

# ML Review

CMSC 491/691  
Robust Machine Learning



**Artificial  
Intelligence**

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

# Course Staff



Instructor: Tejas Gokhale  
Assistant Professor, CSEE

OH:

Wed 2:30 – 3:30 PM ITE 214

[gokhale@umbc.edu](mailto:gokhale@umbc.edu)



TA/Grader ...

???

# Course Website

<https://courses.cs.umbc.edu/graduate/691rml/>



# Class Structure: Overview

- [Wed 08/28]: Today, Class Logistics and Introduction
- [Mon 09/02]: Labor Day, NO CLASS
- [Wed 09/04; Mon 09/09; Wed 09/11]: Machine Learning Review
- Every Week after that:
  - MON: [TEJAS] Overview of a Robustness Challenge
  - WED: [STUDENTS] Quiz; Paper Presentations; Class Discussion



# Announcements

(1) I sent a “Welcome” announcement on Blackboard. No one read it ...

Student Preview

Exit

CMSC491\_2533\_FA2024

CMSC 491 Special Topics in Computer Science (24.2533/CMSC691\_2534) FA2024

Content

Calendar

Announcements

Discussions

Gradebook

Messages


Analytics

Groups

1 Total

Announcement

Posted



**Welcome, Useful Resources, and Logistics**  
Hello everyone, Welcome to the Robust Machine Learning class! We will be using Blackboard for assignment submissions, ...

8/30/24, 5:10 PM

(2) This week and the next, we will do ML review.

(3) Please sign up for presentations and surveys. Link is in Blackboard Announcement.

(4) Start forming project groups / thinking about ideas. Discuss with me in OFFICE HOURS.

# Forcing Function

- [CMSC 491/691 RML Presentation and Survey Signup - Google Sheets](#)

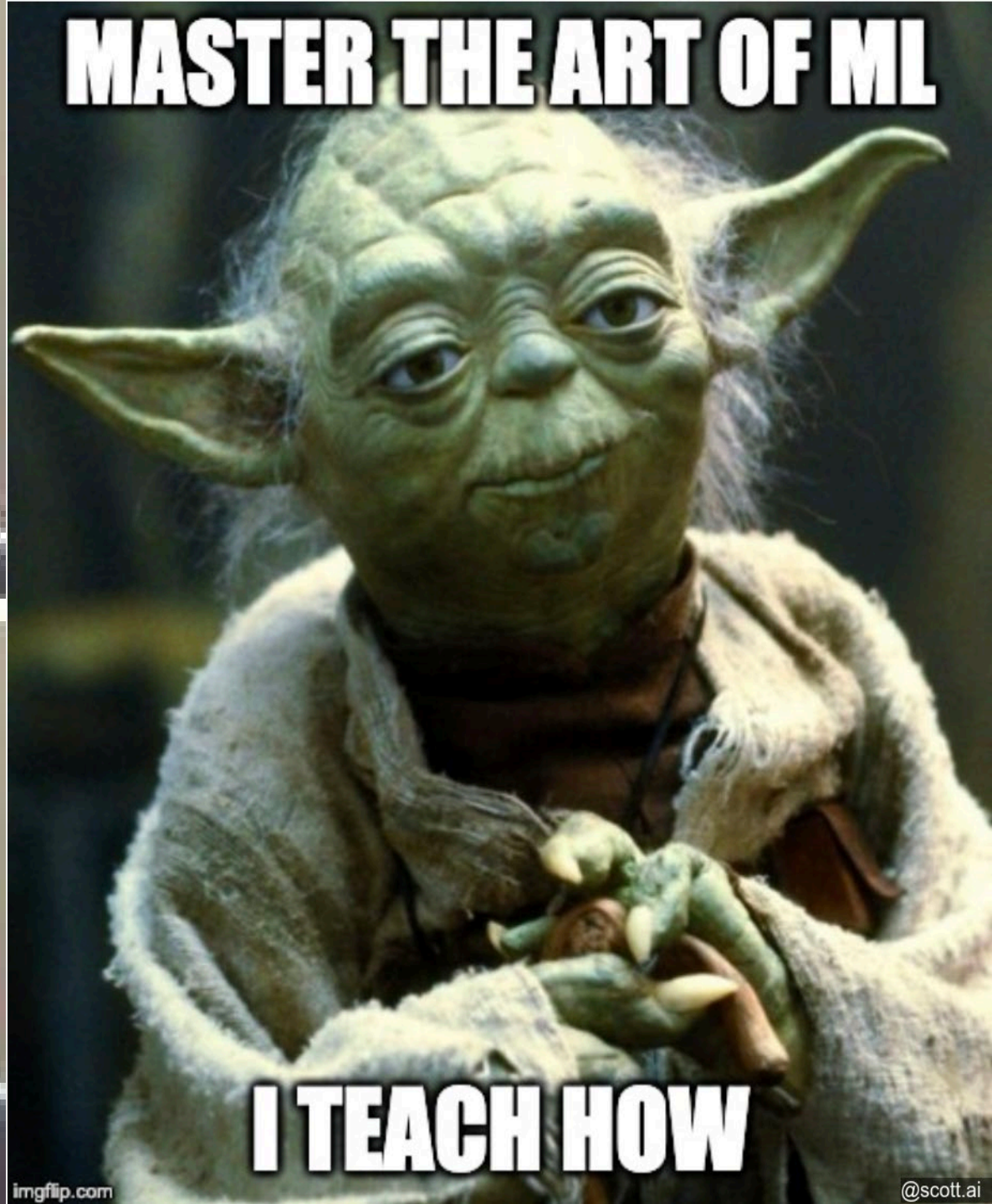


# **Machine Learning Review**

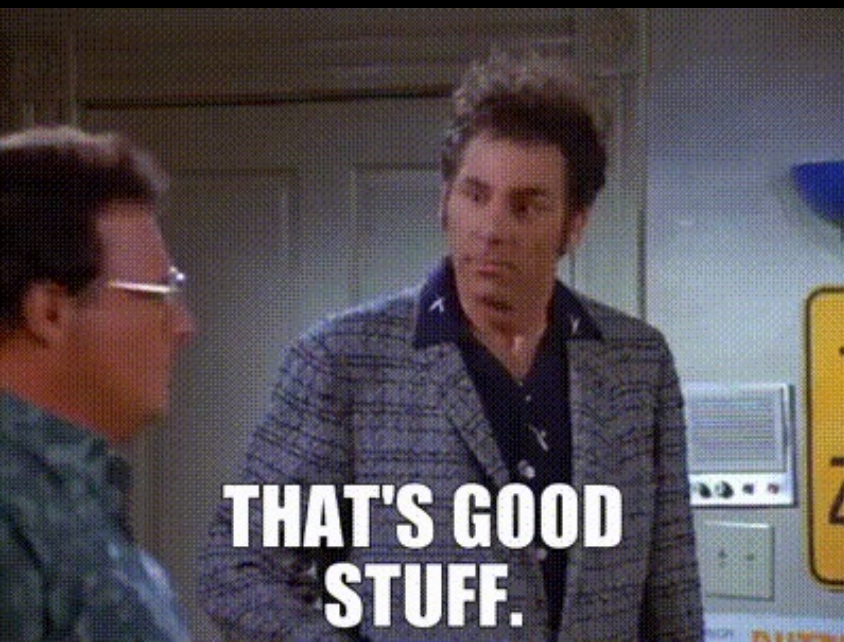
Machine Learning Training and Inference Pipeline

Neural Networks / Deep Learning

ML methods for visual recognition







# Motivating Example: Image Classification



What is this?

{dog, cat, airplane, bus, laptop,  
chair ...}

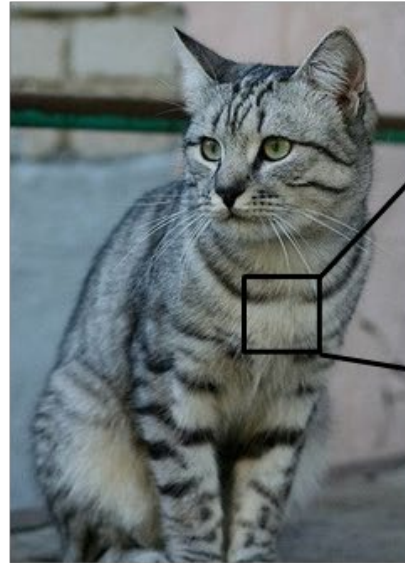
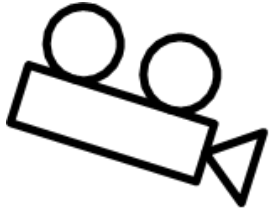
What animal is this ?

{dog, cat, lion, tiger, duck, cow,  
giraffe, ...}

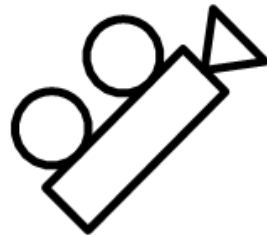
What type of cat is this?

{Cheshire, Siamese, Persian,  
Shorthair, Bombay, ...}

# Challenges: Viewpoint Variation



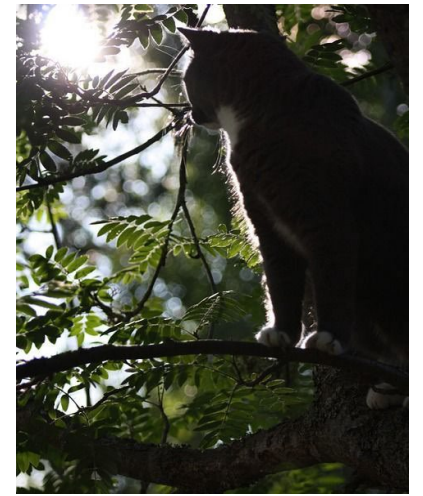
```
[[105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]  
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]  
[ 76 85 98 105 128 105 87 96 95 99 115 112 106 103 99 85]  
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]  
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]  
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]  
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]  
[128 137 144 140 109 95 86 78 62 65 63 63 68 73 86 101]  
[125 133 148 137 110 121 117 94 65 79 80 65 54 64 72 98]  
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]  
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]  
[ 89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 80]  
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]  
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]  
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]  
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]  
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]  
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]  
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]  
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]  
[128 112 96 117 150 144 120 115 104 107 102 93 87 81 72 79]  
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]  
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]  
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]]
```



All pixels change when  
the camera moves!



# Challenges: Illumination





# Challenges: Background Clutter

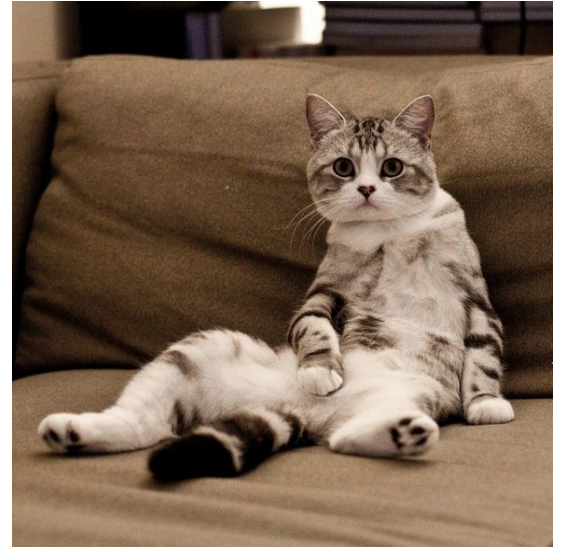
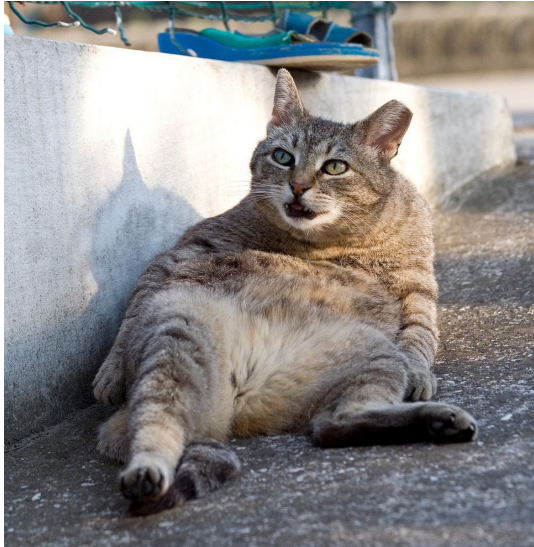




# Challenges: Occlusion



# Challenges: Pose and Deformation (Cat Yoga)





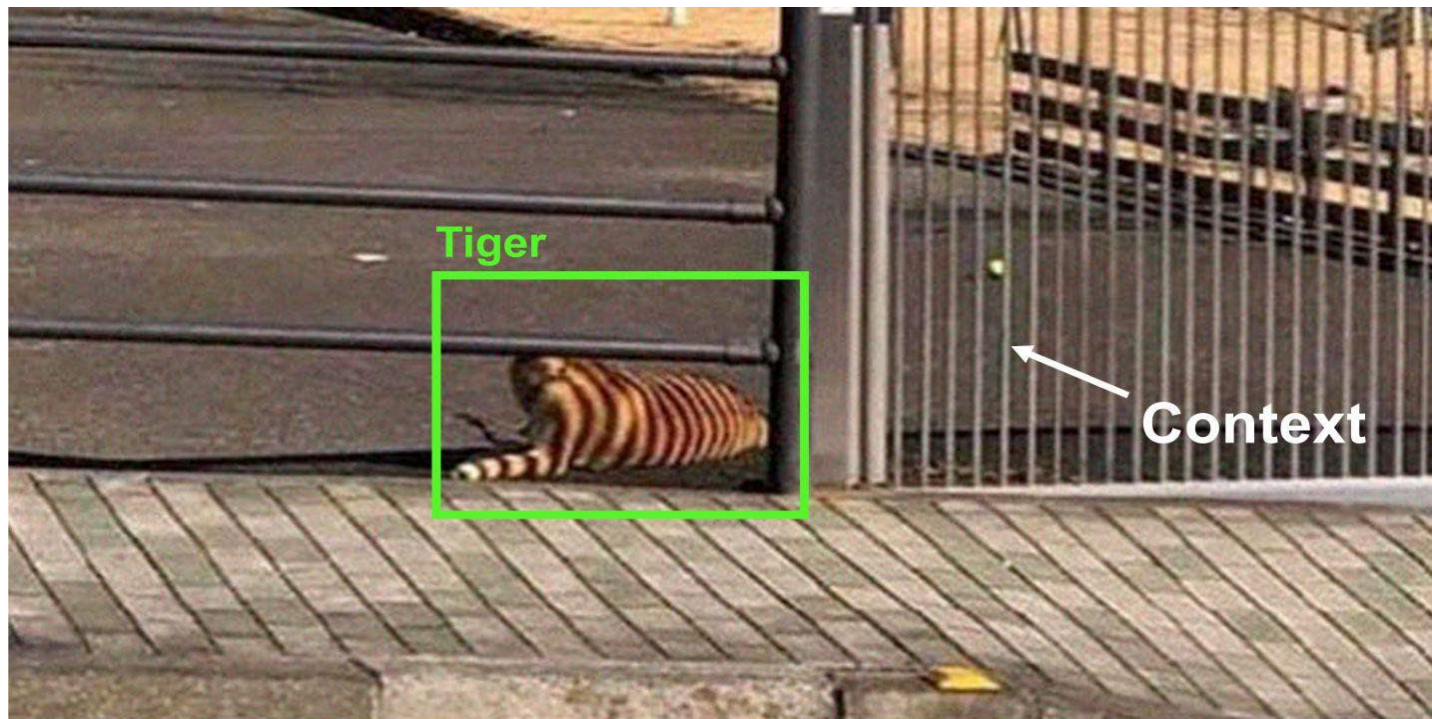
# Challenges: Inter-Class Variation



Slide Credit: Fei-Fei Li



# Challenges: Illusions



# Data !!!



# An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

Unlike e.g. sorting a list of numbers,

**no obvious way to hard-code** the algorithm for recognizing a cat, or other classes.

# Machine Learning

1. Collect a dataset of images and labels
2. Use Machine Learning algorithms to train a classifier
3. Evaluate the classifier on new images

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

Example training set

**airplane**



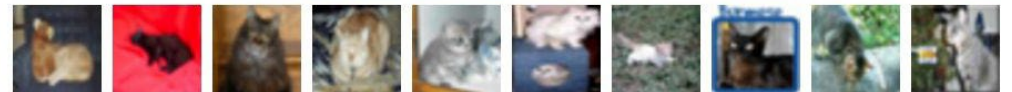
**automobile**



**bird**



**cat**



**deer**





# Nearest Neighbor Classifier

```
def train(images, labels):  
    # Machine learning!  
    return model
```



Memorize all  
data and labels

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```



Predict the label  
of the most similar  
training image

# Nearest Neighbor Classifier

deer



bird



plane



cat



car



Training data with labels



query data



Distance Metric



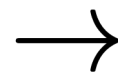
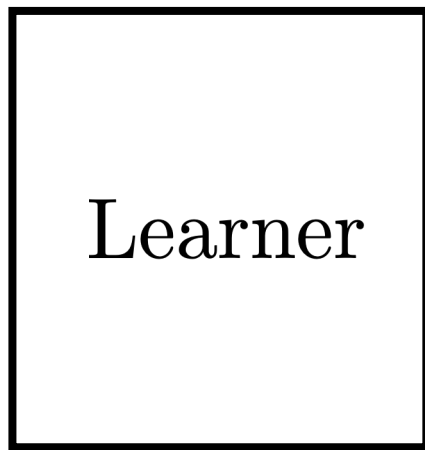
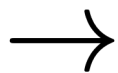
,



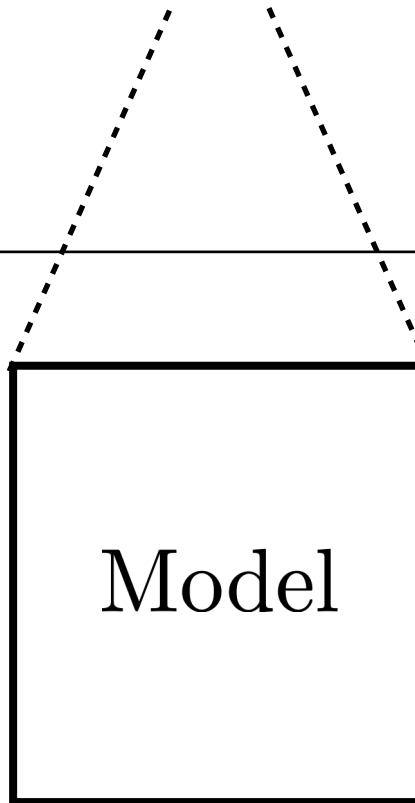
$\rightarrow \mathbb{R}$

Learning

Data

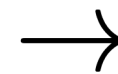
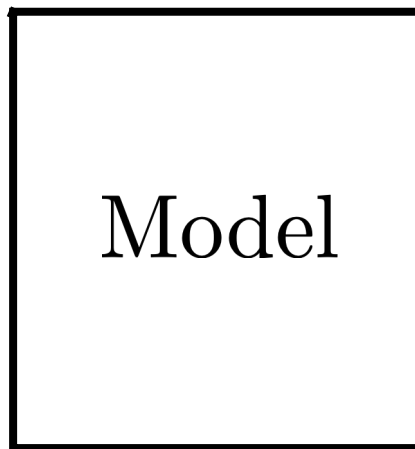
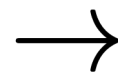


Model



Inference

Input



Output

Learning

Data

```
def train(images, labels):  
    # Machine learning!  
    return model
```

Inference

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

→ Output

The goal of learning is to extract lessons from past experience in order to solve future problems.

**What does ☆ do?**

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

Goal: answer future queries involving ☆

Approach: figure out what ☆ is doing by observing its behavior on examples

**Past experience**

$$2 \star 3 = 36$$

$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$

**Future query**

$$3 \star 5 = ?$$

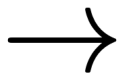
Training

$$2 \star 3 = 36$$

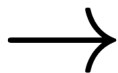
$$7 \star 1 = 49$$

$$5 \star 2 = 100$$

$$2 \star 2 = 16$$



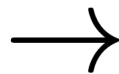
Your  
brain



$$x \star y \rightarrow (xy)^2$$

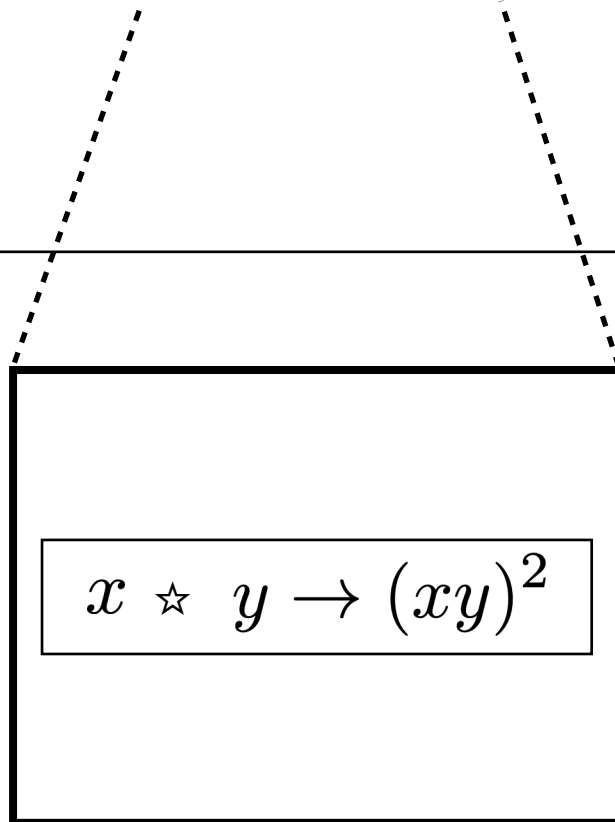
Testing

$$3 \star 5$$



$$x \star y \rightarrow (xy)^2$$

$$\rightarrow 225$$



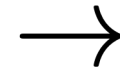
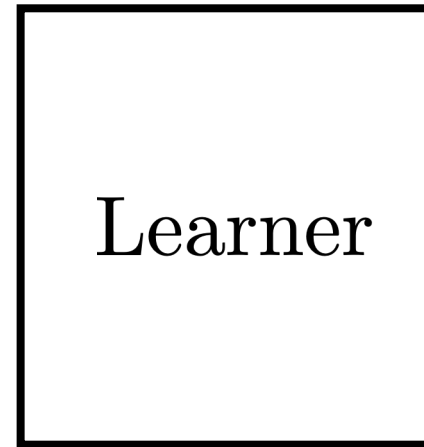
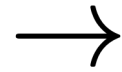


# Learning from examples

(aka **supervised learning**)

Training data

$\{\text{input}: [2, 3], \text{output}: 36\}$   
 $\{\text{input}: [7, 1], \text{output}: 49\}$   
 $\{\text{input}: [5, 2], \text{output}: 100\}$   
 $\{\text{input}: [2, 2], \text{output}: 16\}$



$f$

# Learning from examples

(aka **supervised learning**)

Training data

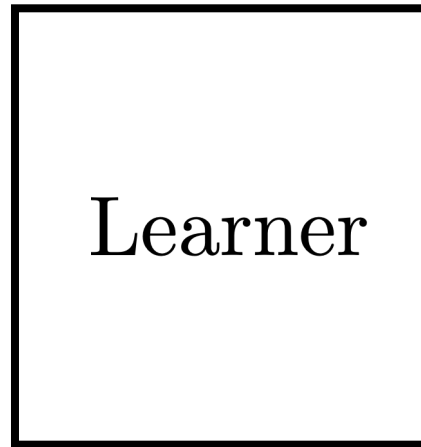
$\{x^{(1)}, y^{(1)}\}$

$\{x^{(2)}, y^{(2)}\}$

$\{x^{(3)}, y^{(3)}\}$

...

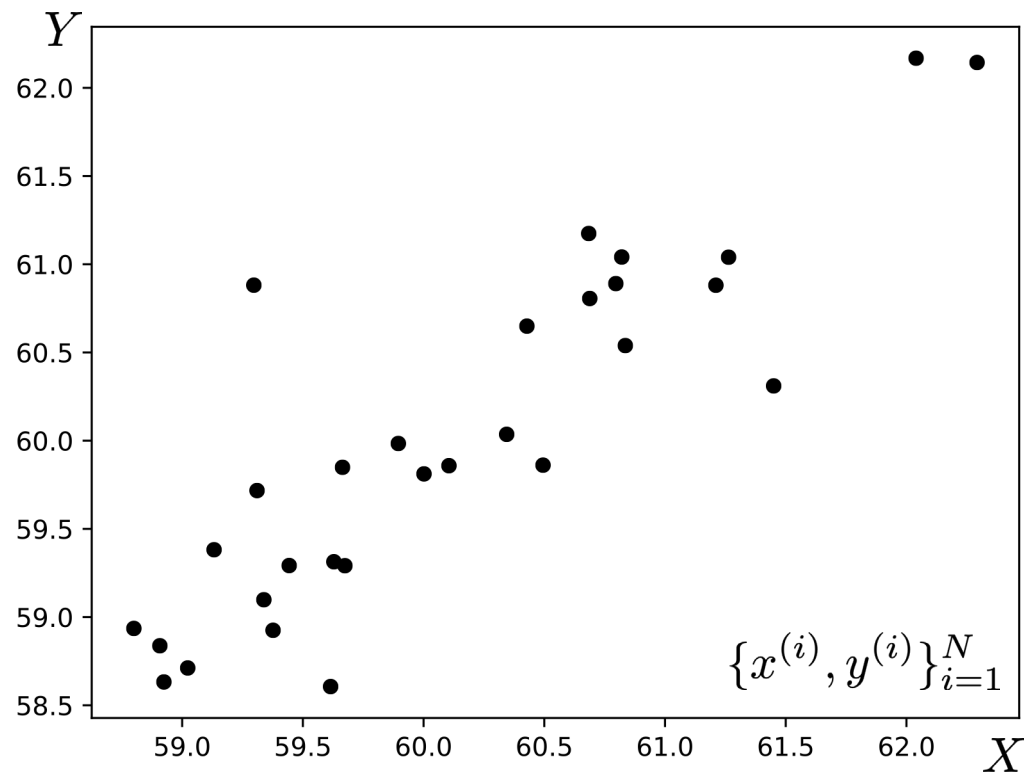
→



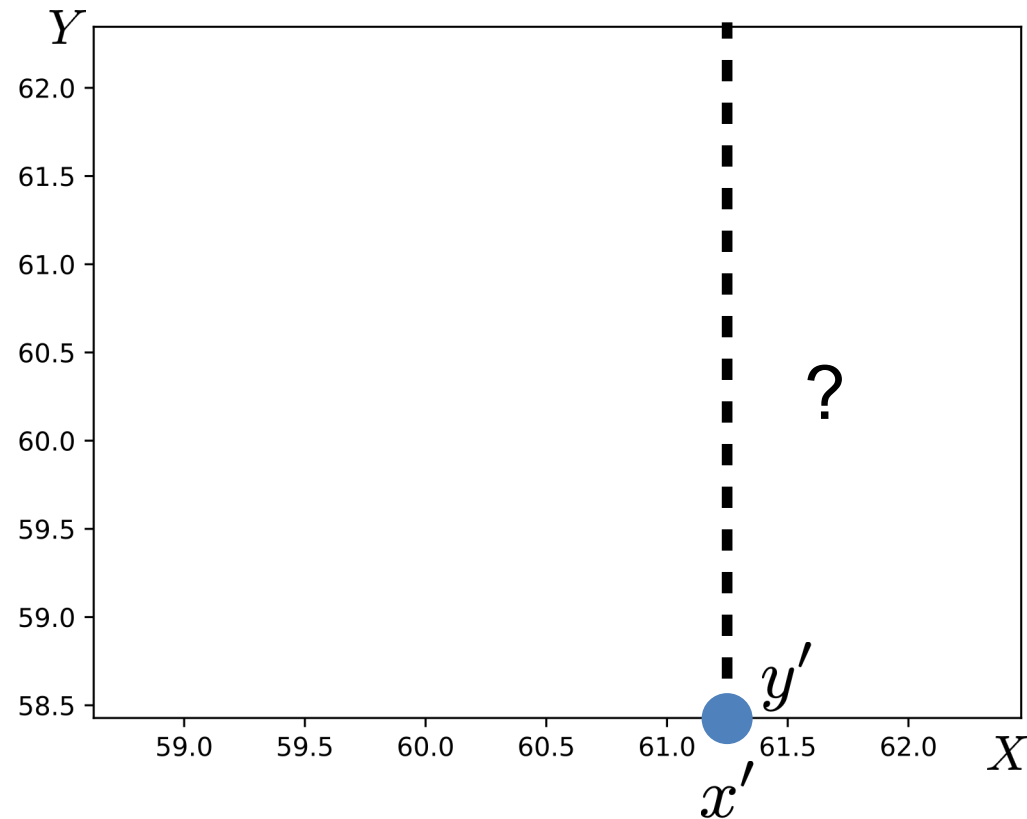
→

$f : X \rightarrow Y$

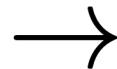
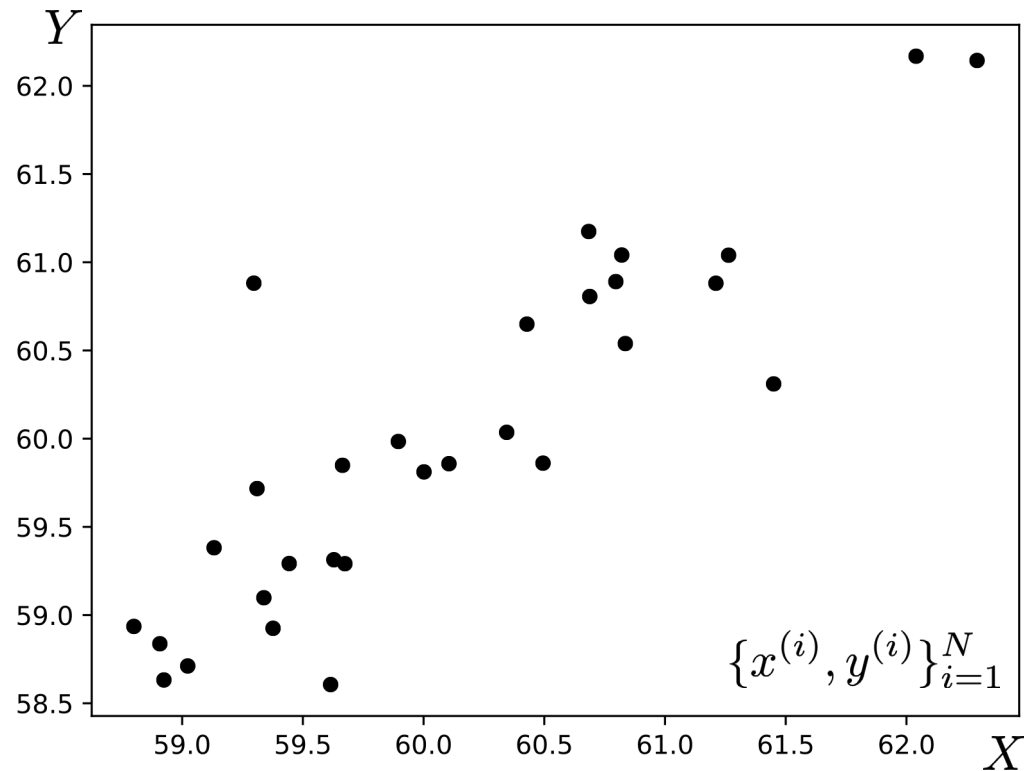
## Training data



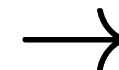
## Test query



## Training data



Learner

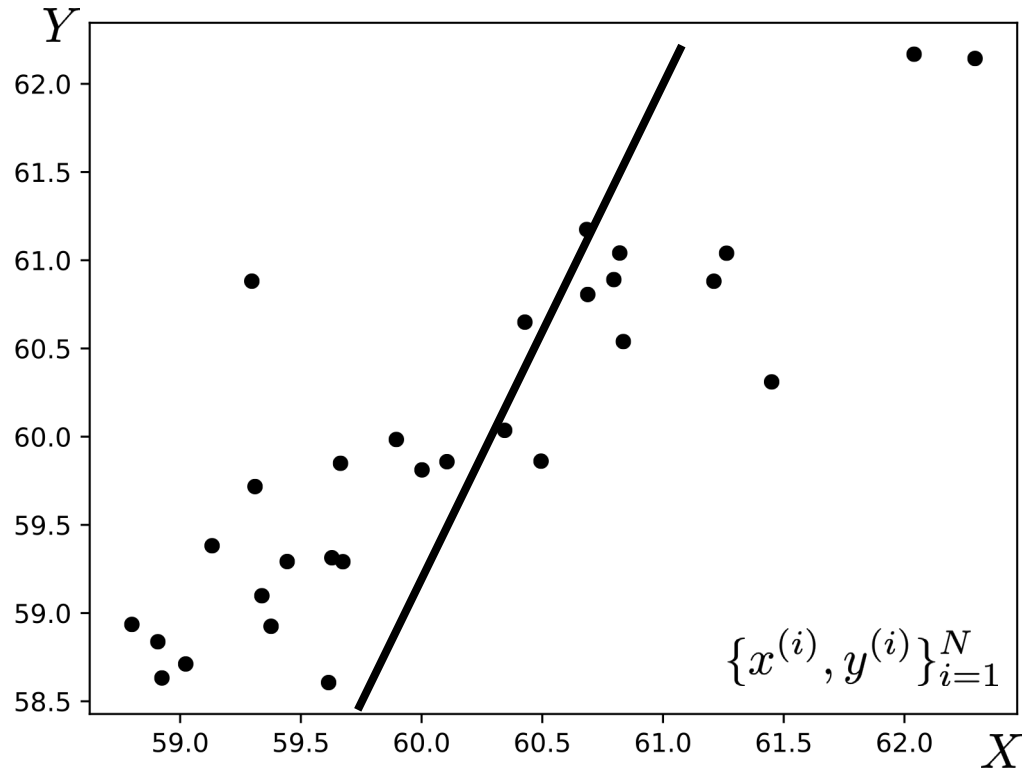


$$f_{\theta}(x) = \theta_1 x + \theta_0$$

## Hypothesis space

The relationship between X and Y is roughly linear:  $y \approx \theta_1 x + \theta_0$

# Training data

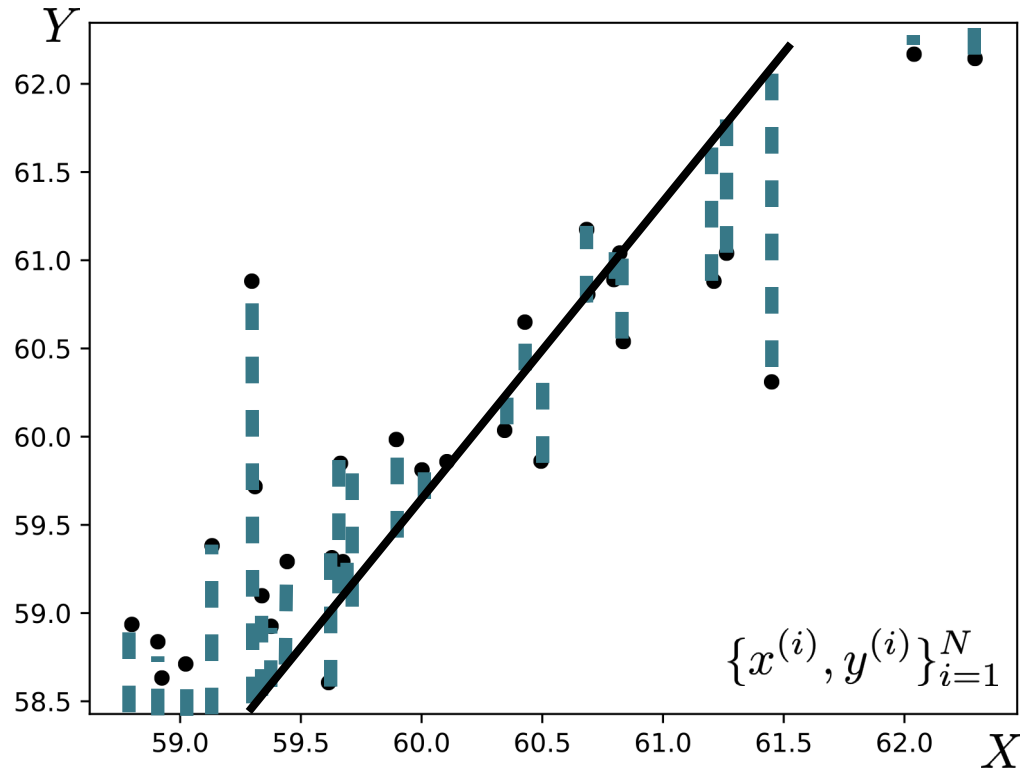


Search for the **parameters**,  $\theta = \{\theta_0, \theta_1\}$ , that best fit the data.

$$f_{\theta}(x) = \theta_1 x + \theta_0$$

**Best fit in what sense?**

# Training data



Search for the **parameters**,  $\theta = \{\theta_0, \theta_1\}$ , that best fit the data.

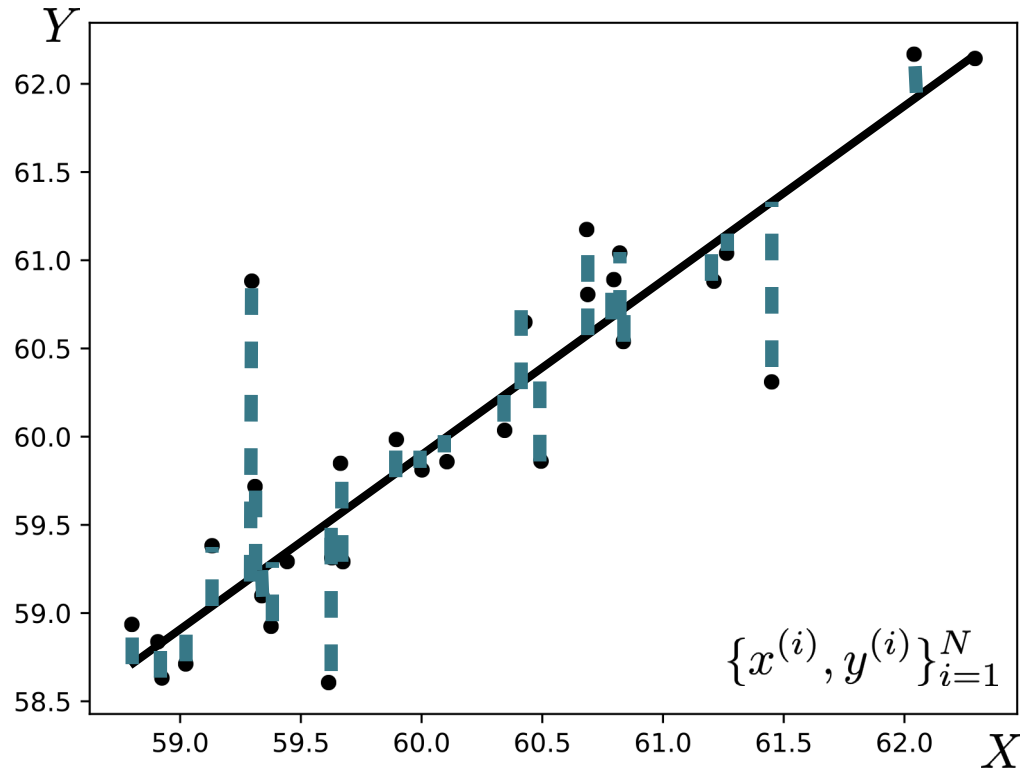
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

**Best fit in what sense?**

The least-squares **objective** (aka **loss**) says the best fit is the function that minimizes the squared error between predictions and target values:

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2 \quad \hat{y} \equiv f_{\theta}(x)$$

# Training data



Search for the **parameters**,  $\theta = \{\theta_0, \theta_1\}$ , that best fit the data.

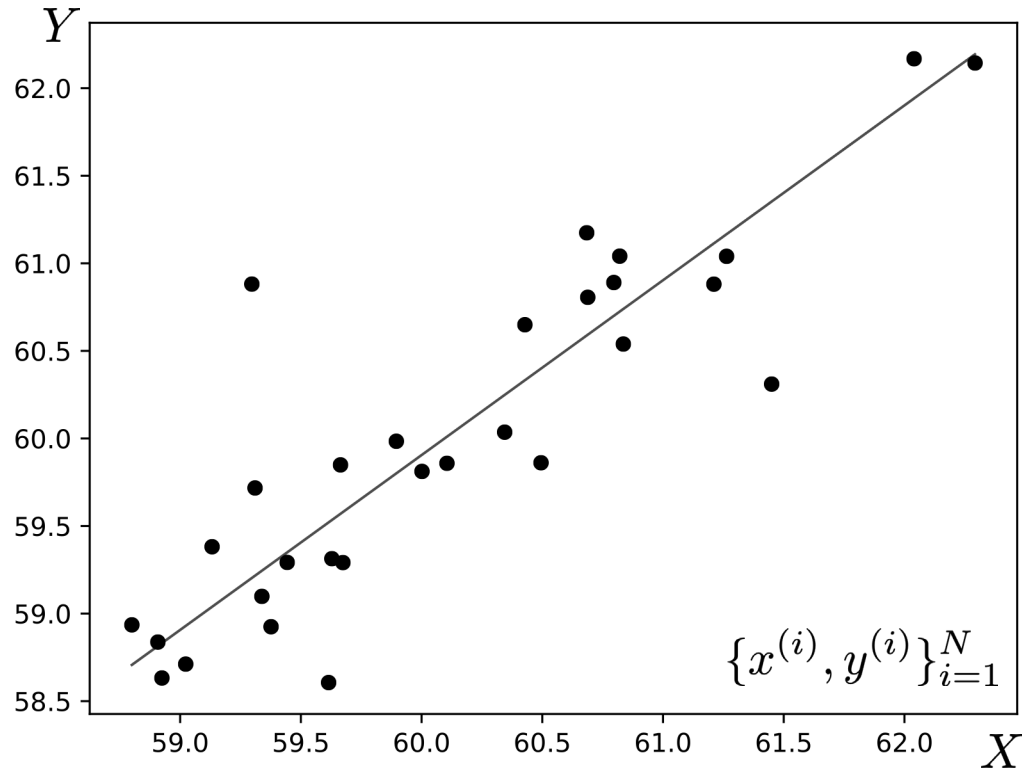
$$f_{\theta}(x) = \theta_1 x + \theta_0$$

**Best fit in what sense?**

The least-squares **objective** (aka **loss**) says the best fit is the function that minimizes the squared error between predictions and target values:

$$\mathcal{L}(\hat{y}, y) = (\hat{y} - y)^2 \quad \hat{y} \equiv f_{\theta}(x)$$

## Training data

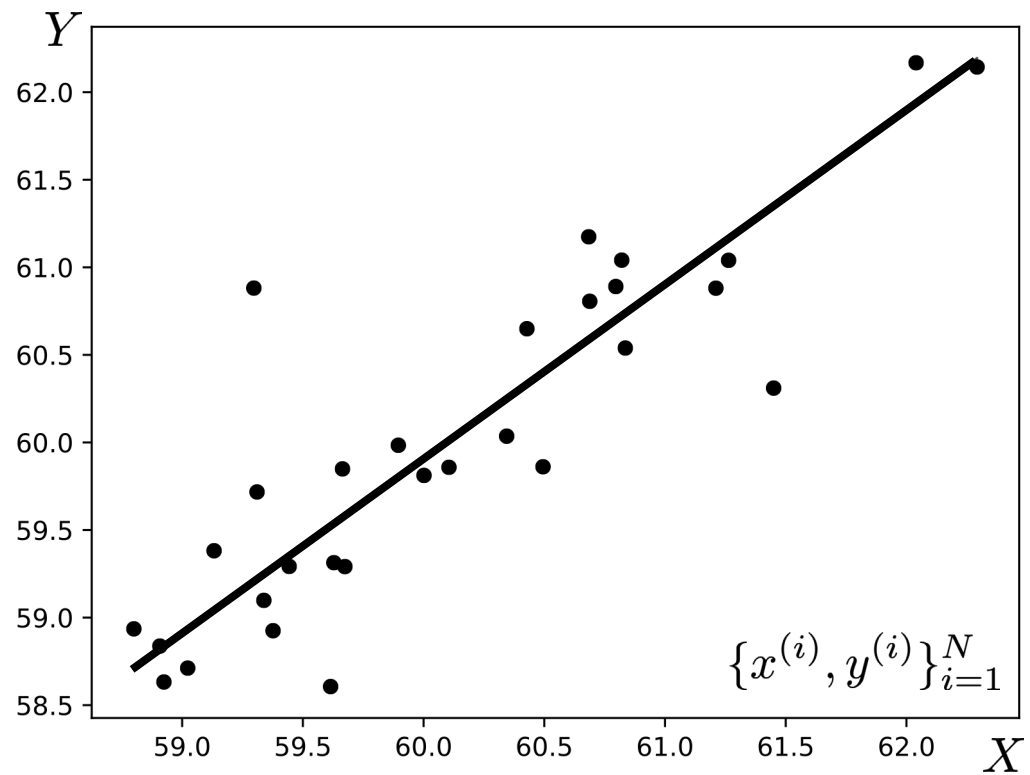


**Complete learning problem:**

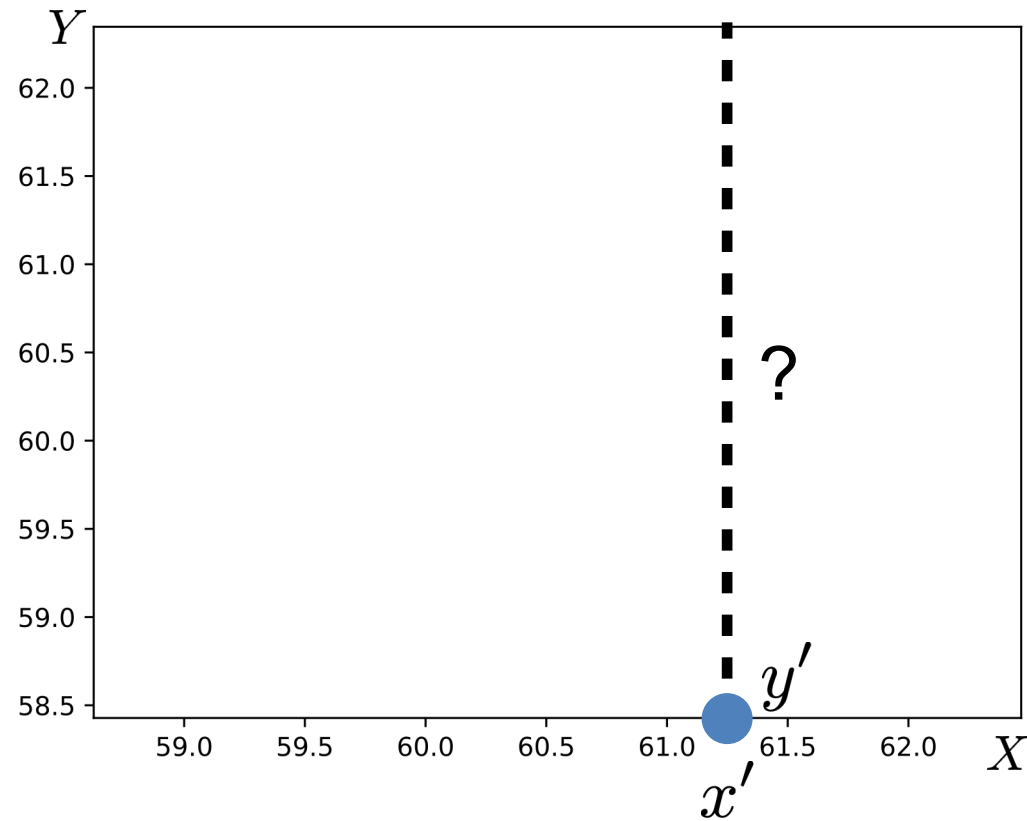
$$\begin{aligned}\theta^* &= \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_0 - y^{(i)})^2\end{aligned}$$



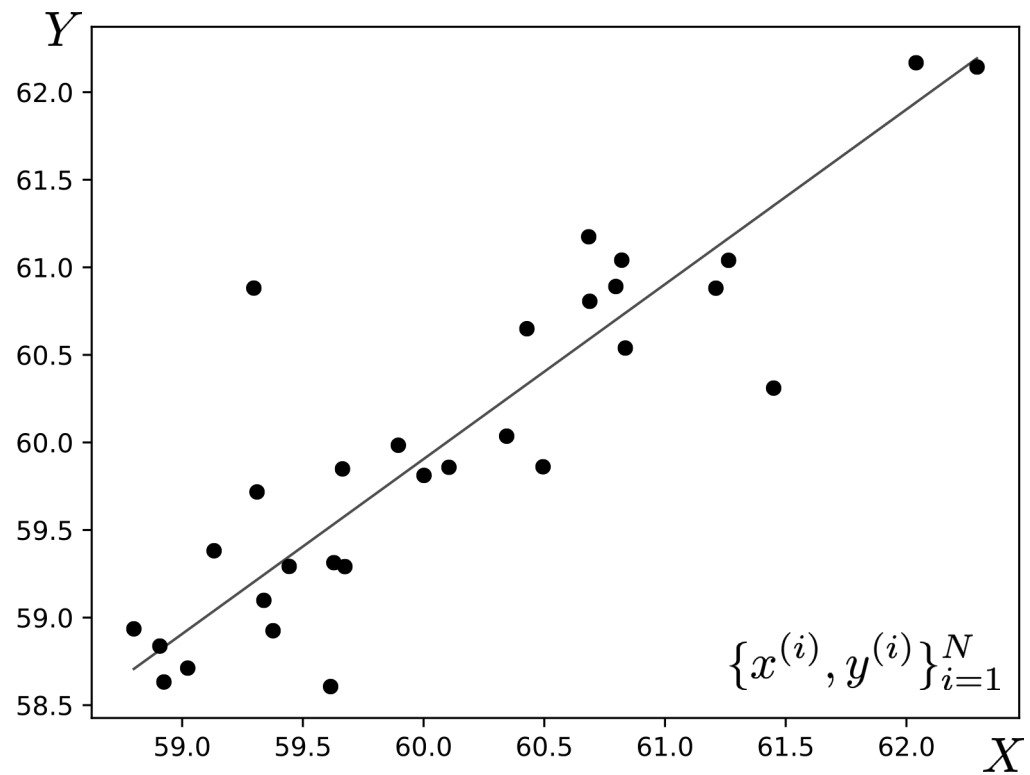
## Training data



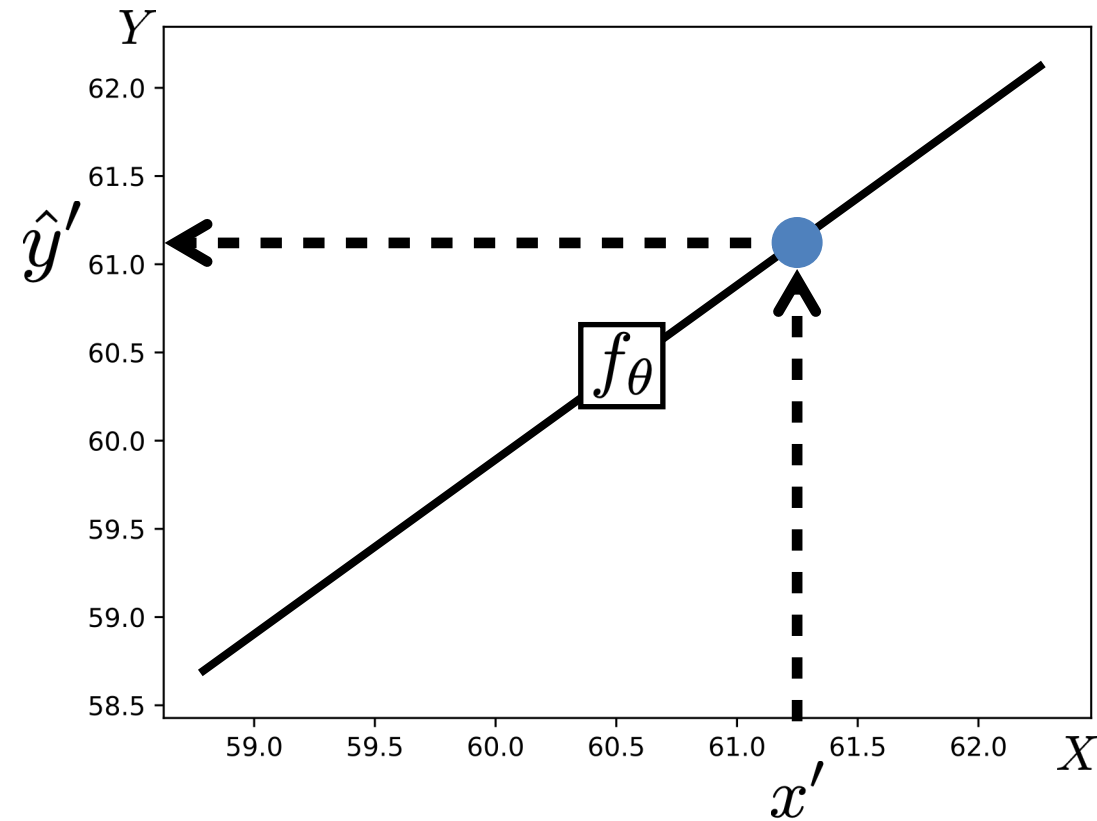
## Test query



Training data



Test query

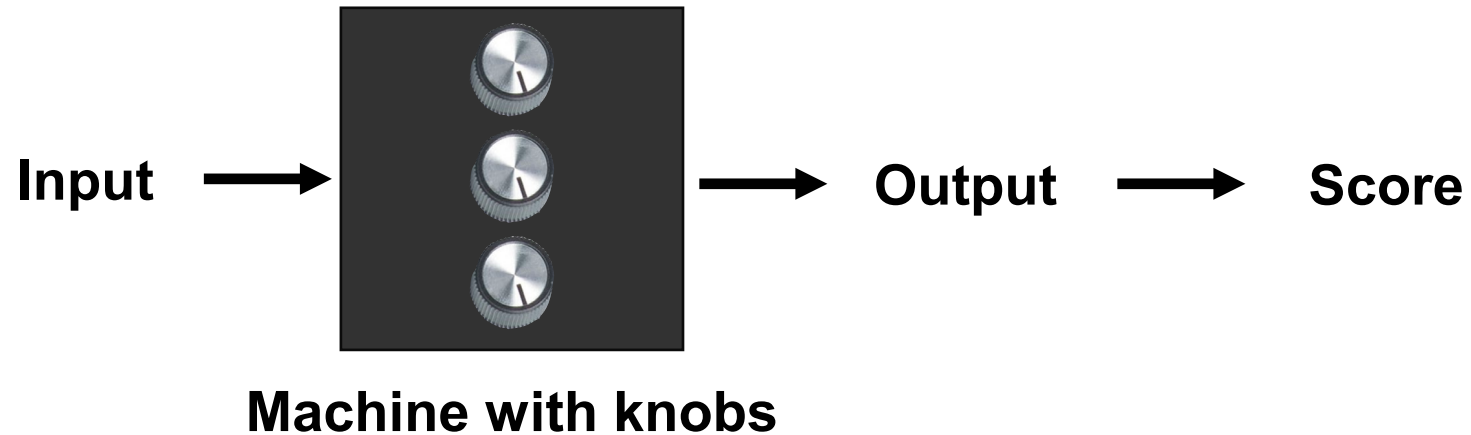


$$x' \dashrightarrow \boxed{f_\theta} \dashrightarrow \hat{y}'$$

How to minimize the objective w.r.t.  $\theta$ ?

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(x^{(i)}) - y^{(i)})^2$$

Use an **optimizer**!



# How to minimize the objective w.r.t. $\theta$ ?

In the linear case:

Learning problem

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_0 - y^{(i)})^2$$

$$\begin{aligned} J(\theta) &= \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_0 - y^{(i)})^2 \\ &= (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta) \end{aligned}$$

$$\mathbf{X} = \begin{pmatrix} x^{(1)} & 1 \\ x^{(2)} & 1 \\ \vdots & \vdots \\ x^{(N)} & 1 \end{pmatrix} \quad \theta = (\theta_1 \quad \theta_0) \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{pmatrix}$$

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

$$\frac{\partial J(\theta)}{\partial \theta} = 0$$

$$\frac{\partial J(\theta)}{\partial \theta} = 2(\mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y})$$

$$2(\mathbf{X}^T \mathbf{X} \theta^* - \mathbf{X}^T \mathbf{y}) = 0$$

$$\mathbf{X}^T \mathbf{X} \theta^* = \mathbf{X}^T \mathbf{y}$$

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Solution

# Empirical Risk Minimization

(formalization of supervised learning)

Linear least squares learning problem

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_0 - y^{(i)})^2$$

# Empirical Risk Minimization

(formalization of supervised learning)

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

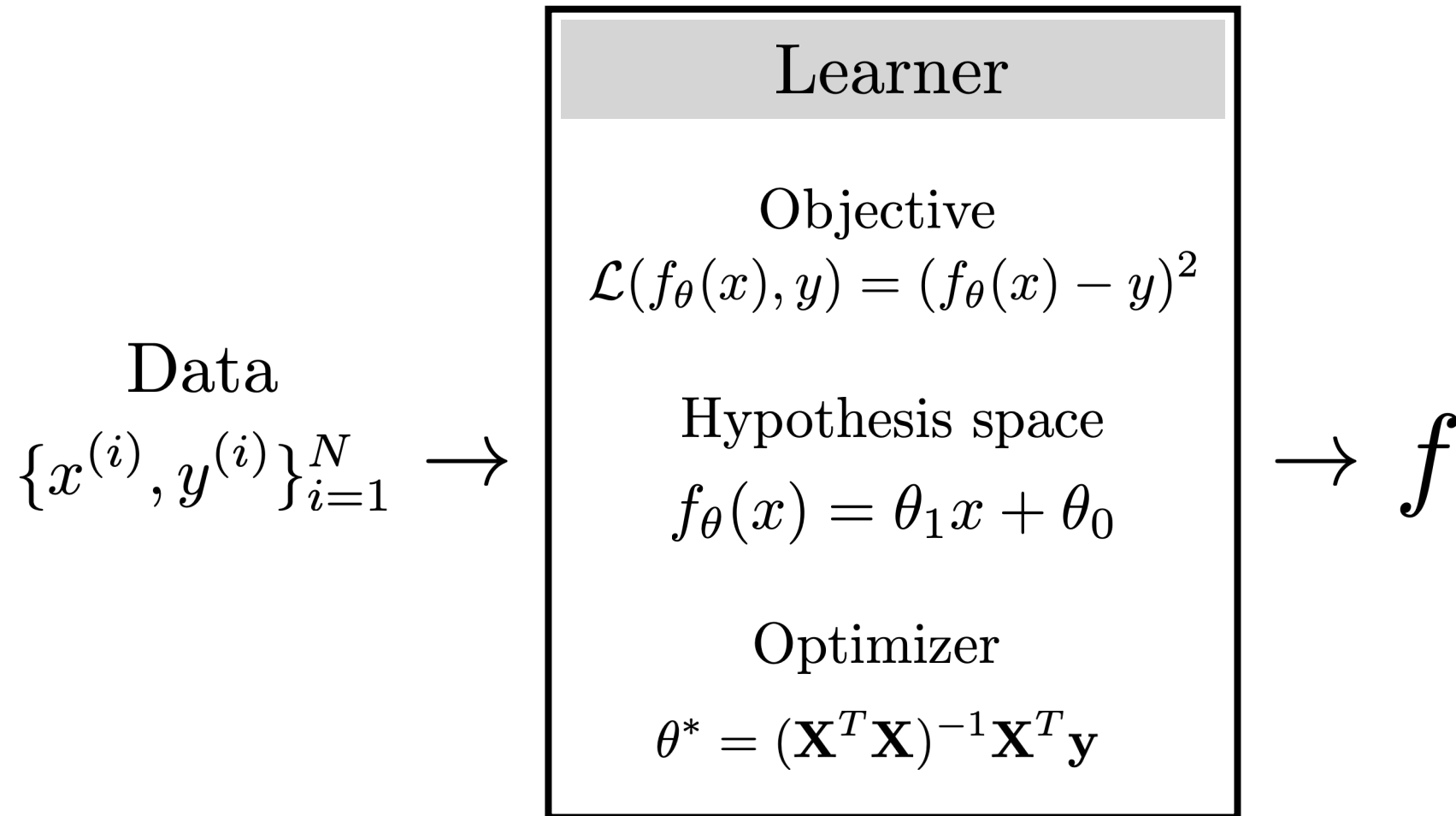
Objective function  
(loss)

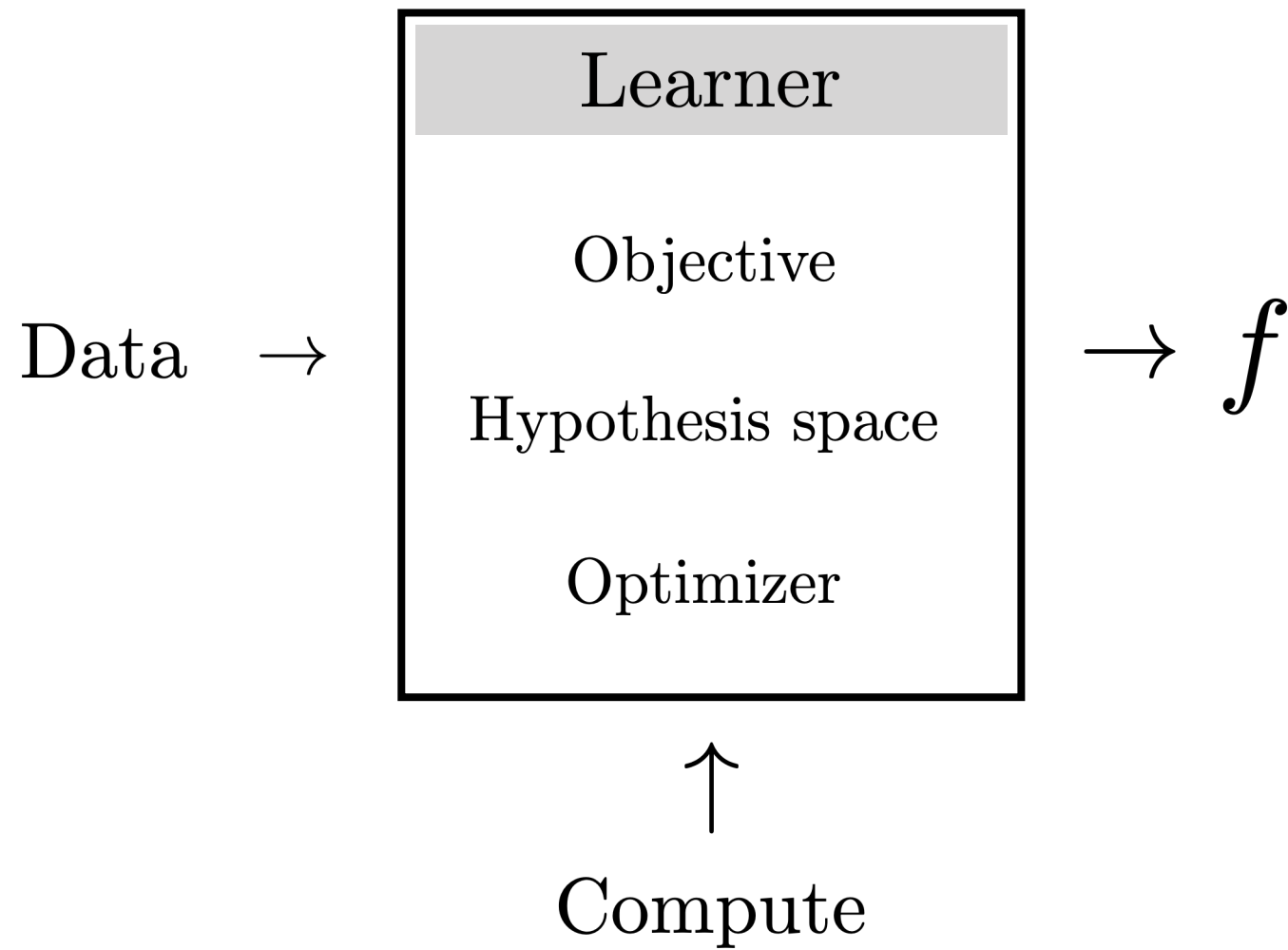
Hypothesis space

Training data

The diagram shows the formula for Empirical Risk Minimization. The text 'Objective function (loss)' has an arrow pointing to the loss function symbol  $\mathcal{L}$ . The text 'Hypothesis space' has an arrow pointing to the set  $\mathcal{F}$  in the minimization argument. The text 'Training data' has two arrows pointing to the input  $\mathbf{x}^{(i)}$  and output  $\mathbf{y}^{(i)}$  of the loss function.

