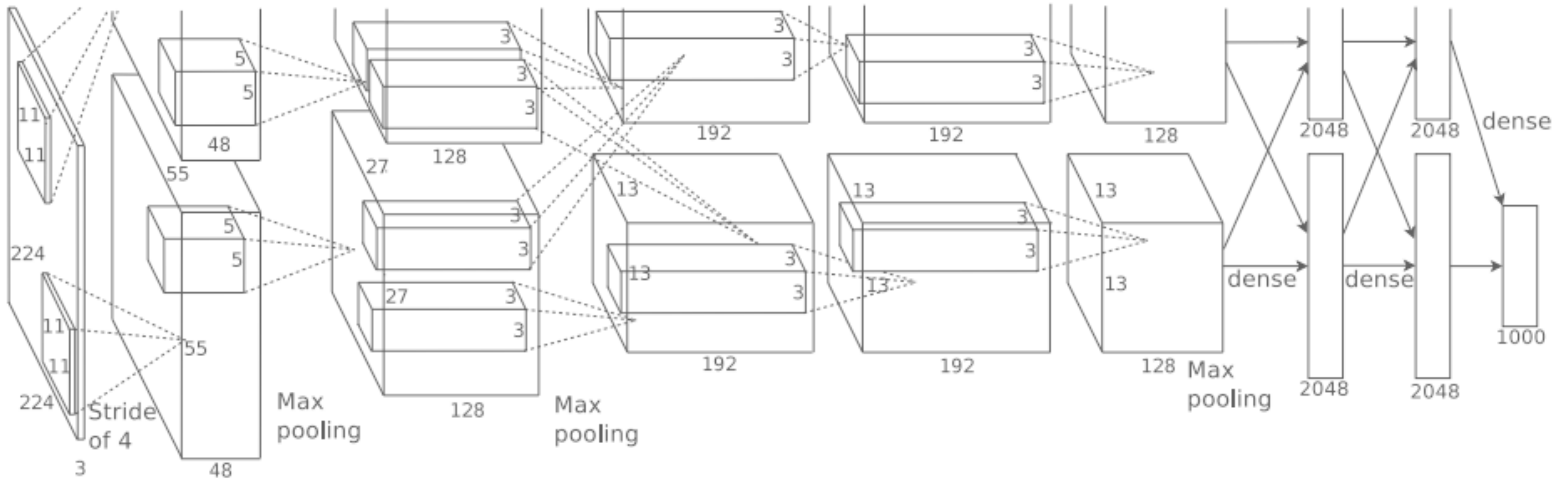


# Lecture 7

## Neural Networks



# 7. Introduction to Deep Learning

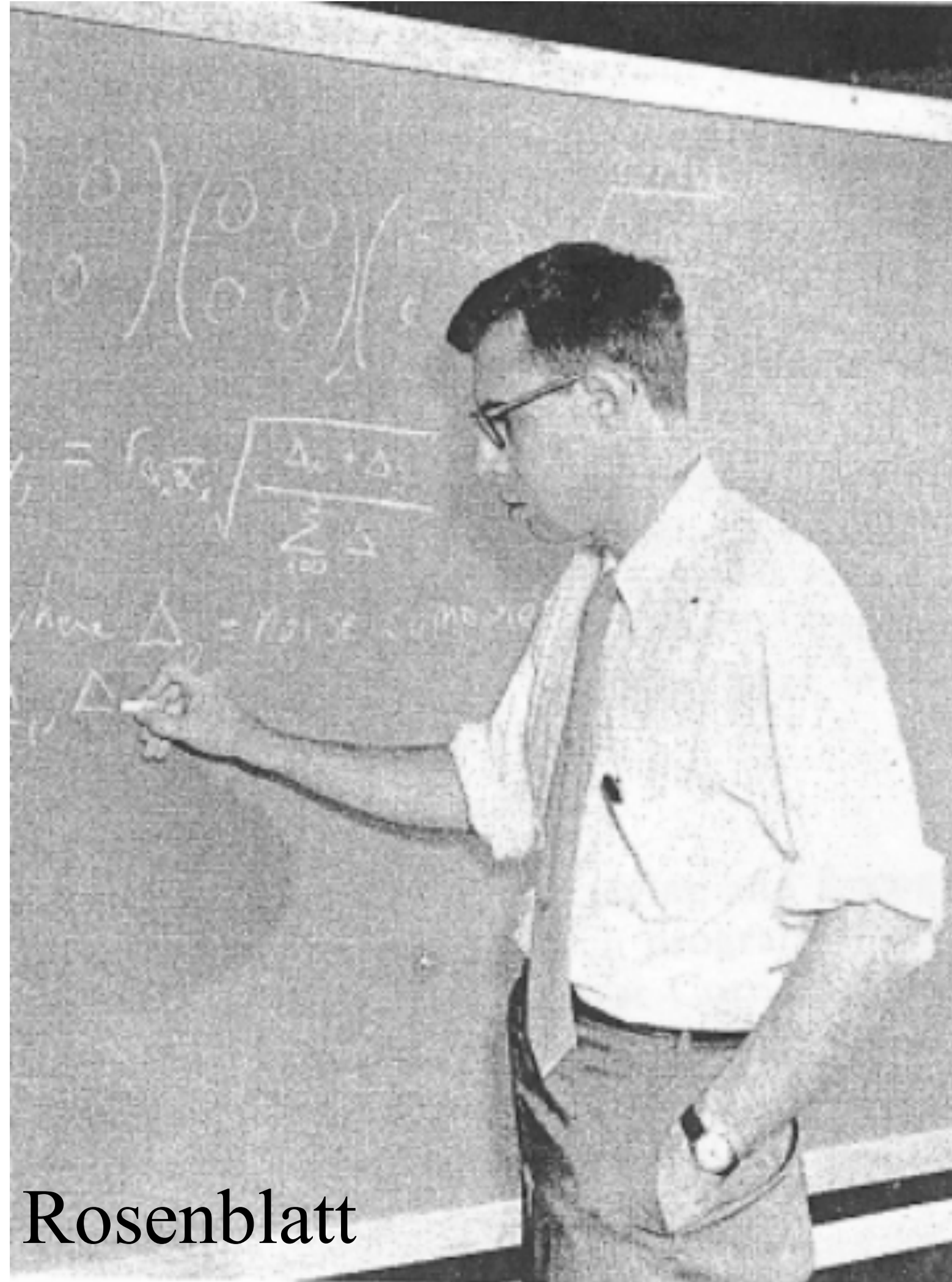
- Brief history
- Basic formulation (*hierarchical processing*)
- Optimization via gradient descent
- Layer types (*Linear, Pointwise non-linearity*)
- Everything is a tensor
- Deep nets as data transformers

# A brief history of Neural Networks



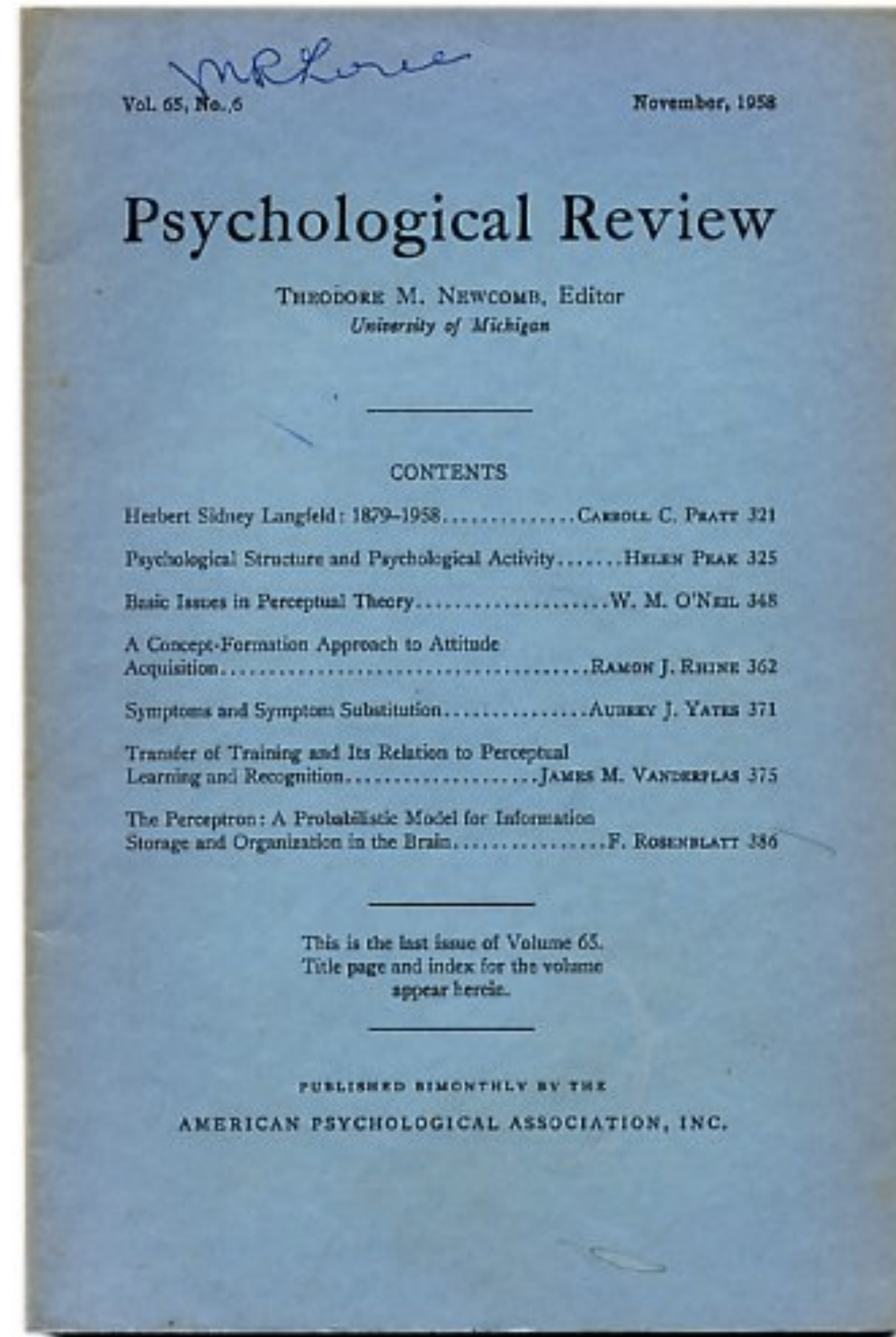


# Perceptrons, 1958



Rosenblatt

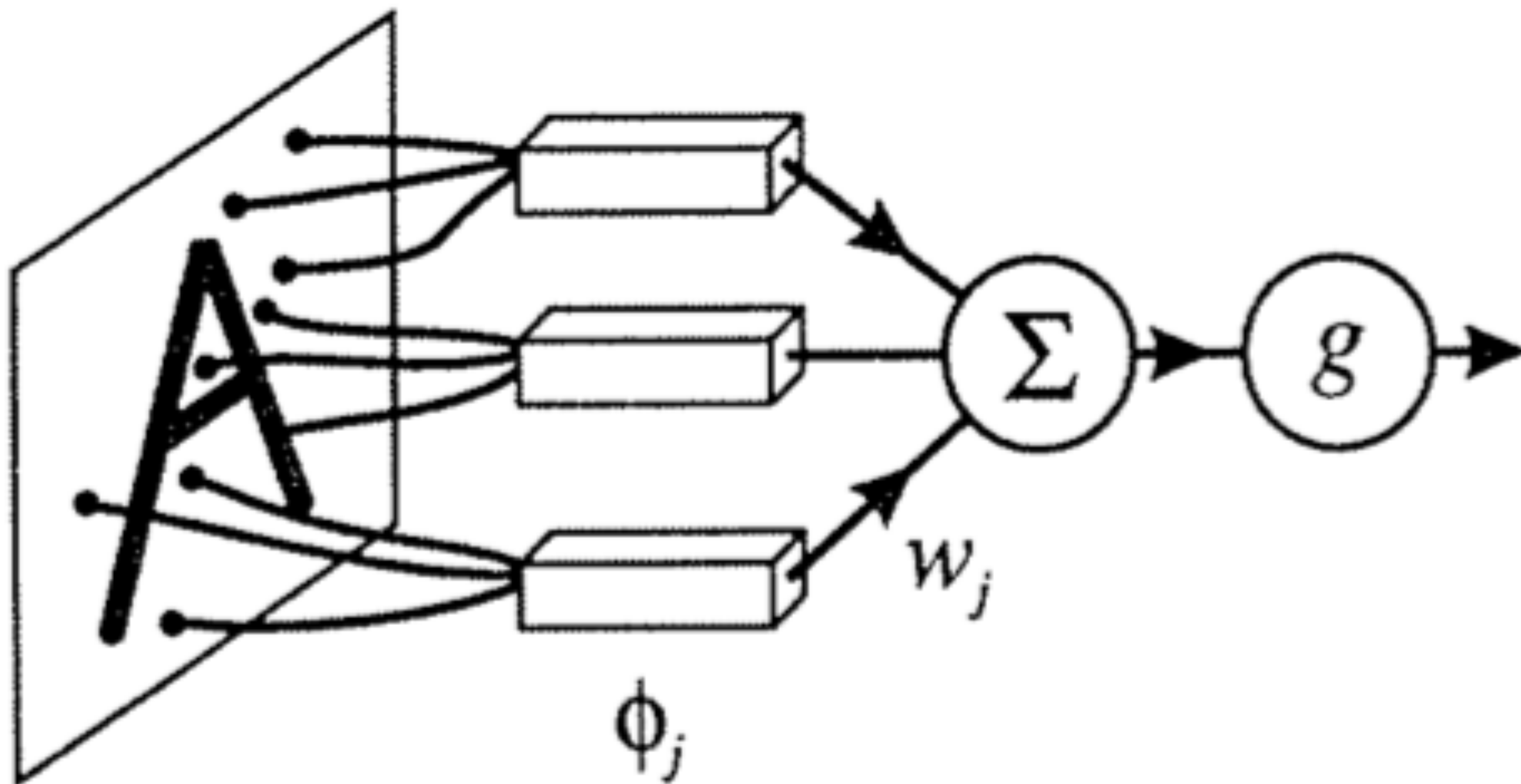
[http://www.ecse.rpi.edu/homepages/nagy/PDF\\_chrono/2011\\_Nagy\\_Pace\\_FR.pdf](http://www.ecse.rpi.edu/homepages/nagy/PDF_chrono/2011_Nagy_Pace_FR.pdf). Photo by George Nagy

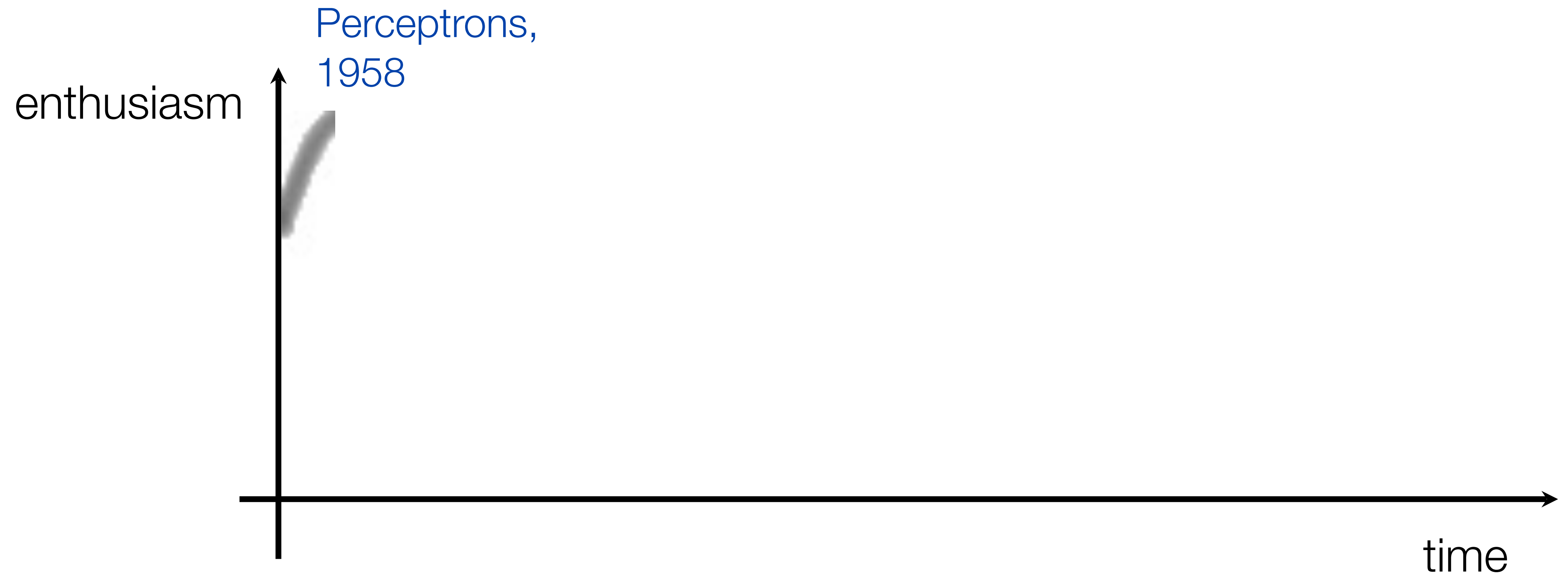


<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.335.3398&rep=rep1&type=pdf>

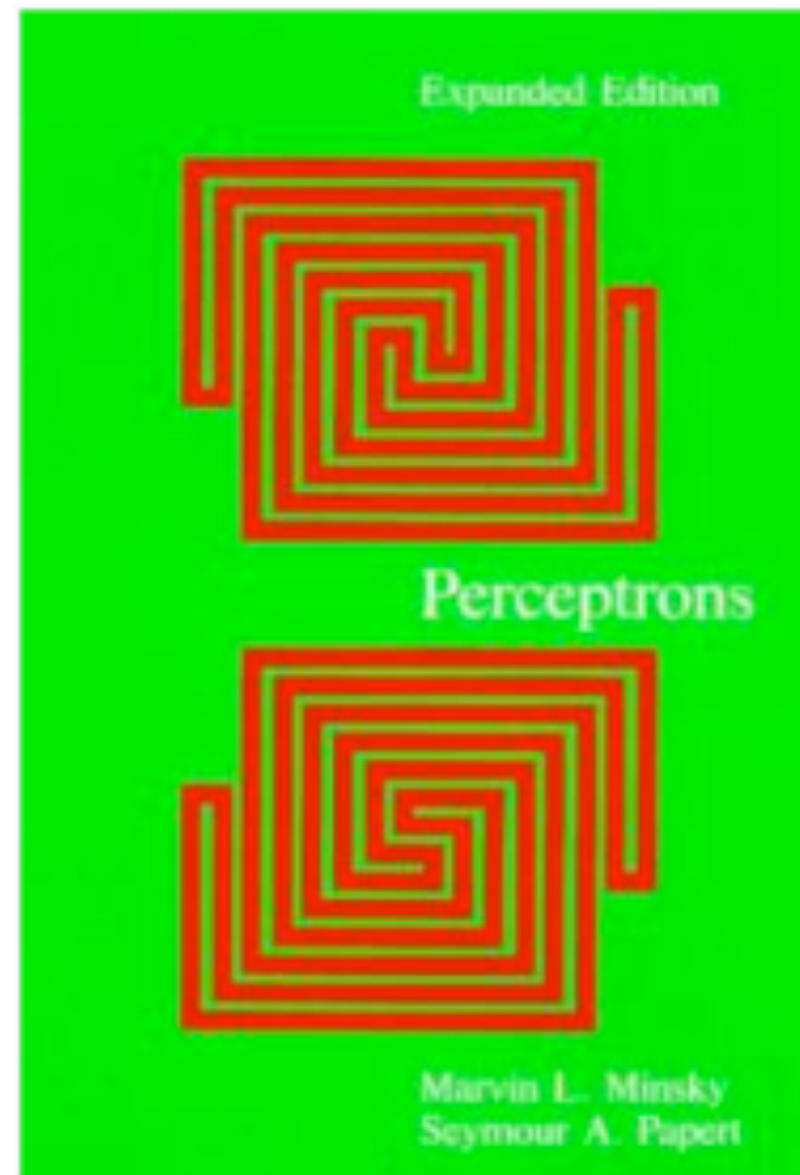


# Perceptrons, 1958





# 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](#)

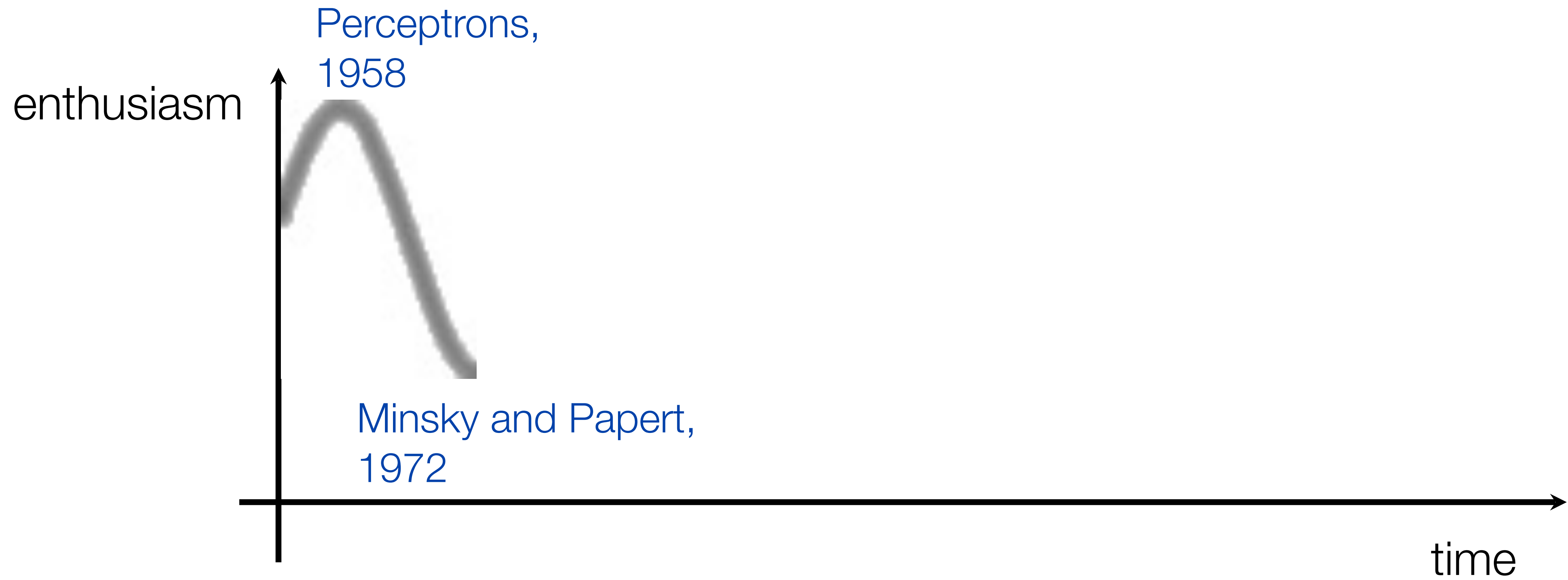
### Overview

*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.

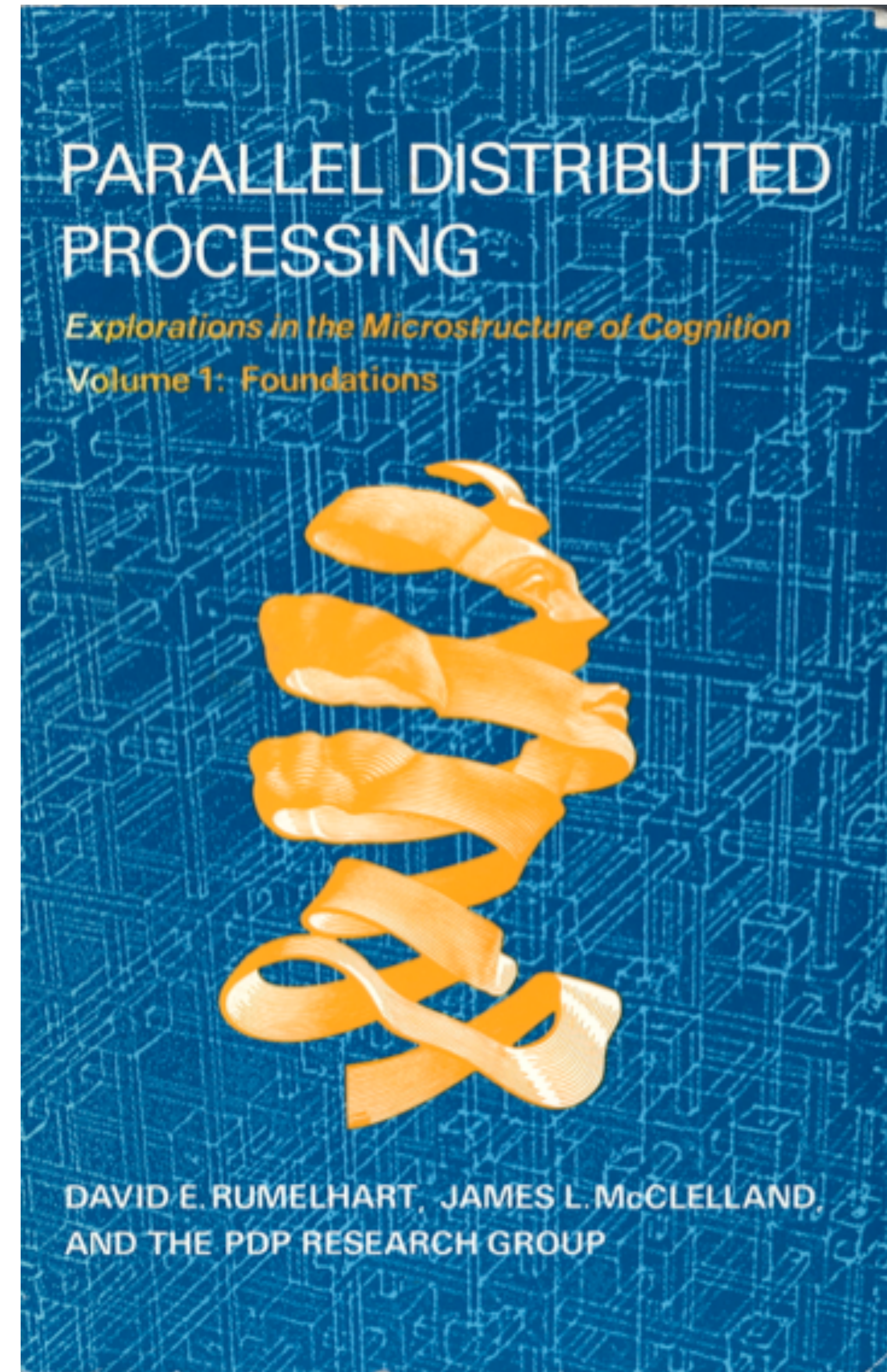
Witnessing the swing of the intellectual pendulum, Minsky and Papert have added a new chapter in which they discuss the current state of parallel computers, review developments since the appearance of the 1972 edition, and identify new research directions related to connectionism. They note a central 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."





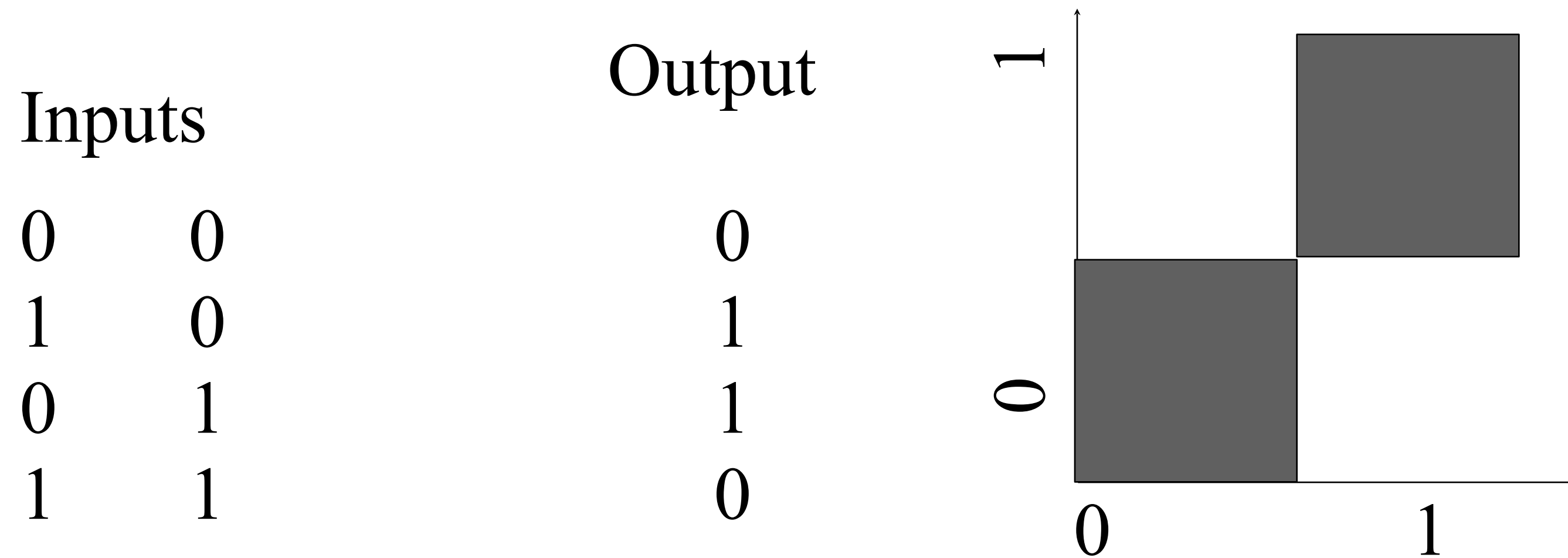


# Parallel Distributed Processing (PDP), 1986



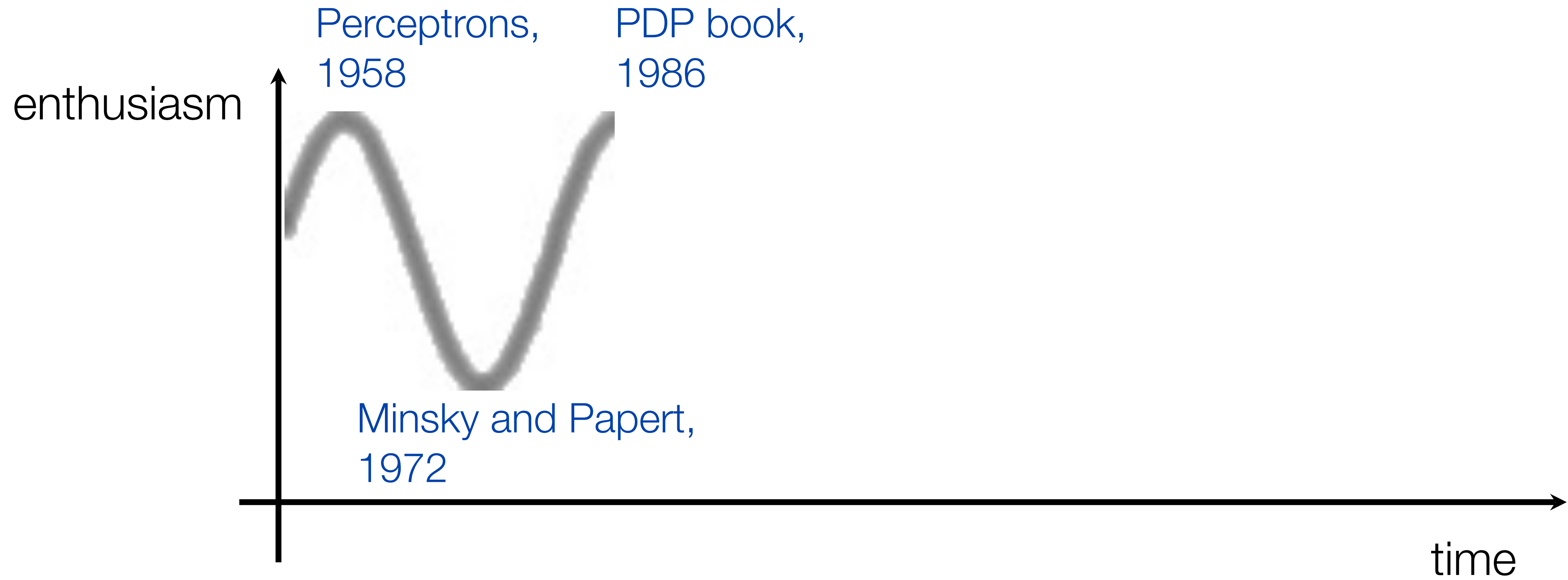


# XOR problem



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

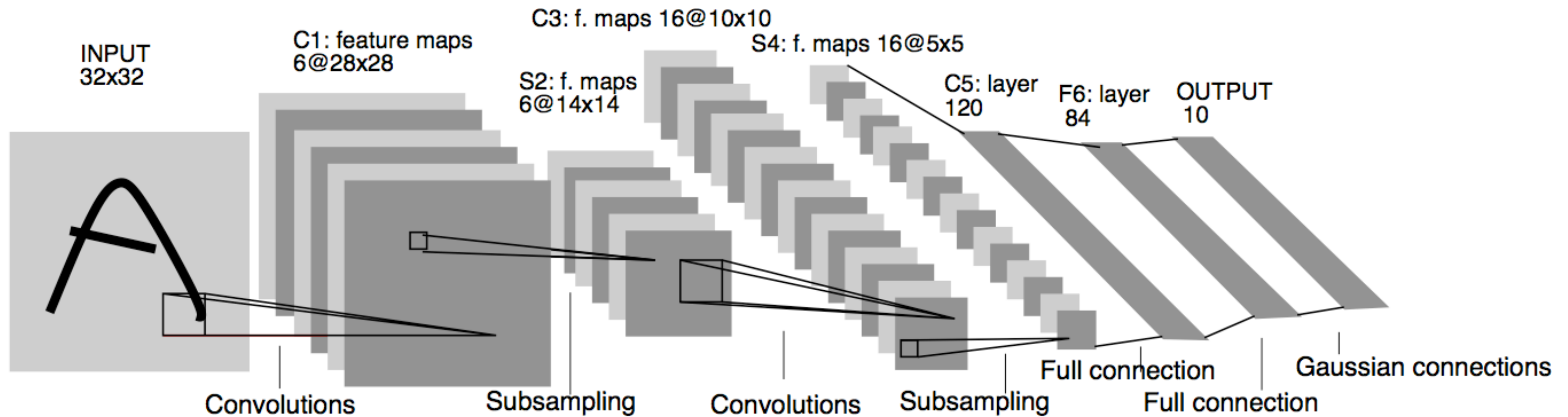


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

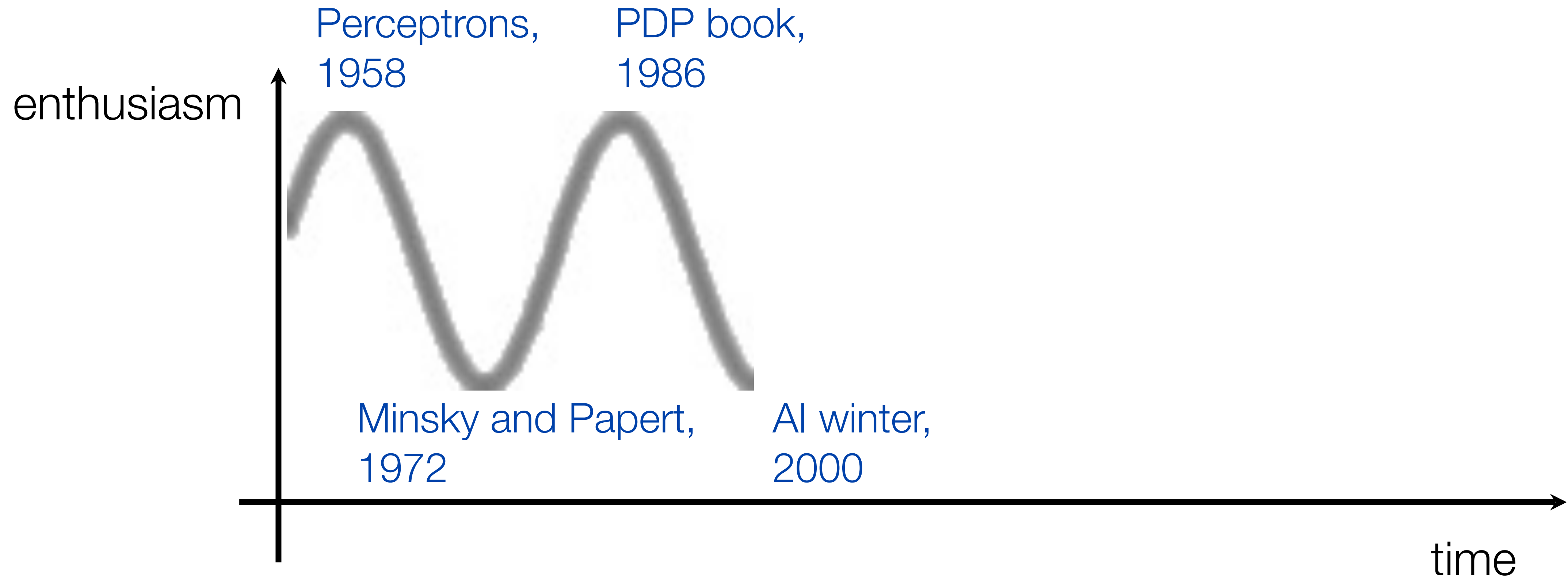
Demos:

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

# Neural Information Processing Systems 2000

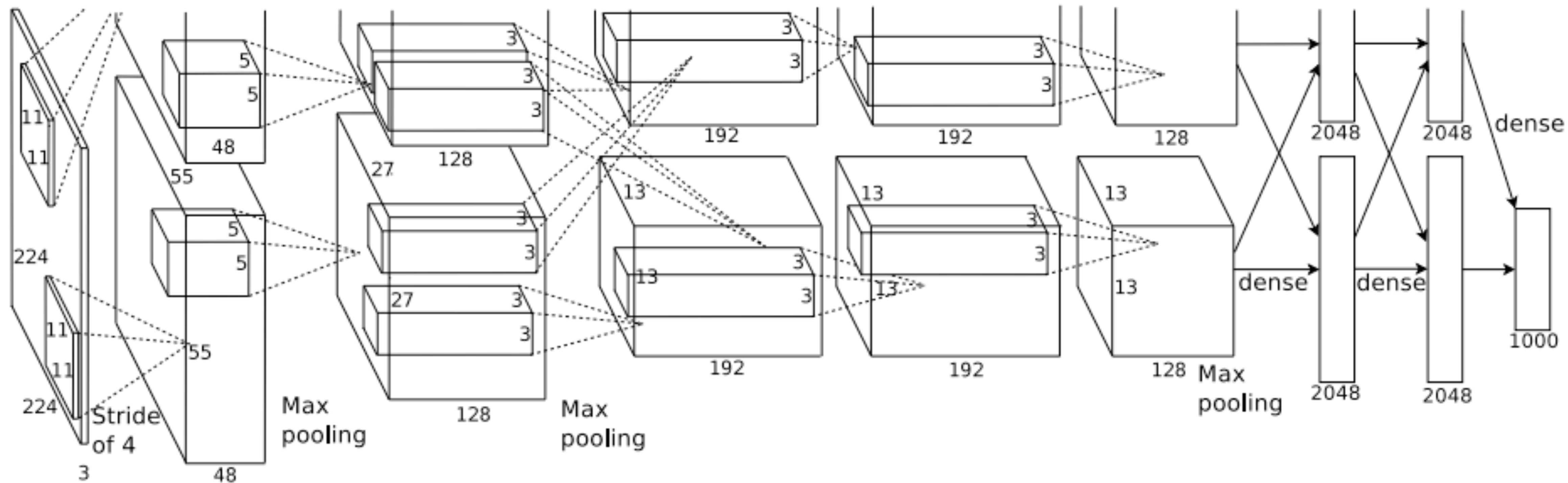
- Neural Information Processing Systems, is the premier conference on machine learning. Evolved from an interdisciplinary conference to a machine learning conference.
- For the 2000 conference:
  - title words predictive of paper acceptance: “Belief Propagation” and “Gaussian”.
  - title words predictive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





# Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

## “Alexnet”

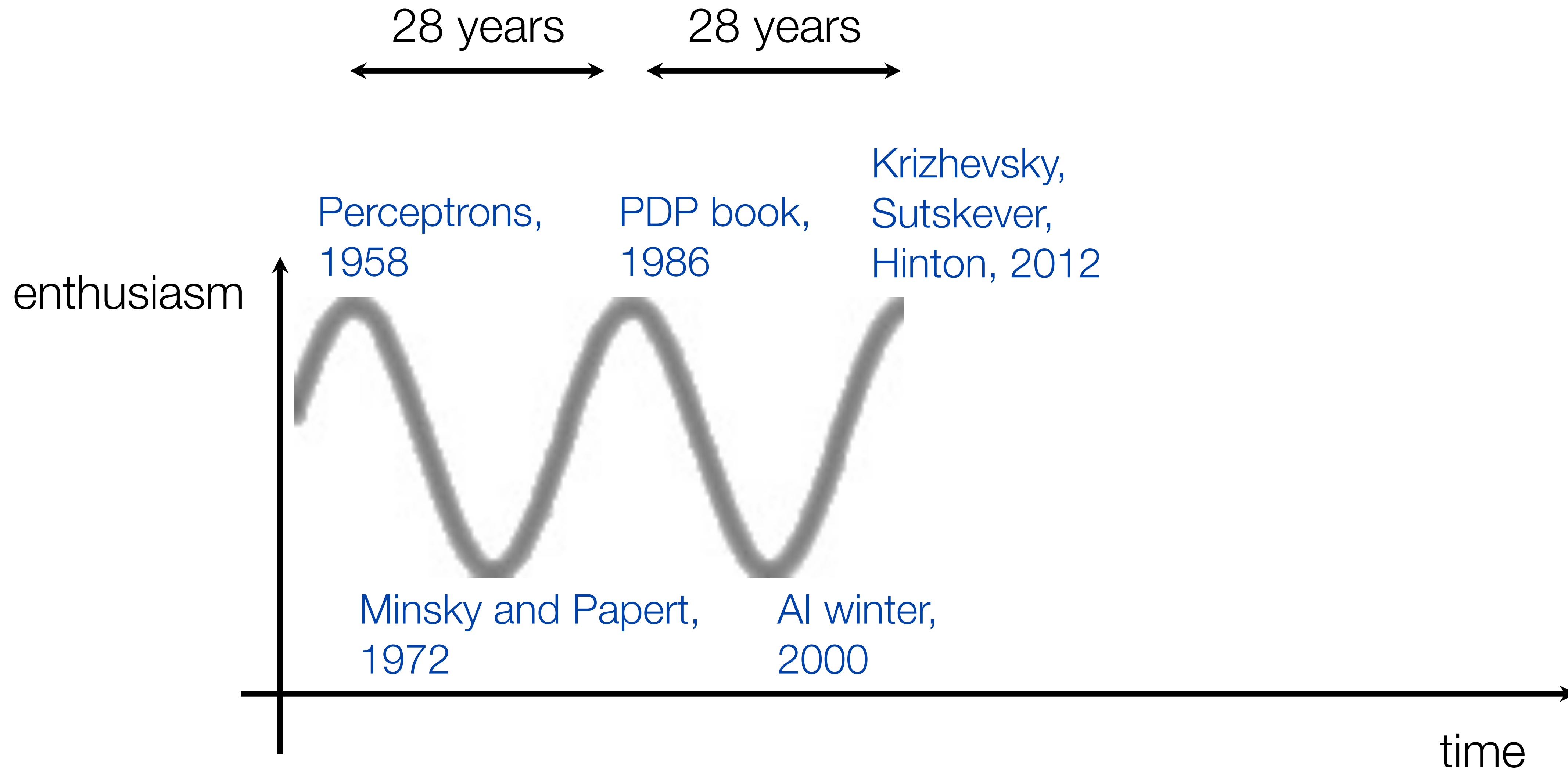




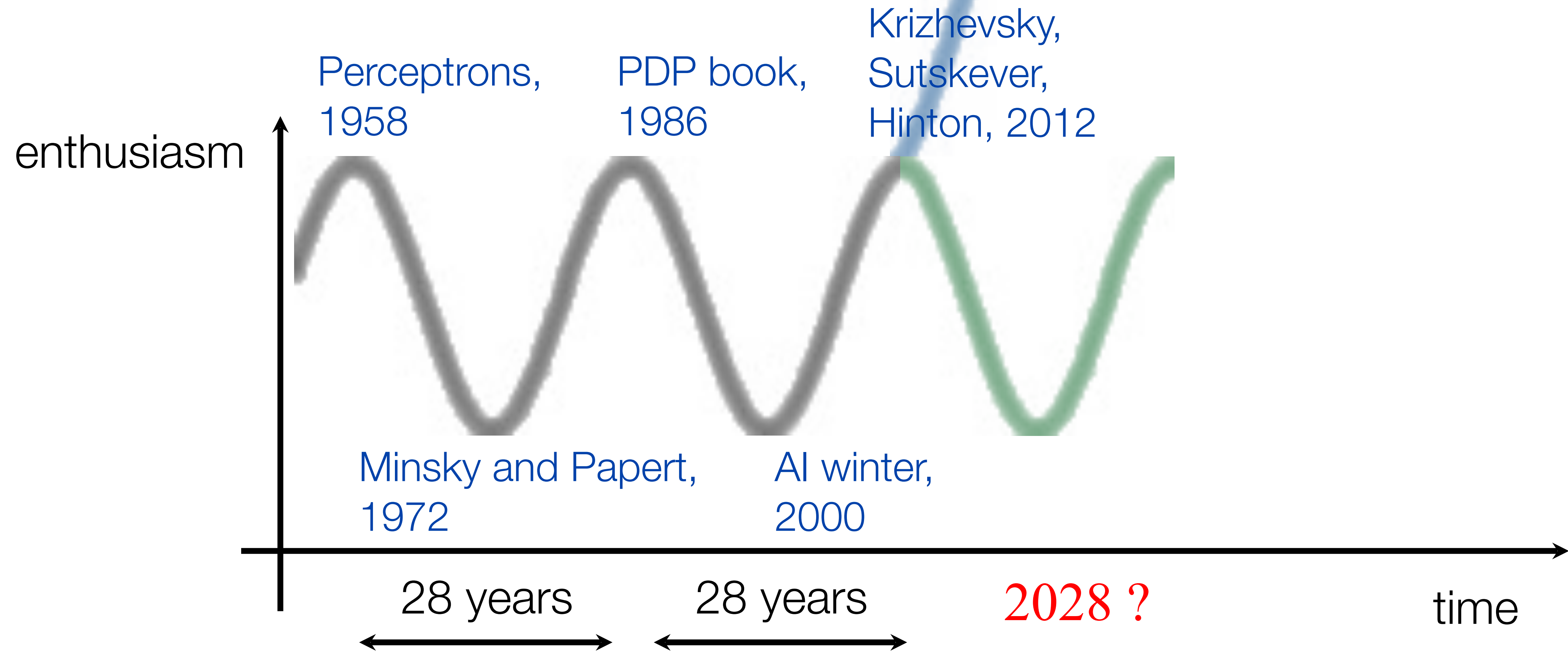
# Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



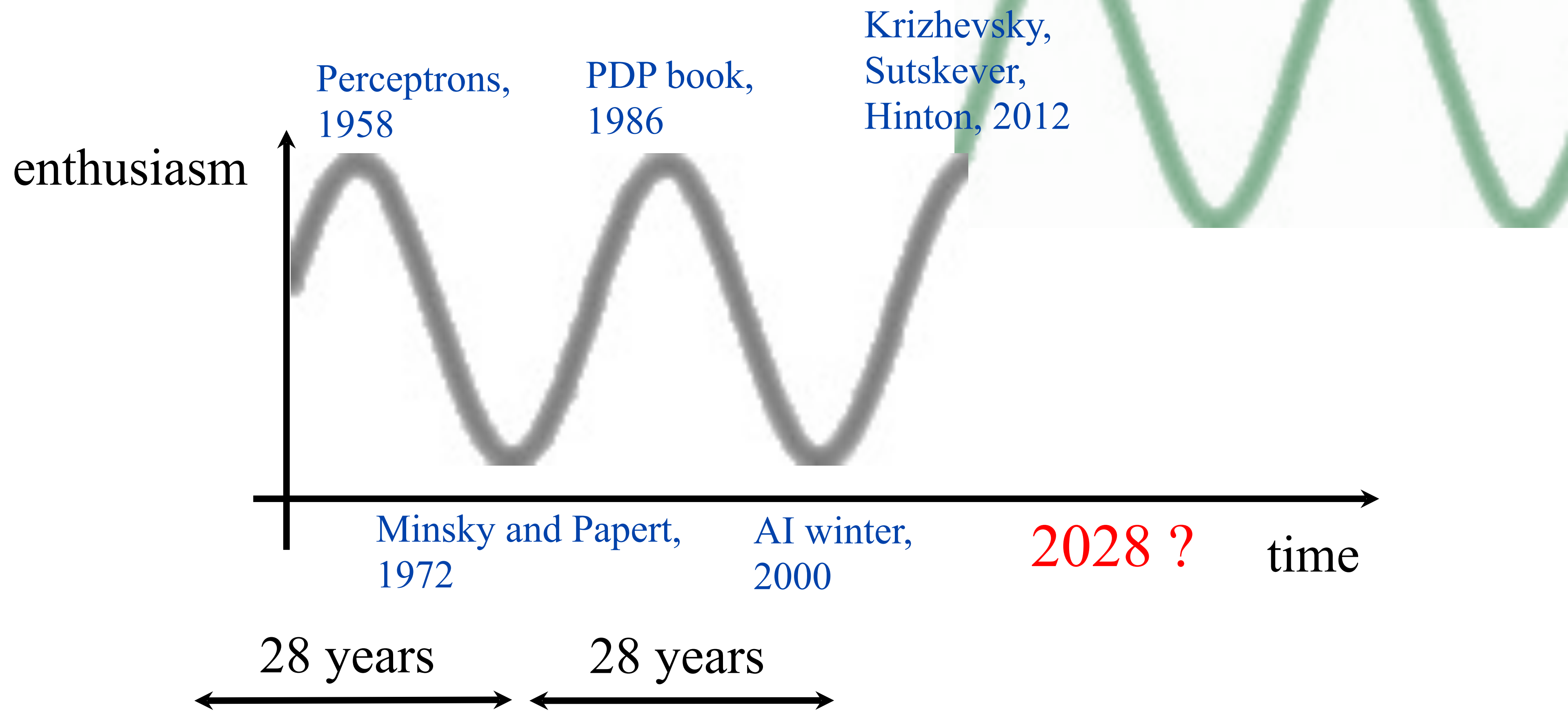




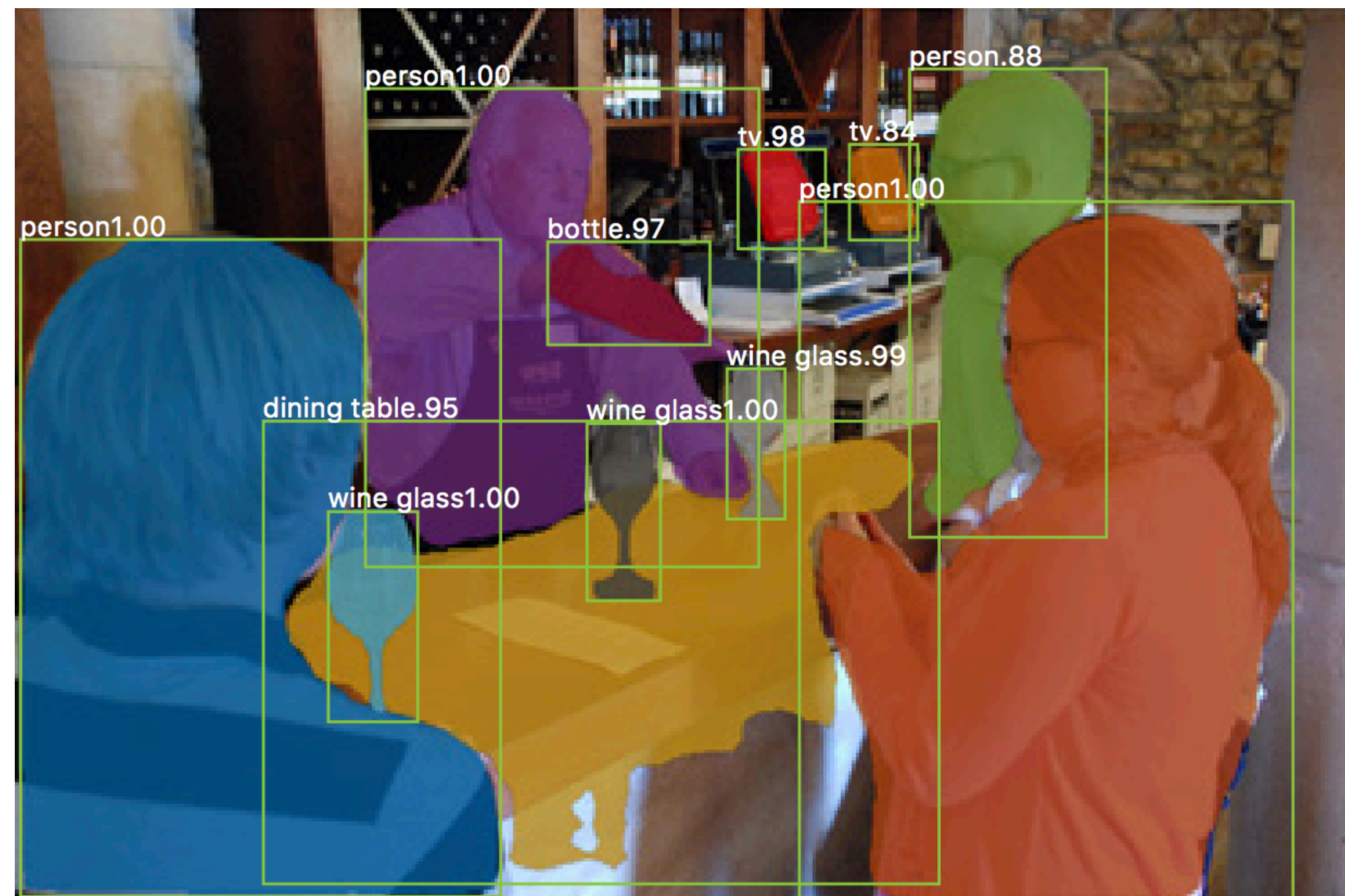
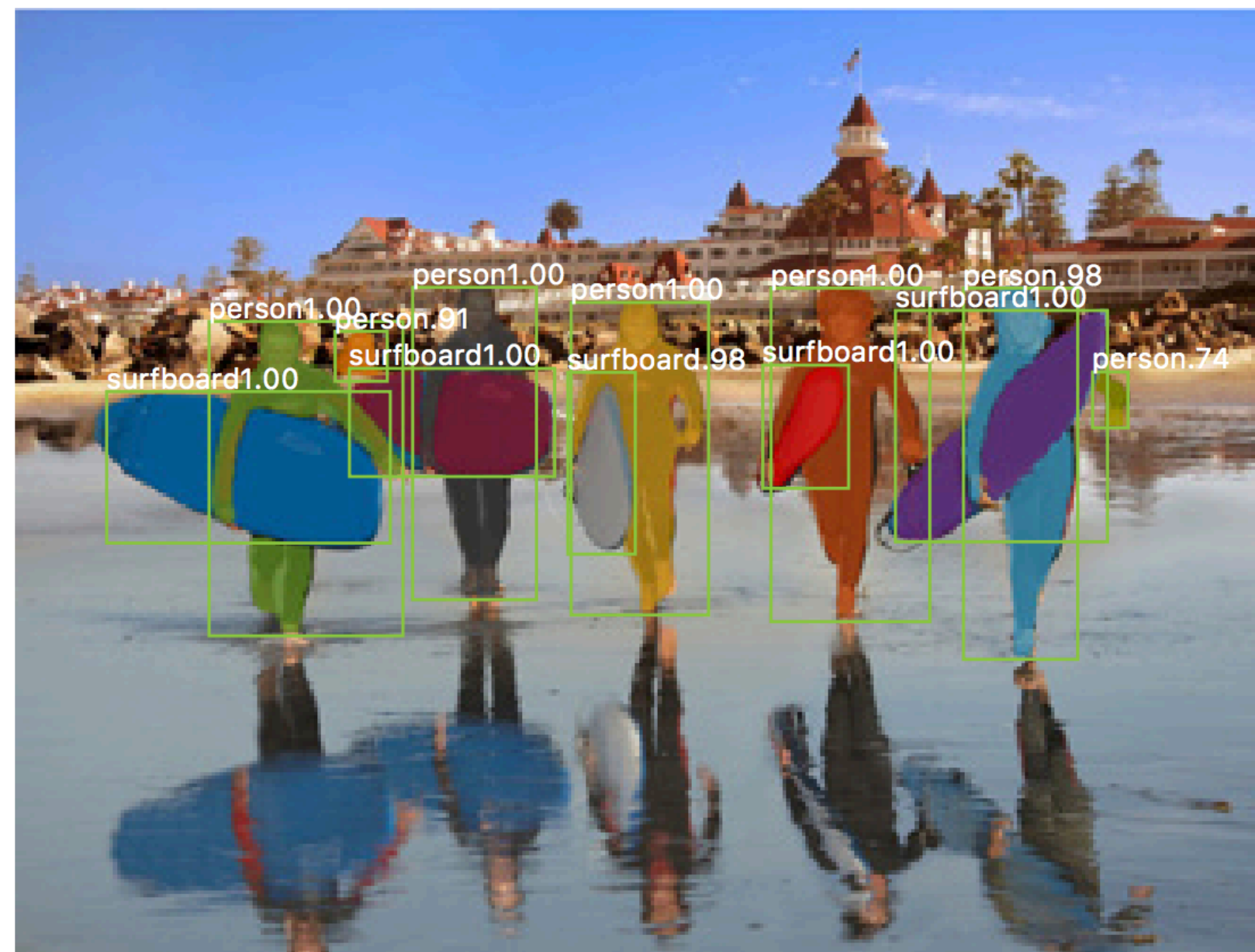
# What comes next?



# What comes next?







["Mask RCNN", He et al. 2017]

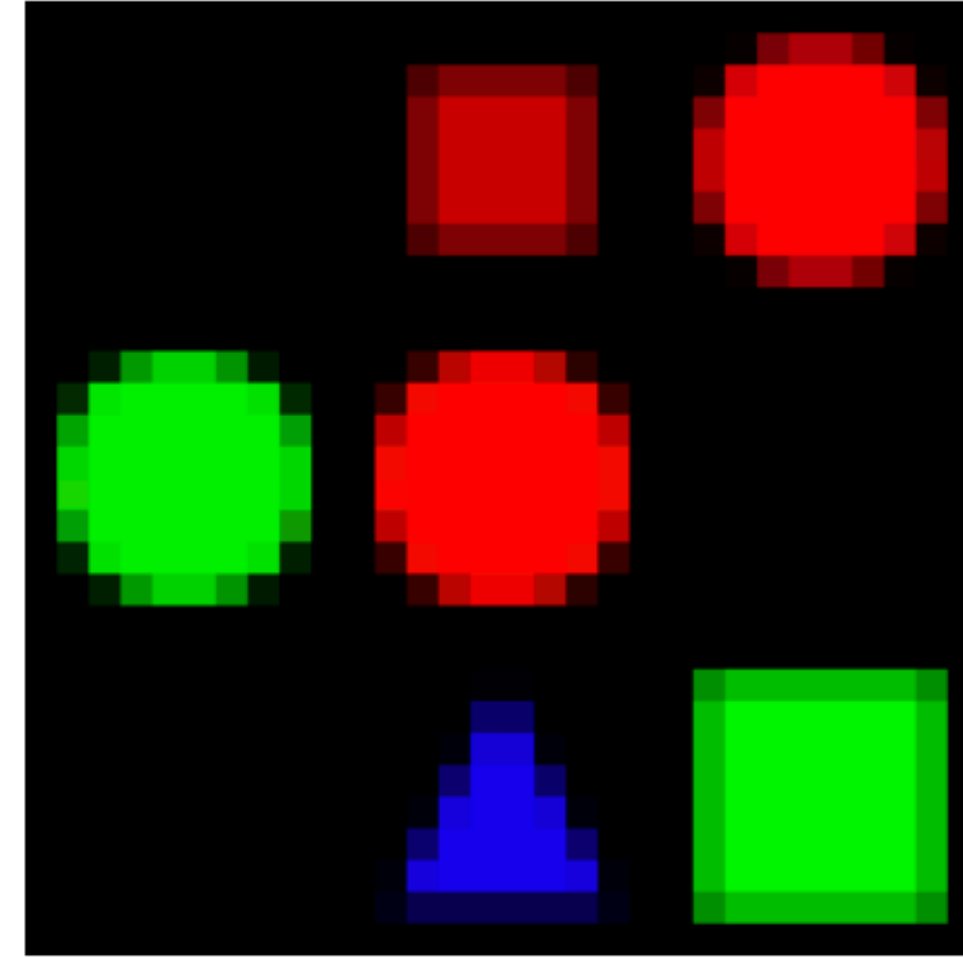




*what color is the vase?*



*is the bus full of passengers?*



*is there a red shape above a circle?*

```
classify[color](
  attend[vase])
```

```
measure[is](
  combine[and](
    attend[bus],
    attend[full]))
```

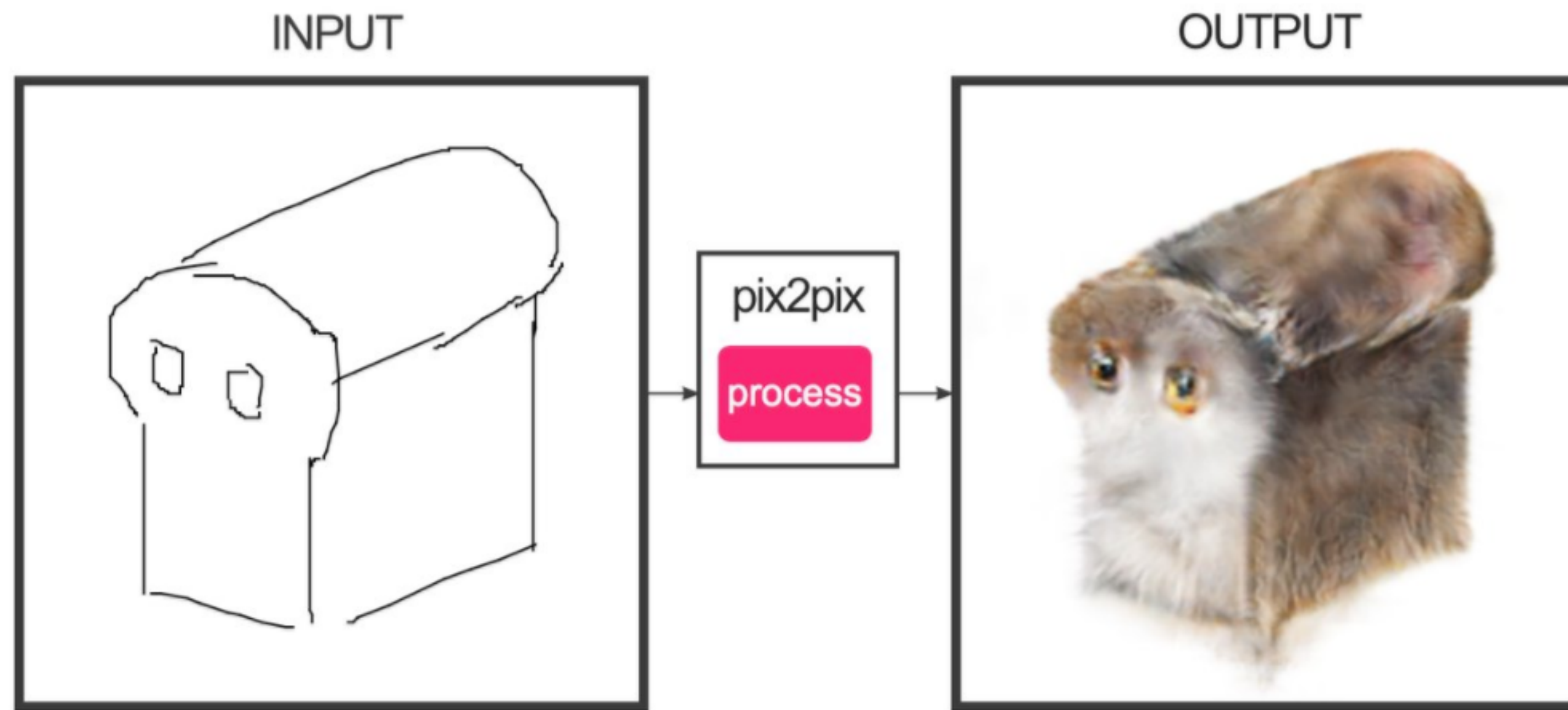
```
measure[is](
  combine[and](
    attend[red],
    re-attend[above](
      attend[circle])))
```

green (green)

yes (yes)

no (no)





Ivy Tasi @ivymyt



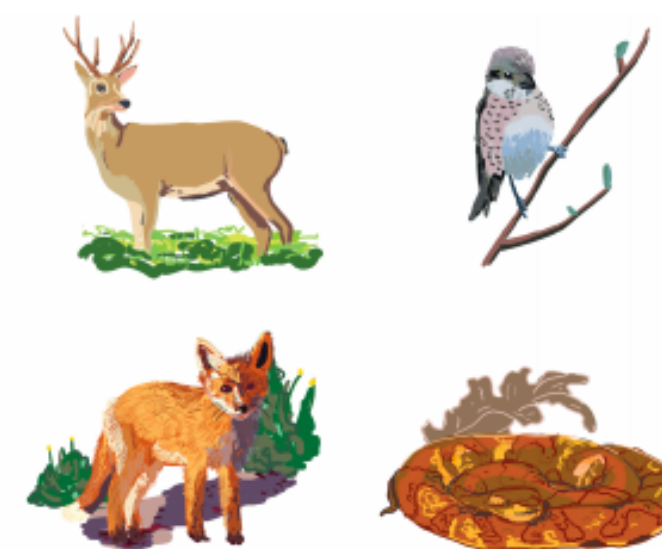
Vitaly Vidmirov @vvid

[“pix2pix”, Isola et al. 2017]





Classification units



PIT/AIT



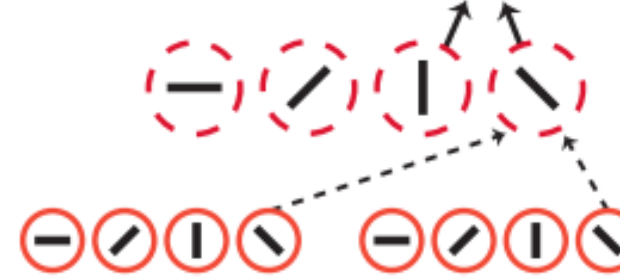
V4/PIT



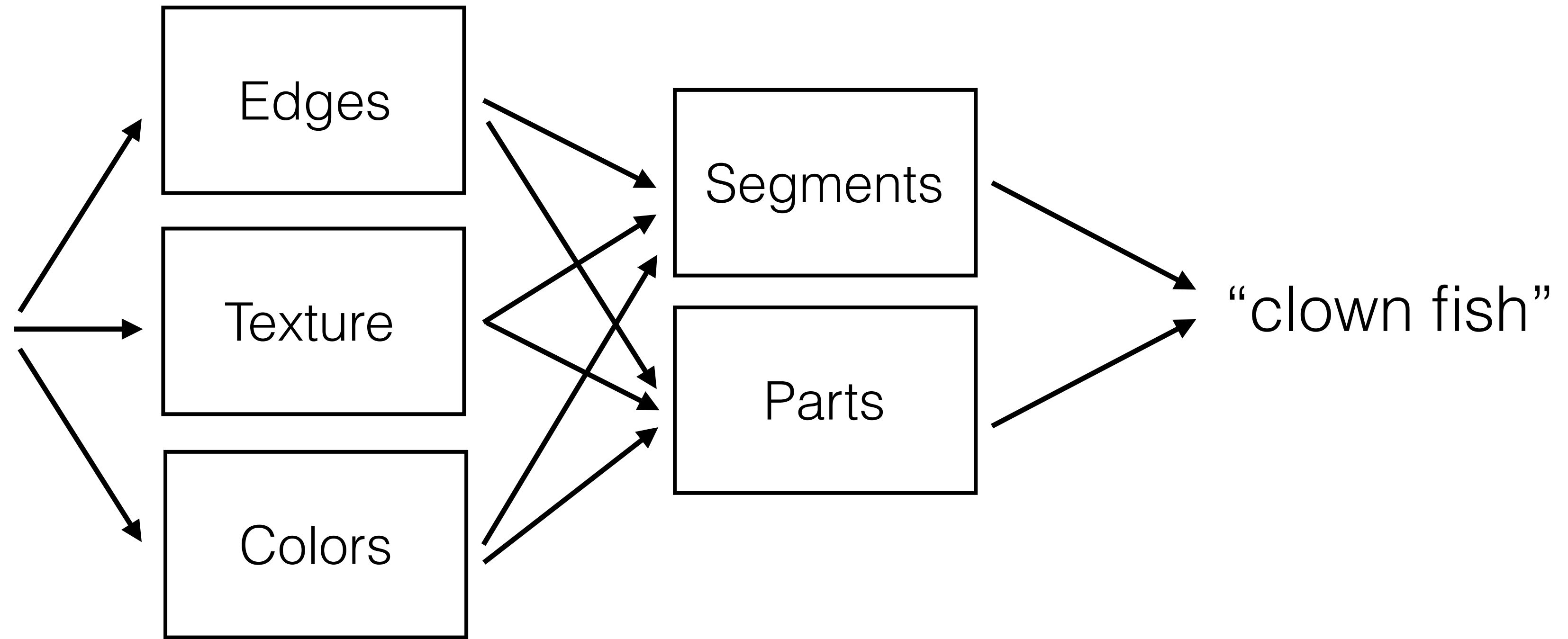
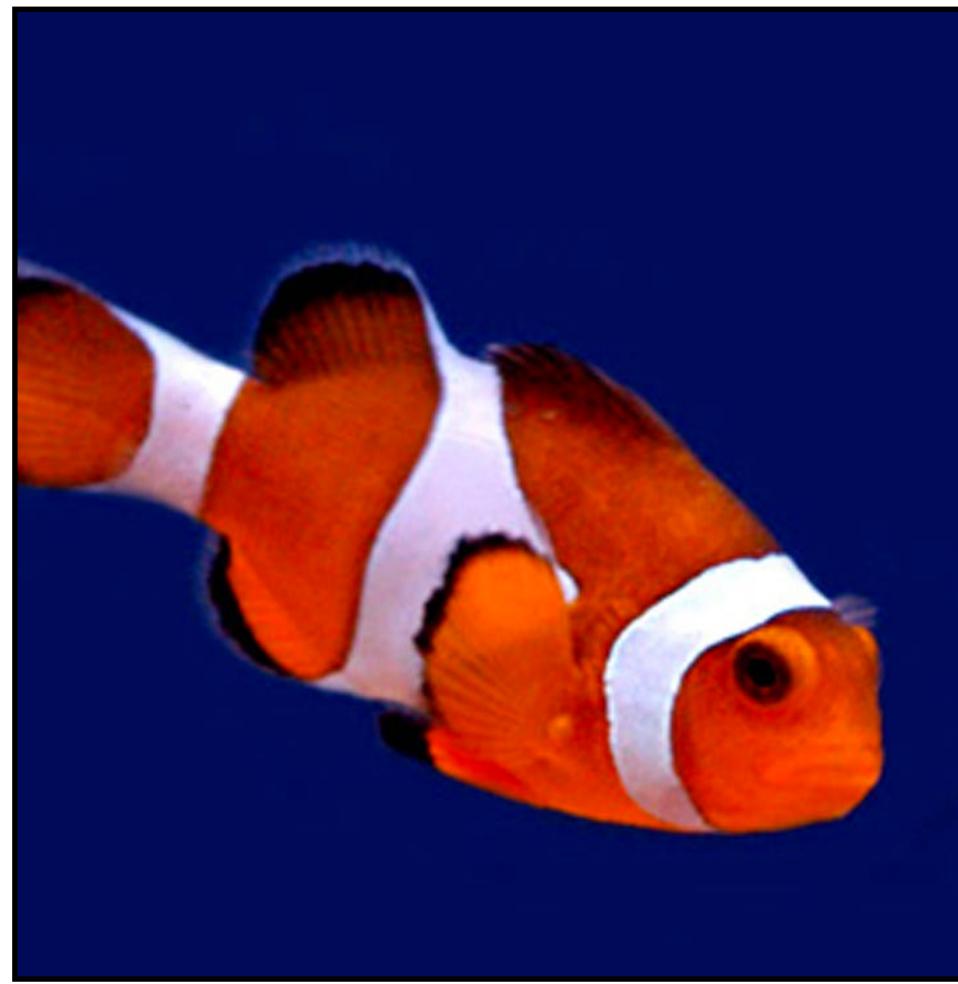
V2/V4



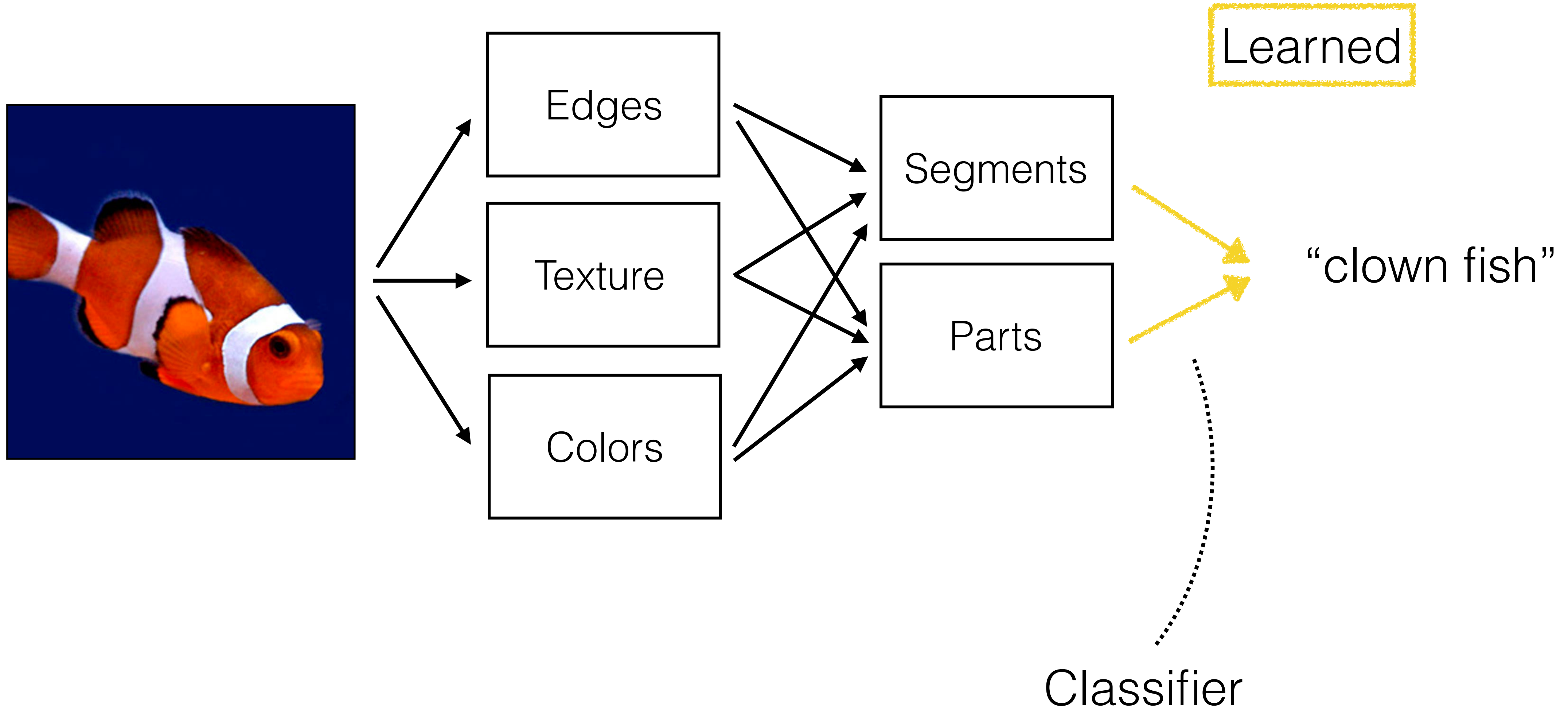
V1/V2



# Image classification

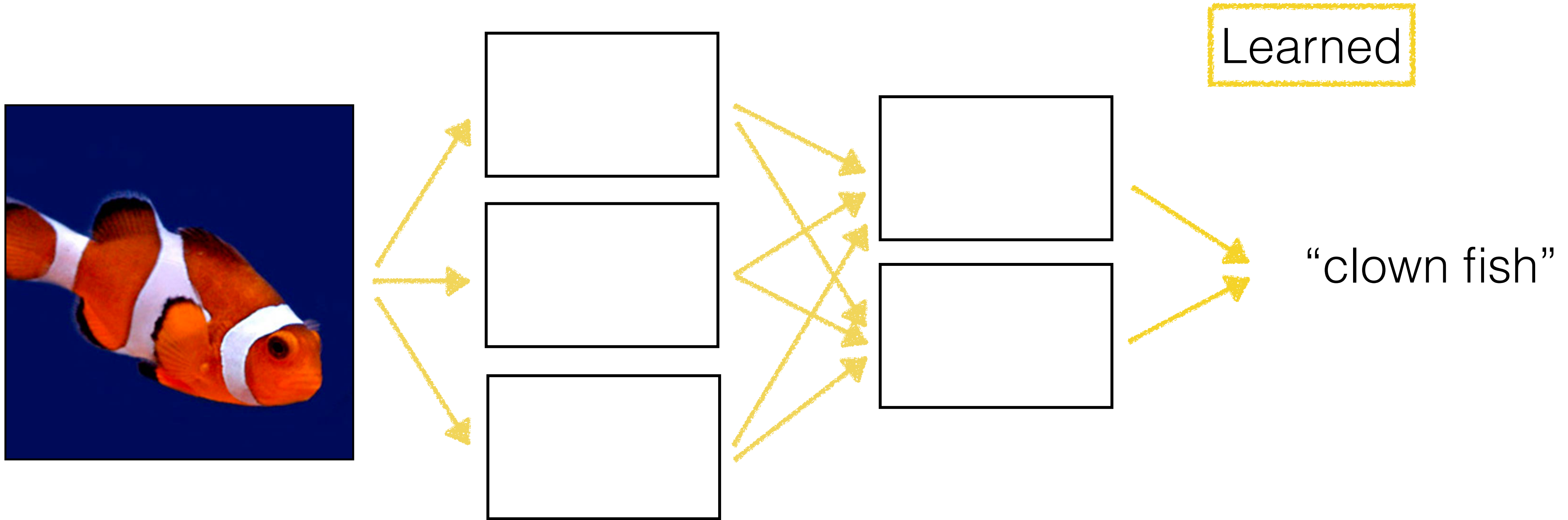


# Image classification

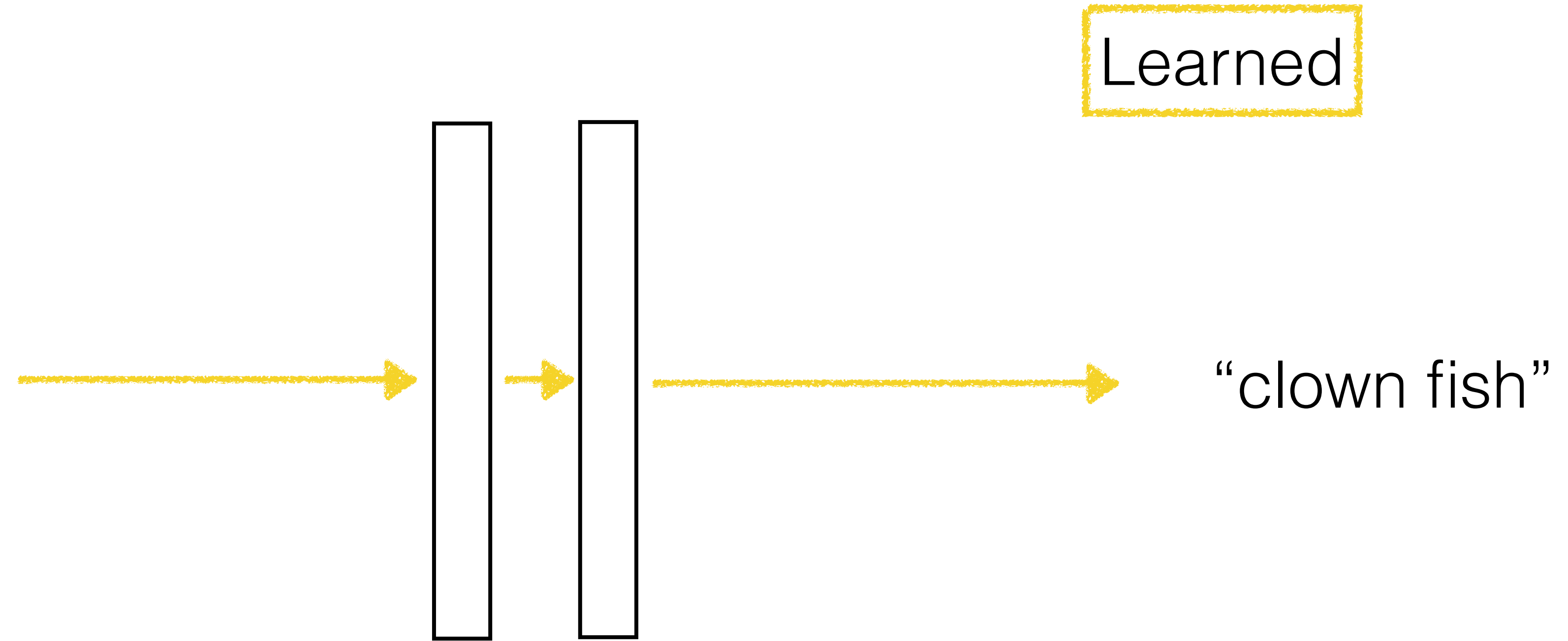
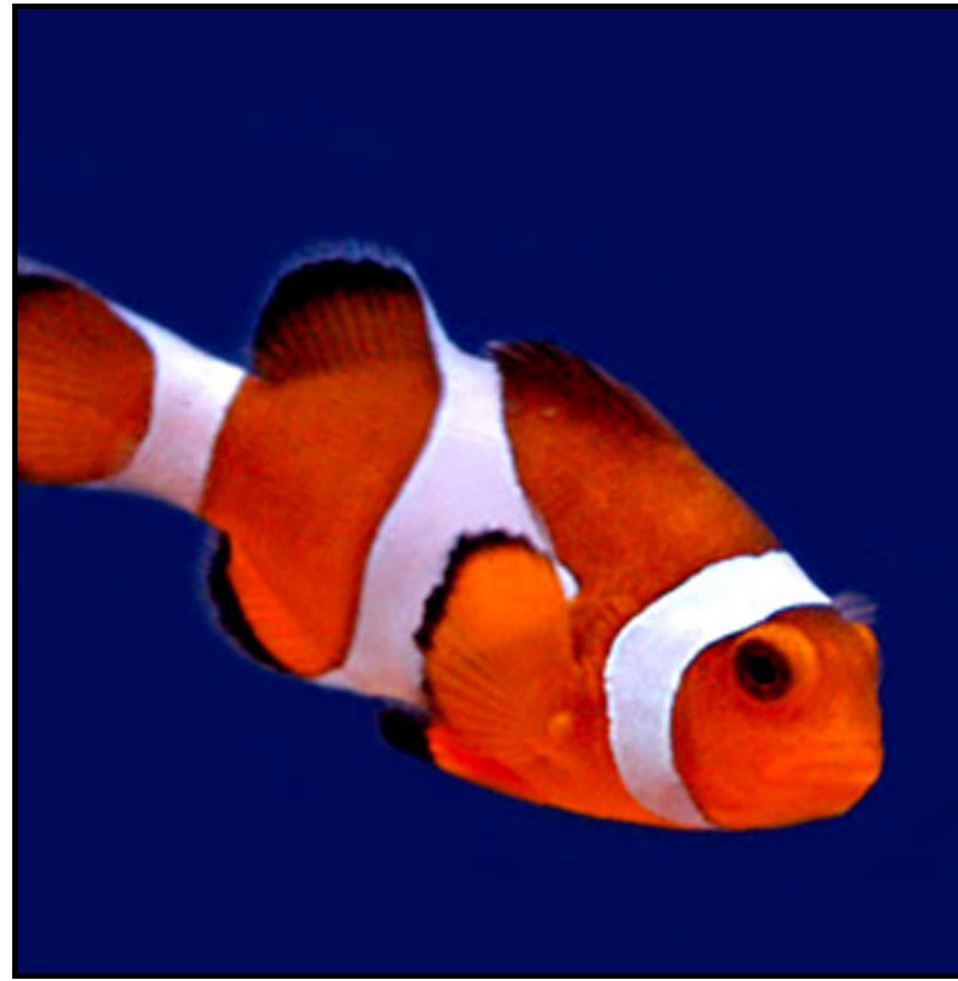




