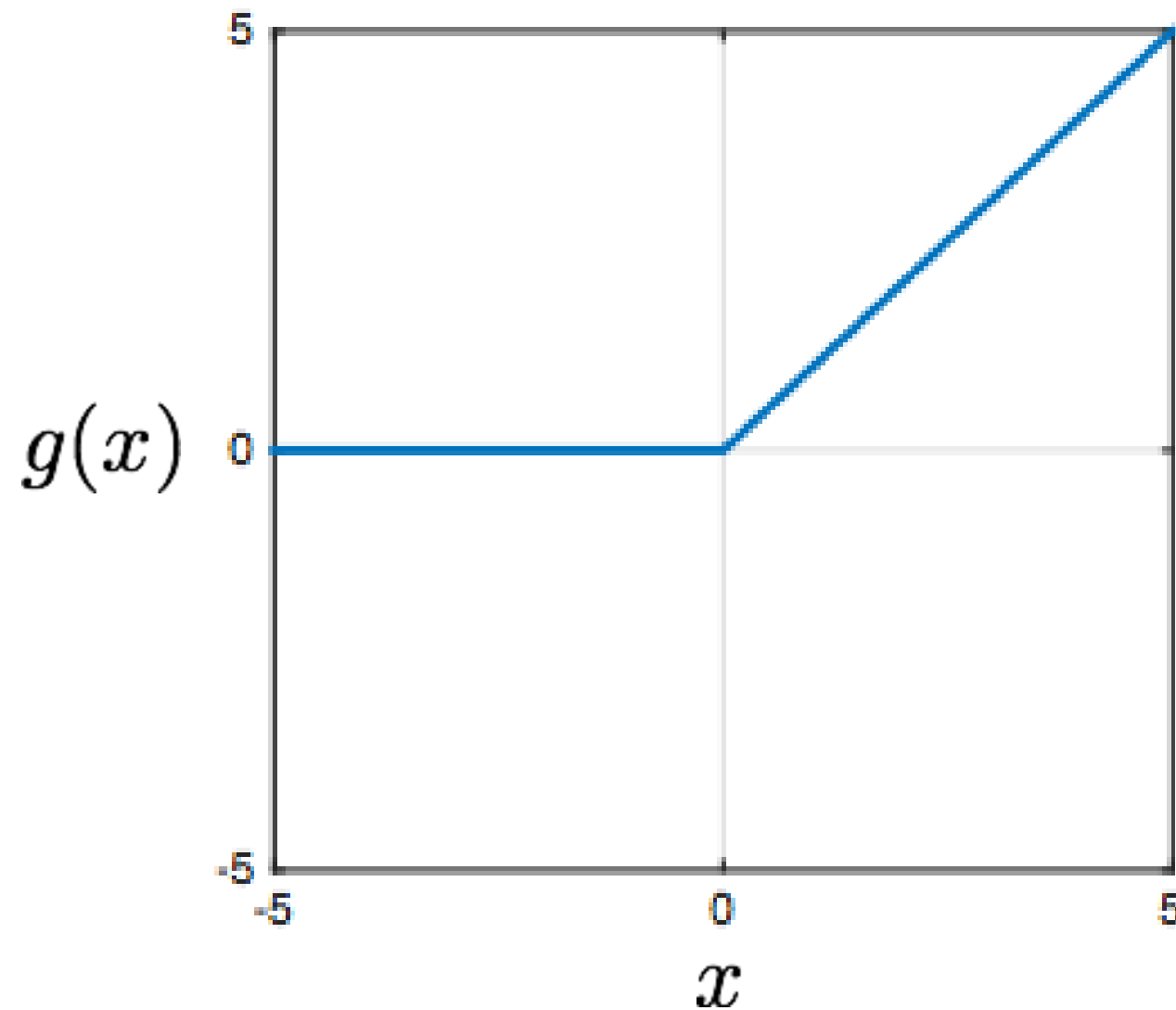


Rectified linear unit (ReLU)



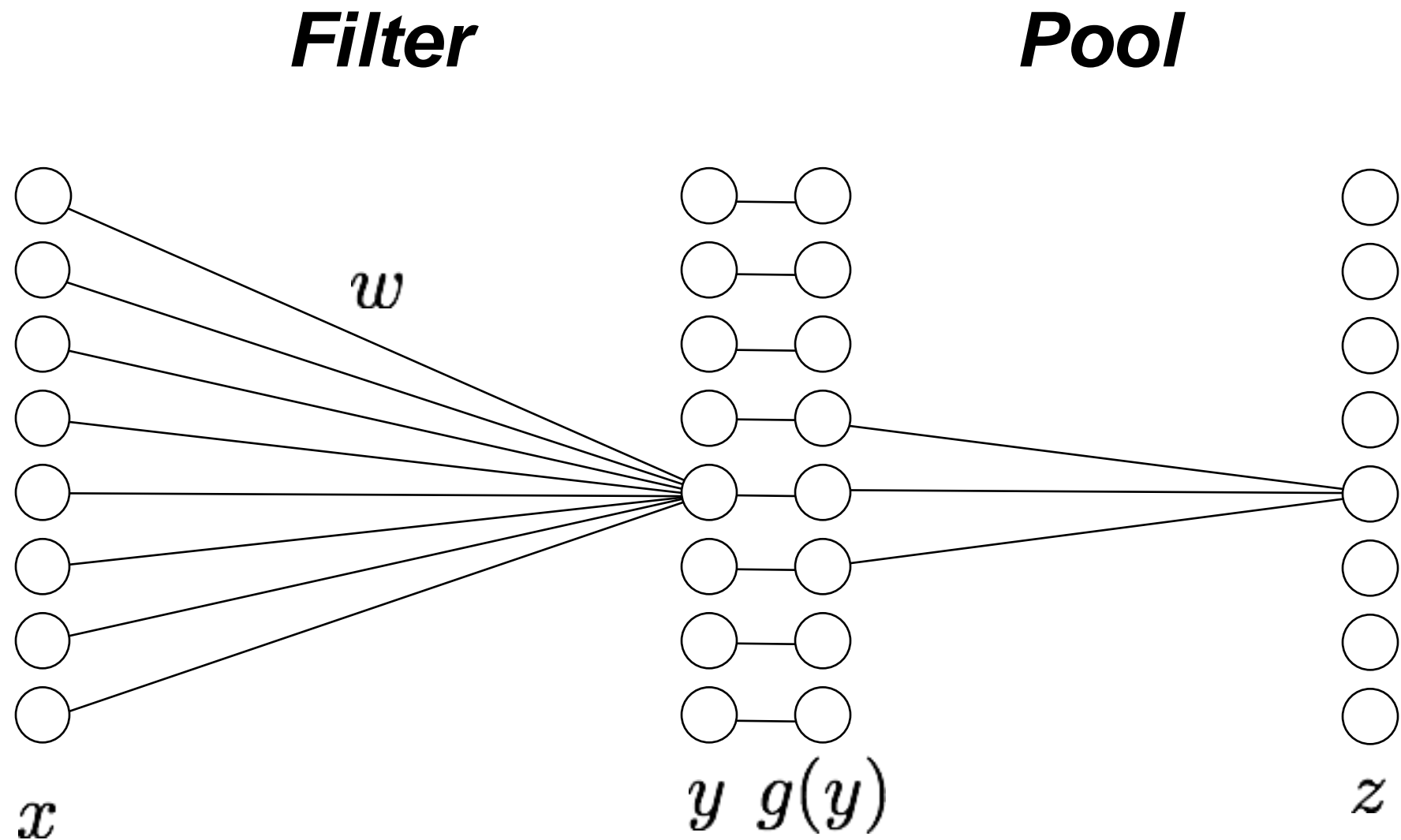
Computation in a neural net

Rectified linear unit (ReLU)



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

Computation in a neural net

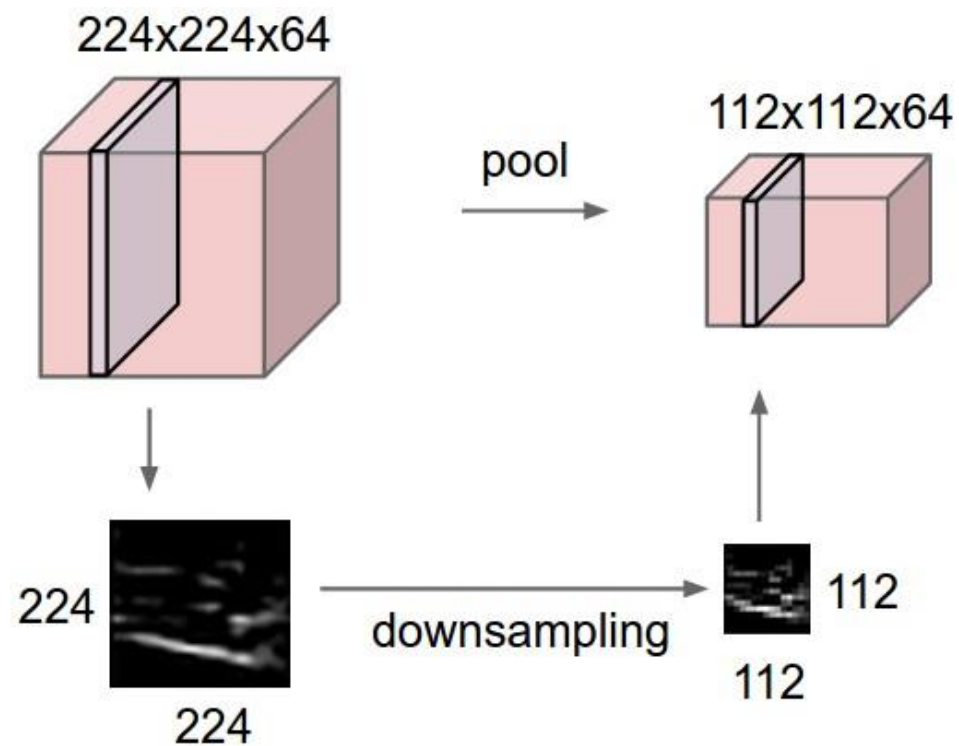
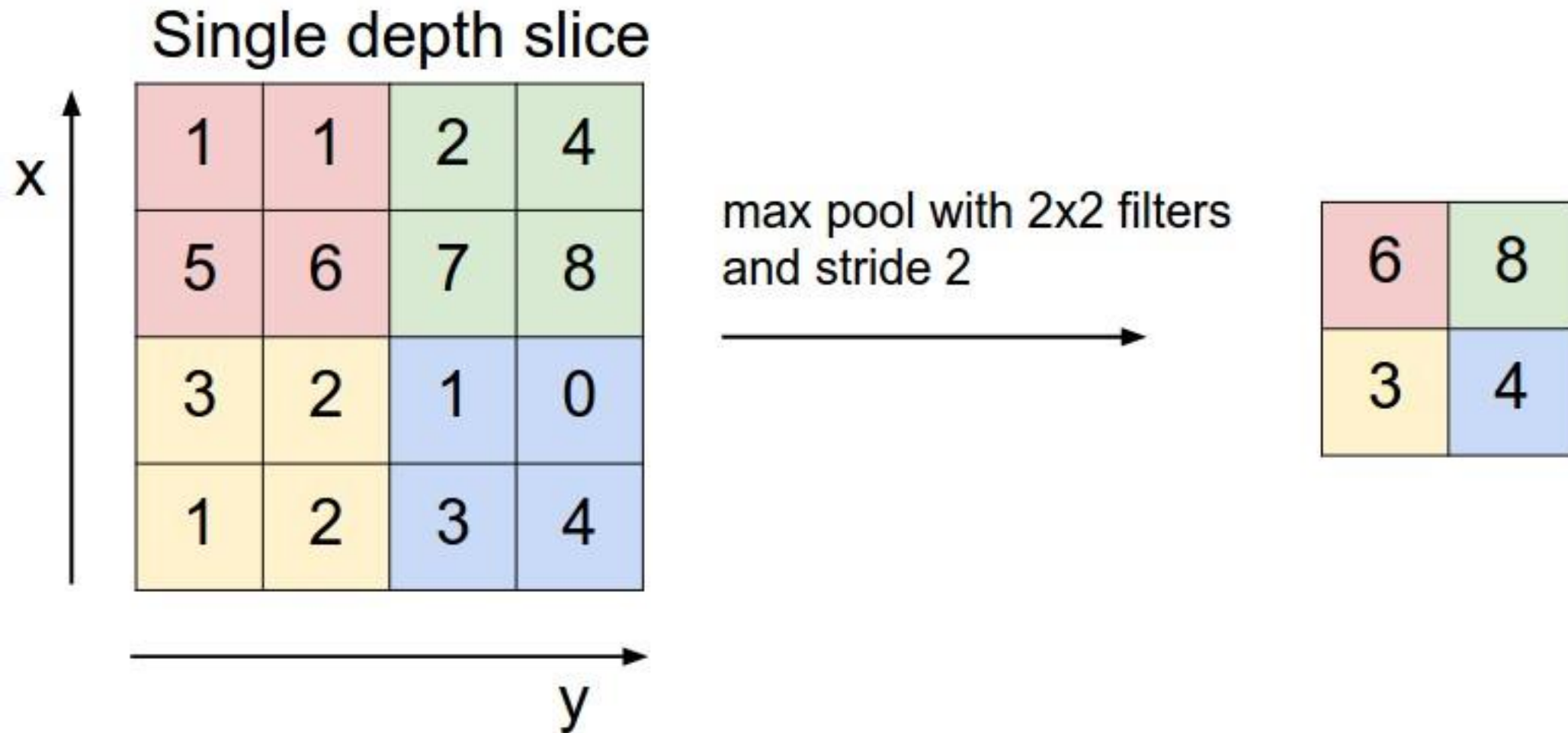


$$y_j = \sum_i w_{ij} x_i$$

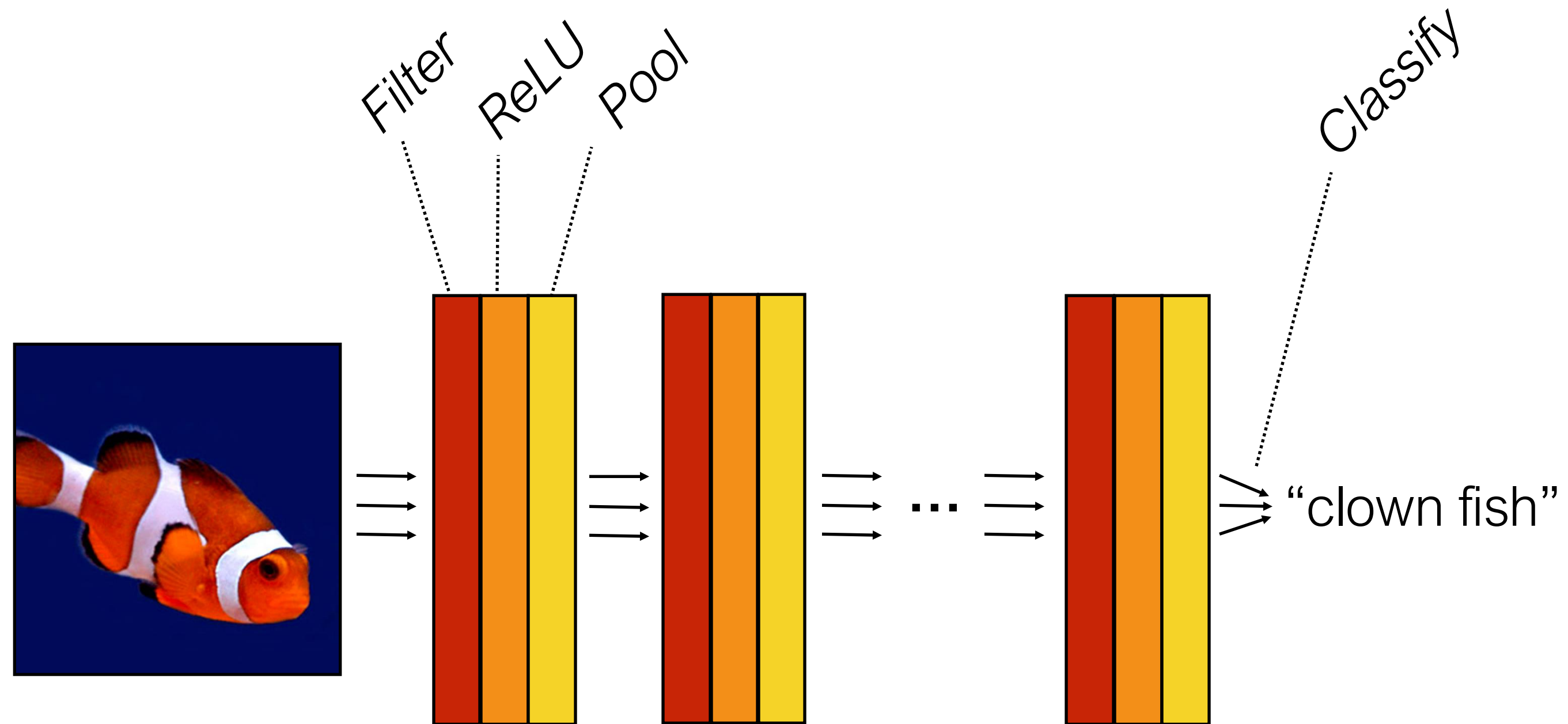
$$z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$

