

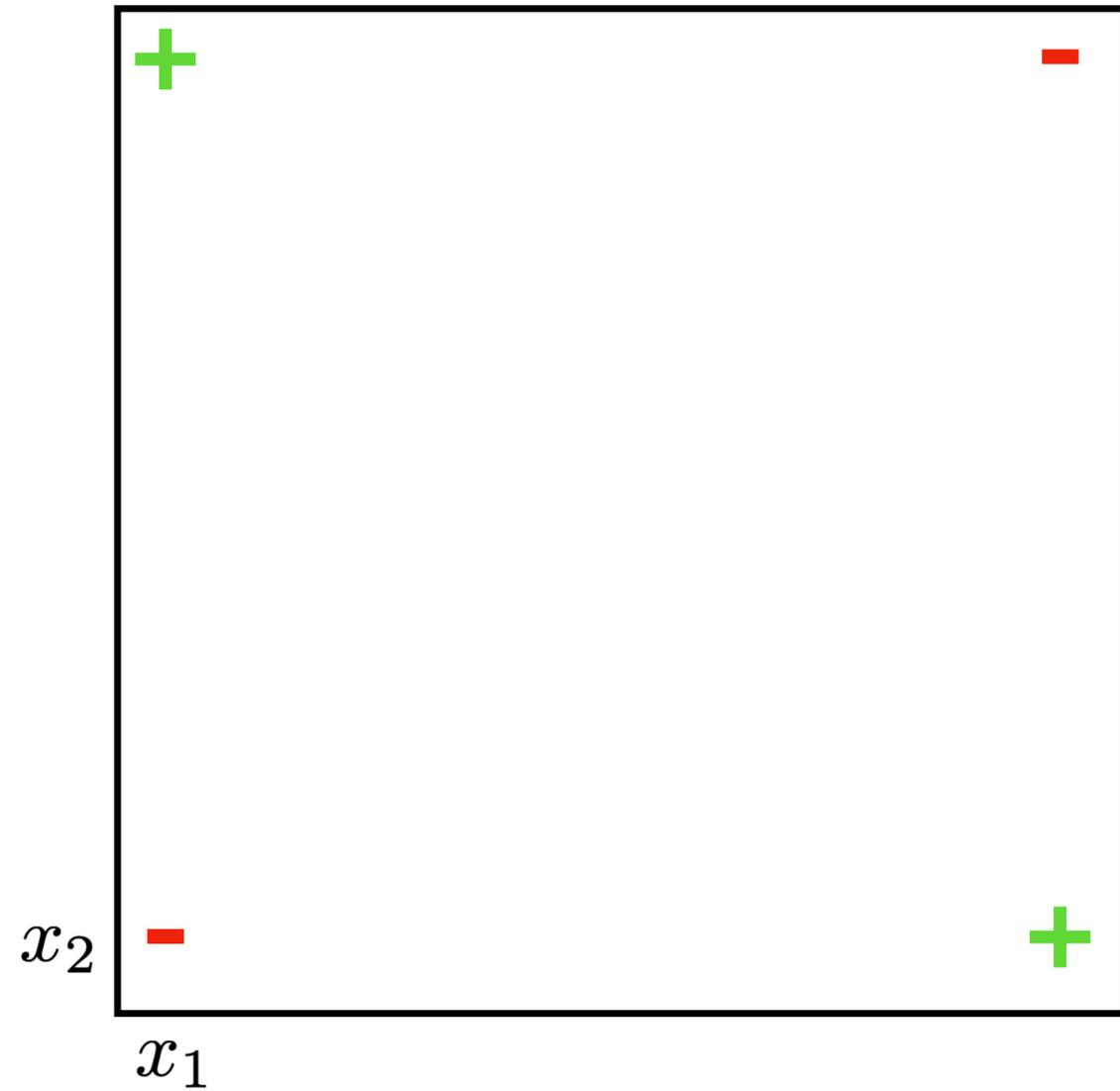
# Lecture 7: Neural networks

- Discussion this week: machine learning
- Reading:
  - Szeliski 5.3
  - Goodfellow Deep Feedforward Networks
- Start thinking about project
- PS2 due today – submit to gradescope and canvas

# Today

- Brief history of neural networks
- Computation in neural networks
- Multi-layer perceptrons (for PS4)
- Estimating gradients (to be continued next class).

# Limitations to linear classifiers



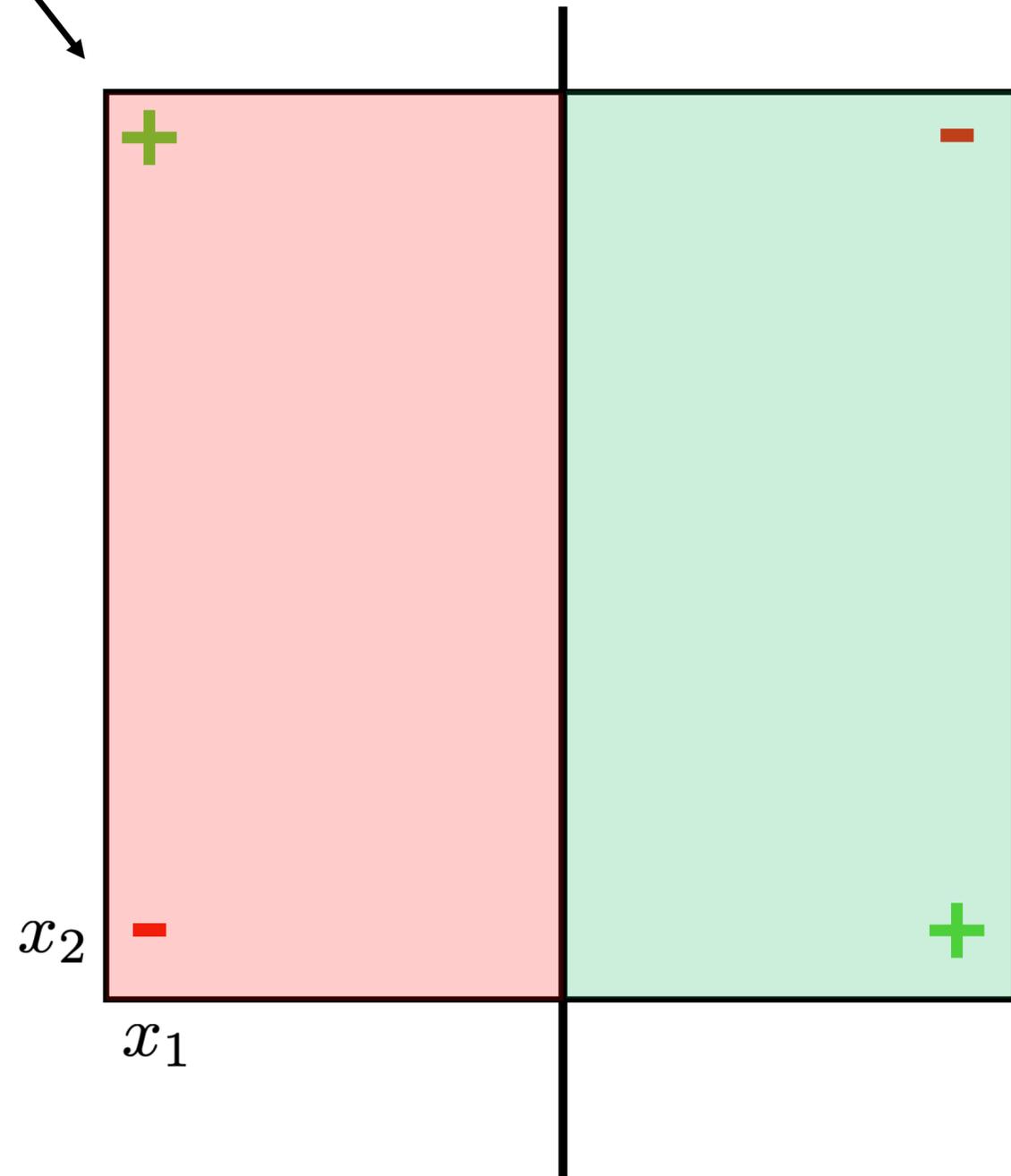
		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

# Limitations to linear classifiers

Wrong!

Wrong!

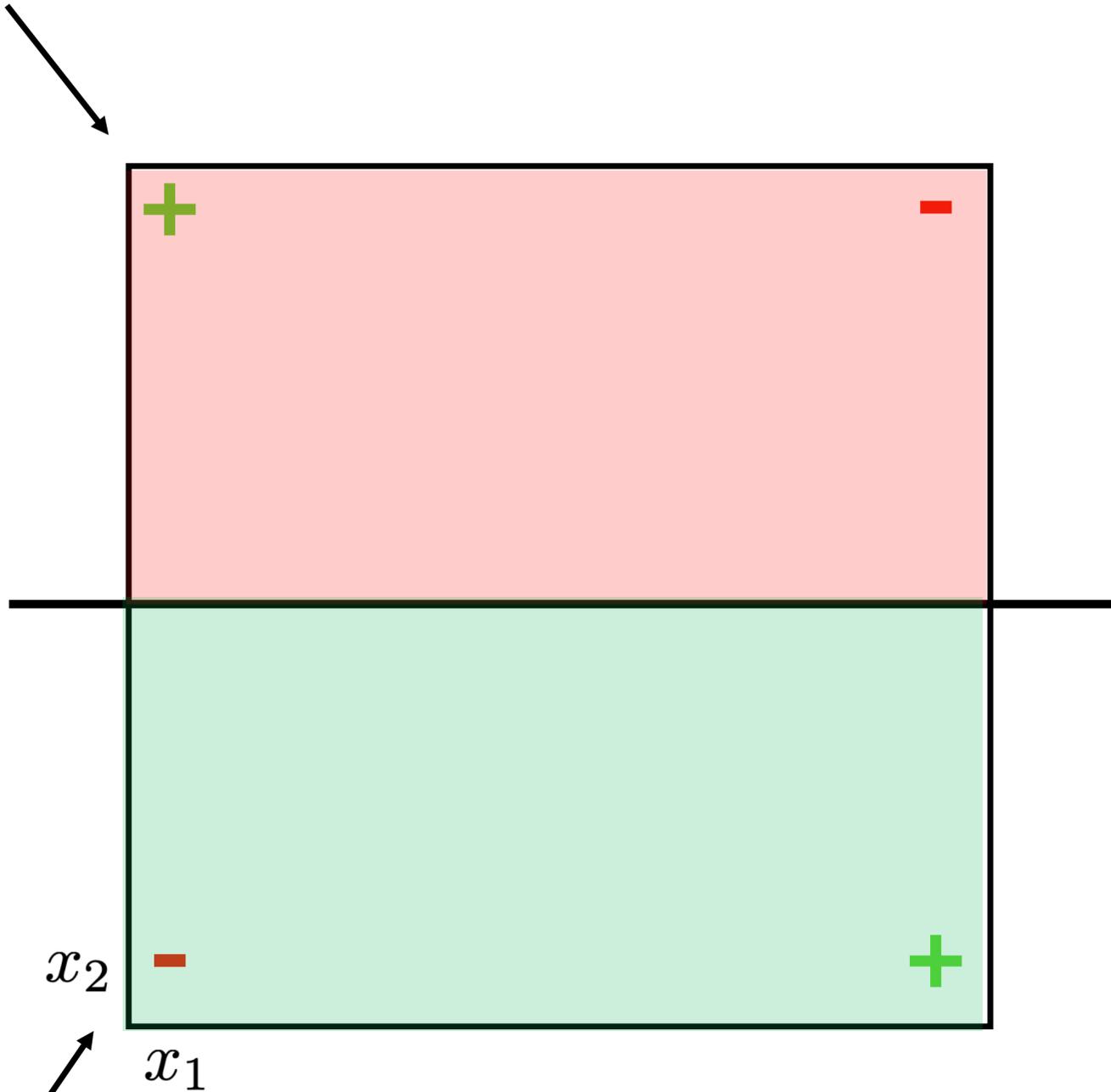


		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

# Limitations to linear classifiers

Wrong!

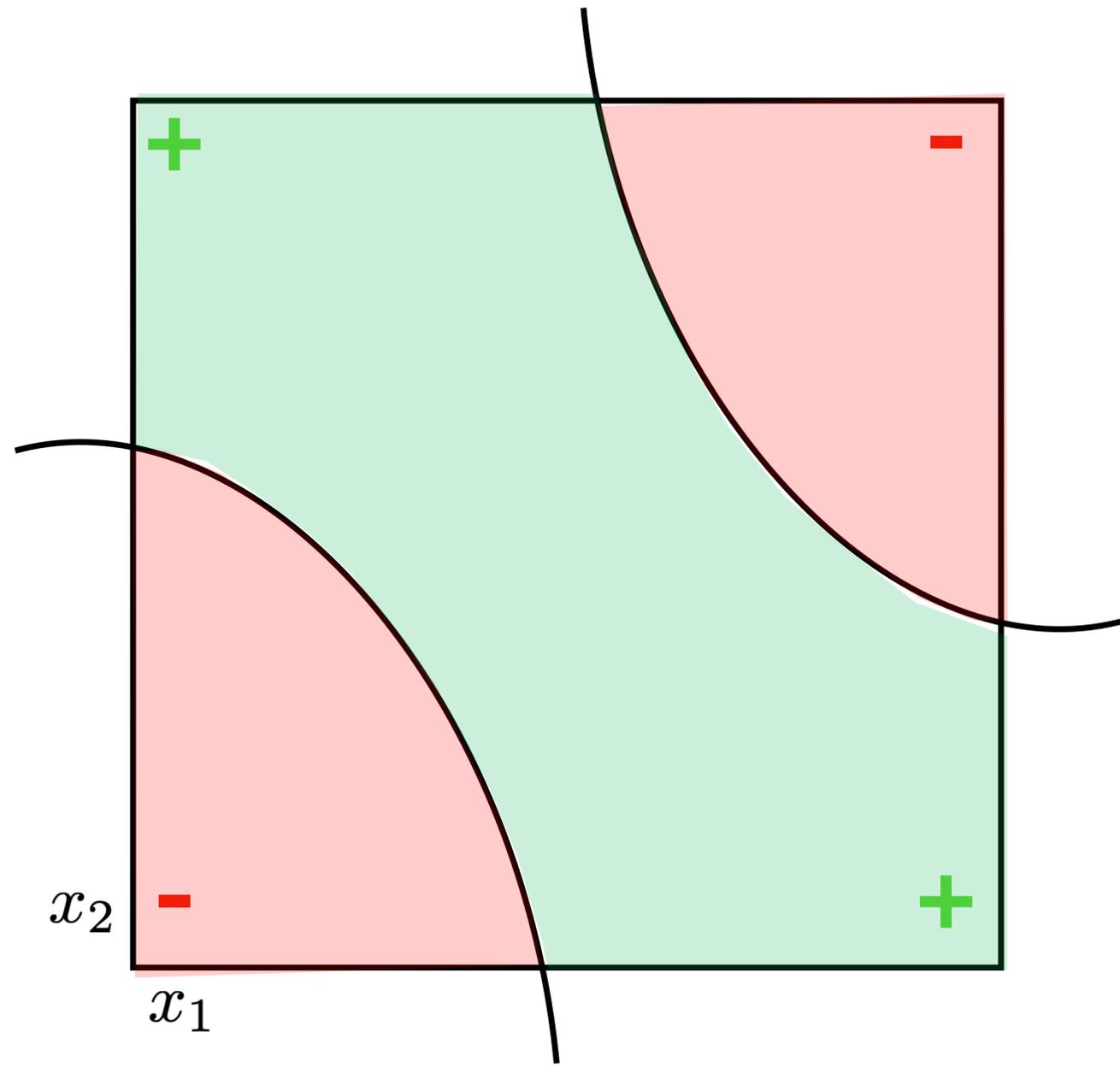


		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

Wrong!

# Goal: Non-linear decision boundary

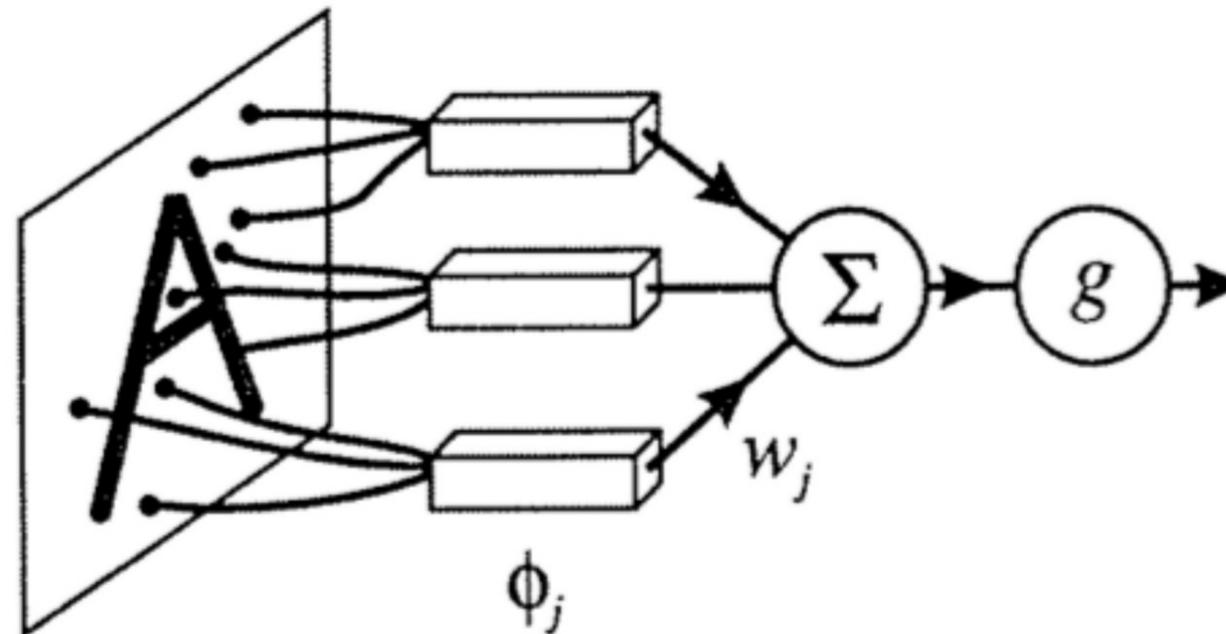


		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

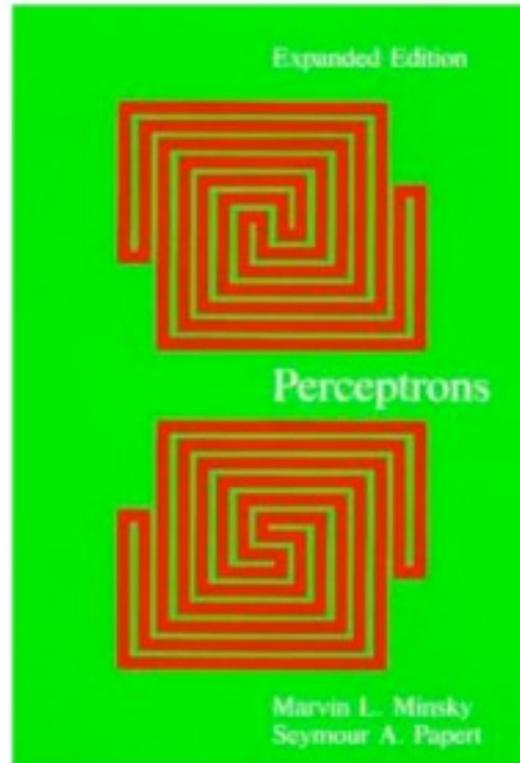
XOR

# Perceptron

- In 1957 Frank Rosenblatt invented the perceptron
- Computers at the time were too slow to run the perceptron, so Rosenblatt built a special purpose machine with adjustable resistors
- New York Times Reported: “The Navy revealed the embryo of an electronic computer that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence”



# Minsky and Papert, Perceptrons, 1972



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

## Perceptrons, expanded edition

An Introduction to Computational Geometry

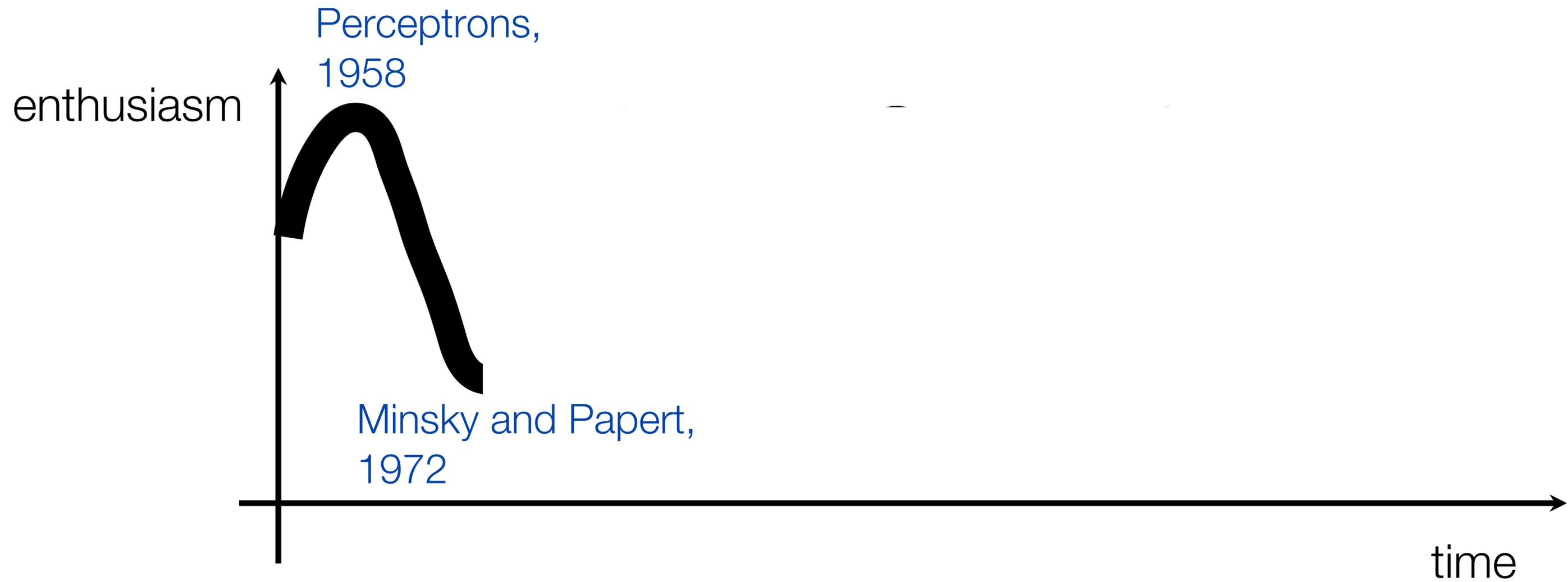
By [Marvin Minsky](#) and [Seymour A. Papert](#)

### Overview

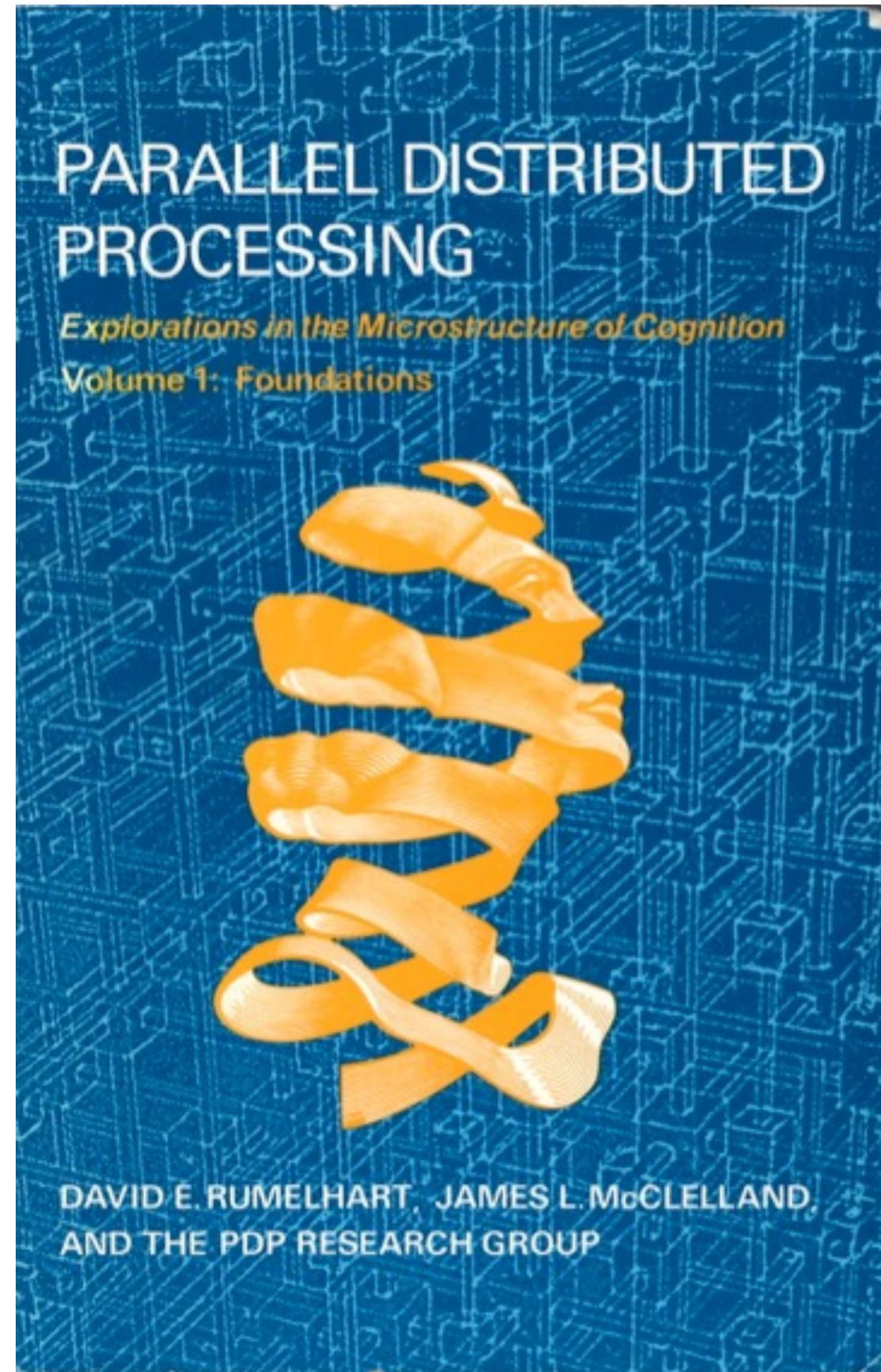
*Perceptrons* - the first systematic study of parallelism in computation - has remained a classical work on threshold automata networks for nearly two decades. It marked a historical turn in artificial intelligence, and it is required reading for anyone who wants to understand the connectionist counterrevolution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

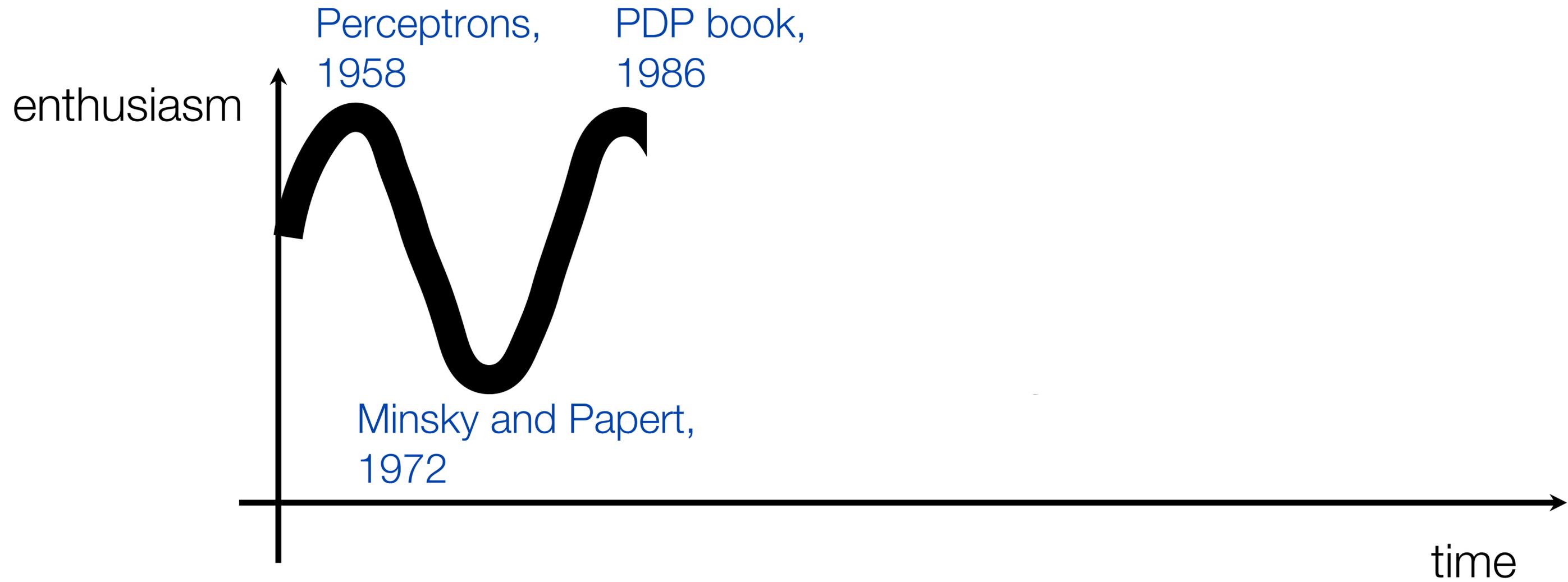
Artificial-intelligence research, which for a time concentrated on the programming of ton Neumann computers, is swinging back to the idea that intelligence might emerge from the activity of networks of neuronlike entities. Minsky and Papert's book was the first example of a mathematical analysis carried 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."