# Case study #1: Linear least squares

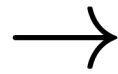
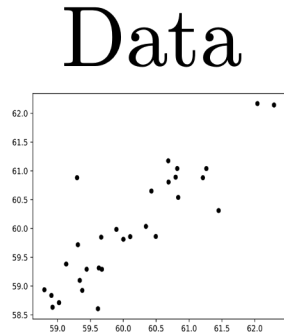




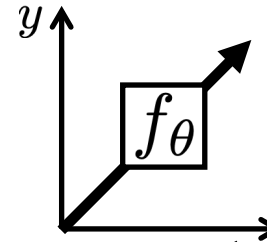
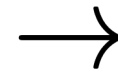


# Example 1: Linear least squares

Training

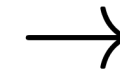
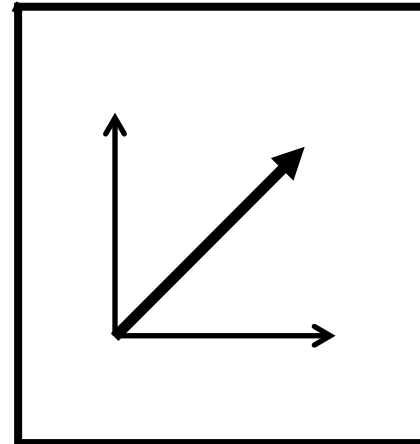
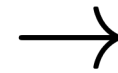


Learner



Testing

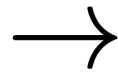
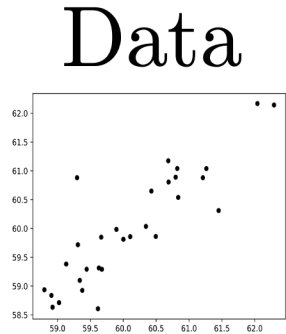
Input



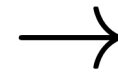
Output

# Example 2: Program induction

Training



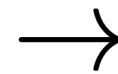
Learner



```
def predict(x):  
    y = 0.8*x + 2  
    return y
```

Testing

Input

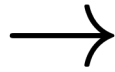
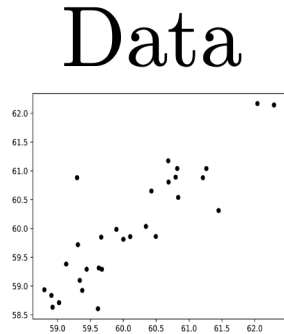


```
def predict(x):  
    y = 0.8*x + 2  
    return y
```

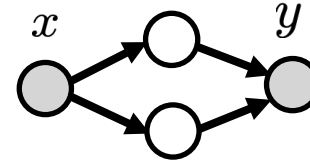
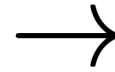
→ Output

# Example 3: Deep learning

Training

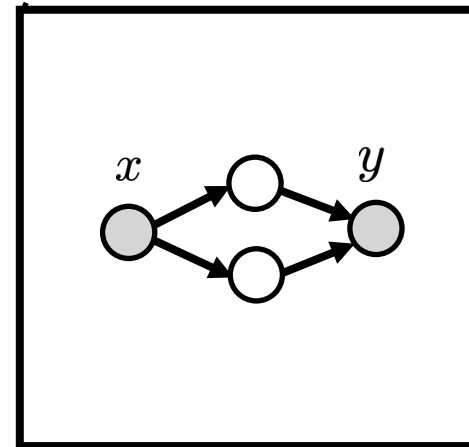
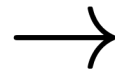


Learner



Testing

Input



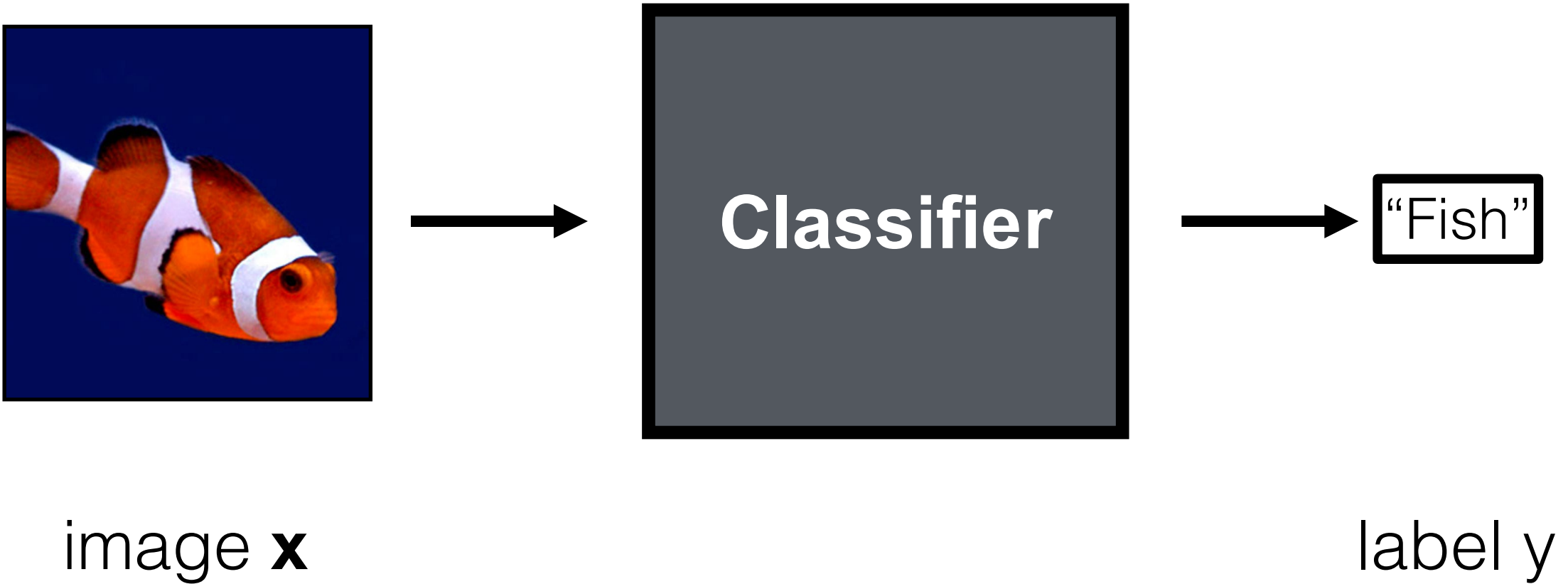
→ Output

# Learning for vision

Big questions:

1. How do you represent the input and output?
2. What is the objective?
3. What is the hypothesis space? (e.g., linear, polynomial, neural net?)
4. How do you optimize? (e.g., gradient descent, Newton's method?)
5. What data do you train on?

# Image classification



# Image classification



image **x**



"Fish"

label **y**



# Image classification



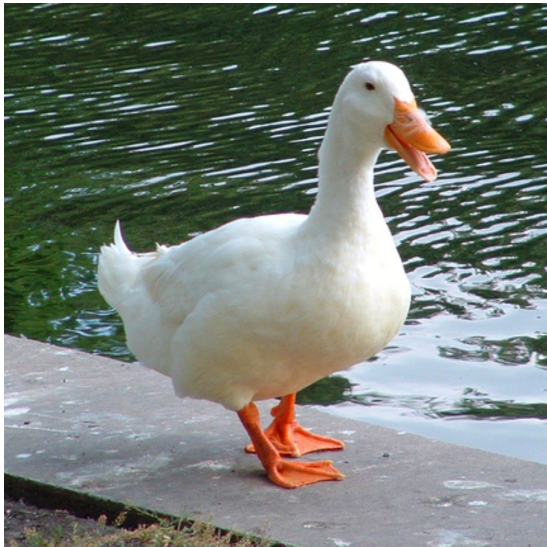
image **x**



"Fish"

label **y**

# Image classification



+

image **x**



**Classifier**



“Duck”

label **y**

$\mathbf{x}$



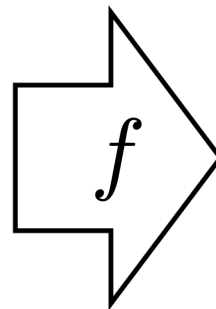
*Training data*

$\mathbf{x}$

$y$



+



$y$

"Fish"

$$\arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}^{(i)}), y^{(i)})$$

# How to represent class labels?

## One-hot vector

*Training data*

$\mathbf{x}$

$y$

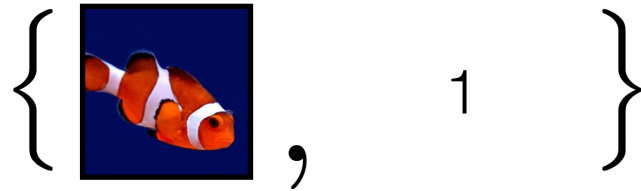


$\vdots$

*Training data*

$\mathbf{x}$

$y$

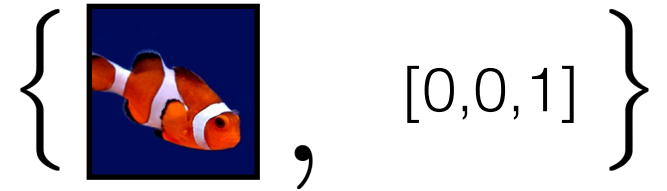


$\vdots$

*Training data*

$\mathbf{x}$

$y$



$\vdots$

# What should the loss be?

**0-1 loss** (number of misclassifications)

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = \mathbb{1}(\hat{\mathbf{y}} \neq \mathbf{y}) \quad \leftarrow \text{discrete, NP-hard to optimize!}$$

**Cross entropy**

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad \leftarrow \begin{array}{l} \text{continuous,} \\ \text{differentiable,} \\ \text{convex} \end{array}$$

Ground truth label    $y$

$x$

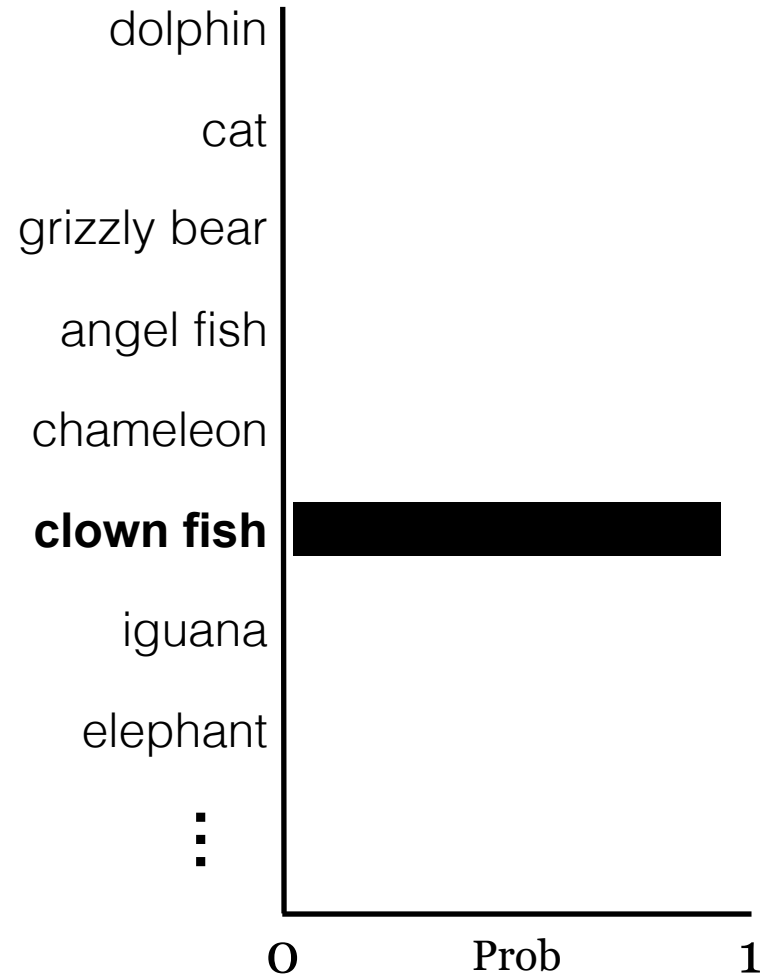


$[0,0,0,0,0,1,0,0,\dots]$

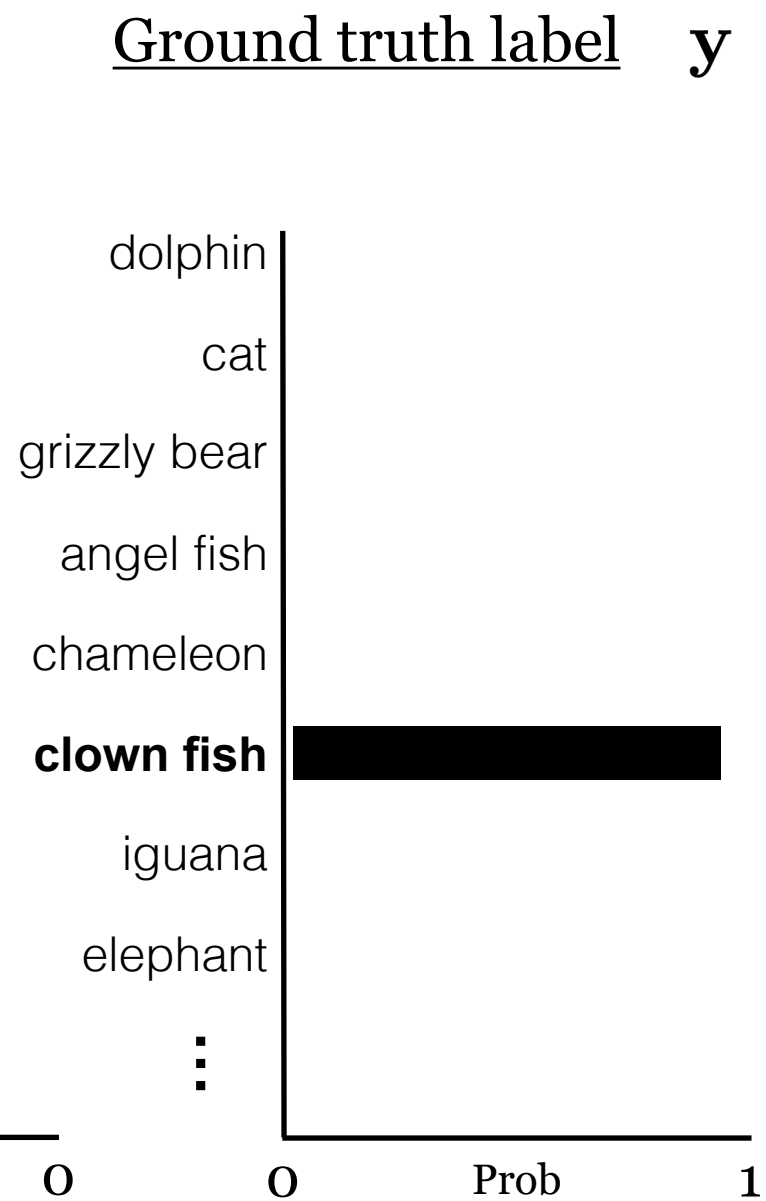
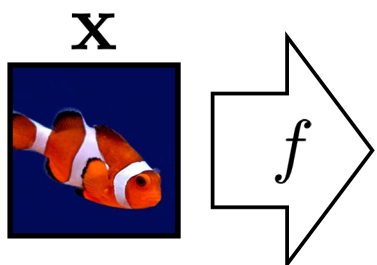


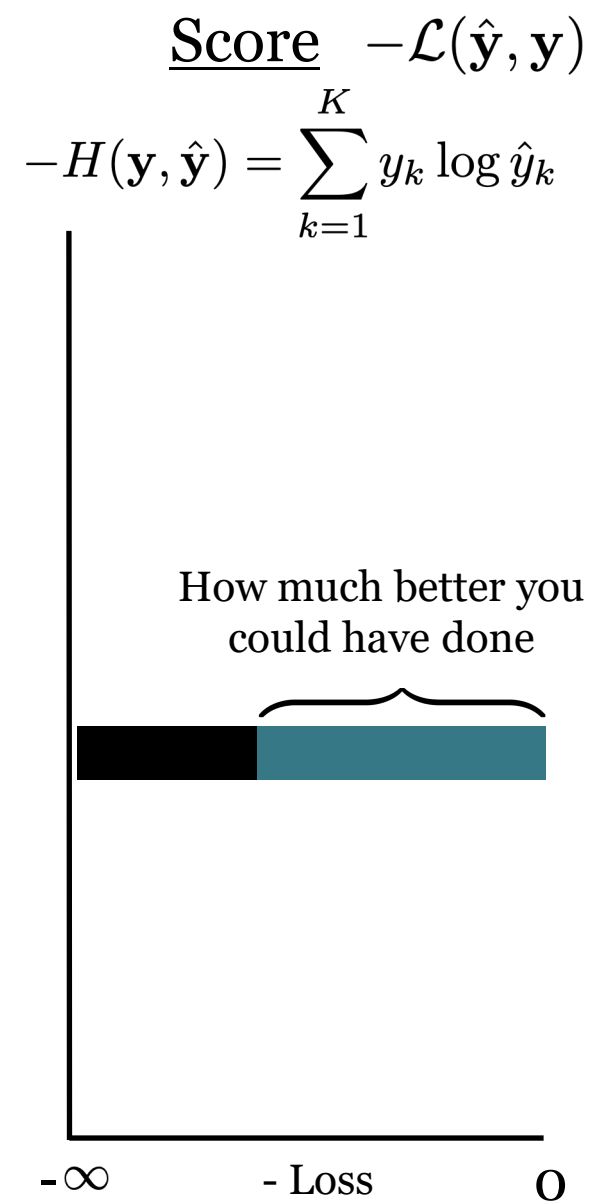
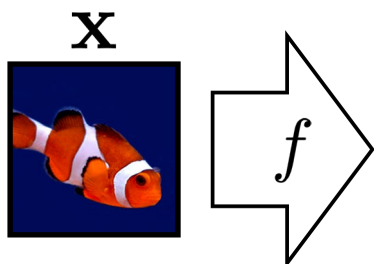


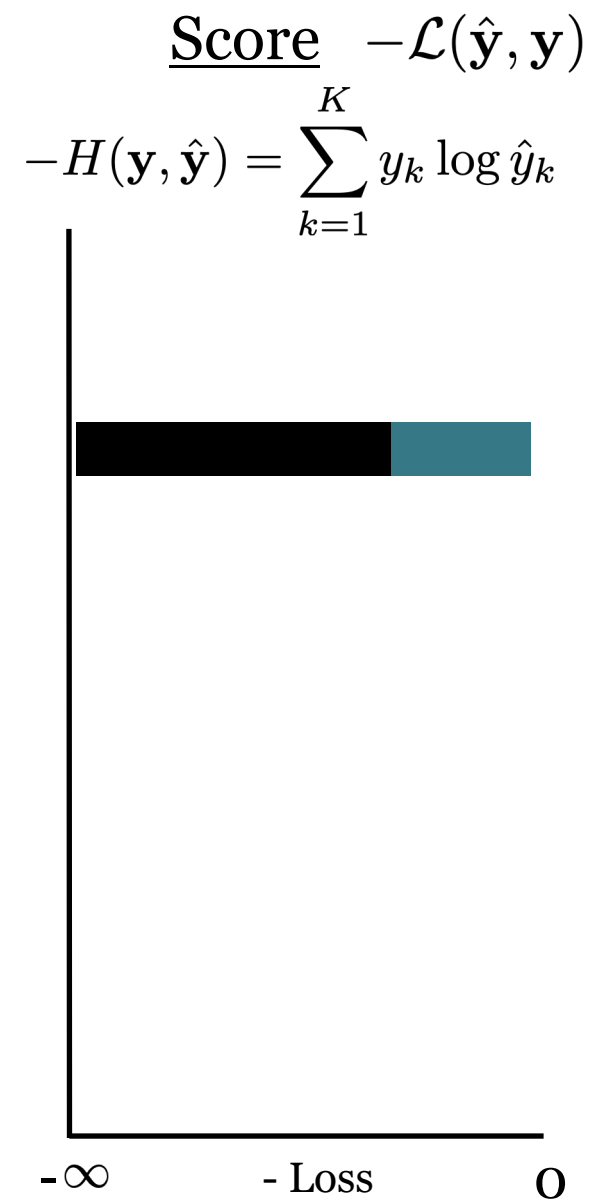
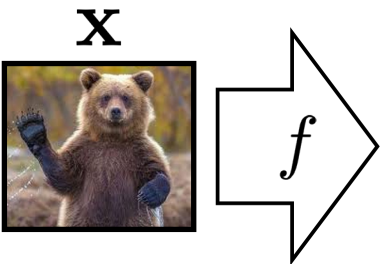
Ground truth label    **y**

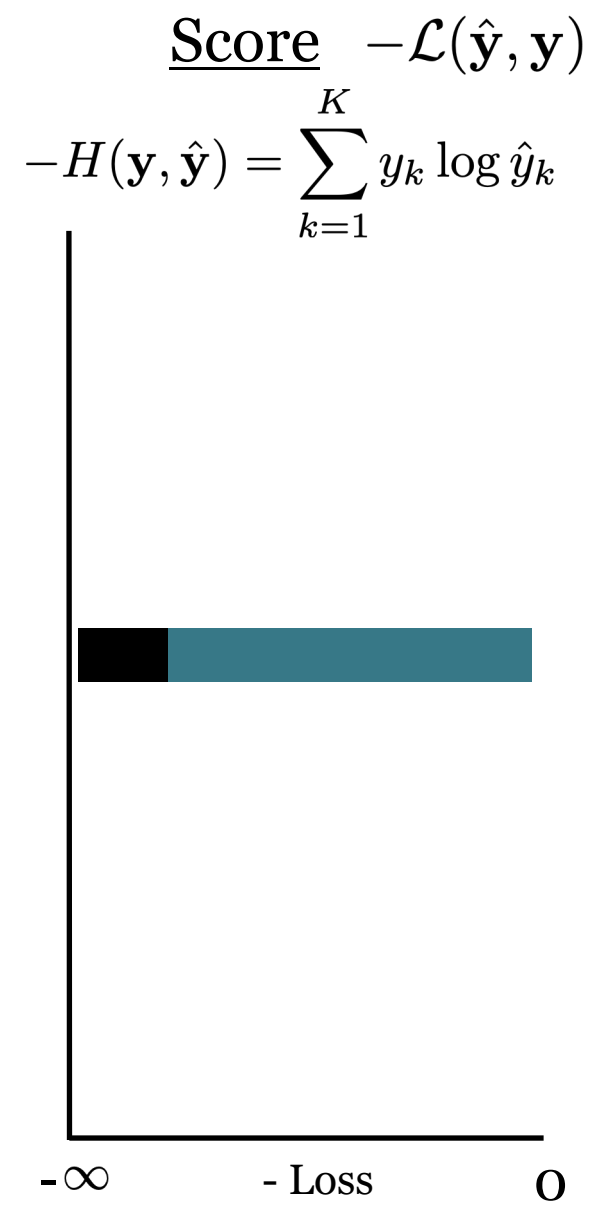
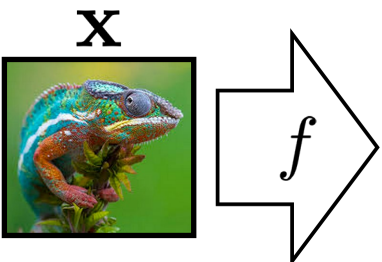


$\infty$









# Softmax regression (a.k.a. multinomial logistic regression)

$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

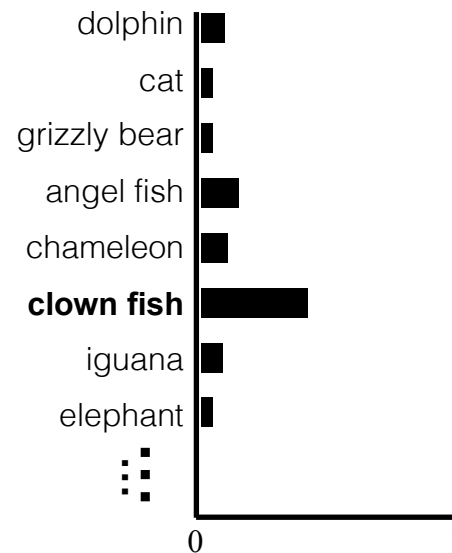
← **logits**: vector of K scores, one for each class

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

← squash into a non-negative vector that sums to 1  
— i.e. **a probability mass function!**

$$\hat{y}_j = \frac{e^{-z_j}}{\sum_{k=1}^K e^{-z_k}}$$

$\hat{\mathbf{y}} =$





## Softmax regression (a.k.a. multinomial logistic regression)

Probabilistic interpretation:

$$\hat{\mathbf{y}} \equiv [P_{\theta}(Y = 1|X = \mathbf{x}), \dots, P_{\theta}(Y = K|X = \mathbf{x})] \quad \leftarrow \text{predicted probability of each class given input } \mathbf{x}$$

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad \leftarrow \text{picks out the -log likelihood of the ground truth class } \mathbf{y} \text{ under the model prediction } \hat{\mathbf{y}}$$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N H(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) \quad \leftarrow \text{max likelihood learner!}$$

## Softmax regression (a.k.a. multinomial logistic regression)

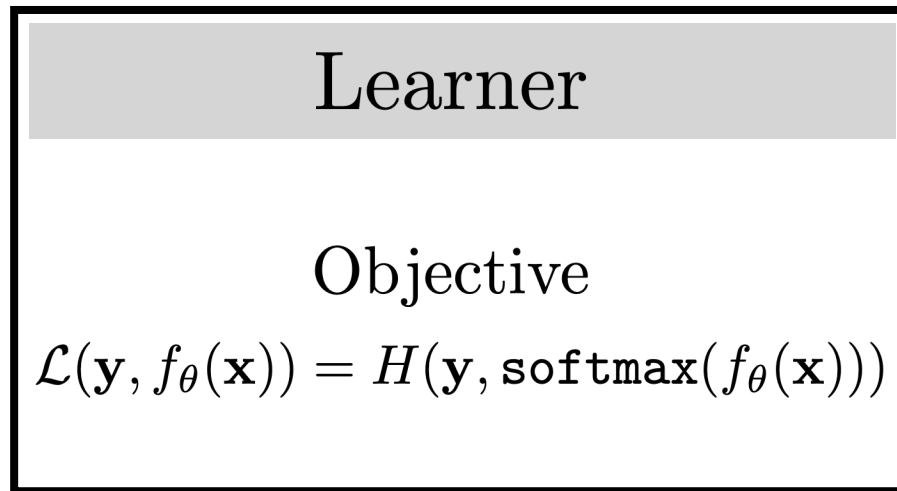
$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

$$\mathbf{z} = f_{\theta}(\mathbf{x})$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

Data

$$\{x^{(i)}, y^{(i)}\}_{i=1}^N \rightarrow$$



$$\rightarrow f$$

# Generalization

“The central challenge in machine learning is that our algorithm must perform well on new, previously unseen inputs—not just those on which our model was trained. The ability to perform well on previously unobserved inputs is called **generalization**.

... [this is what] separates machine learning from optimization.”

— Deep Learning textbook (Goodfellow et al.)

training domain

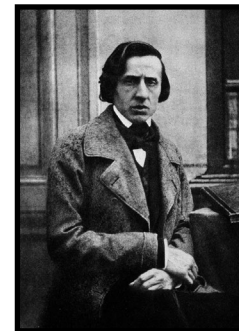
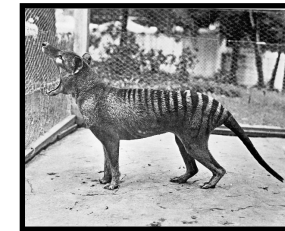
testing domain

(where we actually use our model)

**Domain gap** between  $p_{\text{train}}$  and  $p_{\text{test}}$   
will cause us to fail to generalize.

Space of natural images

Training data



Test data

People telling me AI is going  
to destroy the world

My neural network

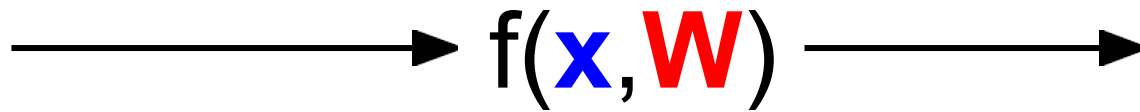


# Parametric Approach

Image



Array of **32x32x3** numbers  
(3072 numbers total)



**10** numbers giving  
class scores

↑  
**W**

parameters  
or weights

# Parametric Approach: Linear Classifier

$$f(x, W) = Wx$$

Image



Array of **32x32x3** numbers  
(3072 numbers total)



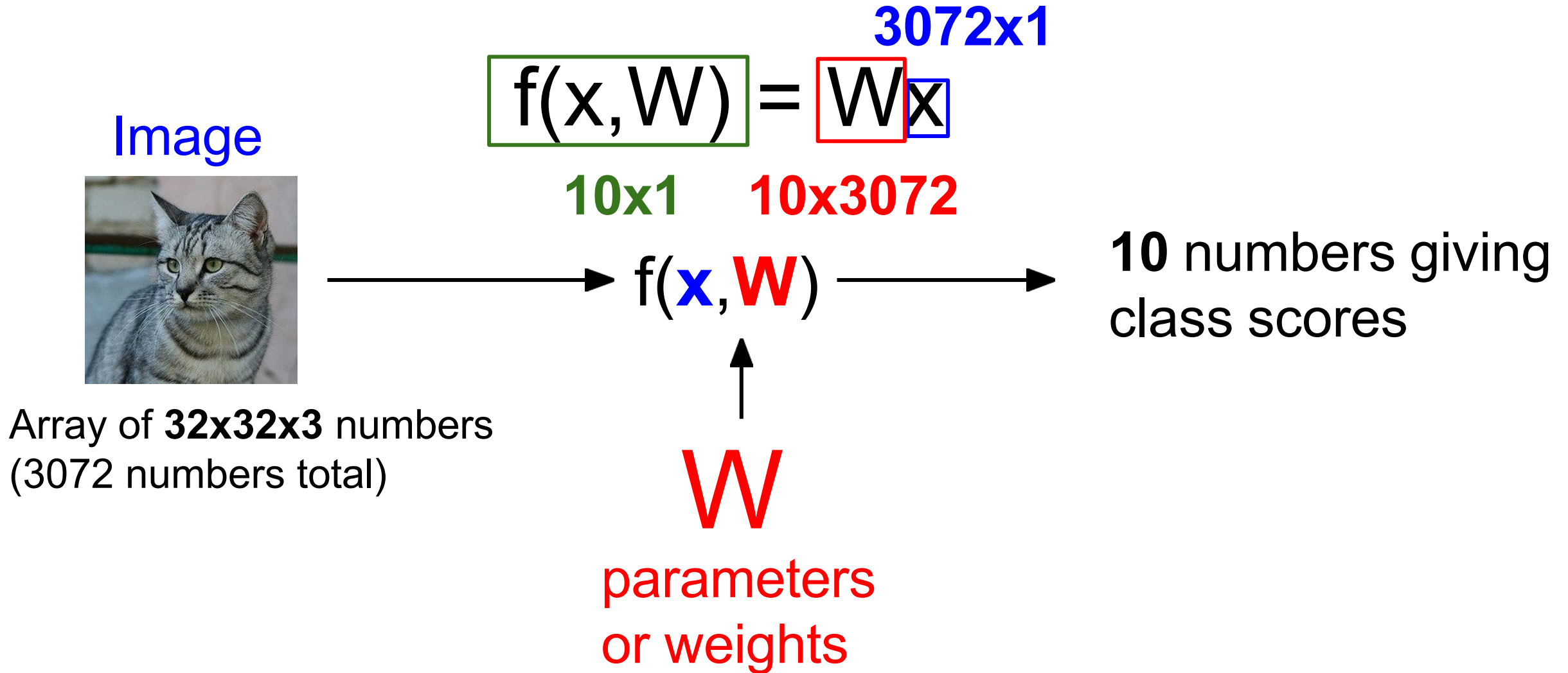
**10** numbers giving  
class scores

**W**

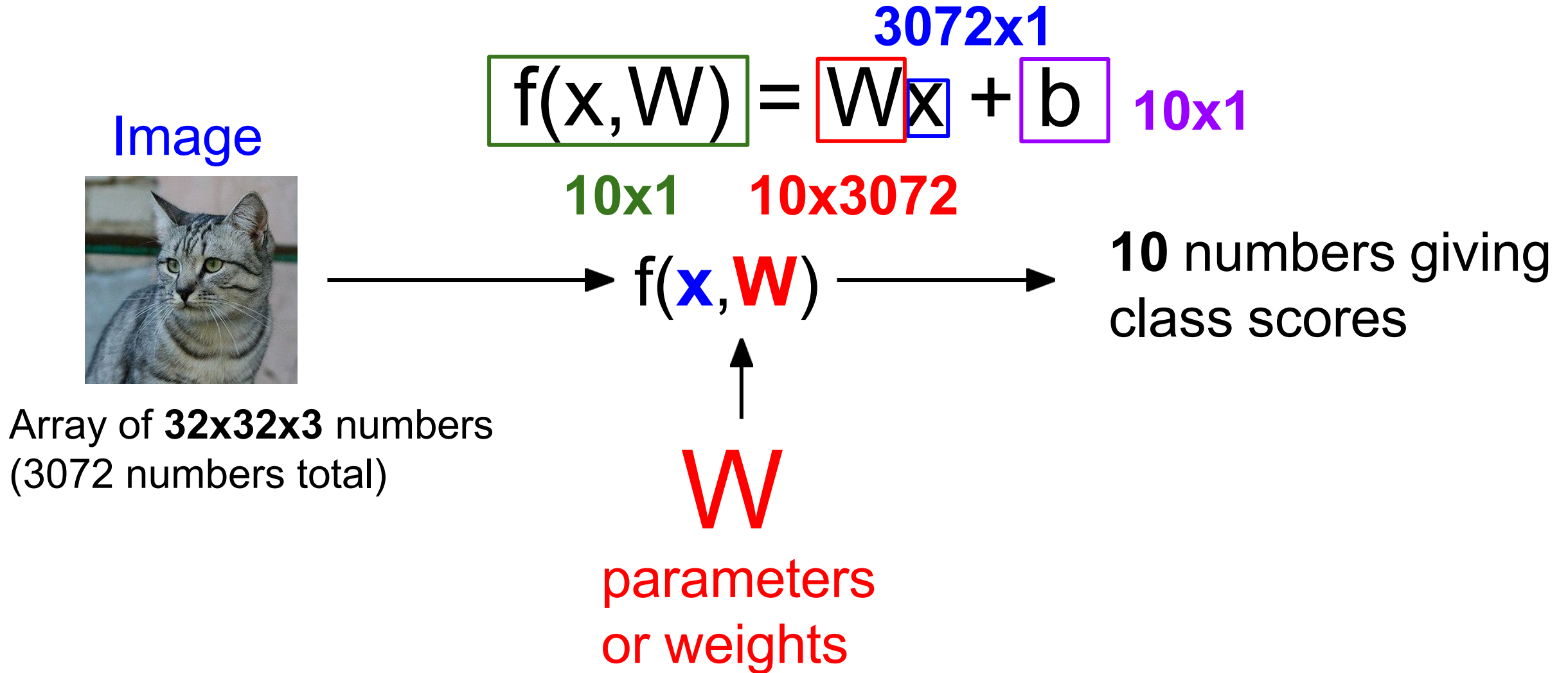
parameters  
or weights



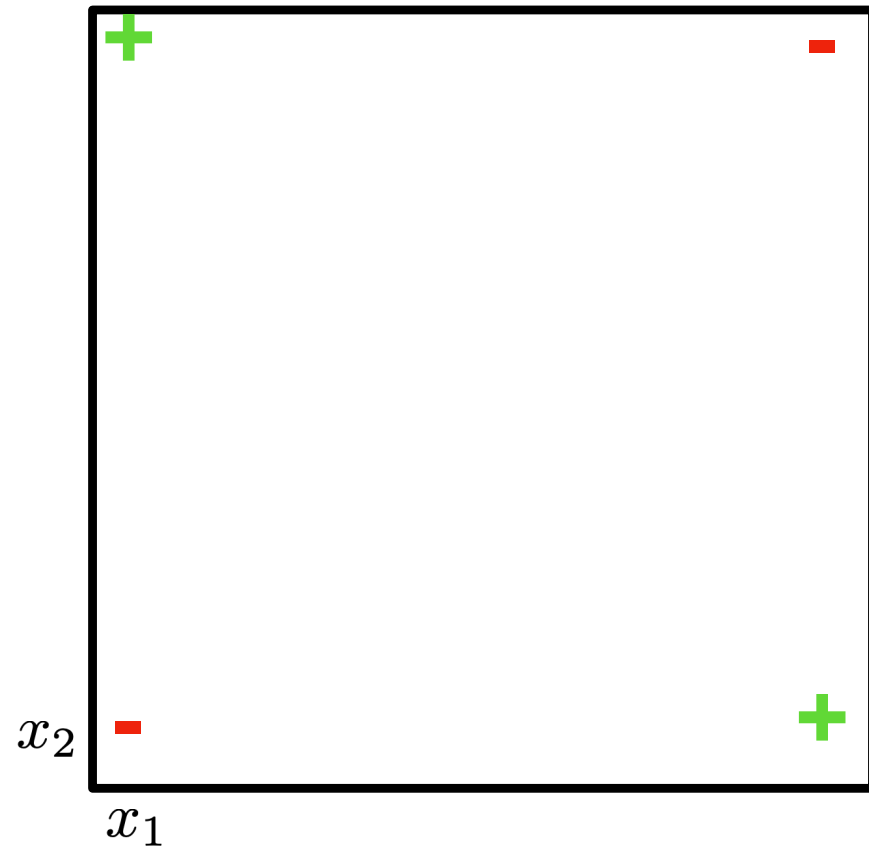
# Parametric Approach: Linear Classifier



# Parametric Approach: Linear Classifier



# Limitations to linear classifiers



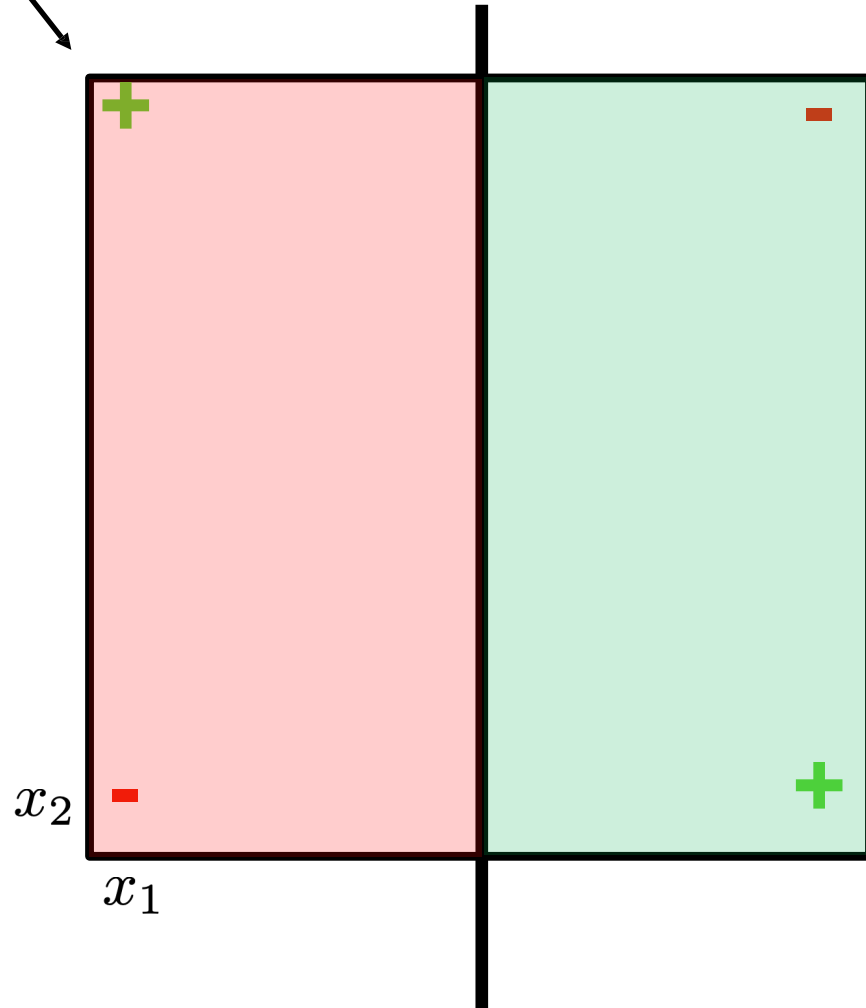
		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

# Limitations to linear classifiers

Wrong!

Wrong!

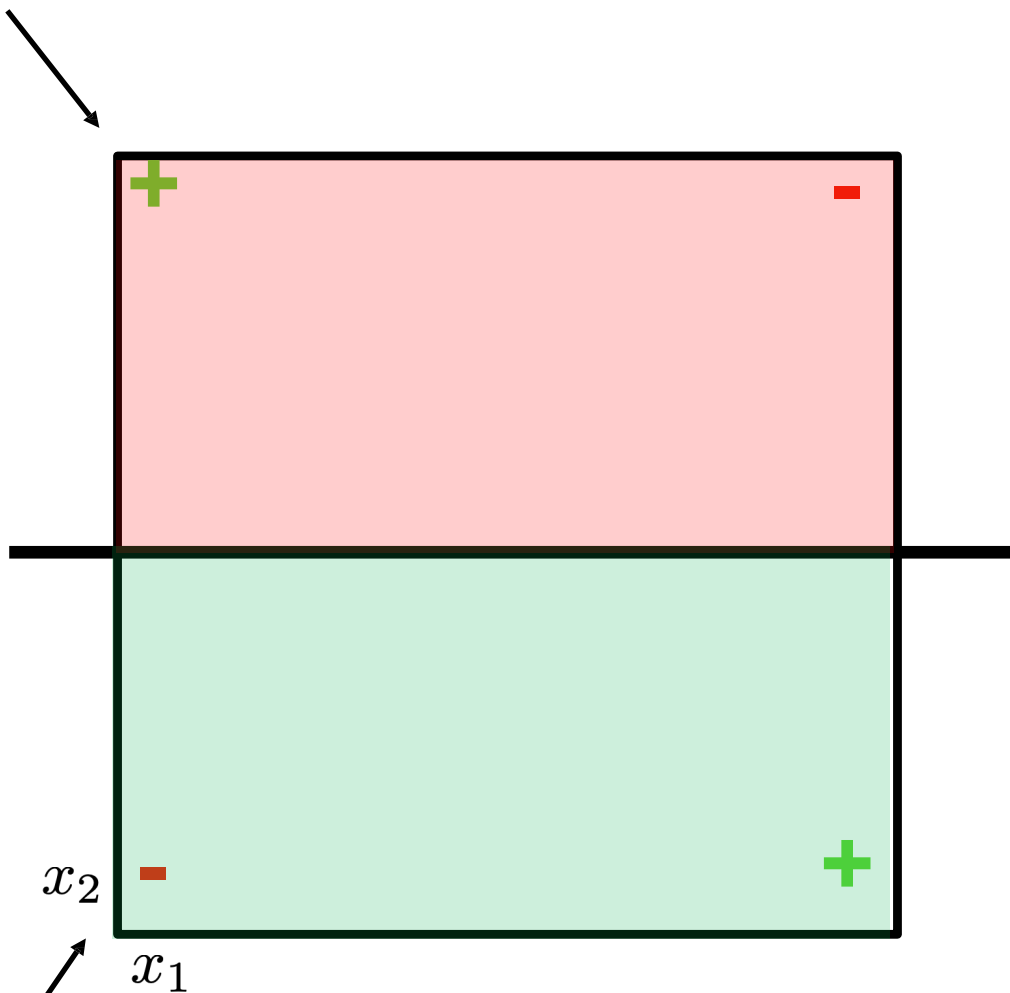


		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

# Limitations to linear classifiers

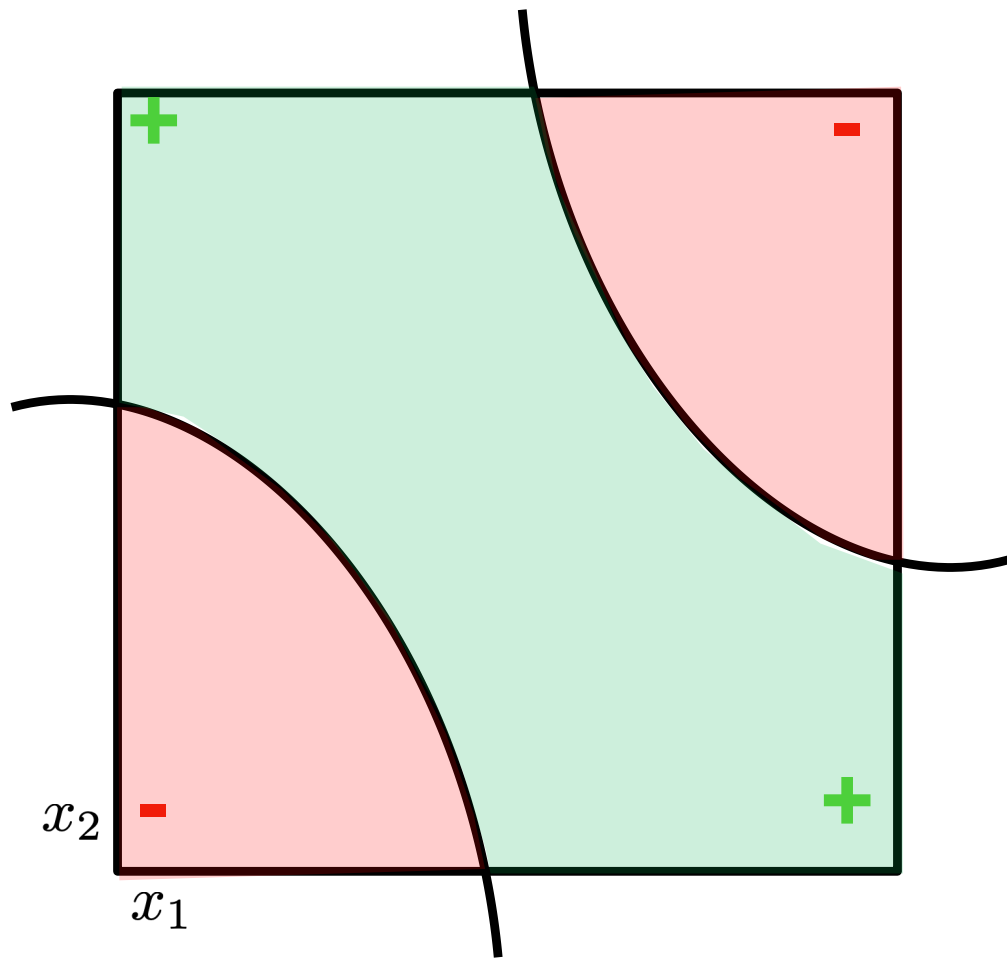
Wrong!



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

XOR

# Goal: Non-linear decision boundary

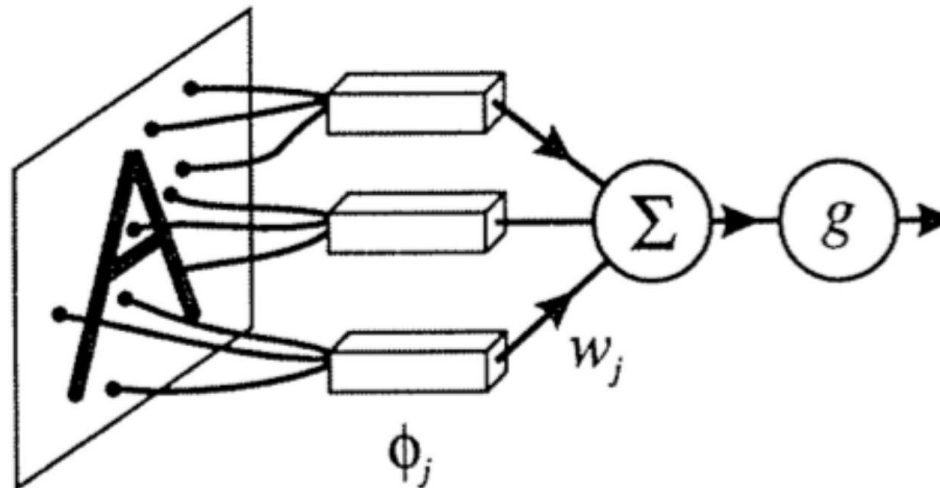


		$x_2$	
		0	1
$x_1$	0	0	1
	1	1	0

XOR

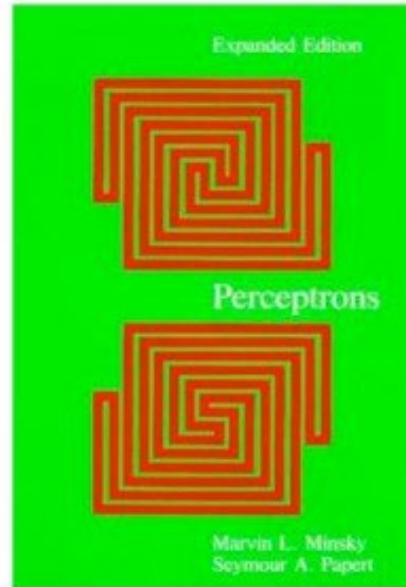
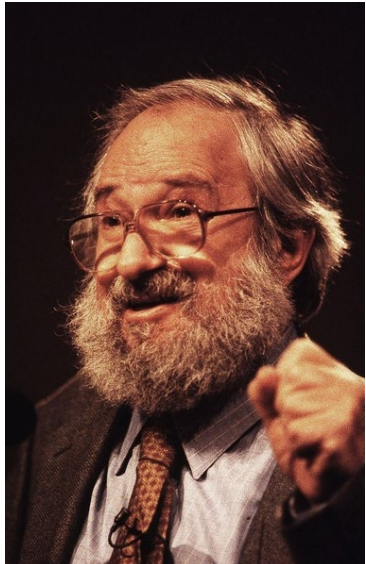
# Perceptron

- 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

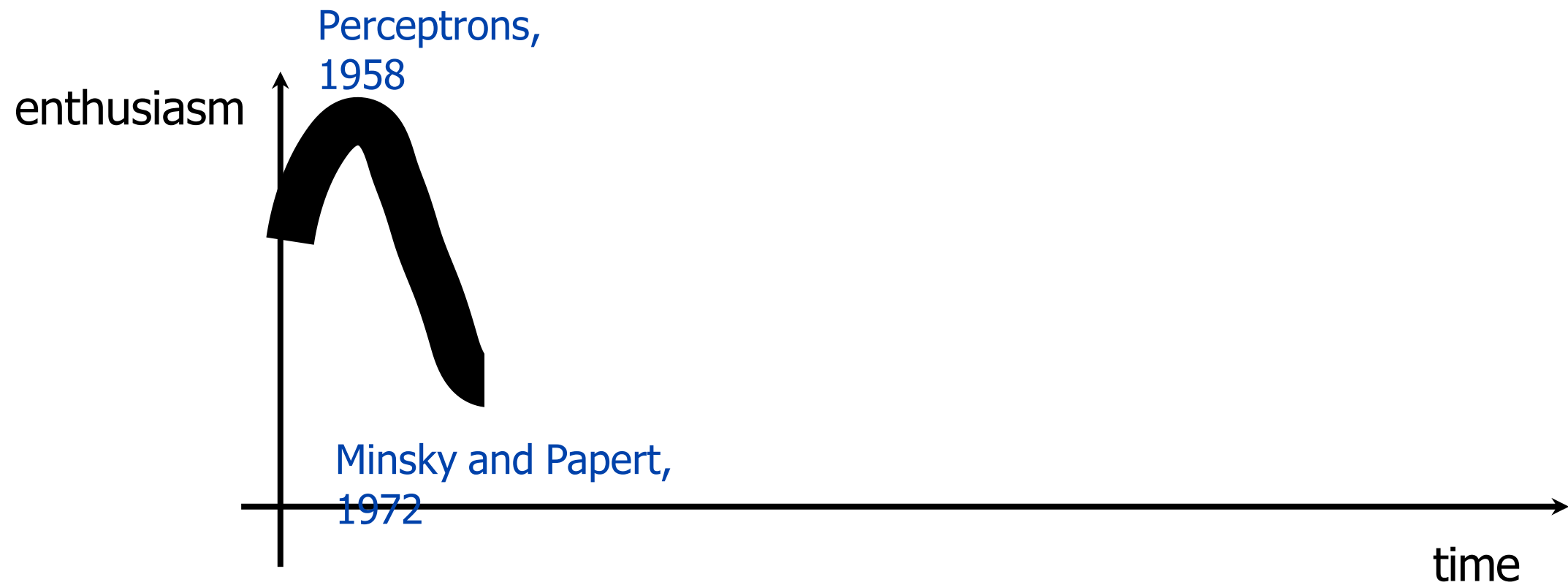
By [Marvin Minsky](#) and [Seymour A. Papert](#)

### 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.

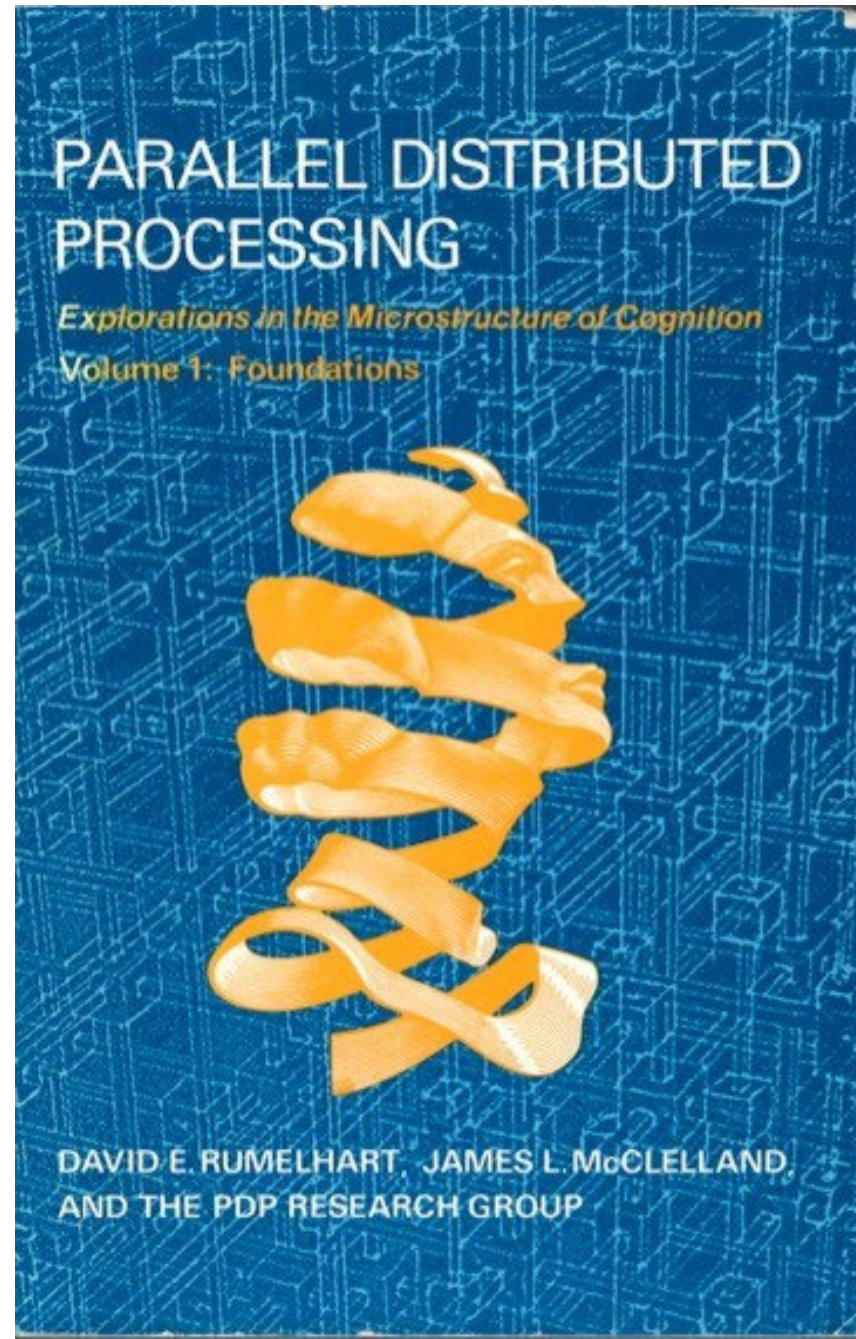
Artificial-intelligence research, which for a time concentrated on the programming of ton Neumann computers, is swinging back to the idea that intelligence might emerge from the activity of networks of neuronlike entities. Minsky and Papert's book was the first example of a mathematical analysis carried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

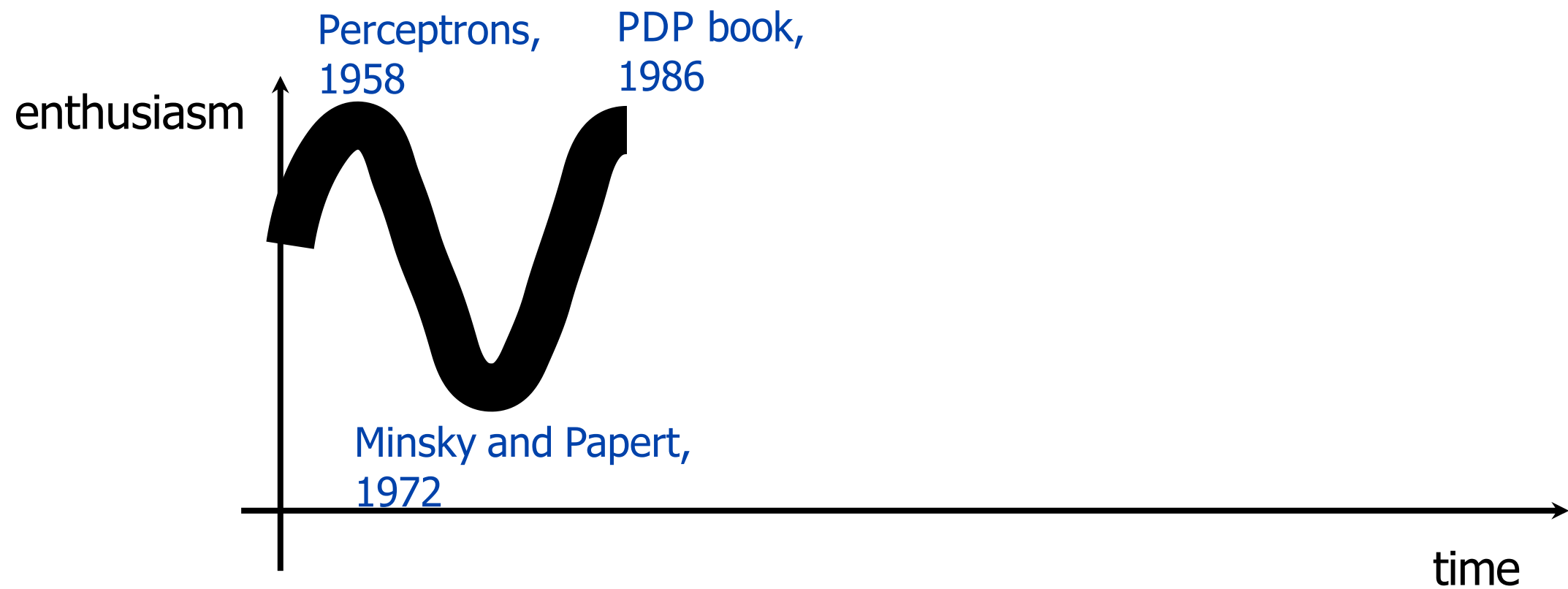




# Parallel Distributed Processing (PDP), 1986



Source: Isola, Torralba, Freeman



# LeCun convolutional neural networks

PROC. OF THE IEEE, NOVEMBER 1998

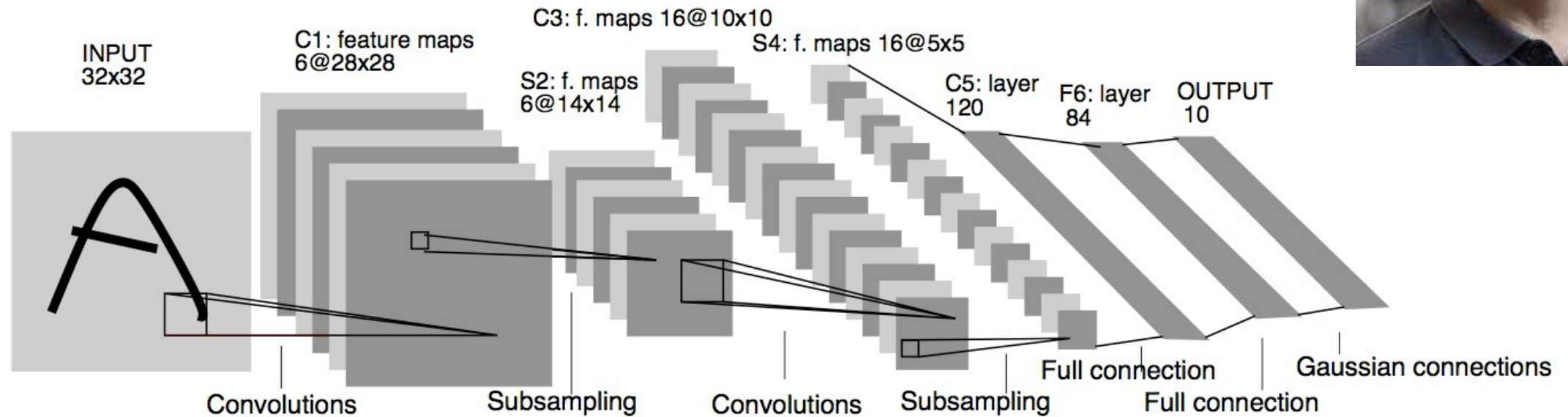


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Demos:

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

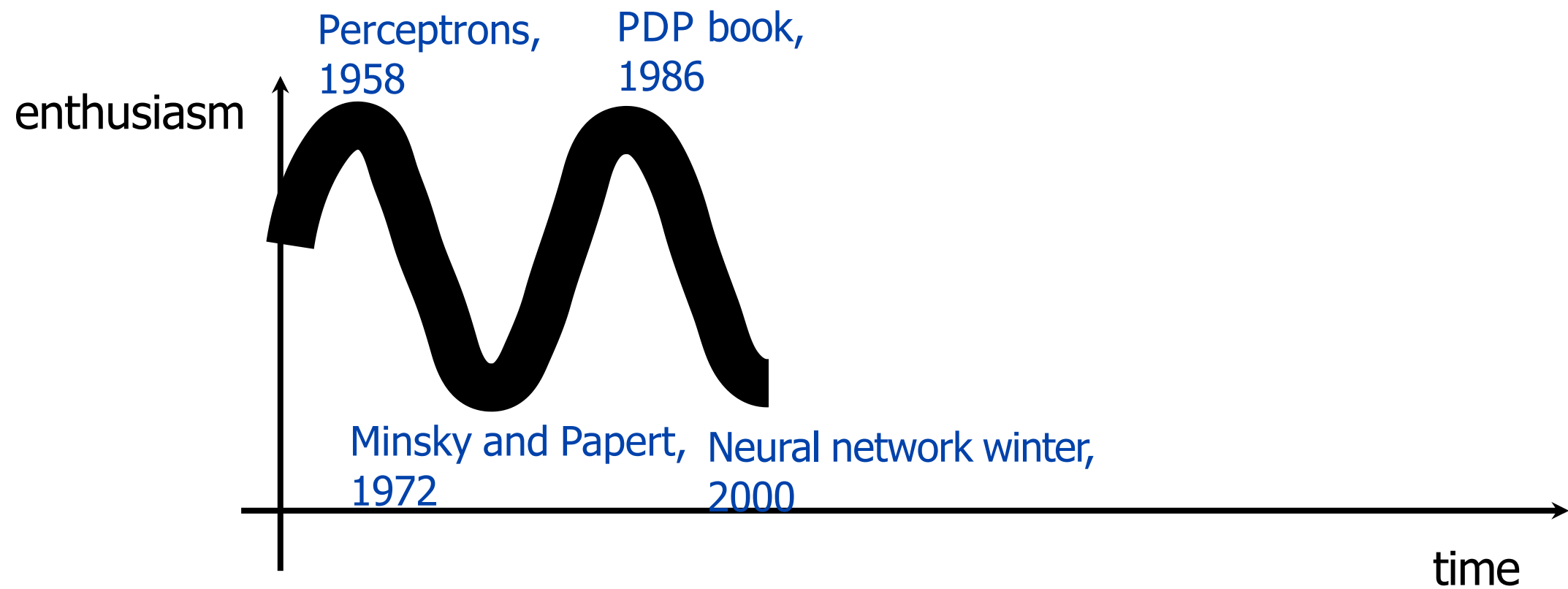


# Yann LeCun

Was at Bell Labs when  
this video was recorded

Now  
Prof @ NYU  
Chief Scientist @ Meta

Turing Award 2018  
(shared with Hinton and  
Bengio)





# ImageNet:

## First (?) large-scale computer vision dataset



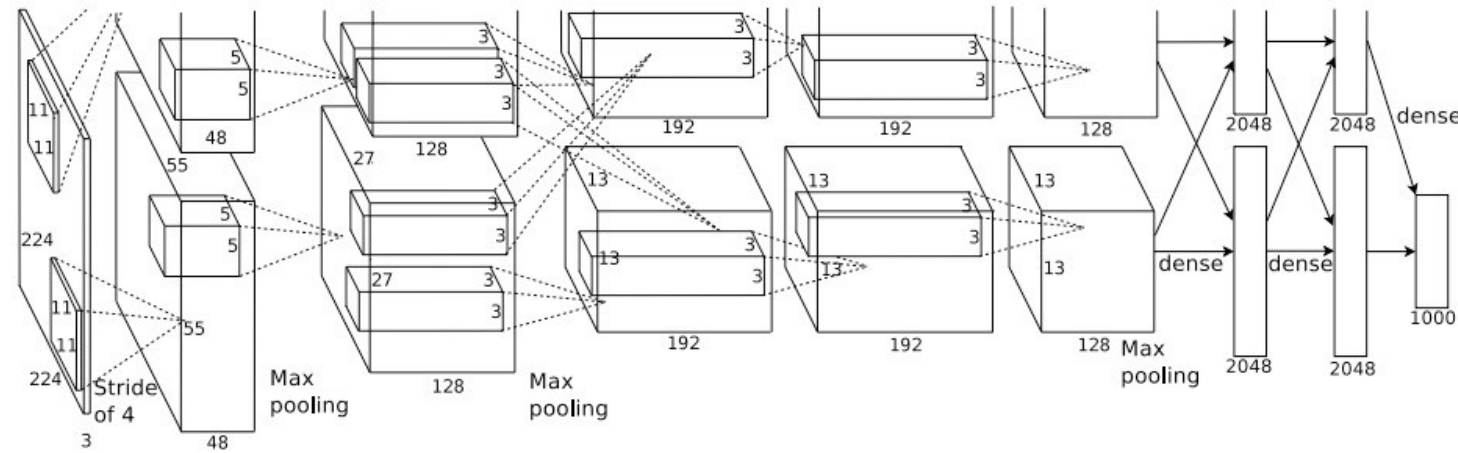
- Millions of images; 1000 categories
- **PI: Fei-Fei Li**
  - Then: Prof, Princeton
  - Now: Prof, Stanford
- 2019 Longuet-Higgins Prize
  - Some argued that Li deserved the 2018 Turing Award along with Hinton, LeCun, Bengio
  - Their work could not have been empirically tested without ImageNet!