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

Computation in a neural net

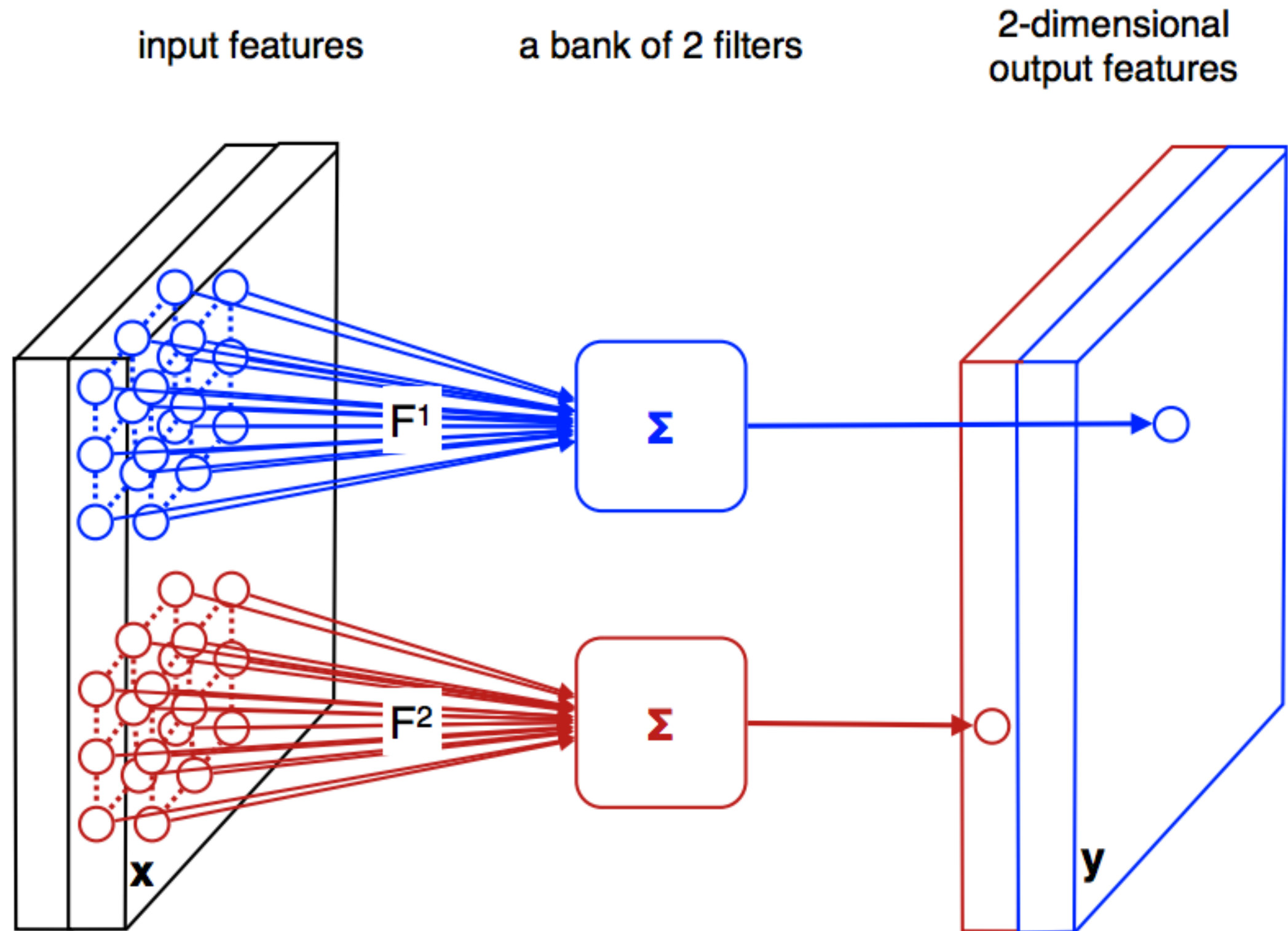


Computation in a neural net



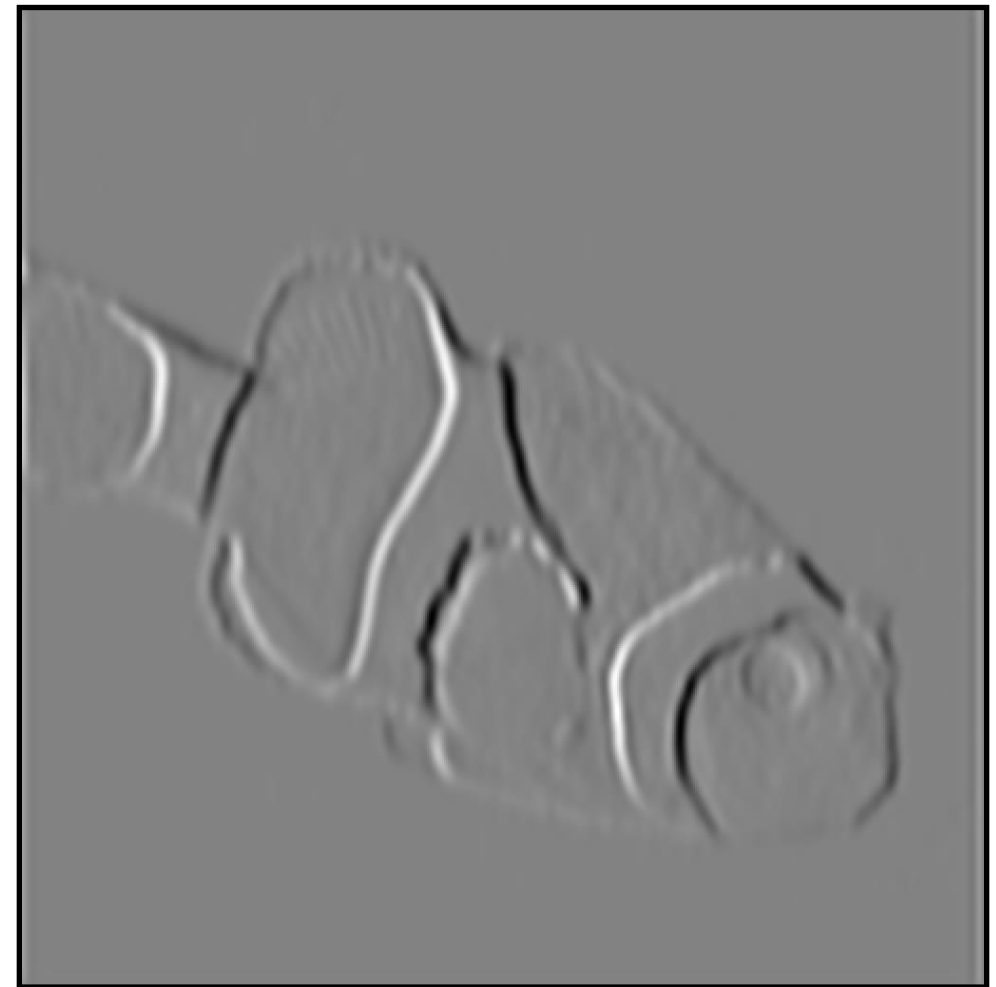
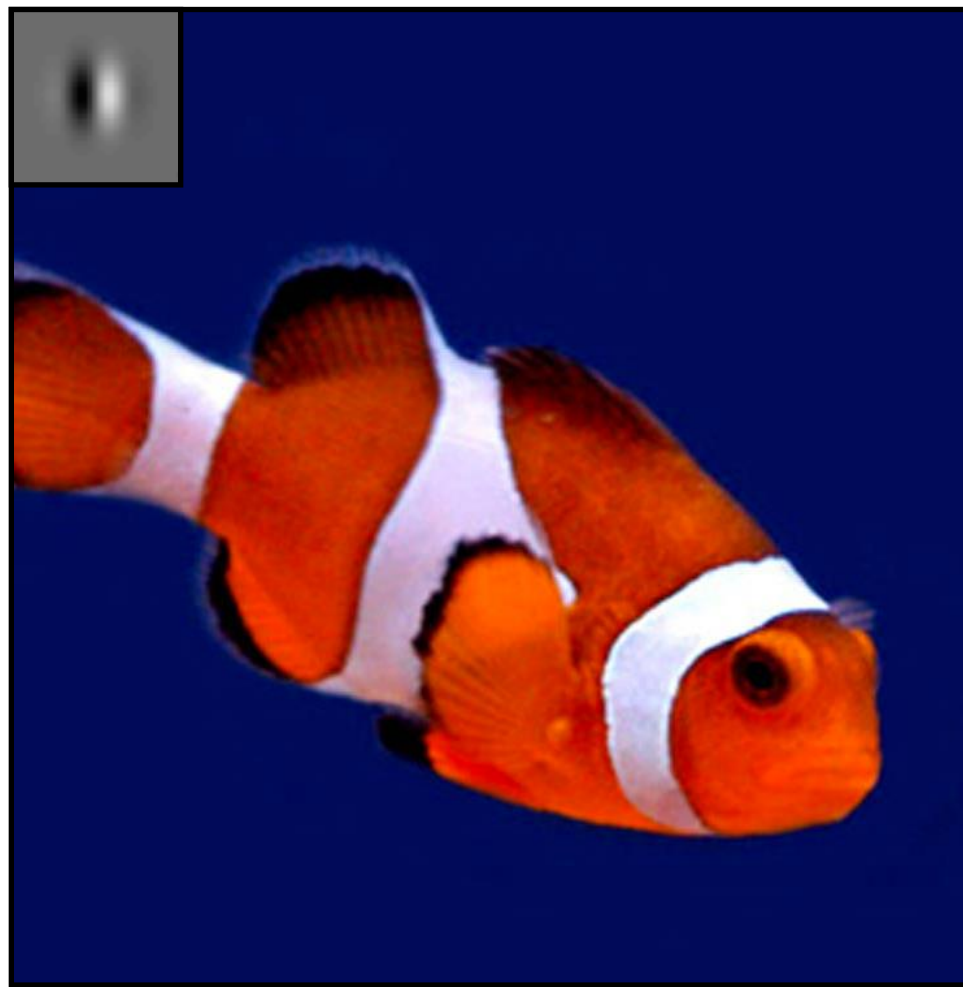
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Convolutional Neural Nets

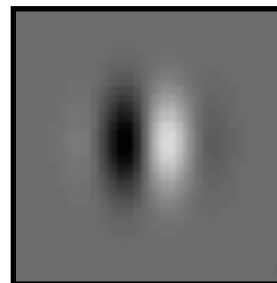


Convolutional Neural Nets

Convolution

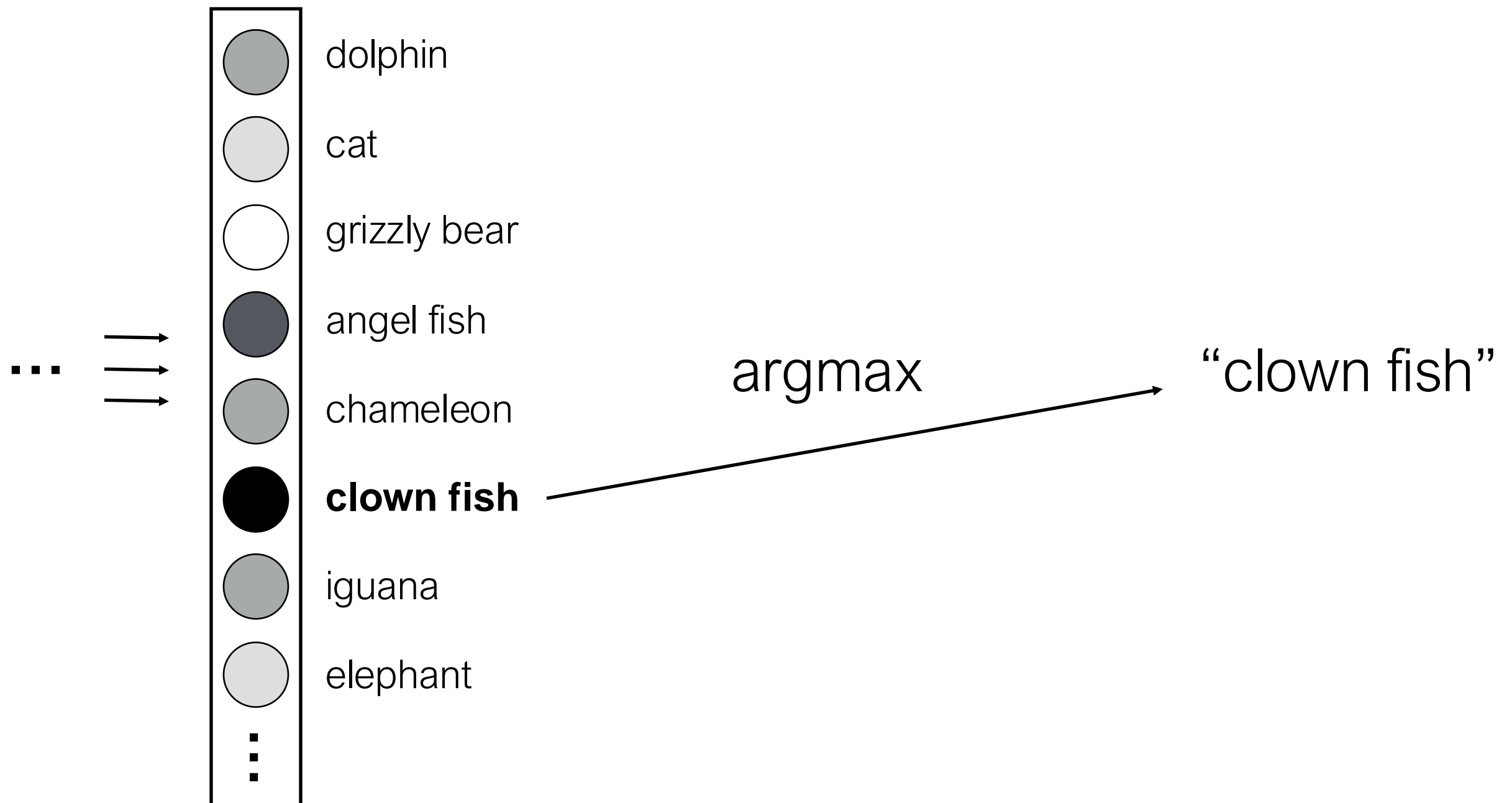


filter



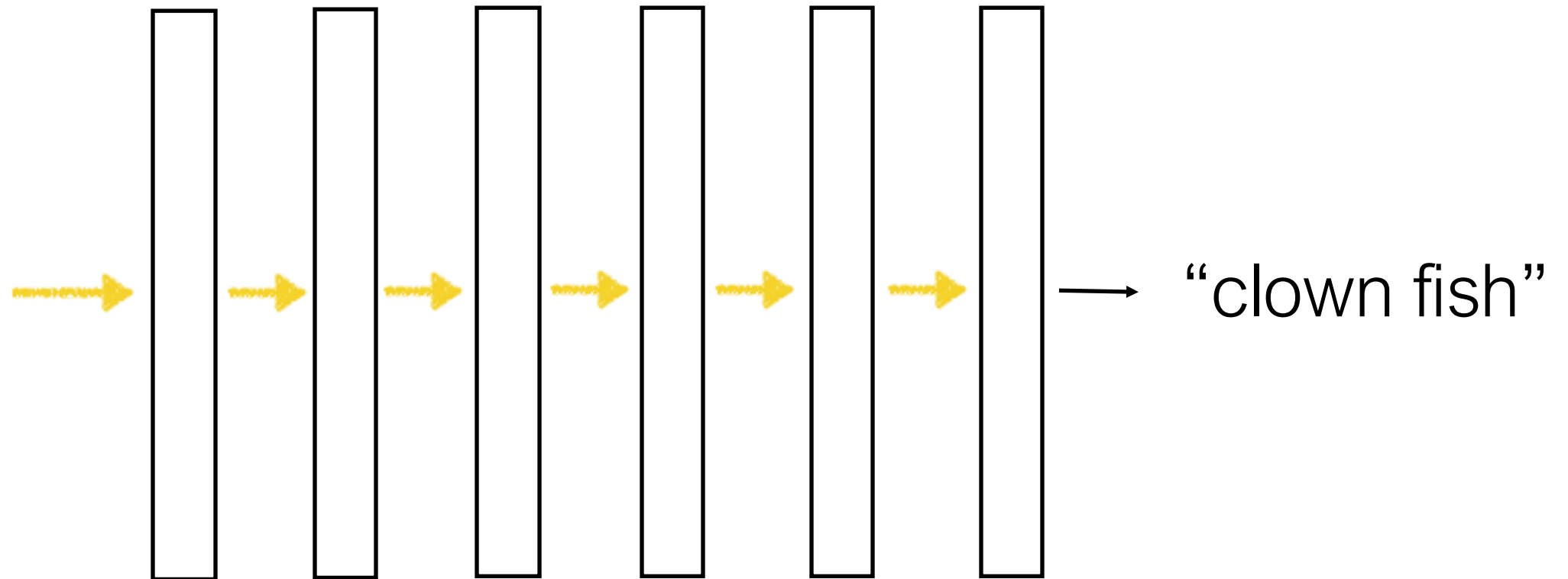
Computation in a neural net

Last layer



Learning with deep nets

Learned



Learning with deep nets



→ “clown fish”