# Image classification

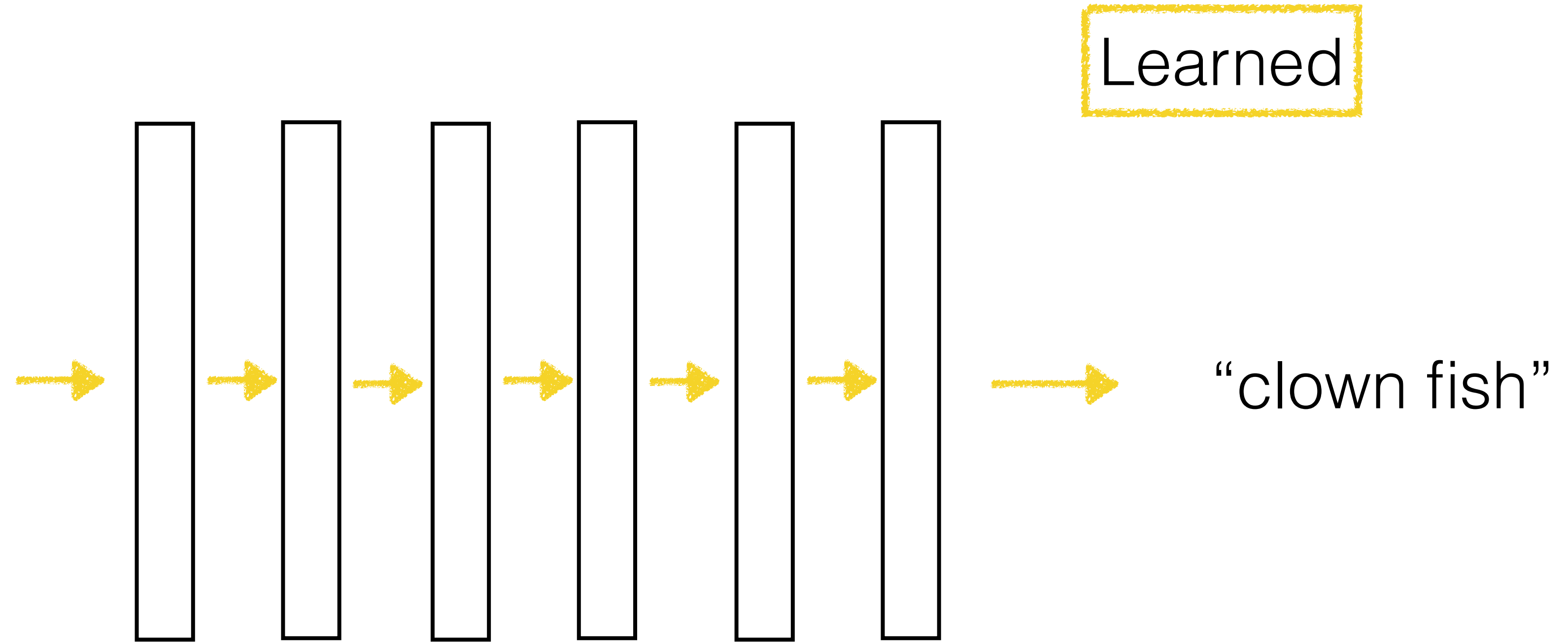


# Image classification



**Neural net**

# Image classification



**Deep neural net**






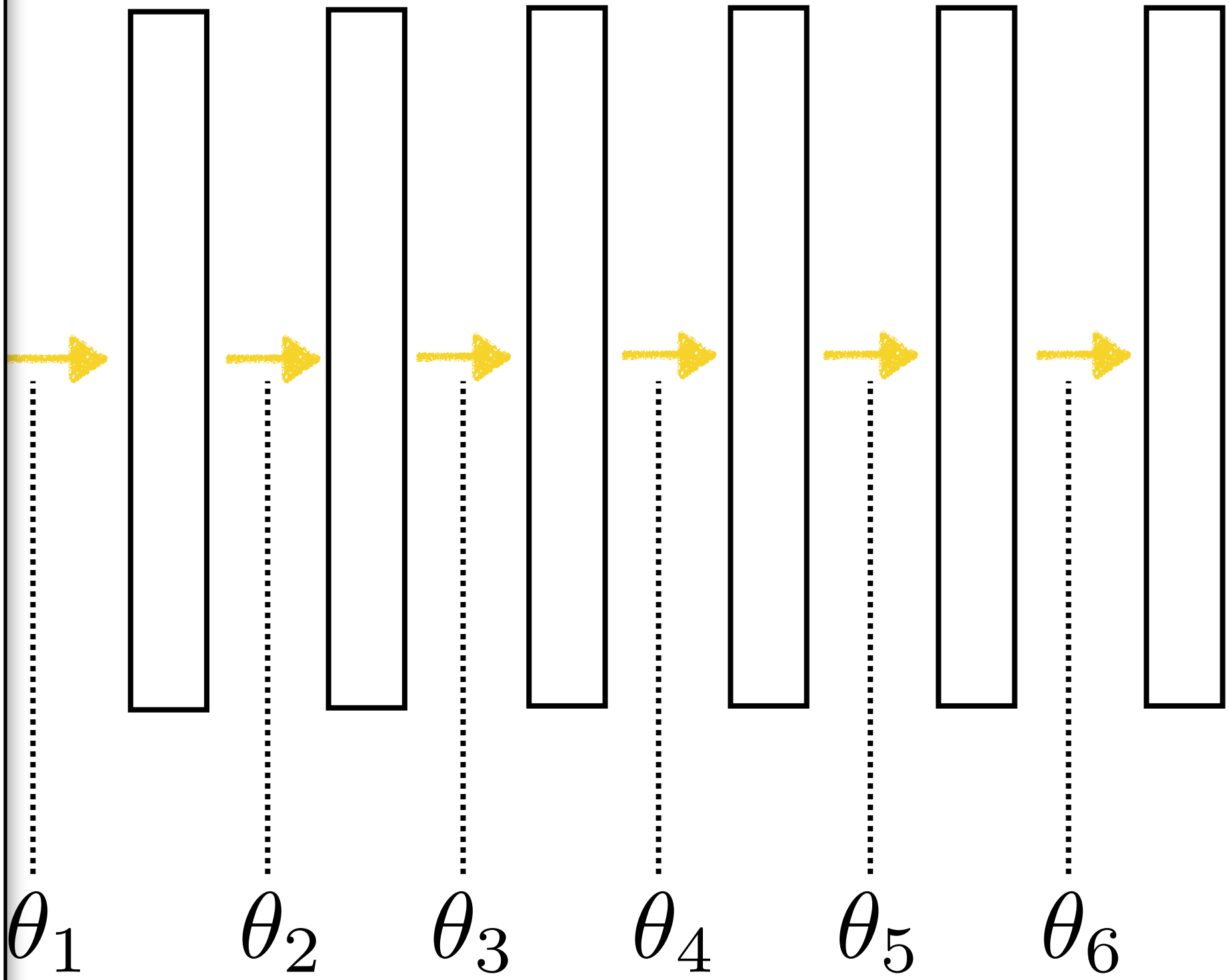
# Deep learning

Learned

$\mathbf{y}^{(i)}$   
“clown fish”

Training data

$\mathbf{x}$	$y$
	“Fish”
	“Grizzly”
	“Chameleon”
⋮	



Loss

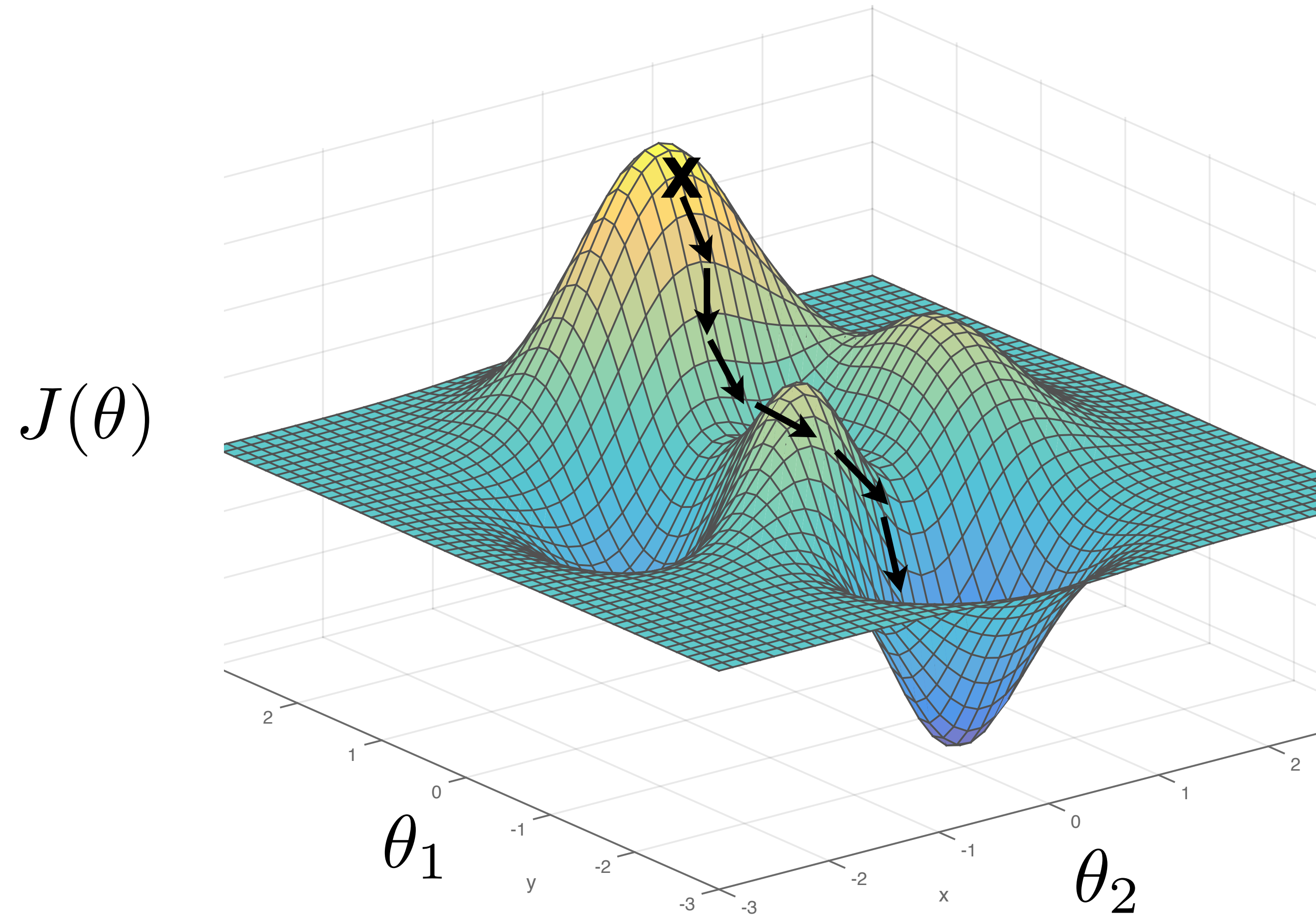
$$\mathcal{L}(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

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

# Gradient descent

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

# Gradient descent



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



# Gradient descent

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

One iteration of gradient descent:

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

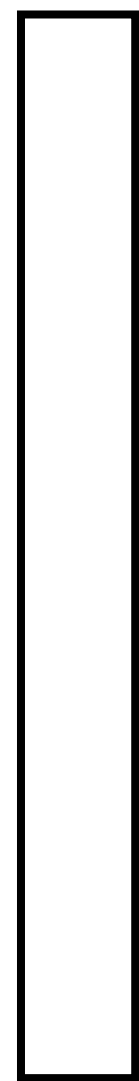
**learning rate**

# Stochastic gradient descent (SGD)

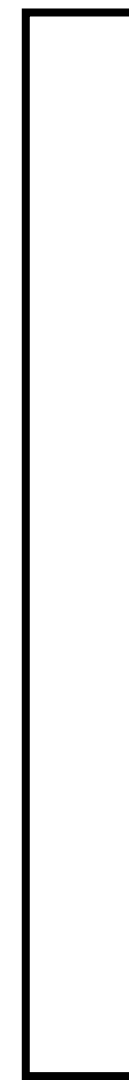
- Want to minimize overall loss function  $\mathbf{J}$ , which is sum of individual losses over each example.
- In Stochastic gradient descent, compute gradient on sub-set (batch) of data.
  - If  $\text{batchsize}=1$  then  $\theta$  is updated after each example.
  - If  $\text{batchsize}=N$  (full set) then this is standard gradient descent.
- Gradient direction is noisy, relative to average over all examples (standard gradient descent).
- Advantages
  - Faster: approximate total gradient with small sample
  - Implicit regularizer
- Disadvantages
  - High variance, unstable updates

# Computation in a neural net

Input  
representation



Output  
representation



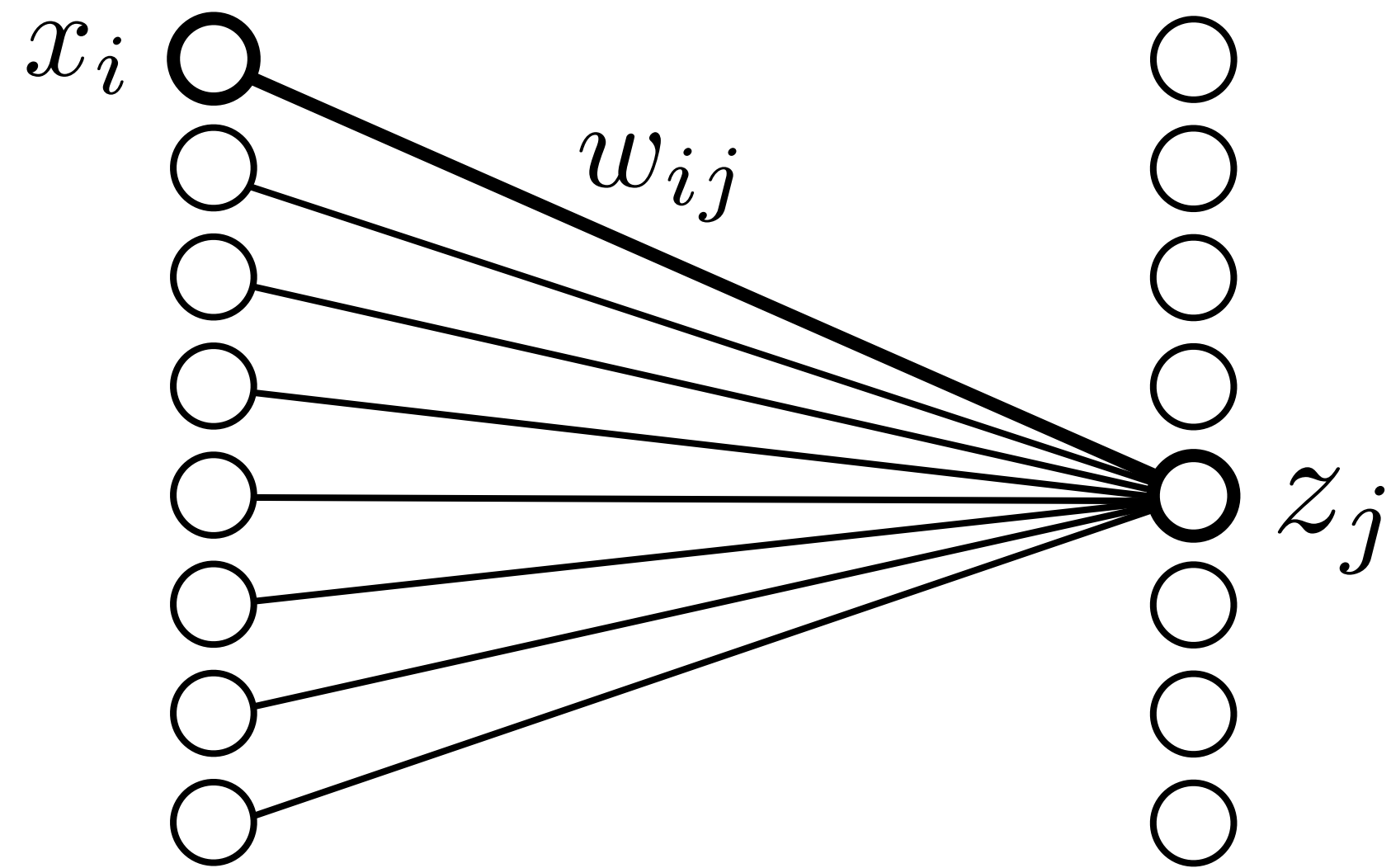


# Computation in a neural net

## Linear layer

Input  
representation

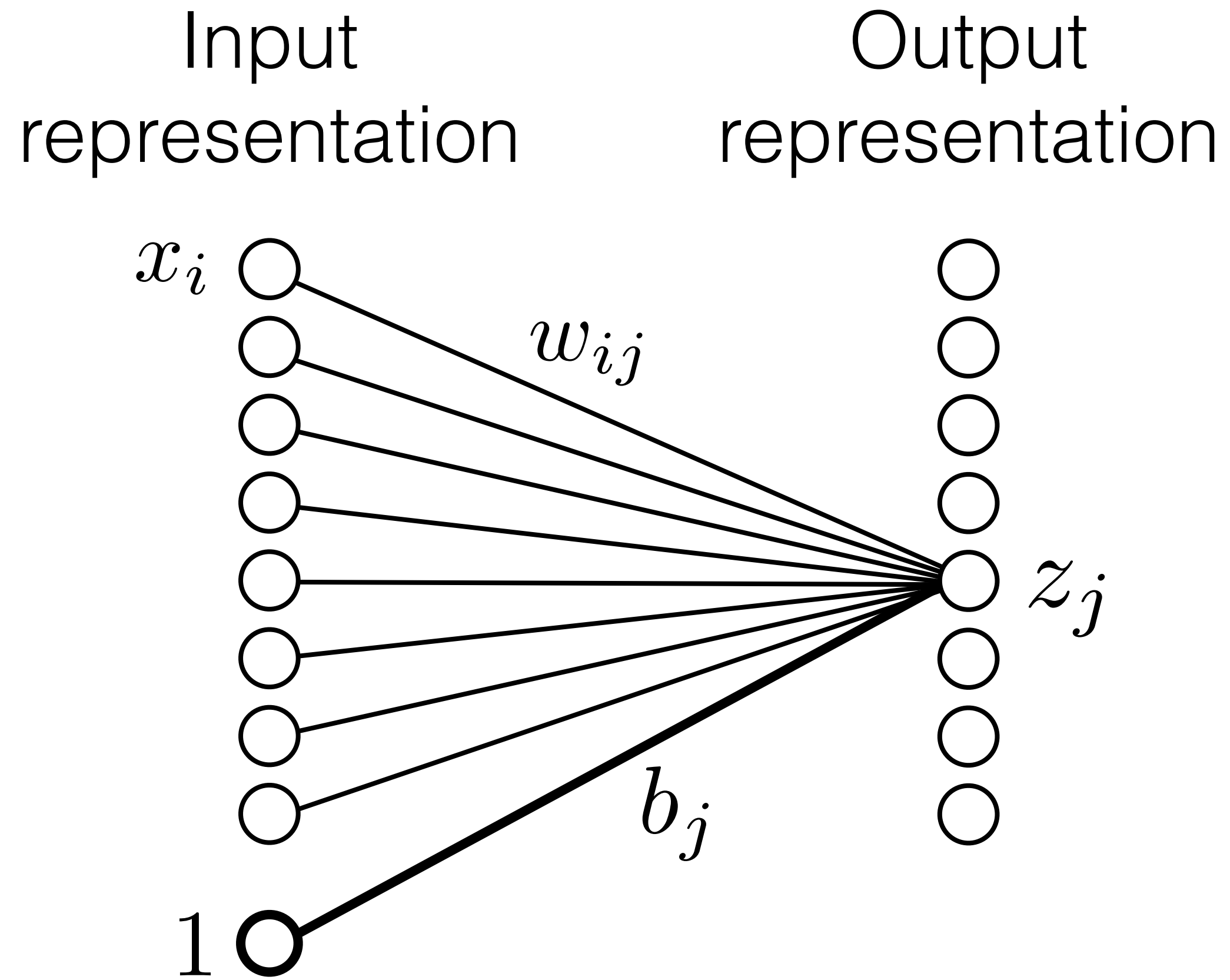
Output  
representation



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

# Computation in a neural net

## Linear layer



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

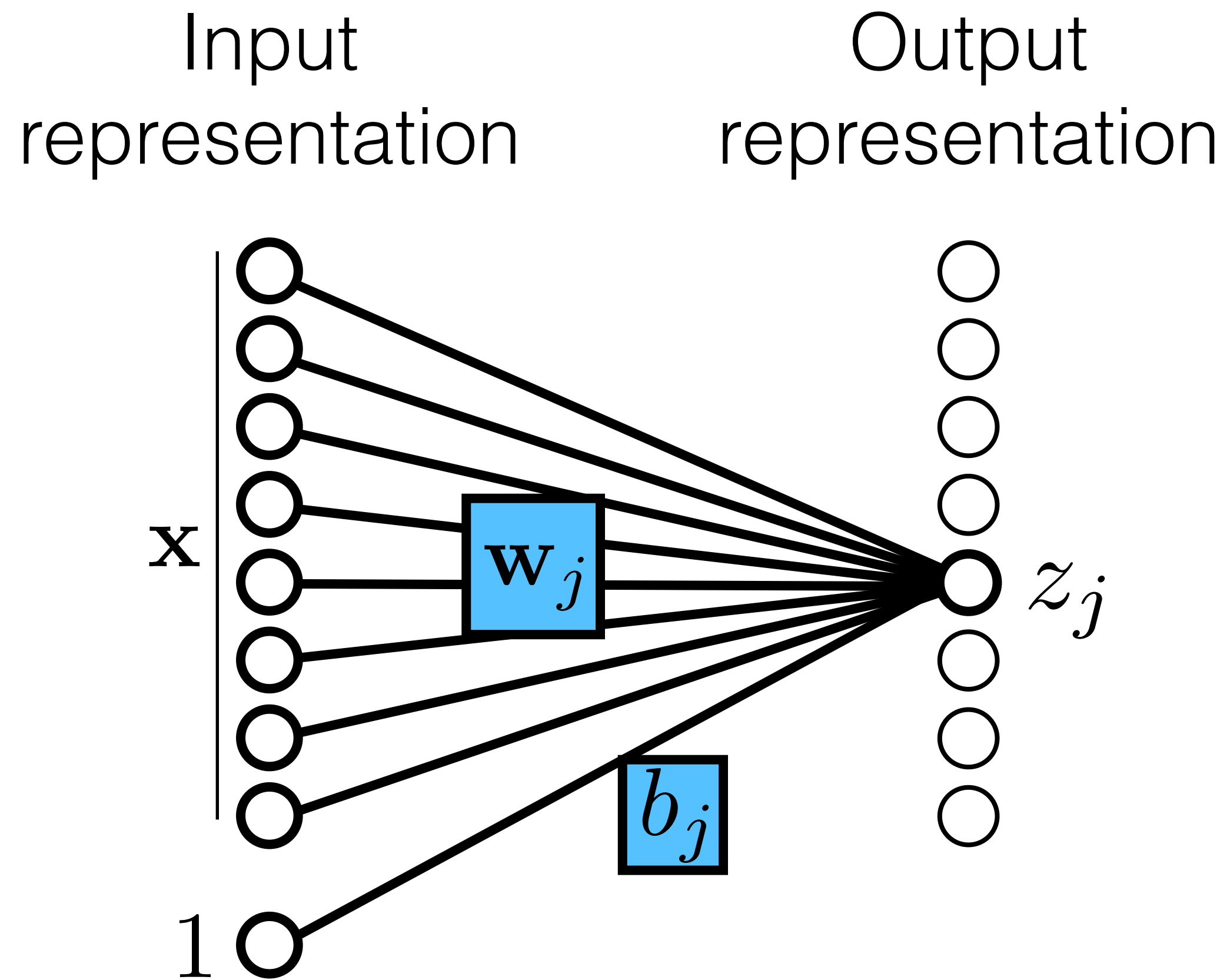
weights

bias

The equation shows the computation of the output  $z_j$  as a weighted sum of the input  $x_i$  plus a bias  $b_j$ . An arrow labeled "weights" points to  $w_{ij}$ , and an arrow labeled "bias" points to  $b_j$ .

# Computation in a neural net

## Linear layer



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

weights

bias

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

parameters of the model

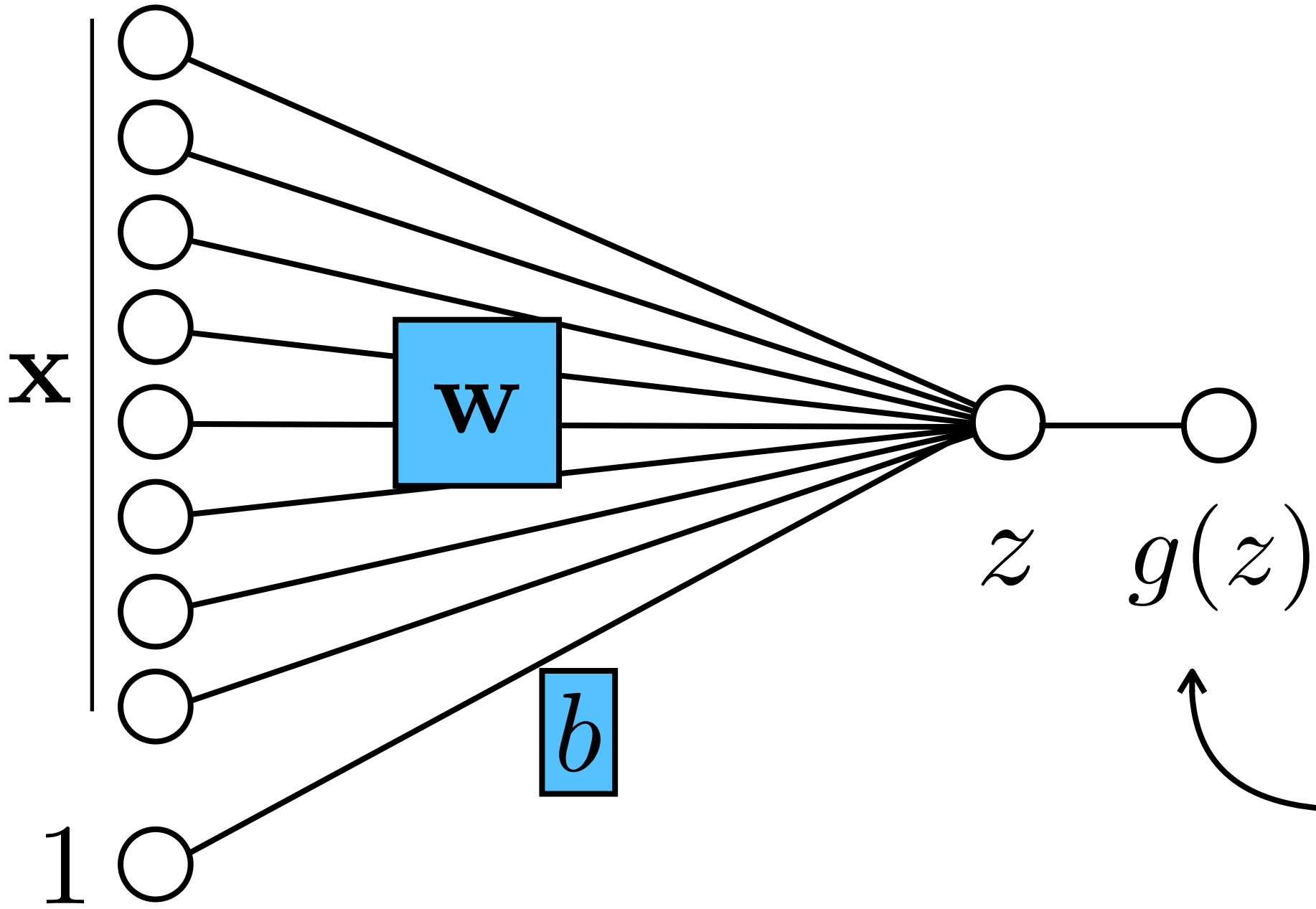


# Computation in a neural net

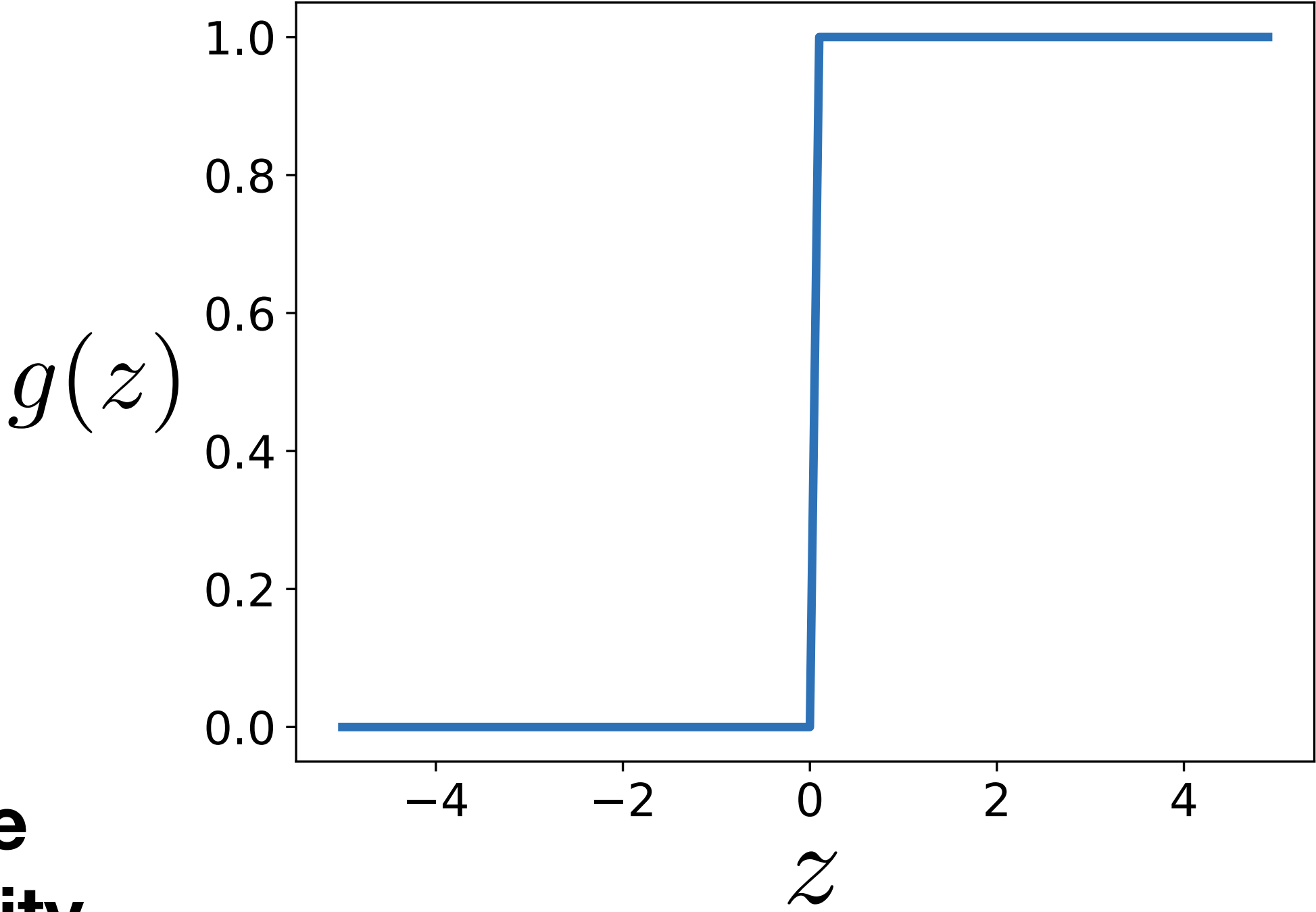
## “Perceptron”

Input representation

Output representation

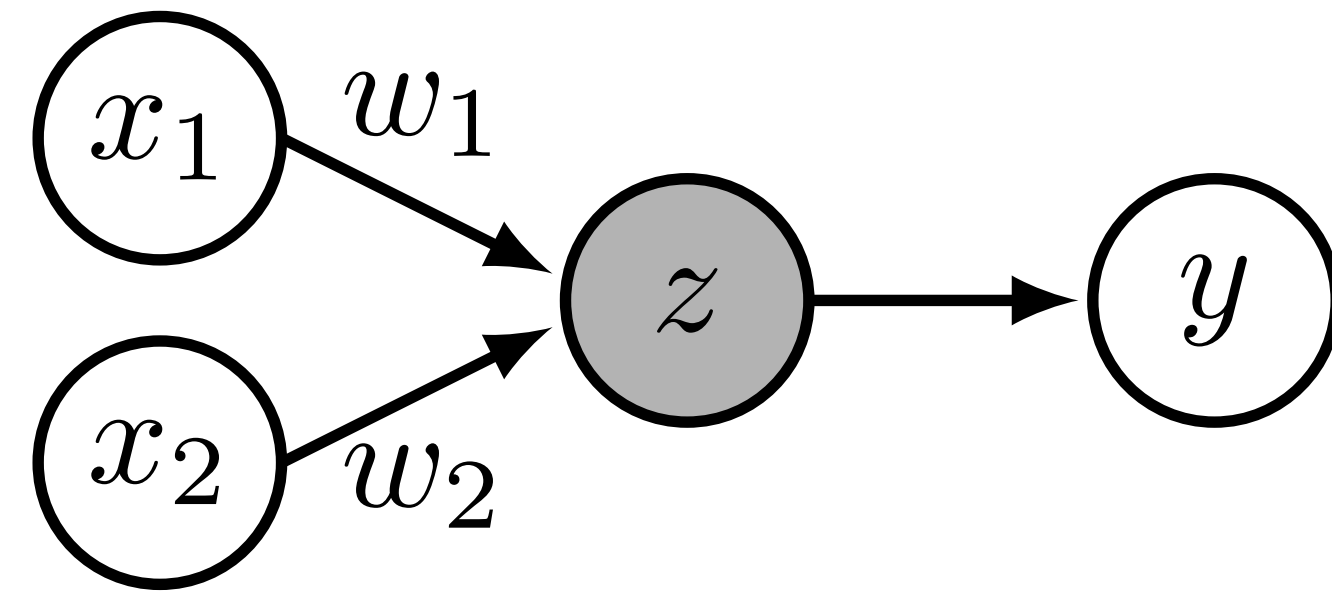


$$g(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{otherwise} \end{cases}$$



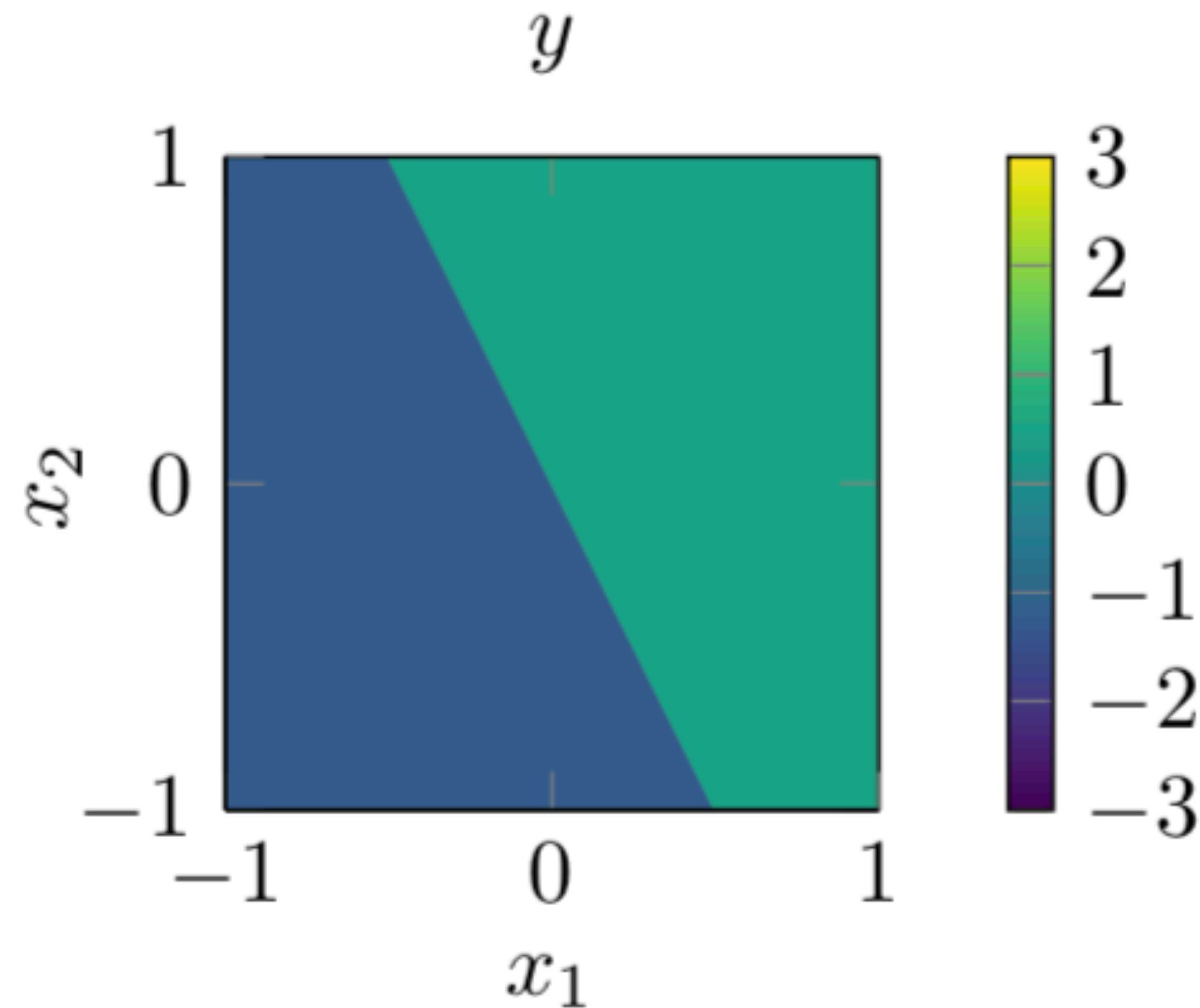
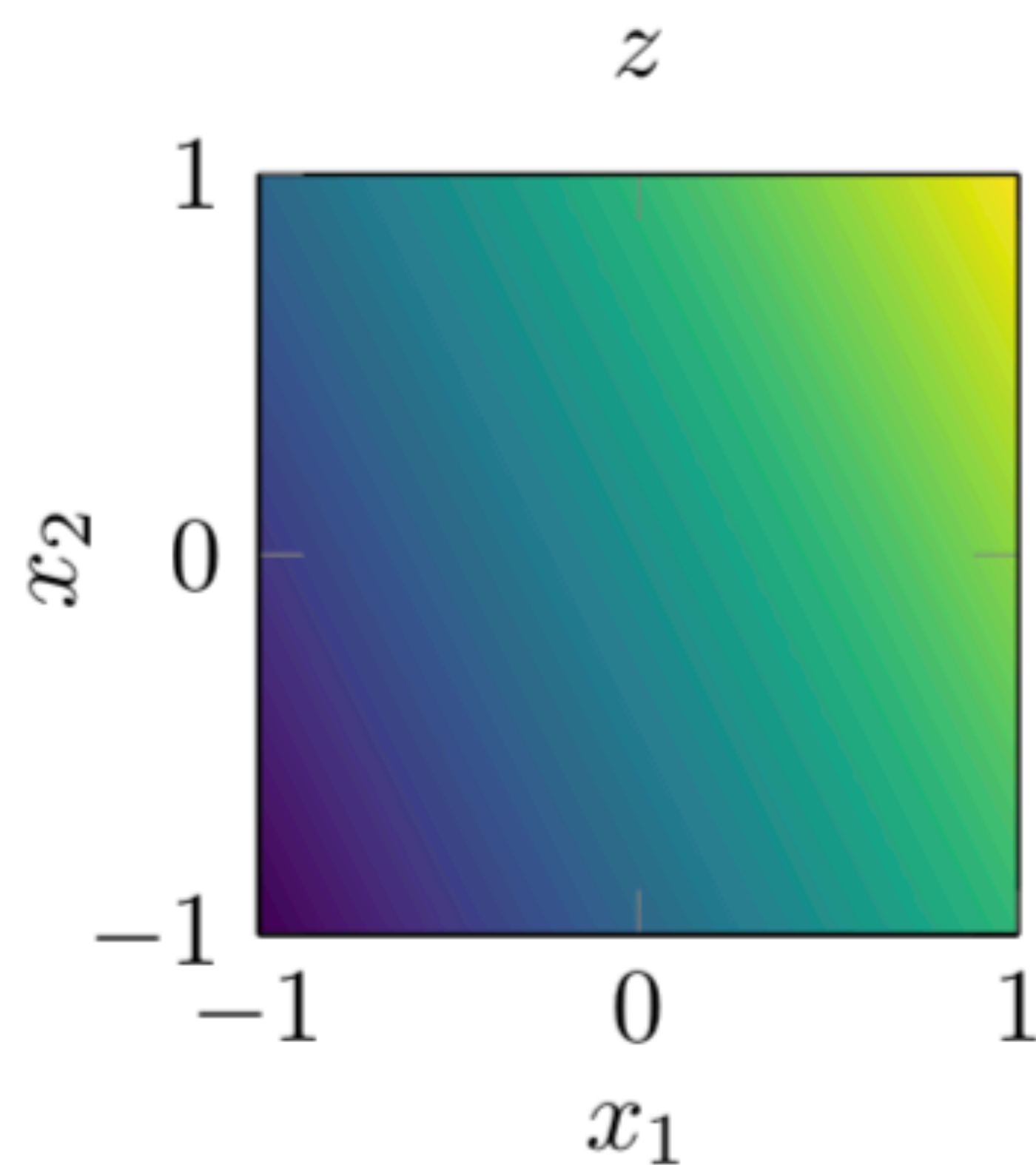
**Pointwise Non-linearity**

