



LUND
UNIVERSITY

Lecture 7

COGNITIVE COMPUTING

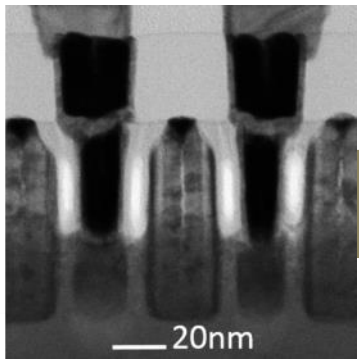


- Cognitive computing
 - Why and what is it?
- Neural network structure
 - The brain
 - Neuron models
 - Network models
- Learning in Neural networks
 - Gradient descent and backpropagation
 - Deep learning
- Hardware for Neural networks

- Hand-in assignment
- Oral exams

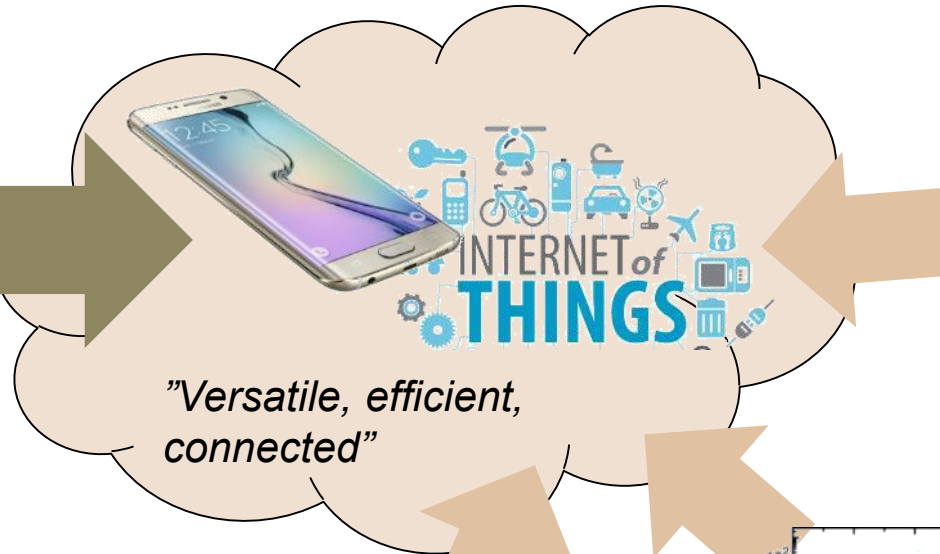
Recap on course

Traditional CMOS

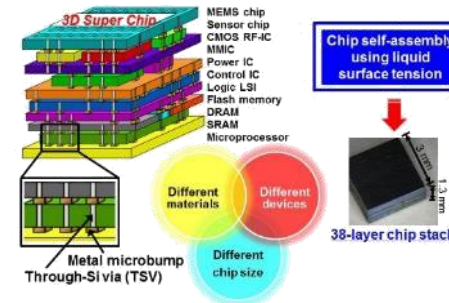


Dennard scaling
"smaller, faster, cheaper"

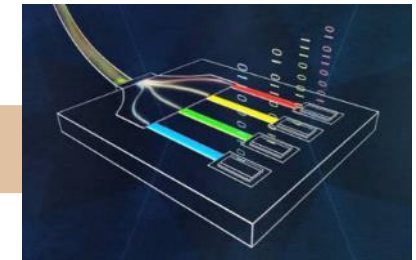
More than Moore applications



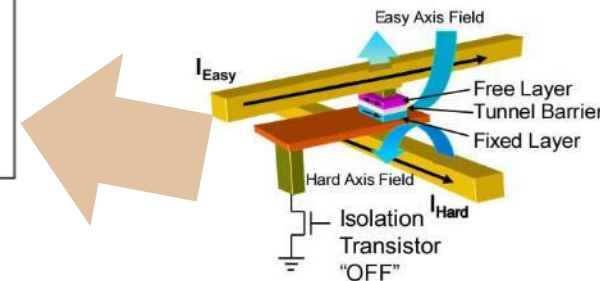
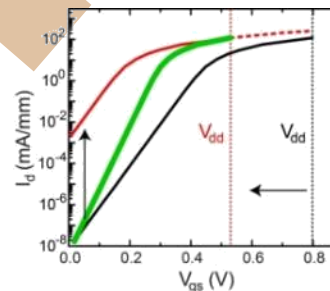
Integrated sensors



