



Western
Science

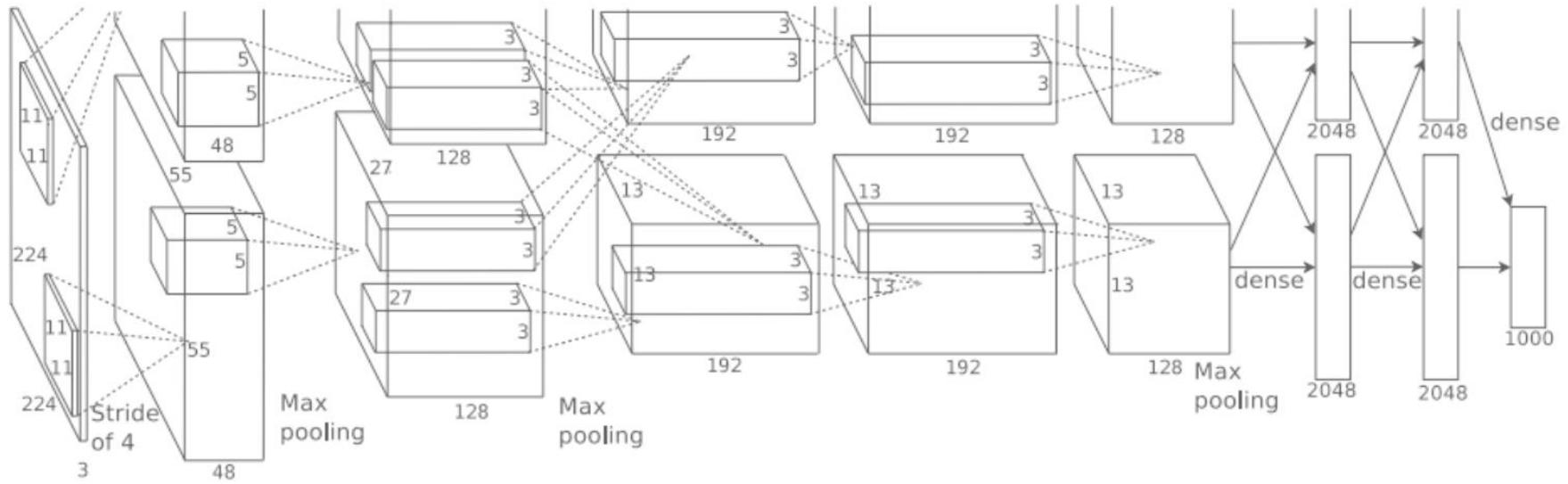
Artificial Intelligence II

Part 2: Lecture 5

Yalda Mohsenzadeh

Winter 2023

Slides are adapted from Antonio Torralba and Philip Isola (MIT)

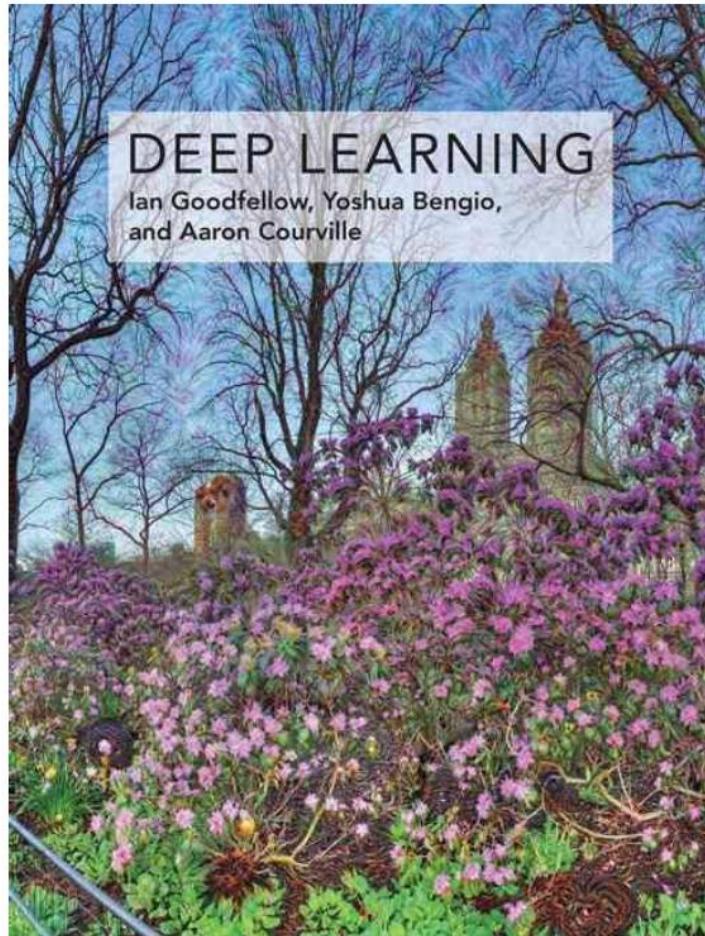


Computer Vision

Neural Networks

Neural Networks

- Brief history
- Basic formulation (hierarchical processing)
- Optimization via gradient descent
- Layer types (Linear, Pointwise, non-linearity)
- Linear classification with perceptron
- Tensor flow
- Regularizers
- Normalization



<http://www.deeplearningbook.org/>

By Ian Goodfellow, Yoshua Bengio and Aaron Courville

November 2016

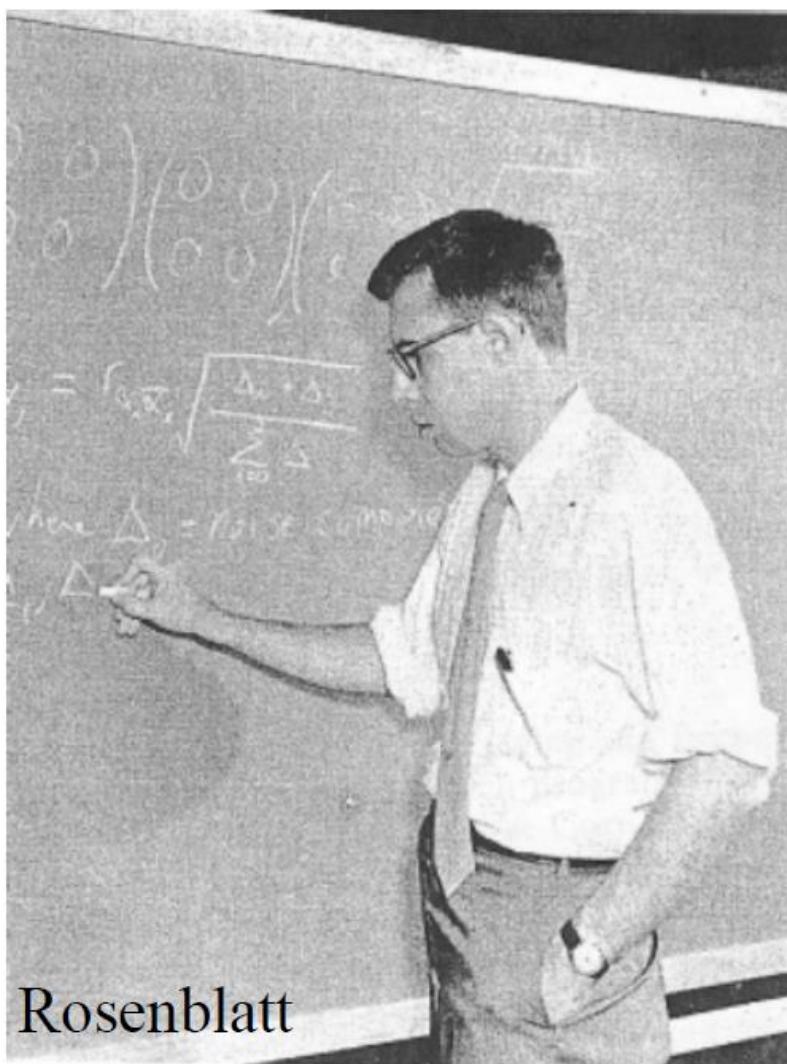
Deep Learning

- Modeling the visual world is incredibly complicated.
We need high capacity models.
- In the past, we didn't have enough data to fit these
models. But now we do!
- We want a class of **high capacity models** that are
easy to optimize.
- **Deep Neural Networks!**

A brief history of neural networks

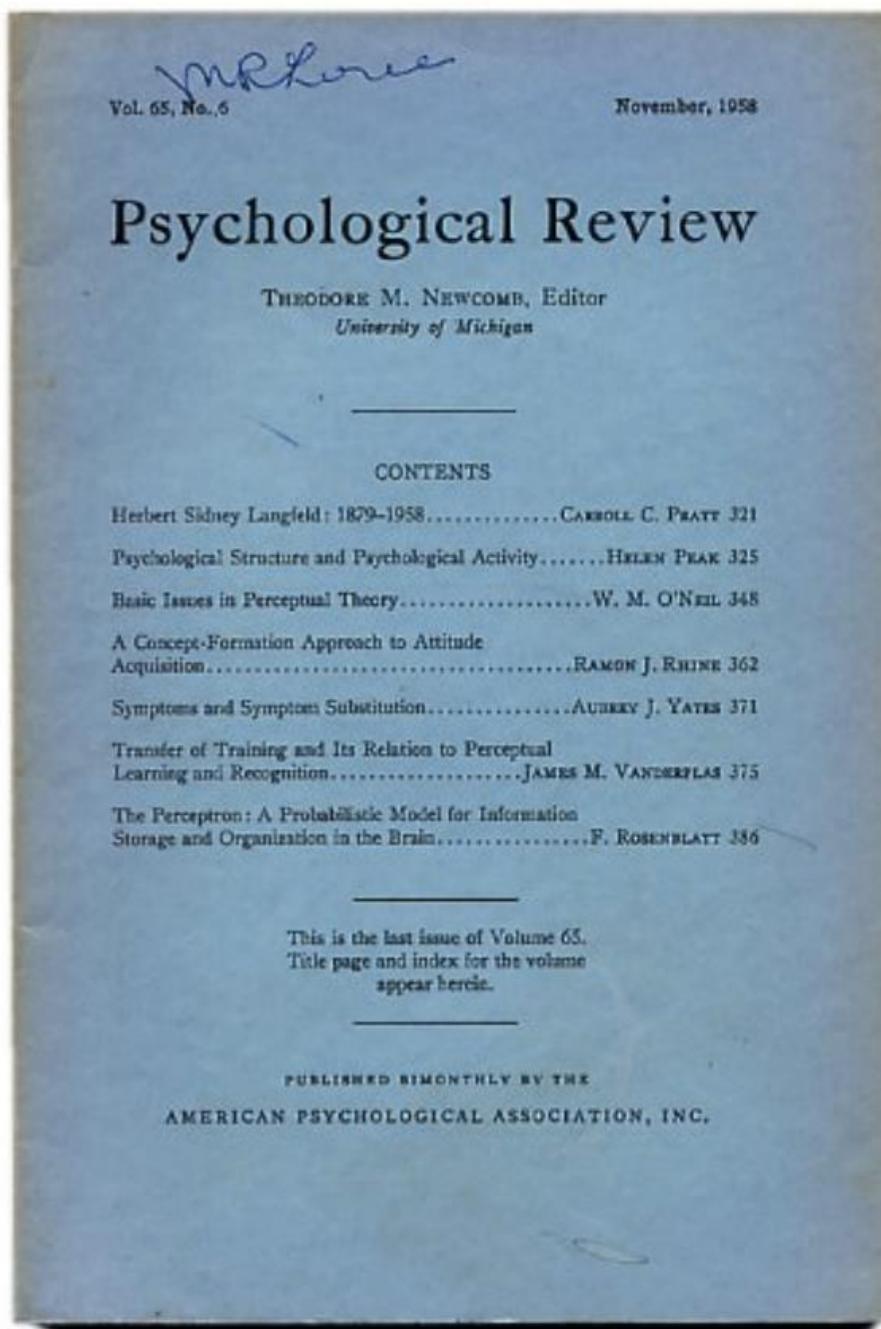


Perceptrons, 1958



Rosenblatt

http://www.ecse.rpi.edu/homepages/nagy/PDF_chrono_2011_Nagy_Pace_FR.pdf. Photo by George Nagy



Vol. 65, No. 6

November, 1958

Psychological Review

THEODORE M. NEWCOMB, Editor
University of Michigan

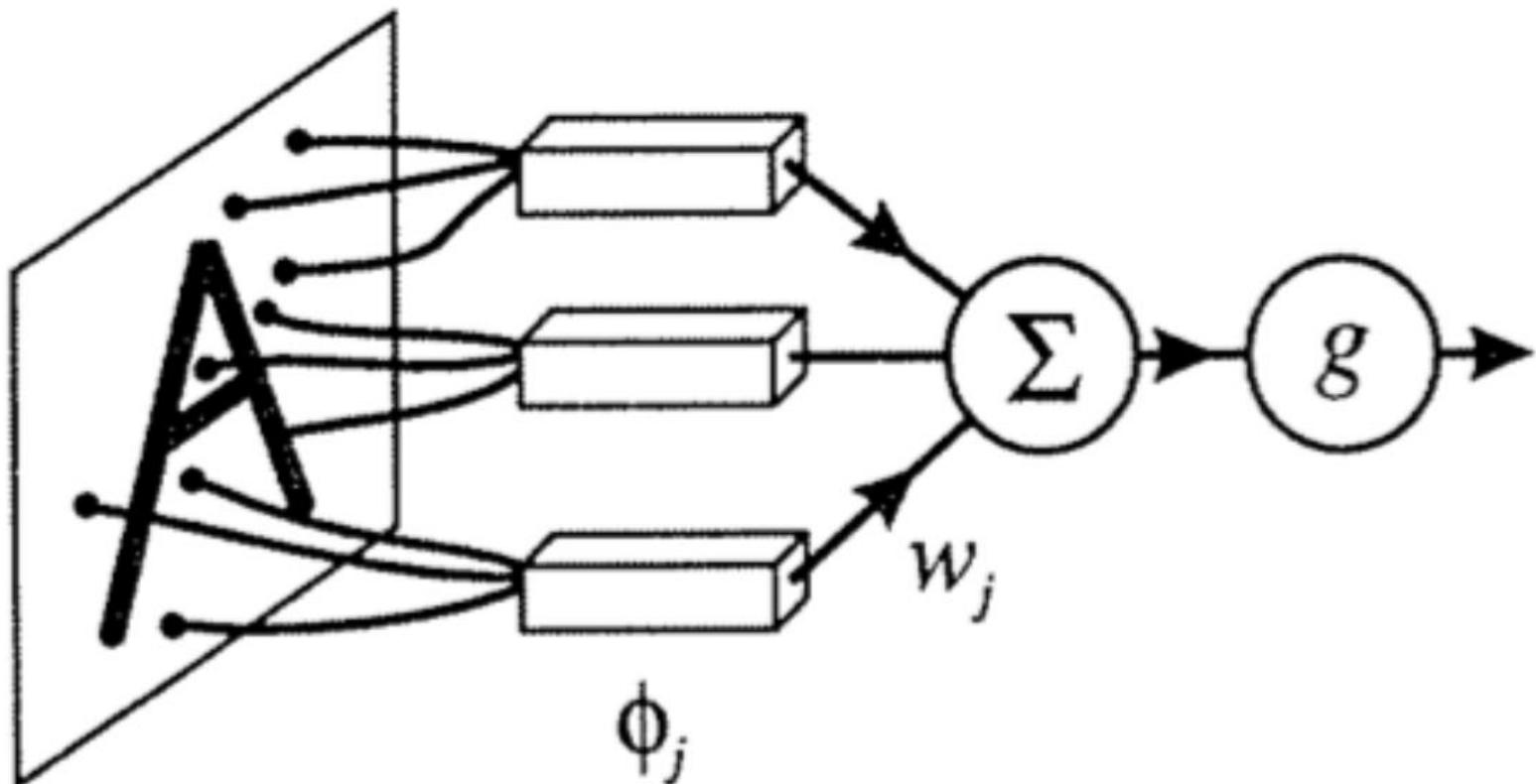
CONTENTS

Herbert Sidney Langfeld: 1879-1958.....	CARROLL C. PRATT 321
Psychological Structure and Psychological Activity.....	HELEN PEAK 325
Basic Issues in Perceptual Theory.....	W. M. O'NEIL 348
A Concept-Formation Approach to Attitude Acquisition.....	RAMON J. RHINE 362
Symptoms and Symptom Substitution.....	AUBREY J. YATES 371
Transfer of Training and Its Relation to Perceptual Learning and Recognition.....	JAMES M. VANDERPLAS 375
The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain.....	F. ROSENBLATT 386

This is the last issue of Volume 65.
Title page and index for the volume
appear herein.

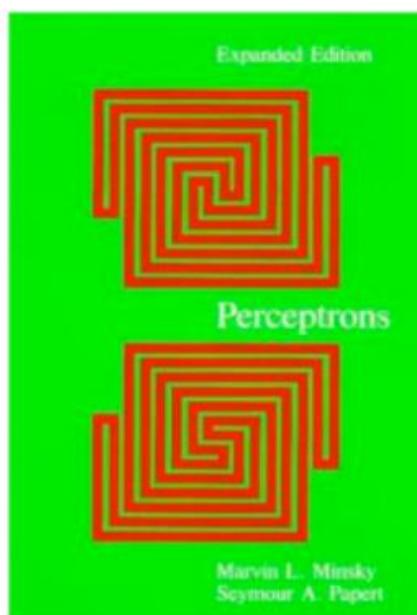
PUBLISHED BIMONTHLY BY THE
AMERICAN PSYCHOLOGICAL ASSOCIATION, INC.

Perceptron, 1958





Minsky and Papert, Perceptrons, 1972



Perceptrons, expanded edition

An Introduction to Computational Geometry

By Marvin Minsky and Seymour A. Papert

Overview

Perceptrons - the first systematic study of parallelism in computation - has remained a classical work on threshold automata networks for nearly two decades. It marked a historical turn in artificial intelligence, and it is required reading for anyone who wants to understand the connectionist counterrevolution that is going on today.

Artificial-intelligence research, which for a time concentrated on the programming of von Neumann computers, is swinging back to the idea that intelligence might emerge from the activity of networks of neuronlike entities. Minsky and Papert's book was the first example of a mathematical analysis carried far enough to show the exact limitations of a class of computing machines that could seriously be considered as models of the brain. Now the new developments in mathematical tools, the recent interest of physicists in the theory of disordered matter, the new insights into and psychological models of how the brain works, and the evolution of fast computers that can simulate networks of automata have given *Perceptrons* new importance.

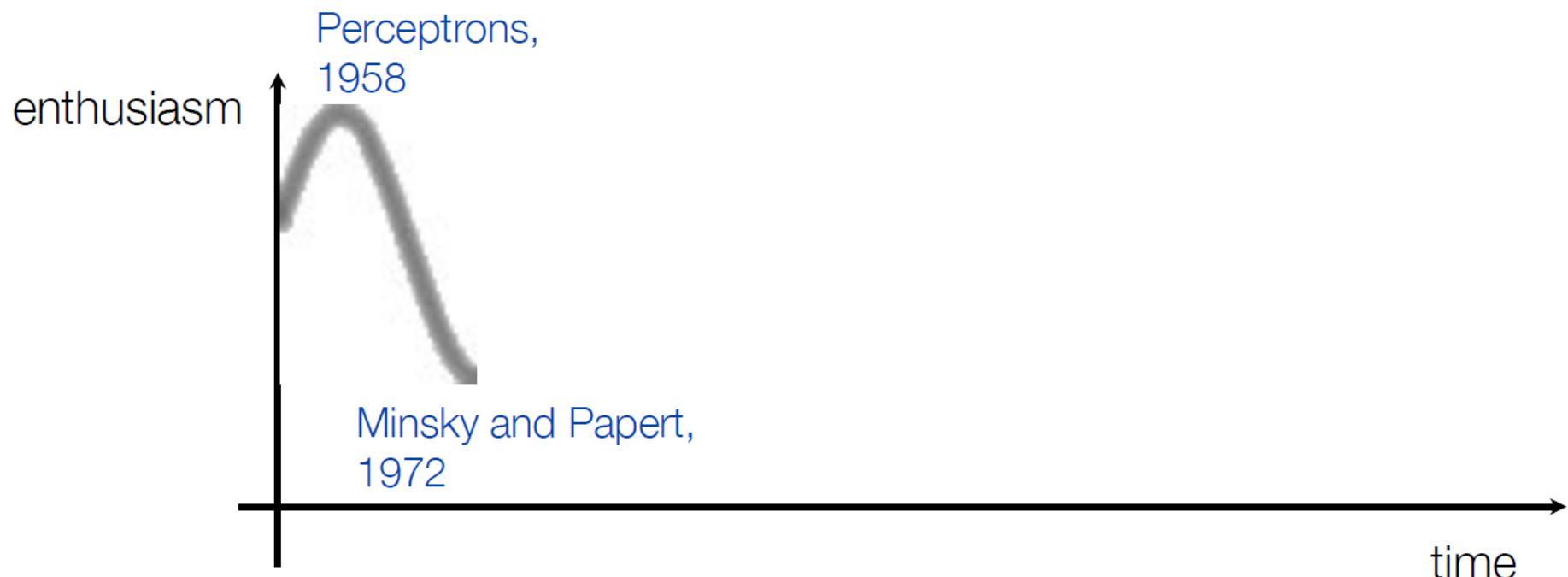
Witnessing the swing of the intellectual pendulum, Minsky and Papert have added a new chapter in which they discuss the current state of parallel computers, review developments since the appearance of the 1972 edition, and identify new research directions related to connectionism. They note a central theoretical challenge facing connectionism: the challenge to reach a deeper understanding of how "objects" or "agents" with individuality can emerge in a network. Progress in this area would link connectionism with what the authors have called "society theories of mind."



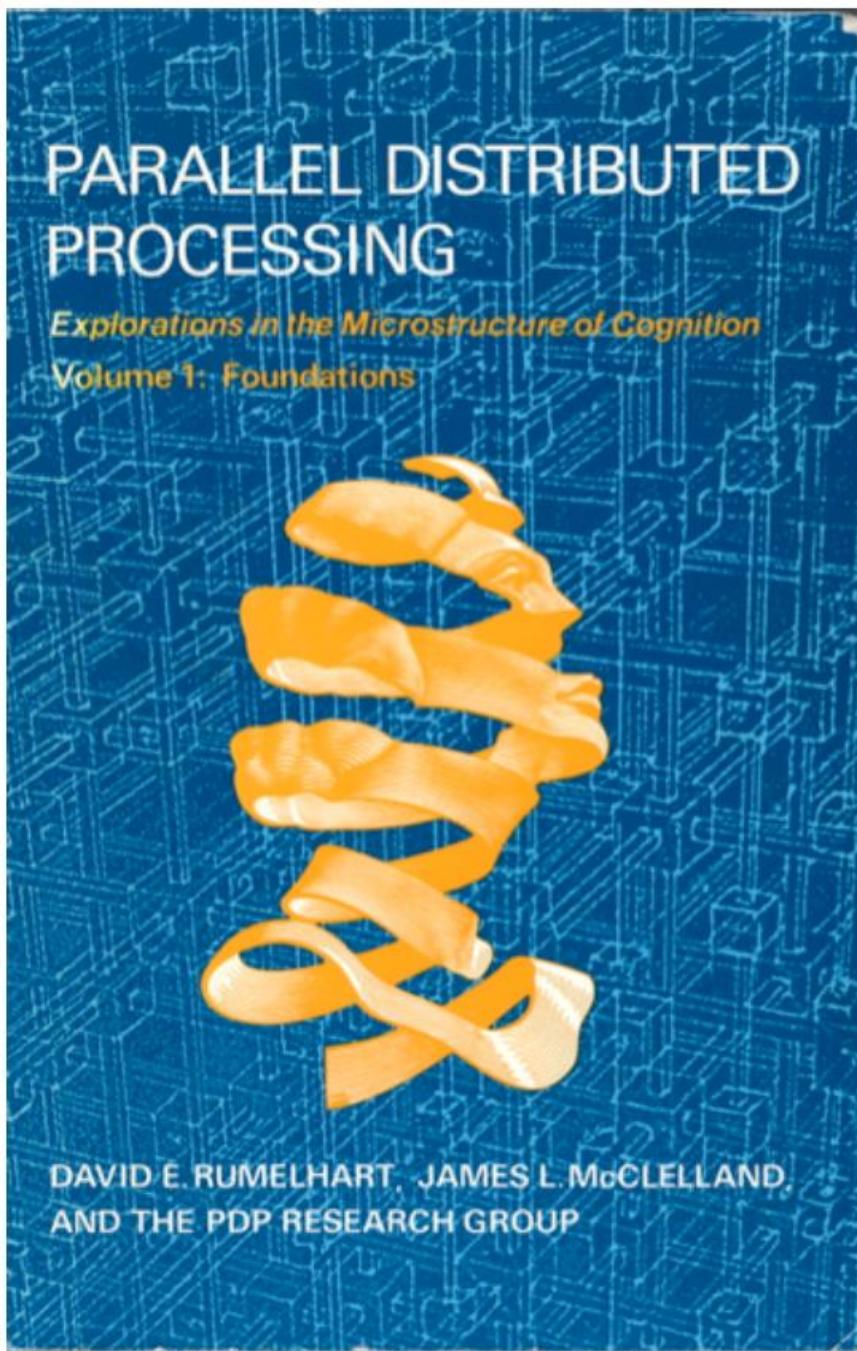
FOR BUYING OPTIONS, START HERE

Select Shipping Destination ▾

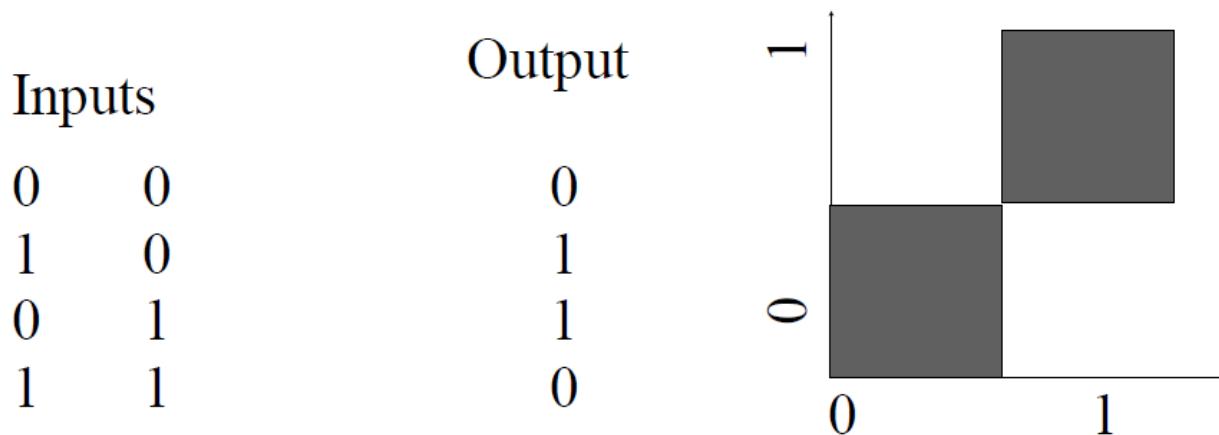
Paperback | \$35.00 Short | £24.95 |
ISBN: 9780262631112 | 308 pp. | 6 x
8.9 in | December 1987



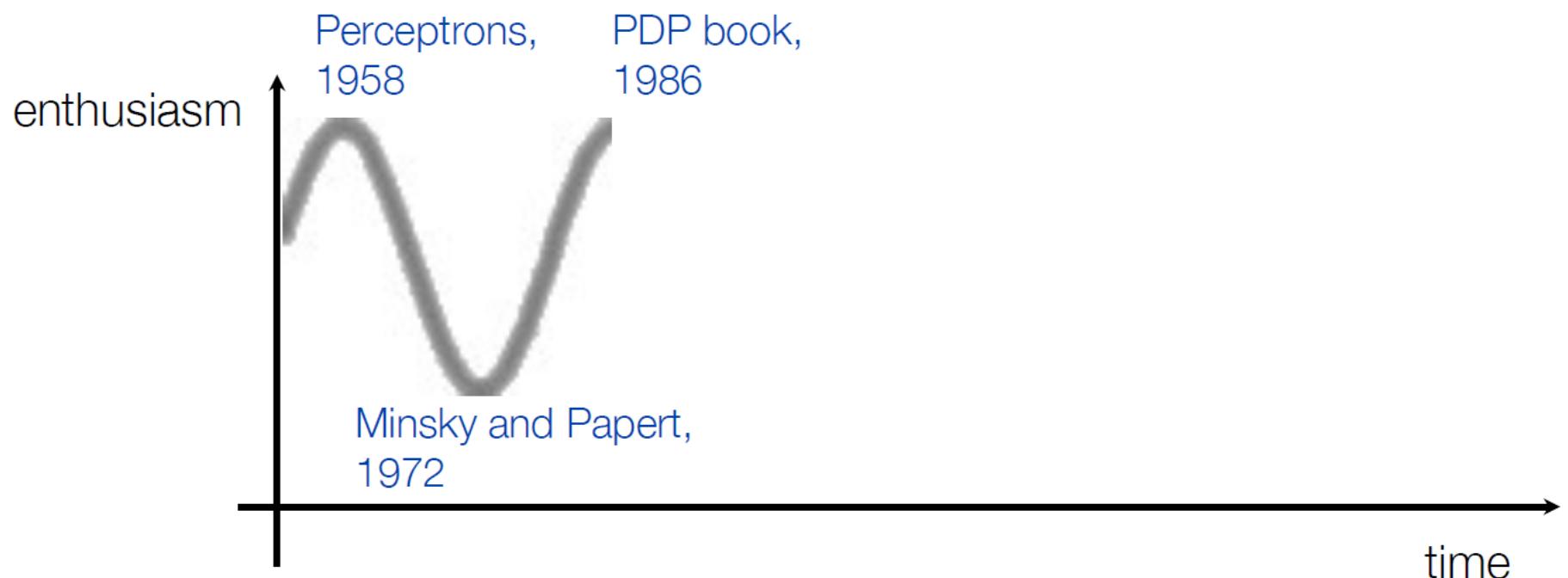
Parallel Distributed Processing (PDP), 1986



XOR problem



PDP authors pointed to the backpropagation algorithm as a breakthrough, allowing multi-layer neural networks to be trained. Among the functions that a multi-layer network can represent but a single-layer network cannot: the XOR function.



LeCun conv nets, 1998

PROC. OF THE IEEE, NOVEMBER 1998

7

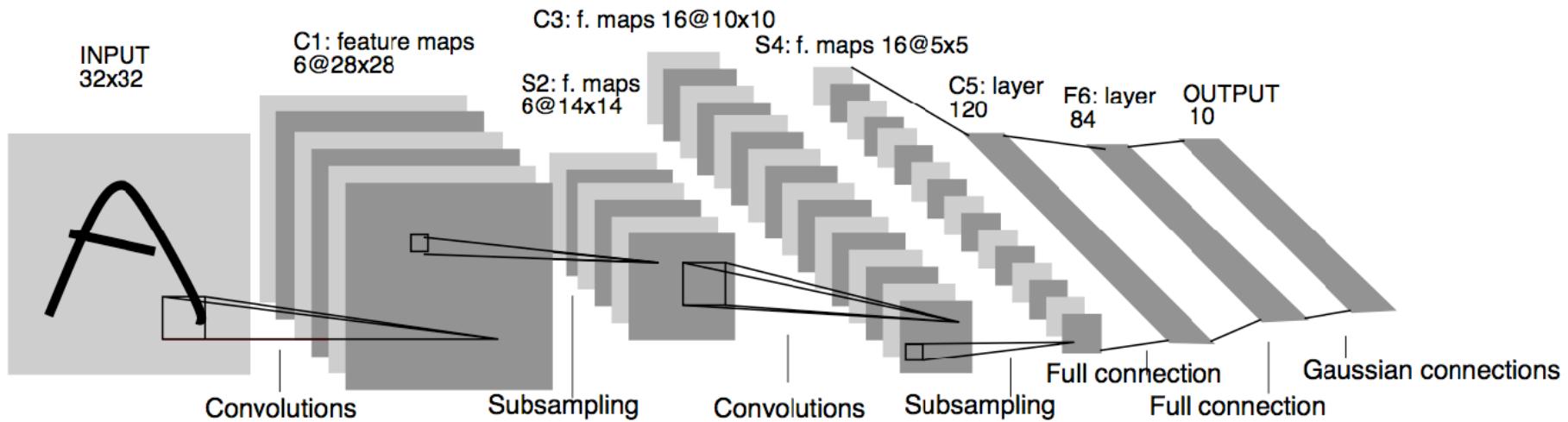


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Demos:

<http://yann.lecun.com/exdb/lenet/index.html>

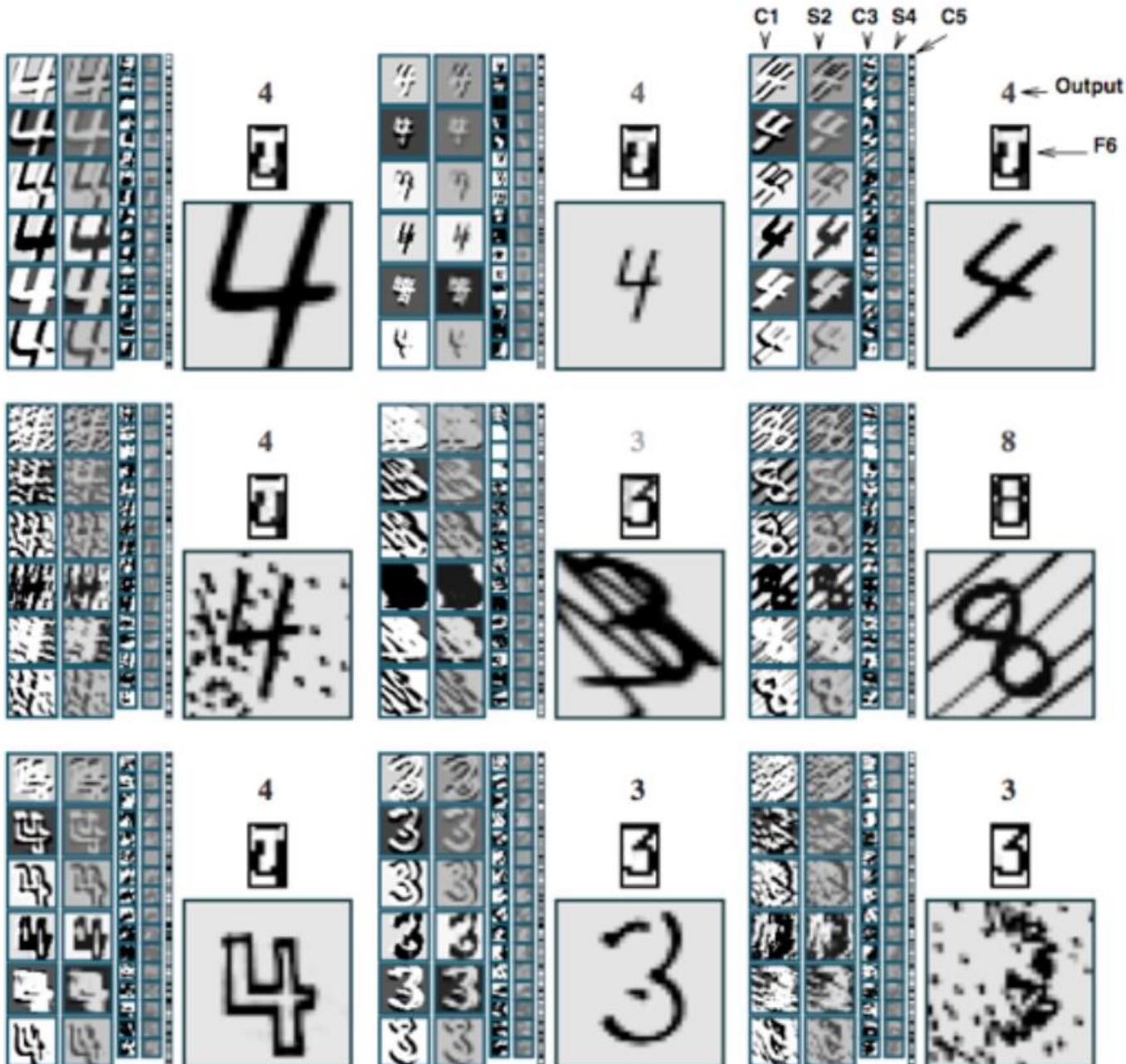
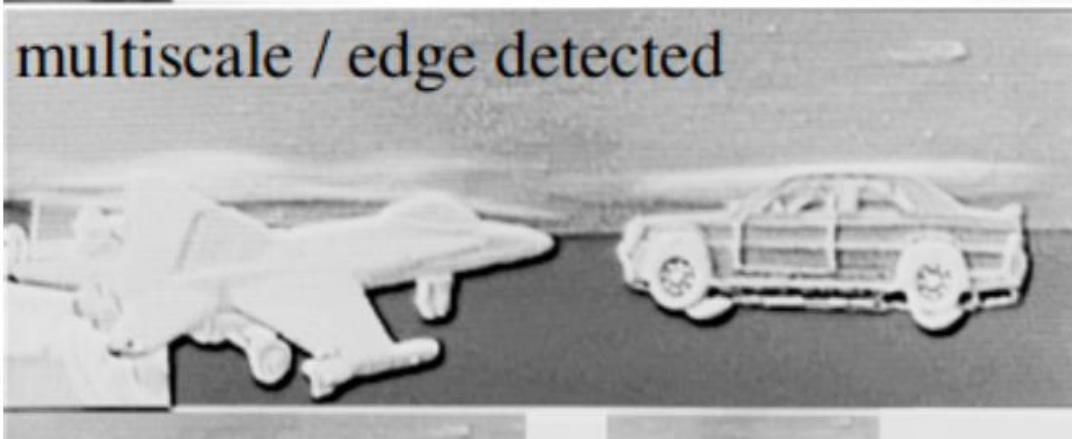


Fig. 13. Examples of unusual, distorted, and noisy characters correctly recognized by LeNet-5. The grey-level of the output label represents the penalty (lighter for higher penalties).