# 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

	mite
	black widow
	cockroach
	tick
	starfish

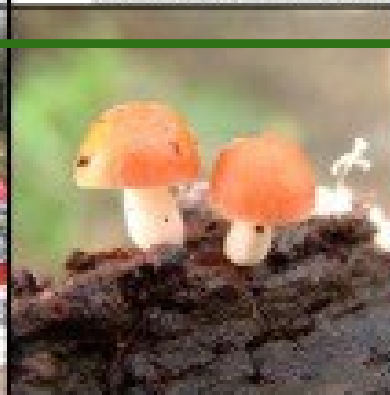
	container ship
	lifeboat
	amphibian
	fireboat
	drilling platform

	motor scooter
	go-kart
	moped
	bumper car
	golfcart

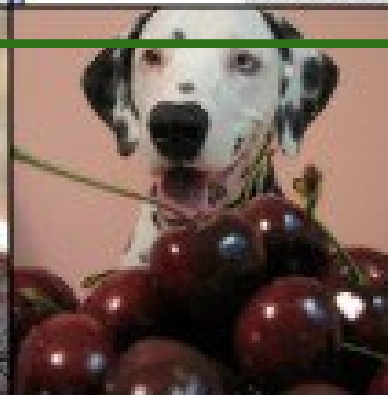
	leopard
	jaguar
	cheetah
	snow leopard
	Egyptian cat



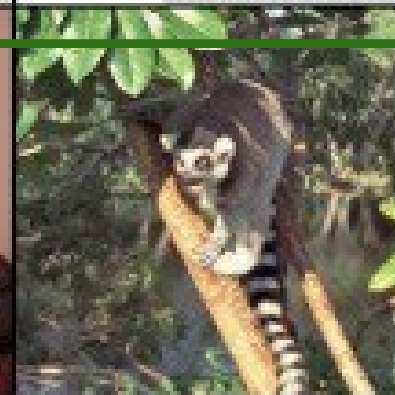
grille



mushroom



cherry



Madagascar cat

	convertible
	grille
	pickup
	beach wagon
	fire engine

	agaric
	mushroom
	jelly fungus
	gill fungus
	dead-man's-fingers

	dalmatian
	grape
	elderberry
	ffordshire bullterrier
	currant

	squirrel monkey
	spider monkey
	titi
	indri
	howler monkey

# Krizhevsky, Sutskever, and Hinton, NeurIPS 2012



mite



container ship



motor scooter



leopard

	mite
	black widow
	cockroach
	tick
	starfish

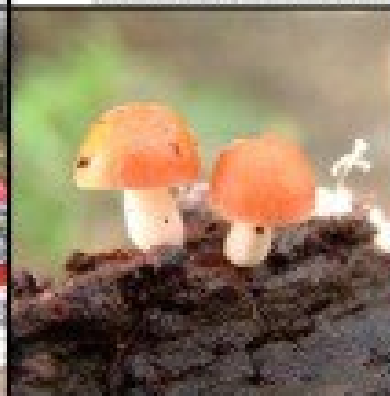
	container ship
	lifeboat
	amphibian
	fireboat
	drilling platform

	motor scooter
	go-kart
	moped
	bumper car
	golfcart

	leopard
	jaguar
	cheetah
	snow leopard
	Egyptian cat



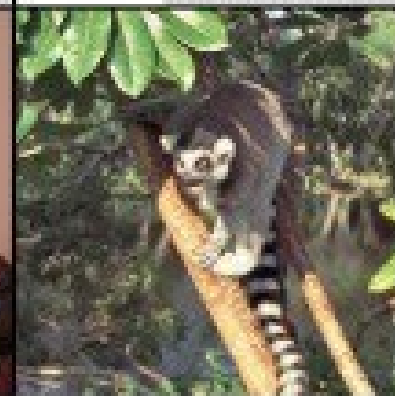
grille



mushroom



cherry



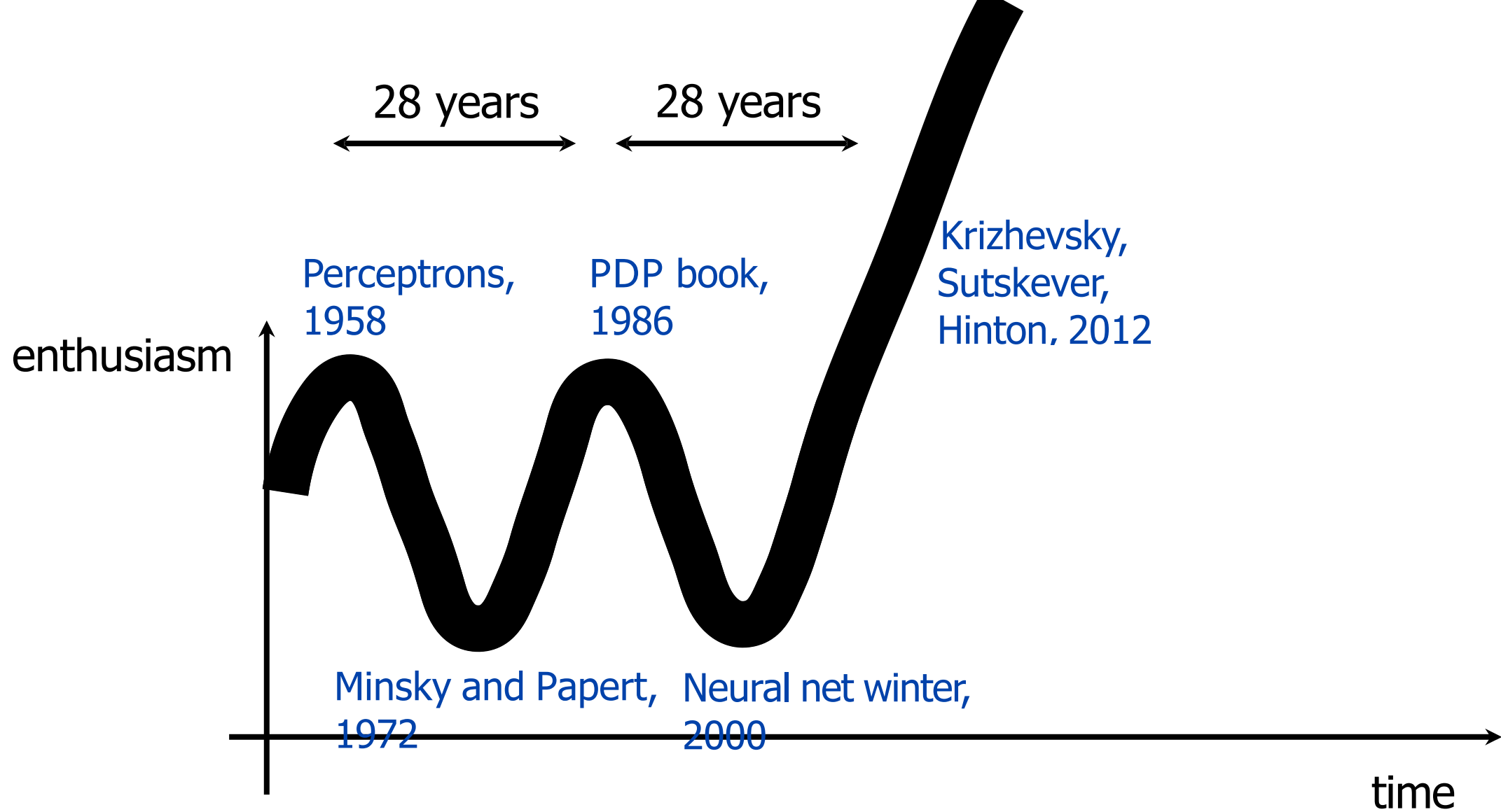
Madagascar cat

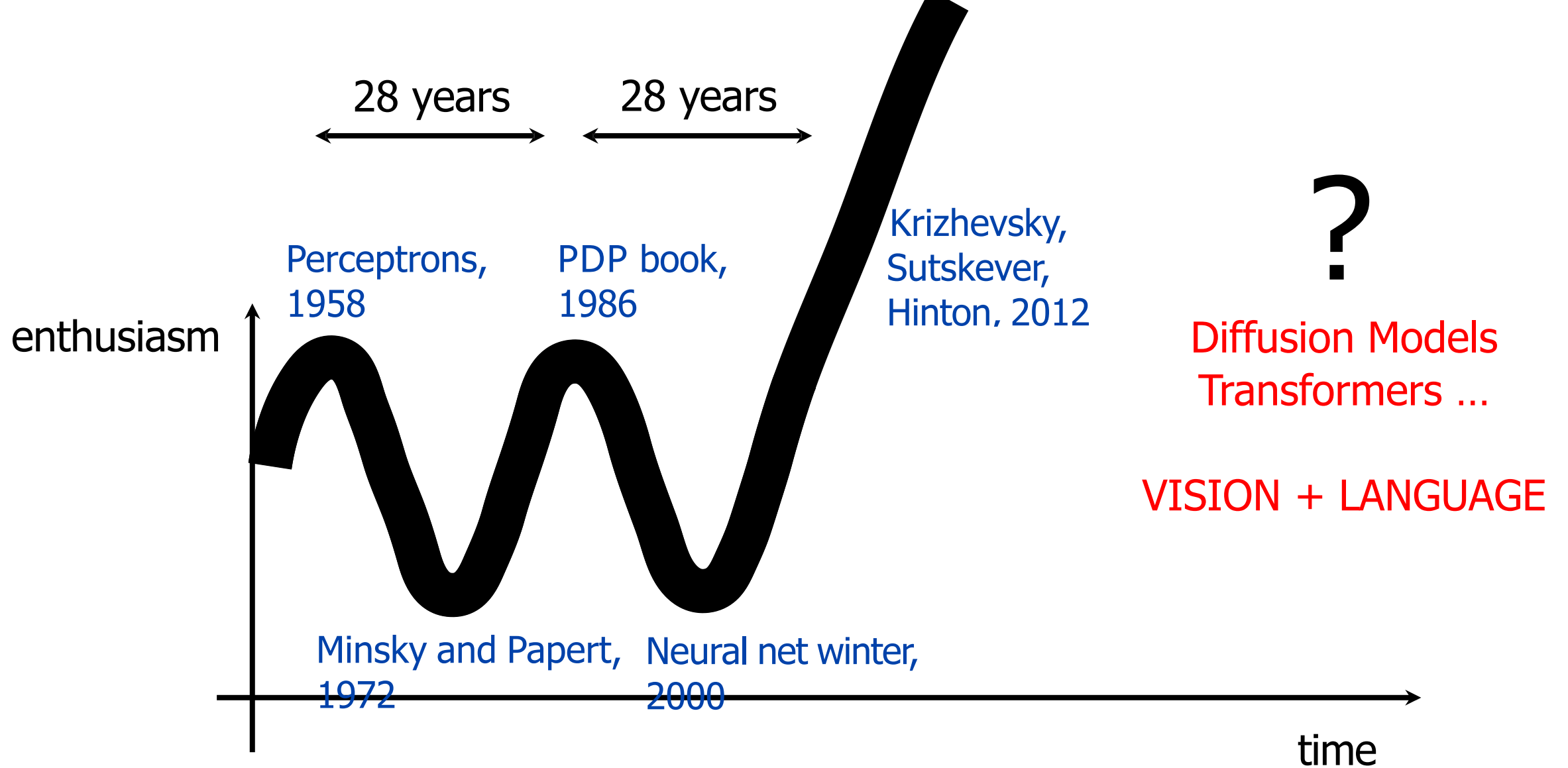
	convertible
	grille
	pickup
	beach wagon
	fire engine

	agaric
	mushroom
	jelly fungus
	gill fungus
	dead-man's-fingers

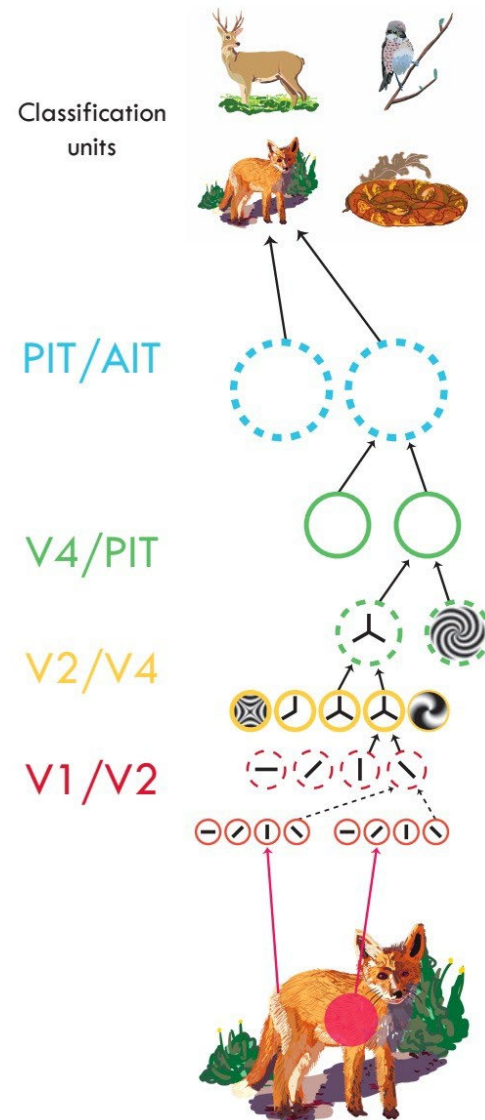
	dalmatian
	grape
	elderberry
	ffordshire bullterrier
	currant

	squirrel monkey
	spider monkey
	tit
	indri
	howler monkey





# Inspiration: Hierarchical Representations



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

# Object recognition

Pixel 1



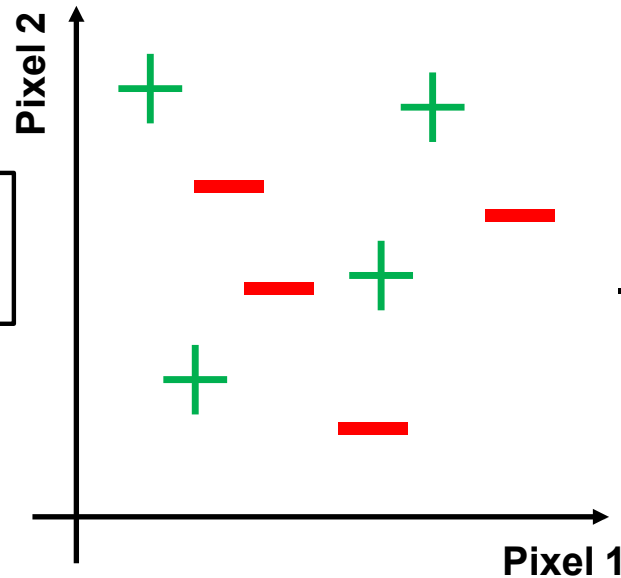
Neural Network

Is dog?

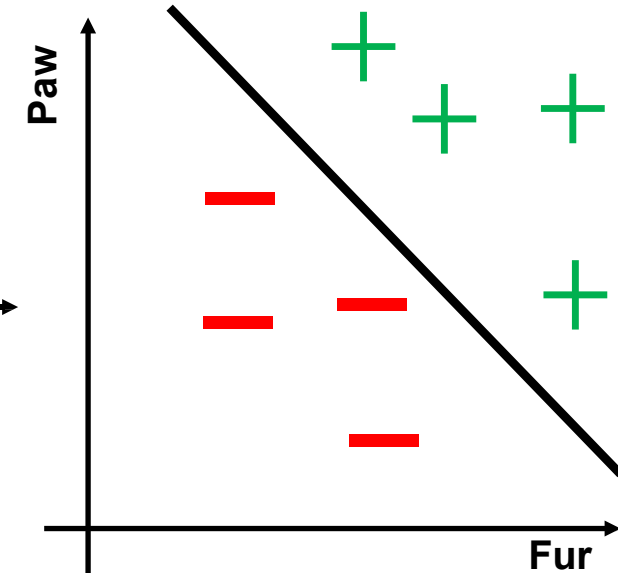
Pixel 2

Input Space

Feature Space



$f(x)$

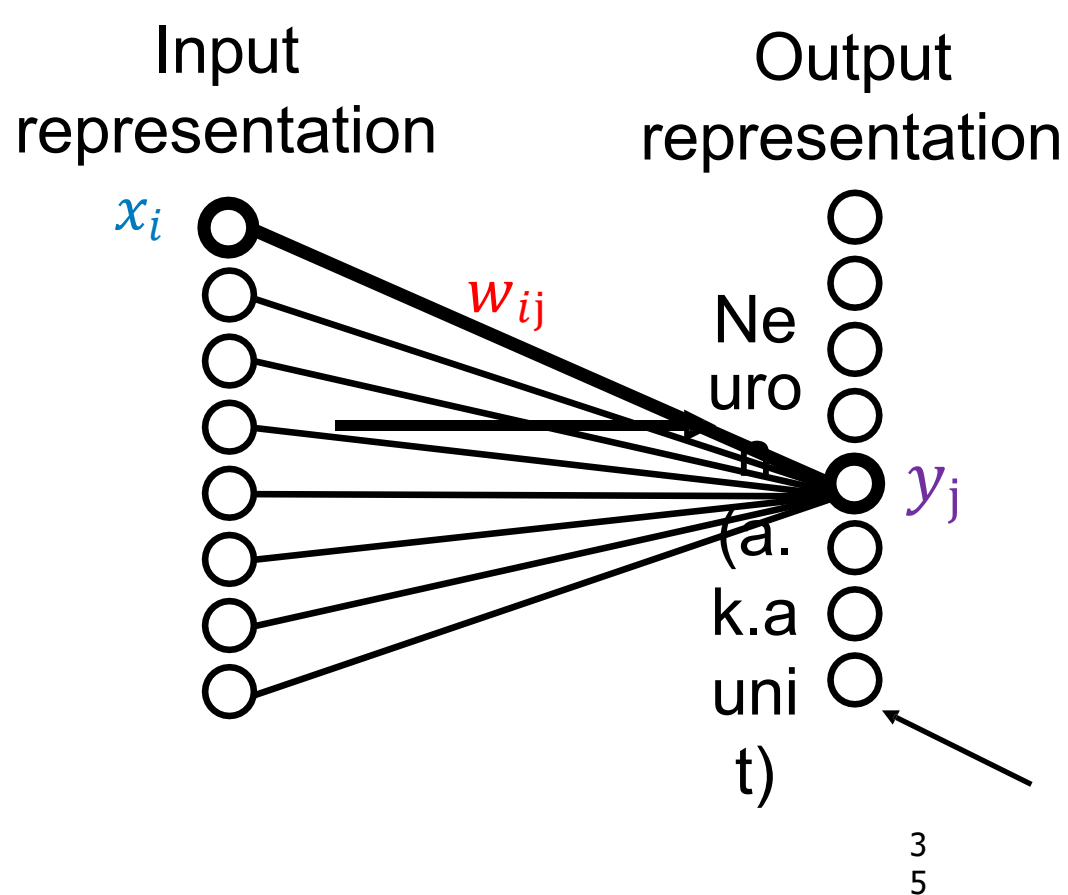


Goal: automatically learn a function that maps data from the input space to a feature space, i.e., "feature learning", rather than use hand-crafted features

# Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space:

## Linear layer



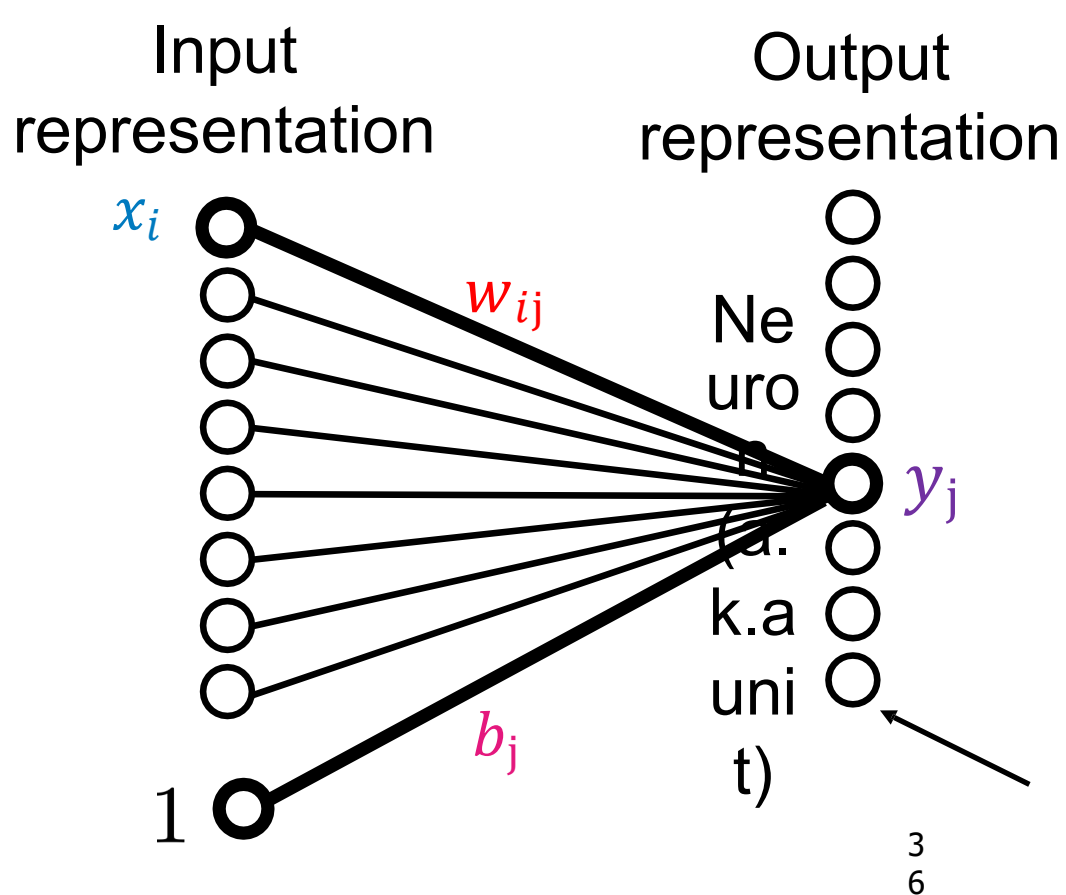
$$y_j = \sum_i \overset{\text{weights}}{w_{ij}} x_i$$



# Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space

## Linear layer



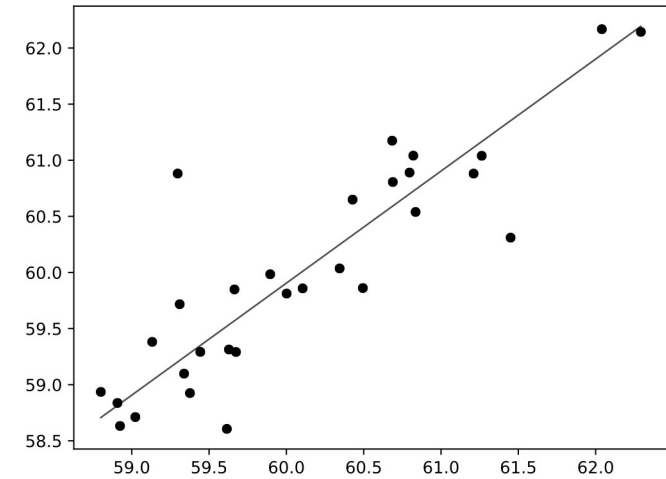
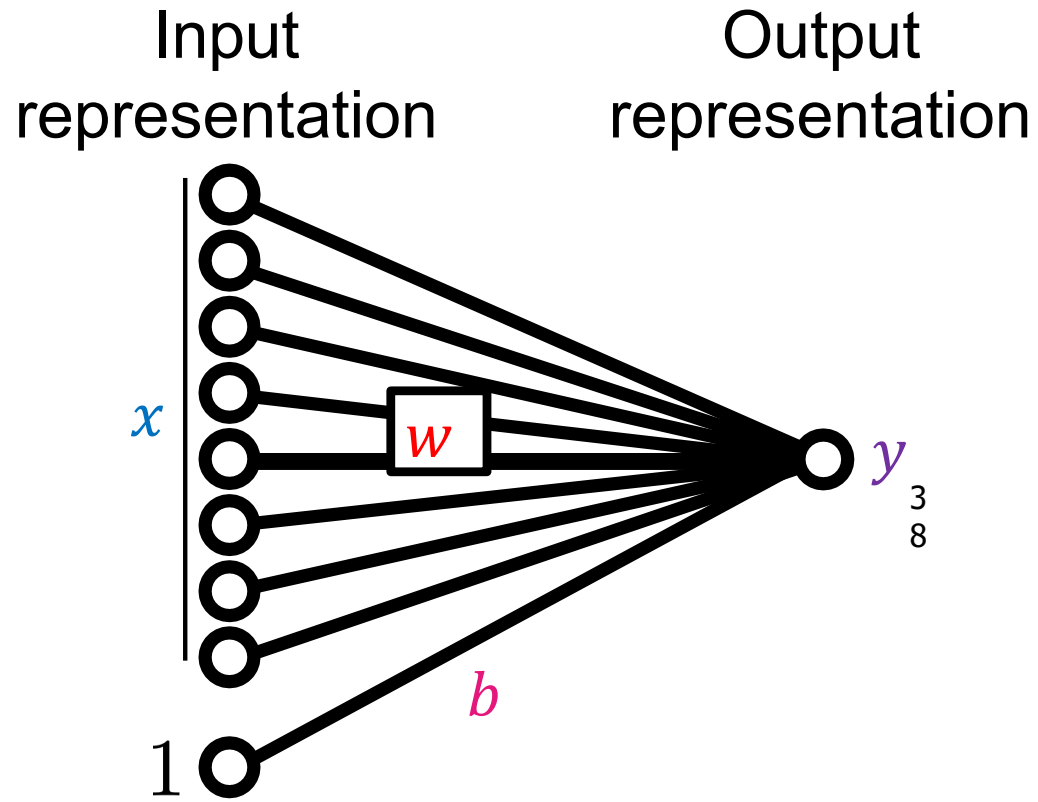
$$y_j = \sum_i w_{ij} x_i + b_j$$

weights

bias

# Example: Linear Regression

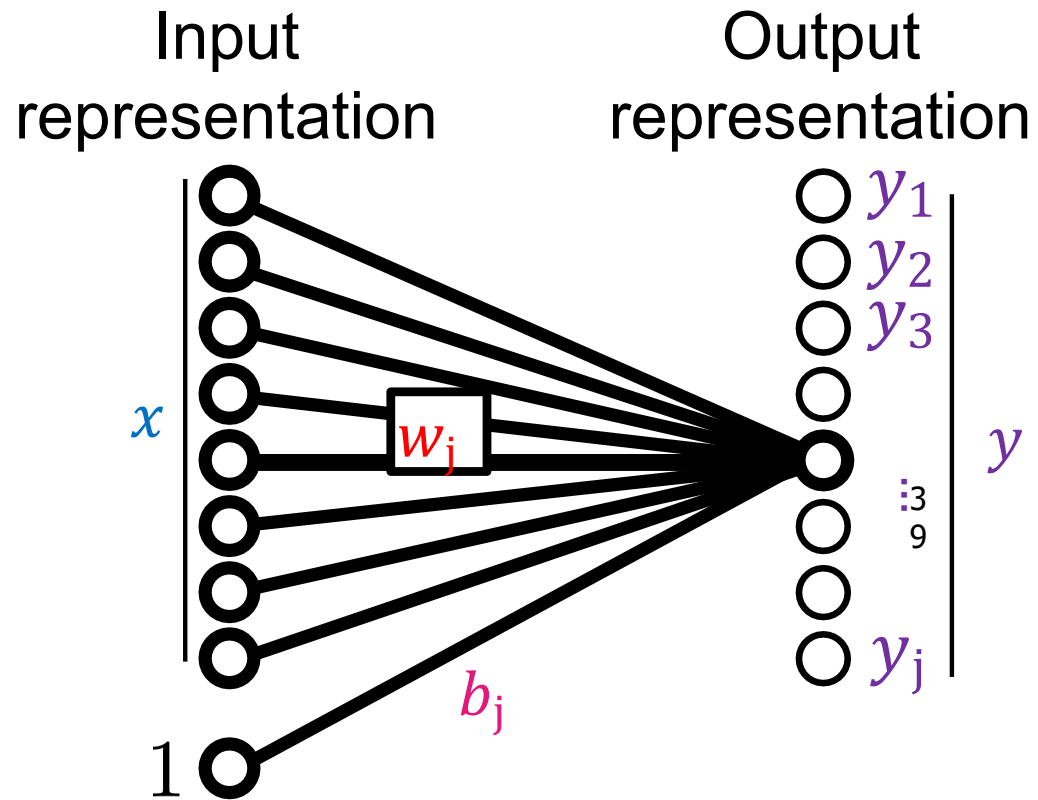
## Linear layer



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

# Computation in a neural net – Full Layer

## Linear layer



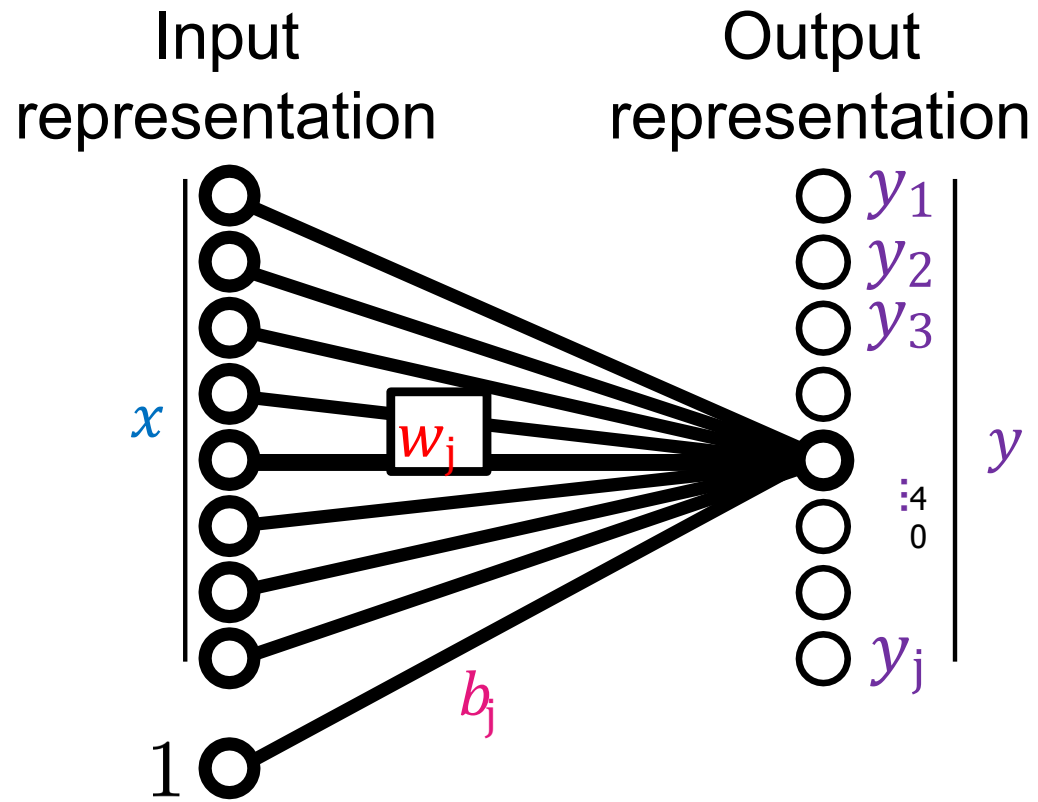
$$y = Wx + b$$

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

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

# Computation in a neural net – Full Layer

## Linear layer



## Full layer

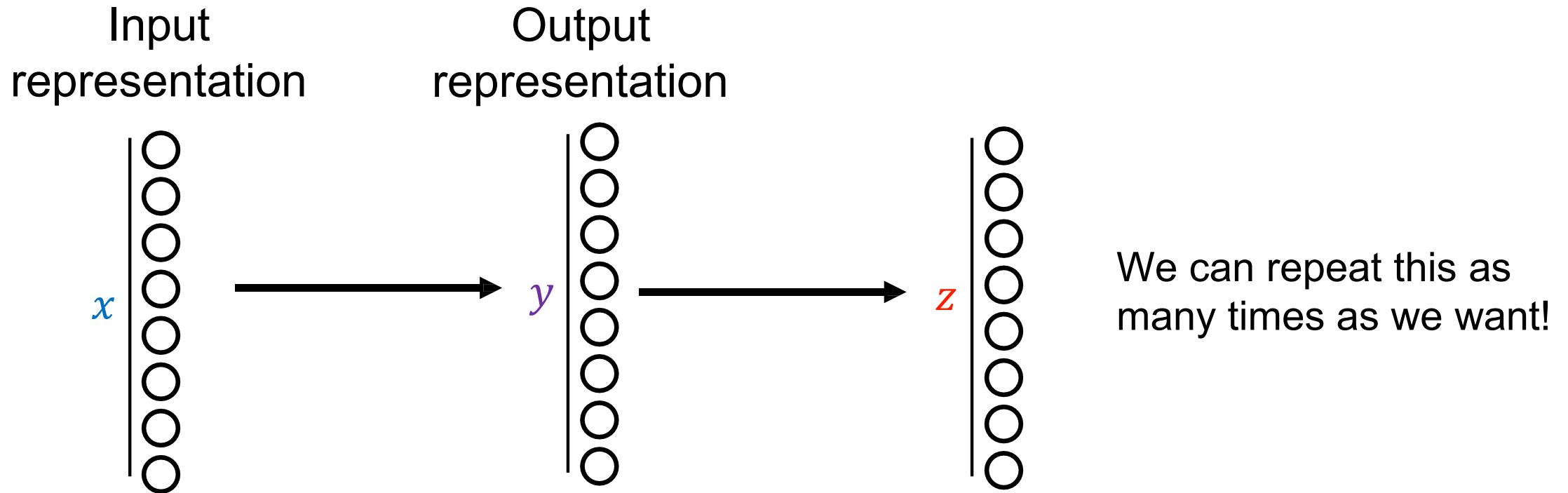
$$y = Wx + b$$

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

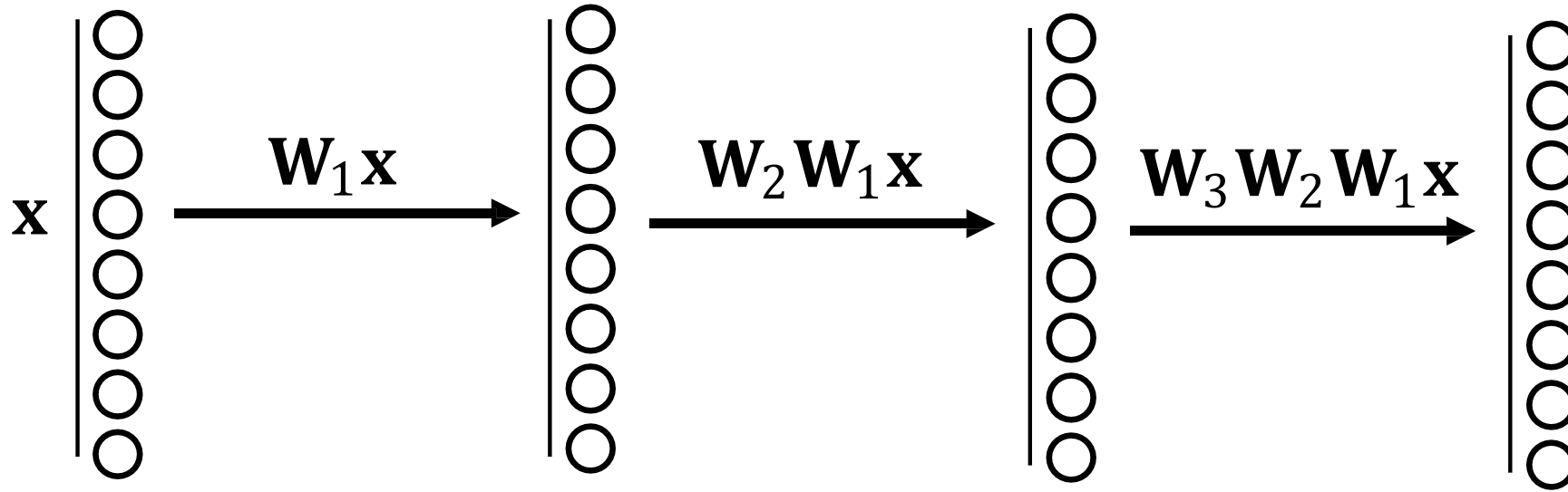
Can again simplify notation by appending a 1 to  $\mathbf{x}$

# 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:



# What is the problem with this idea?



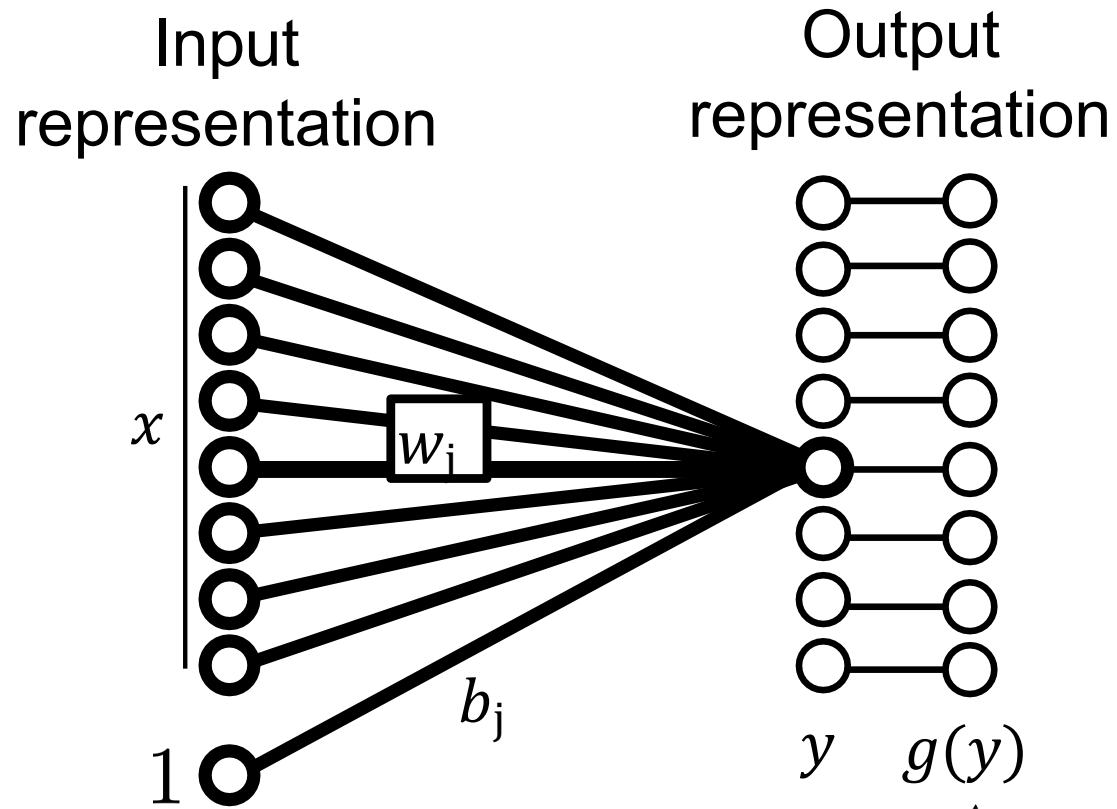
Can be expressed as single linear layer!

$$\begin{pmatrix} \mathbf{G} & \mathbf{W}_i \\ i \end{pmatrix} \mathbf{x} = \hat{\mathbf{W}} \mathbf{x}$$

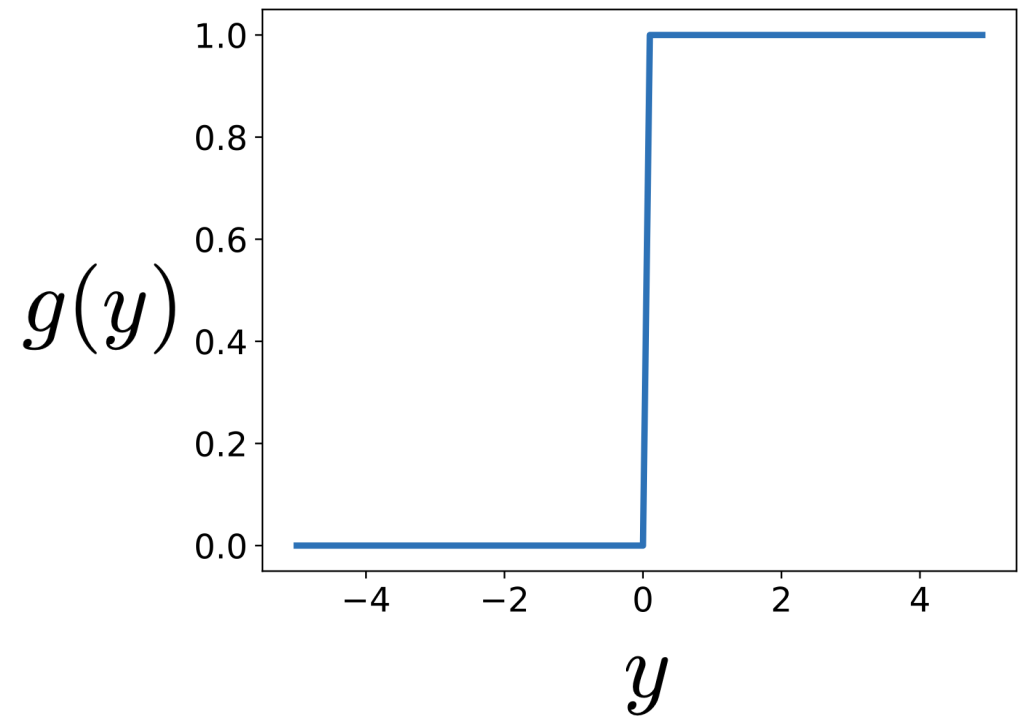
Limited power: can't solve XOR ☹️

# Solution: simple nonlinearity

## Linear layer

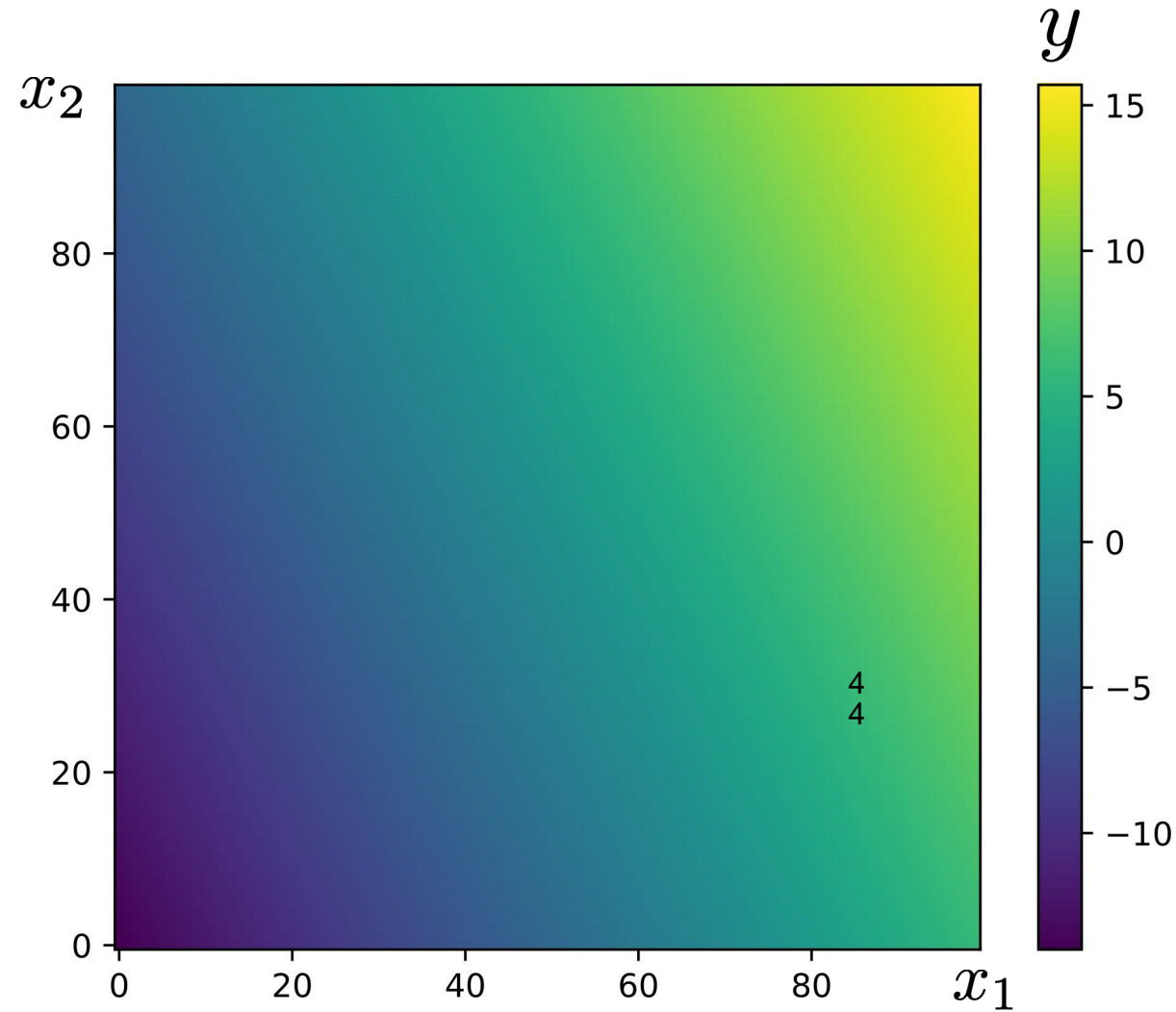


$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



Pointwise  
Non-linearity

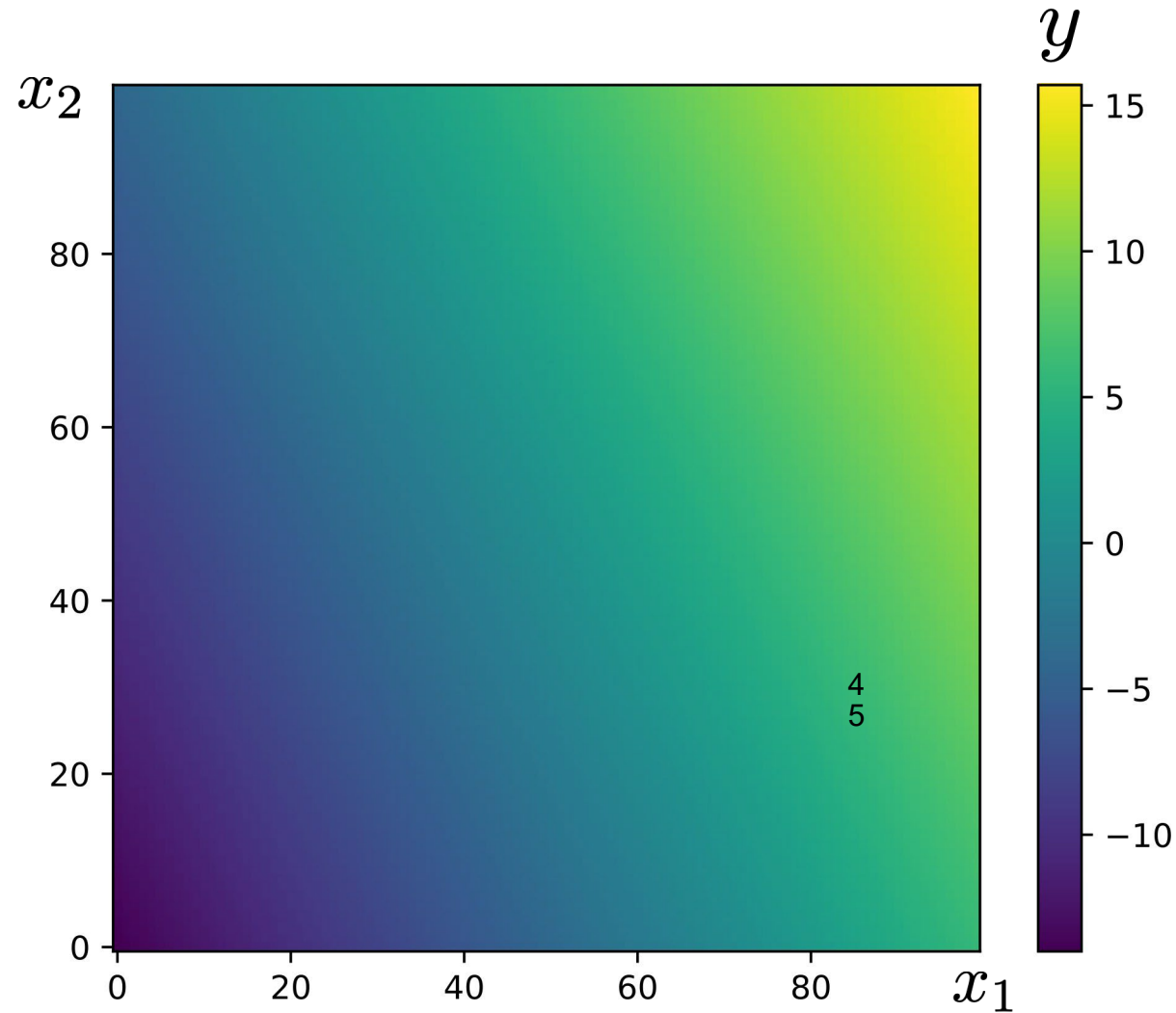
# Example: linear classification with a perceptron



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



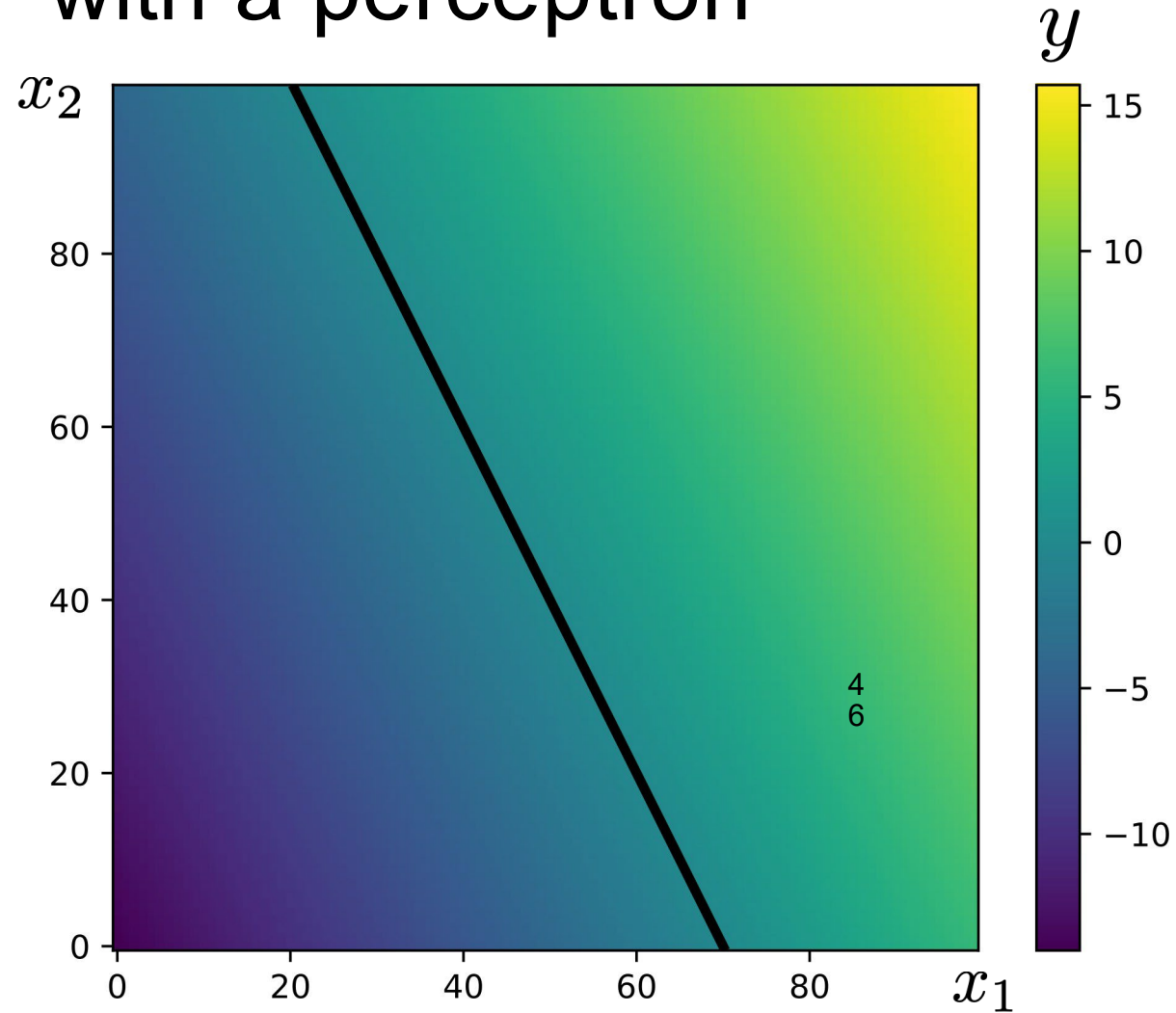
# Example: linear classification with a perceptron



