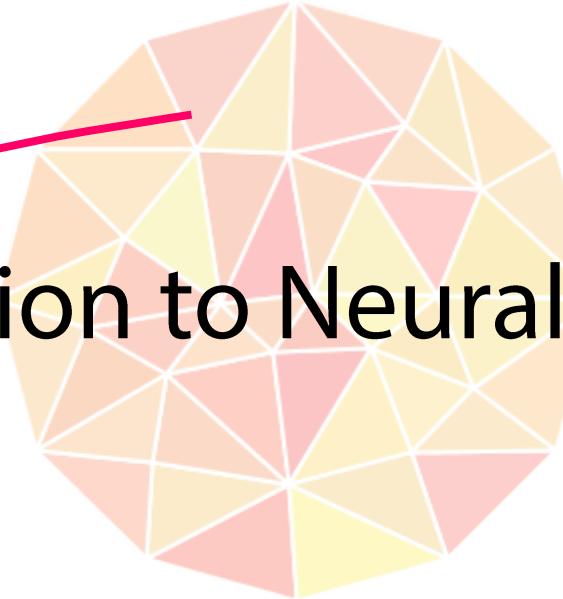


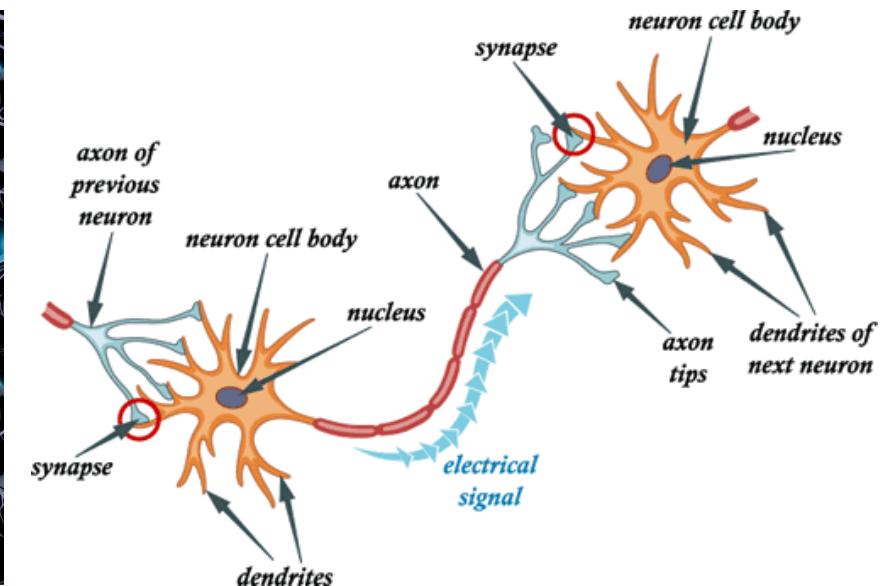
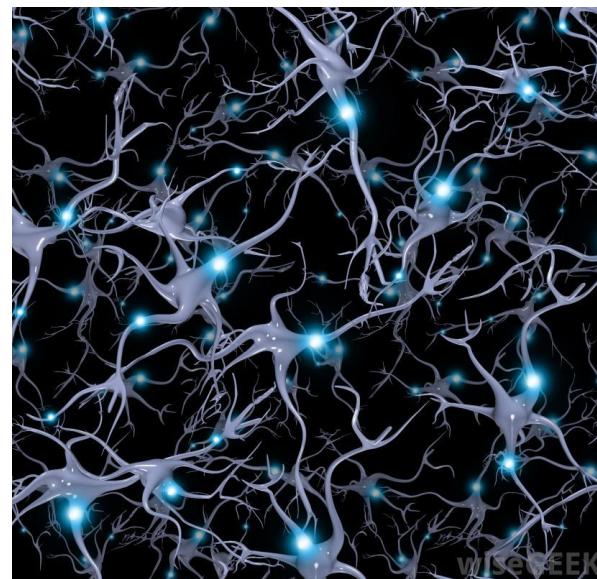
A Quick and Dirty

# Introduction to Neural Networks

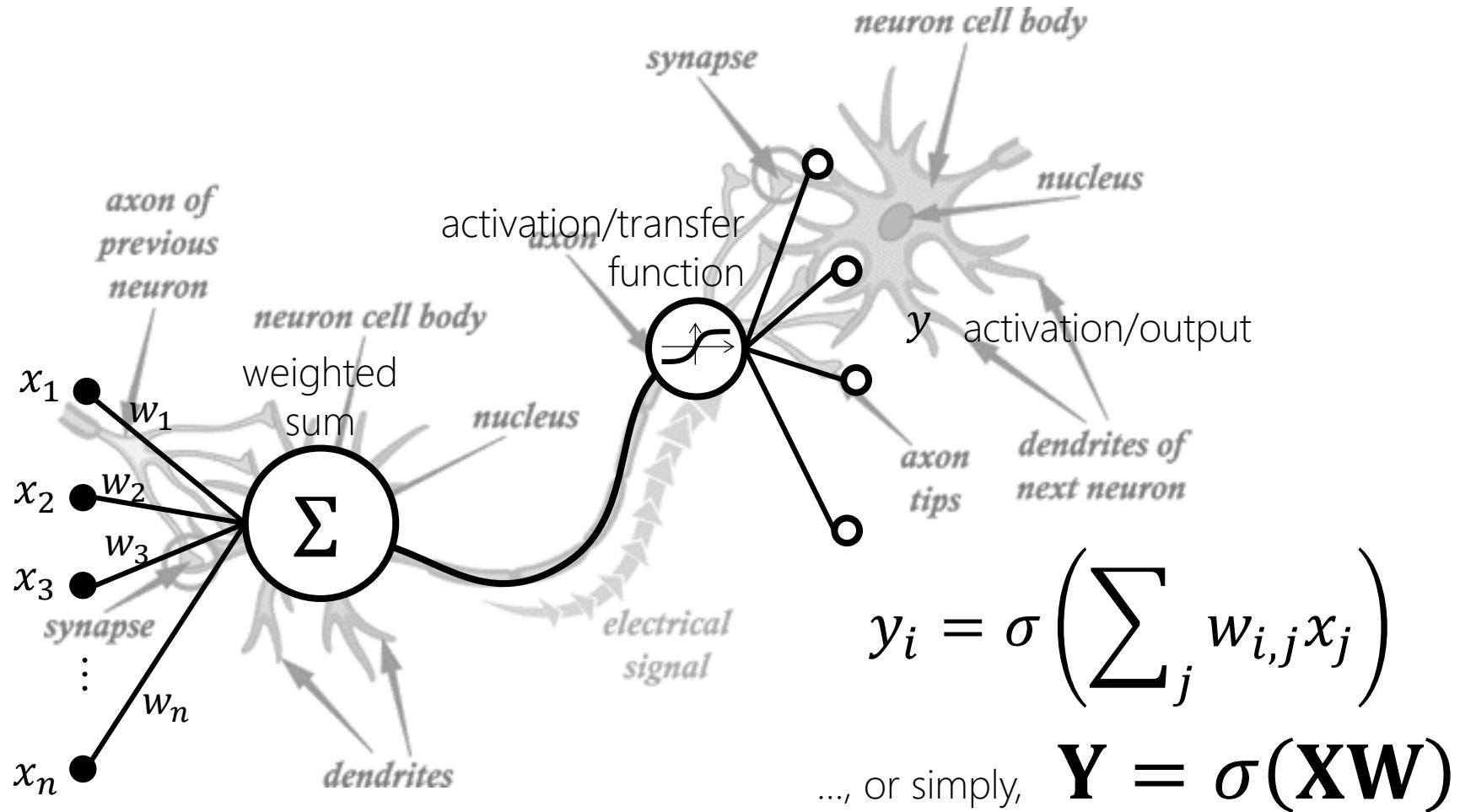


Stephen Baek

# Biological Neural Networks



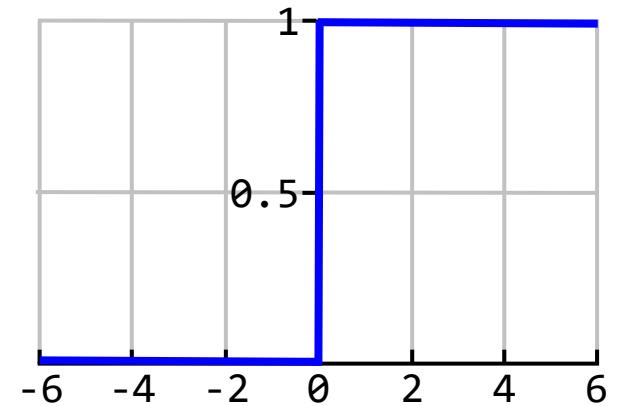
# Artificial Neural Networks (ANN)



# Activation/Transfer Functions

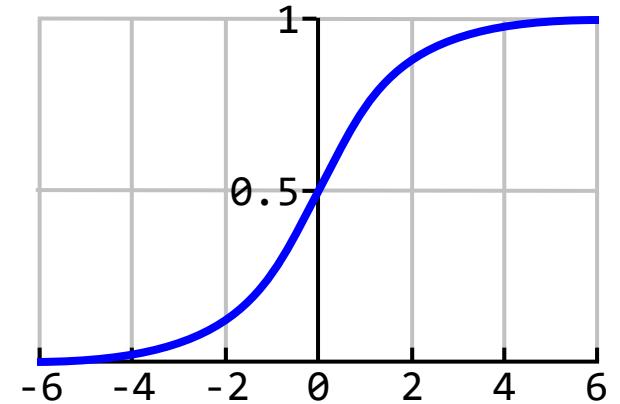
- “Threshold” in biological systems
- Step Function:

$$y = \sigma(z) = \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



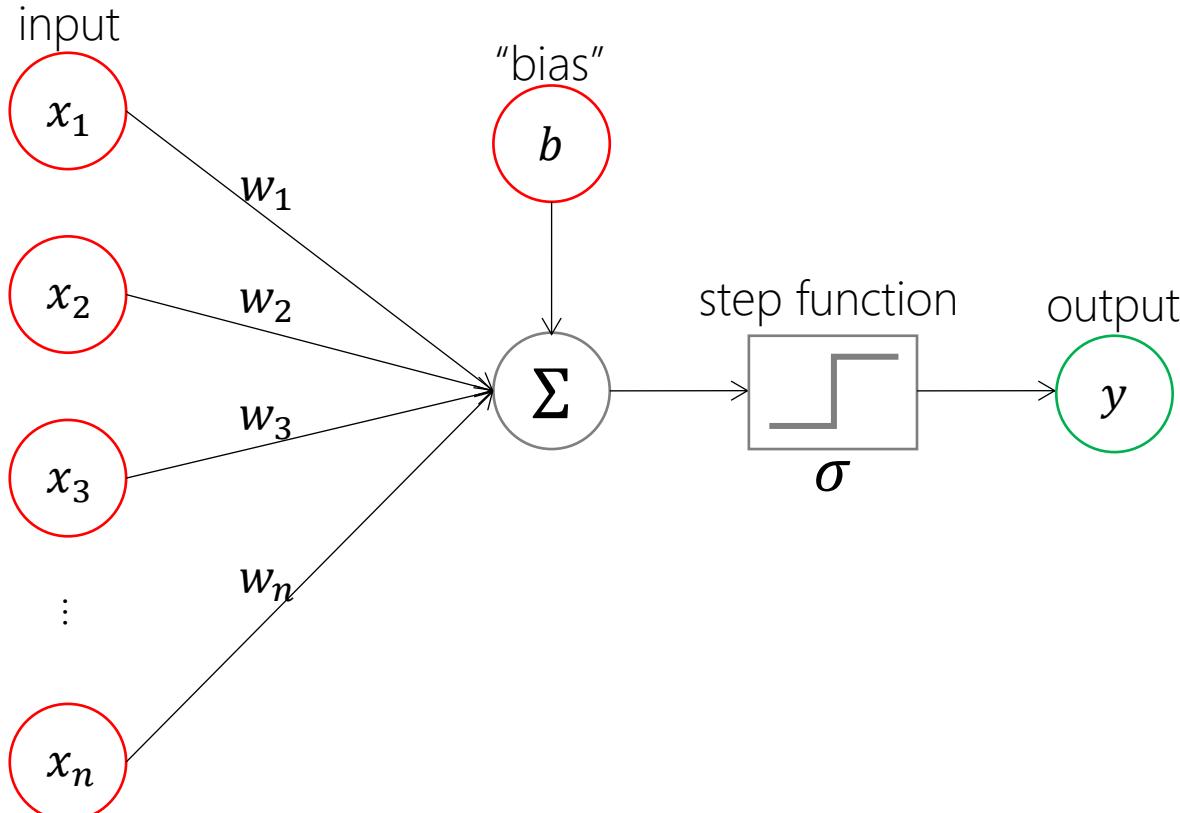
- Sigmoid Function:

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



# The simplest model: the “Perceptron”

- Frank Rosenblatt (1957)



$$y = \sigma \left( \sum_j w_j x_j + b \right)$$

For each example  $(x^{(t)}, \hat{y}^{(t)})$ , try to minimize:

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} (y^{(t)} - \hat{y}^{(t)})^2 \\ \frac{\partial \mathcal{L}}{\partial w_j} &= (y^{(t)} - \hat{y}^{(t)}) \frac{\partial y^{(t)}}{\partial w_j} \\ &= (y^{(t)} - \hat{y}^{(t)}) x_j^{(t)}\end{aligned}$$

Weight update scheme:

$$w(t+1) = w(t) - r(y^{(t)} - \hat{y}^{(t)}) x_j^{(t)}$$

learning rate

# The simplest model: the “Perceptron”

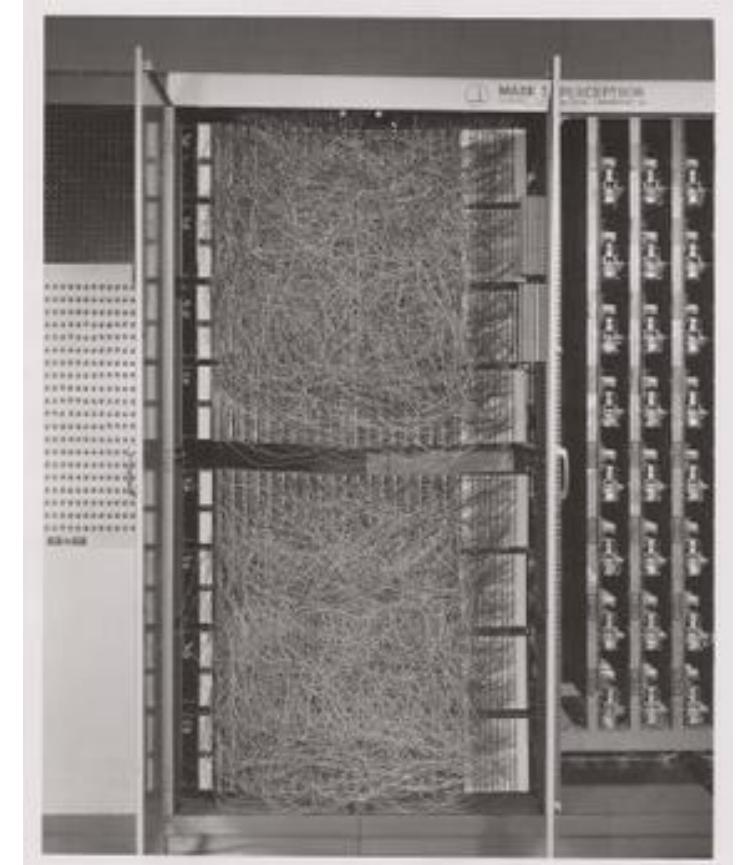
- Frank Rosenblatt (1957)
  - The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.
  - Designed for image recognition.
  - Connected to a camera that used a 20x20 cadmium sulfide photocell array to produce a 400-pixel image.



<http://www.rutherfordjournal.org/article040101.html>

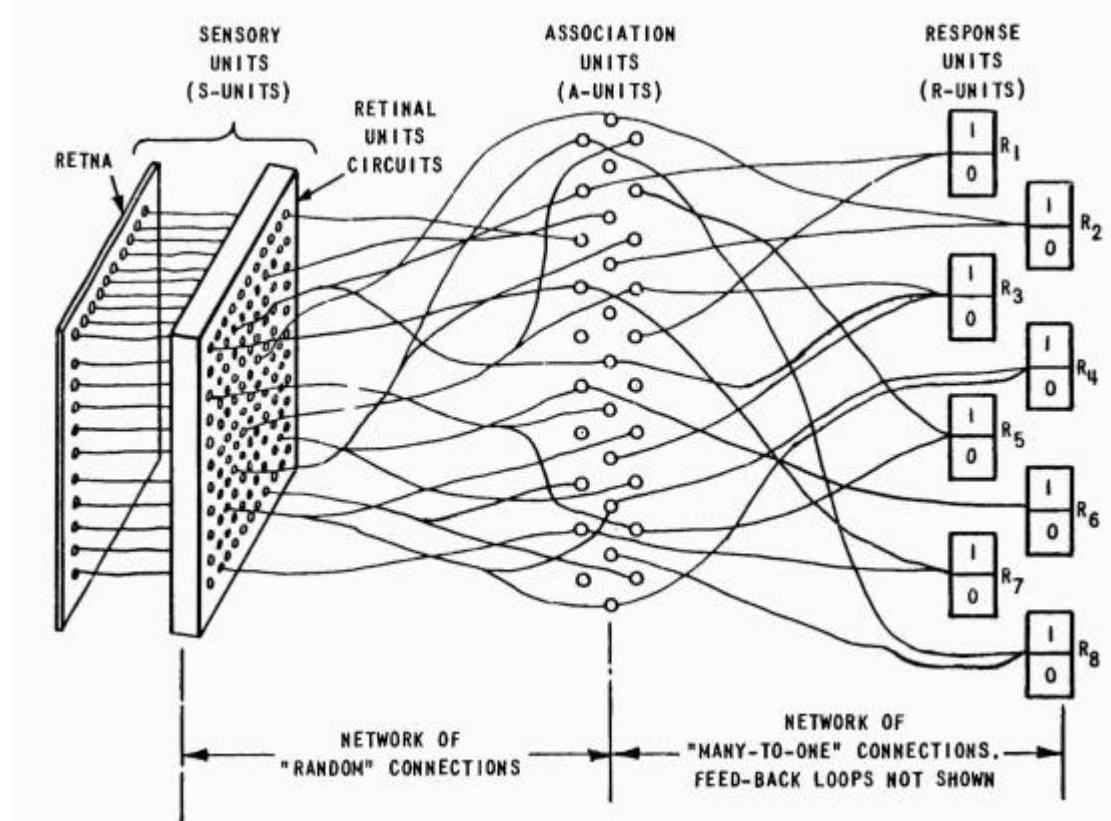
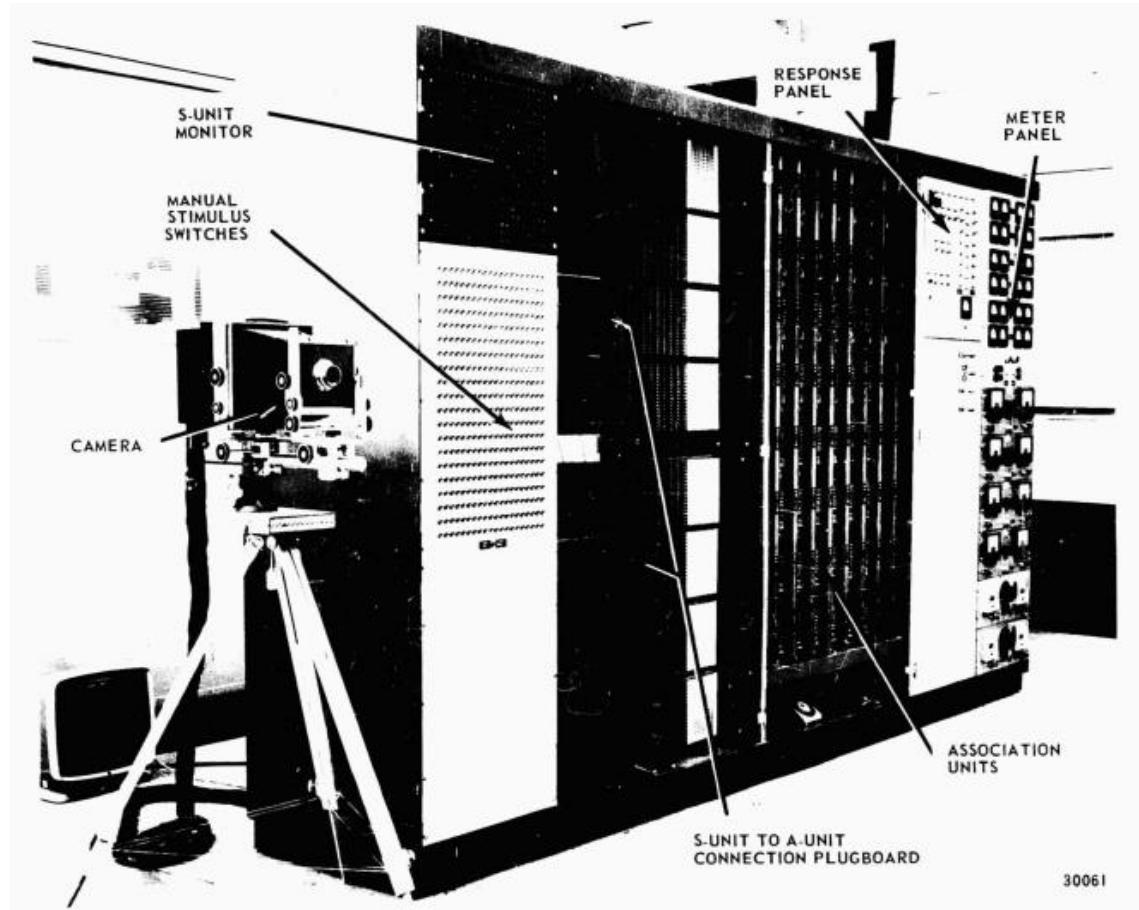
# The simplest model: the “Perceptron”

- Frank Rosenblatt (1957)
  - The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.
  - Designed for image recognition.
  - Connected to a camera that used a 20x20 cadmium sulfide photocell array to produce a 400-pixel image.
  - Weights were encoded in potentiometers.
  - Weight updates were performed by electric motors.