# Example: linear classification with a perceptron

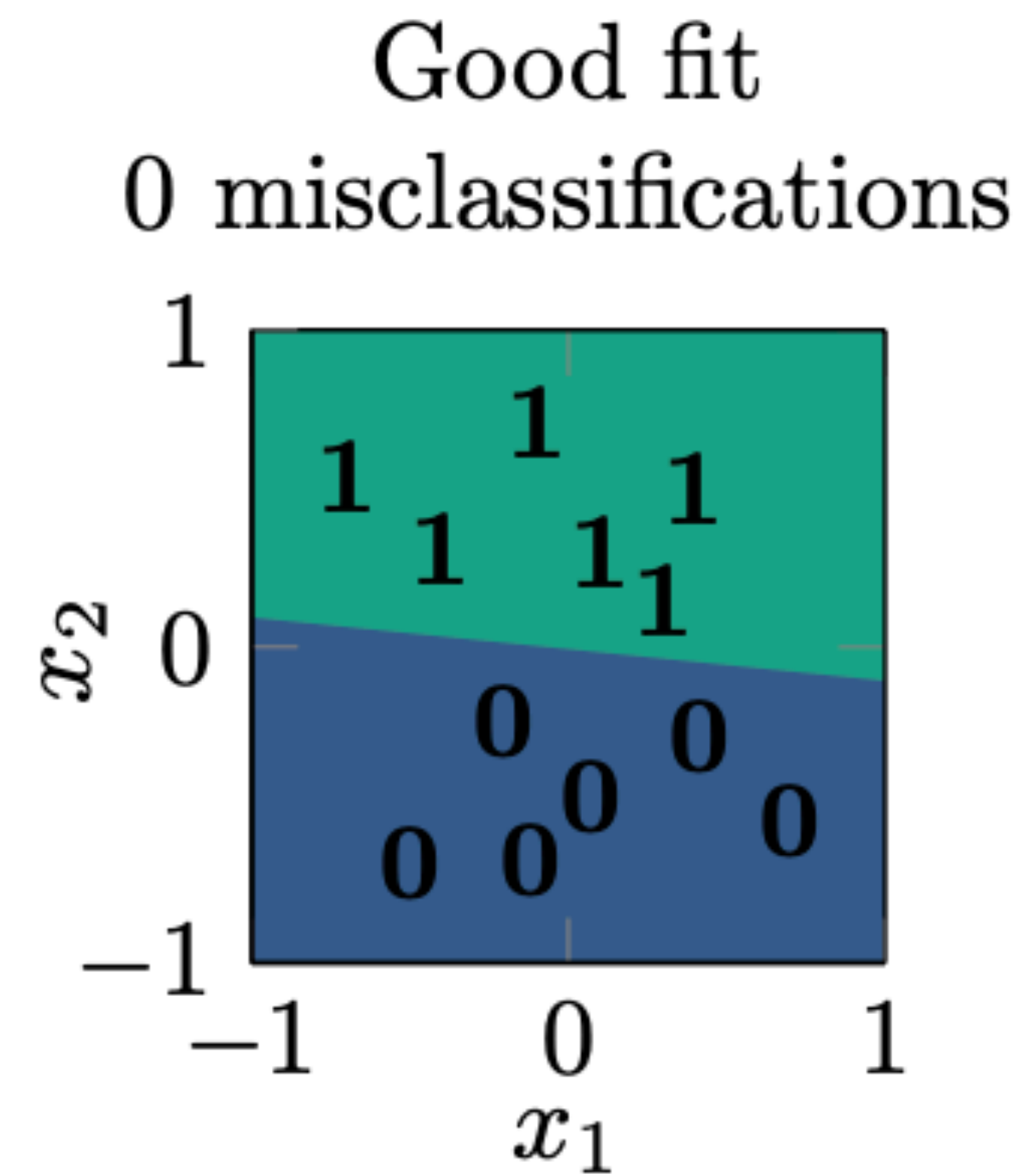
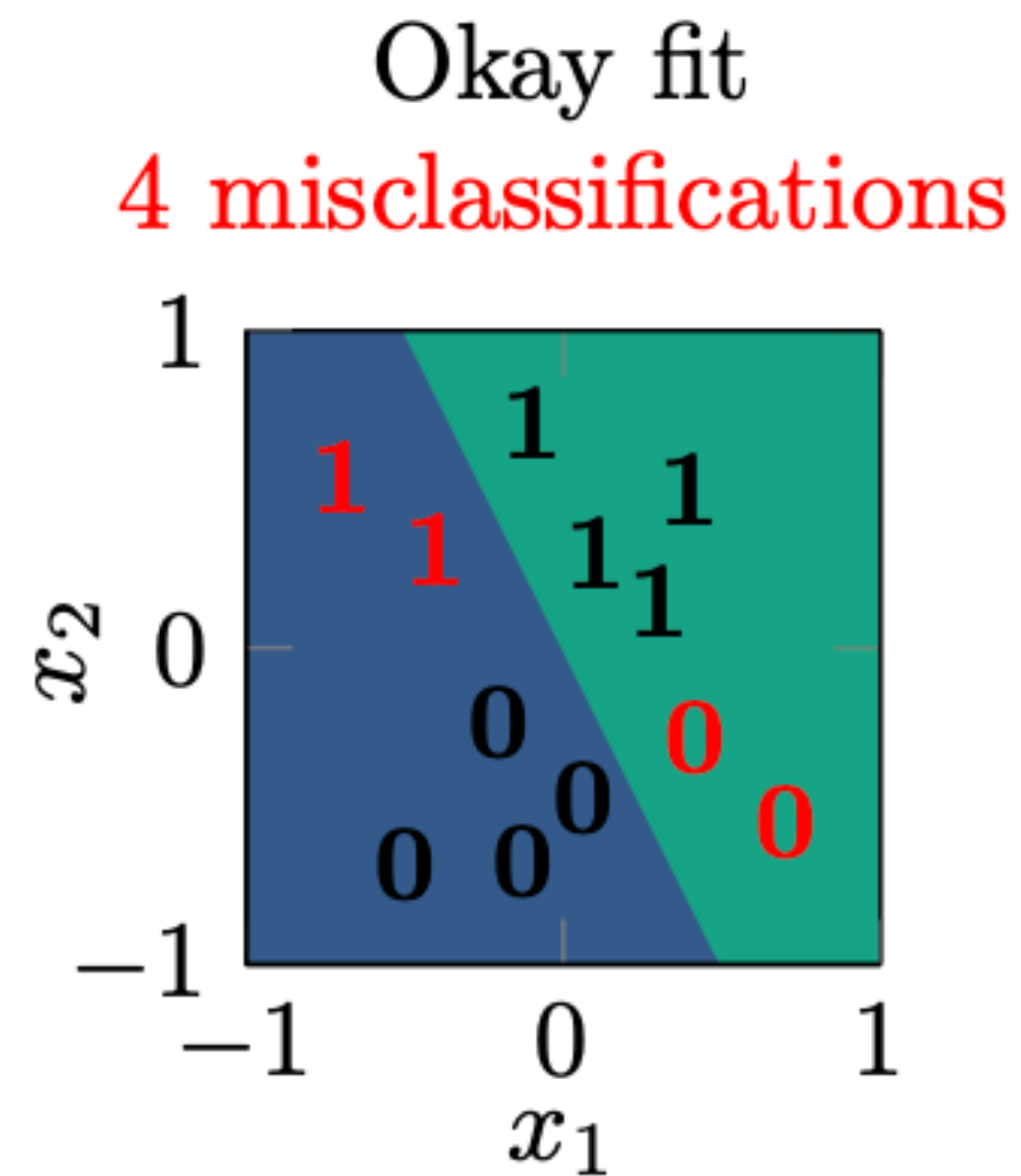
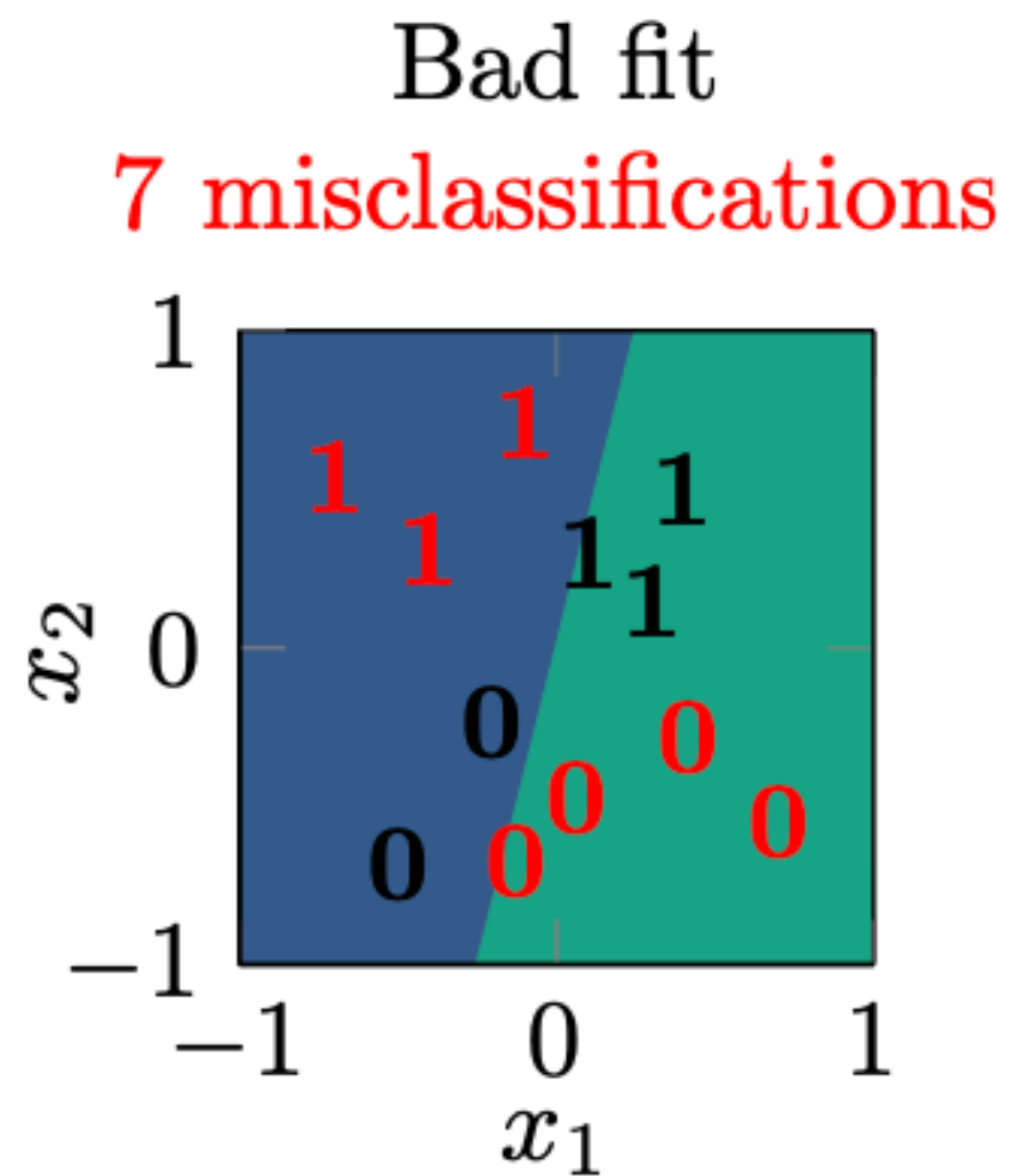
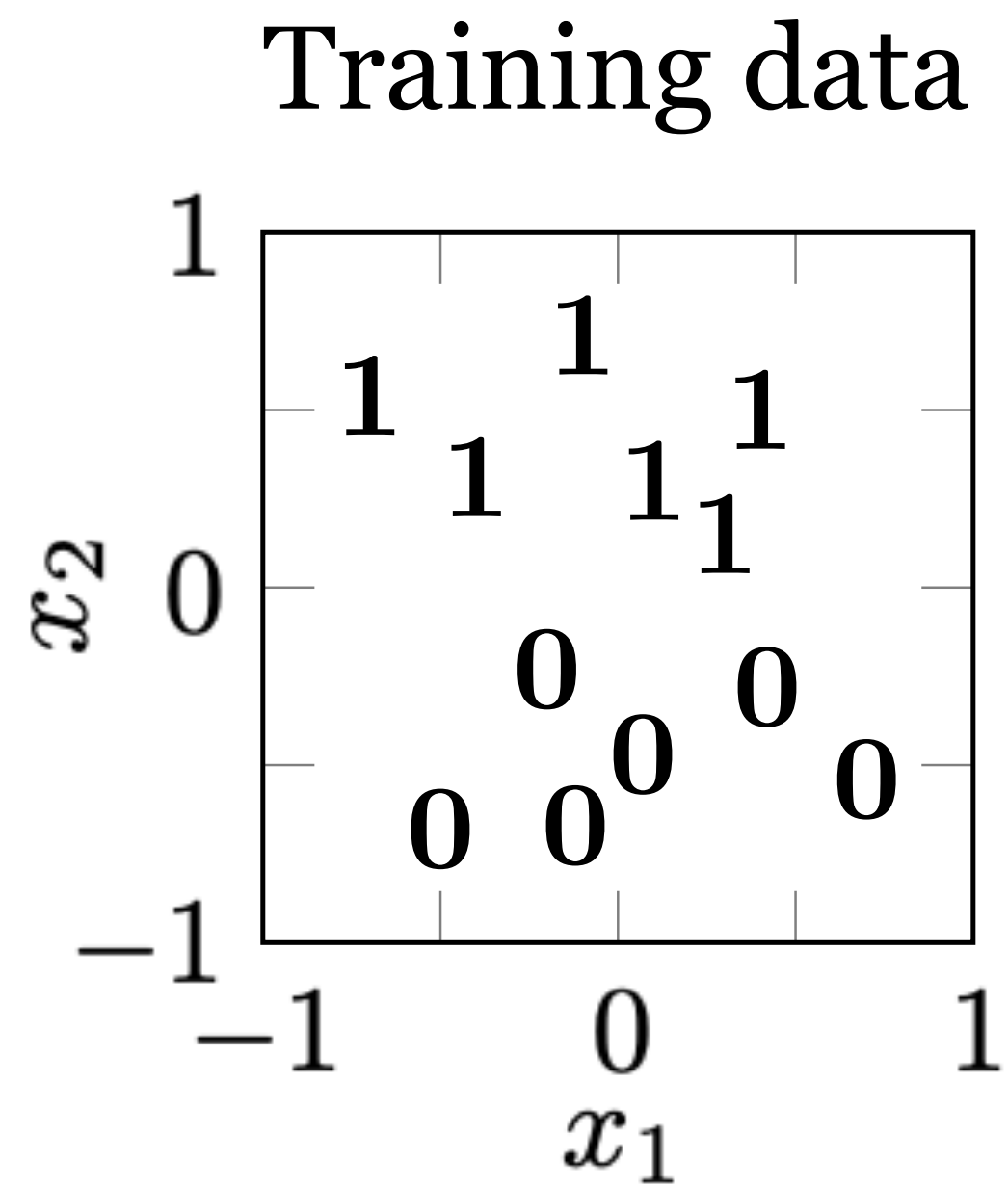


$$z = \mathbf{x}^T \mathbf{w} + b$$

$$y = g(z)$$



One layer neural net (perceptron) can perform linear classification!



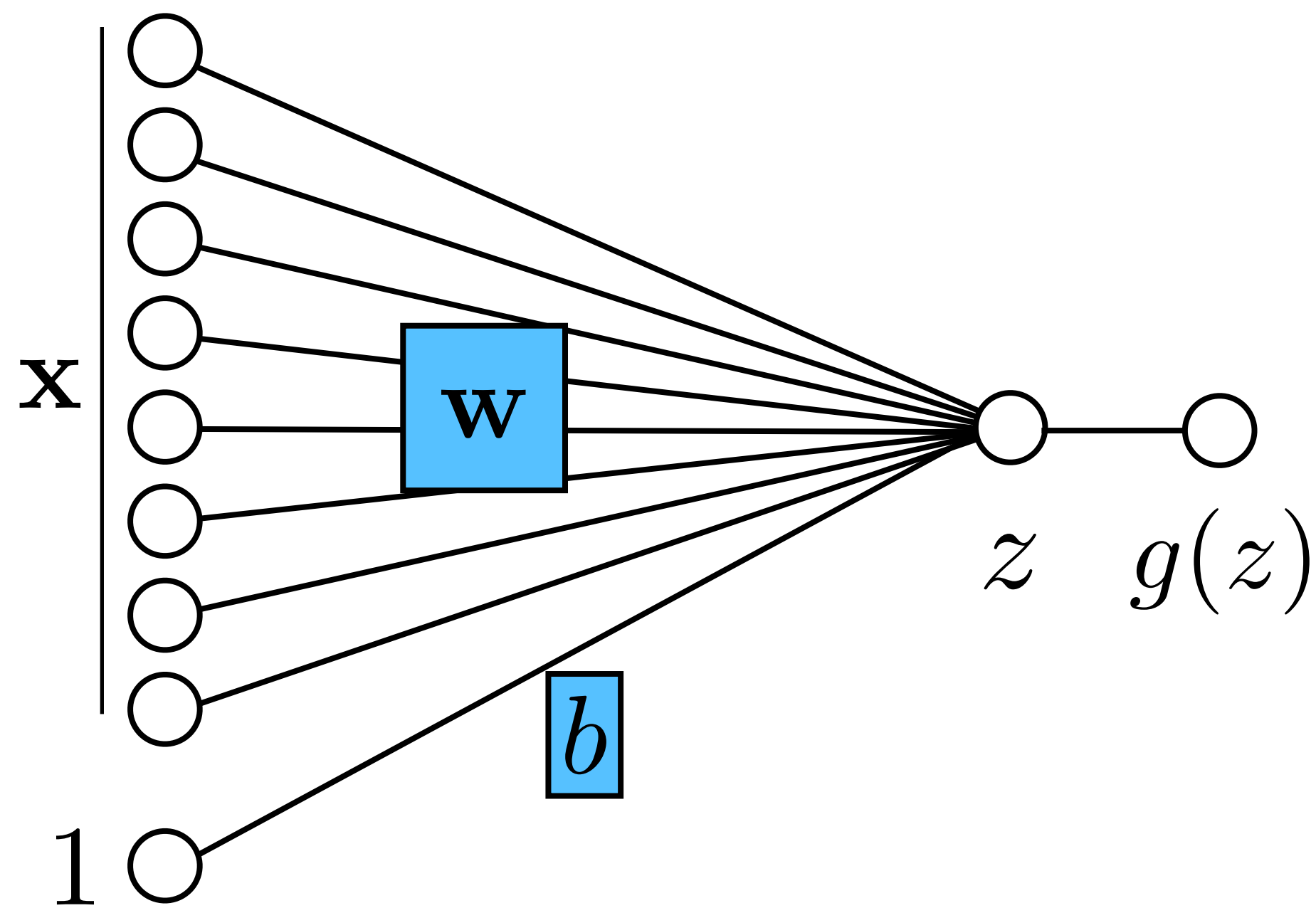
$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \sum_{i=1}^N \mathcal{L}(g(z^{(i)}), y^{(i)})$$



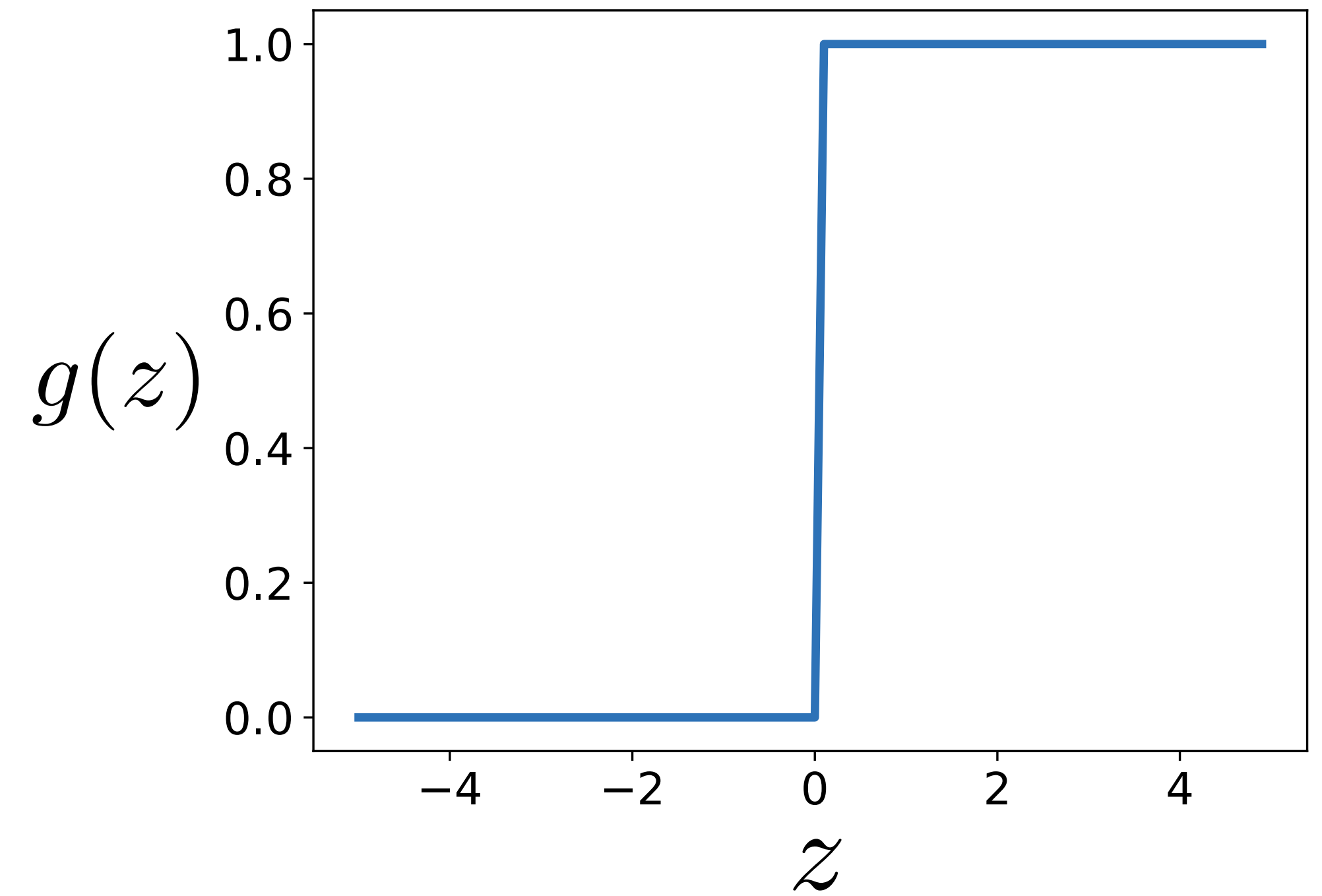
# Computation in a neural net

Input representation

Output representation



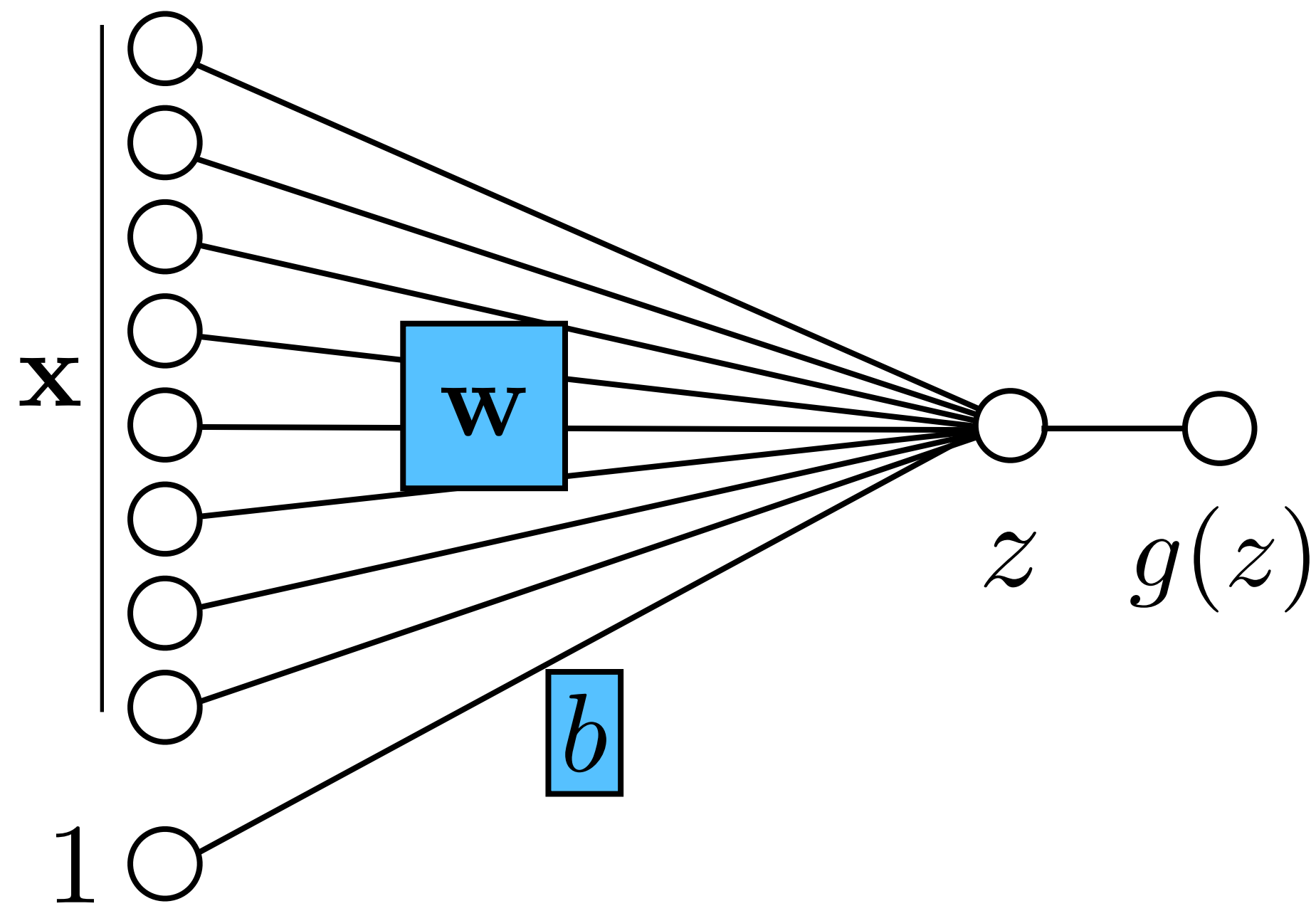
$$g(z) = \begin{cases} 1, & \text{if } z > 0 \\ 0, & \text{otherwise} \end{cases}$$



# Computation in a neural net — nonlinearity

Input representation

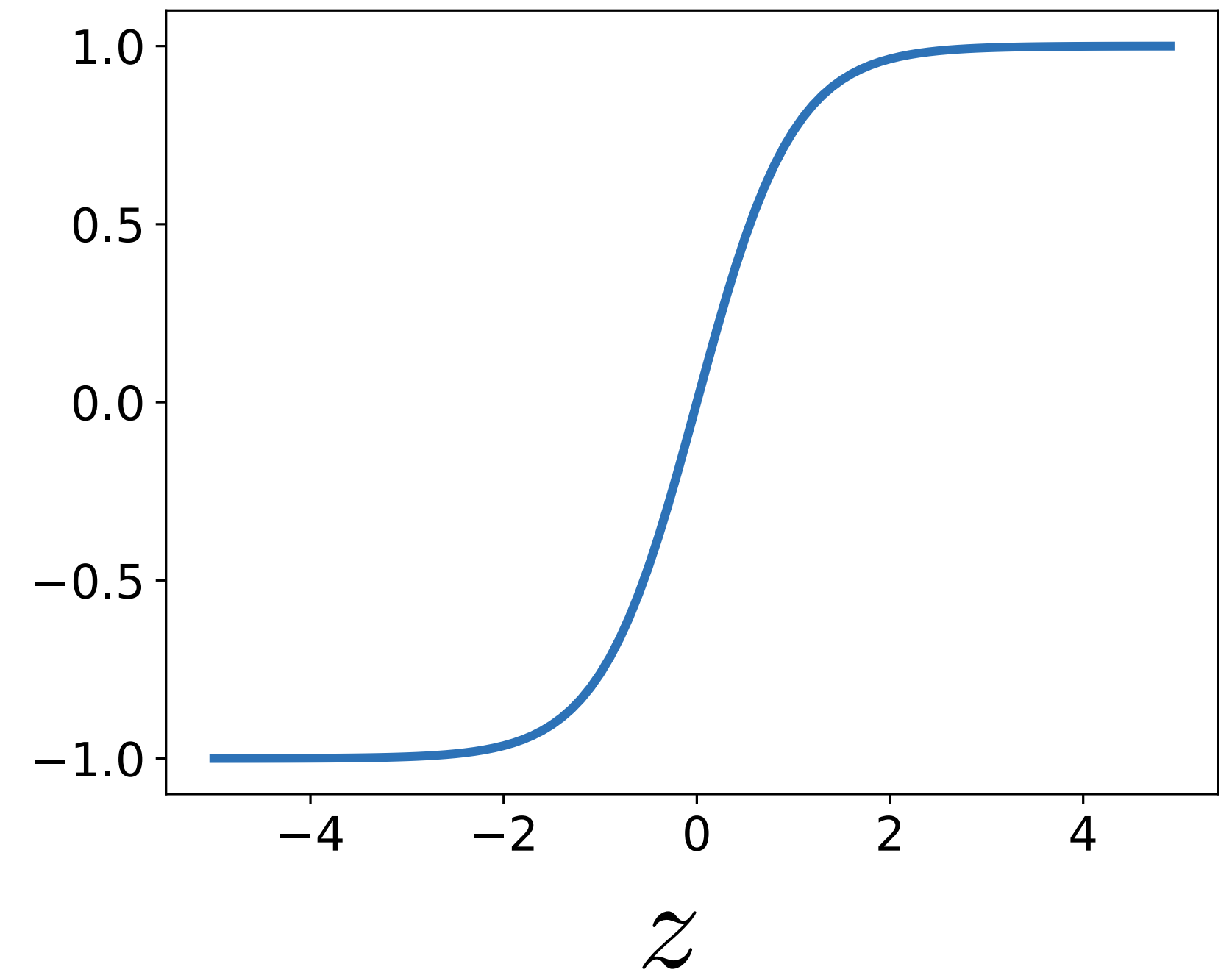
Output representation



**Tanh**

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

$g(z)$

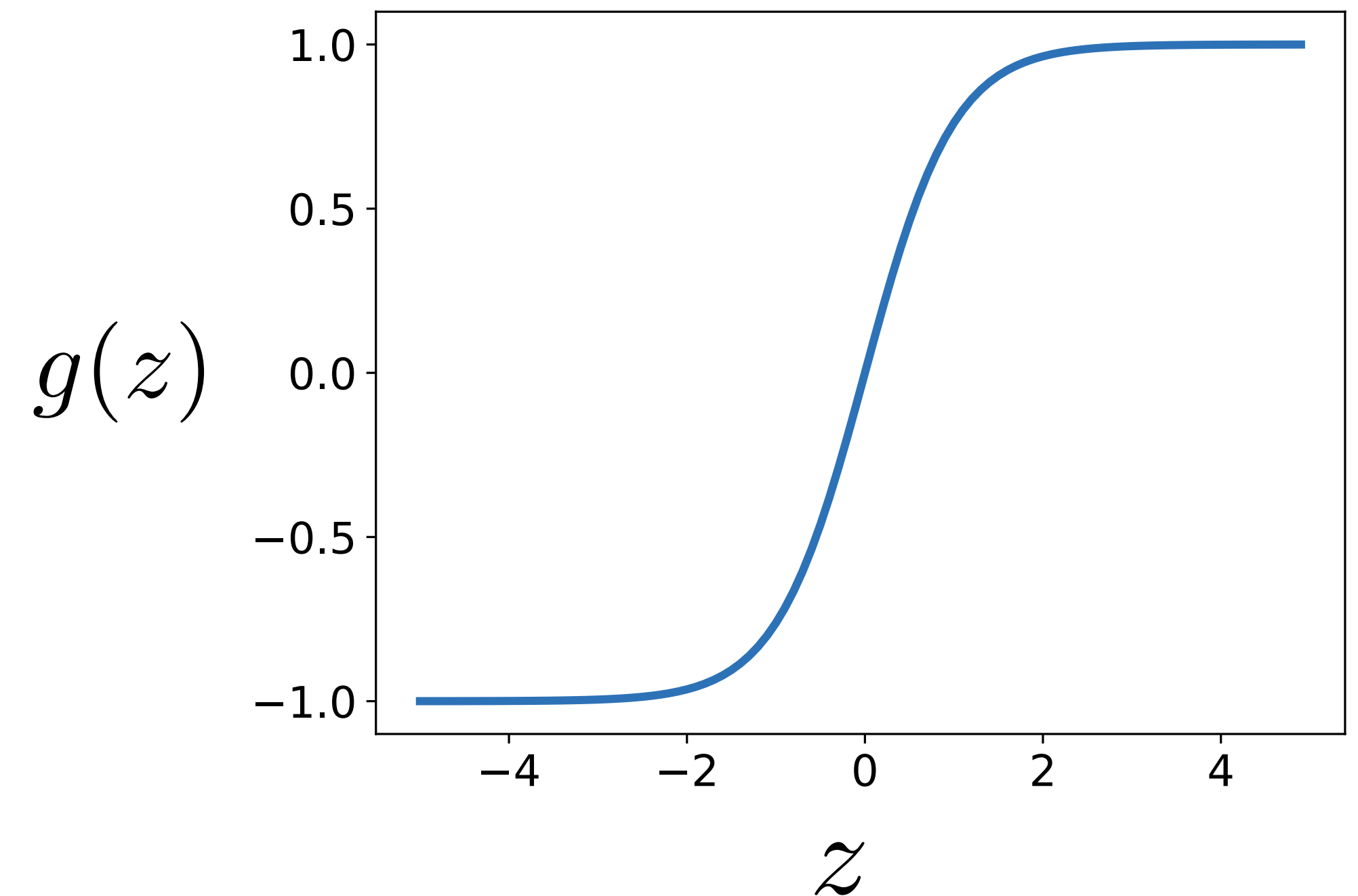


# Computation in a neural net — nonlinearity

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

**Tanh**

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



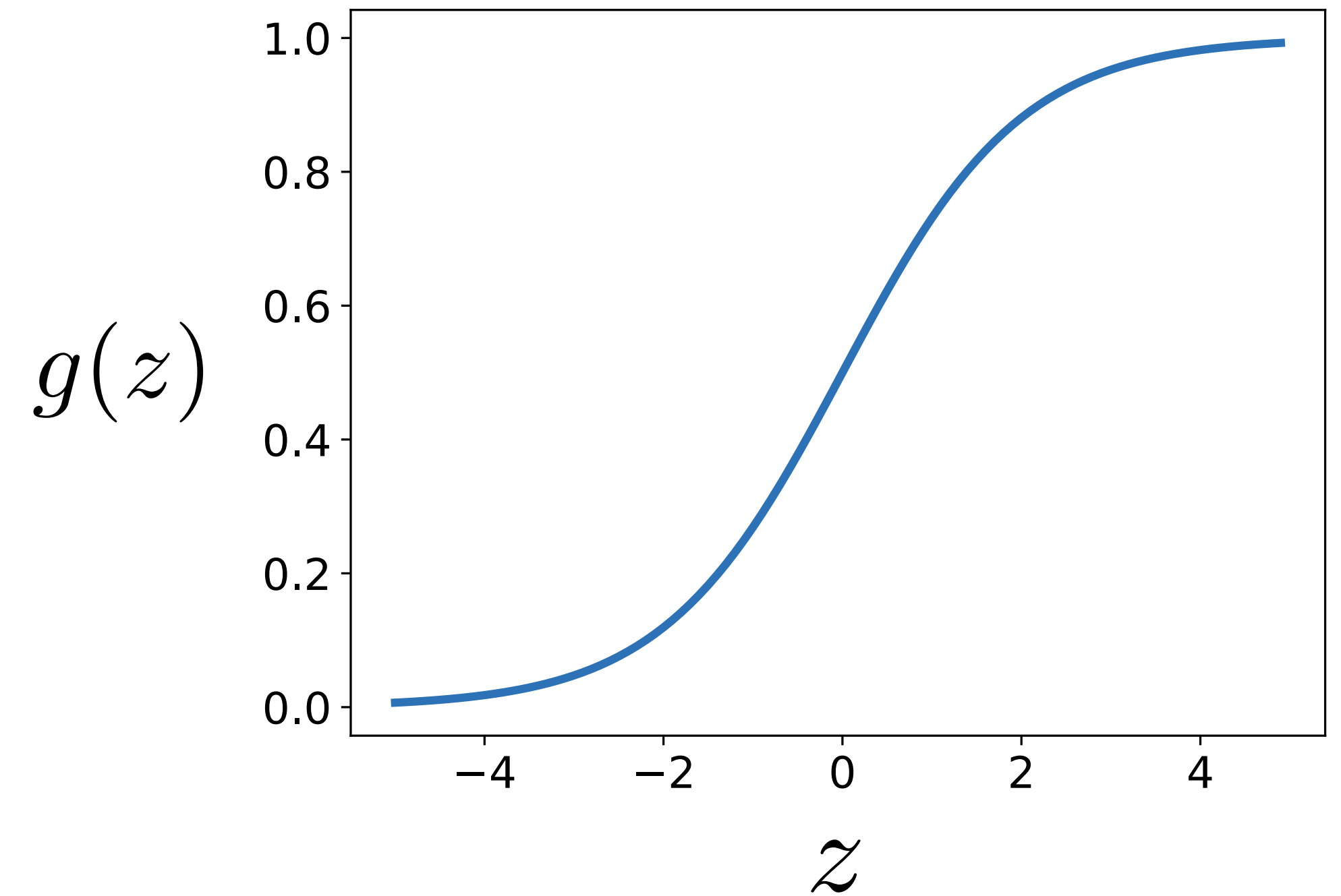


# Computation in a neural net — nonlinearity

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

## Sigmoid

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



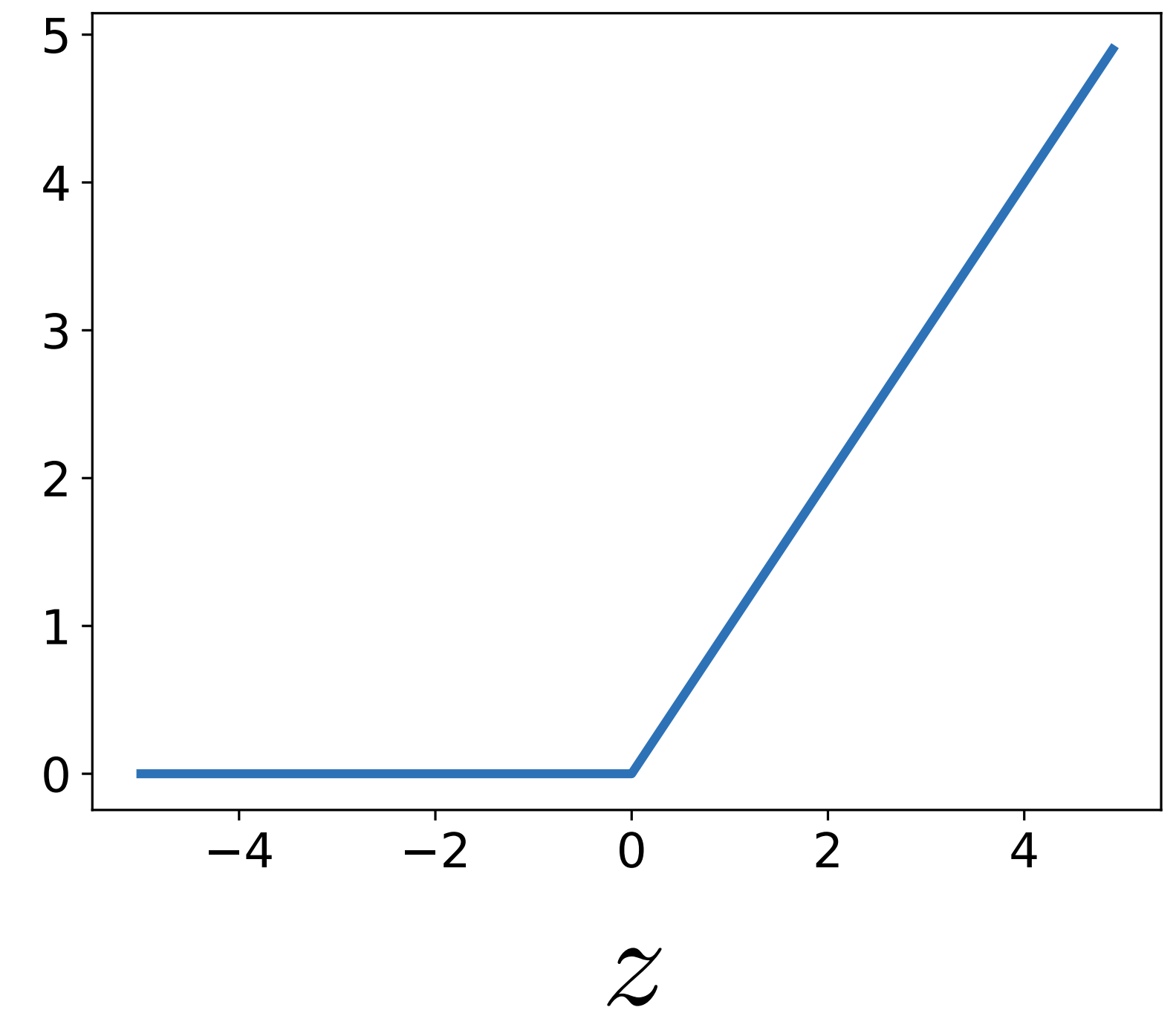
# Computation in a neural net — nonlinearity

- Unbounded output (on positive side)
- Efficient to implement:  $\frac{\partial g}{\partial z} = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{if } z \geq 0 \end{cases}$
- Also seems to help convergence (see 6x speedup vs tanh in [Krizhevsky et al.])
- Drawback: if strongly in negative region, unit is dead forever (no gradient).
- Default choice: widely used in current models.