3D circuits and heterogeneous materials



Optical interconnects



New types of Memory

Why cognitive computing?

- Peta-bytes of constant data generation
 - Mostly unstructured
- Ubiquitous computing power



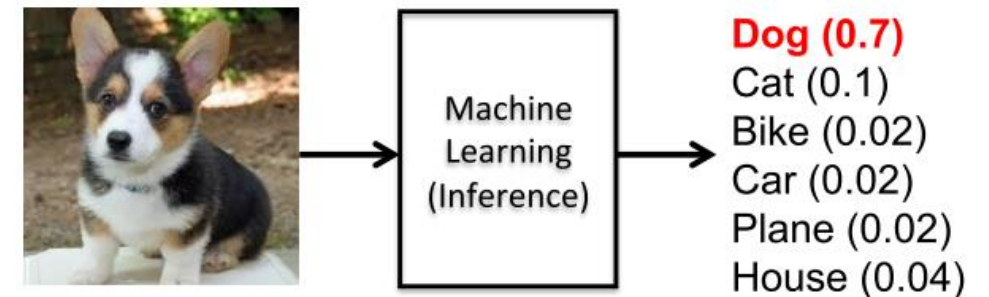
How to handle the data flow and utilize this power?

We need cognitive systems that:

1. Identify data patterns
2. Identifies anomalies in data
3. Can find optimal decision based on data

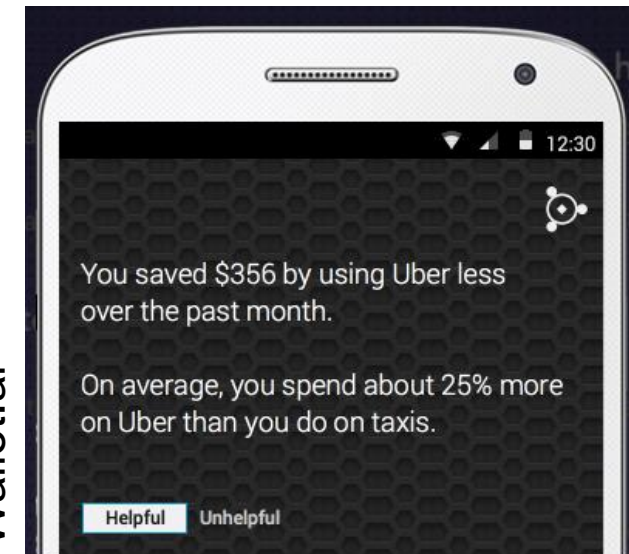
Pattern recognition

- **Image recognition**
 - Computer vision → data from images, surveillance, autonomous cars,...
- **Speech recognition**
 - Personal assistants, surveillance
- **Language**
 - Data mining, translation
- Unstructured → Structured data



Decision making

- **Health care**
 - Diagnosis
 - Identifying individual treatments for patients
 - Drugs: Choosing optimal trials to reduce time to market
- **Finance**
 - Stock analysis:
 - Analyse earnings statements, news reports and regulatory filings looking for clues on how to view a stock.
 - Personal Financial Advisor
 - Analyse what impacts personal economy and suggest solutions



Wallet.ai

Monitor processes and detect issues/trends early on and alert/make decision

- **Finding illegal behavior**
 - Observe traffic (monetary, wares, transactions)
- **Failures**
 - Monitor sensor status and detect irregularities before breakdown
- **Health care**
 - Monitoring of ECG/EEG signals can detect cardiovascular deceases, early onset of seizure in epilepsy



Who is working on it?

