

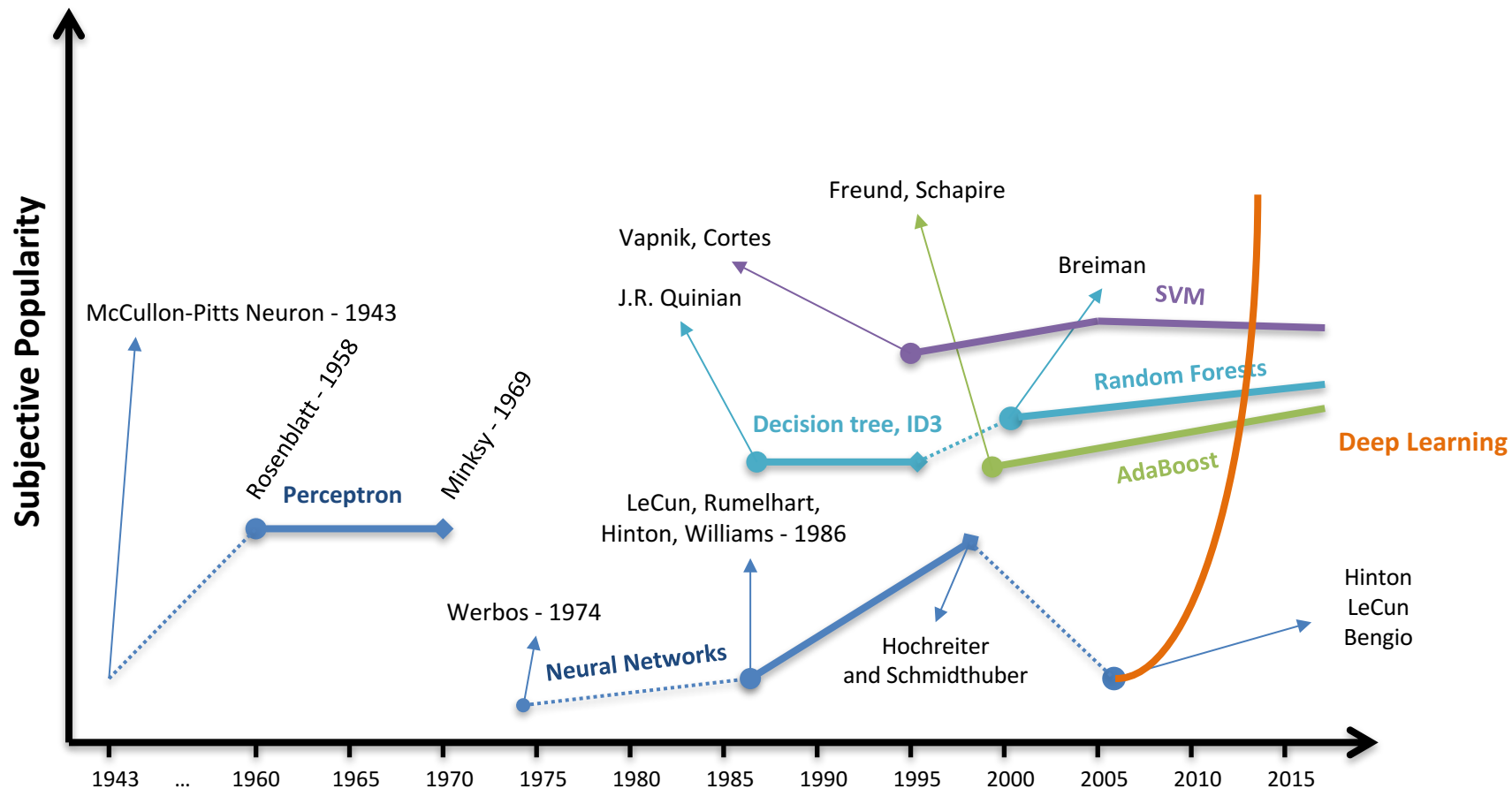
A Concise History of Neural Networks

Dr. Haojin Yang

Why Learn History?

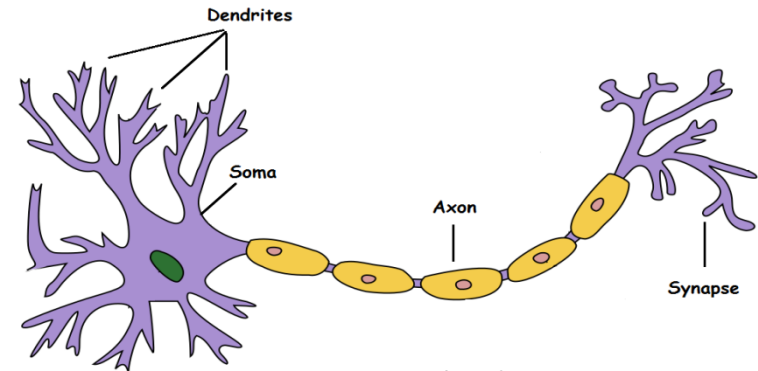
“... the best reason to learn history: not in order to predict the future, but to free yourself of the past and imagine alternative destinies. Of course this is not total freedom - we cannot avoid being shaped by the past. But some freedom is better than none.”

- *Homo Deus: A Brief History of Tomorrow* (Yuval Noah Harari, 2015)



A Biological Neuron

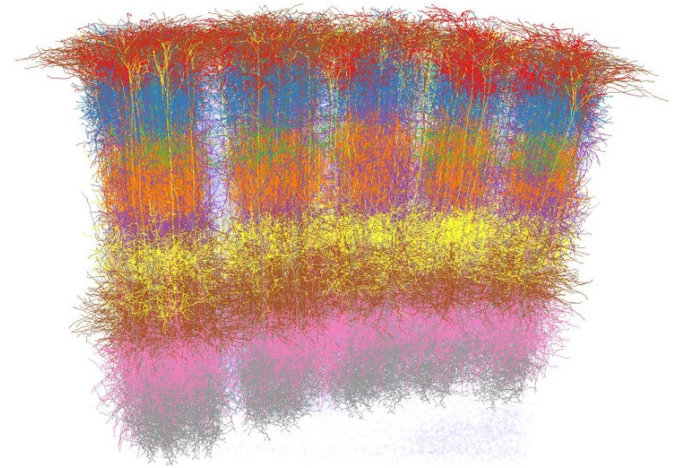
- Dendrite: Receives signals from other neurons
- Soma: Processes the information
- Axon: Transmits the output of this neuron
- Synapse: Point of connection to other neurons



Source: Wikipedia

Cortical Columns

- Neurons in cortex arranged into many stacks, or “columns” that process information in parallel

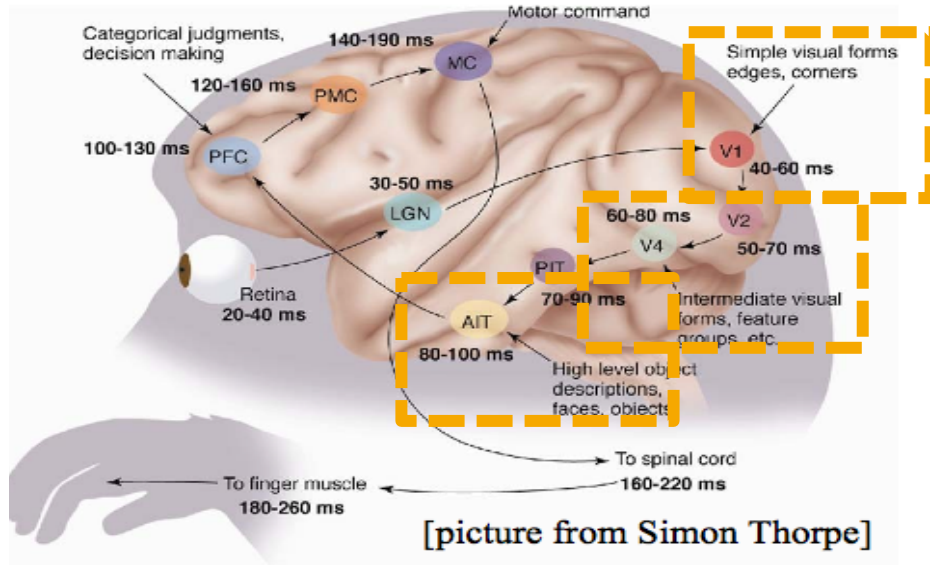


credit: Marcel Oberlaender et al.

Enlightenment: The Mammalian Visual Cortex is Hierarchical

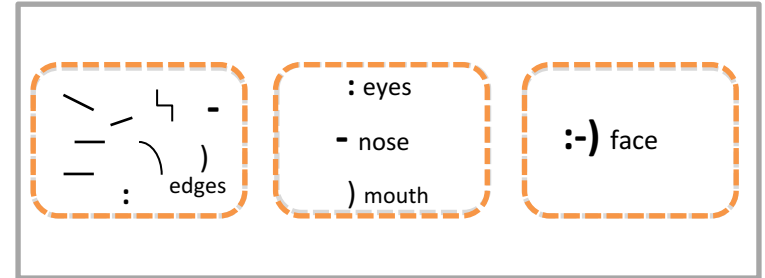
The ventral (recognition) pathway in the visual cortex has multiple stages:

- Retina - LGN - V1 - V2 - V4 - PIT - AIT ..., lots of intermediate representations



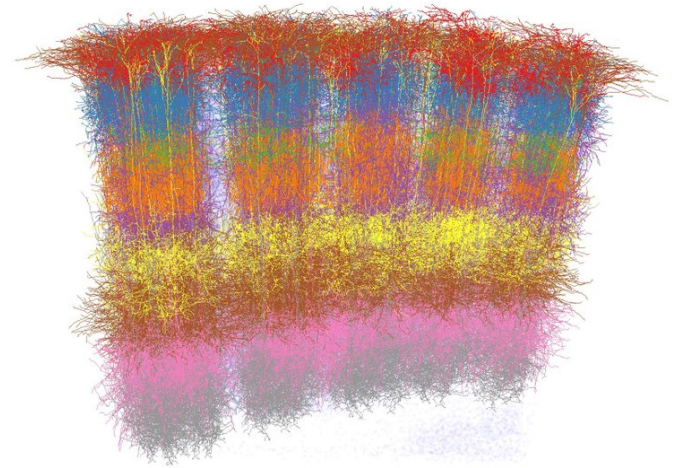
[picture from Simon Thorpe]

[Gallant & Van Essen]



Cortical Columns

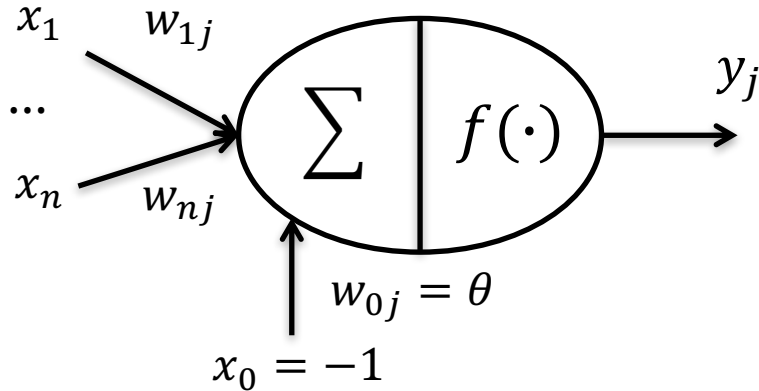
- “mini-columns” of around 100 neurons are organized into larger “hyper-columns”.
- There are 100 million mini-columns in your cortex
- mini-column is statistically meaningful, but what is the functional meaning?



credit: Marcel Oberlaender et al.