# Parallel Distributed Processing (PDP), 1986





# LeCun convolutional neural networks

PROC. OF THE IEEE, NOVEMBER 1998

7

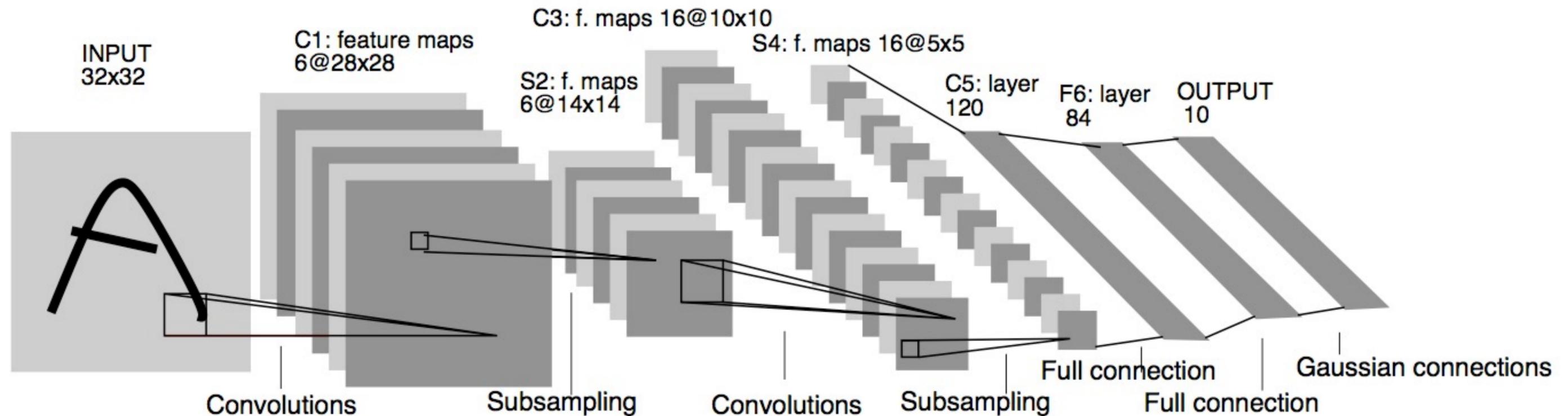


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Demos:

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

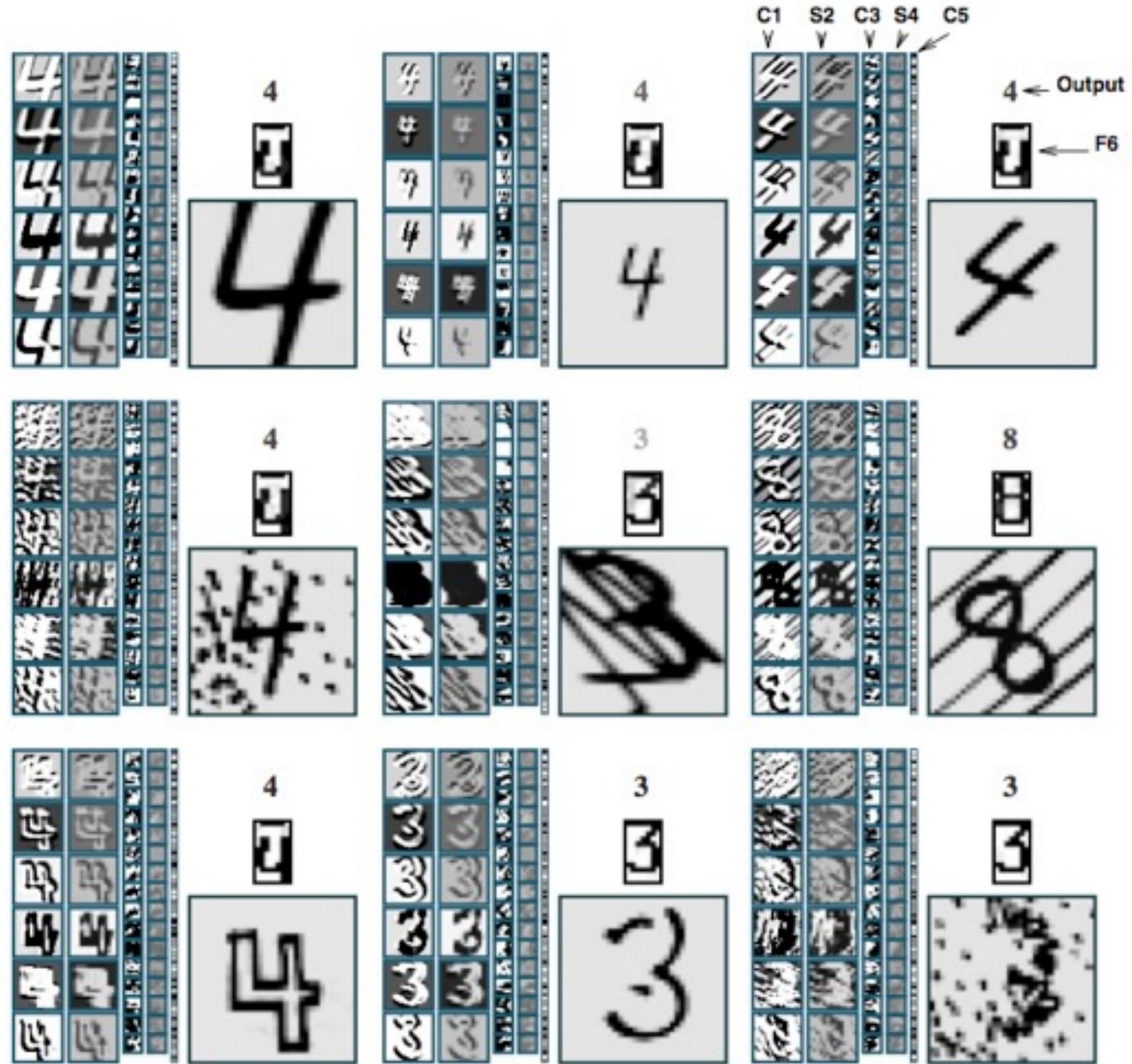
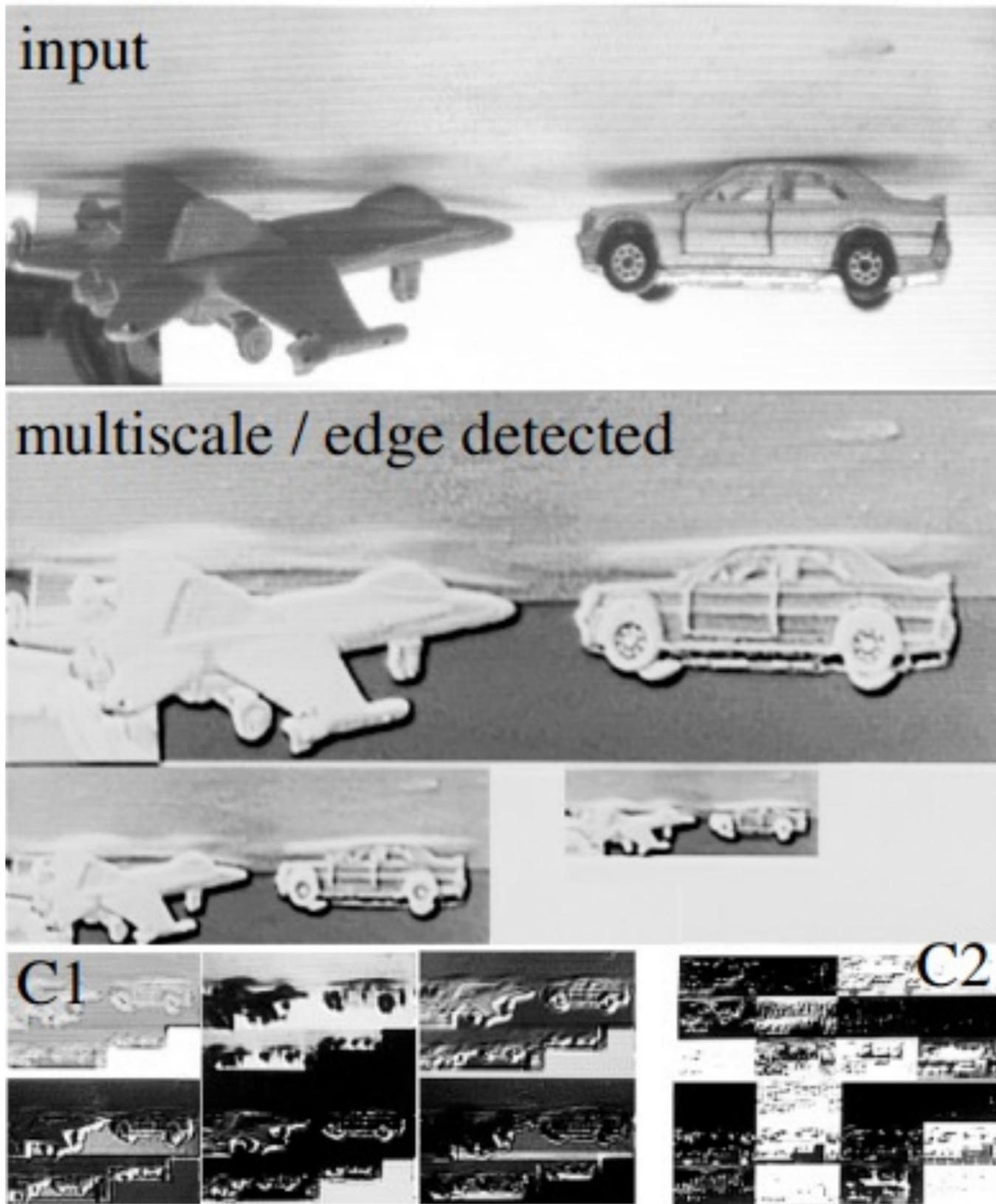


Fig. 13. Examples of unusual, distorted, and noisy characters correctly recognized by LeNet-5. The grey-level of the output label represents the penalty (lighter for higher penalties).



Neural networks to recognize handwritten digits and human faces?

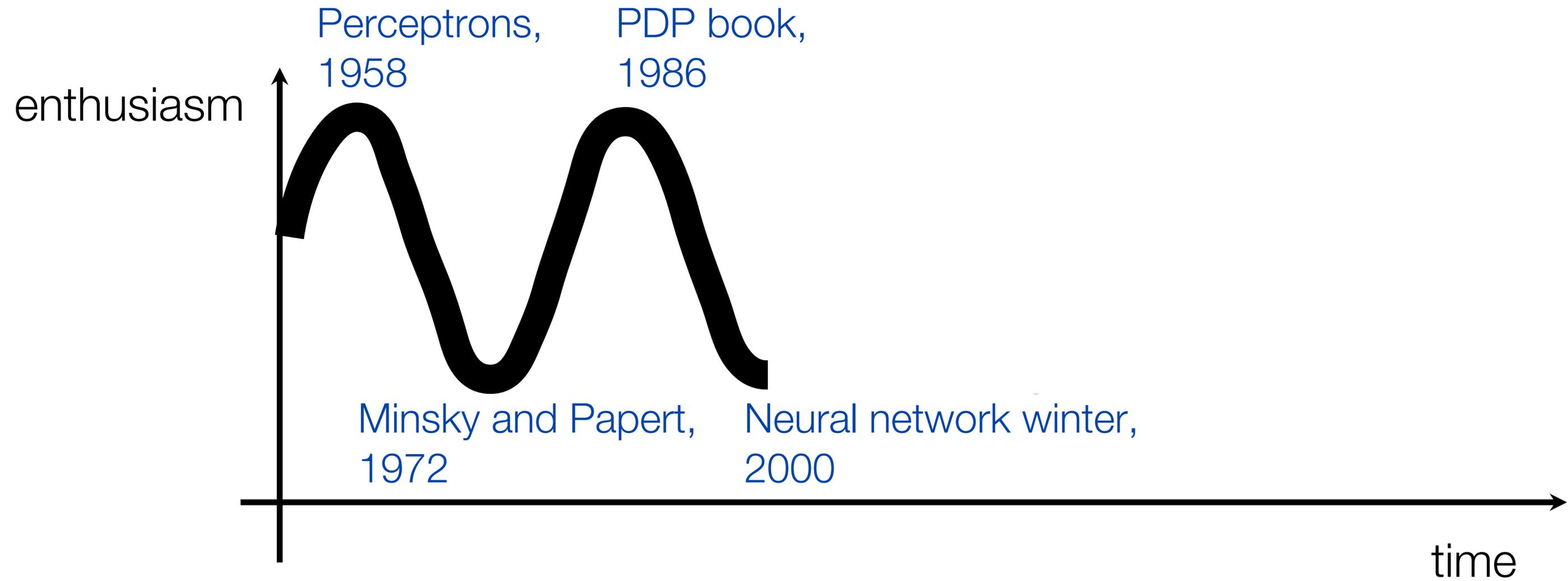
yes

Neural networks for tougher problems?

not really

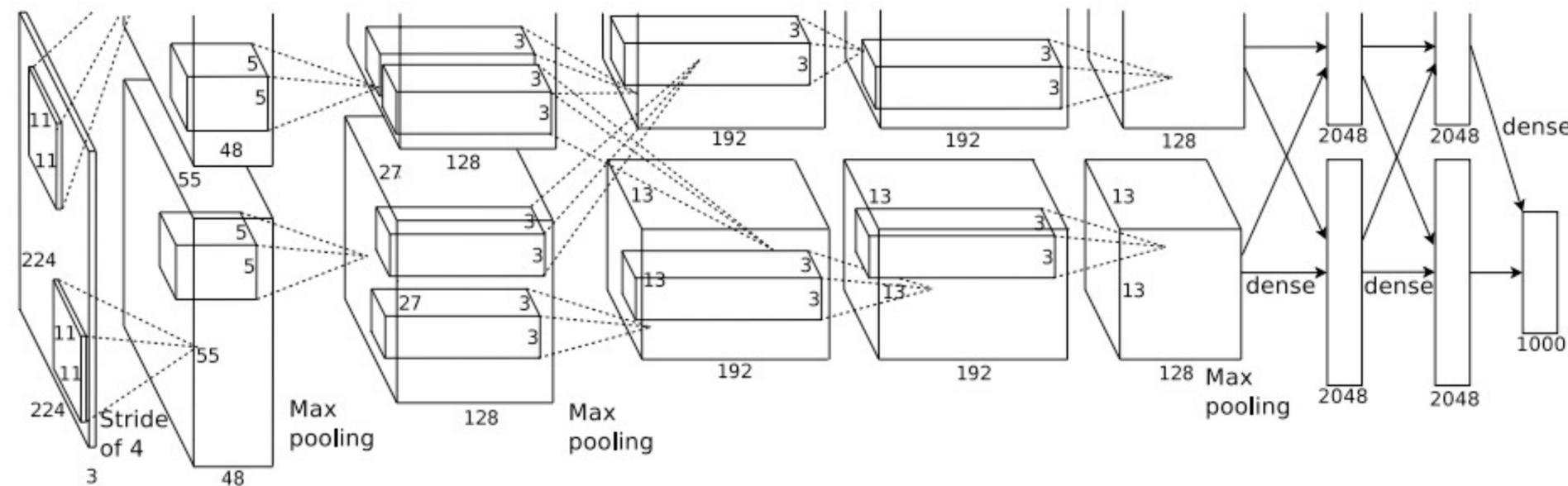
# Machine learning circa 2000

- Neural Information Processing Systems (NeurIPS), is a top conference on machine learning.
- 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”



Got all the “pieces” right, e.g.,

- Trained on ImageNet
- 8 layer architecture (for reference: today we have architectures with 100+ layers)
- Allowed for multi-GP training



# Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



**mite**

**container ship**

**motor scooter**

**leopard**

--	--	--	--



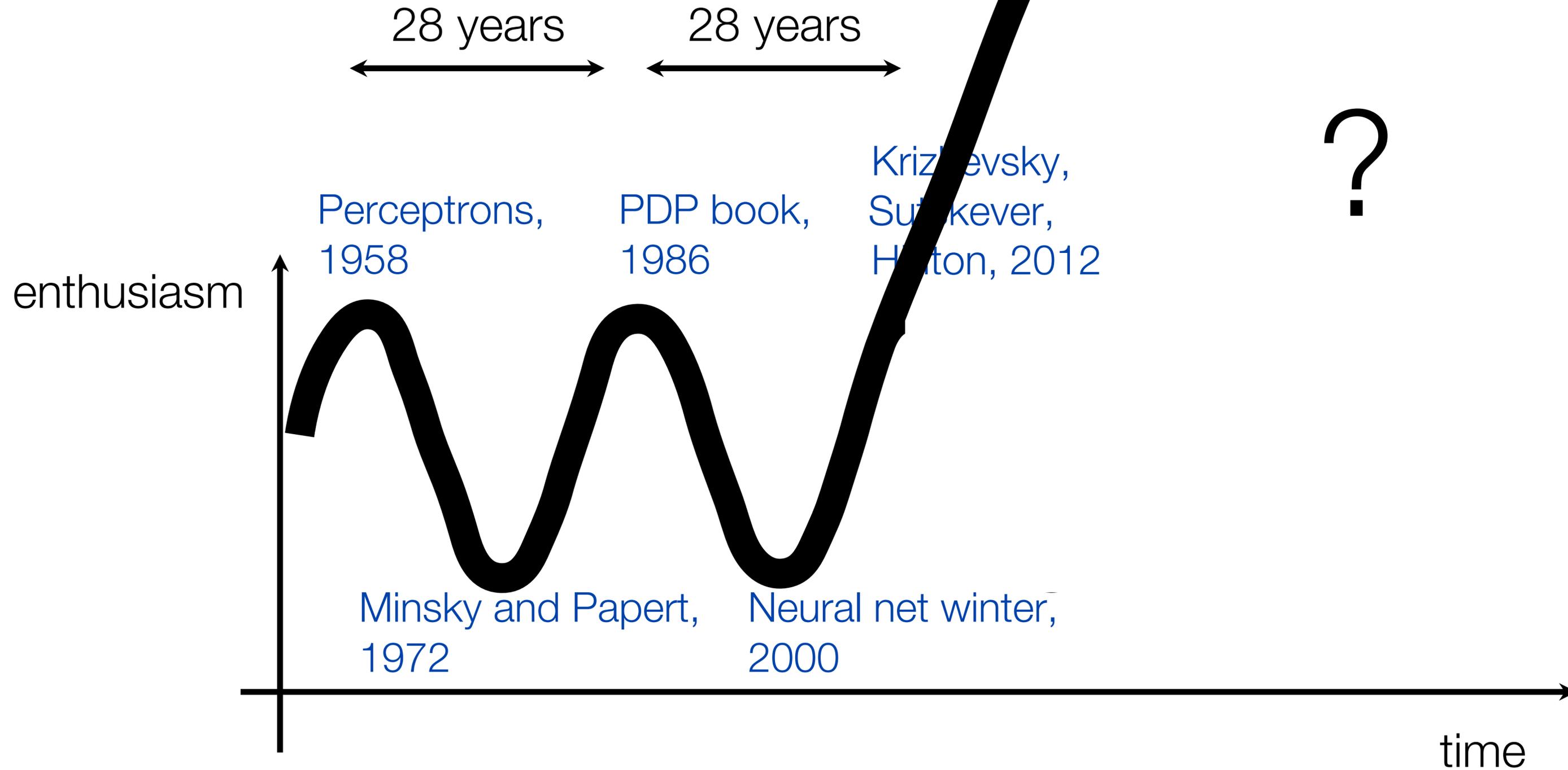
**grille**

**mushroom**

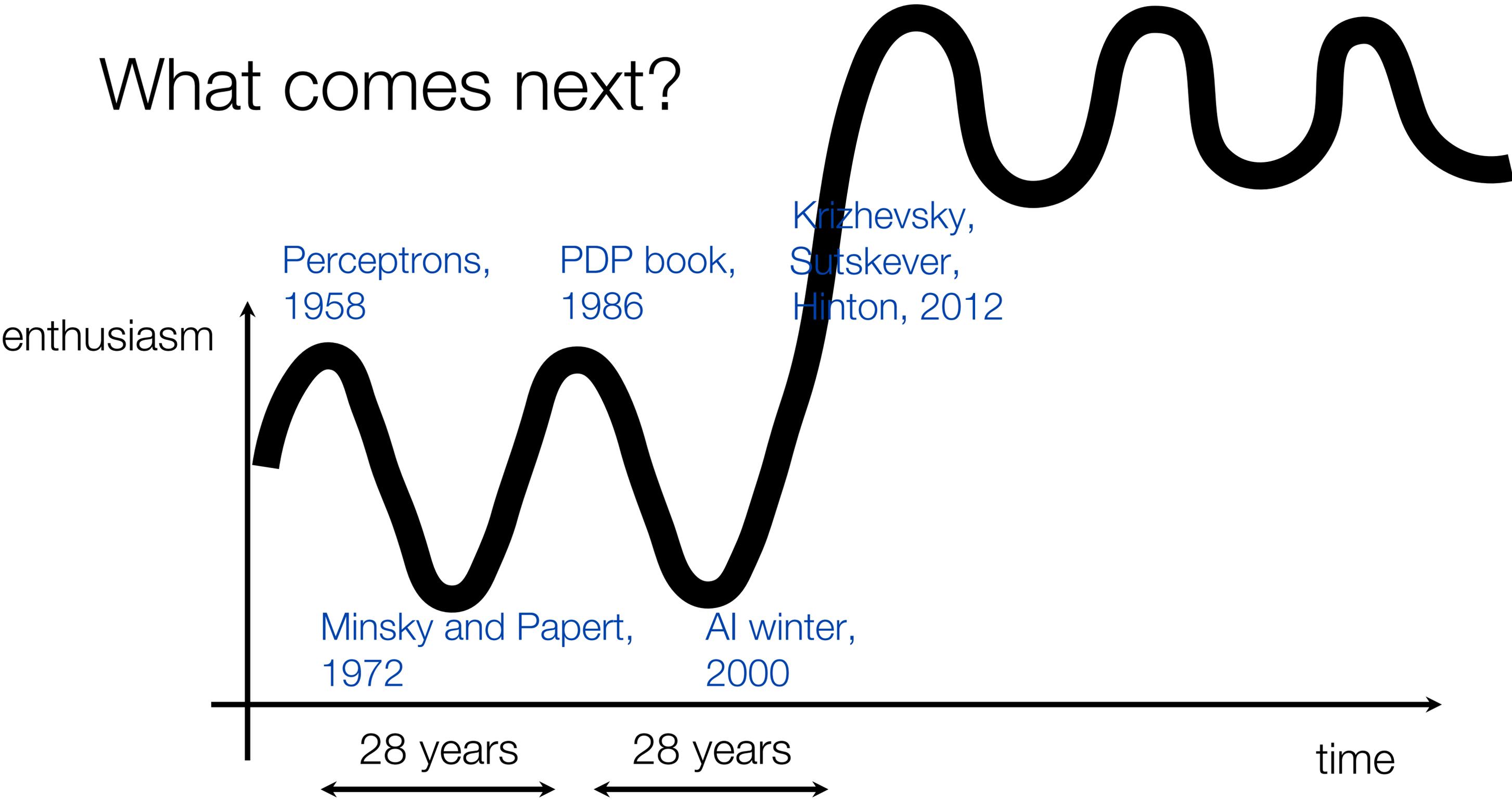
**cherry**

**Madagascar cat**

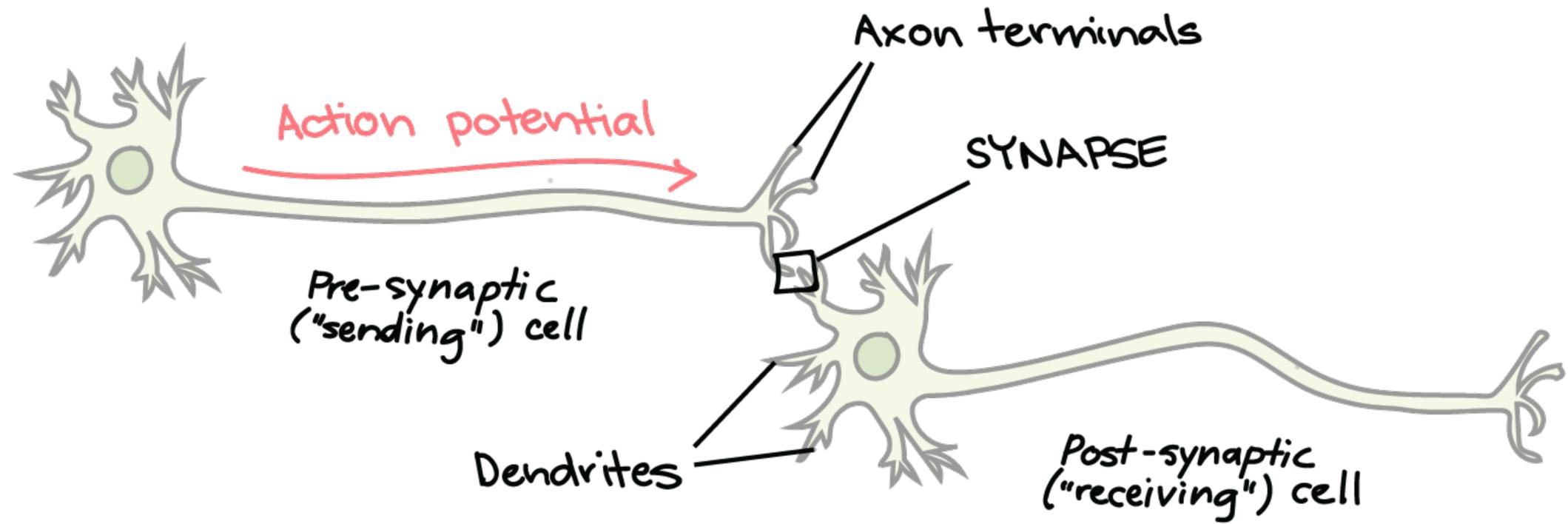
--	--	--	--



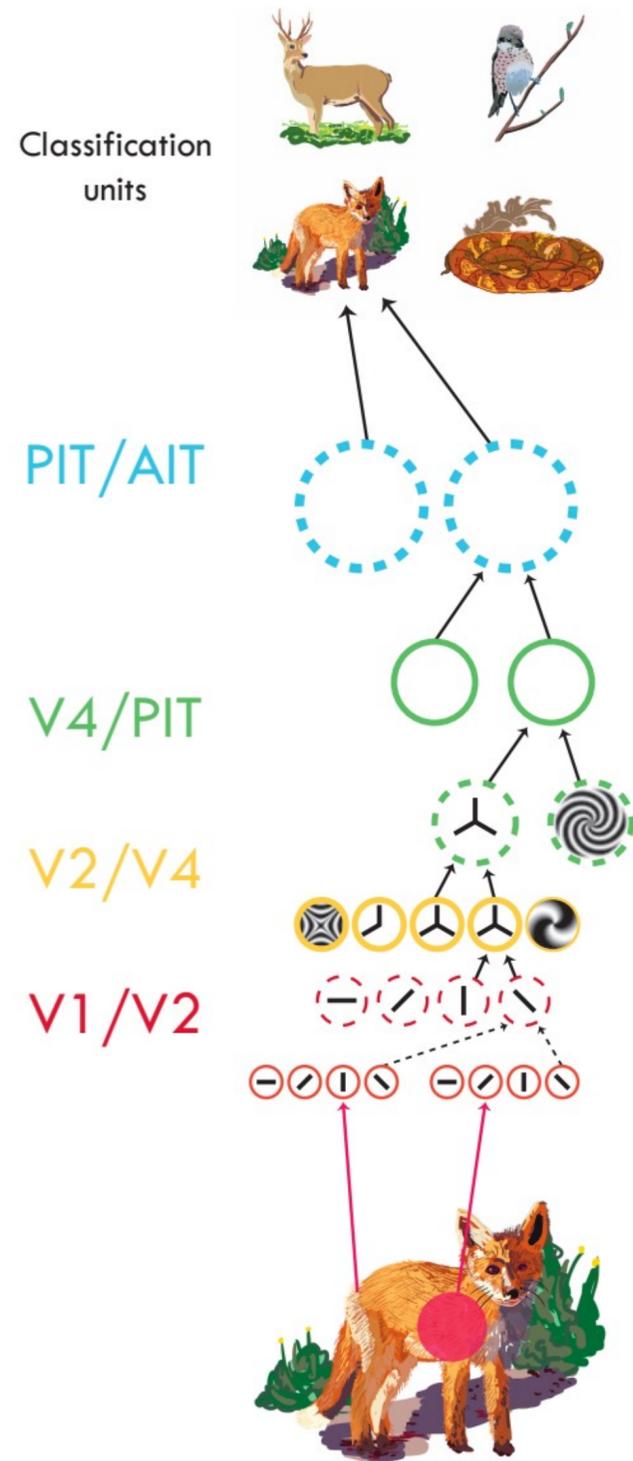
# What comes next?



# Inspiration: Neurons



# Inspiration: Hierarchical Representations



Best to treat as *inspiration*. The neural nets we'll talk about aren't very biologically plausible.

[Serre, 2014]

# Object recognition

Pixel 1



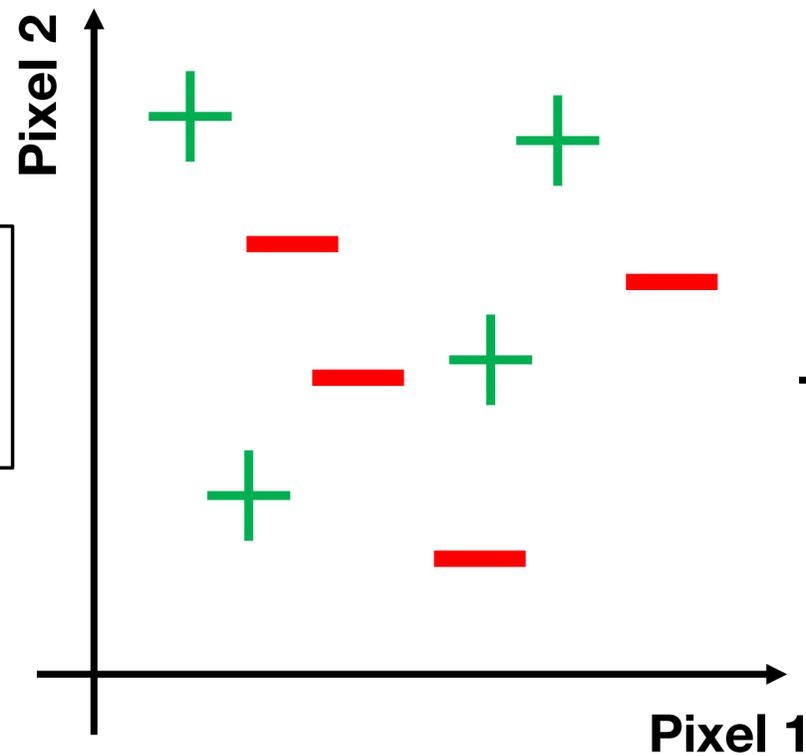
Neural Network

Is dog?

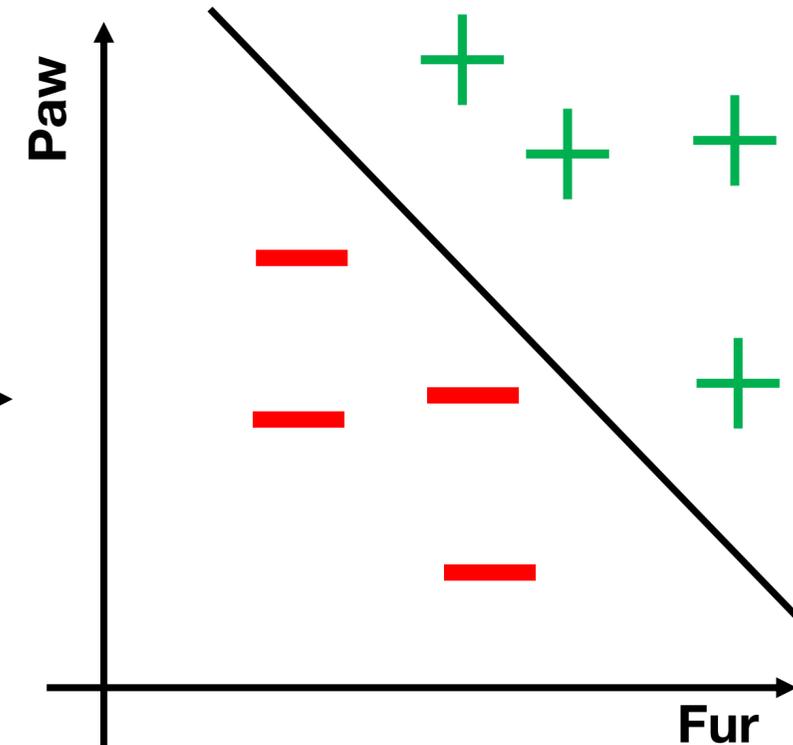
Pixel 2

Input Space

Feature Space



$f(x)$

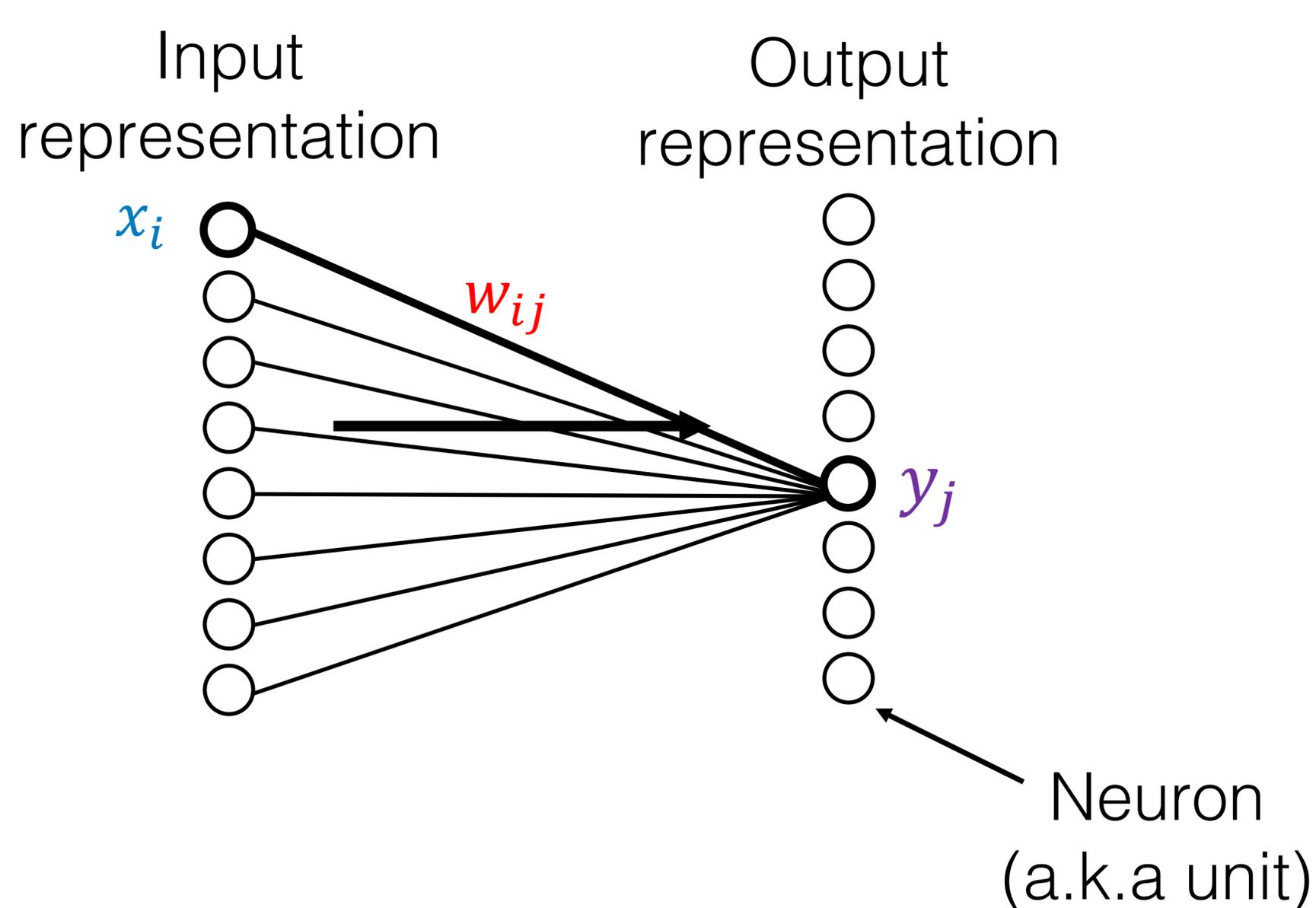


Goal: automatically learn a function that maps data from the input space to a feature space, i.e., "feature learning", rather than use hand-crafted features

# Computation in a neural net

Lets say we have some 1D input that we want to convert to some new feature space:

## Linear layer



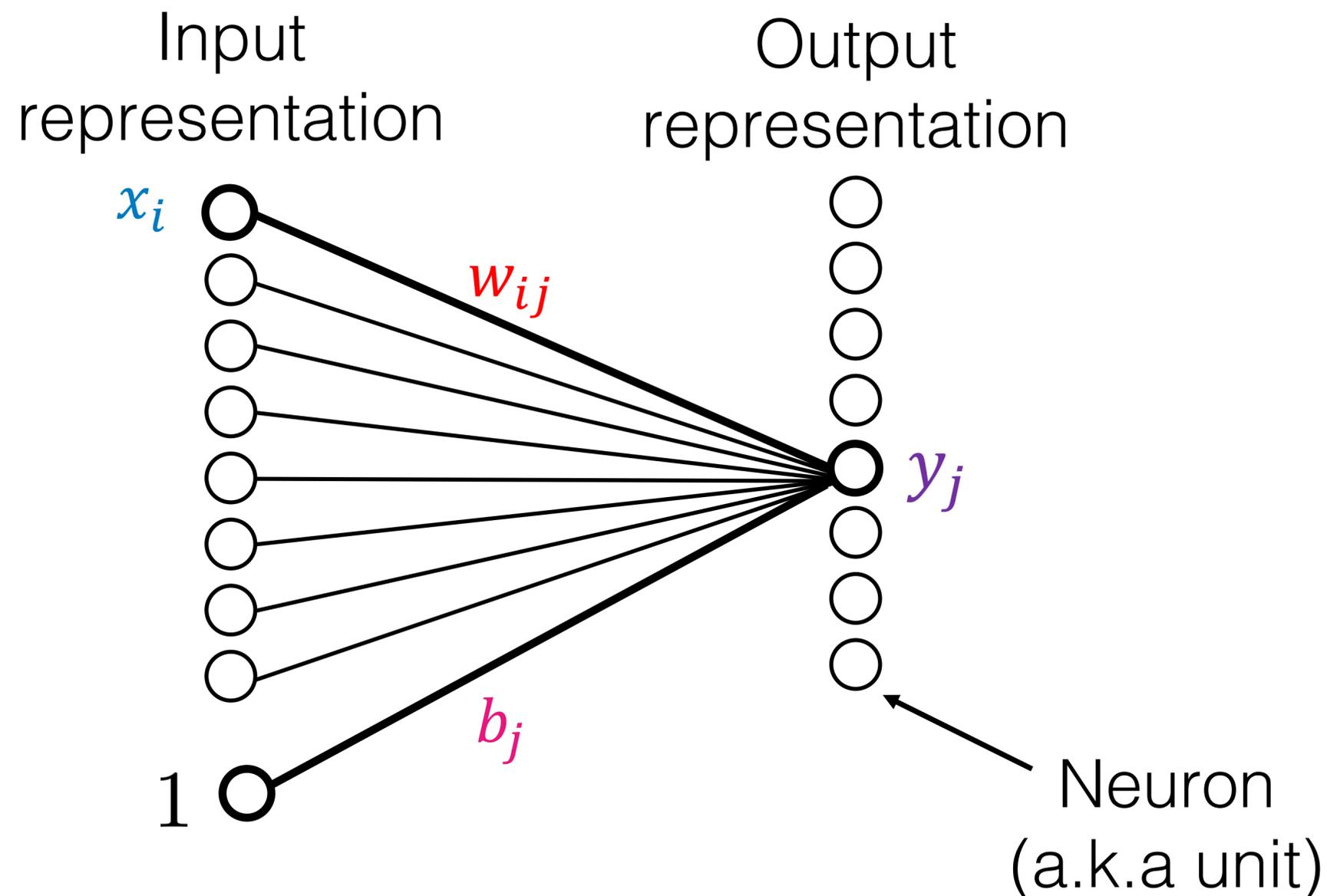
weights

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

# Computation in a neural net

Lets say we have some 1D input that we want to convert to some new feature space:

## Linear layer



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

The equation shows the computation of the output  $y_j$ . The term  $w_{ij} x_i$  is labeled "weights" with an arrow pointing to  $w_{ij}$ . The term  $b_j$  is labeled "bias" with an arrow pointing to  $b_j$ .

# Computation in a neural net – Matrix Multiplication

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

$$\sum_i w_{ij} x_i = \mathbf{x} \cdot \mathbf{w}_j = \mathbf{x}^T \mathbf{w}_j$$

↓

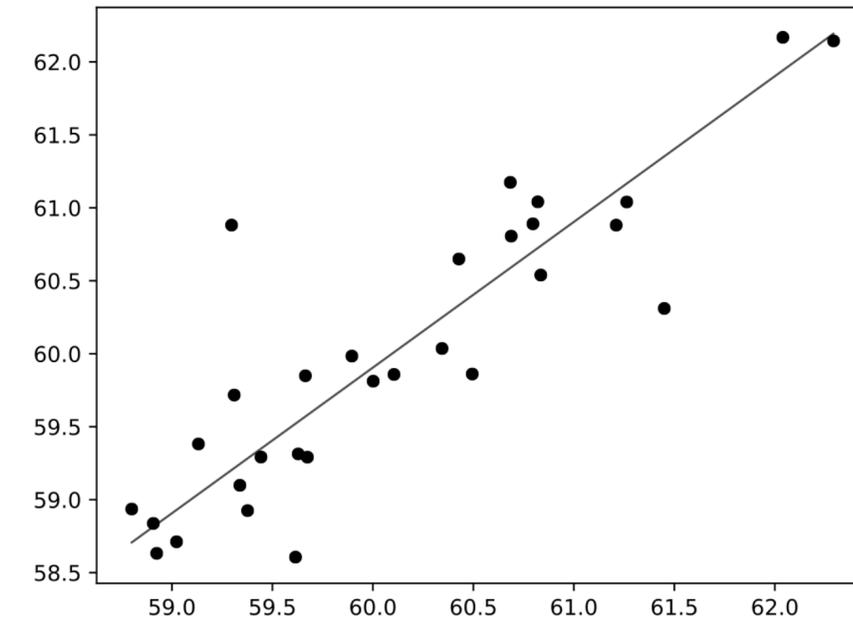
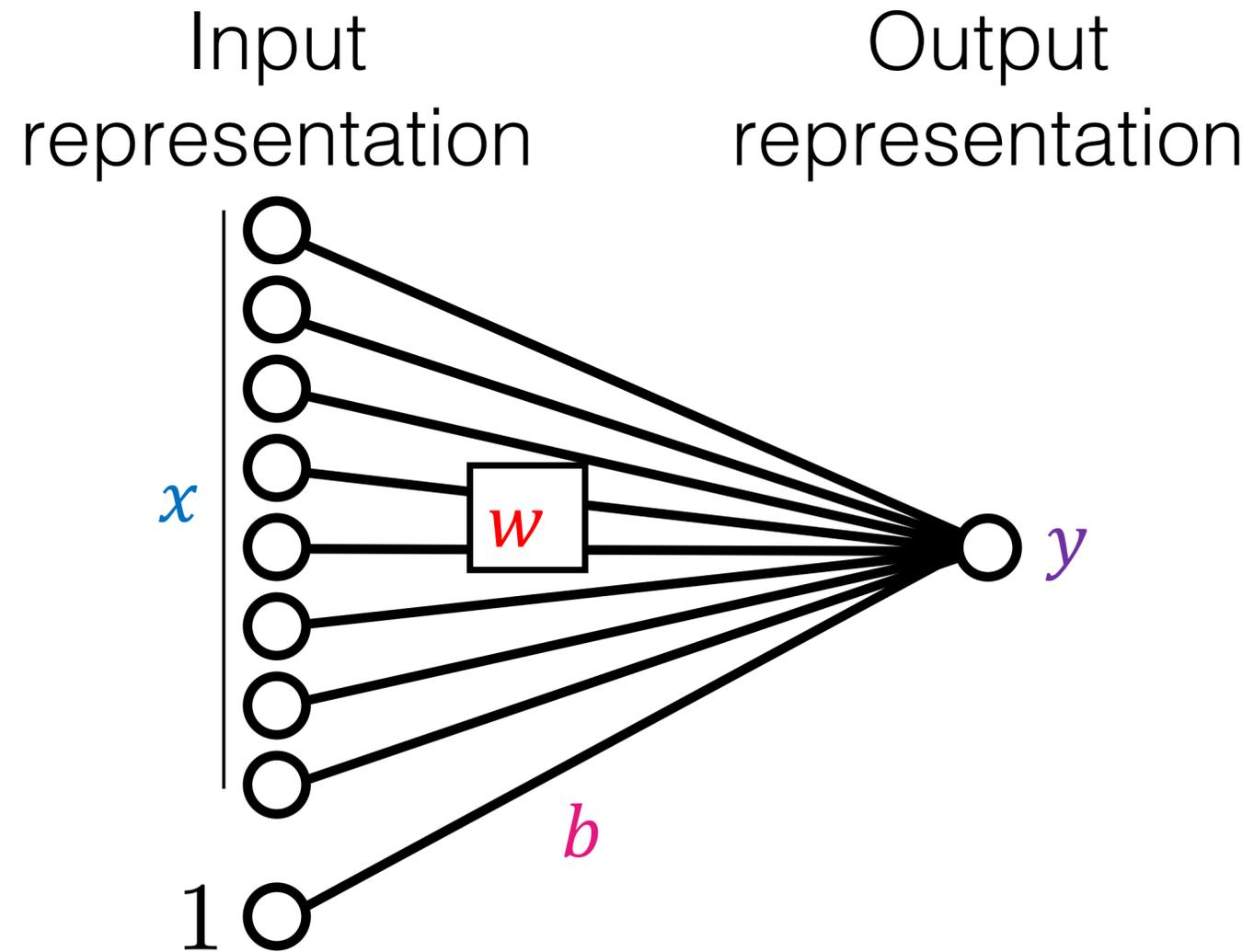
$$y_j = \mathbf{x}^T \mathbf{w}_j + b_j$$

Vector of all input units →  $\mathbf{x}^T$        $\mathbf{w}_j$  ← Vector of weights

$$[x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_n \end{bmatrix} + b_j = [x_1 \ x_2 \ \dots \ x_n \ 1] \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_n \\ b_j \end{bmatrix}$$