Wikipedia – Perceptron

# The simplest model: the “Perceptron”

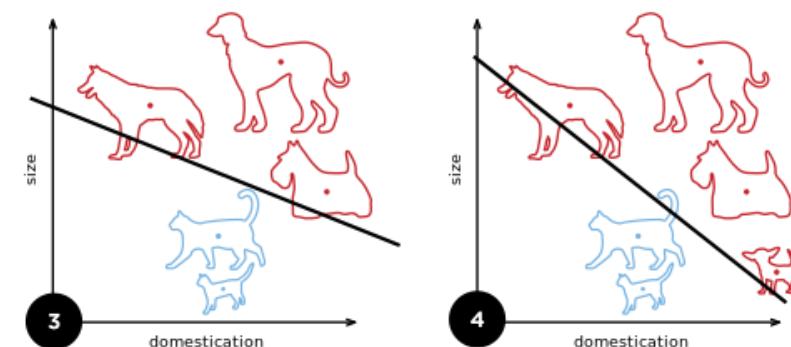
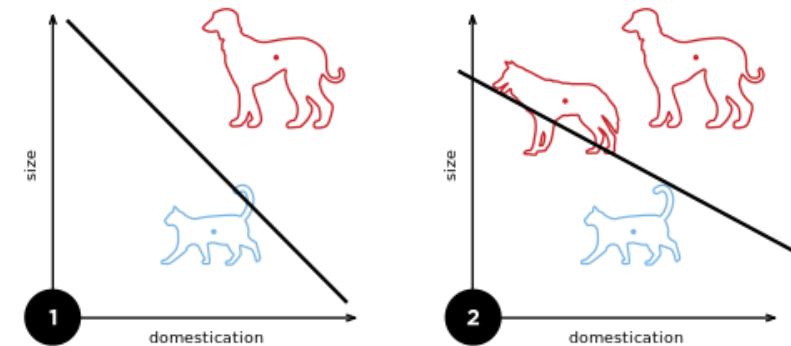
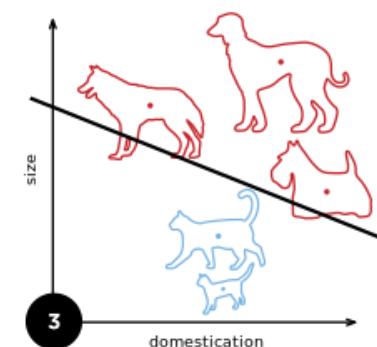
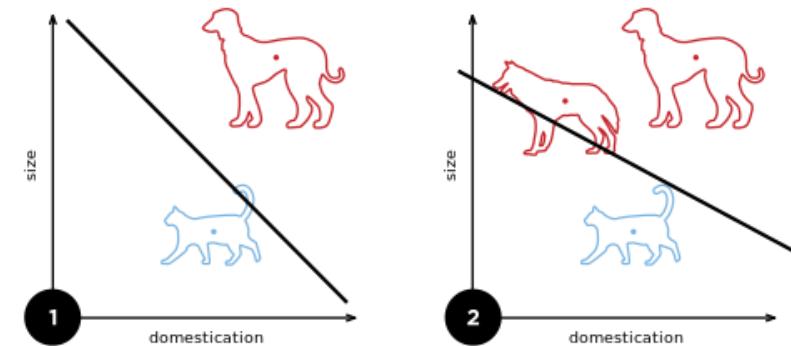
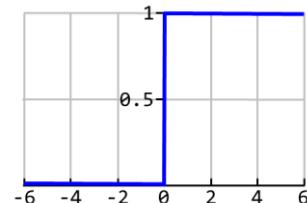


# Perceptron as Linear Binary Classification

- Perceptron returns **1** iff.  $\sum_j w_j x_j + b \geq 0$  and **0** otherwise.
- Essentially, this is equivalent to linear binary classification problem

$$z = \sum_j w_j x_j + b$$

$$y = \sigma(z) = \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



# The XOR Problem

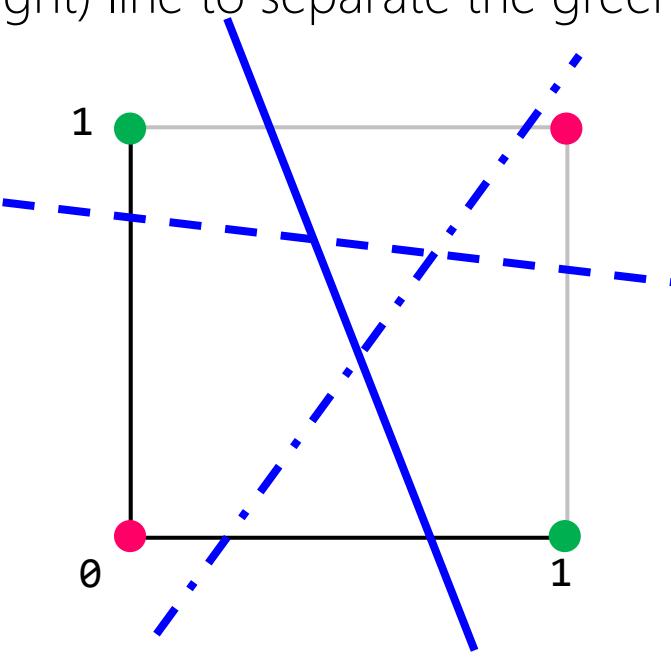
- XOR: Exclusive OR:

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

- Q. Can you implement the XOR gate with the perceptron model?

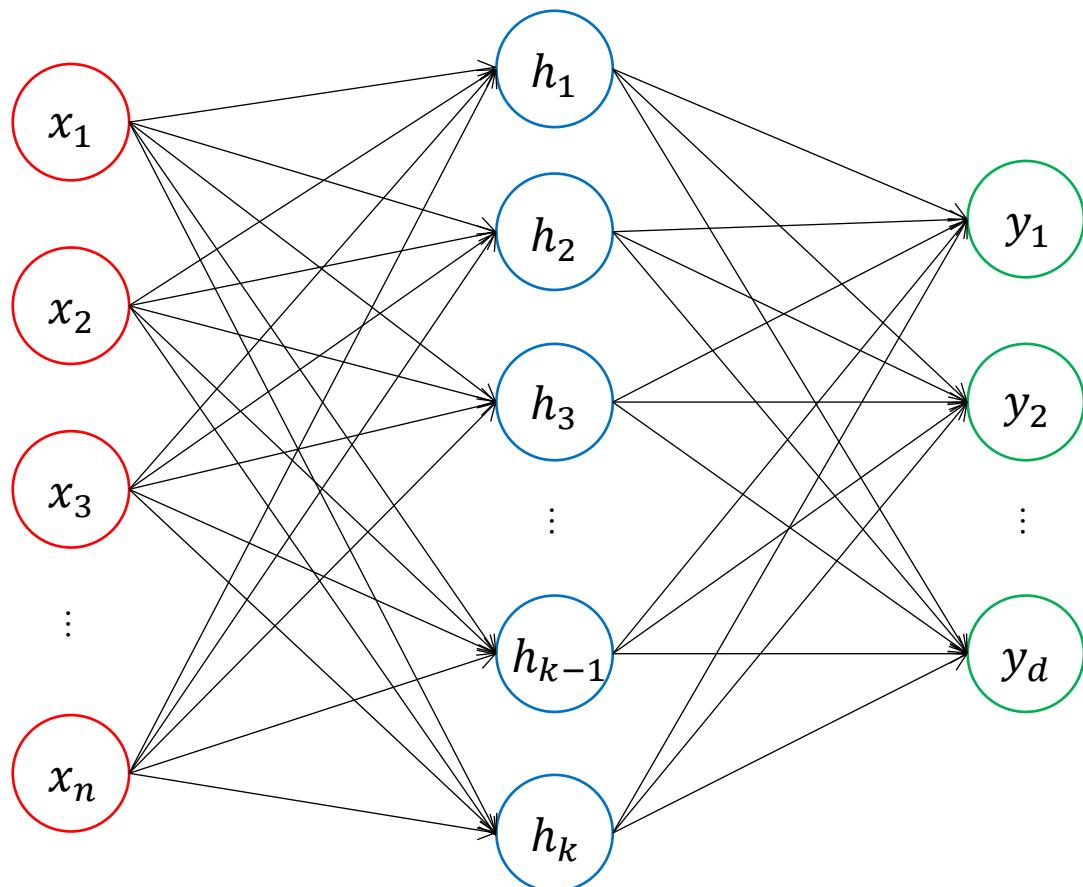
# The XOR Problem

- The quick answer is NO.
- Why?
  - a perceptron is a linear classifier.
  - Can you draw a single (straight) line to separate the green and red dots?



# Multi-Layer Perceptron (MLP)

- Minsky & Papert. (1969)

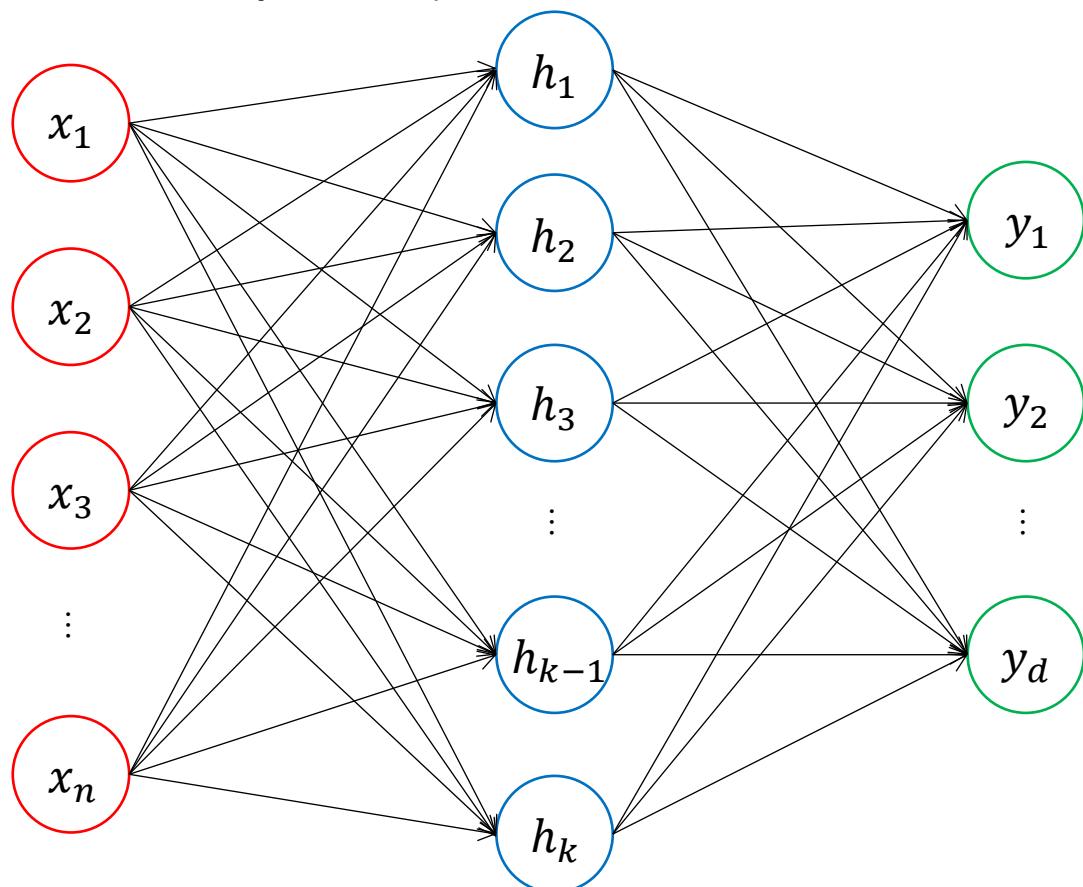


The perceptron has shown itself worthy of study despite (and even because of!) its severe limitations. It has many features to attract attention: its linearity; its intriguing learning theorem; its clear paradigmatic simplicity as a kind of parallel computation. There is no reason to suppose that any of these virtues carry over to the many-layered version. Nevertheless, we consider it to be an important research problem to elucidate (or reject) our intuitive judgement that the extension to multilayer systems is sterile.

*Minsky & Papert (1969, pp. 231-232)*

# Multi-Layer Perceptron (MLP)

- Minsky & Papert. (1969)



$$h_i = \sigma_h \left( \sum_j w_{h,i,j} x_j + b_h \right) \quad y_k = \sigma_o \left( \sum_i w_{o,k,i} h_i + b_o \right)$$

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

BRACE YOURSELVES

$$\frac{\partial \mathcal{L}}{\partial w} = ?$$



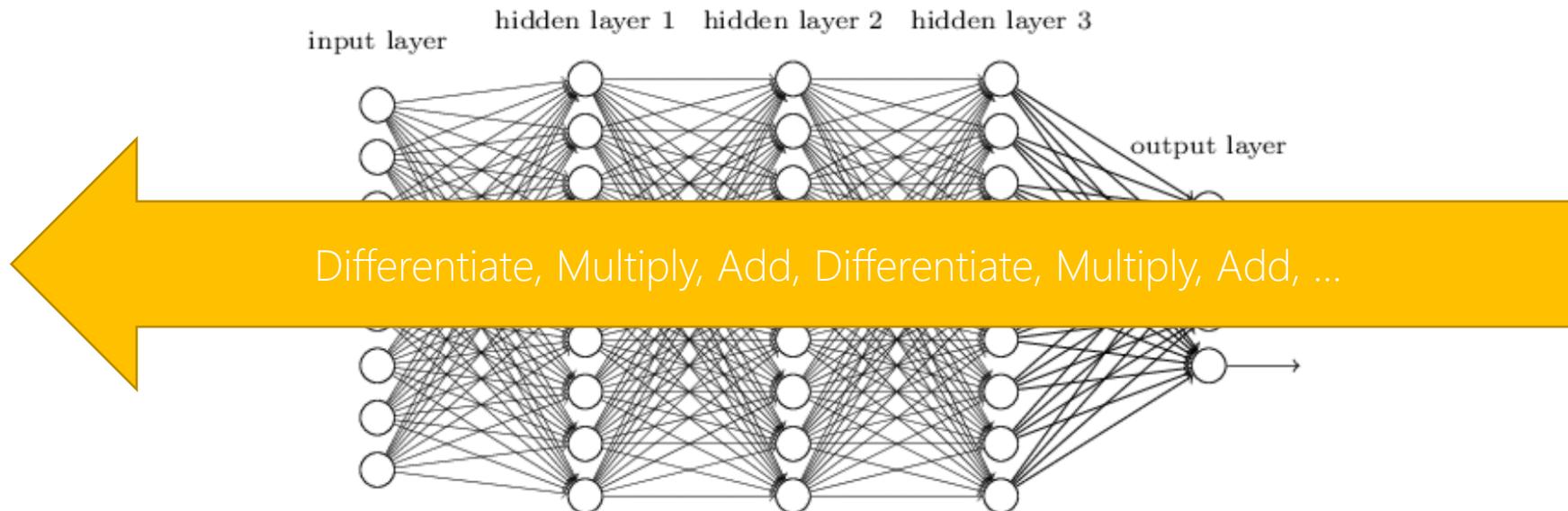
- M&P's criticism on (single-layer) perceptrons
- No practical way to train MLPs

→ The First AI Winter (1970s)

A FEW  
MOMENTS LATER

# Backpropagation

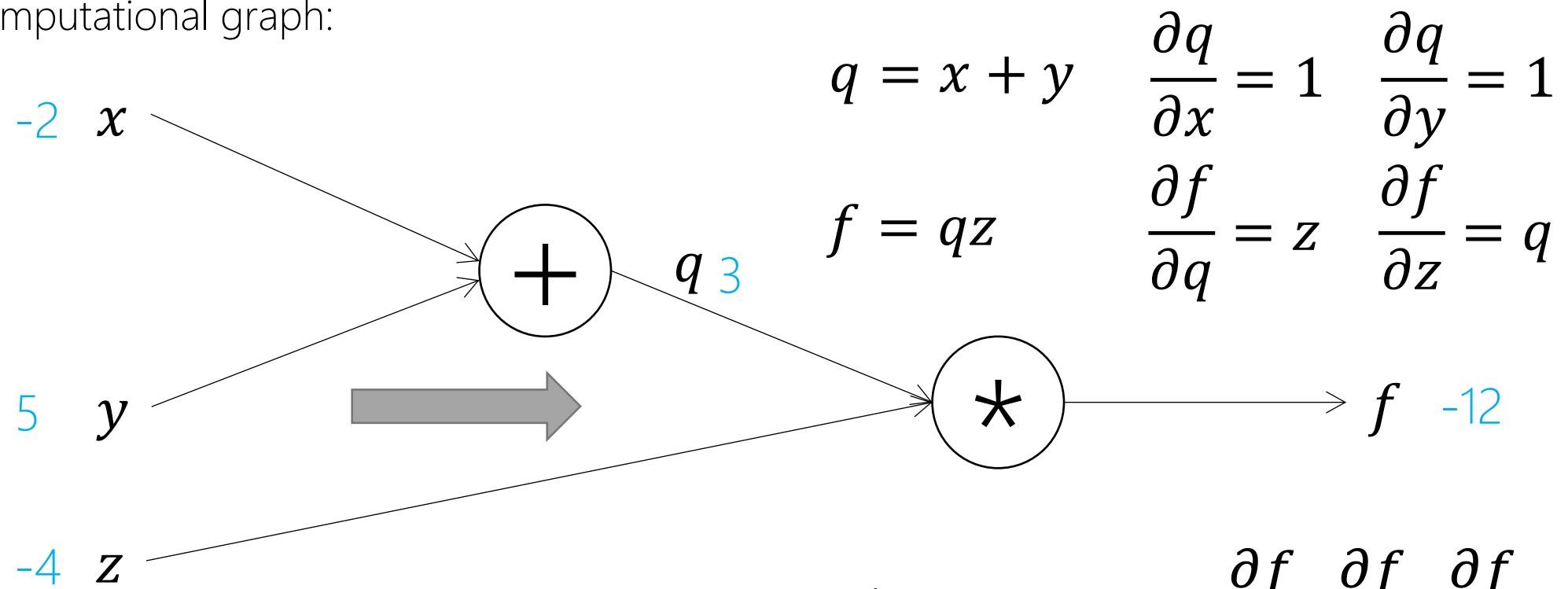
- Rumelhart, Hinton, and Williams. (1986).
- A popular training method for neural nets
- Propagate what? “the gradient of the current error”



# Backpropagation

- A simple example:  $f(x, y, z) = (x + y)z$

- Computational graph:

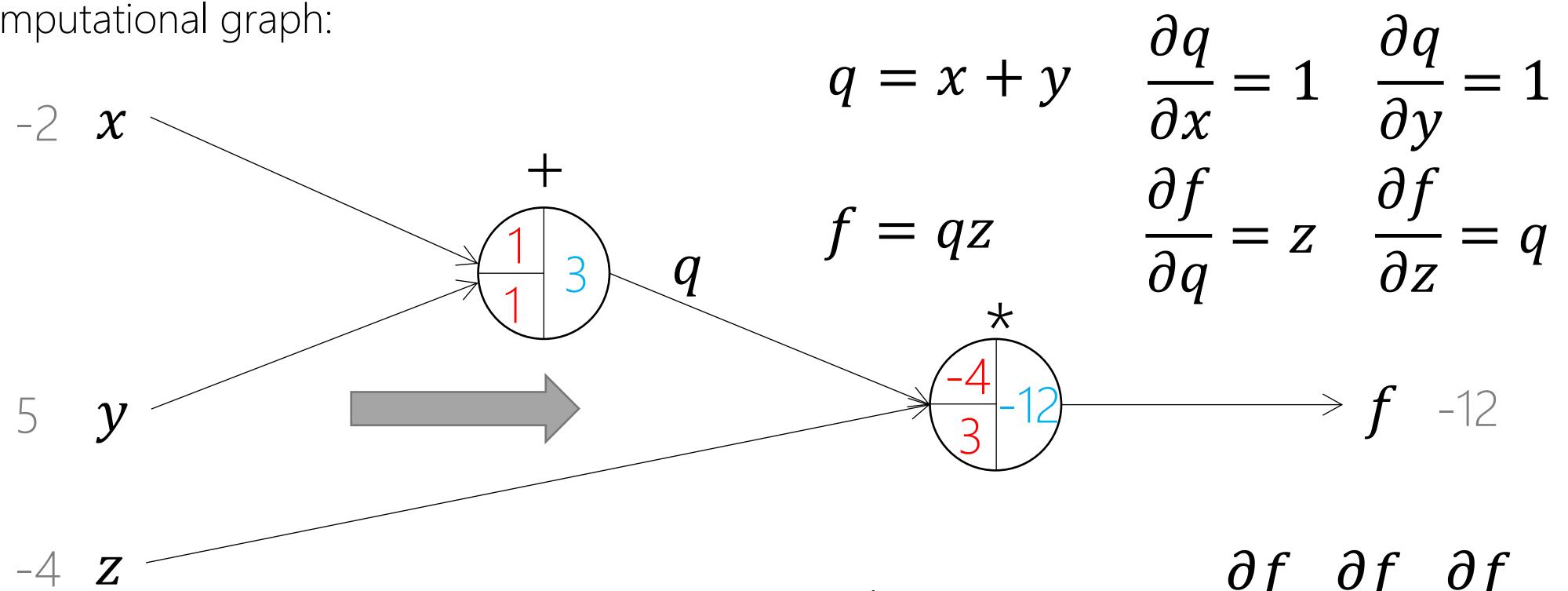


What we want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

# Backpropagation

- A simple example:  $f(x, y, z) = (x + y)z$

- Computational graph:

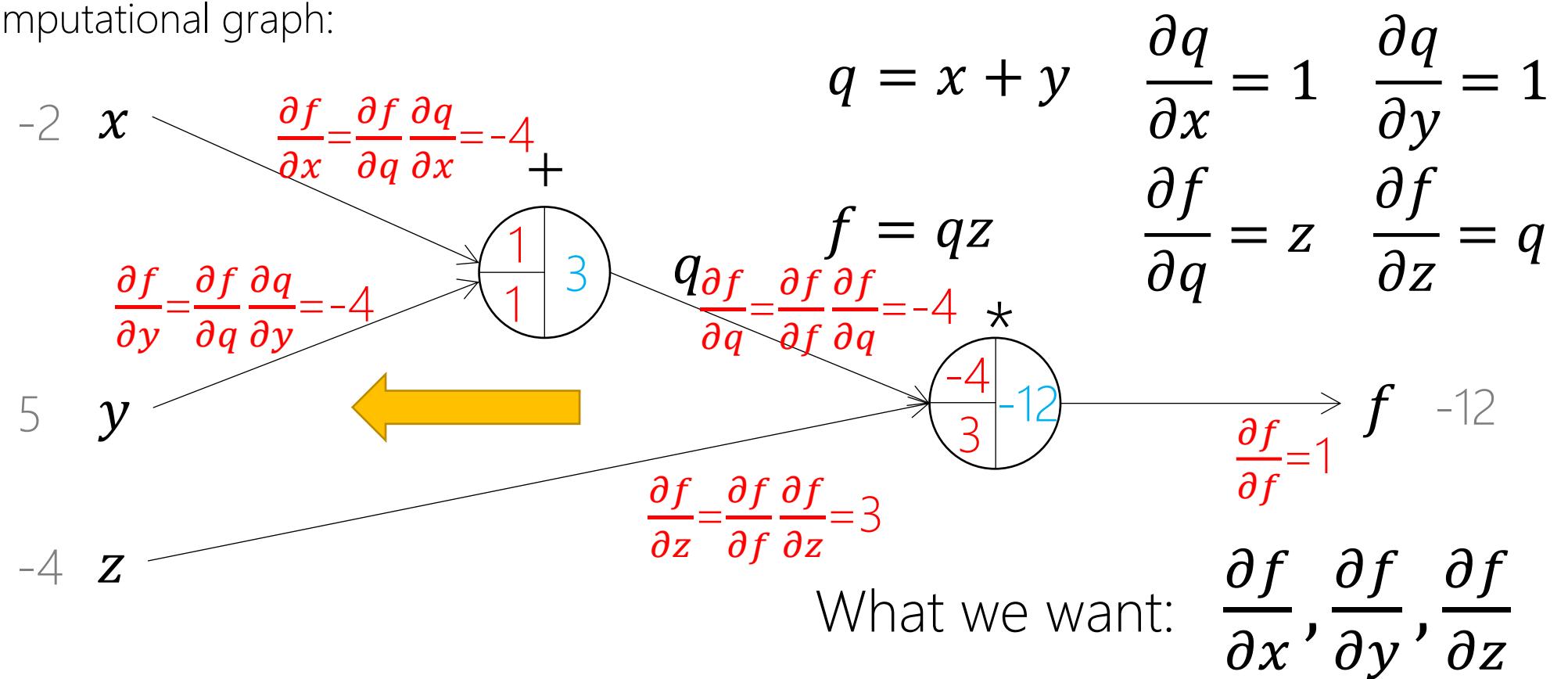


What we want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

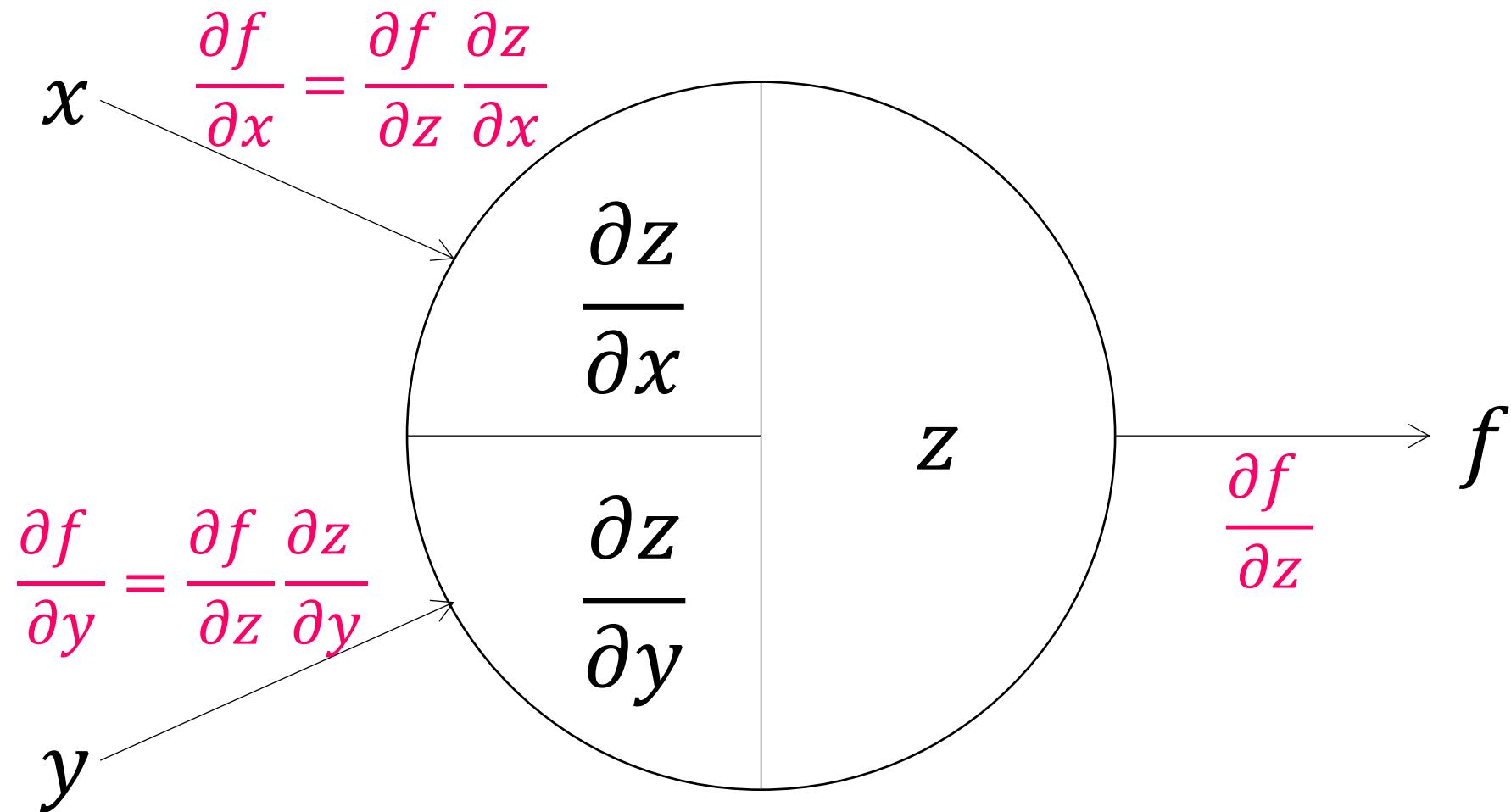
# Backpropagation

- A simple example:  $f(x, y, z) = (x + y)z$

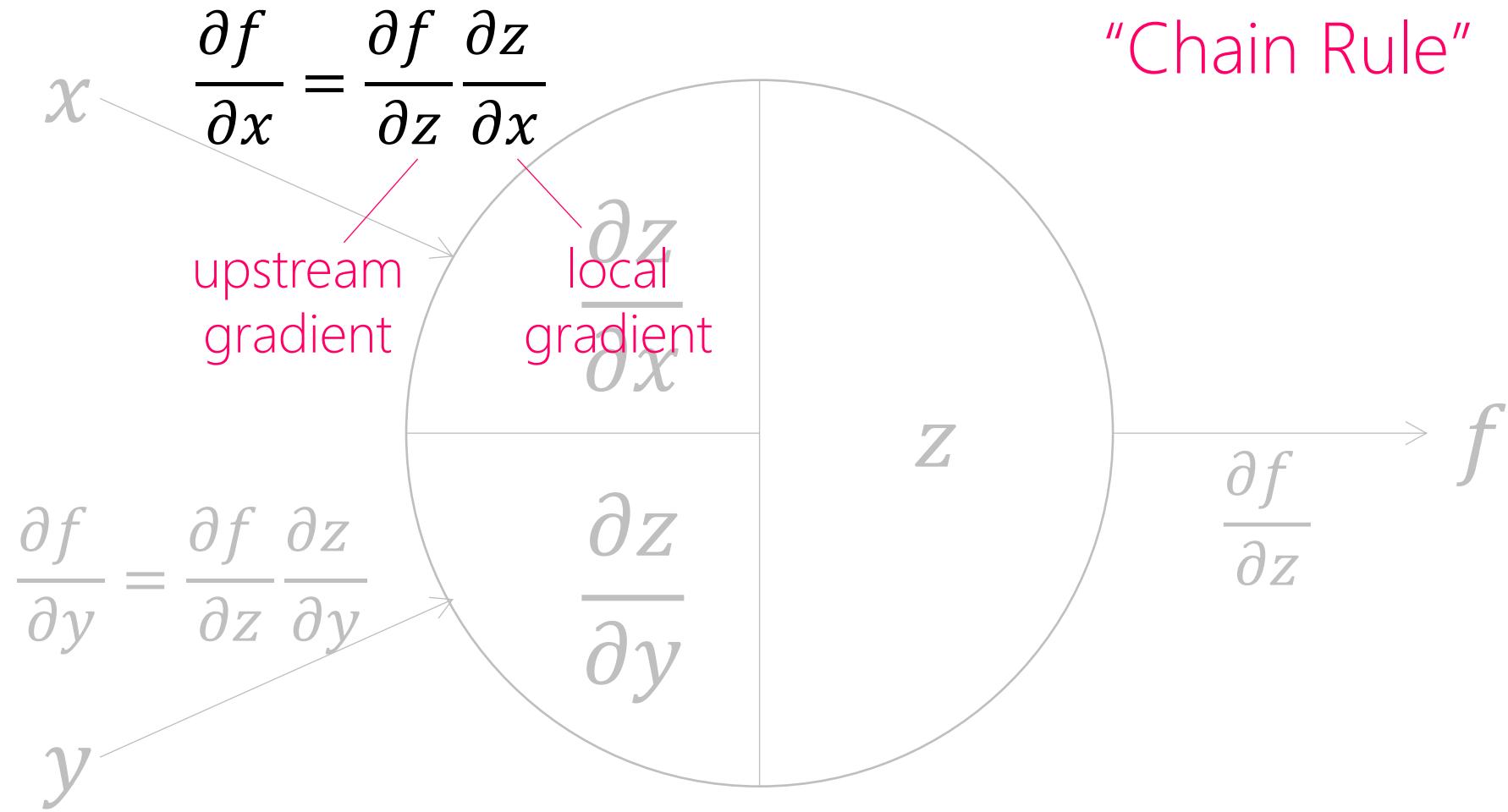
- Computational graph:



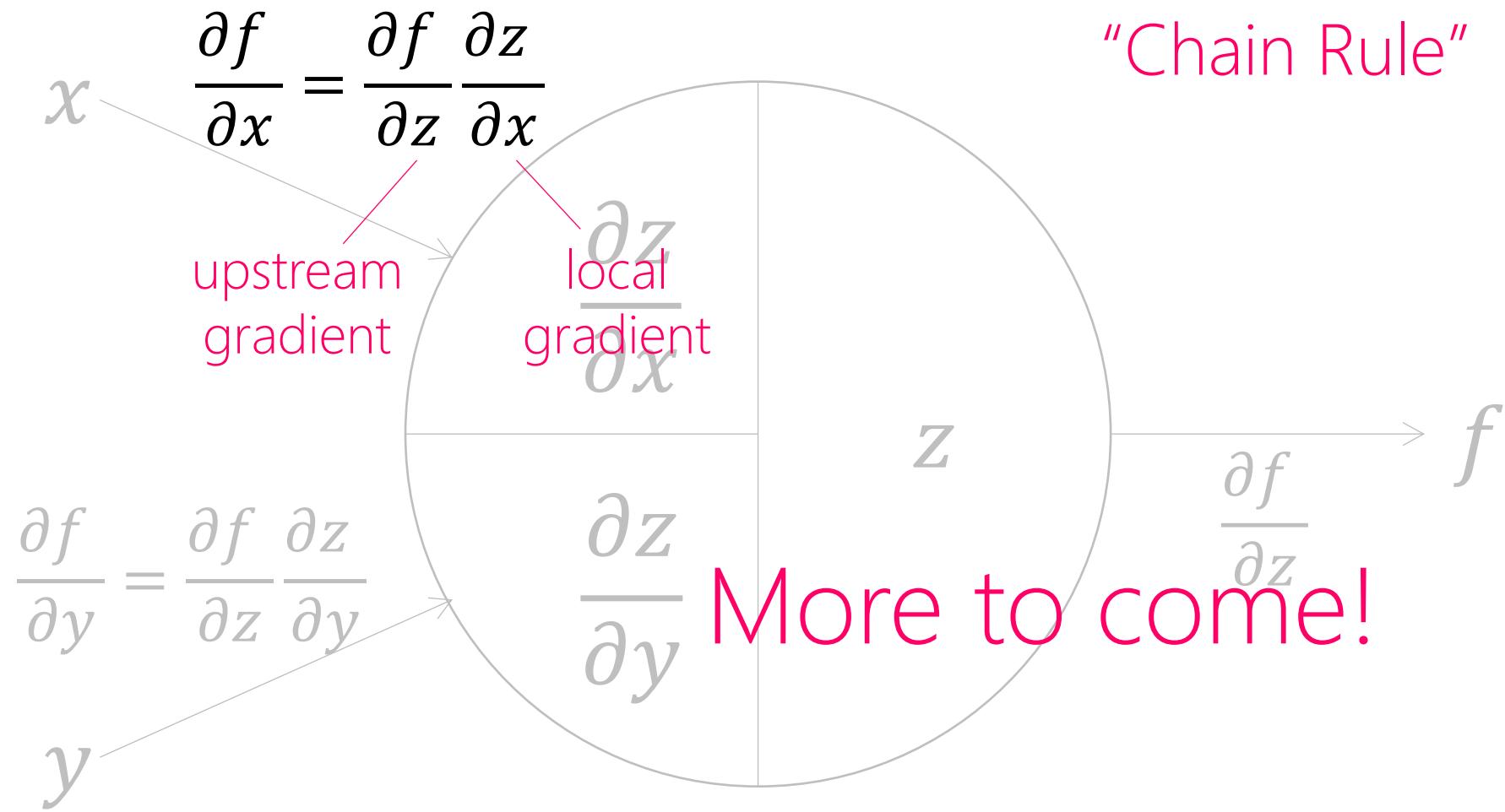
# Backpropagation



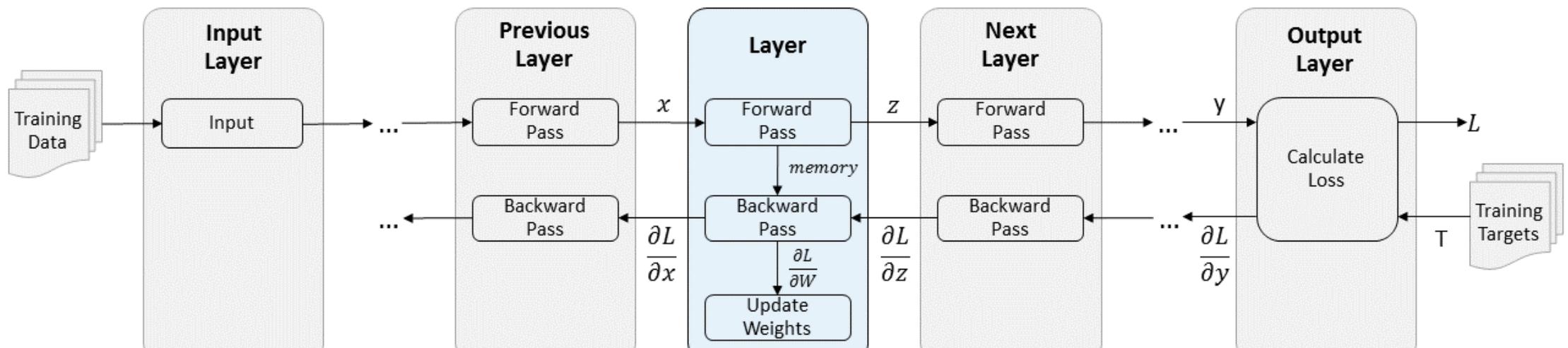
# Backpropagation



# Backpropagation



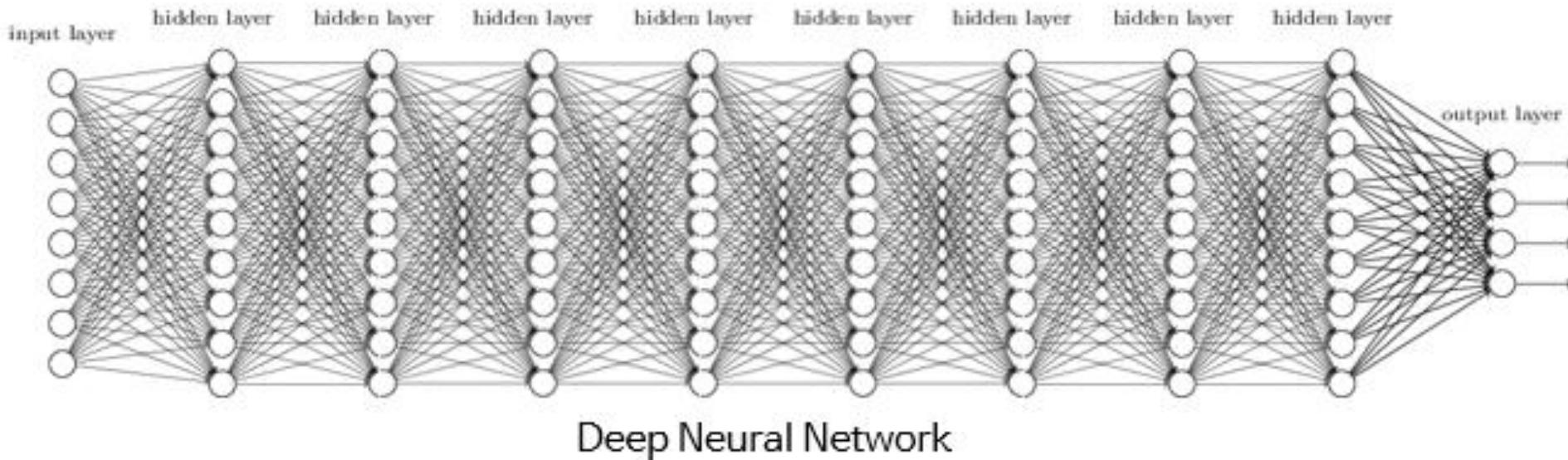
# Putting them all together



<https://www.mathworks.com/help/nnet/ug/define-custom-deep-learning-layers.html>

# Putting them all together

- Linear model:  $f = Wx + b$
- Neural network with a single hidden layer:  $f = W_2 \max(0, W_1 x + b_1) + b_2$
- Neural network with two hidden layers:  
$$f = W_3 \max(0, W_2 \max(0, W_1 x + b_1) + b_2) + b_3$$
- ... stacking up hidden layers → “deep” neural networks





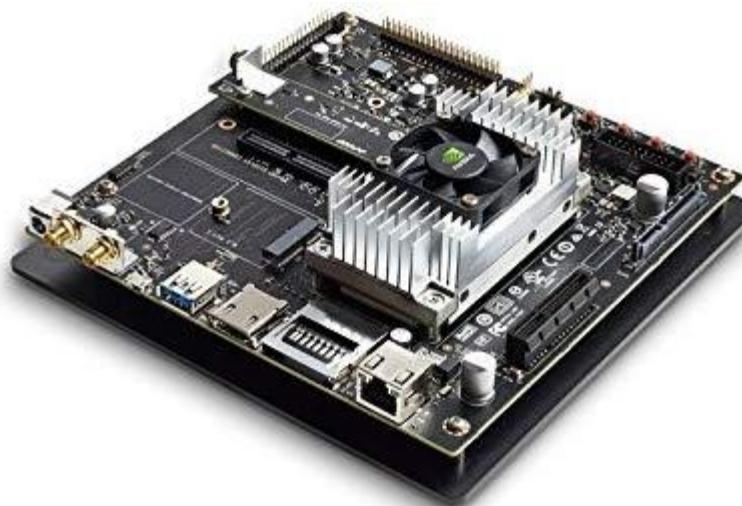
5MB Hard Drive Being Shipped by IBM, 1956

A FEW  
MOMENTS LATER

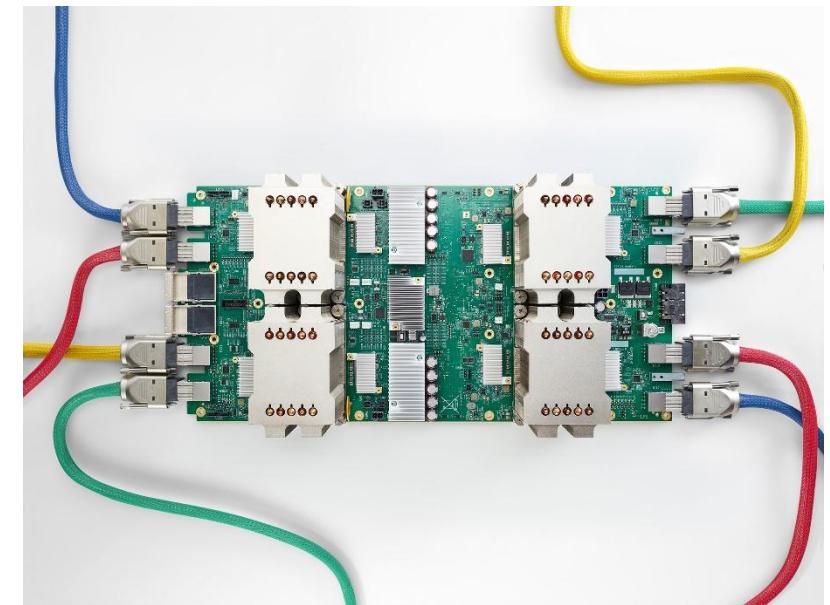
# 40 years later...