→ “grizzly bear”



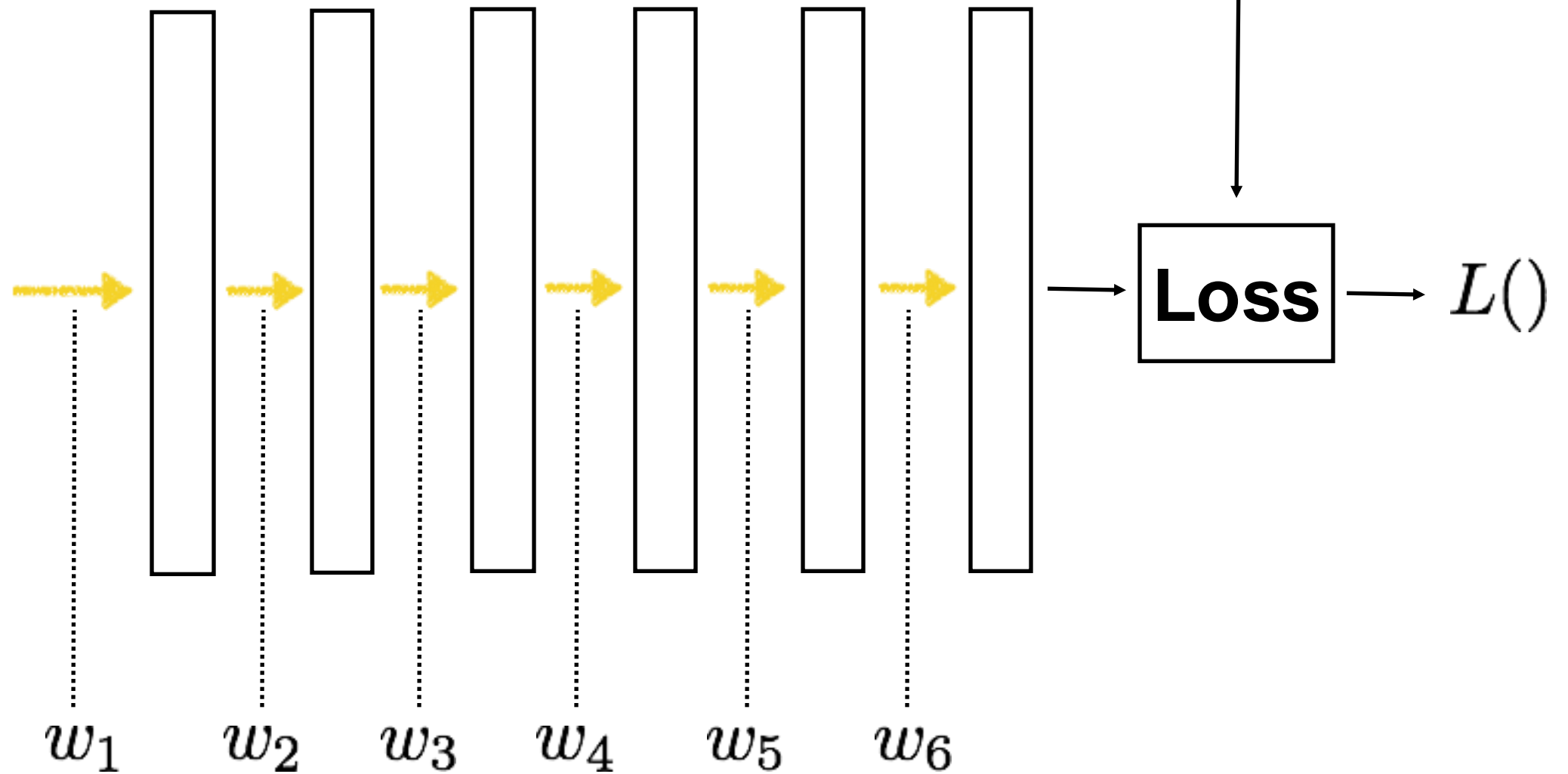
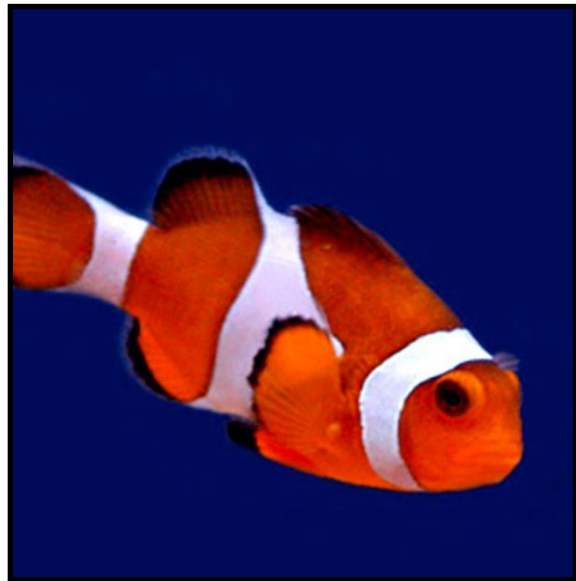
→ “chameleon”

Train network to
associate the right label
with each image

Learning with deep nets

Learned

“clown fish”