IBM (Watson)



Data mining
Health
IoT
Business
...



Nvidia
GPU AI Hardware...



Personal Assistant
Smarter search
Autonomous driving
...



Amazon Alexa
"Personal assistant"

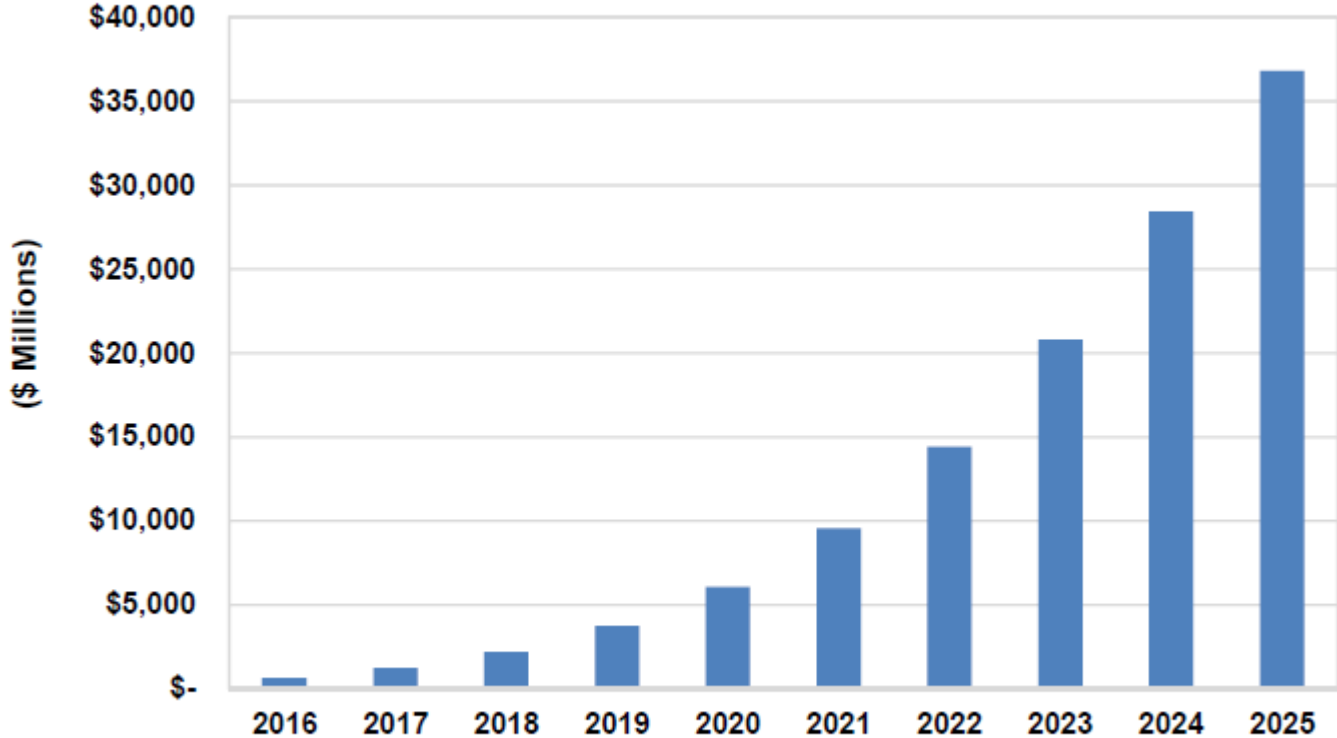


Visual, Speech, Text
recognition
Smarter search
Knowledge mining

Giant new industry

AI has applications and use cases in almost every industry vertical and is considered the **next big technological shift**, similar to past shifts like the industrial revolution, the computer age, and the smartphone revolution