NVIDIA GeForce Quadro GV100



NVIDIA Jetson TX2



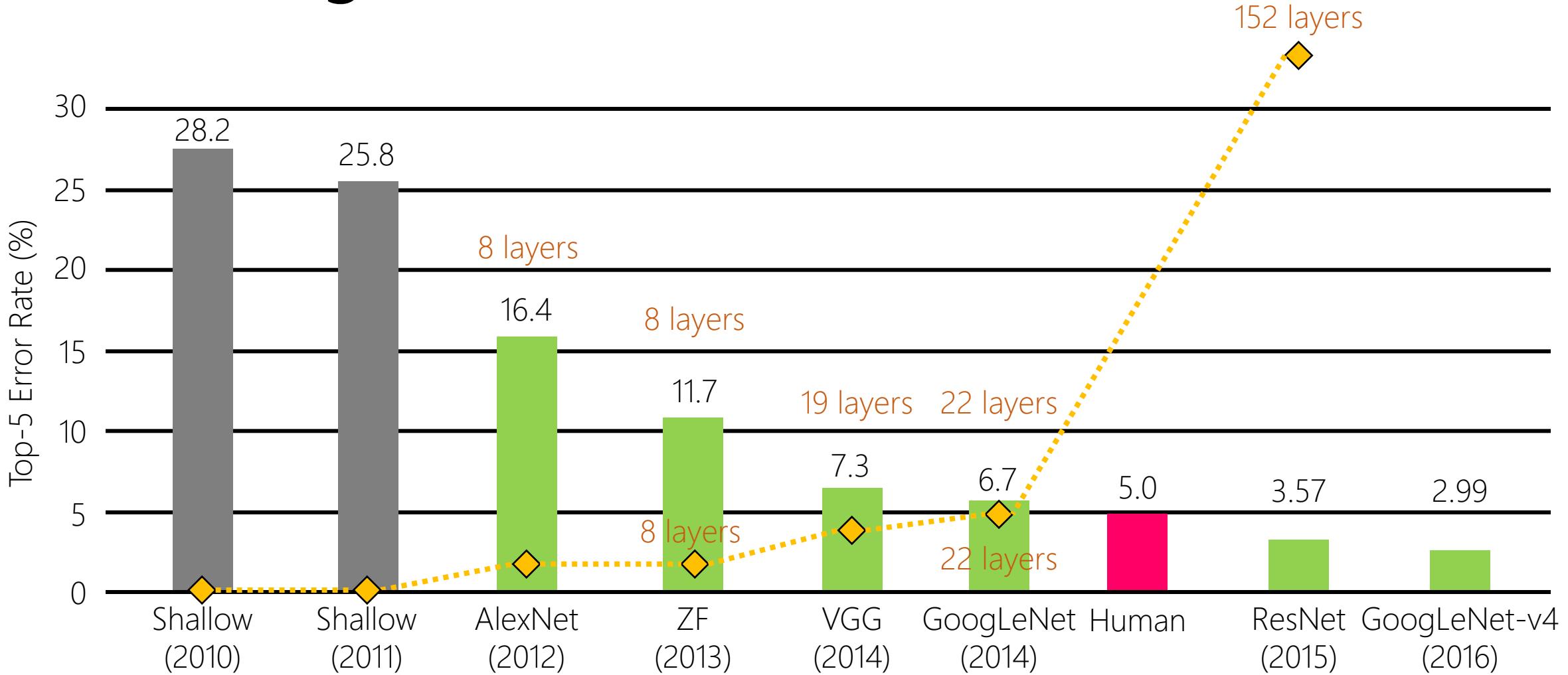
Google Tensor Processing Unit (TPU)



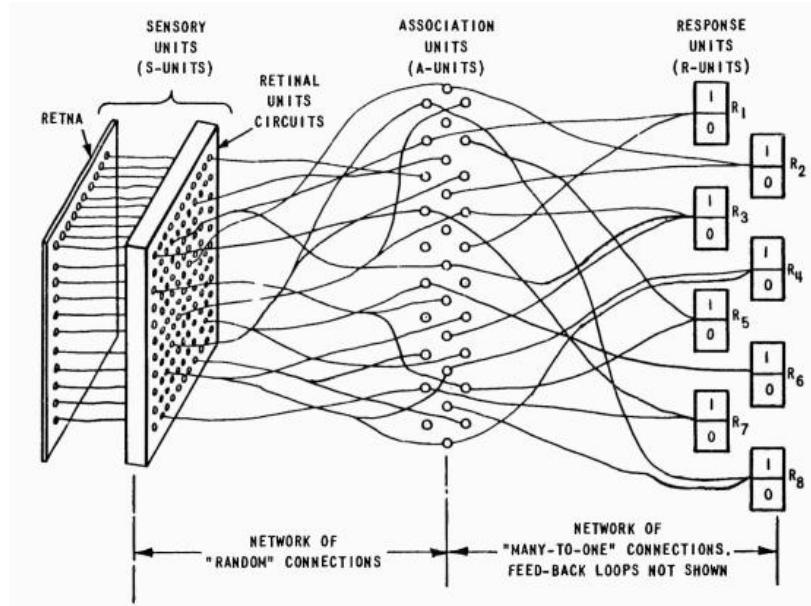
# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

- ImageNet (First appearance as a poster at CVPR 2009)
  - Over 14 million images with hand-annotated labels (crowdsourced via M-turk)
  - Over 20 thousand categories
  - 1M+ images with bounding boxes.
- ImageNet Challenge
  - Since 2010
  - Uses a trimmed set of thousand non-overlapping classes.
  - Humans error is about 5%\*
  - Playground of computer vision researchers.
  - Marked the start of the current AI boom.

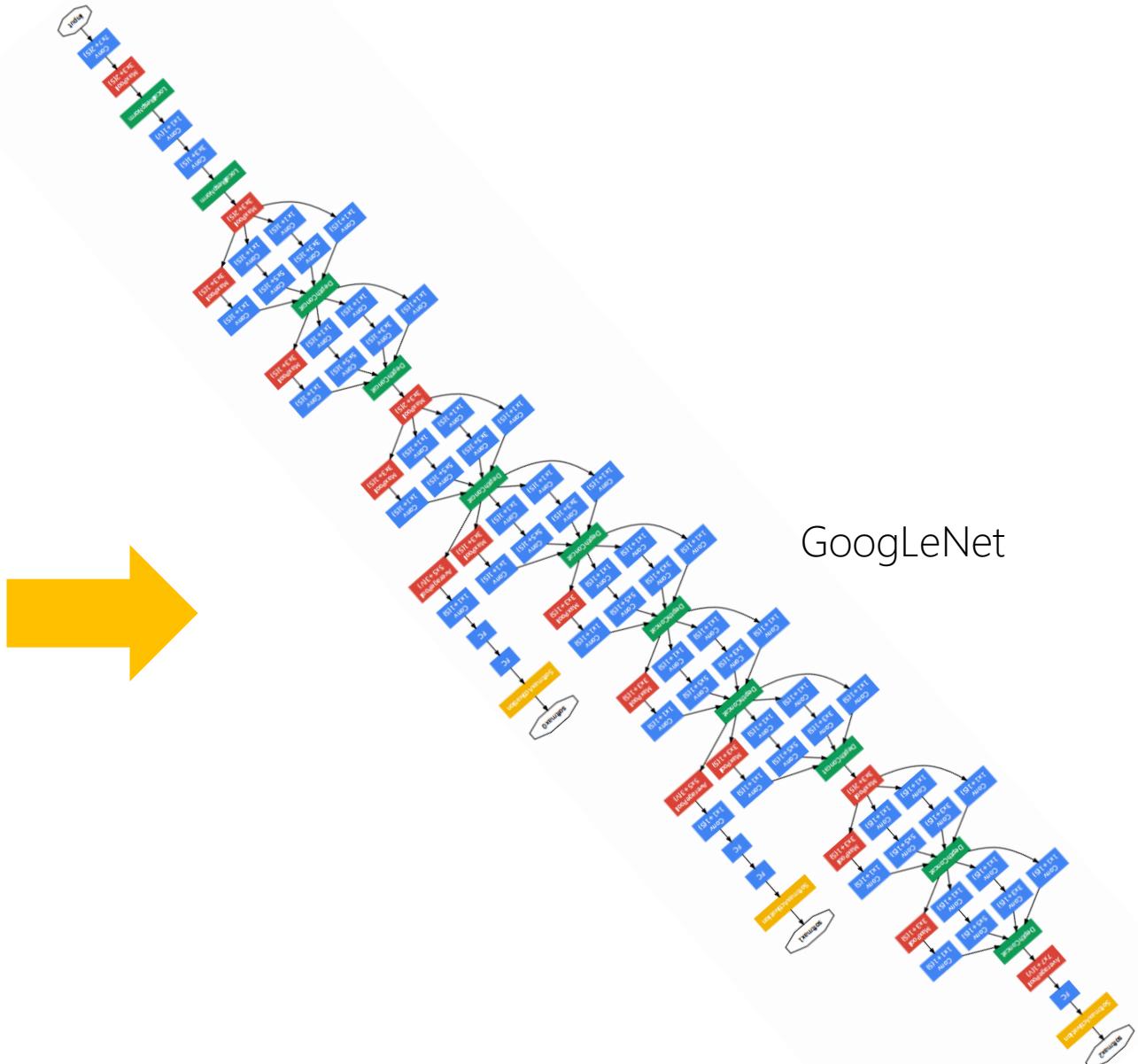
# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

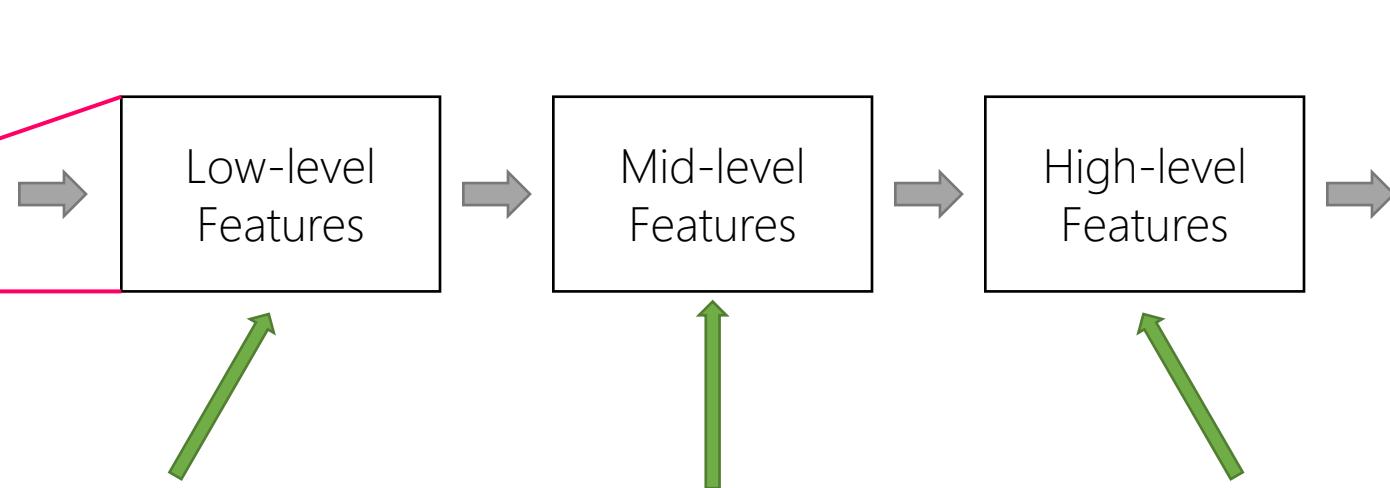
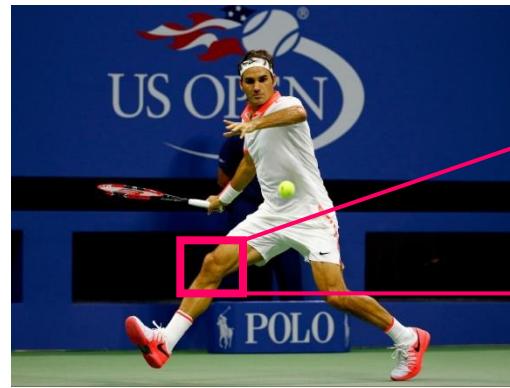


# Going Deeper



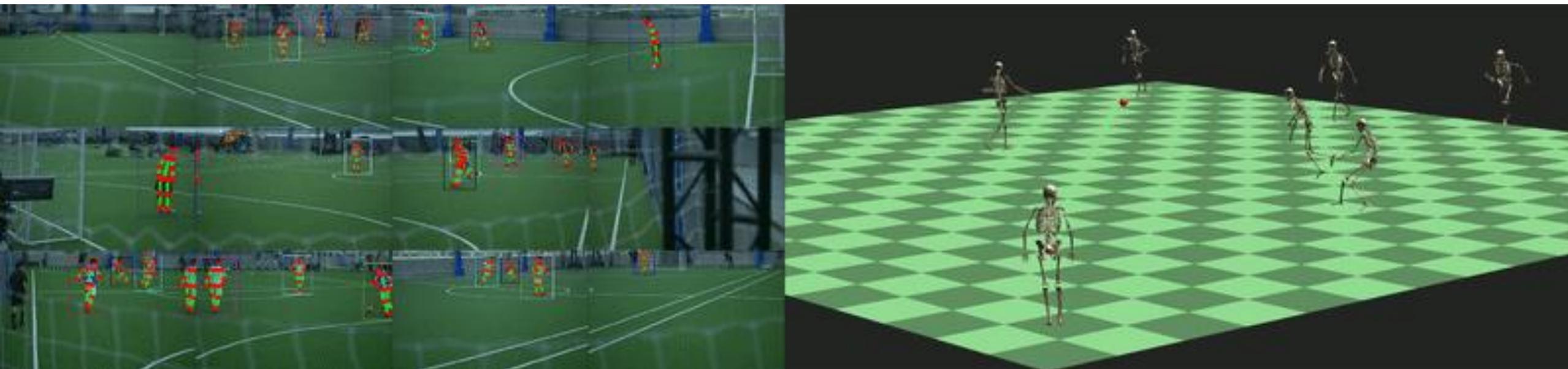
Perceptron



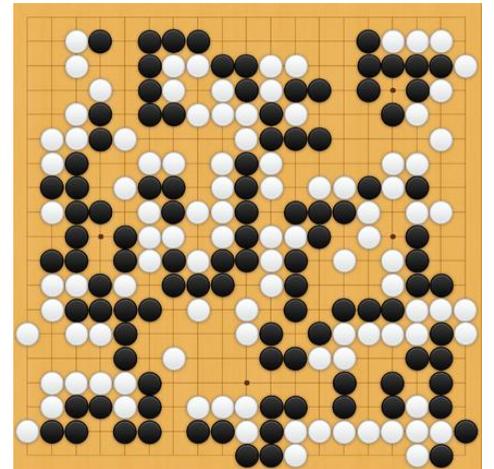
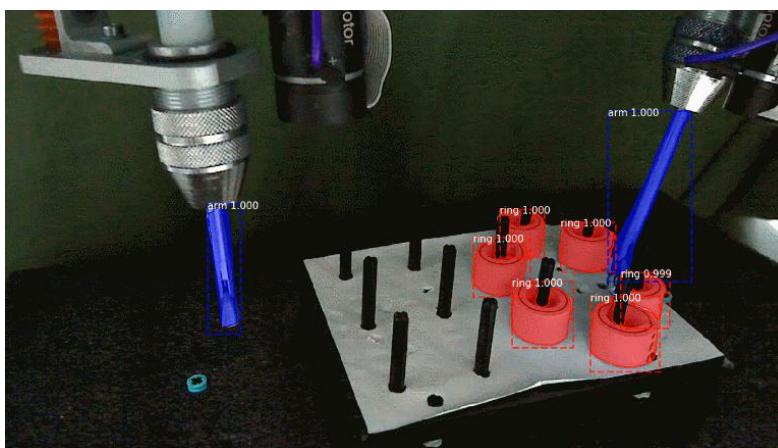
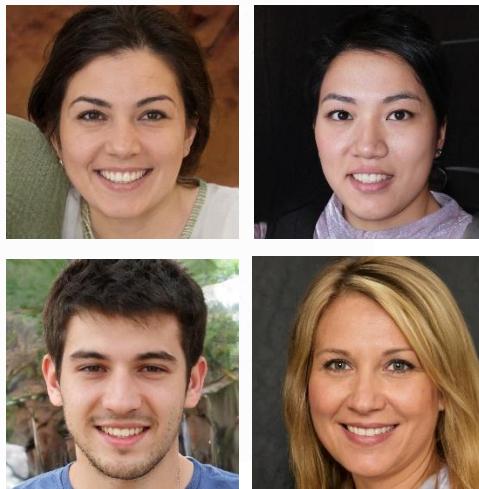
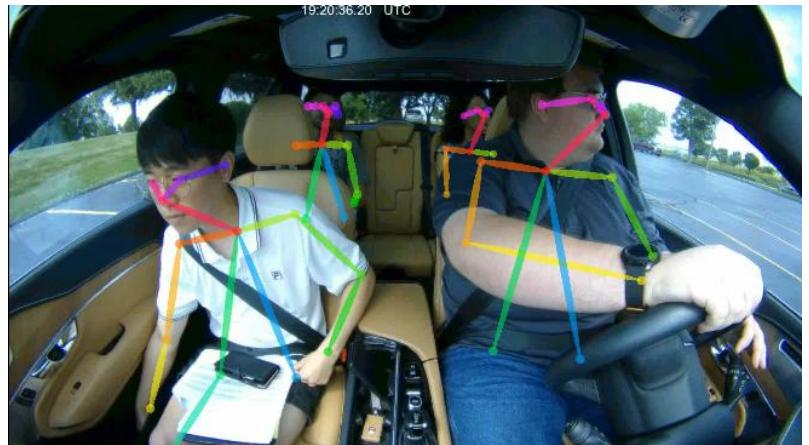
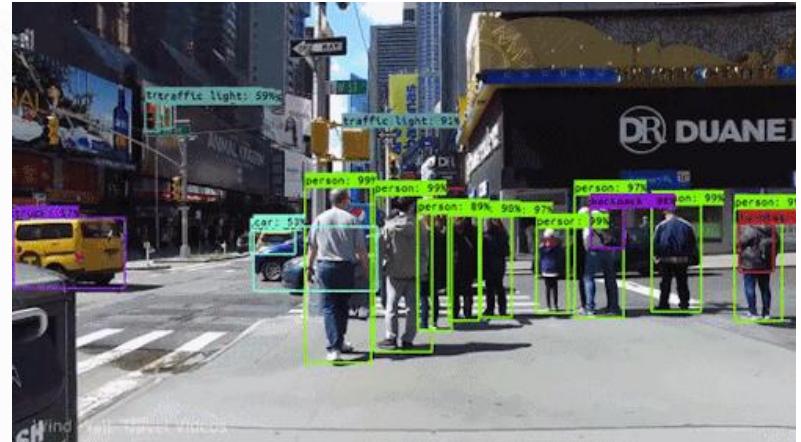


Shoulder:	0.01
Elbow:	0.12
Wrist:	0.01
<b>Knee:</b>	<b>0.74</b>
Ankle:	0.03
...	

# e.g. Pose Estimation

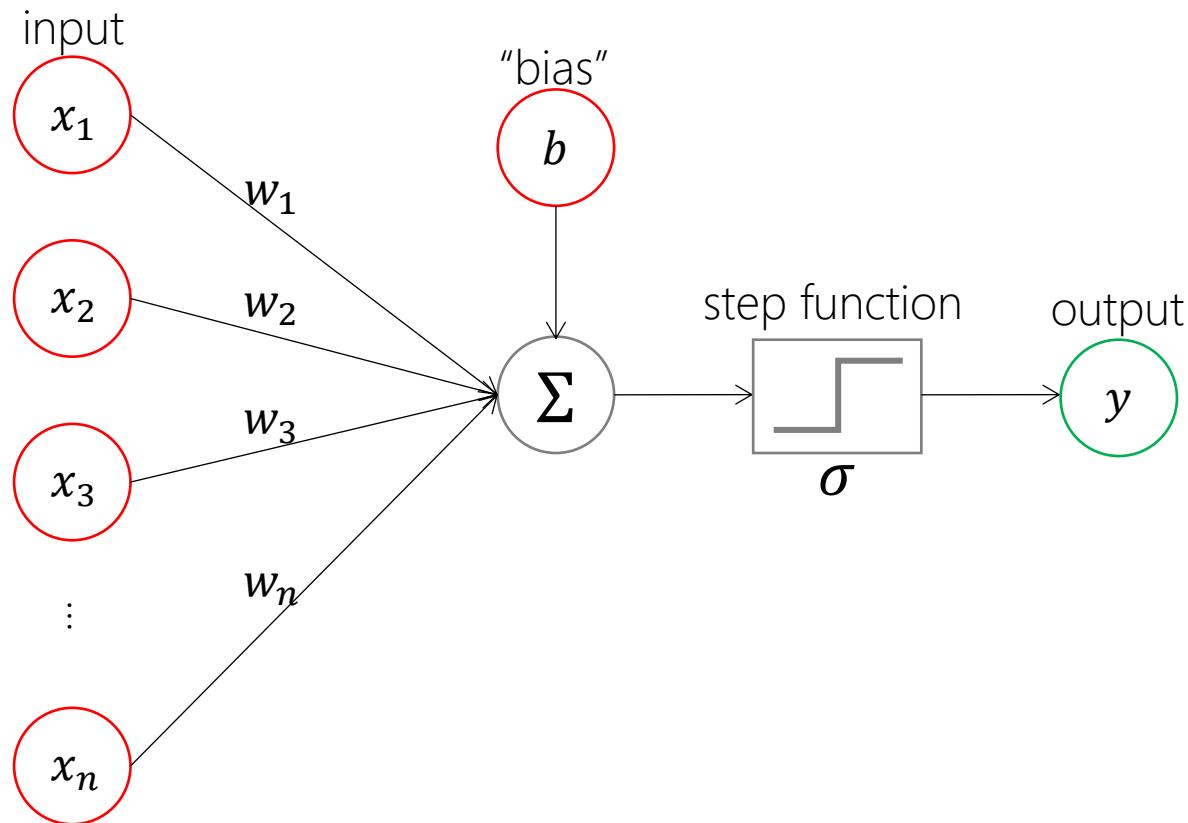


...and many others



# Closer Look into Neural Nets

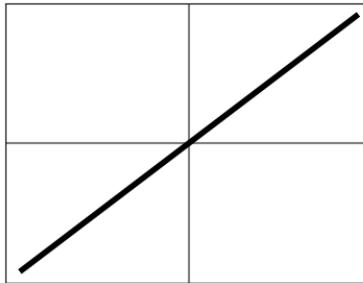
- Back to perceptron...



$$z = \sum_j w_j x_j + b$$
$$y = \sigma(z)$$

Linear (No Activation)

$$\sigma(z) = z$$



Leaky ReLU

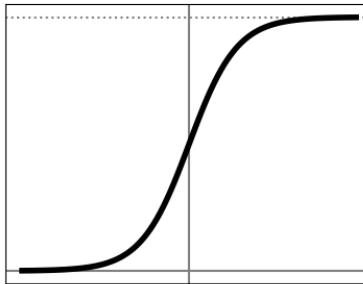
$$\sigma(z) = \max(0.1z, z)$$

Parametric ReLU (PReLU)

$$\sigma(z) = \max(\alpha z, z)$$

Sigmoid

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

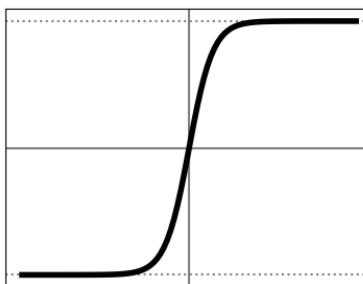


Exponential ReLU (ELU)

$$\sigma(z) = \begin{cases} z, & z \geq 0 \\ \alpha(e^z - 1), & z < 0 \end{cases}$$

Hyperbolic Tangent

$$\sigma(z) = \tanh(z)$$



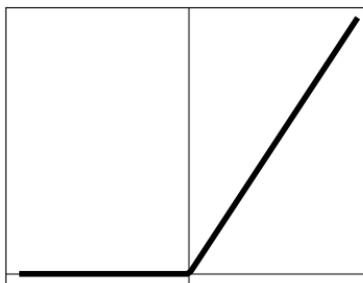
Swish

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

$$= z * \text{sigmoid}(z)$$

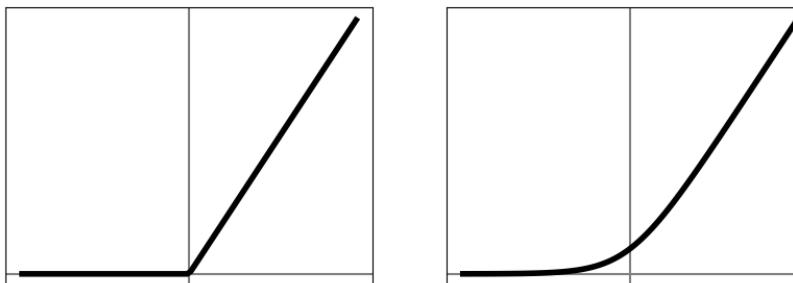
Rectified Linear Unit  
(ReLU)