## Rectified linear unit (ReLU)

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

$g(z)$



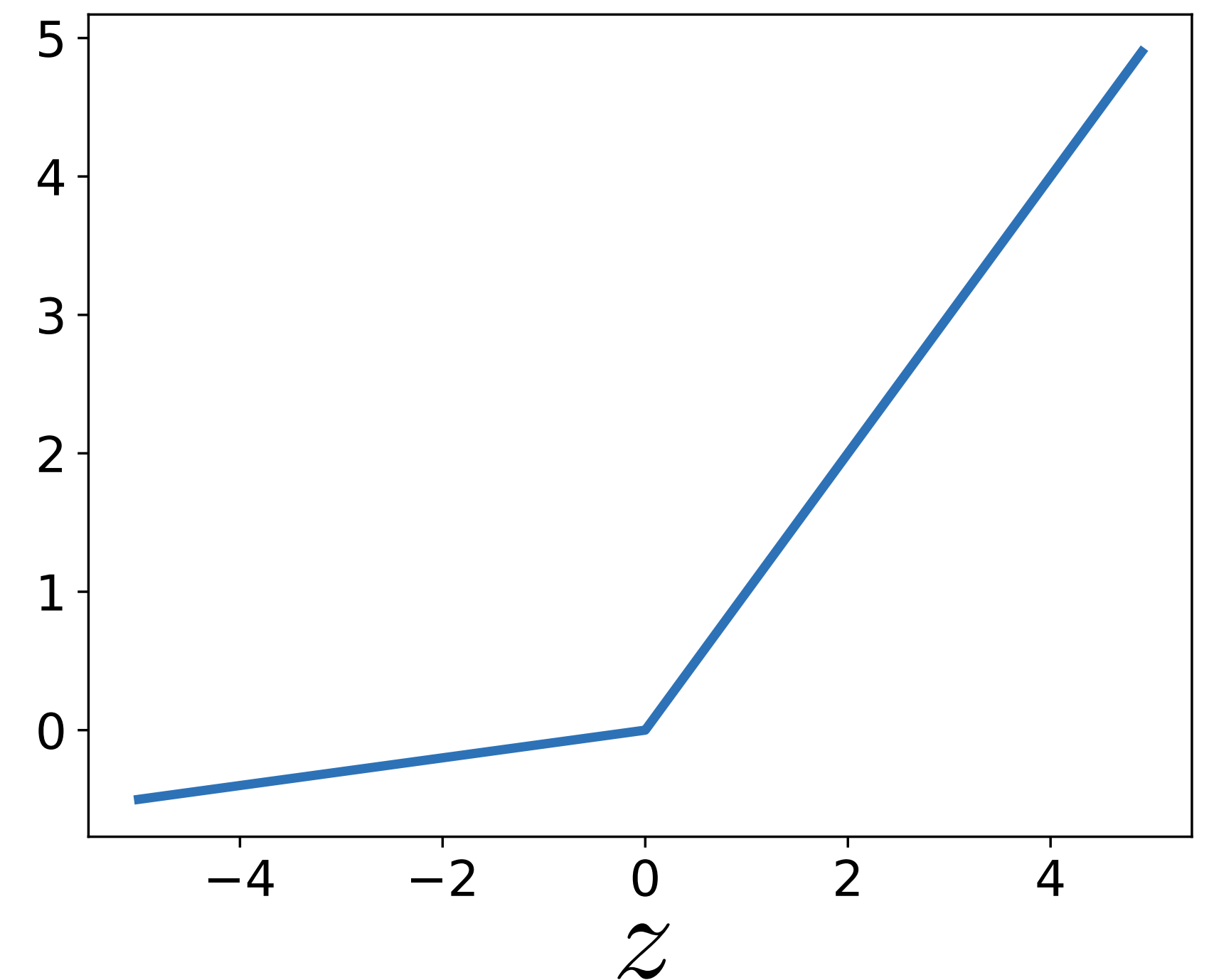
# Computation in a neural net — nonlinearity

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

## Leaky ReLU

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

$g(z)$



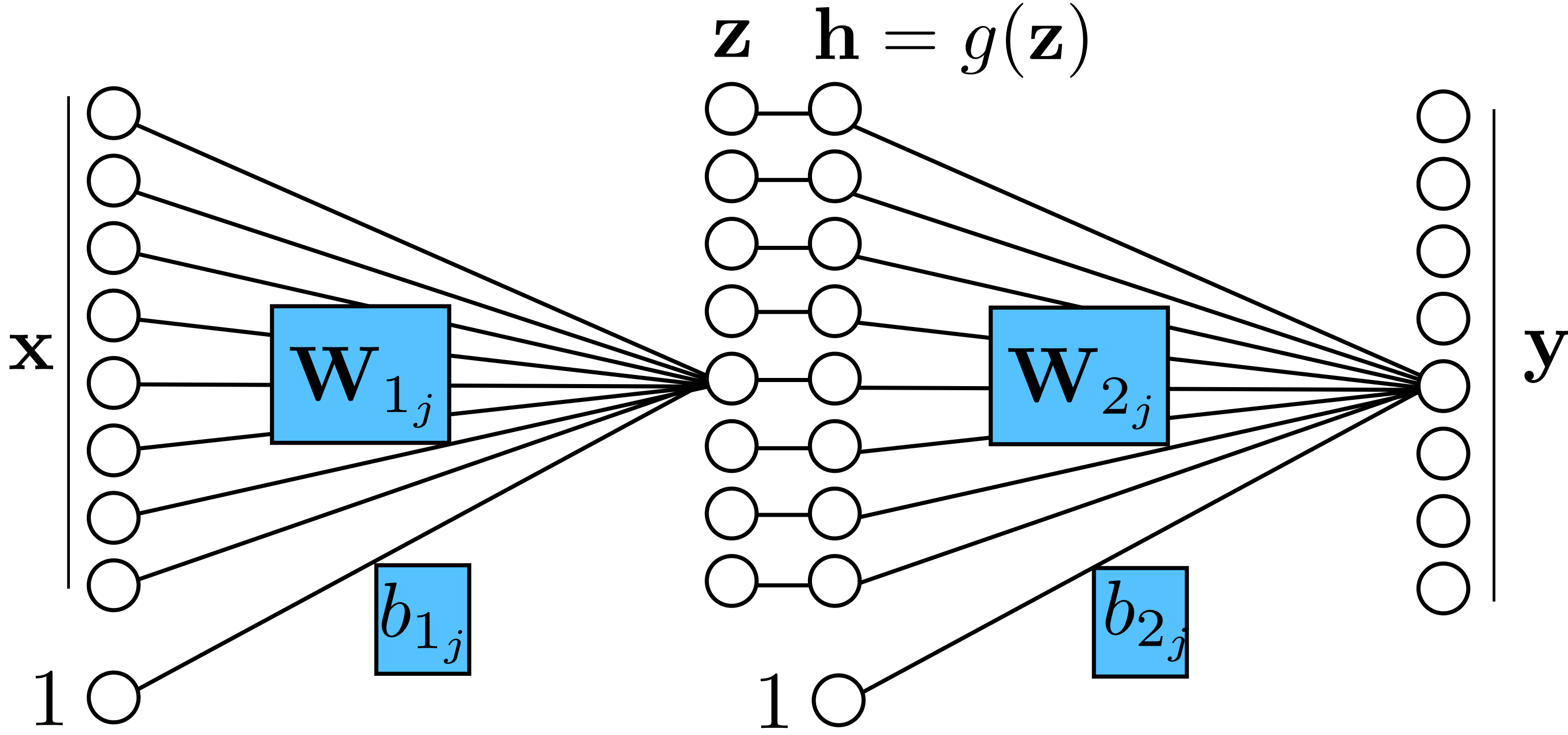


# Stacking layers

Input representation

Intermediate representation

Output representation



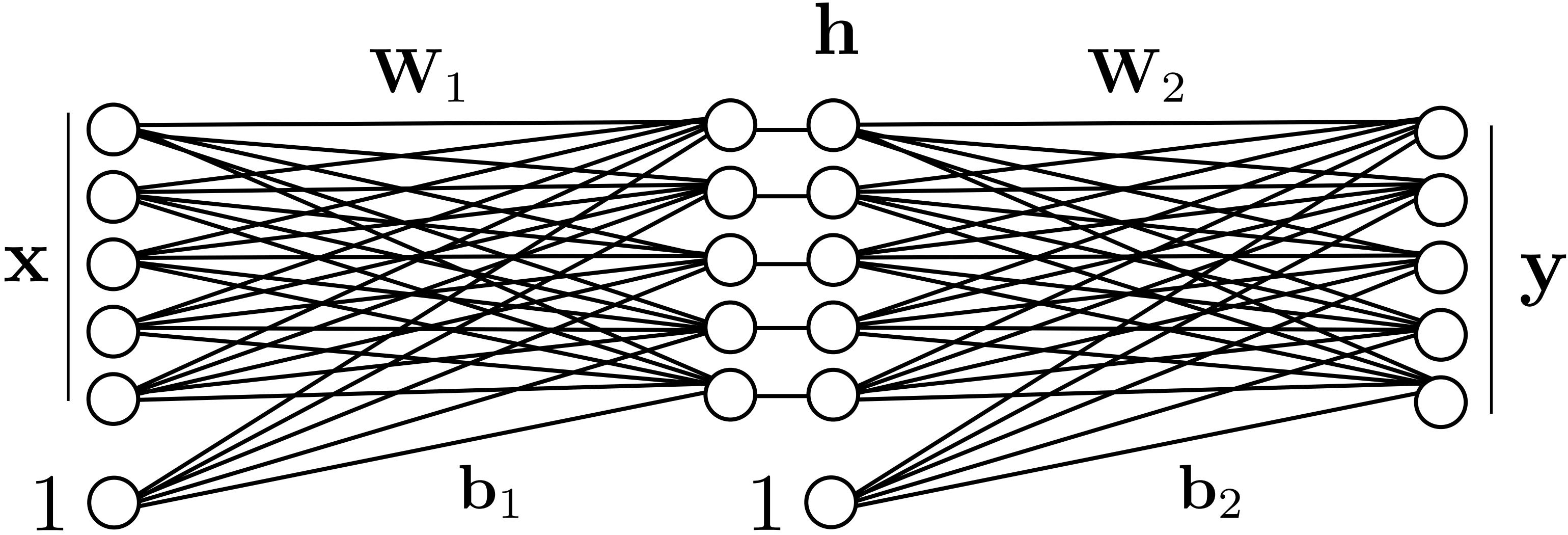
$z, h =$  "hidden units"

# Stacking layers

Input  
representation

Intermediate  
representation

Output  
representation



$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

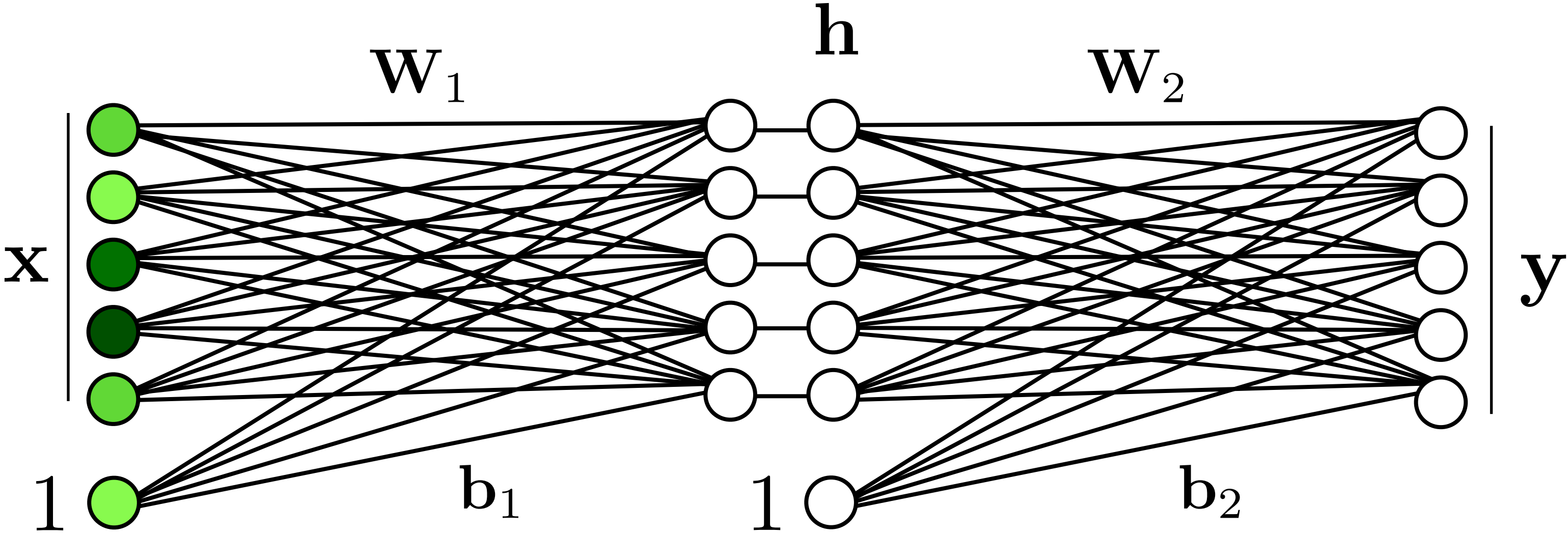
$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers

Input  
representation

Intermediate  
representation

Output  
representation



$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

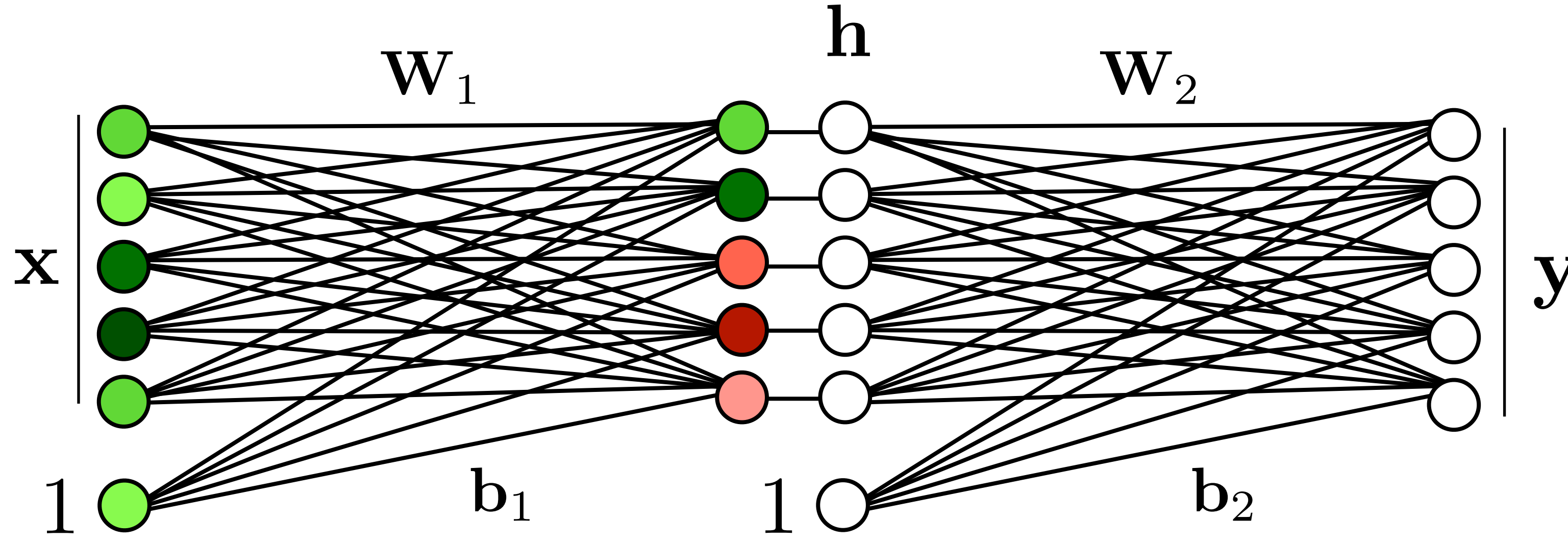


# Stacking layers

Input representation

Intermediate representation

Output representation



$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

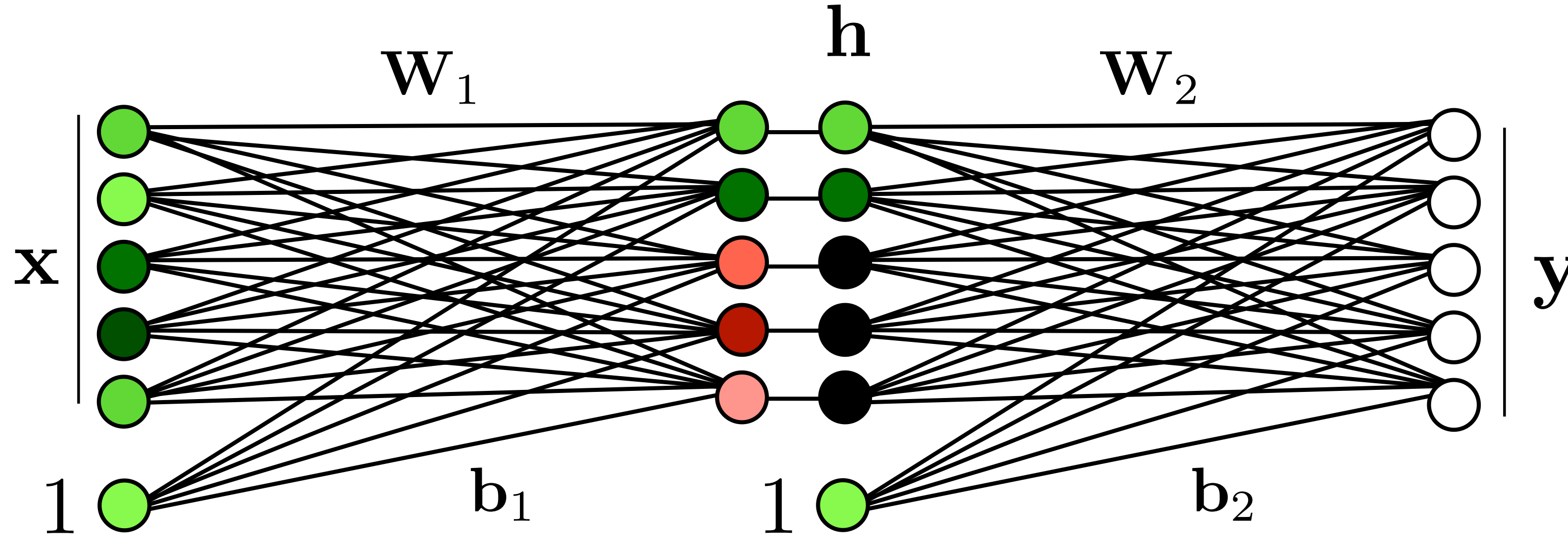
$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers

Input representation

Intermediate representation

Output representation



positive

negative

$$\mathbf{h} = g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{y} = g(\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

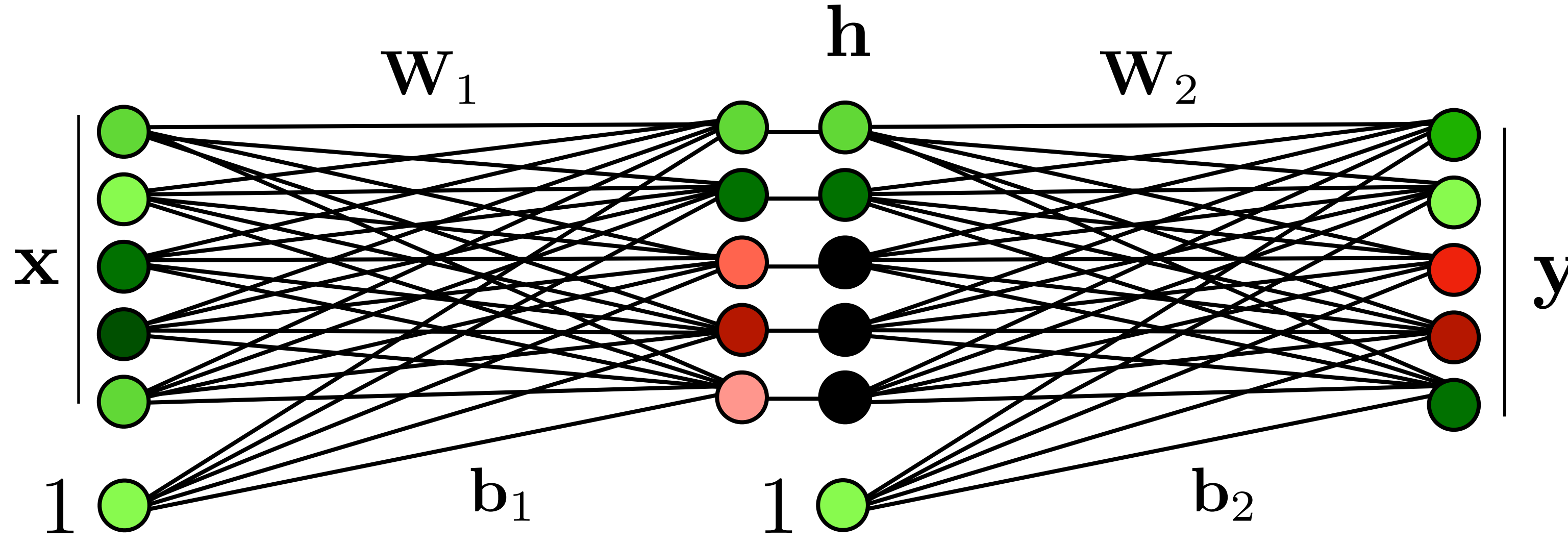
$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Stacking layers

Input representation

Intermediate representation

Output representation



positive

negative

$$\mathbf{h} = g(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)$$

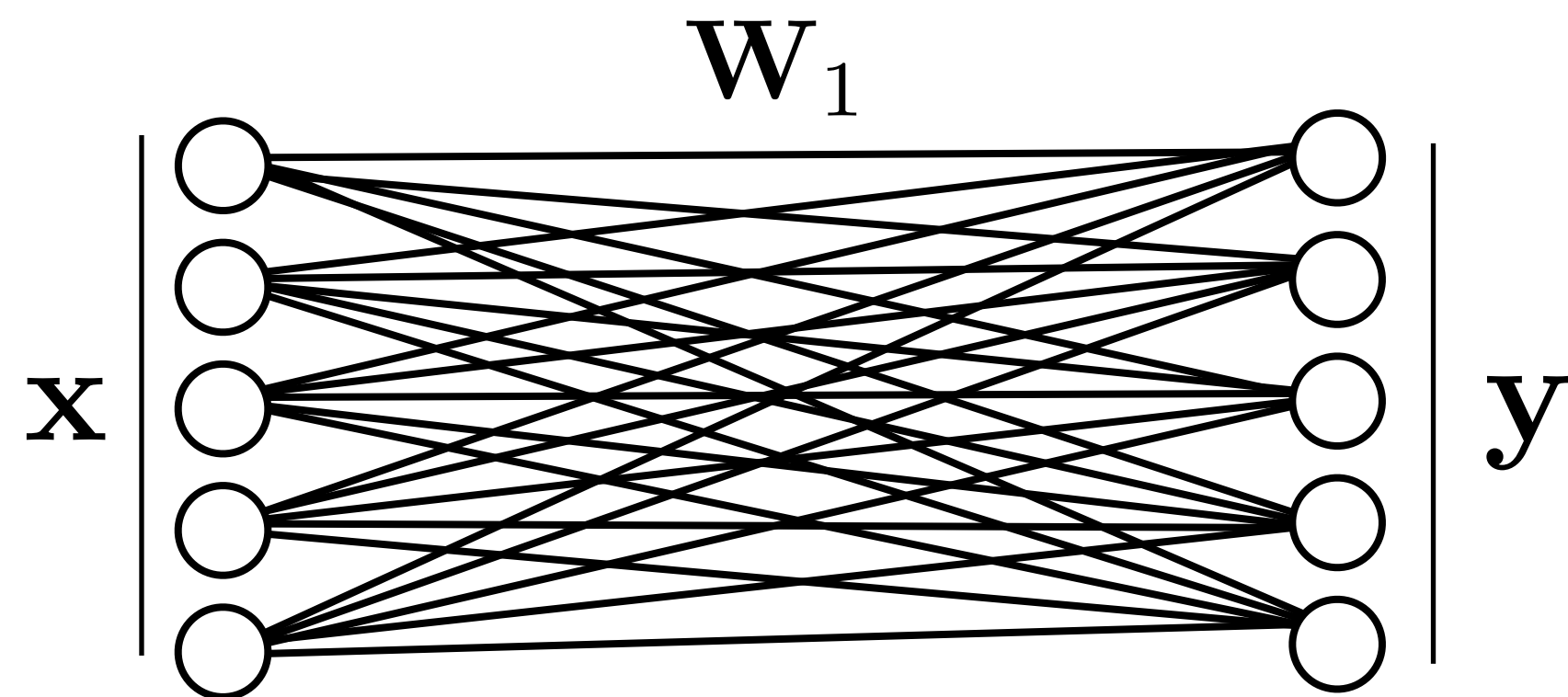
$$\mathbf{y} = g(\mathbf{W}_2\mathbf{h} + \mathbf{b}_2)$$

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

# Connectivity patterns

Input  
representation

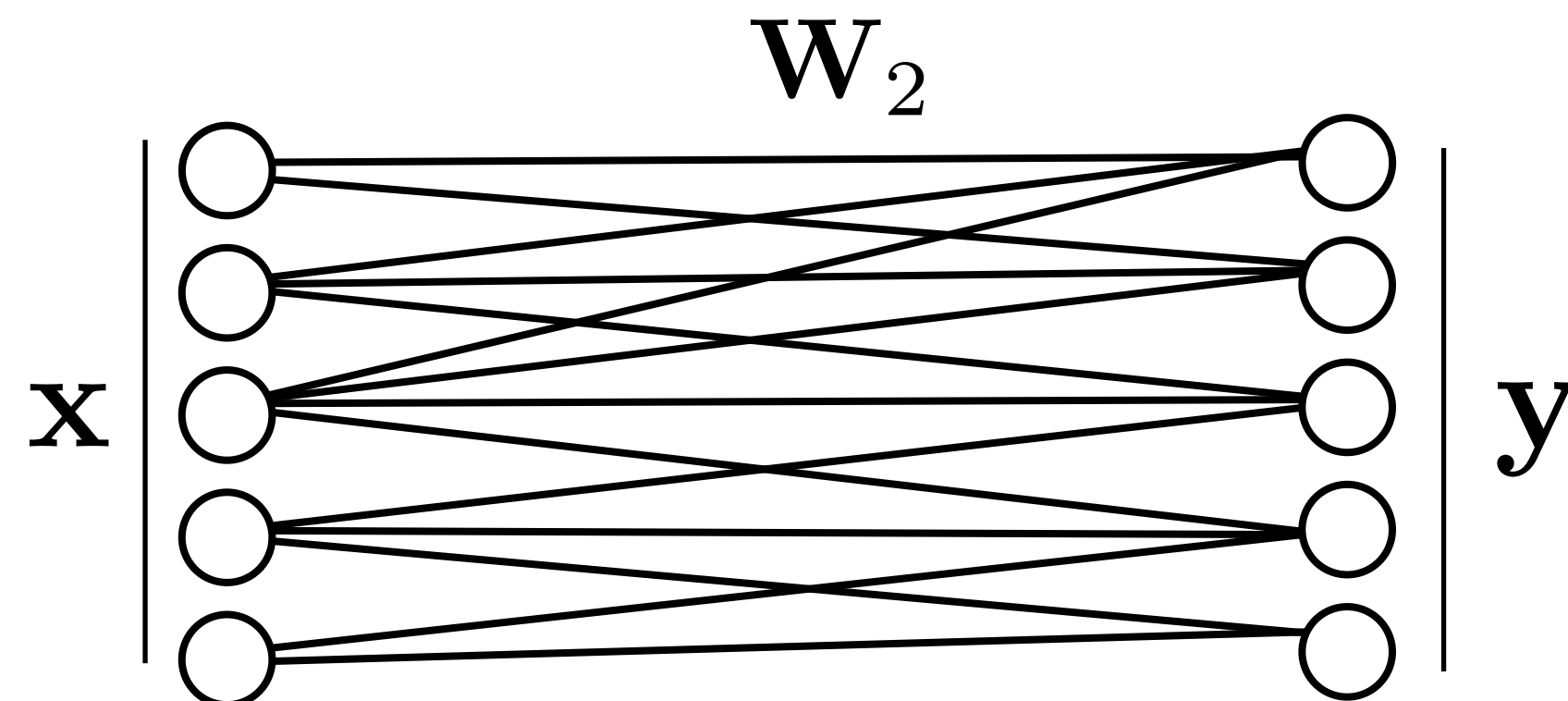
Output  
representation



*Fully connected layer*

Input  
representation

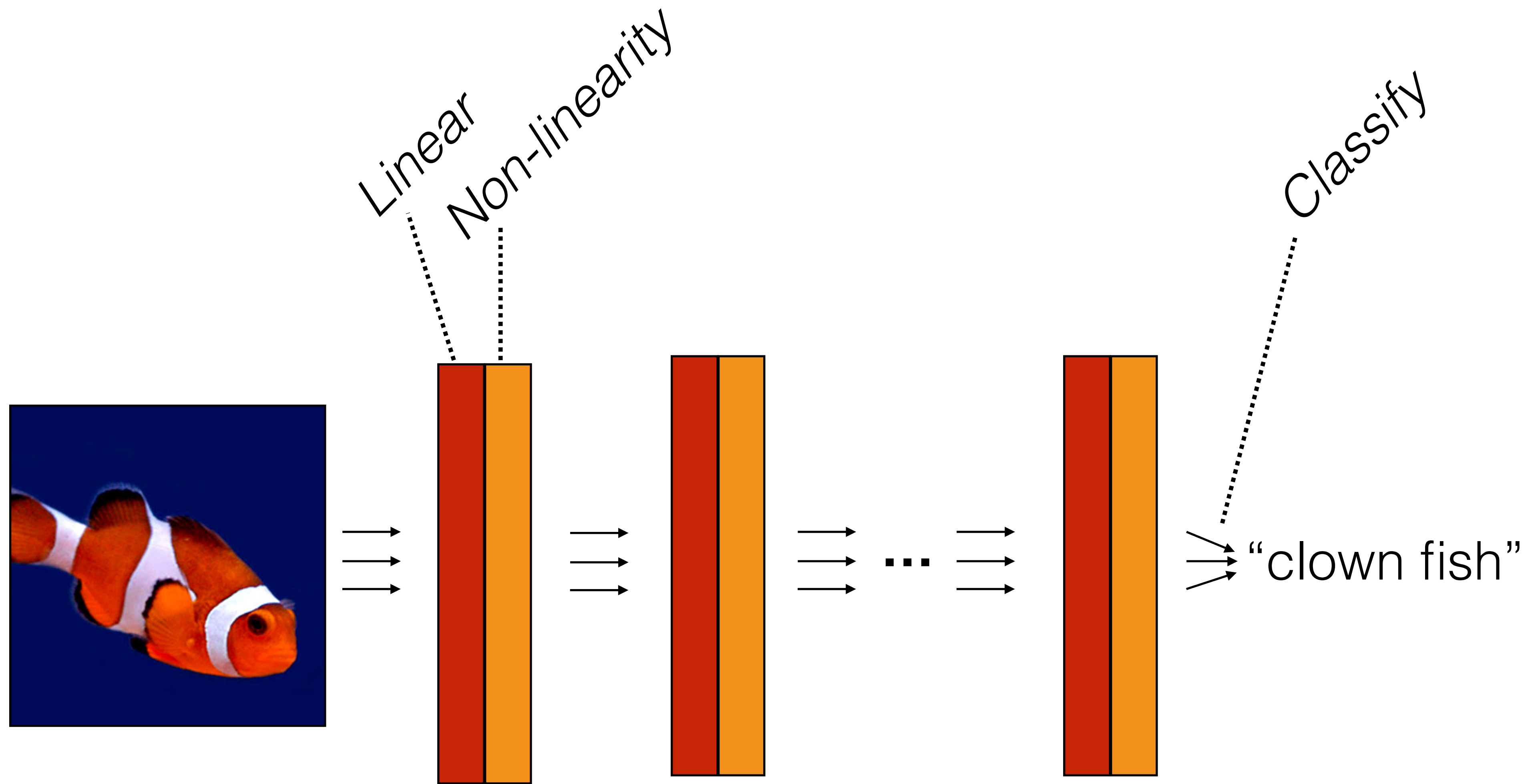
Output  
representation



*Locally connected layer  
(Sparse  $W$ )*

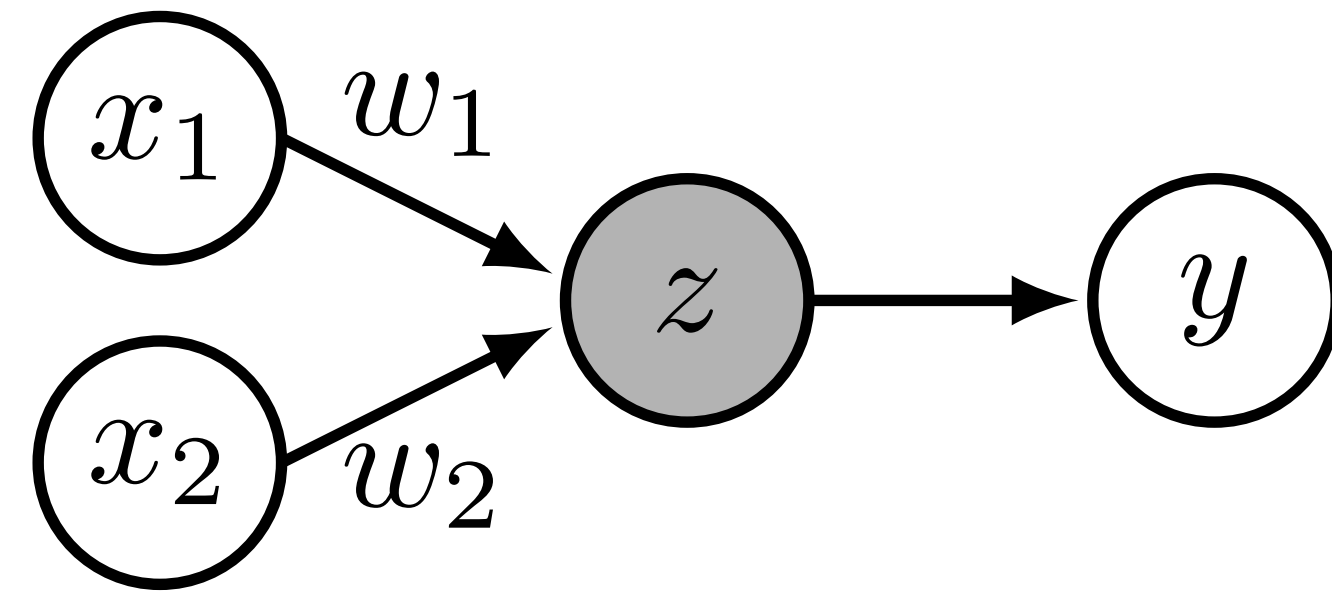


# Deep nets



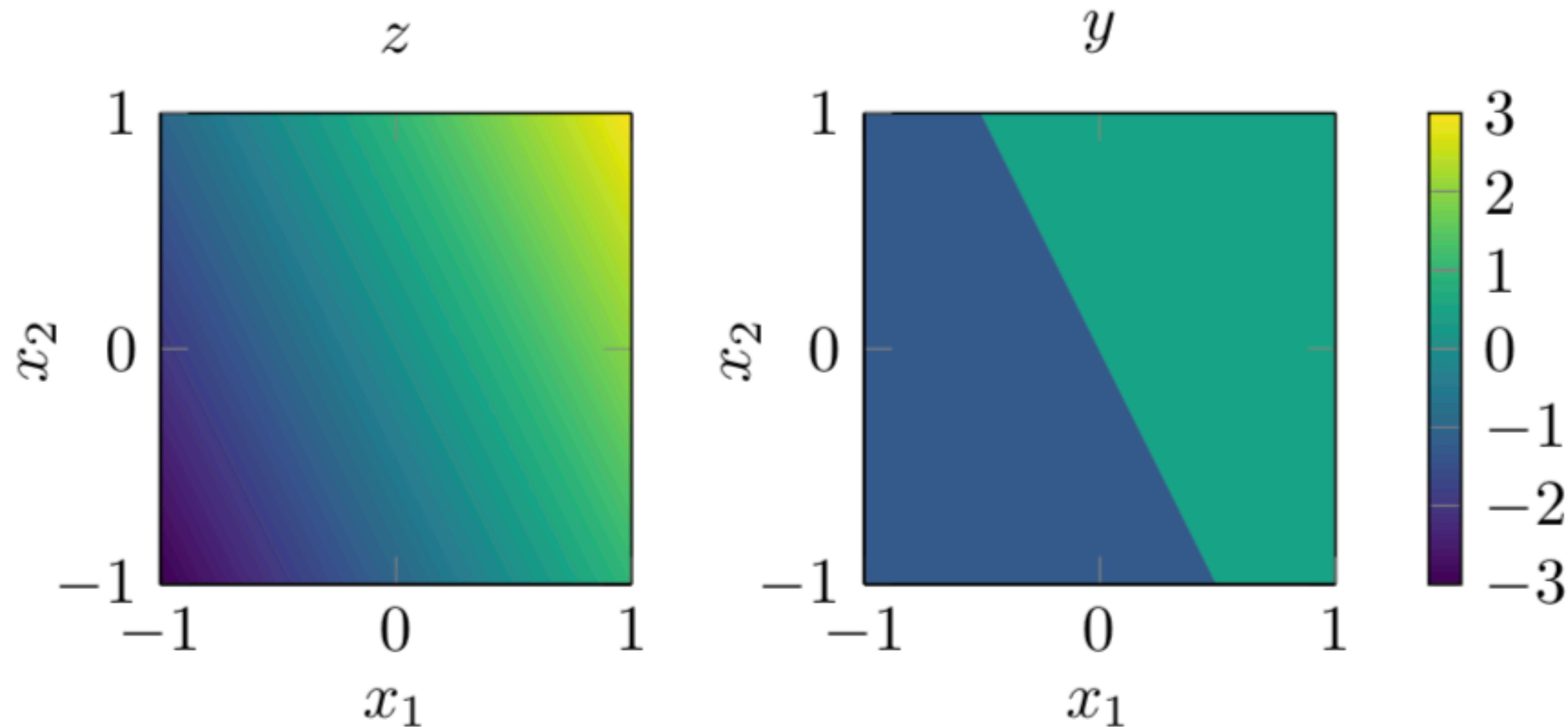
$$f(\mathbf{x}) = f_L(f_{L-1}(\dots f_2(f_1(\mathbf{x}))))$$

# Example: linear classification with a perceptron



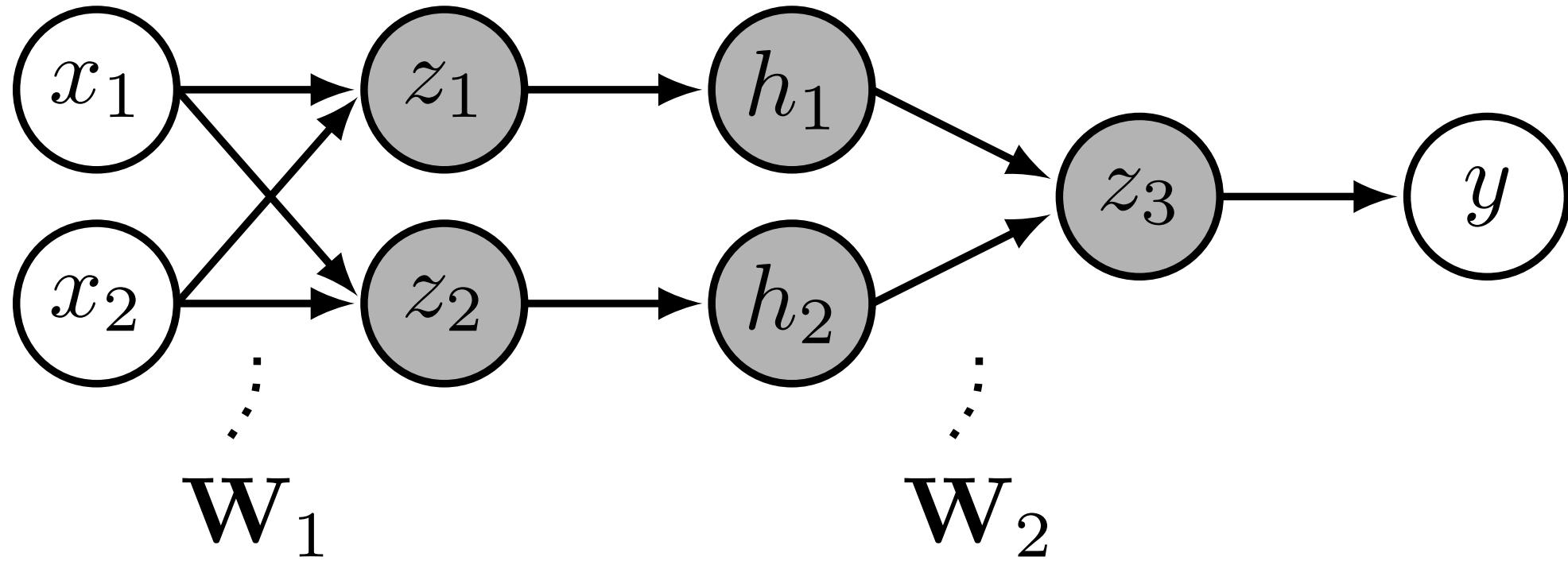
$$z = \mathbf{x}^T \mathbf{w} + b$$

$$y = g(z)$$

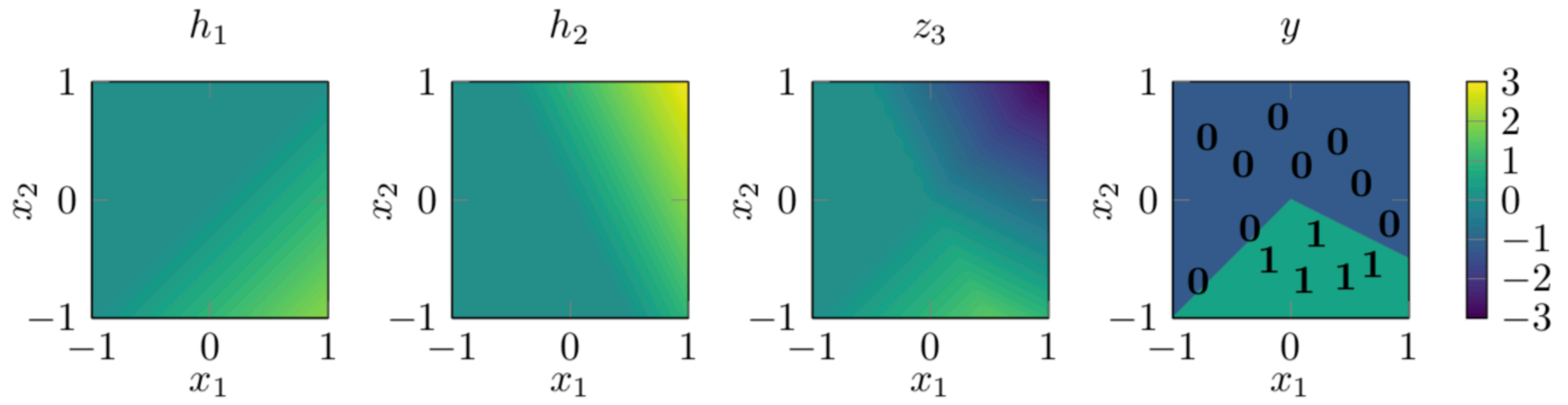


One layer neural net (perceptron) can perform linear classification!

# Example: nonlinear classification with a deep net net



$$\mathbf{z} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$
$$\mathbf{h} = g(\mathbf{z})$$
$$z_3 = \mathbf{W}_2 \mathbf{h} + b_2$$
$$y = 1(z_3 > 0)$$

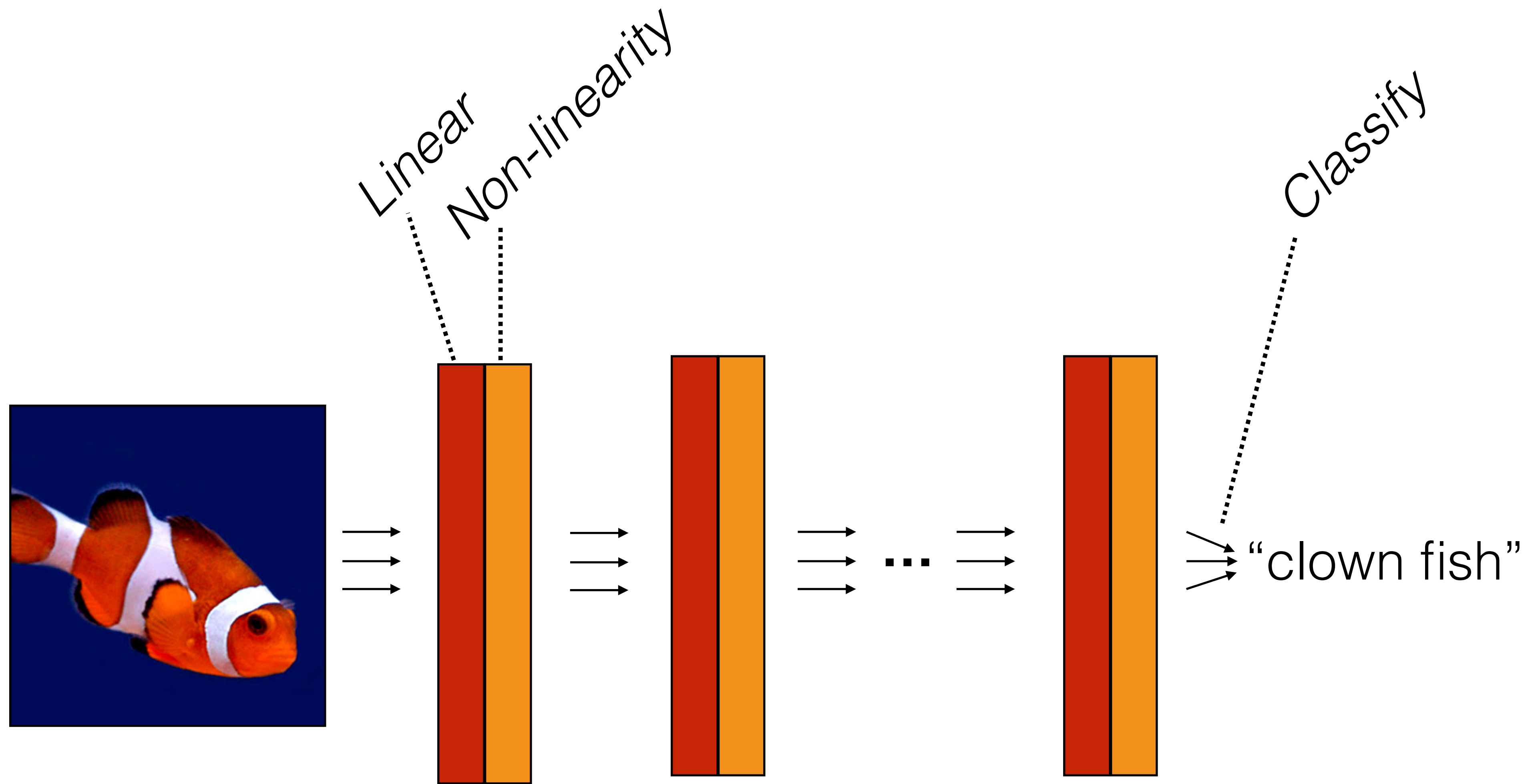


# Representational power

- 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/ftml.pdf>
  - Bengio, Courville, Goodfellow book  
<http://www.deeplearningbook.org/contents/mlp.html>
  - Simple proof by M. Neilsen  
<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.