M-P Neuron

A Logical Calculus of the Ideas Immanent in Nervous Activity, Warren McCulloch and Walter Pitts (1943)



$$y_j = \sum_{i=1}^n w_{ij} x_i - \theta_j$$

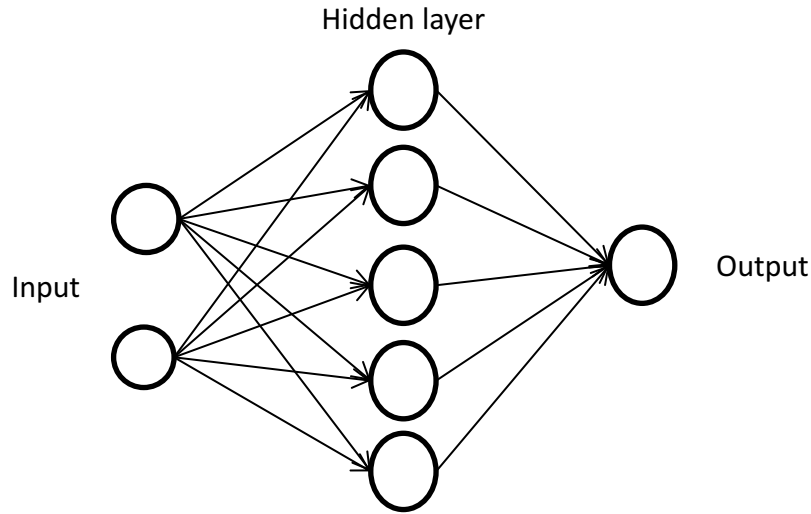
Biological	M-P
neuron	j
inputs	x_i
weights	w_{ij}
output	y_j
sum	Σ
membrane potential	$\sum_{i=1}^n w_{ij} x_i$
threshold	θ



Warren McCulloch & Walter Pitts

Perceptron

- Single layer neural network



Frank Rosenblatt
(1928-1971)

$$y_j = \begin{cases} 0 & \text{if } \sum_{i=1}^n w_{ij}x_i \leq \theta_j \\ 1 & \text{if } \sum_{i=1}^n w_{ij}x_i > \theta_j \end{cases}$$

Perceptron

- Follows the Hebbian theory

$$w_{ij} = x_j \rightarrow x_i$$

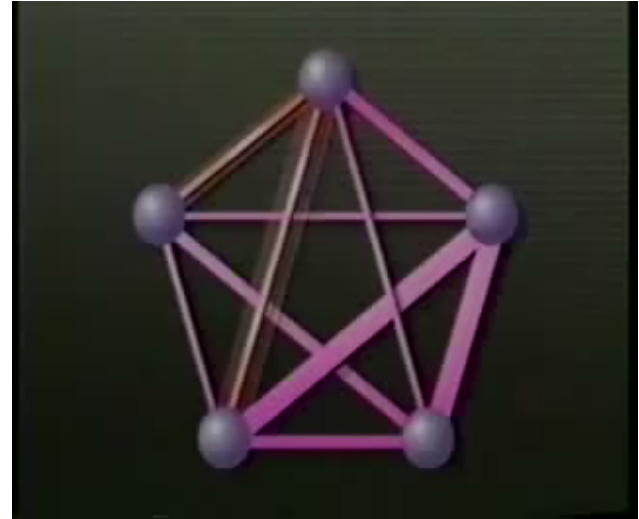
$$\Delta w_{ij} = \eta \cdot a_i \cdot y_j$$

Δw_{ij} : the changes in the i _th synaptic weight

η : learning rate

a_i : activation value of i _th synapse

y_j : output of j _th synapse



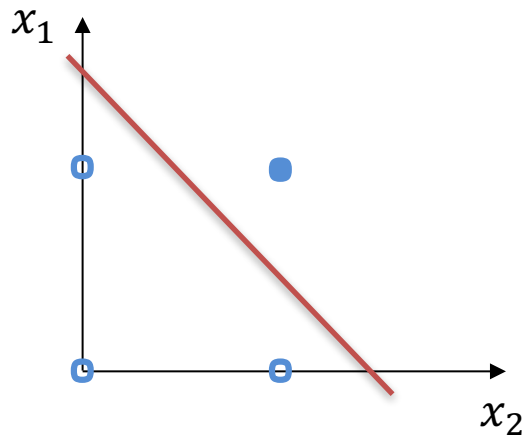
source: youtube

Perceptron

Perceptrons, Marvin Minsky and Seymour Papert (1969)

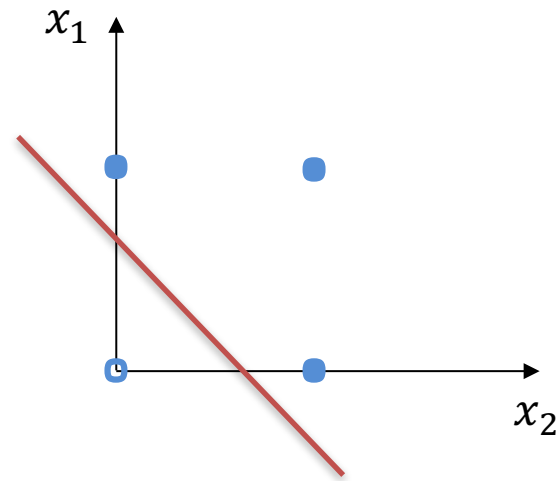
x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

AND



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

OR

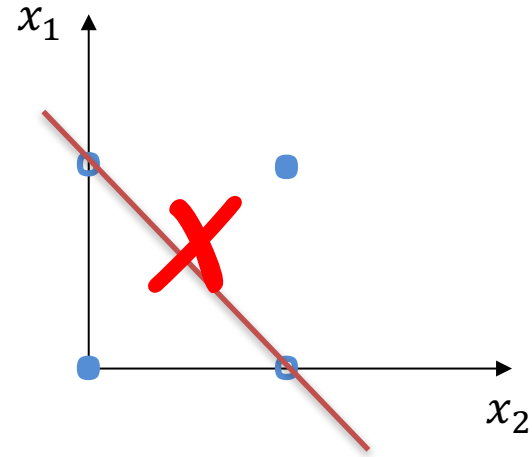


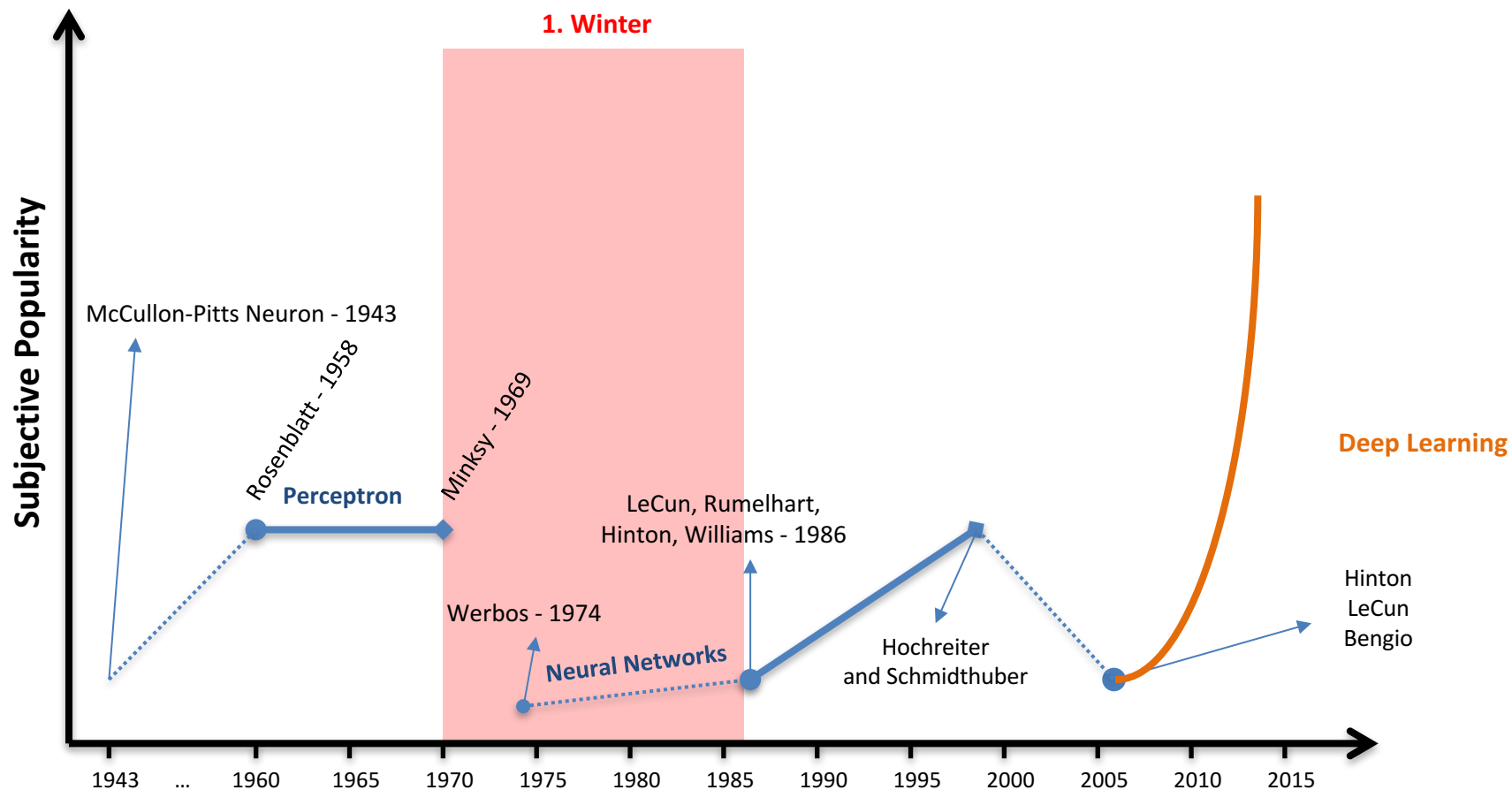
Perceptron

Perceptrons, Marvin Minsky and Seymour Papert (1969)

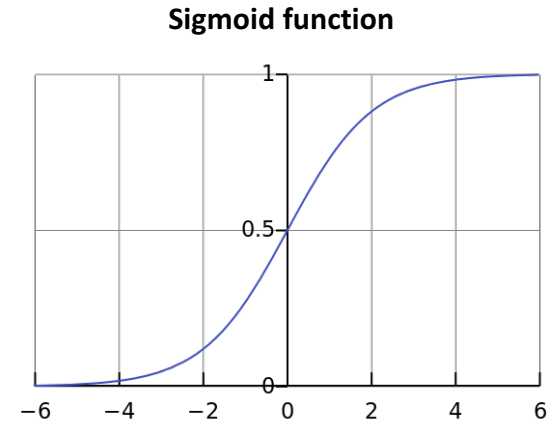
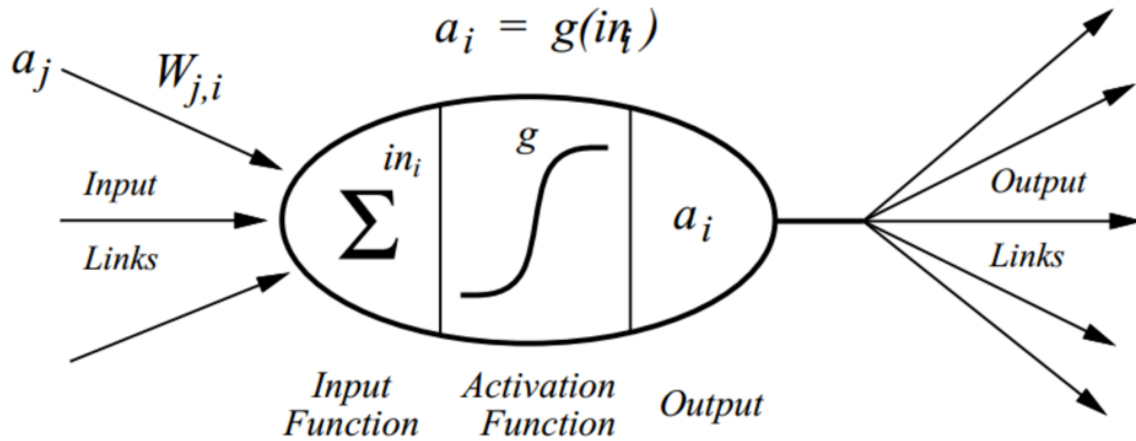
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

XOR





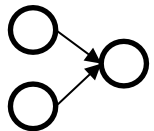

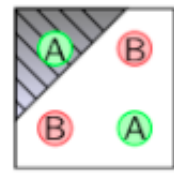
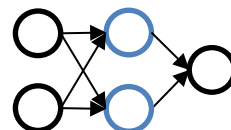

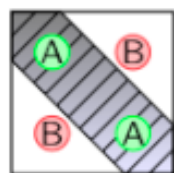
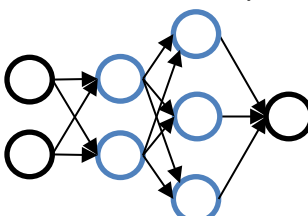

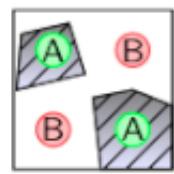
Non-linear Activation



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

Multi-Layer Perceptron

Perceptrons - Expanded Edition, Marvin Minsky (1987)

Structure	Decision area	Area shape	XOR problem
<p>linear activation</p> 	<p>divided by one hyperplane</p>		
<p>1 hidden layer, non-linear</p> 	<p>open convex area and closed convex area</p>		
<p>double hidden layer</p> 	<p>arbitrary shape (complexity depending on the number of units)</p>		

Backpropagation

- ~1986 **Backpropagation**: “Learning representations by back-propagating errors”
- BP enables training multi-layer artificial neural networks
 - Update weights using Gradient Descent
 - BP is very efficient for calculating $\nabla_{\theta} \text{Loss}$