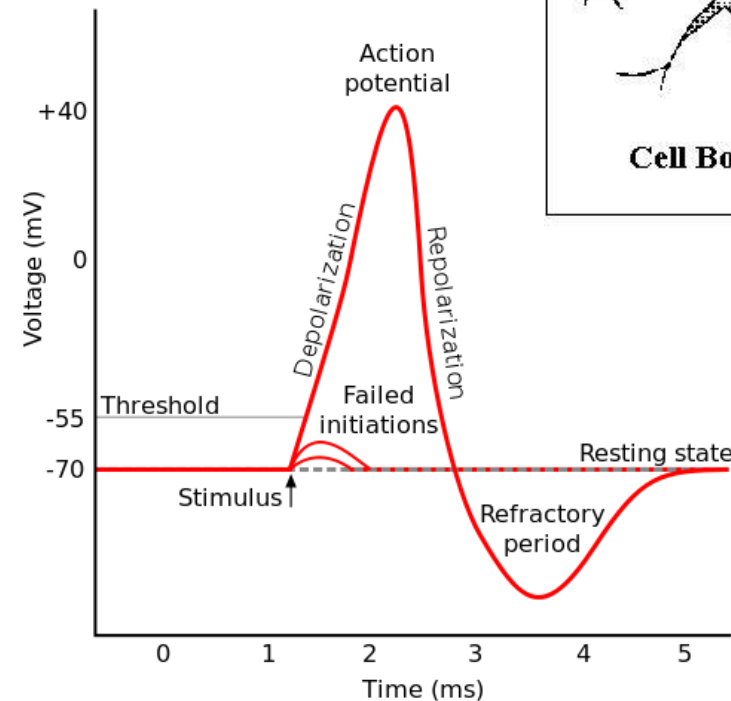
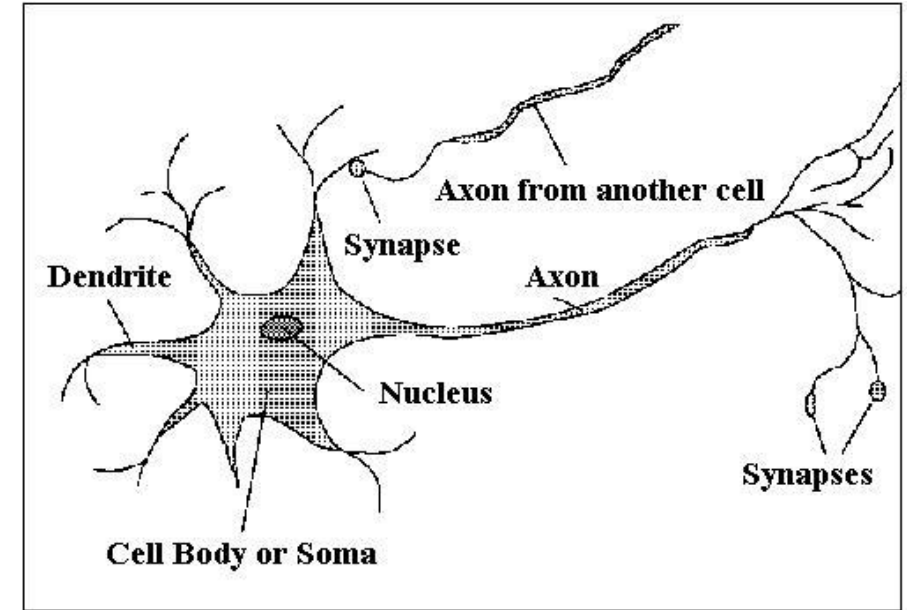
Chart 1.1 Artificial Intelligence Revenue, World Markets: 2016-2025



(Source: Tractica)