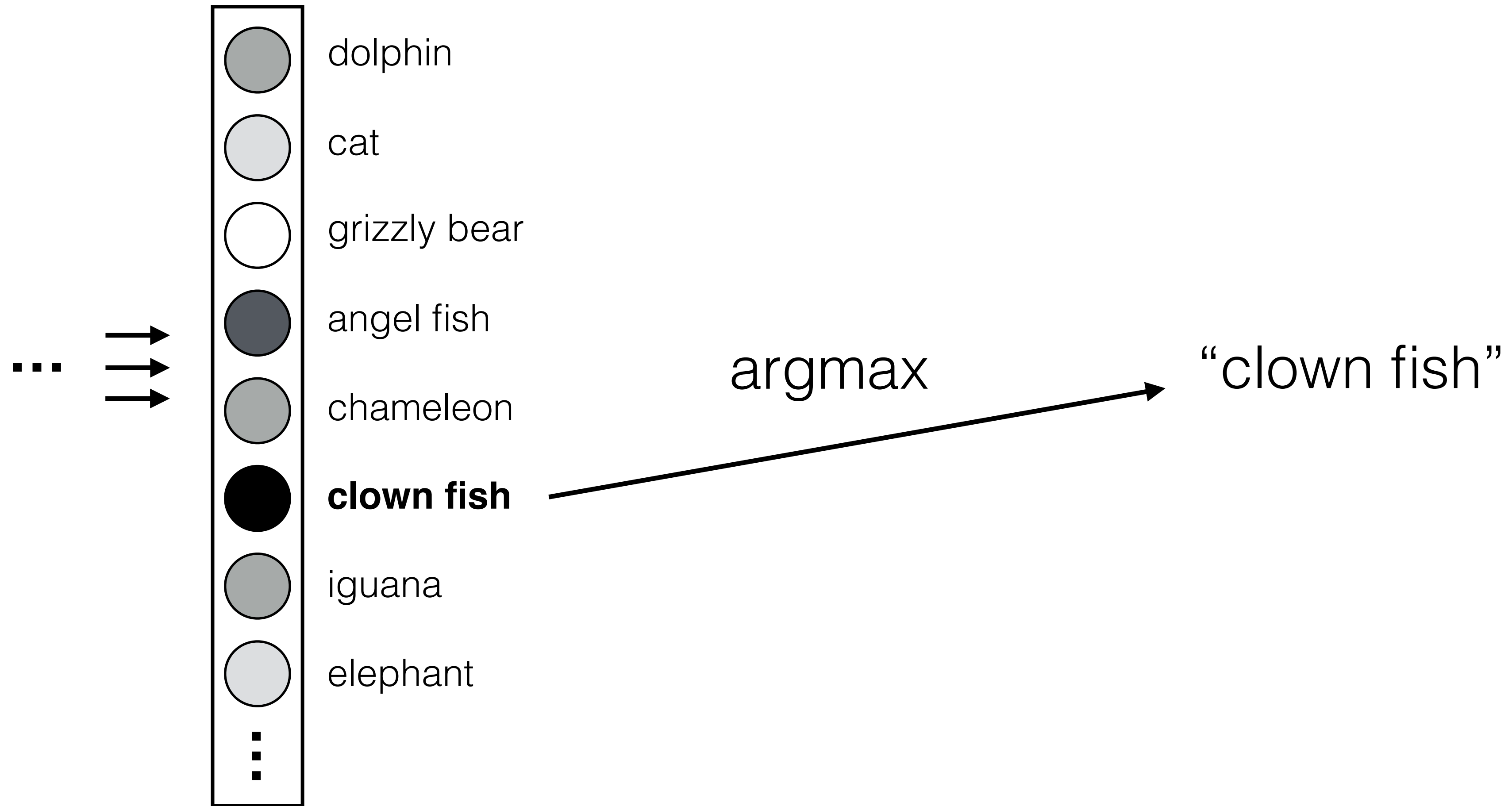
# Deep nets



$$f(\mathbf{x}) = f_L(f_{L-1}(\dots f_2(f_1(\mathbf{x}))))$$

# Classifier layer

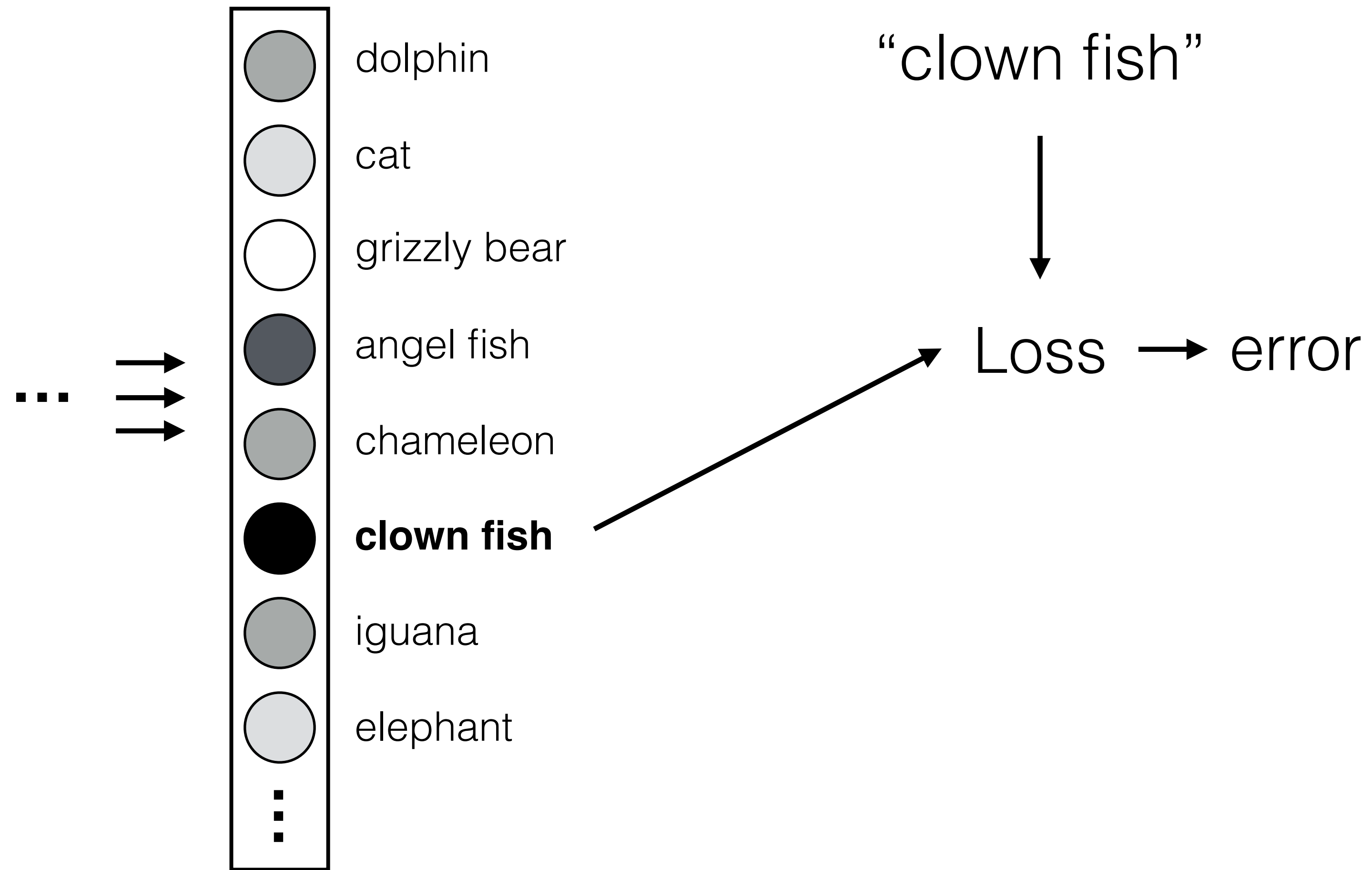
Last layer



# Loss function

Network output

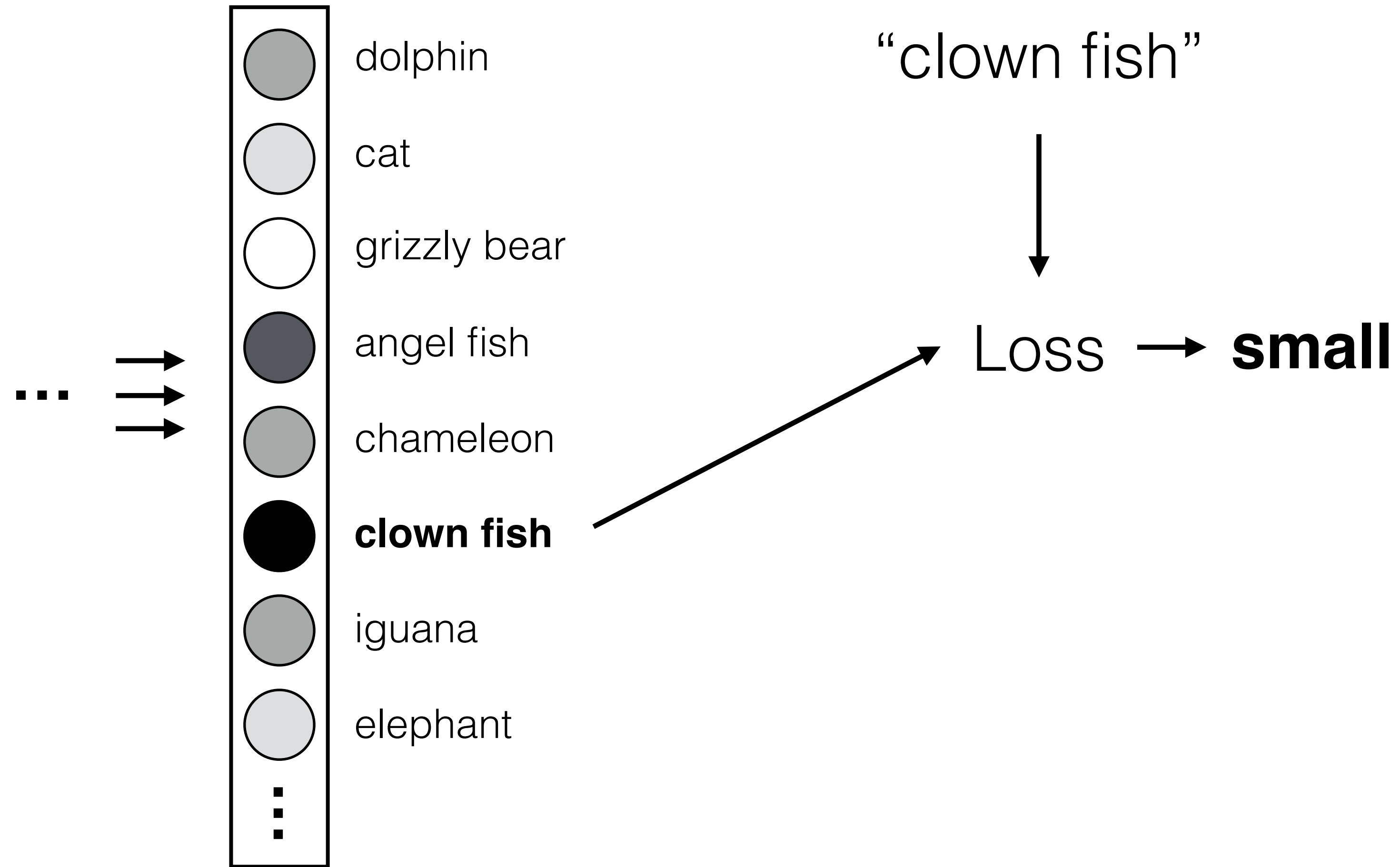
Ground truth label



# Loss function

Network output

Ground truth label





# Loss function

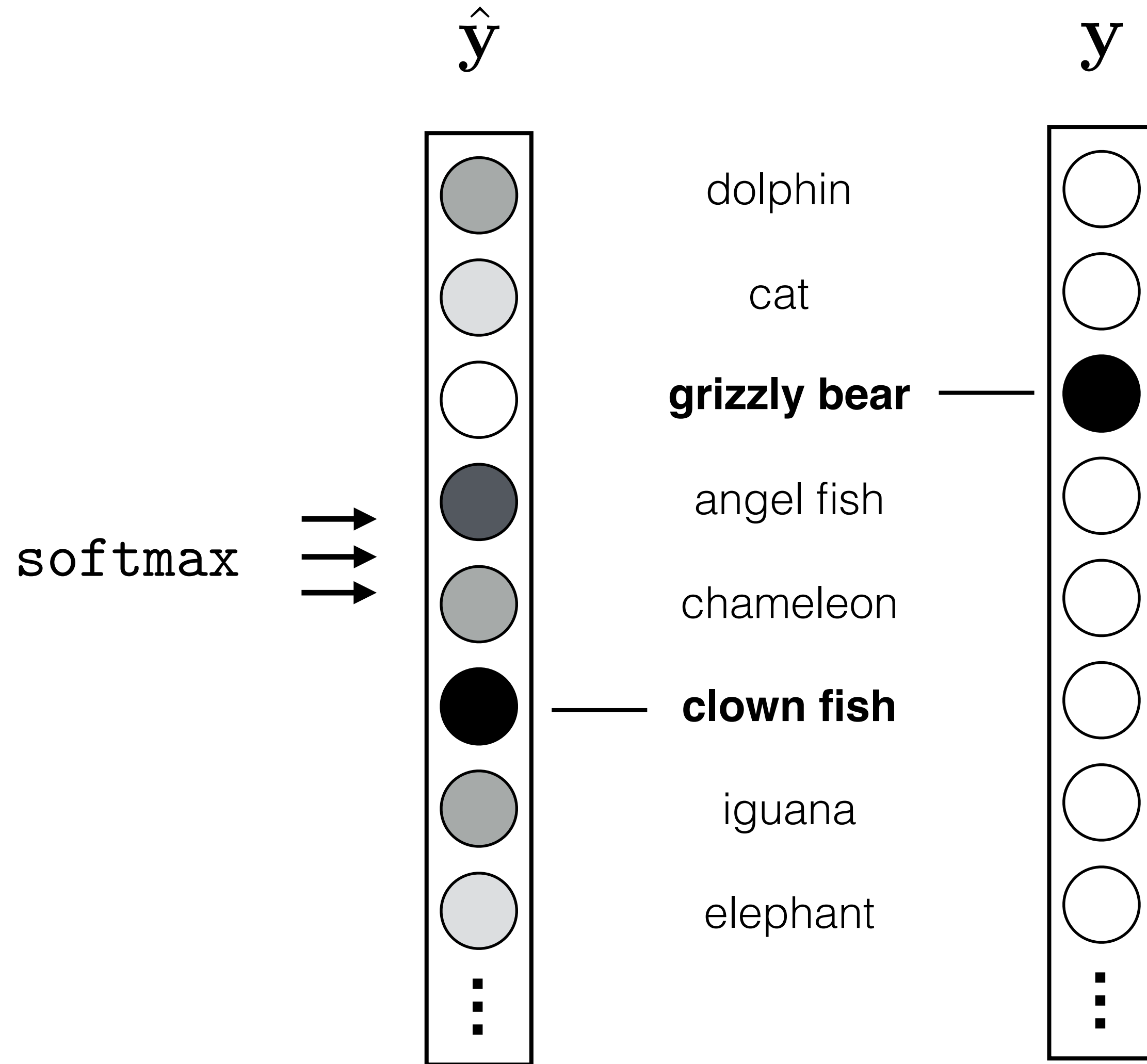
Network output

Ground truth label



Network output

Ground truth label

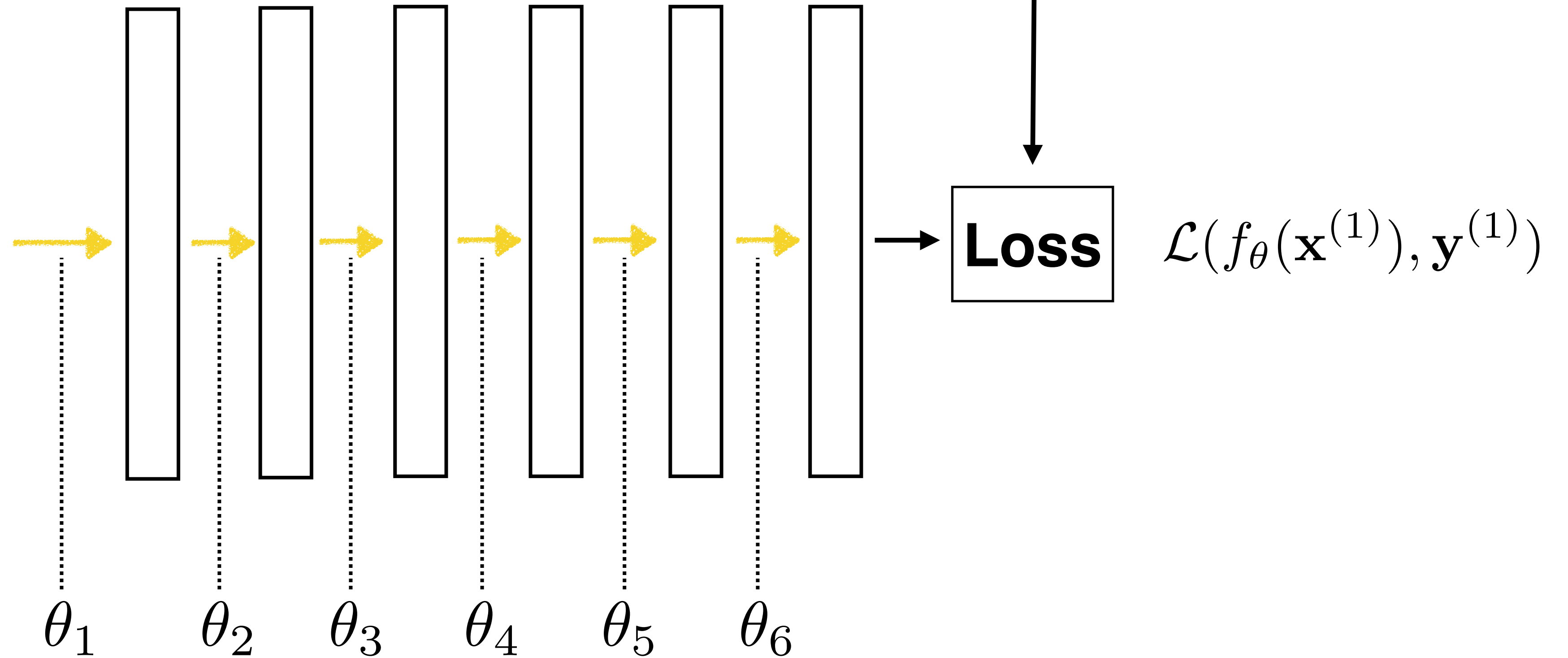
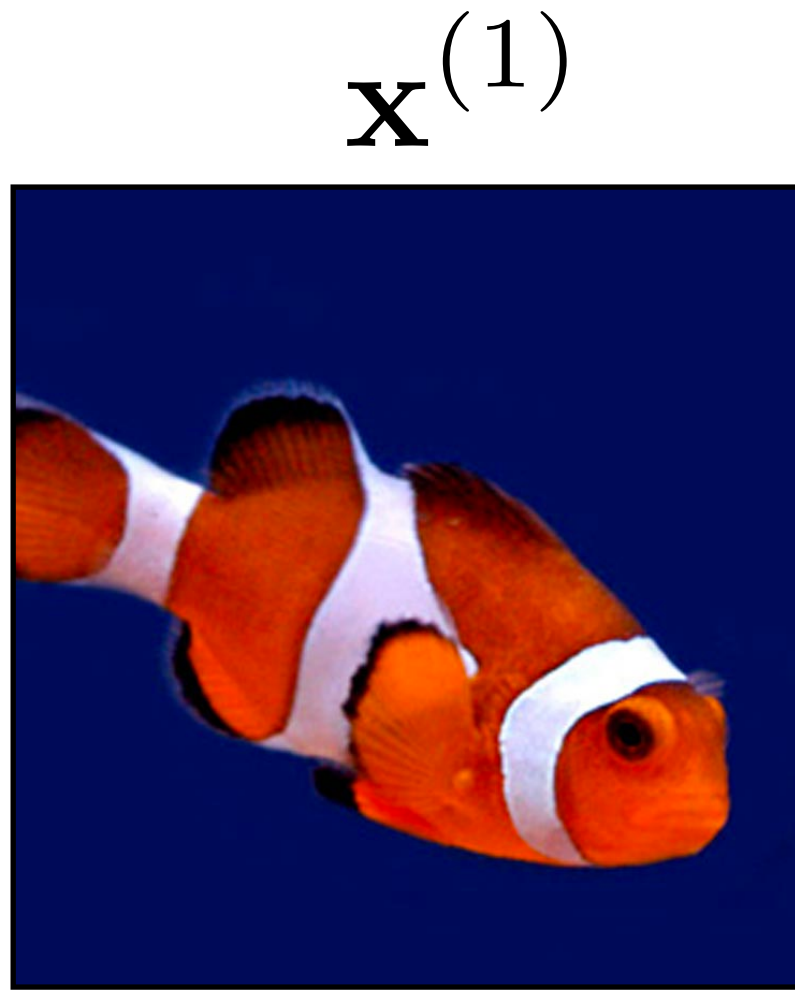


Probability of the observed data under the model

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

# Deep learning

$\mathbf{y}^{(1)}$   
“clown fish”

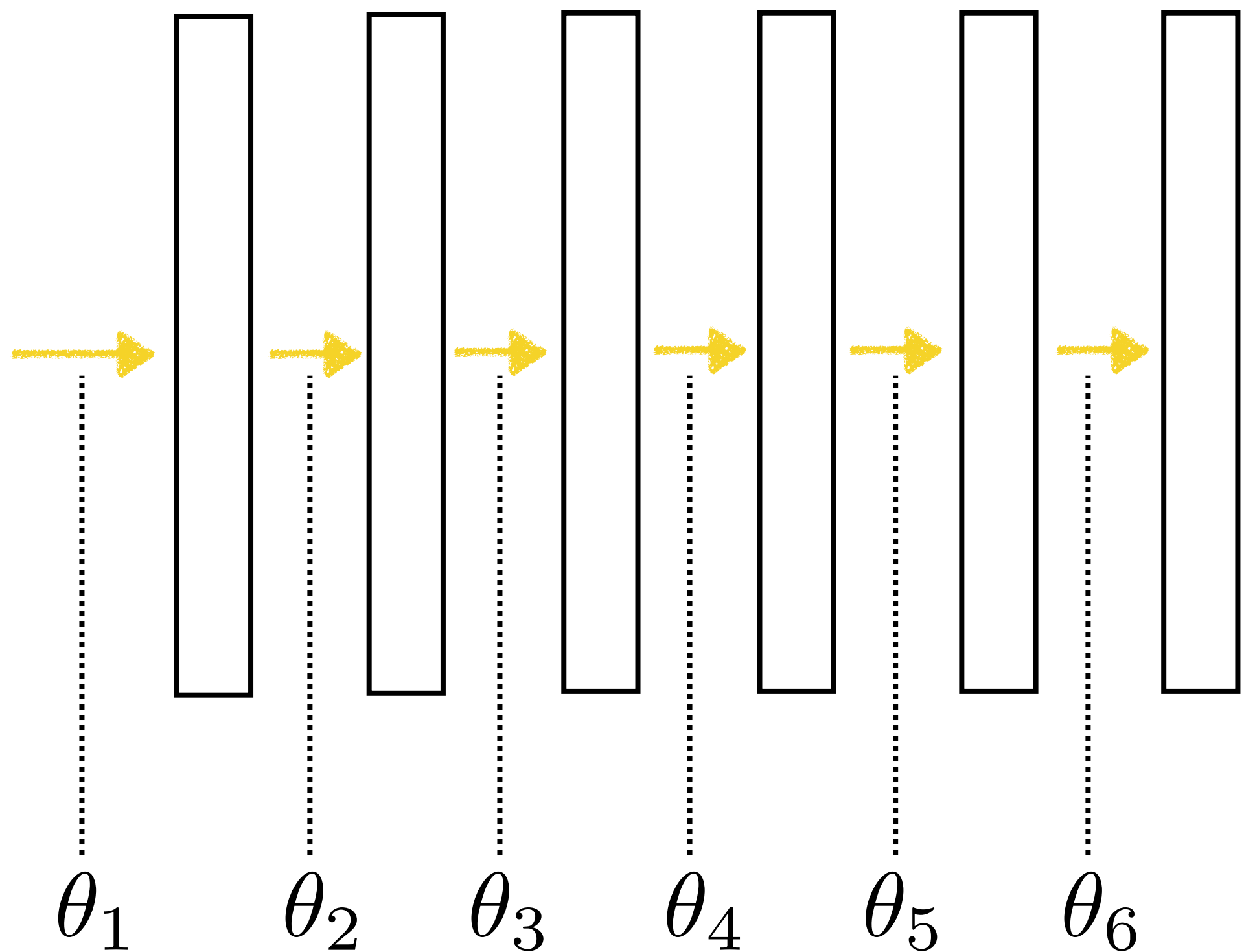


$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

# Deep learning

$\mathbf{y}^{(2)}$   
“grizzly bear”

Learned



Loss

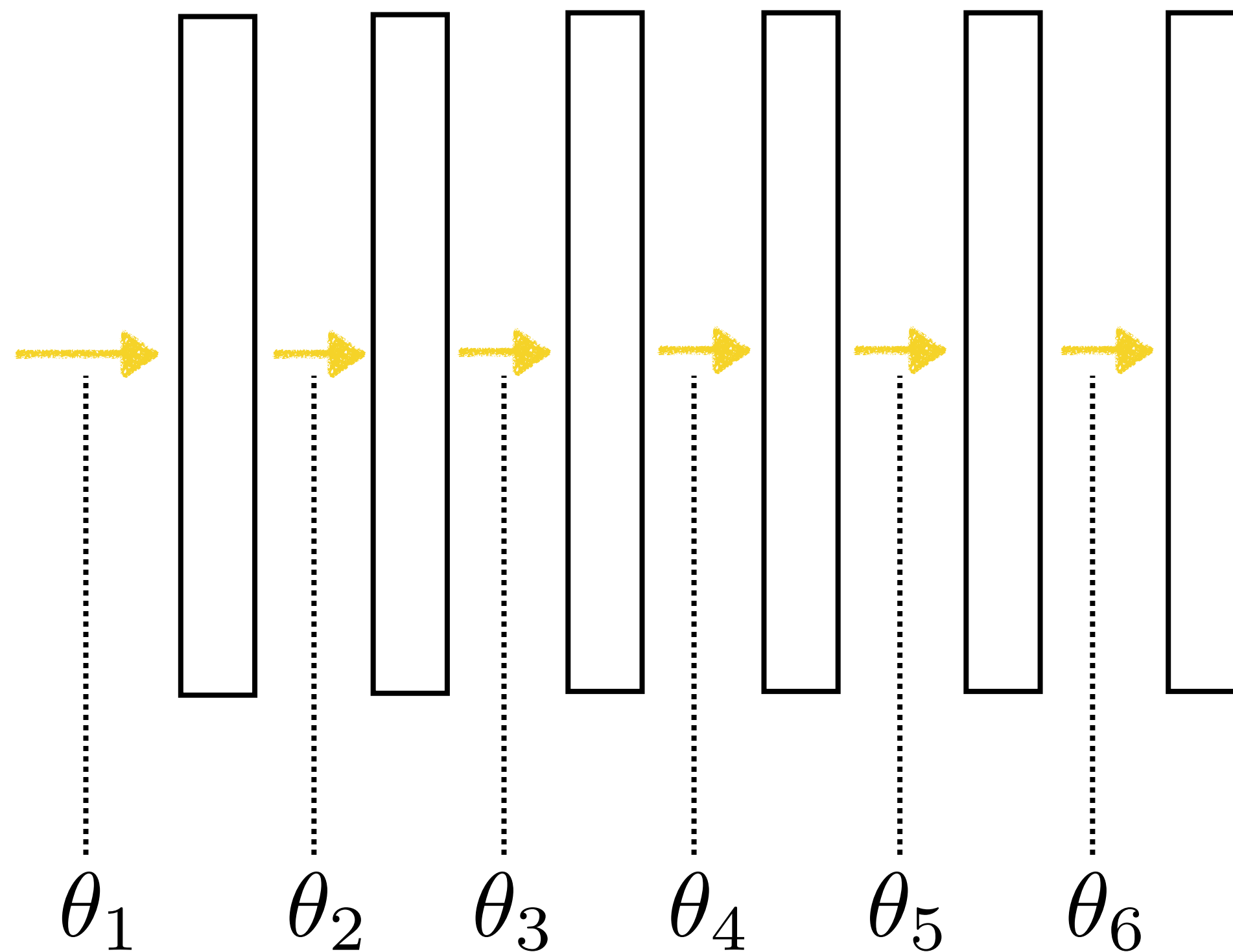
$$\mathcal{L}(f_{\theta}(\mathbf{x}^{(2)}), \mathbf{y}^{(2)})$$

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



# Deep learning

$\mathbf{y}^{(i)}$   
“chameleon”



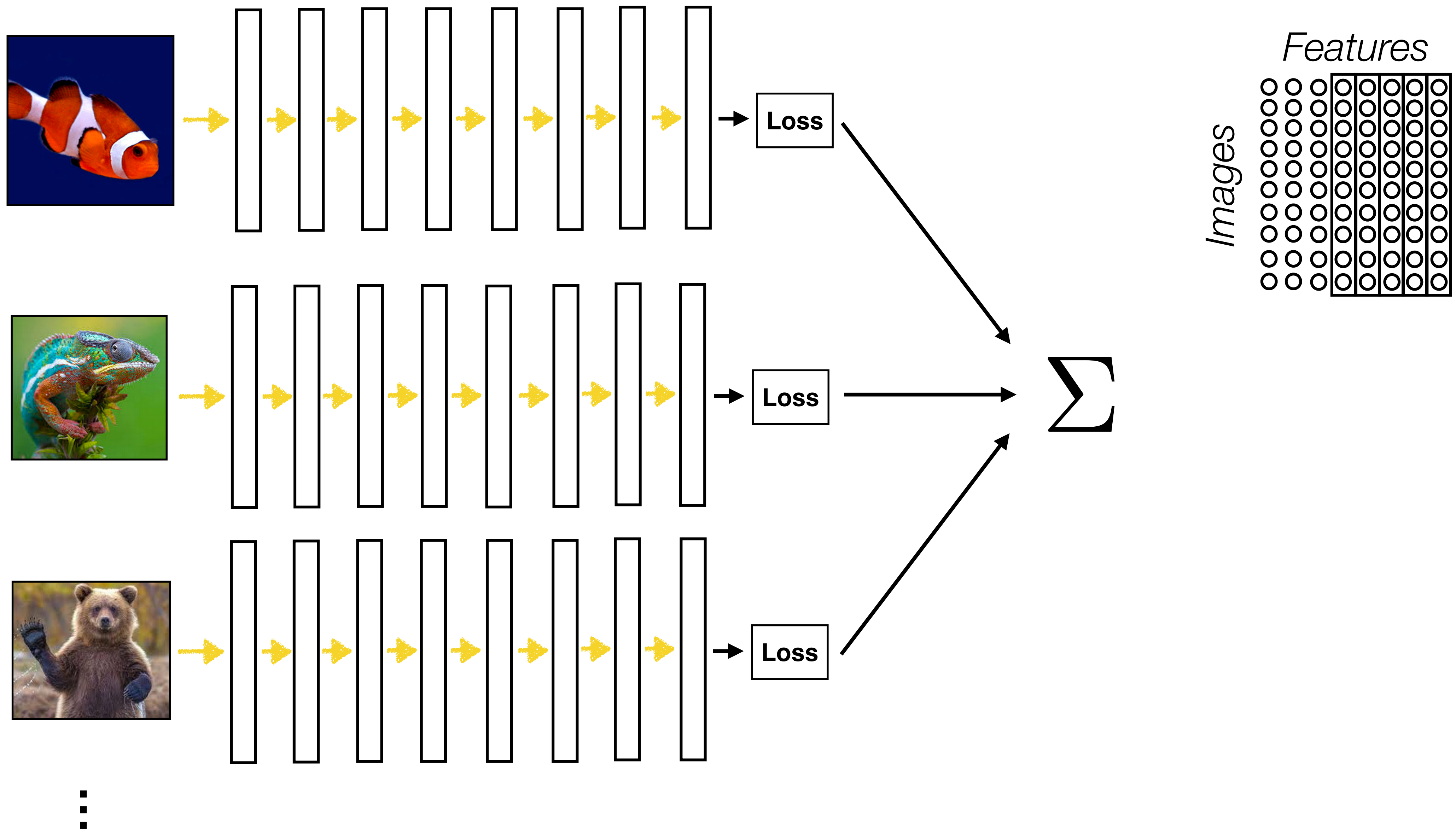
Learned

**Loss**

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

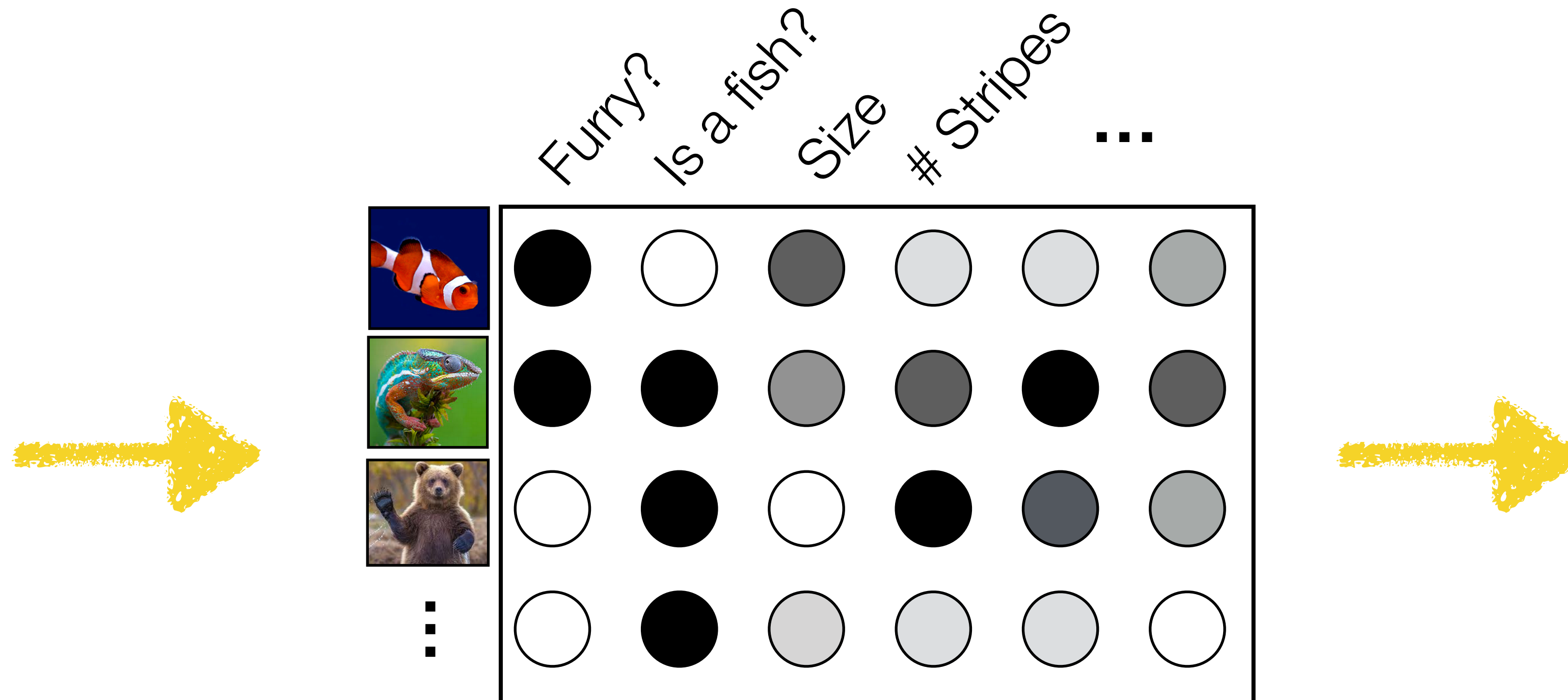
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

# Batch (parallel) processing



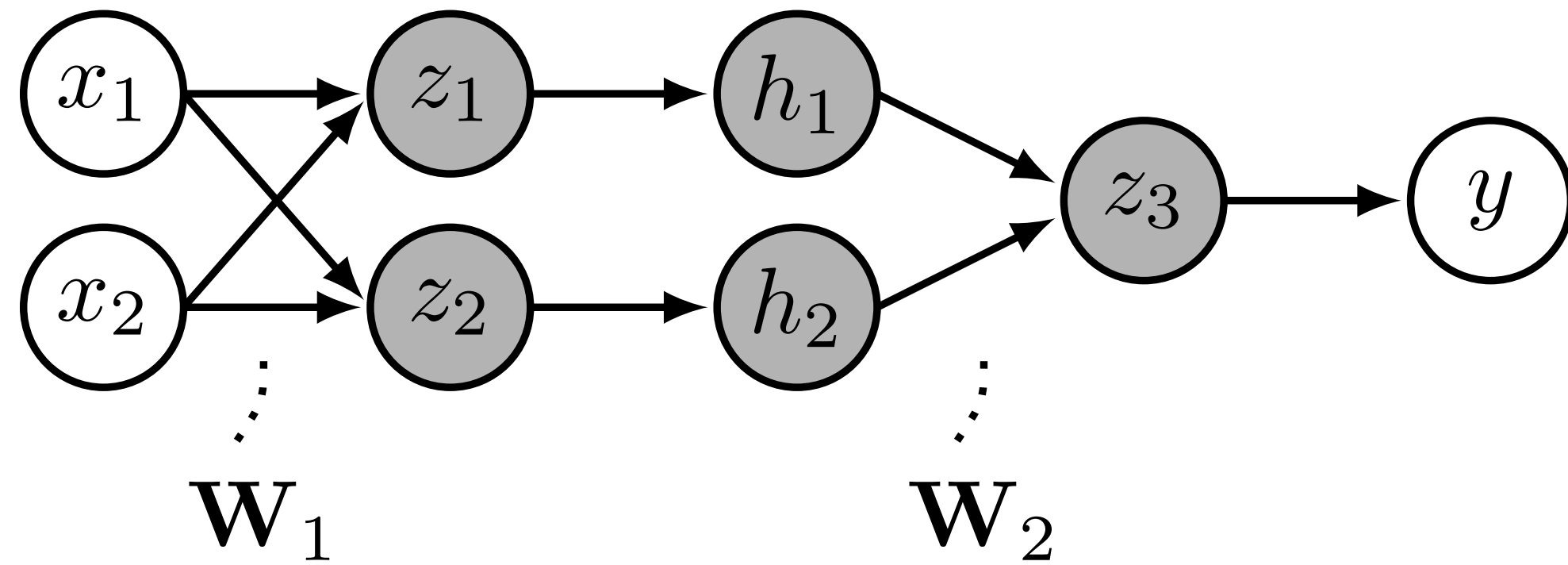
# Tensors

(multi-dimensional arrays)



*Each layer is a representation of the data*

# Everything is a tensor



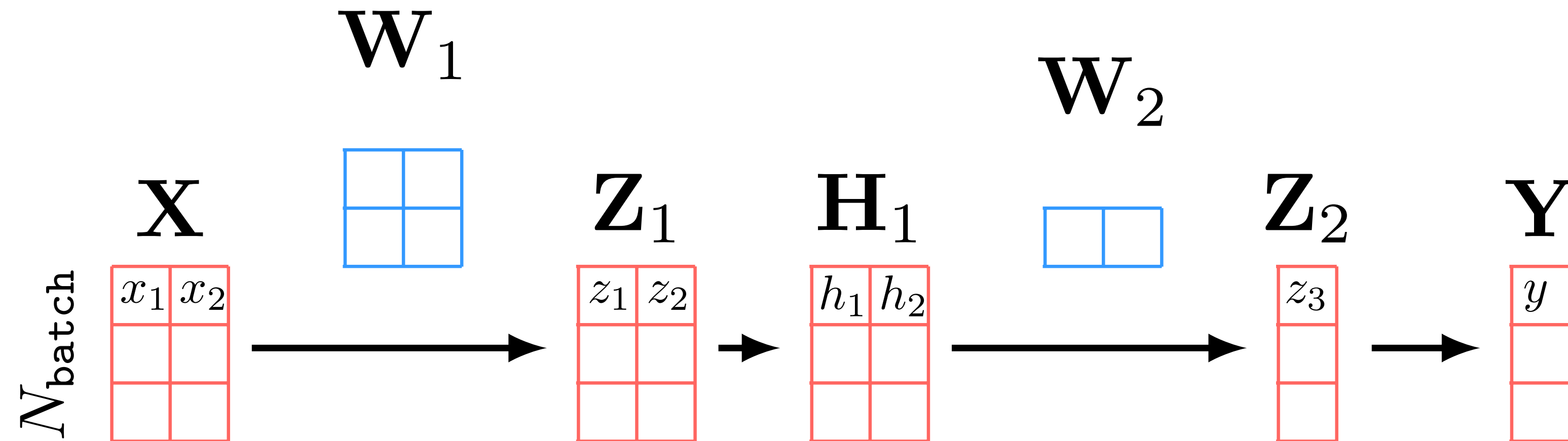
$$\mathbf{z} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

$$\mathbf{h} = g(\mathbf{z})$$

$$z_3 = \mathbf{W}_2 \mathbf{h} + b_2$$

$$y = 1(z_3 > 0)$$

Tensor processing with batch size = 3:





# 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 prevent the network from overfitting?

1. Fewer neurons, fewer layers
2. Weight decay and other regularizers
3. Normalization layers
4. ...

# Recall: regularized least squares

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

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

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

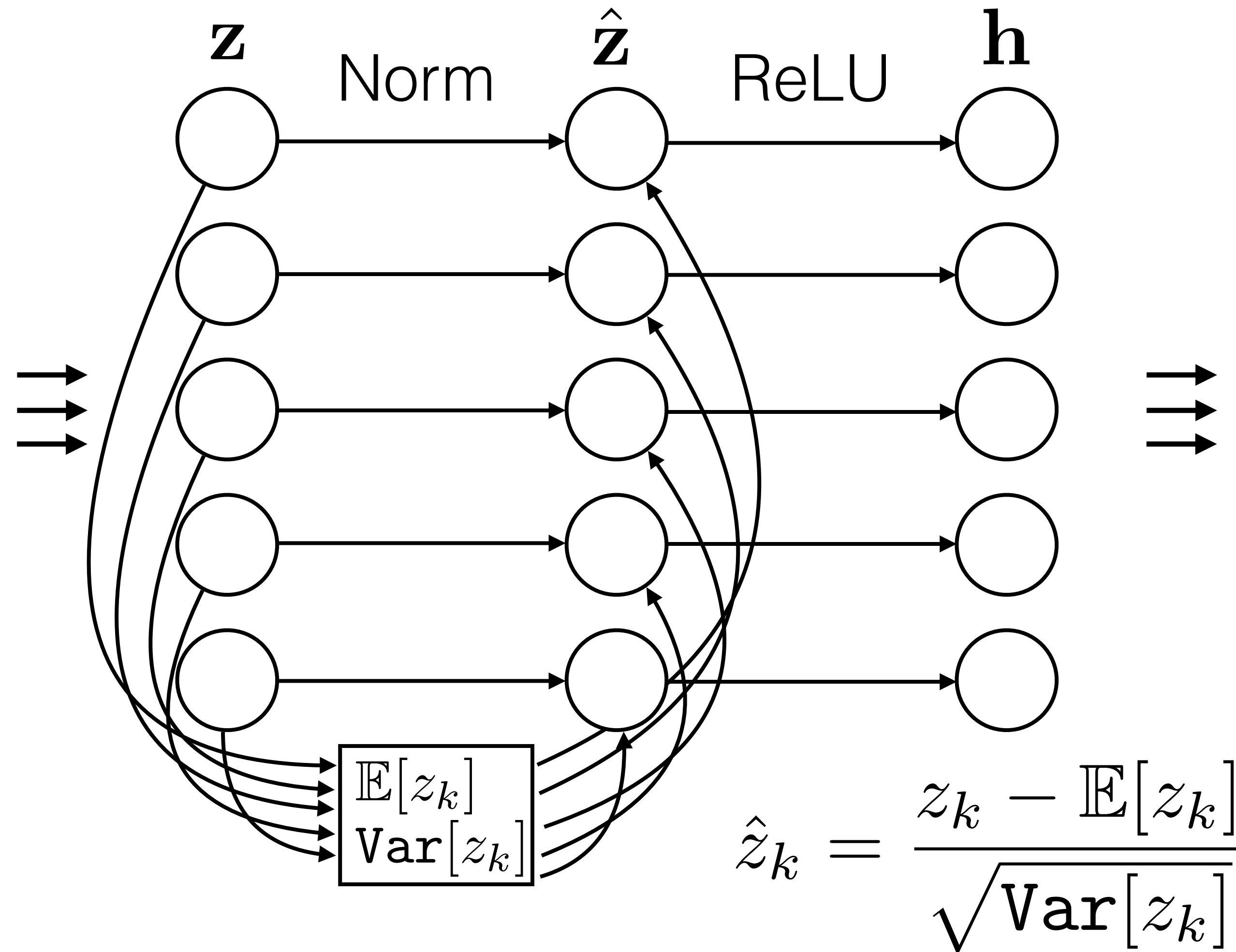
# Regularizing the weights in a neural net

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

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

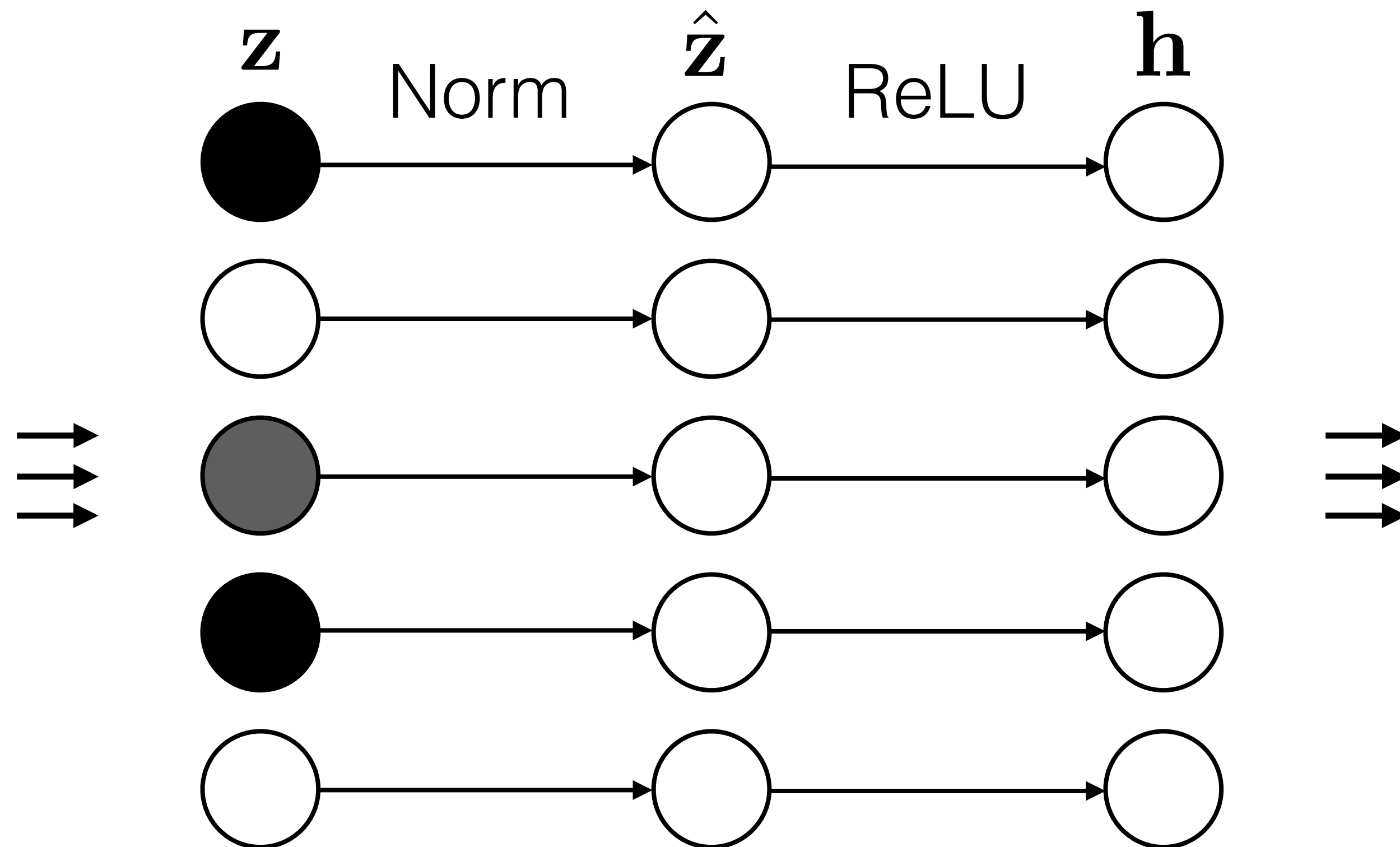
“We prefer to keep weights small.”