Geoffrey Hinton

NATURE VOL. 323 9 OCTOBER 1986

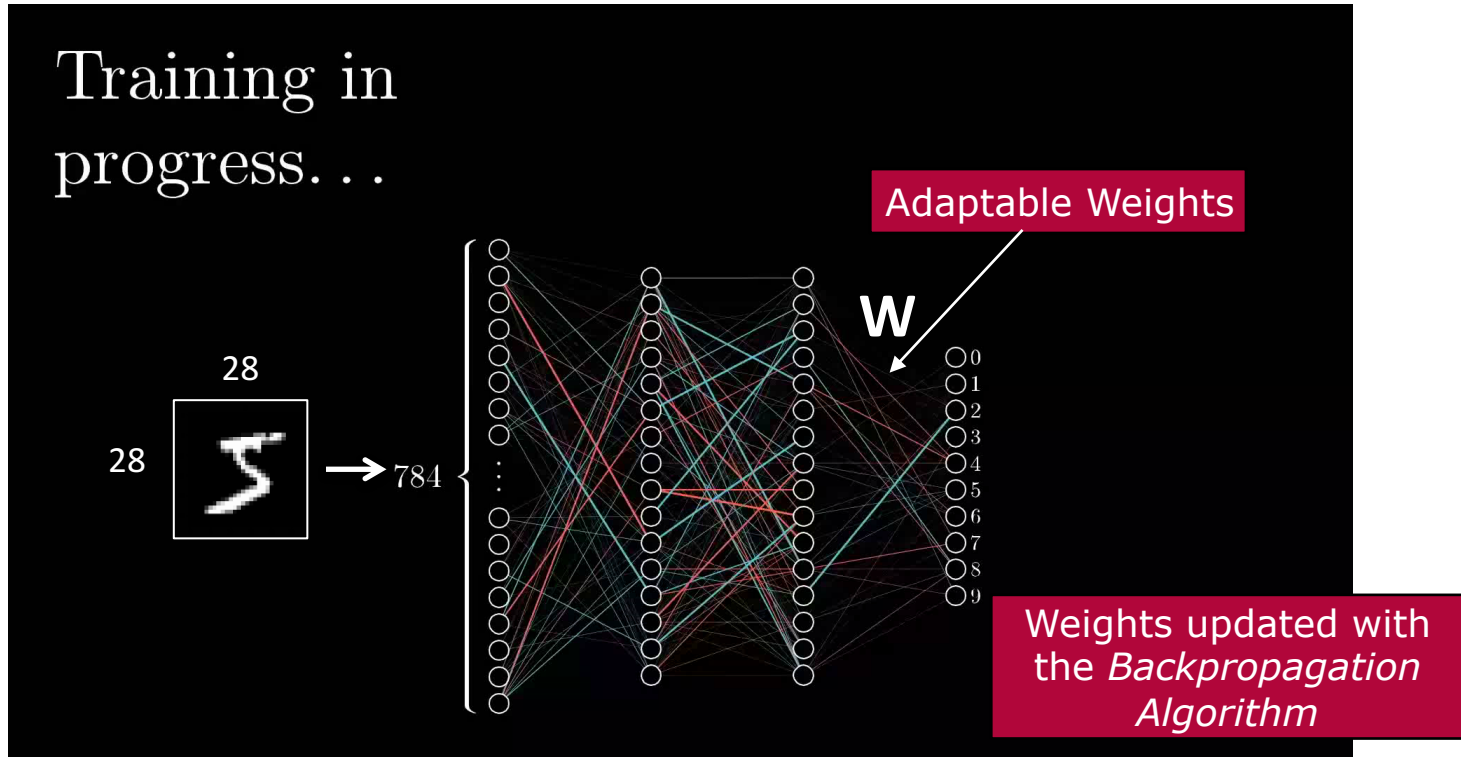
Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

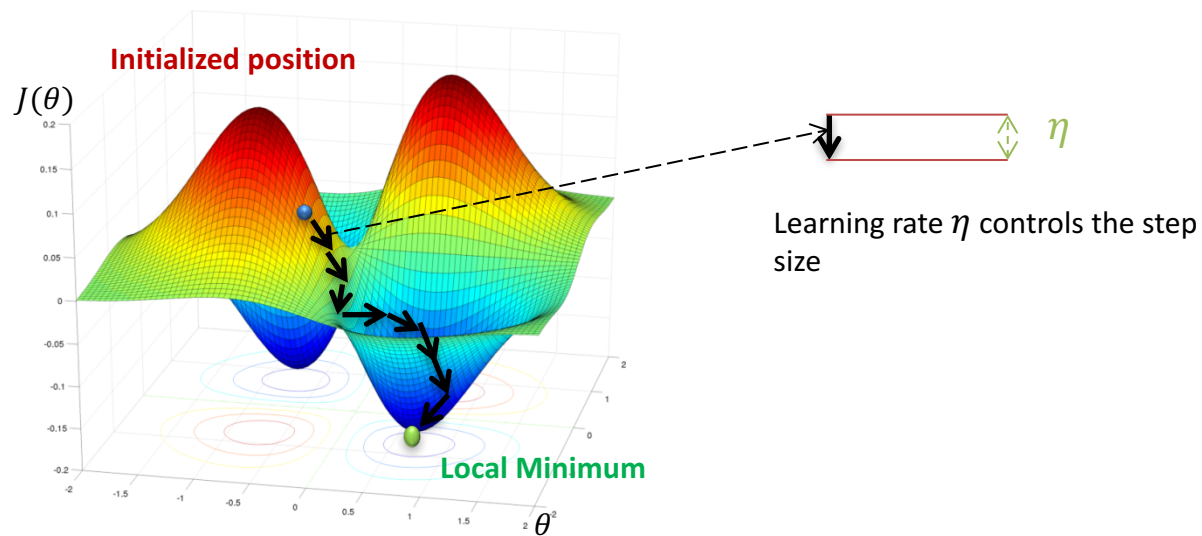
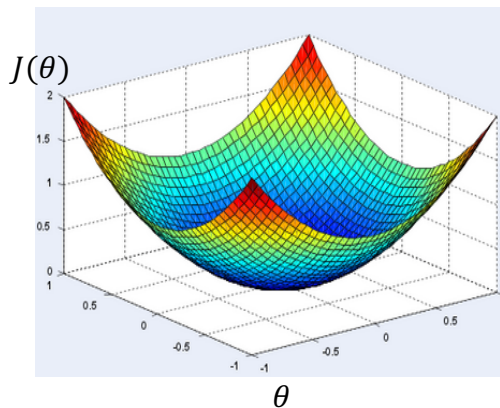
† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

Artificial Neural Networks



Gradient Descent

$$\theta := \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)})$$

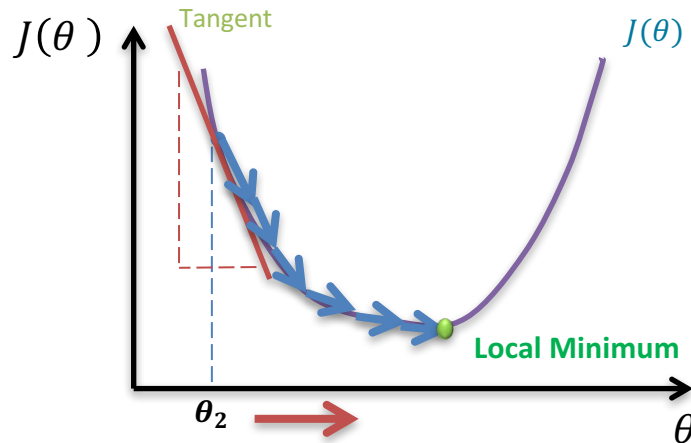


Gradient Descent

- Optimization objective: $\min J(\theta; x^{(i:i+n)}; y^{(i:i+n)}), \theta \in \mathbb{R}$

$$\theta := \theta - \eta \cdot \frac{\partial J(\theta; x^{(i:i+n)}; y^{(i:i+n)})}{\partial \theta}$$

Learning rate Derivative

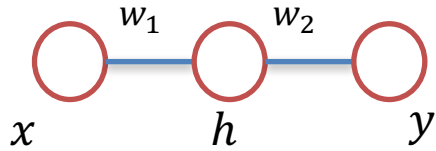


$$\theta_2 := \theta_2 - \eta \cdot \frac{\partial J(\theta_2)}{\partial \theta_2}$$

< 0

$$\theta_2 := \theta_2 - \eta \cdot (\text{negative num})$$

Chaining

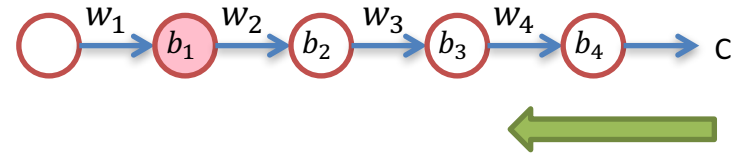
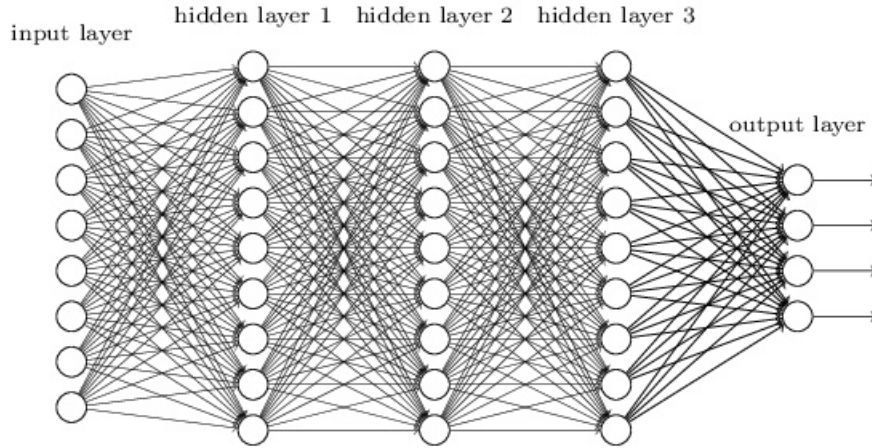


$$h = x * w_1$$
$$\frac{\partial h}{\partial w_1} = x$$

$$y = h * w_2$$
$$\frac{\partial y}{\partial h} = w_2$$

$$y = x * w_1 * w_2$$
$$\frac{\partial y}{\partial w_1} = w_2 * x$$
$$\frac{\partial y}{\partial w_1} = \frac{\partial y}{\partial h} \frac{\partial h}{\partial w_1}$$

Chaining

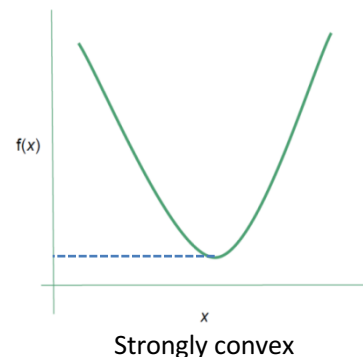
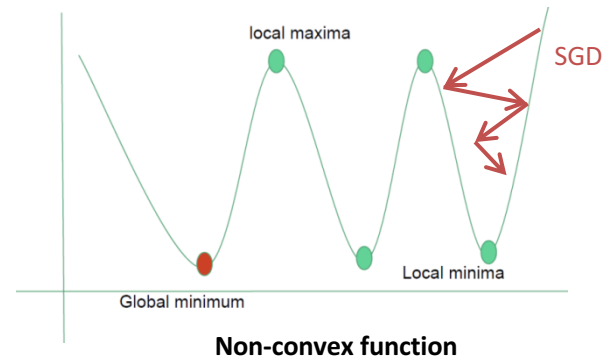