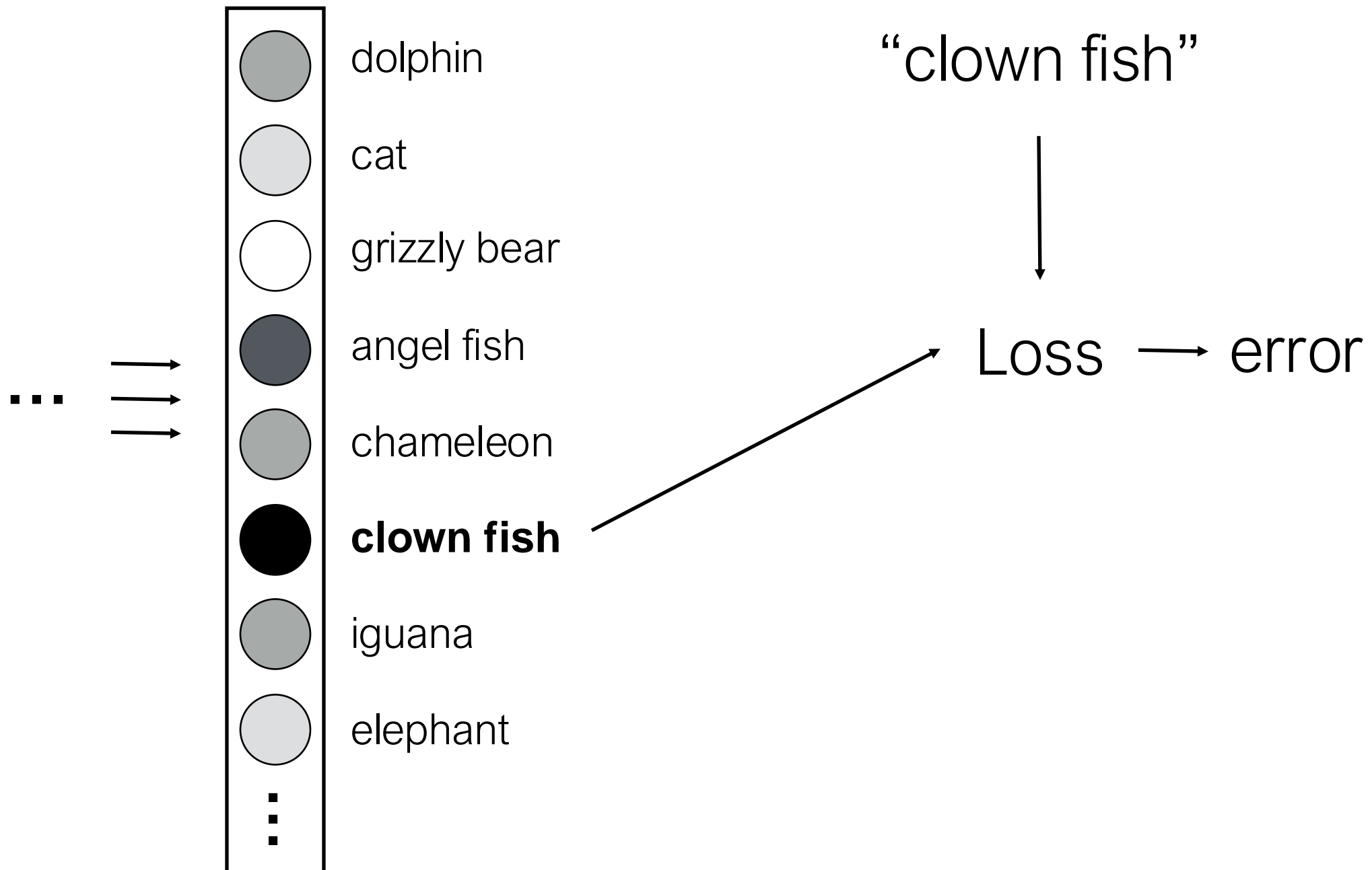


$$\underset{\mathbf{w}}{\operatorname{argmin}} L(w_1, \dots, w_6)$$

Loss function

Network output

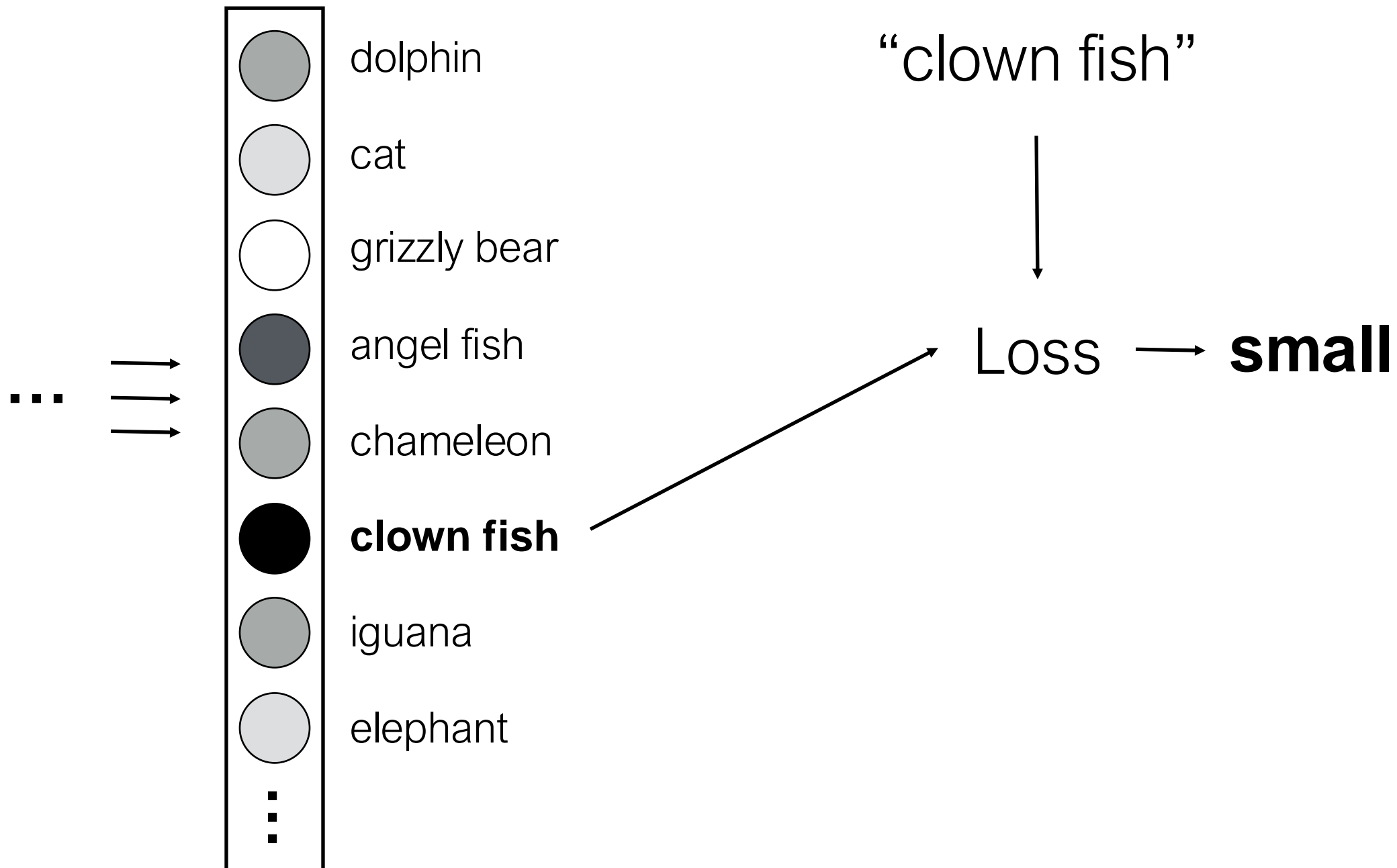
Ground truth label



Loss function

Network output

Ground truth label



Loss function

Network output

Ground truth label



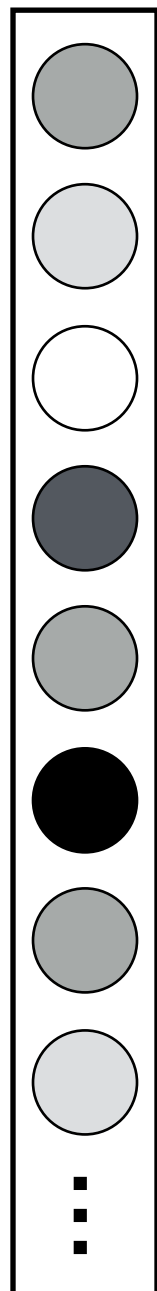
Loss function for classification

Network output

Ground truth label

$\hat{\mathbf{z}}$

\mathbf{z}



dolphin

cat

grizzly bear

angel fish

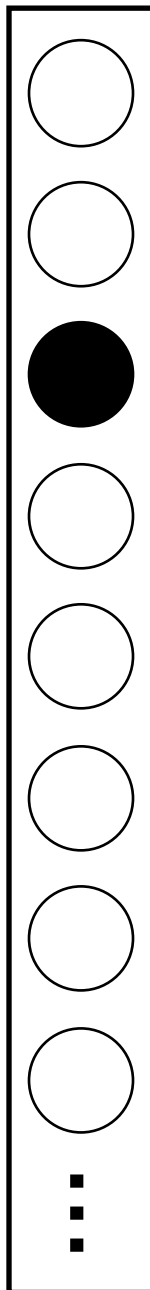
chameleon

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

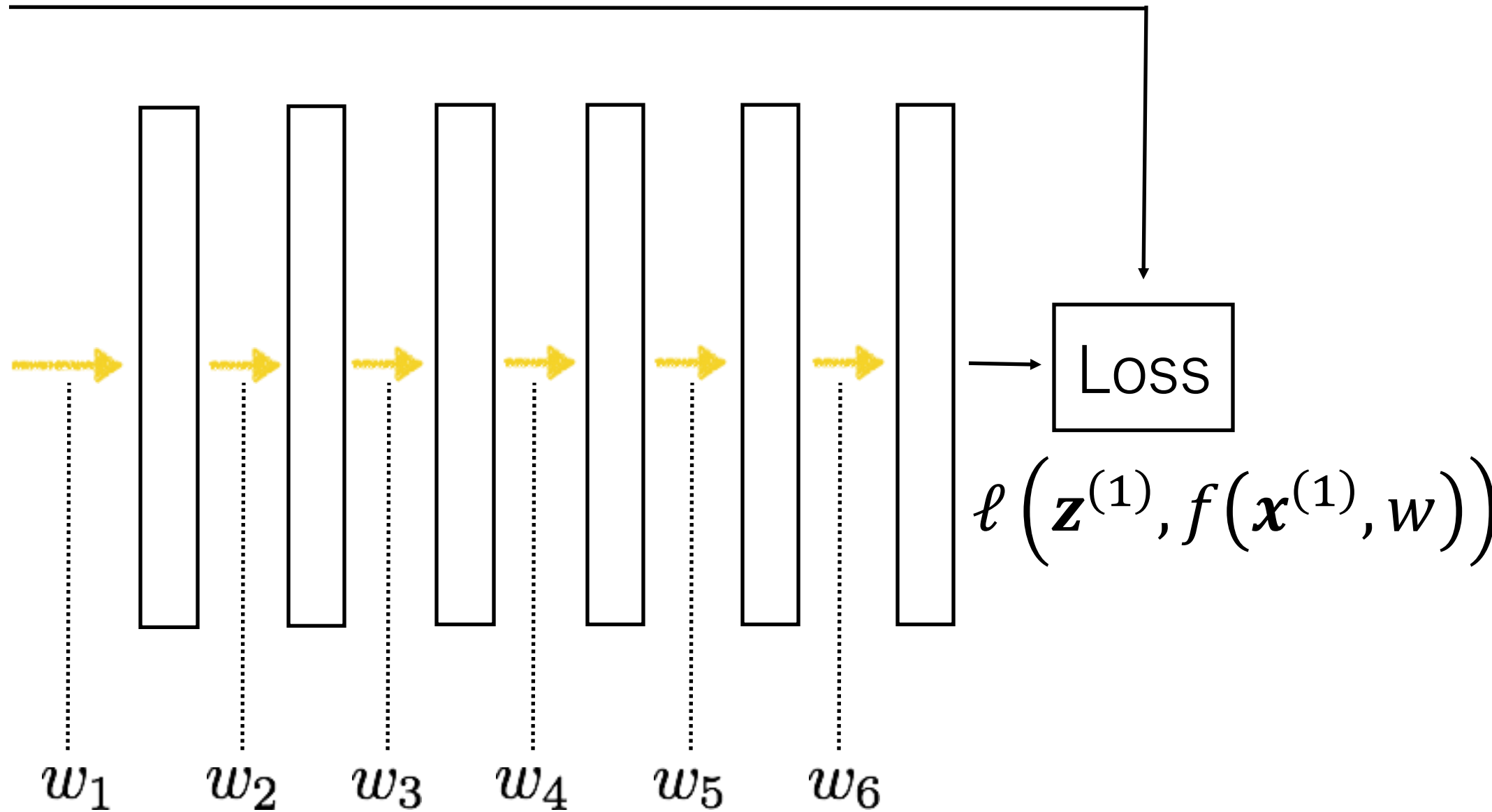
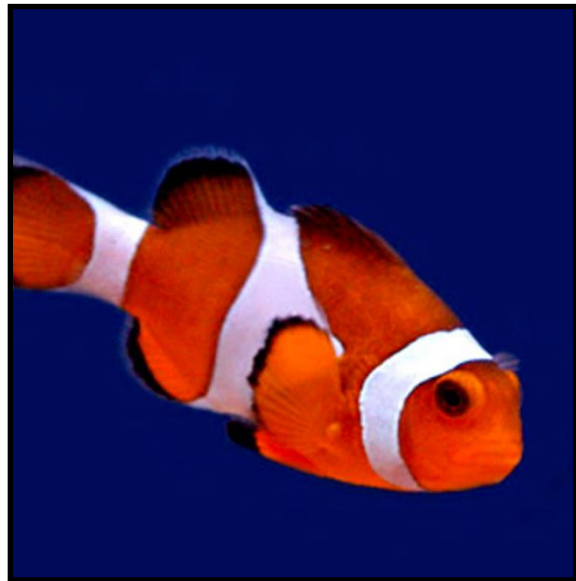
c is the c^{th} class in the output

Learning with deep nets

Learned

$\mathbf{z}^{(1)}$
“clown fish”

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



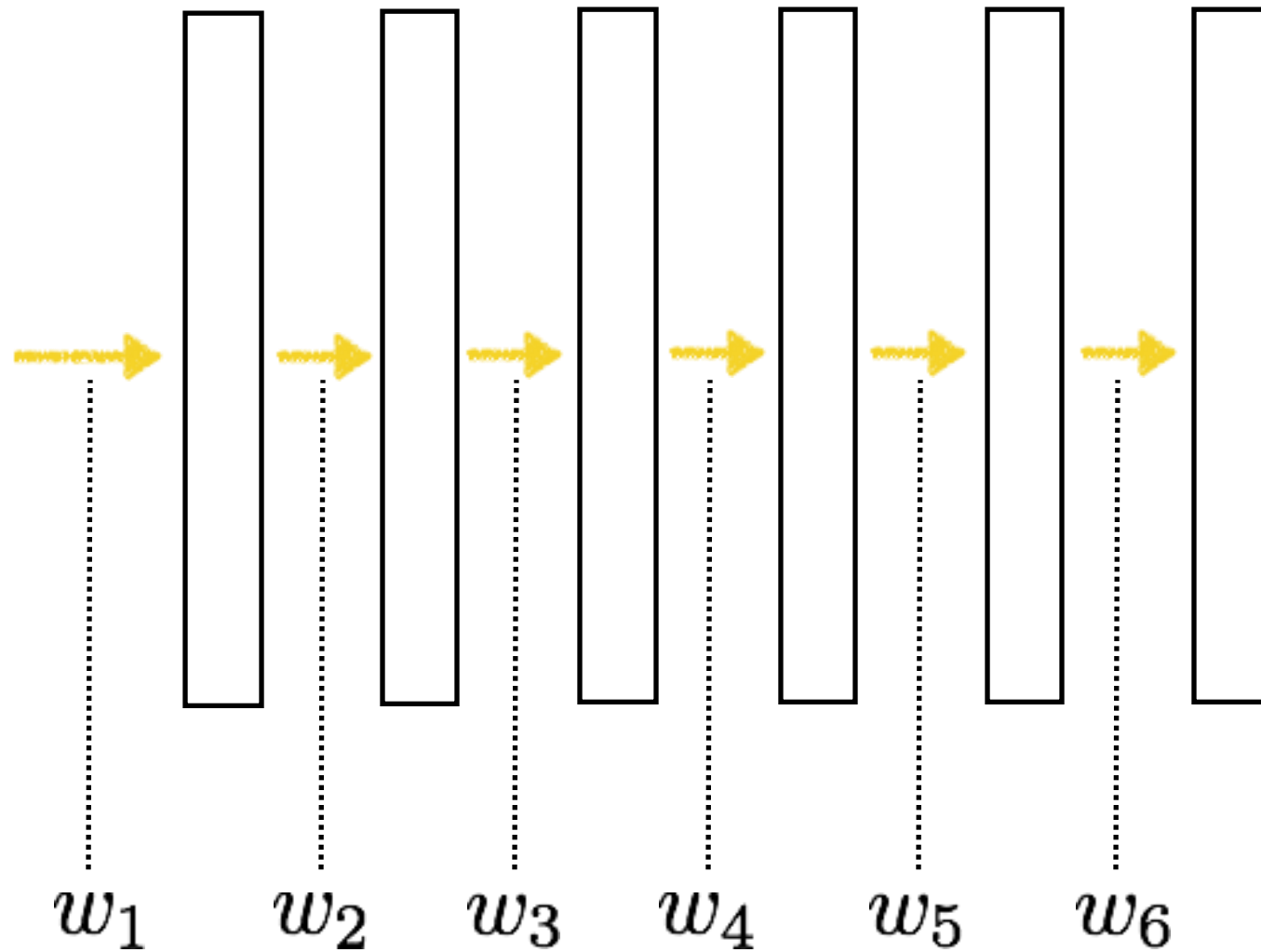
$\mathbf{x}^{(1)}, \mathbf{z}^{(1)}$ is the input and label
of the 1st training image

Learning with deep nets

Learned

$\mathbf{z}^{(2)}$
“grizzly bear”

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

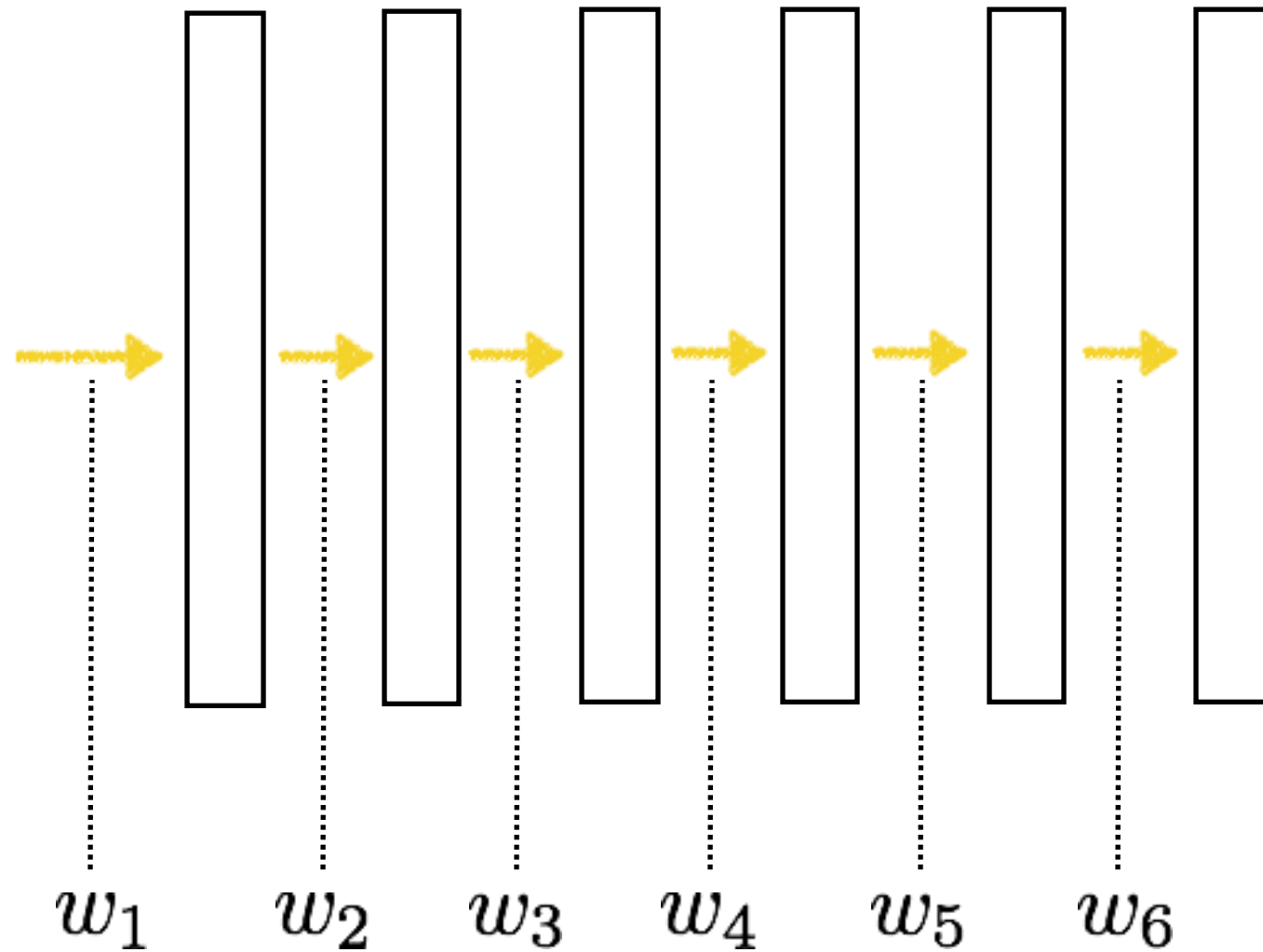
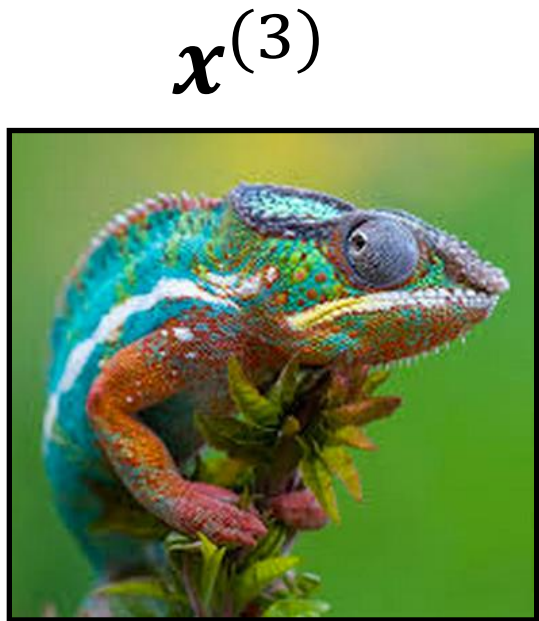


$\mathbf{x}^{(2)}, \mathbf{z}^{(2)}$ is the input and label
of the 2nd training image

Learning with deep nets

Learned

$\mathbf{z}^{(3)}$
“chameleon”



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

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

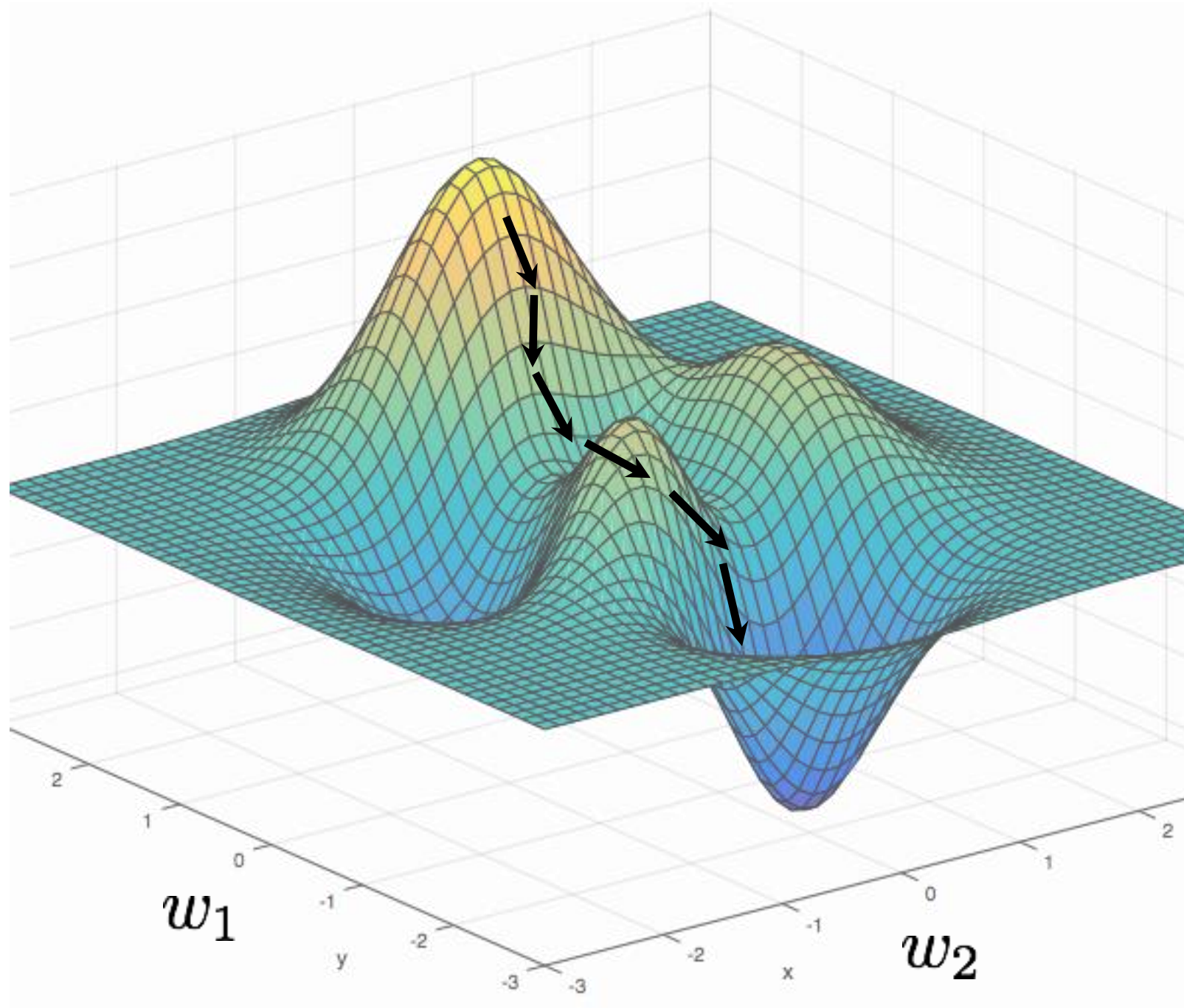
One iteration of gradient descent:

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

learning rate

Gradient descent

$L(\mathbf{w})$



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



mite

container ship

motor scooter

leopard



grille



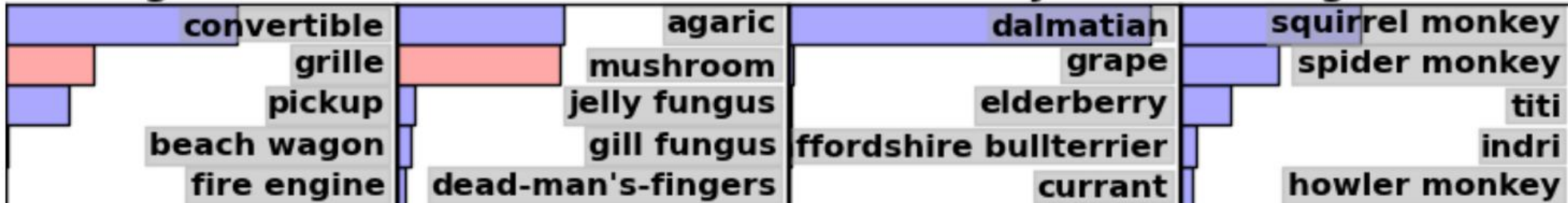
mushroom



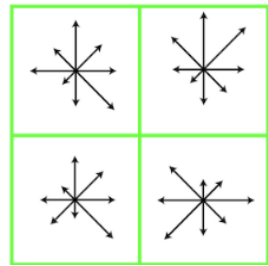
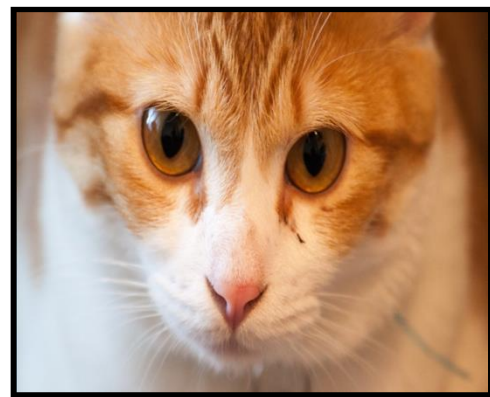
cherry



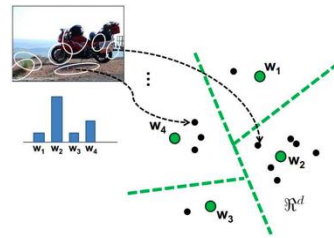
Madagascar cat



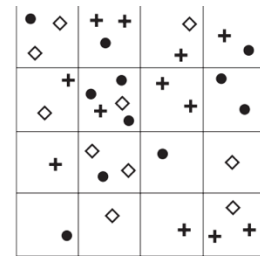
Computer Vision before 2012



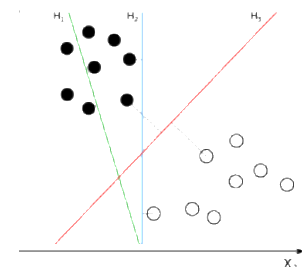
Features



Clustering



Pooling

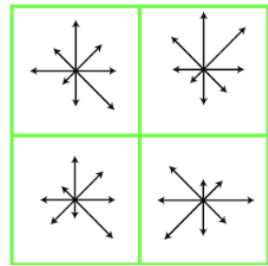
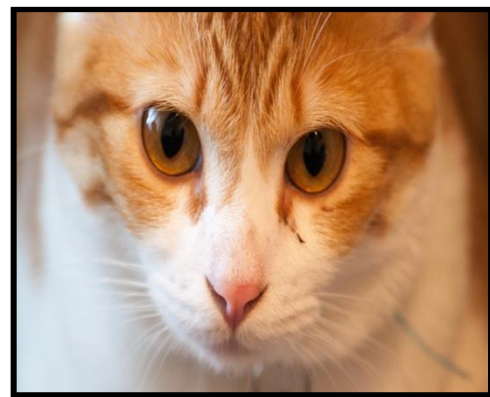


Classification

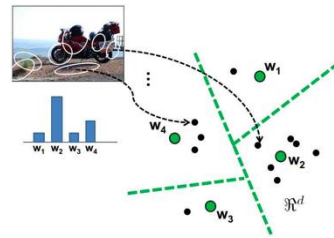


Cat

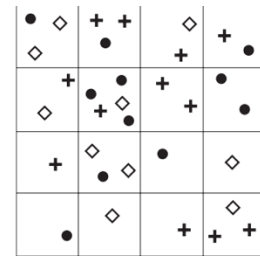
Computer Vision Now



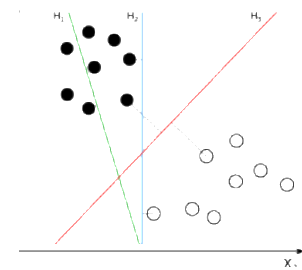
Features



Clustering



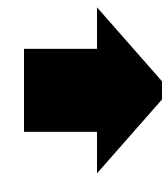
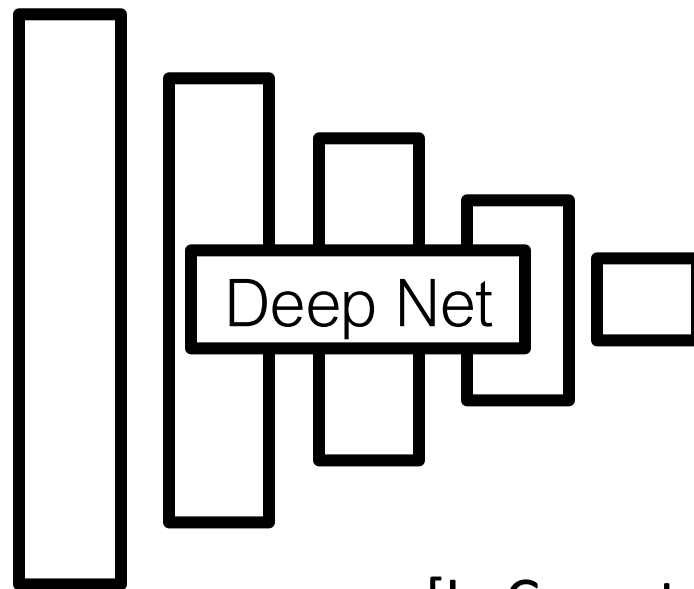
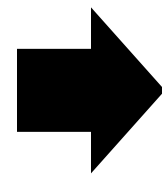
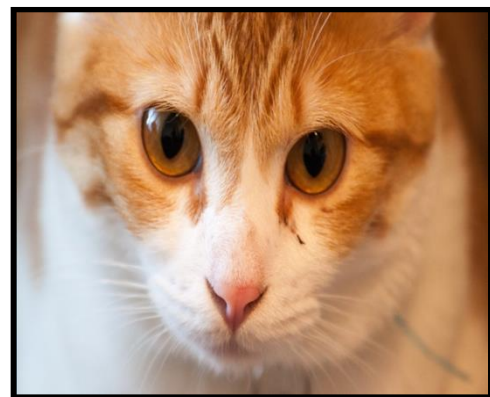
Pooling



Classification



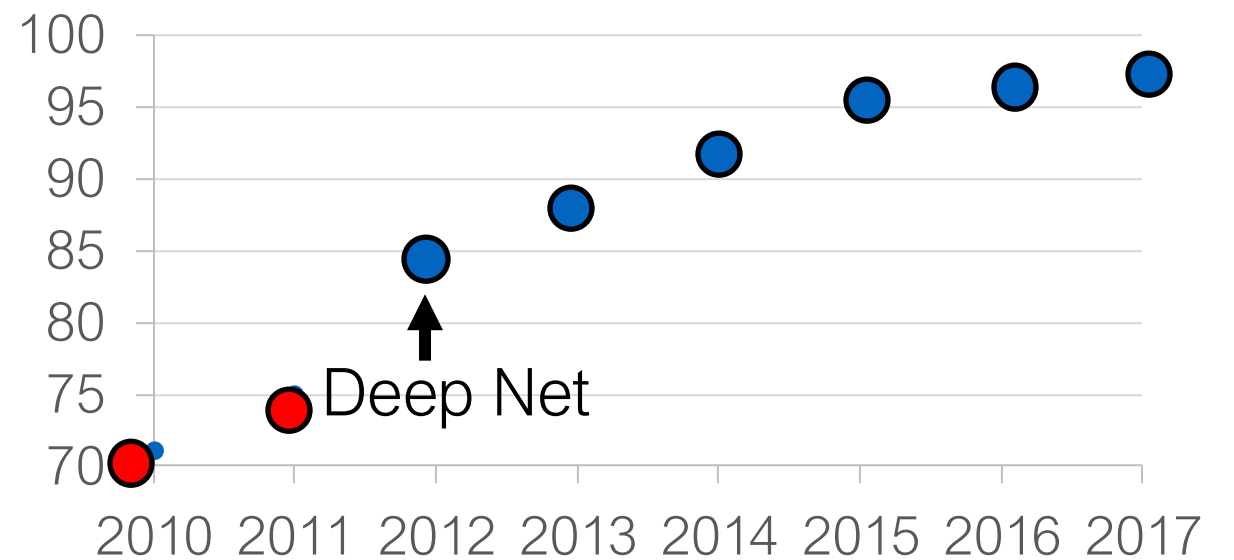
Cat



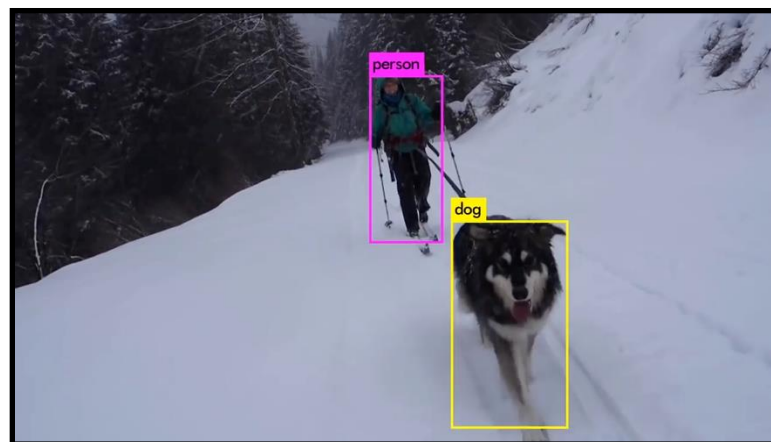
Cat

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

Deep Learning for Computer Vision

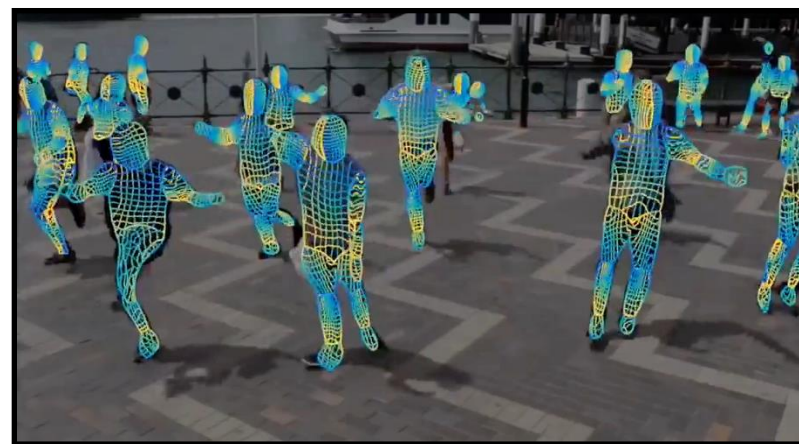


Top 5 **accuracy** on ImageNet benchmark



[Redmon et al., 2018]

Object detection



[Güler et al., 2018]

Human understanding



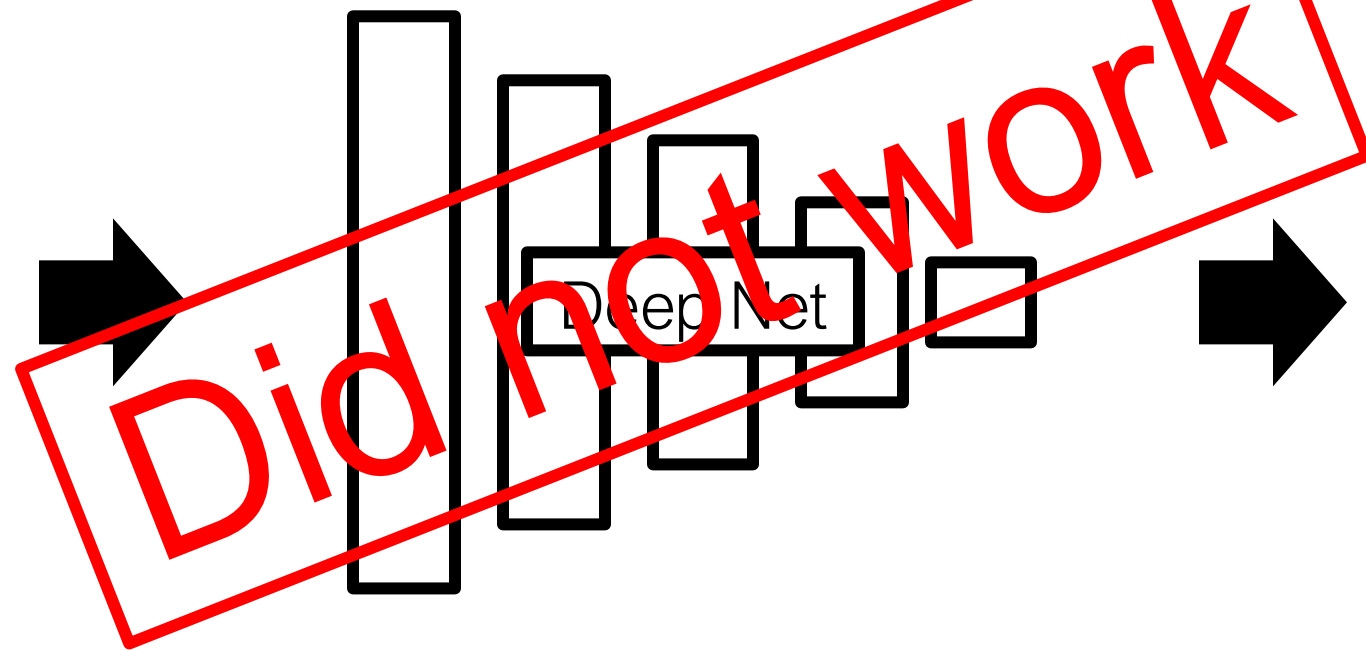
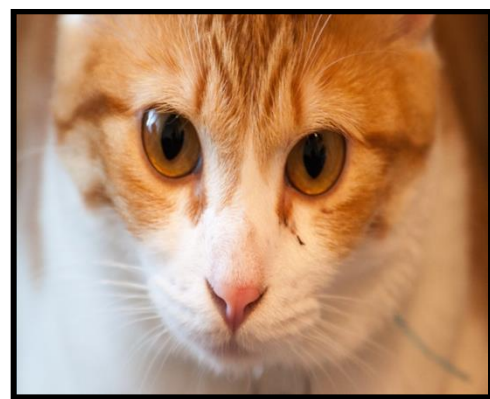
[Zhao et al., 2017]

Autonomous driving

Can Deep Learning Help Graphics?

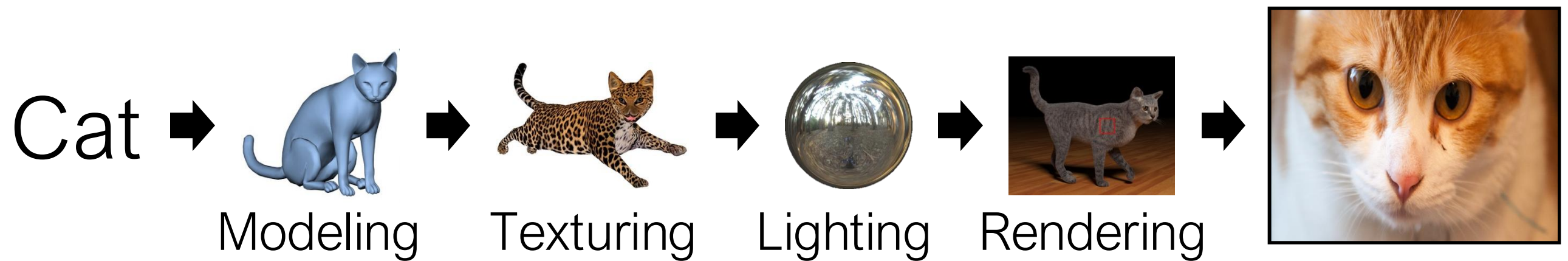


Can Deep Learning Help Graphics?