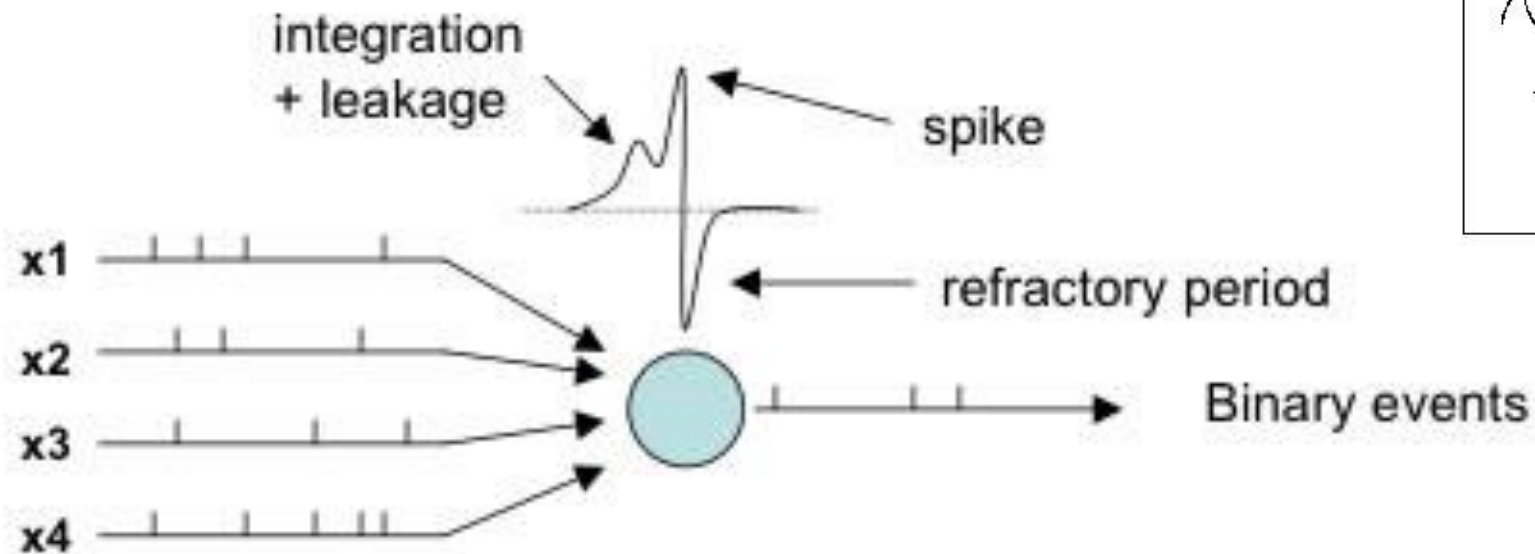
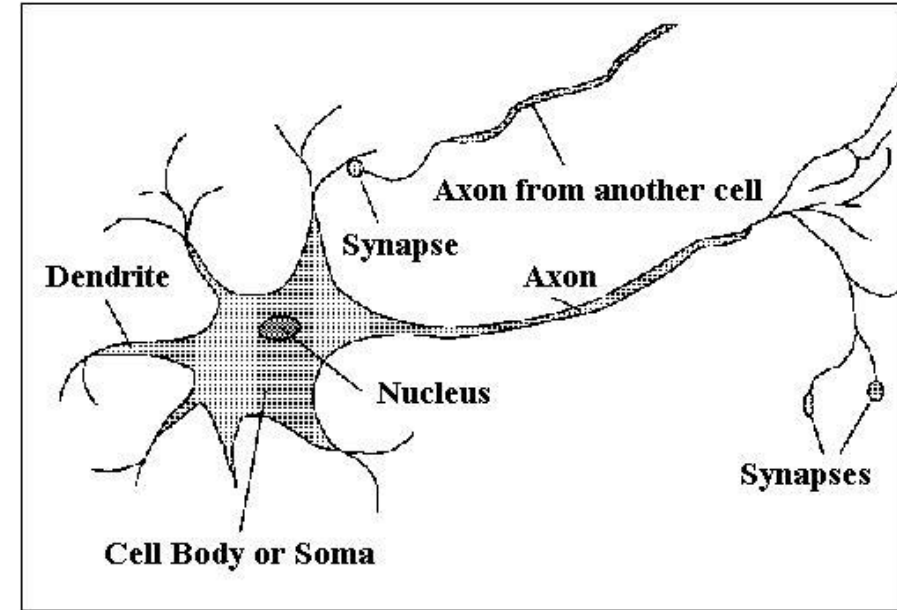
How the brain works ("much simplified")