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

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

# Example: linear classification with a perceptron



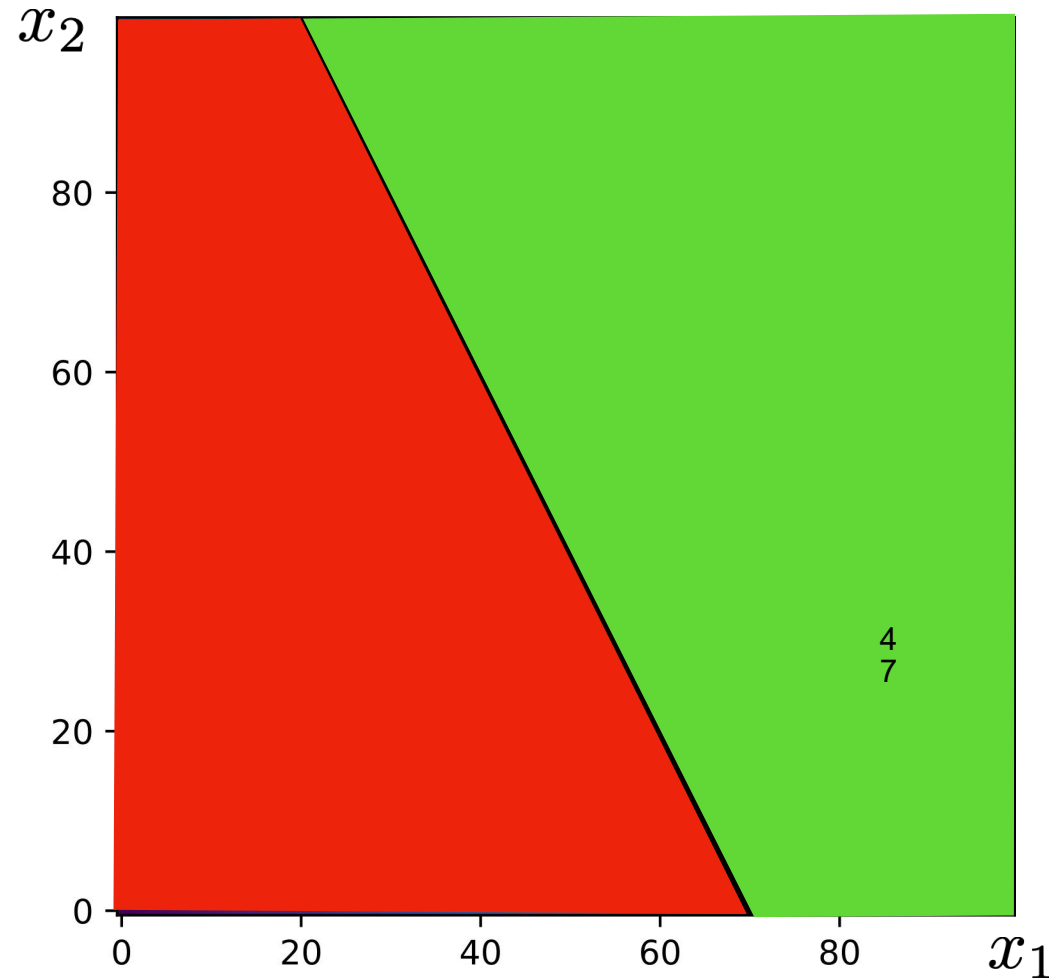
$$y = \mathbf{x}^T \mathbf{w} + b$$

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

“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

$g(y)$



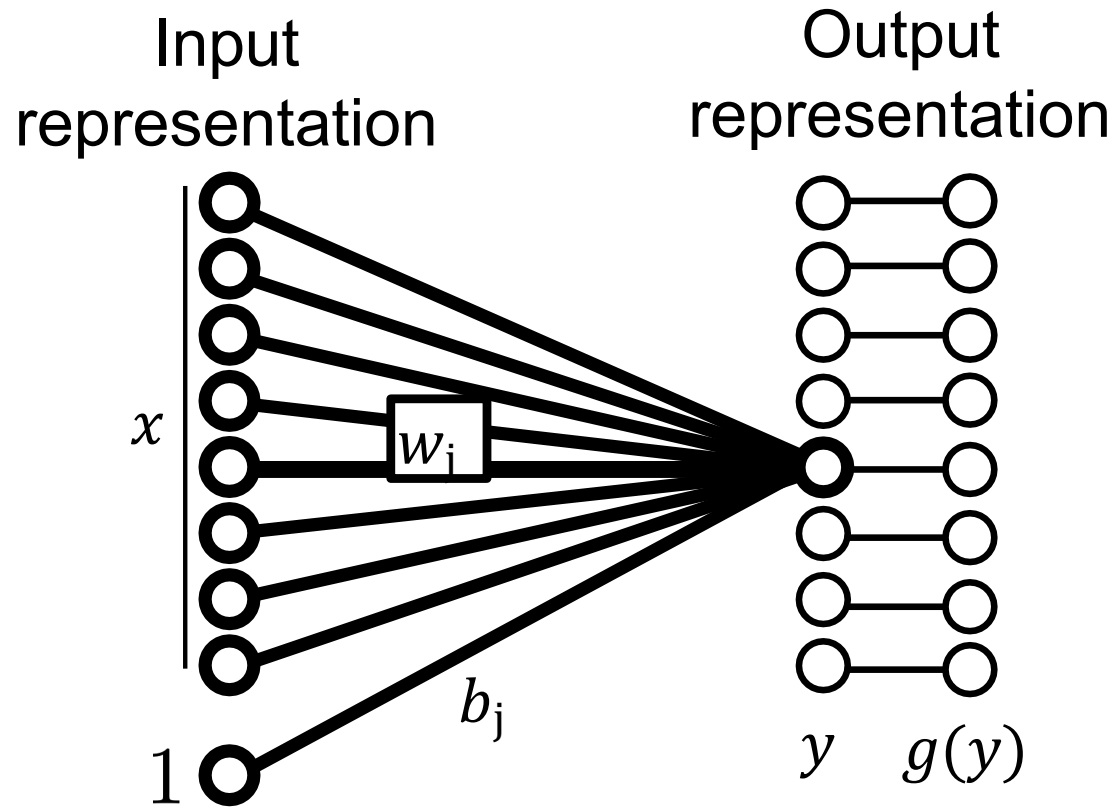
$$y = \mathbf{x}^T \mathbf{w} + b$$

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

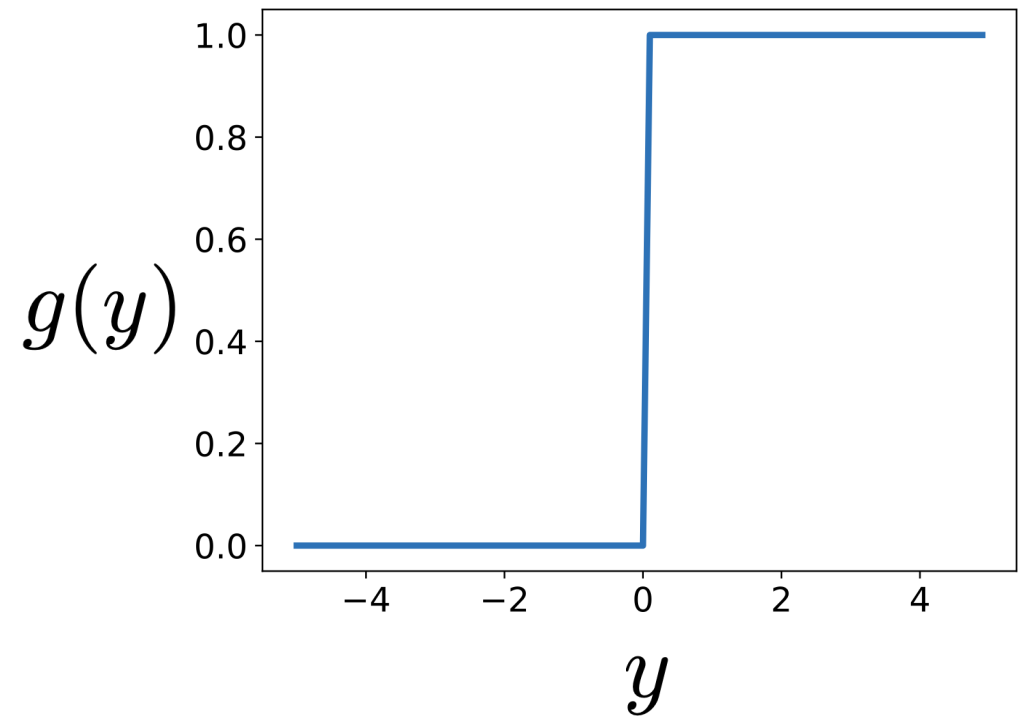
“when  $y$  is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)”

# Computation in a neural net - nonlinearity

## Linear layer



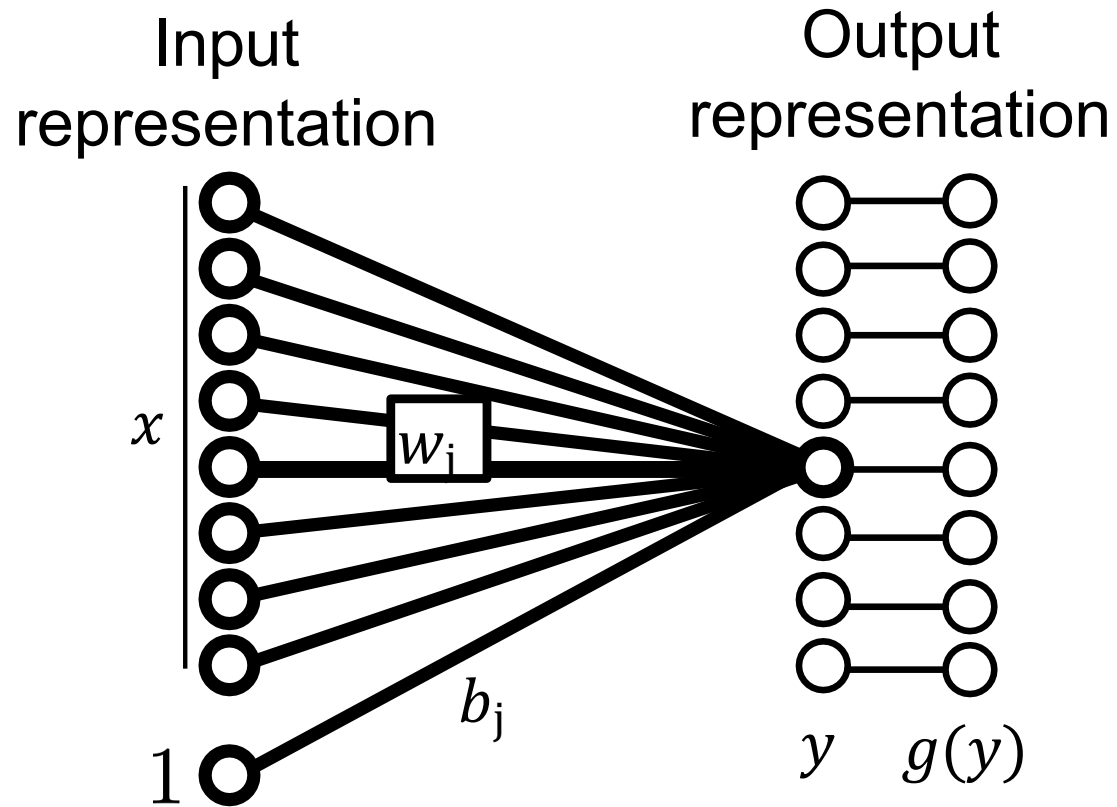
$$g(y) = \begin{cases} 1, & \text{if } y > 0 \\ 0, & \text{otherwise} \end{cases}$$



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

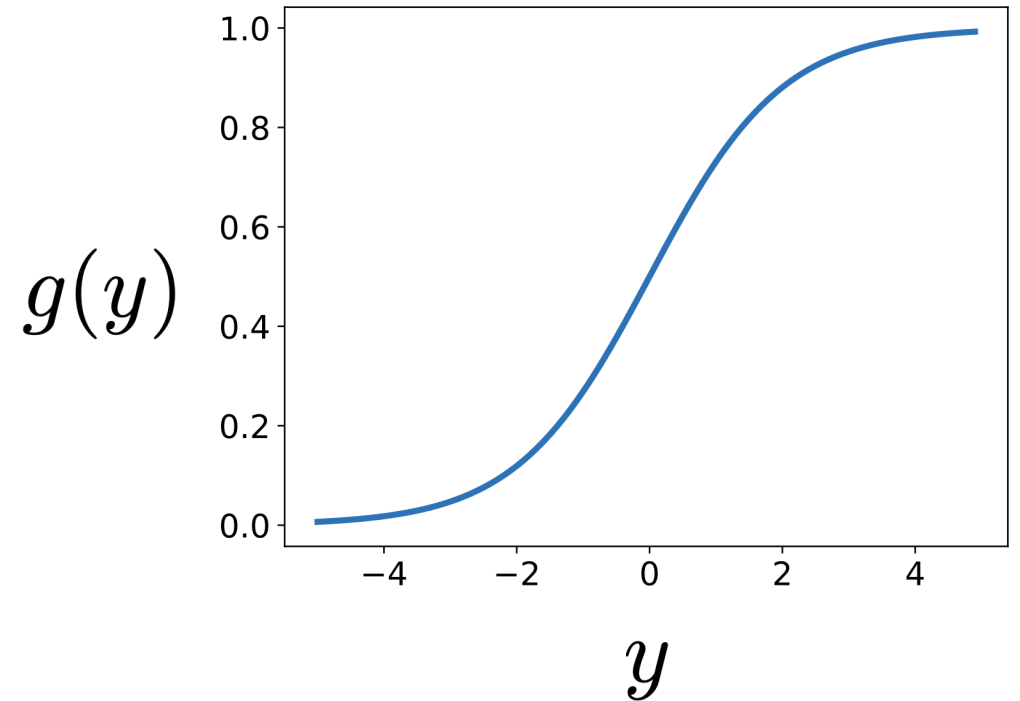
# Computation in a neural net - nonlinearity

## Linear layer



## Sigmoid

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

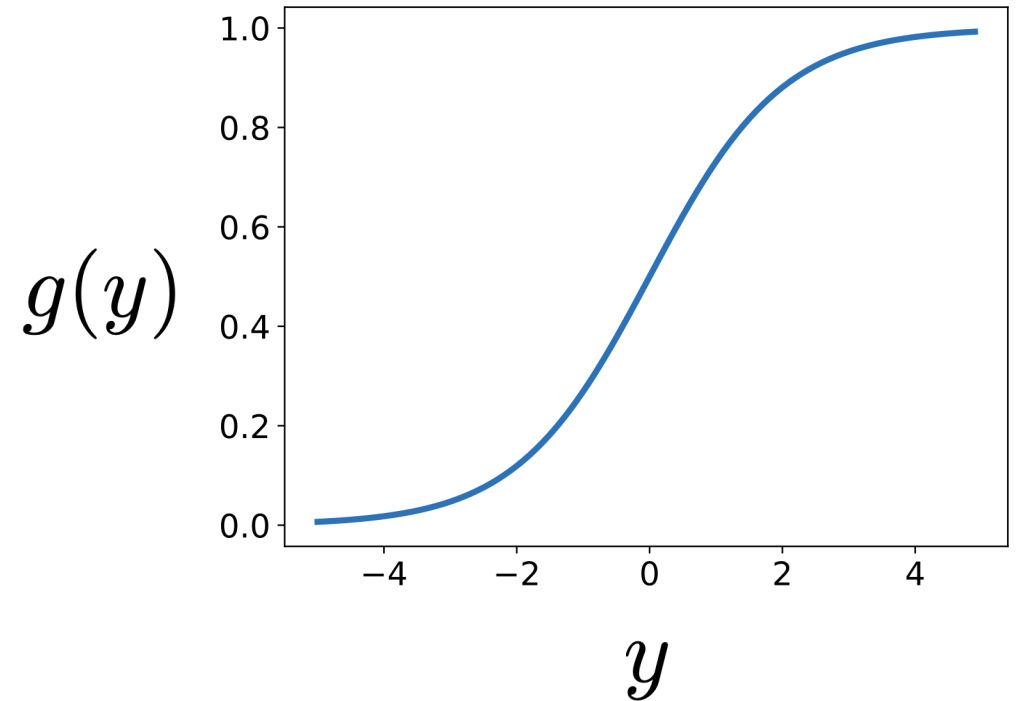


# Computation in a neural net - nonlinearity

- Bounded between [0,1]
- Saturation for large +/- inputs
- Gradients go to zero

**Sigmoid**

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

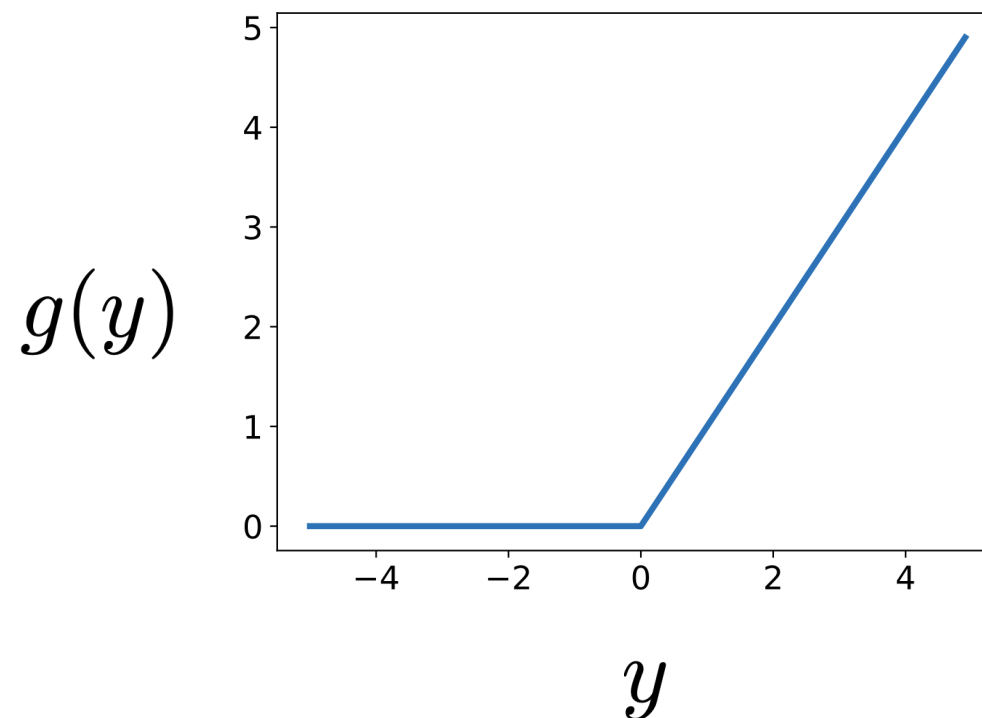


# Computation in a neural net — nonlinearity

- 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)

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



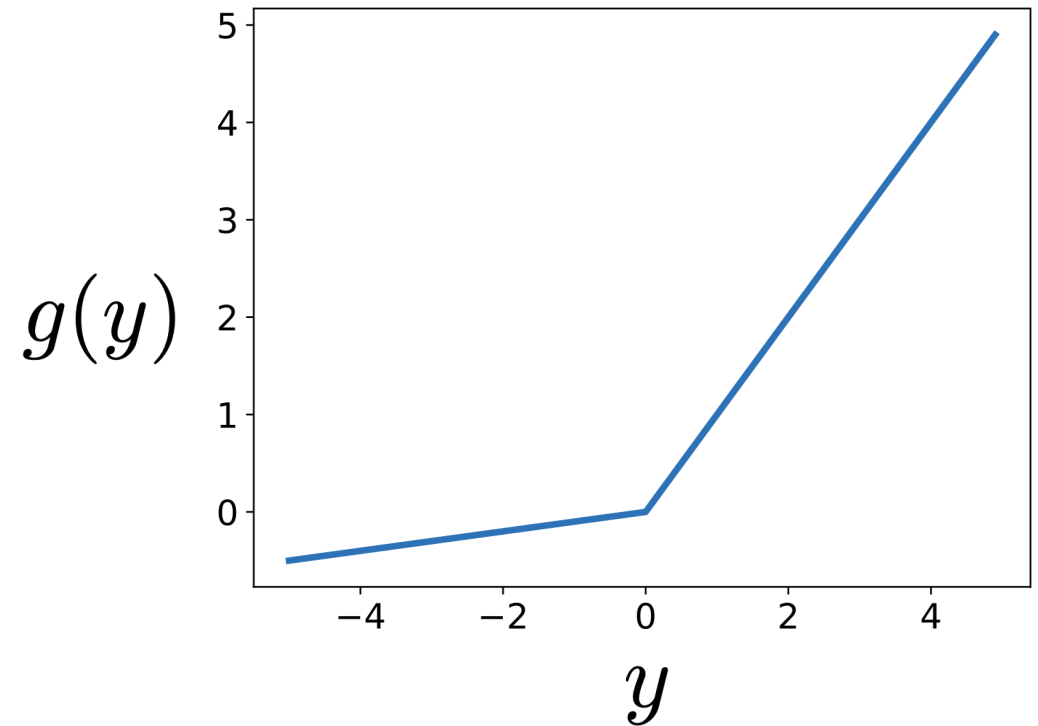
# Computation in a neural net — nonlinearity

- where  $a$  is small (e.g., 0.02)
- Efficient to implement:
- Has non-zero gradients everywhere (unlike ReLU)

$$\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if } y < 0 \\ 1, & \text{if } y \geq 0 \end{cases}$$

## Leaky ReLU

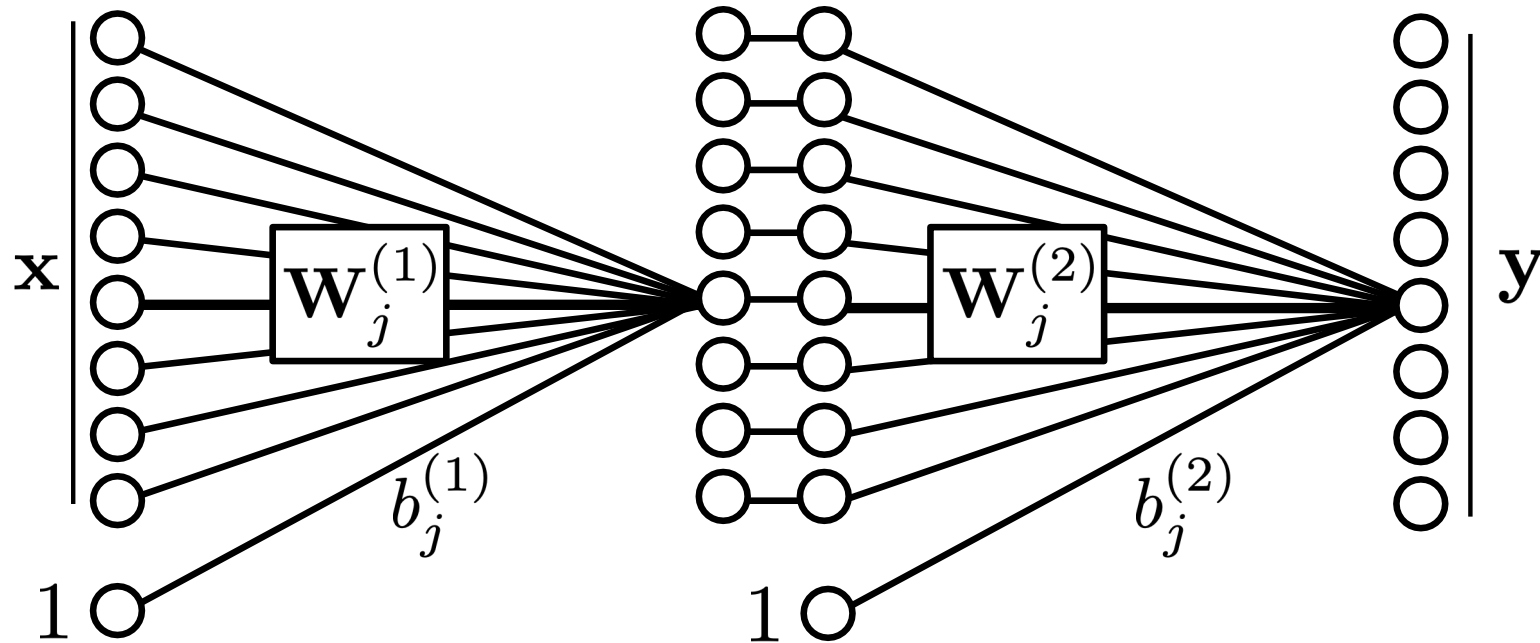
$$g(y) = \begin{cases} \max(0, y), & \text{if } y \geq 0 \\ a \min(0, y), & \text{if } y < 0 \end{cases}$$





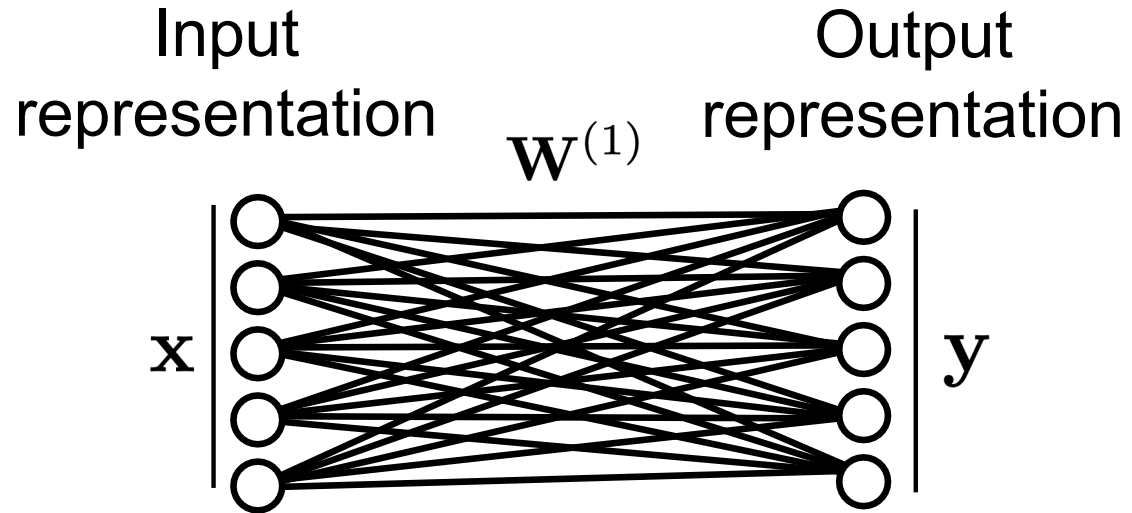
# Stacking layers

Input representation      Intermediate representation      Output representation

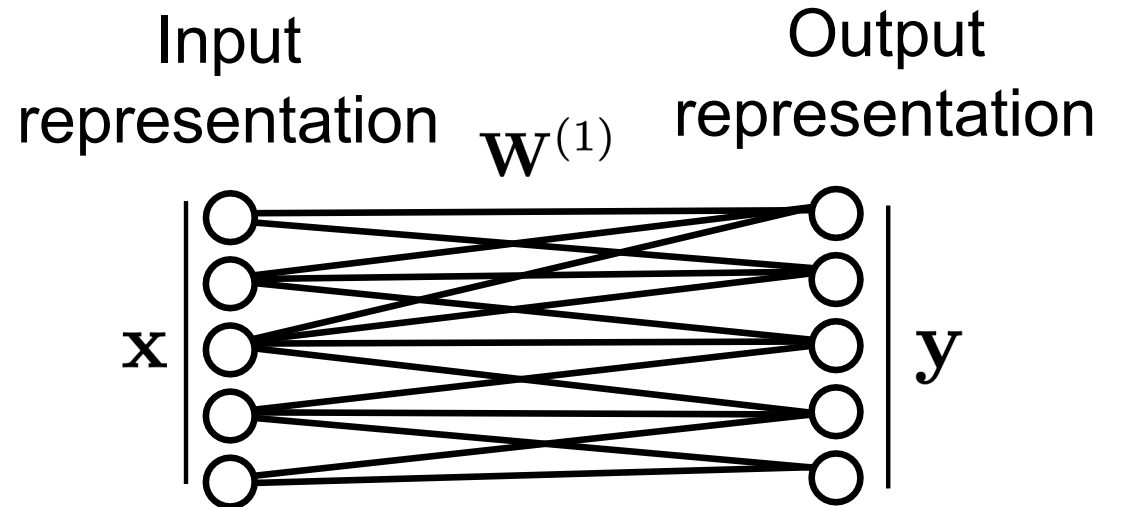


**h** = "hidden units"

# Connectivity patterns

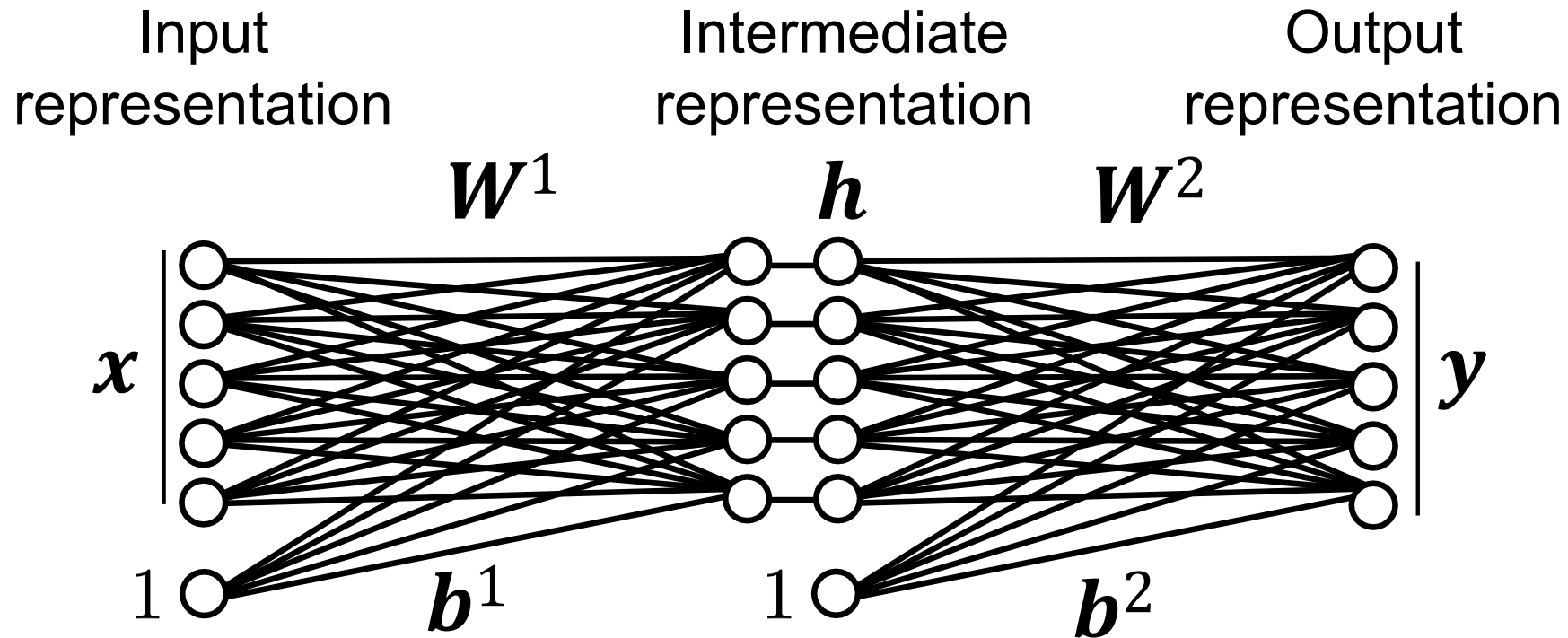


*Fully connected layer*



*Locally connected layer  
(Sparse  $\mathbf{W}$ )*

# Stacking layers

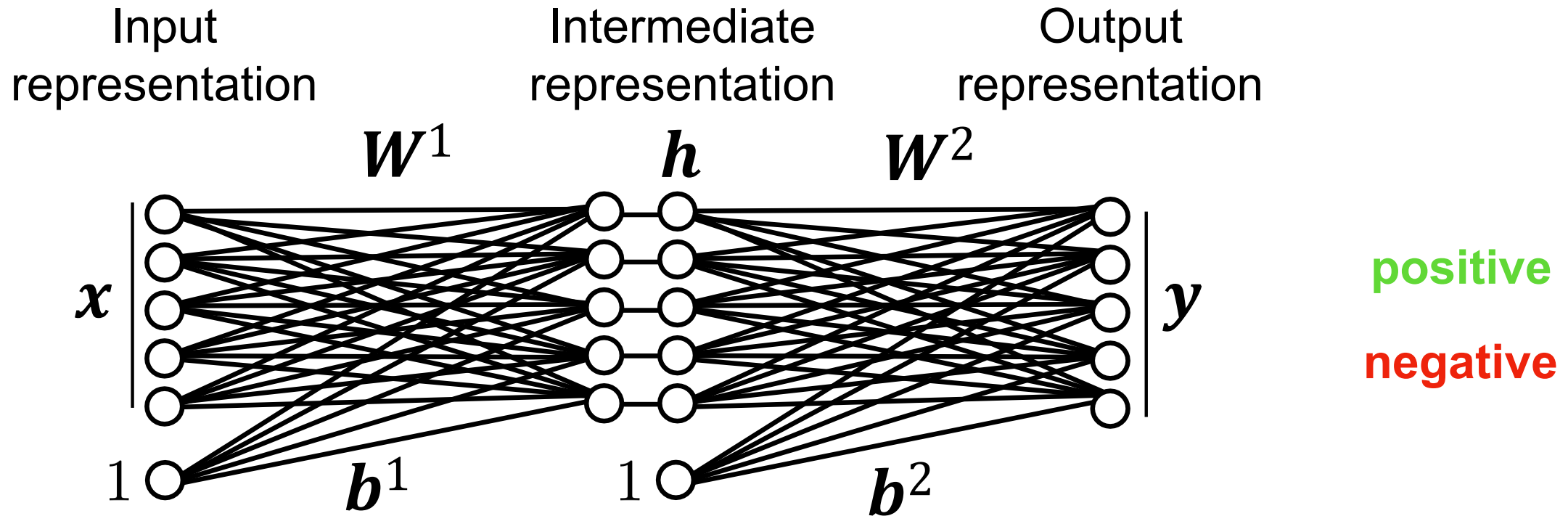


$$\mathbf{h} = g(W^1 x + b^1) \quad y = g(W^2 \mathbf{h} + b^2)$$

ReLU  $\nearrow$

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

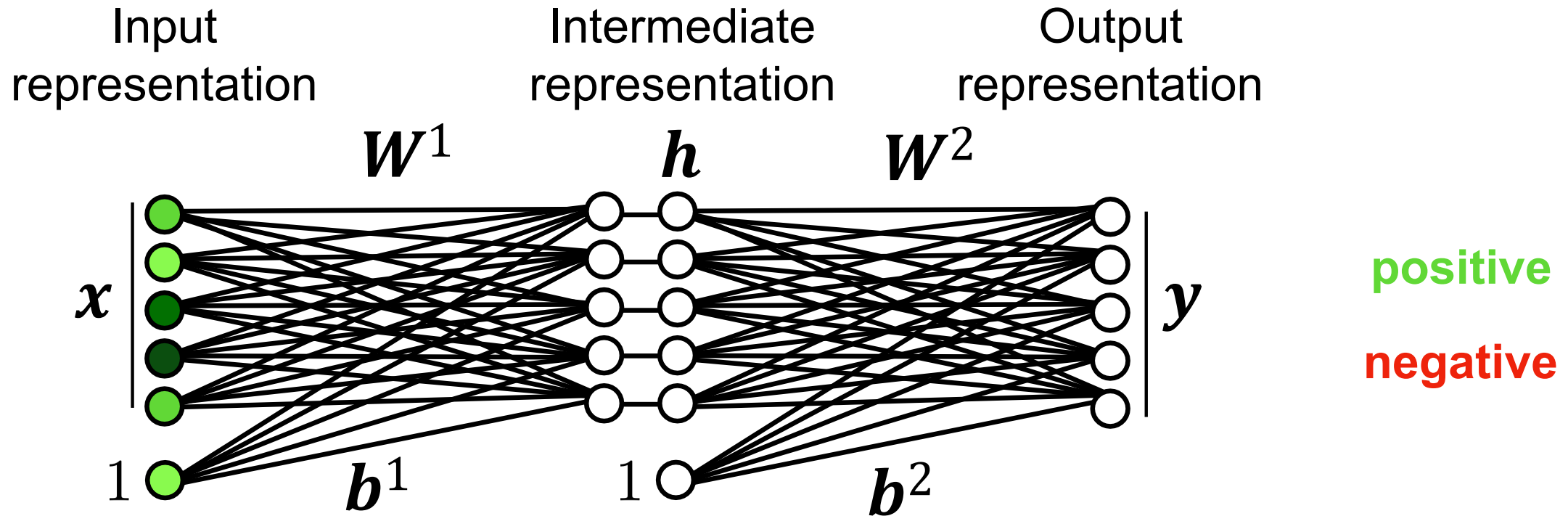
# Stacking layers



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

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

# Stacking layers

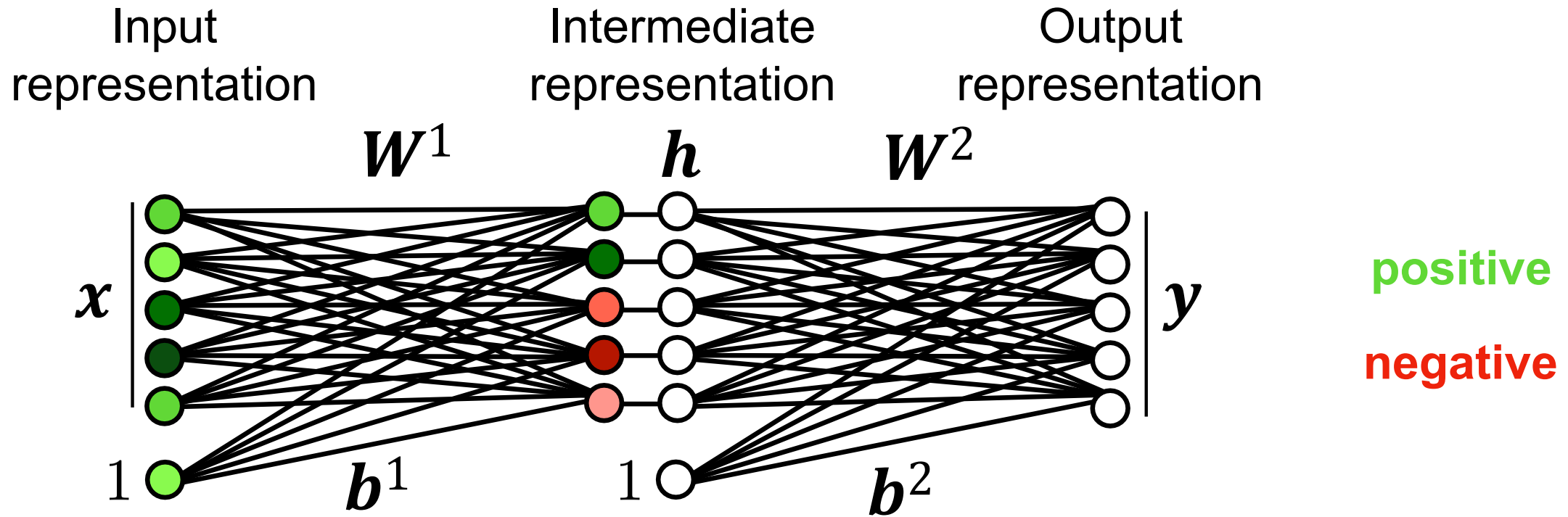


$$\mathbf{h} = g(W^1 \mathbf{x} + \mathbf{b}^1) \quad \mathbf{y} = g(W^2 \mathbf{h} + \mathbf{b}^2)$$

ReLU  $\nearrow$

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

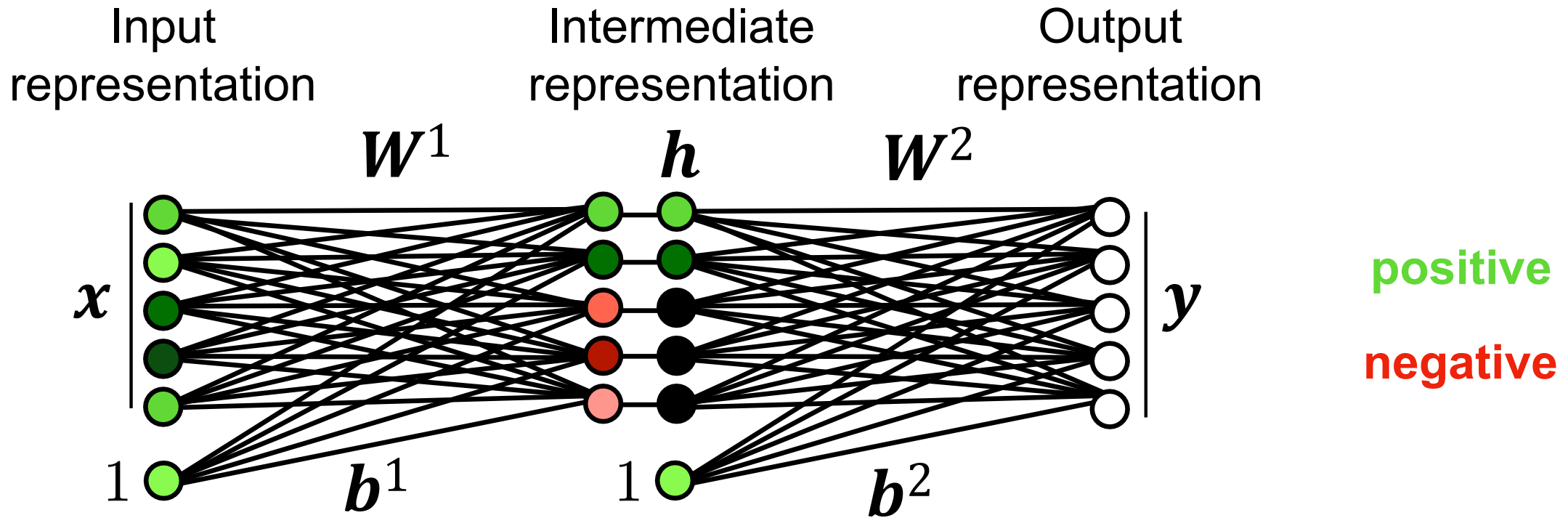
# Stacking layers



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

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

# Stacking layers

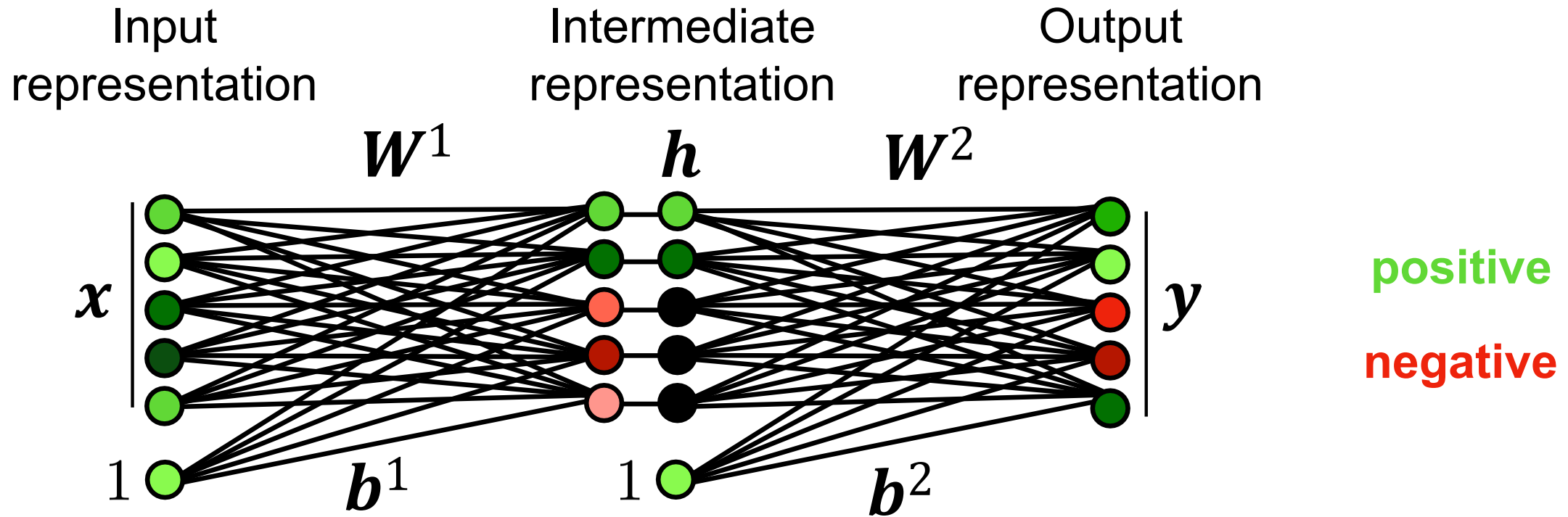


$$\mathbf{h} = g(W^1 x + b^1) \quad y = g(W^2 \mathbf{h} + b^2)$$

ReLU  $\nearrow$

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

# Stacking layers



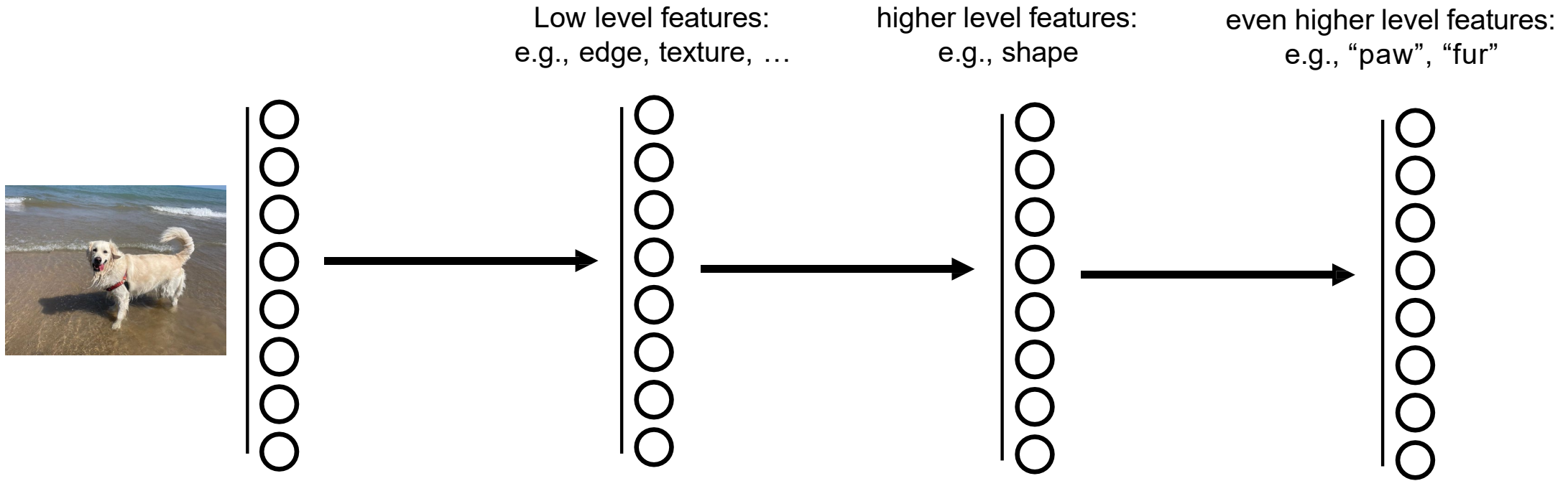
$$\mathbf{h} = g(W^1 x + b^1) \quad y = g(W^2 \mathbf{h} + b^2)$$

ReLU  $\nearrow$

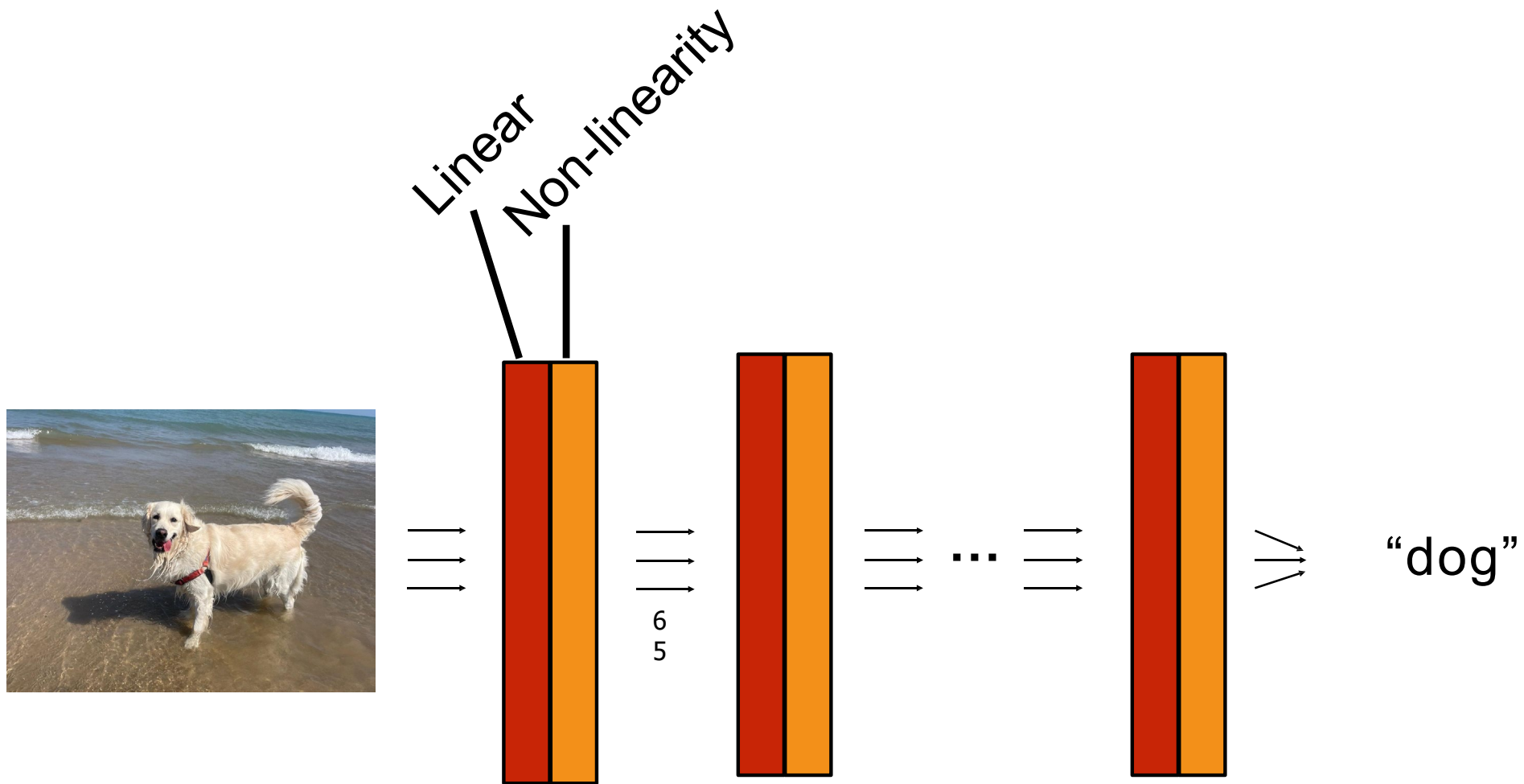
$$\theta = \{W^1, \dots, W^L, b^1, \dots, b^L\}$$



# Stacking layers - What's actually happening?

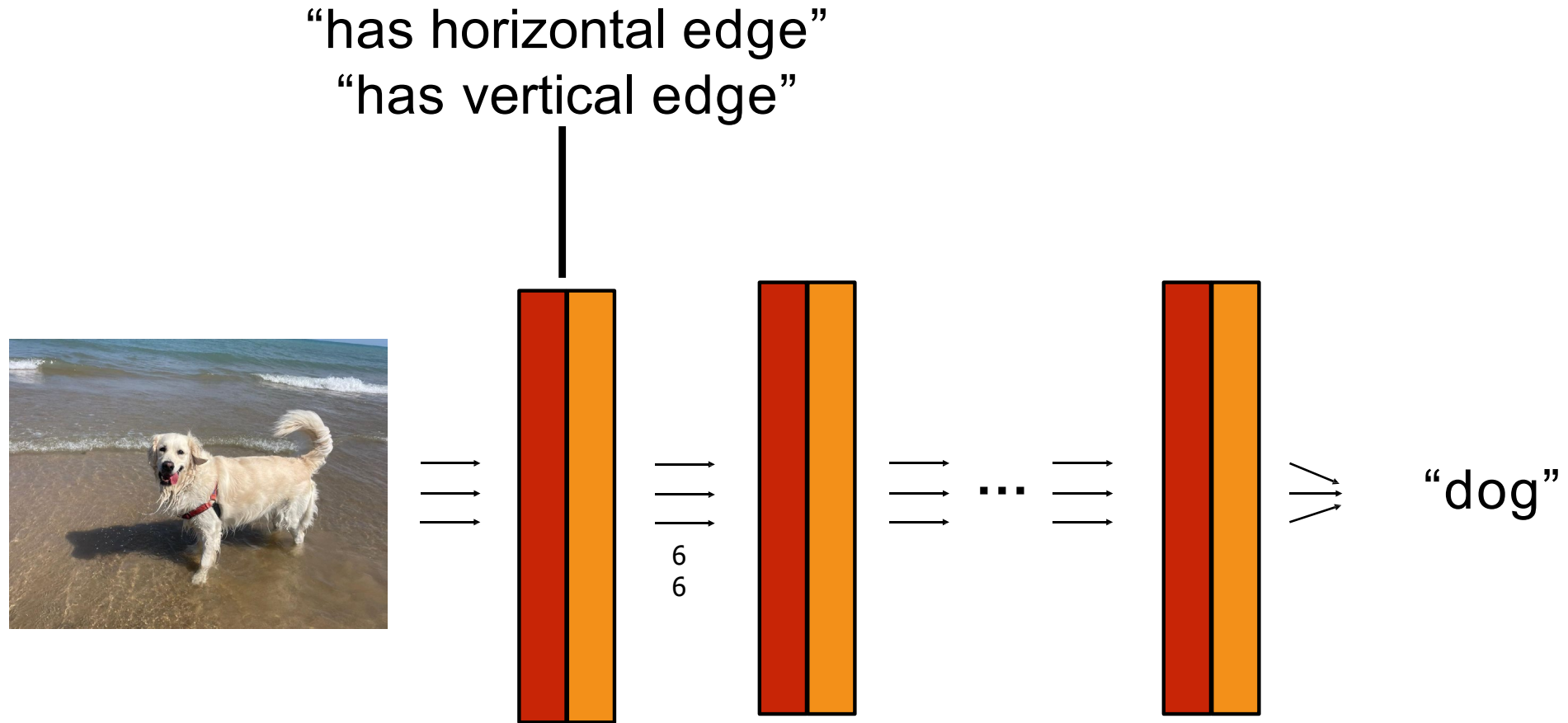


# Deep nets

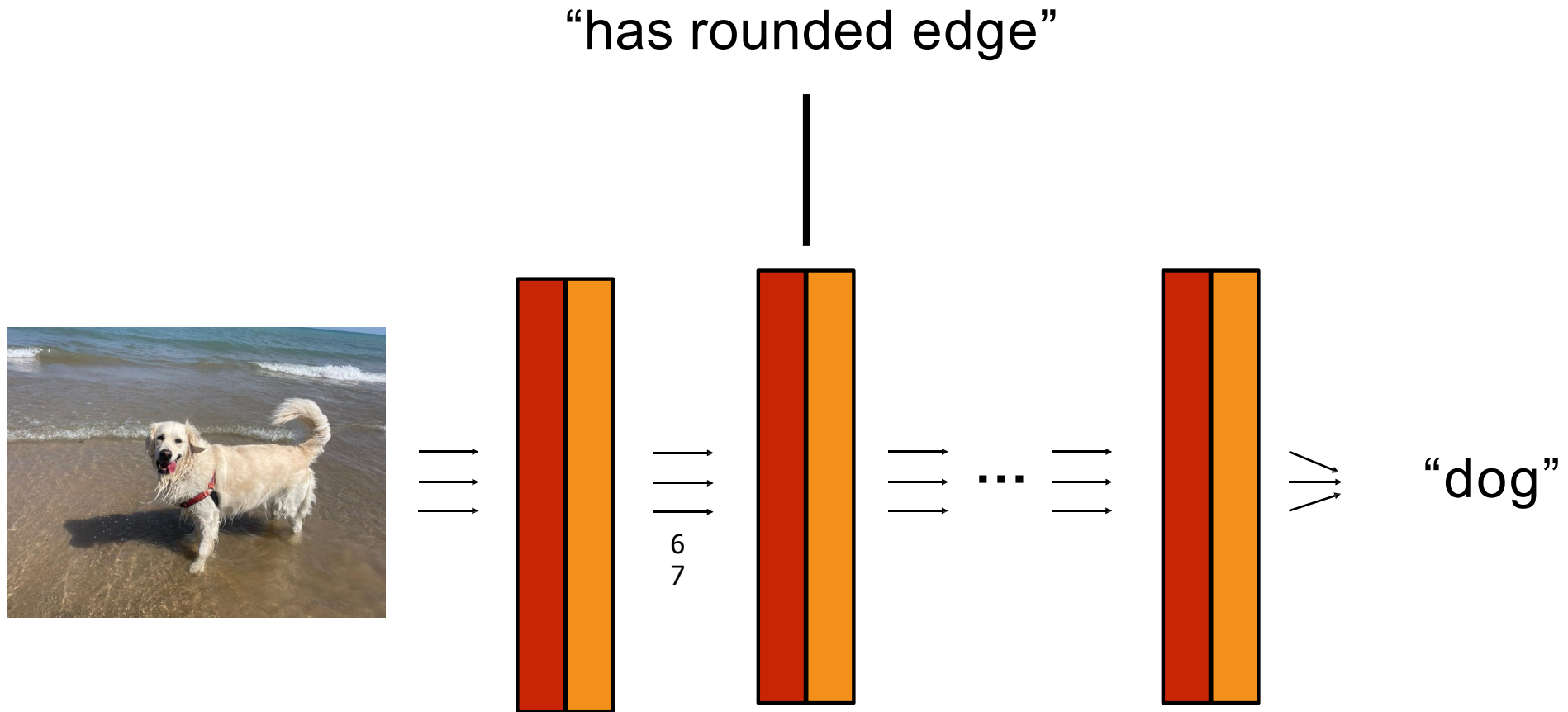


$$f(x) = f_L( \dots f_3(f_2(f_1(x))) )$$

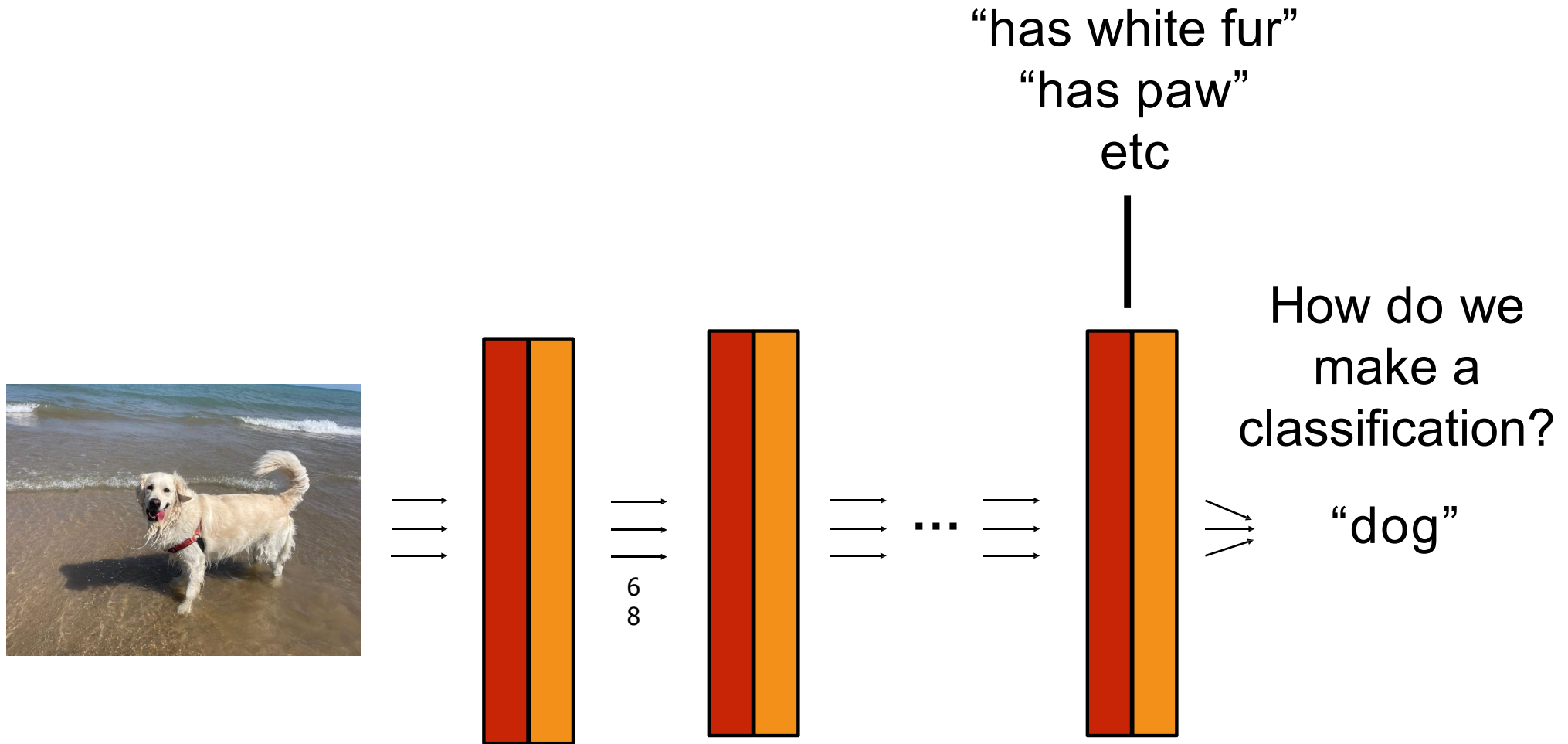
# Deep nets - Intuition



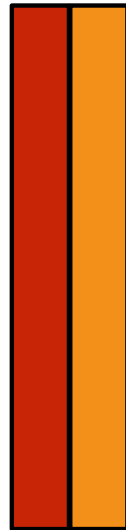
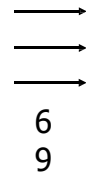
# Deep nets - Intuition



# Deep nets - Intuition



# Deep nets - Intuition

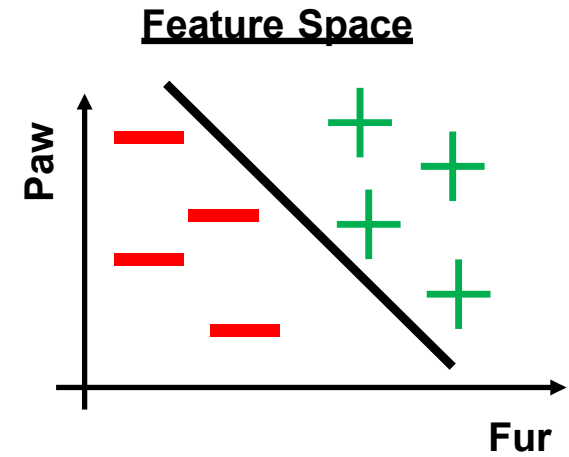


“dog”

Classify

“has white fur”  
“has paw”  
etc

Recall:



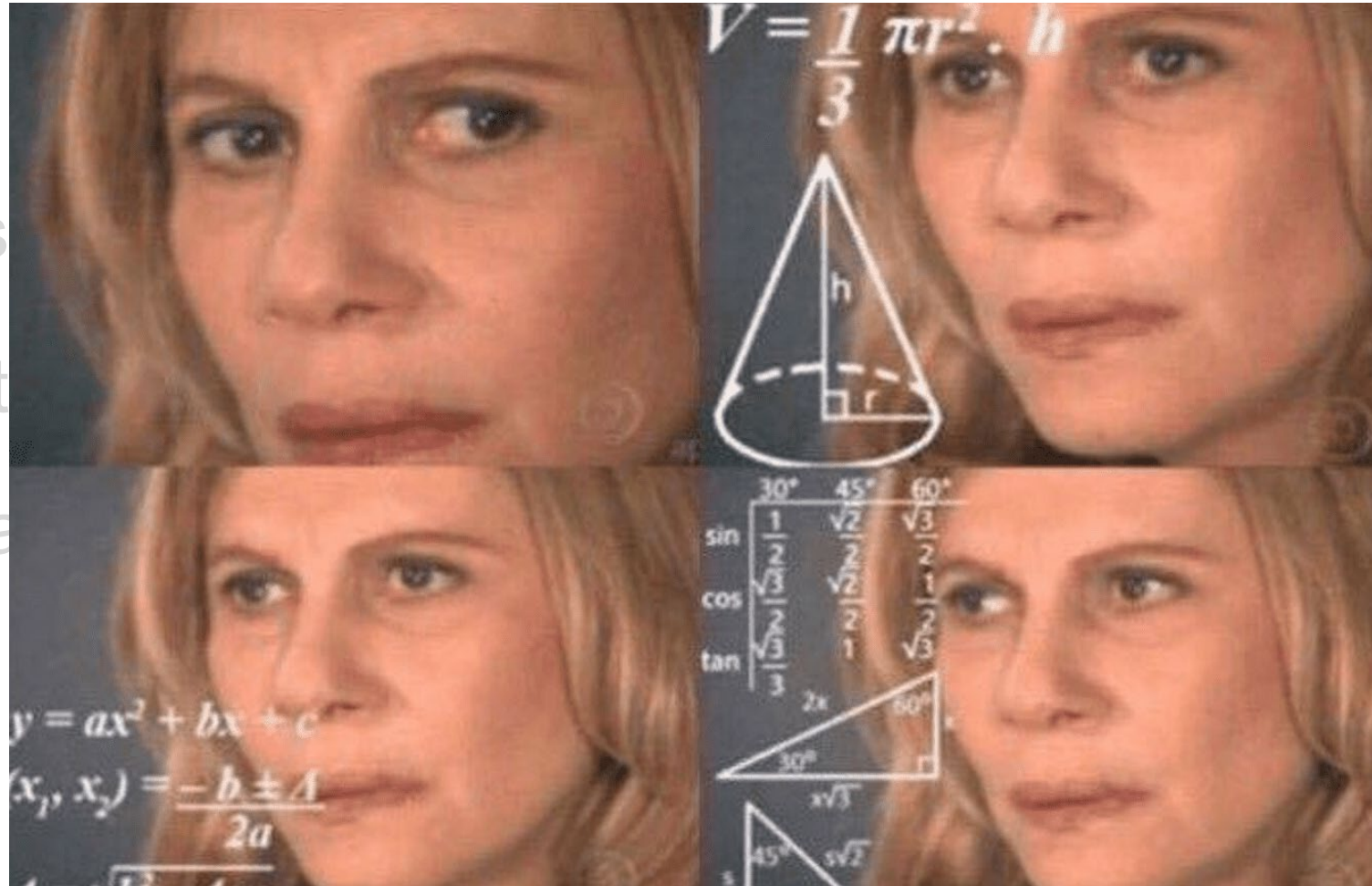
# Computation has a simple form

- 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?