# Example: Linear Regression

## Linear layer

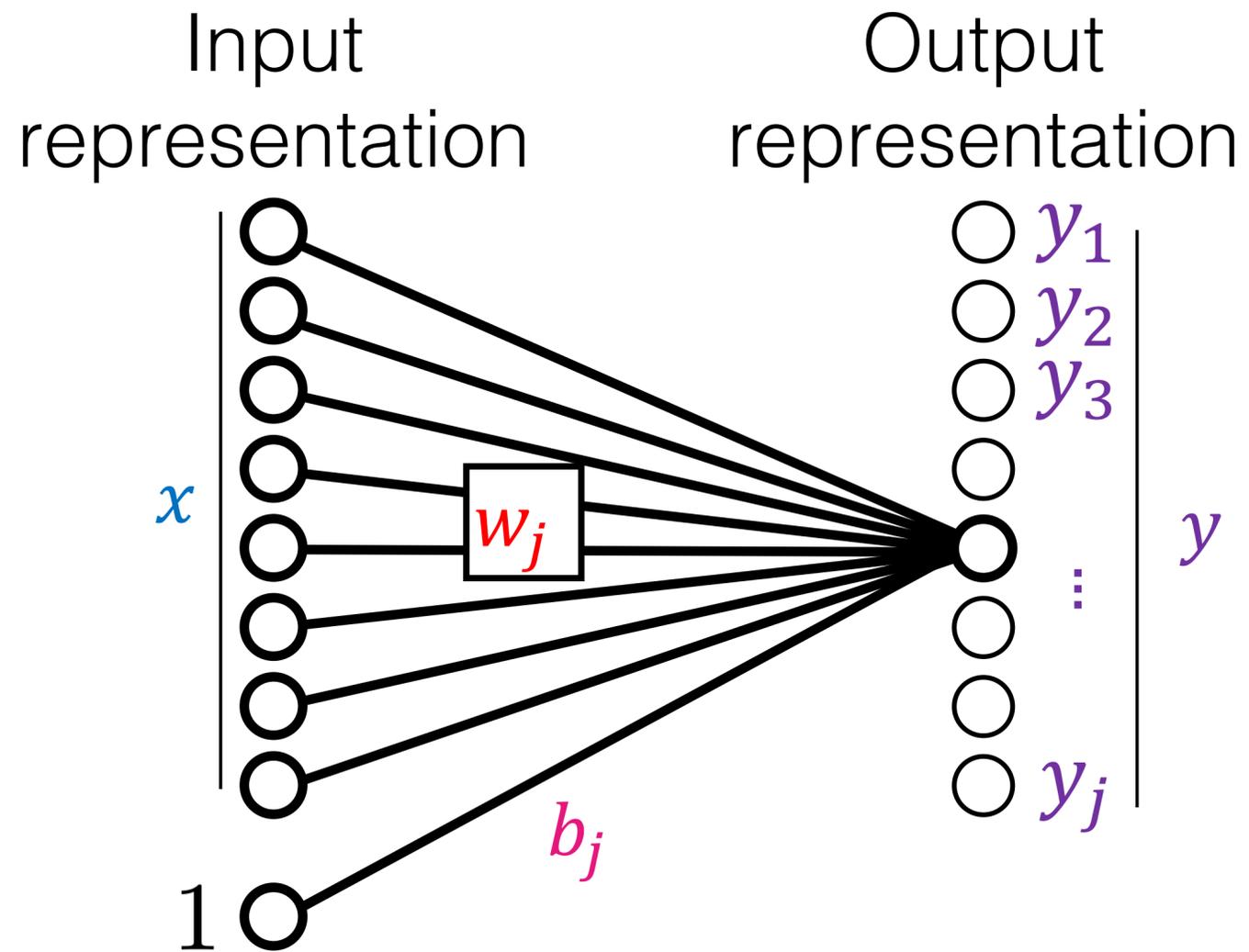


$$f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$$

# Computation in a neural net – Full Layer

## Linear layer

$$y = Wx + b$$

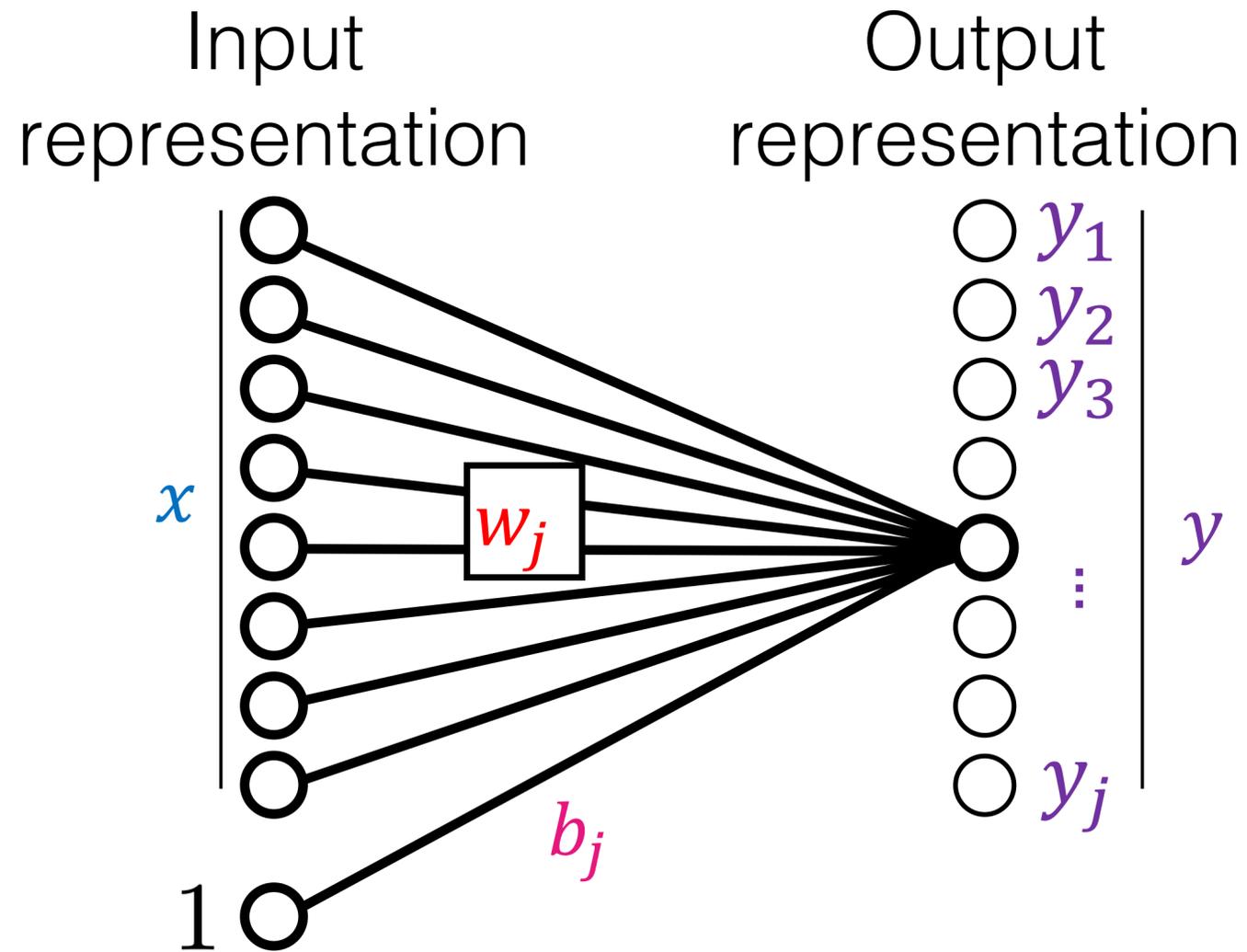


$$\begin{bmatrix} W_{11} & \cdots & W_{1n} \\ \vdots & \ddots & \vdots \\ W_{j1} & \cdots & W_{jn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \cdots \\ b_j \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_j \end{bmatrix}$$

parameters of the model:  $\theta = \{W, b\}$

# Computation in a neural net – Full Layer

## Linear layer



## Full layer

$$y = Wx + b$$

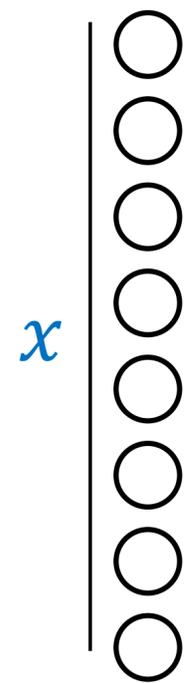
$$\begin{bmatrix} w_{11} & \cdots & w_{jn} & b_1 \\ \vdots & \ddots & \vdots & \vdots \\ w_{j1} & \cdots & w_{jn} & b_j \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \\ 1 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_j \end{bmatrix}$$

Can again simplify notation by appending a 1 to  $\mathbf{x}$

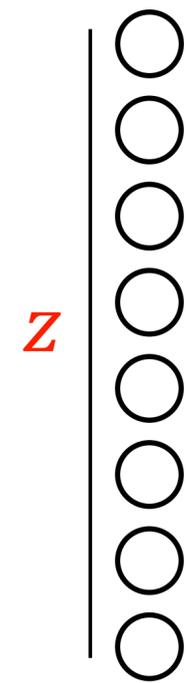
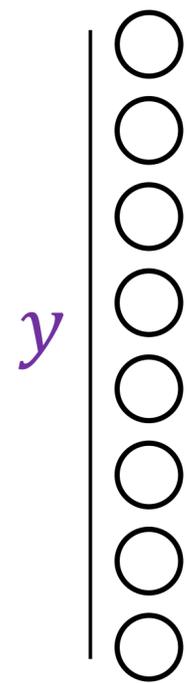
# Computation in a neural net – Recap

We can now transform our input representation vector into some output representation vector using a bunch of linear combinations of the input:

Input  
representation

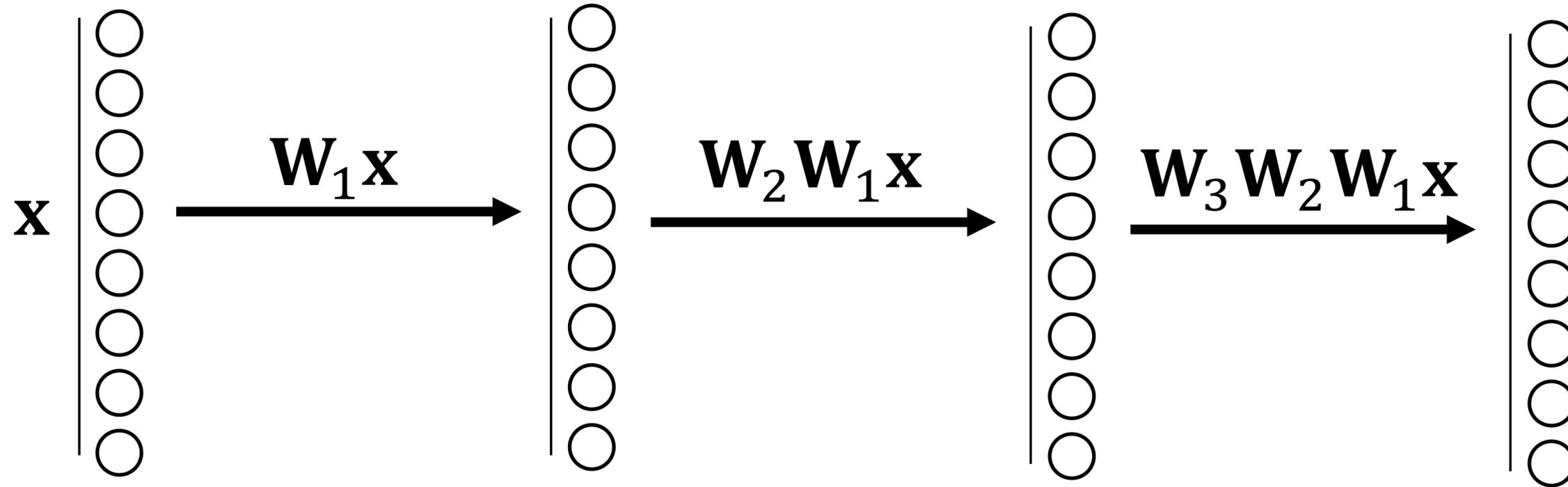


Output  
representation



We can repeat this as  
many times as we want!

What is the problem with this idea?



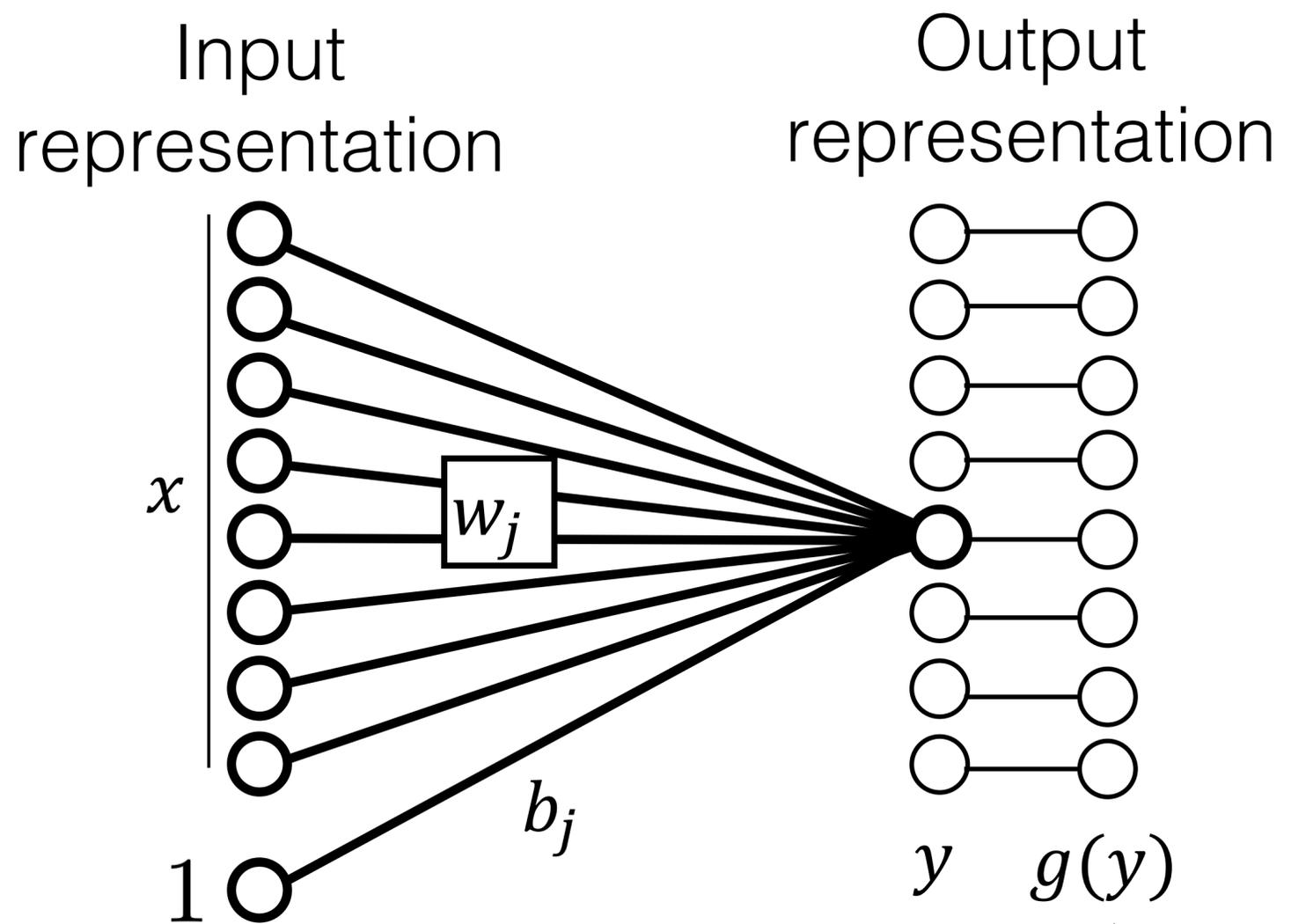
Can be expressed as single linear layer!

$$\left( \prod_i \mathbf{W}_i \right) \mathbf{x} = \hat{\mathbf{W}} \mathbf{x}$$

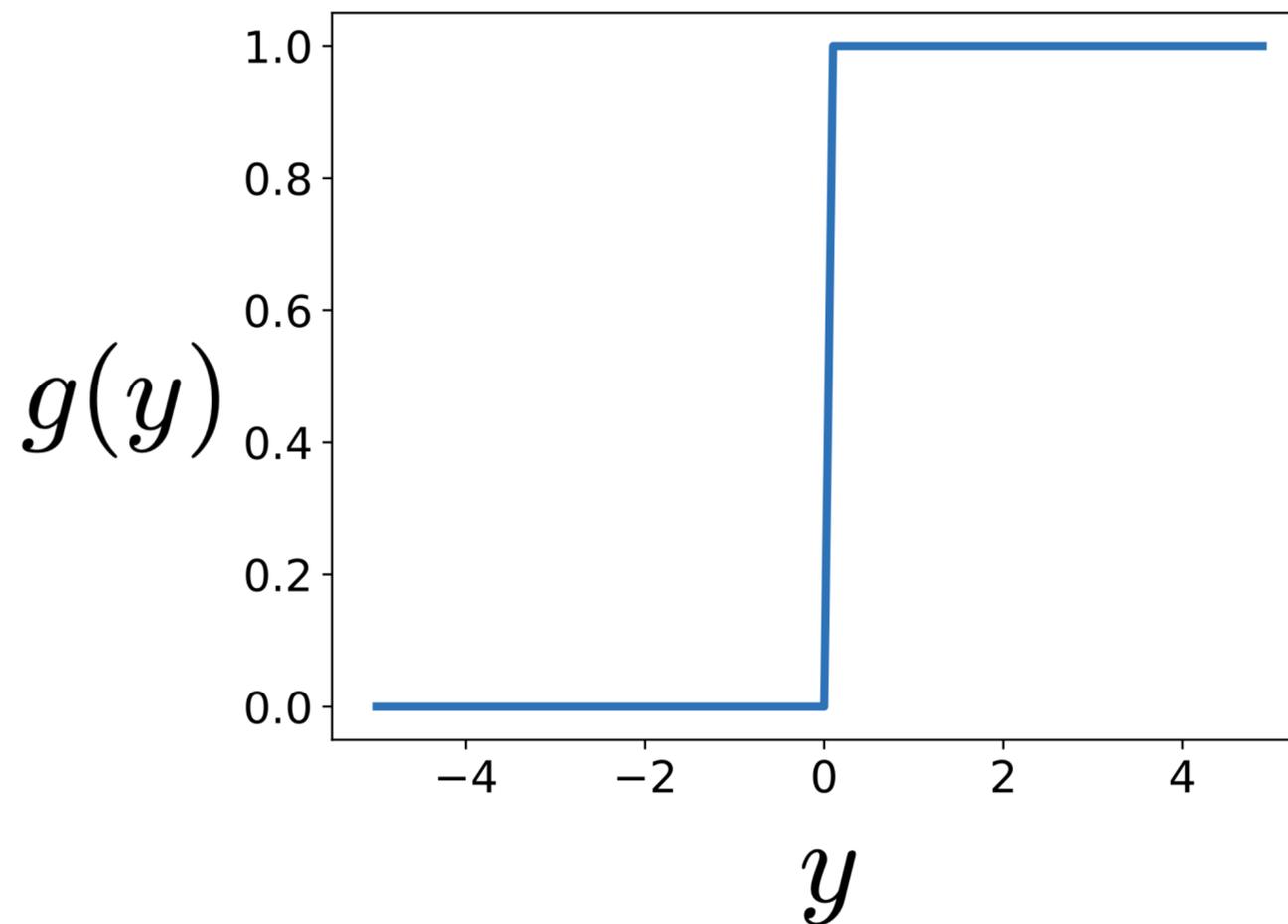
Limited power: can't solve XOR :(

# Solution: simple nonlinearity

## Linear layer

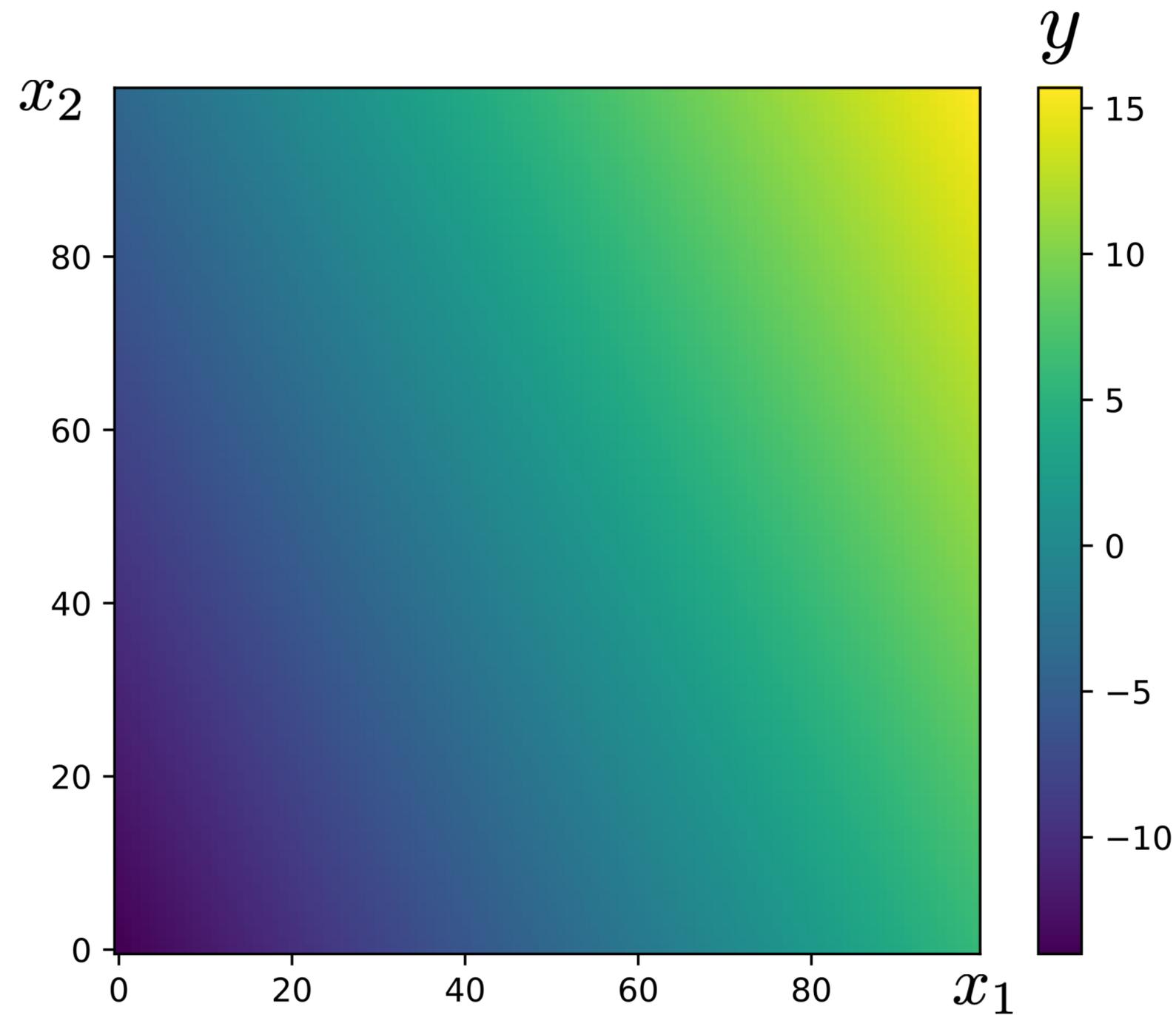


$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



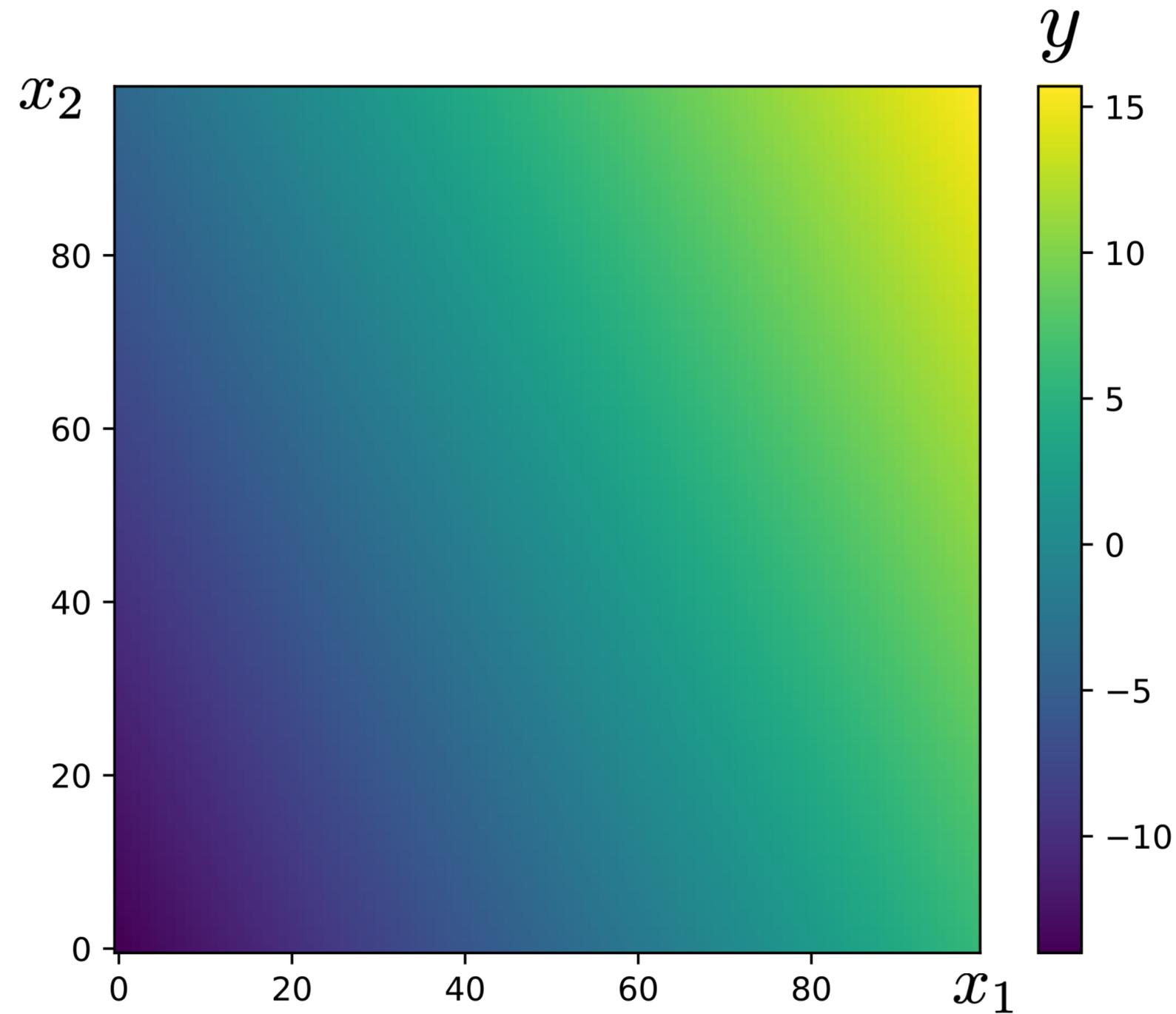
**Pointwise  
Non-linearity**

# Example: linear classification with a perceptron



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

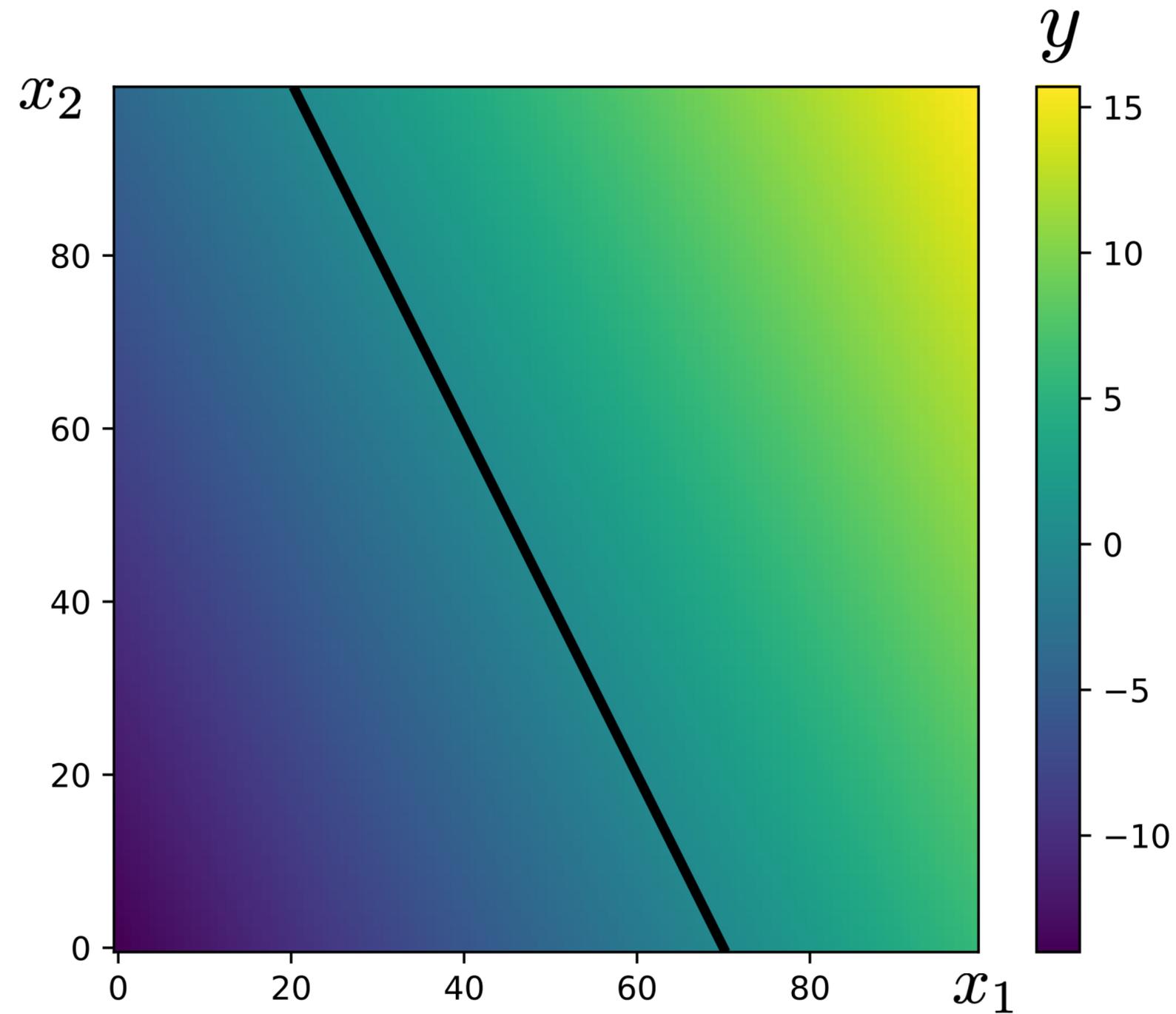
# Example: linear classification with a perceptron