- Neurons $\sim 10^{11}$
- Connected by synapses with varying "resistance"/"weight"
 $\sim 10^{15}$ synapses
- Electrical stimuli above threshold close in time causes them to fire a signal
- Signals propagate through network
- Connections encode logic/memory
- Hierarchical "layered" structure



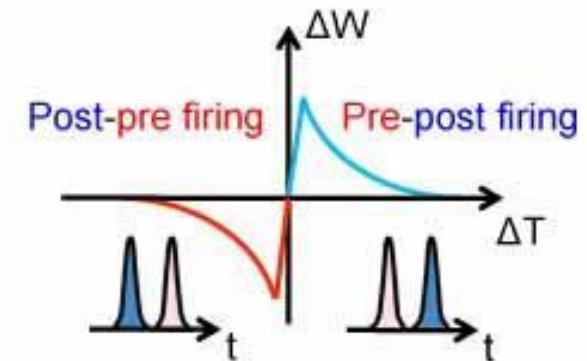
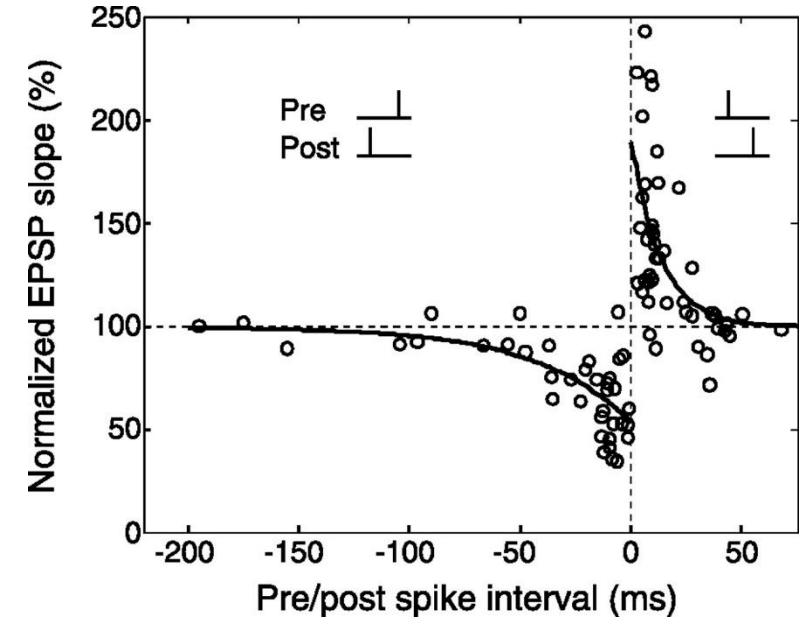
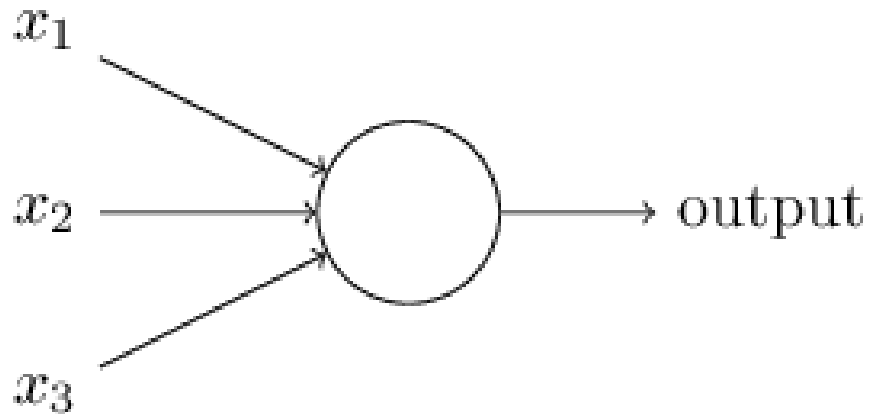
Spiking neurons