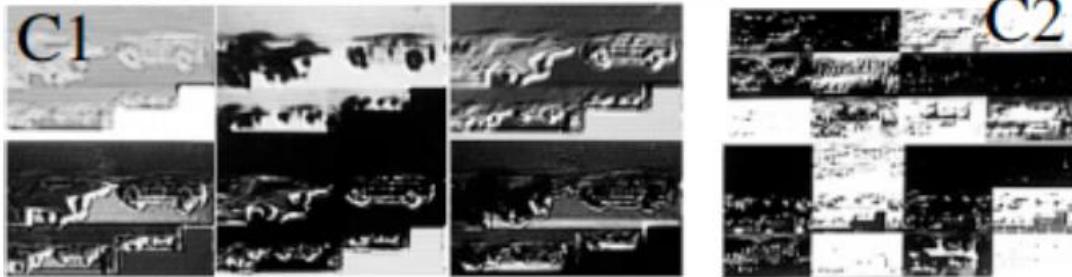
input



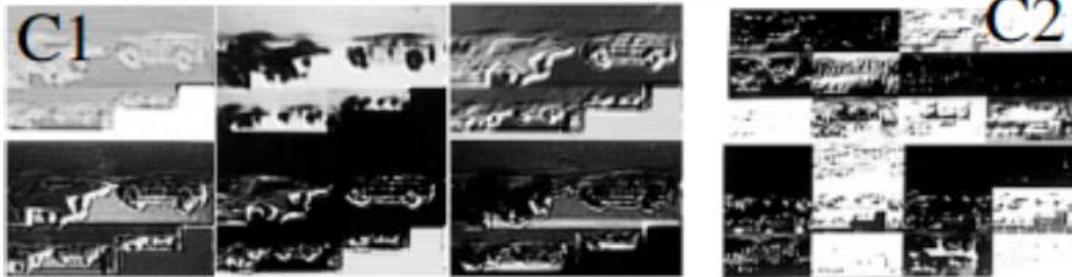
multiscale / edge detected



C1



C2



Neural networks to
recognize handwritten
digits?

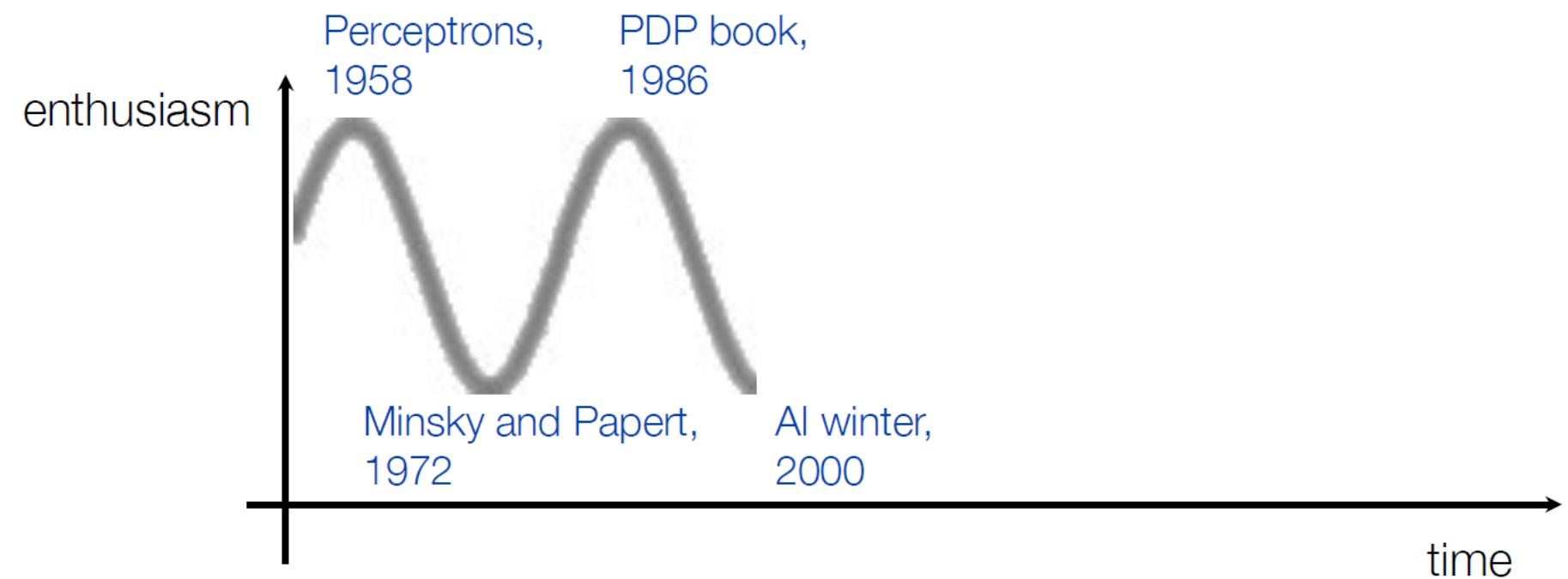
Yes

Neural networks for
tougher problems?

Not really

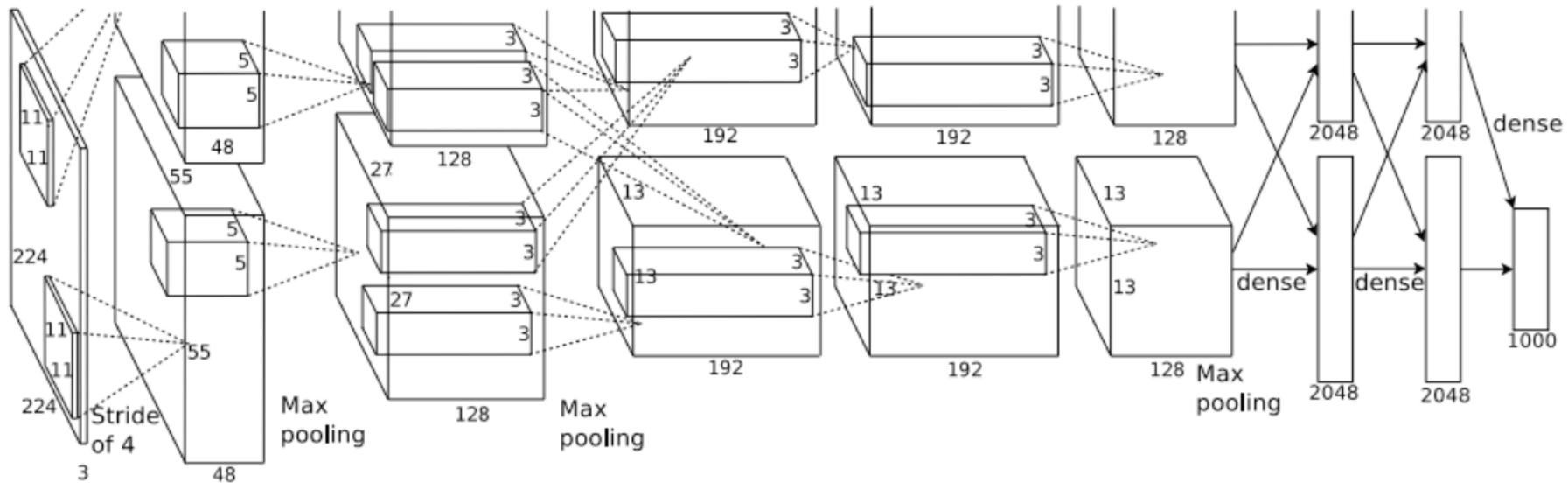
Neural Information Processing Systems 2000

- Neural Information Processing Systems, is the premier conference on machine learning. Evolved from an interdisciplinary conference to a machine learning conference.
- For the 2000 conference:
 - title words predictive of paper acceptance: “Belief Propagation” and “Gaussian”.
 - title words predictive of paper rejection: “Neural” and “Network”.



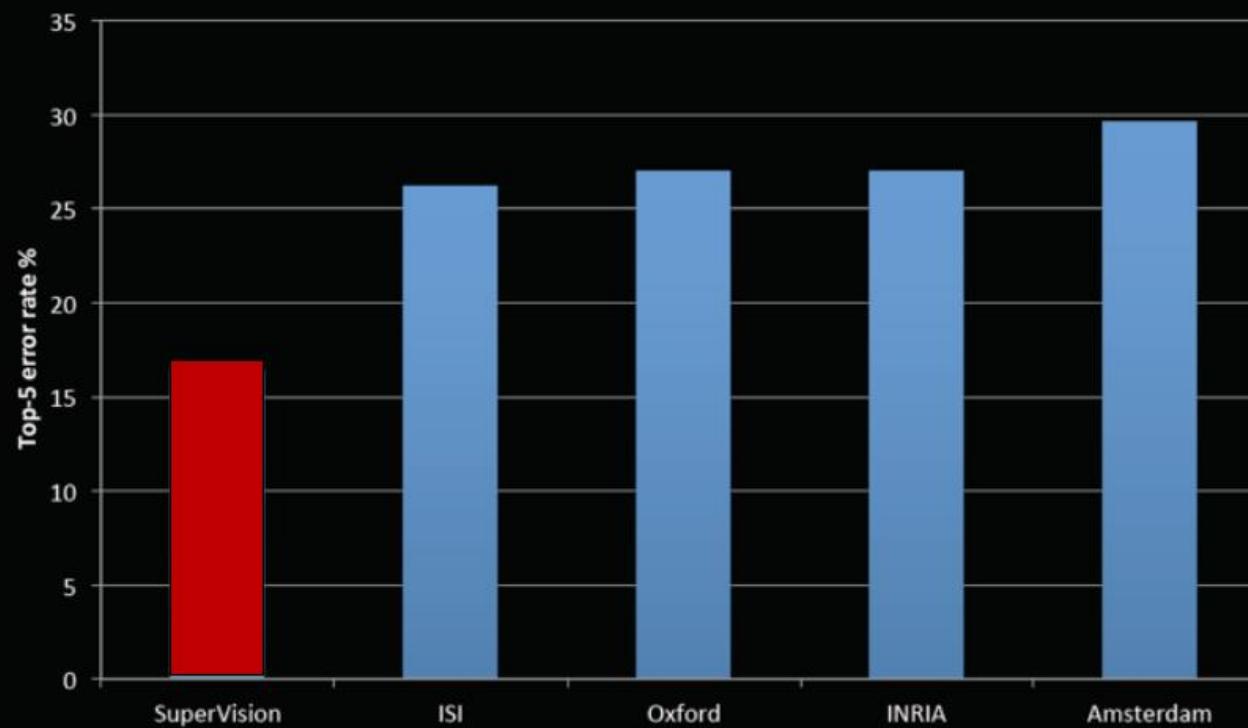
Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

“Alexnet”



ImageNet Classification 2012

- Krizhevsky et al. -- 16.4% error (top-5)
- Next best (non-convnet) – 26.2% error



Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



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



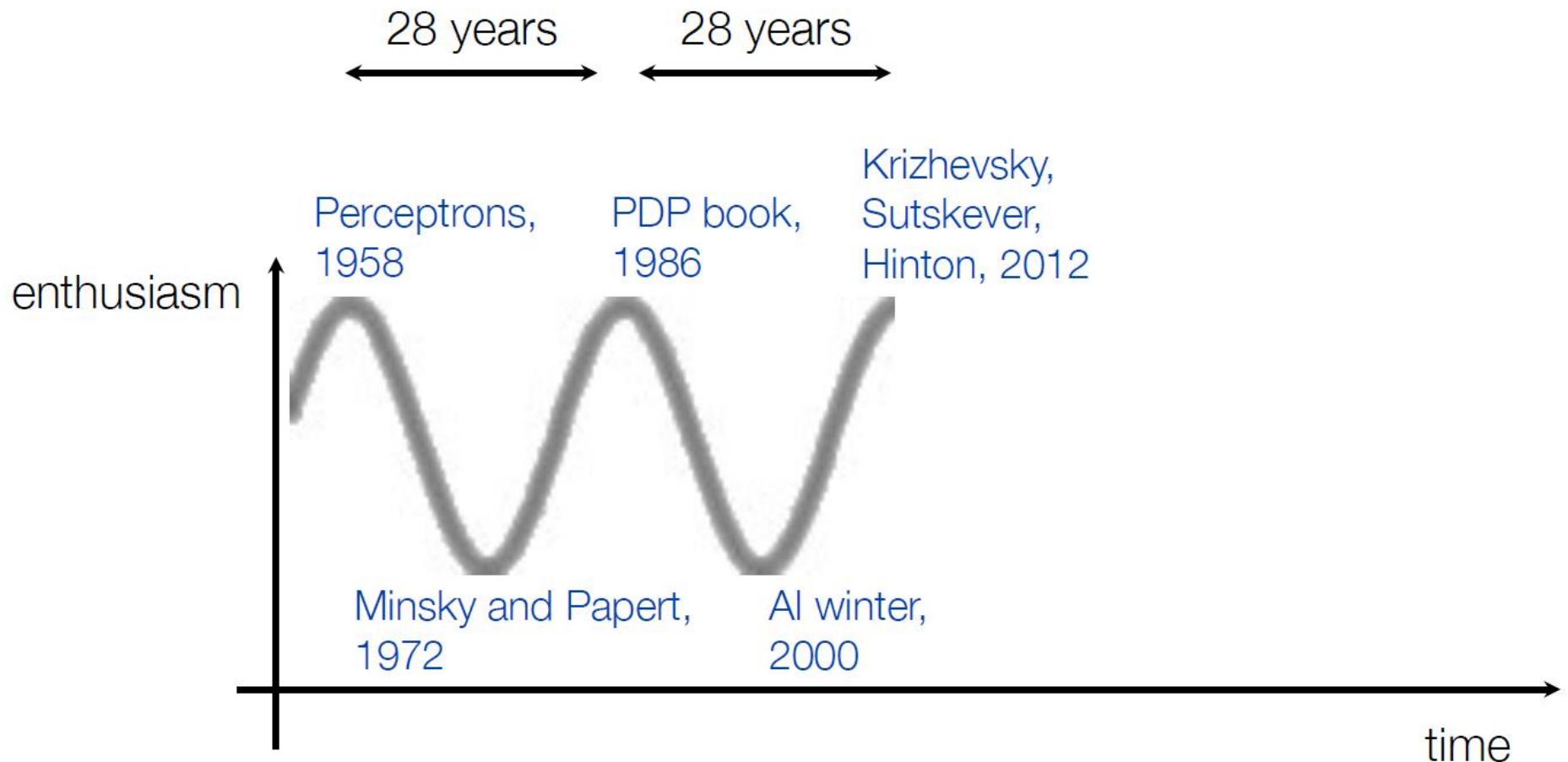
grille

mushroom

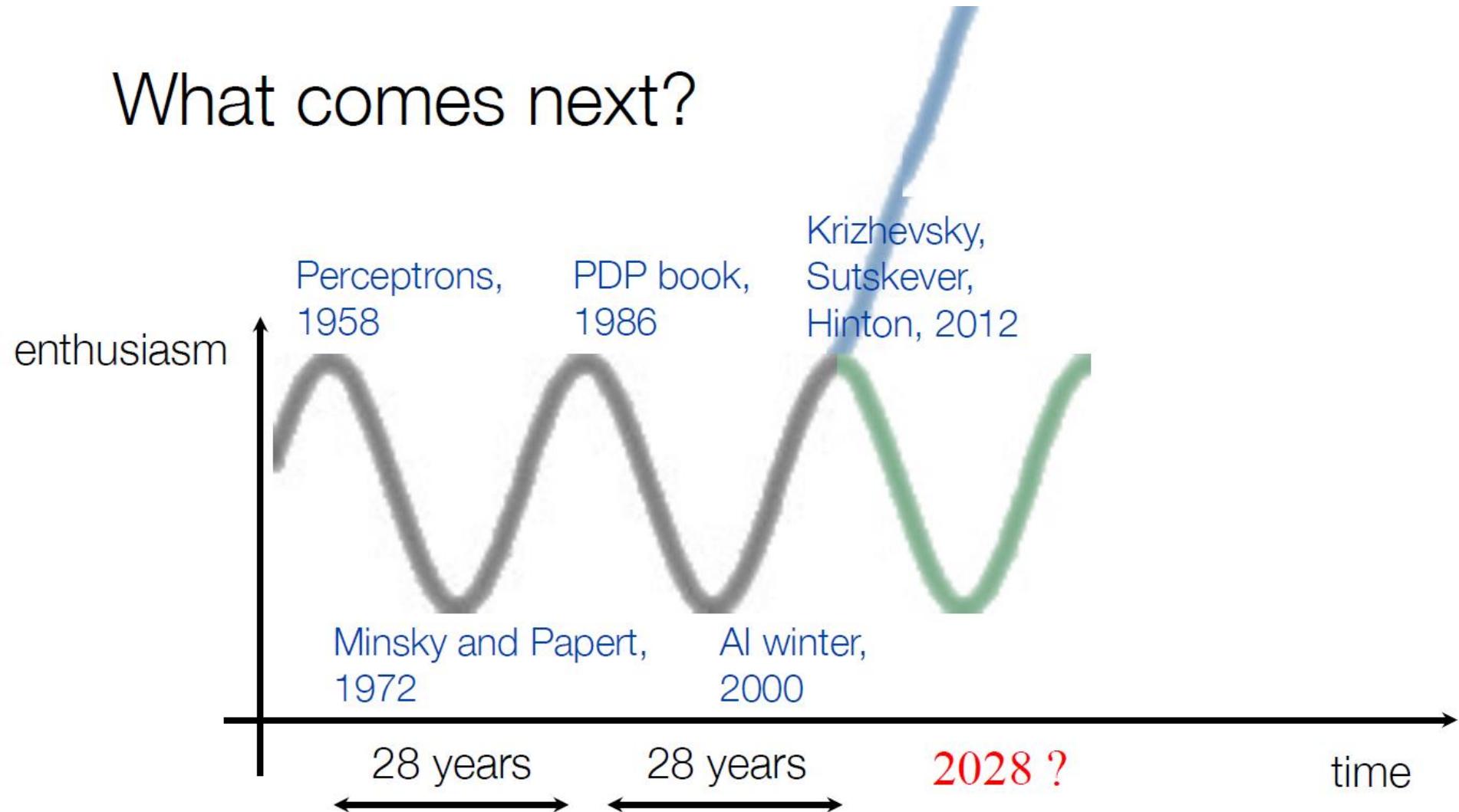
cherry

Madagascar cat

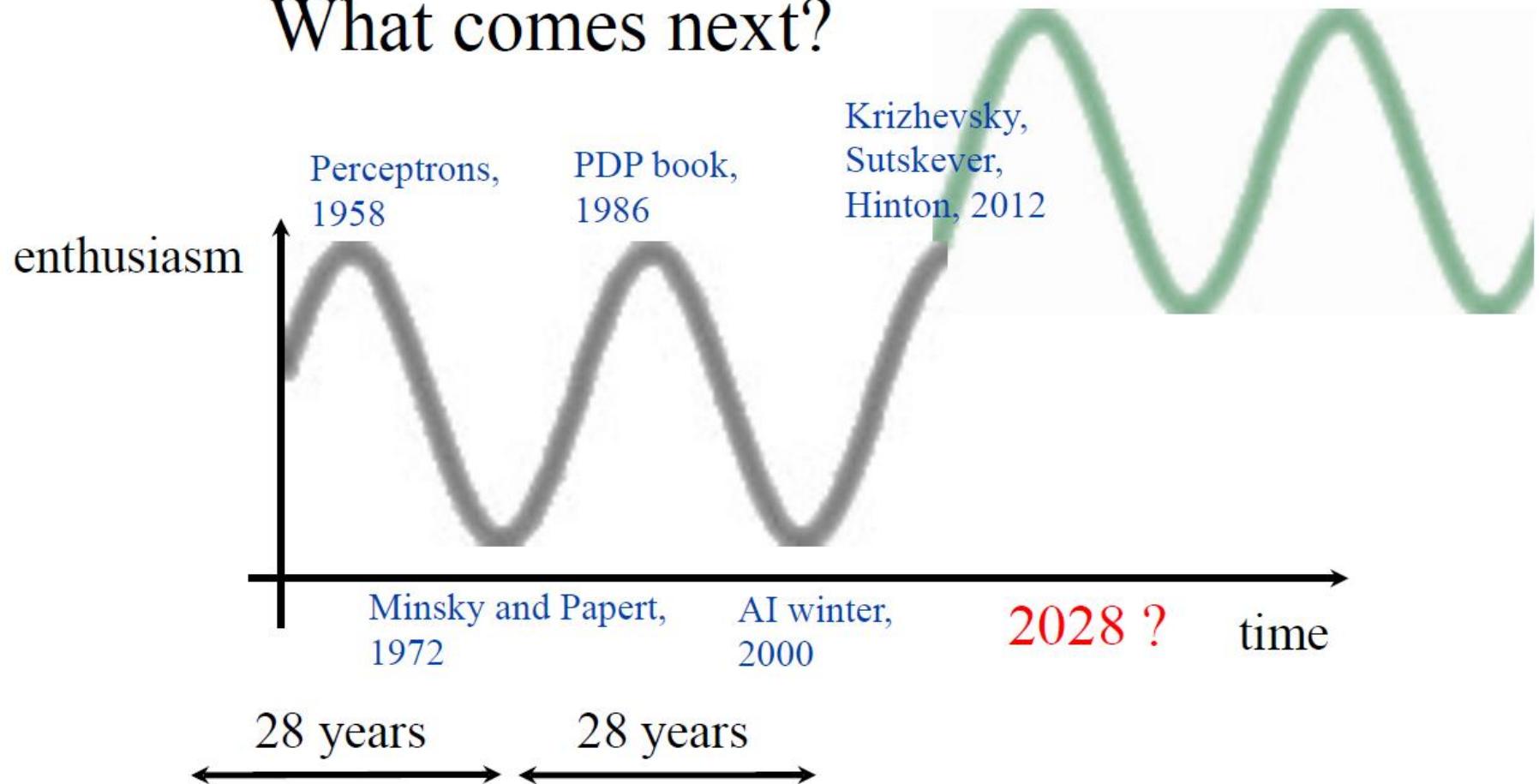
convertible	agaric	dalmatian	squirrel monkey
grille	mushroom	grape	spider monkey
pickup	jelly fungus	elderberry	titi
beach wagon	gill fungus	ffordshire bullterrier	indri
fire engine	dead-man's-fingers	currant	howler monkey



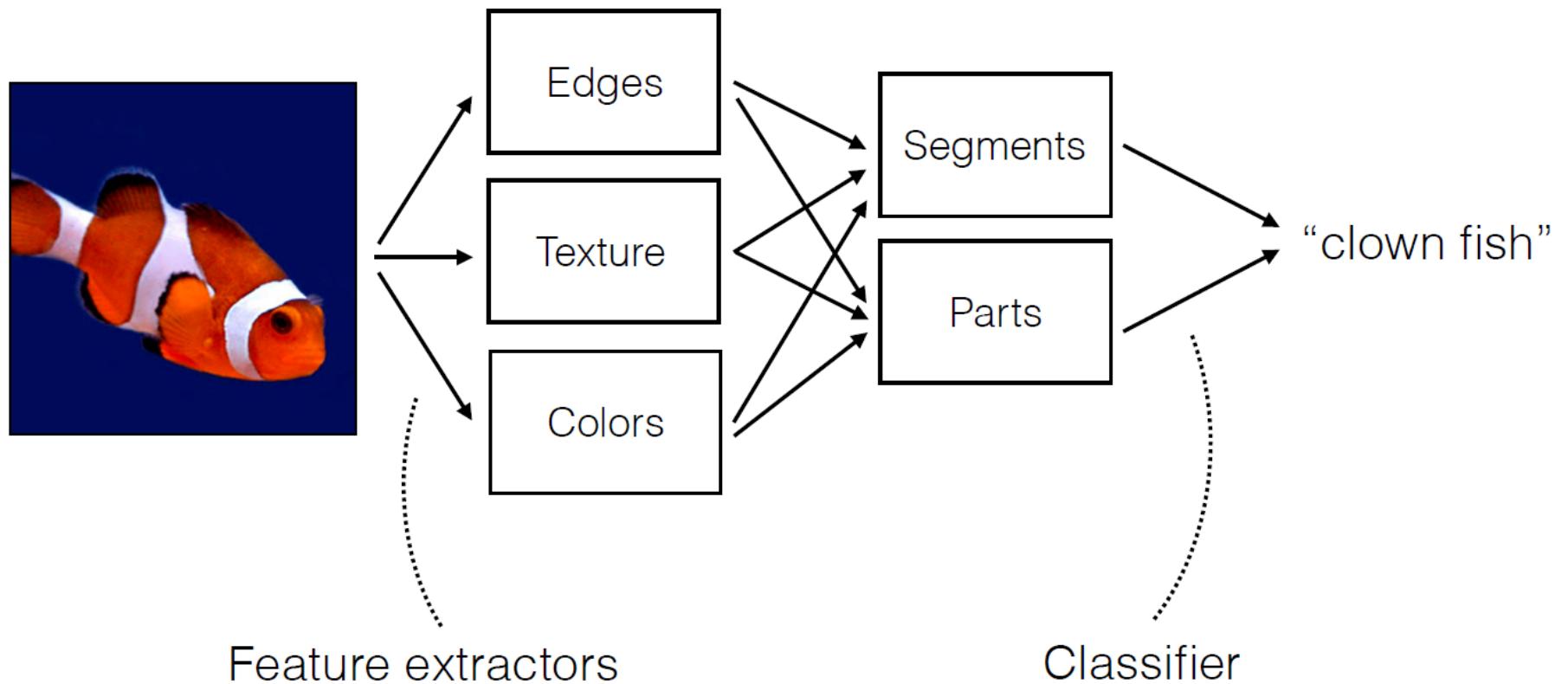
What comes next?



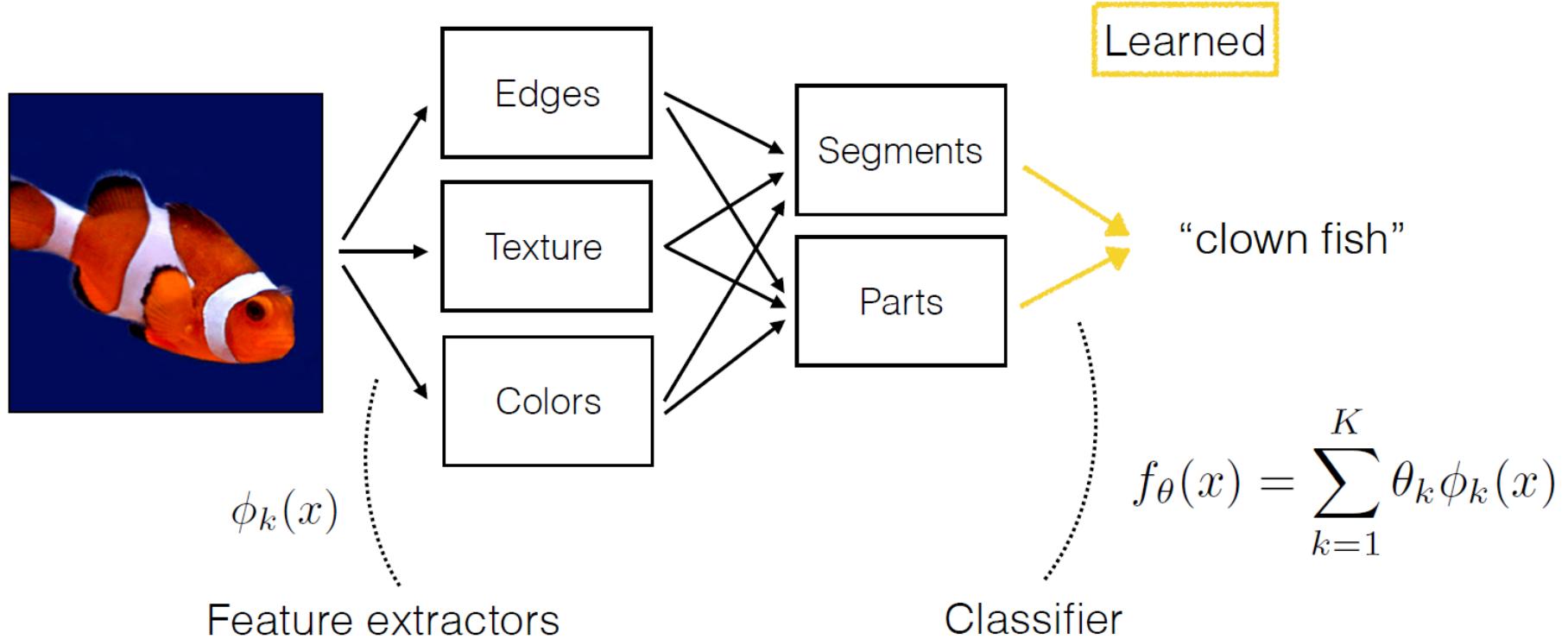
What comes next?