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

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

# Example: linear classification with a perceptron



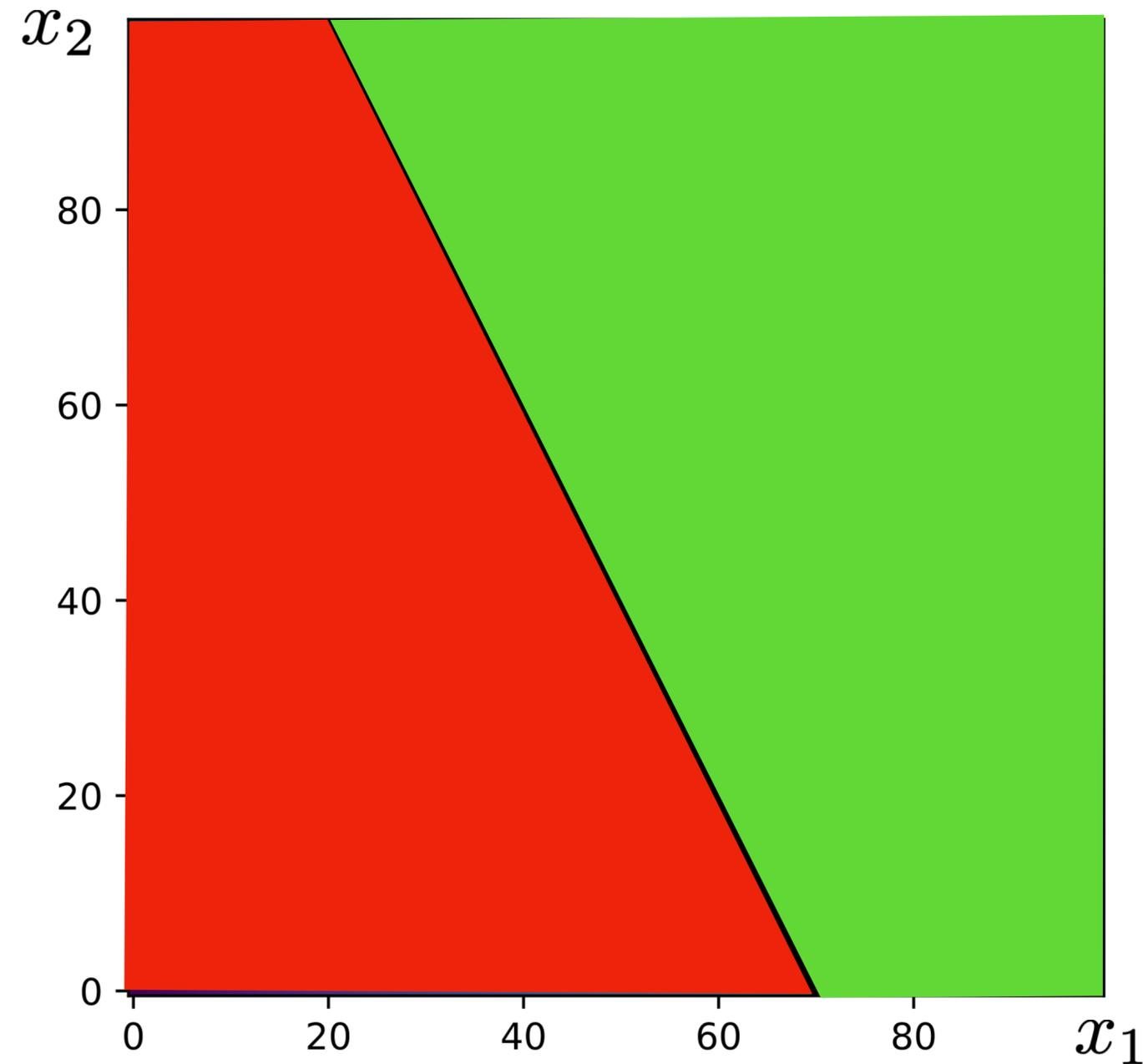
$$y = \mathbf{x}^T \mathbf{w} + b$$

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

“when  $y$  is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)”

# Example: linear classification with a perceptron

$$g(y)$$



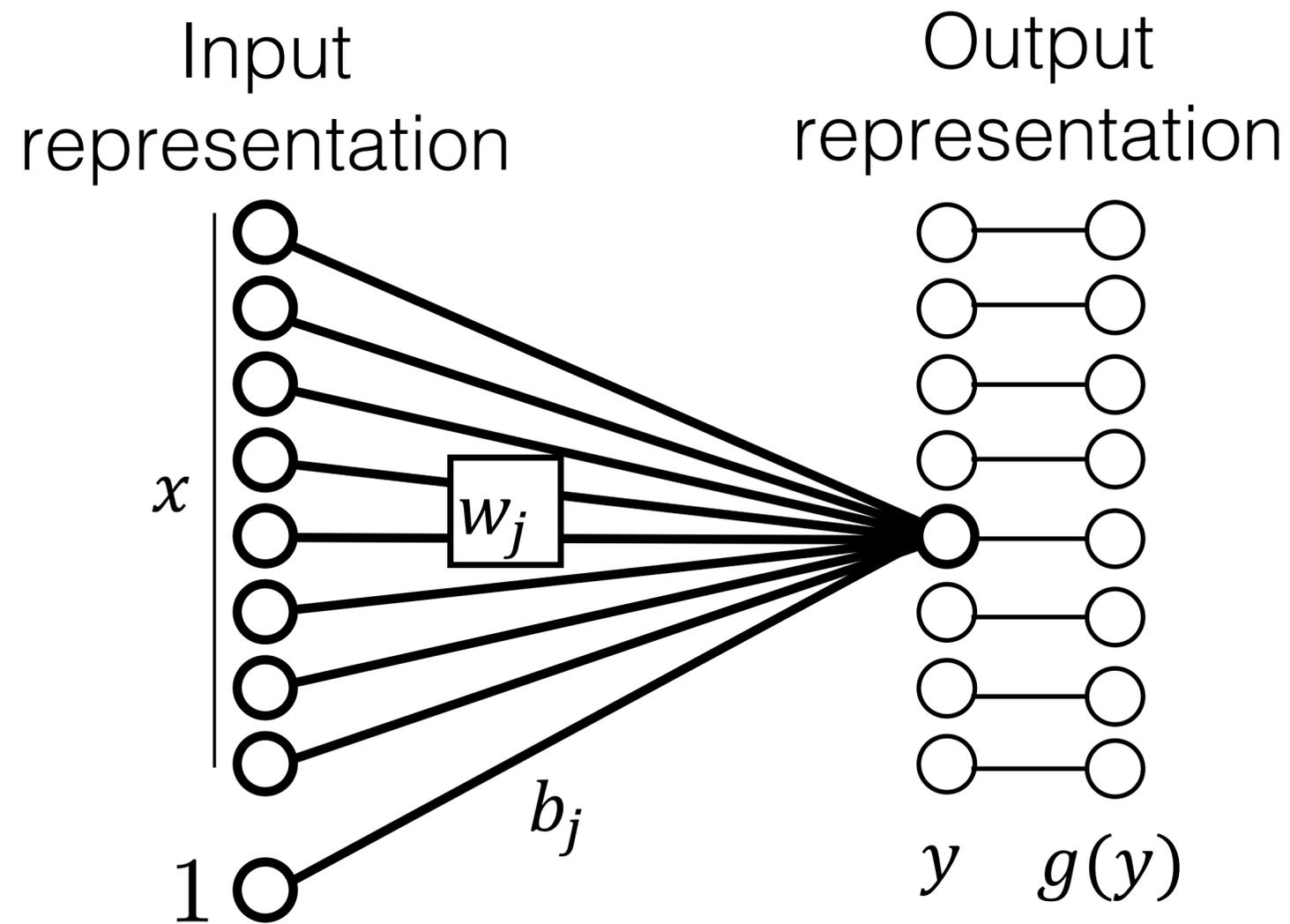
$$y = \mathbf{x}^T \mathbf{w} + b$$

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

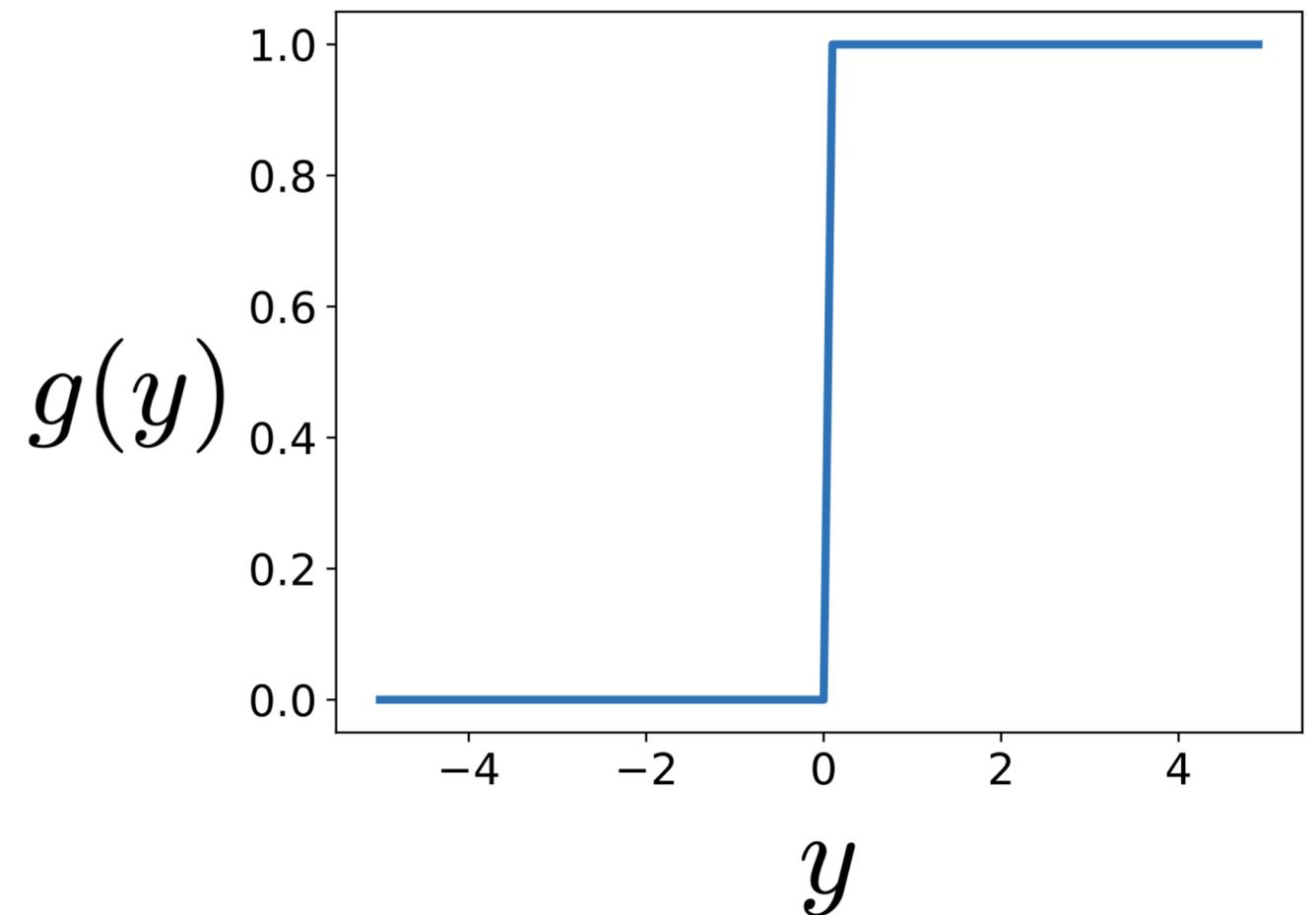
“when  $y$  is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)”

# Computation in a neural net - nonlinearity

## Linear layer



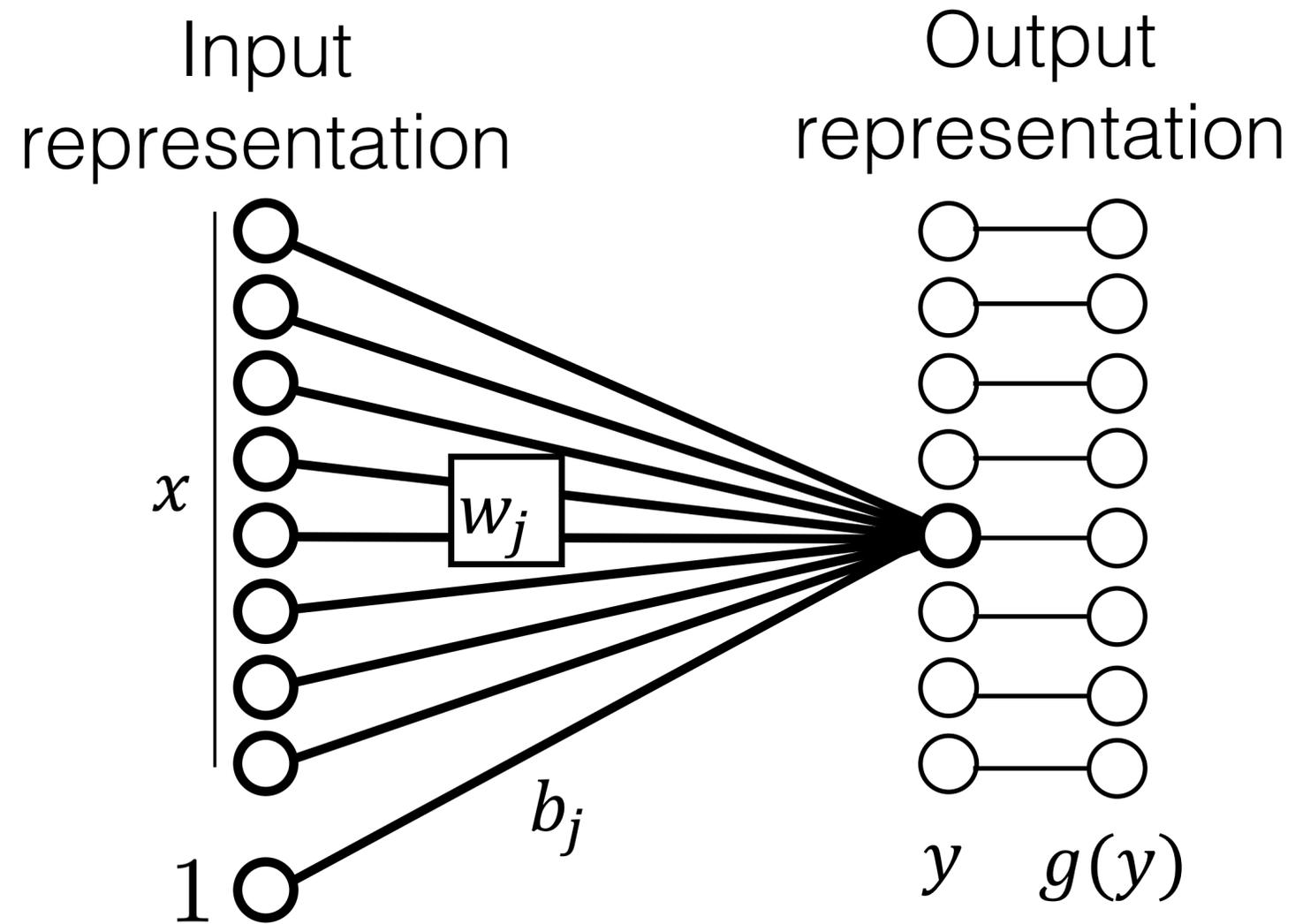
$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



Can't use with gradient descent,  $\frac{\partial}{\partial y} g = 0$

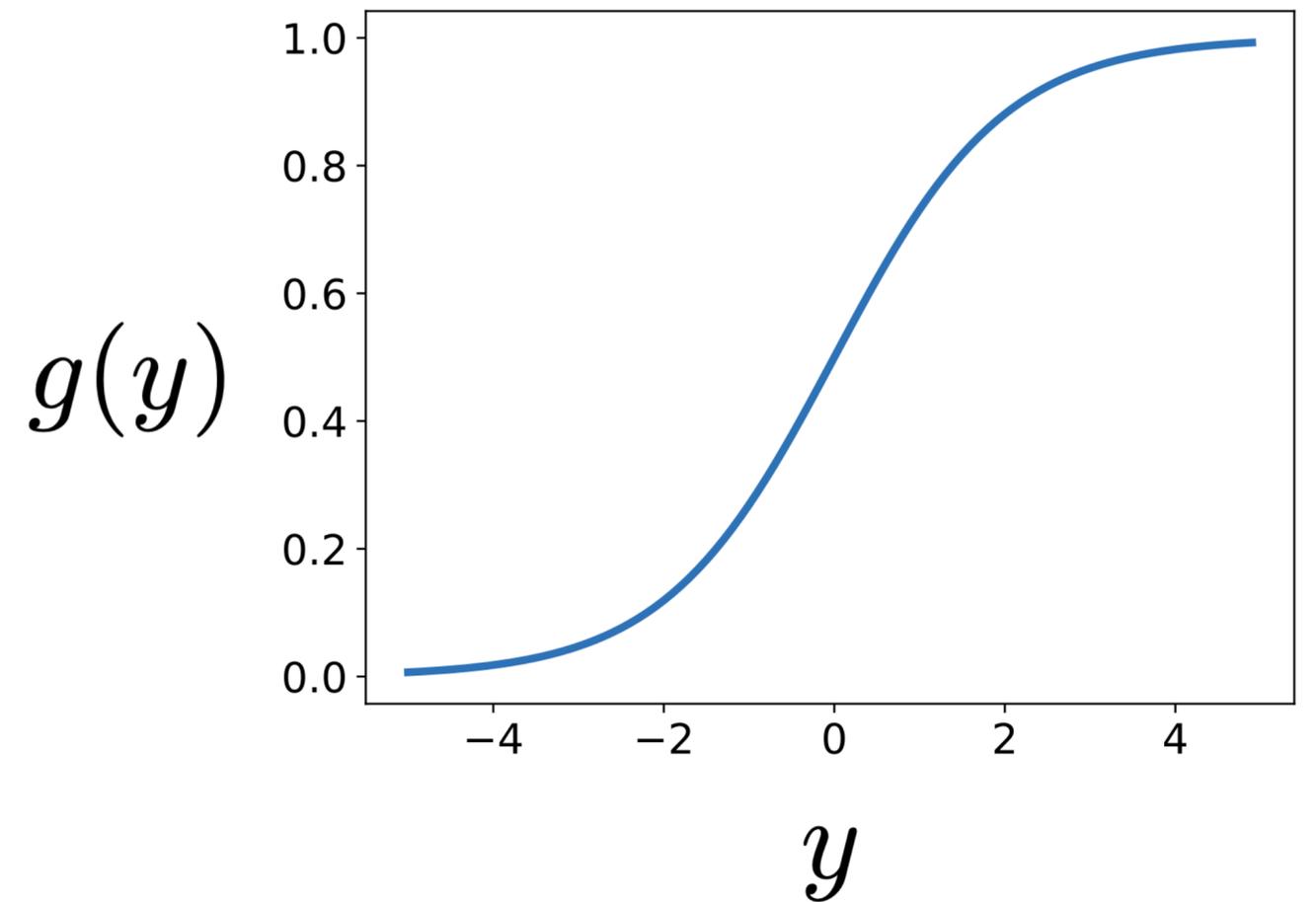
# Computation in a neural net - nonlinearity

## Linear layer



## Sigmoid

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

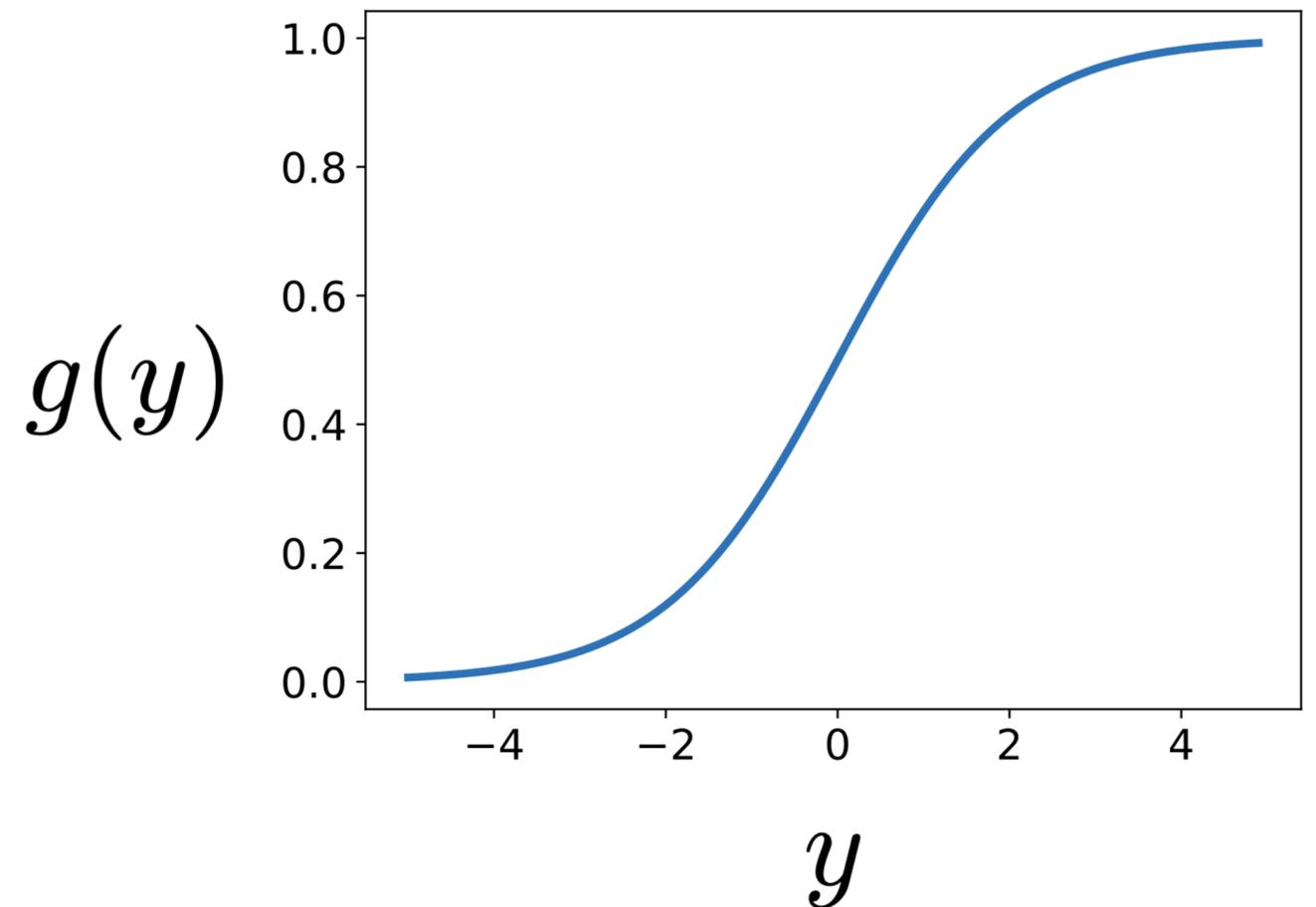


# Computation in a neural net - nonlinearity

- Bounded between  $[0, 1]$
- Saturation for large +/- inputs
- Gradients go to zero
- Better in practice to use:  $\tanh(y)$   
 $= 2g(y) - 1$

## Sigmoid

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



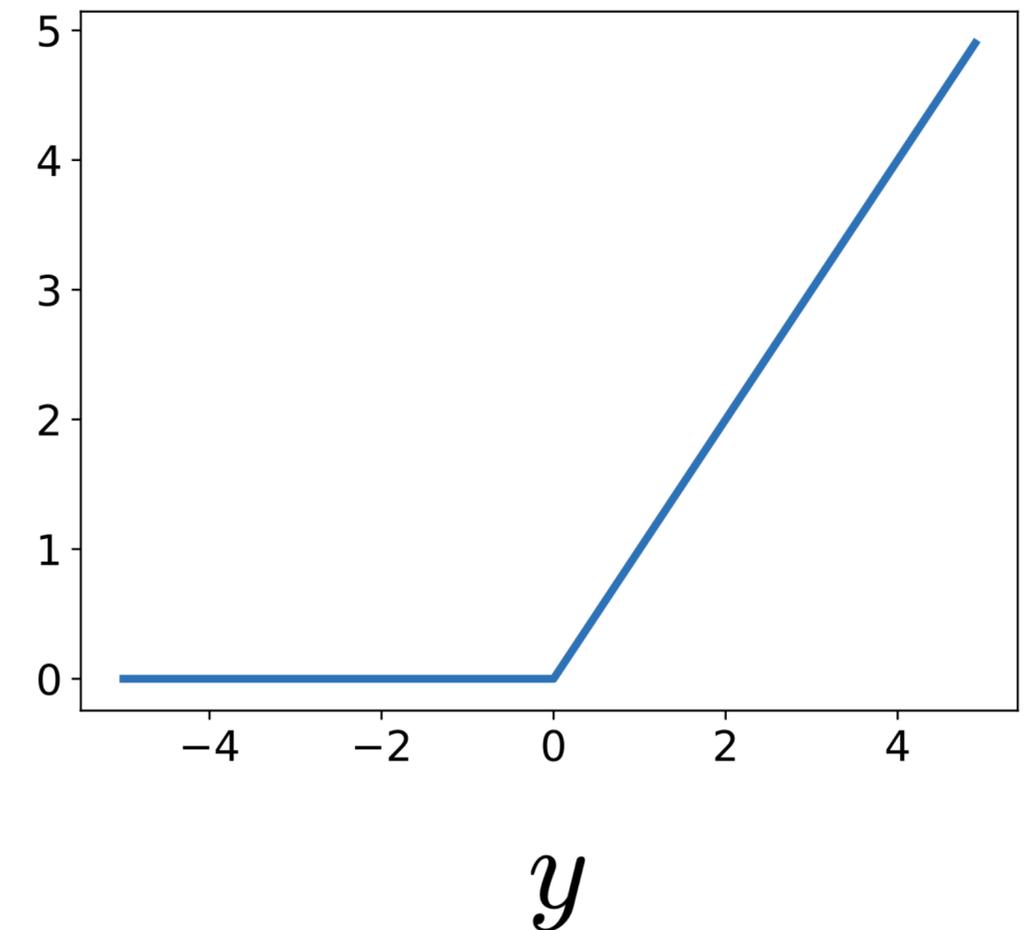
# Computation in a neural net — nonlinearity

- 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 (6x speedup vs. tanh in [Krizhevsky et al. 2012])
- Drawback: if strongly in negative region, unit is dead forever (no gradient).
- Default choice: widely used in current models!