# Computation has a simple form

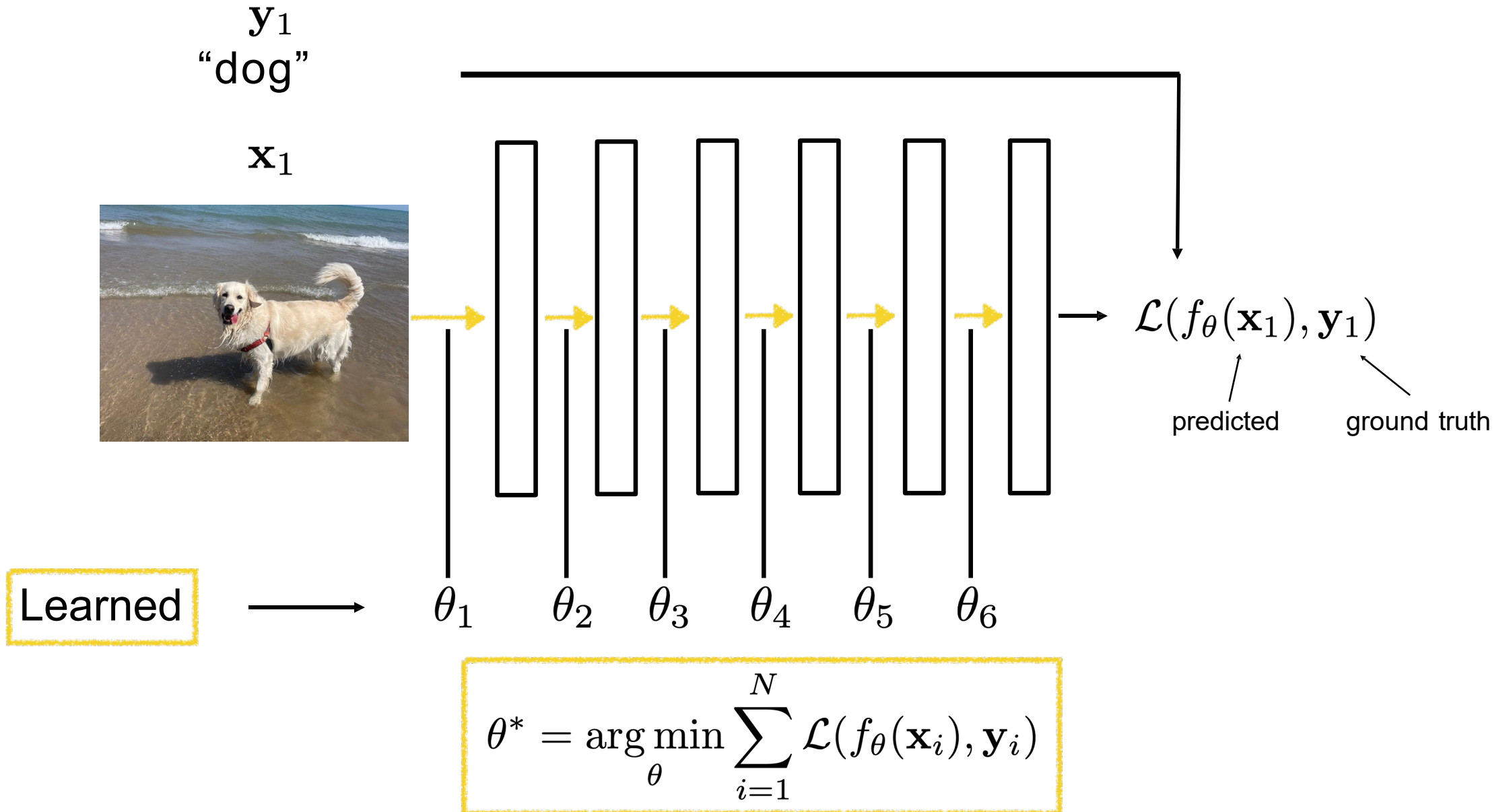
- Compos
  - E.g. mat
  - Do a ma
- n between  
afterwards  
o 0, repeat



But where do we get the weights from?



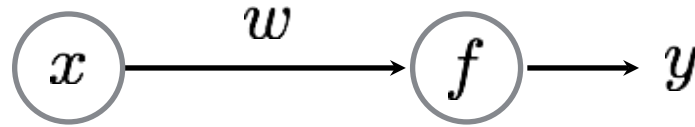
# How would we learn the parameters?





Let's start easy

world's smallest neural network!  
(a “perceptron”)



$$y = wx$$

(a.k.a. line equation, linear regression)

# Training a Neural Network

Given a set of samples and a Perceptron

$$\{x_i, y_i\}$$

$$y = f_{\text{PER}}(x; w)$$

Estimate the parameter of the Perceptron

$$w$$

Given training data:

$x$	$y$
10	10.1
2	1.9
3.5	3.4
1	1.1

*What do you think the weight parameter is?*

$$y = wx$$

Given training data:

$x$	$y$
10	10.1
2	1.9
3.5	3.4
1	1.1

*What do you think the weight parameter is?*

$$y = wx$$

not so obvious as the network gets more complicated so we use ...

# An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$



# An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$

Modify weight  $w$  such that  $\hat{y}$  gets **'closer'** to  $y$

# An Incremental Learning Strategy

(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

and a perceptron

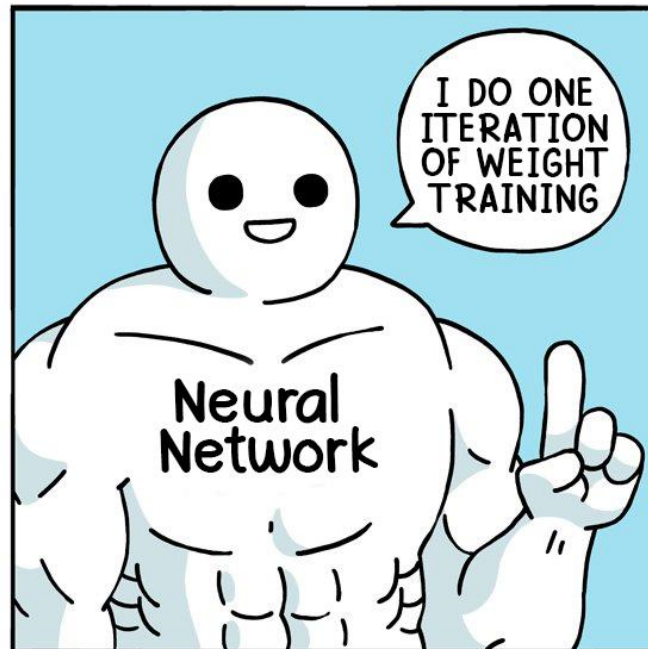
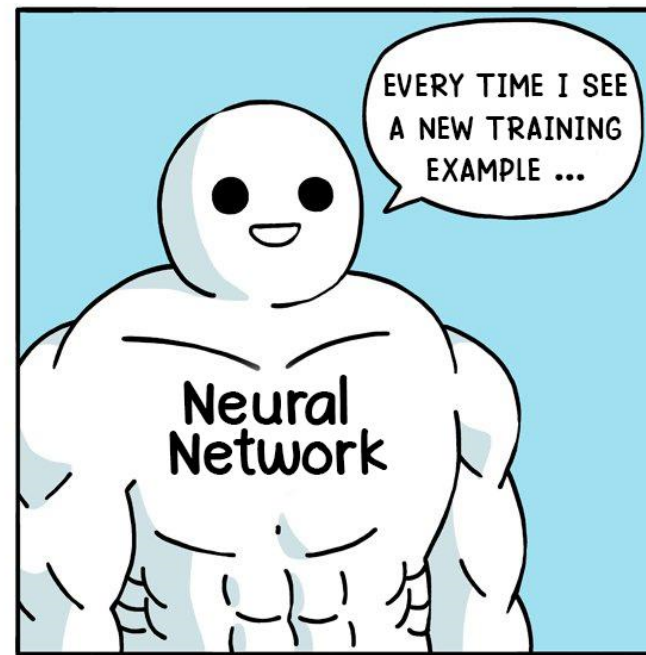
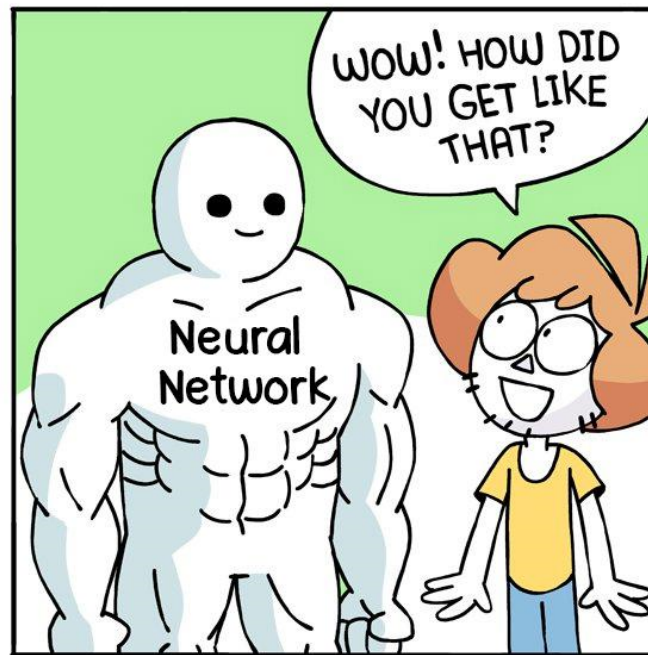
$$\hat{y} = wx$$

Modify weight  $w$  such that  $\hat{y}$  gets **'closer'** to  $y$

perceptron  
parameter

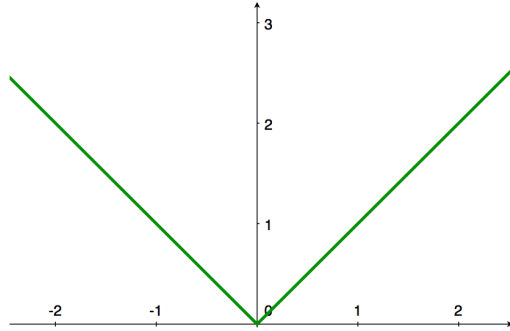
perceptron  
output

true  
label



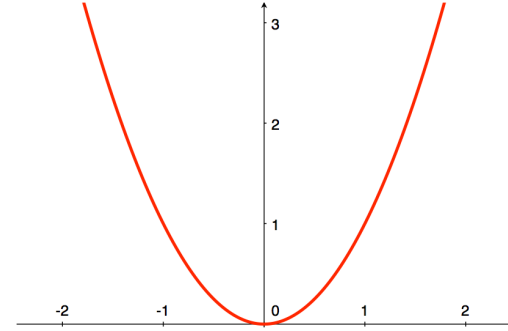
## L1 Loss

$$\ell(\hat{y}, y) = |\hat{y} - y|$$



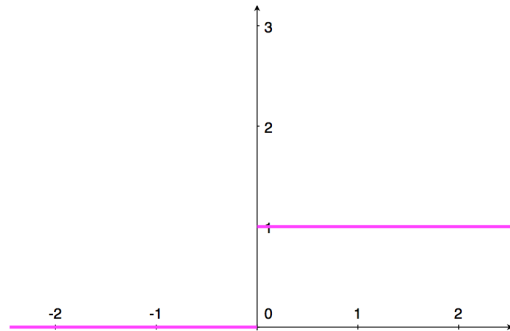
## L2 Loss

$$\ell(\hat{y}, y) = (\hat{y} - y)^2$$



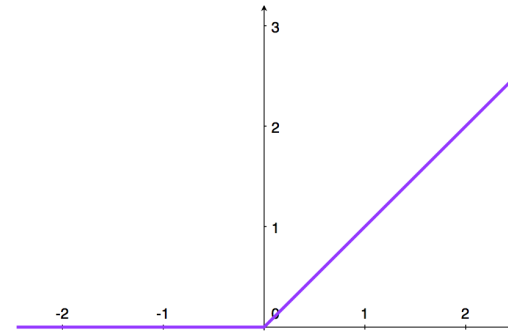
## Zero-One Loss

$$\ell(\hat{y}, y) = \mathbf{1}[\hat{y} \neq y]$$

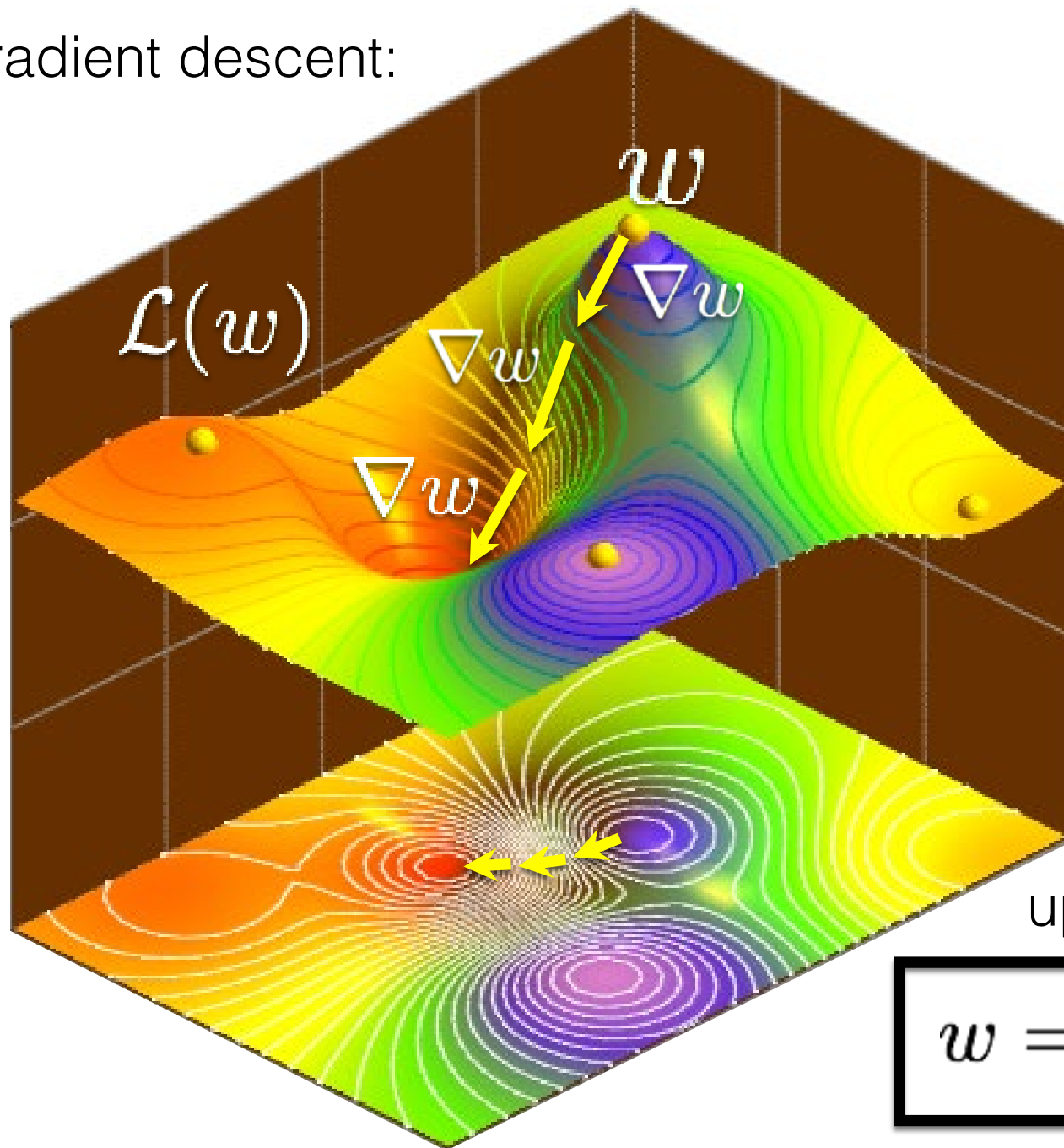


## Hinge Loss

$$\ell(\hat{y}, y) = \max(0, 1 - y \cdot \hat{y})$$



Gradient descent:



update rule:

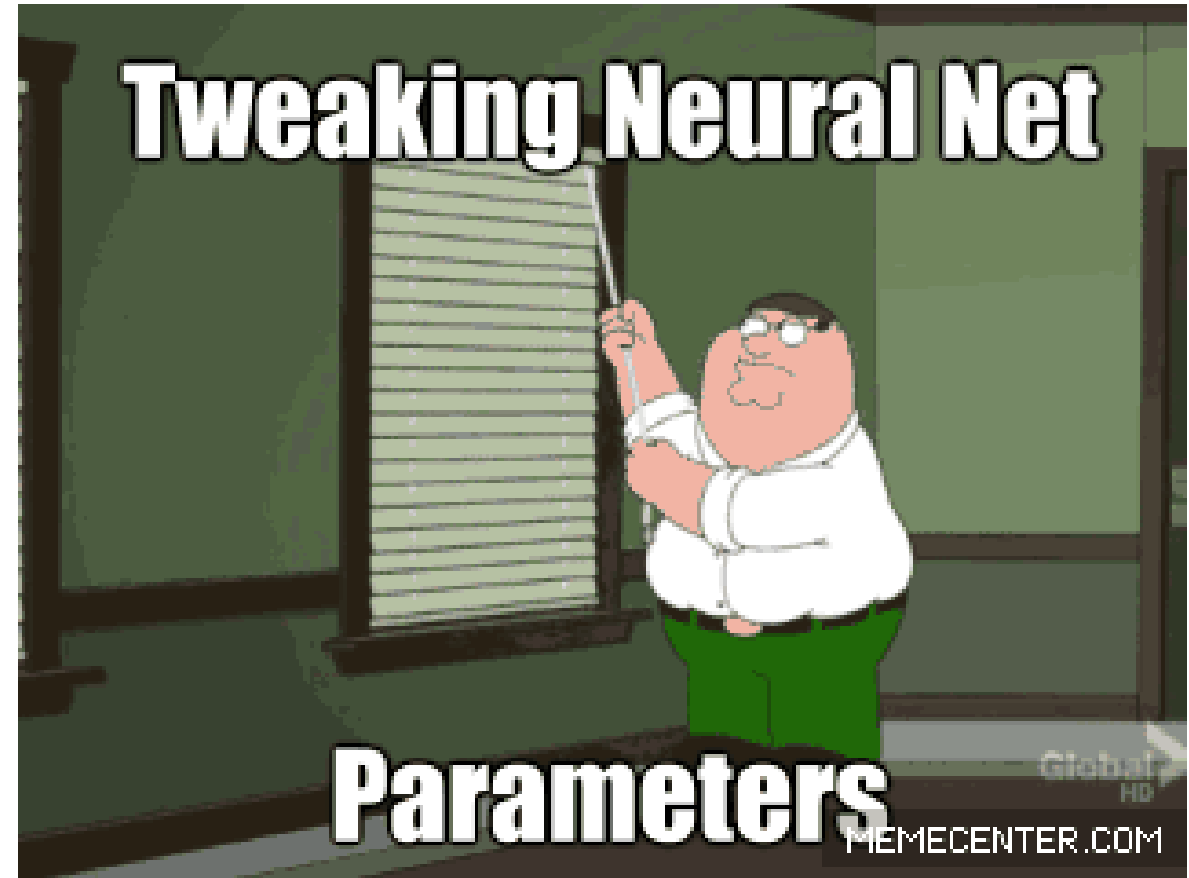
$$w = w - \nabla w$$

# Backpropagation

Geoff Hinton after writing the paper on backprop in 1986



Backpropagation





$\frac{d\mathcal{L}}{dw}$  ...is the rate at which **this** will change...

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

the loss function

... per unit change of **this**

$$y = wx$$

the weight parameter

Let's compute the derivative...



Compute the derivative

$$\begin{aligned}\frac{d\mathcal{L}}{dw} &= \frac{d}{dw} \left\{ \frac{1}{2} (y - \hat{y})^2 \right\} \\ &= -(y - \hat{y}) \frac{dw x}{dw} \\ &= -(y - \hat{y}) x = \nabla w\end{aligned}$$

That means the weight update for **gradient descent** is:

$$\begin{aligned}w &= w - \nabla w \quad \text{move in direction of negative gradient} \\ &= w + (y - \hat{y}) x\end{aligned}$$

## Gradient Descent (world's smallest perceptron)

For each sample

$$\{x_i, y_i\}$$

1. Predict

a. Forward pass

$$\hat{y} = wx_i$$

b. Compute Loss

$$\mathcal{L}_i = \frac{1}{2}(y_i - \hat{y})^2$$

2. Update

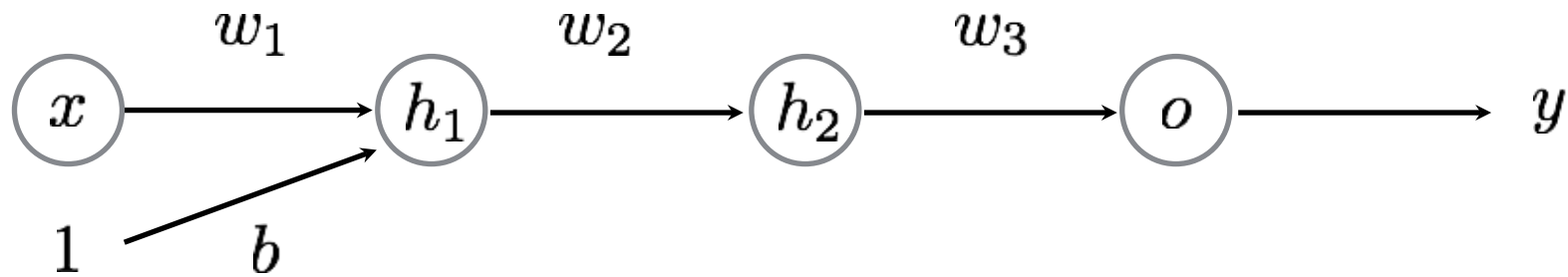
a. Back Propagation

$$\frac{d\mathcal{L}_i}{dw} = -(y_i - \hat{y})x_i = \nabla w$$

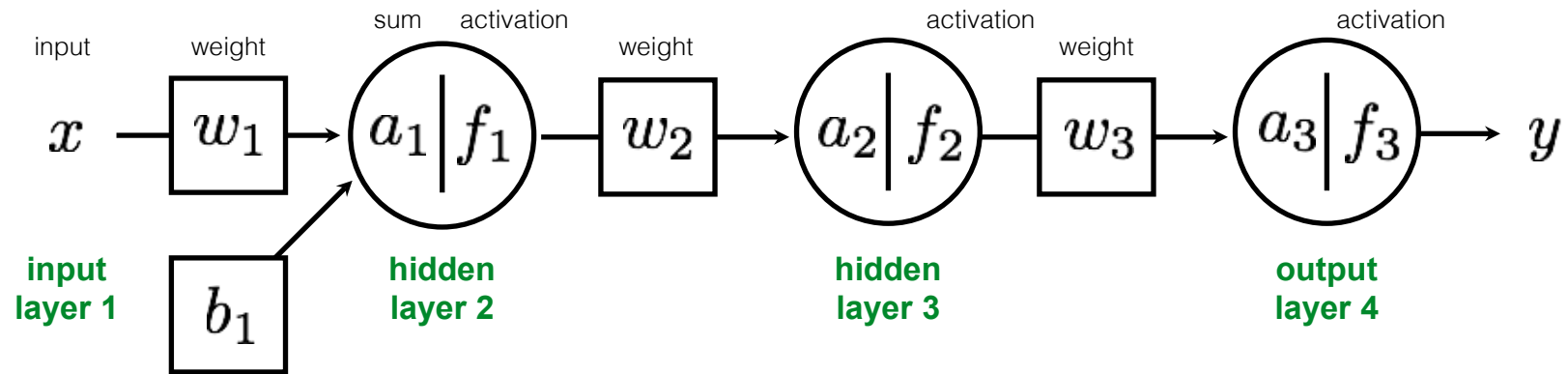
b. Gradient update

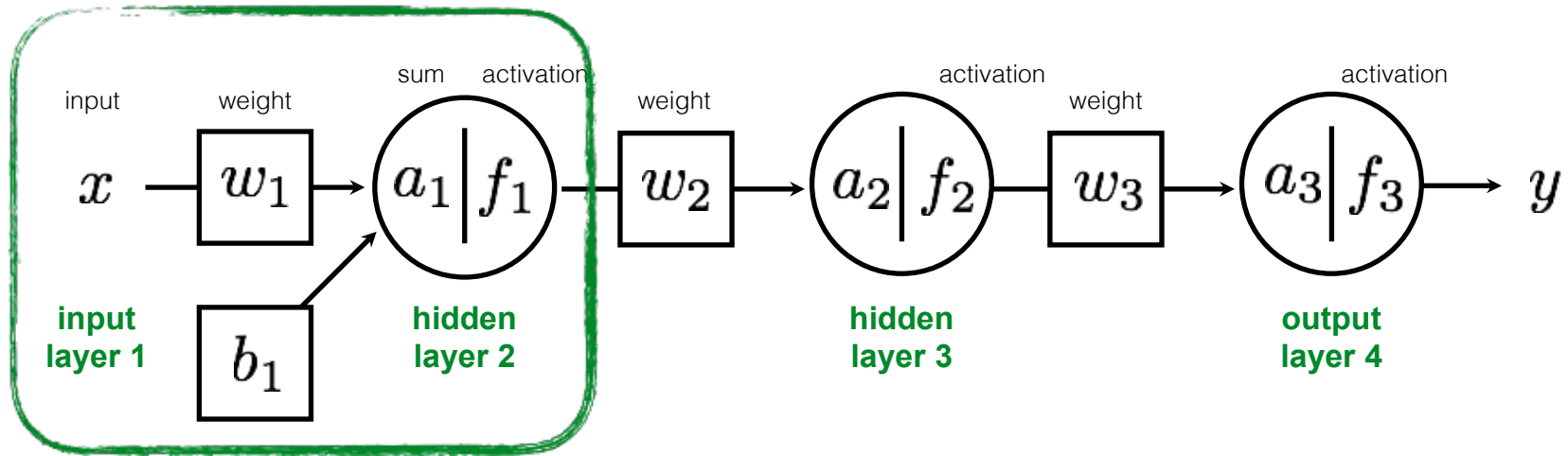
$$w = w - \nabla w$$

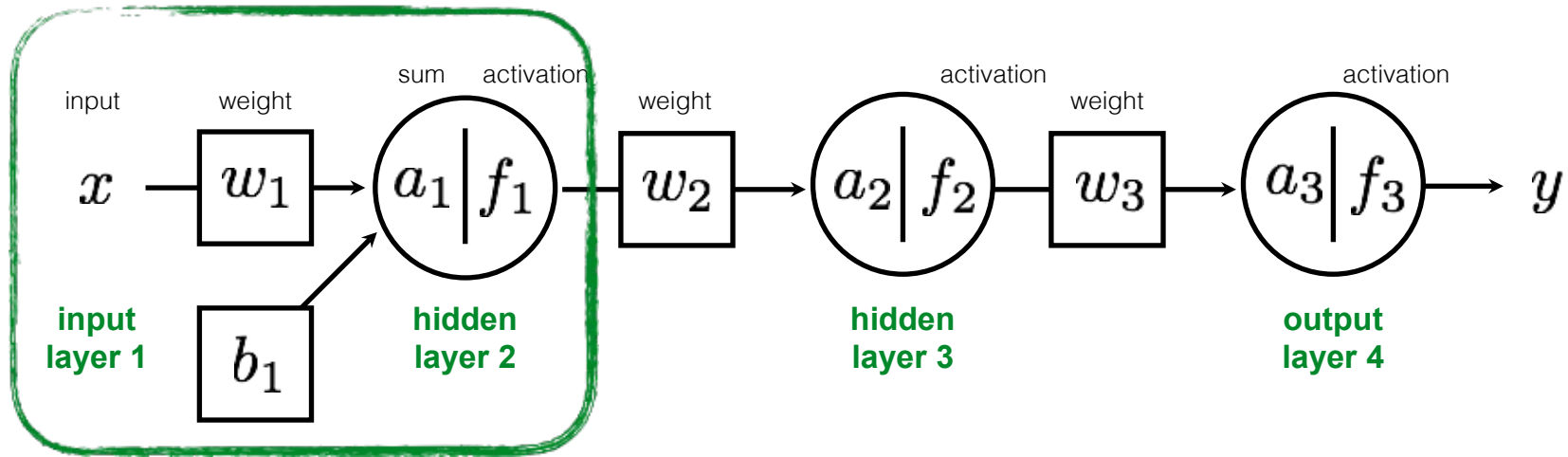
# multi-layer perceptron



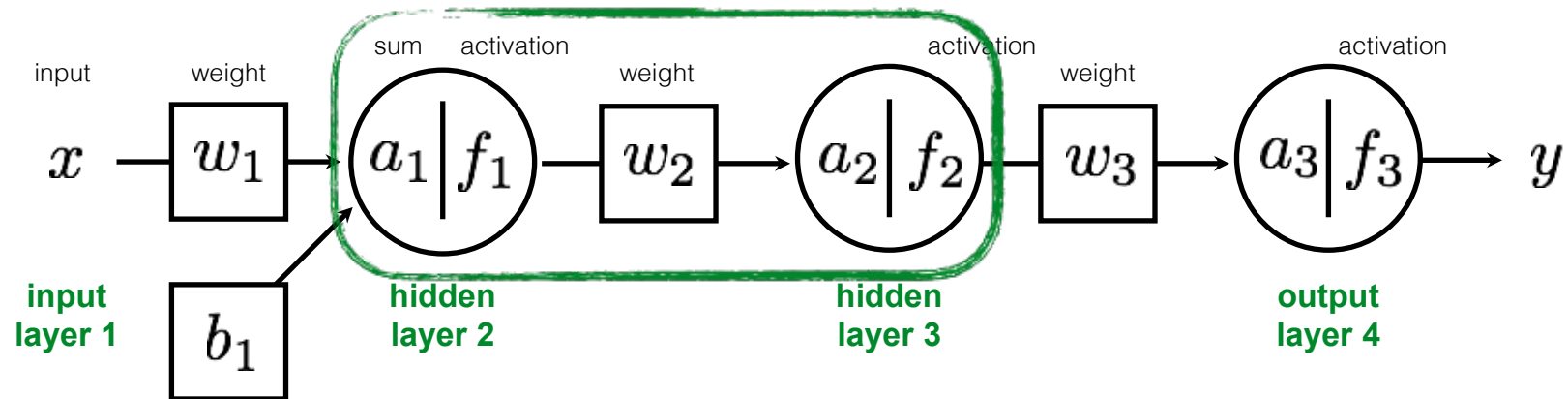
function of **FOUR** parameters and **FOUR** layers!



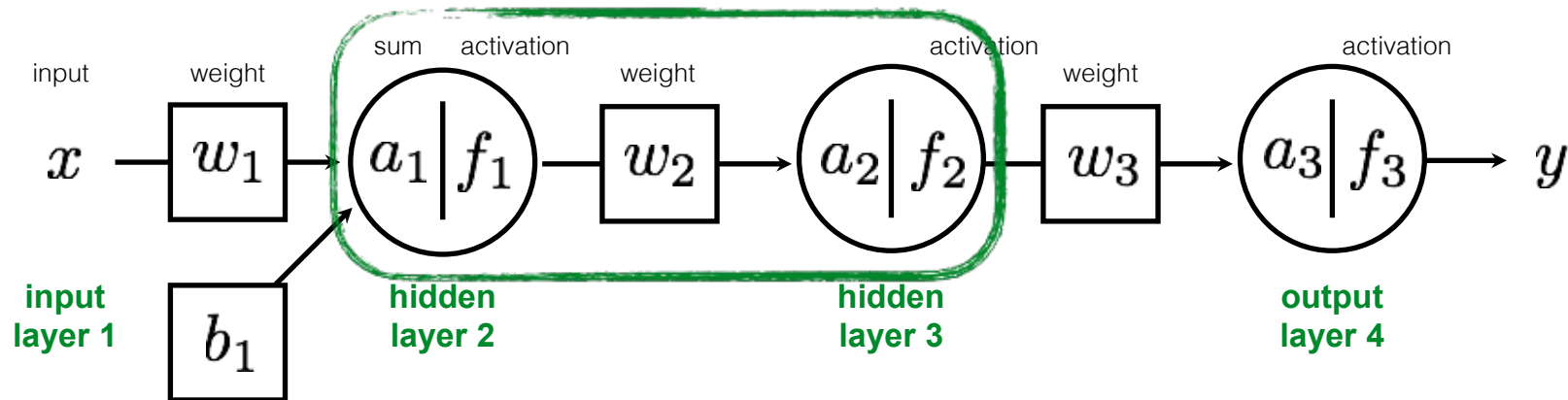




$$a_1 = w_1 \cdot x + b_1$$



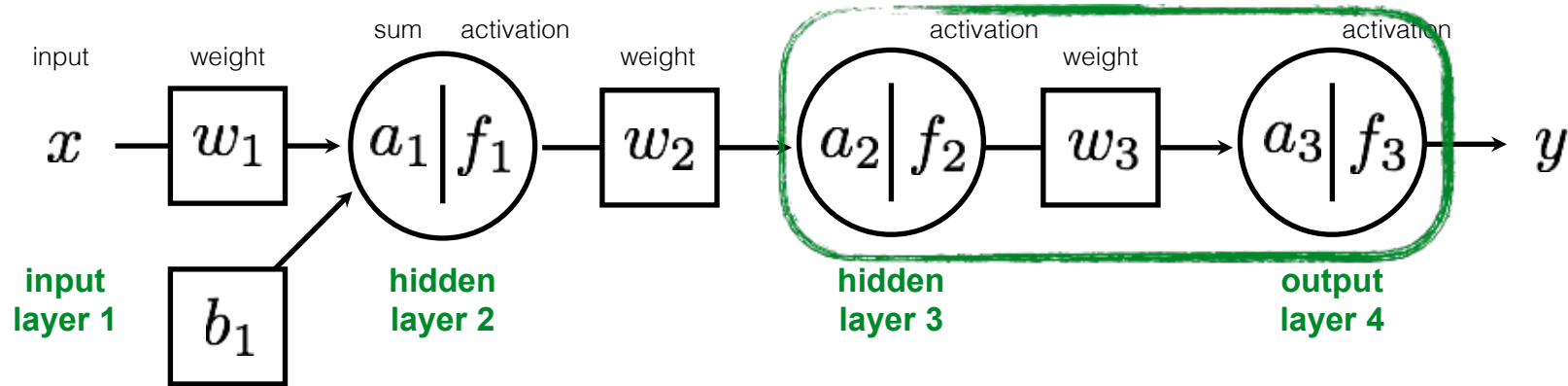
$$a_1 = w_1 \cdot x + b_1$$



$$a_1 = w_1 \cdot x + b_1$$

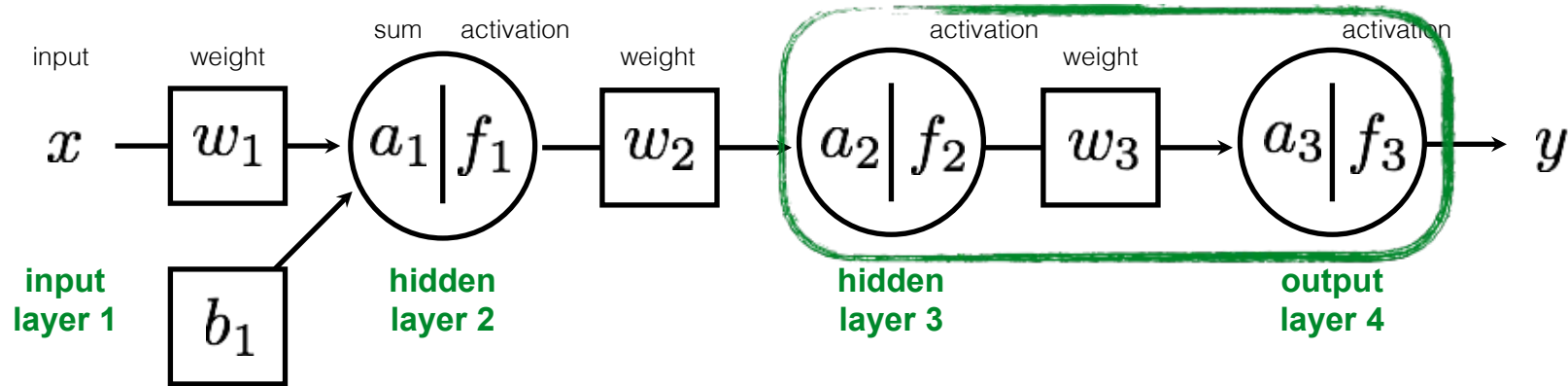
$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$





$$a_1 = w_1 \cdot x + b_1$$

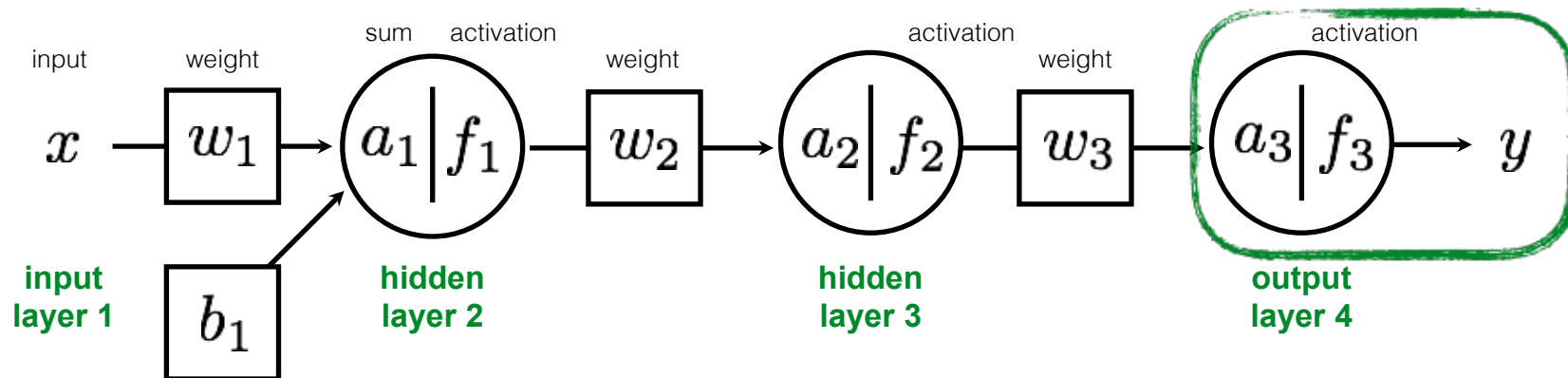
$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

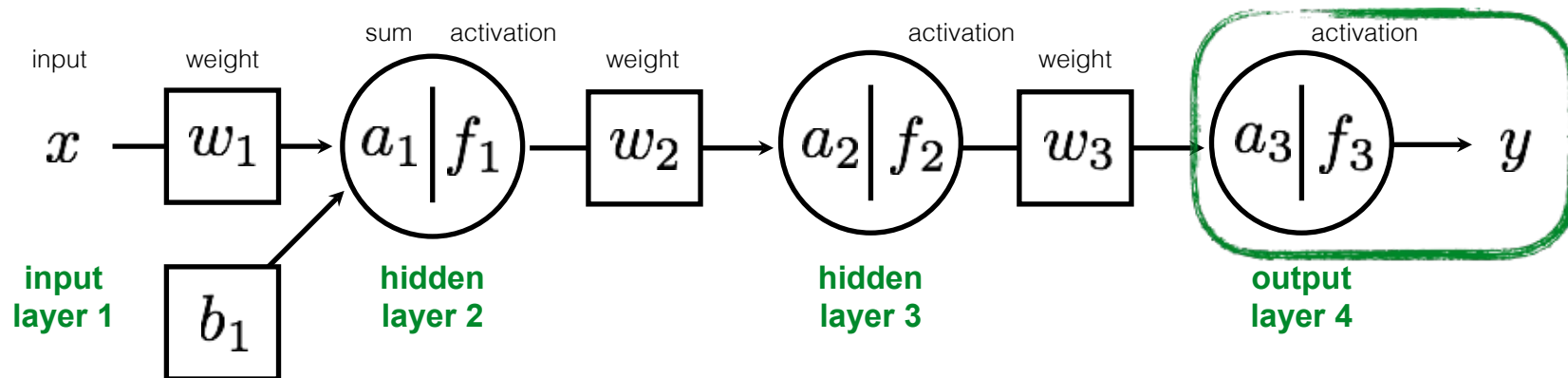
$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$



$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

We need to train the network:

*What is known? What is unknown?*

Entire network can be written out as a long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$



known

We need to train the network:

*What is known? What is unknown?*

Entire network can be written out as a long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

activation function  
sometimes has  
parameters

**unknown**



We need to train the network:

*What is known? What is unknown?*

# Learning an MLP

Given a set of samples and a MLP

$$\{x_i, y_i\}$$

$$y = f_{\text{MLP}}(x; \theta)$$

Estimate the parameters of the MLP

$$\theta = \{f, w, b\}$$



## Gradient Descent

For each **random** sample  $\{x_i, y_i\}$

1. Predict

a. Forward pass  $\hat{y} = f_{\text{MLP}}(x_i; \theta)$

b. Compute Loss

2. Update

a. Back Propagation

$$\frac{\partial \mathcal{L}}{\partial \theta}$$

vector of parameter partial derivatives

b. Gradient update

$$\theta \leftarrow \theta - \eta \nabla \theta$$

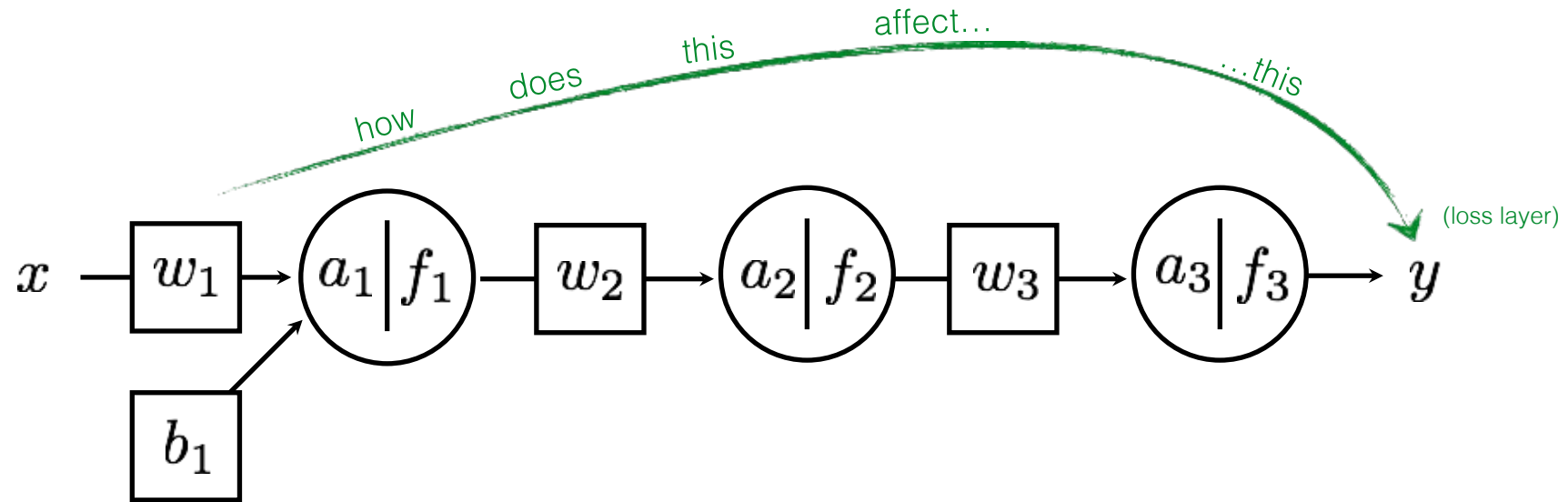
vector of parameter update equations

So we need to compute the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \left[ \frac{\partial \mathcal{L}}{\partial w_3} \frac{\partial \mathcal{L}}{\partial w_2} \frac{\partial \mathcal{L}}{\partial w_1} \frac{\partial \mathcal{L}}{\partial b} \right]$$

Remember,

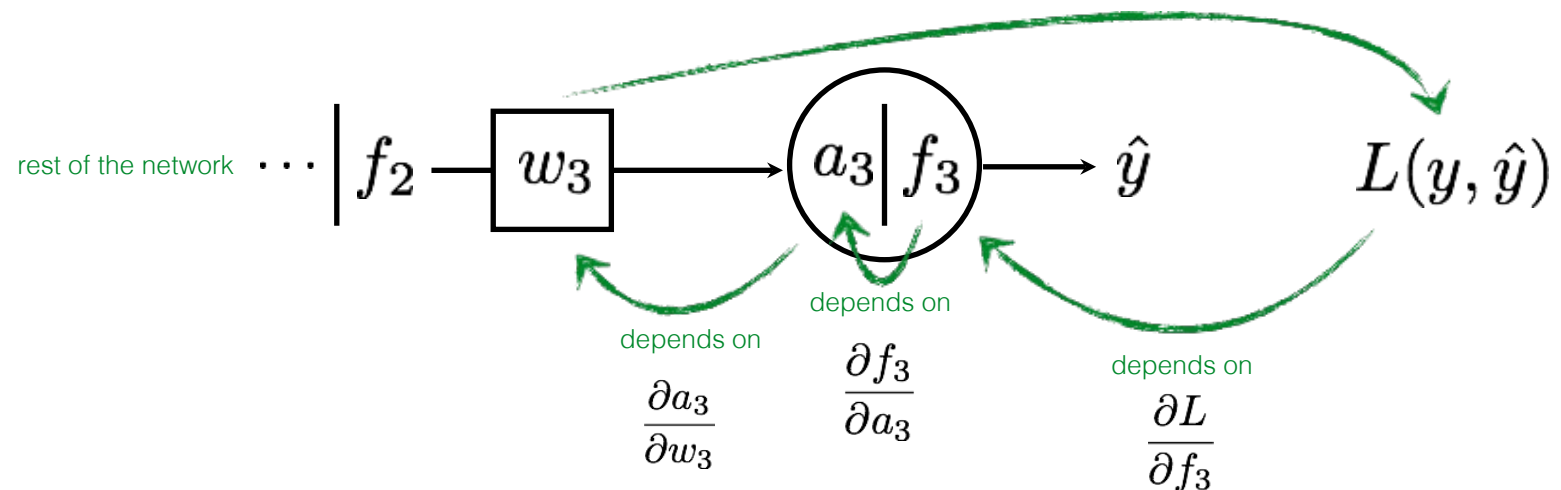
Partial derivative  $\frac{\partial L}{\partial w_1}$  describes...

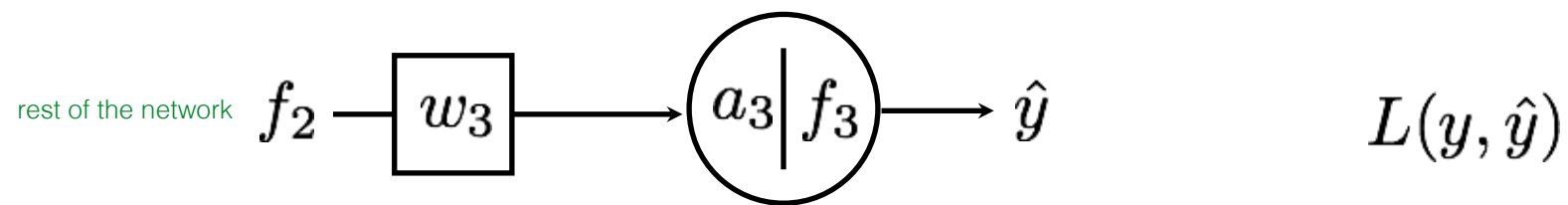


According to the chain rule...

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

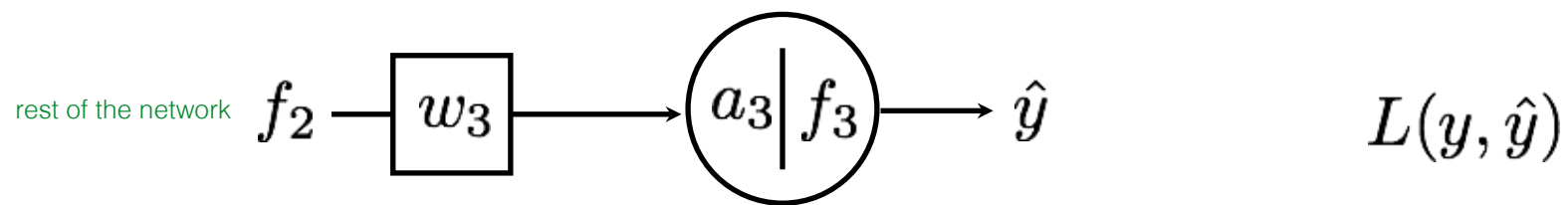
Intuitively, the effect of weight on loss function :  $\frac{\partial L}{\partial w_3}$





$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

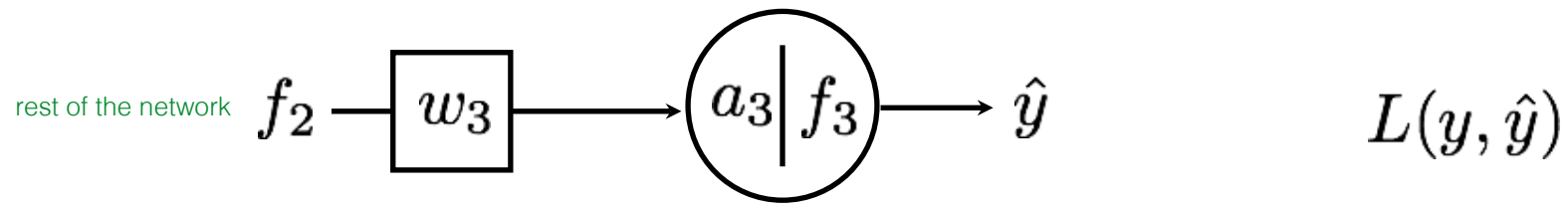
Chain Rule!



$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

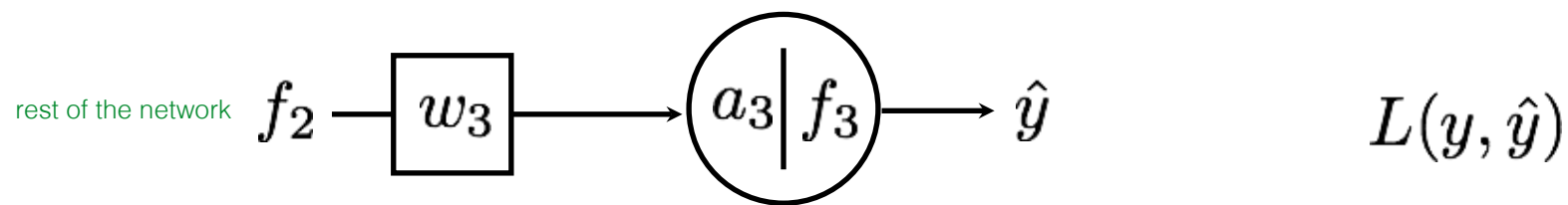
Just the partial  
derivative of L2 loss



$$\begin{aligned}\frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ &= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}\end{aligned}$$

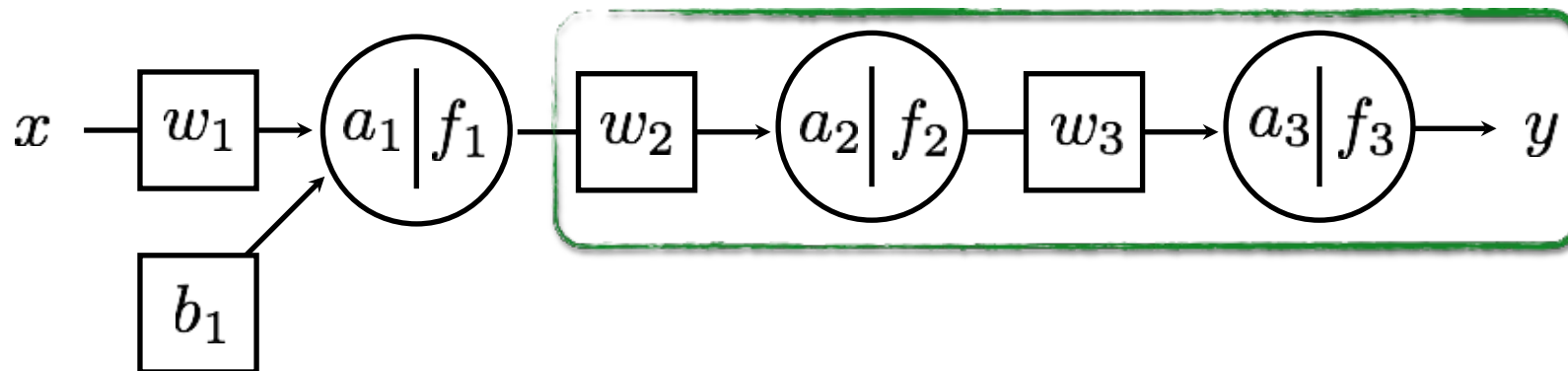
Let's use a Sigmoid function

$$\frac{ds(x)}{dx} = s(x)(1 - s(x))$$

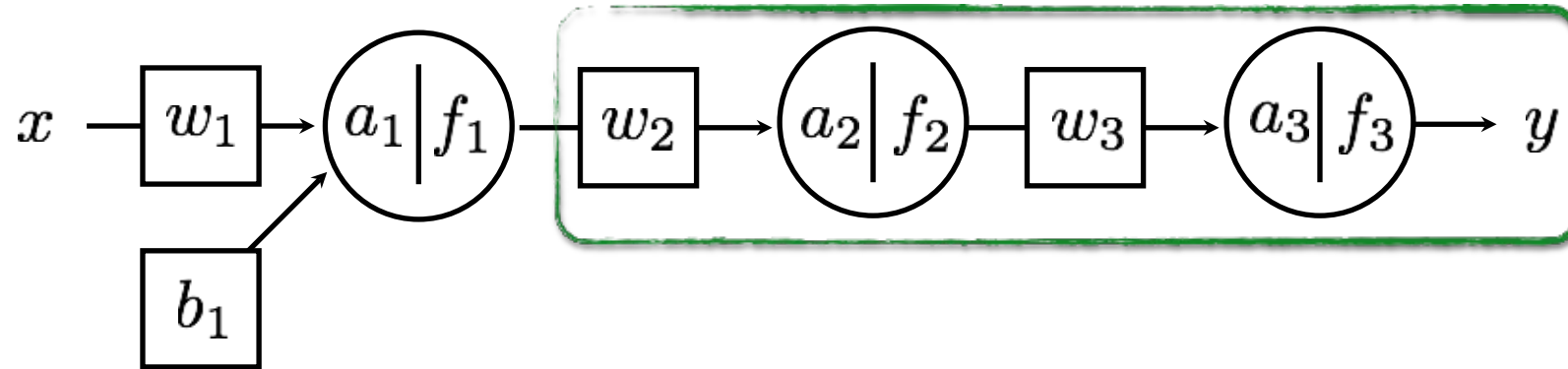


$$\begin{aligned}
 \frac{\partial L}{\partial w_3} &= \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) f_3(1 - f_3) \frac{\partial a_3}{\partial w_3} \\
 &= -\eta(y - \hat{y}) f_3(1 - f_3) f_2
 \end{aligned}$$





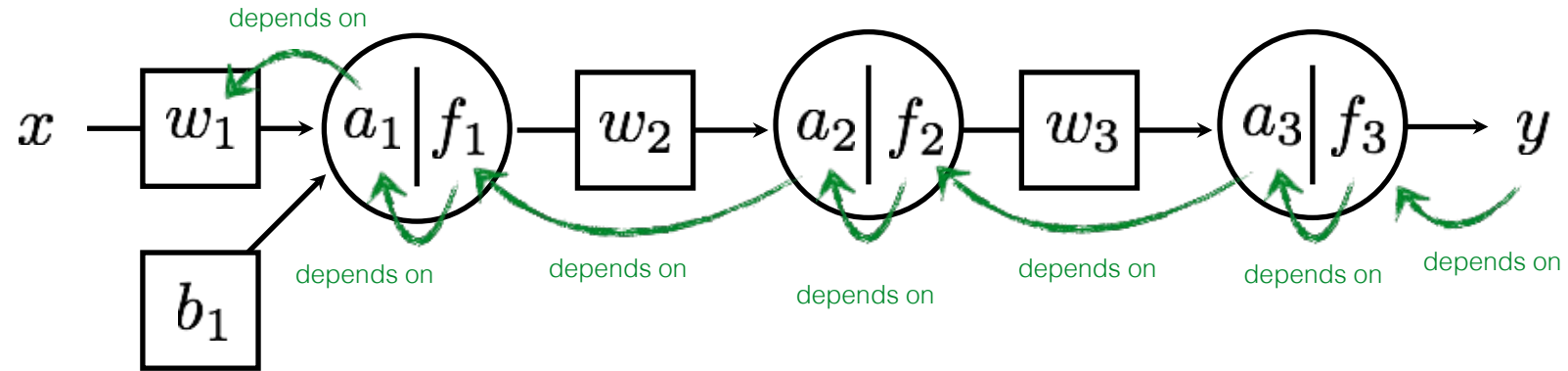
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

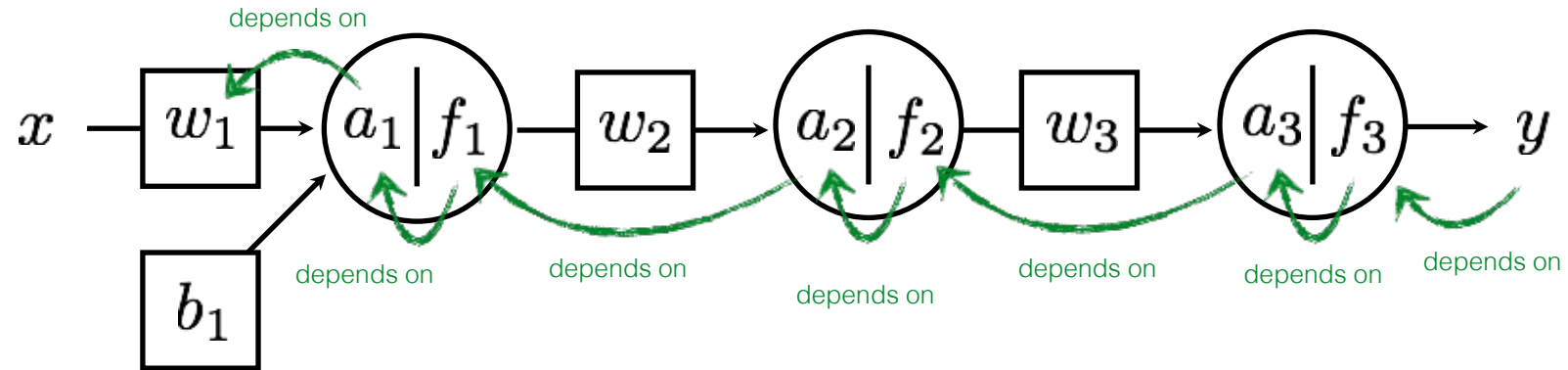
already computed.  
re-use (propagate)!

The chain rule says...