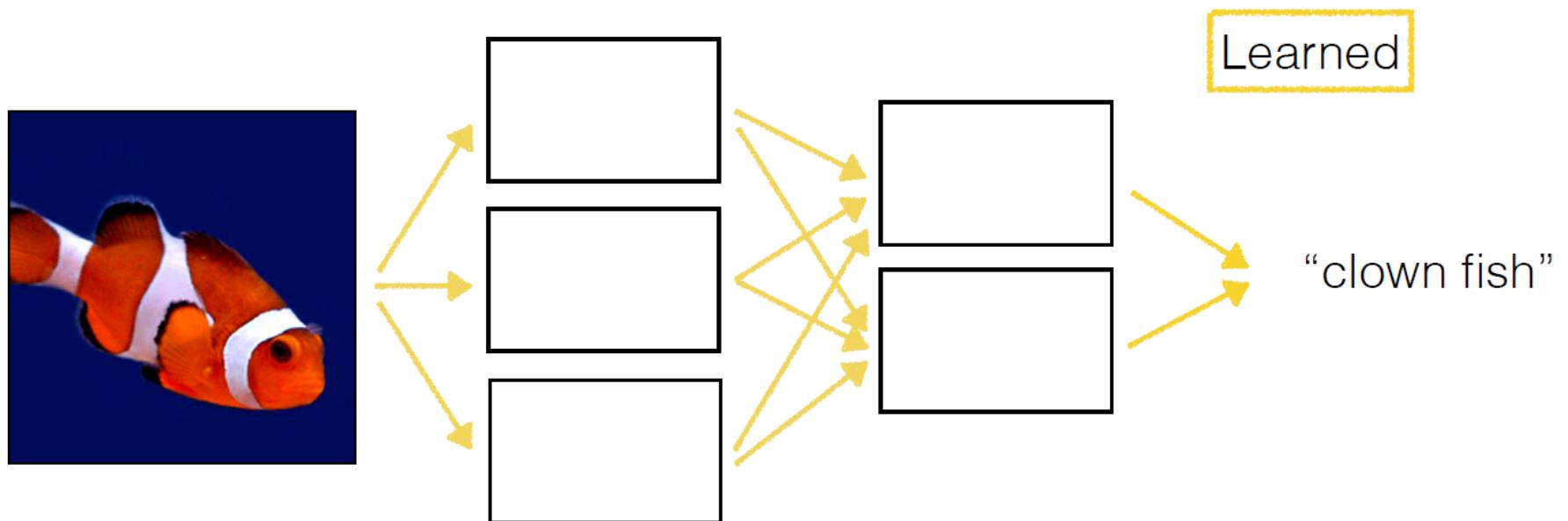
Object Recognition



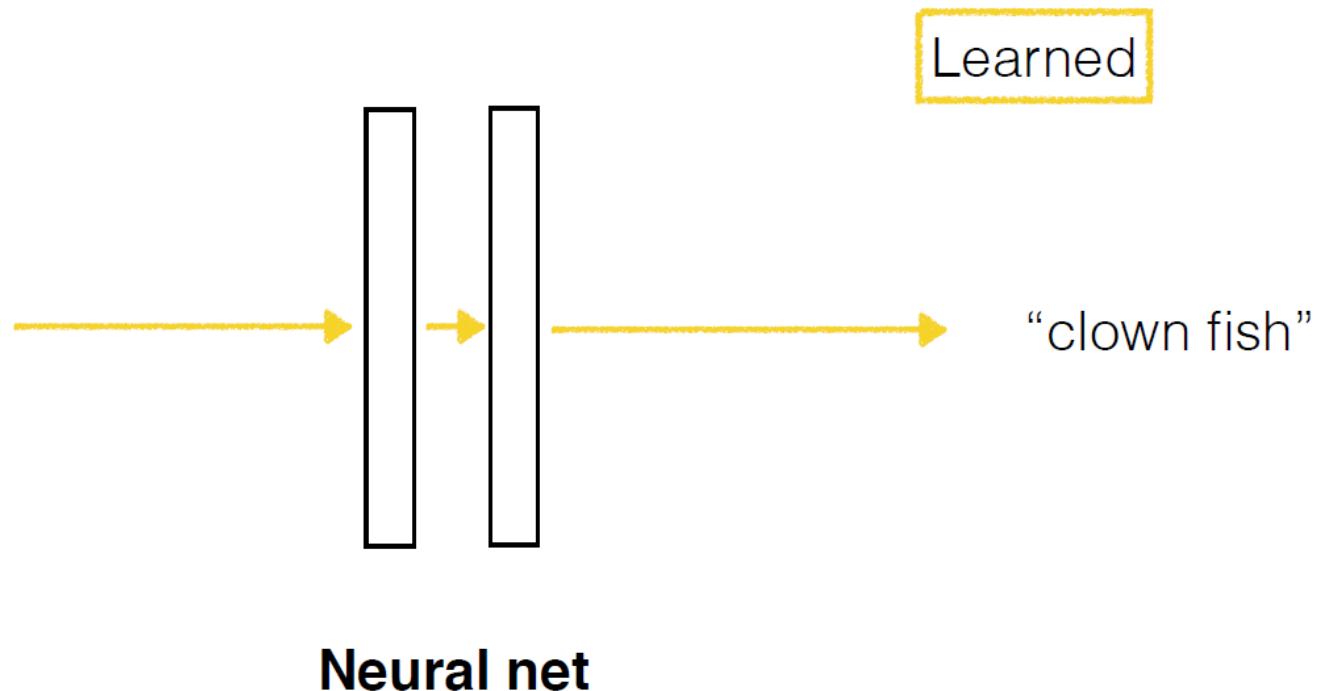
Object Recognition



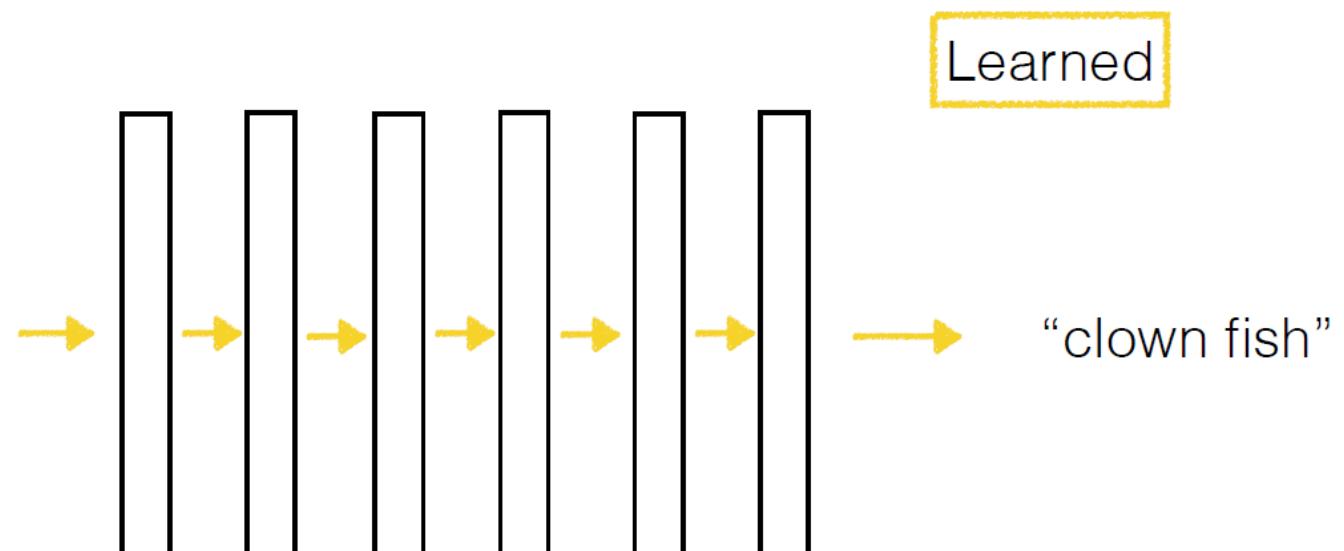
Object Recognition



Object Recognition



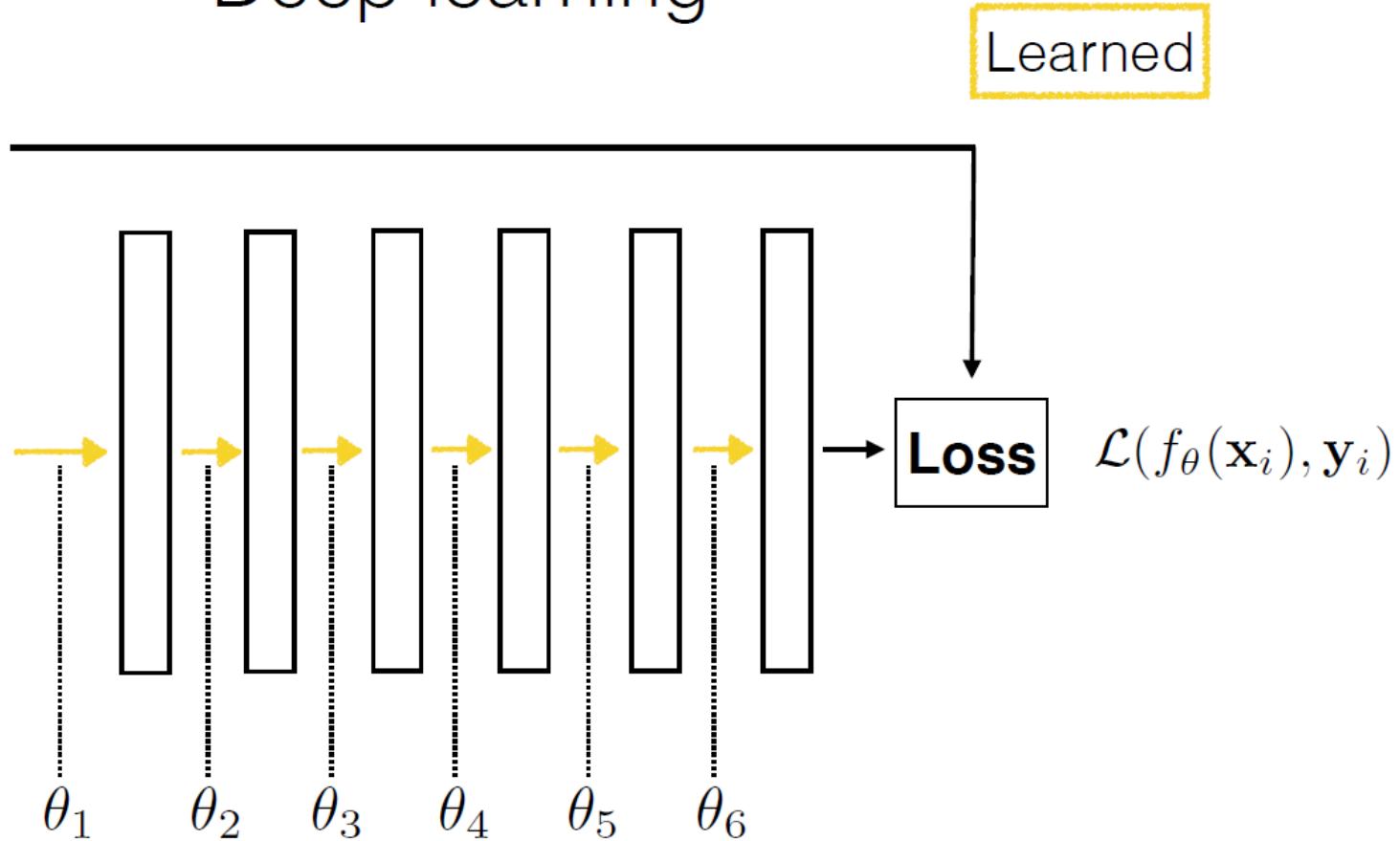
Object Recognition



Deep neural net

Deep learning

\mathbf{y}_i
“clown fish”



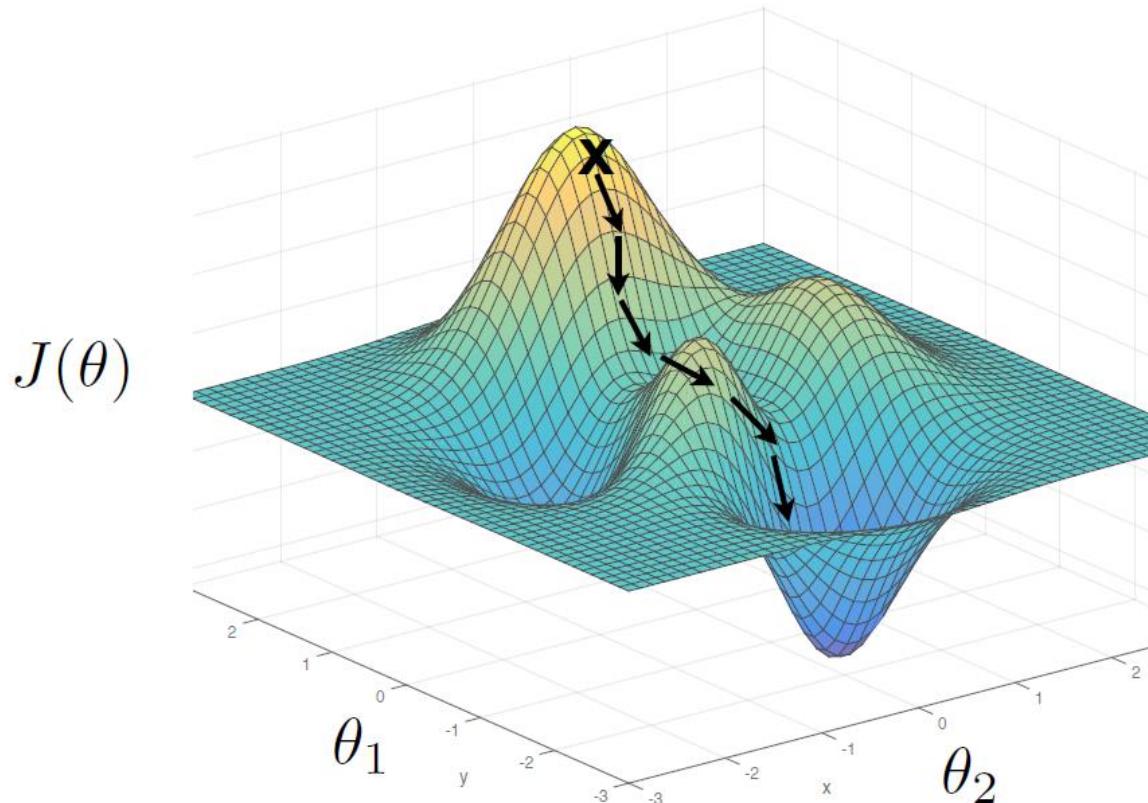
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

Gradient descent

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

$$\overbrace{\hspace{10em}}^{J(\theta)}$$

Gradient Descent



$$\theta^* = \arg \min_{\theta} J(\theta)$$

Gradient Descent

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

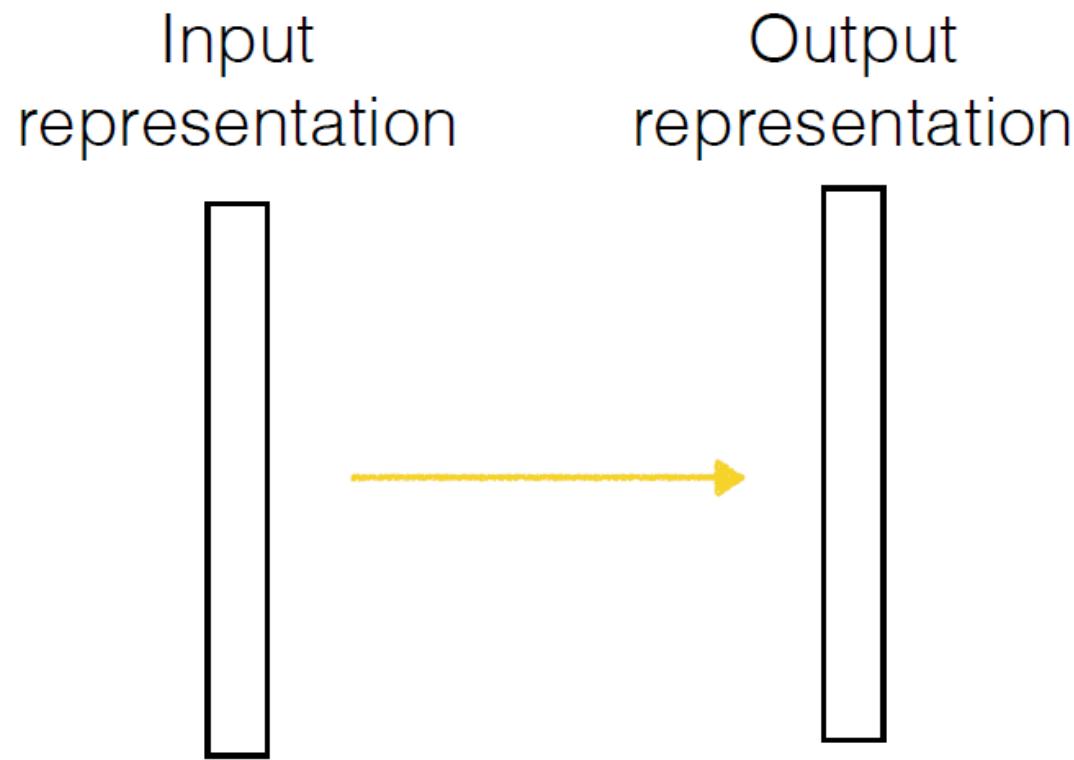
$\overbrace{\hspace{10em}}$
 $J(\theta)$

One iteration of gradient descent:

$$\theta^{t+1} = \theta^t - \eta_t \frac{\partial J(\theta)}{\partial \theta} \Big|_{\theta=\theta^t}$$

(learning rate)

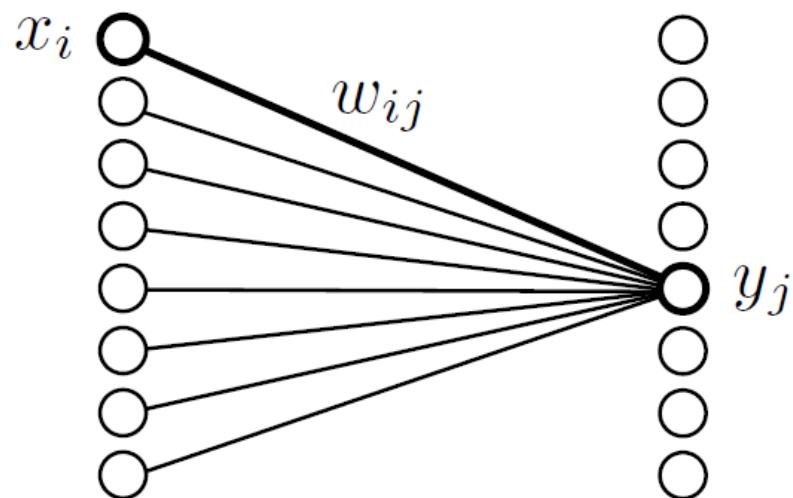
Computation in a neural network



Computation in a neural network

Linear layer

Input representation Output representation

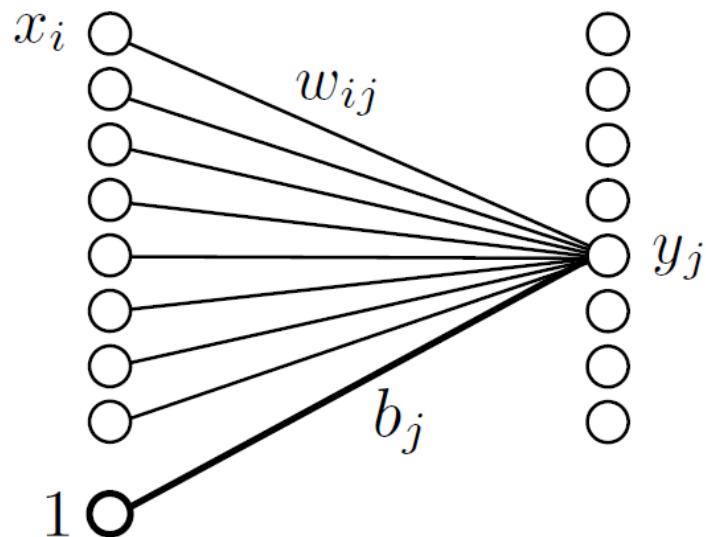


$$y_j = \sum_i w_{ij} x_i$$

Computation in a neural network

Linear layer

Input representation Output representation

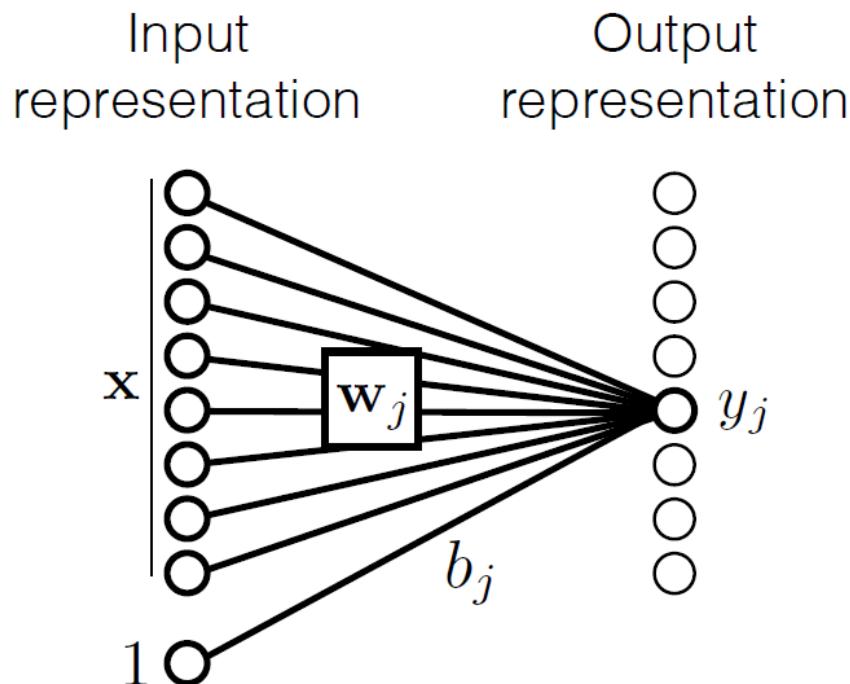


$$y_j = \sum_i w_{ij} x_i + b_j$$

weights
bias

Computation in a neural network

Linear layer



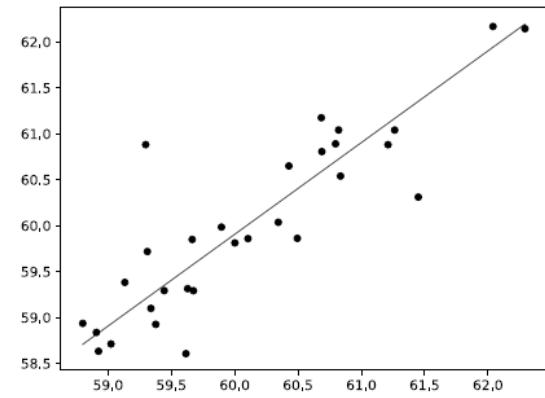
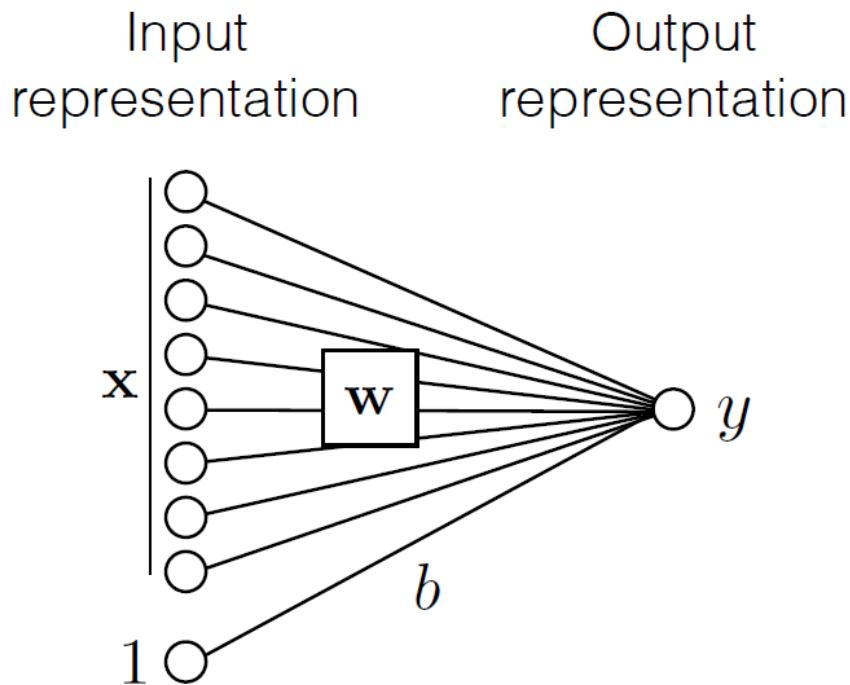
$$y_j = \mathbf{x}^T \mathbf{w}_j + b_j$$

weights
bias
parameters of the model

$\theta = \{\mathbf{W}, \mathbf{b}\}$

Example: linear regression with a neural net

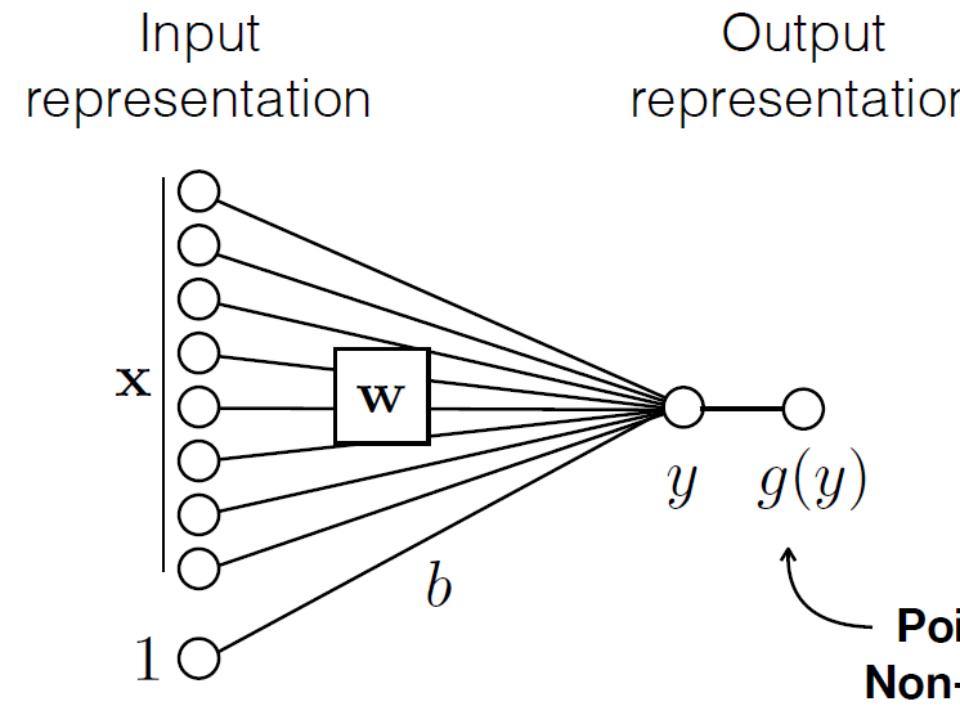
Linear layer



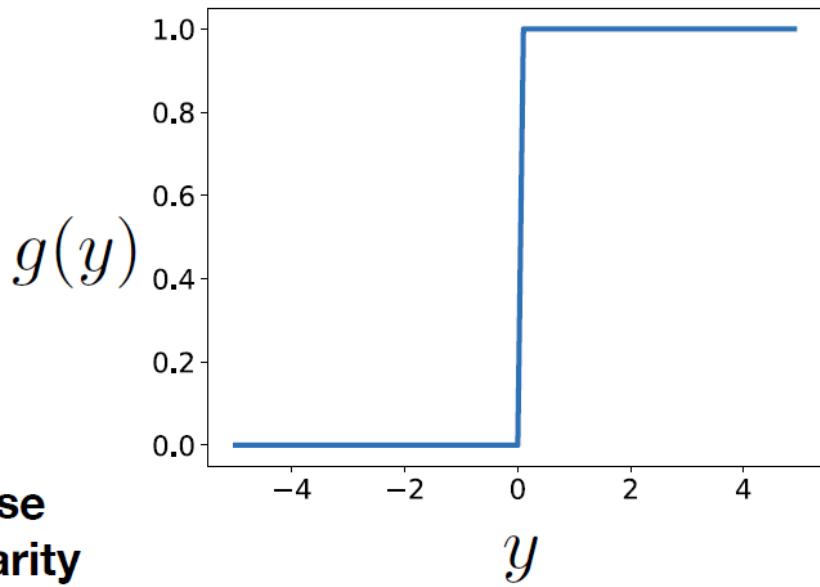
$$f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$$

Computation in a neural network

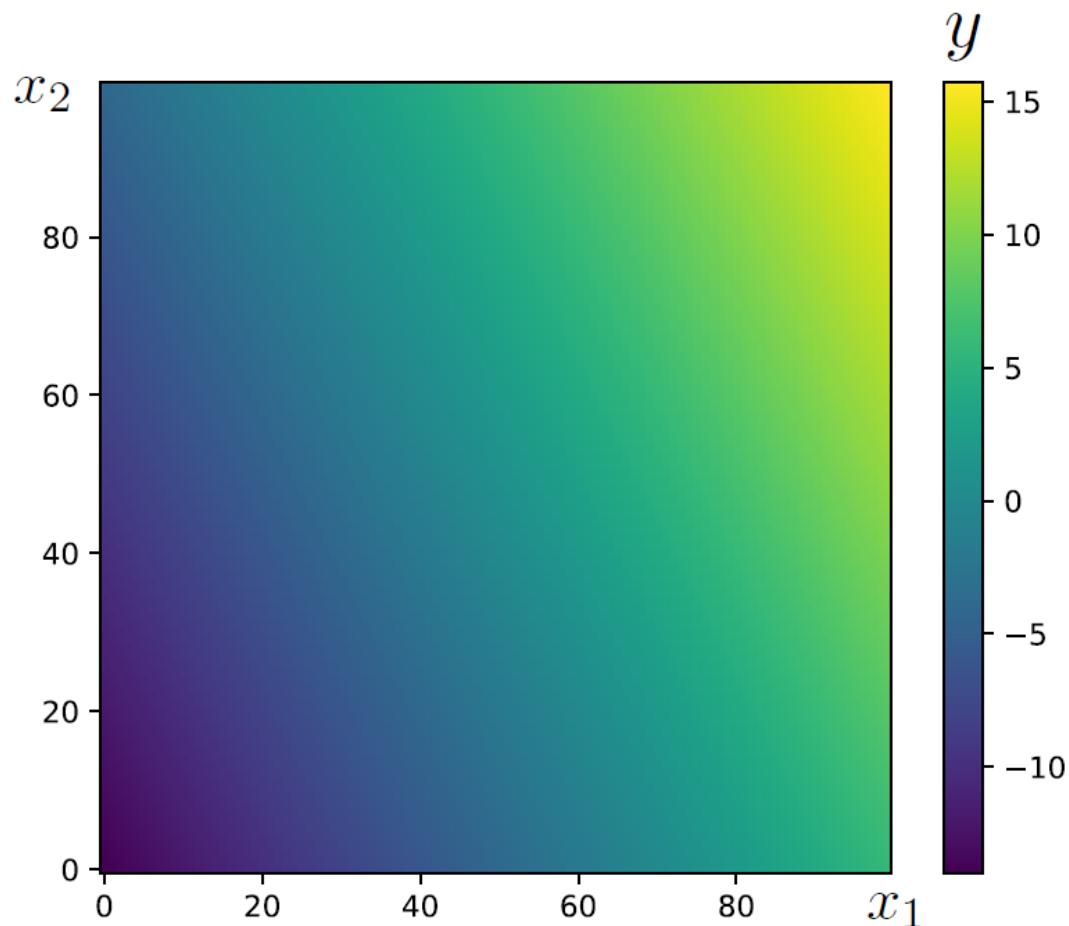
“Perceptron”



$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

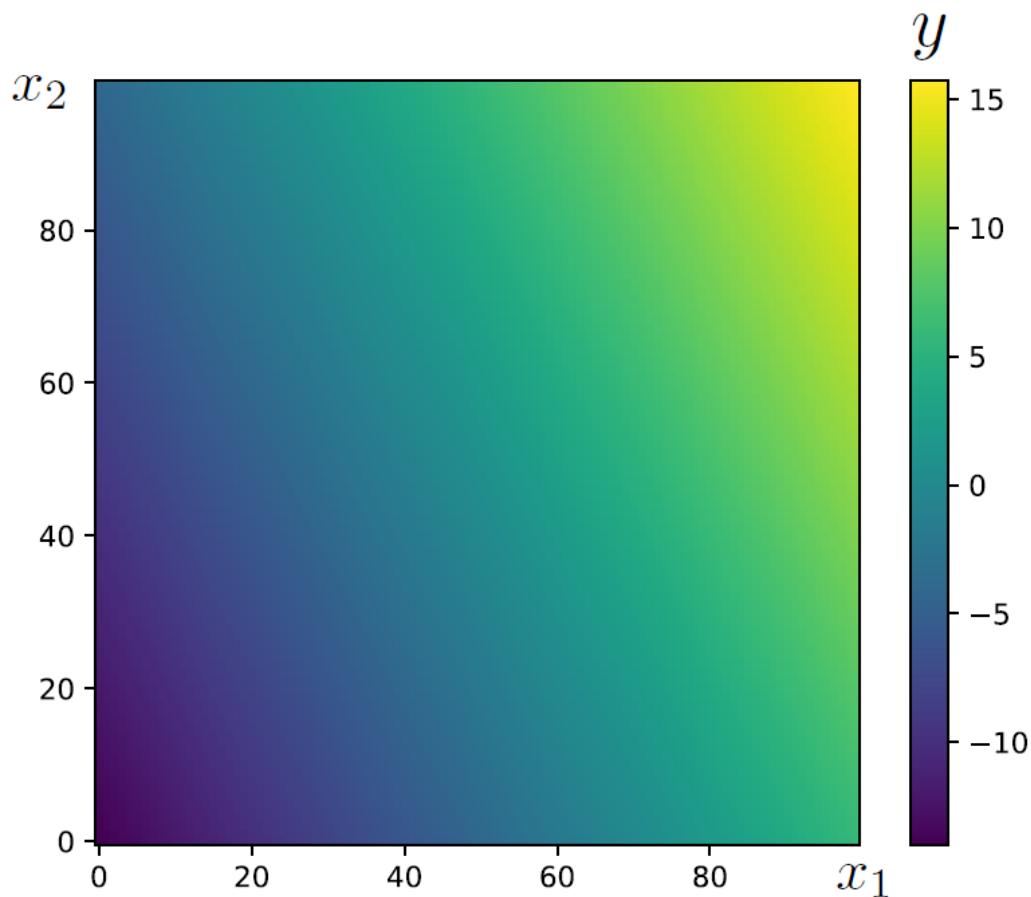


Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

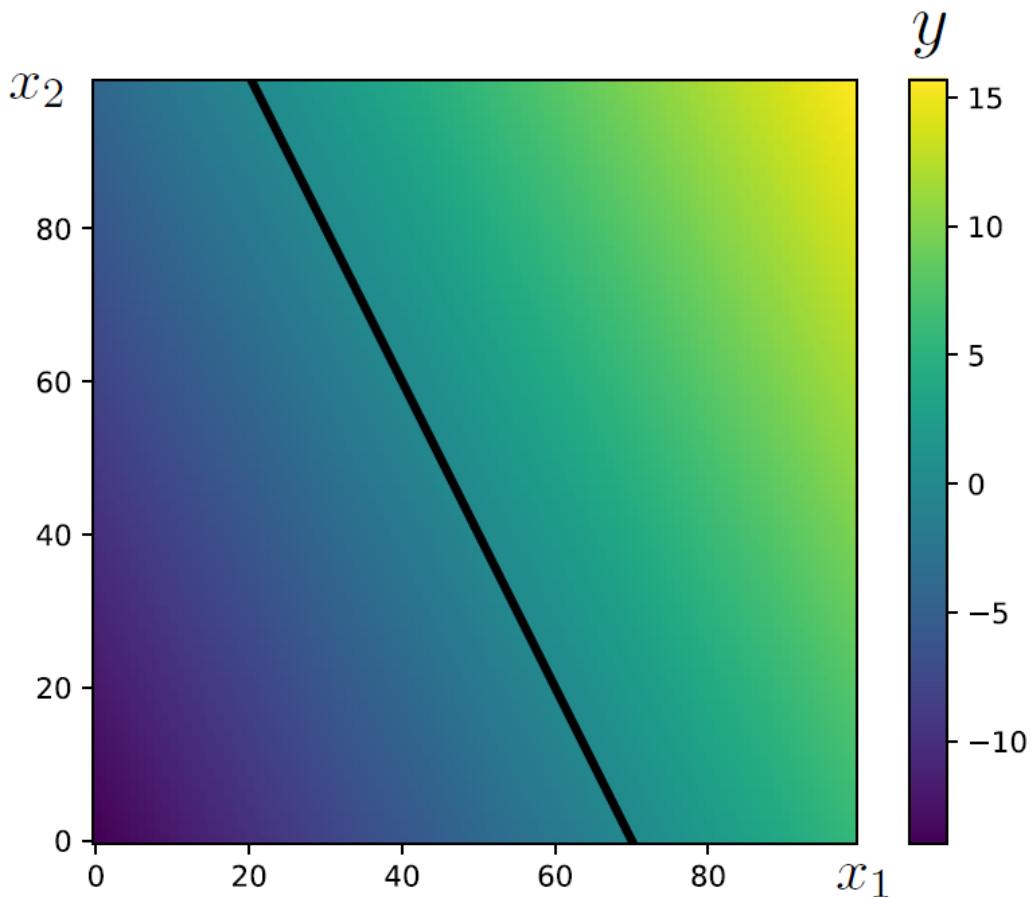
Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