## Rectified linear unit (ReLU)

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

$g(y)$



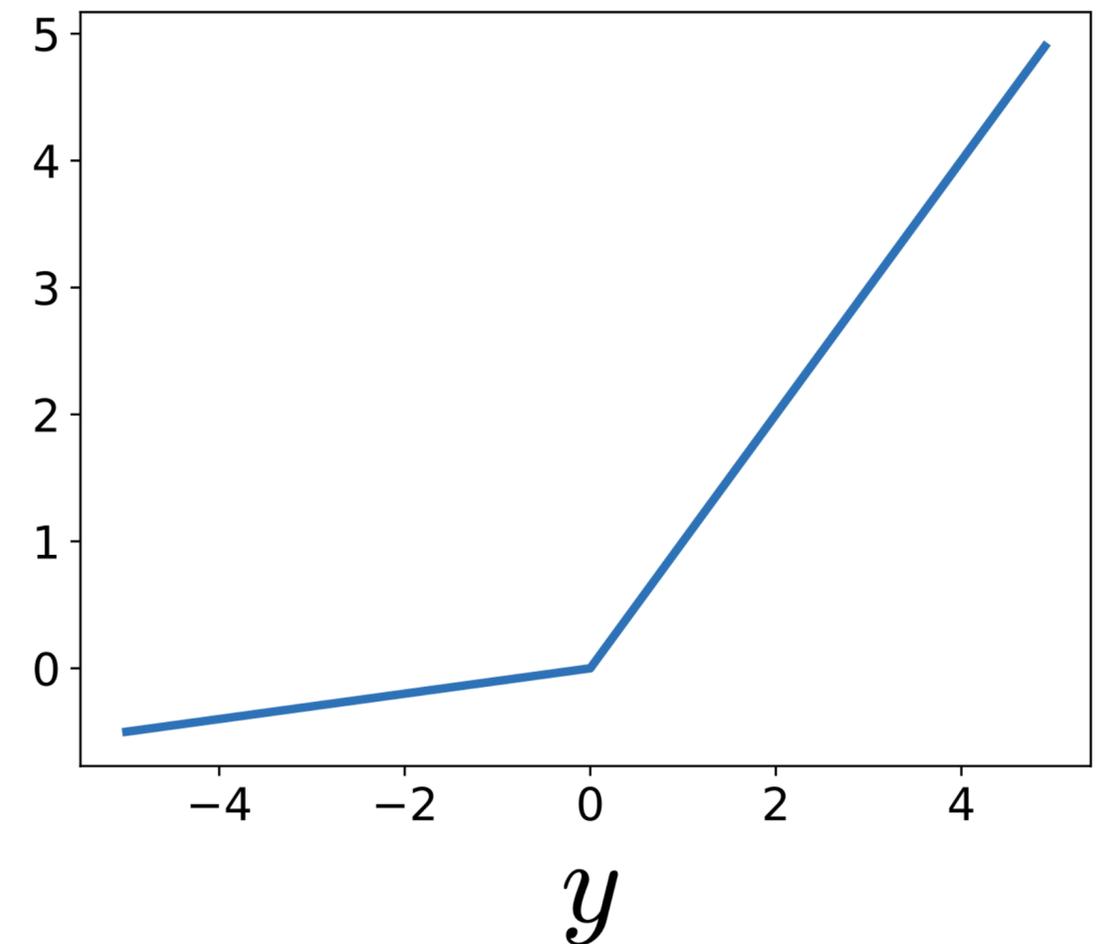
# Computation in a neural net — nonlinearity

- where  $a$  is small (e.g., 0.02)
- Efficient to implement:  $\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$
- Has non-zero gradients everywhere (unlike ReLU)

## Leaky ReLU

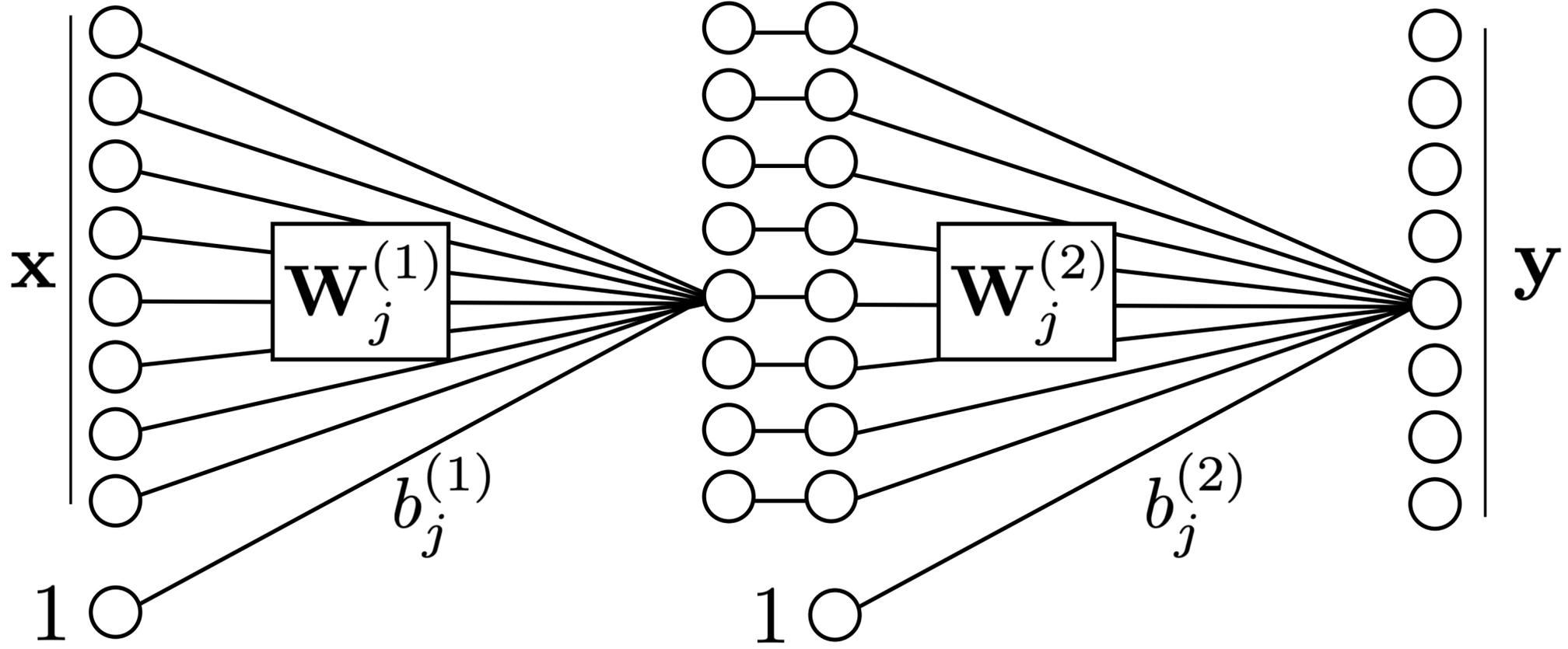
$$g(y) = \begin{cases} \max(0, y), & \text{if } y \geq 0 \\ a \min(0, y), & \text{if } y < 0 \end{cases}$$

$g(y)$



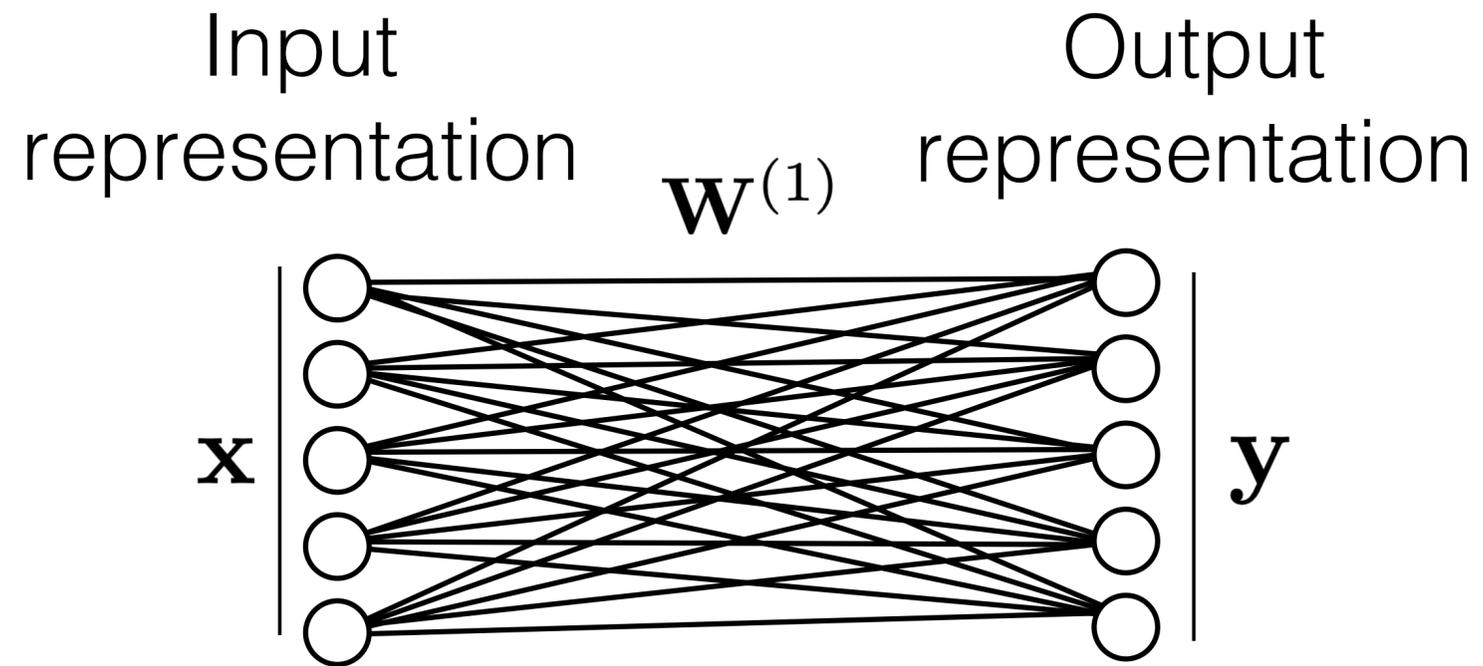
# Stacking layers

Input representation      Intermediate representation      Output representation

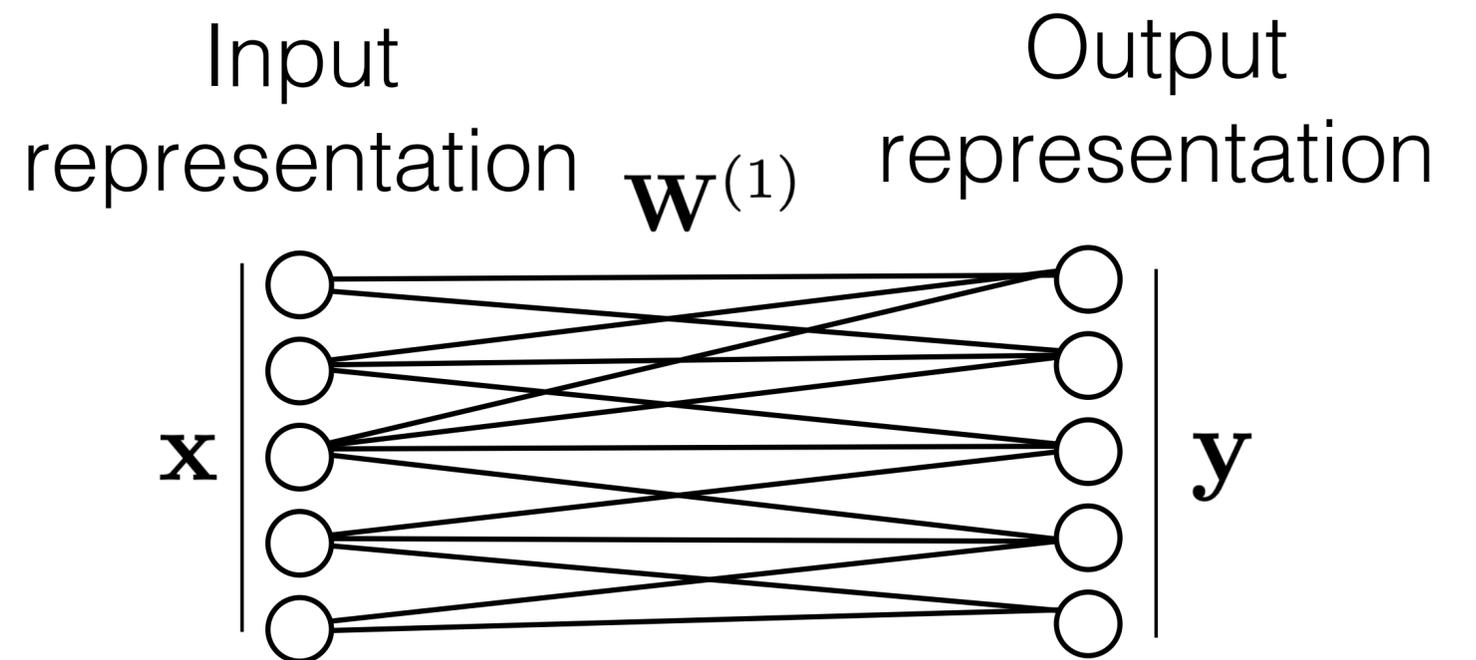


$\mathbf{h}$  = "hidden units"