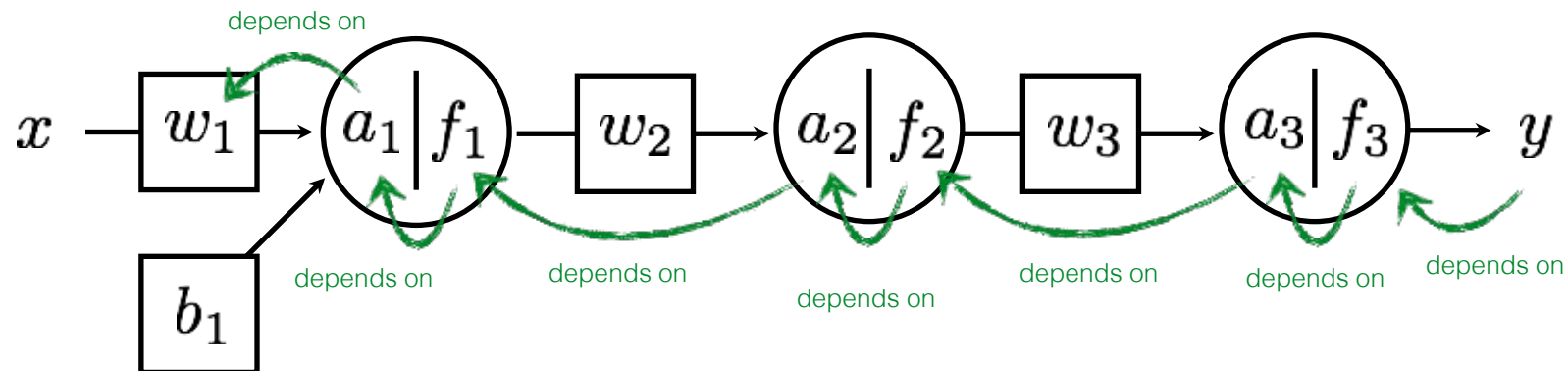
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

The chain rule says...

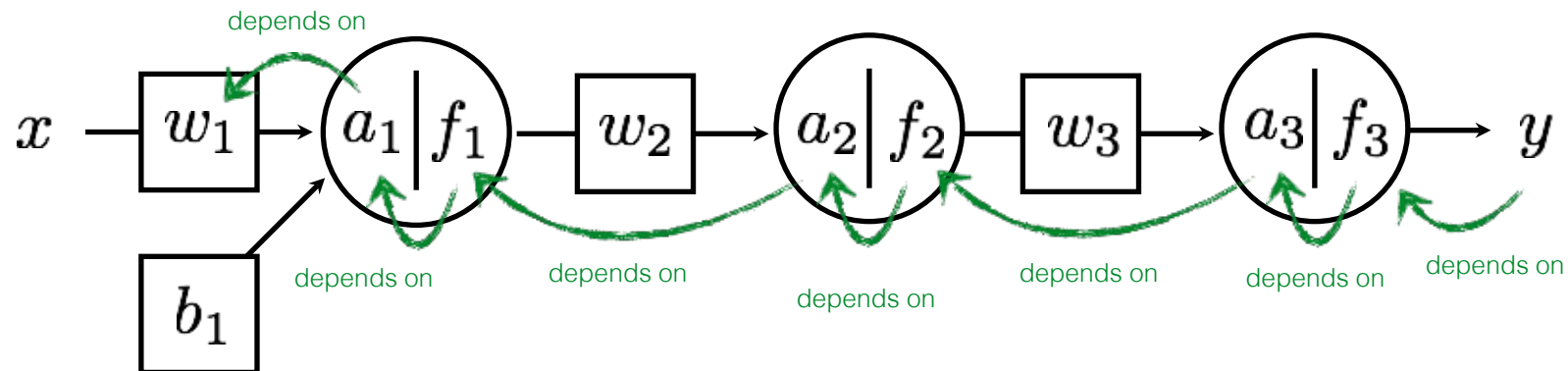


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

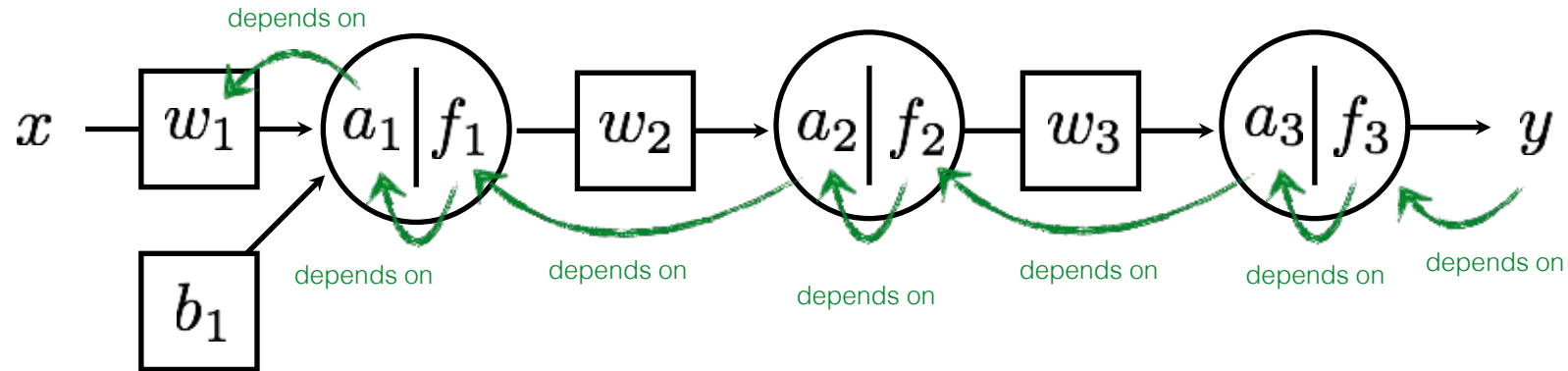
already computed.  
re-use (propagate)!



$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\
 \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\
 \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}
 \end{aligned}$$



$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\
 \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\
 \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\
 \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}
 \end{aligned}$$



$$\frac{\partial \mathcal{L}}{\partial w_3} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2}$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}$$

# Gradient Descent

For each example sample  $\{x_i, y_i\}$

1. Predict

a. Forward pass  $\hat{y} = f_{\text{MLP}}(x_i; \theta)$

b. Compute Loss  $\mathcal{L}_i$

2. Update

a. Back Propagation

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_3} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3} \\ \frac{\partial \mathcal{L}}{\partial w_2} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} \\ \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial \mathcal{L}}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial f_2} \frac{\partial f_2}{\partial a_2} \frac{\partial a_2}{\partial f_1} \frac{\partial f_1}{\partial a_1} \frac{\partial a_1}{\partial b}\end{aligned}$$

b. Gradient update

$$\begin{aligned}w_3 &= w_3 - \eta \nabla w_3 \\ w_2 &= w_2 - \eta \nabla w_2 \\ w_1 &= w_1 - \eta \nabla w_1 \\ b &= b - \eta \nabla b\end{aligned}$$



# Gradient Descent

For each example sample  $\{x_i, y_i\}$

1. Predict

a. Forward pass

$$\hat{y} = f_{\text{MLP}}(x_i; \theta)$$

b. Compute Loss

$$\mathcal{L}_i$$

2. Update

a. Back Propagation

$$\frac{\partial \mathcal{L}}{\partial \theta}$$

vector of parameter partial derivatives

b. Gradient update

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

vector of parameter update equations

# Gradient Descent

$$\theta = \begin{bmatrix} \frac{\partial L}{\partial w} & \frac{\partial L}{\partial b} \end{bmatrix} \quad \theta = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$

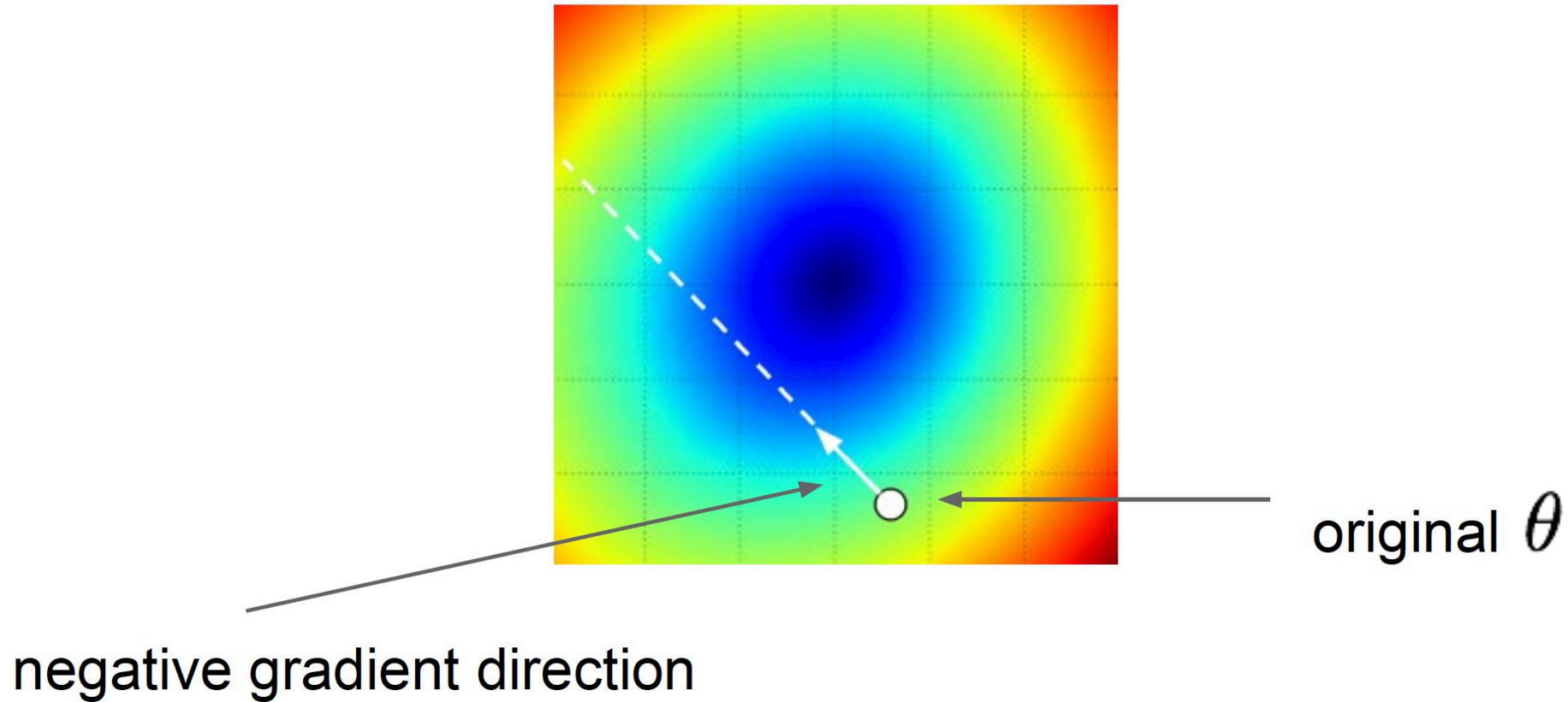
$$\theta = \theta - \eta \theta$$

$$\theta = \theta - \eta \theta$$

$$\theta = \theta - \eta \theta$$



# Learning rates



$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

Step size: learning rate

Too big: will miss the minimum

Too small: slow convergence

# Learning rate scheduling

- Use different learning rate at each iteration.
- Most common choice:

$$\eta_t = \frac{\eta_0}{\sqrt{t}}$$

Need to select initial learning rate  $\eta_0$

More modern choice: Adaptive learning rates.

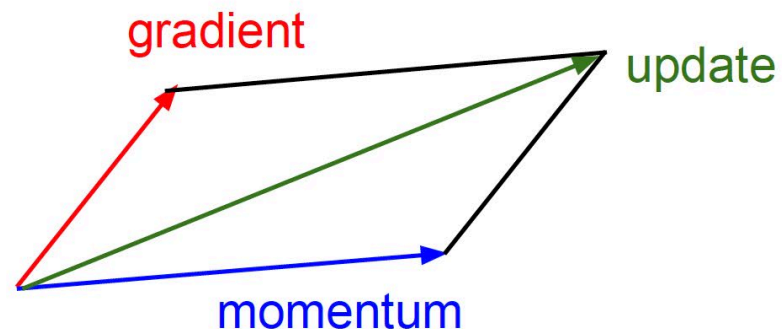
$$\eta_t = G \left( \left\{ \frac{\partial L}{\partial \theta} \right\}_{i=0}^t \right)$$

Many choices for G (Adam, Adagrad, Adadelata).

# Momentum Update

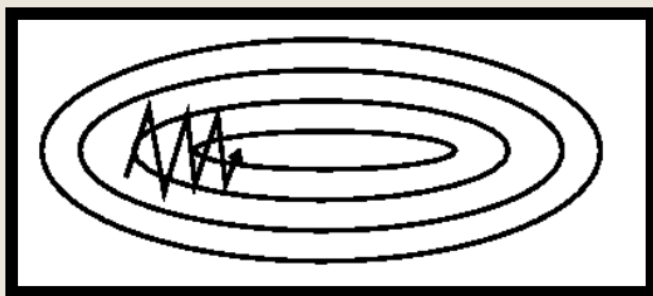
$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

$$\Delta\theta \leftarrow w \frac{\partial L}{\partial \theta} + (1 - w)\Delta\theta$$

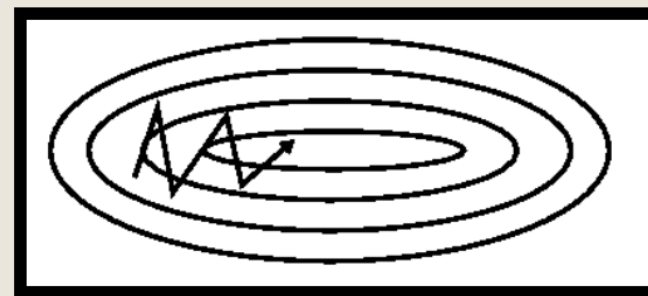


Take direction history  
into account!

```
weights_grad = evaluate_gradient(loss_fun, data, weights)
vel = vel * 0.9 + step_size * weights_grad
weights += vel
```



(Fig. 2a)



(Fig. 2b)

## Many other ways to perform optimization...

- Second order methods that use the Hessian (or its approximation): BFGS, **LBFGS**, etc.
- Currently, the lesson from the trenches is that well-tuned SGD+Momentum is very hard to beat for CNNs.
- No consensus on Adam etc.: Seem to give faster performance to worse local minima.

**CONVOLUTIONAL NEURAL  
NETWORKS**



**IS THERE ANYTHING THEY CAN'T  
DO?**

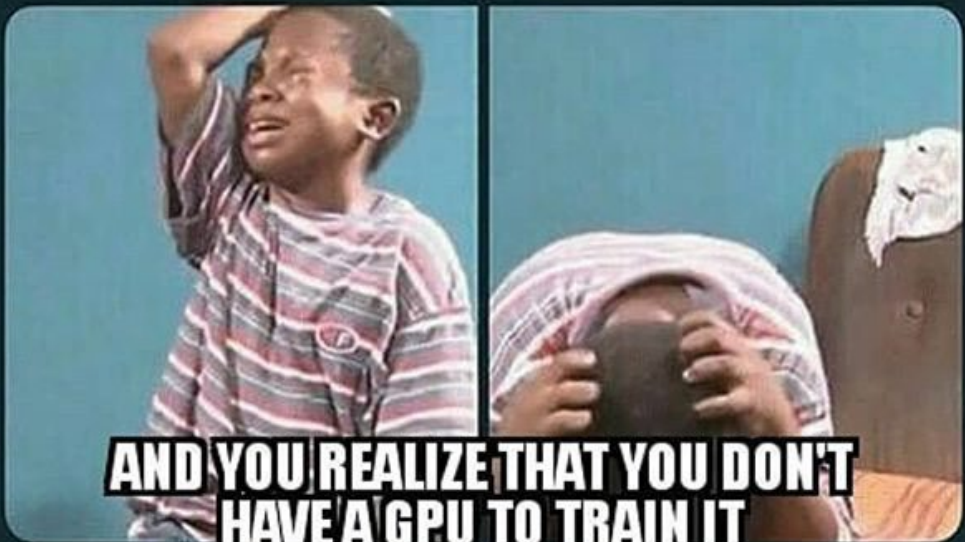
memegenerator.net

**I KNOW**

**CONVOLUTIONAL  
NEURAL NETWORKS**

imgflip.com

**WHEN YOU CREATE A  
CONVOLUTIONAL NEURAL NETWORK**



**AND YOU REALIZE THAT YOU DON'T  
HAVE A GPU TO TRAIN IT**



**DEEP  
LEARNING**



**MATH,  
MATH MATH**



[HOME](#)[MENU](#)[CONNECT](#)[THE LATEST](#)[POPULAR](#)[MOST SHARED](#)

MIT  
Technology  
Review

# 10 BREAKTHROUGH TECHNOLOGIES 2013

[Introduction](#)[The 10 Technologies](#)[Past Years](#)

## Deep Learning

With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.



## Temporary Social Media

Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.



## Prenatal DNA Sequencing

Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?



## Additive Manufacturing

Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.



## Baxter: The Blue-Collar Robot

Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.



## Memory Implants

## Smart Watches

## Ultra-Efficient Solar

## Big Data from

## Supergrids

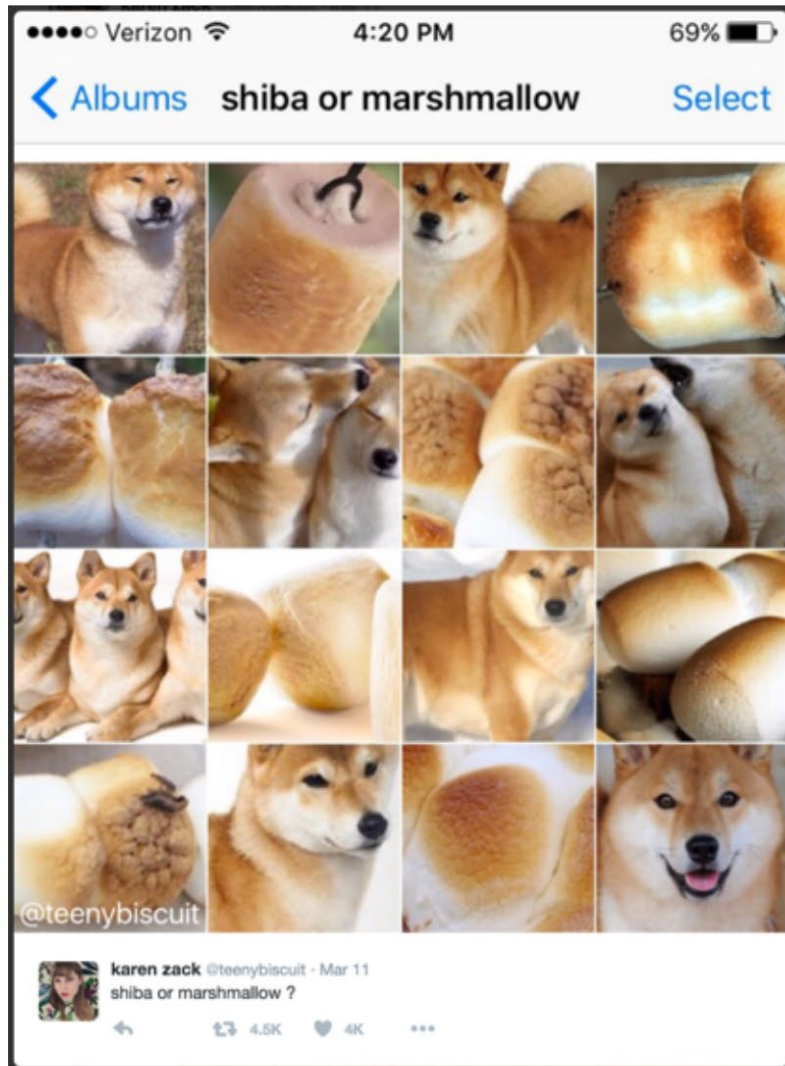


# (Unrelated) Dog vs Food



[Karen Zack, @teenybiscuit]

# (Unrelated) Dog vs Food



[Karen Zack, @teenybiscuit]



# CNNs in 2012: “SuperVision” (aka “AlexNet”)

**“AlexNet”** — Won the ILSVRC2012 Challenge

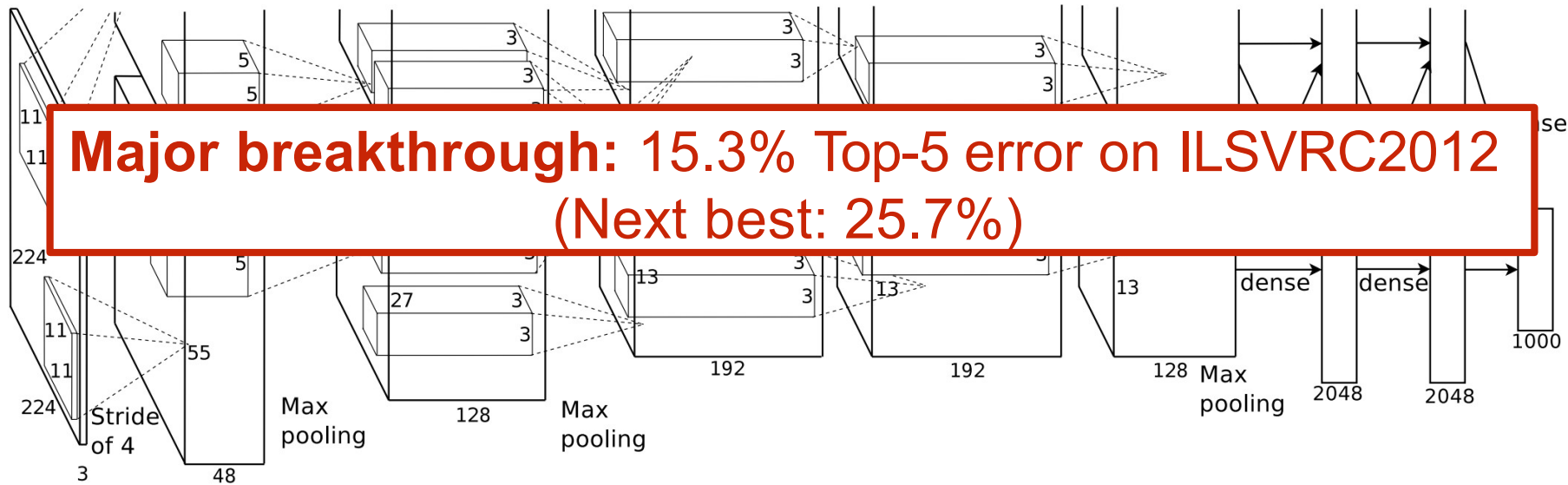
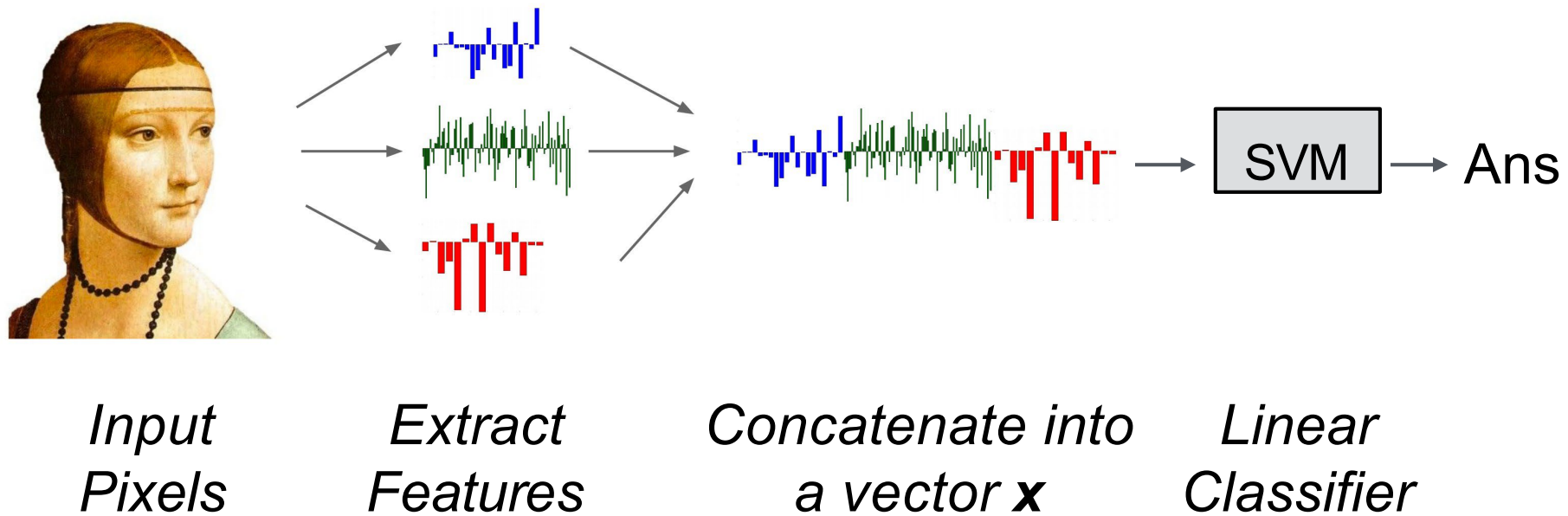


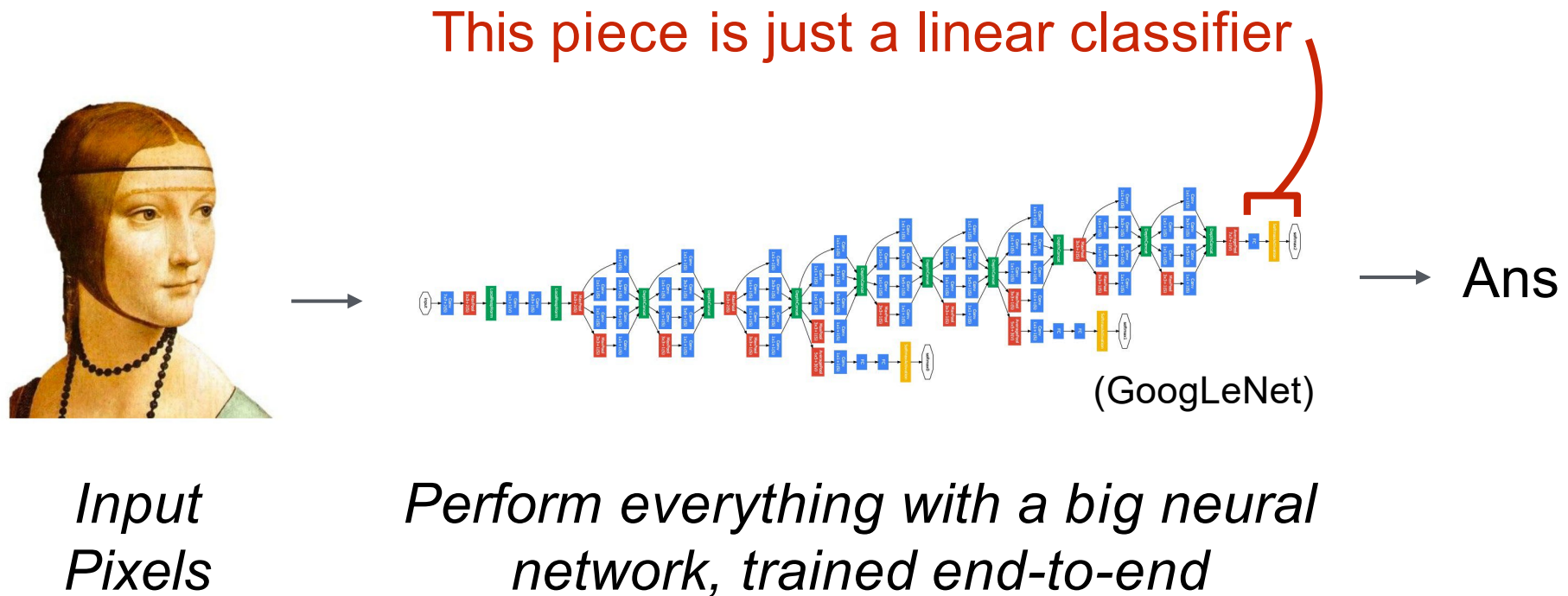
Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network’s input is 150,528-dimensional, and the number of neurons in the network’s remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

[Krizhevsky, Sutskever, Hinton. NIPS 2012]

# Recap: Before Deep Learning



# The last layer of (most) CNNs are linear classifiers

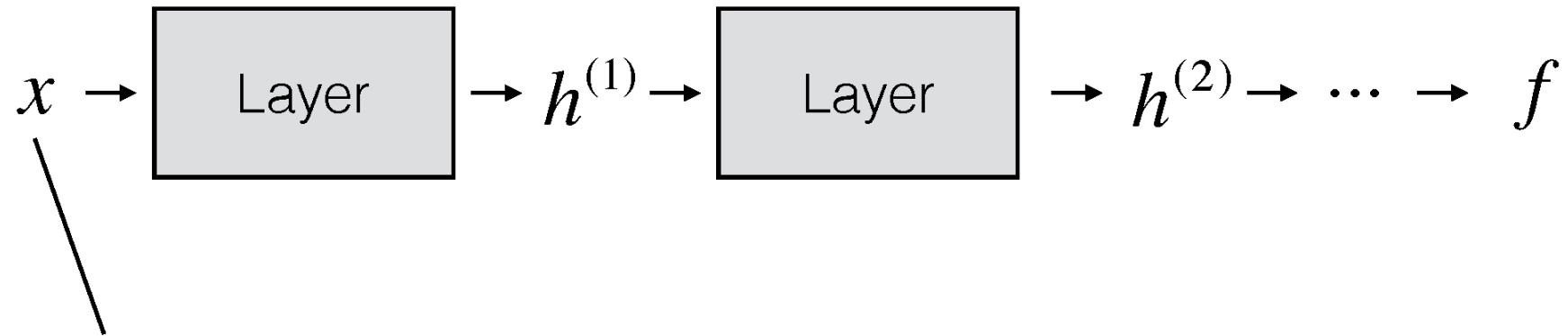


**Key:** perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

# ConvNets

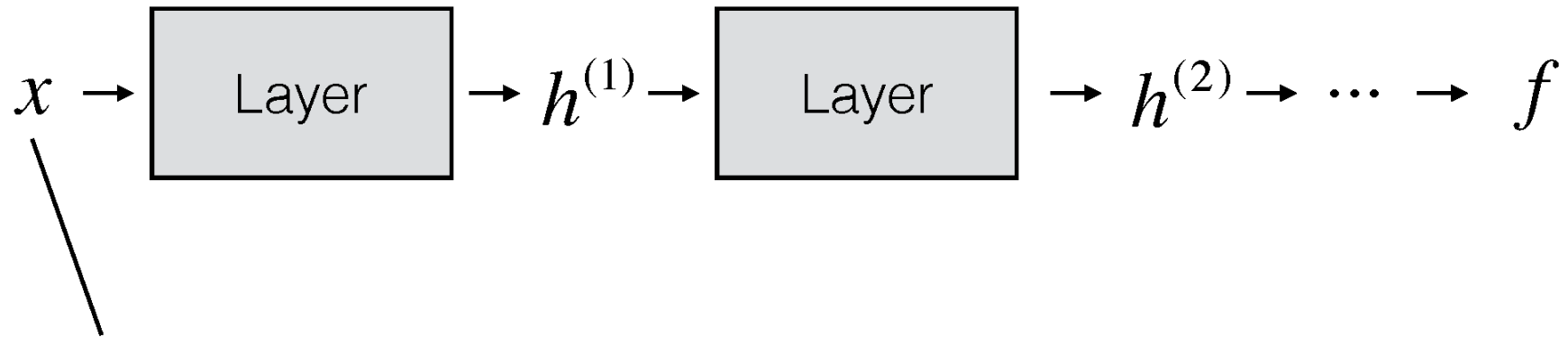
They're just neural networks with  
3D activations and weight sharing

# What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)

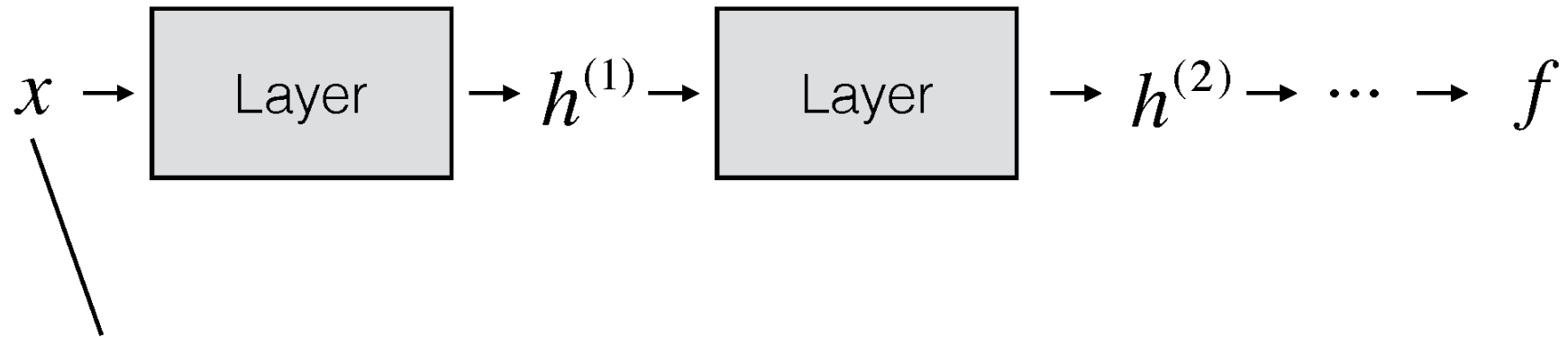
# What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure



# What shape should the activations have?



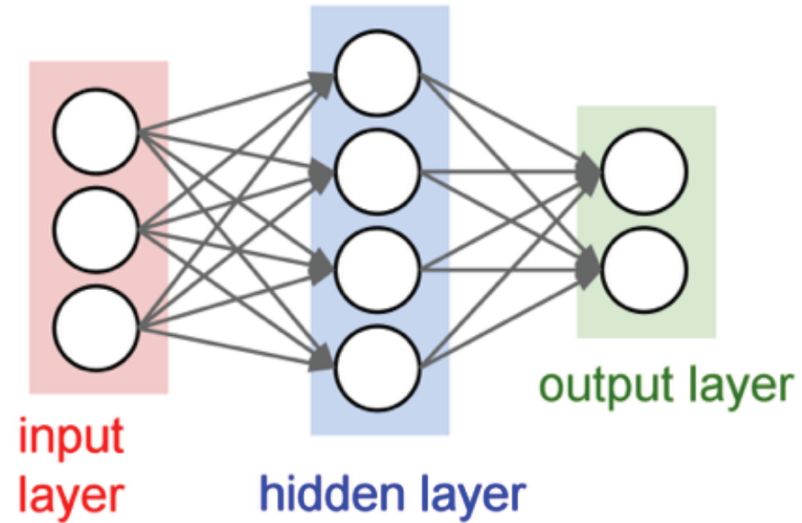
- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure
- What about keeping everything in 3D?

# ConvNets

They're just neural networks with  
3D activations and weight sharing

# 3D Activations

before:



**(1D vectors)**

# 3D Activations

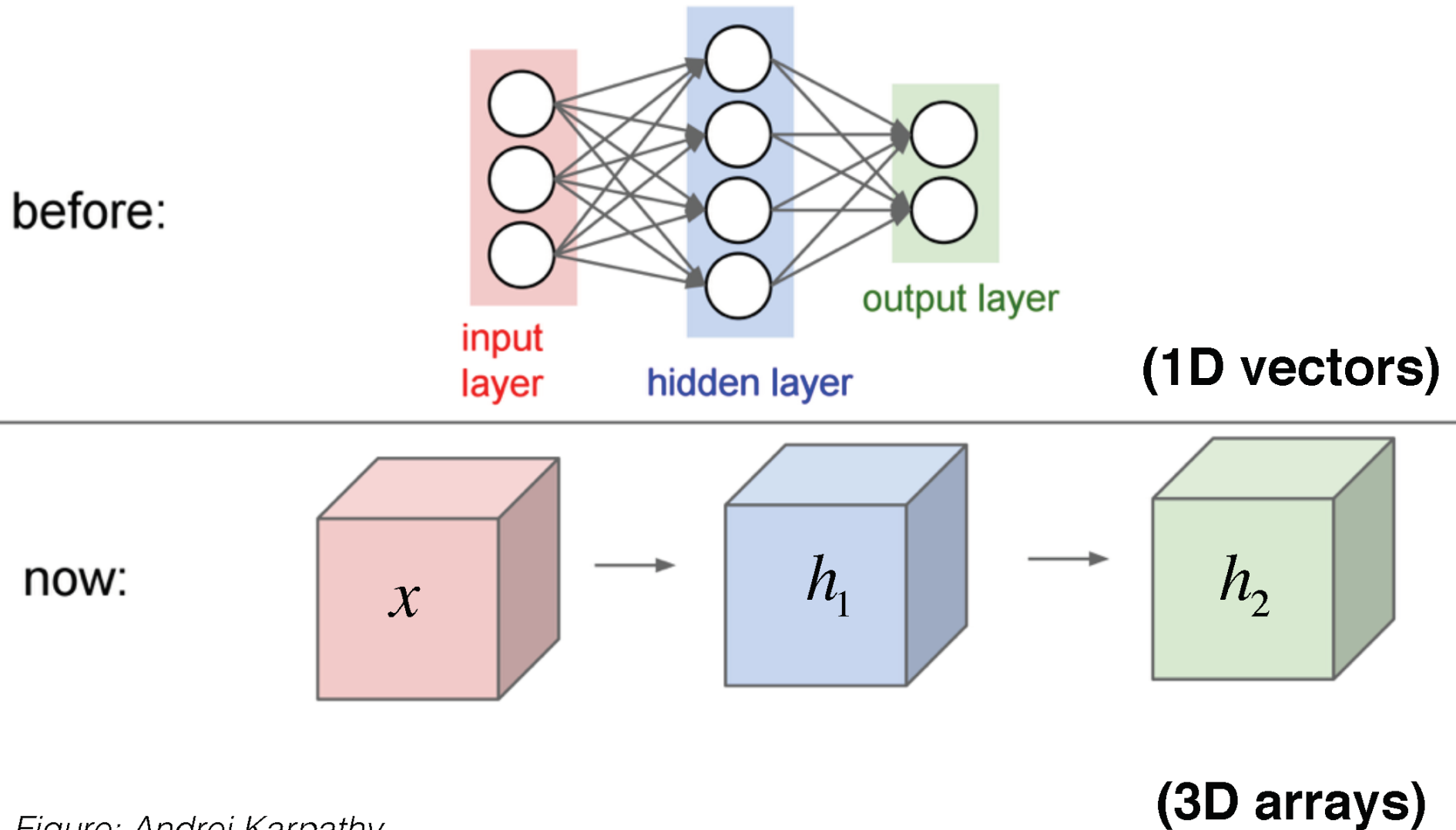
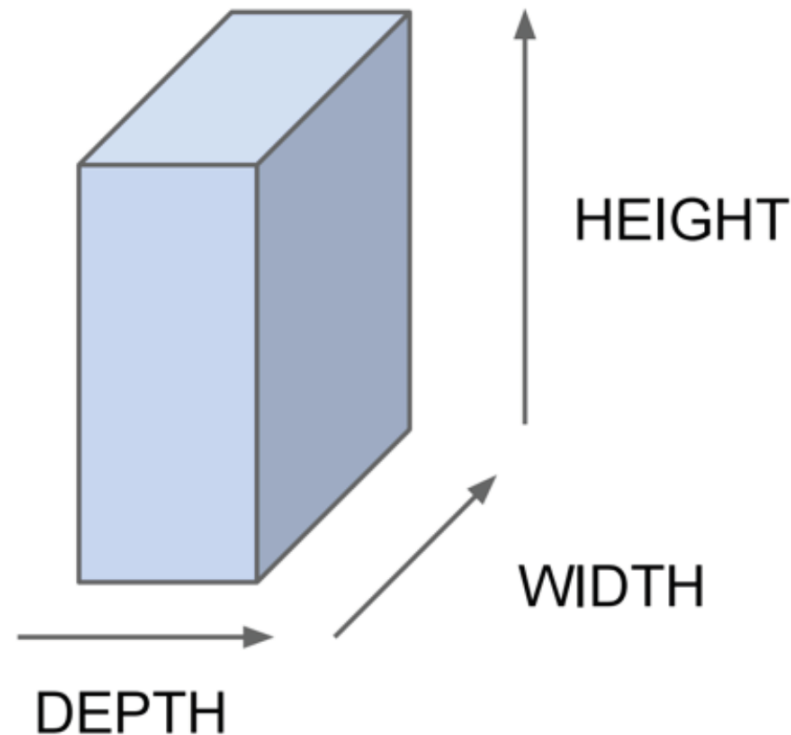


Figure: Andrej Karpathy

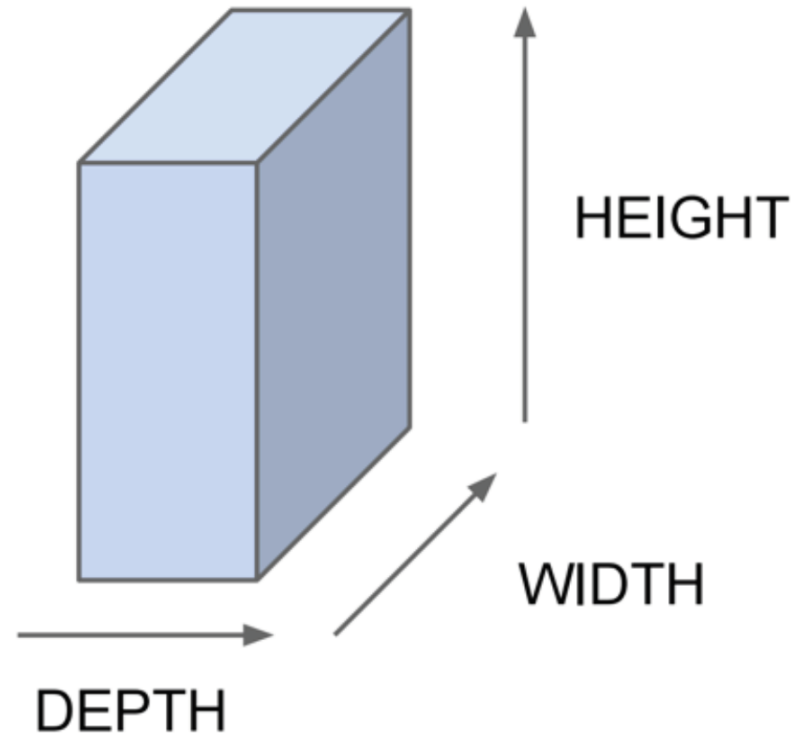
# 3D Activations

All Neural Net  
activations  
arranged in **3  
dimensions:**



# 3D Activations

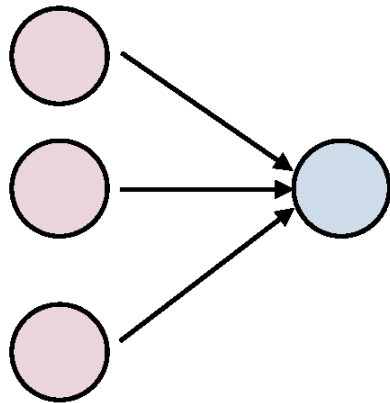
All Neural Net  
activations  
arranged in **3  
dimensions**:



For example, a CIFAR-10 image is a  $3 \times 32 \times 32$  volume  
(3 depth — RGB channels, 32 height, 32 width)

# 3D Activations

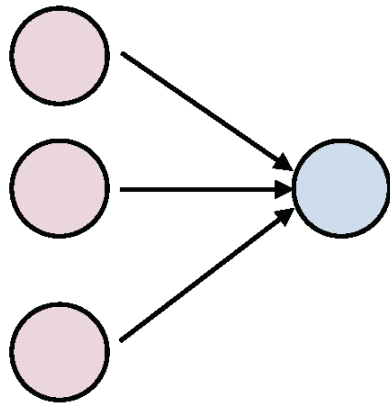
**1D Activations:**



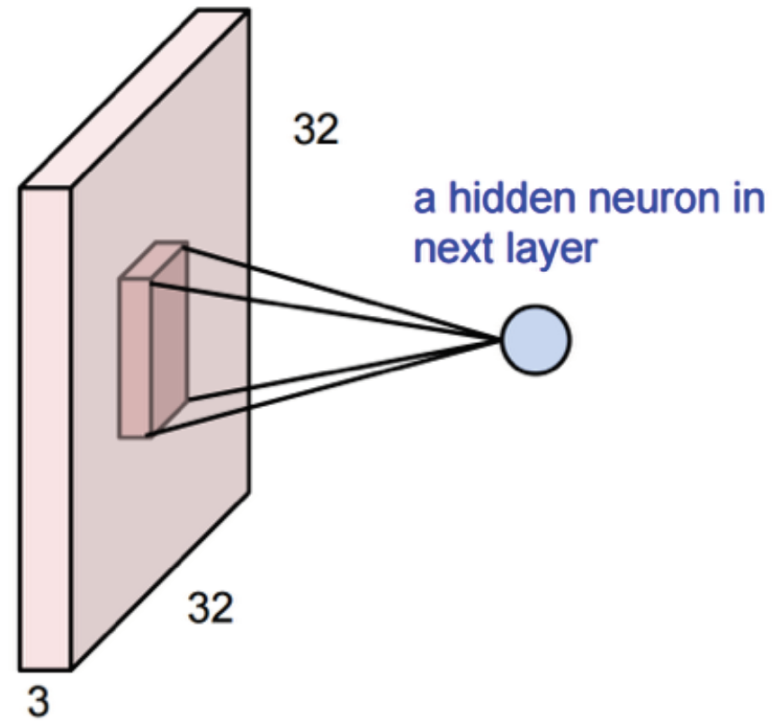
*Figure: Andrej Karpathy*

# 3D Activations

**1D Activations:**



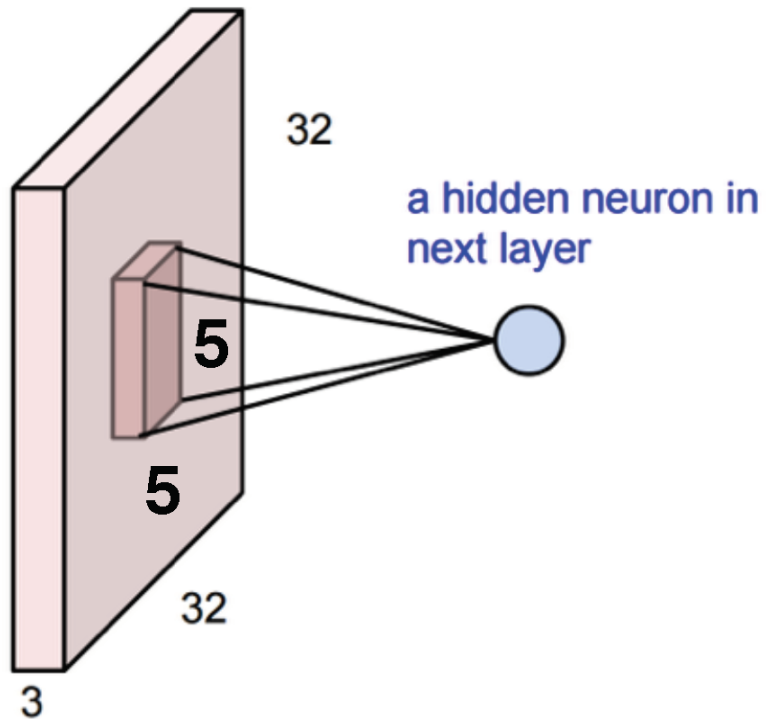
**3D Activations:**



*Figure: Andrej Karpathy*



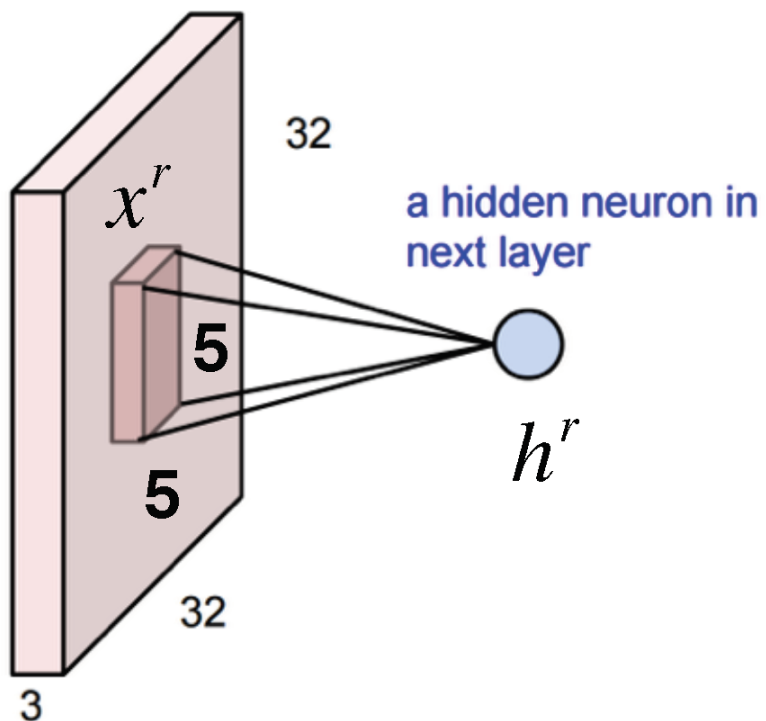
# 3D Activations



- The input is  $3 \times 32 \times 32$
- This neuron depends on a  $3 \times 5 \times 5$  chunk of the input
- The neuron also has a  $3 \times 5 \times 5$  set of weights and a bias (scalar)

Figure: Andrej Karpathy

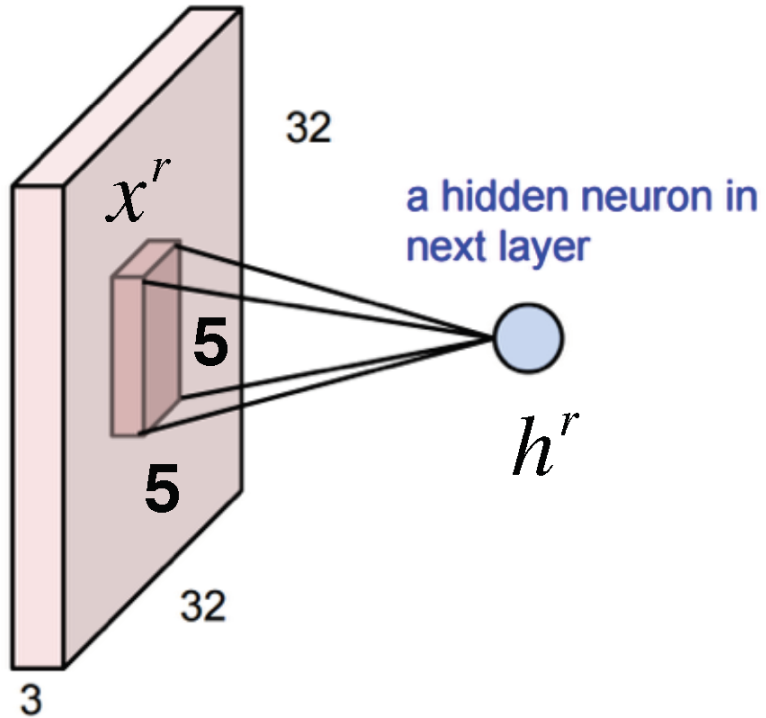
# 3D Activations



Example: consider the region of the input " $x^r$ "

With output neuron  $h^r$

# 3D Activations



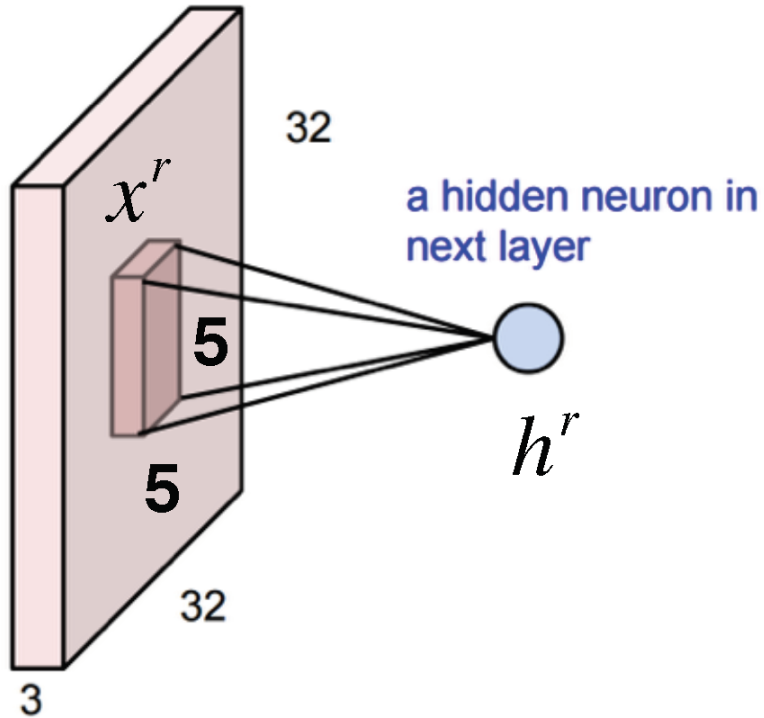
Example: consider the region of the input “ $x^r$ ”

With output neuron  $h^r$

Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

# 3D Activations



Example: consider the region of the input “ $x^r$ ”

With output neuron  $h^r$

Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

Sum over 3 axes

# 3D Activations

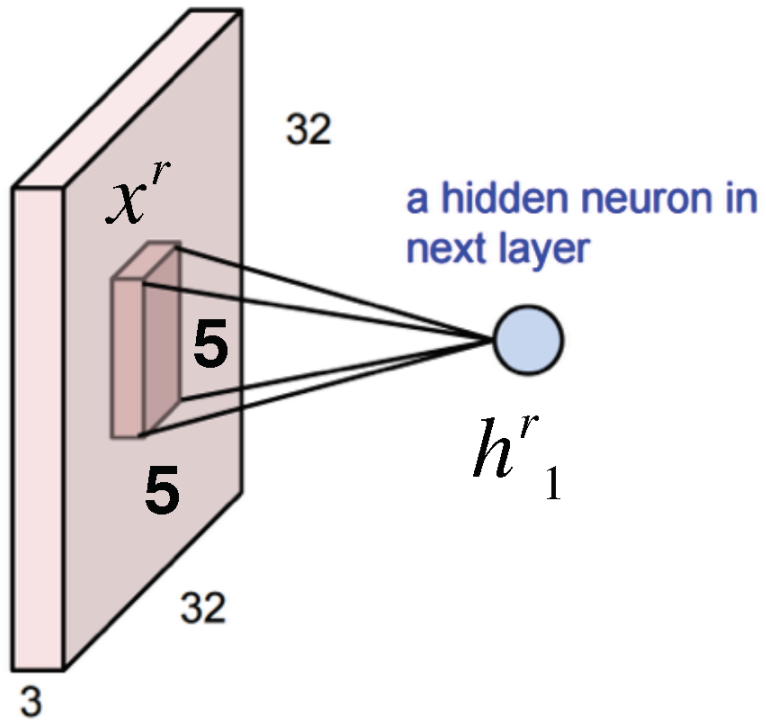


Figure: Andrej Karpathy

# 3D Activations

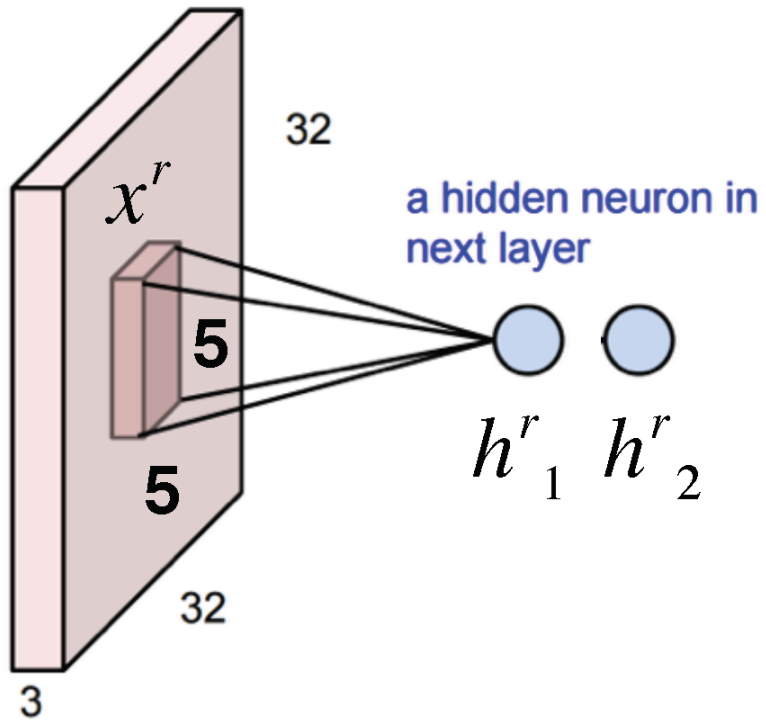
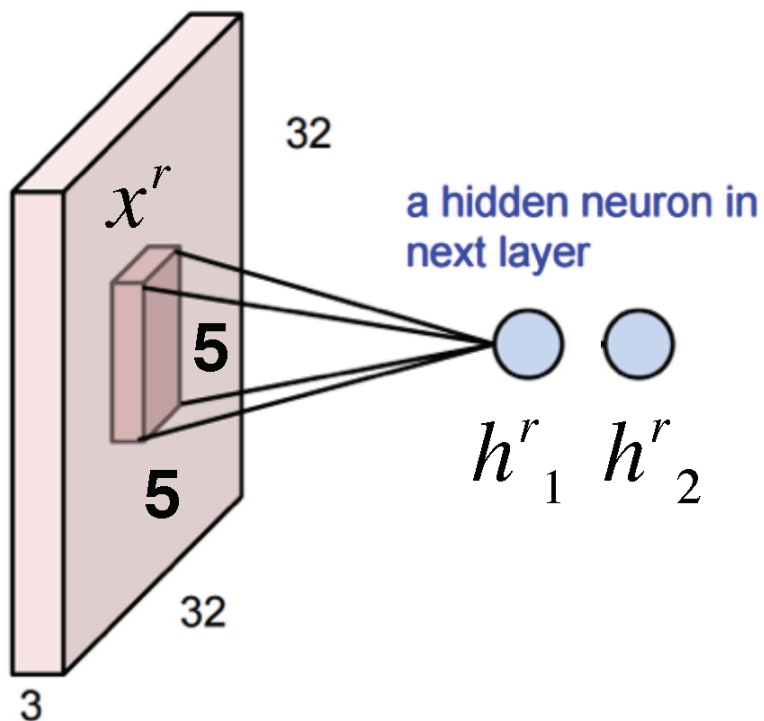


Figure: Andrej Karpathy

# 3D Activations

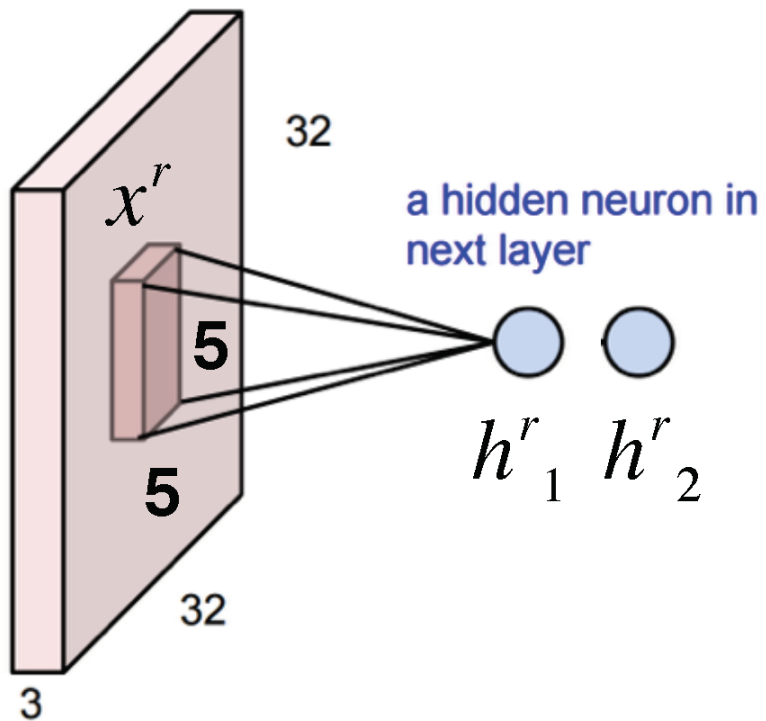


With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_1$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$

# 3D Activations



With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_{1}$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_{2}$$