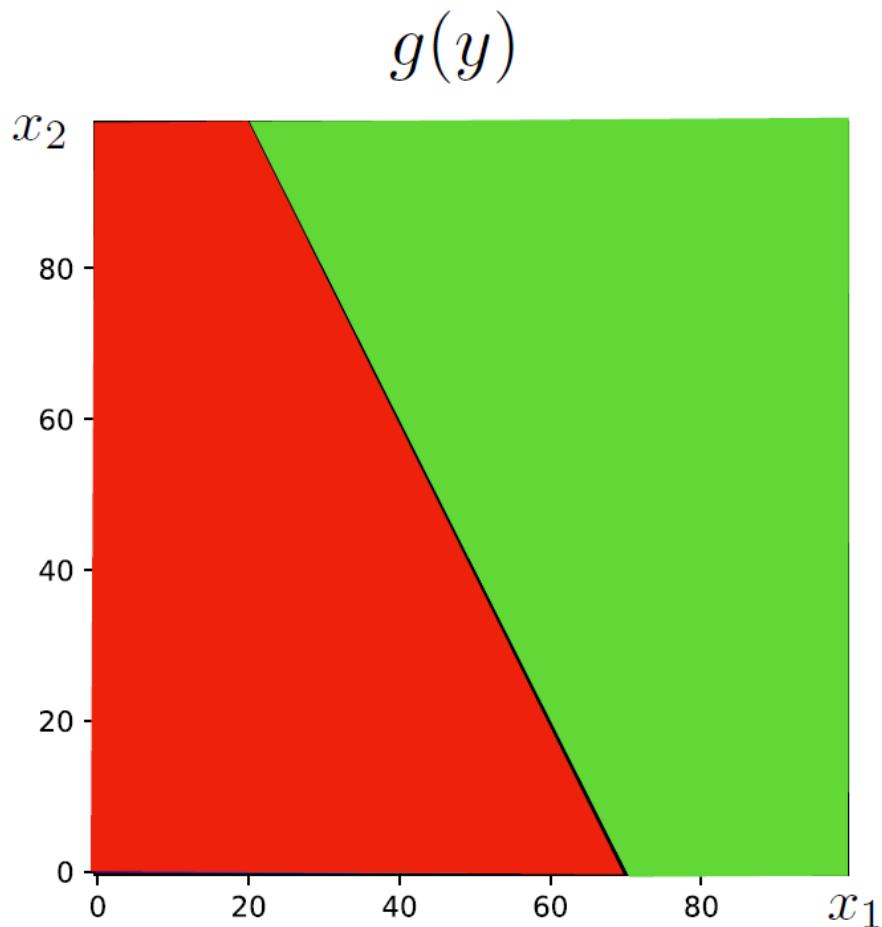
Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

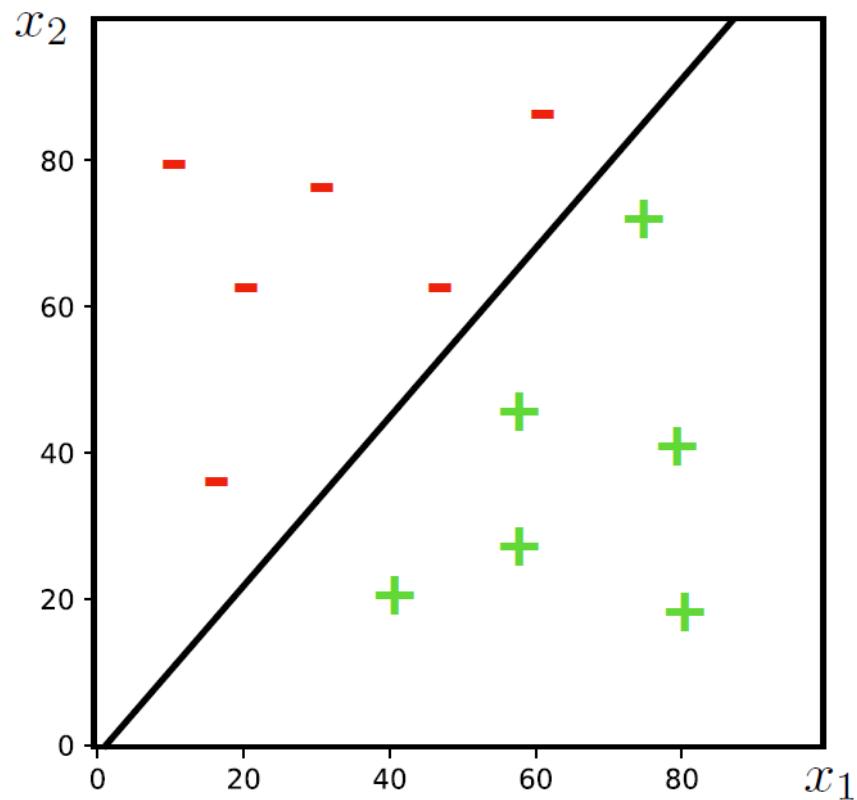
Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$

Example: linear classification with a perceptron

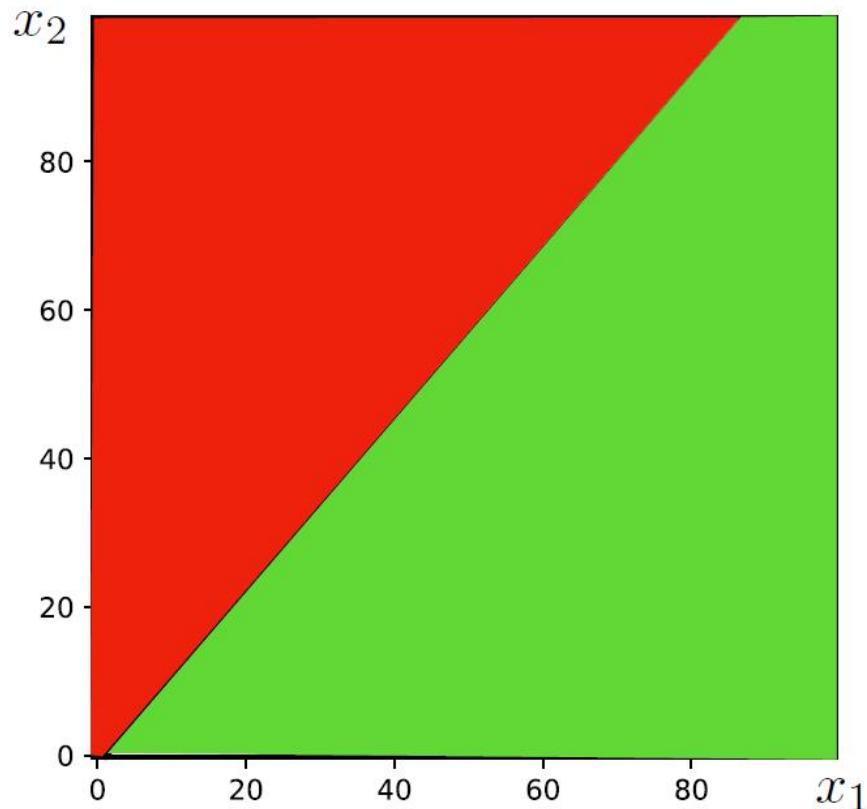


$$\hat{y} = \mathbf{x}^T \mathbf{w} + b$$

$$g(\hat{y}) = \begin{cases} 1, & \text{if } \hat{y} > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \mathcal{L}(g(\hat{y}), y_i)$$

Example: linear classification with a perceptron



$$\hat{y} = \mathbf{x}^T \mathbf{w} + b$$

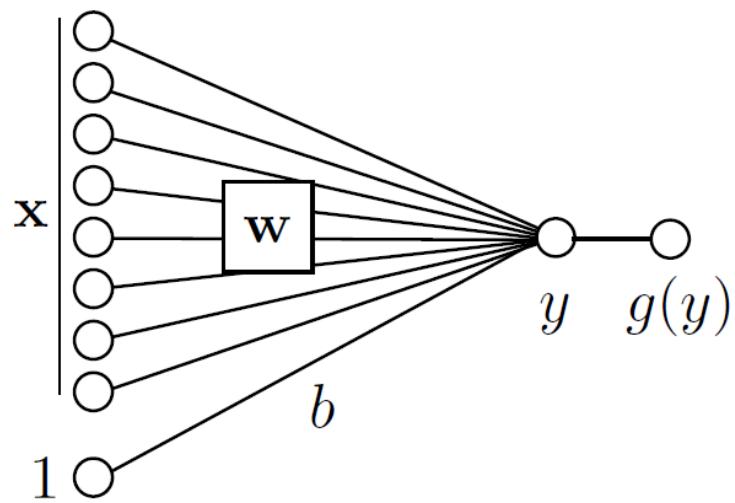
$$g(\hat{y}) = \begin{cases} 1, & \text{if } \hat{y} > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \mathcal{L}(g(\hat{y}), y_i)$$

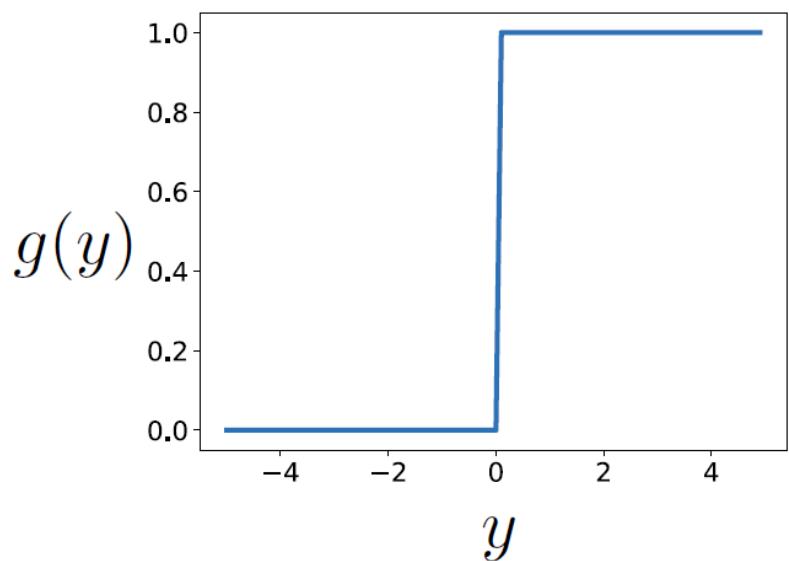
Computation in a neural network

Input representation

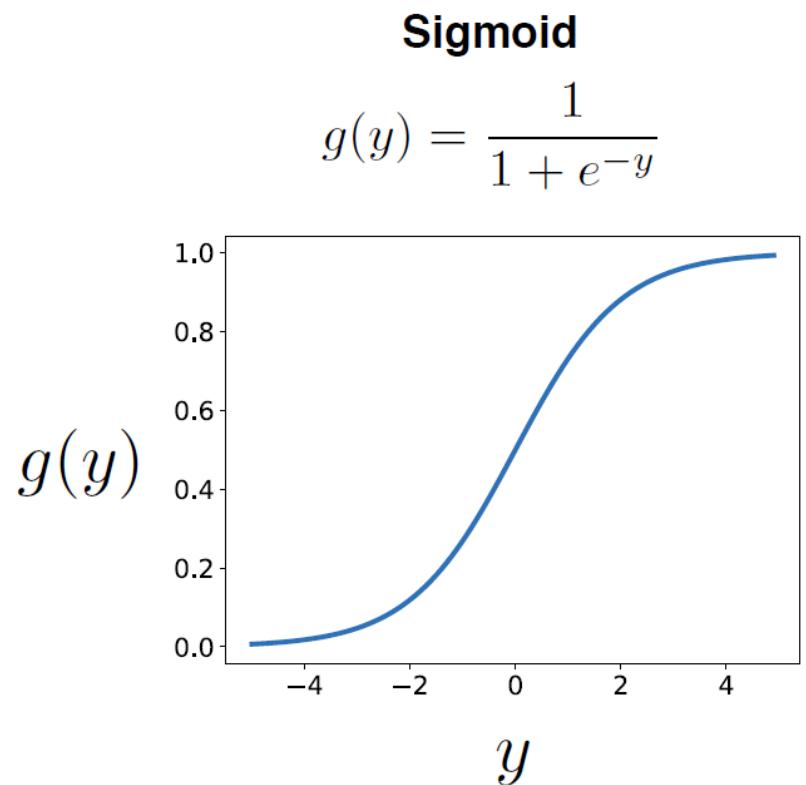
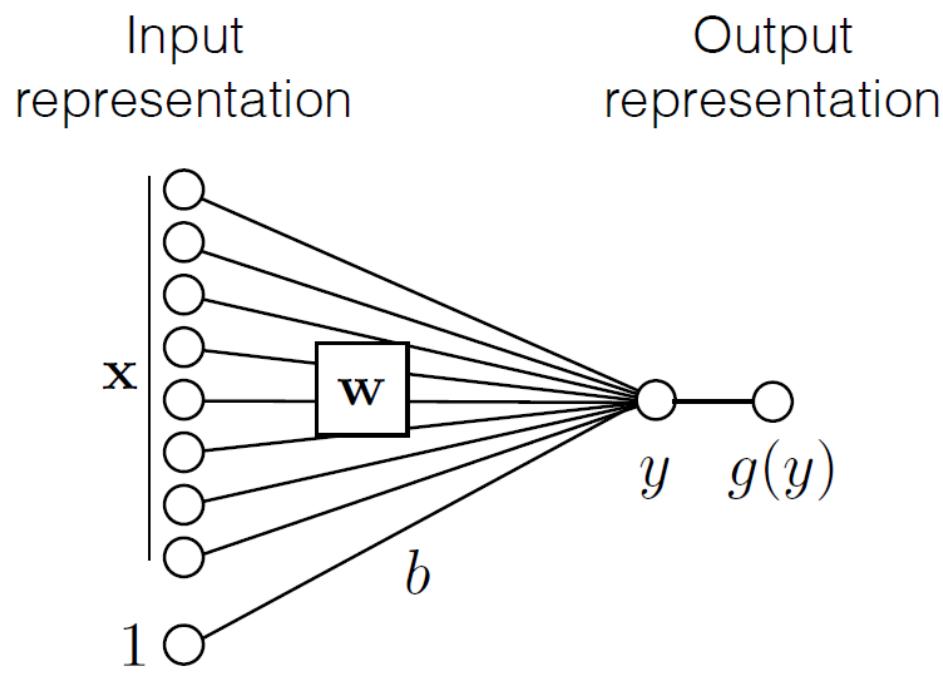
Output representation



$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



Computation in a neural network - nonlinearity

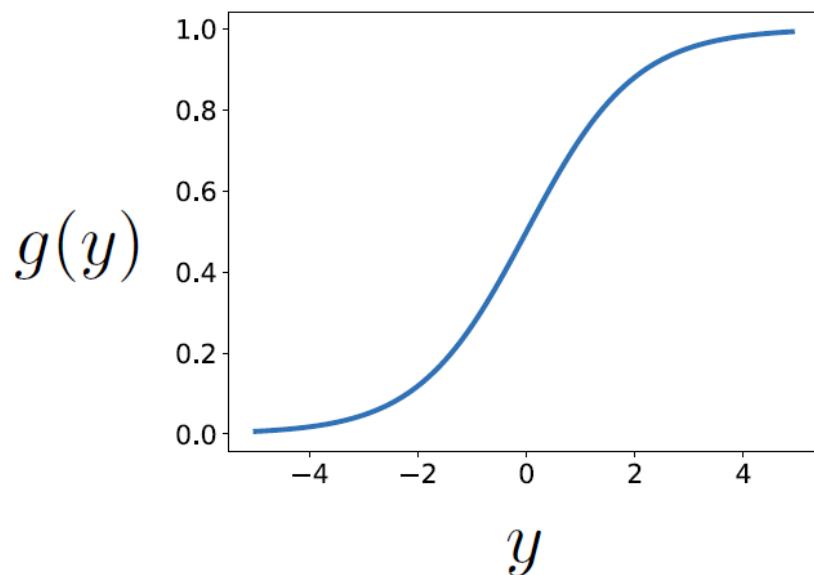


Computation in a neural network - nonlinearity

- Interpretation as firing rate of neuron
- Bounded between [0,1]
- Saturation for large +/- inputs
- Gradients go to zero
- Output centered at 0.5 (poor conditioning)
- Not used in practice

Sigmoid

$$g(y) = \frac{1}{1 + e^{-y}}$$

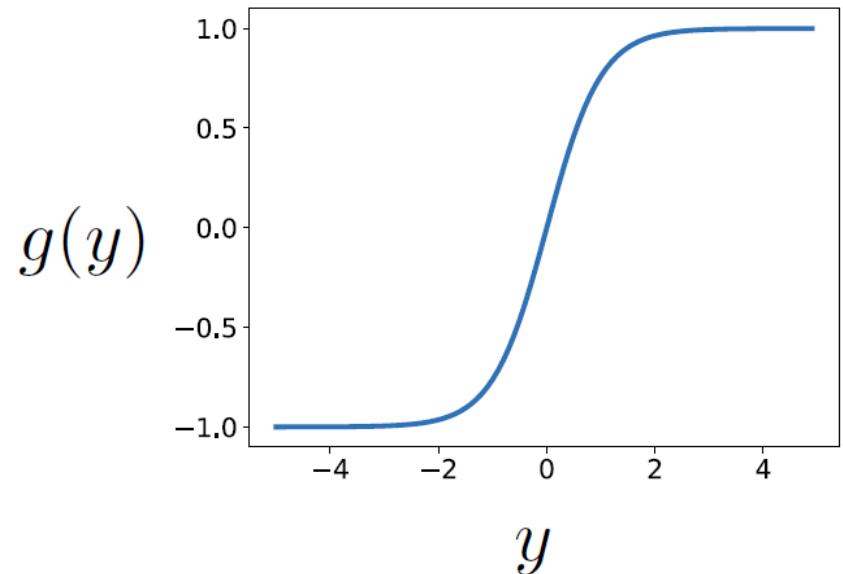


Computation in a neural network - nonlinearity

- Bounded between [-1,+1]
- Saturation for large +/- inputs
- Gradients go to zero
- Output centered at 0
- Preferable to sigmoid
- $\tanh(x) = 2 \text{ sigmoid}(2x) - 1$

Tanh

$$g(y) = \frac{e^y - e^{-y}}{e^y + e^{-y}}$$



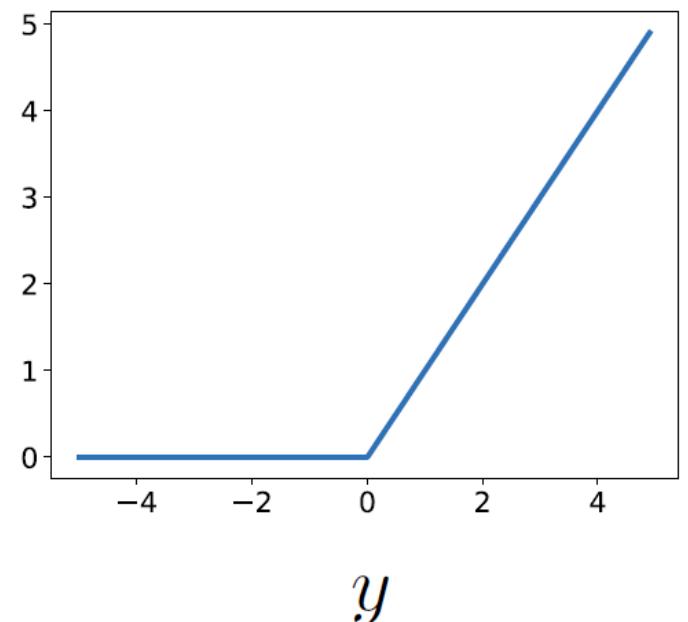
Computation in a neural network - nonlinearity

- unbounded output on positive side
- Efficient to implement
 - $\frac{\partial g}{\partial y} = \begin{cases} 0, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$
- Also seems to help convergence
- Drawback: if strongly in negative region, unit is dead forever (no gradient)
- Default choice: widely used in current models.

Rectified linear unit (ReLU)

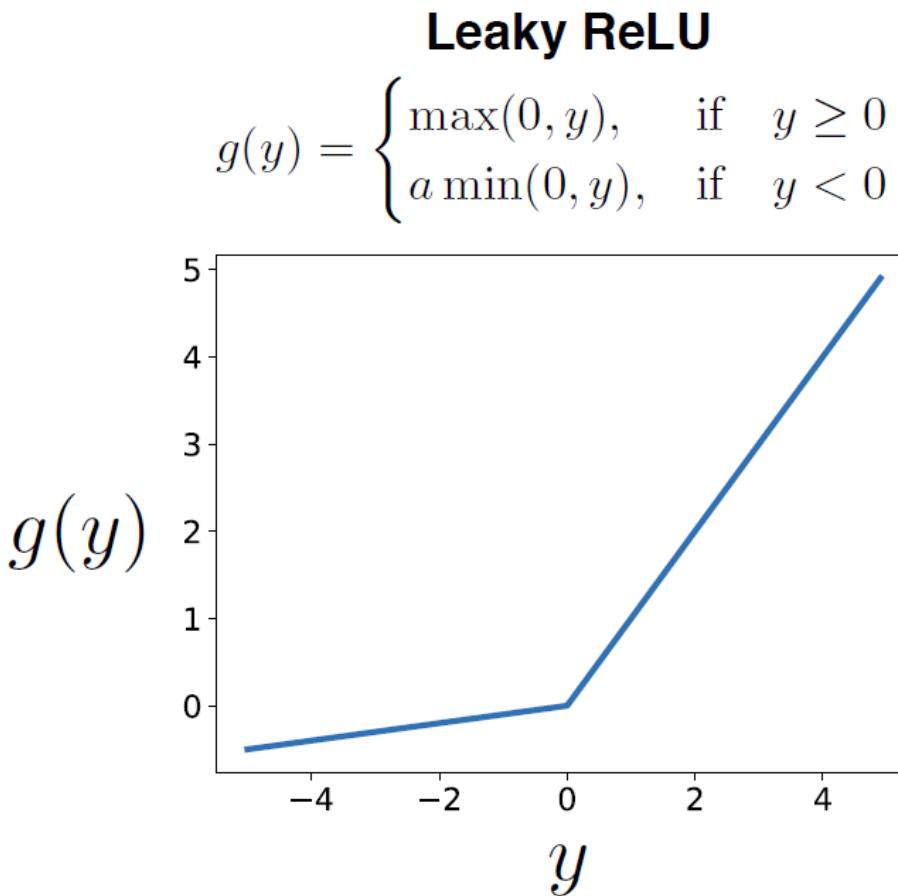
$$g(y) = \max(0, y)$$

$$g(y)$$

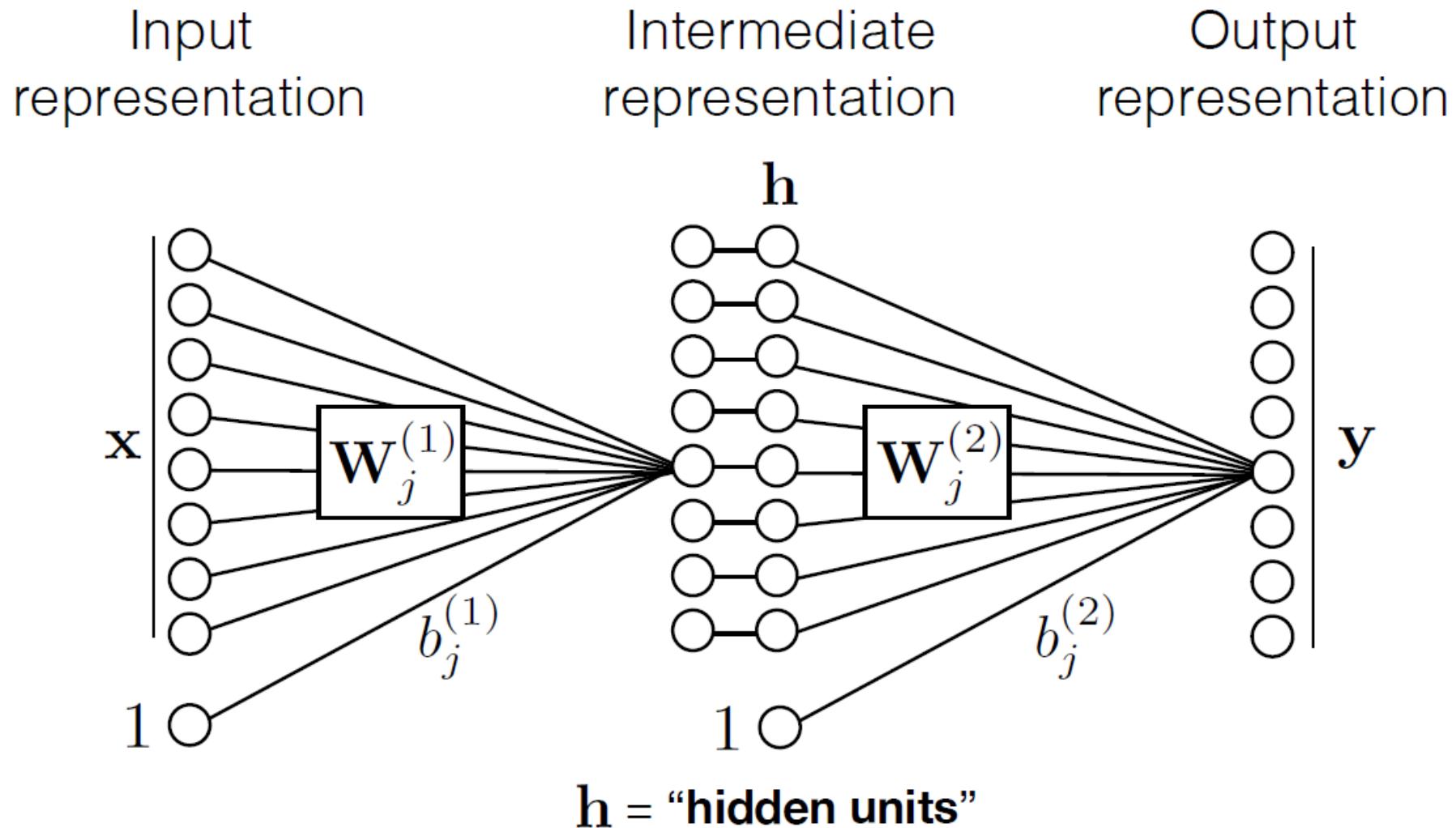


Computation in a neural network - nonlinearity

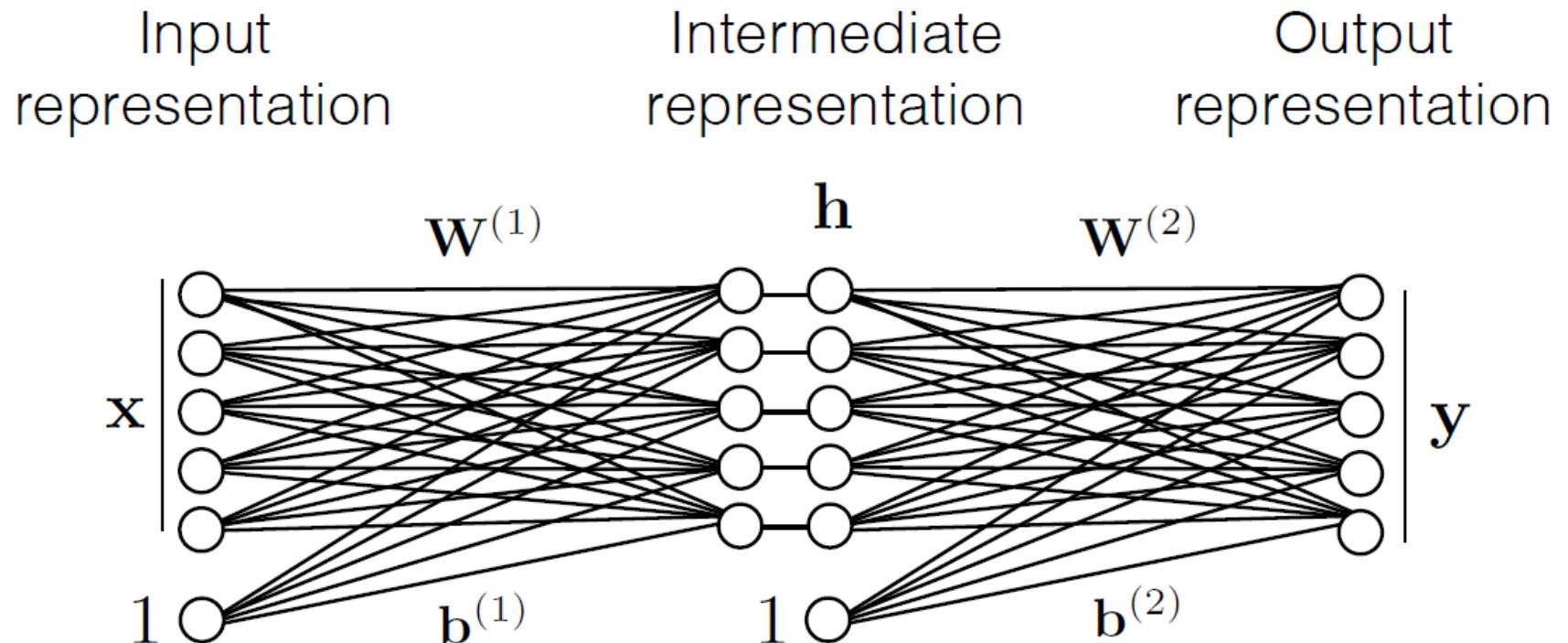
- Where a is small (e.g. 0.02)
- Efficient to implement
 - $\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$
- Also known as probabilistic Relu (PReLU)
- Has non-zero gradients everywhere (unlike ReLU)
- a can also be learned (see Kaiming He et al. 2015)



Stacking layers



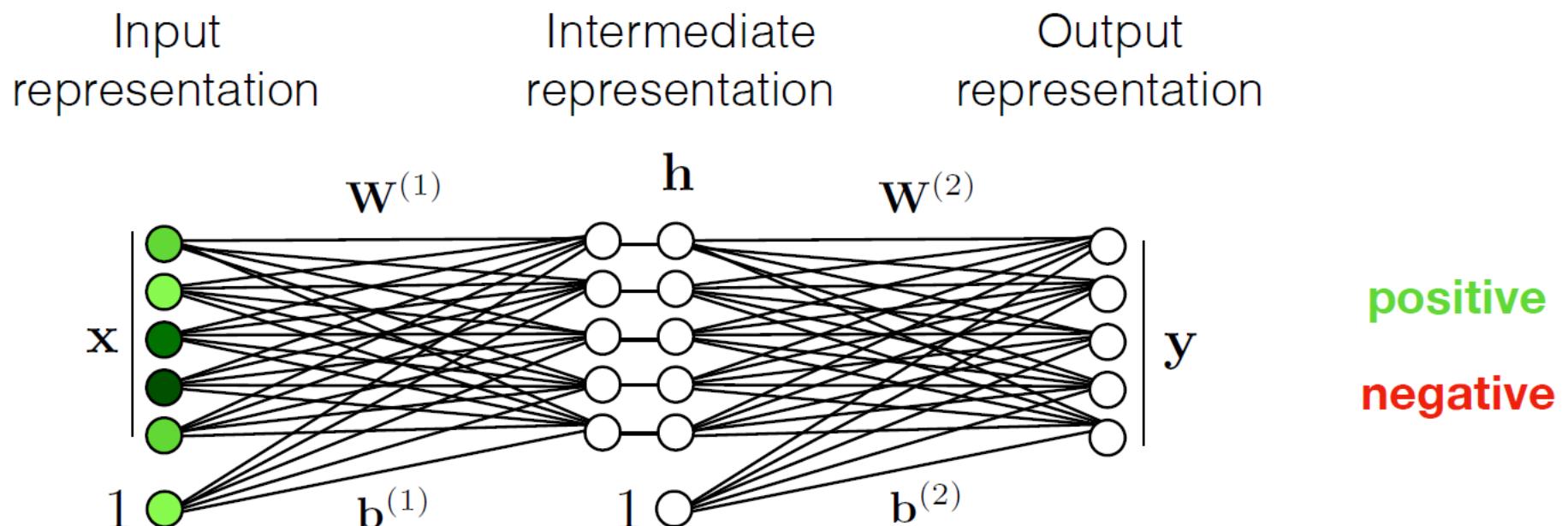
Stacking layers



$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \quad \mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