# Connectivity patterns

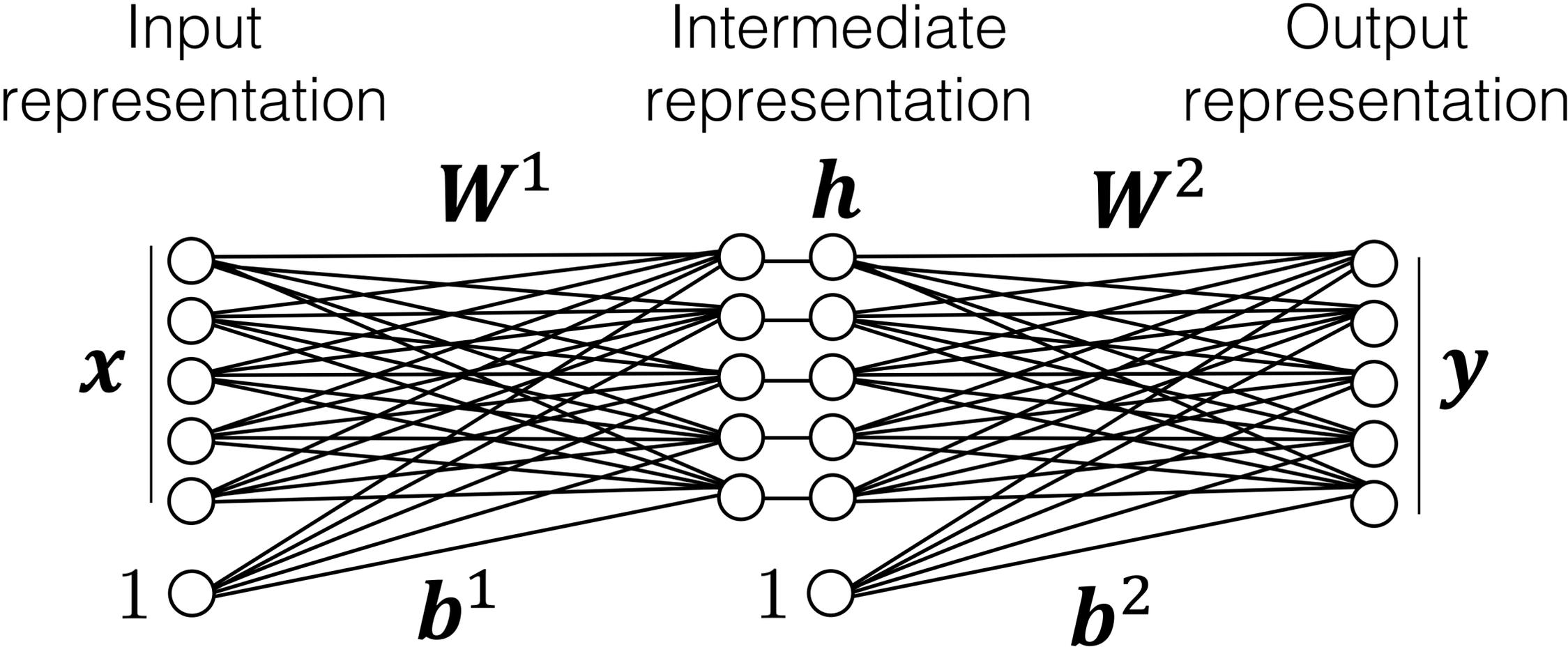


*Fully connected layer*



*Locally connected layer  
(Sparse  $W$ )*

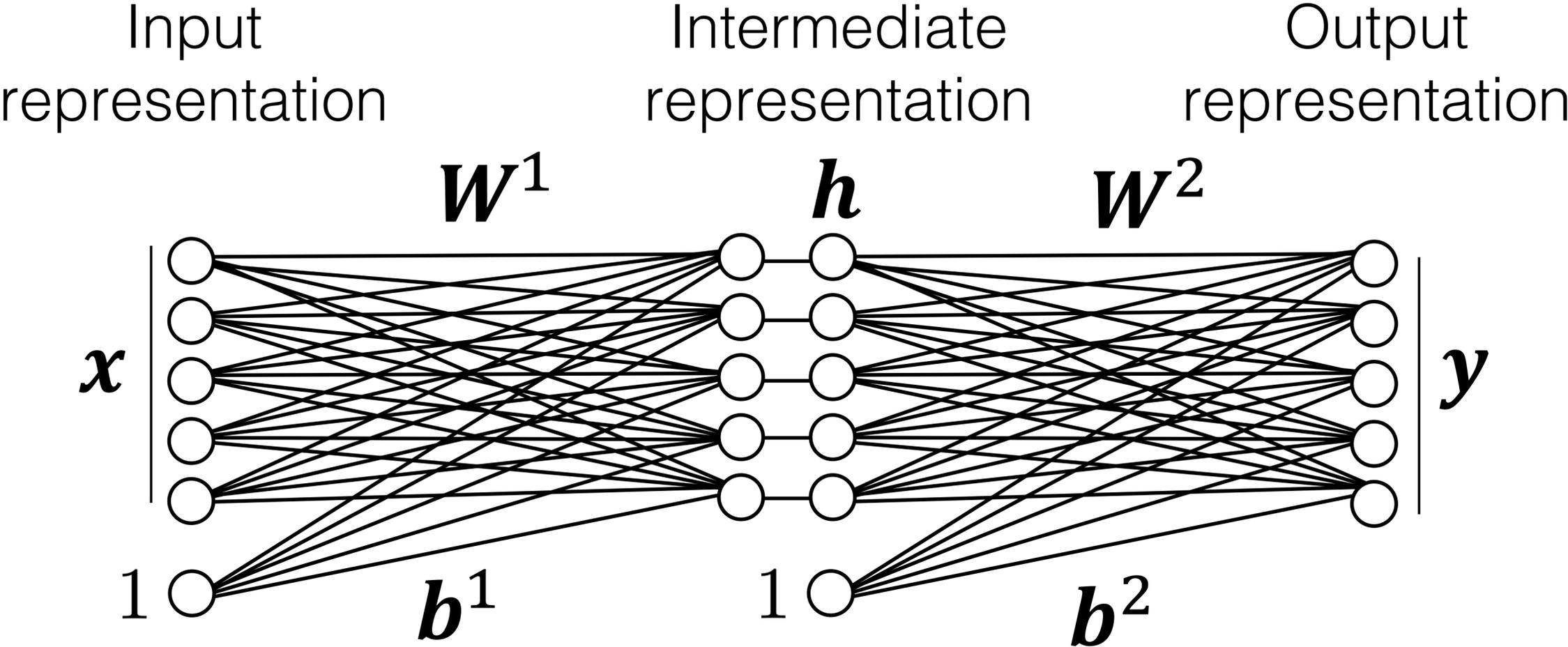
# Stacking layers



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

ReLU  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

# Stacking layers

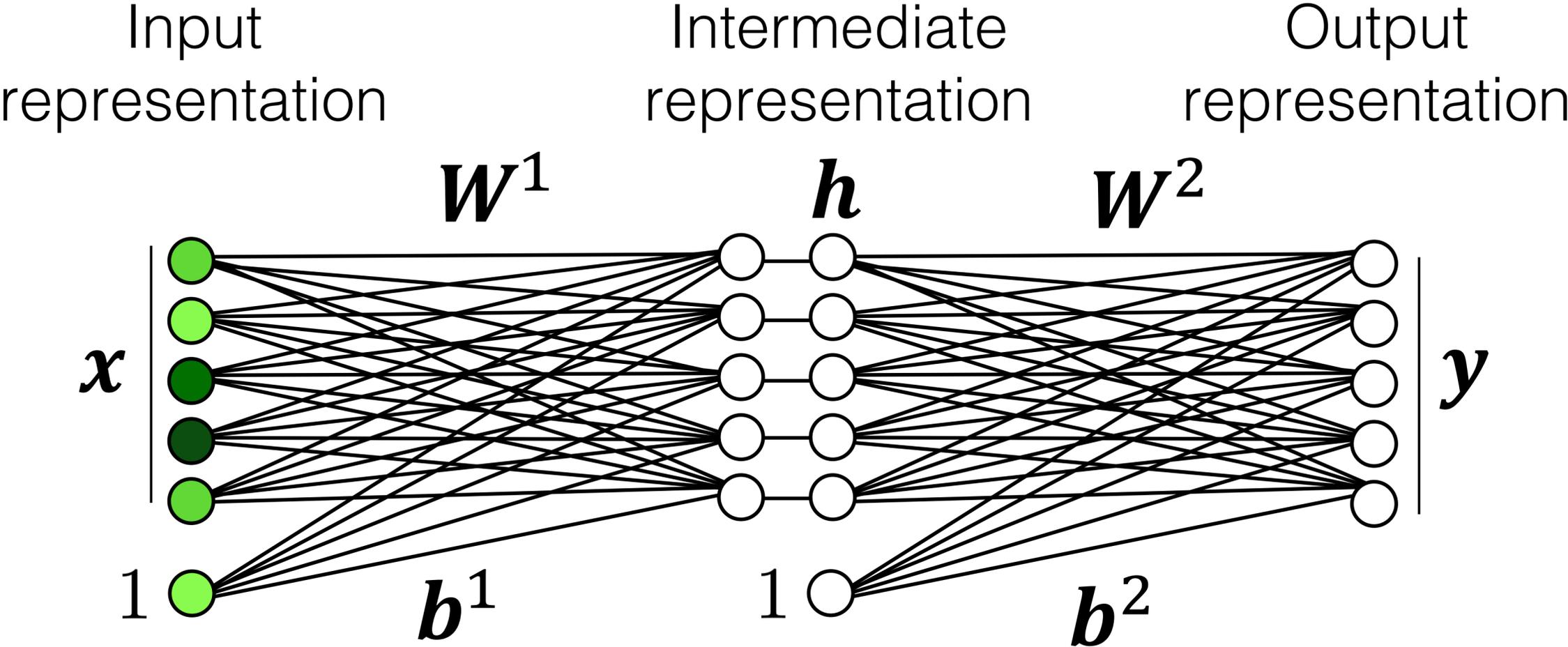


positive  
negative

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

ReLU  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

# Stacking layers

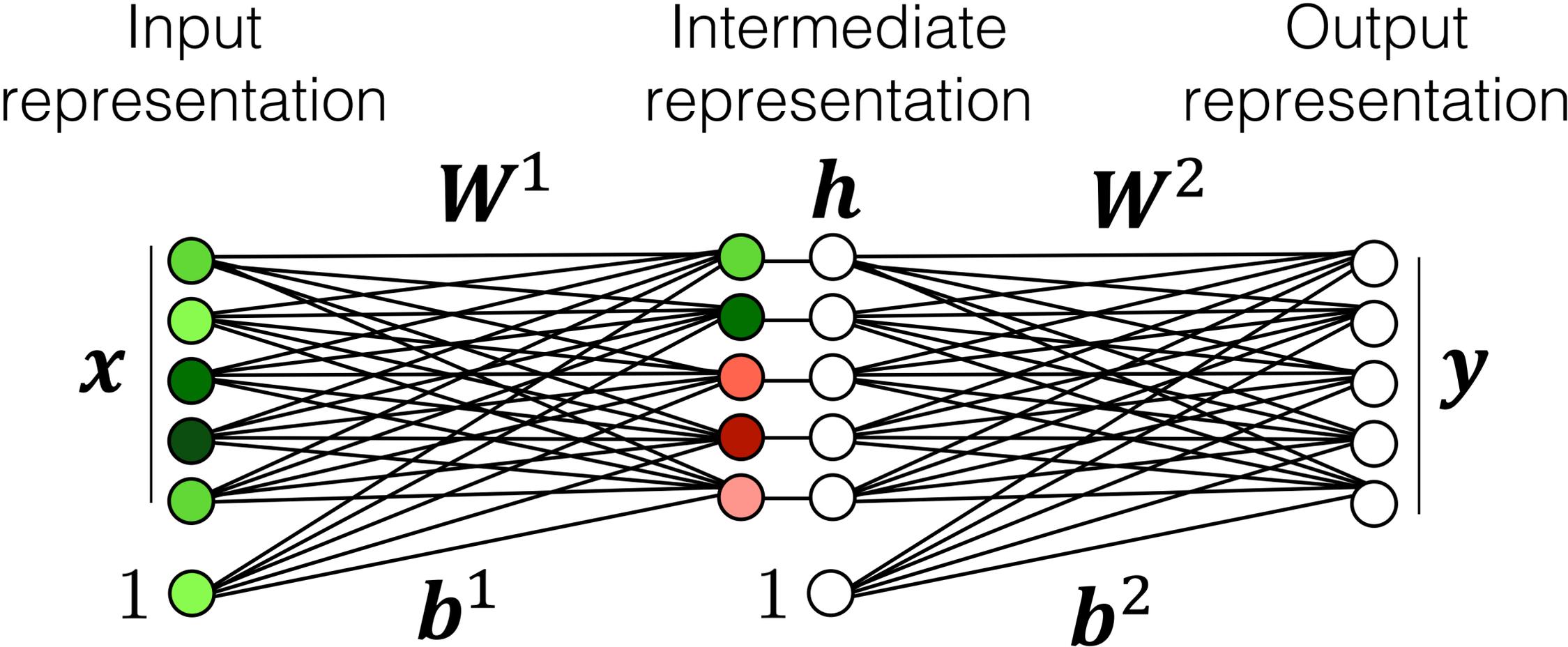


positive  
negative

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

ReLU  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

# Stacking layers

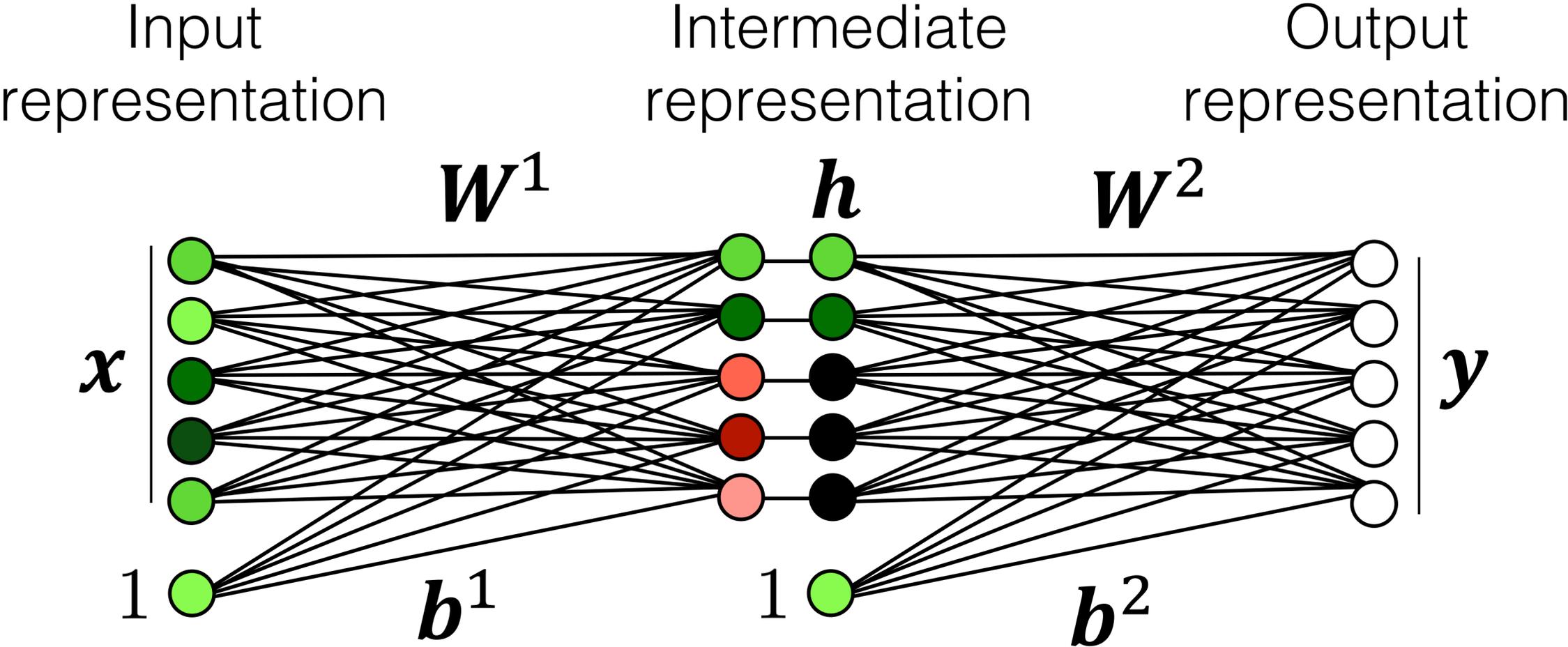


positive  
negative

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

ReLU →  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

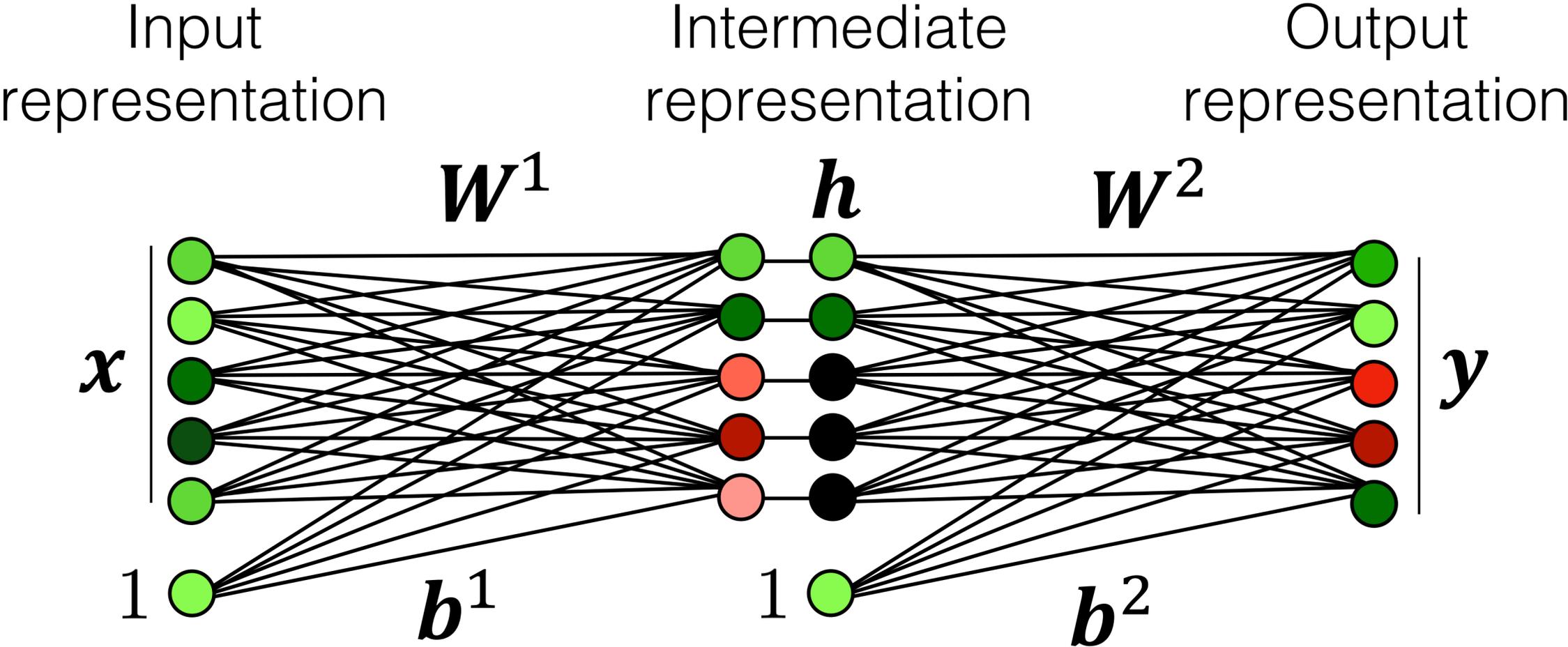
# Stacking layers



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

ReLU →  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

# Stacking layers

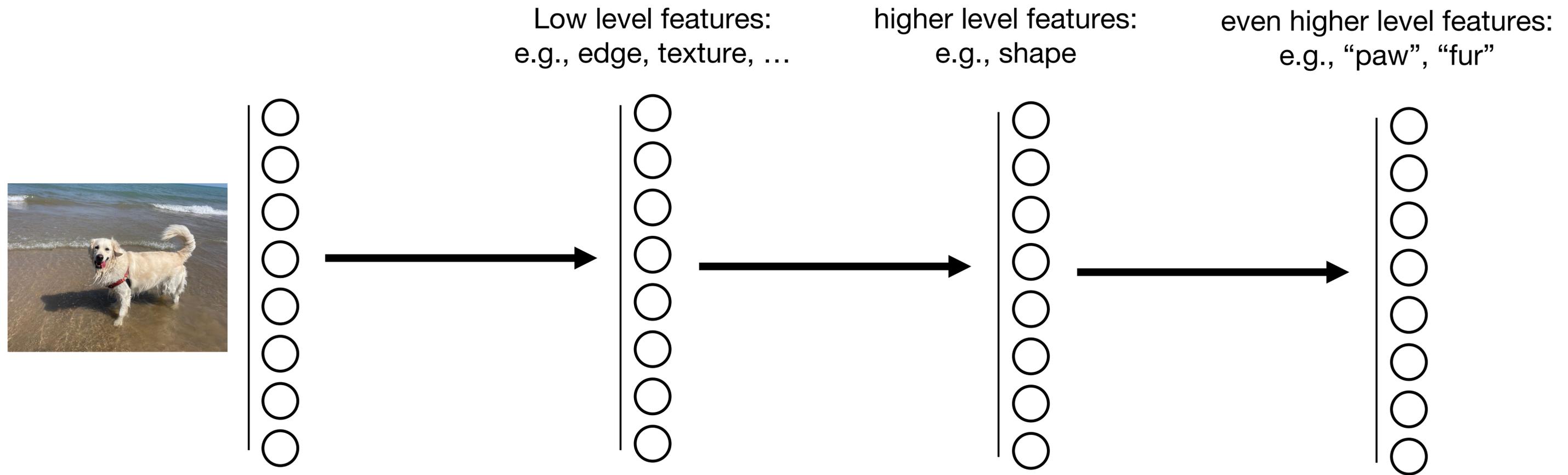


positive  
negative

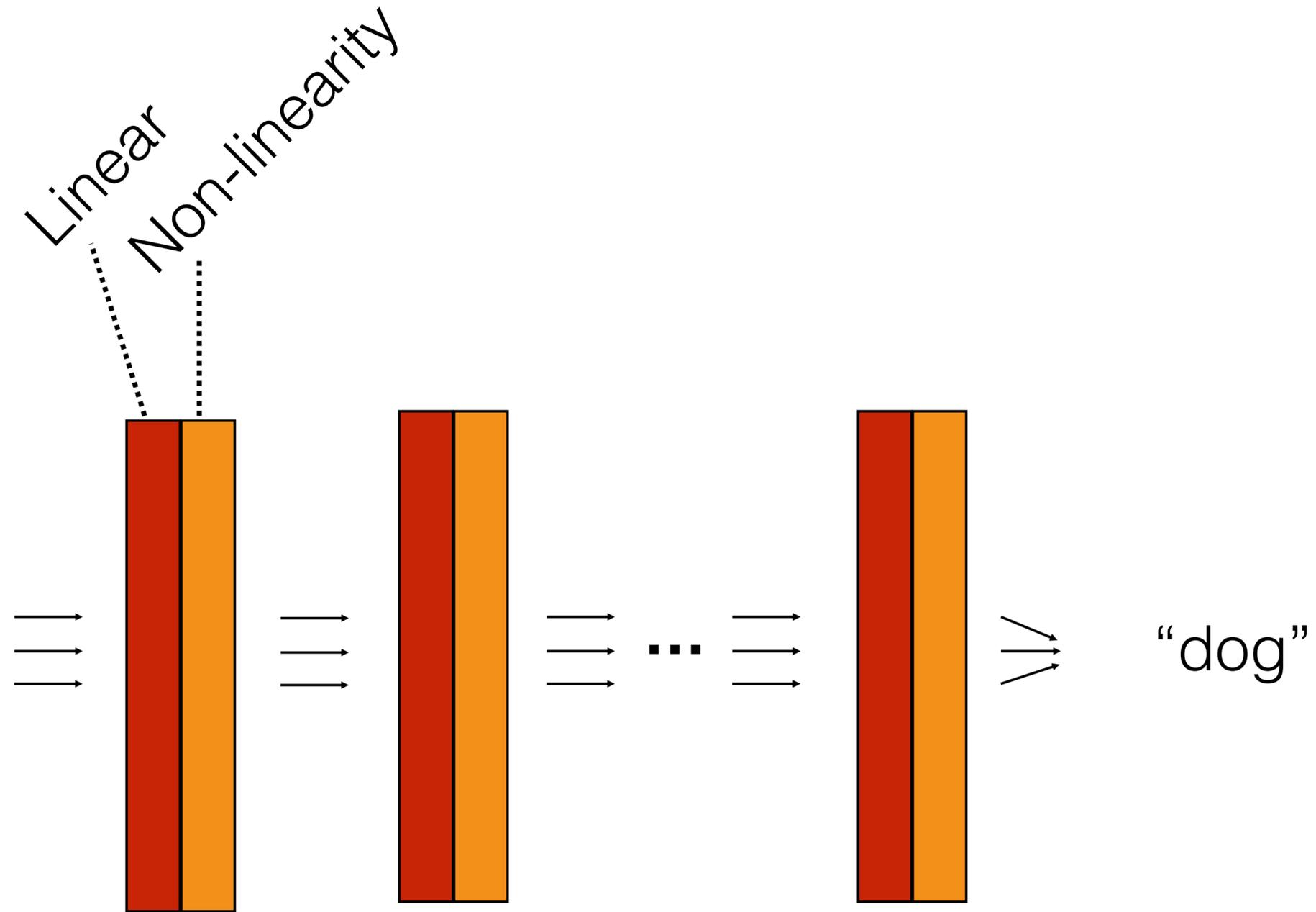
$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU  $\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$

# Stacking layers - What's actually happening?



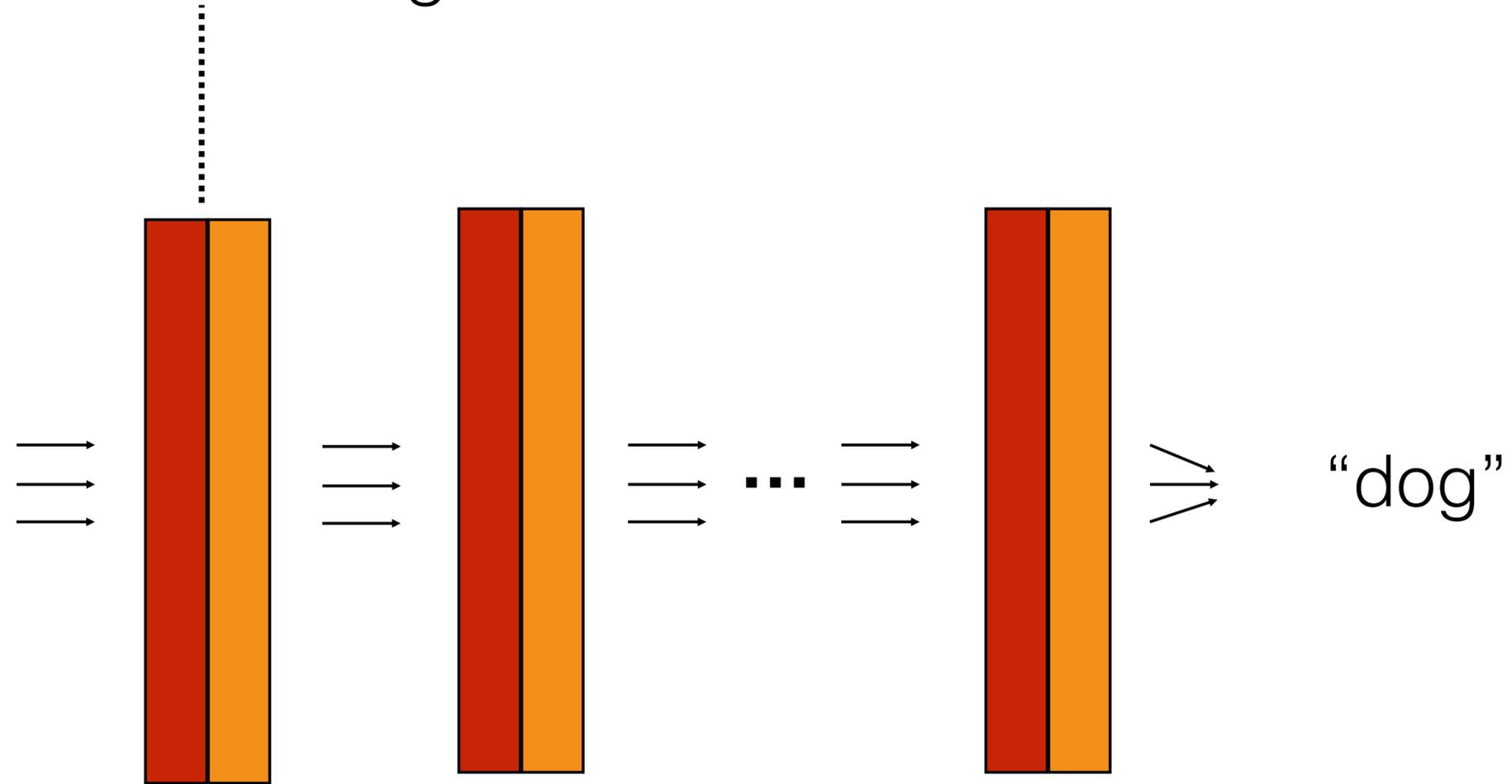
# Deep nets



$$f(x) = f_L(\dots f_3(f_2(f_1(x))))$$

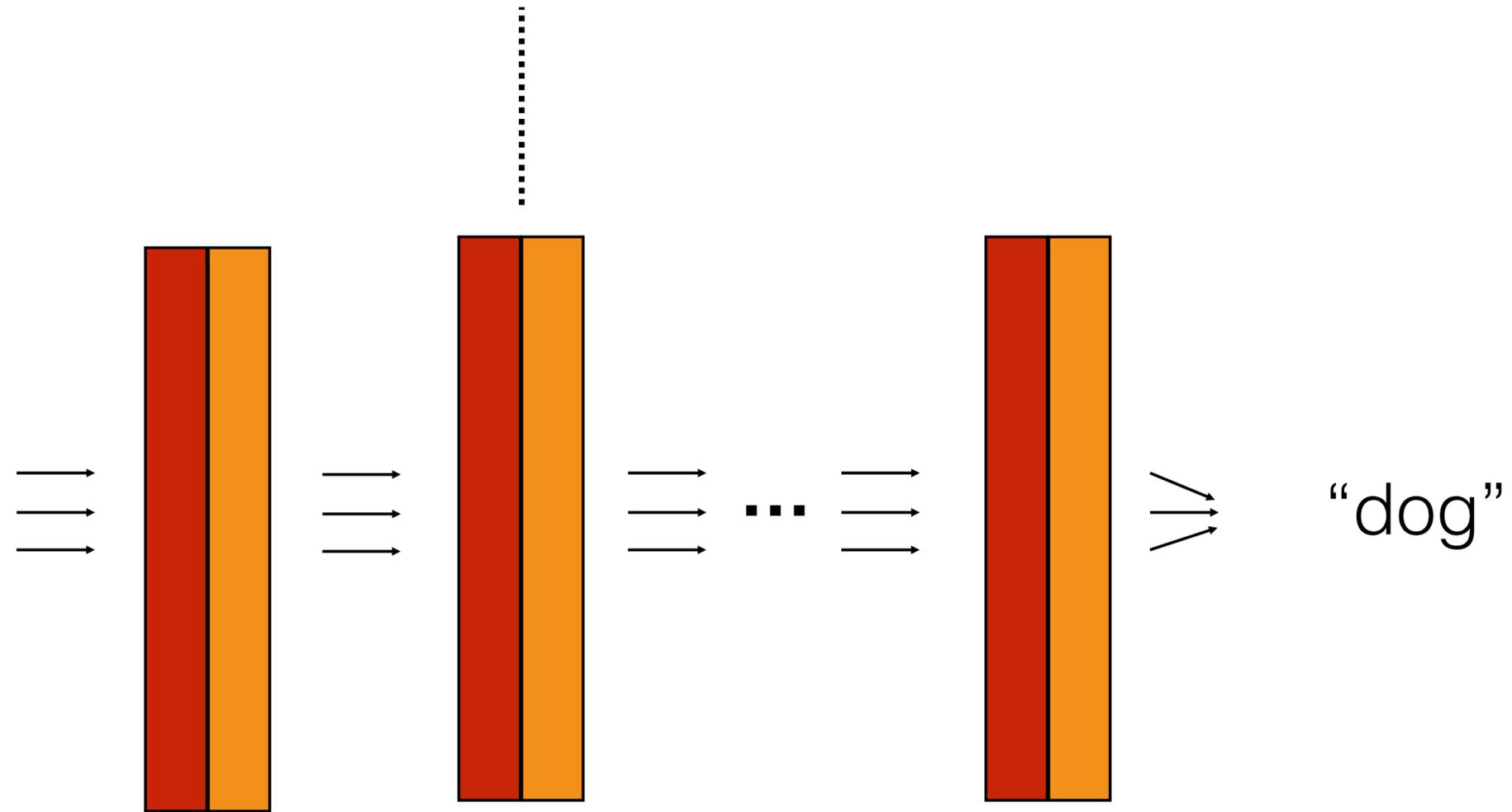
# Deep nets - Intuition

“has horizontal edge”  
“has vertical edge”

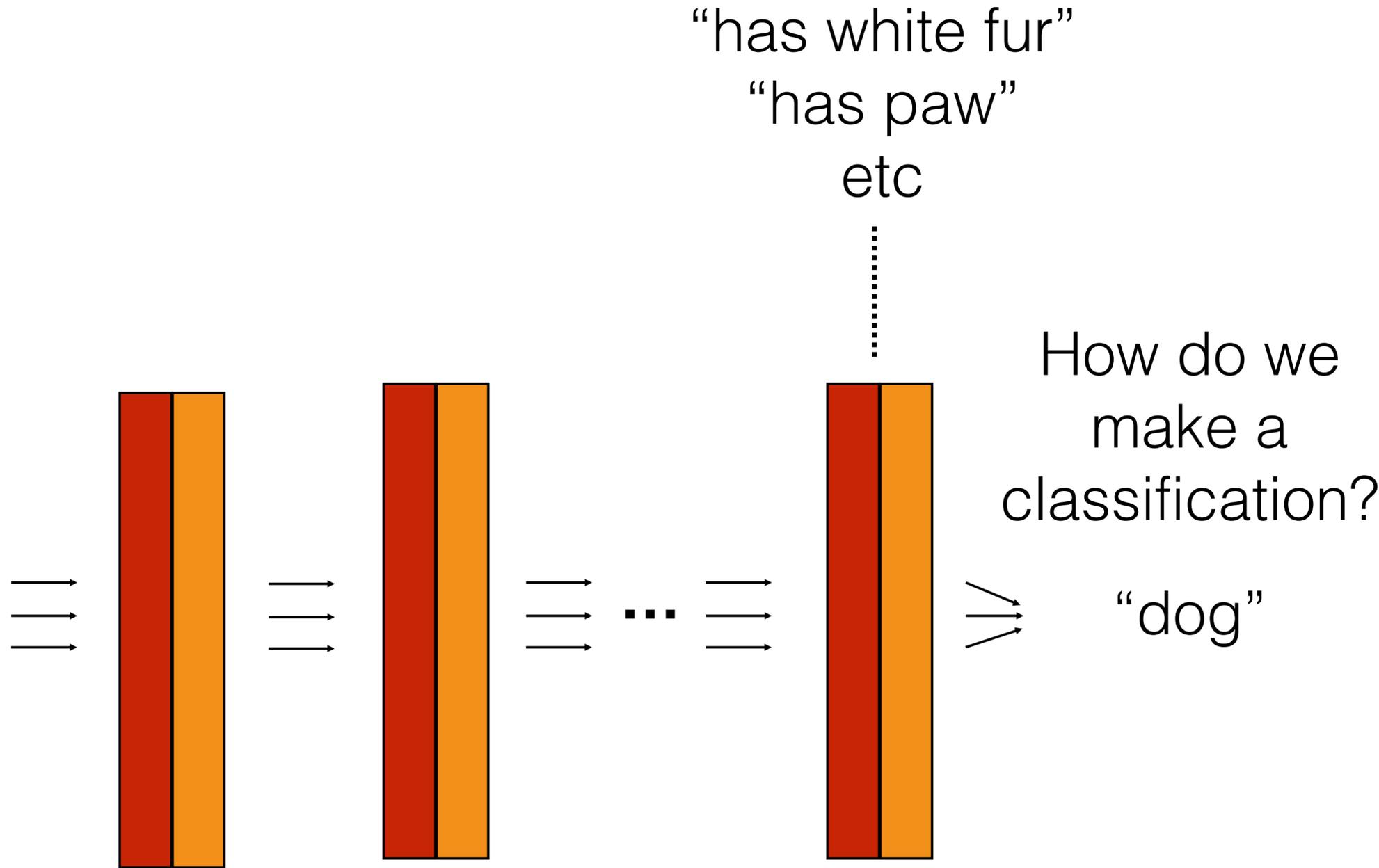


# Deep nets - Intuition

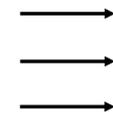
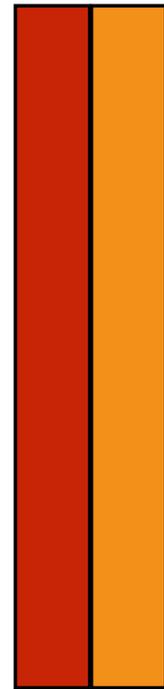
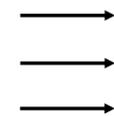
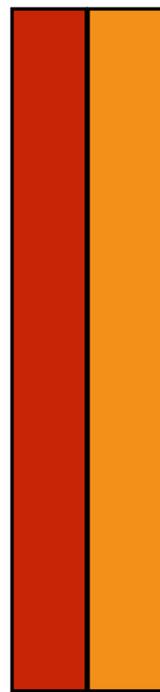
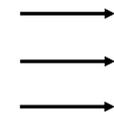
“has rounded edge”



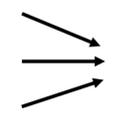
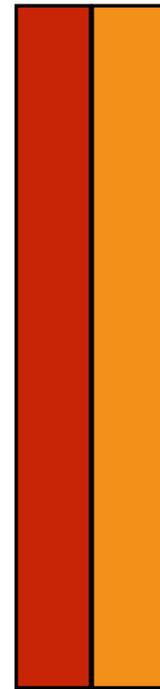
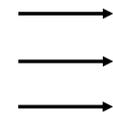
# Deep nets - Intuition



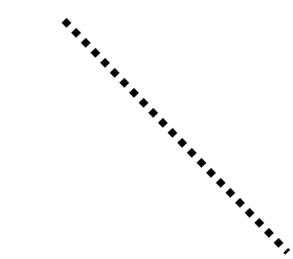
# Deep nets - Intuition



...