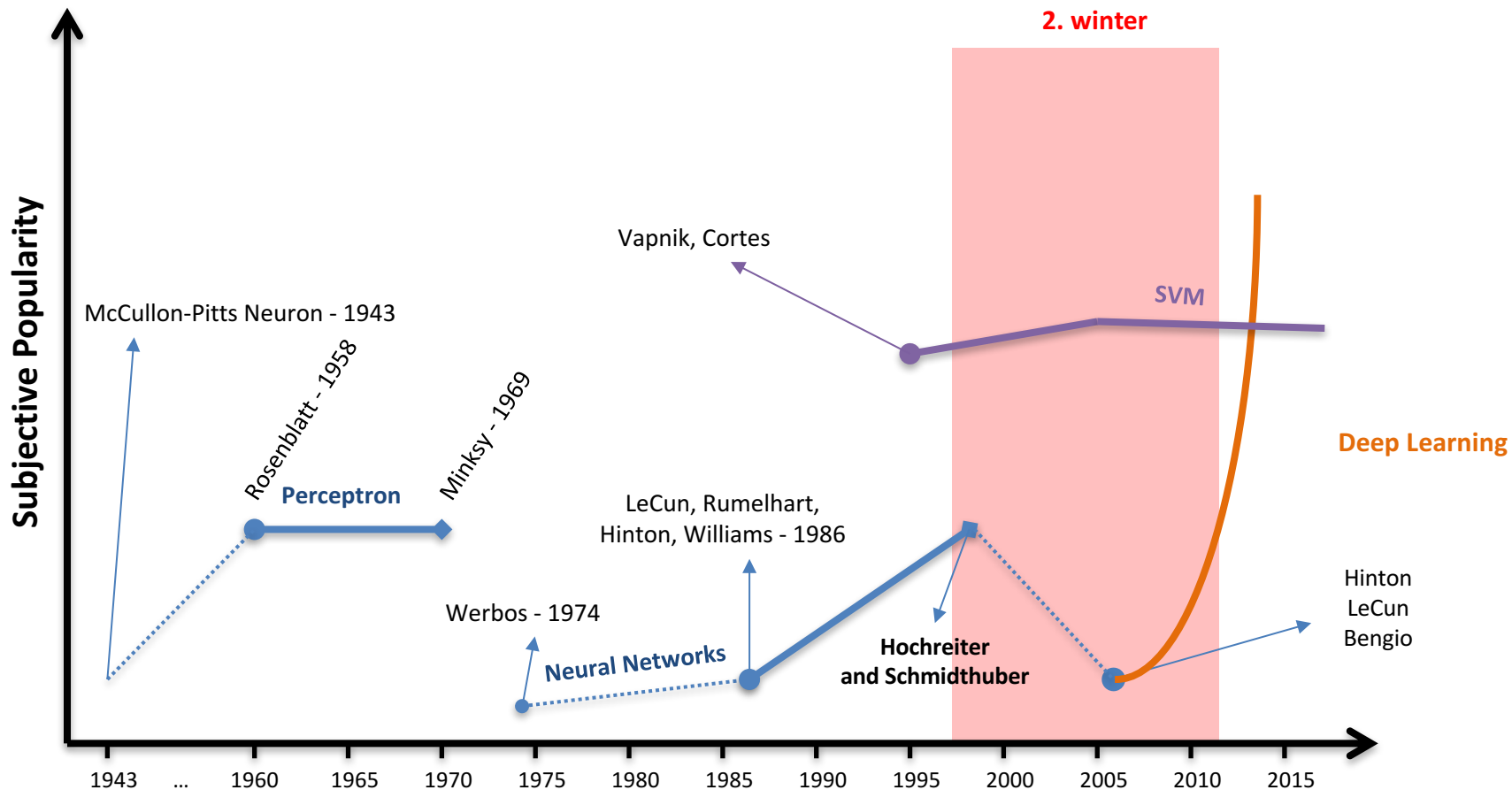
$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \frac{\partial x_2}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

Backpropagation

- BP is biologically difficult to explain
- BP based on Gradient Descent algorithms → **non-convex** optimization problem
- Other methods were better e.g., SVM
- **Unable to successfully train deeper networks**

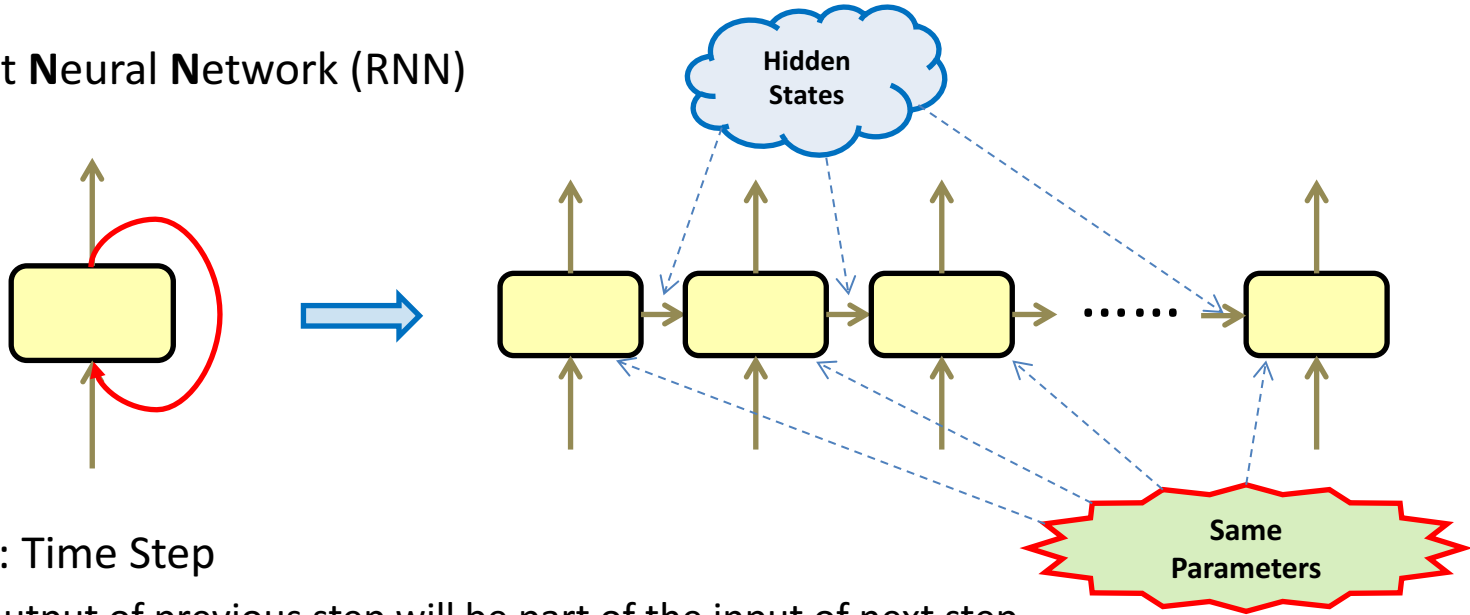




Why can't it be deeper?

Untersuchungen zu dynamischen neuronalen Netzen. Sepp Hochreiter, Diploma thesis, TU Munich (1991)

- **Recurrent Neural Network (RNN)**



- **Keyword: Time Step**

- The output of previous step will be part of the input of next step.
 - “Hidden State”
- Parameters in the layer are kept static through continuous steps.

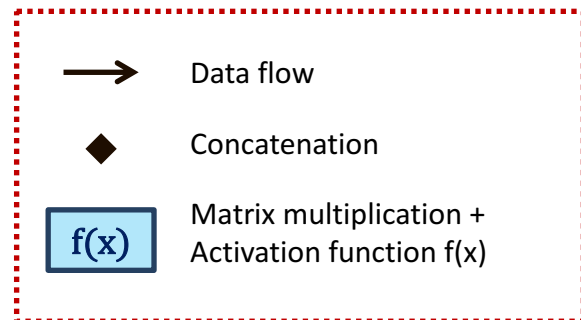
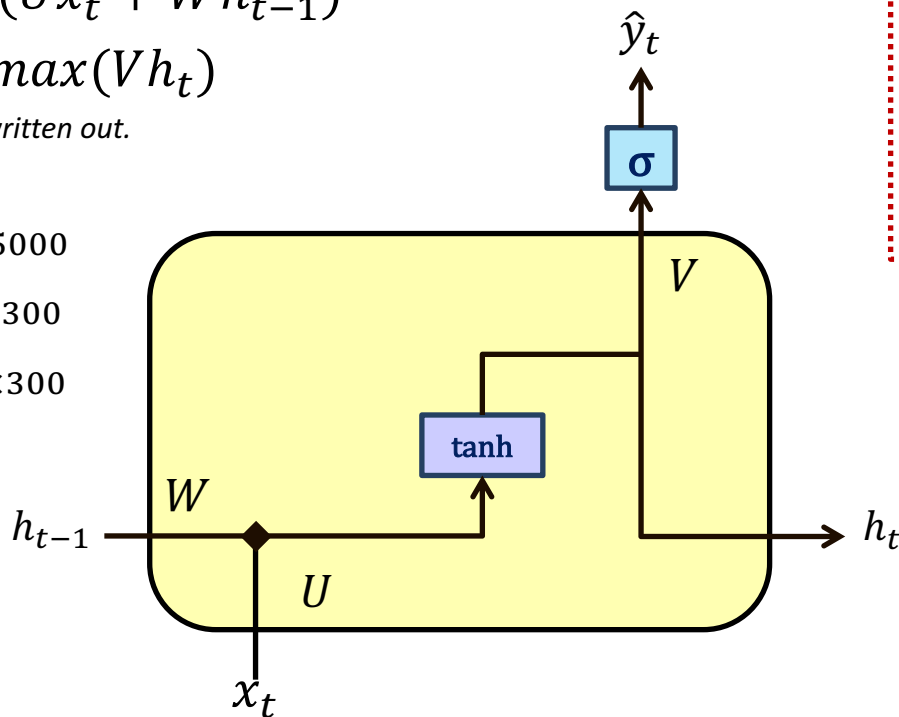
Gradient Vanishing/Exploding

- When the chain of back propagation is very long, so many times of **multiplication** would be applied, somehow repeated (e.g. in W^T).
 - **Gradient Vanishing**
 - ... when the values are smaller than 1.
 - ... the model only learns how to keep “short memory”.
 - ... no good solution found with vanilla RNN.
 - **Gradient Exploding**
 - ... when the values are greater than 1.
 - ... could be partially solved by setting a threshold to “clip” the gradient.

Vanilla RNN

- $h_t = \tanh(Ux_t + Wh_{t-1})$
- $\hat{y}_t = \text{softmax}(Vh_t)$
- * Bias 'b' is NOT written out.

- $U \in \mathbb{R}^{300 \times 5000}$
- $W \in \mathbb{R}^{300 \times 300}$
- $V \in \mathbb{R}^{5000 \times 300}$

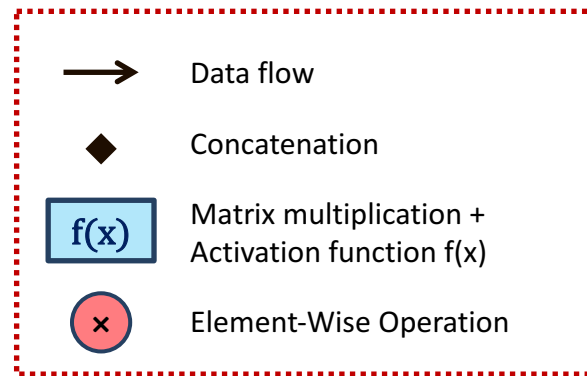
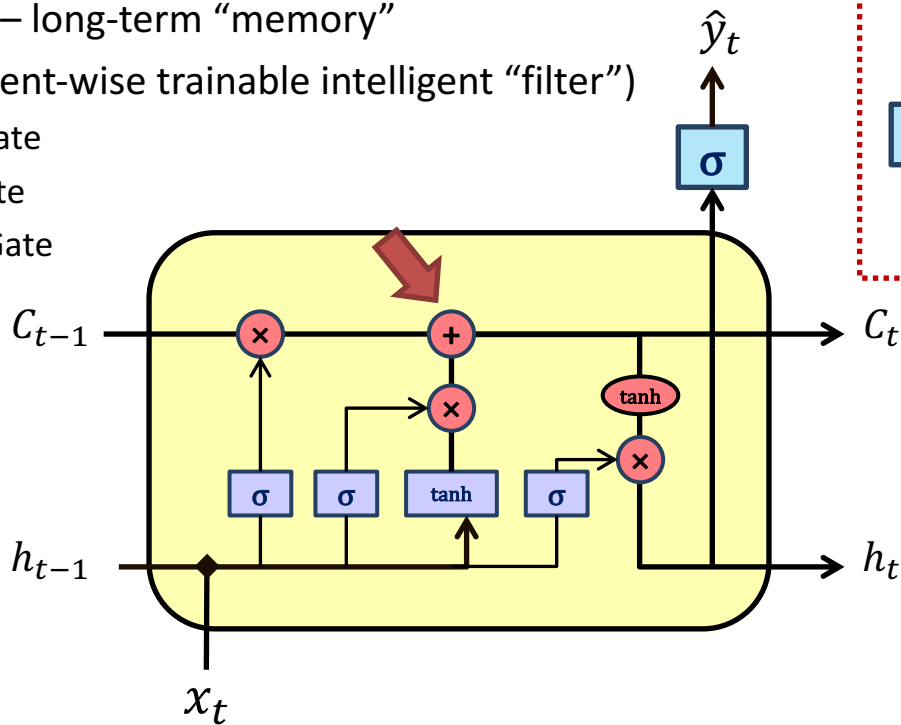


LSTM

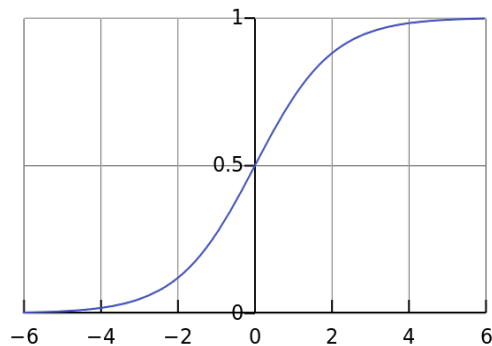
Hochreiter & Schmidhuber (1997)

- **Long Short Term Memory**

- Cell State C_t – long-term “memory”
- Gates: (element-wise trainable intelligent “filter”)
 - Forget Gate
 - Input Gate
 - Output Gate