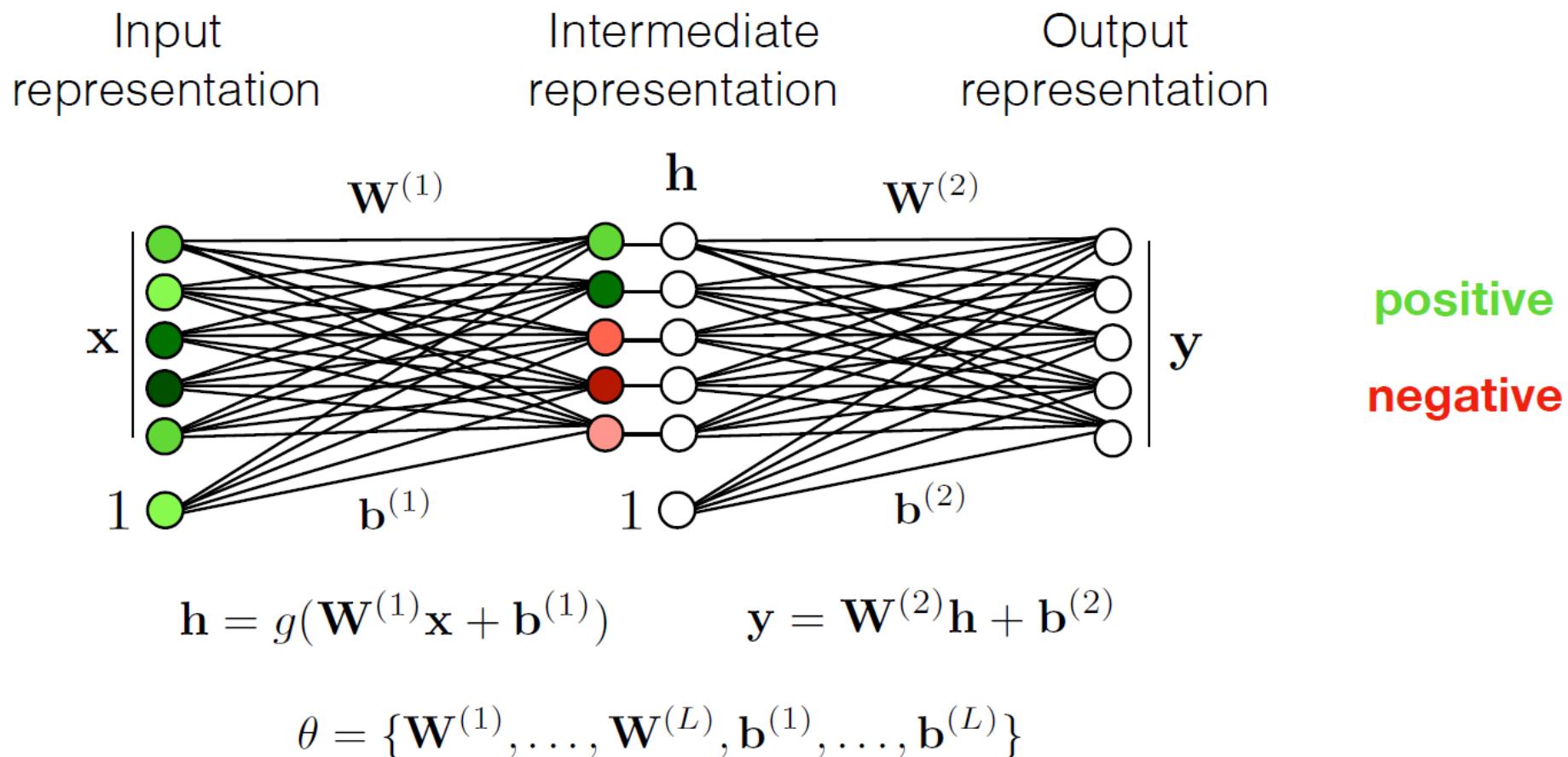
Stacking layers



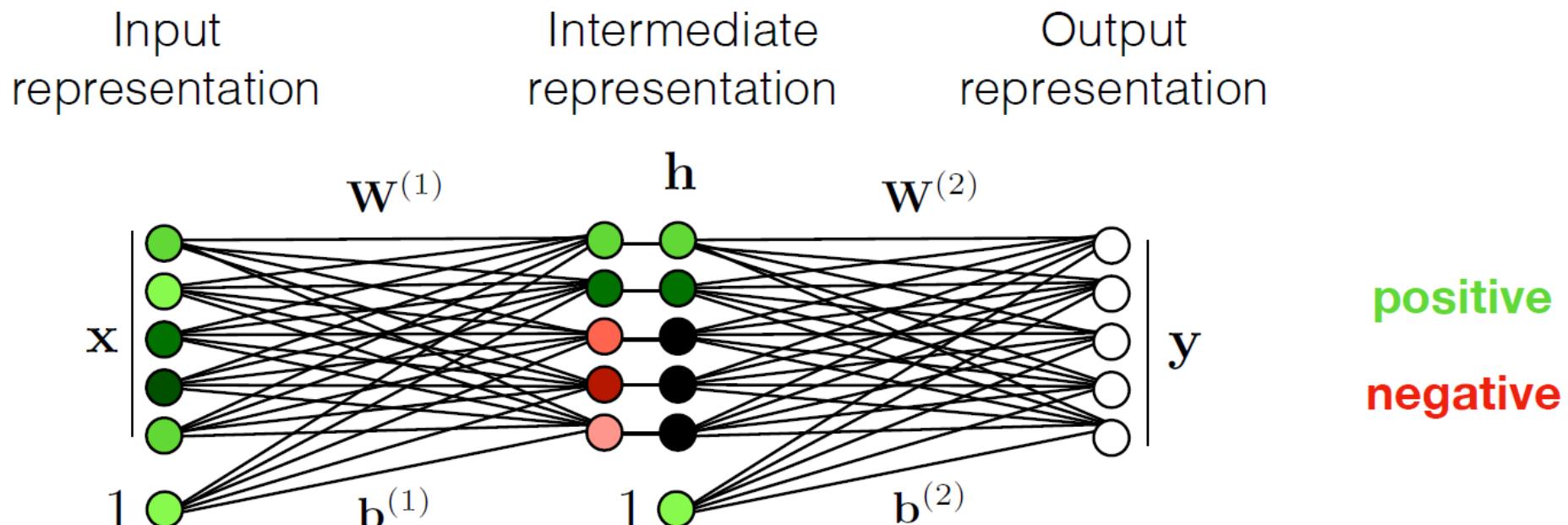
$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \quad \mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Stacking layers



Stacking layers

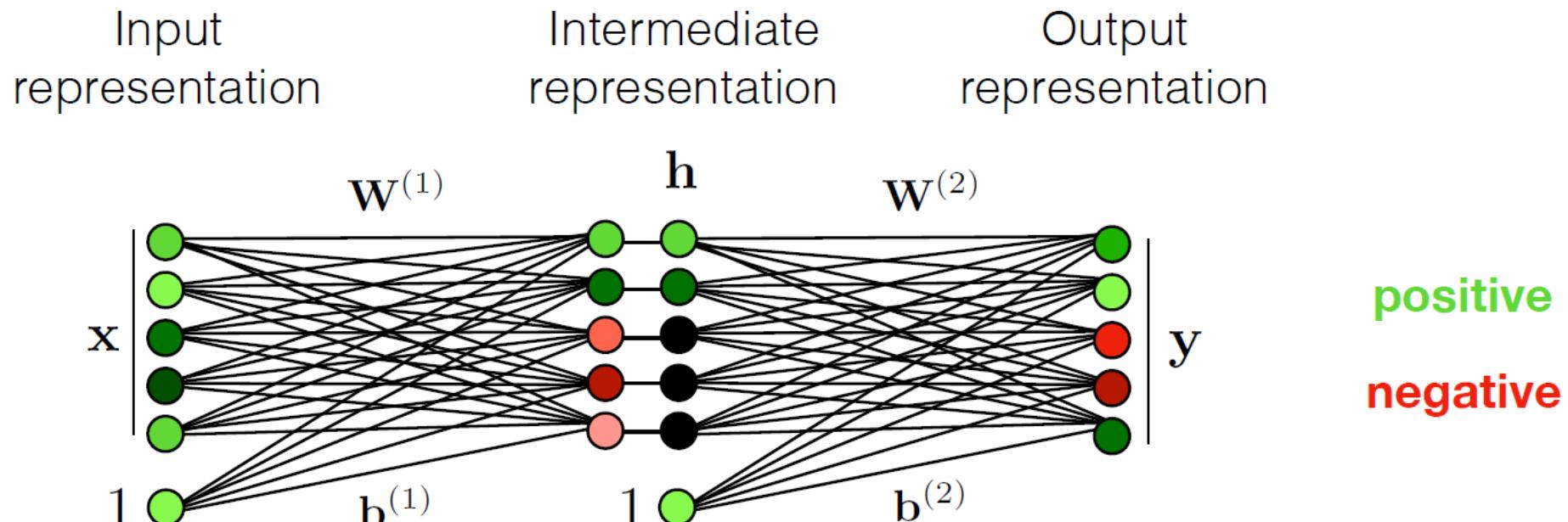


$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Stacking layers

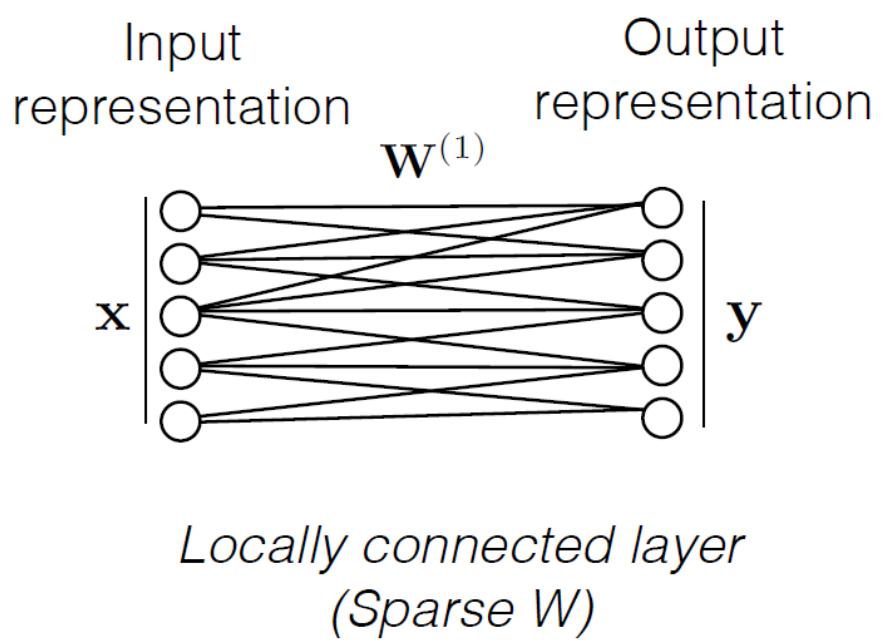
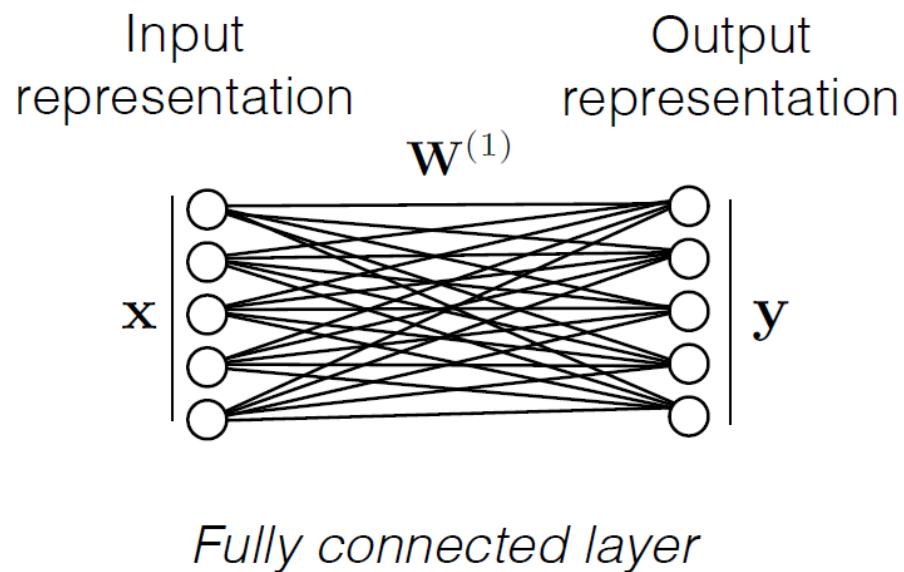


$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

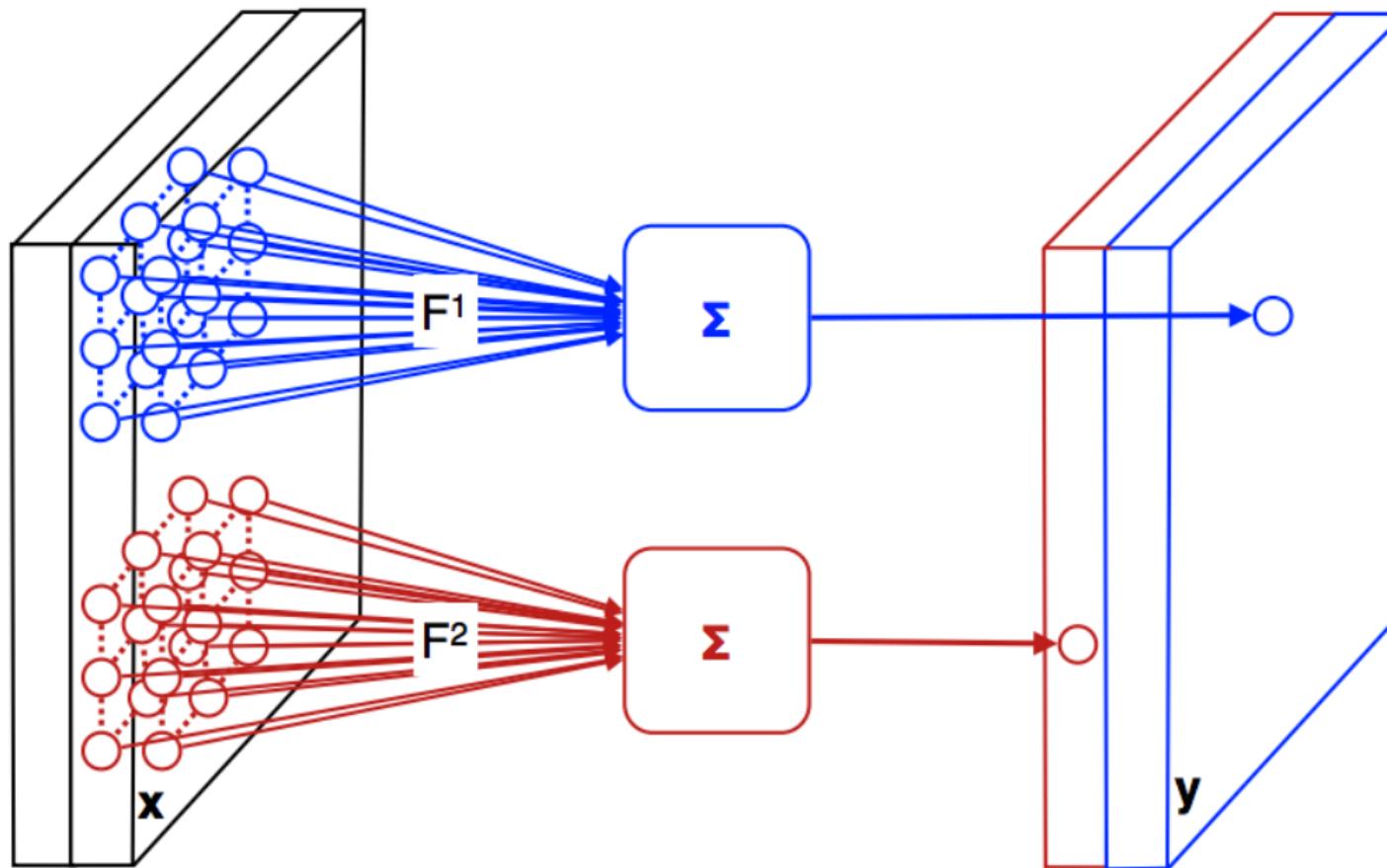
$$\mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Connectivity patterns



2-dimensional
input representation

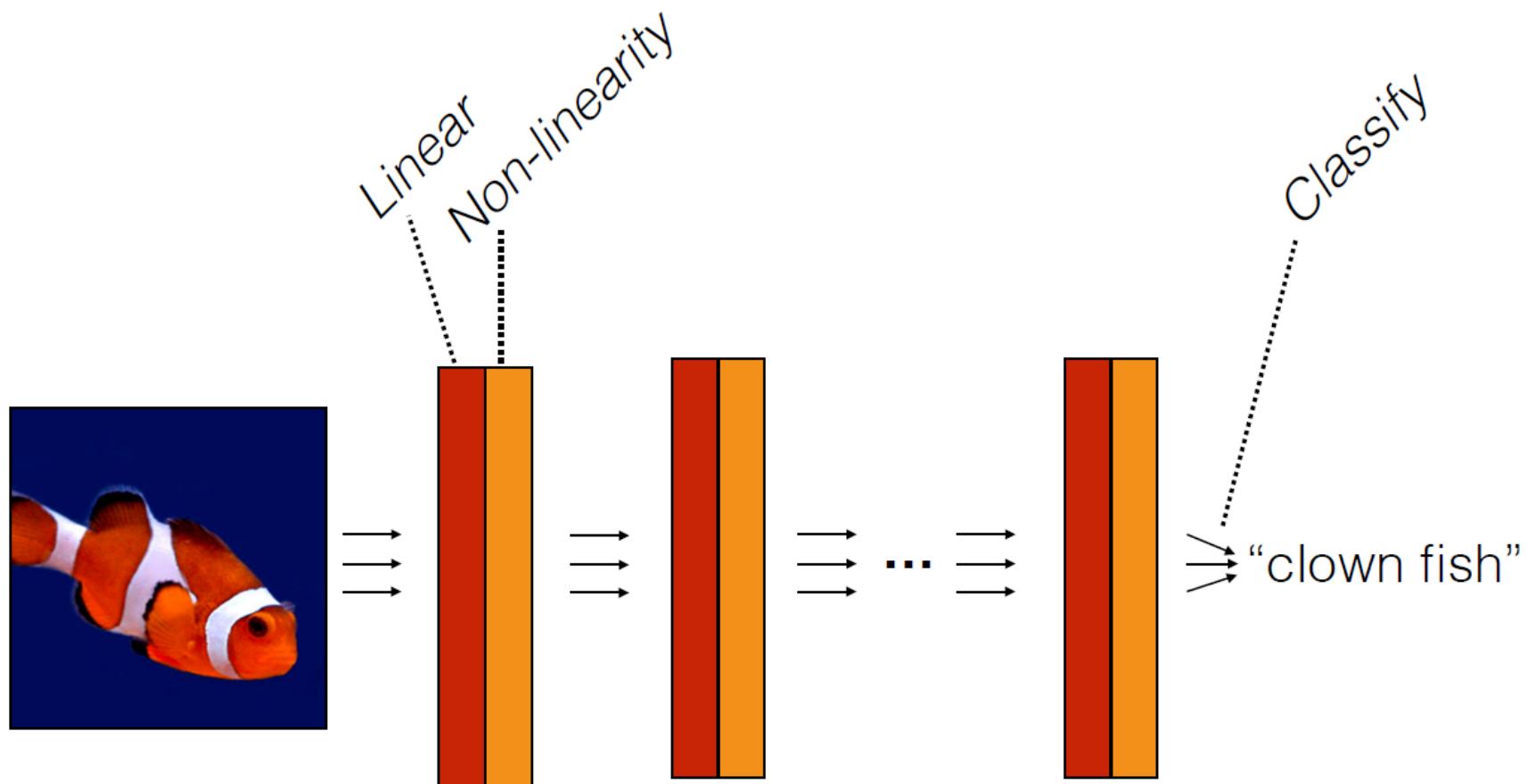


2-dimensional output
representation

$$\mathbb{R}^{H \times W \times C^{(l)}} \rightarrow \mathbb{R}^{H \times W \times C^{(l+1)}}$$

[Figure from Andrea Vedaldi]

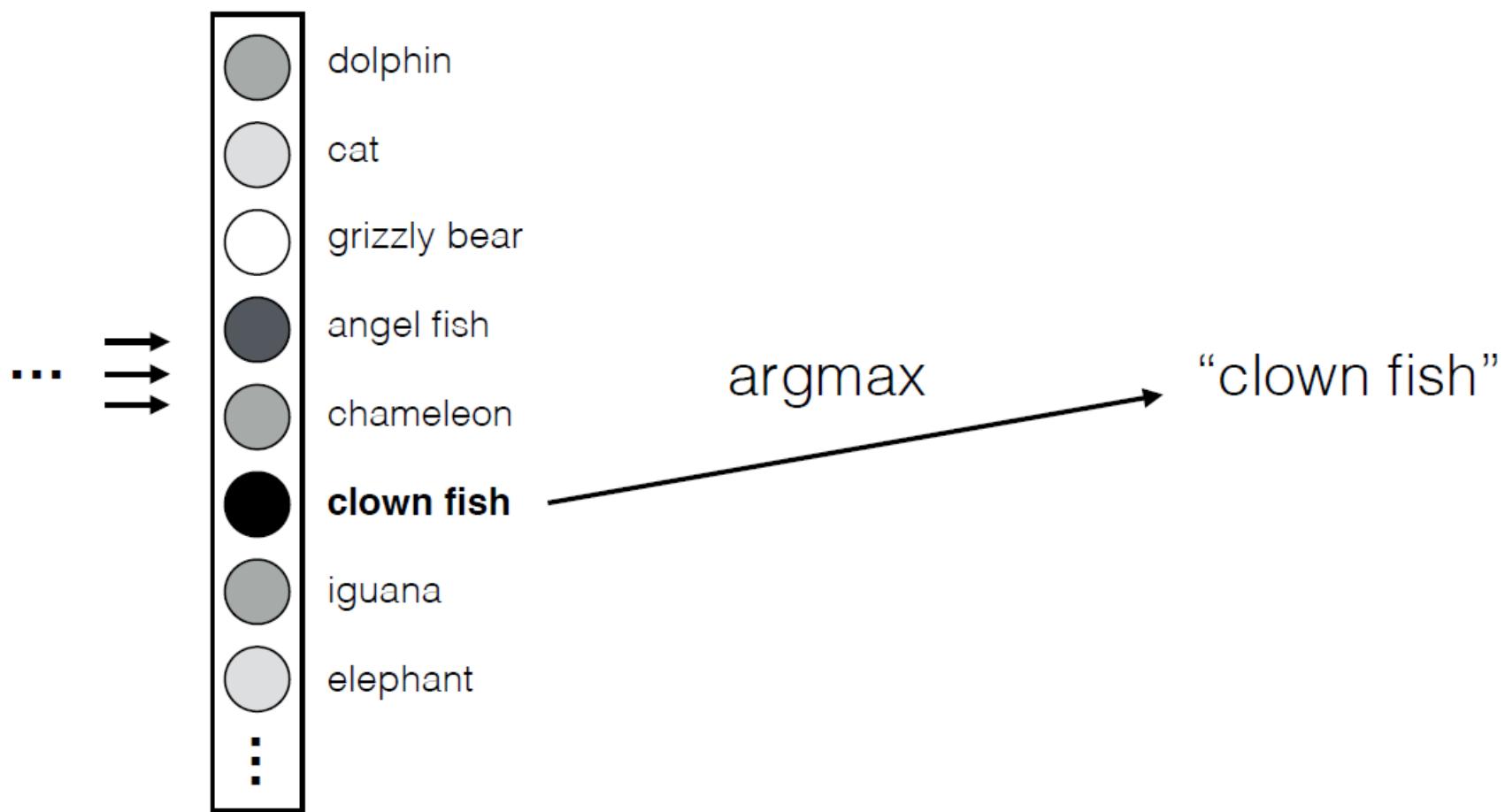
Deep nets



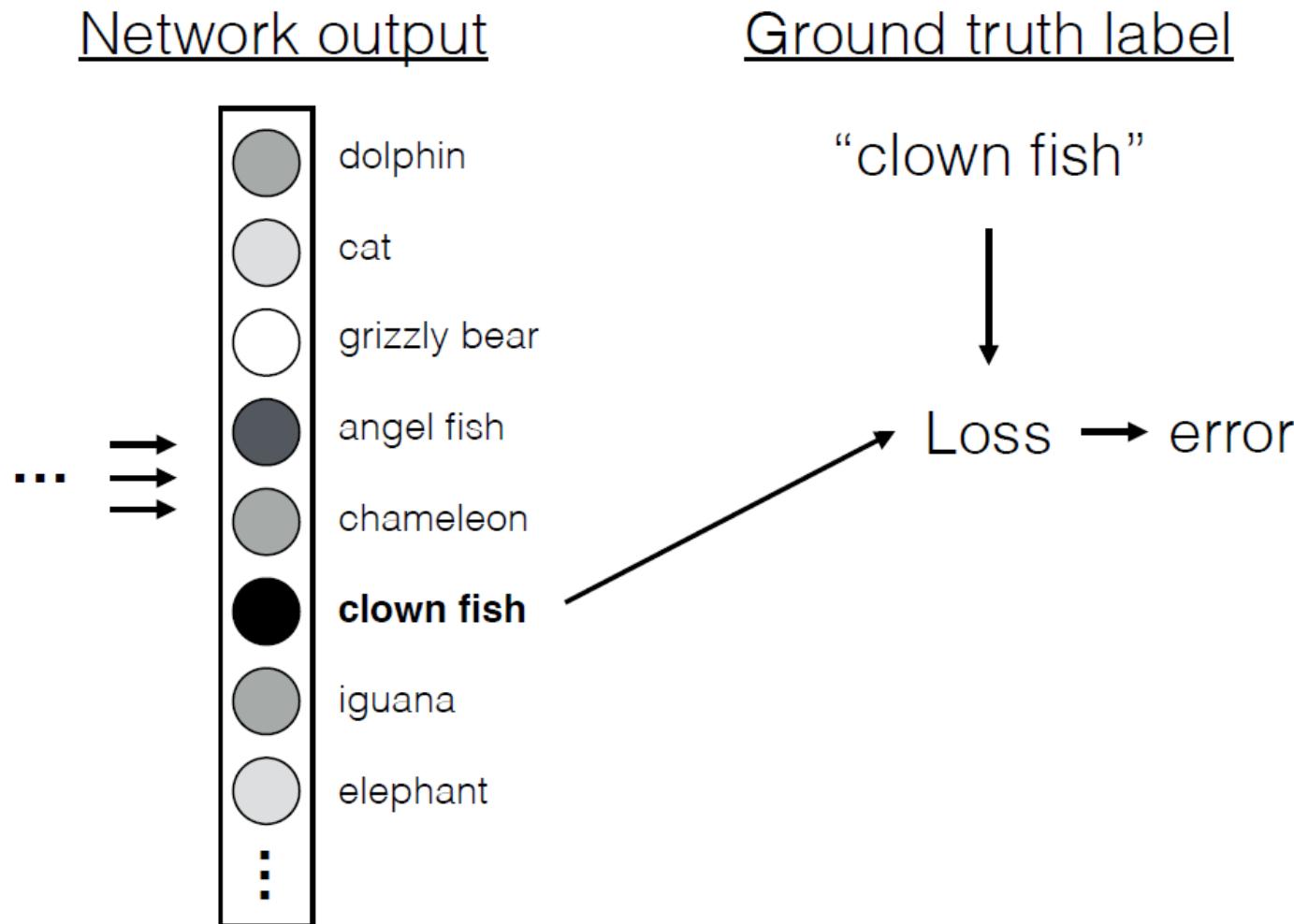
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Classifier layer

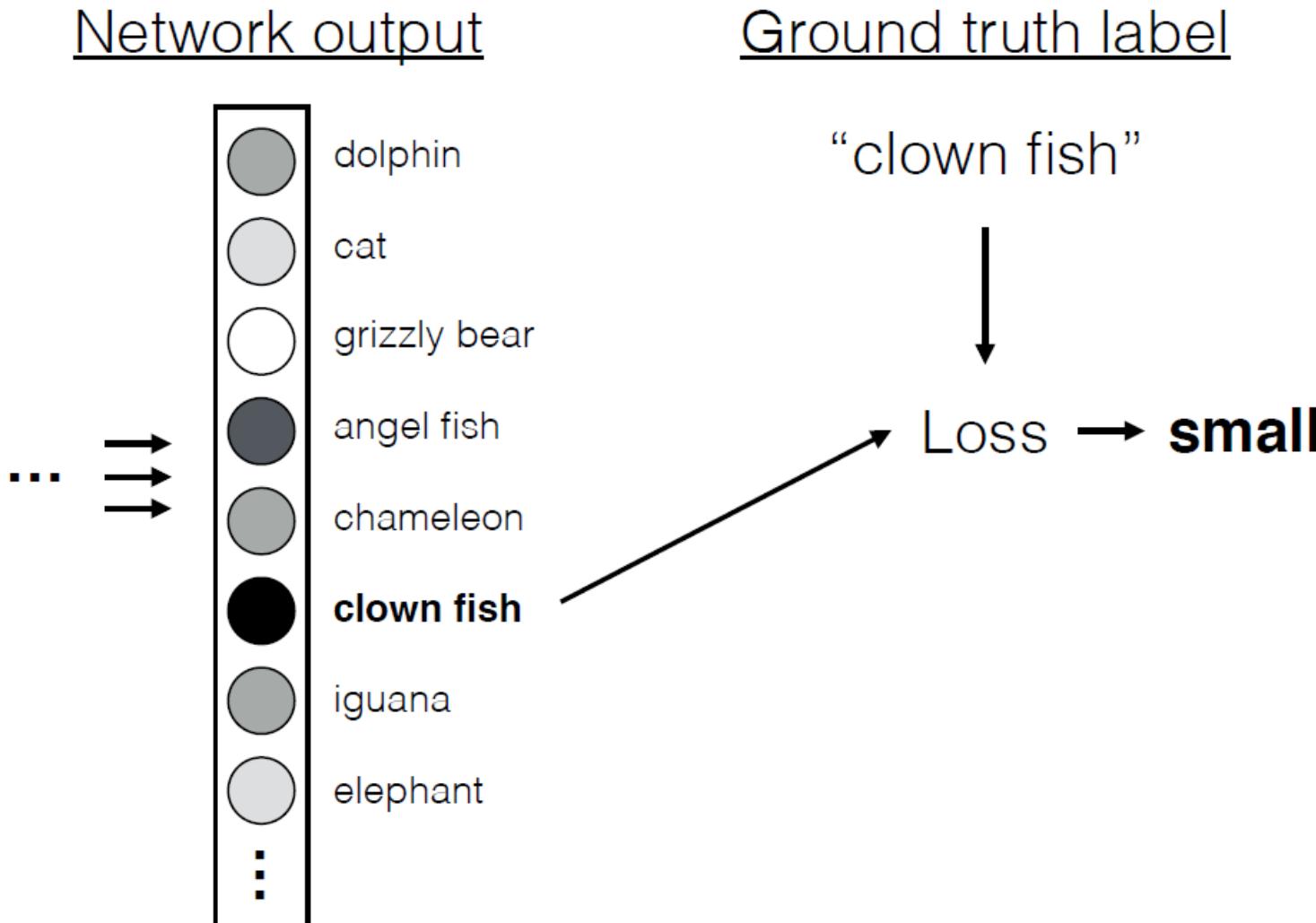
Last layer



Loss function



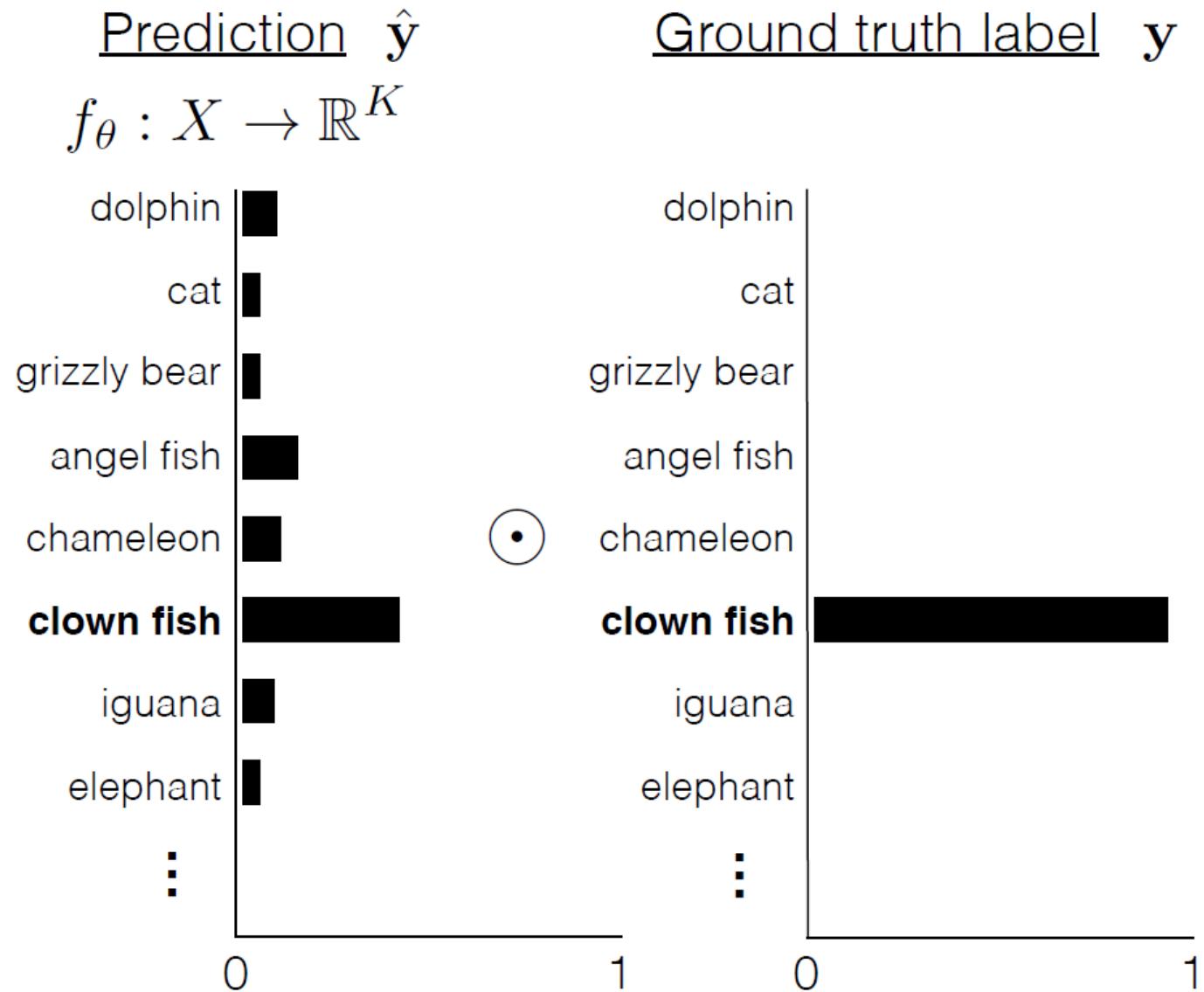
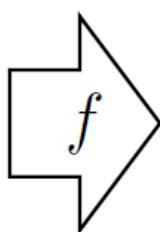
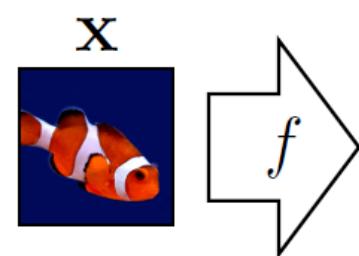
Loss function