“dog”



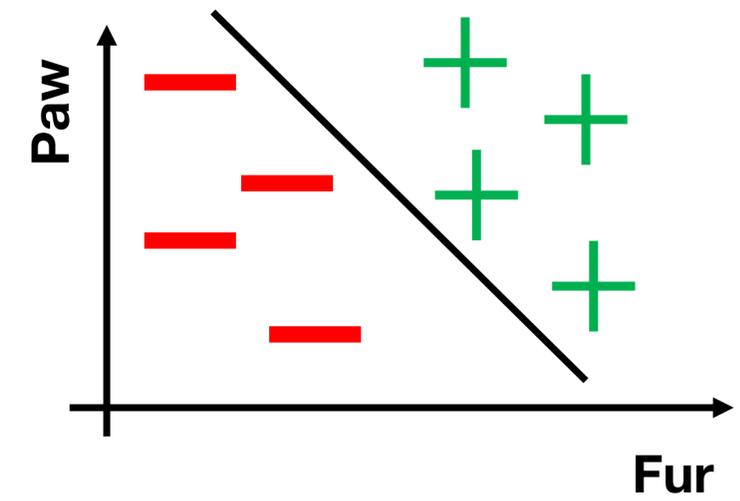
Classify

“has white fur”  
“has paw”  
etc



Recall:

Feature Space



# Computation has a simple form

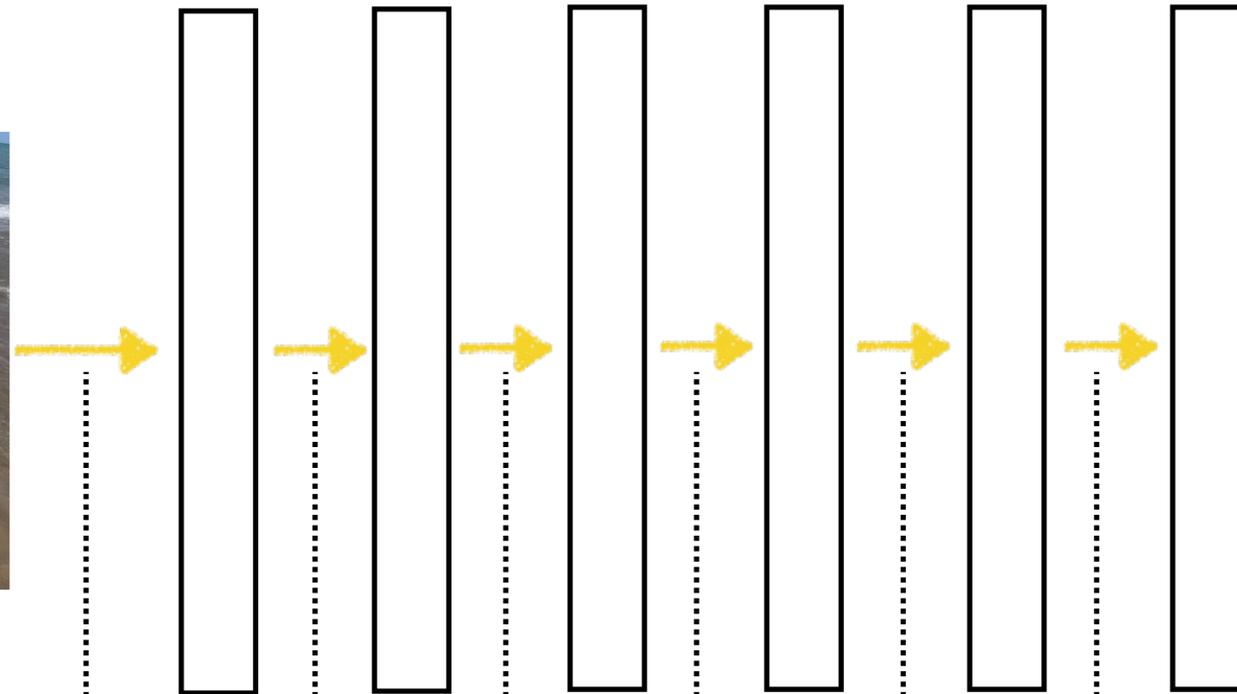
- Composition of linear functions with nonlinearities in between
- E.g. matrix multiplications with ReLU,  $\max(0, \mathbf{x})$  afterwards
- Do a matrix multiplication, set all negative values to 0, repeat

**But where do we get the weights from?**

# How do we learn the parameters?

$y_1$   
“dog”

$x_1$



$\mathcal{L}(f_{\theta}(\mathbf{x}_1), \mathbf{y}_1)$

predicted

ground truth

Learned

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

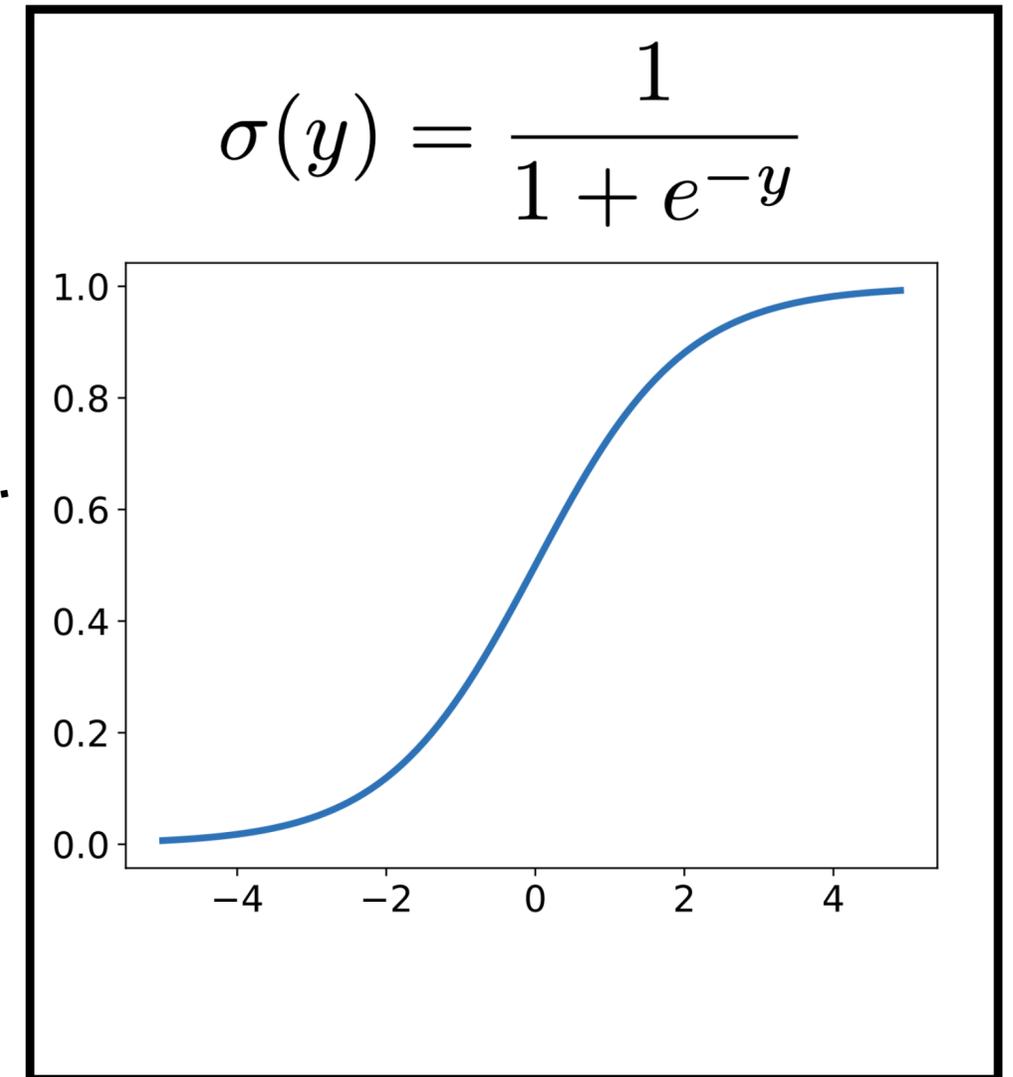
# Learning parameters

Squared loss with single-variable network:

$$L = \frac{1}{2} (y - f(x))^2$$

$$L = \frac{1}{2} (y - \sigma(wx + b))^2$$

Want: derivatives  $\frac{\partial L}{\partial w}$ ,  $\frac{\partial L}{\partial b}$



# Computing derivatives with the chain rule

Given:  $L = \frac{1}{2} (y - \sigma(wx + b))^2$

Writing out the layers explicitly:

$$z = wx + b$$

$$t = \sigma(z)$$

$$L = \frac{1}{2} (y - t)^2$$

$$\rightarrow \frac{\partial L}{\partial w} = \frac{\partial L}{\partial t} \frac{\partial t}{\partial z} \frac{\partial z}{\partial w}$$

# Computing derivatives with the chain rule

$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w} \left[ \frac{1}{2} (y - \sigma(wx + b))^2 \right]$$

$$= (y - \sigma(wx + b)) \frac{\partial}{\partial w} (y - \sigma(wx + b))$$

$$= (y - \sigma(wx + b)) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b)$$

$$= (y - \sigma(wx + b)) \sigma'(wx + b) x$$

$$\frac{\partial L}{\partial b} = \frac{\partial}{\partial b} \left[ \frac{1}{2} (y - \sigma(wx + b))^2 \right]$$

$$= (y - \sigma(wx + b)) \frac{\partial}{\partial b} (y - \sigma(wx + b))$$

$$= (y - \sigma(wx + b)) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b)$$

$$= (y - \sigma(wx + b)) \sigma'(wx + b)$$

Note: For each of these derivatives, you'll have to compute many things multiple times!

# Limitations to this approach

- Inefficient! Lots of redundant computation
- We'll also need to extend this to multivariable functions
- **Next lecture:** backpropagation

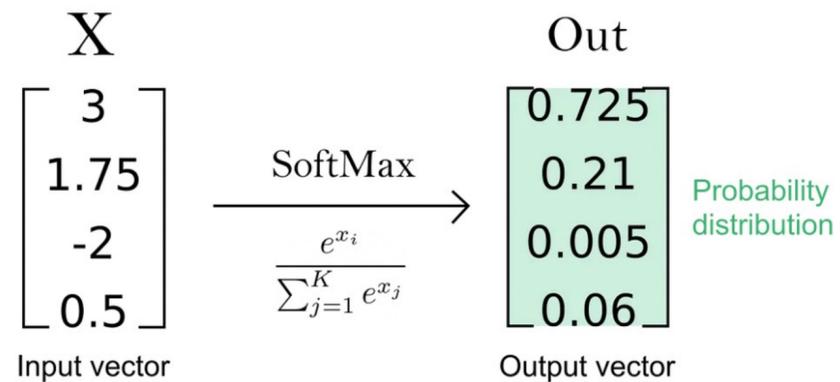
# Representational power

- 1 layer? Linear decision surface.
- 2+ layers? In theory, can represent any function! (if it was infinitely wide with infinite data)
  - Simple proof by M. Nielsen  
<http://neuralnetworksanddeeplearning.com/chap4.html>
- But issue is efficiency: very wide two layers vs narrow deep model? In practice, more layers helps.

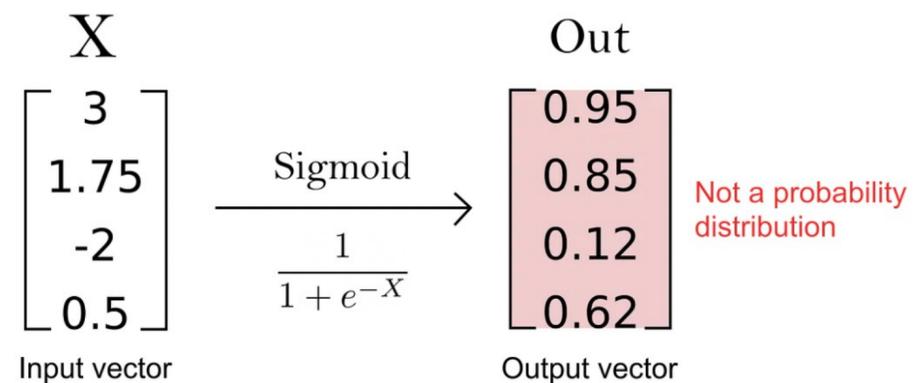
# Backup Slides

# Sigmoid vs. Softmax

Softmax outputs a probability distribution over all predicted classes:

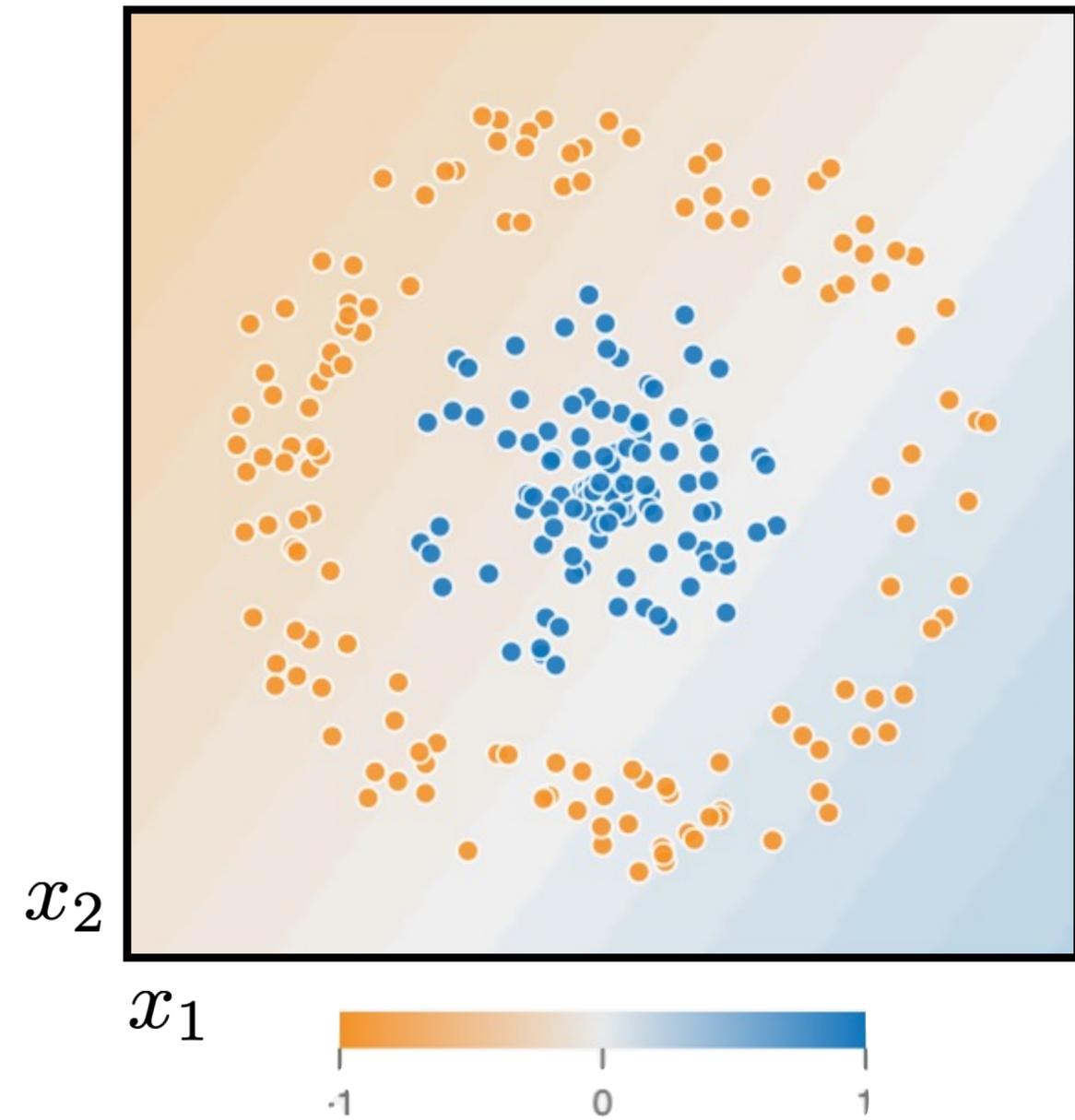
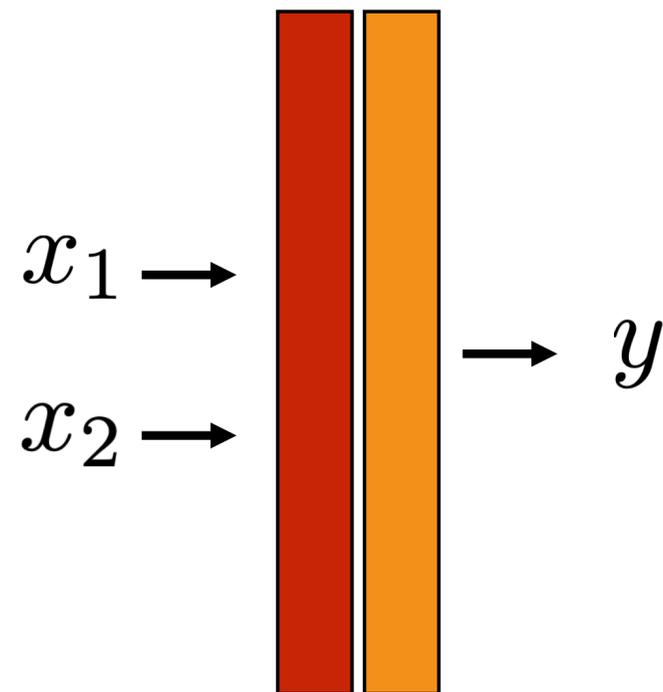


Sigmoid – not a probability distribution:



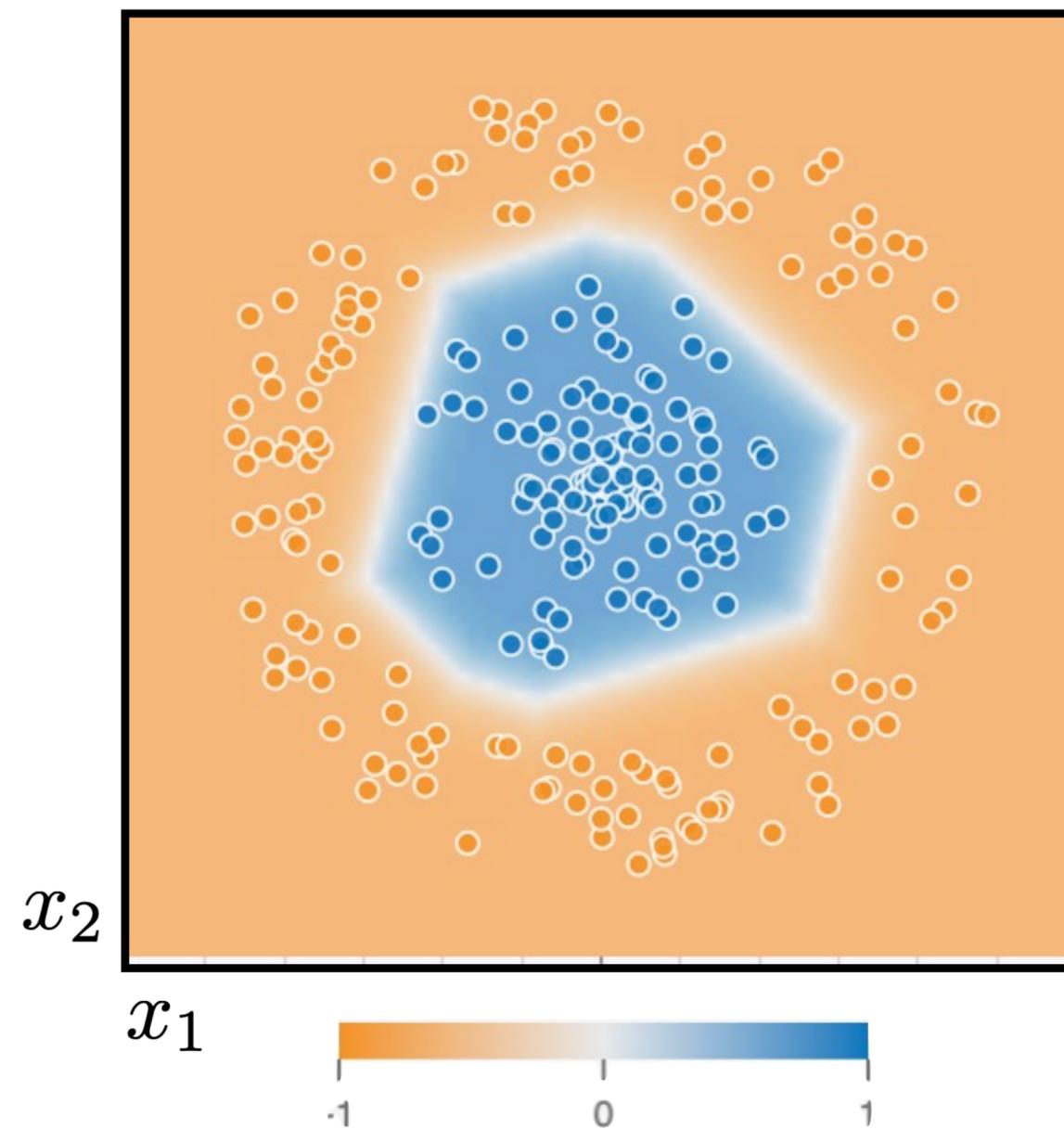
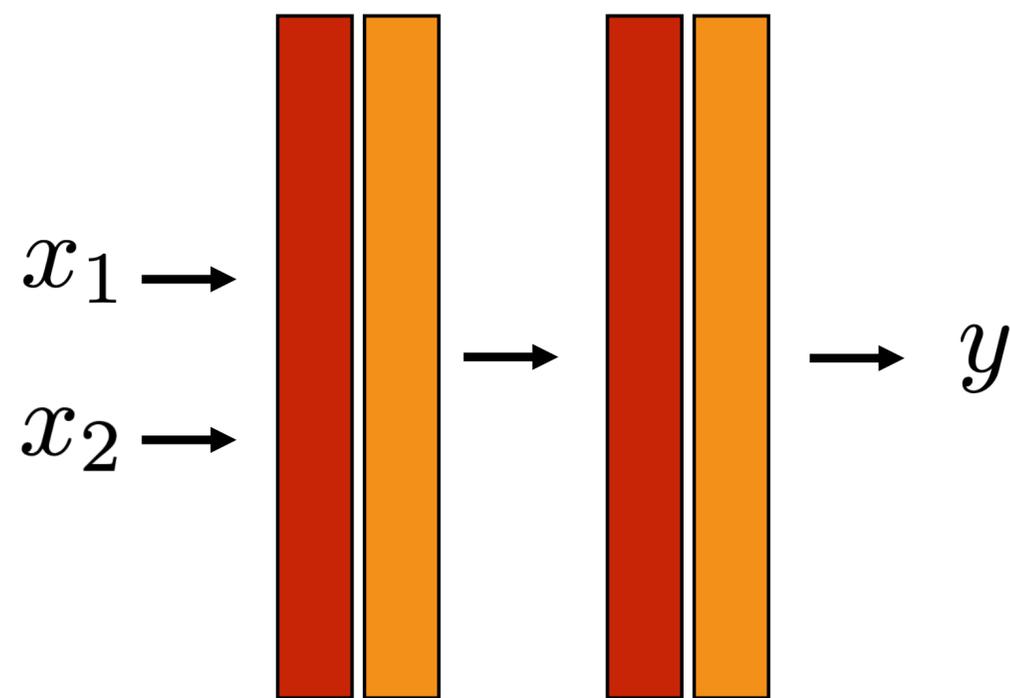
# Example: perceptron

$$\mathbf{y} = \sigma(\mathbf{W}^{(1)}\mathbf{x})$$

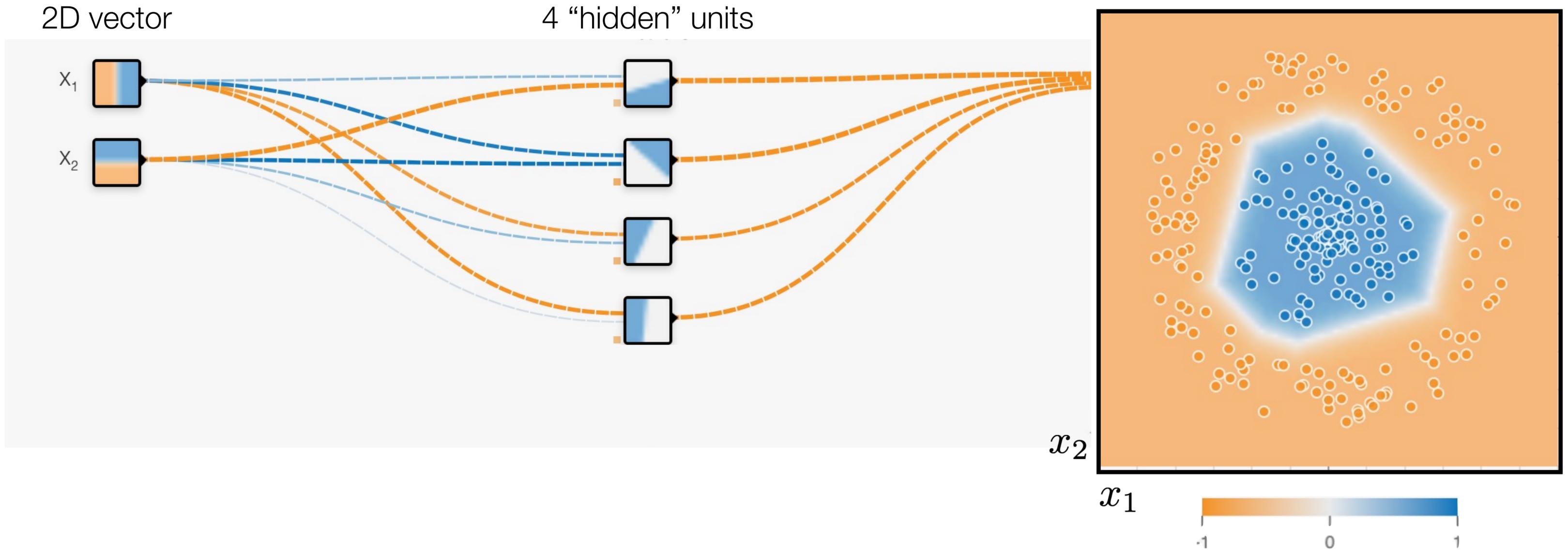


# Example: multilayer perceptron (MLP)

$$\mathbf{y} = \sigma(\mathbf{W}^{(2)} \max(0, \mathbf{W}^{(1)} \mathbf{x}))$$



# Example: multilayer perceptron (MLP)



Example source: <http://playground.tensorflow.org>

# Example: multilayer perceptron (MLP)