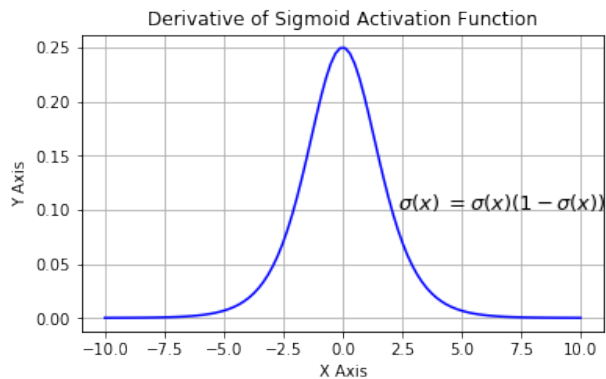


Sigmoid Function



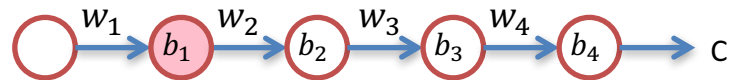
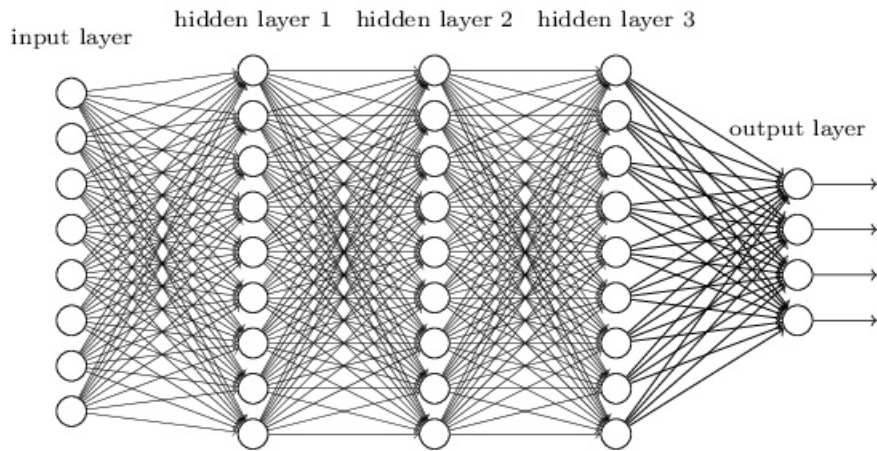
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

```
def sigmoid(x):  
    return 1.0 / (1.0 + math.exp(-x))
```



$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

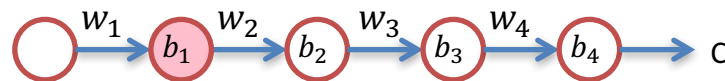
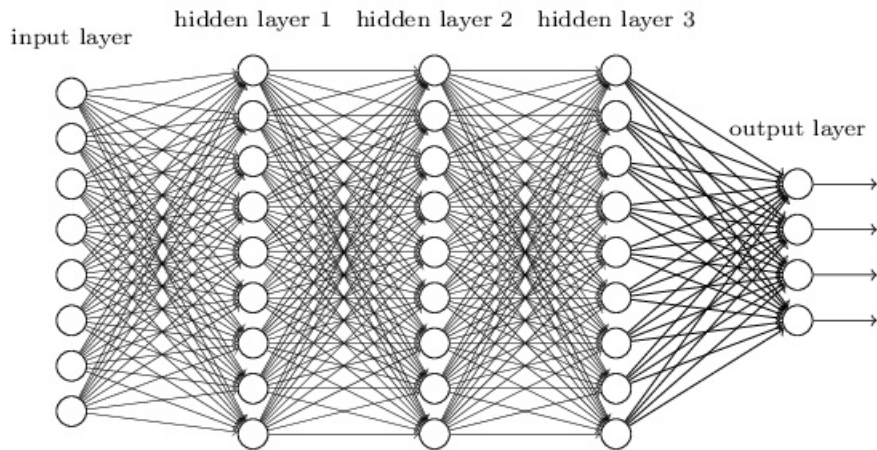
Chaining



$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \\ &= \frac{\partial C}{\partial y_4} \end{aligned}$$

Chaining

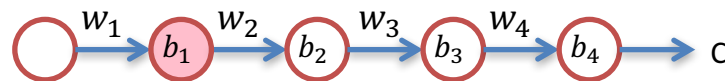
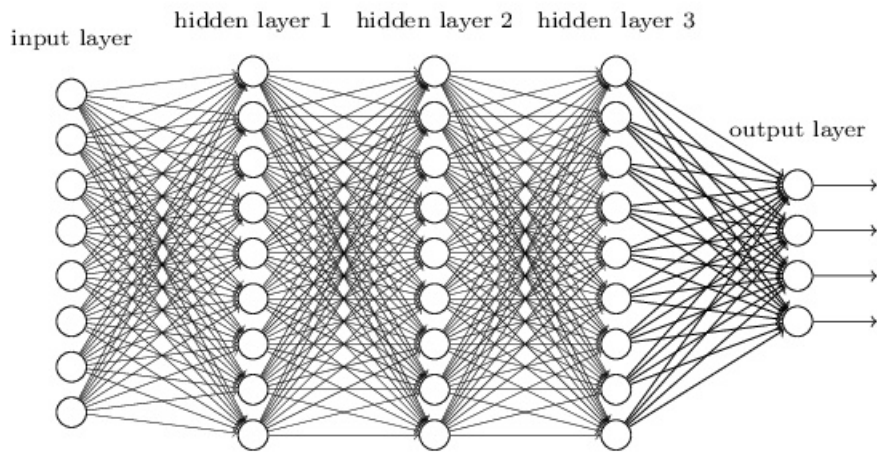


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) \end{aligned}$$

$$y_4 = \sigma(z_4)$$

Chaining

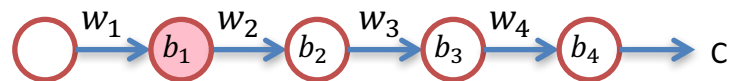
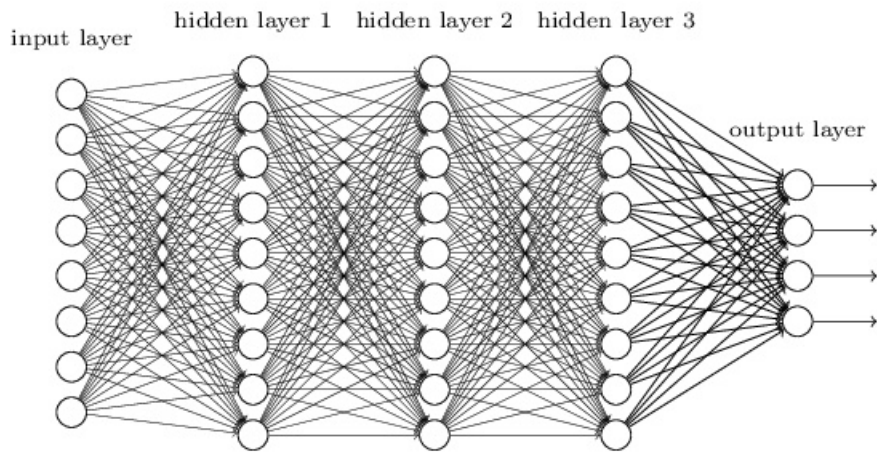


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \end{aligned}$$

$$\frac{\partial z_4}{\partial x_4} = \frac{\partial (w_4 x_4 + b_4)}{\partial x_4}$$

Chaining

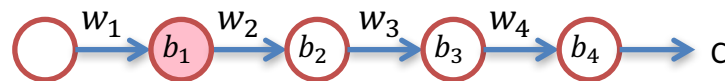
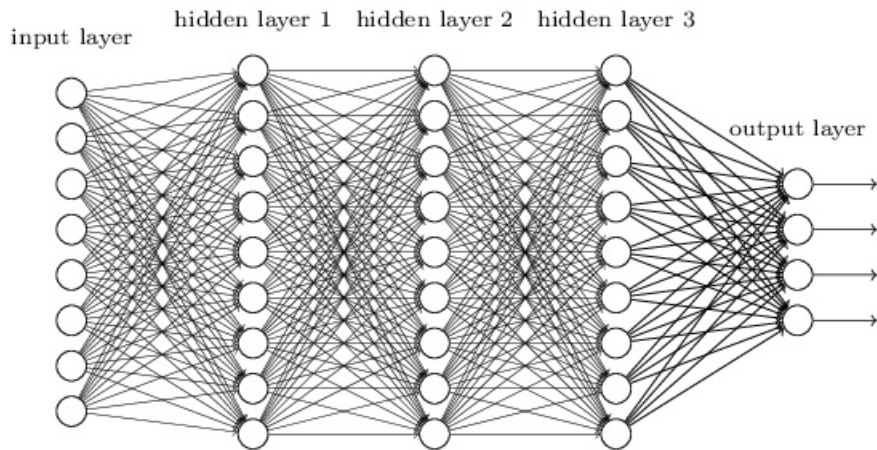


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) \end{aligned}$$

$$x_4 = \sigma(z_3)$$

Chaining

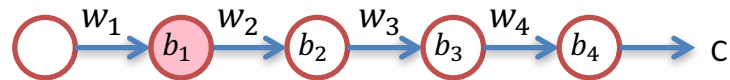
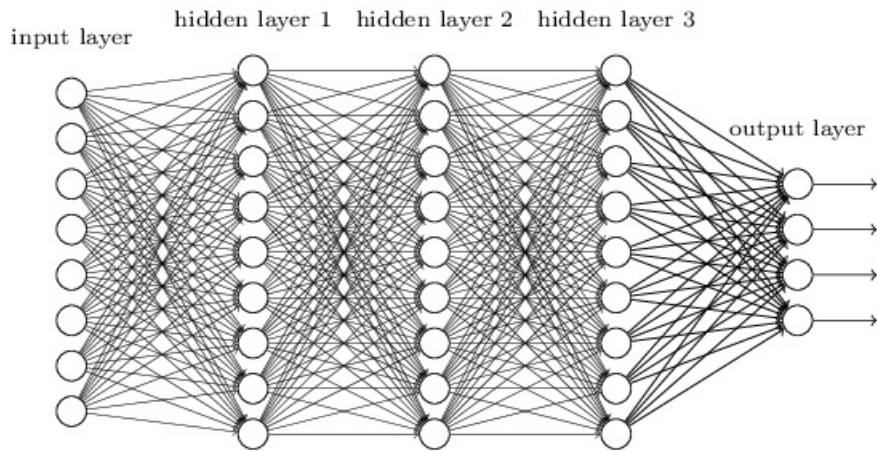


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \end{aligned}$$

$$\frac{\partial z_3}{\partial x_3} = \frac{\partial (w_3 x_3 + b_3)}{\partial x_3}$$

Chaining

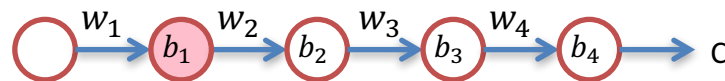
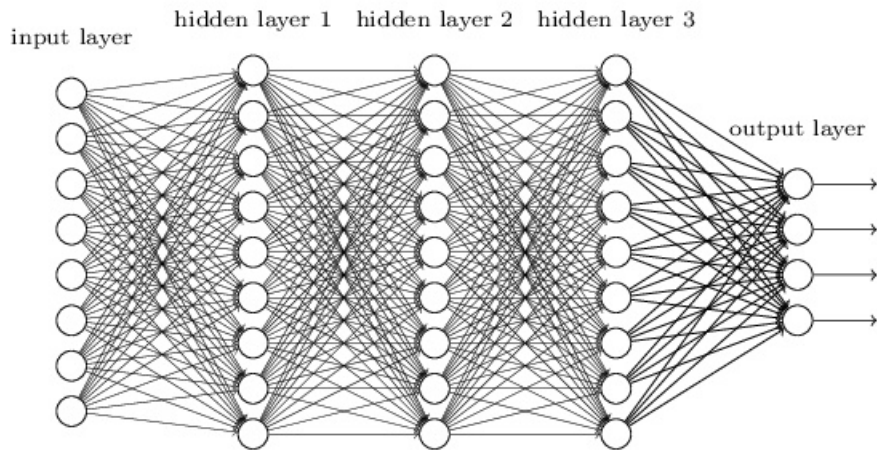


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) \end{aligned}$$

$$x_3 = \sigma(z_2)$$

Chaining

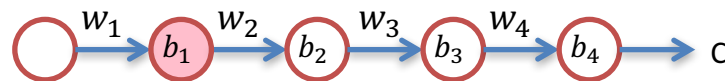
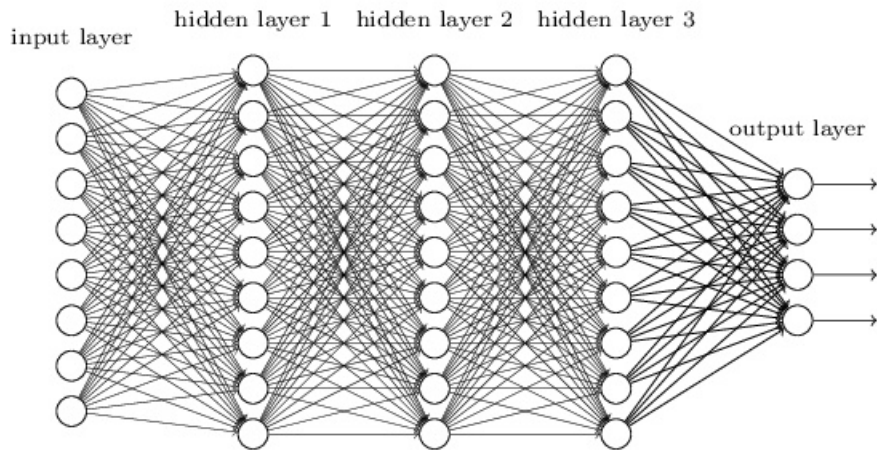