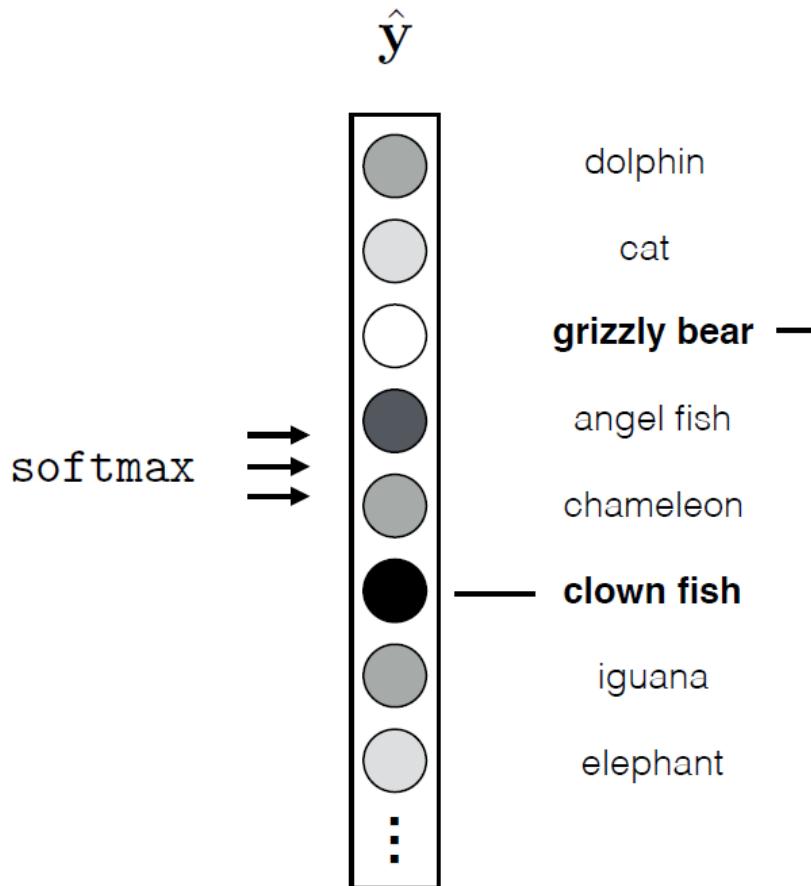
Loss function

Network output





Network output



Ground truth label

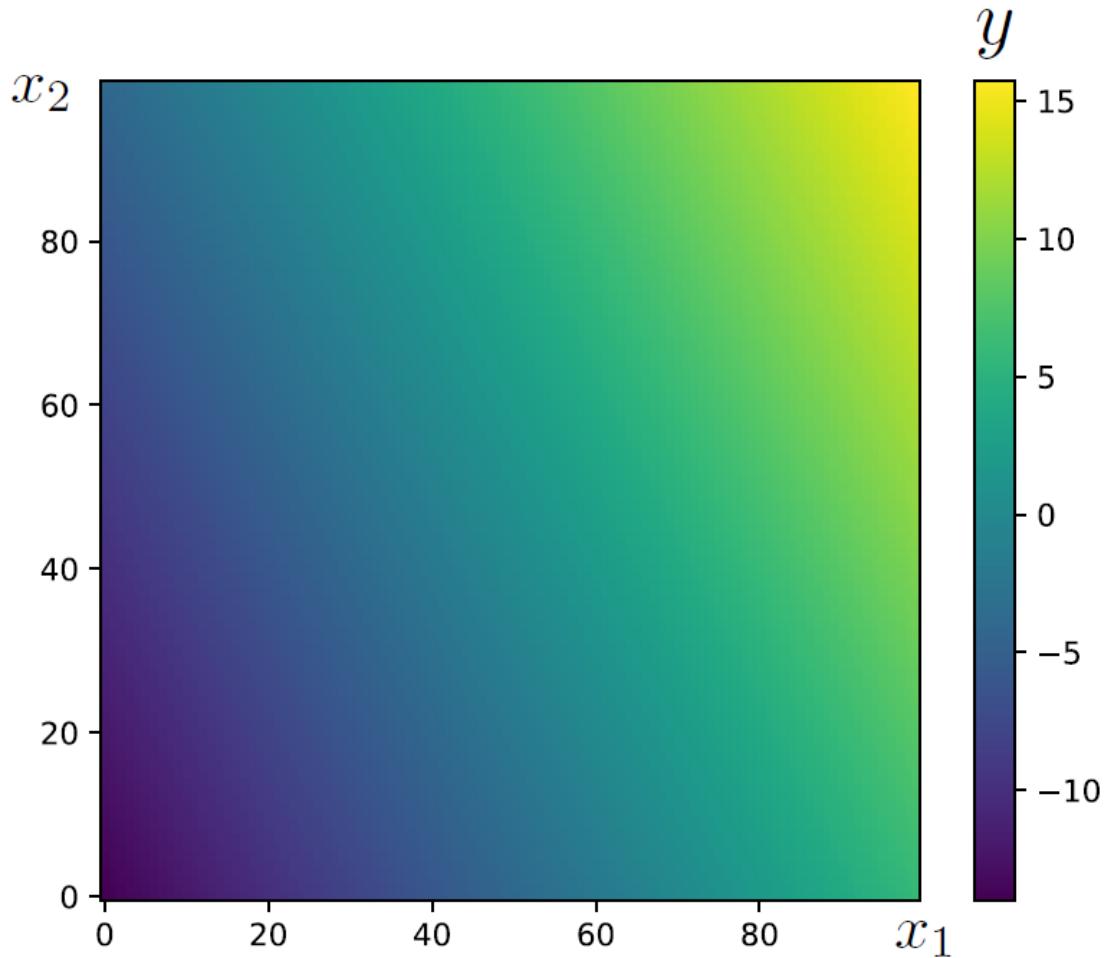
Probability of the observed data under the model

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

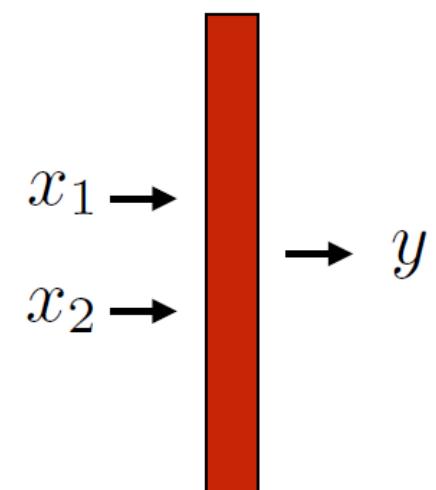
Representational power

- 1 layer? Linear decision surface.
- 2+ layers? In theory, can represent any function.
Assuming non-trivial non-linearity.
 - Bengio 2009,
<http://www.iro.umontreal.ca/~bengioy/papers/fml.pdf>
 - Bengio, Courville, Goodfellow book
<http://www.deeplearningbook.org/contents/mlp.html>
 - Simple proof by M. Nielsen
<http://neuralnetworksanddeeplearning.com/chap4.html>
 - D. Mackay book
<http://www.inference.phy.cam.ac.uk/mackay/itprnn/ps/482.491.pdf>
- But issue is efficiency: very wide two layers vs narrow deep model? In practice, more layers helps.

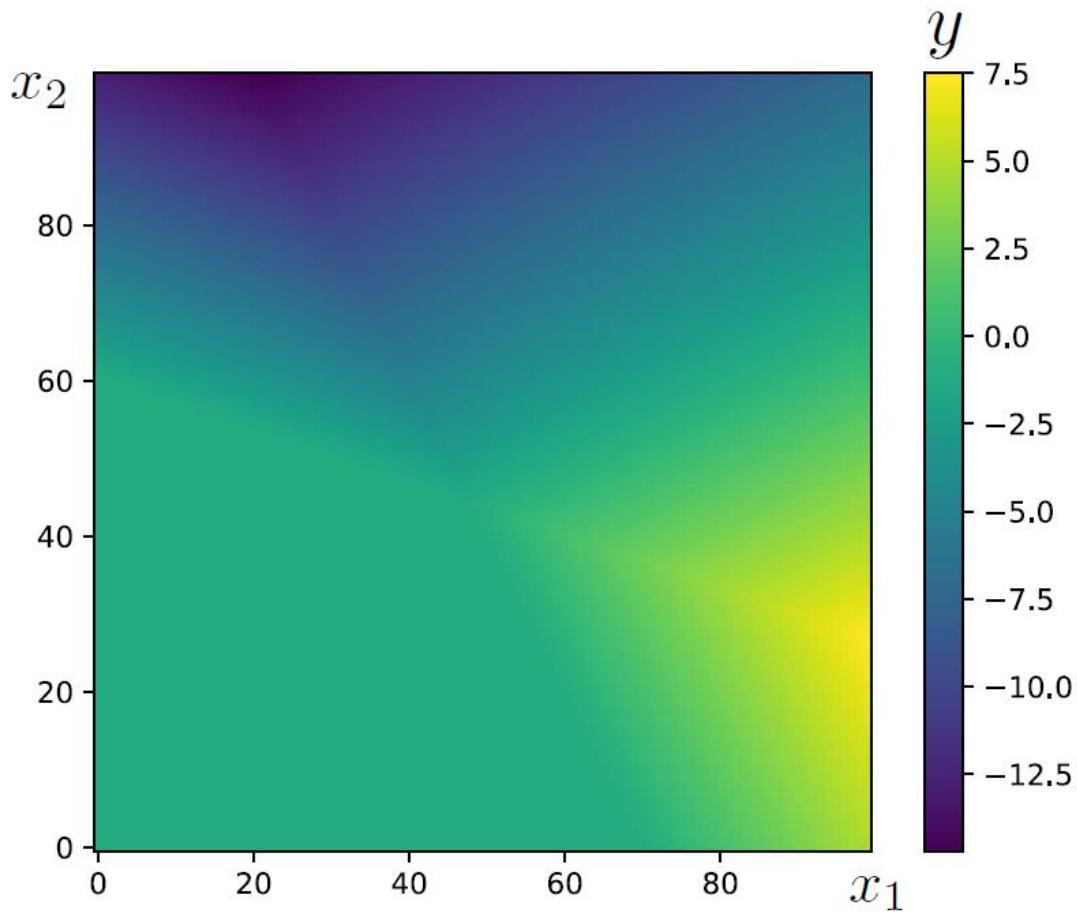
Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

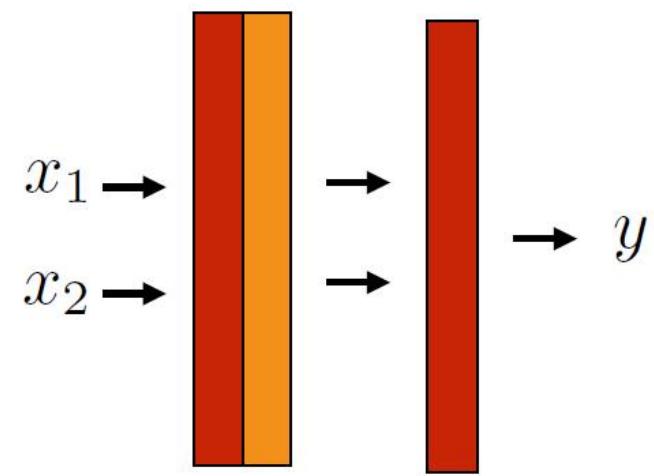


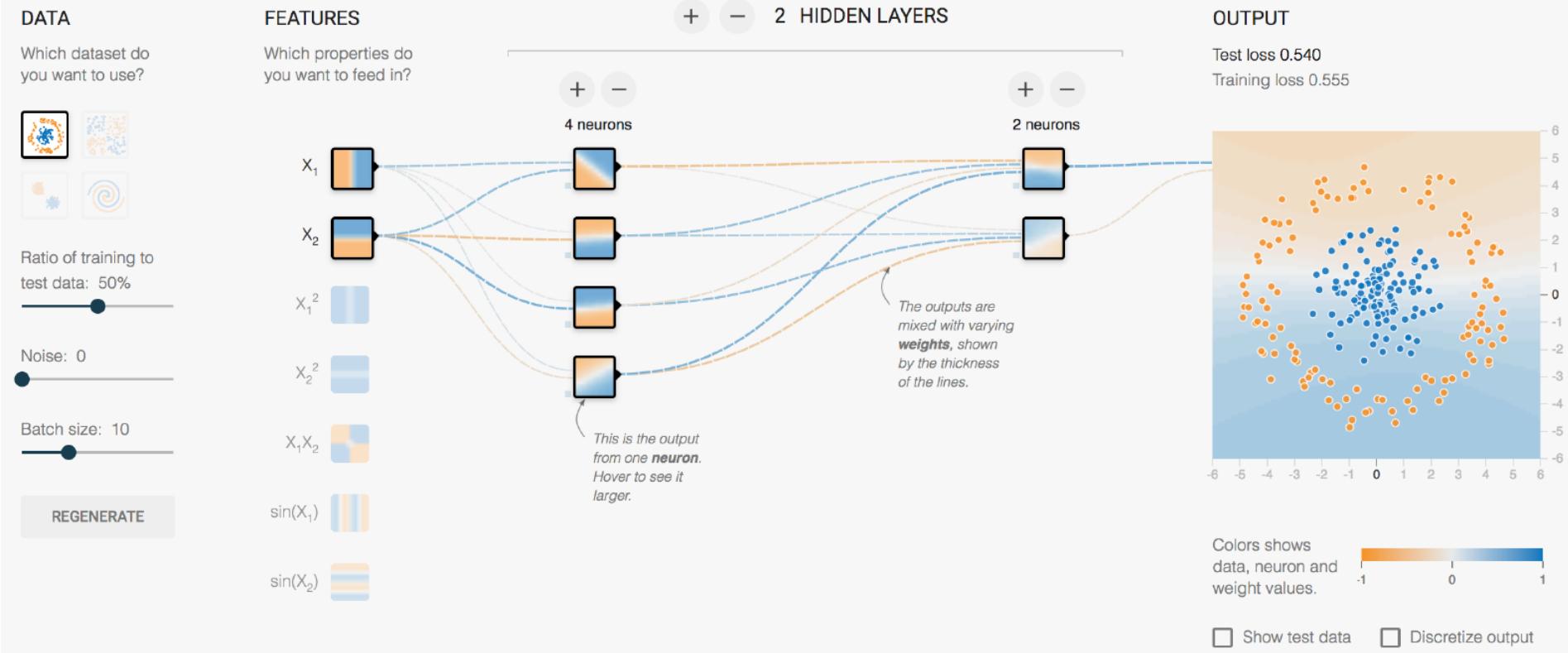
Example: nonlinear classification with a deep net



$$\mathbf{h} = g(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})$$

$$y = \mathbf{W}^{(2)} \mathbf{h} + b^{(2)}$$

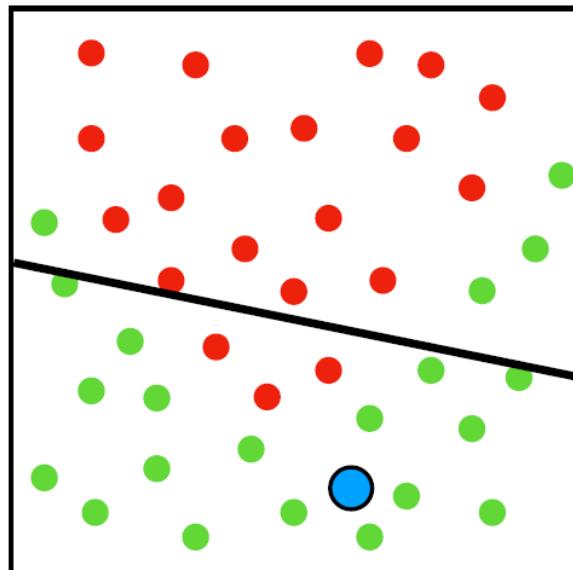




[<http://playground.tensorflow.org>]

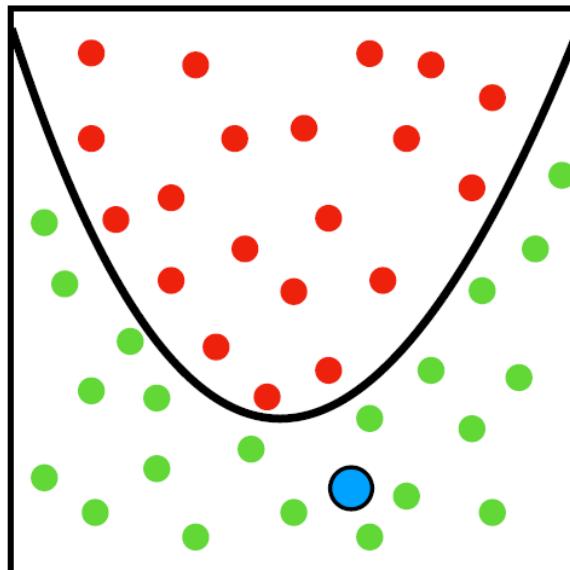
Example: nonlinear classification with a deep net

What class is ● ?



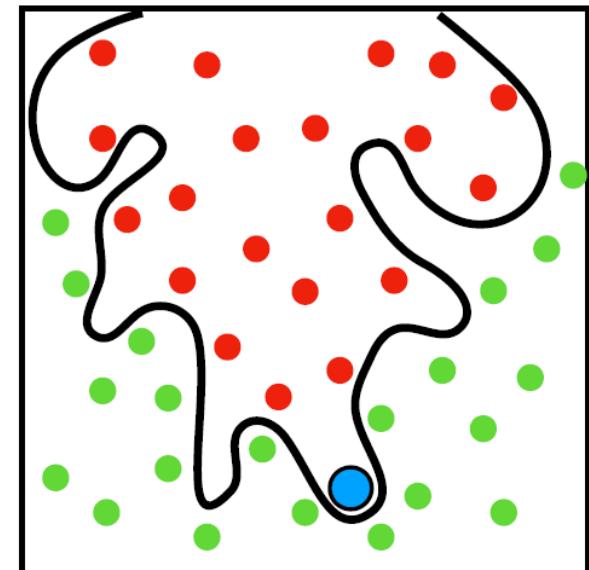
Answer: ●

Underfitting



Answer: ●

Appropriate model

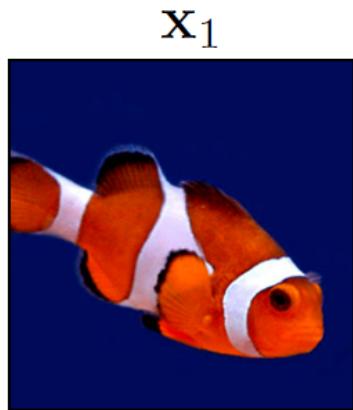


Answer: ●

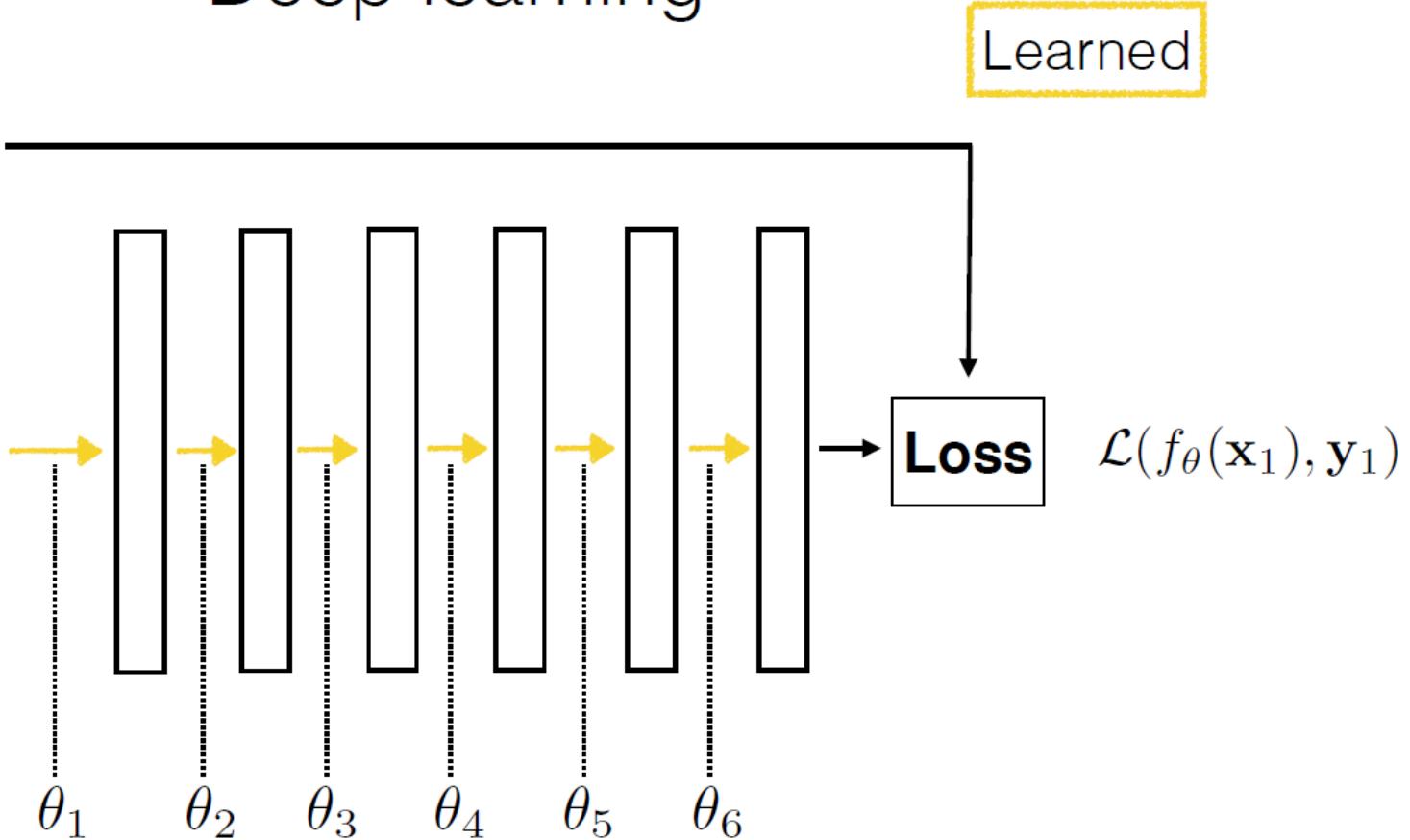
Overfitting

Deep learning

y_1
“clown fish”



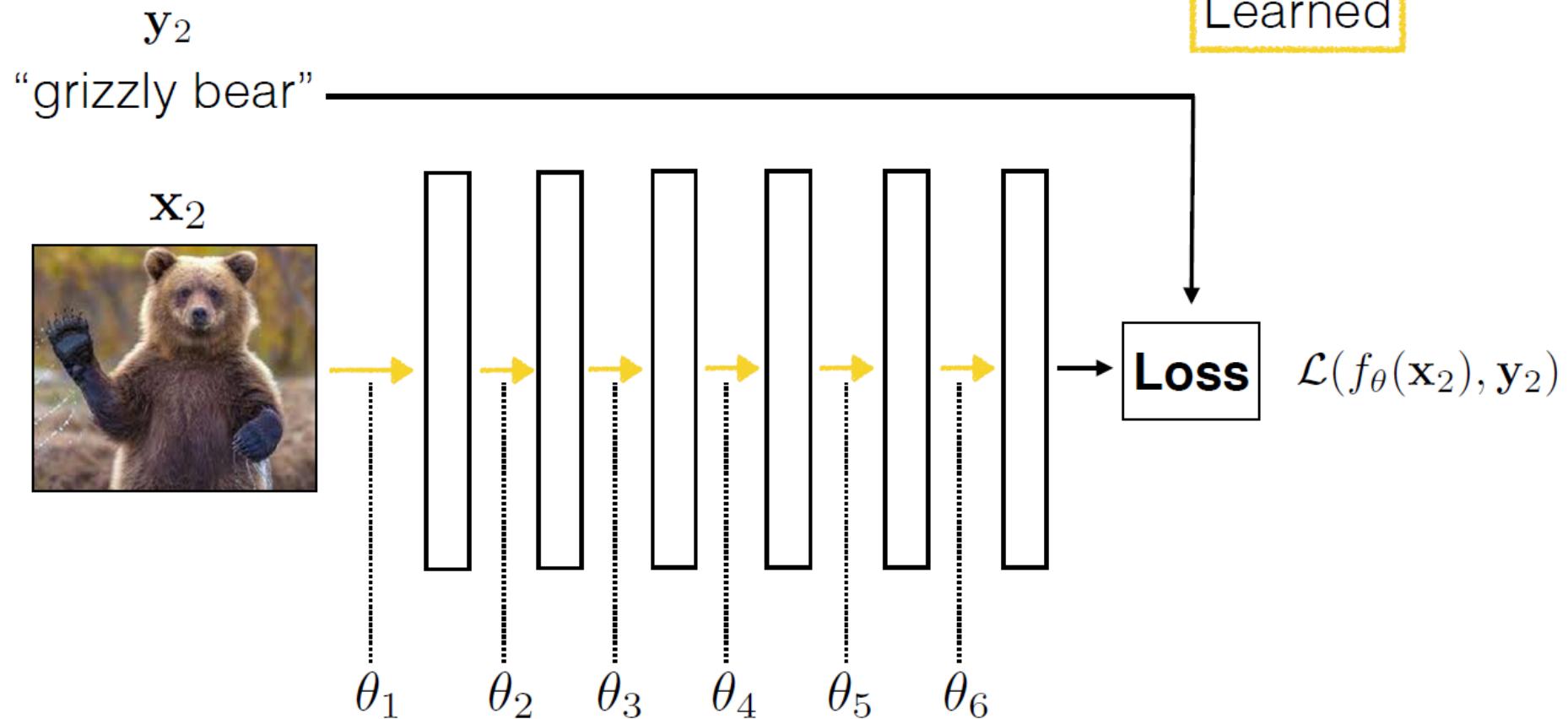
x_1



Learned

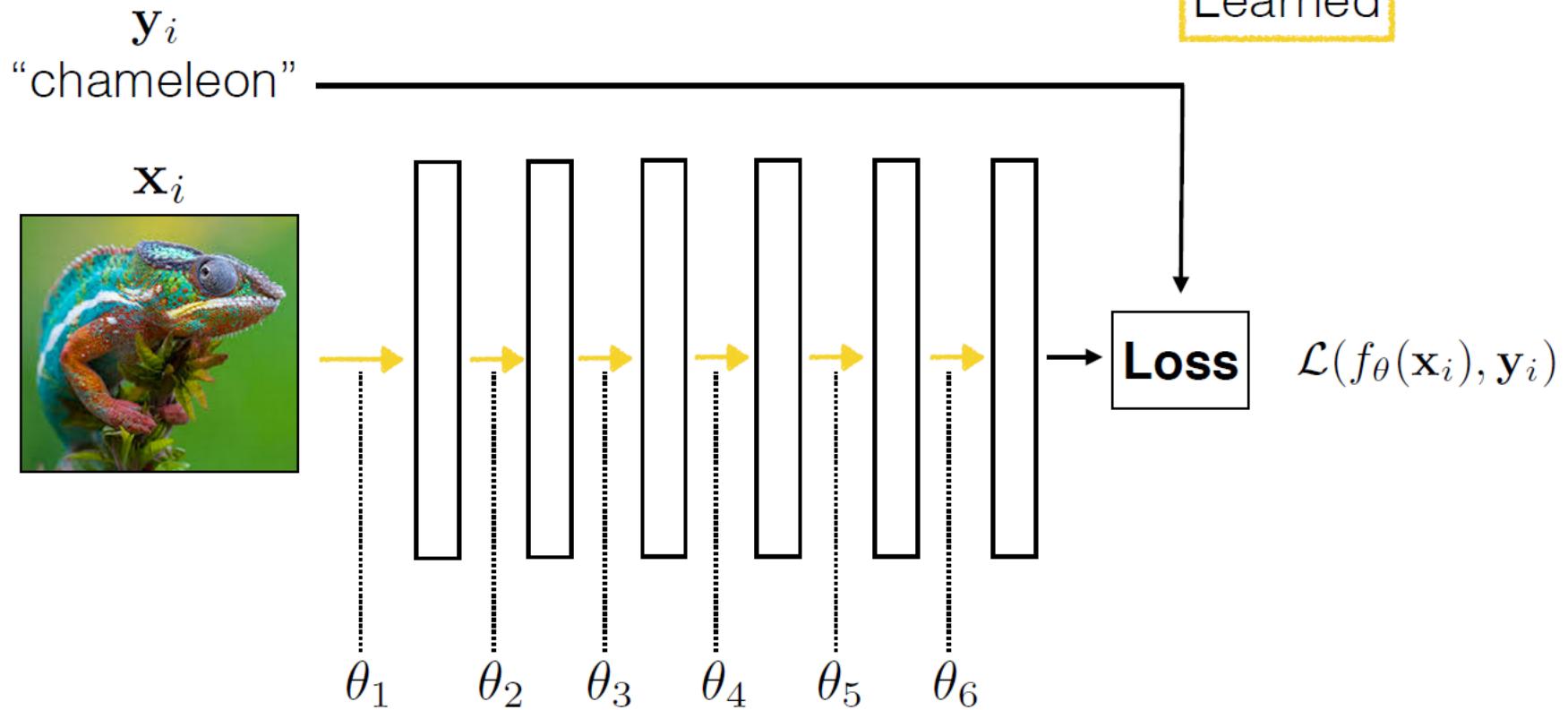
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

Deep learning



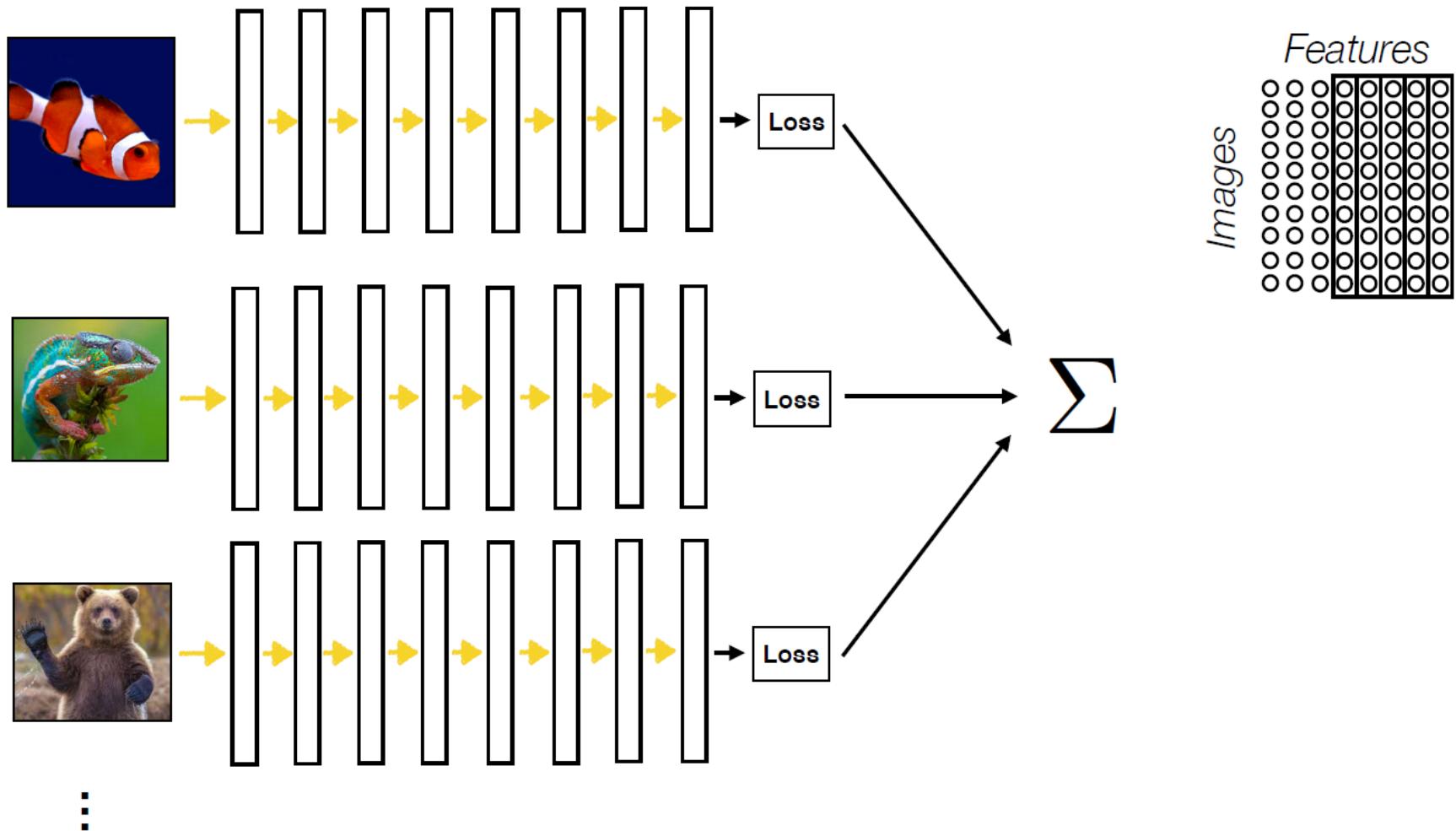
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