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



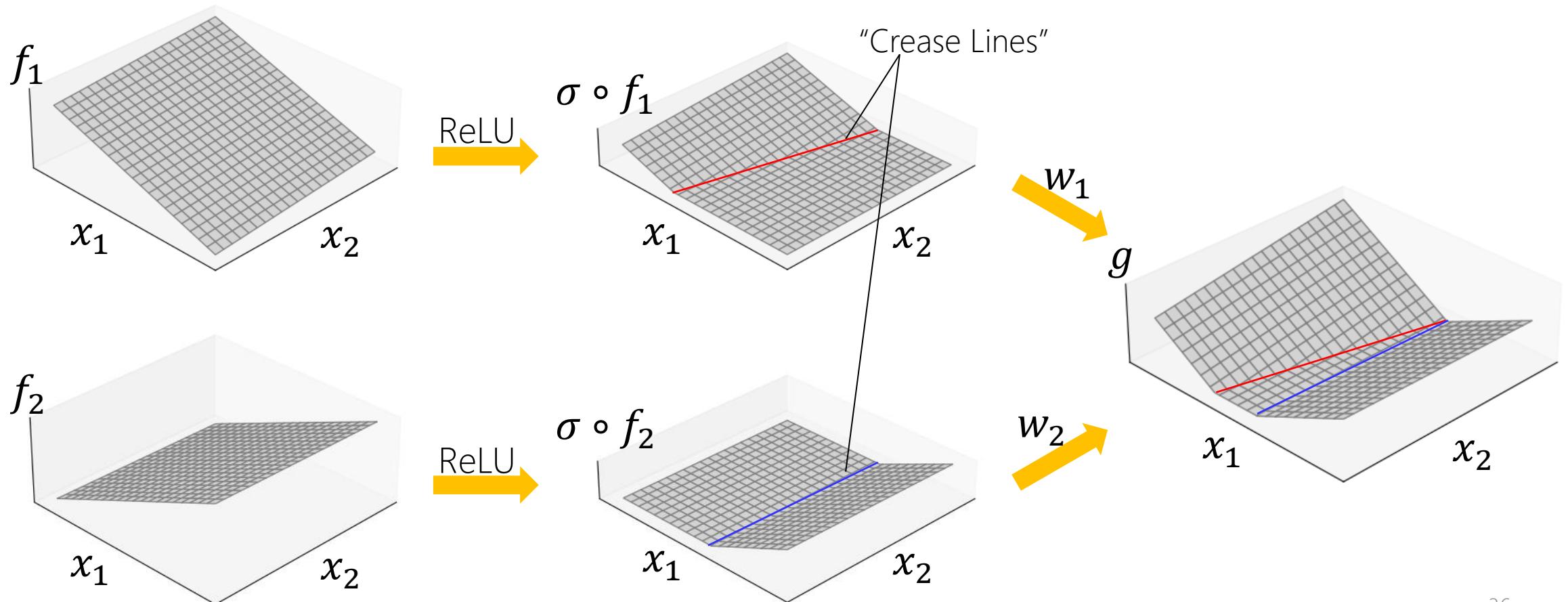
SoftPlus

$$\sigma(z) = \log(1 + e^z)$$

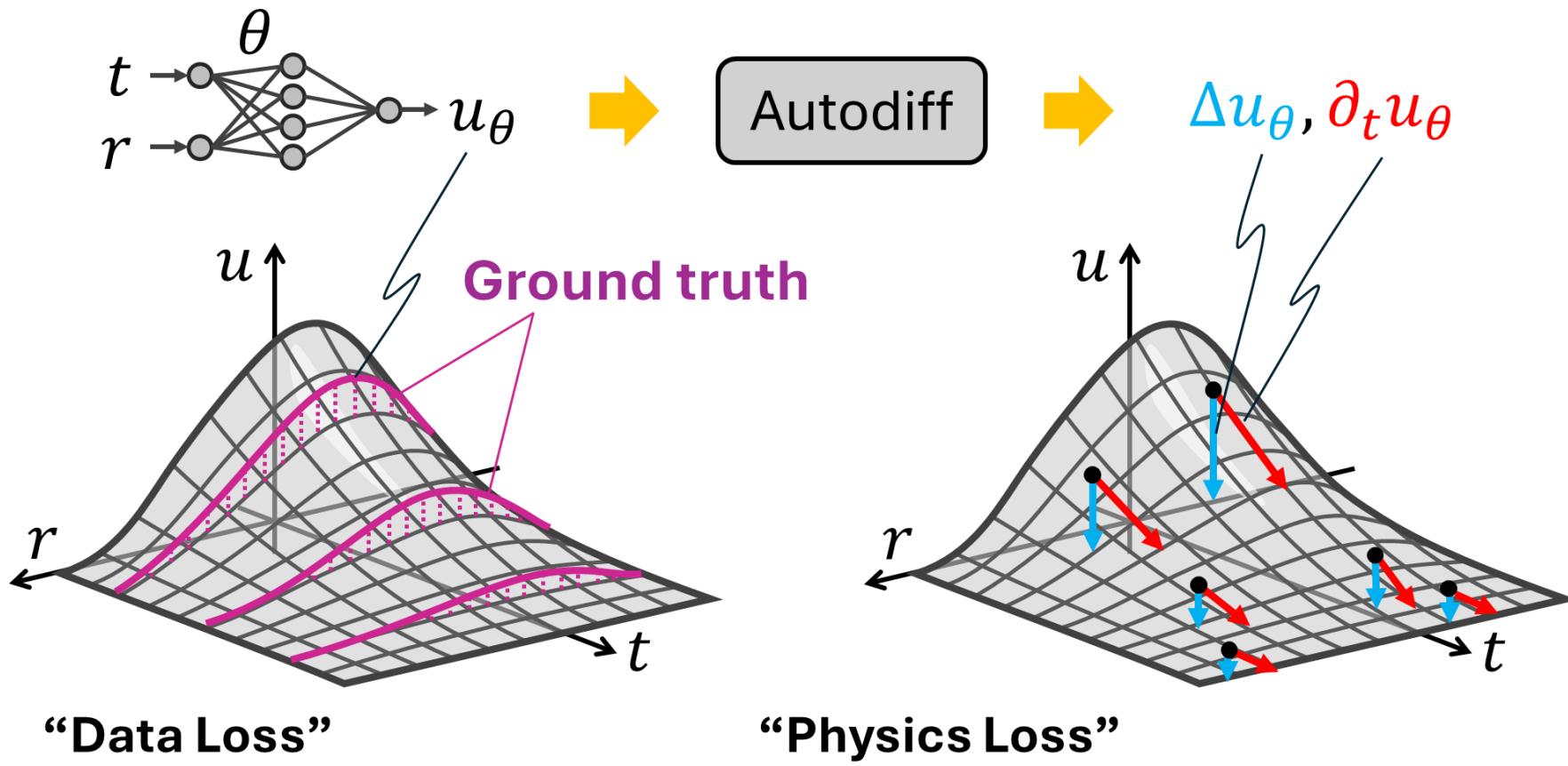


# Closer Look into Neural Nets

- Geometric view of perceptron?



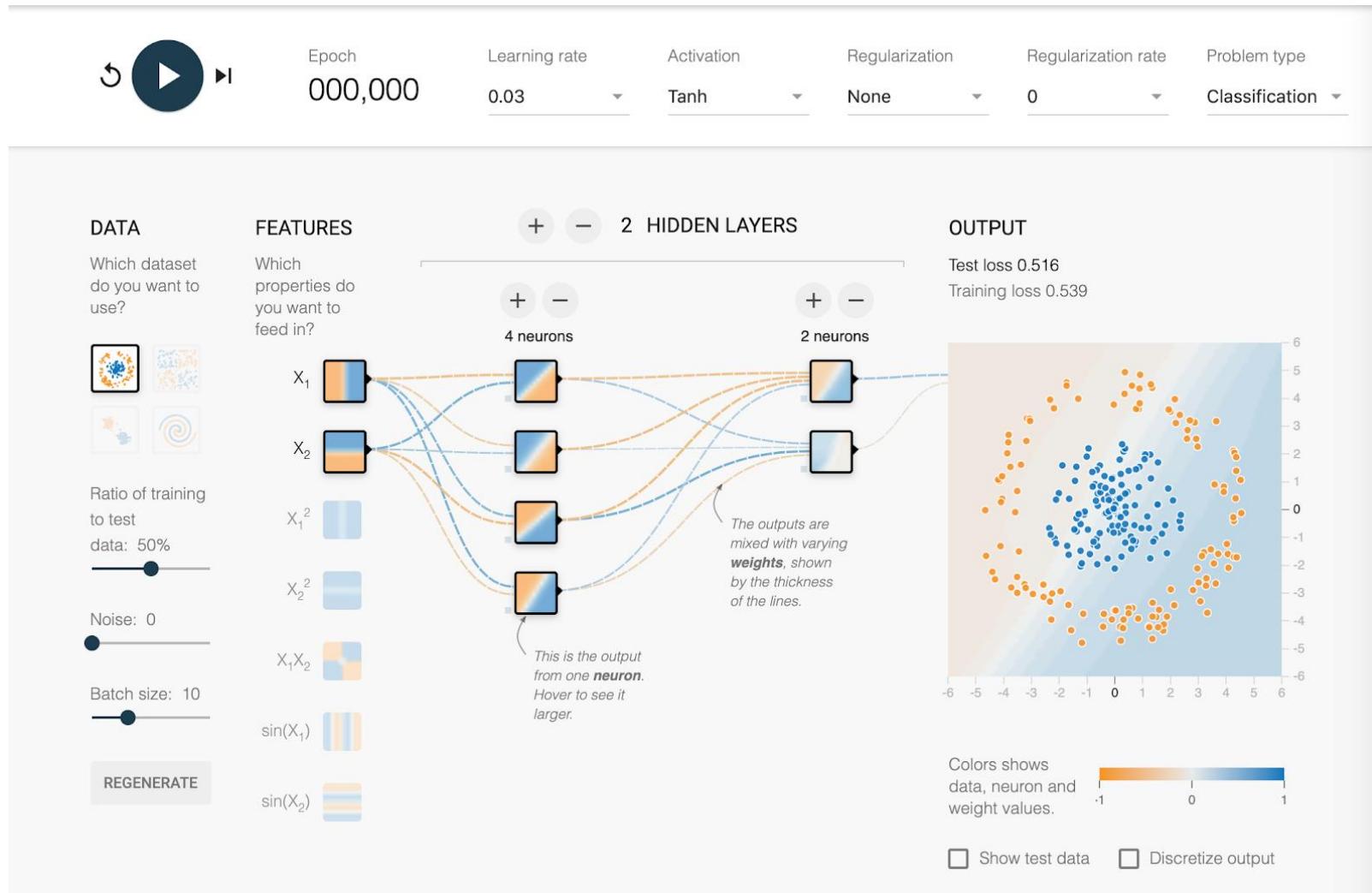
# Peek on Physics-Informed Neural Nets



# Universal Approximation Theorem (UAT)

- In a nutshell: "

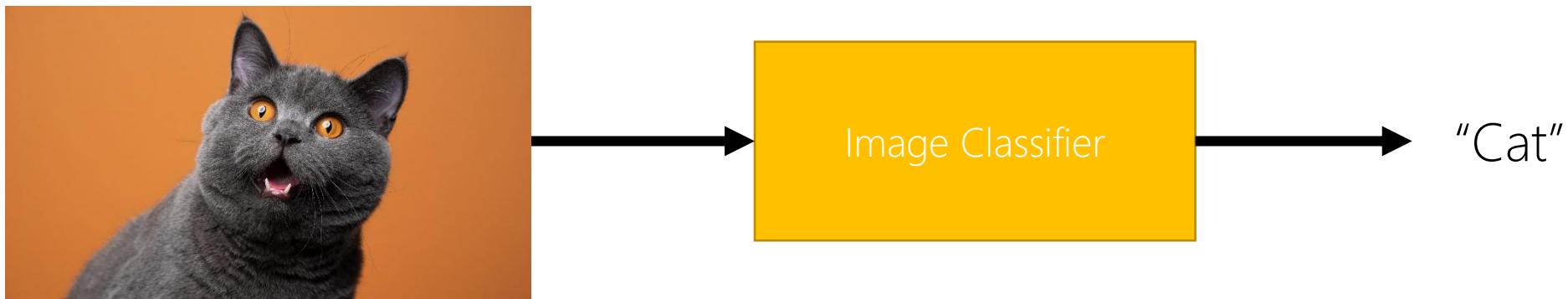
# TensorFlow Playground



# Formulating NN Problems

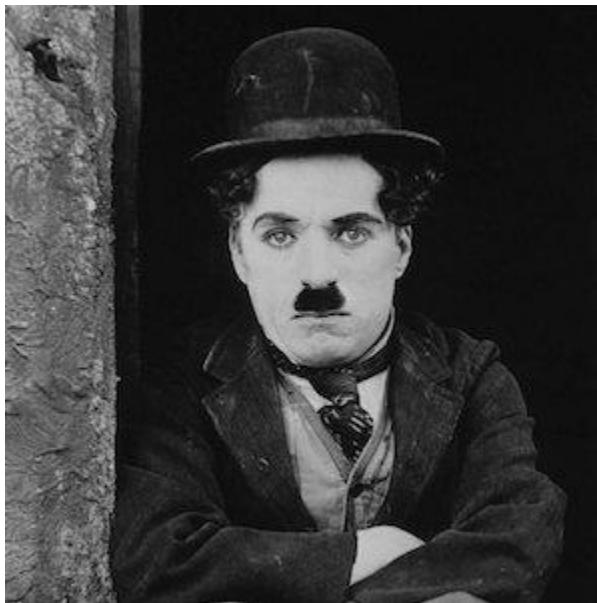
# Image Classification

- Images are more straightforward and intuitive subject to study deep networks
- Easy to generalize to other forms of data
- Major breakthroughs in deep learning research were in computer vision
- Image classification = core task in computer vision



# Image Classification is Challenging!

- Problem 1. Semantic Challenge



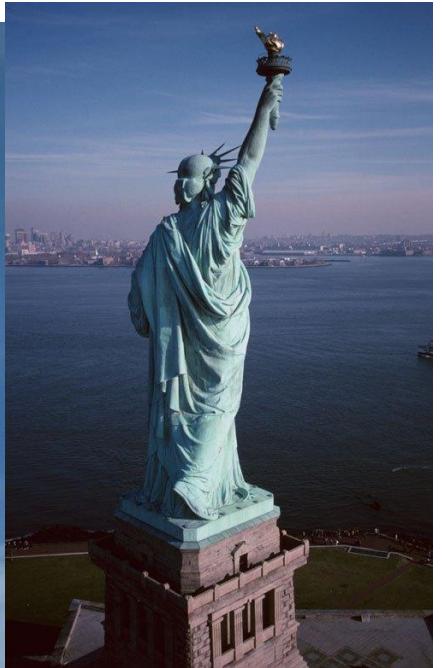
What we see

127	123	101	40	12	13	14	22	26	18	13	13	13	13	13	15
114	81	98	38	13	14	28	36	34	29	17	13	13	13	13	15
142	102	87	33	13	15	31	38	34	36	22	13	13	14	15	17
132	112	87	24	14	22	34	52	46	37	26	20	14	15	16	18
120	103	83	19	15	18	57	158	161	131	24	19	14	15	15	17
117	110	84	18	14	21	114	163	171	128	49	17	15	15	15	17
117	112	89	21	14	26	155	163	174	157	117	15	16	16	16	17
117	103	84	21	15	13	120	174	106	182	87	11	15	15	16	17
117	103	86	21	16	16	56	158	125	156	55	19	21	16	16	17
112	102	94	23	17	29	33	86	97	80	44	38	40	26	15	16
107	91	95	25	19	33	39	98	77	105	65	42	35	40	17	16
98	95	97	26	25	29	39	101	108	70	70	39	29	46	20	16
111	92	91	34	34	23	27	64	144	89	54	23	30	42	31	16
113	96	90	36	32	28	21	31	92	83	47	21	35	42	42	17
119	109	87	34	31	33	45	51	45	134	128	36	33	40	35	17
98	90	78	29	28	34	28	29	39	70	51	65	64	35	23	16

What computers see  
= just a grid of numbers (e.g. 1024 x 768 x 3)

# Image Classification is Challenging!

- Problem 2. Varying Perspectives



# Image Classification is Challenging!

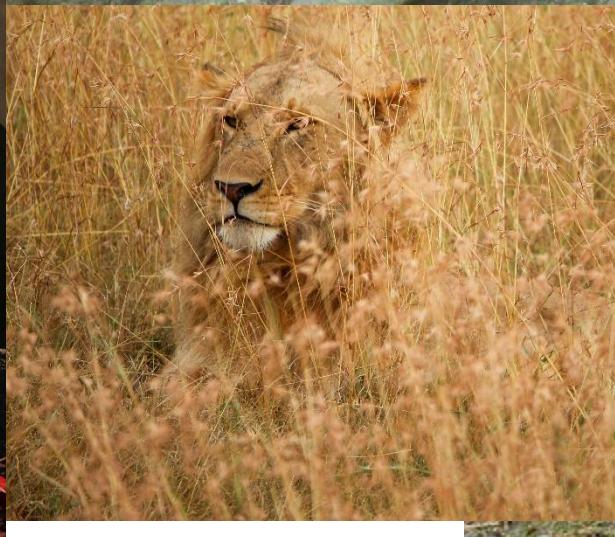
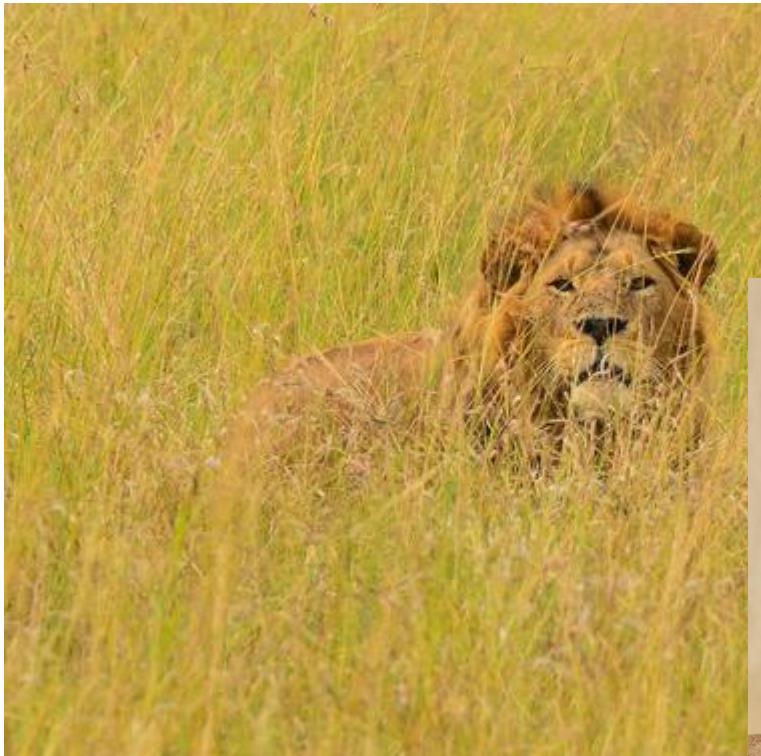
- Problem 3. Illumination



<http://www.ifp.illinois.edu/~tchen5/vlfr.html>

# Image Classification is Challenging!

- Problem 4. Occlusion



# Image Classification is Challenging!

- Problem 5. Deformation, Shape & Posture Change



The same guy in front of you in 2011

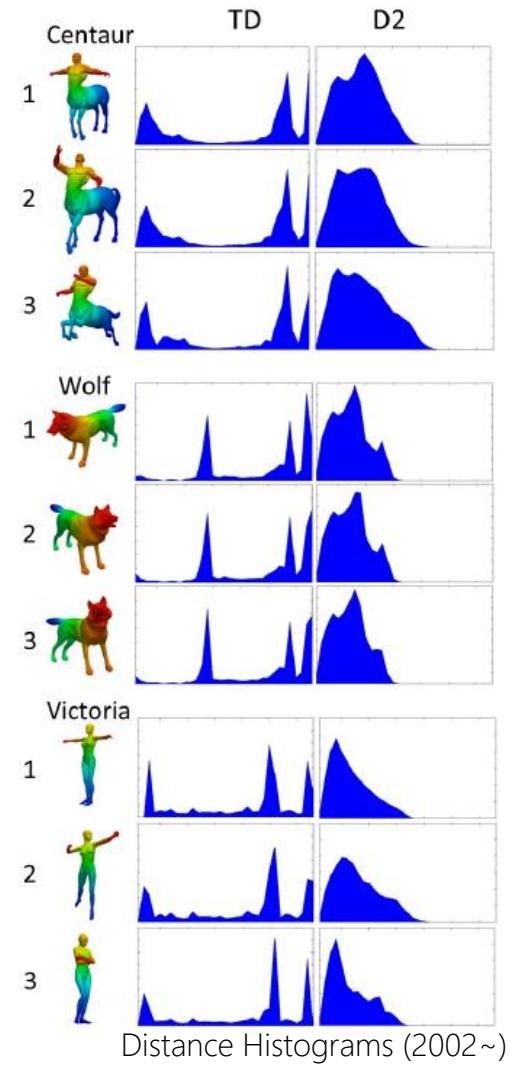
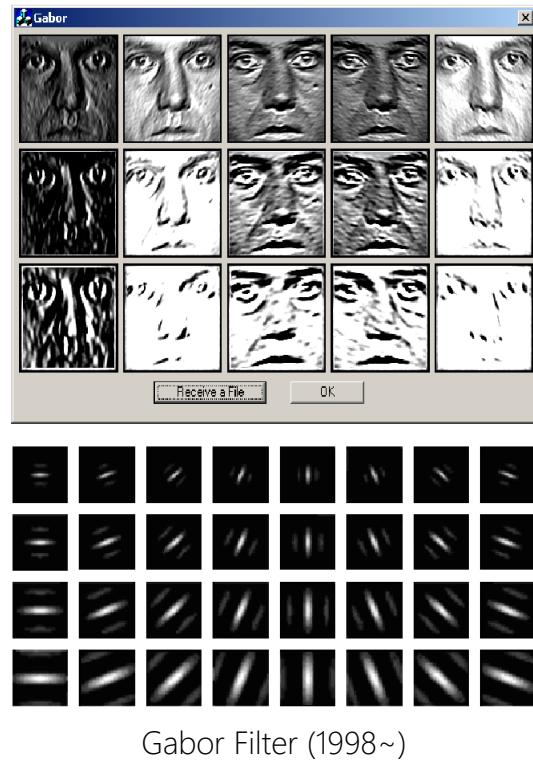
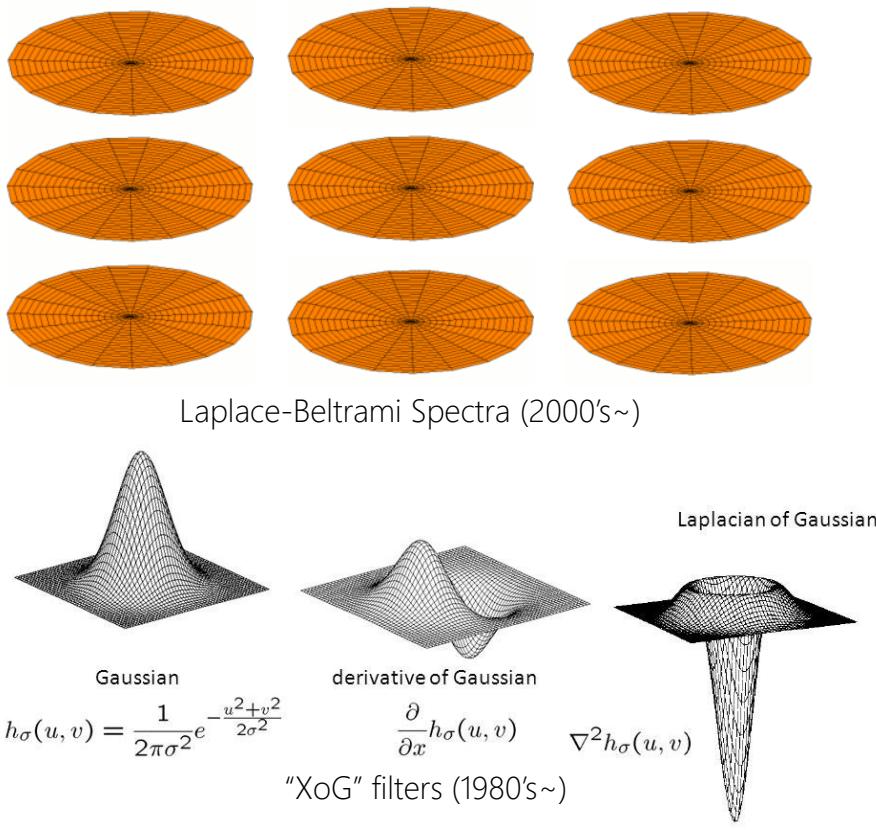
# Image Classification is Challenging!