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \end{aligned}$$

$$\frac{\partial z_2}{\partial x_2} = \frac{\partial (w_2 x_2 + b_2)}{\partial x_2}$$

Chaining

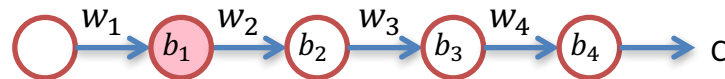
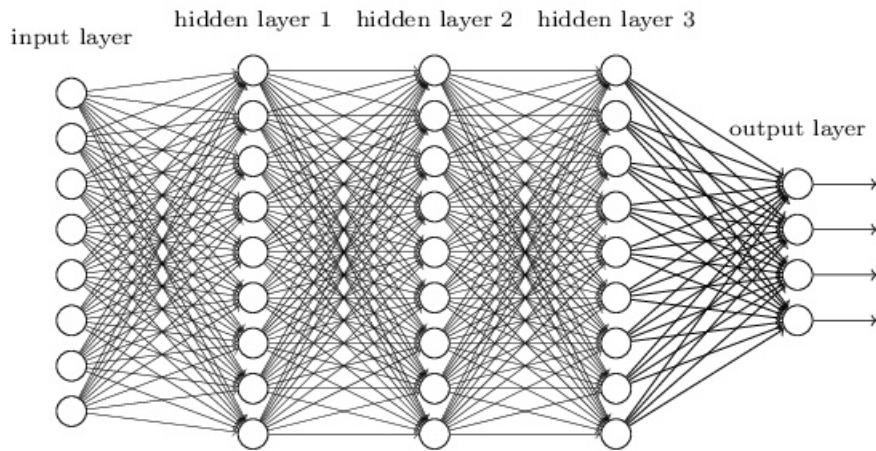


$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \frac{\partial x_2}{\partial z_1} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) \end{aligned}$$

$$x_2 = \sigma(z_1)$$

Chaining



$$y = x_{i+1} = \sigma(z_i) = \sigma(w_i x_i + b_i), \text{ where } z_i = w_i x_i + b_i$$

$$\begin{aligned} \frac{\partial C}{\partial b_1} &= \frac{\partial C}{\partial y_4} \frac{\partial y_4}{\partial z_4} \frac{\partial z_4}{\partial x_4} \frac{\partial x_4}{\partial z_3} \frac{\partial z_3}{\partial x_3} \frac{\partial x_3}{\partial z_2} \frac{\partial z_2}{\partial x_2} \frac{\partial x_2}{\partial z_1} \frac{\partial z_1}{\partial b_1} \\ &= \frac{\partial C}{\partial y_4} \sigma'(z_4) w_4 \sigma'(z_3) w_3 \sigma'(z_2) w_2 \sigma'(z_1) * 1 \end{aligned}$$

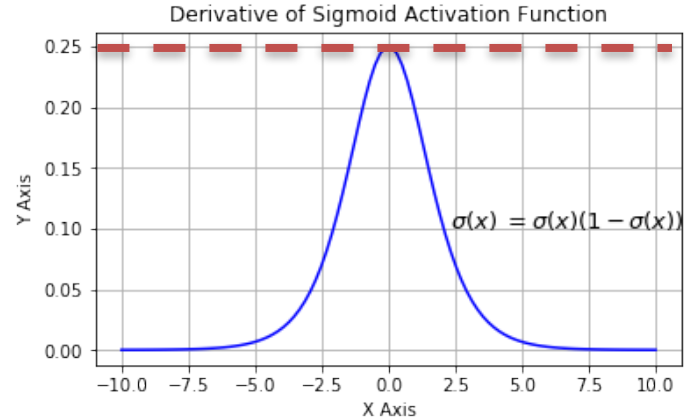
$$\frac{\partial z_1}{\partial b_1} = \frac{\partial (w_1 x_1 + b_1)}{\partial b_1}$$

Gradient Vanishing Problem



- Maximum of σ' is 0.25 and $|w_{init}| < 1$,
- Thus $|w \sigma'(z_i)| < 0.25$

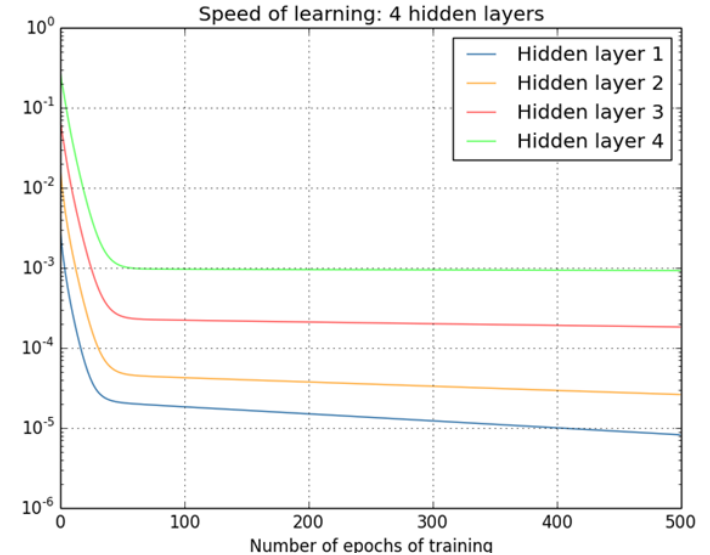
$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial y_4} \underbrace{\sigma'(z_4)w_4}_{<0.25} \underbrace{\sigma'(z_3)w_3}_{<0.25} \underbrace{\sigma'(z_2)w_2}_{<0.25} \sigma'(z_1)$$



Gradient Vanishing Problem

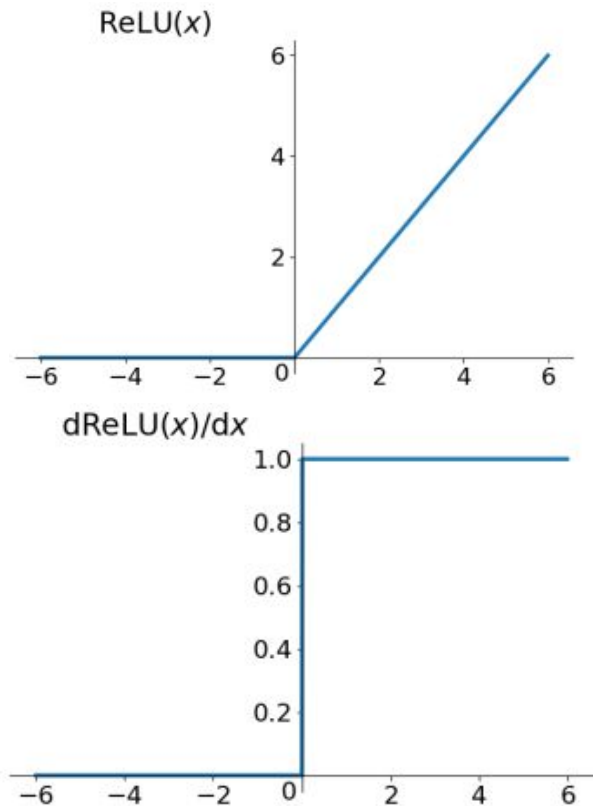


- The more front of the layer, the slower the gradient changes
- We were not able to train a very deep network
- How to solve it?



ReLU (Rectified Linear Unit)

Rectified linear units improve restricted boltzmann machines, Nair & Hinton (2010)

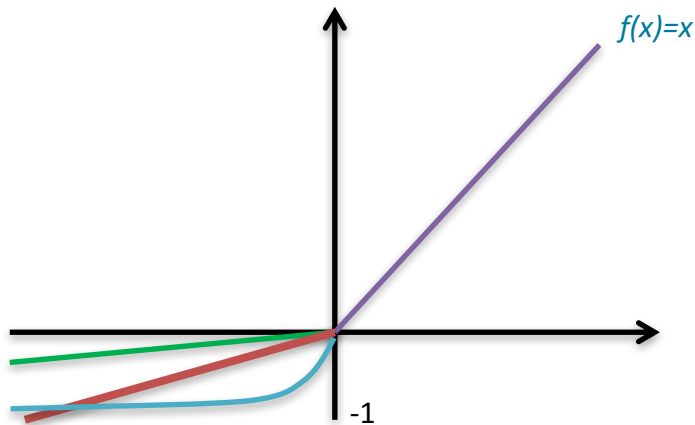


$$f(x) = \max(x, 0)$$

```
# ReLU
def relu(x):
    return np.maximum(x, 0)
```

$$f'(x) = \begin{cases} 0, & \text{if } x < 0 \\ 1, & \text{if } x \geq 0 \end{cases}$$

Leaky ReLU (LReLU), Parametric ReLU (PReLU), Exponential Linear Unit (ELU)



Leaky ReLU:

$$f(x) = \begin{cases} 0.01x, & \text{if } x < 0 \\ x, & \text{if } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} 0.01, & \text{if } x < 0 \\ 1, & \text{if } x \geq 0 \end{cases}$$

PReLU:

$$f(x) = \begin{cases} ax, & \text{if } x < 0 \\ x, & \text{if } x \geq 0 \end{cases}$$

$$f'(x) = \begin{cases} a, & \text{if } x < 0 \\ 1, & \text{if } x \geq 0 \end{cases}$$


ELU:

$$f(x) = \begin{cases} x, & \text{if } x \geq 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \end{cases}$$

$$f'(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ f(x) + \alpha, & \text{if } x < 0 \end{cases}$$


ImageNet Challenge 2012

Given an image, classify what is depicted



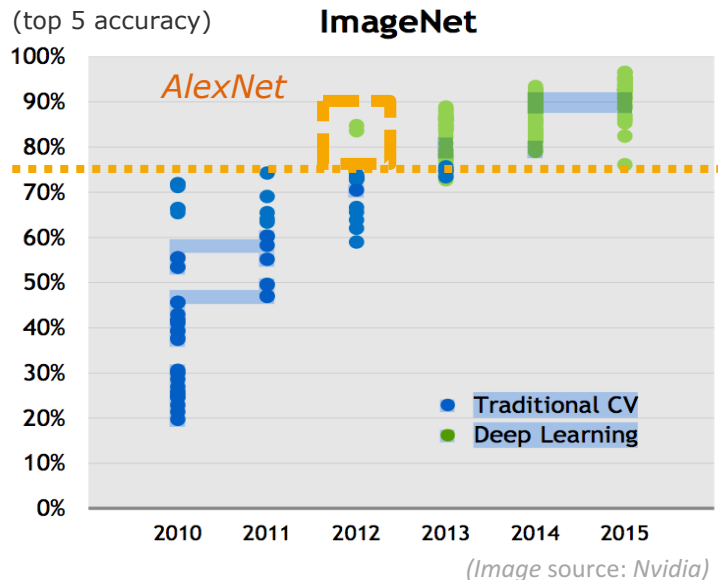
IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images



Output: Scale T-shirt <u>Steel drum</u> Drumstick Mud turtle	✓	Output: Scale T-shirt Giant panda Drumstick Mud turtle	✗
--	---	--	---

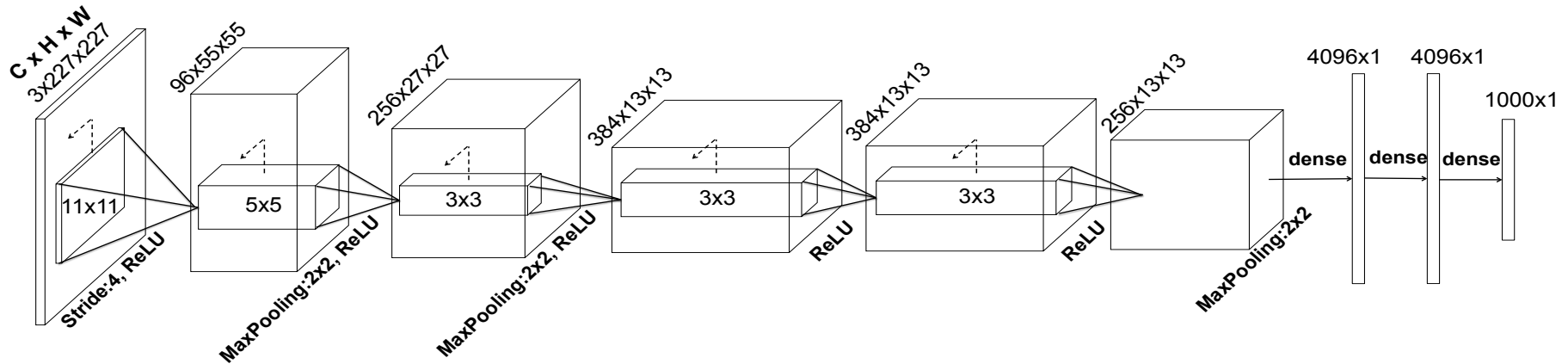
(Deng, J. et al. 2014)