Deep learning

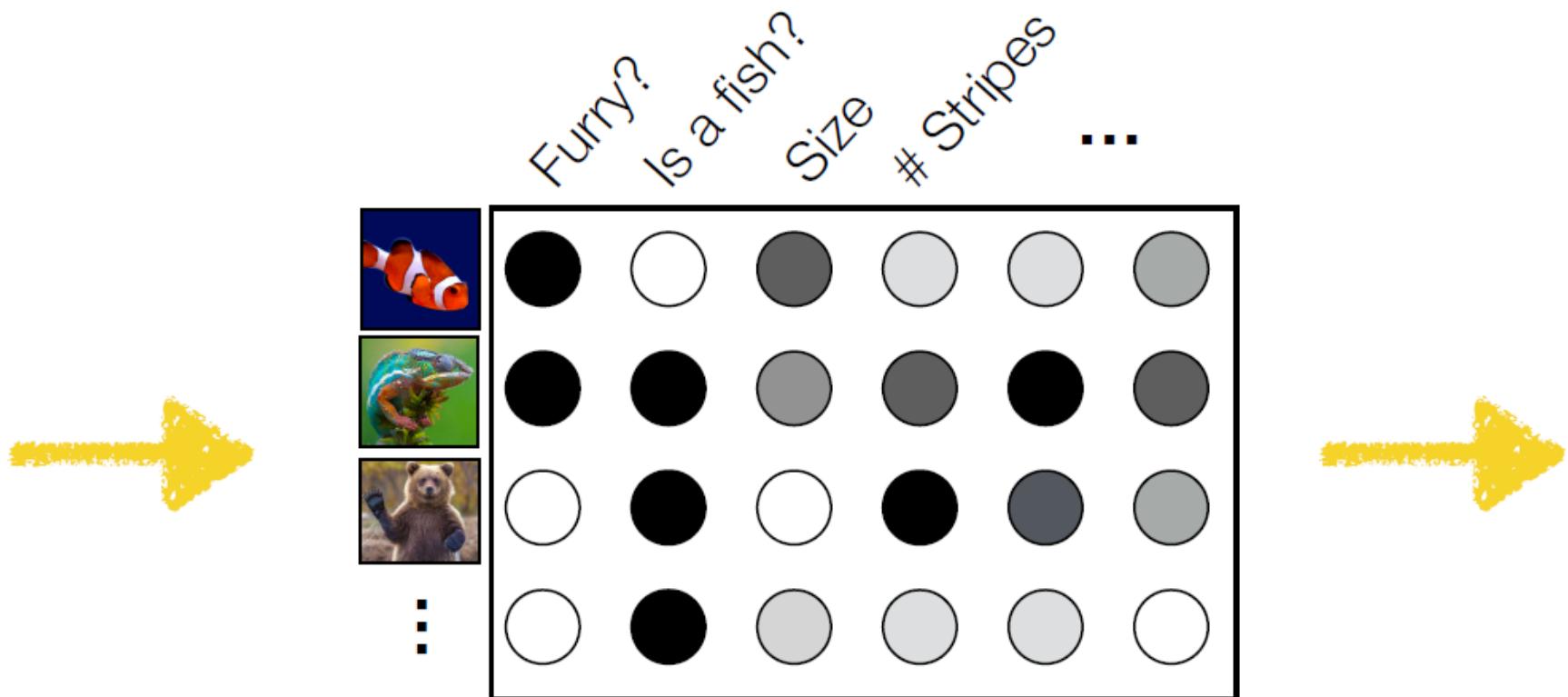


$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

Batch (parallel) processing



Tensors (multi-dimensional arrays)



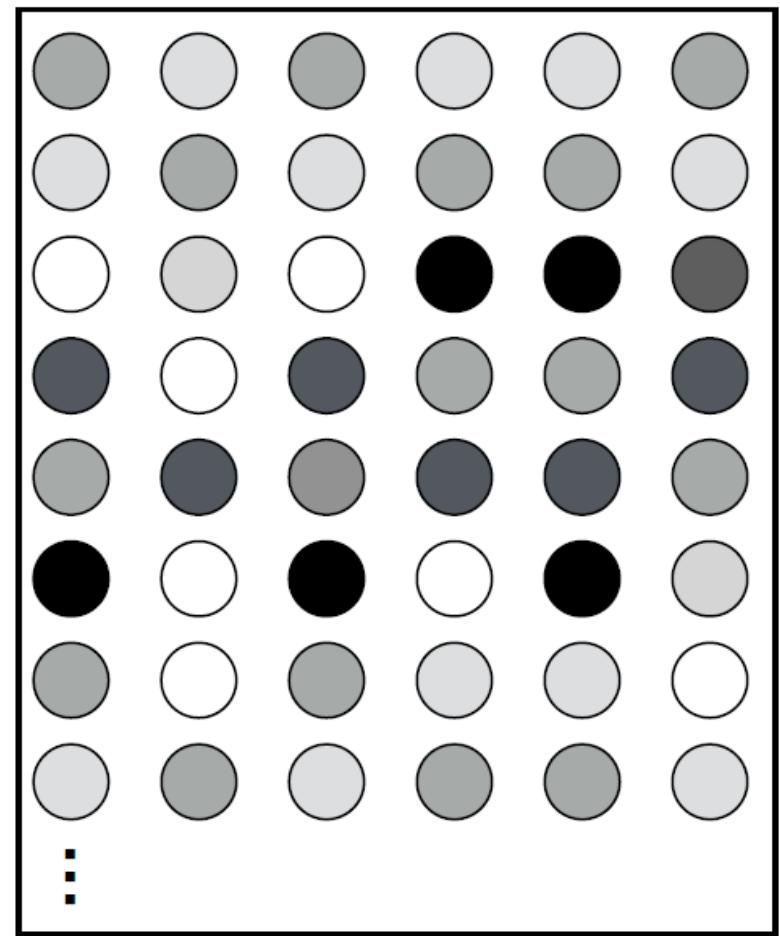
Each layer is a representation of the data

Tensors (multi-dimensional arrays)

$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\text{batch}} \times C^{(1)}}$

neurons
features
units
“channels”

N_{batch}



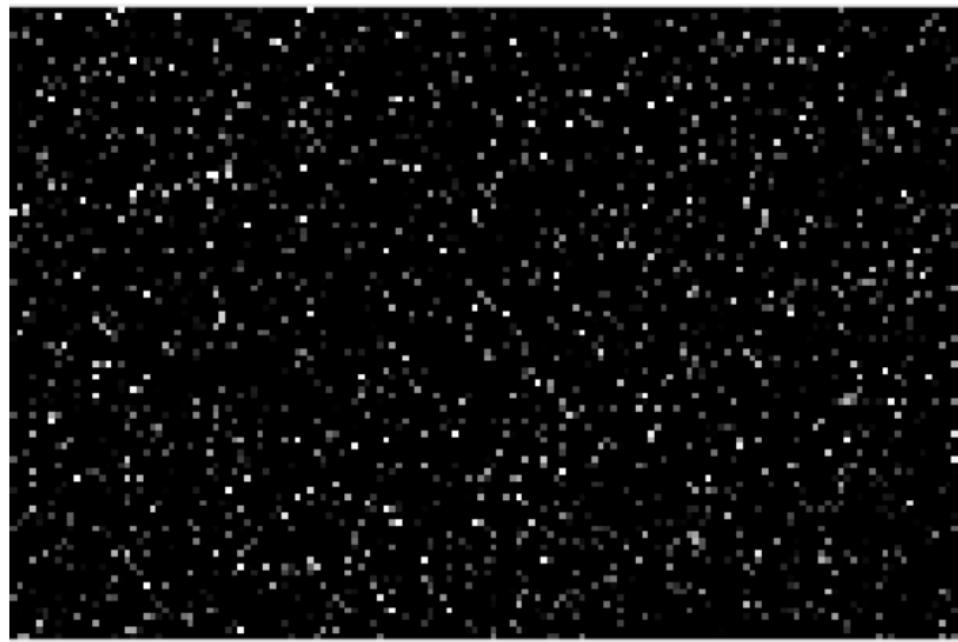
Tensors (multi-dimensional arrays)

$$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\text{batch}} \times C^{(1)}}$$

neurons
features
units
“channels”

N_{batch}

$$C^{(1)}$$

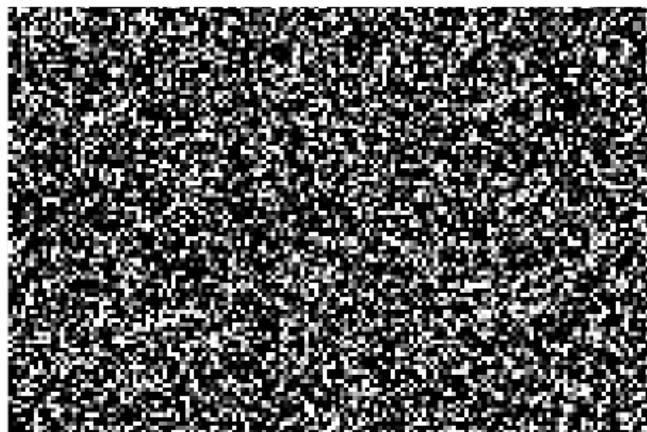


“Tensor flow”

$$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\text{batch}} \times C^{(1)}}$$

N_{batch}

$C^{(1)}$

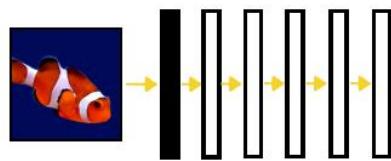


$$\mathbf{h}^{(2)} \in \mathbb{R}^{N_{\text{batch}} \times C^{(2)}}$$

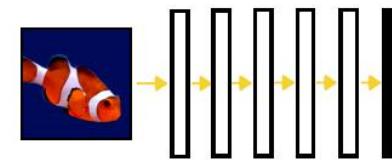
N_{batch}

$C^{(2)}$

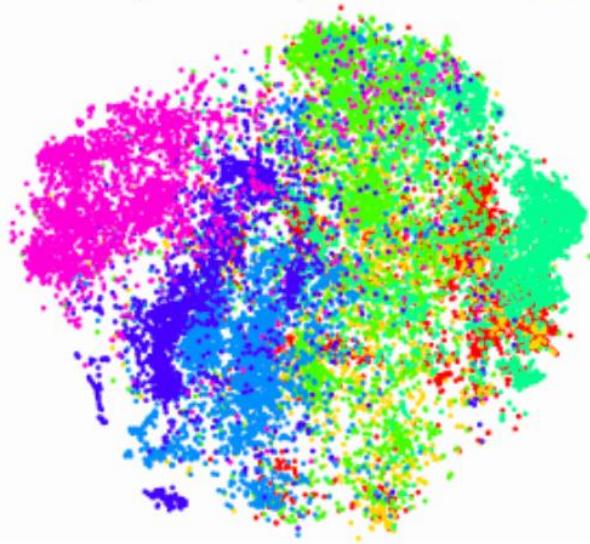
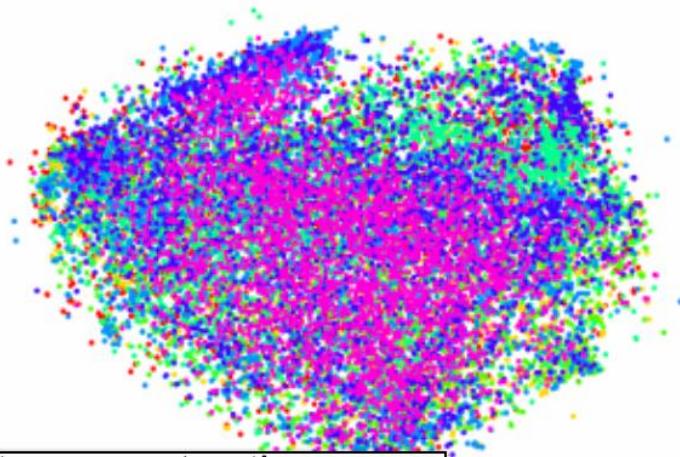




Layer 1 representation



Layer 6 representation



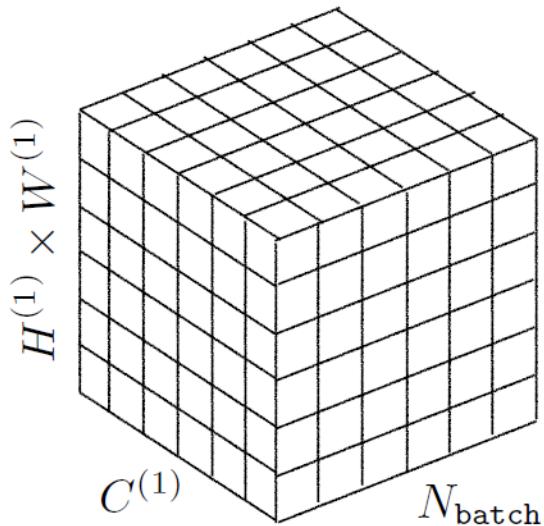
- structure, construction
- covering
- commodity, trade good, good
- conveyance, transport
- invertebrate
- bird
- hunting dog

[DeCAF, Donahue, Jia, et al. 2013]

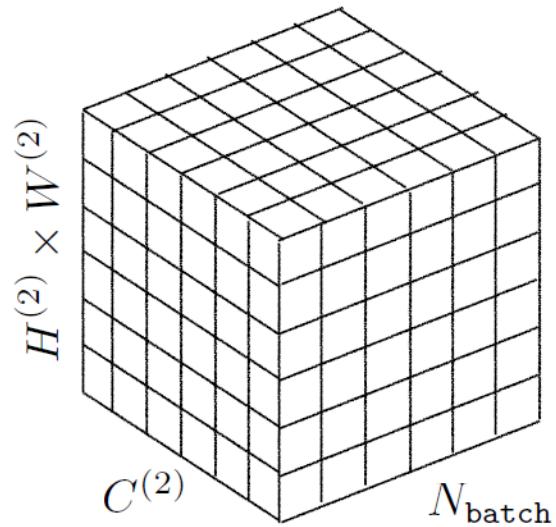
[Visualization technique : t-SNE, van der Maaten & Hinton, 2008]

“Tensor flow”

$$\mathbf{h}^{(1)} \in \mathbb{R}^{N_{\text{batch}} \times H^{(1)} \times W^{(1)} \times C^{(1)}}$$



$$\mathbf{h}^{(2)} \in \mathbb{R}^{N_{\text{batch}} \times H^{(2)} \times W^{(2)} \times C^{(2)}}$$



Regularization

Regularizing deep nets

Deep nets have millions of parameters!

On many datasets, it is easy to overfit — we may have more free parameters than data points to constrain them.

How can we regularize to prevent the network from overfitting?

1. Fewer neurons, fewer layers
2. Weight decay
3. Dropout
4. Normalization layers
5. ...

Regularized least squares

$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

$$R(\theta) = \lambda \|\theta\|_2^2 \leftarrow \text{Only use polynomial terms if you really need them! Most terms should be zero}$$

ridge regression, a.k.a., **Tikhonov regularization**

Probabilistic interpretation: R is a Gaussian **prior** over values of the parameters.

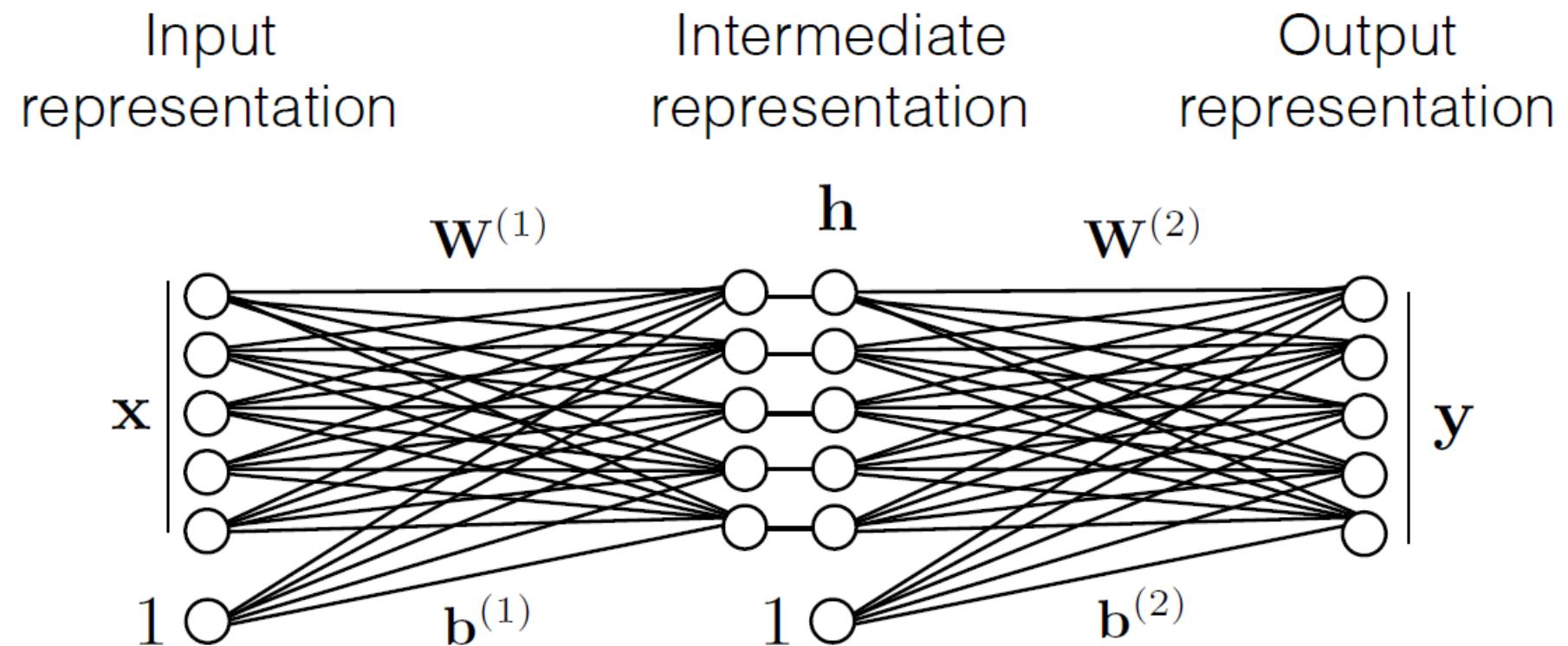
Regularizing the weights in a neural net

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + R(\theta)$$

$$R(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 \quad \longleftarrow \quad \text{weight decay}$$

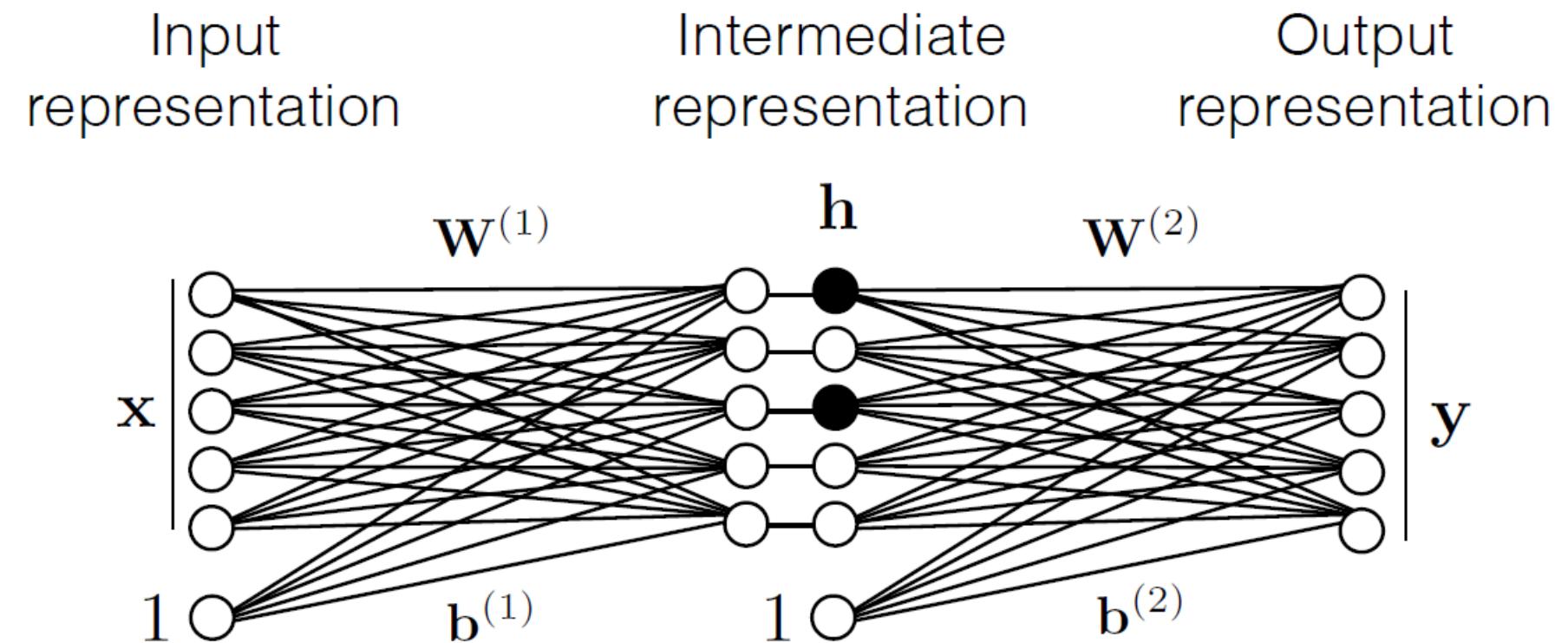
“We prefer to keep weights small.”

Dropout



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout



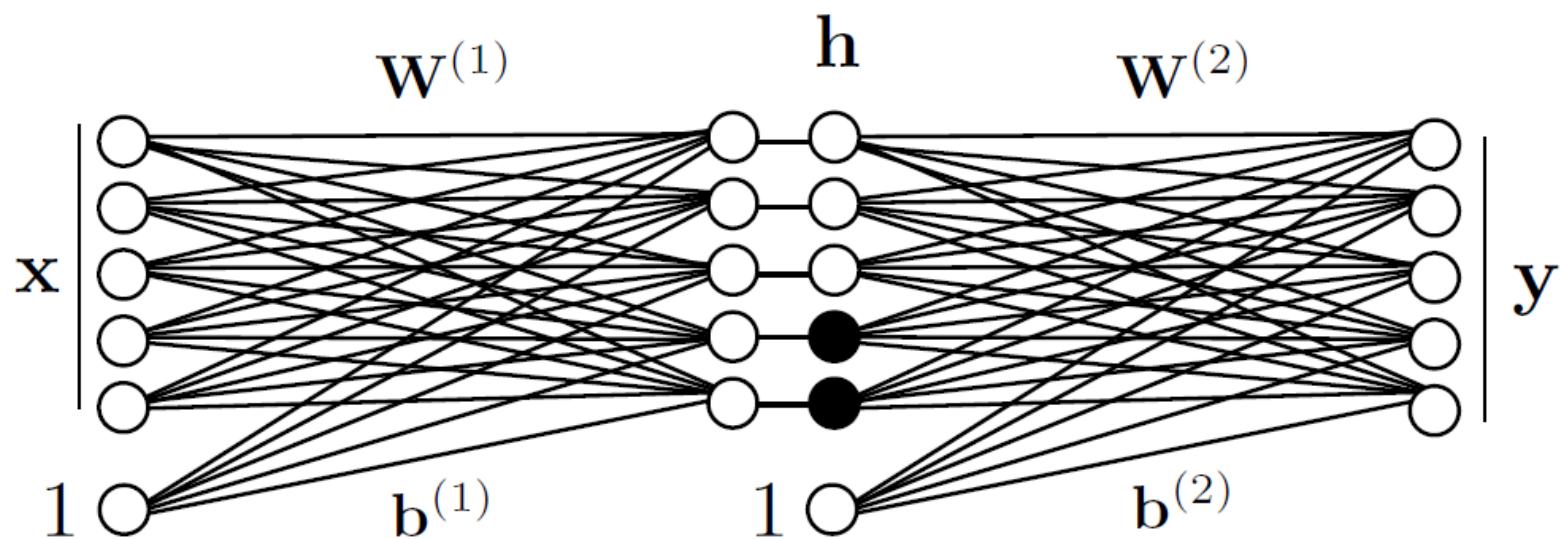
$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

Input
representation

Intermediate
representation

Output
representation



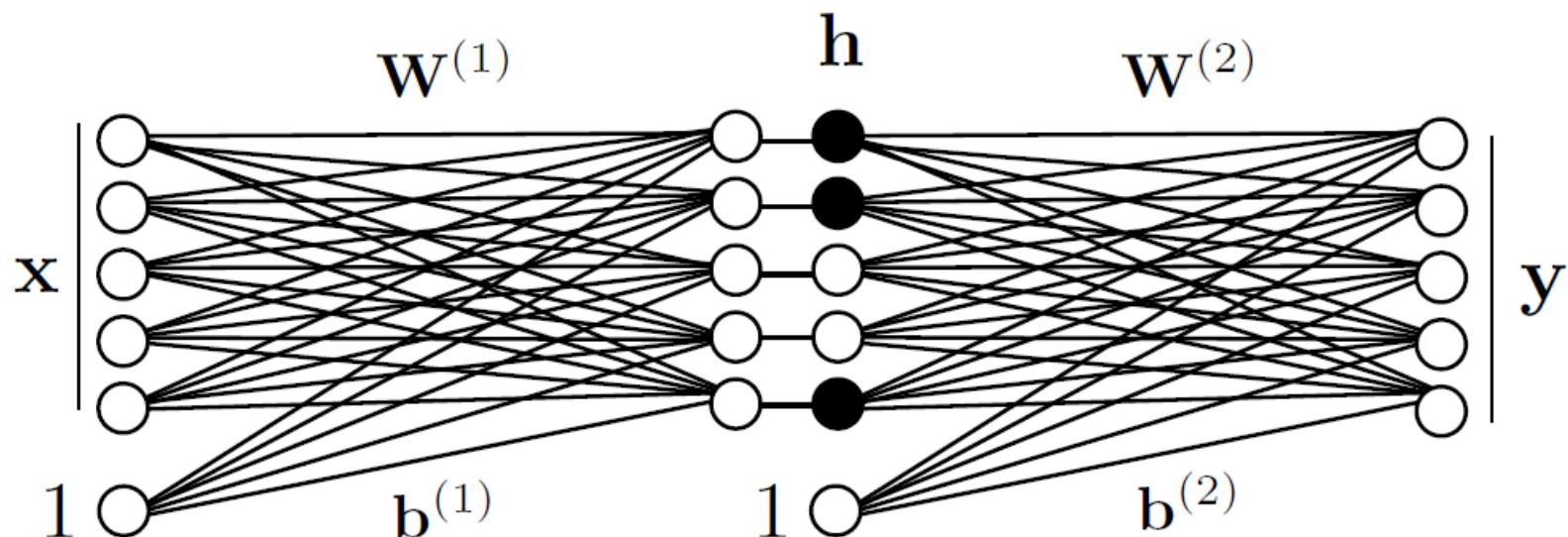
$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

Input
representation

Intermediate
representation

Output
representation



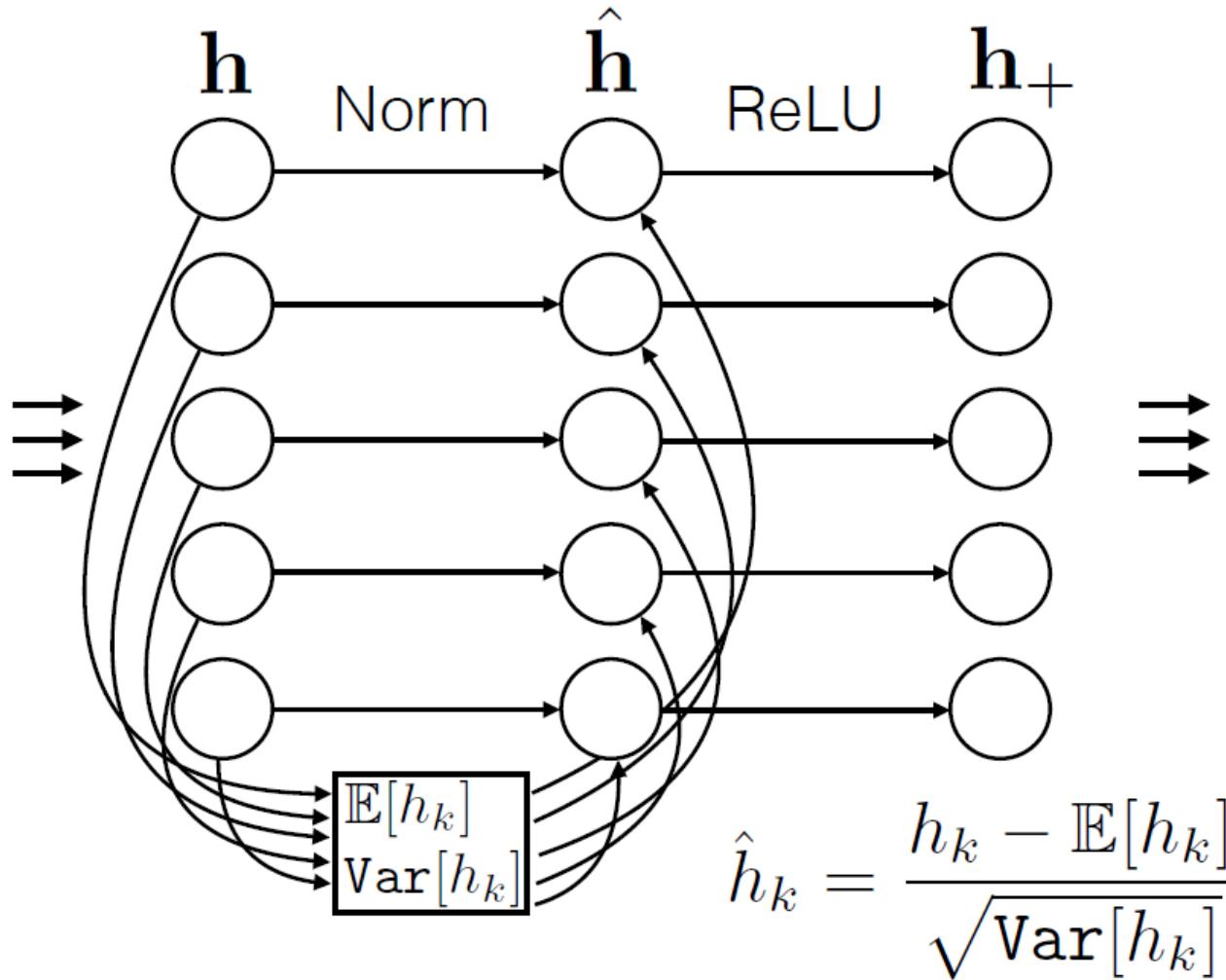
$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

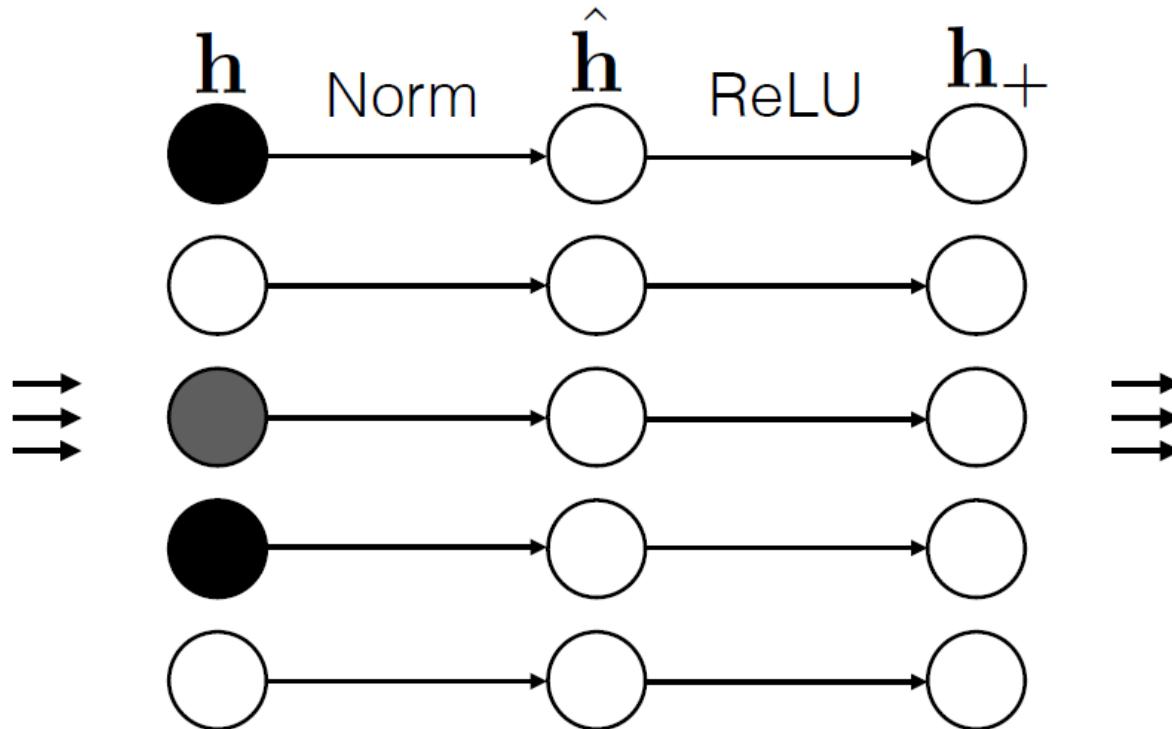
- Randomly zero out hidden units
- Prevents network from relying too much on spurious correlations between different hidden units.
- Can be understood as averaging over an exponential ensemble of subnetworks. This averaging smooths the function, thereby reducing the effective capacity of the network.