- Problem 6. Intraclass Variation

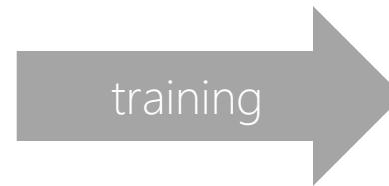
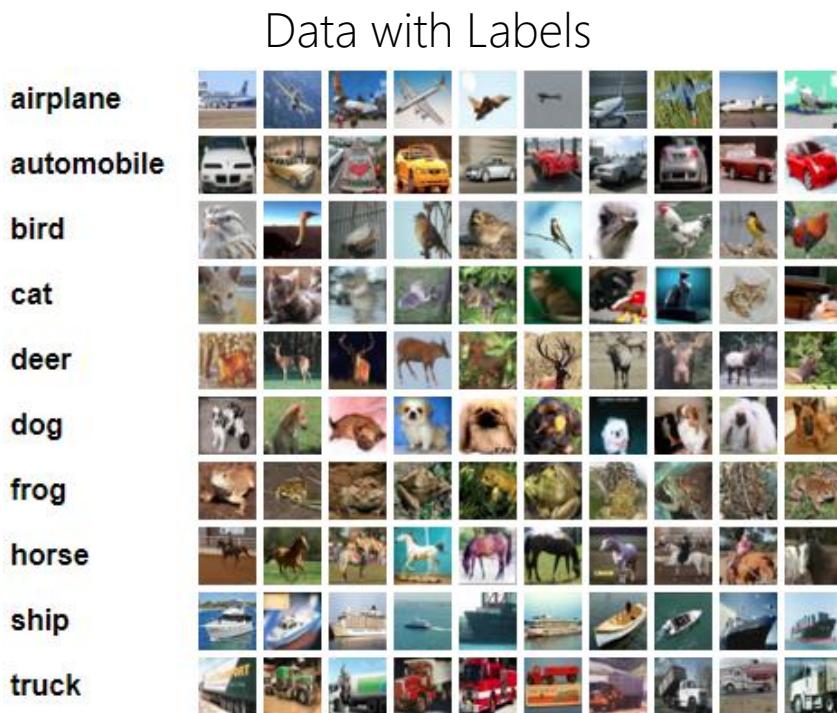


Arman, Alarm Clocks (Réveils), 1960, MCA Chicago

# Conventional Approach (Model-Driven)



# Machine Learning (Data-Driven)



# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.

# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.
- Our model " $f$ " can be written formally as  $f(\mathbf{x} | \boldsymbol{\theta}) := \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\boldsymbol{\theta} := (\mathbf{W}, \mathbf{b})$ .

# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.
- Our model " $f$ " can be written formally as  $f(\mathbf{x} | \boldsymbol{\theta}) := \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\boldsymbol{\theta} := (\mathbf{W}, \mathbf{b})$ .
- Goal: Find  $\boldsymbol{\theta}$  ( $\mathbf{W}$  and  $\mathbf{b}$ ) that minimize the "*difference*" between the model prediction  $f(\mathbf{x} | \boldsymbol{\theta})$  and the ground truth  $y$

# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.
- Our model " $f$ " can be written formally as  $f(\mathbf{x} | \boldsymbol{\theta}) := \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\boldsymbol{\theta} := (\mathbf{W}, \mathbf{b})$ .
- Goal: Find  $\boldsymbol{\theta}$  ( $\mathbf{W}$  and  $\mathbf{b}$ ) that minimize the "*difference*" between the model prediction  $f(\mathbf{x} | \boldsymbol{\theta})$  and the ground truth  $y$ 
  - In English: We want the predicted label of our model  $f(\mathbf{x} | \boldsymbol{\theta})$  to be as close as possible to the true label  $y$  and we're going to accomplish that by carefully "tuning"  $\mathbf{W}$  and  $\mathbf{b}$ .

# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.
- Our model " $f$ " can be written formally as  $f(\mathbf{x} | \boldsymbol{\theta}) := \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\boldsymbol{\theta} := (\mathbf{W}, \mathbf{b})$ .
- Goal: Find  $\boldsymbol{\theta}$  ( $\mathbf{W}$  and  $\mathbf{b}$ ) that minimize the "*difference*" between the model prediction  $f(\mathbf{x} | \boldsymbol{\theta})$  and the ground truth  $y$ 
  - In English: We want the predicted label of our model  $f(\mathbf{x} | \boldsymbol{\theta})$  to be as close as possible to the true label  $y$  and we're going to accomplish that by carefully "tuning"  $\mathbf{W}$  and  $\mathbf{b}$ .
  - Based on what, though?

# Image Classification with Perceptron

- Idea: Formulate the image classification problem as  $y = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\mathbf{x}$  is an image and  $y$  is the image label.
- Our model " $f$ " can be written formally as  $f(\mathbf{x} | \boldsymbol{\theta}) := \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$  where  $\boldsymbol{\theta} := (\mathbf{W}, \mathbf{b})$ .
- Goal: Find  $\boldsymbol{\theta}$  ( $\mathbf{W}$  and  $\mathbf{b}$ ) that minimize the "*difference*" between the model prediction  $f(\mathbf{x} | \boldsymbol{\theta})$  and the ground truth  $y$ 
  - In English: We want the predicted label of our model  $f(\mathbf{x} | \boldsymbol{\theta})$  to be as close as possible to the true label  $y$  and we're going to accomplish that by carefully "tuning"  $\mathbf{W}$  and  $\mathbf{b}$ .
  - Based on what, though? Lots of labeled data!  $(\mathbf{x}_i, y_i)_{i=1,\dots,N}$

# Image Classifier: Step-by-Step

1. Collect a bunch of data  $(\mathbf{x}_i, y_i)_{i=1,\dots,N}$
2. Split them into two sets: training set and test set
3. Solve  $\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^{N_{\text{train}}} d(y_i, f(\mathbf{x}_i | \boldsymbol{\theta}))$
4. See if  $d(y, f(\mathbf{x} | \boldsymbol{\theta}))$  is “close enough to zero” on the test set
  - If it does: profit!
  - If it doesn’t: sit and cry (or build a better model, collect more data, etc.: Next class!)

# Example: Cat/Dog Classifier

- The cats and dogs dataset: <https://www.microsoft.com/en-us/download/details.aspx?id=54765>
- 12,500 cat images & 12,500 dog images from internet



# Example: Cat/Dog Classifier

$$(\text{cat image}, -1)$$
$$(\text{cat image}, -1)$$
$$(\text{cat image}, -1)$$

...

$$(\text{dog image}, 1)$$
$$(\text{puppies image}, 1)$$
$$(\text{two dogs image}, 1)$$

...

# Example: Cat/Dog Classifier



width = 400

height = 500

channels = 3

# Example: Cat/Dog Classifier



gray + resize



64 x 64

width = 400  
height = 500  
channels = 3

# Example: Cat/Dog Classifier



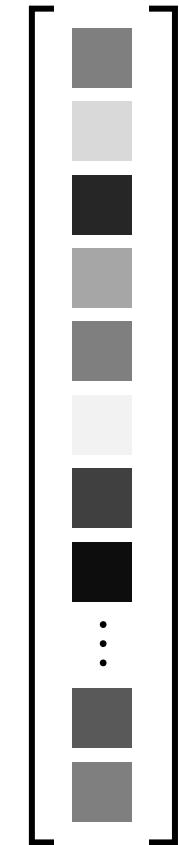
width = 400  
height = 500  
channels = 3

gray + resize



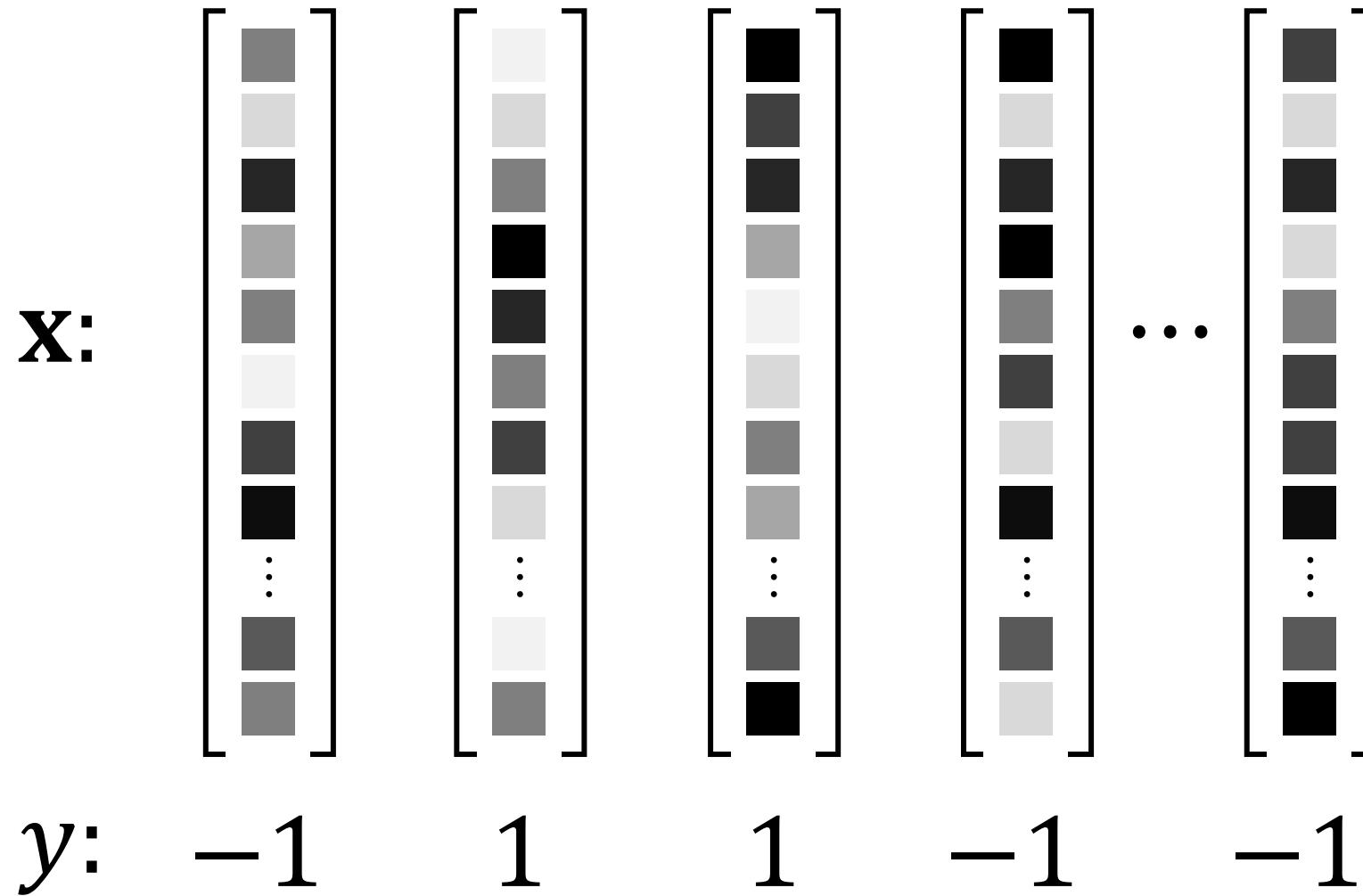
64 x 64

flatten



4096 x 1

# Example: Cat/Dog Classifier



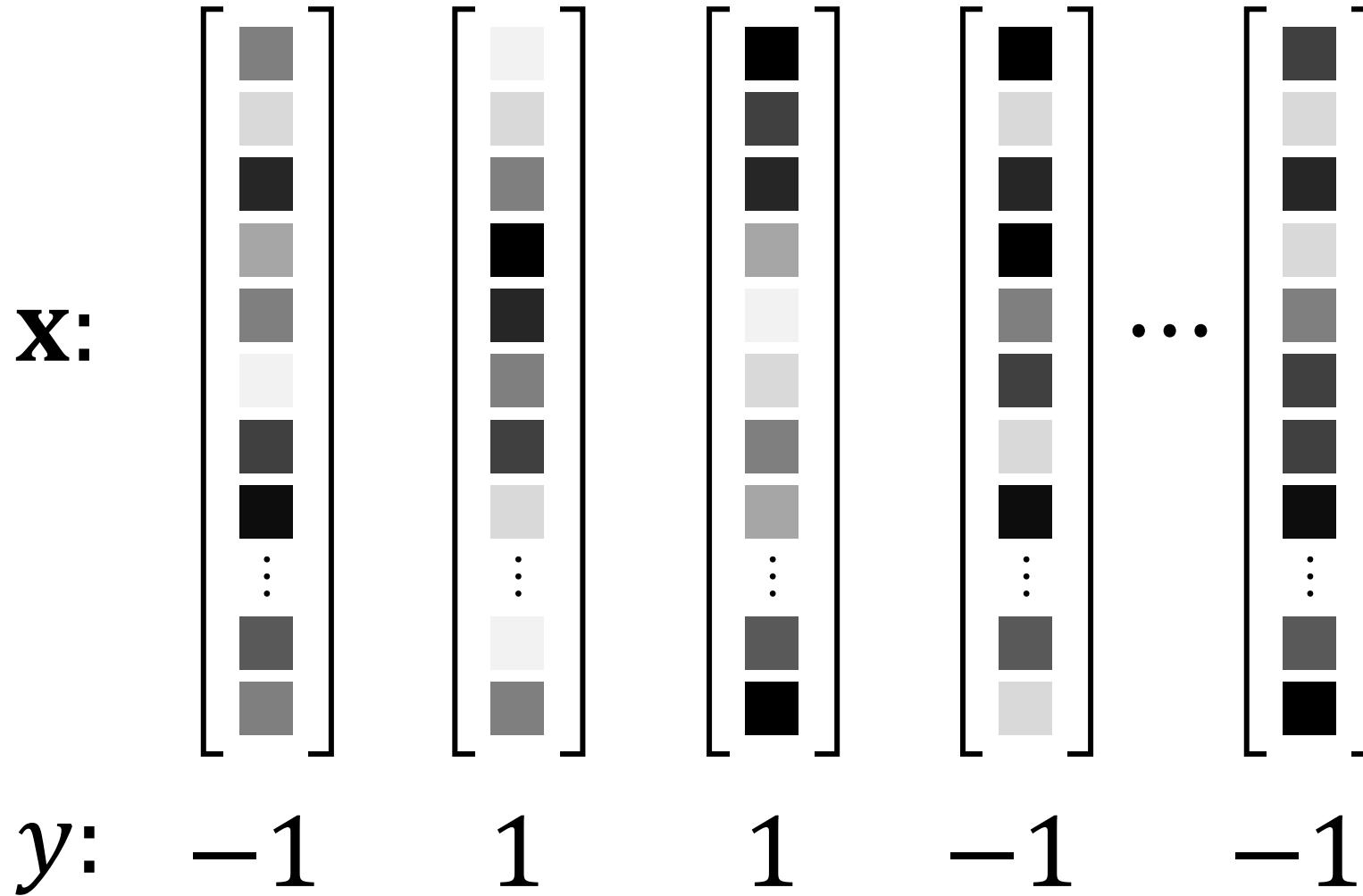
# Example: Cat/Dog Classifier

X:	$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$	$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$	$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$	$\dots$	$\begin{bmatrix} \dots \\ \vdots \\ \dots \end{bmatrix}$
y:	-1	1	1	-1	-1

$$\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N \|\mathbf{W}\mathbf{x}_i + \mathbf{b} - y_i\|^2$$

1x4096 / 1x1

# Example: Cat/Dog Classifier



$$\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^N |\mathbf{W}\mathbf{x}_i + \mathbf{b} - y_i|^2$$



$$\mathbf{W}^*, \mathbf{b}^*$$

# Example: Cat/Dog Classifier

(Online)



(\*Never seen by the model)

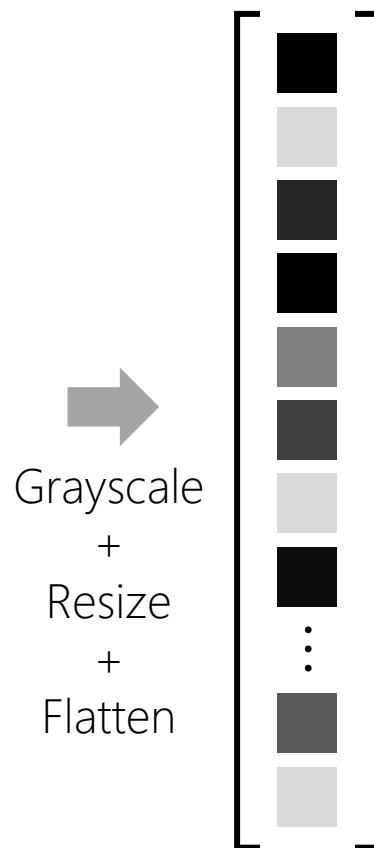
Query

# Example: Cat/Dog Classifier

(Online)



Query

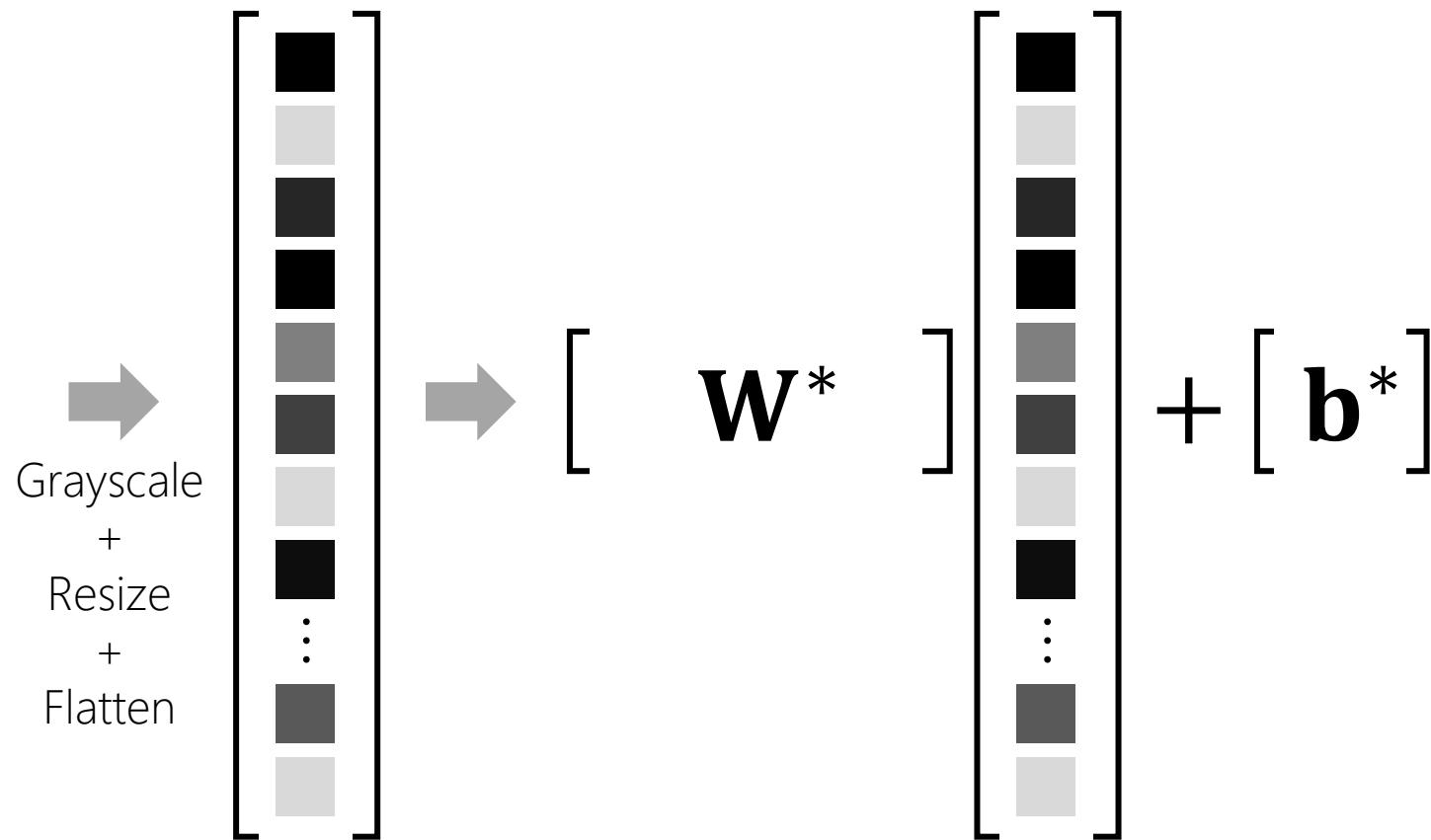


# Example: Cat/Dog Classifier

(Online)