AlexNet

Imagenet classification with deep convolutional neural network, Krizhevsky, Sutskever, Hinton. (NIPS'12)

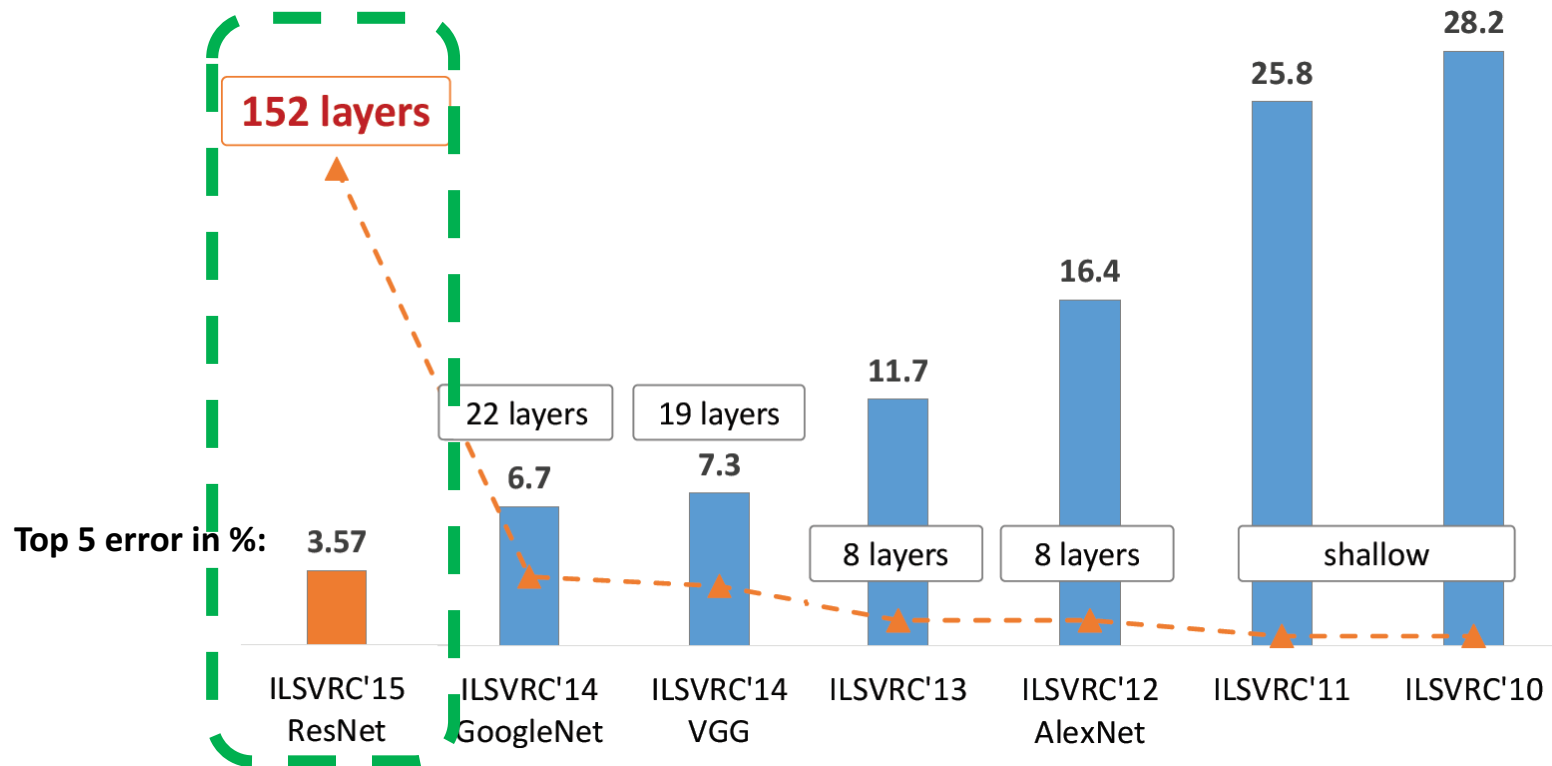
- 62.3 million parameters
 - Conv: 3.7 million (6%) , FC: 58.6 million (94%)
- Computation: Conv (95%) , FC (5%)
- The network takes 90 epochs in **6** days to train on two GTX 580 GPUs



A close-up shot from the movie Inception showing Leonardo DiCaprio and Matt Damon. DiCaprio is on the left, looking slightly down and to the right with a serious expression. Damon is on the right, seen in profile, looking towards DiCaprio. The lighting is dramatic, with strong highlights and deep shadows. The background is blurred, showing what appears to be a window or architectural structure.

WE NEED TO GO DEEPER

ImageNet Challenge Winners



ResNet

Deep residual learning for image recognition, He et al. (CVPR'16, best paper award)

- What happens when we continue stacking layers on a “plain” convolutional neural network like AlexNet?

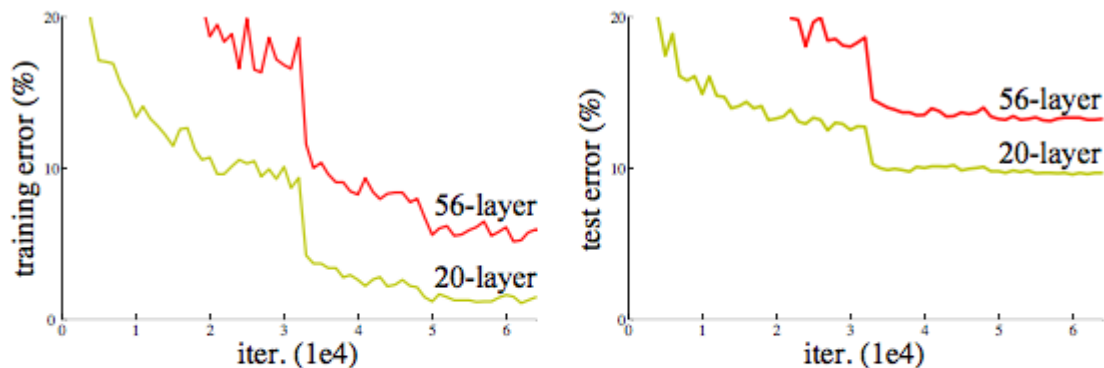


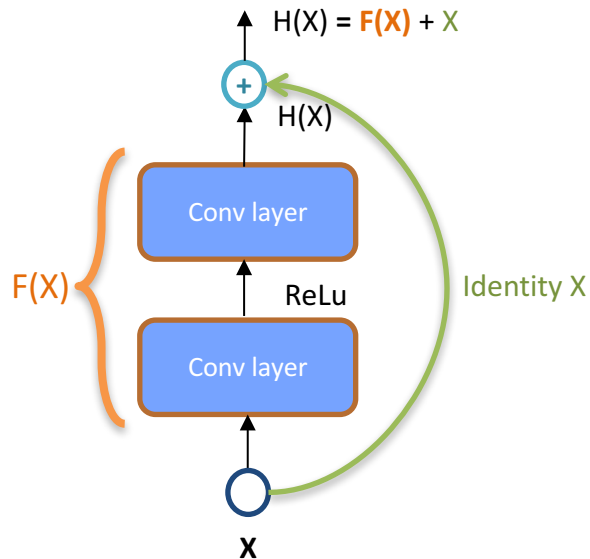
Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

- The deeper model performs worse, overfitting?
- **Network degradation**

ResNet

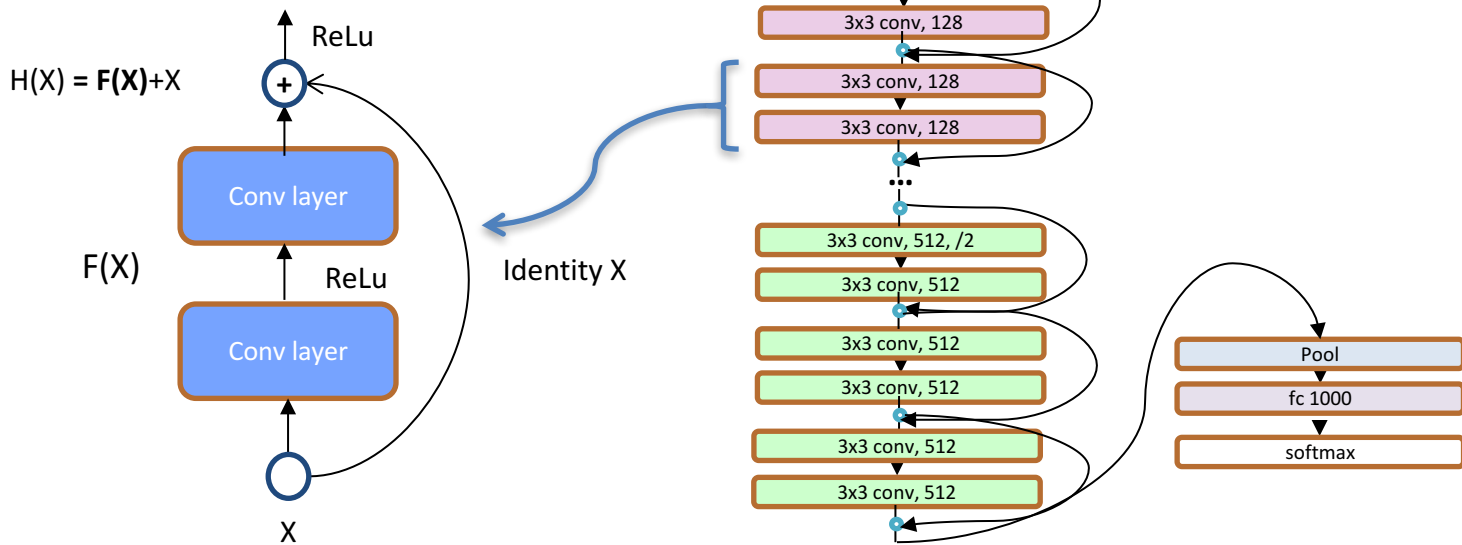
Deep residual learning for image recognition, He et al. (CVPR'16)

- Intuition: learn $F(X)$ not $H(X)$ in the **residual block**
 - If $F(X) = 0$, then $H(X)$ is an identity mapping
 - Use layers in residual block to learn the residual function $F(X) = H(X) - X$ instead of learning $H(X)$
 - Learning $F(X)$ should be easier
 - **Better gradient flow**



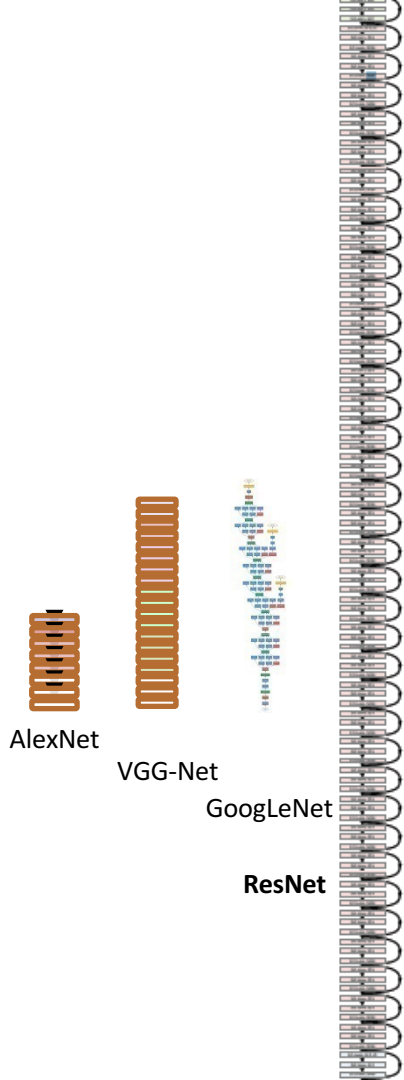
ResNet

- Stacked residual blocks
- At least two conv layers in each residual block
- Periodically, double filters and downsample spatially using stride 2