Normalization

Normalization layers

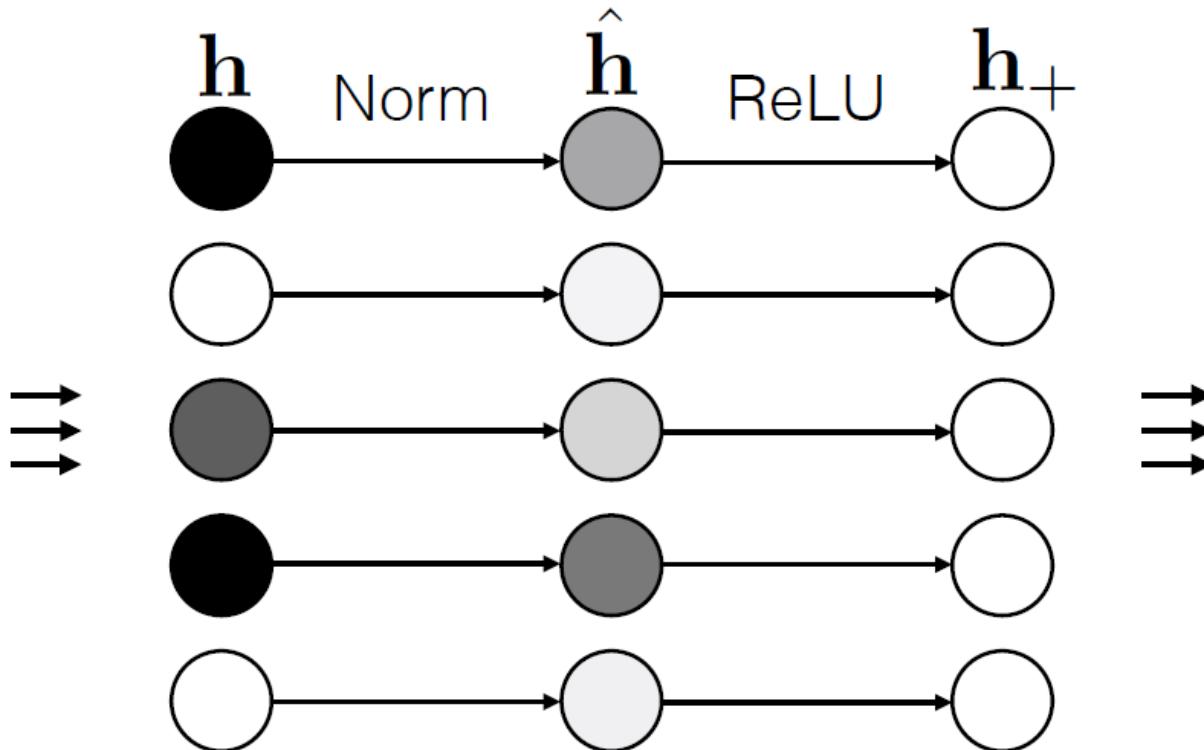


Normalization layers



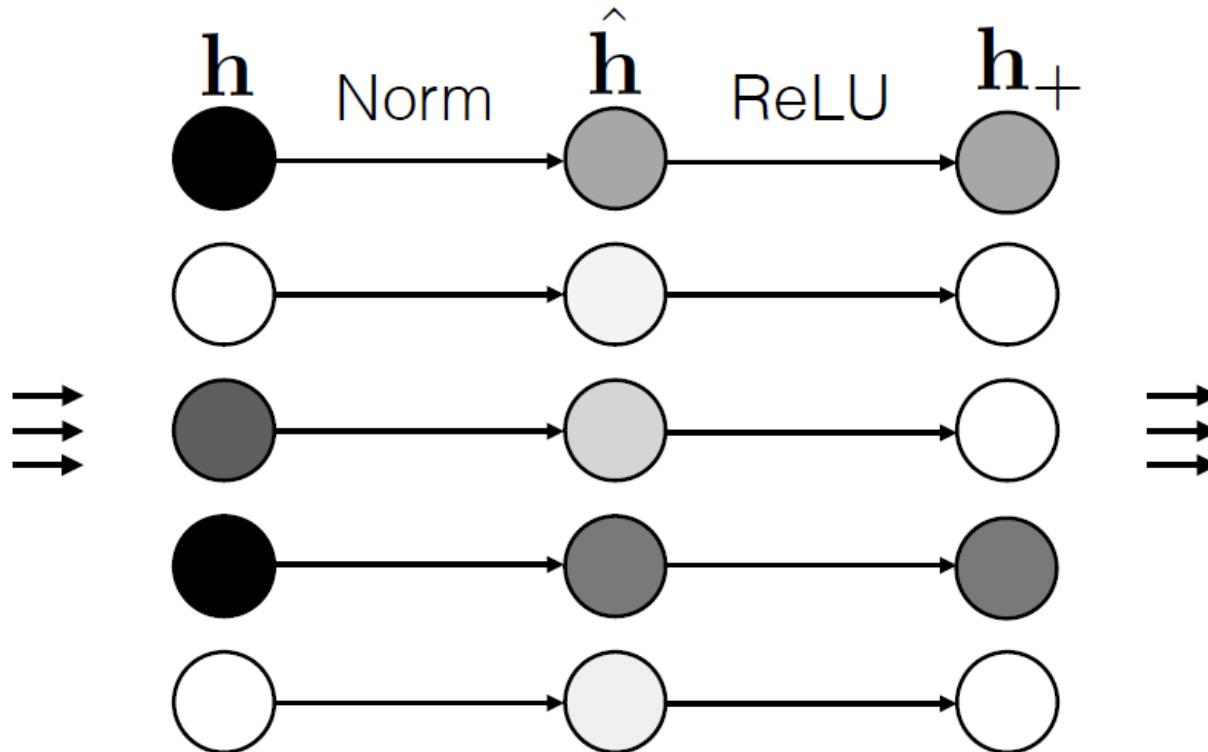
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers

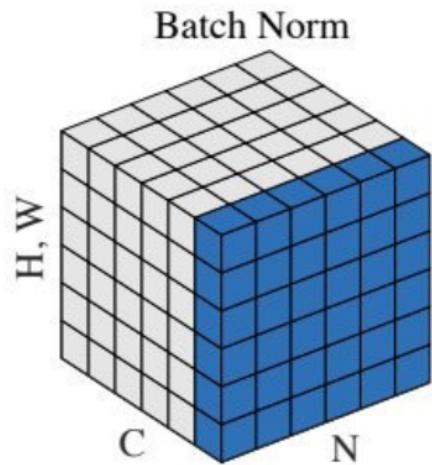


$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers

- Keep track of mean and variance of a unit (or a population of units) over time.
- Standardize unit activations by subtracting mean and dividing by variance.
- Squashes units into a **standard range**, avoiding overflow.
- Also achieves **invariance** to mean and variance of the training signal.
- Both these properties reduce the effective capacity of the model, i.e. regularize the model.

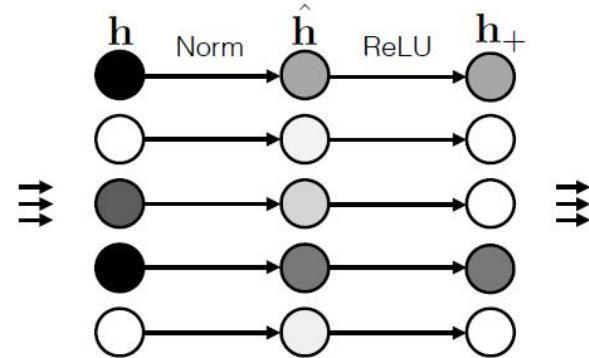
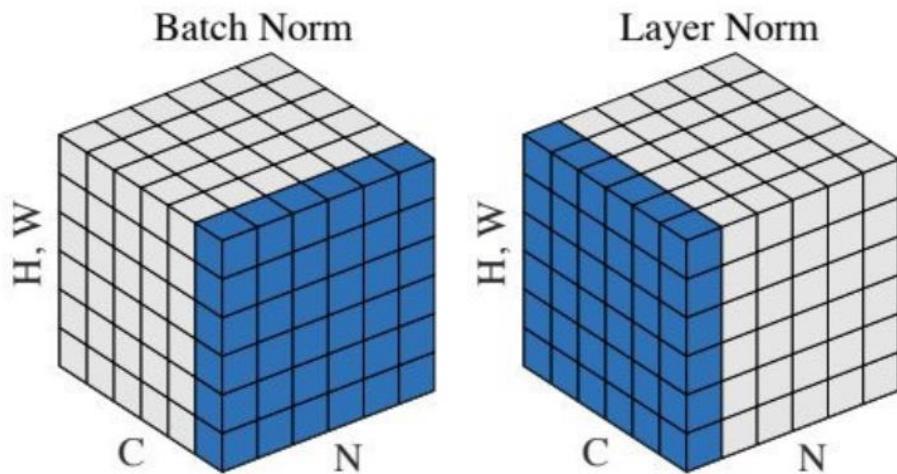
Normalization layers



Normalize w.r.t. a single hidden unit's pattern of activation over training examples (a batch of examples).

[Figure from Wu & He, arXiv 2018]

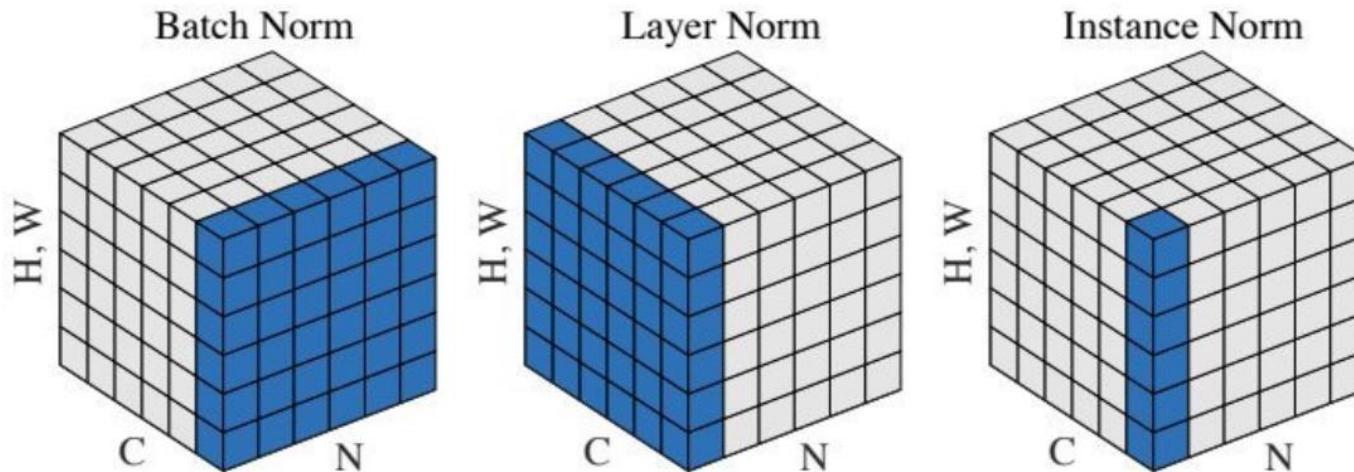
Normalization layers



Normalize w.r.t. the mean and variance of the activations of all the hidden units (neurons) on this layer (c).

[Figure from Wu & He, arXiv 2018]

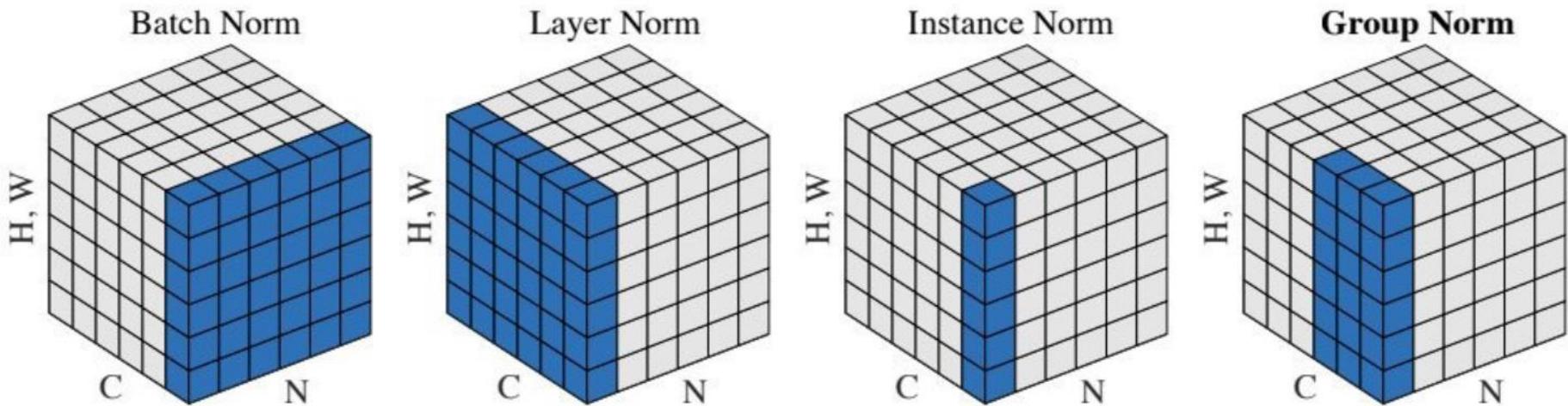
Normalization layers



Normalize w.r.t. the mean and variance of the activations of all the hidden units (neurons) on this layer (c) that process this **particular example**.

[Figure from Wu & He, arXiv 2018]

Normalization layers

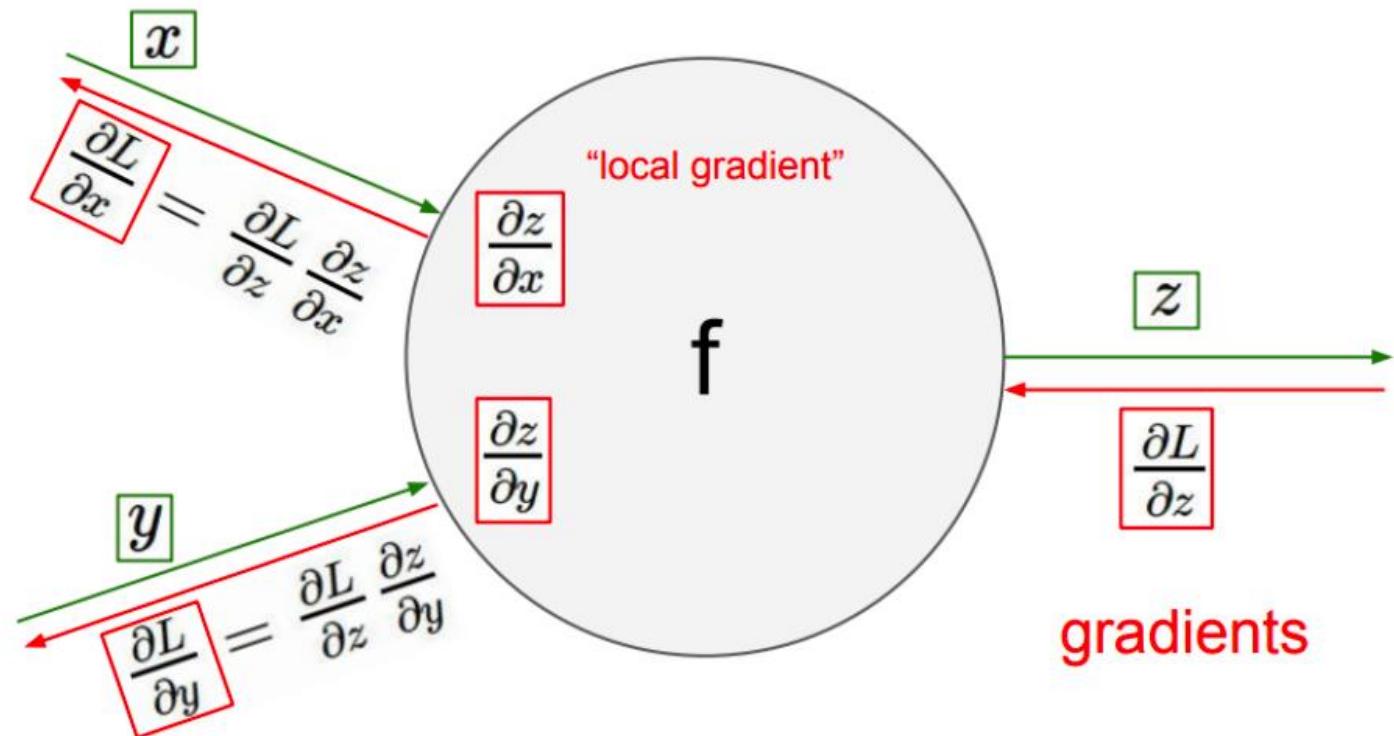


Might as well...

[Figure from Wu & He, arXiv 2018]

Next up:

Optimization of NNs (backprop)



Taxonomy of Learning Problems

Supervised learning: learning from examples

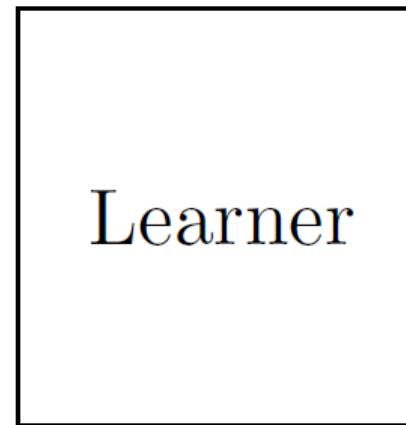
Training data

$$\{x_1, y_1\}$$

$$\{x_2, y_2\} \rightarrow$$

$$\{x_3, y_3\}$$

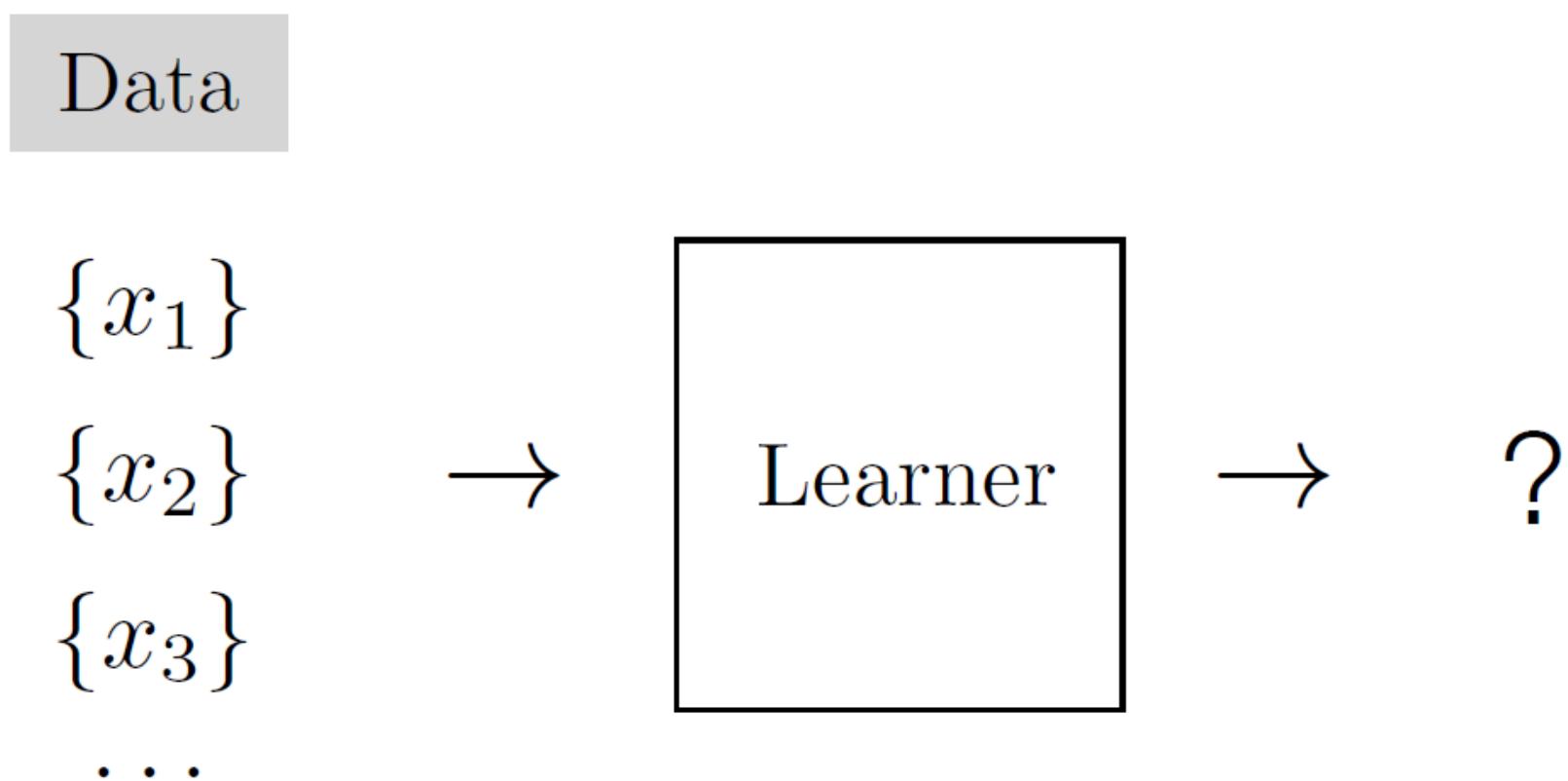
...



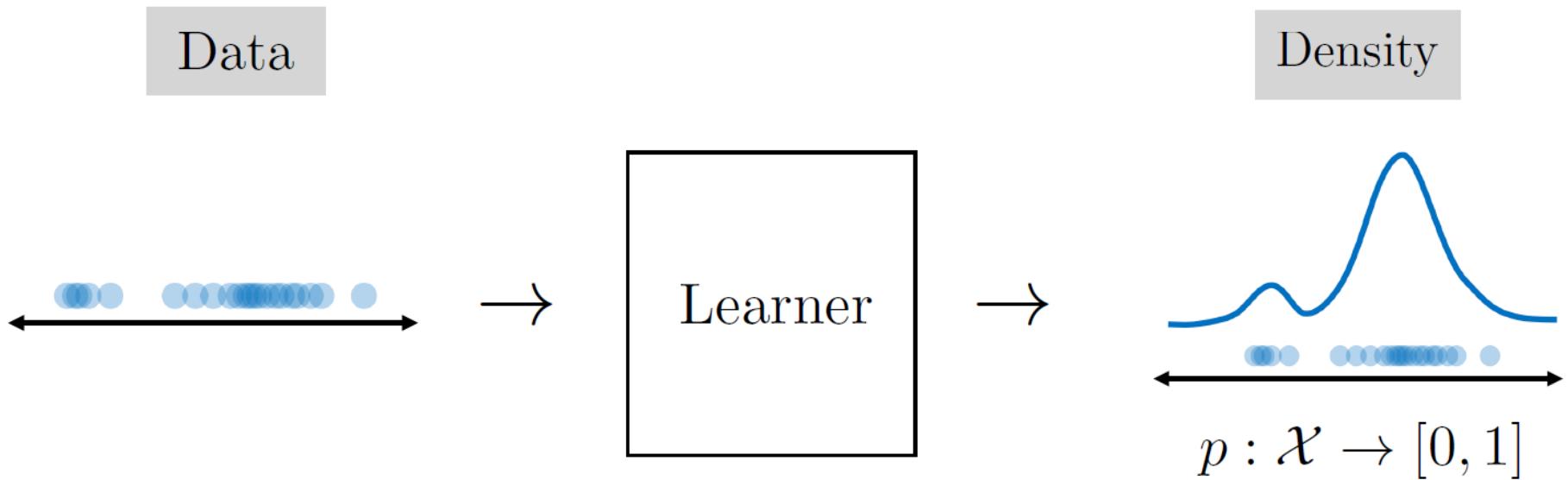
$$\rightarrow f : X \rightarrow Y$$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(x_i), y_i)$$

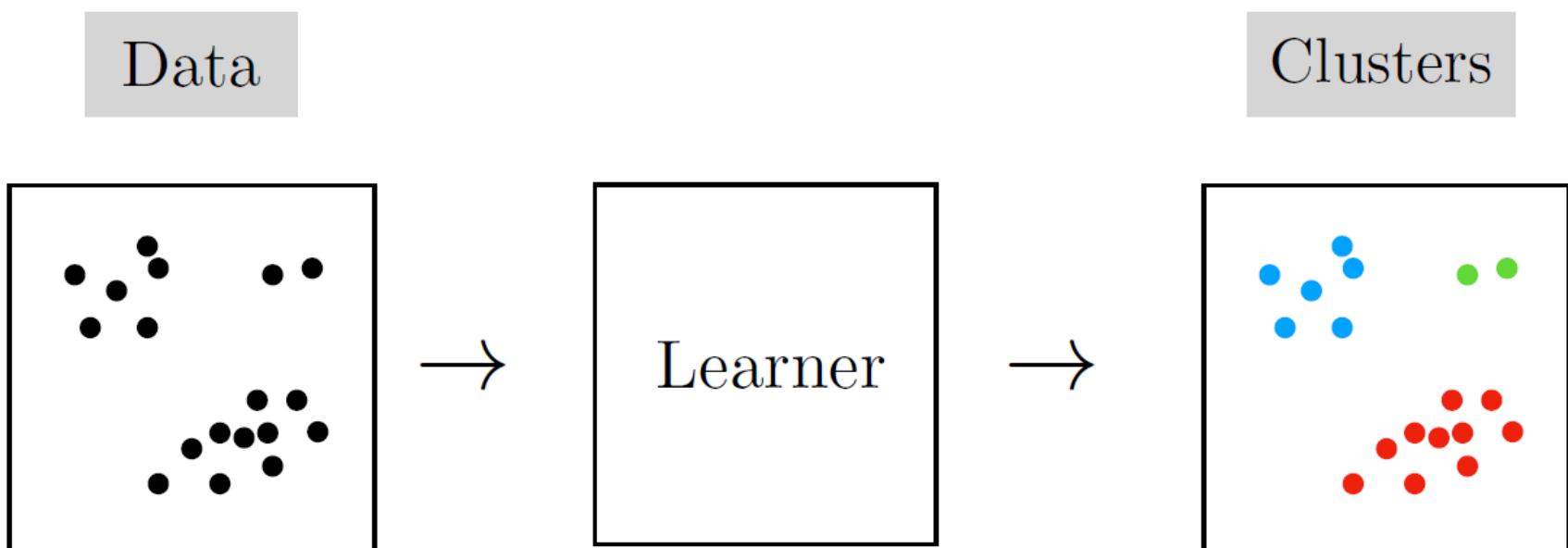
Unsupervised learning: Learning without examples



Density modeling



Clustering



$$f : \mathcal{X} \rightarrow \{1, \dots, k\}$$

Representation Learning

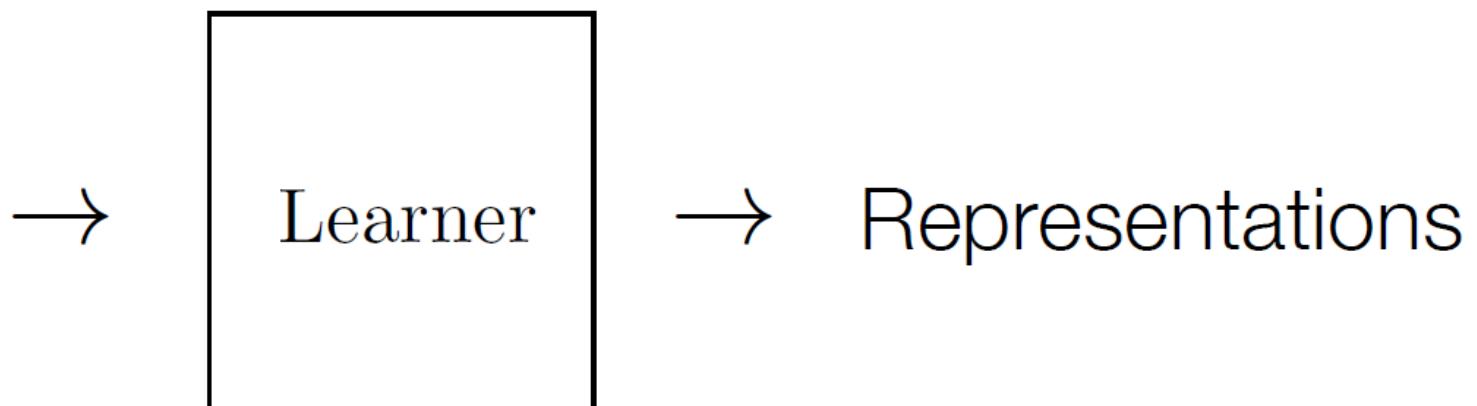
Data

$\{x_1\}$

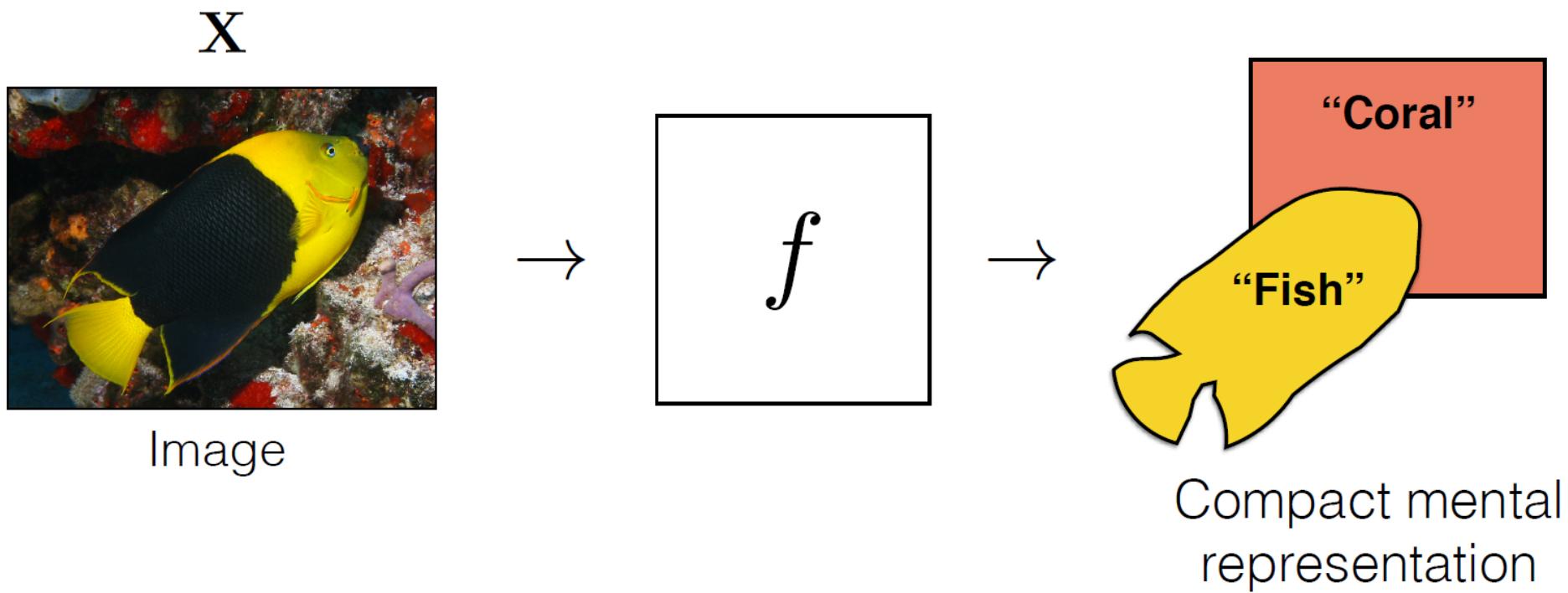
$\{x_2\}$

$\{x_3\}$

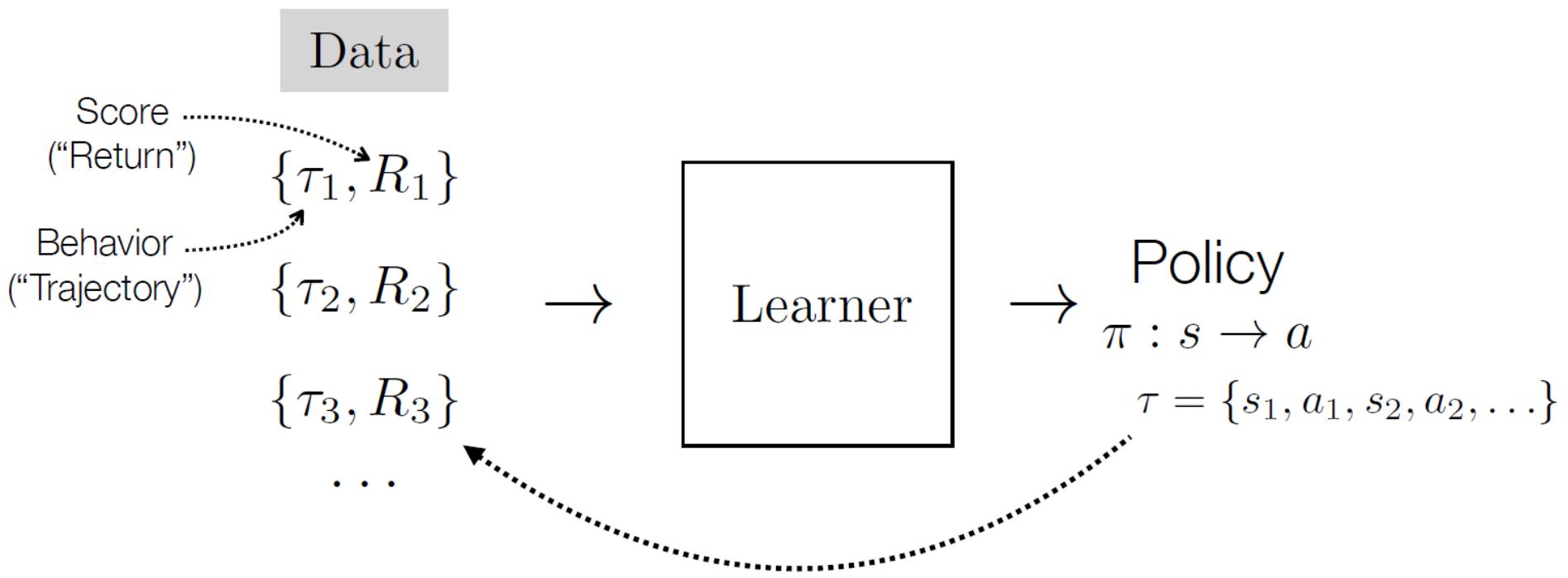
...



Representation Learning



Reinforcement learning



Differentiable programming

Deep nets are popular for a few reasons:

1. High capacity
2. Easy to optimize (differentiable)
3. Compositional “block based programming”

An emerging term for general models with these properties is **differentiable programming**.



Yann LeCun

January 5 ·



OK, Deep Learning has outlived its usefulness as a buzz-phrase.
Deep Learning est mort. Vive Differentiable Programming!



Thomas G. Dietterich

@tdietterich

Following

DL is essentially a new style of programming--"differentiable programming"--and the field is trying to work out the reusable constructs in this style. We have some: convolution, pooling, LSTM, GAN, VAE, memory units, routing units, etc. 8/

8:02 AM - 4 Jan 2018

65 Retweets 194 Likes



6 65 194