# Normalization layers



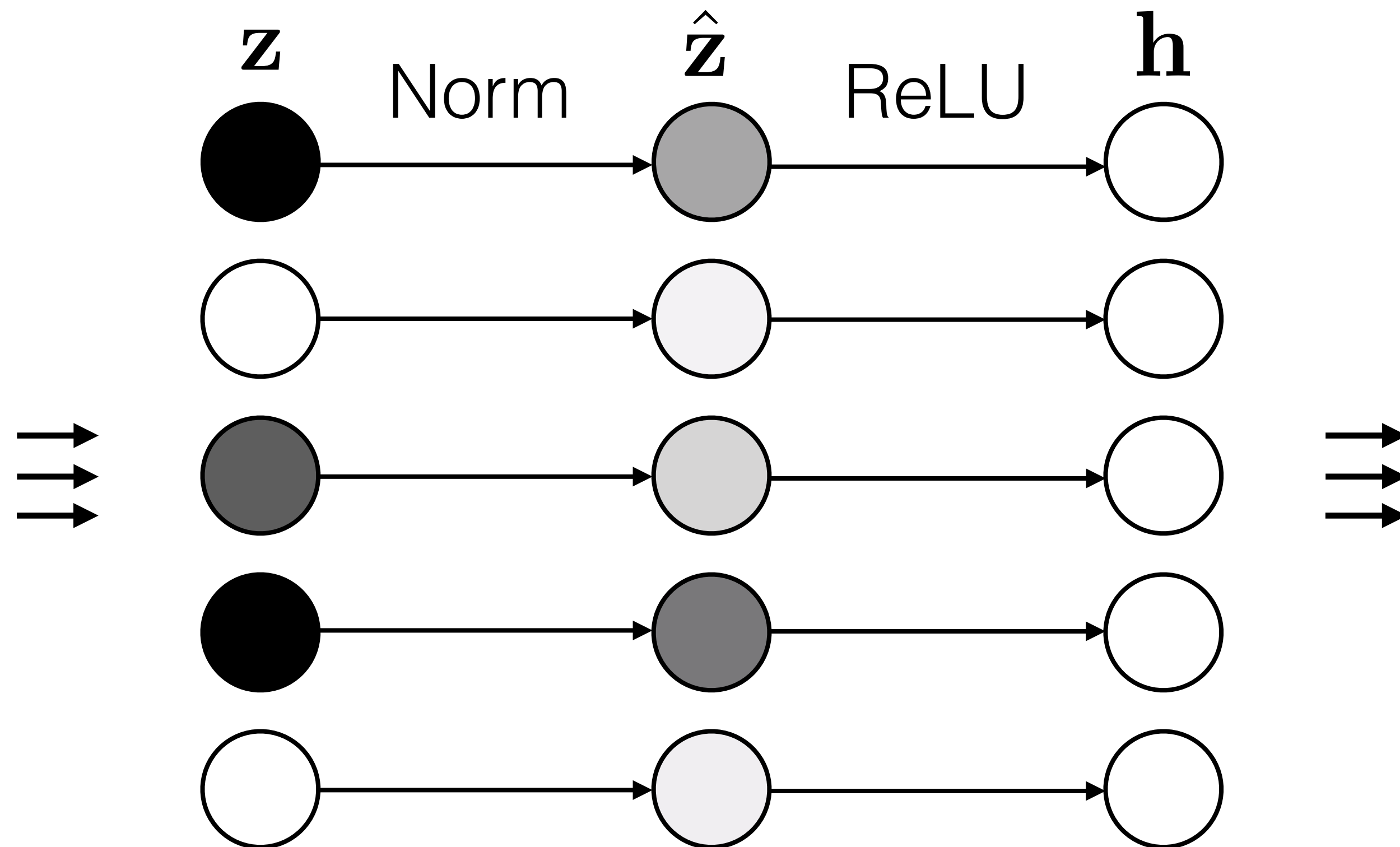


# Normalization layers



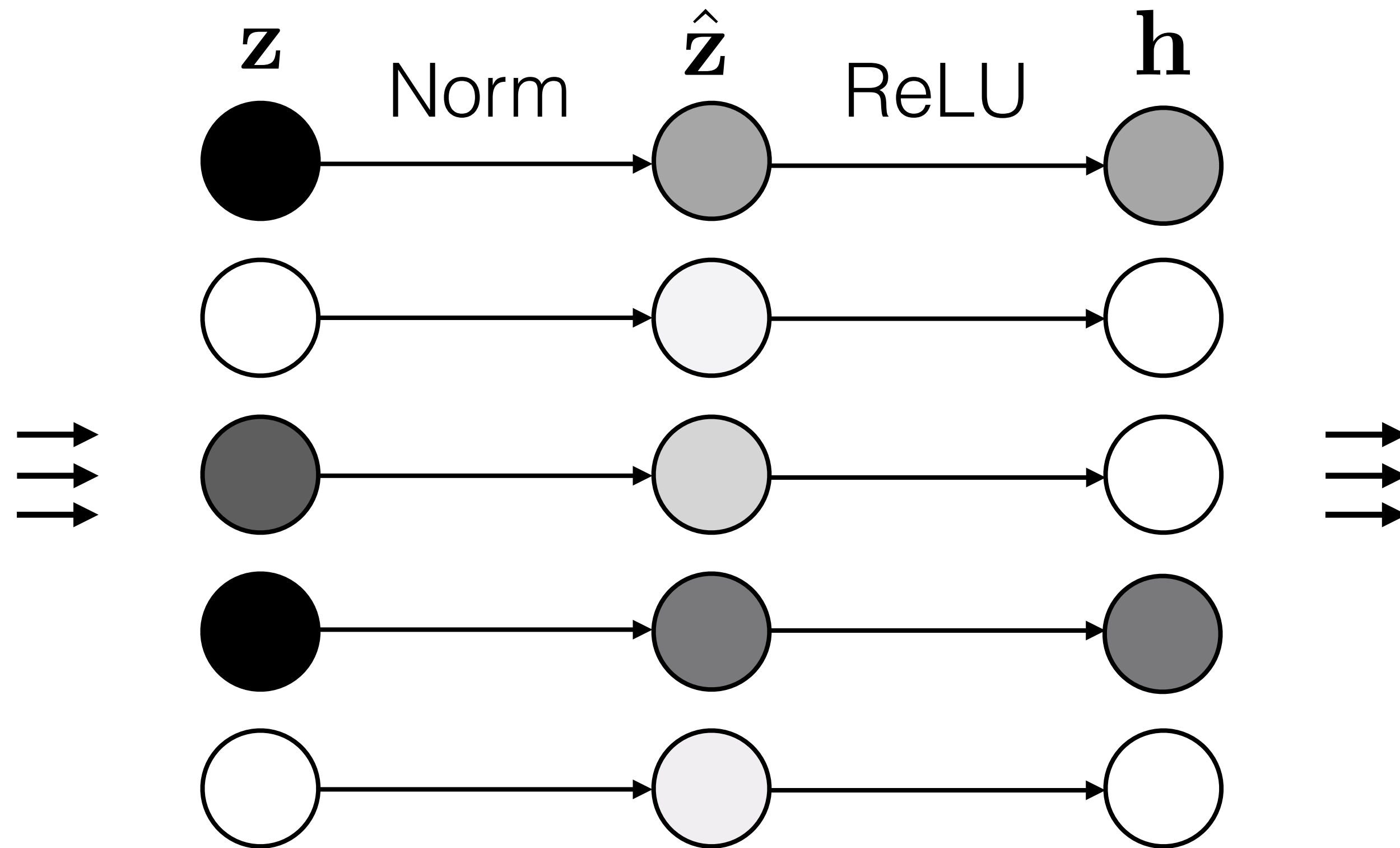
$$\hat{z}_k = \frac{z_k - \mathbb{E}[z_k]}{\sqrt{\text{Var}[z_k]}}$$

# Normalization layers



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

# Normalization layers



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

# 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.

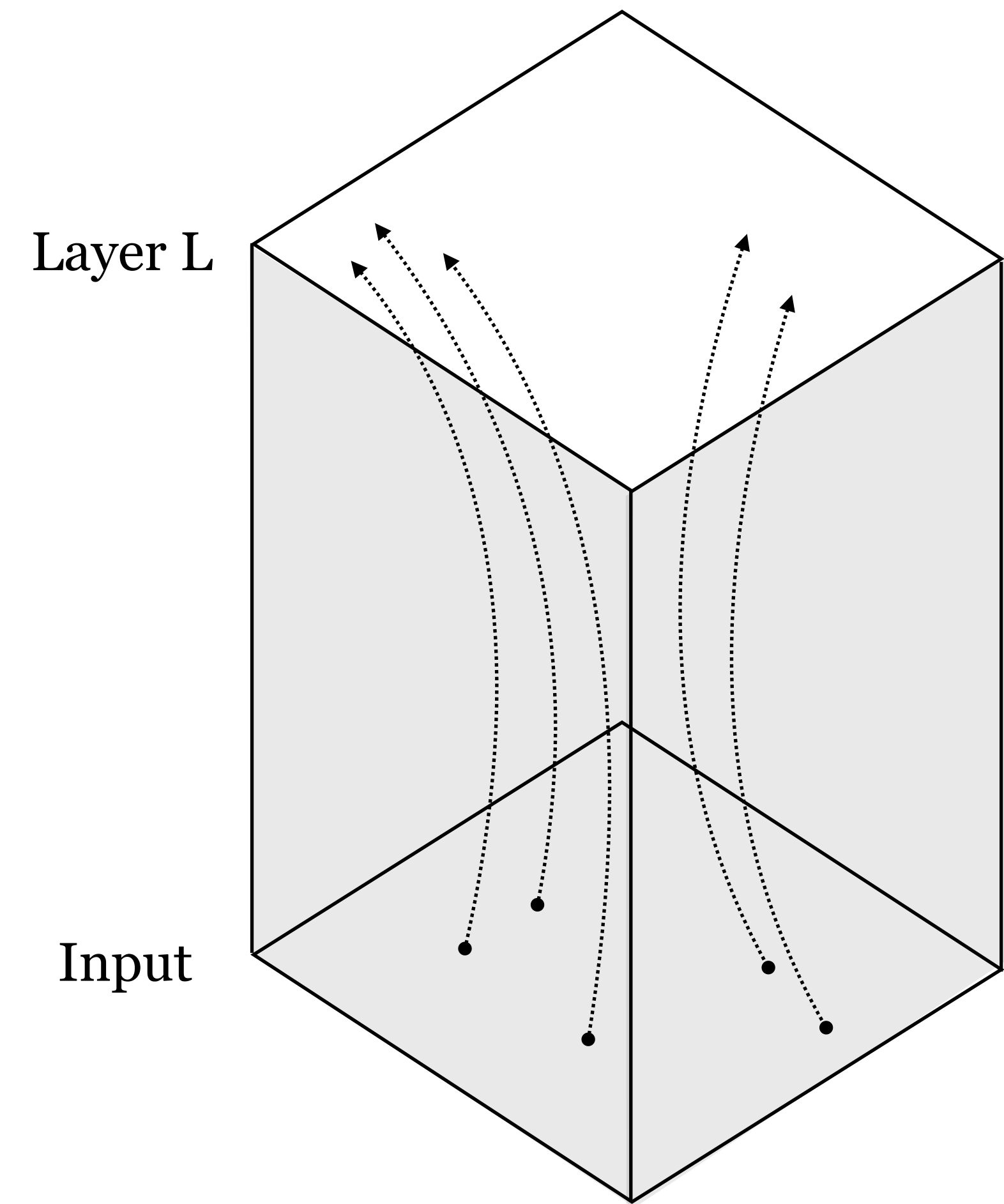
Squashes units into a **standard range**, avoiding overflow.

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.

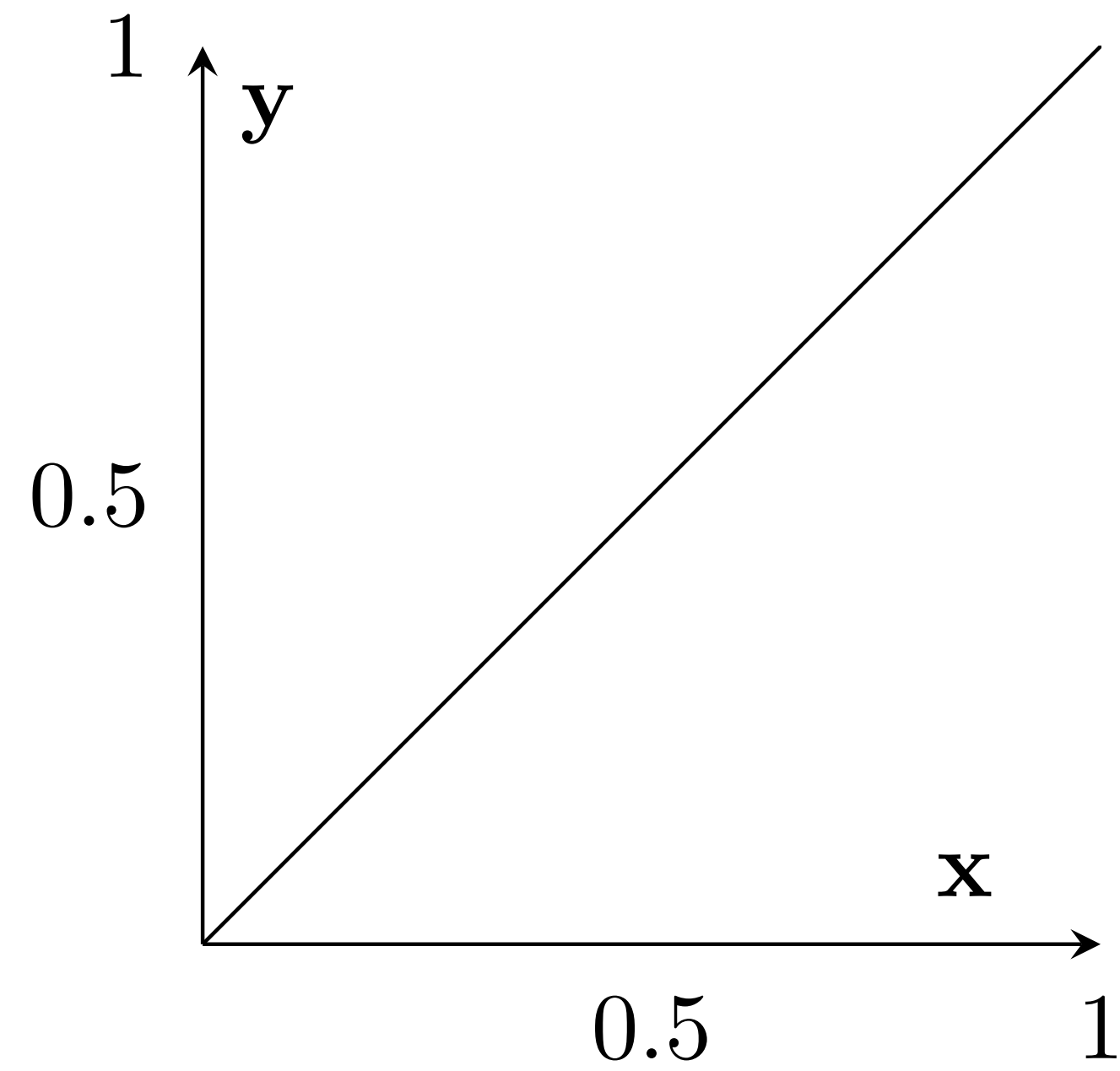
# Deep nets are data transformers

- Deep nets transform datapoints, layer by layer
- Each layer is a different *representation* of the data
- We call these representations **embeddings**

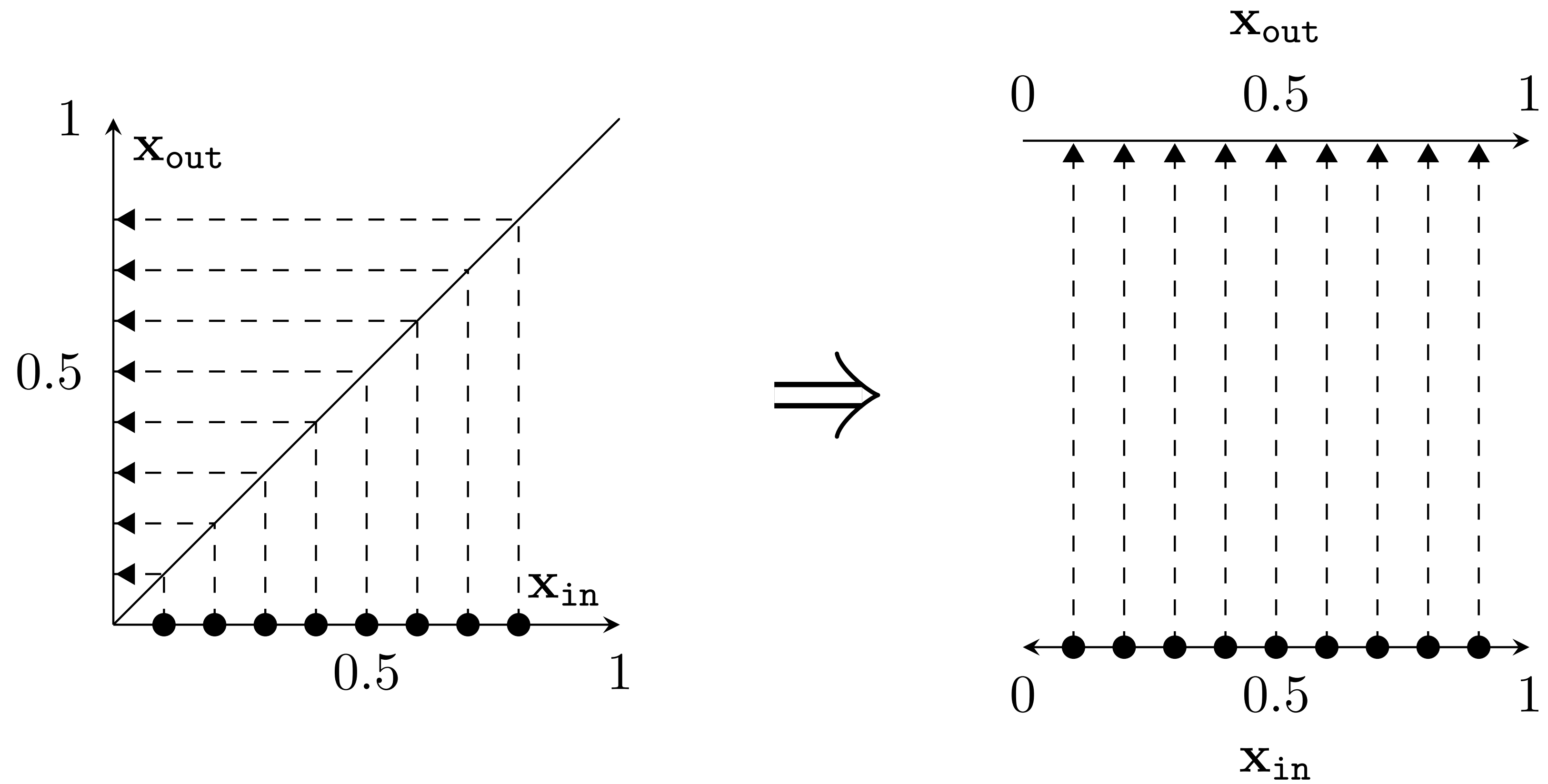




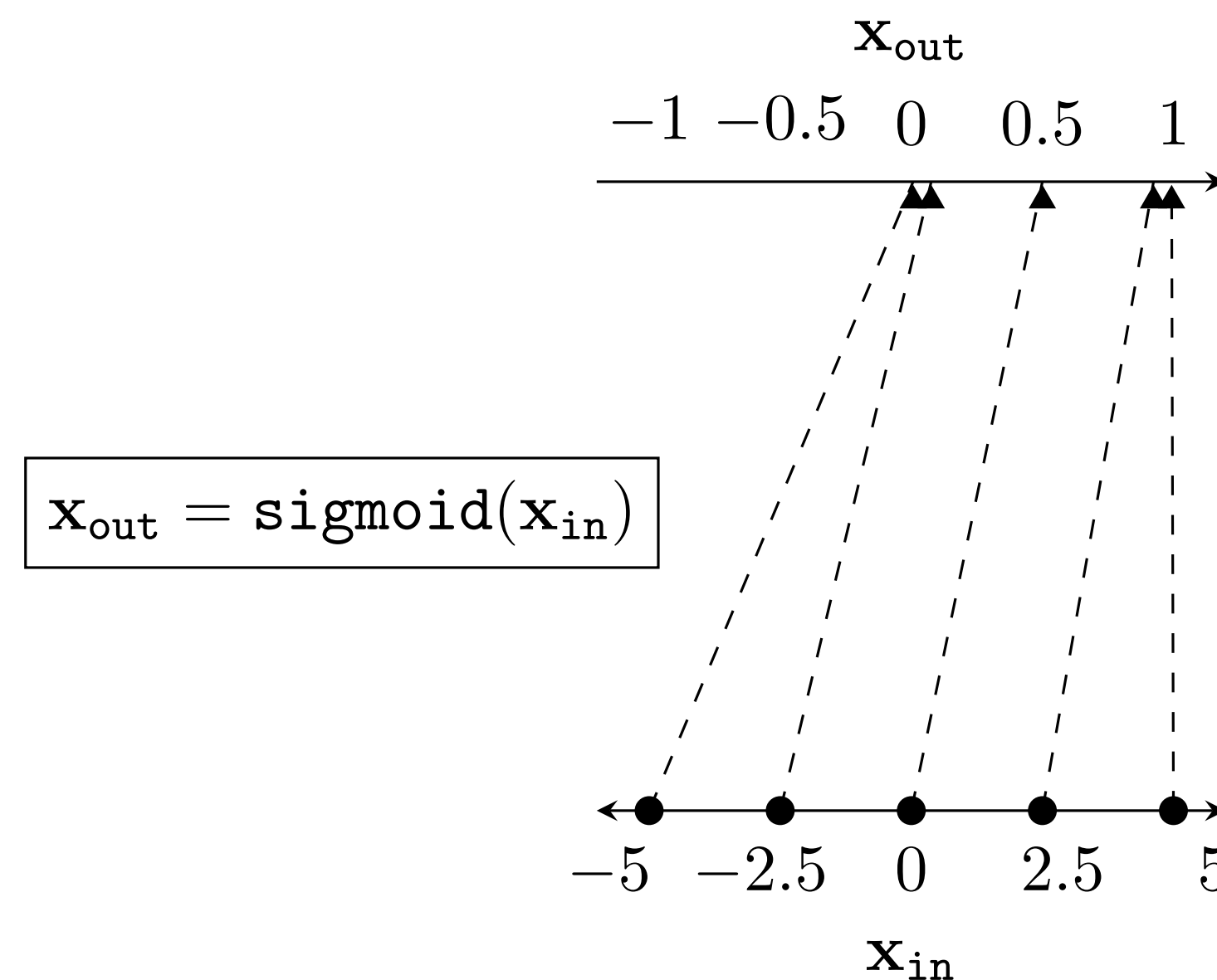
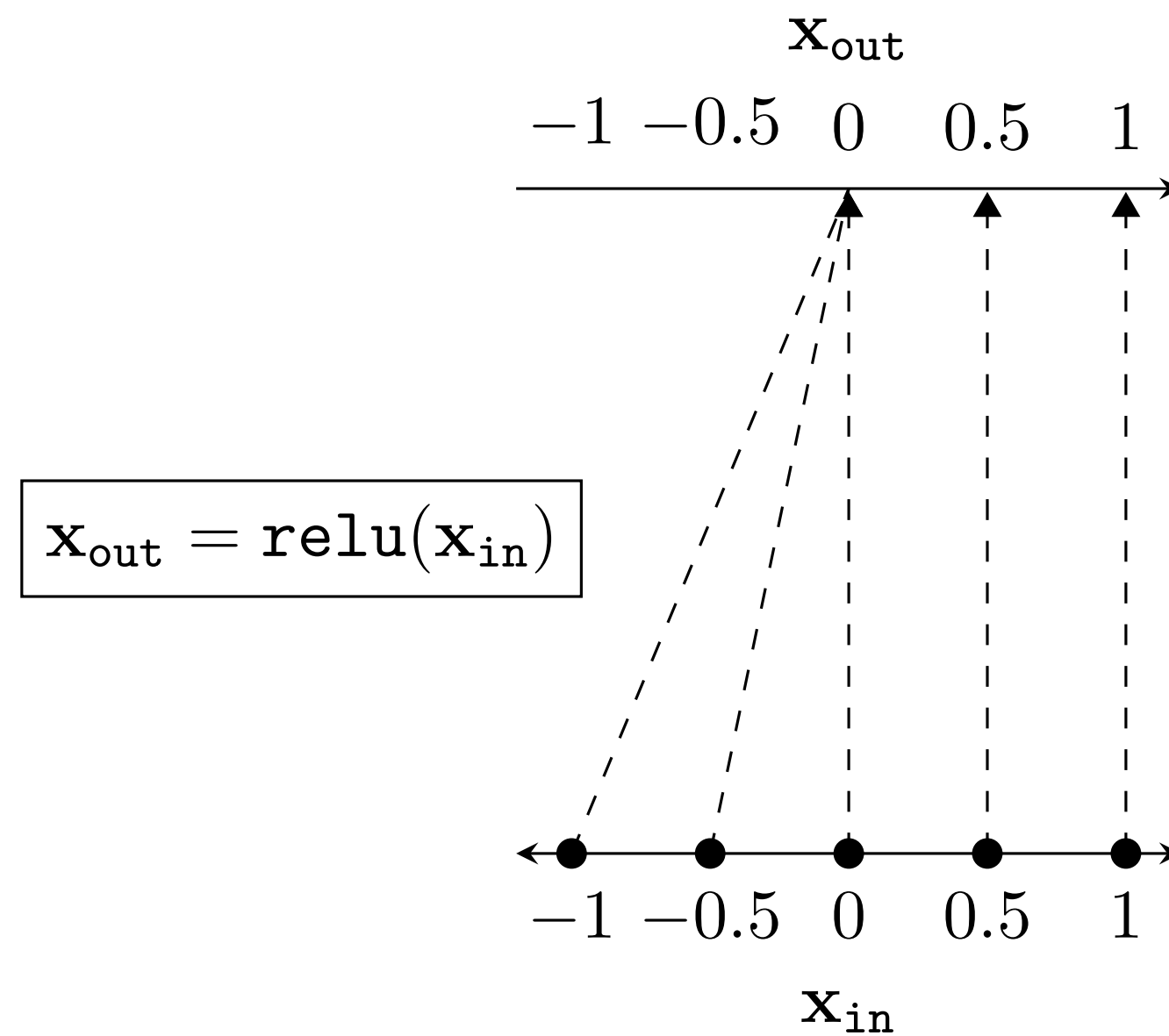
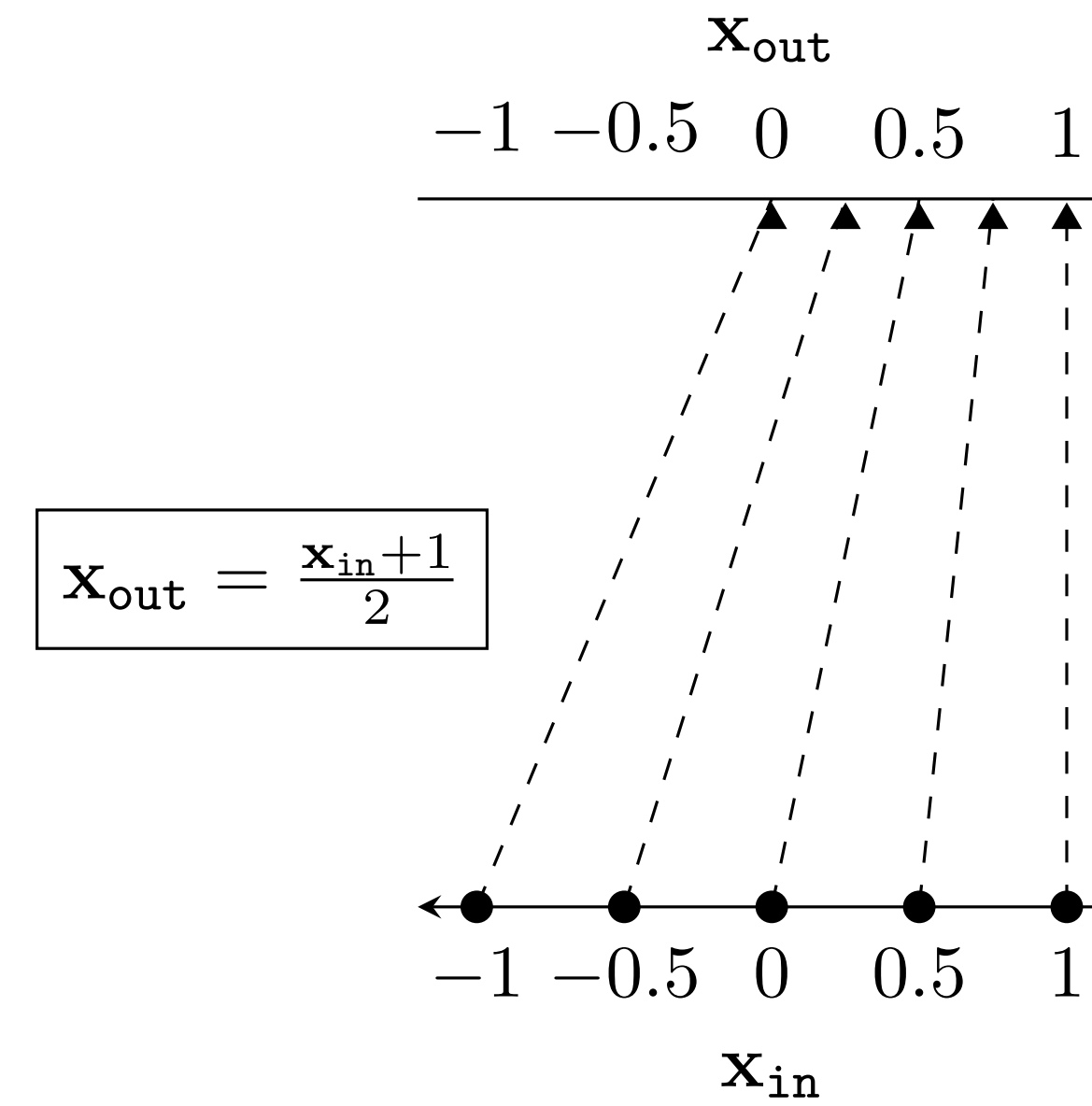
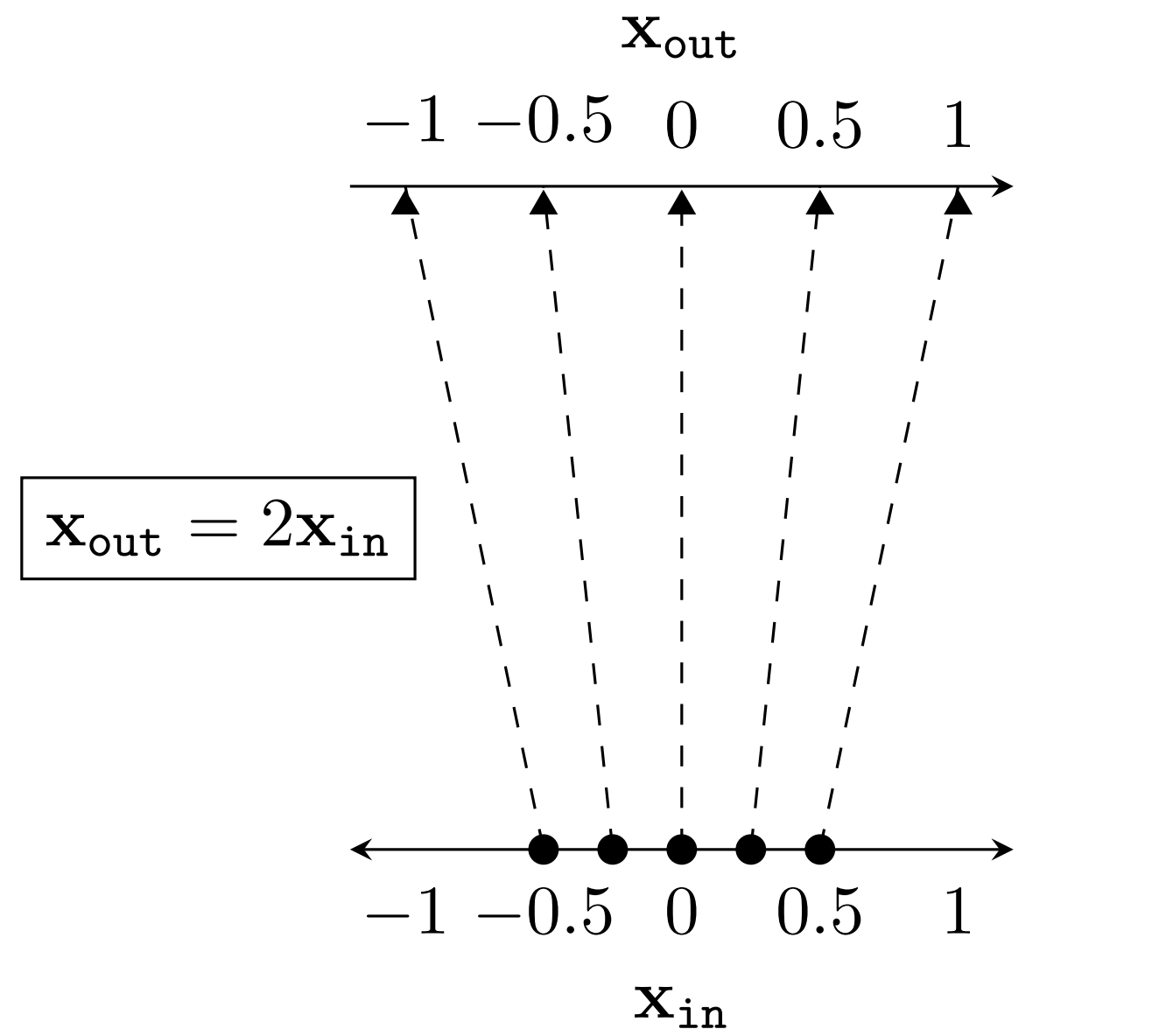
# Two different ways to represent a function