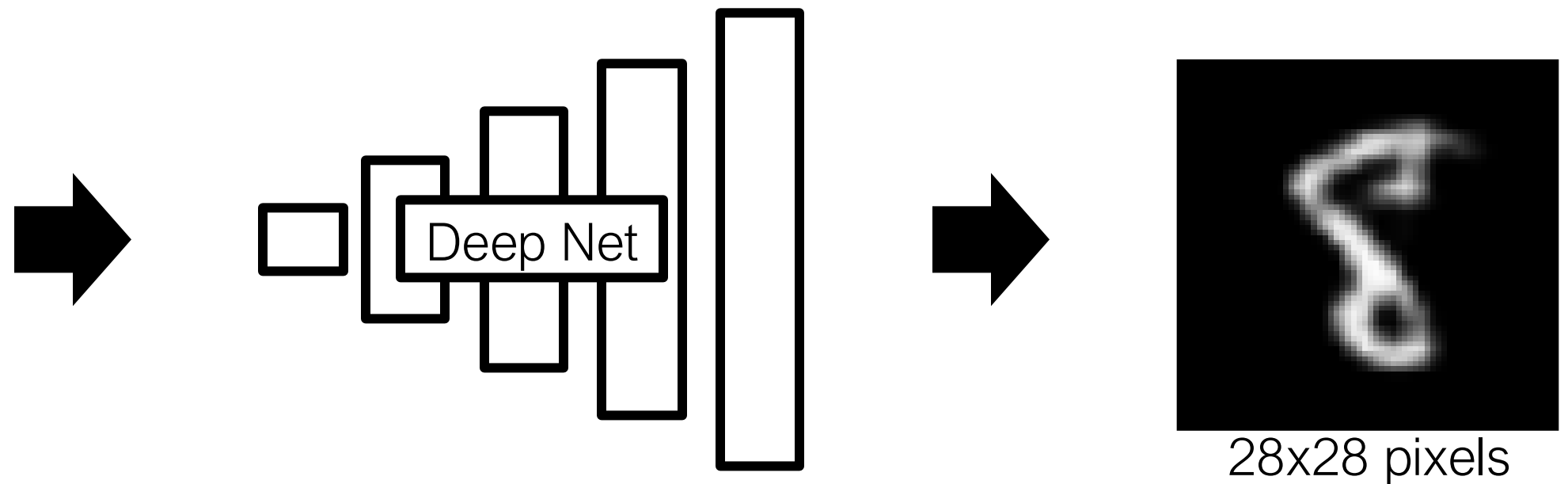


Cat

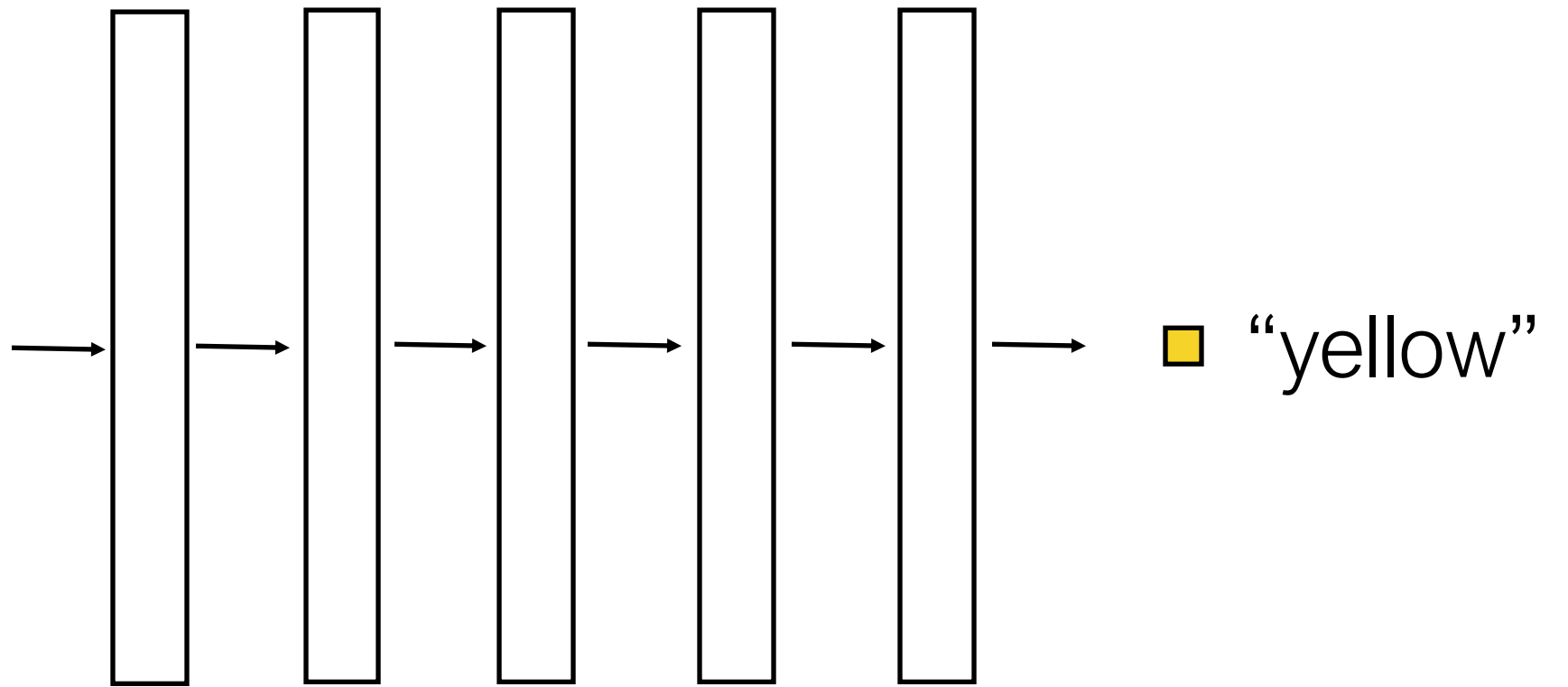
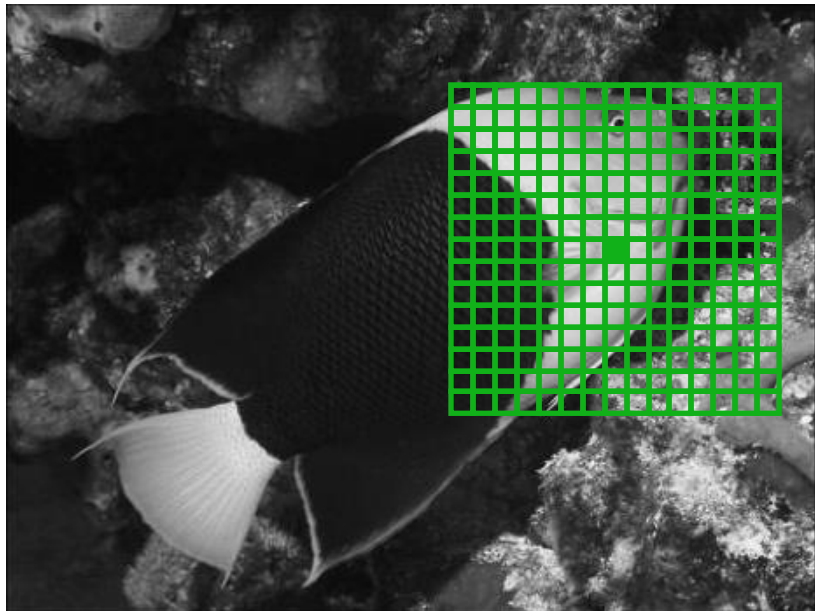
Generating images is hard!