# 3D Activations

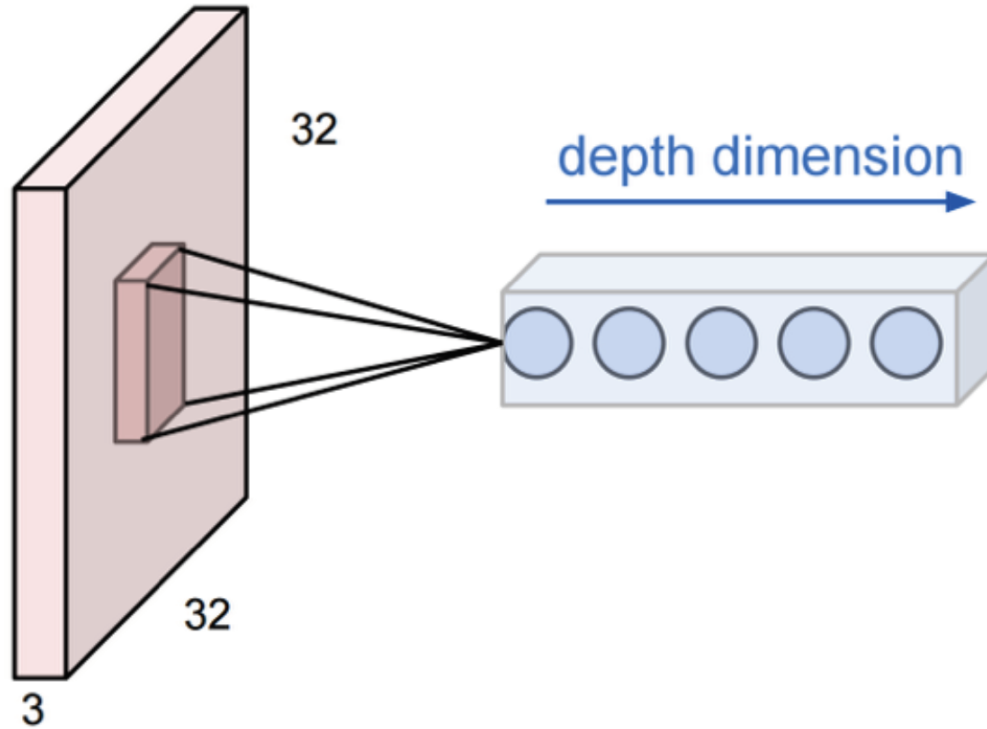
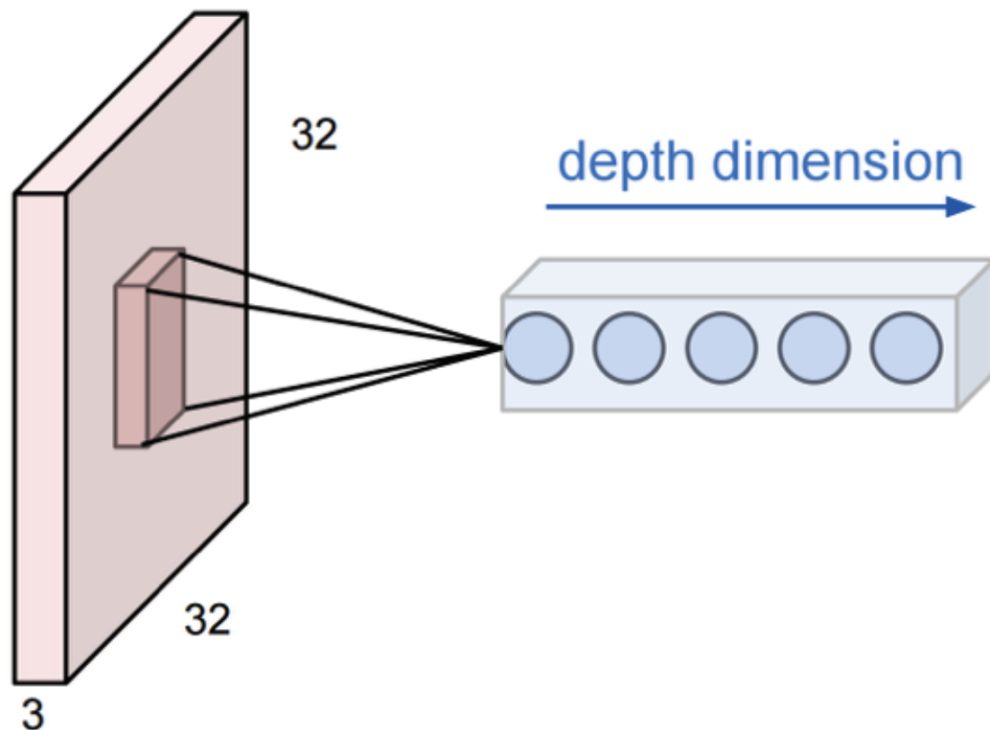


Figure: Andrej Karpathy

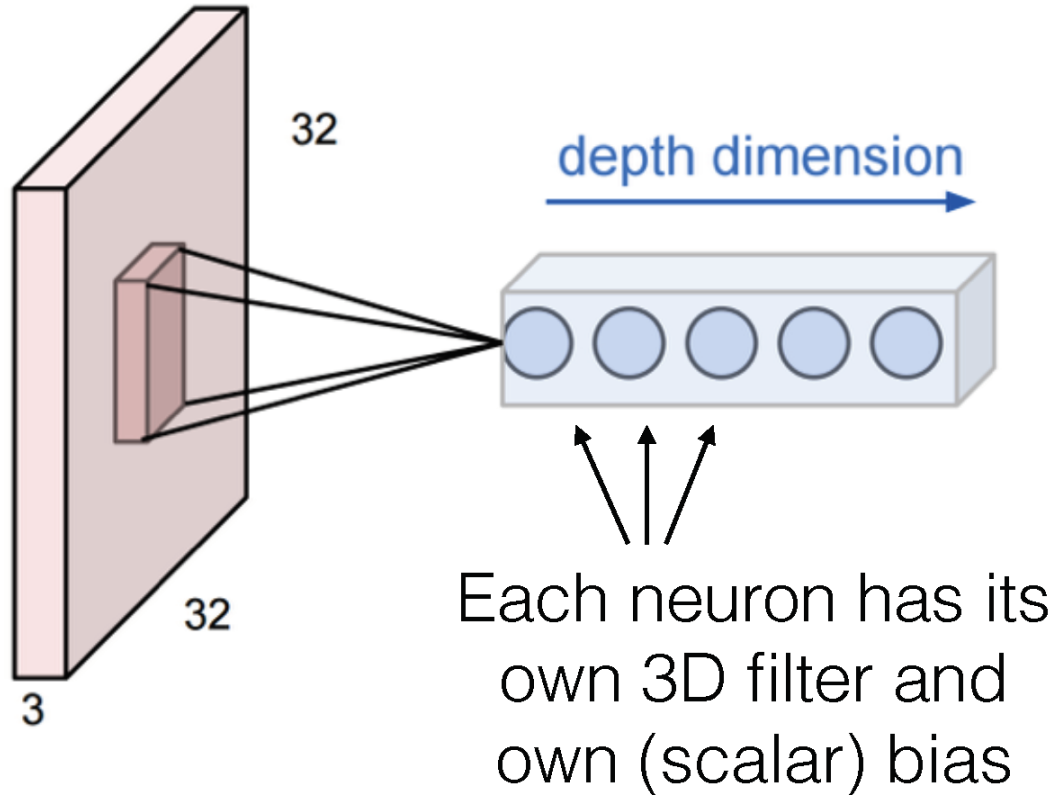
# 3D Activations



We can keep adding more outputs

These form a column in the output volume:  
[depth x 1 x 1]

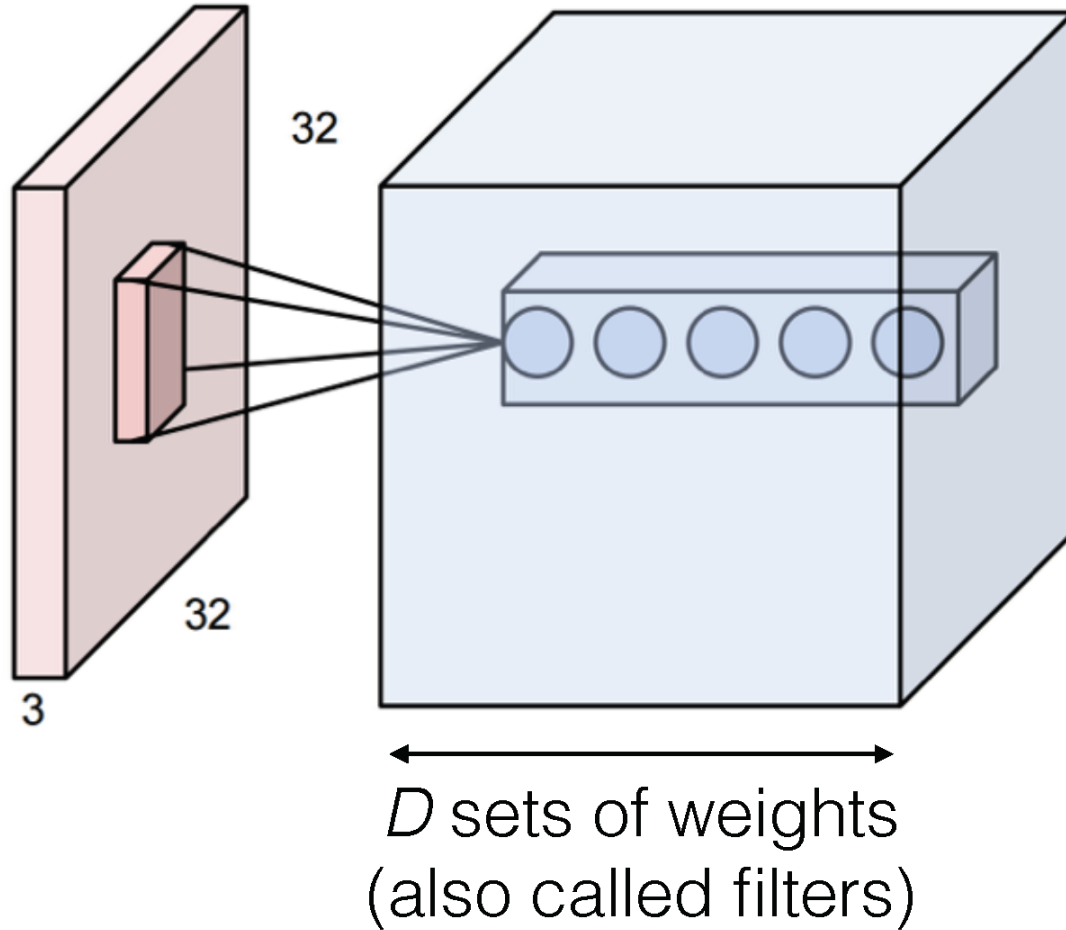
# 3D Activations



We can keep adding more outputs

These form a column in the output volume:  
[depth x 1 x 1]

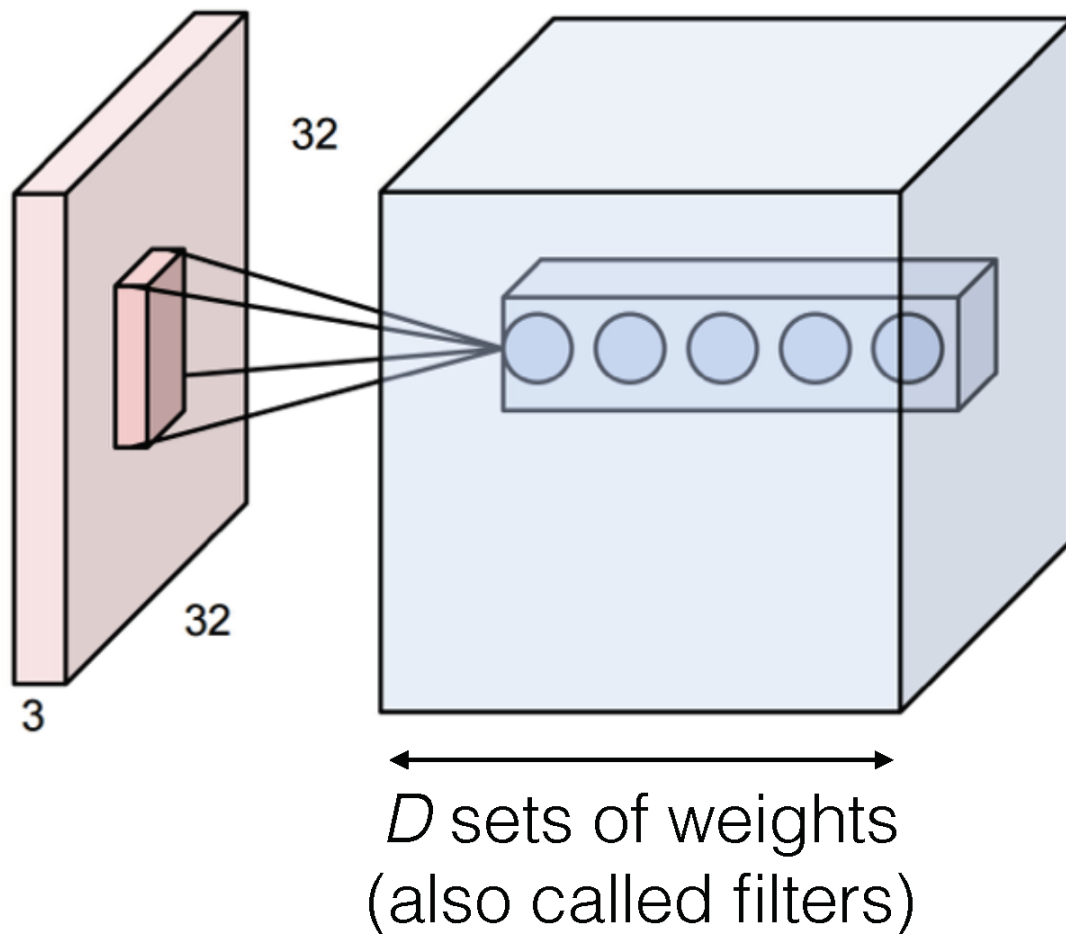
# 3D Activations



Now repeat this  
across the input

Figure: Andrej Karpathy

# 3D Activations



Now repeat this  
across the input

## **Weight sharing:**

Each filter shares  
the same weights  
(but each depth  
index has its own  
set of weights)

Figure: Andrej Karpathy

# 3D Activations

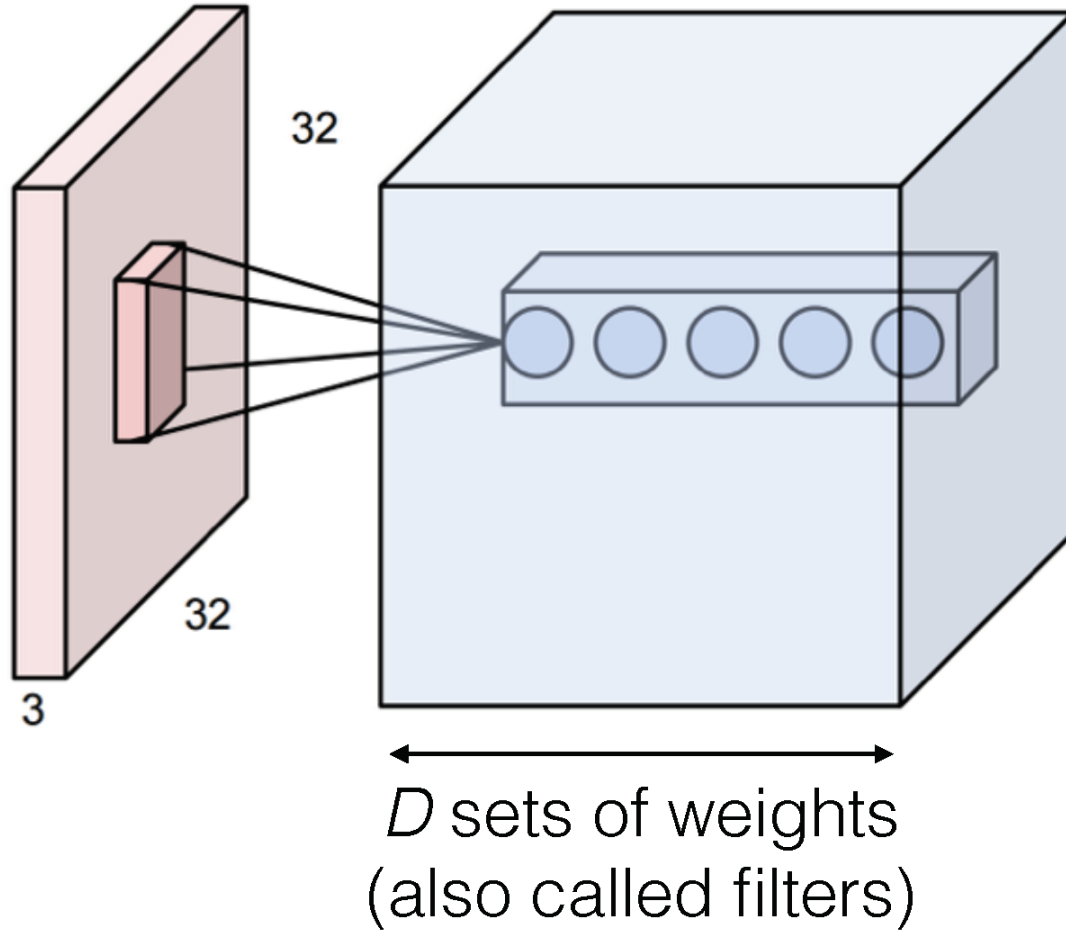
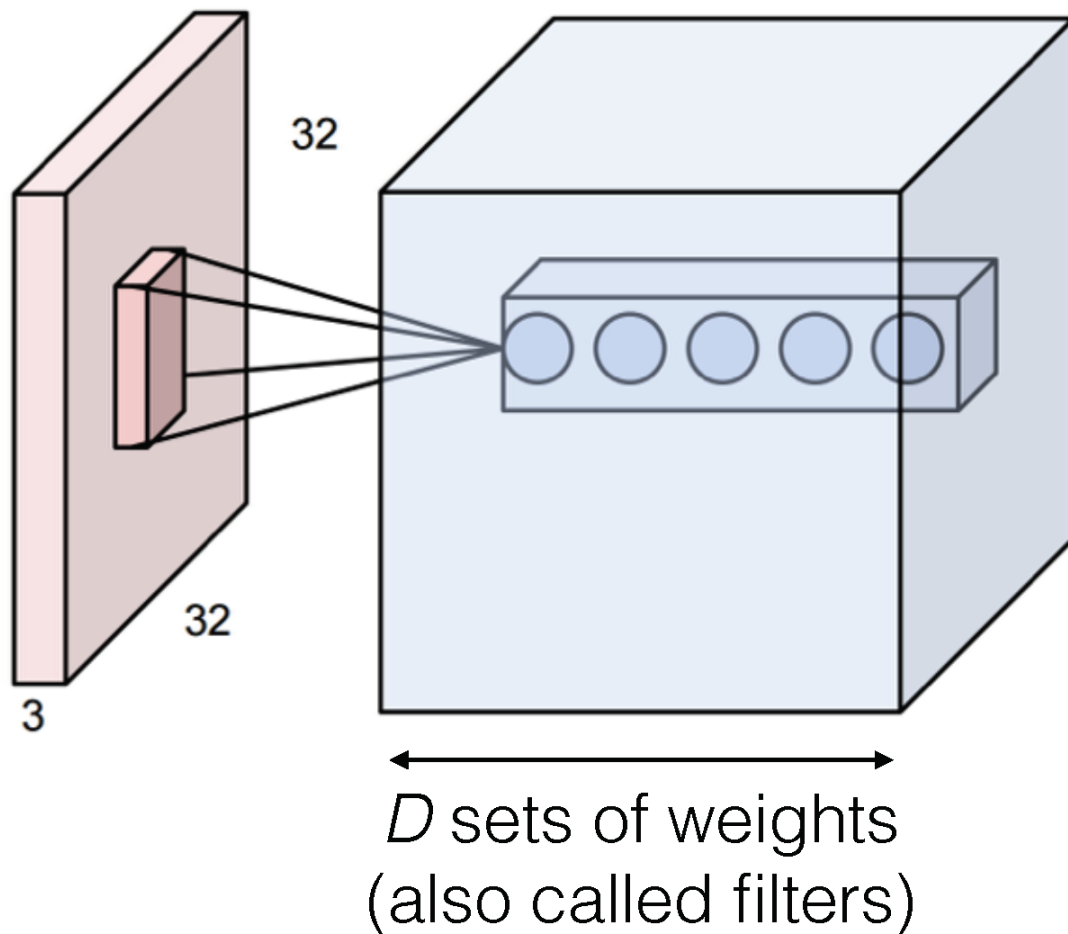


Figure: Andrej Karpathy

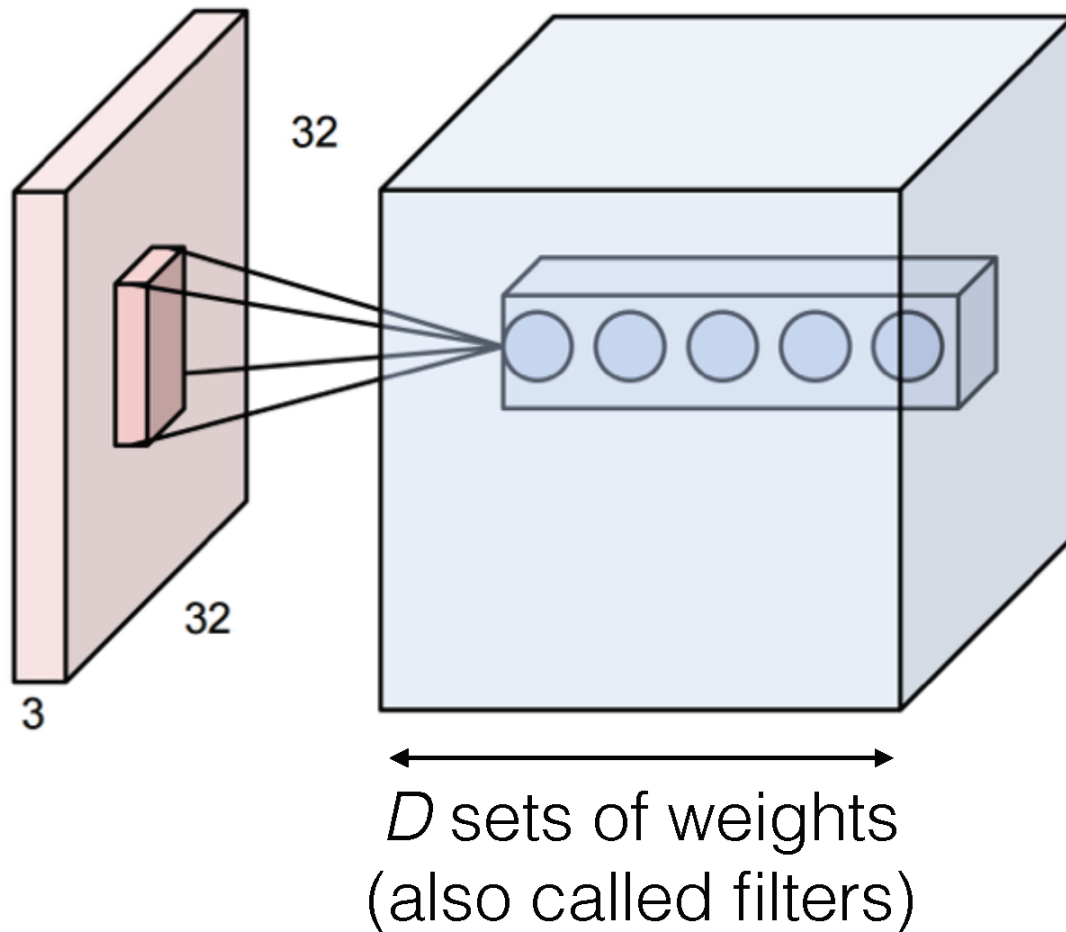
# 3D Activations



With weight sharing,  
this is called **convolution**

Figure: Andrej Karpathy

# 3D Activations



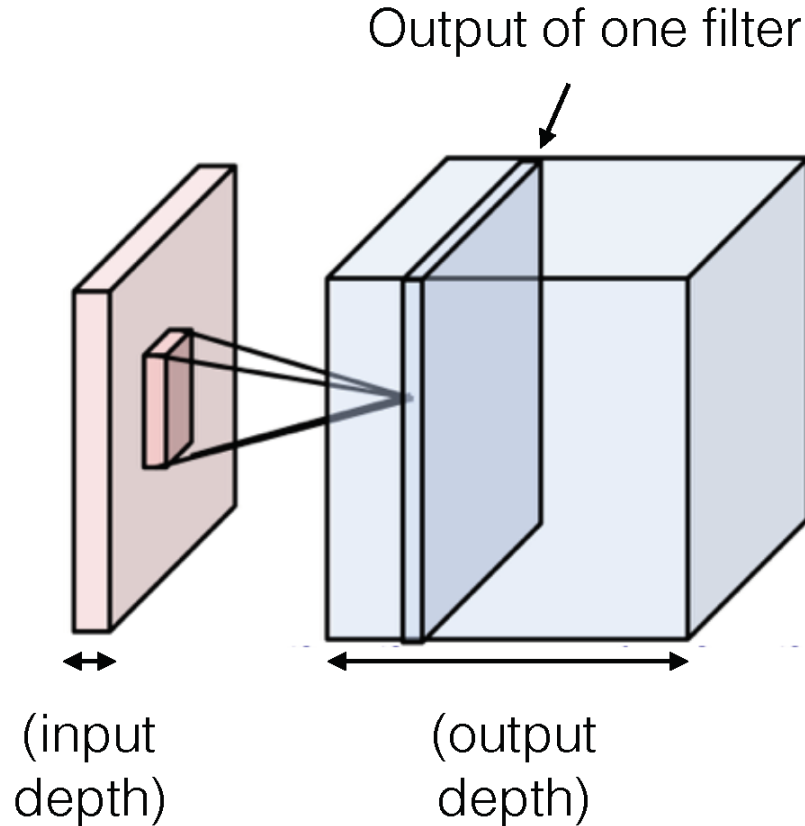
With weight sharing,  
this is called  
**convolution**

Without weight sharing,  
this is called a  
**locally connected layer**

Figure: Andrej Karpathy



# 3D Activations

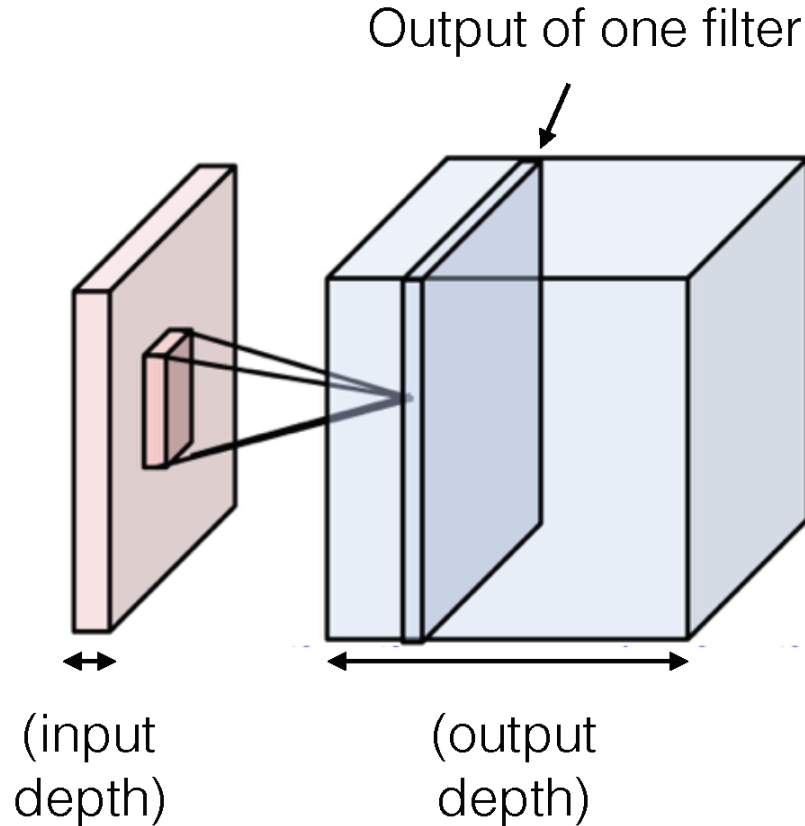


One set of weights gives one slice in the output

To get a 3D output of depth  $D$ , use  $D$  different filters

In practice, ConvNets use many filters ( $\sim 64$  to  $1024$ )

# 3D Activations



One set of weights gives one slice in the output

To get a 3D output of depth  $D$ , use  $D$  different filters

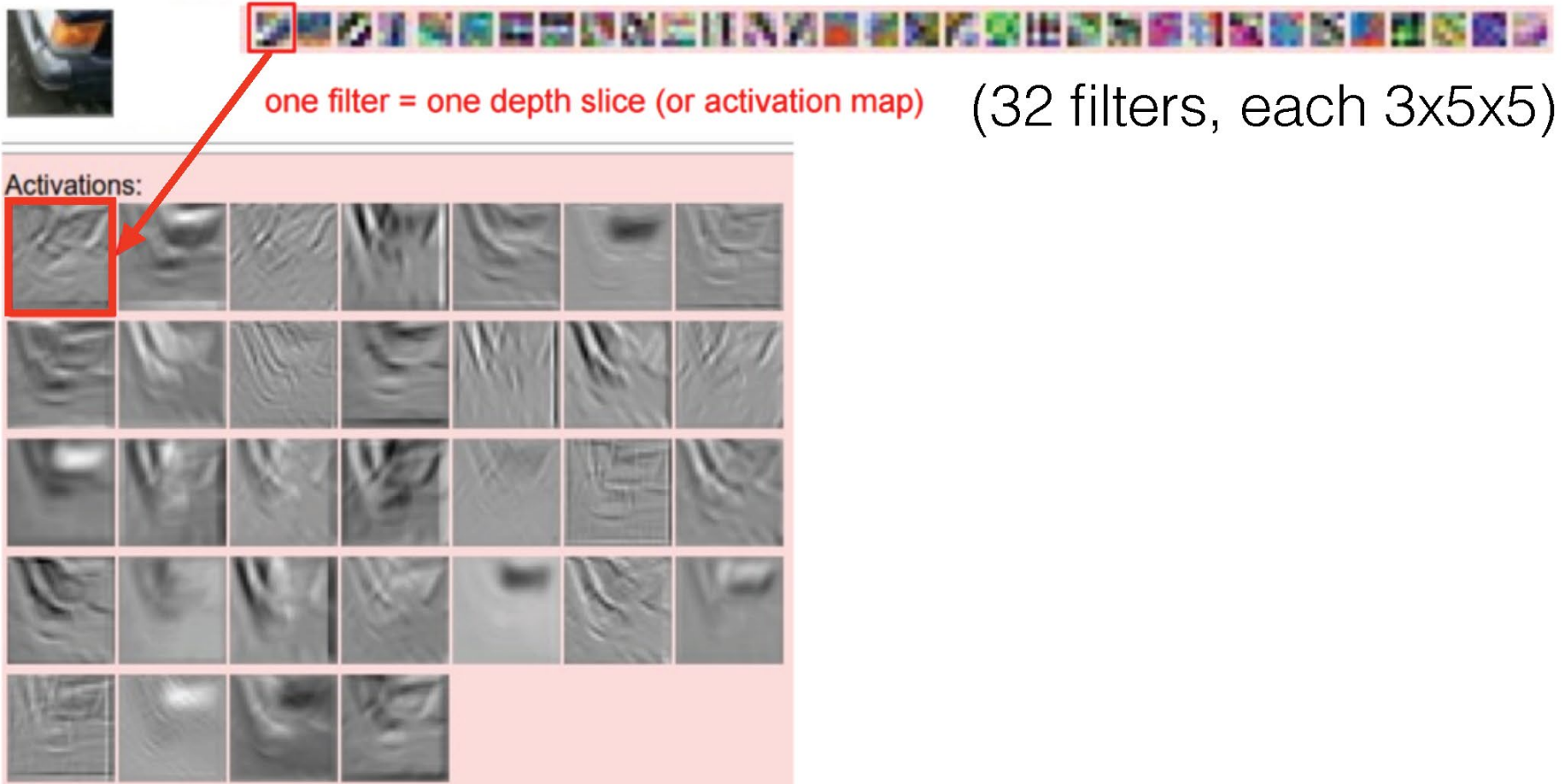
In practice, ConvNets use many filters ( $\sim 64$  to  $1024$ )

All together, the weights are **4** dimensional:  
(output depth, input depth, kernel height, kernel width)

# 3D Activations

**We can unravel the 3D cube and show each layer separately:**

(Input)



*Figure: Andrej Karpathy*

# 3D Activations

**We can unravel the 3D cube and show each layer separately:**

(Input)

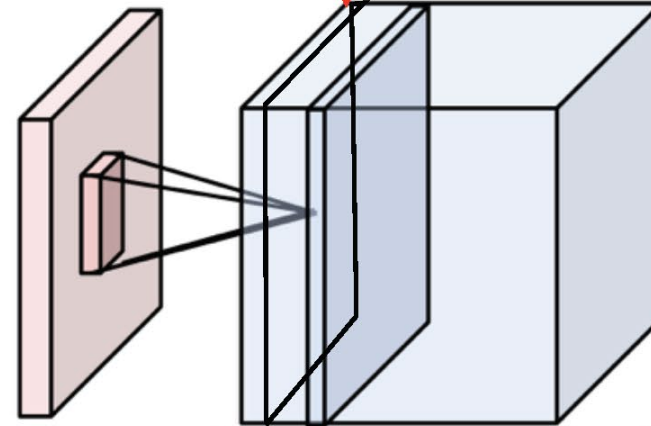
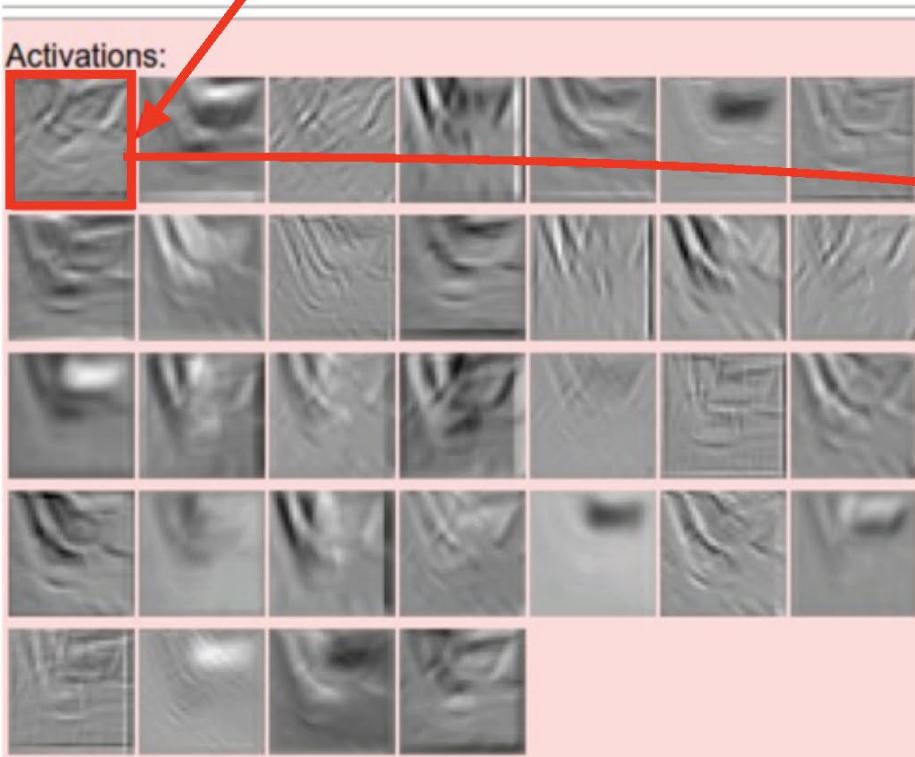


Figure: Andrej Karpathy

# 3D Activations

**We can unravel the 3D cube and show each layer separately:**

(Input)

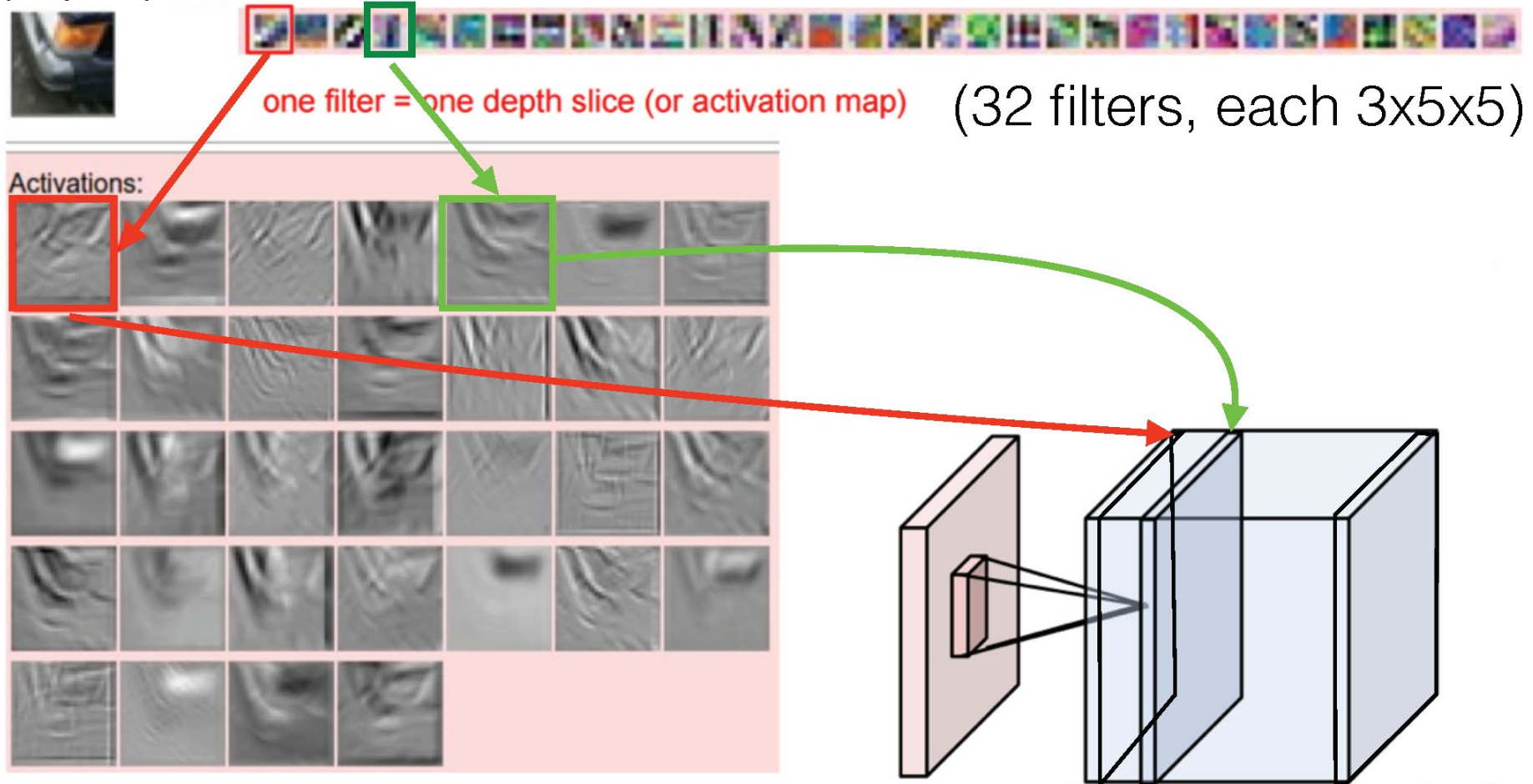


Figure: Andrej Karpathy



# 3D Activations

**We can unravel the 3D cube and show each layer separately:**

(Input)

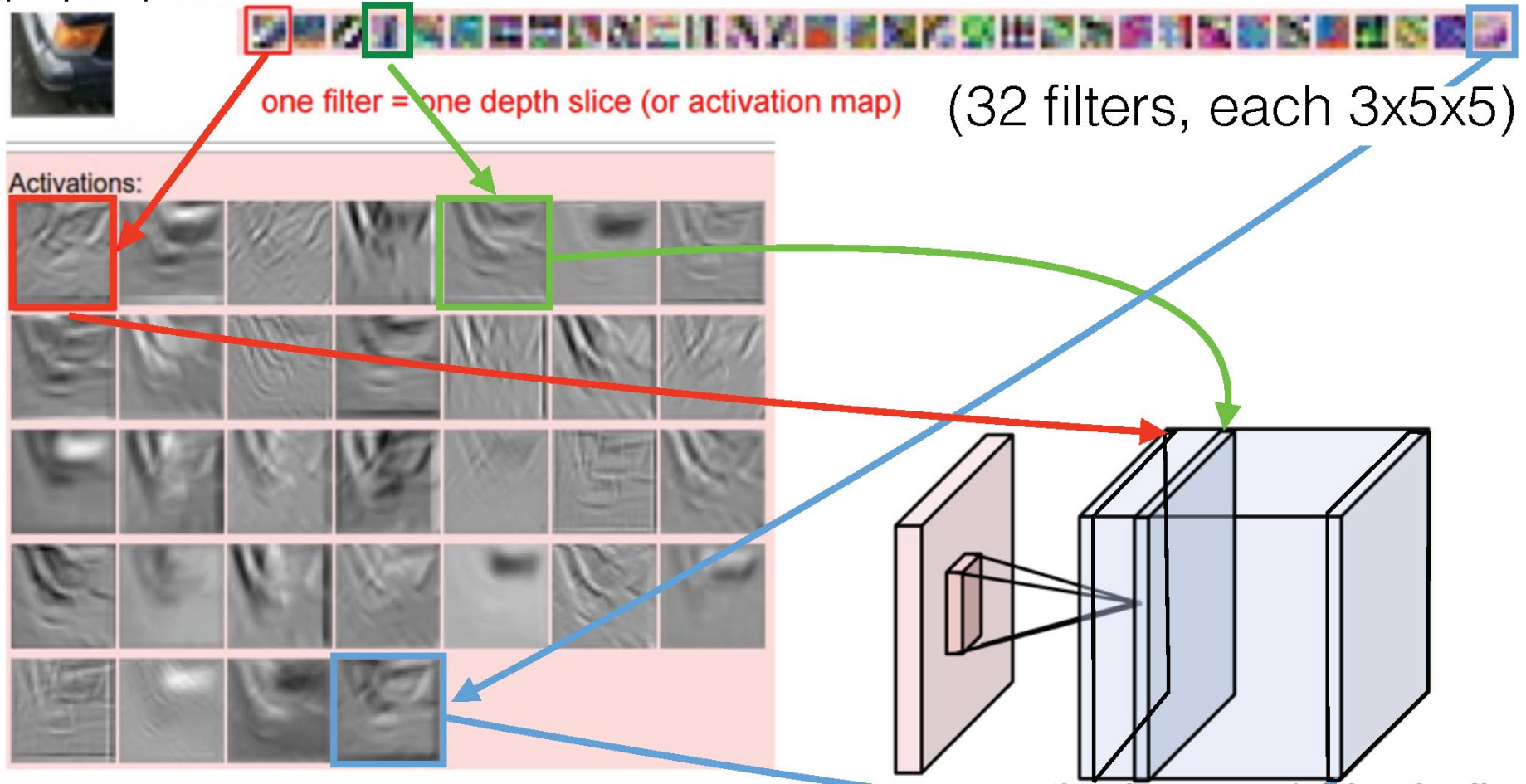
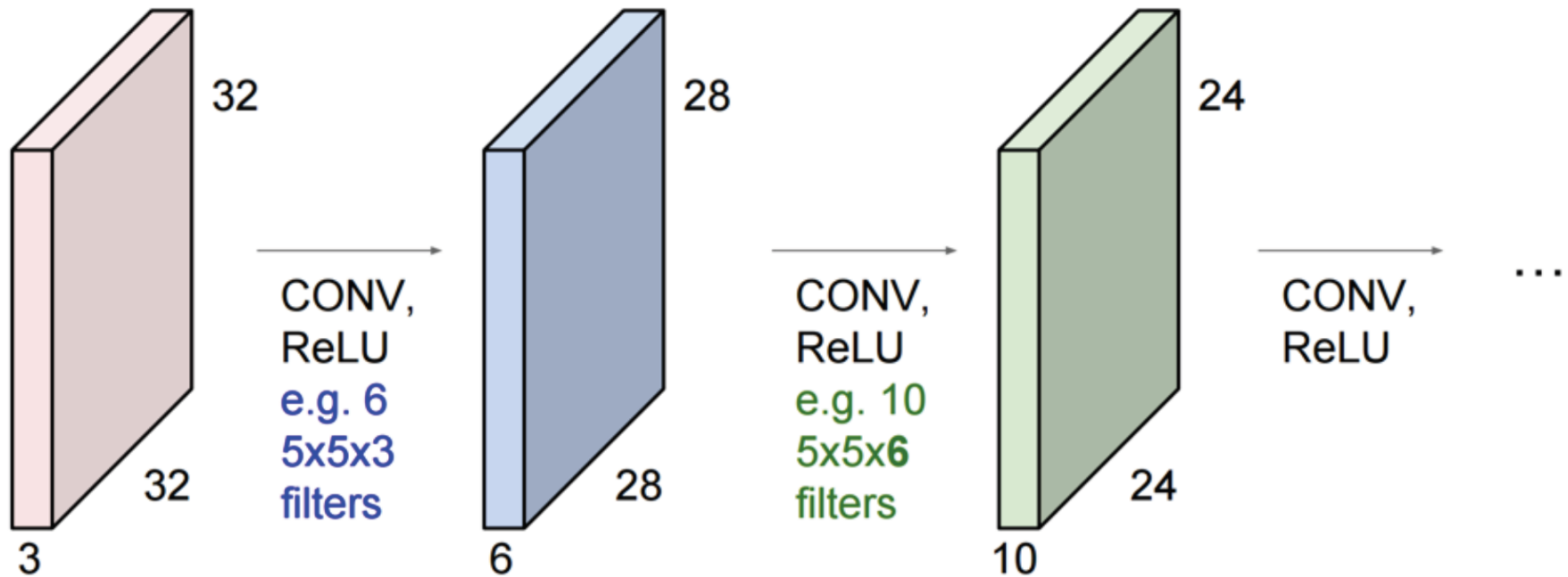


Figure: Andrej Karpathy

# (Recap)

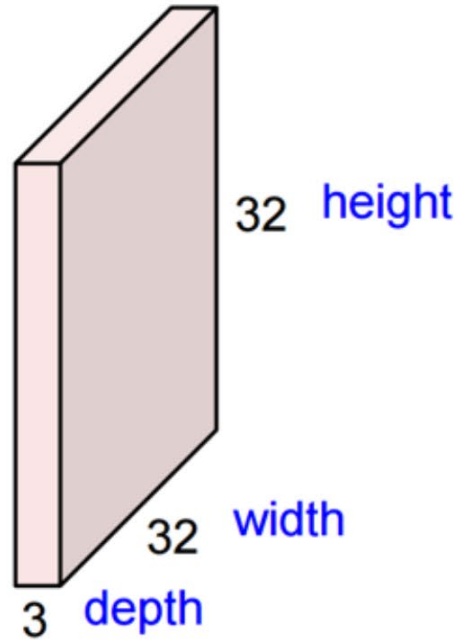
A **ConvNet** is a sequence of convolutional layers, interspersed with activation functions (and possibly other layer types)



(Recap)

## Convolution Layer

32x32x3 image

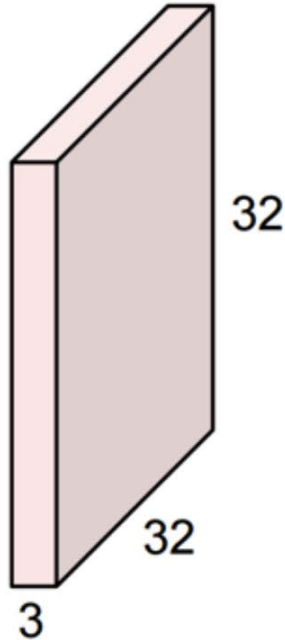




# (Recap)

## Convolution Layer

32x32x3 image



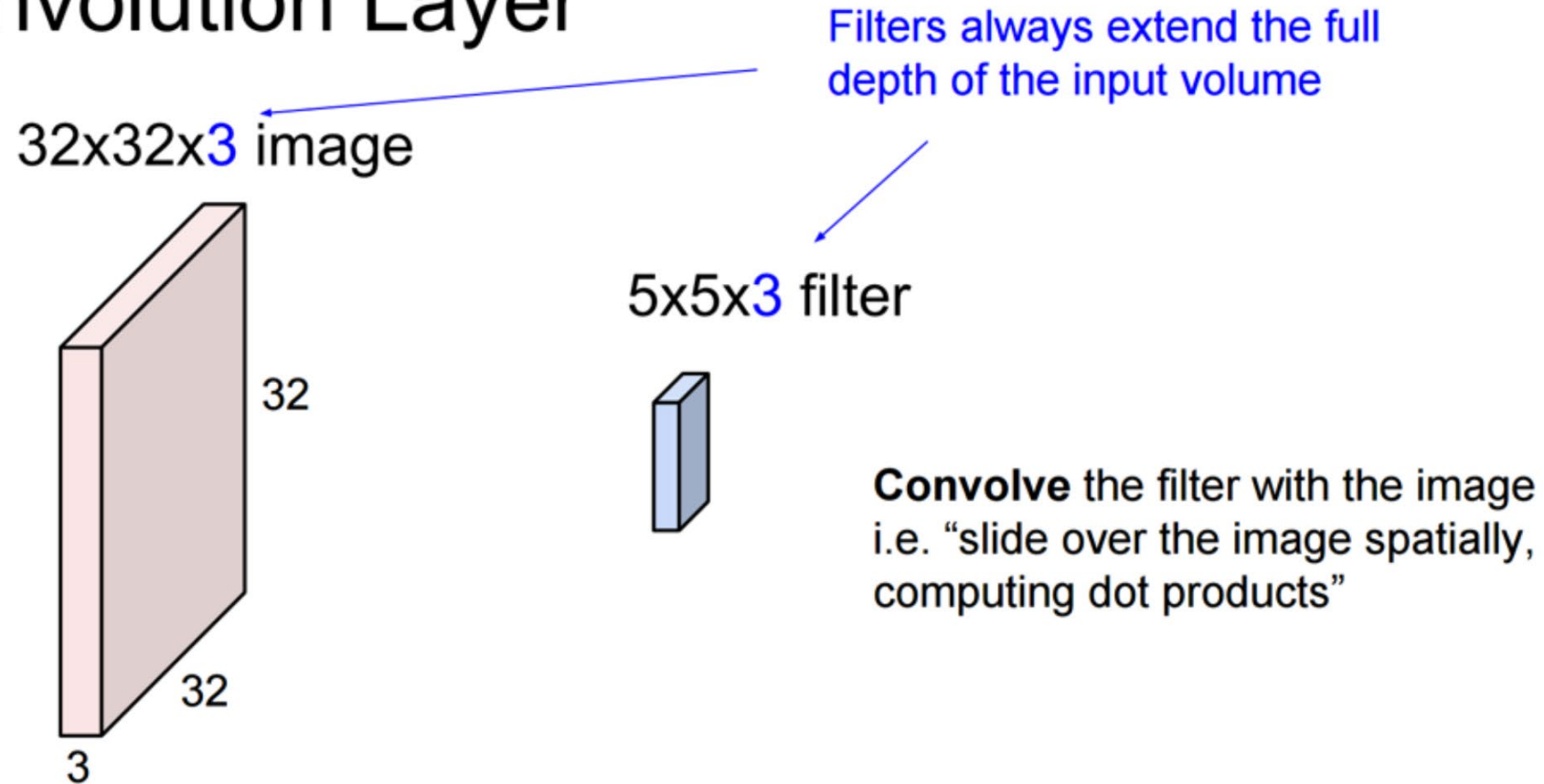
5x5x3 filter



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

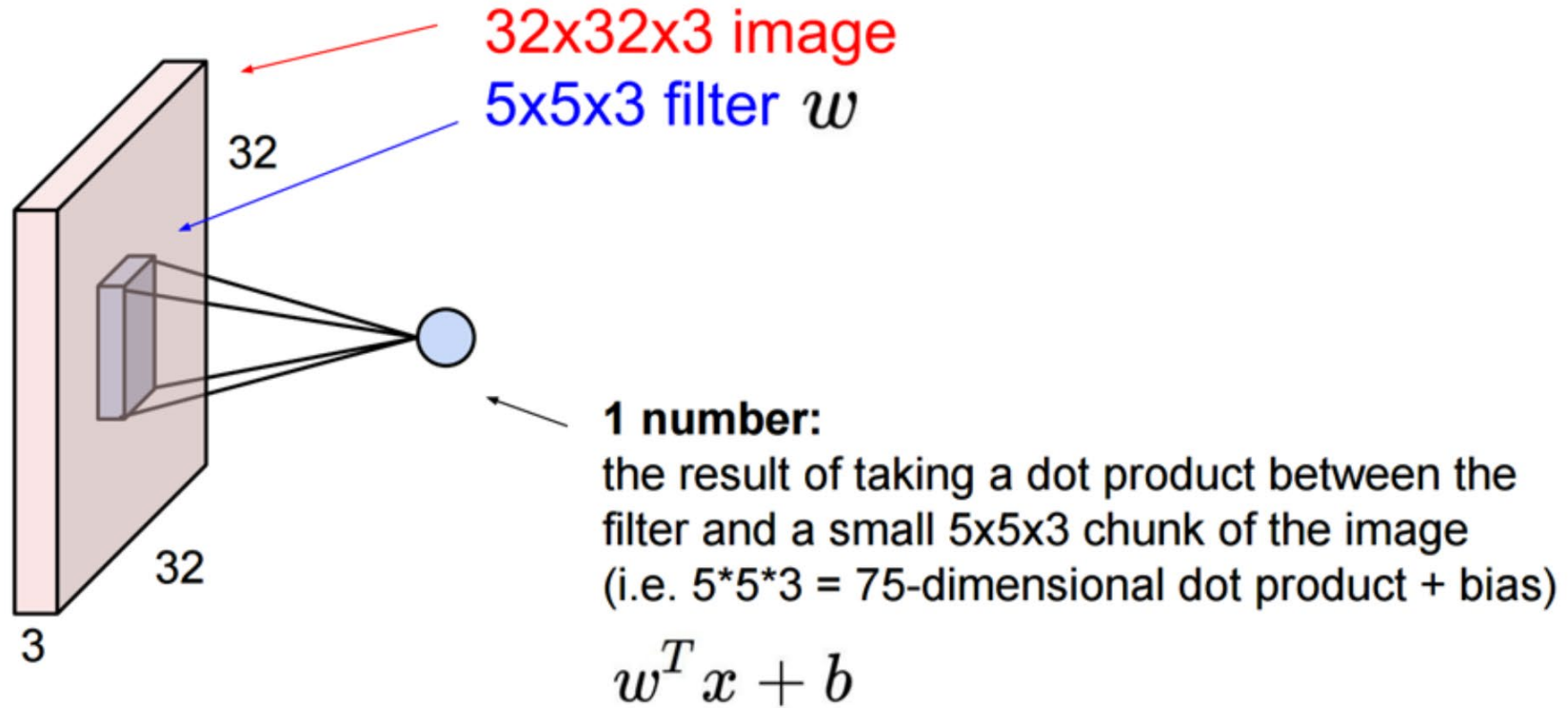
# (Recap)

## Convolution Layer



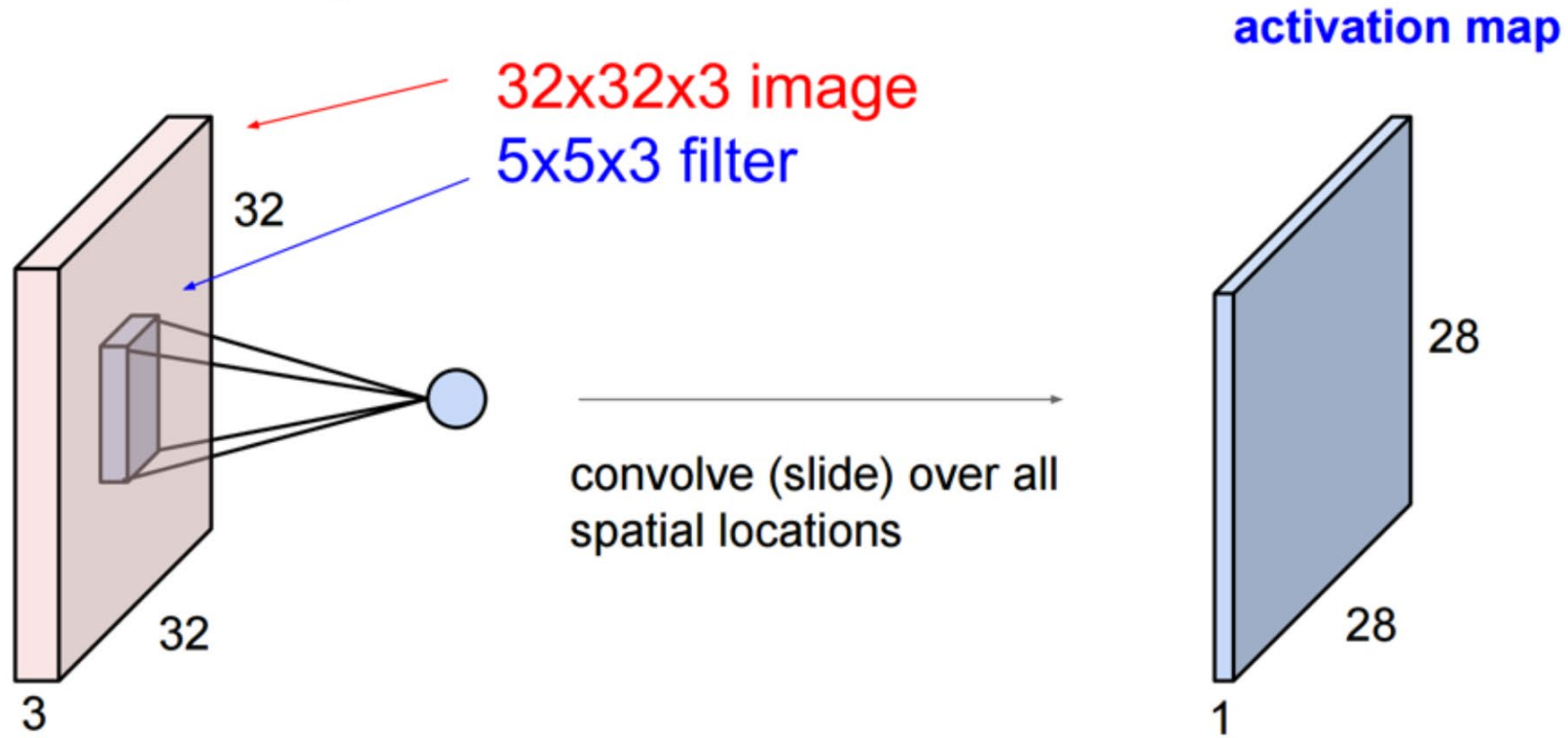
# (Recap)