# Two different ways to represent a function



# Data transformations for a variety of neural net layers

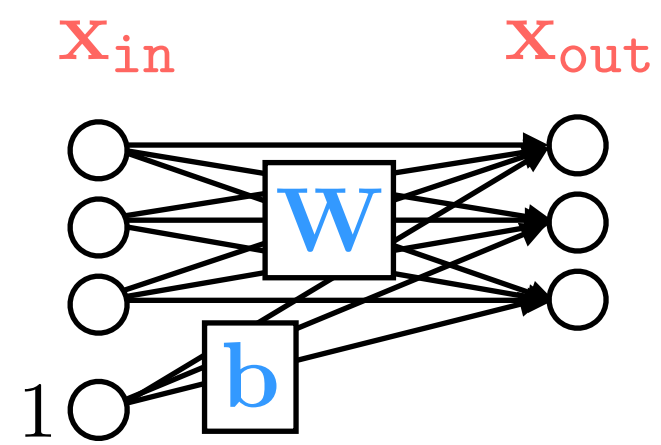


Activations

Parameters

Wiring graph

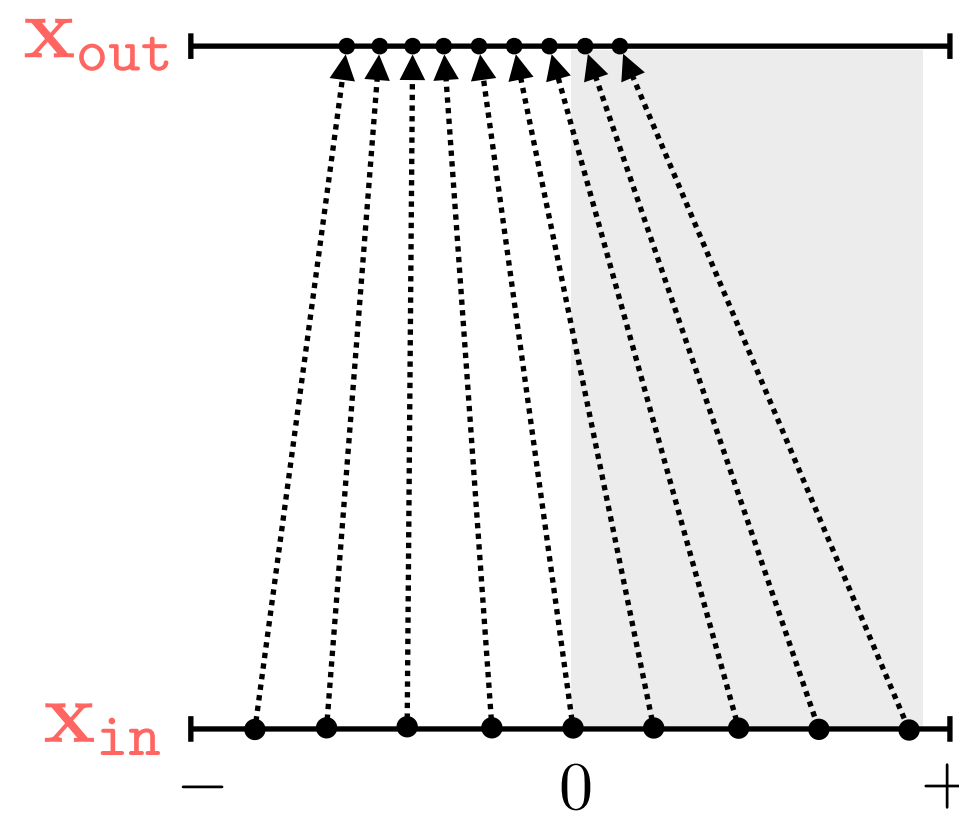
Linear



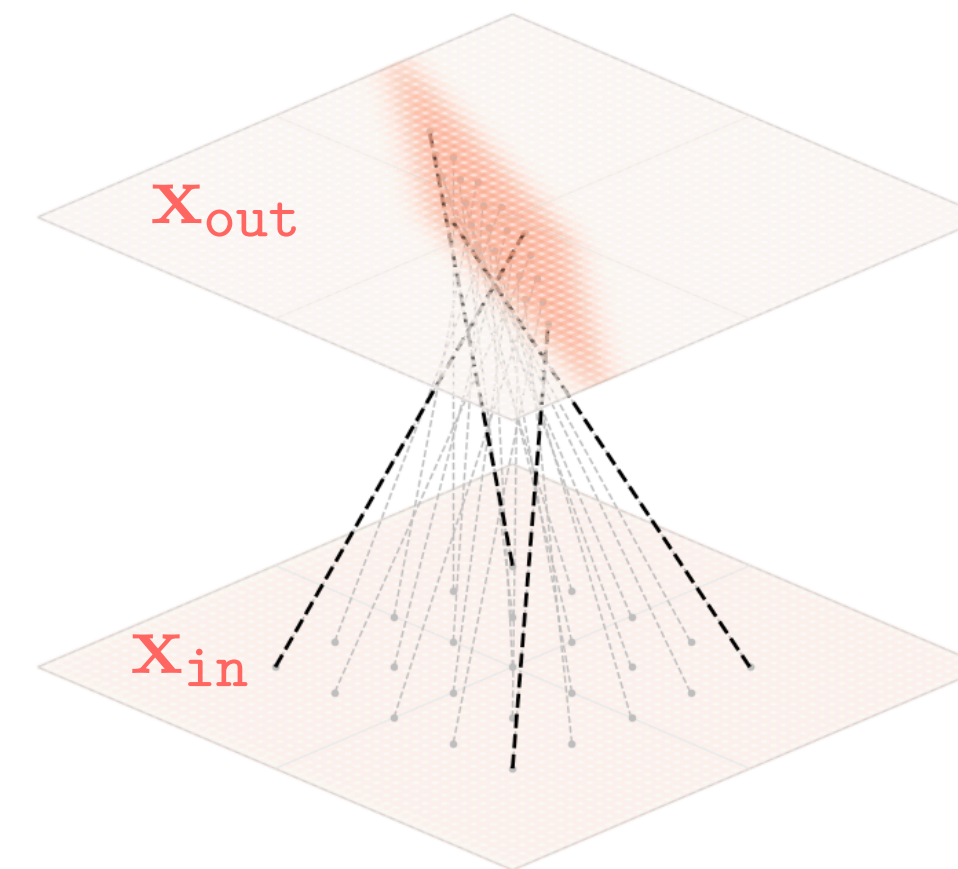
Equation

$$\mathbf{x}_{out} = \mathbf{W}\mathbf{x}_{in} + \mathbf{b}$$

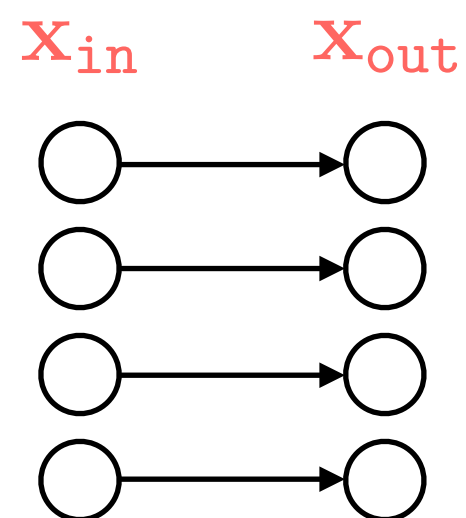
Mapping 1D



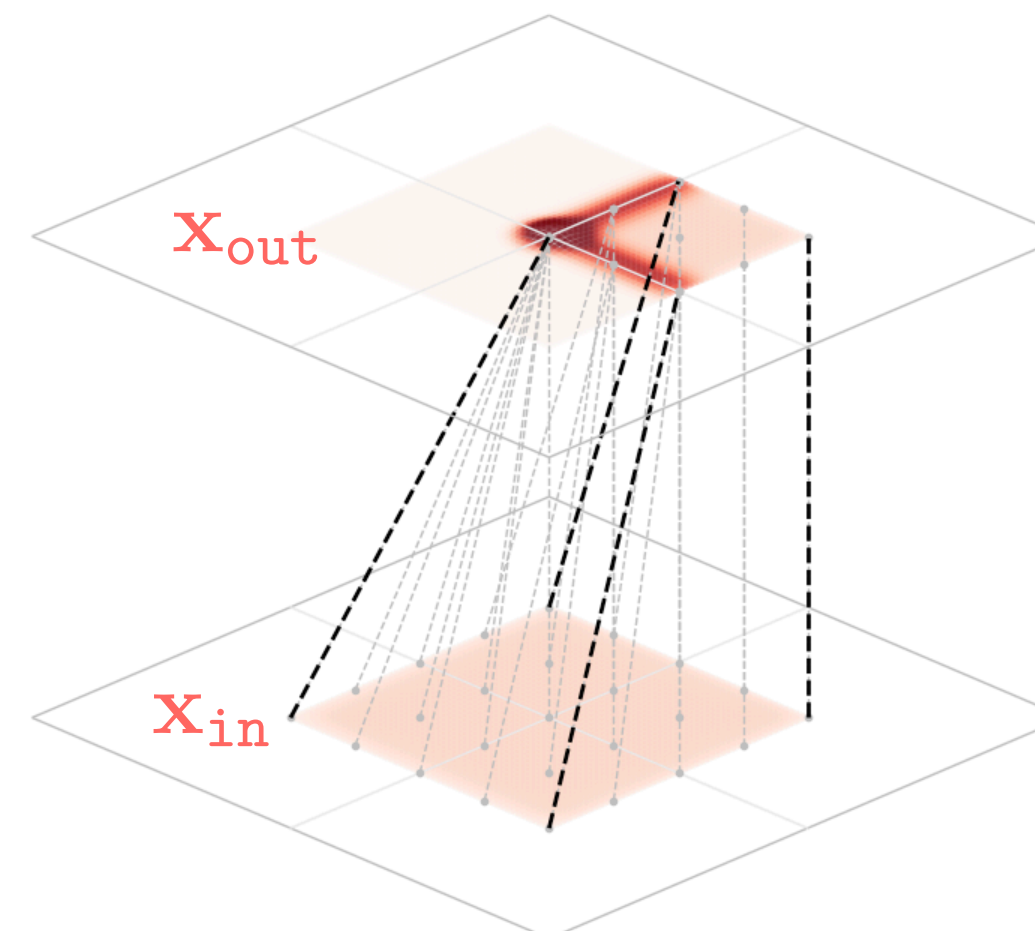
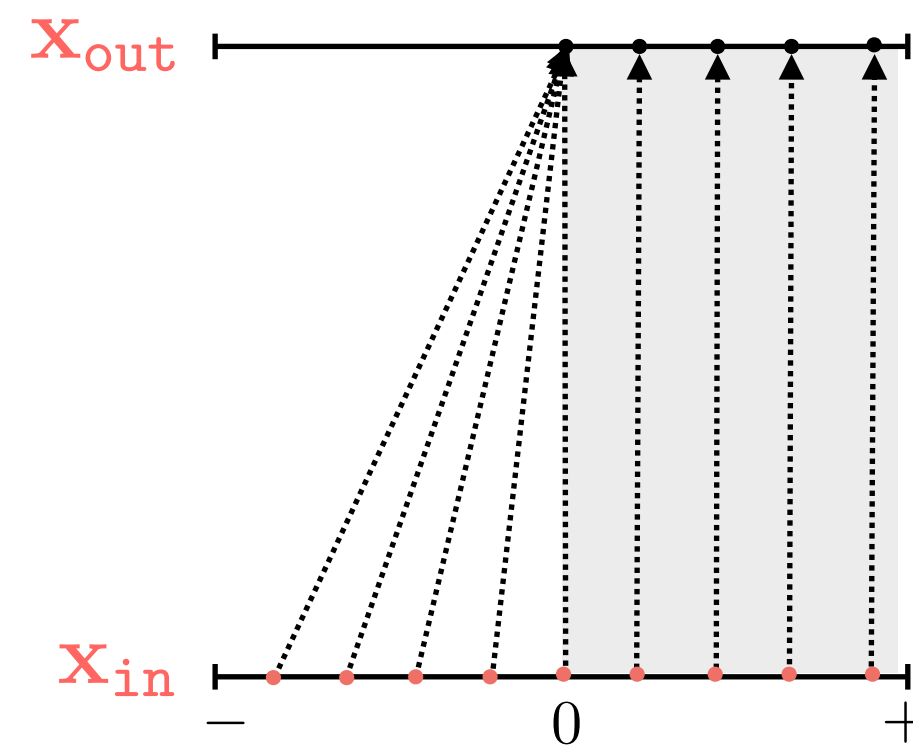
Mapping 2D

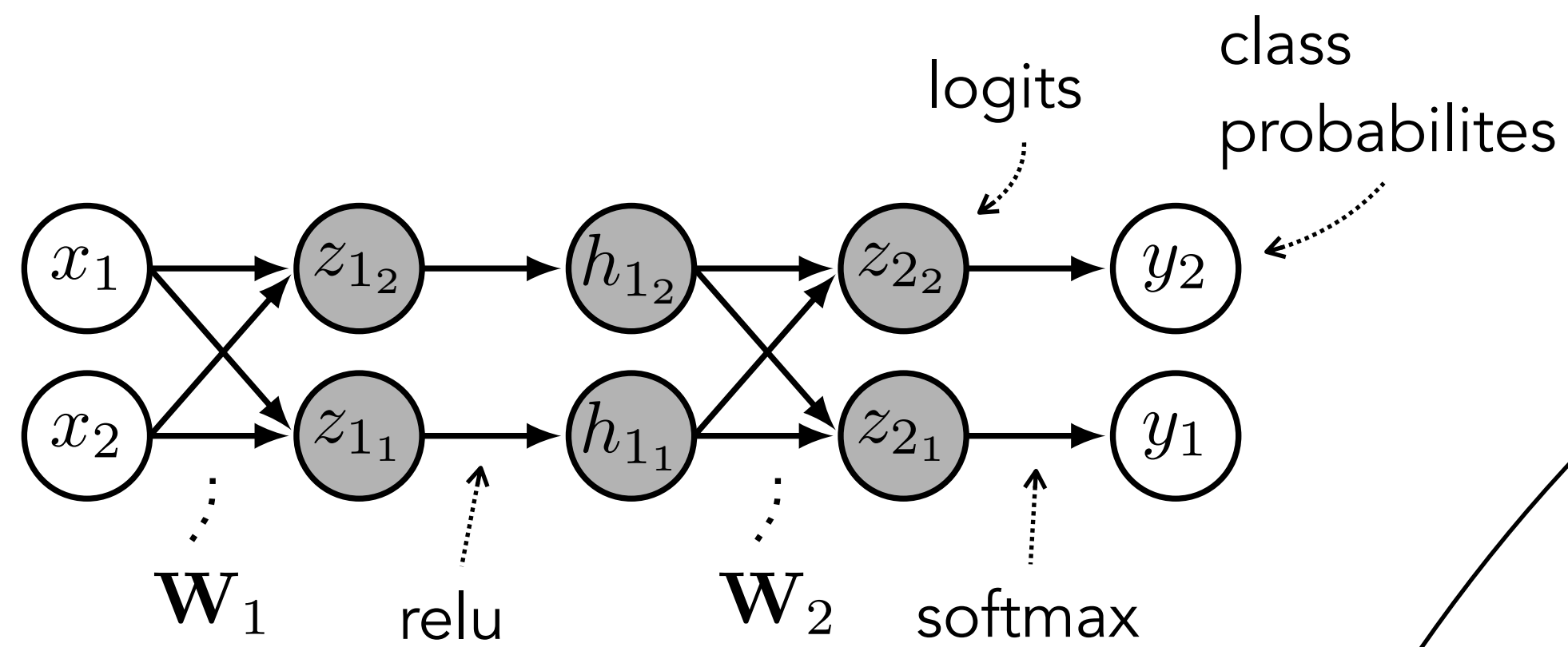


relu

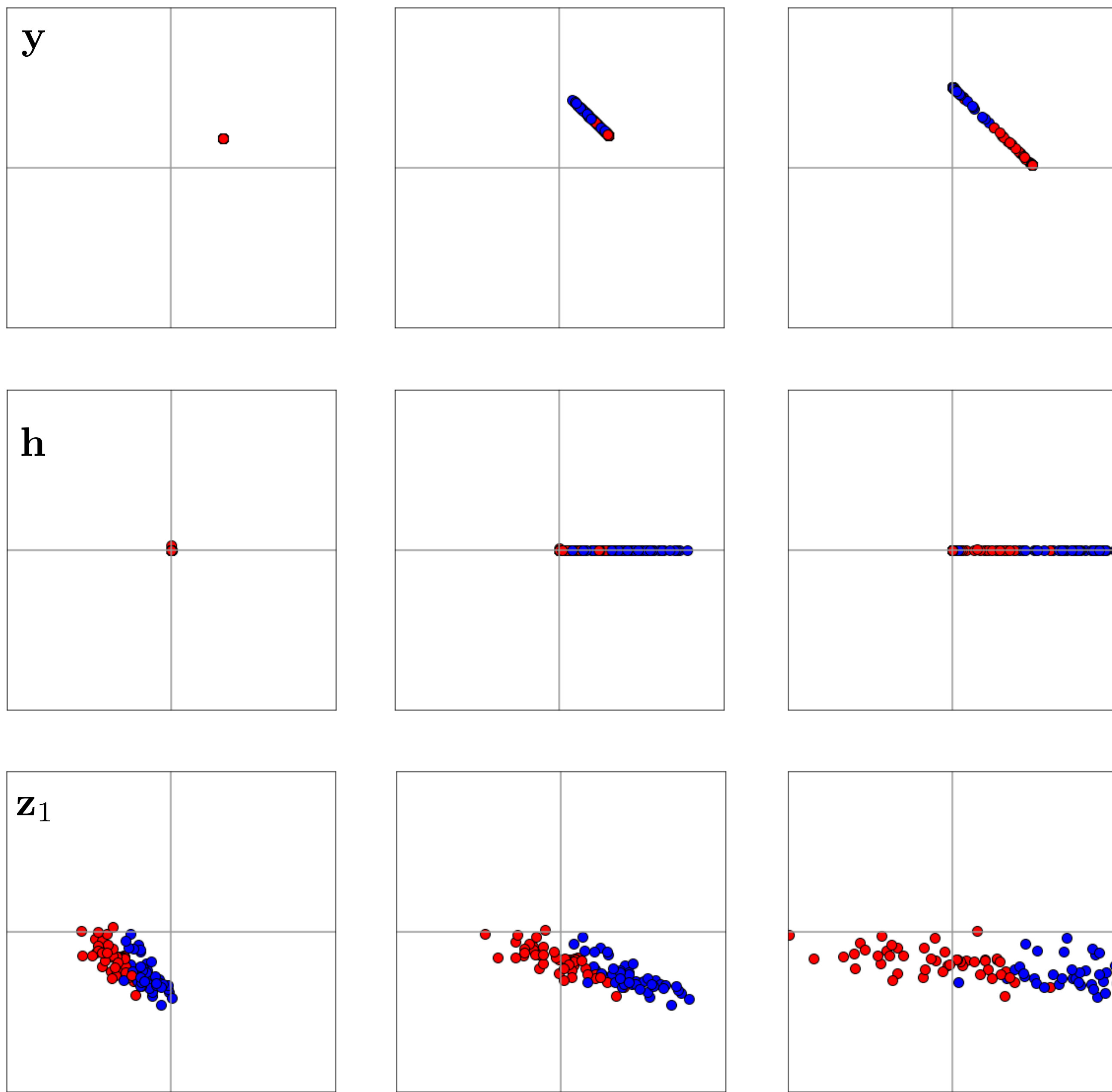
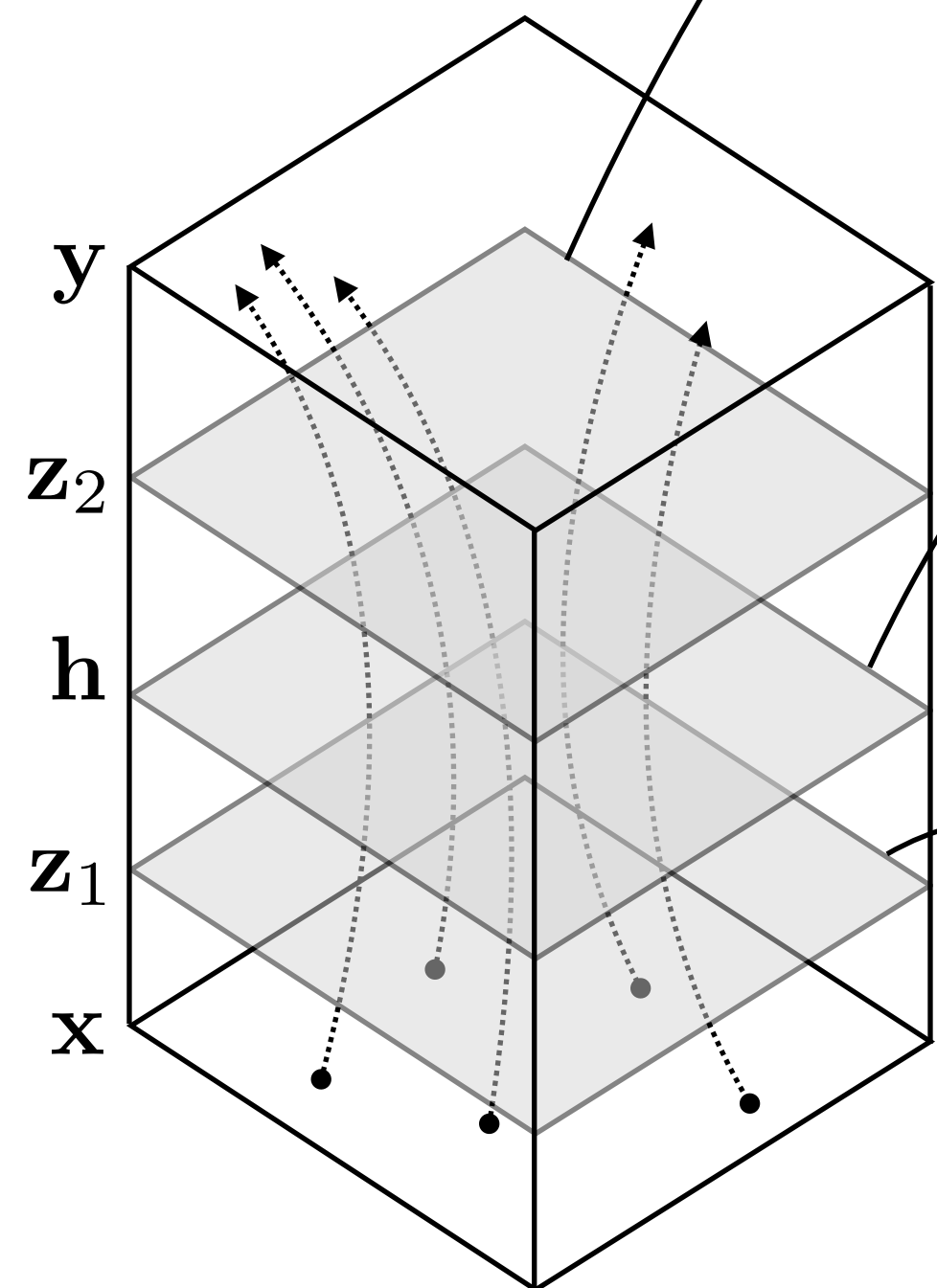
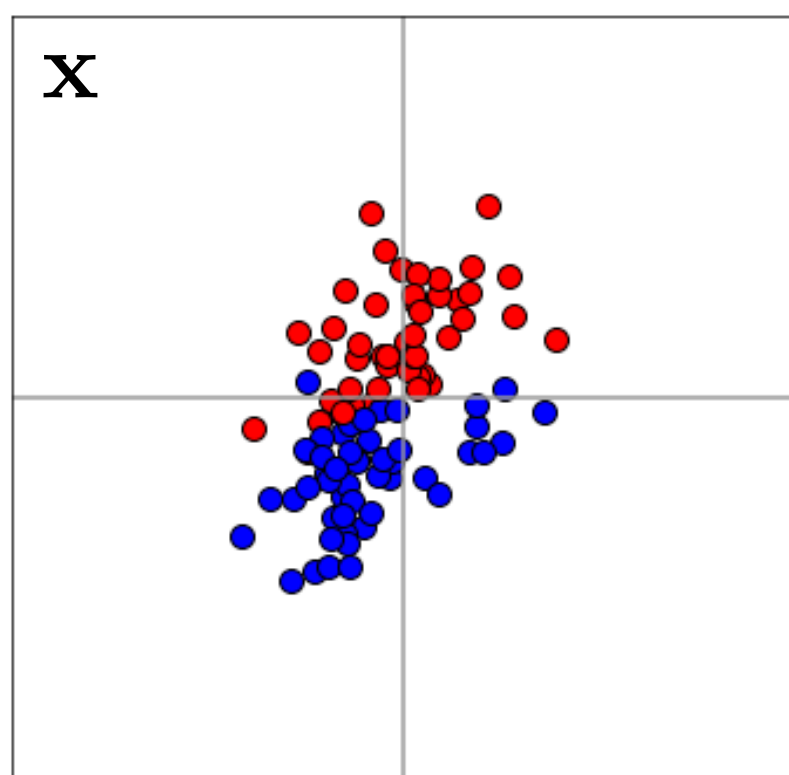


$$x_{out_i} = \max(x_{in_i}, 0)$$



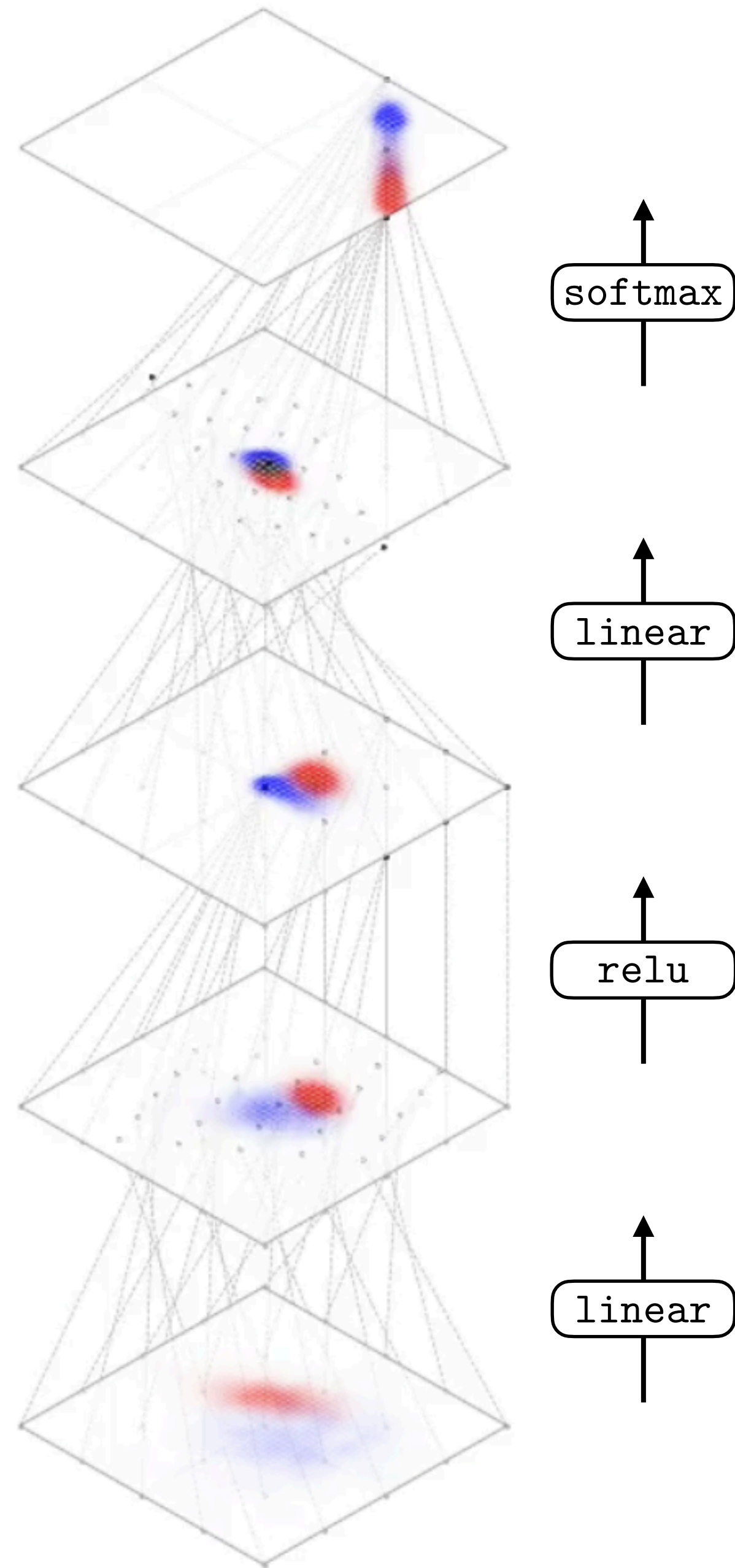


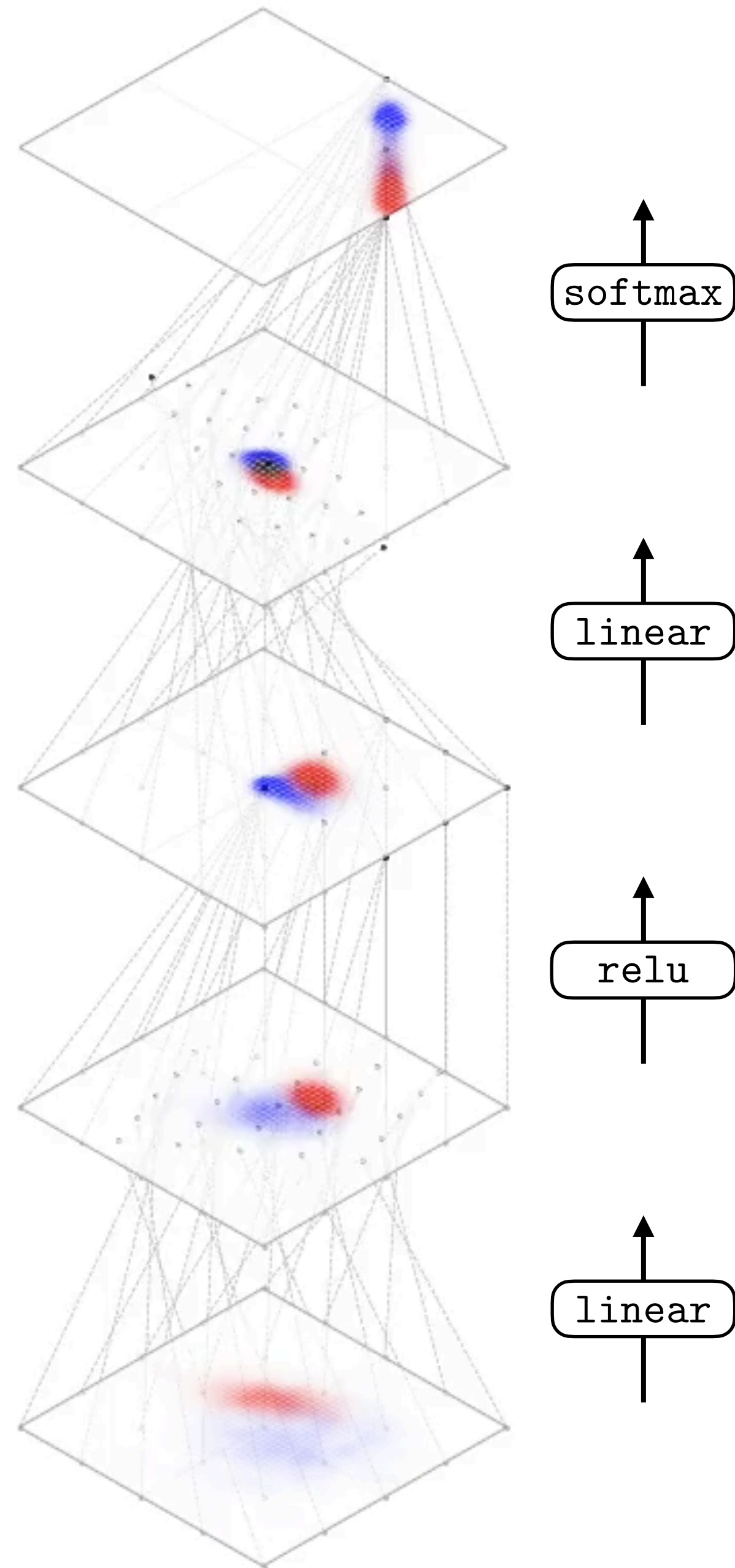
Training data

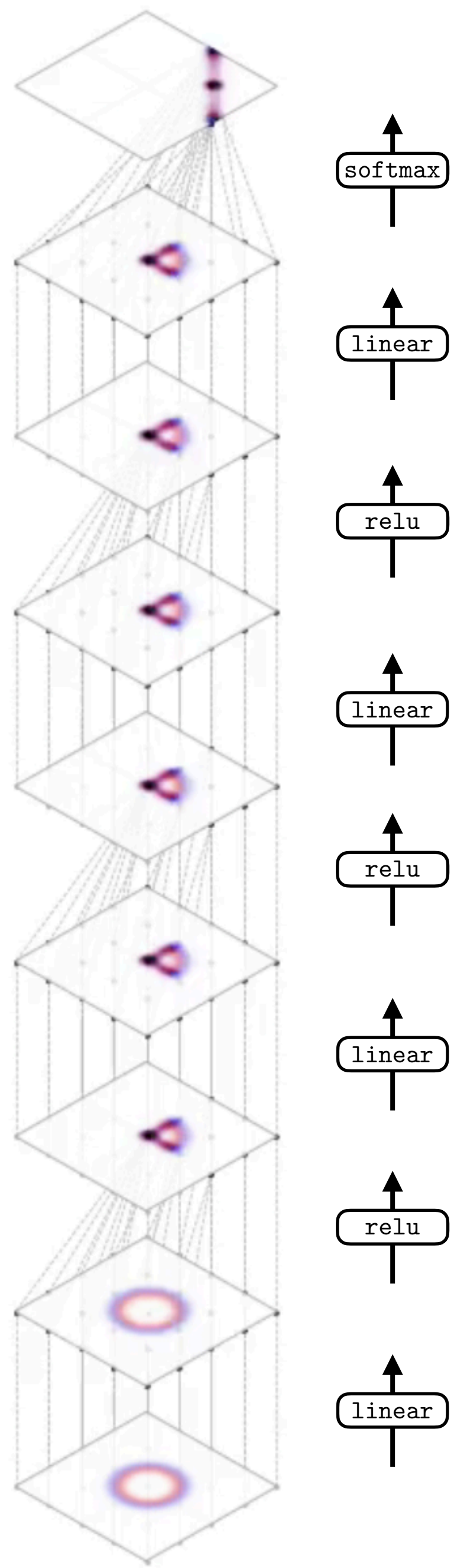


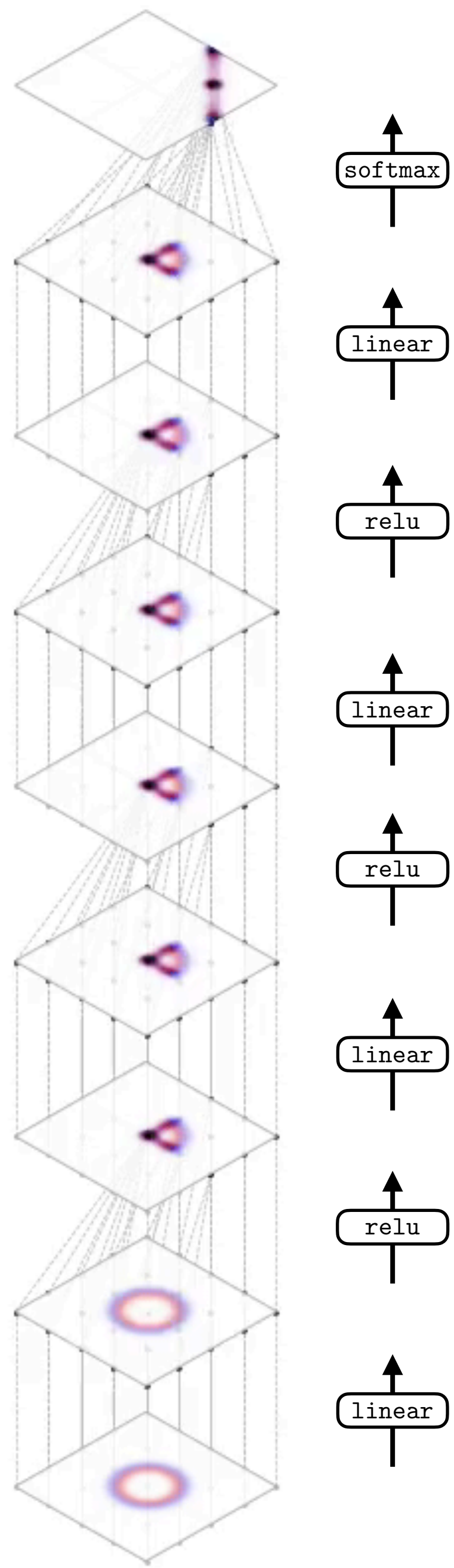
Training iteration  $\longrightarrow$



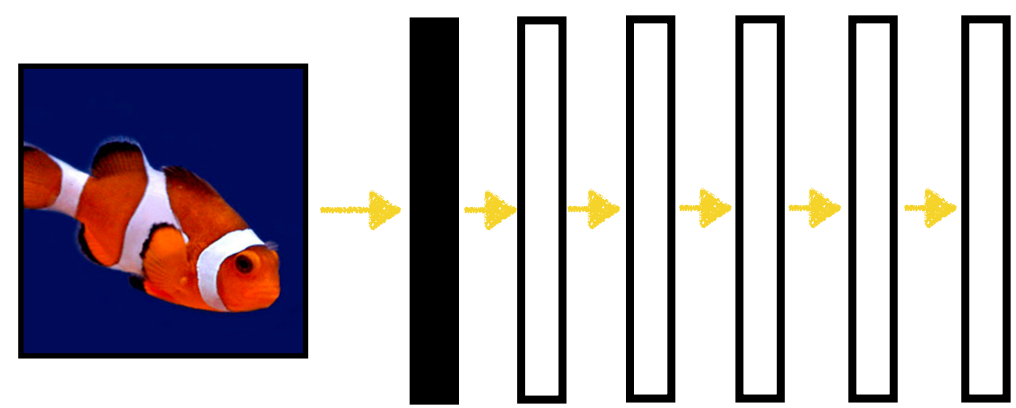




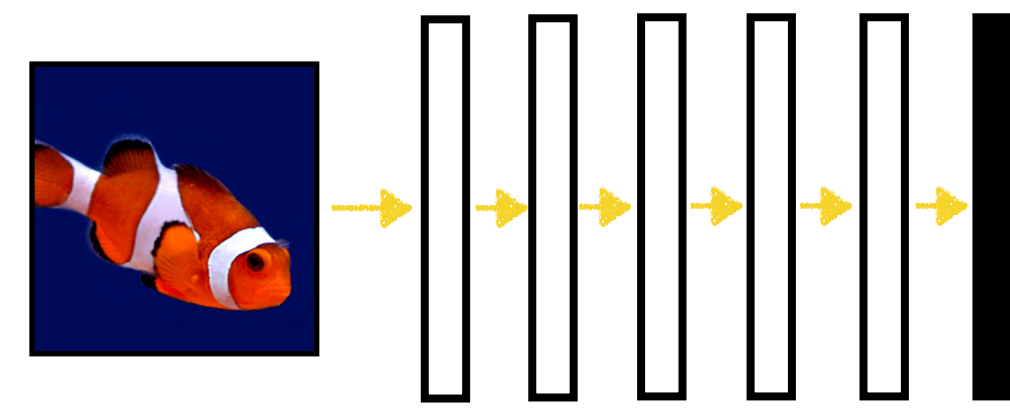




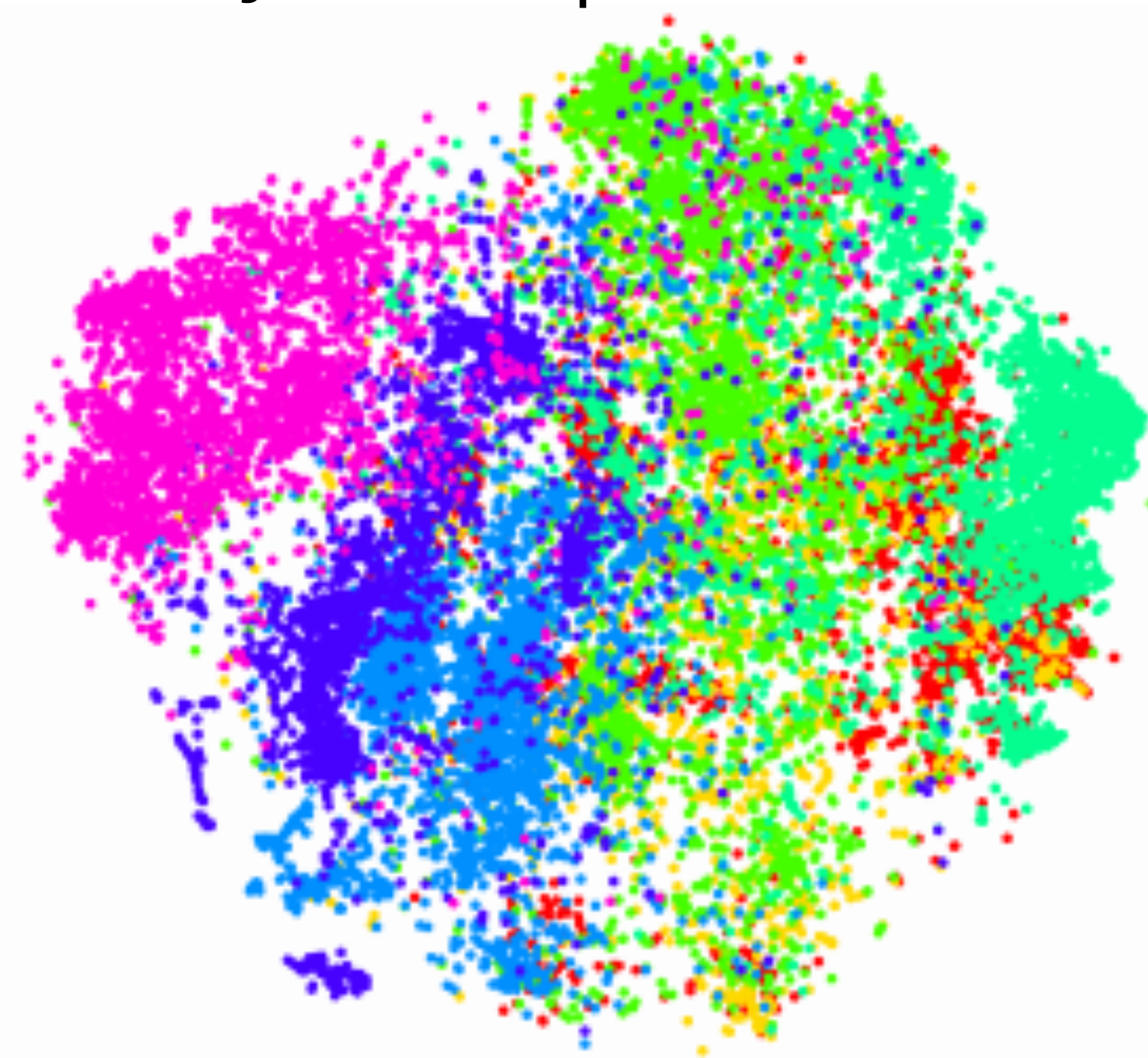
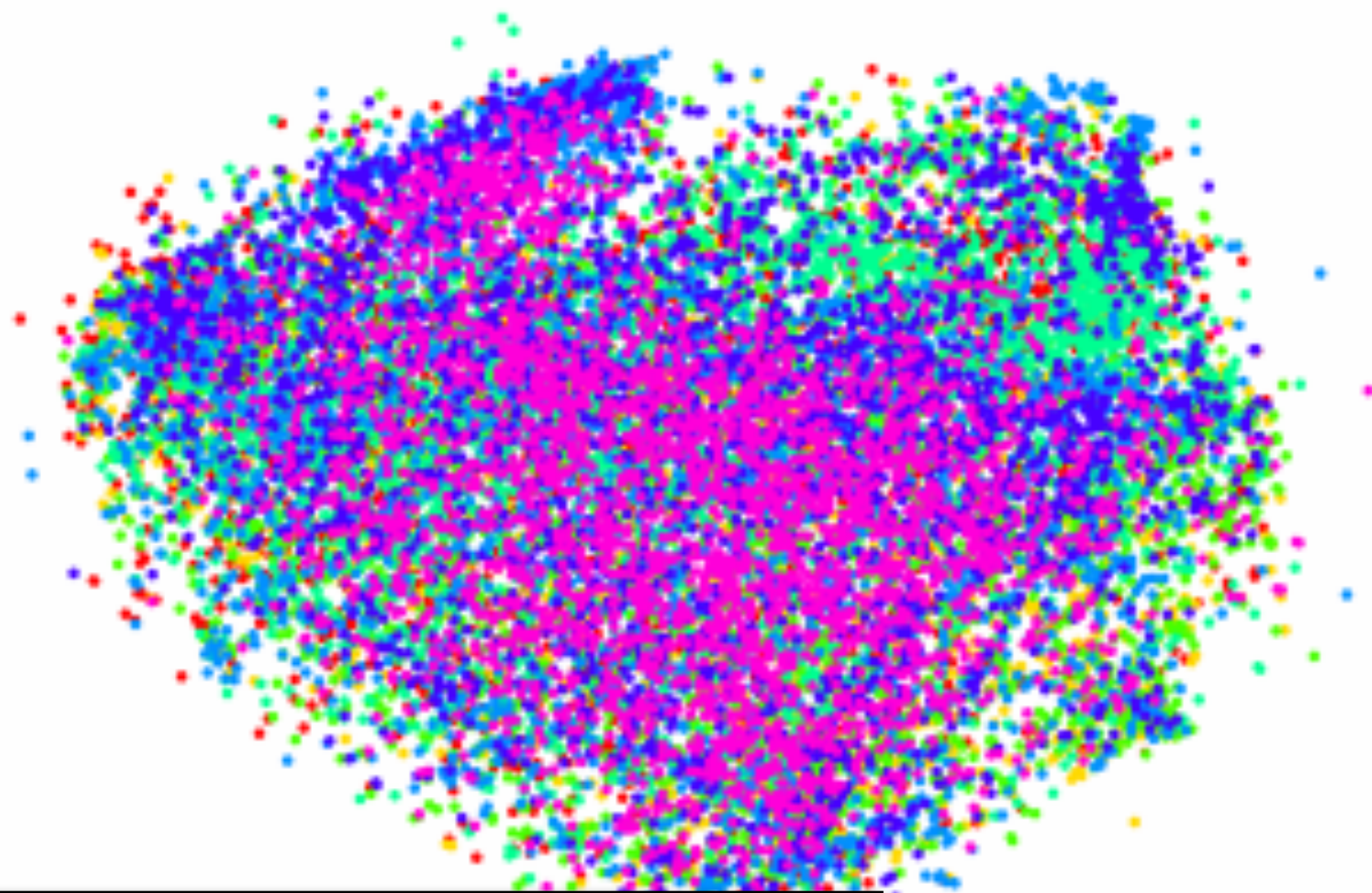




Layer 1 representation



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]



# 7. Introduction to Deep Learning

- Brief history
- Basic formulation (*hierarchical processing*)
- Optimization via gradient descent
- Layer types (*Linear, Pointwise non-linearity*)
- Everything is a tensor
- Deep nets as data transformers