Query



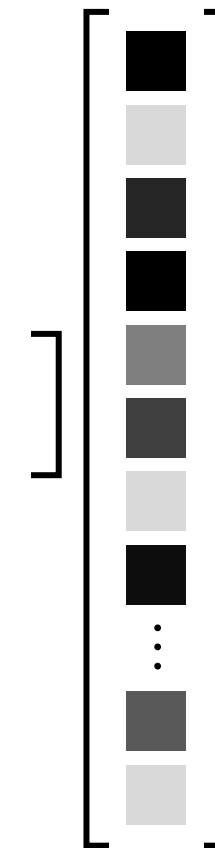
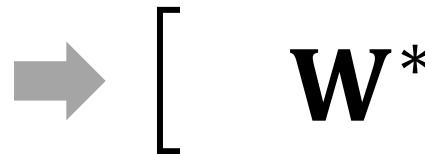
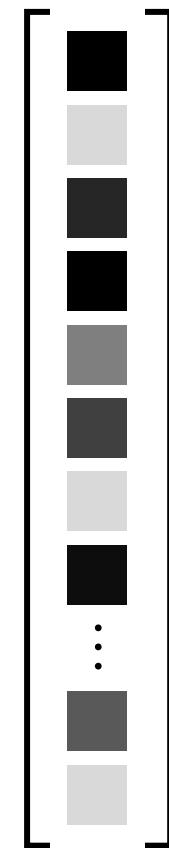
# Example: Cat/Dog Classifier

(Online)



## Query

Grayscale  
+  
Resize  
+  
Flatten



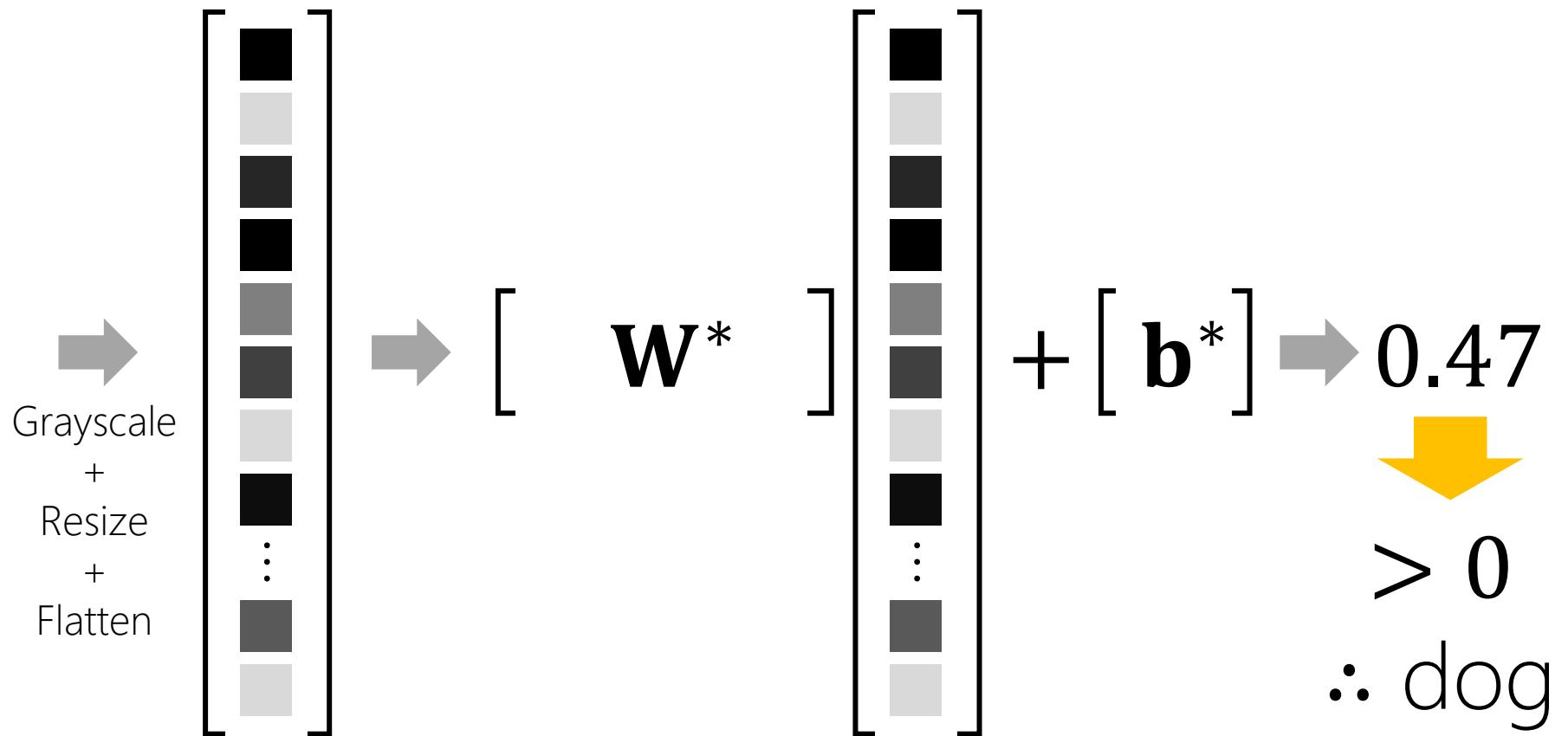
$$+ [b^*] \rightarrow 0.47$$

# Example: Cat/Dog Classifier

(Online)



Query



# Example: Cat/Dog Classifier

(Online)



Query

→  
Grayscale  
+  
Resize  
+  
Flatten

$$\begin{bmatrix} \text{Grayscale Row 1} \\ \text{Grayscale Row 2} \\ \vdots \\ \text{Grayscale Row N} \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} \mathbf{W}^* \end{bmatrix}$$

$$\begin{bmatrix} \text{Weight Row 1} \\ \text{Weight Row 2} \\ \vdots \\ \text{Weight Row N} \end{bmatrix}$$

$$+ \begin{bmatrix} \mathbf{b}^* \end{bmatrix} \rightarrow -0.7$$



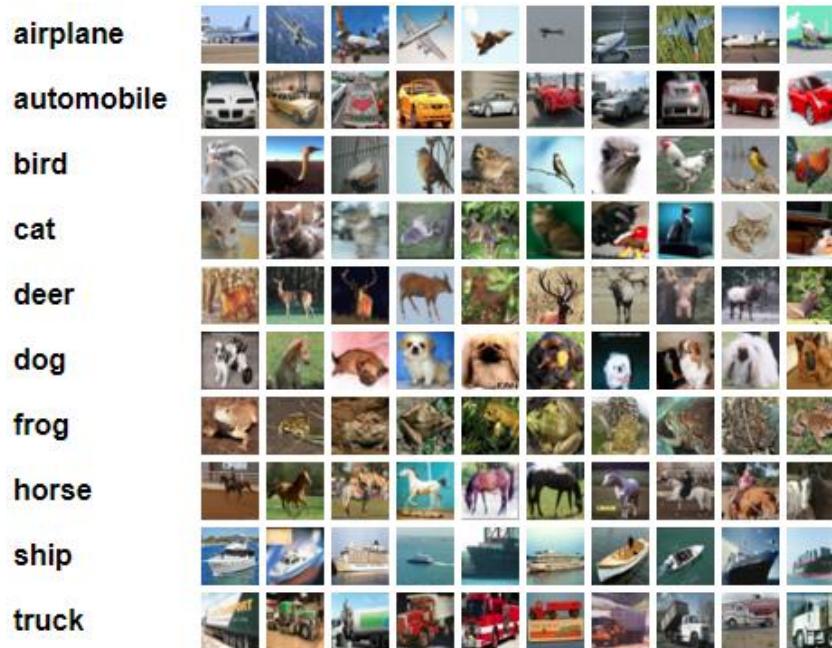
< 0

∴ cat

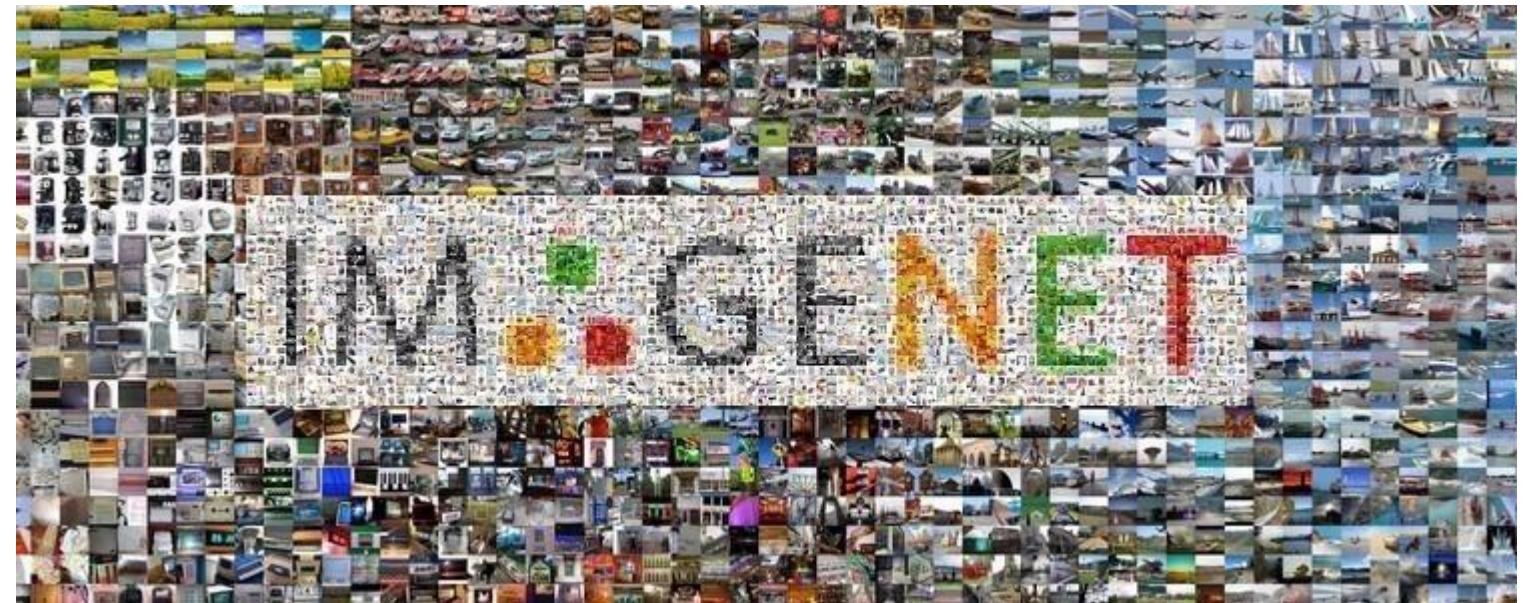
# Example: 10-way Classifier (CIFAR-10 Dataset)

- What if there were more than two classes (e.g. CIFAR, ImageNet)?

CIFAR-10: 60,000 images in 10 classes



ImageNet: 14,197,122 images in 21,841 classes  
ILSVRC: 1.4M images in 1,000 classes



# Probability Mass Function

$$\text{airplane} \quad \text{automobile} \quad \text{bird} \quad \text{cat} \quad \text{deer} \quad \text{dog} \quad \text{frog} \quad \text{horse} \quad \text{ship} \quad \text{truck}$$

→

$$y = \begin{bmatrix} P_{\text{airplane}} \\ P_{\text{automobile}} \\ P_{\text{bird}} \\ P_{\text{cat}} \\ P_{\text{deer}} \\ P_{\text{dog}} \\ P_{\text{frog}} \\ P_{\text{horse}} \\ P_{\text{ship}} \\ P_{\text{truck}} \end{bmatrix}$$
$$0 \leq P_i \leq 1, \forall i$$
$$\sum_i P_i = 1$$

# “One-hot” Encoding



True PMF

$y$

1	airplane
0	automobile
0	bird
0	cat
0	deer
0	dog
0	frog
0	horse
0	ship
0	truck

# “One-hot” Encoding



True PMF

$y$

1	airplane
0	automobile
0	bird
0	cat
0	deer
0	dog
0	frog
0	horse
0	ship
0	truck

What ML model predicts

$f$

0.7
0.1
0.05
0
0
0
0
0
0
0.15
0

# Cross-entropy

$$-\sum_i y_i \log(f_i)$$

In English: How similar are the two distributions?  
(We will re-visit this later to learn how and why.)

True PMF $y$	What ML model predicts $f$
1 airplane	0.7
0 automobile	0.1
0 bird	0.05
0 cat	0
0 deer	0
0 dog	0
0 frog	0
0 horse	0
0 ship	0.15
0 truck	0

# Example: 10-way Classifier (CIFAR-10 Dataset)



$$\mathbf{x}_1 = \begin{matrix} \text{airplane} \\ \text{automobile} \\ \text{bird} \\ \text{cat} \\ \text{deer} \\ \text{dog} \\ \text{frog} \\ \text{horse} \\ \text{ship} \\ \text{truck} \end{matrix} \dots \quad y_1 = [1 \quad 0 \quad 0]$$



$$\mathbf{x}_2 = \begin{matrix} \text{airplane} \\ \text{automobile} \\ \text{bird} \\ \text{cat} \\ \text{deer} \\ \text{dog} \\ \text{frog} \\ \text{horse} \\ \text{ship} \\ \text{truck} \end{matrix} \dots \quad y_2 = [0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0]$$

⋮

⋮

⋮

# Example: 10-way Classifier (CIFAR-10 Dataset)



$$\mathbf{x}_1 = \begin{matrix} \text{airplane} \\ \text{automobile} \\ \text{bird} \\ \text{cat} \\ \text{deer} \\ \text{dog} \\ \text{frog} \\ \text{horse} \\ \text{ship} \\ \text{truck} \end{matrix} \quad \mathbf{y}_1 = [1 \quad 0 \quad 0]$$



$$\text{Minimize } - \sum_i y_i \log(f(\mathbf{x}_i | \mathbf{W}, \mathbf{b})) \rightarrow \mathbf{W}^*, \mathbf{b}^*$$

$$\mathbf{x}_2 = \begin{matrix} \text{airplane} \\ \text{automobile} \\ \text{bird} \\ \text{cat} \\ \text{deer} \\ \text{dog} \\ \text{frog} \\ \text{horse} \\ \text{ship} \\ \text{truck} \end{matrix} \quad \mathbf{y}_2 = [0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0]$$

⋮

⋮

⋮

# Example: 10-way Classifier (CIFAR-10 Dataset)

(online)



Query

$$\begin{bmatrix} \text{10 x 4096} \\ \mathbf{W}^* \end{bmatrix} + \begin{bmatrix} \text{10 x 1} \\ \mathbf{b}^* \end{bmatrix}$$

4096 x 1

The diagram illustrates the computation of a query vector for a 10-way classifier. It shows the addition of two vectors:  $\mathbf{W}^*$  (Weight matrix, 10 x 4096) and  $\mathbf{b}^*$  (Bias vector, 10 x 1). The result is a 4096 x 1 vector. The weight matrix  $\mathbf{W}^*$  is represented by a vertical column of 10 horizontal bars, each corresponding to one of the 10 classes. The bias vector  $\mathbf{b}^*$  is represented by a single horizontal bar above the addition sign.

# Example: 10-way Classifier (CIFAR-10 Dataset)

(online)



Query

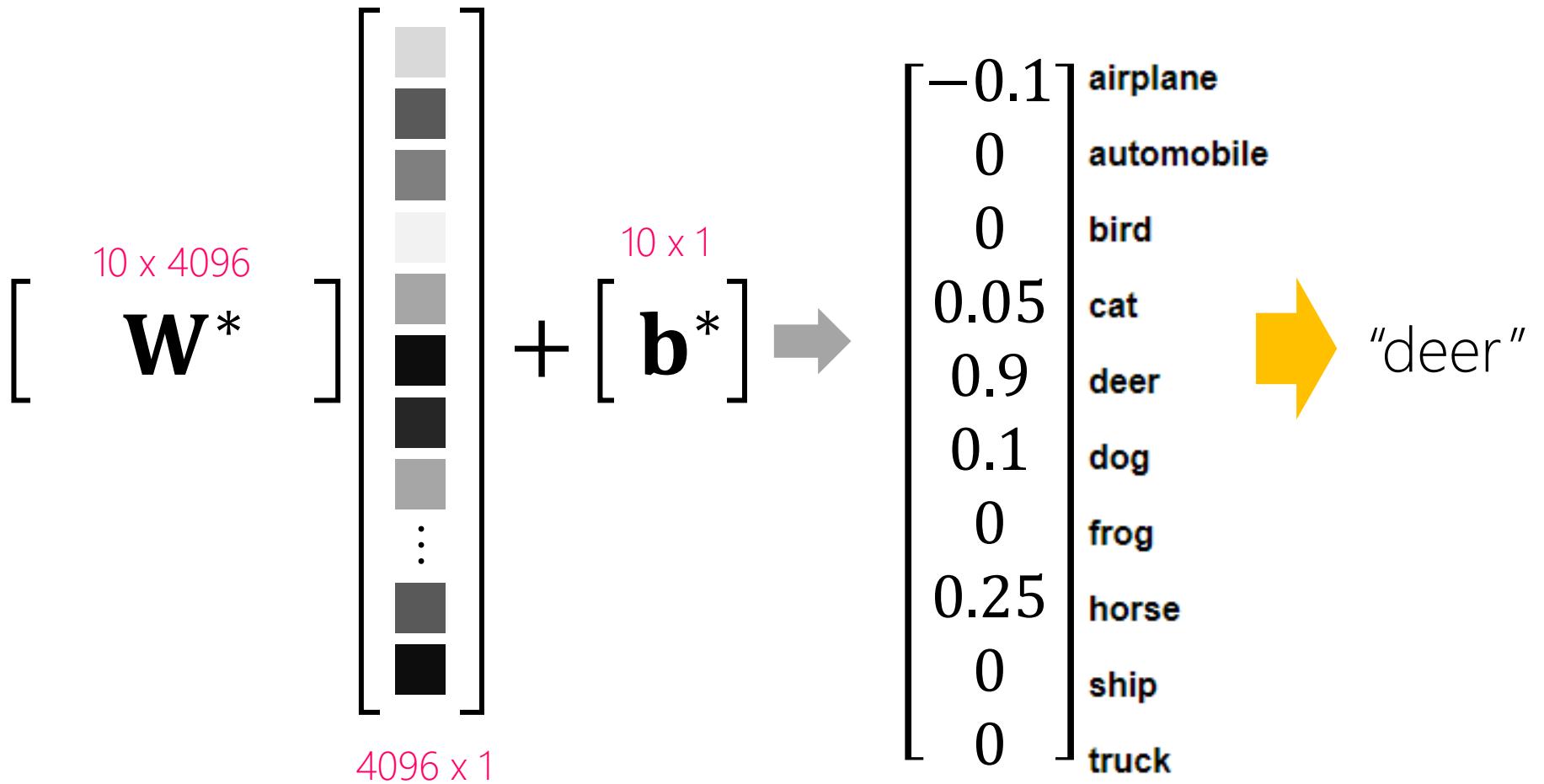
$$\begin{matrix} \text{(online)} \\ \text{Query} \\ \text{---} \\ \text{[ } \mathbf{W}^* \text{ ]} \end{matrix} \begin{matrix} 10 \times 4096 \\ \vdots \\ 4096 \times 1 \end{matrix} + \begin{matrix} 10 \times 1 \\ \text{[ } \mathbf{b}^* \text{ ]} \end{matrix} \rightarrow \begin{matrix} -0.1 \\ 0 \\ 0 \\ 0.05 \\ 0.9 \\ 0.1 \\ 0 \\ 0.25 \\ 0 \\ 0 \end{matrix}$$

The diagram illustrates the computation of a 10-way classification score vector. It starts with a query image of a deer, which is processed by a weight matrix  $\mathbf{W}^*$  (dimensions  $10 \times 4096$ ) and a bias vector  $\mathbf{b}^*$  (dimensions  $10 \times 1$ ). The result is a vector of scores for ten categories: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. The scores are numerical values ranging from -0.1 to 0.25, with 'deer' having the highest positive score of 0.9.