## Convolution Layer



# (Recap)

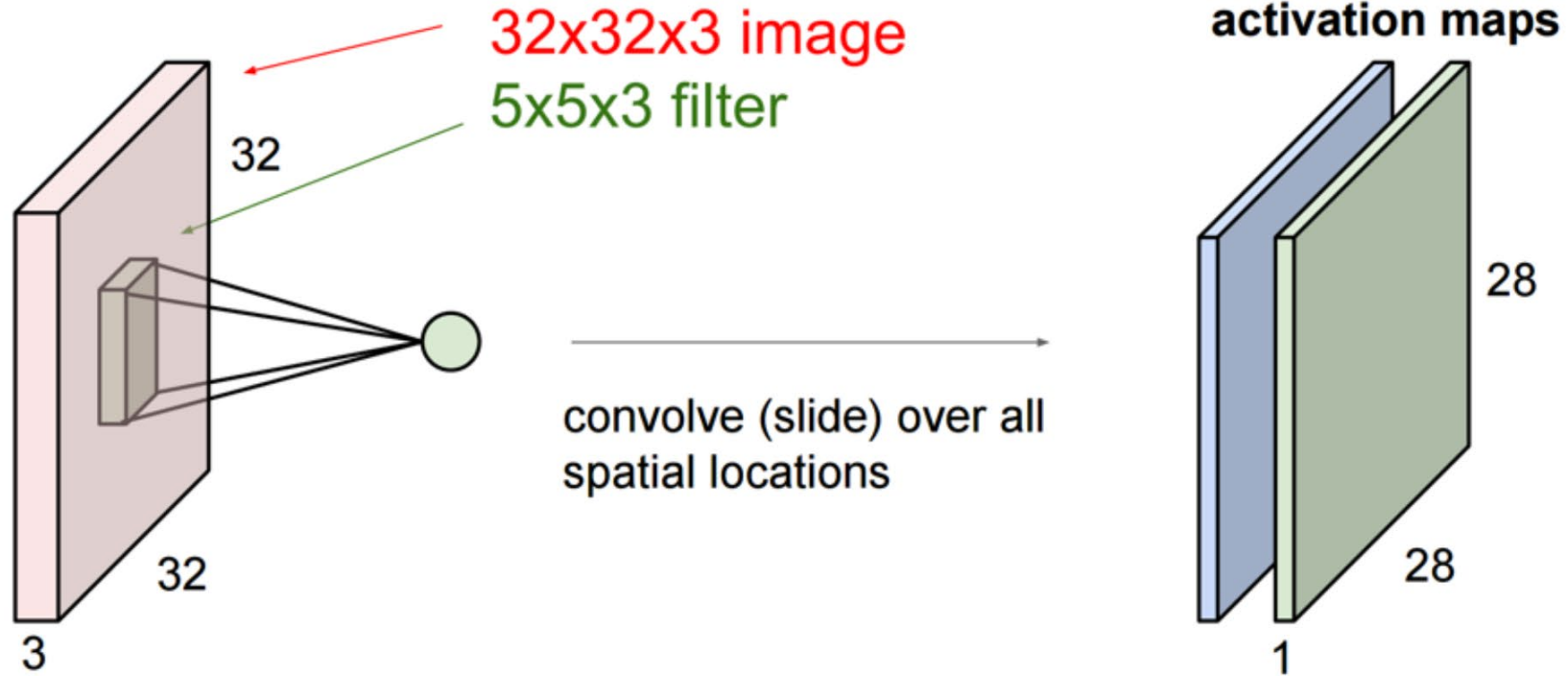
## Convolution Layer



# (Recap)

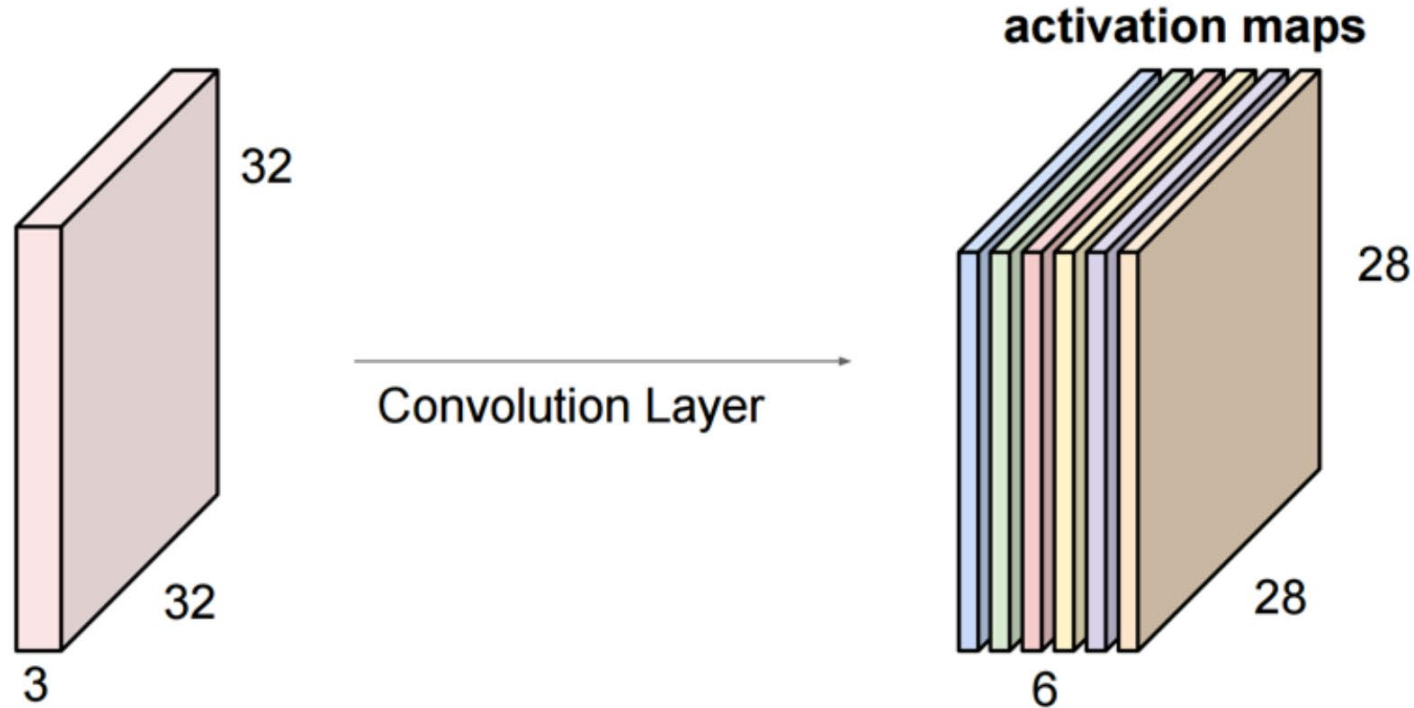
## Convolution Layer

consider a second, **green** filter



# (Recap)

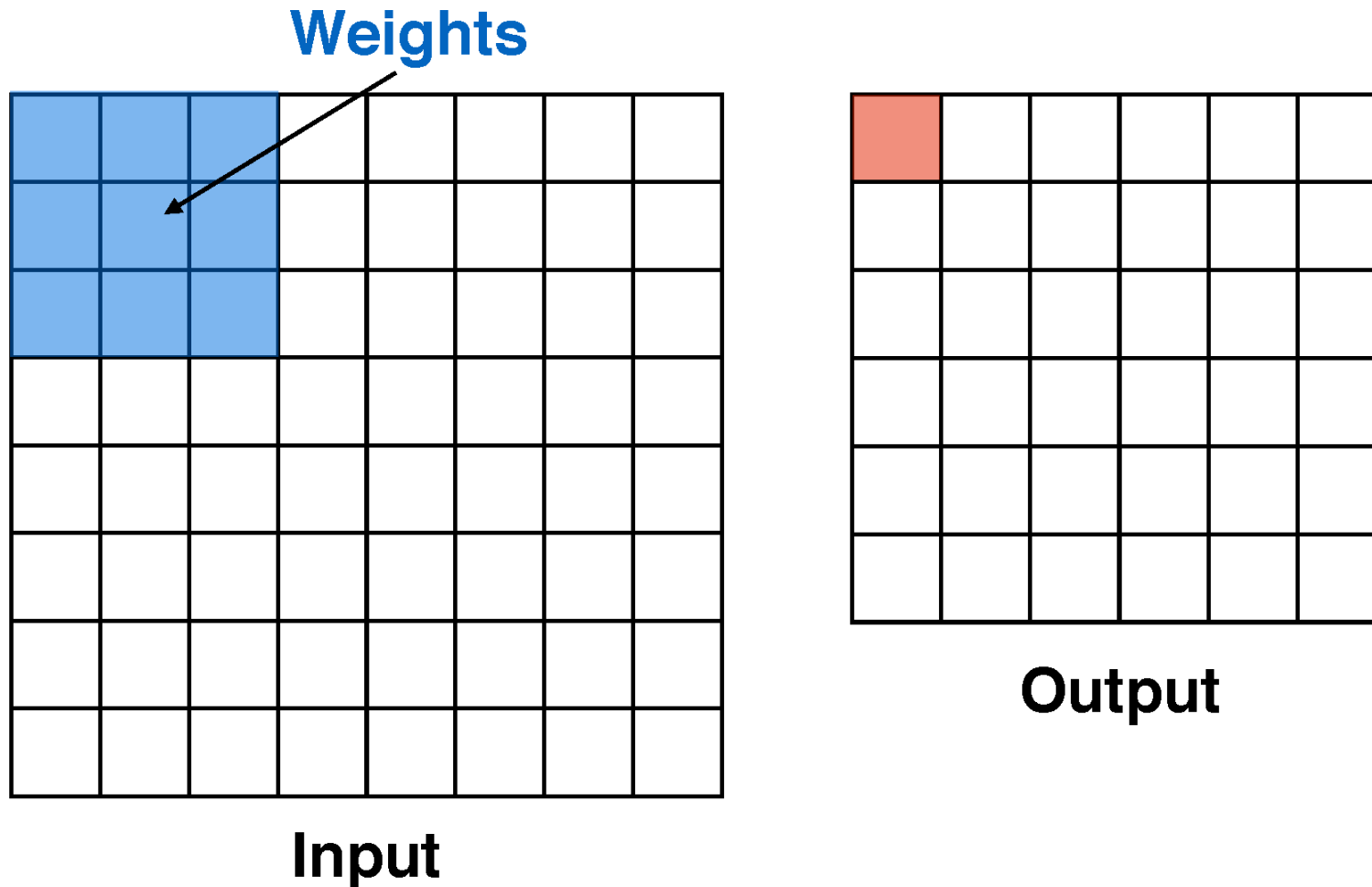
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

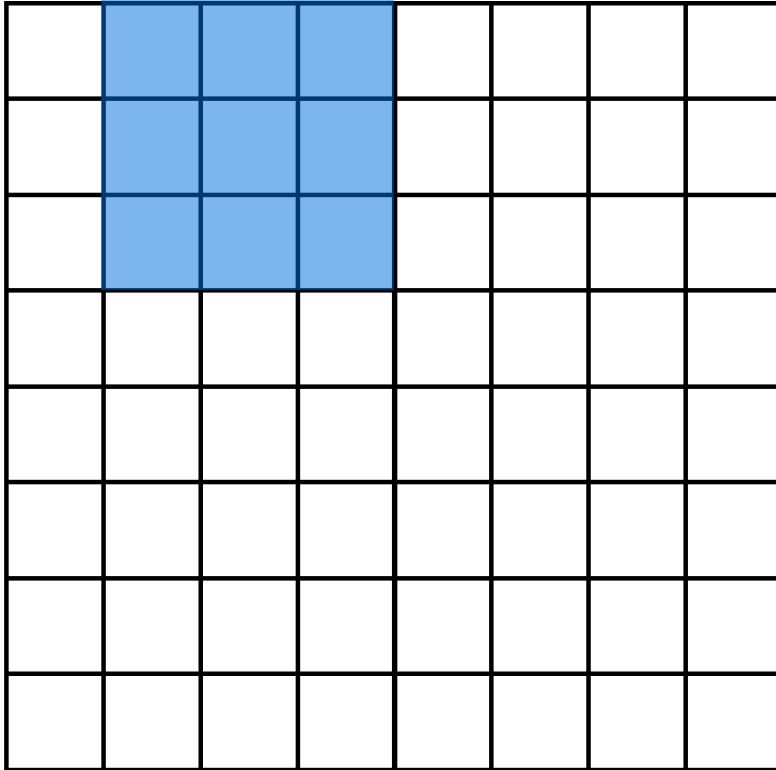
# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

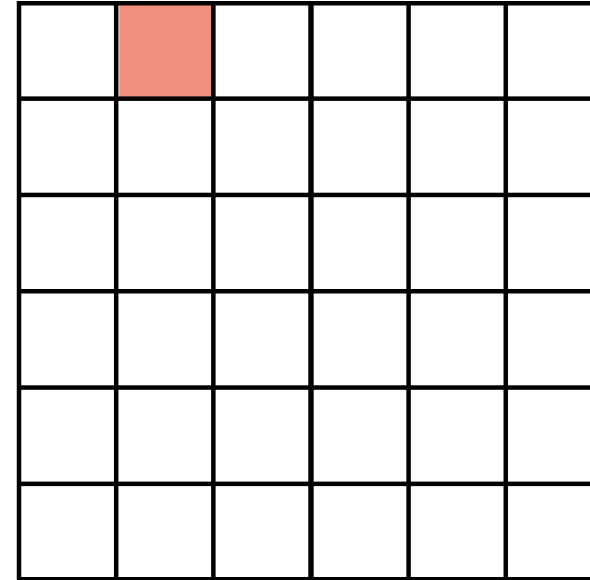


# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



**Input**

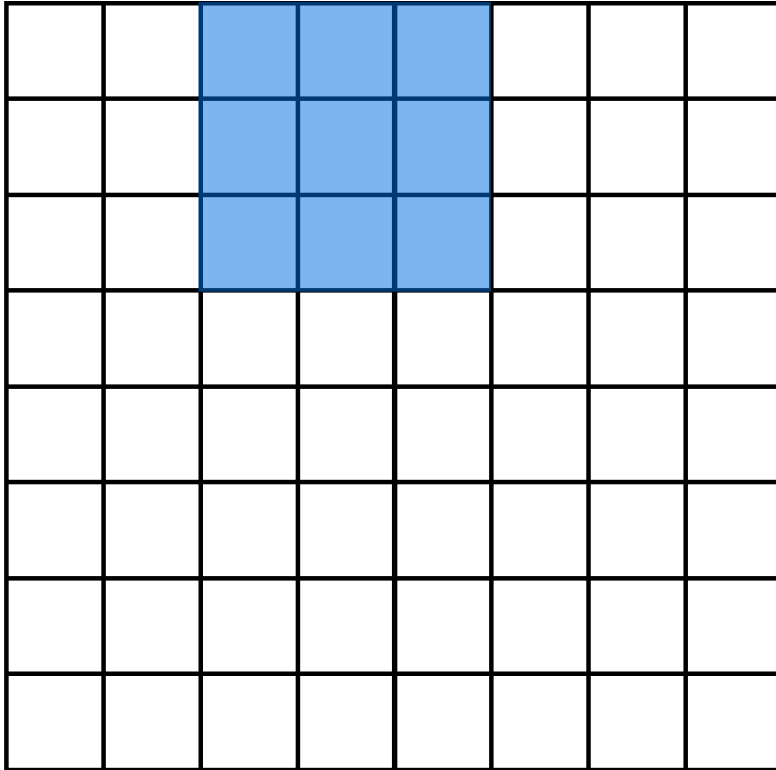


**Output**

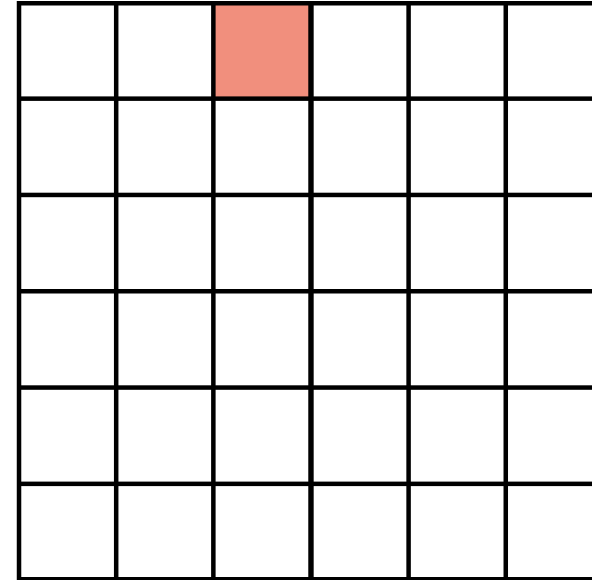


# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



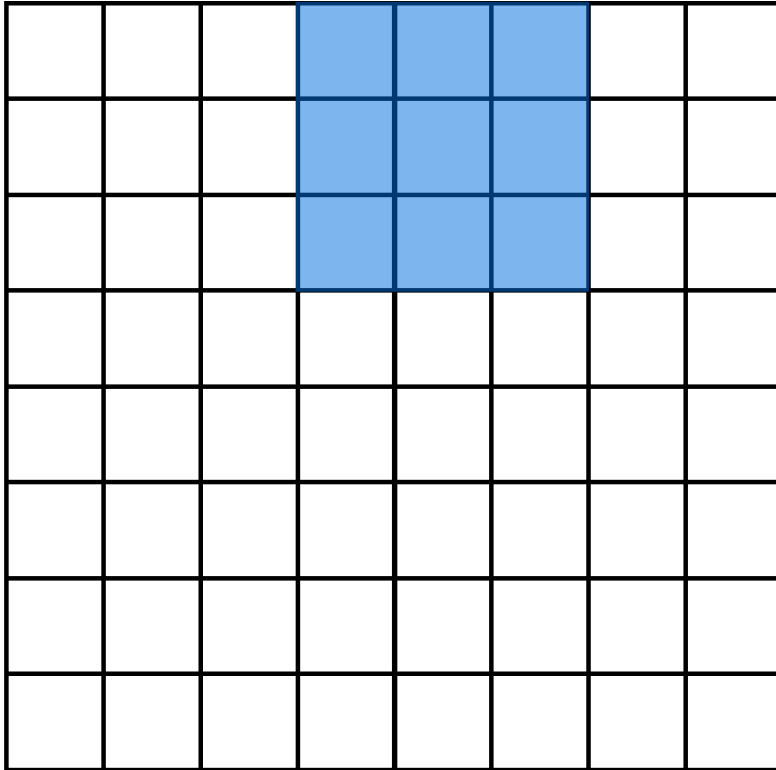
**Input**



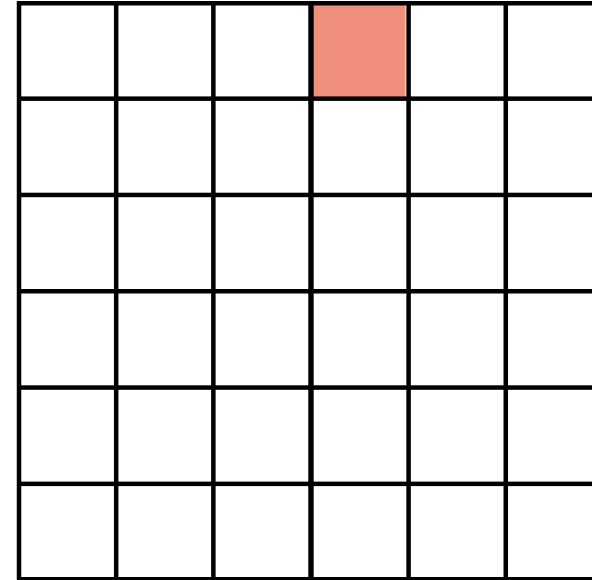
**Output**

# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



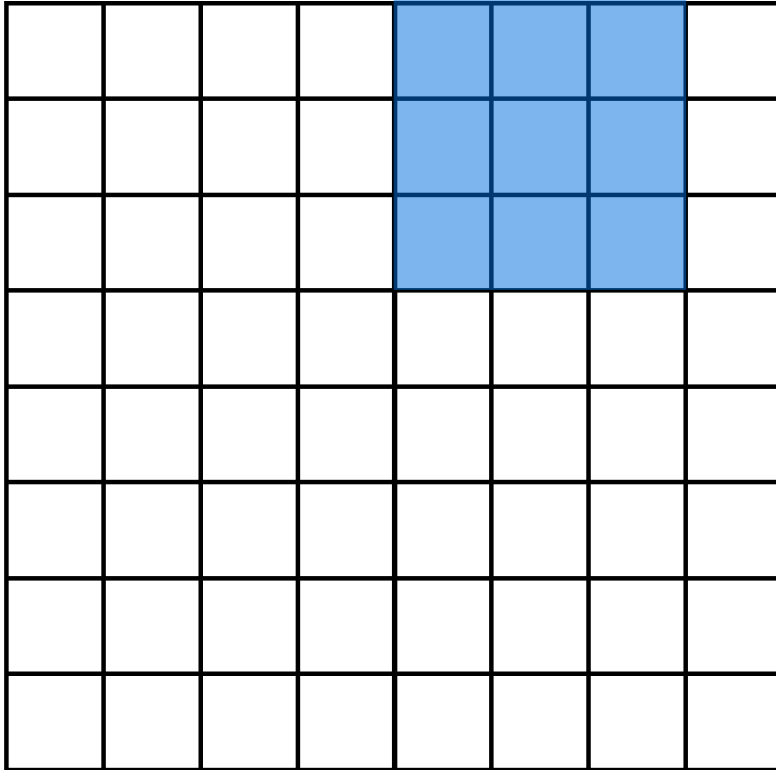
**Input**



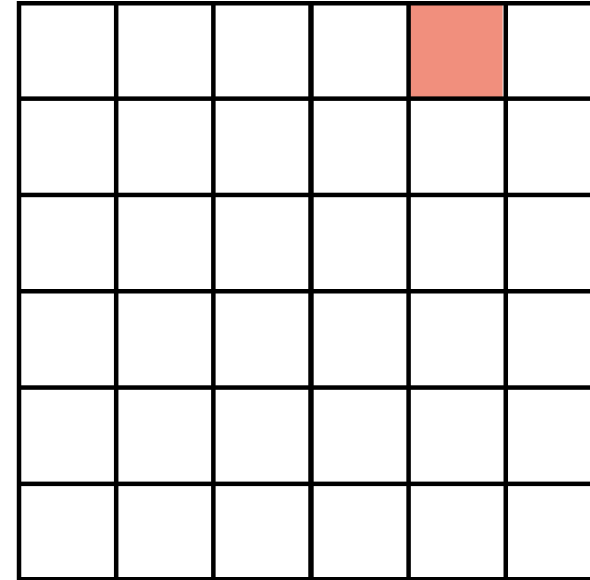
**Output**

# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



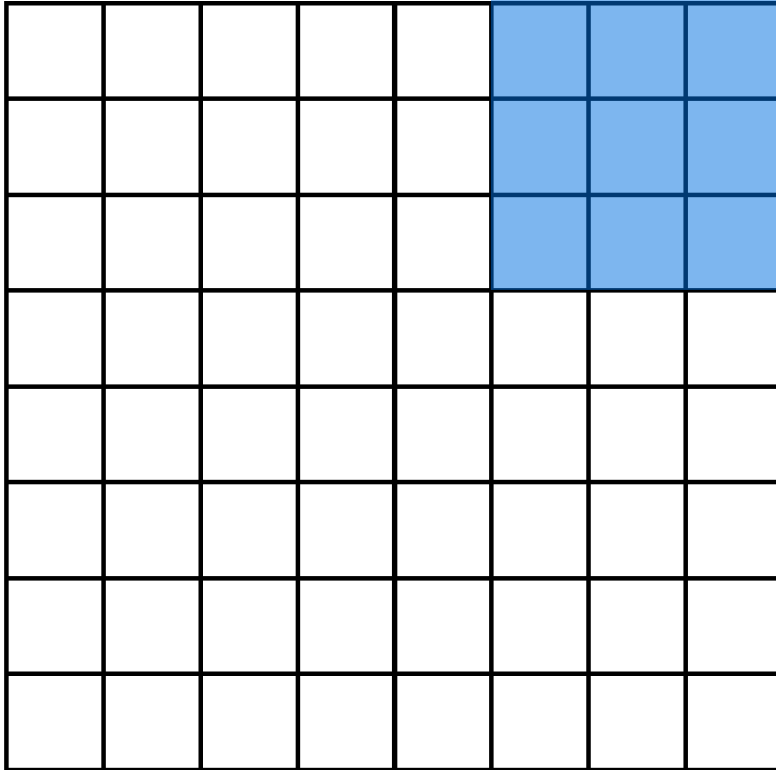
**Input**



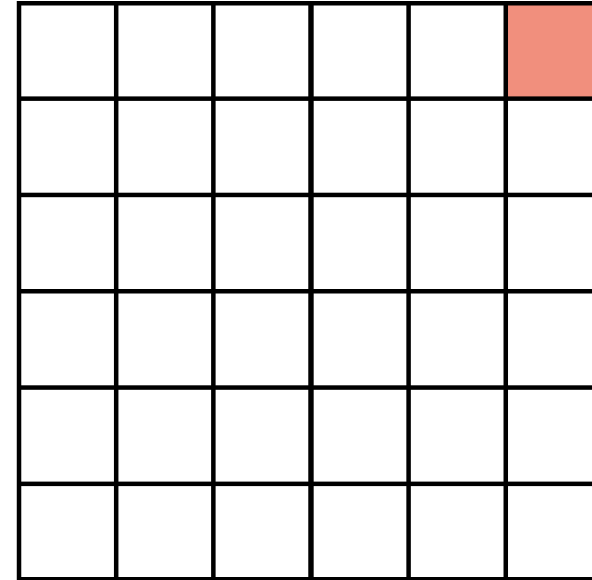
**Output**

# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



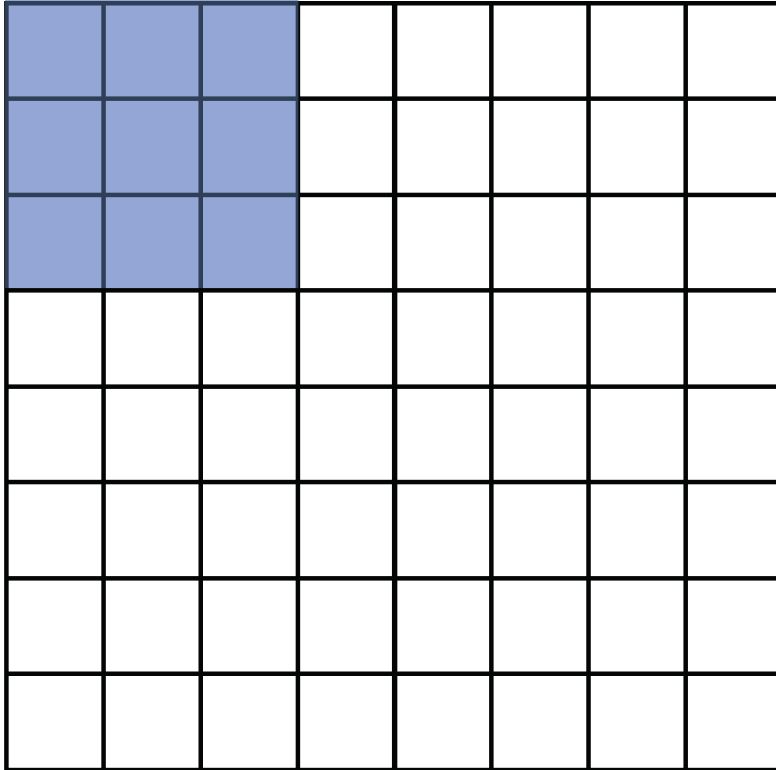
**Input**



**Output**

# Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



**Input**

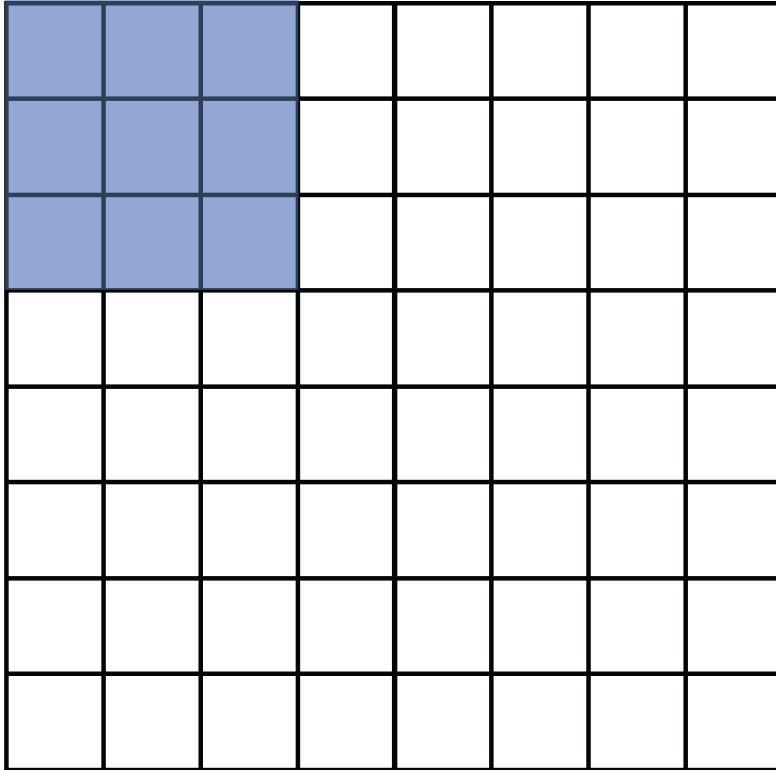
Recall that at each position, we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

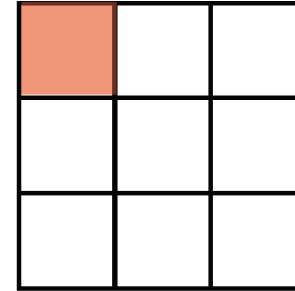
*(channel, row, column)*

# Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



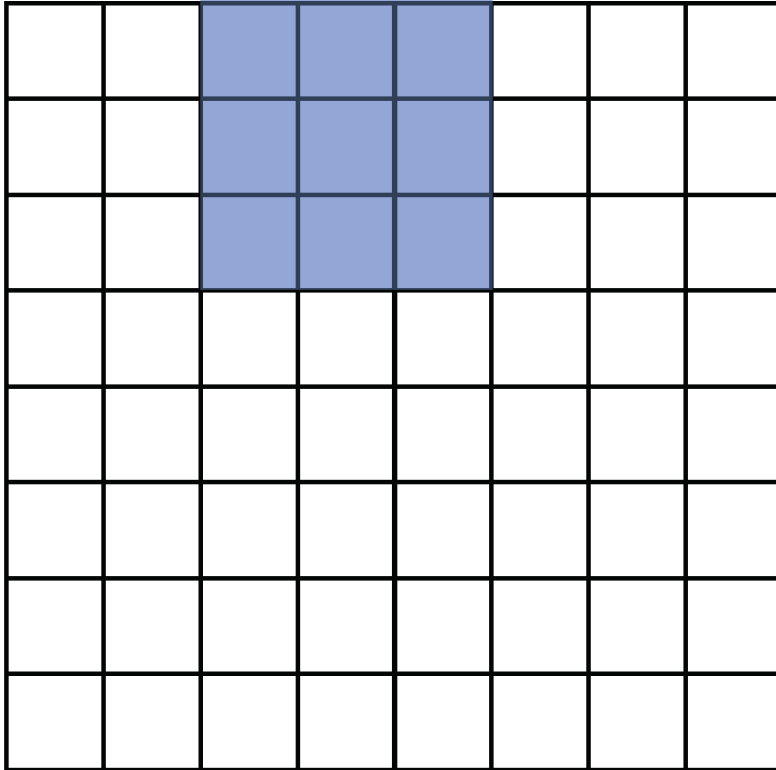
**Input**



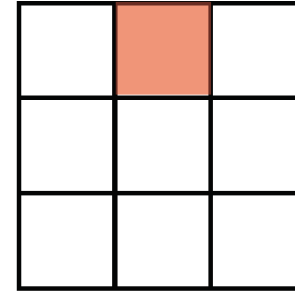
**Output**

# Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



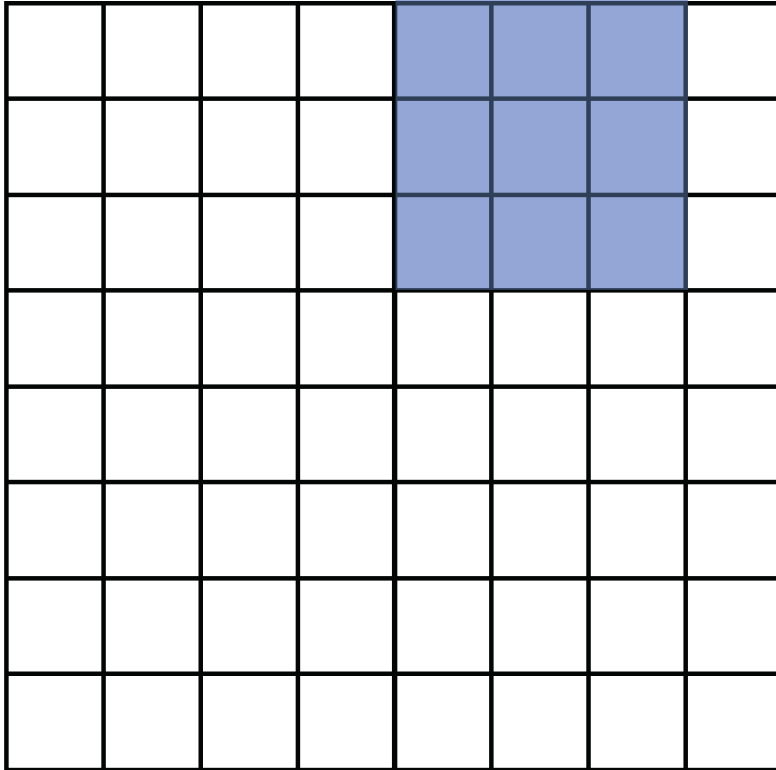
**Input**



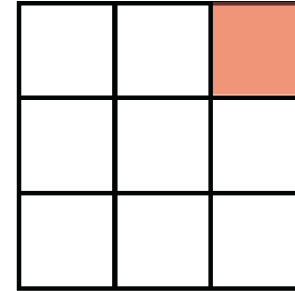
**Output**

# Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



**Input**

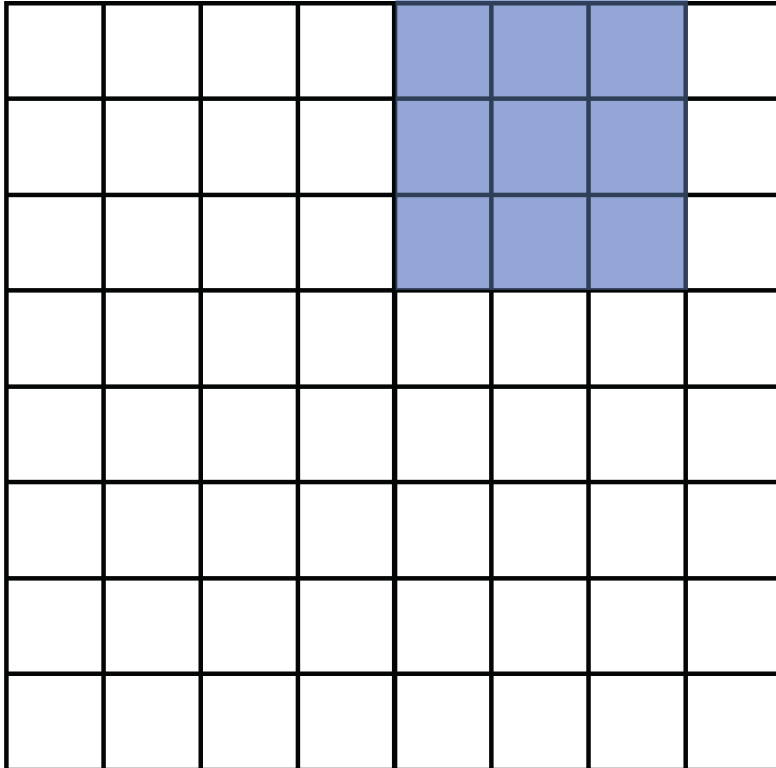


**Output**

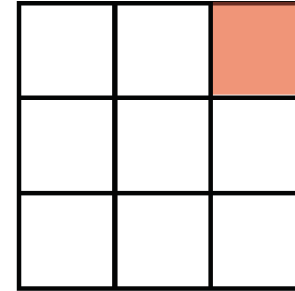


# Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



**Input**



**Output**

- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

# Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

**Input**


**Output**

# Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

**Input**


**Output**

# Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

**Input**


**Output**

# Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

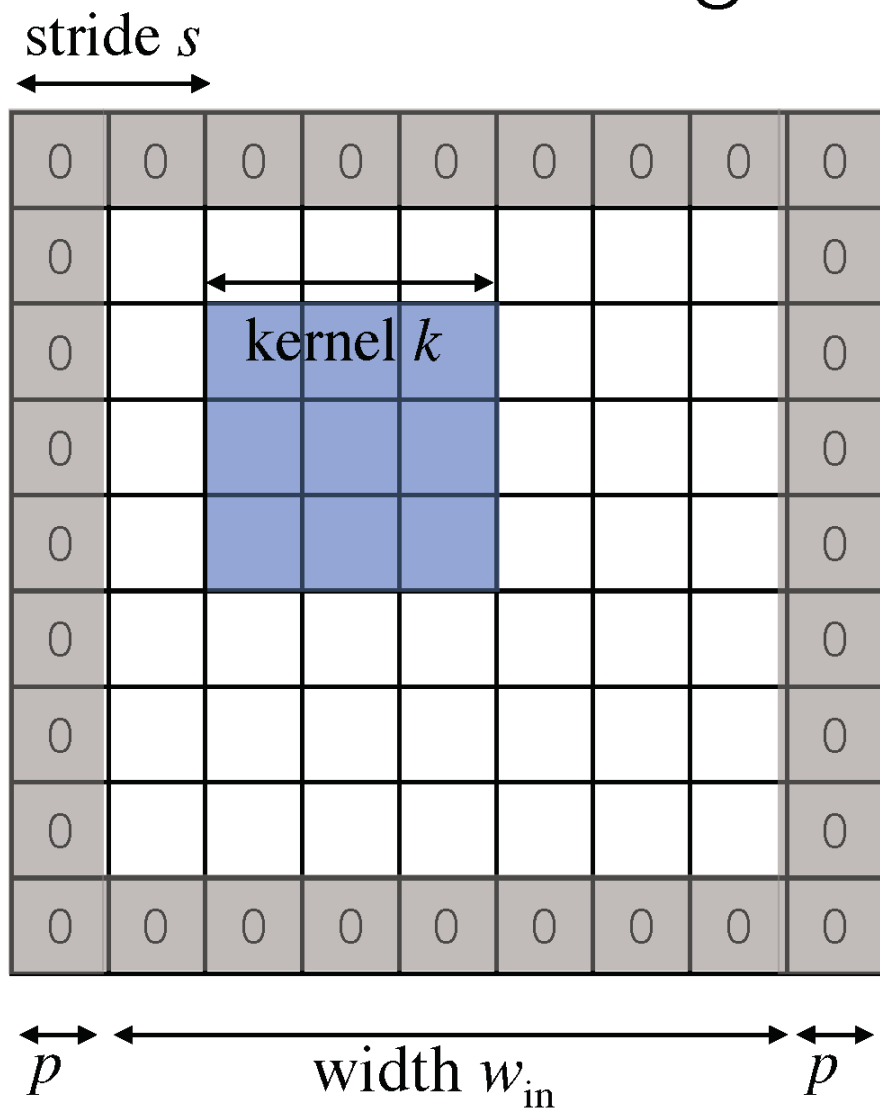
0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

**Input**


**Output**

# Convolution:

## How big is the output?

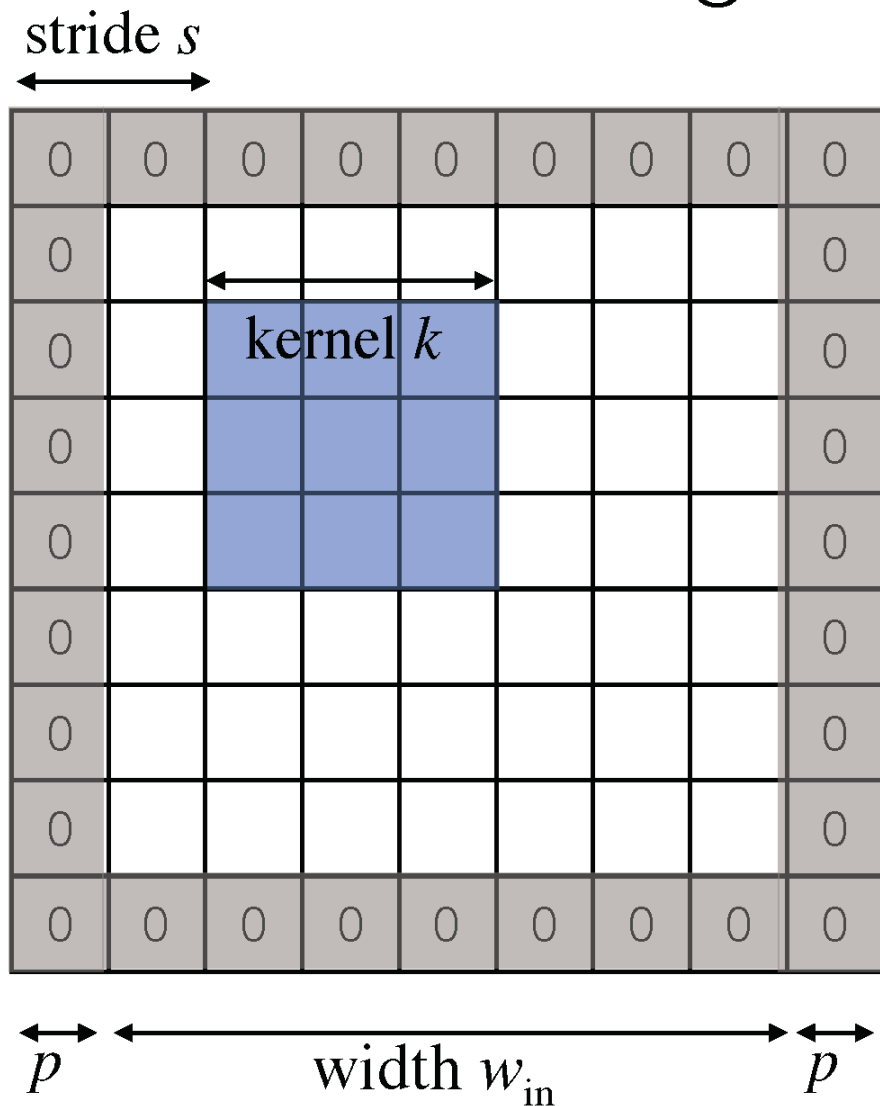


In general, the output has size:

$$w_{\text{out}} = \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1$$

# Convolution:

## How big is the output?



**Example:**  $k=3$ ,  $s=1$ ,  $p=1$

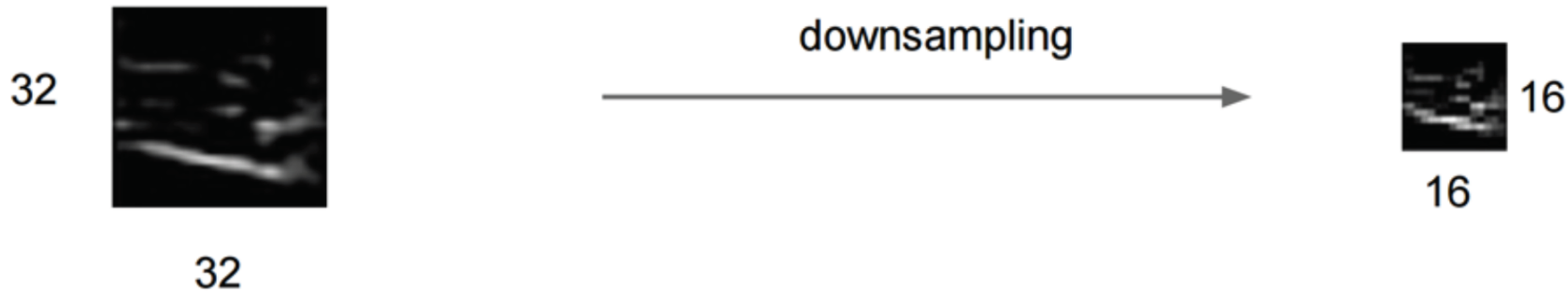
$$\begin{aligned}w_{out} &= \left\lfloor \frac{w_{in} + 2p - k}{s} \right\rfloor + 1 \\&= \left\lfloor \frac{w_{in} + 2 - 3}{1} \right\rfloor + 1 \\&= w_{in}\end{aligned}$$

VGGNet [Simonyan 2014]  
uses filters of this shape

# Pooling

For most ConvNets, **convolution** is often followed by **pooling**:

- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common
- Why might “avg” be a poor choice?

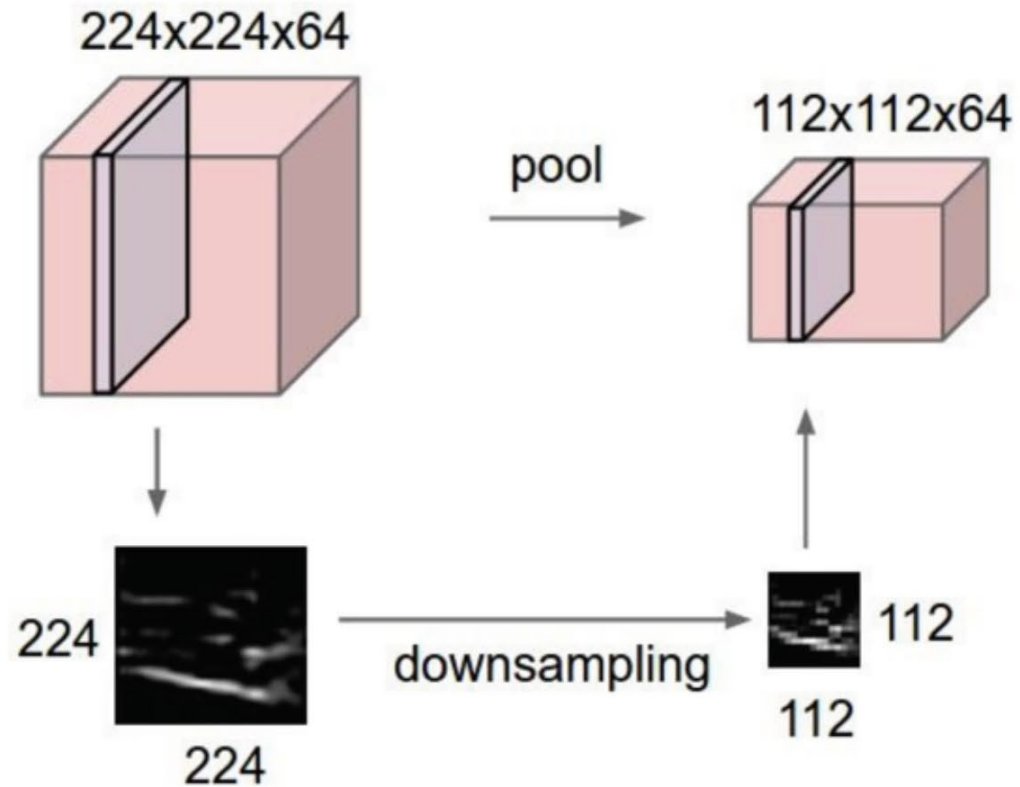


*Figure: Andrej Karpathy*

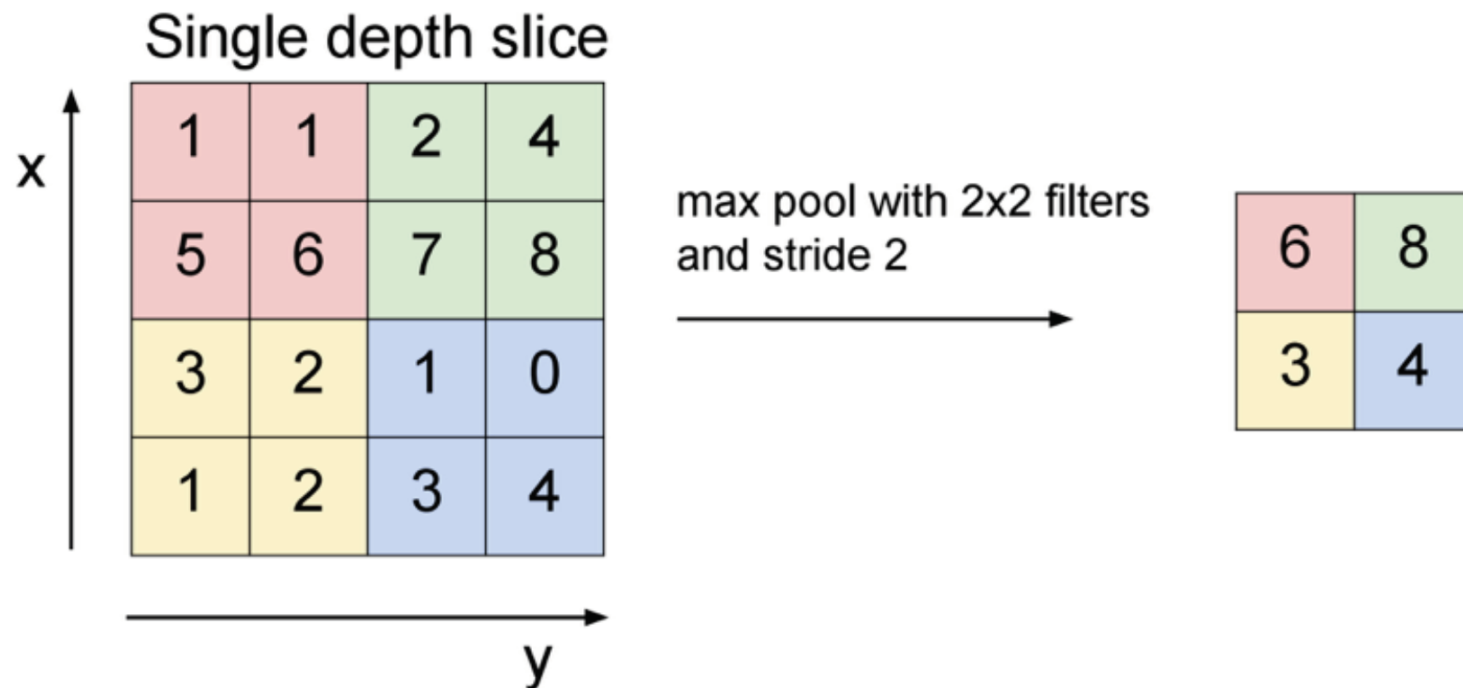


# Pooling

- makes the representations smaller and more manageable
- operates over each activation map independently:

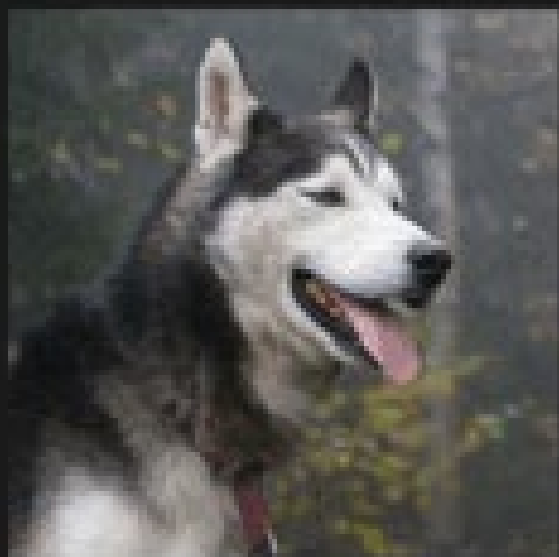


# Max Pooling



What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index



# Example ConvNet



Figure: Andrej Karpathy

# Example ConvNet

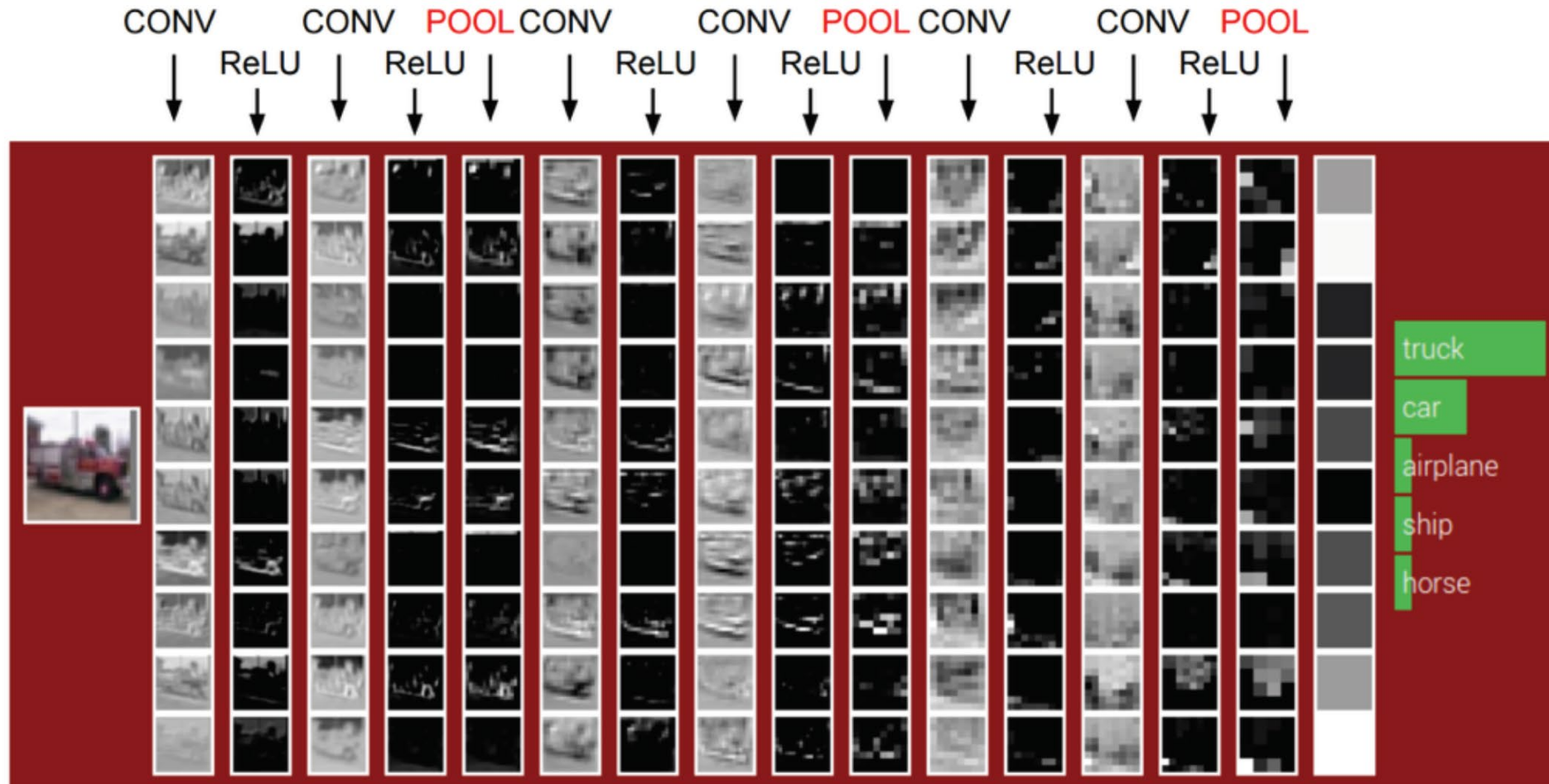


Figure: Andrej Karpathy

# Example ConvNet

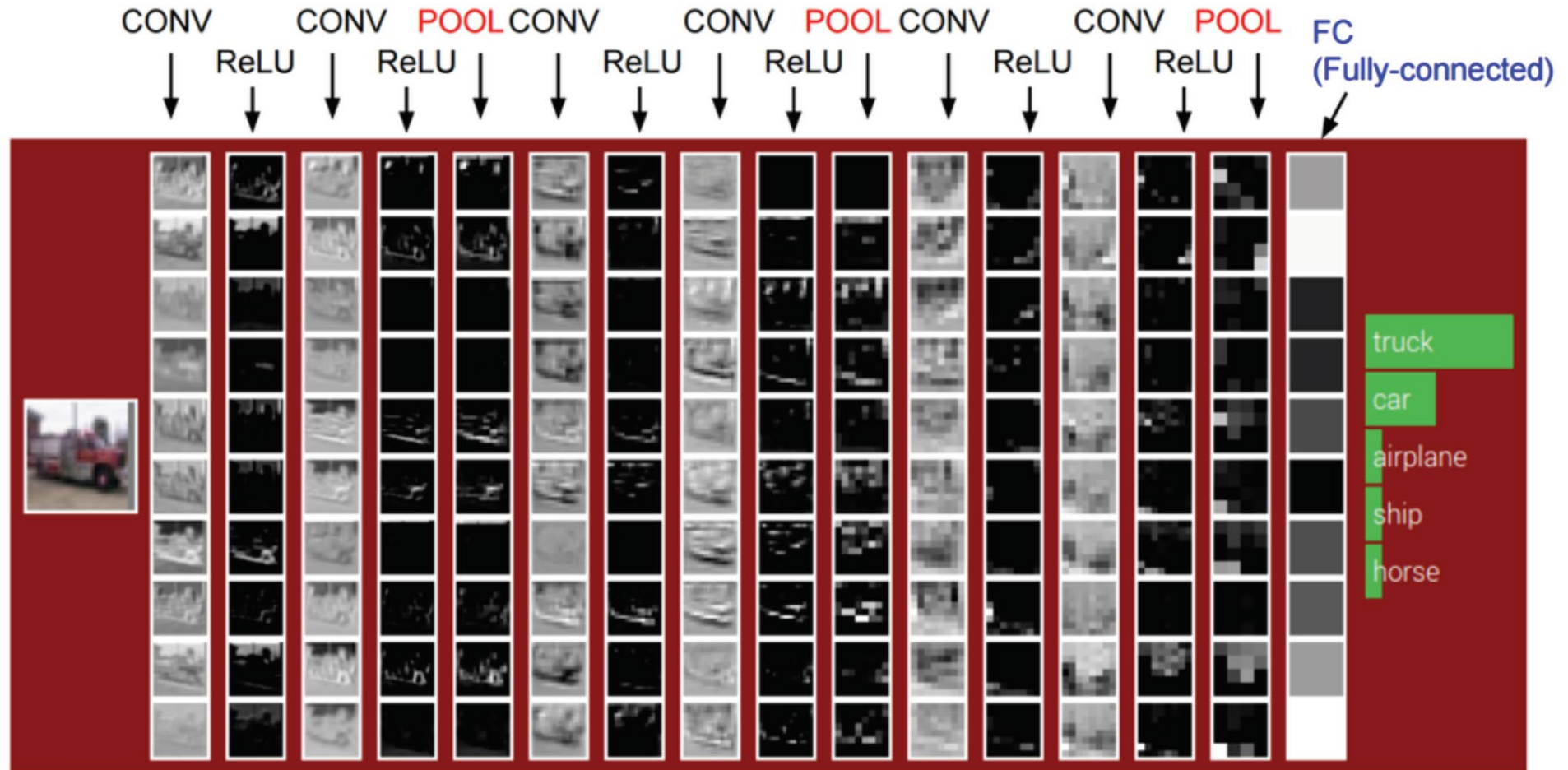
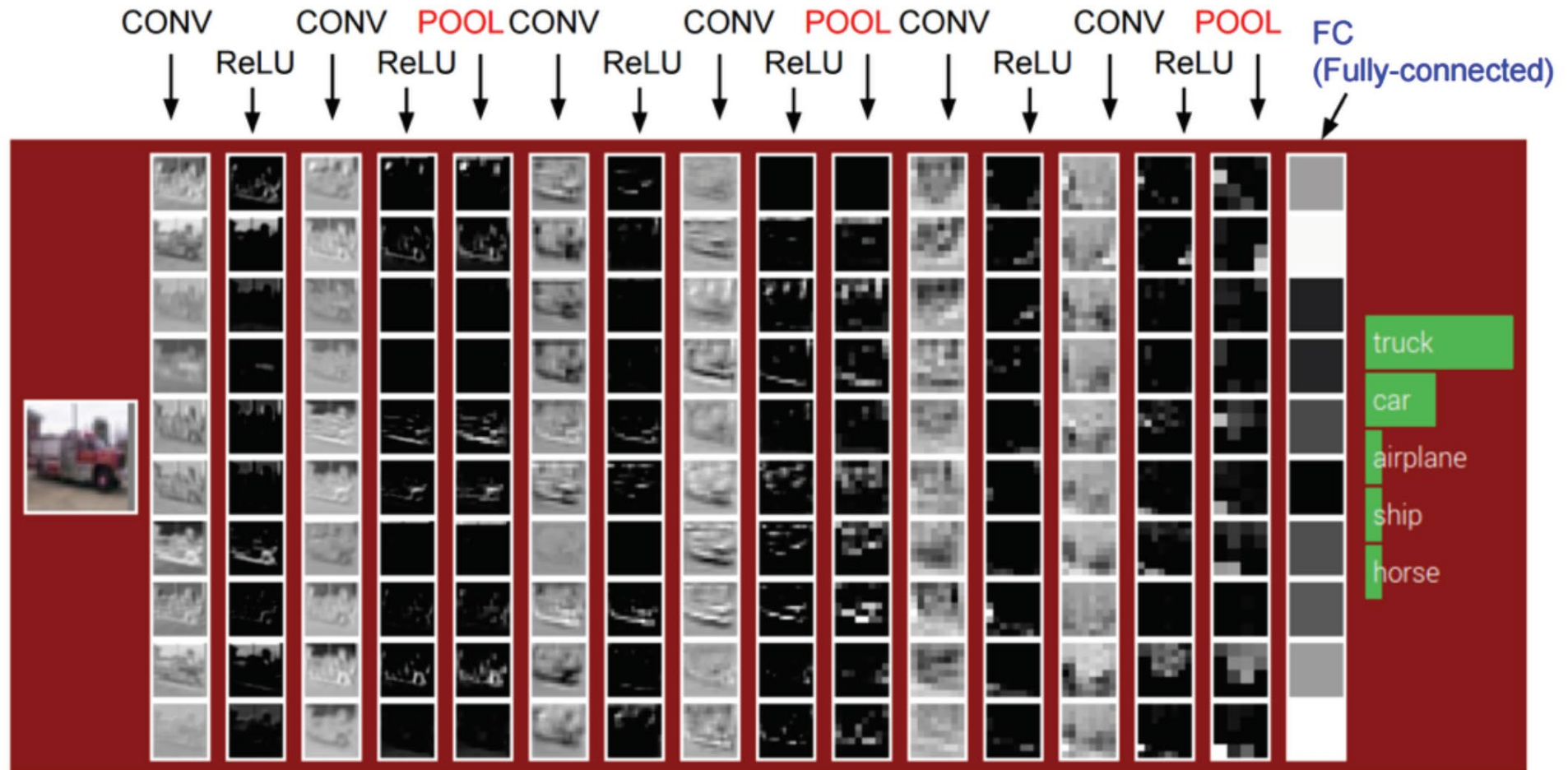


Figure: Andrej Karpathy



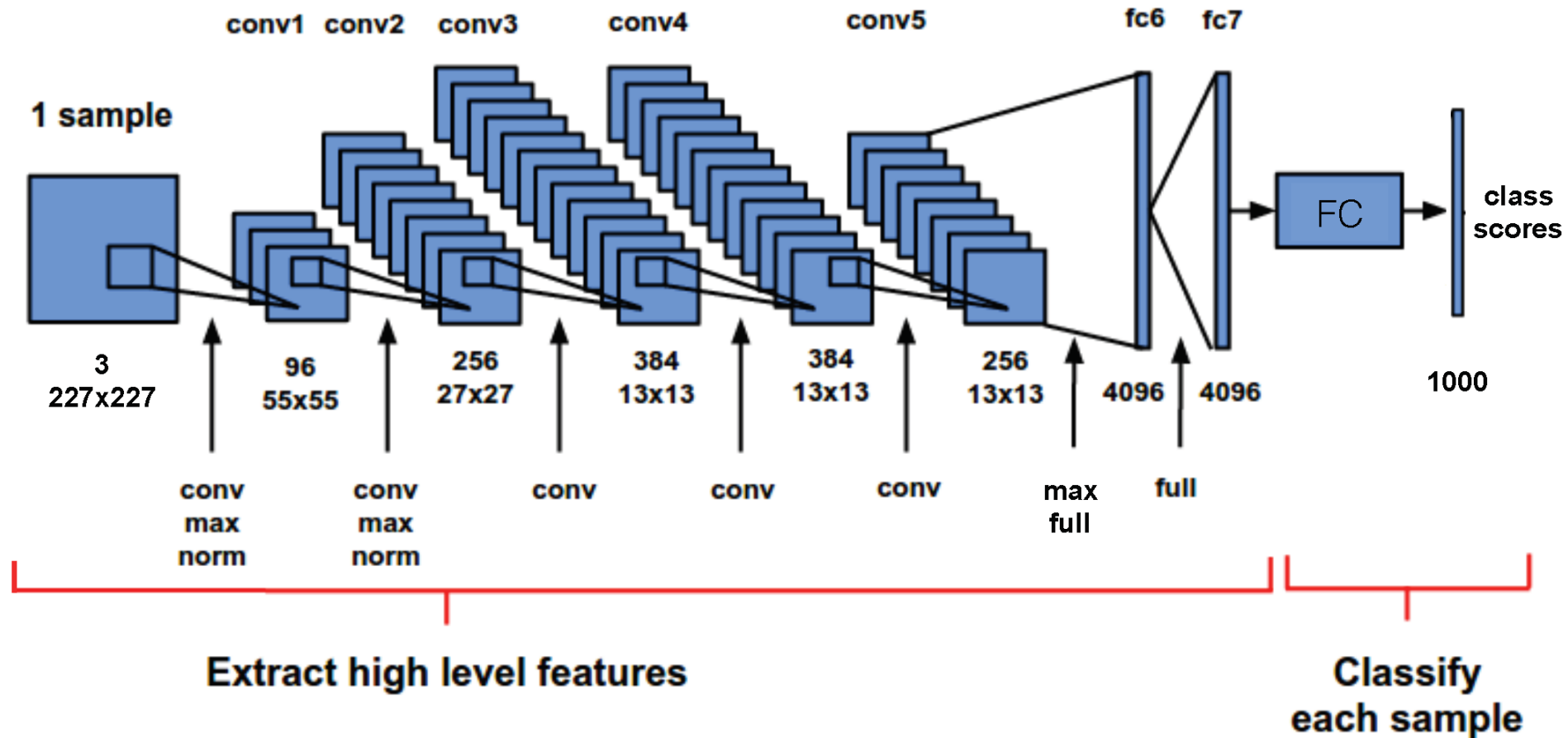
# Example ConvNet



10x3x3 conv filters, stride 1, pad 1  
2x2 pool filters, stride 2

Figure: Andrej Karpathy

# Example: AlexNet [Krizhevsky 2012]



“max”: max pooling

“norm”: local response normalization

“full”: fully connected

Figure: [Karnowski 2015] (with corrections)