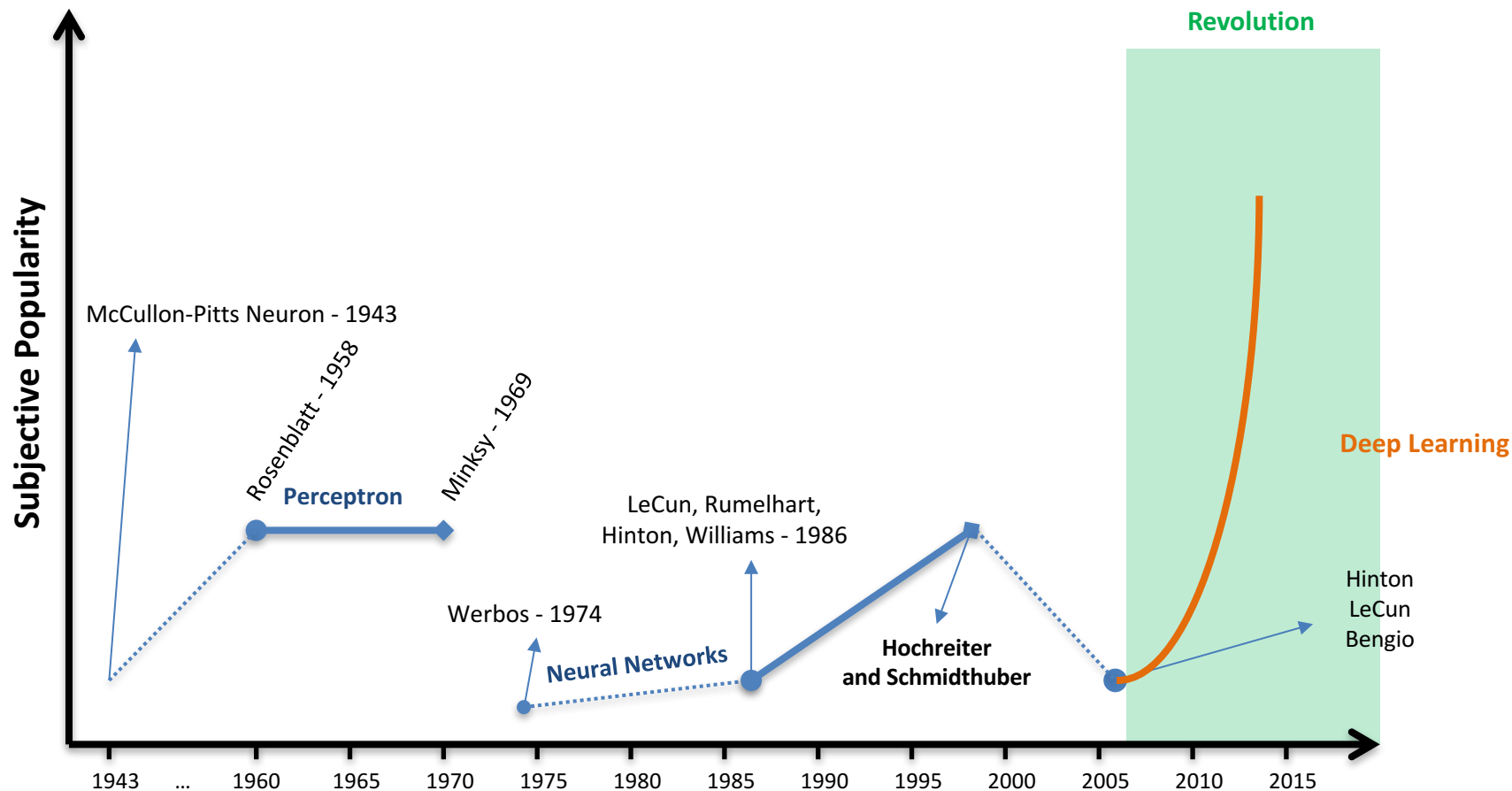
ResNet

- 152-layer on ImageNet, 1202-layer on CIFAR data set
- Swept 1st place in all ILSVRC and COCO 2015 competitions
 - Including *image classification*, *object detection* and *object segmentation*
- ILSVRC'15 classification winner (**3.57%** top 5 error) → better than “human performance!” (**5.1%** reported in *Russakovsky 2014*)
- Contribution
 - Preventing network degradation problem with residual connections
 - We are able to train extreme deeper networks

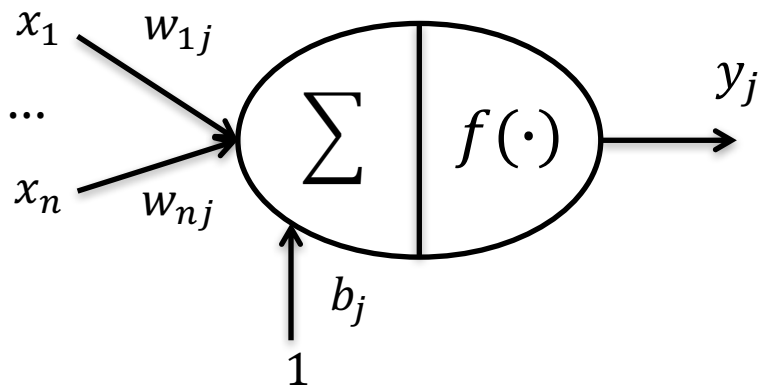


Recent Efforts

- **Recent efforts more based on Information/Gradient Flow**
 - Sigmoid saturation, gradient vanishing problem : **ReLU** (6003 citations)
 - If $x < 0$, $\text{ReLU} = 0$: **LeakyReLU, PReLU, ELU** etc. (5334 citations in total)
 - Network too “deep”, gradient vanishing : **highway** (parameterized shortcut connection between layers) (779 citations)
 - Simplified highway : **ResNet** (shortcut connection without params) (**23397** citations)
 - Forced stability of the mean and variance of activations: **BatchNorm** (10717)
 - Adding noises in gradient flow : **Dropout** (12387)
 - RNN with gradient exploding and vanishing : **LSTM** (adding more gated controls) (18907)
 - Simplified LSTM : **GRU** (gated recurrent unit) (2664)
 - etc.

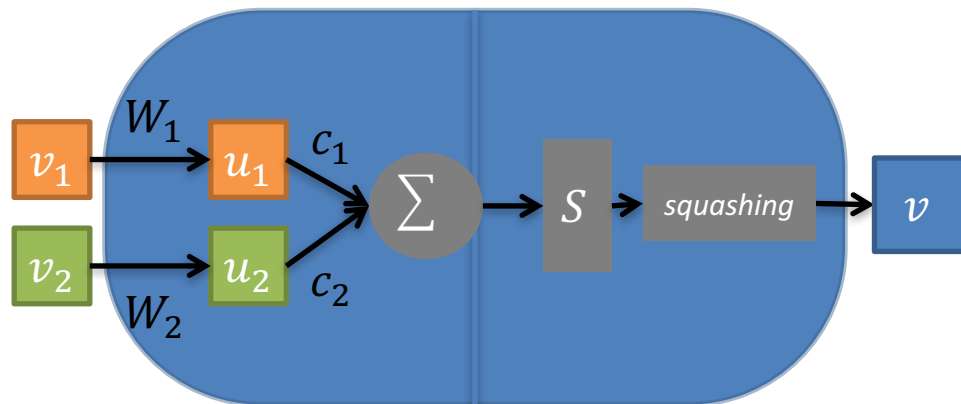


Next week: Capsule vs. Classical NN



Classical neuron: **scalar to scalar**

$$y_j = \sum_{i=1}^n w_{ij} x_i + b_j$$

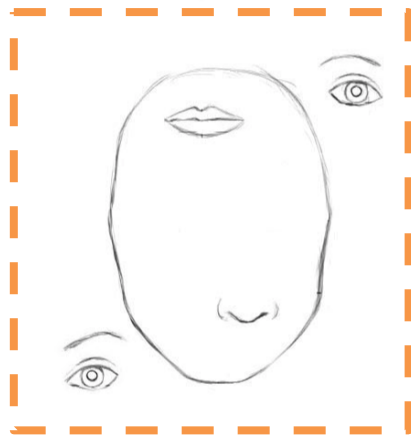
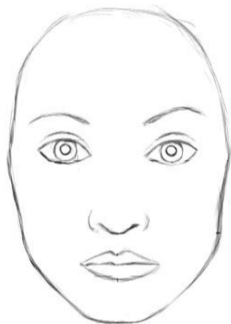


Capsule neuron: **vector to vector**

$$u_i = W_i v_i \quad s = \sum_{i=1}^n c_i u_i$$

$$v = \text{Squashing}(s) = \frac{\|s\|^2}{1 + \|s\|^2} \frac{s}{\|s\|}$$

Limitations of ConvNets



person 0.88



reddish orange color 0.78



light brown color 0.78



starlet 0.66



entertainer 0.66



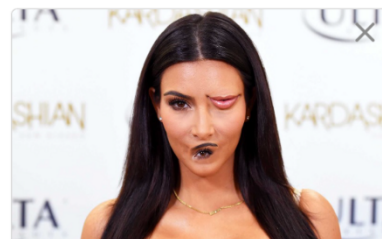
female 0.60



woman 0.59



young lady (heroine) 0.59



person 0.90



light brown color 0.84



starlet 0.77



entertainer 0.77



female 0.65



woman 0.64



young lady (heroine) 0.64



reddish orange color 0.64



newsreader 0.50



coal black color 0.79



hairpiece (hair) 0.71



dress 0.71



maroon color 0.71



person 0.58



toupee (hairpiece) 0.58



woman 0.56



Earrings 0.55

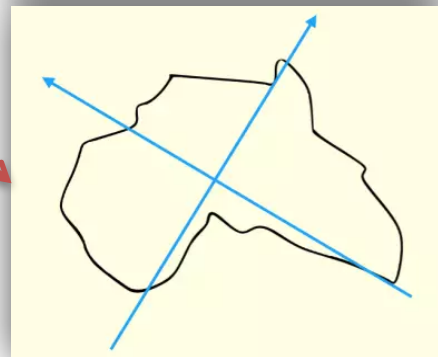
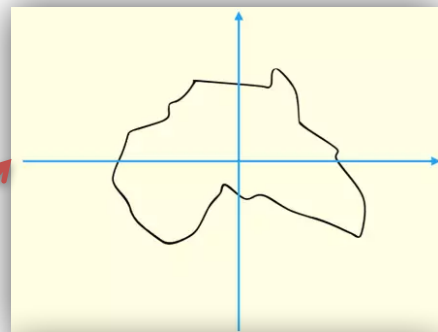
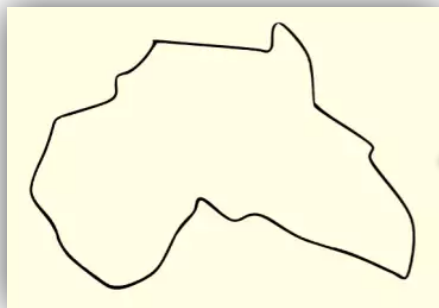
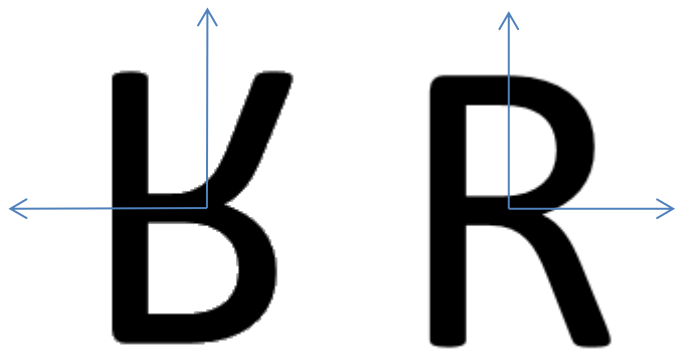


female 0.50



Example

- Human vision has **coordinate frame**, which influences the perception system.



Capsule

Sabour, Frosst & Hinton (NIPS'17, ICLR'18)

Keypoints

- Using vector (matrix) instead of scalar
- Encapsulate (visual) entity with its pattern
- Equivariance instead of Invariance:
 - $\text{Transform}(\text{Represent}(x)) = \text{Represent}(\text{Transform}(x))$
 - Invariance e.g., CNN: $\text{Represent}(x) = \text{Represent}(\text{Transform}(x))$
- Routing by agreement instead of pooling
- Using Matrix to represent the relationships of visual entities

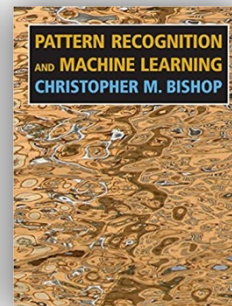
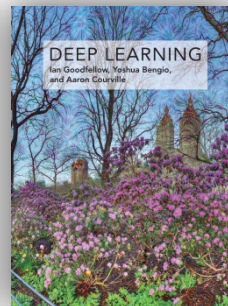
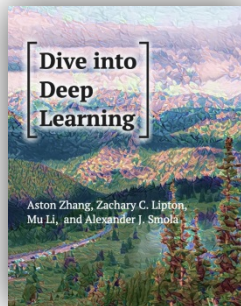
Literature


Books

- Alex Smola et al., [“Dive into deep learning”](#), 2019
- Ian J. Goodfellow and Yoshua Bengio and Aaron Courville, [“Deep Learning”](#), 2016
- Christopher M. Bishop “Pattern Recognition and Machine Learning” (google it)

Online courses:

- cs231n: Convolutional Neural Networks for Visual Recognition (Feifei Li et al.)
- Deep Learning courses at Coursera (Andrew Ng et al.)



A colorful, stylized illustration of a building with a green sign that says "HPI" and "Hochschule für Angewandte Wissenschaften". The building is yellow and orange, with green windows. The sign is green with a yellow square containing a stylized 'H' and the text "HPI" and "Hochschule für Angewandte Wissenschaften".

Thank you for your Attention!

Dr. Haojin Yang

Office: H-1.22

Email: haojin.yang@hpi.de