Category	Score
airplane	-0.1
automobile	0
bird	0
cat	0.05
deer	0.9
dog	0.1
frog	0
horse	0.25
ship	0
truck	0

# Example: 10-way Classifier (CIFAR-10 Dataset)

(online)

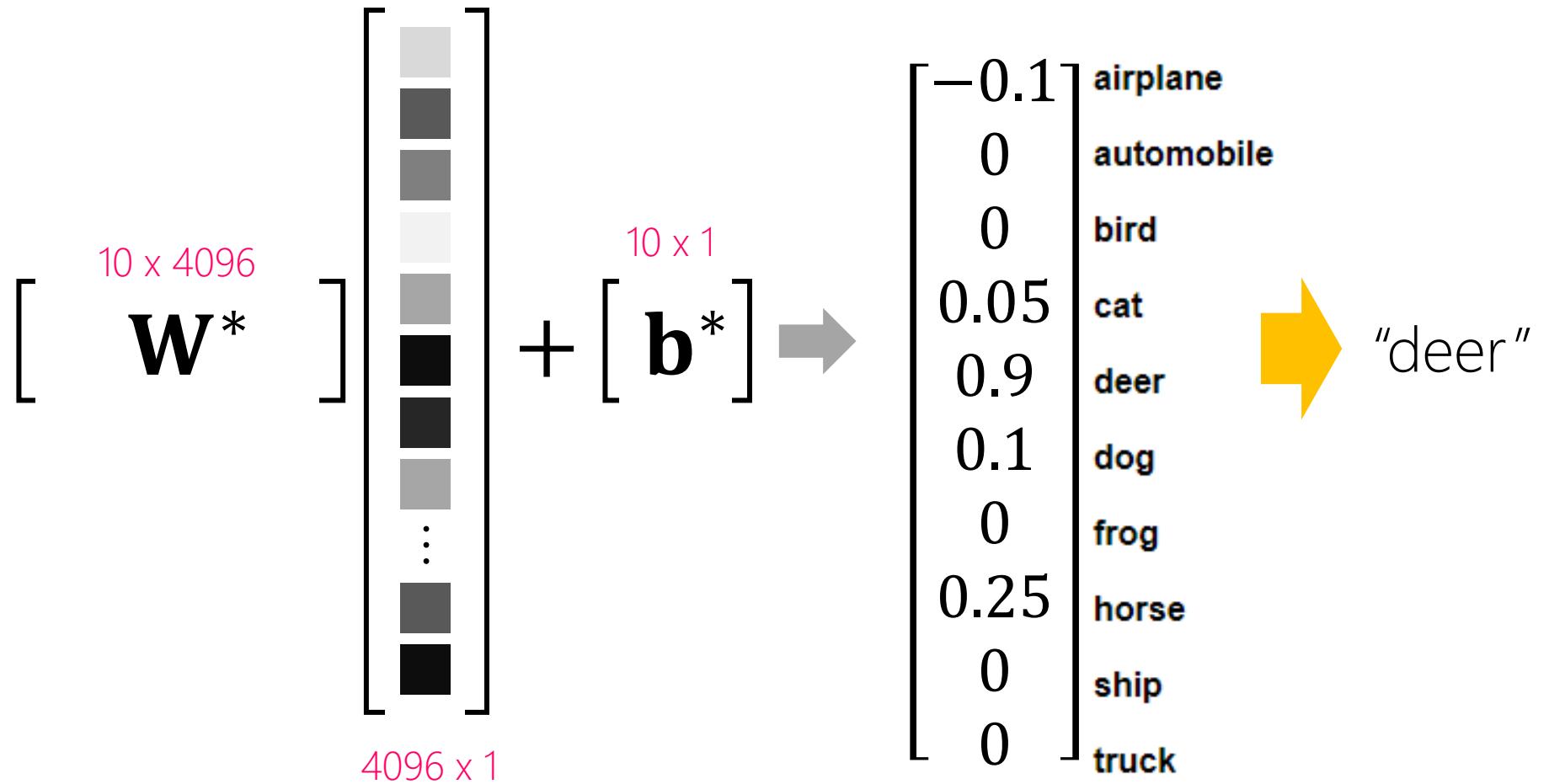


# Example: 10-way Classifier (CIFAR-10 Dataset)

(online)



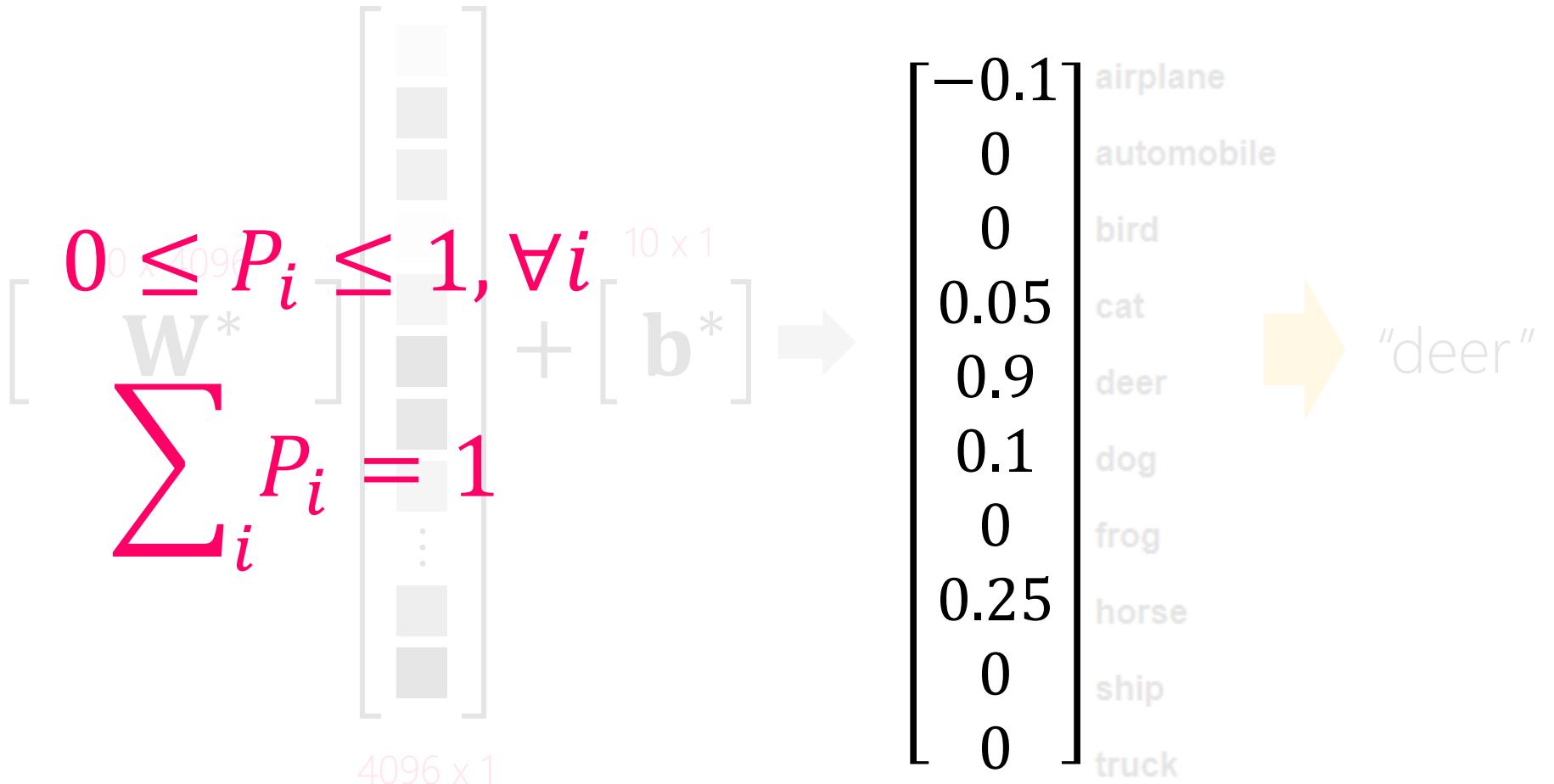
Query



Wait a second, anything wrong?

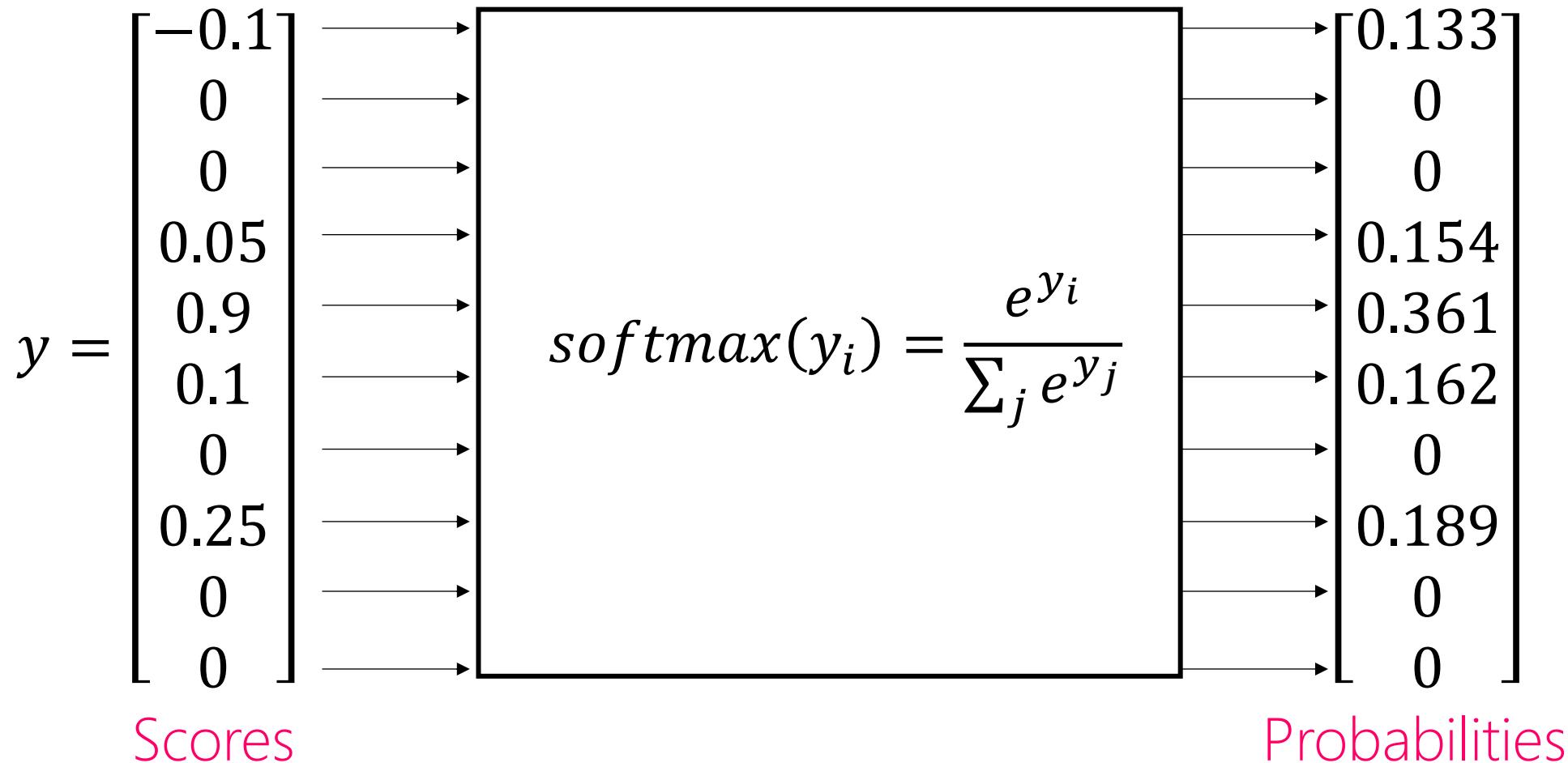
# Example: 10-way Classifier (CIFAR-10 Dataset)

(online)



Wait a second, anything wrong?

# Softmax!

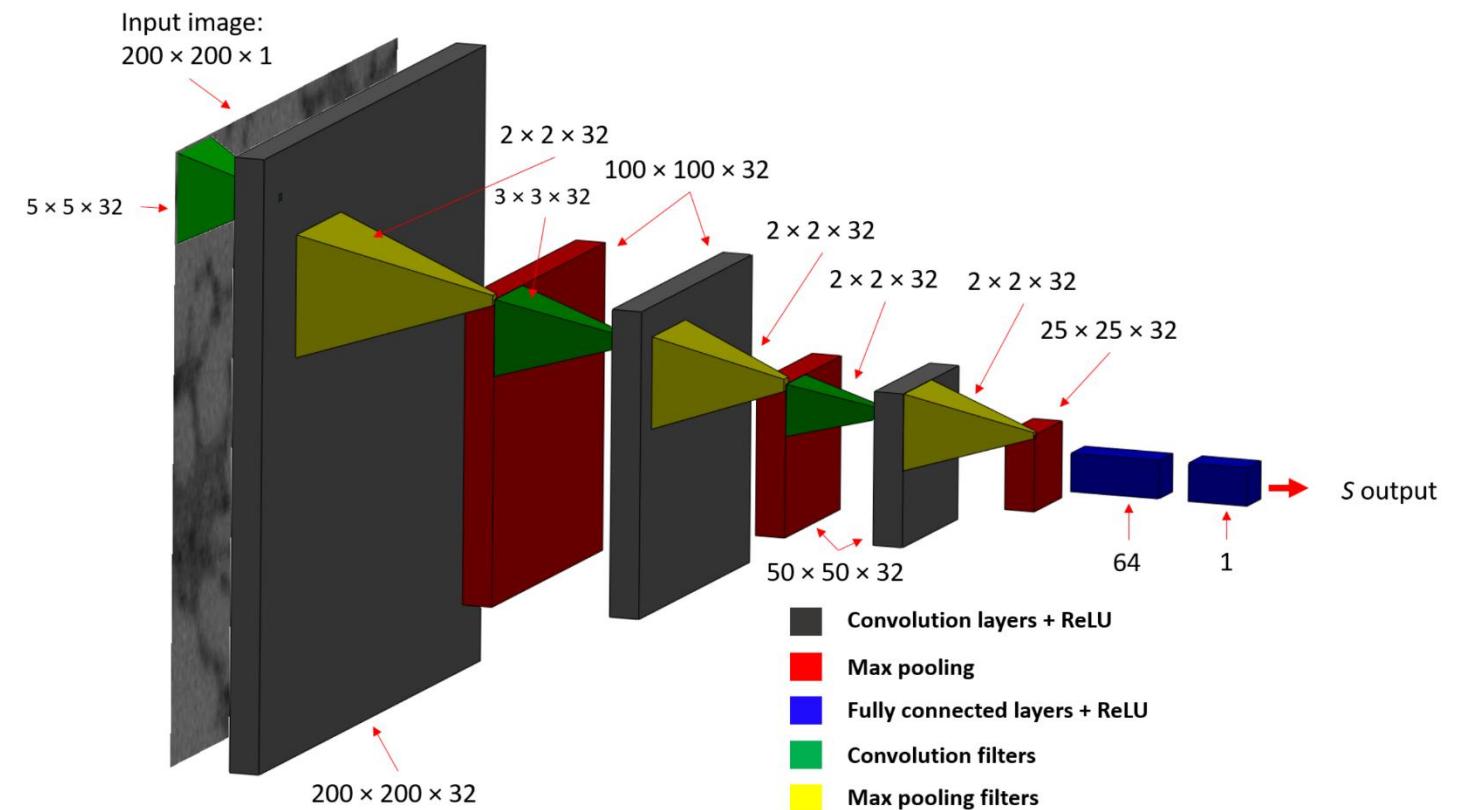
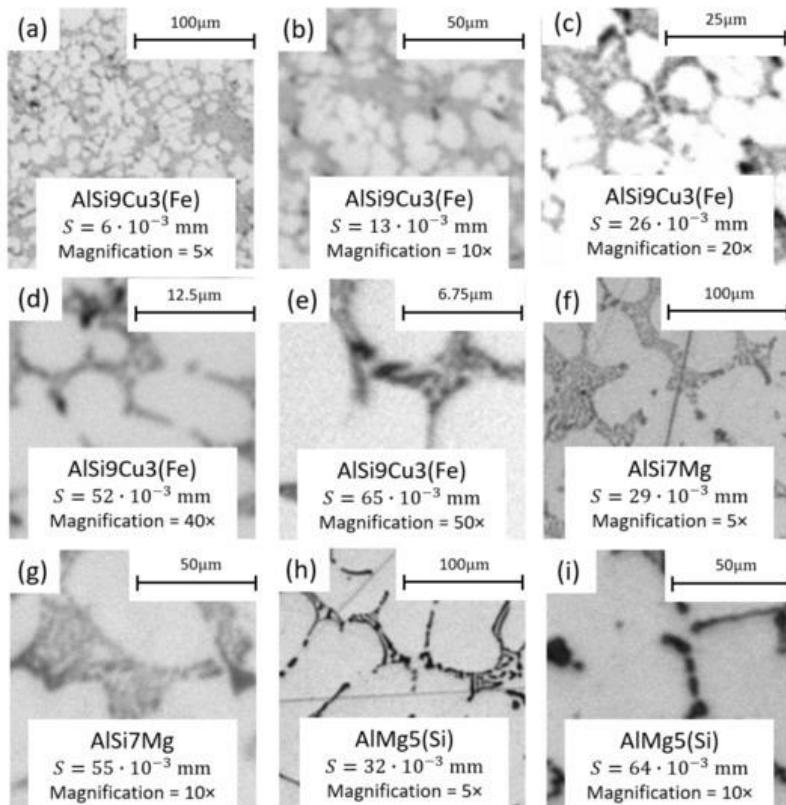


# n-way Softmax Classifier: Step-by-Step

1. Collect a bunch of data  $(\mathbf{x}_i, y_i)_{i=1,\dots,N}$
2. Encode  $y_i$  in one-hot encoding
3. Split the data into two sets: training set and test set
4. Solve  $\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^{N_{\text{train}}} -y_i \log(f_i)$  where  $f_i = \text{softmax}(\sigma(\mathbf{W}\mathbf{x}_i + \mathbf{b}))$
5. Test the trained model  $f(\mathbf{x} | \mathbf{W}^*, \mathbf{b}^*) = \text{softmax}(\sigma(\mathbf{W}^*\mathbf{x} + \mathbf{b}^*))$  on the test set

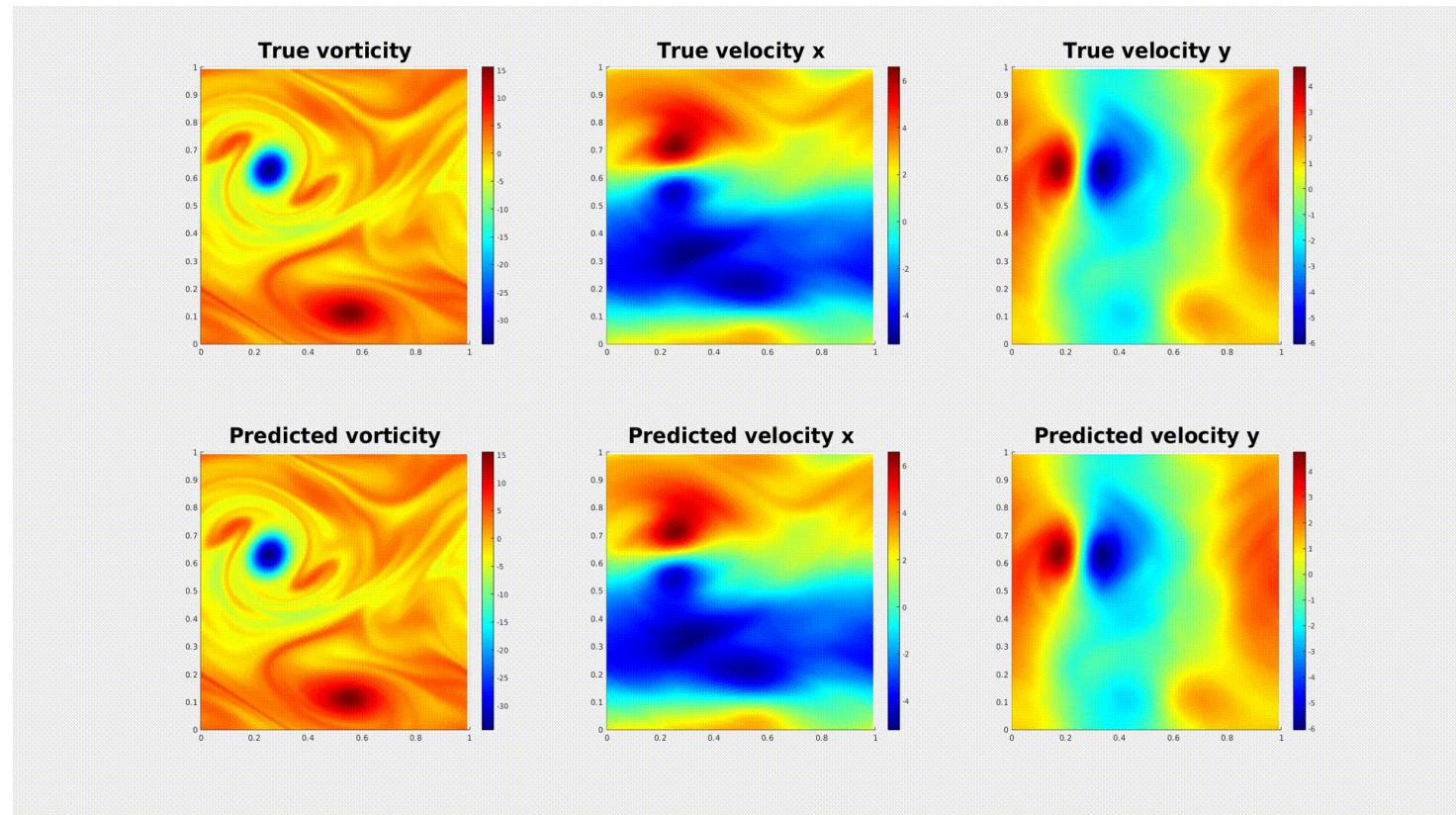
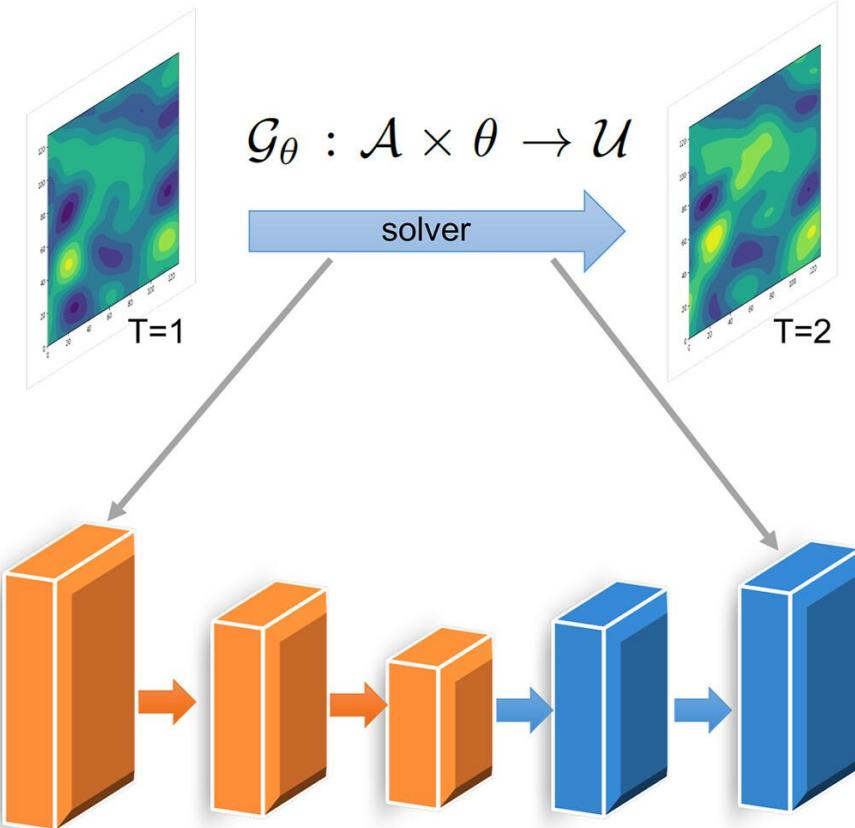
# Generalizing this formulation...

- Instead of  $y$  being a class probability...



# Generalizing this formulation...

- Instead of  $y$  being a class probability...



# More details?

- Coming up:
  - Lecture 3: Training neural networks: gradient descent, backpropagation, and more
  - Lab 1: Introduction to PyTorch—Build your first neural network
- Fun videos to watch:
  - <https://www.3blue1brown.com/topics/neural-networks>