8

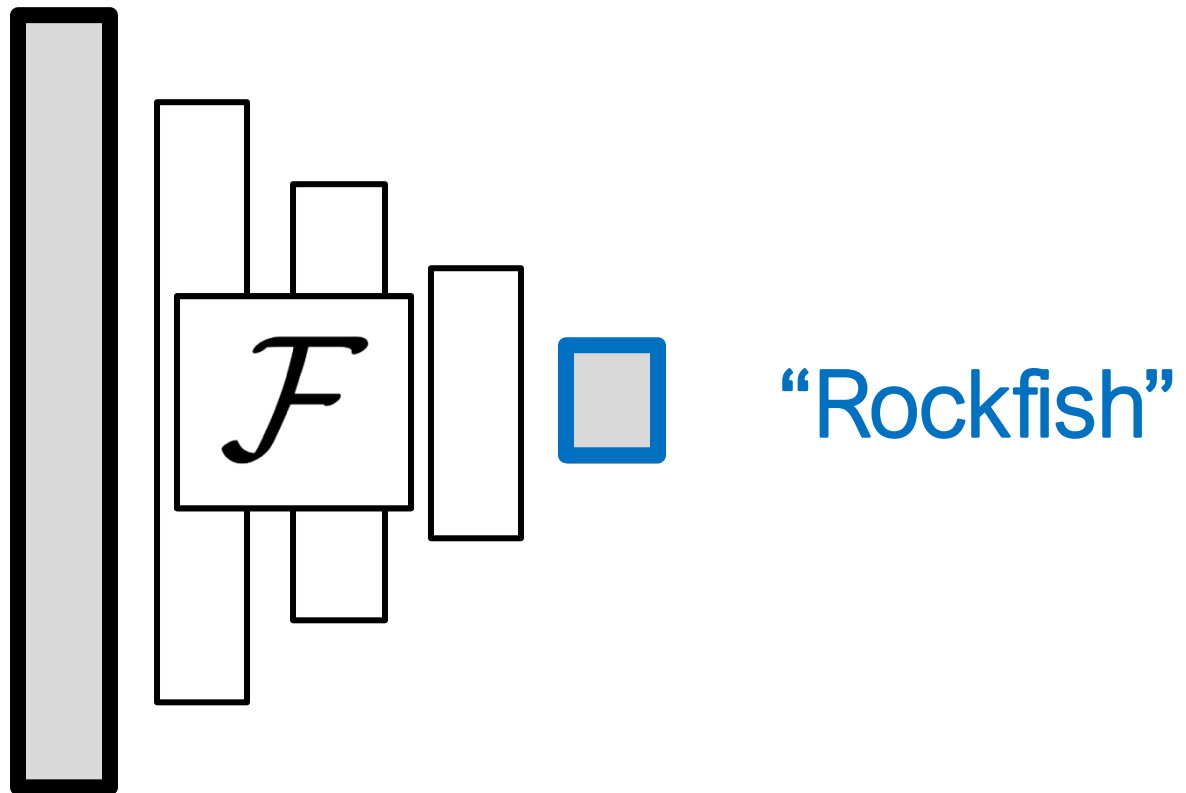


from Classification to Generation

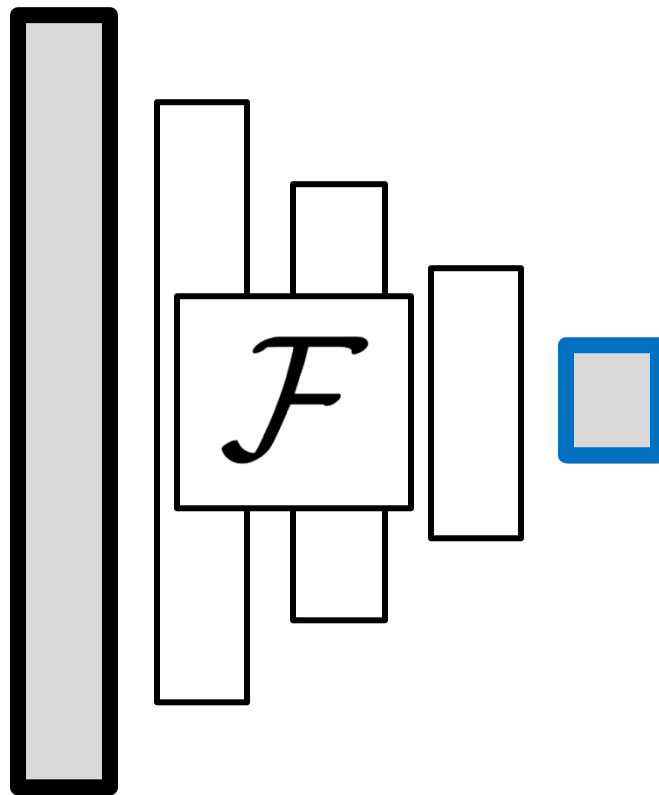


Predicting the color value of an output pixel given a patch ⁴³

Discriminative Deep Networks

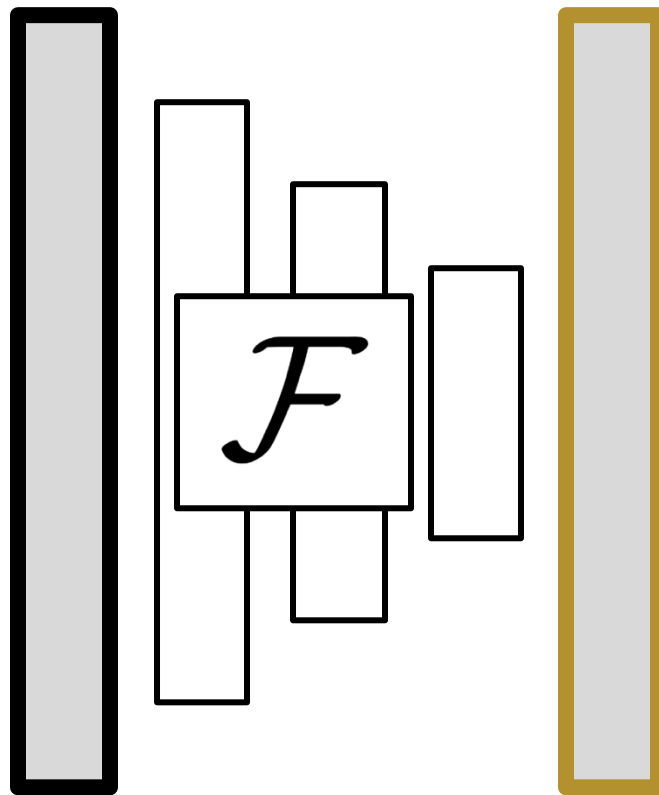


Discriminative Deep Networks



Raw, Unlabeled Pixels

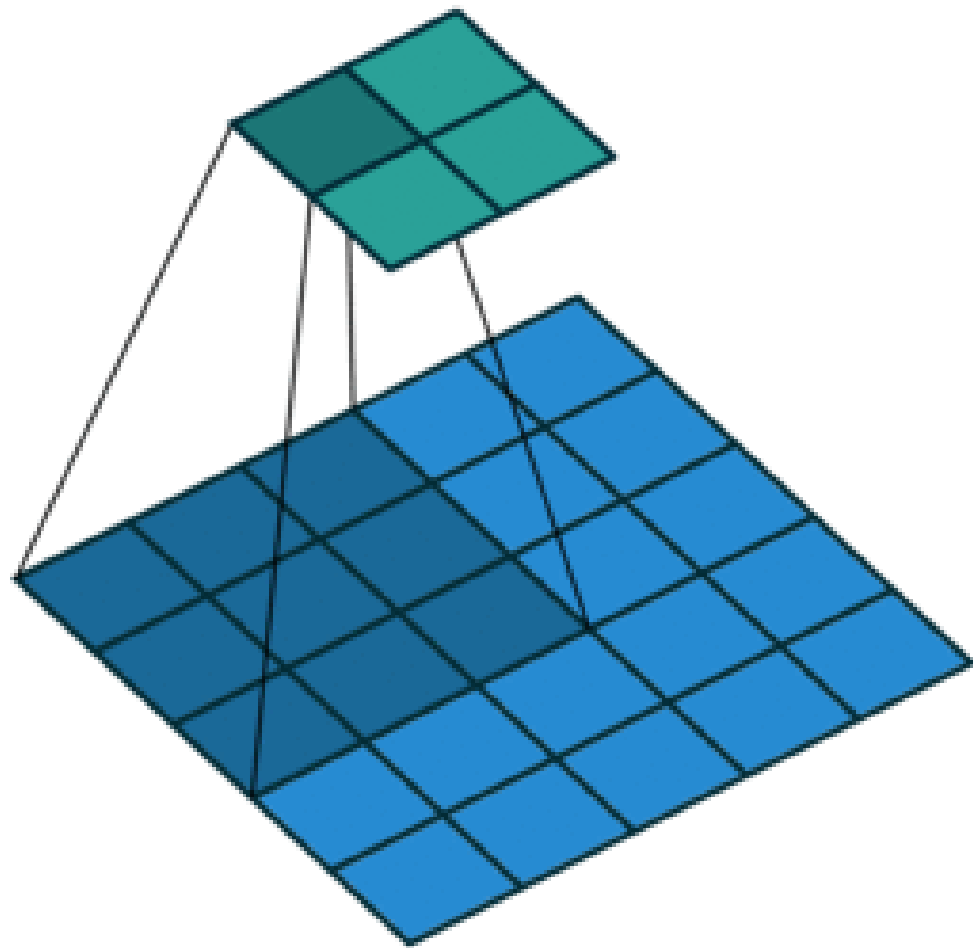
Generative Deep Networks



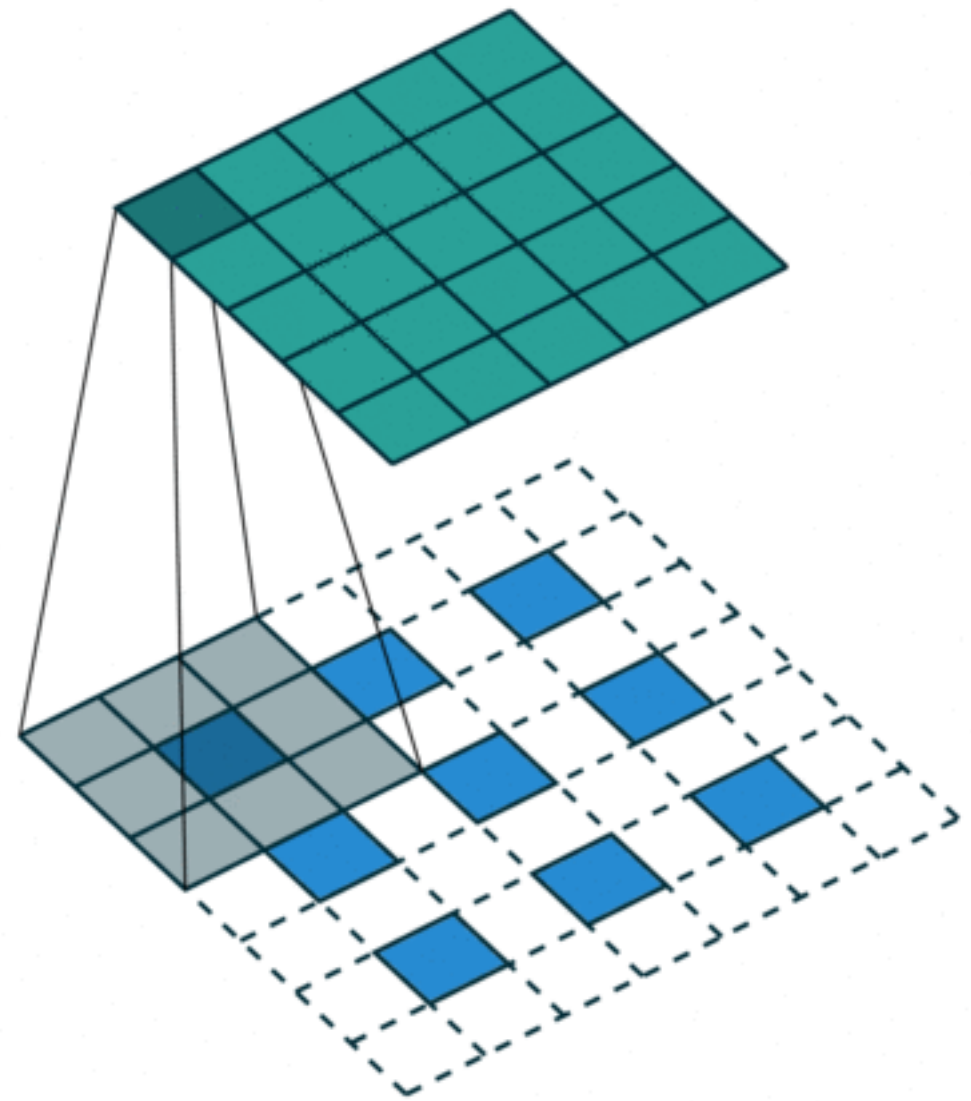
Raw, Unlabeled Pixels

Better Architectures

Fractionally-strided Convolution

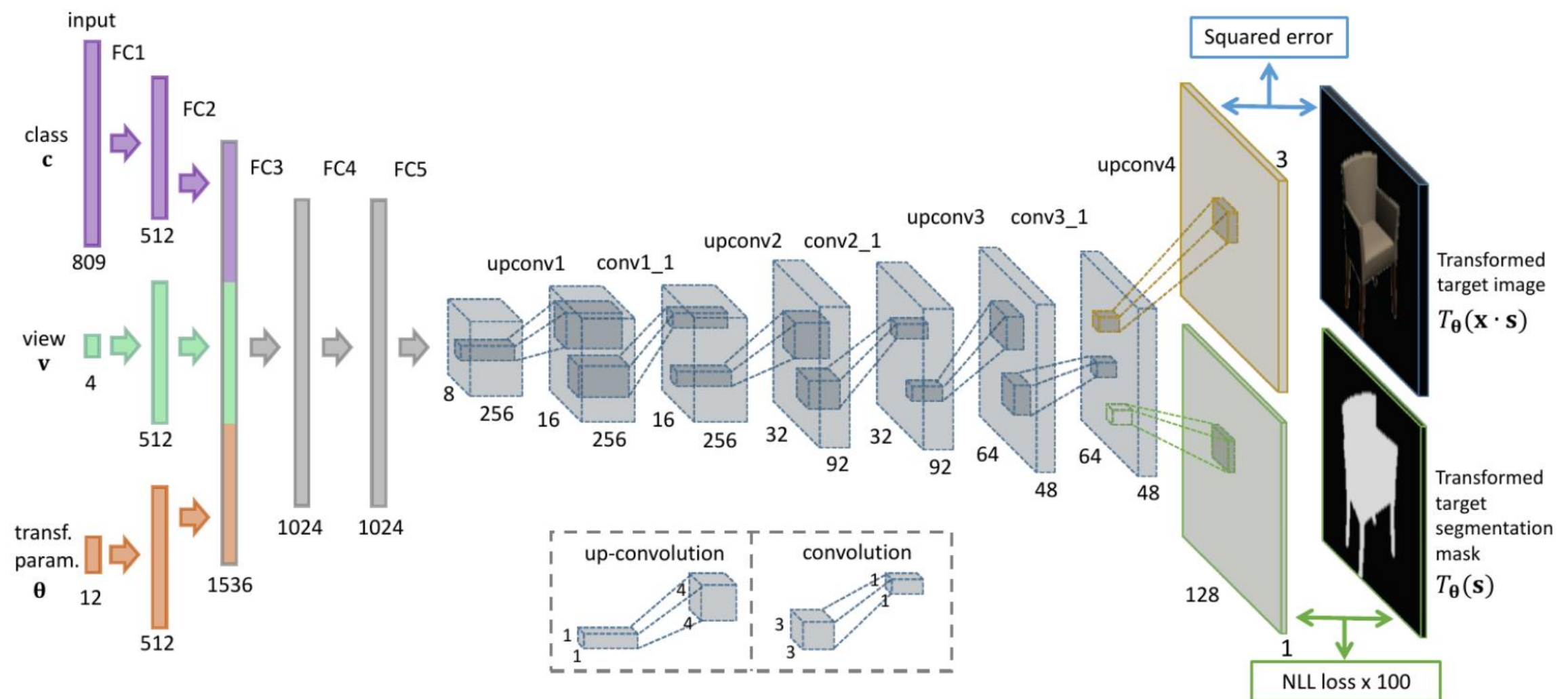


Regular conv (no padding)



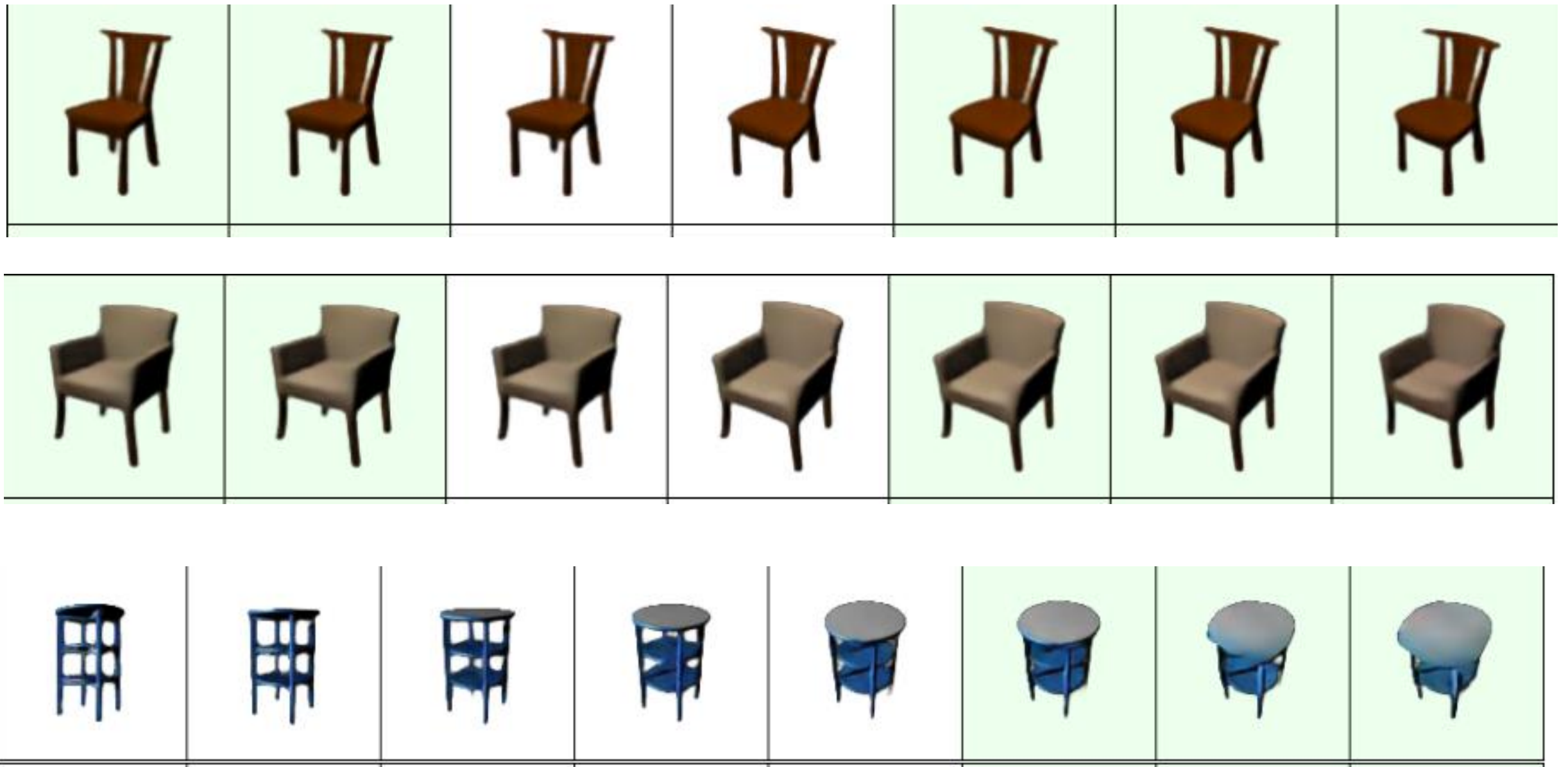
Fractionally-strided conv

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



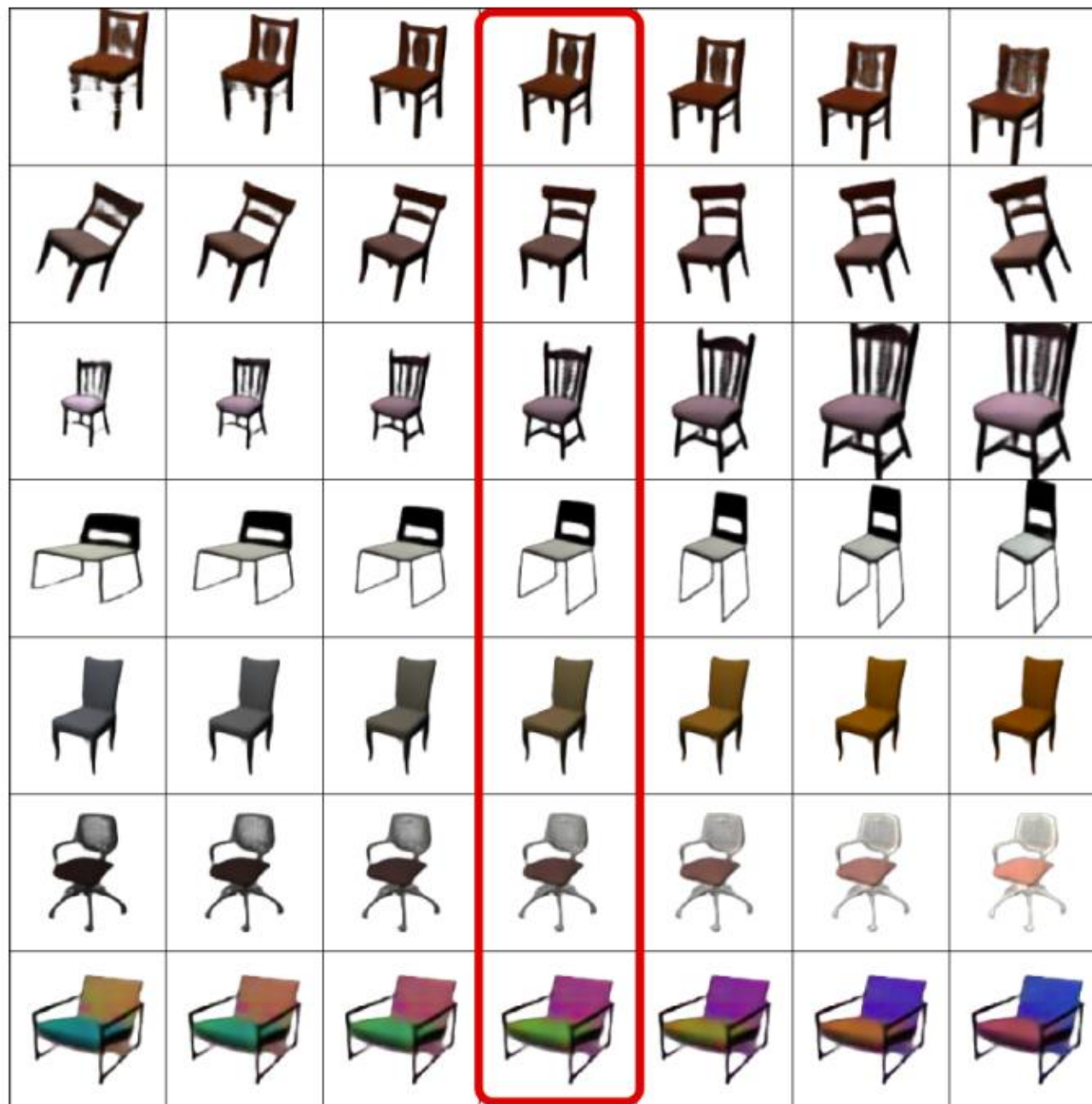
Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks
PAMI 2017 (CVPR 2015)

With Varying Viewpoints

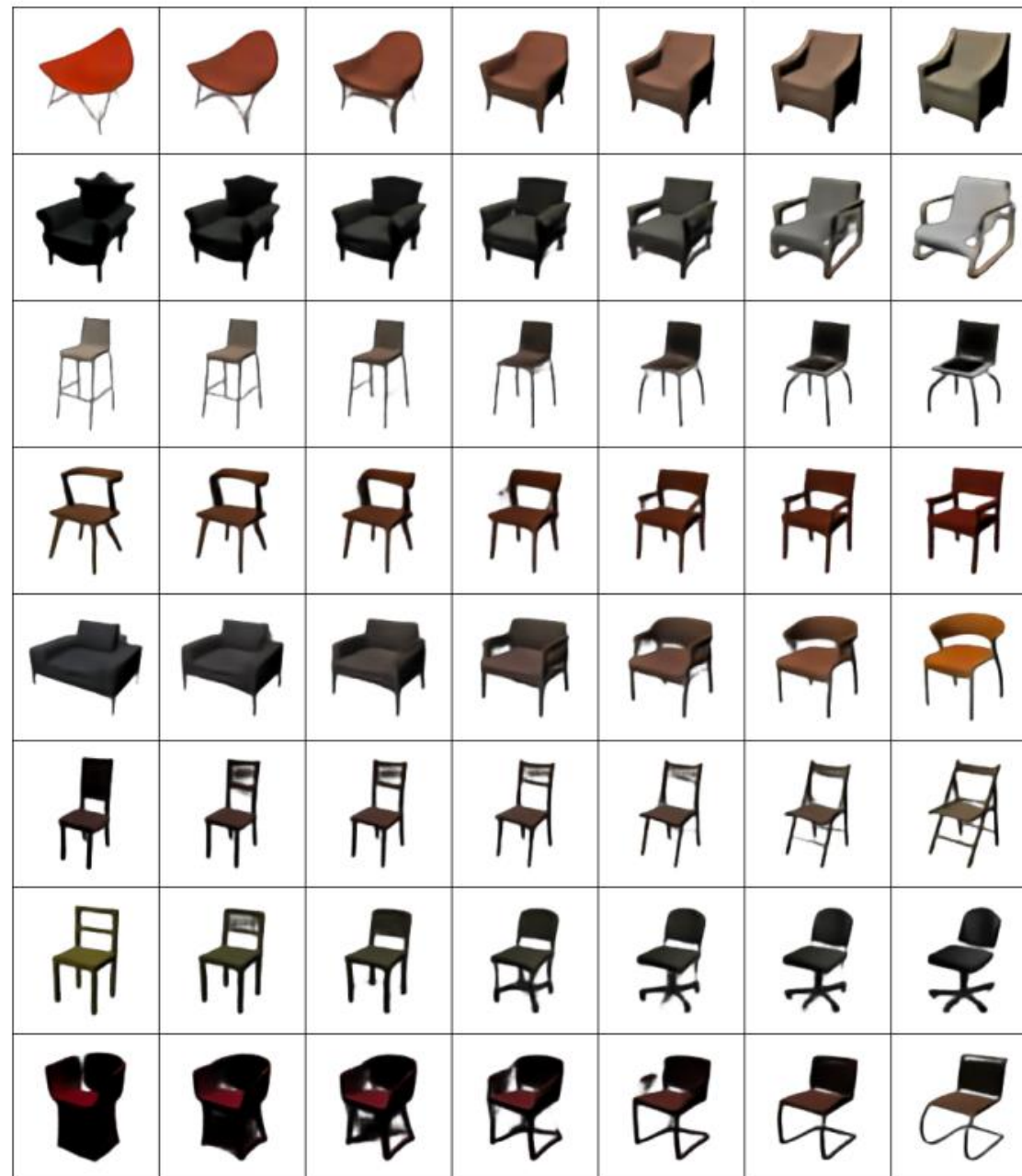


Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks
PAMI 2017 (CVPR 2015)

With Varying Transformation Parameters



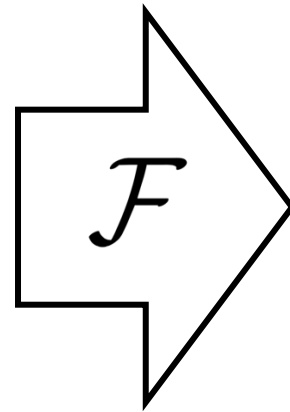
Interpolation between Two Chairs



Better Loss Functions



Ansel Adams. *Yosemite Valley Bridge.*

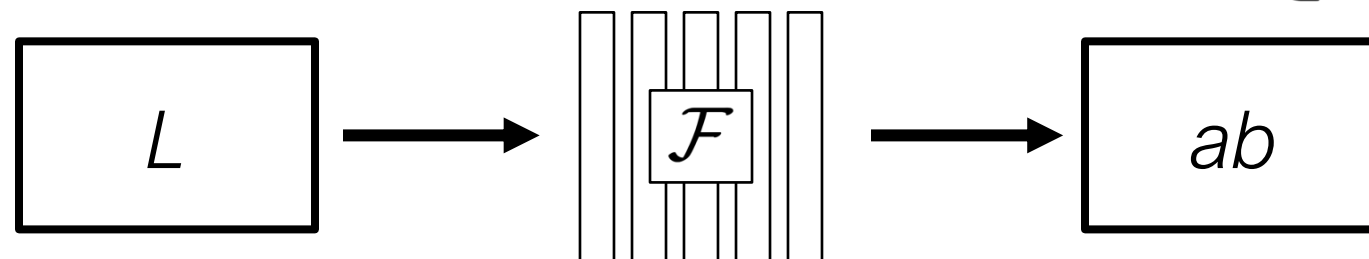


Grayscale image: L channel

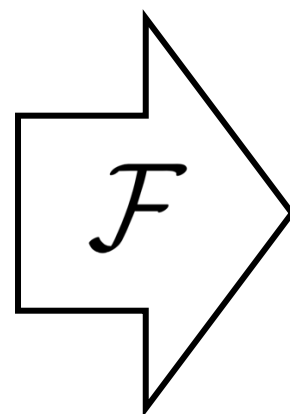
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: ab channels

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



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

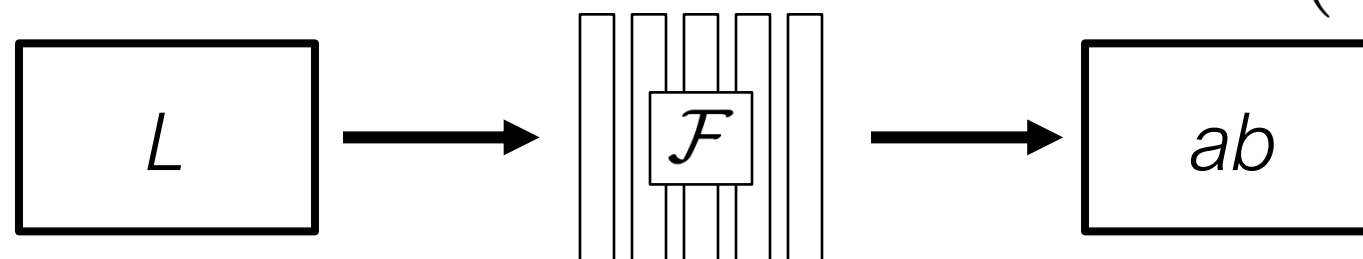


Grayscale image: L channel

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

Concatenate (L, ab) channels

$$(\mathbf{X}, \hat{\mathbf{Y}})^{56}$$



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

Simple L2 regression doesn't work ☹️

Input



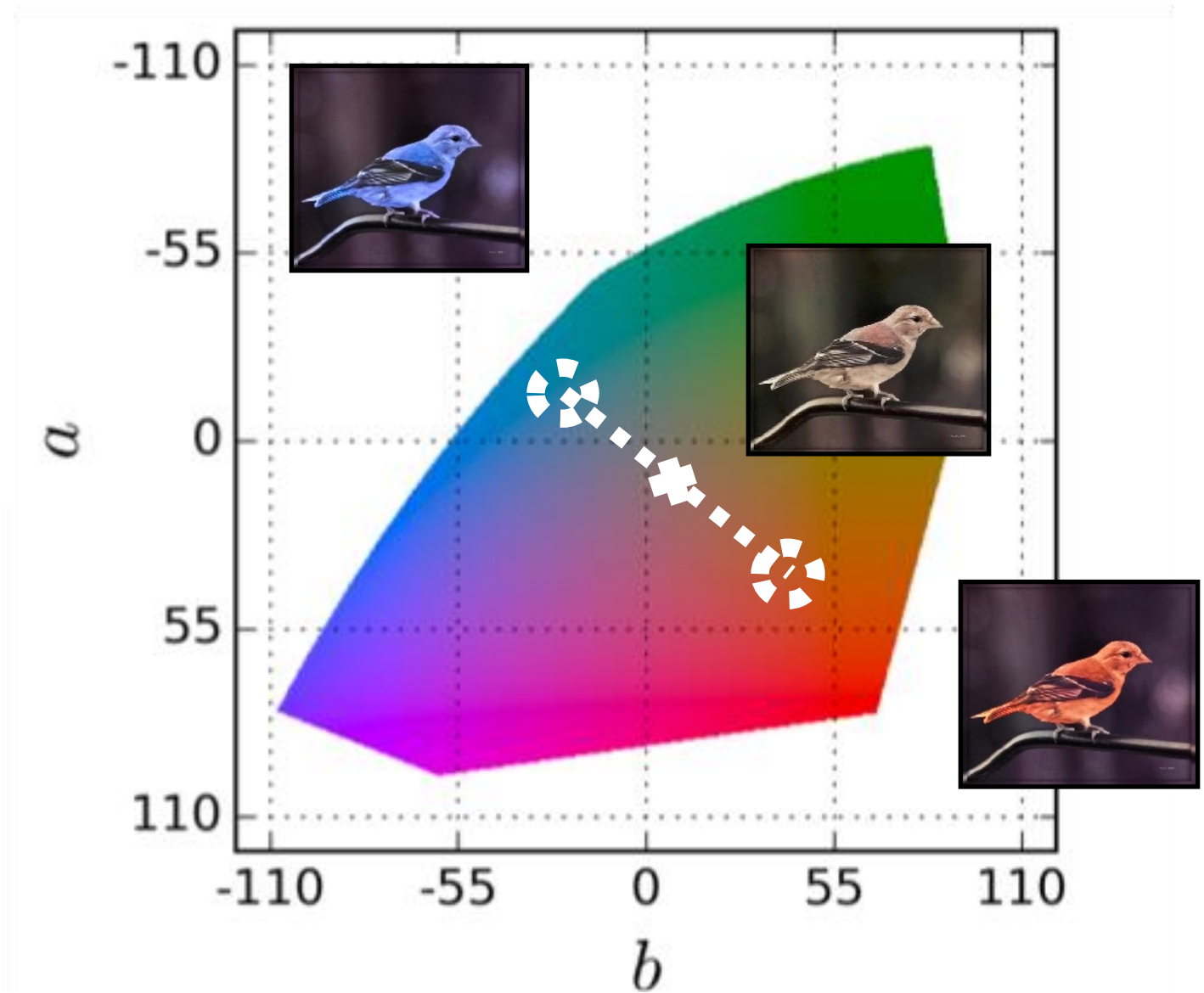
Output



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

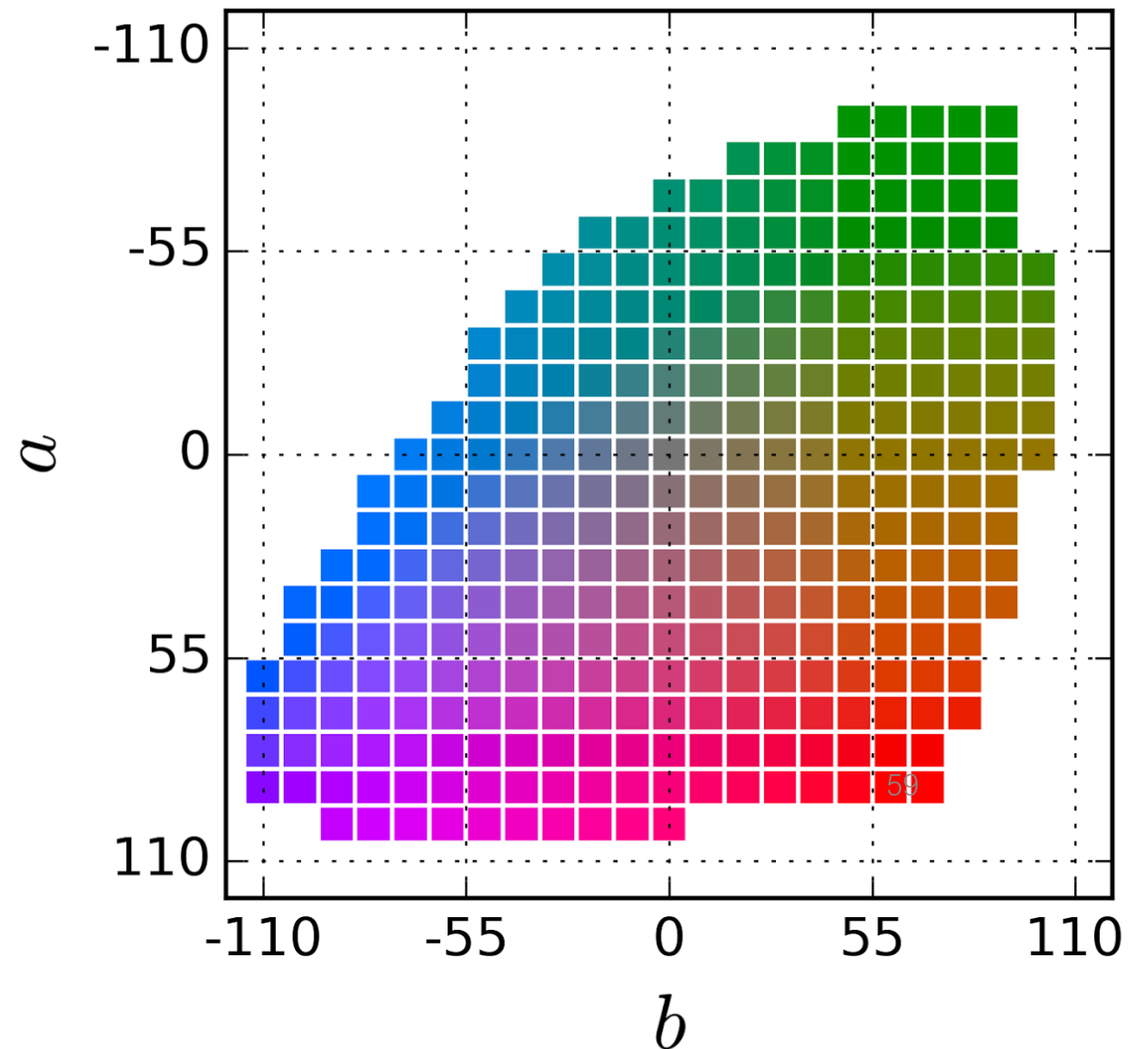
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

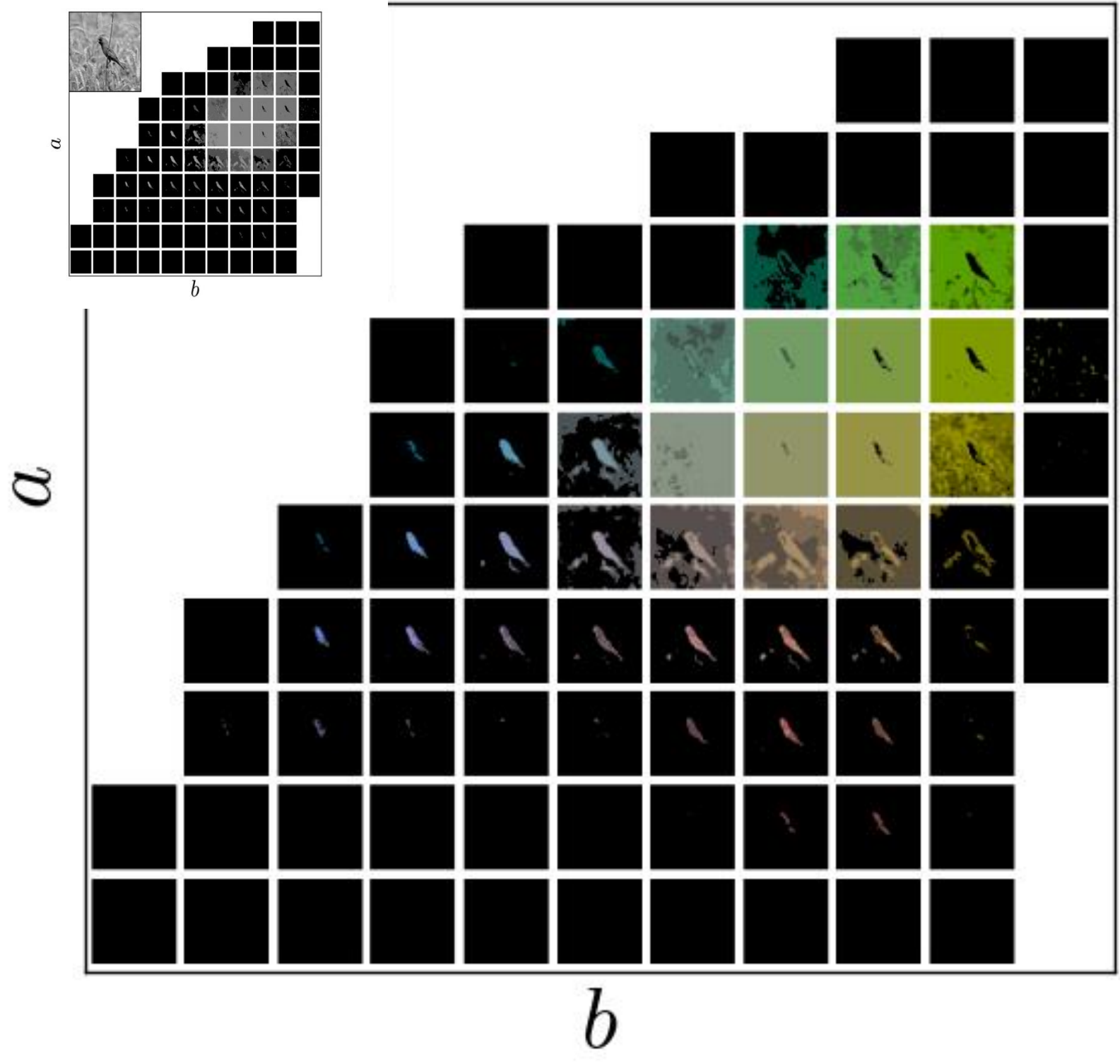
- Use per-pixel multinomial classification

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

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 2025

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