- Input signals in dendrites are integrated in the neuron
- Many inputs in short time interval
 → a threshold is overcome
 → the neuron will fire a signal into the axon



Spike-timing dependent plasticity

- How the brain "learns"
- Inputs prior to neuron fires are strengthened
- Inputs after the neuron fires are weakened
- Neurons tend to fire when many inputs occur at the same time (threshold) \rightarrow Subset of correlated inputs remain

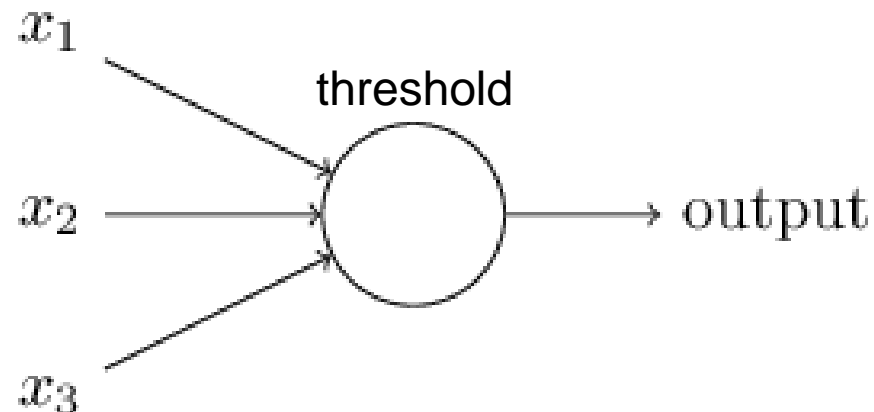


"Neurons that fire together wire together"

The perceptron

- A simplified representation of a neuron
 - A number of inputs, one output
 - Threshold determines fire/not.
 - Also called *bias*

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases}$$

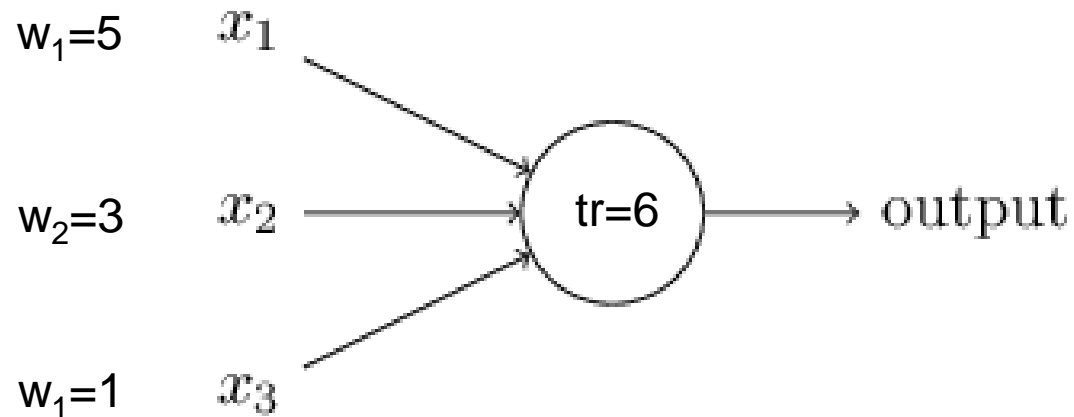


Perceptron example

Cheese festival in another town and you love cheese

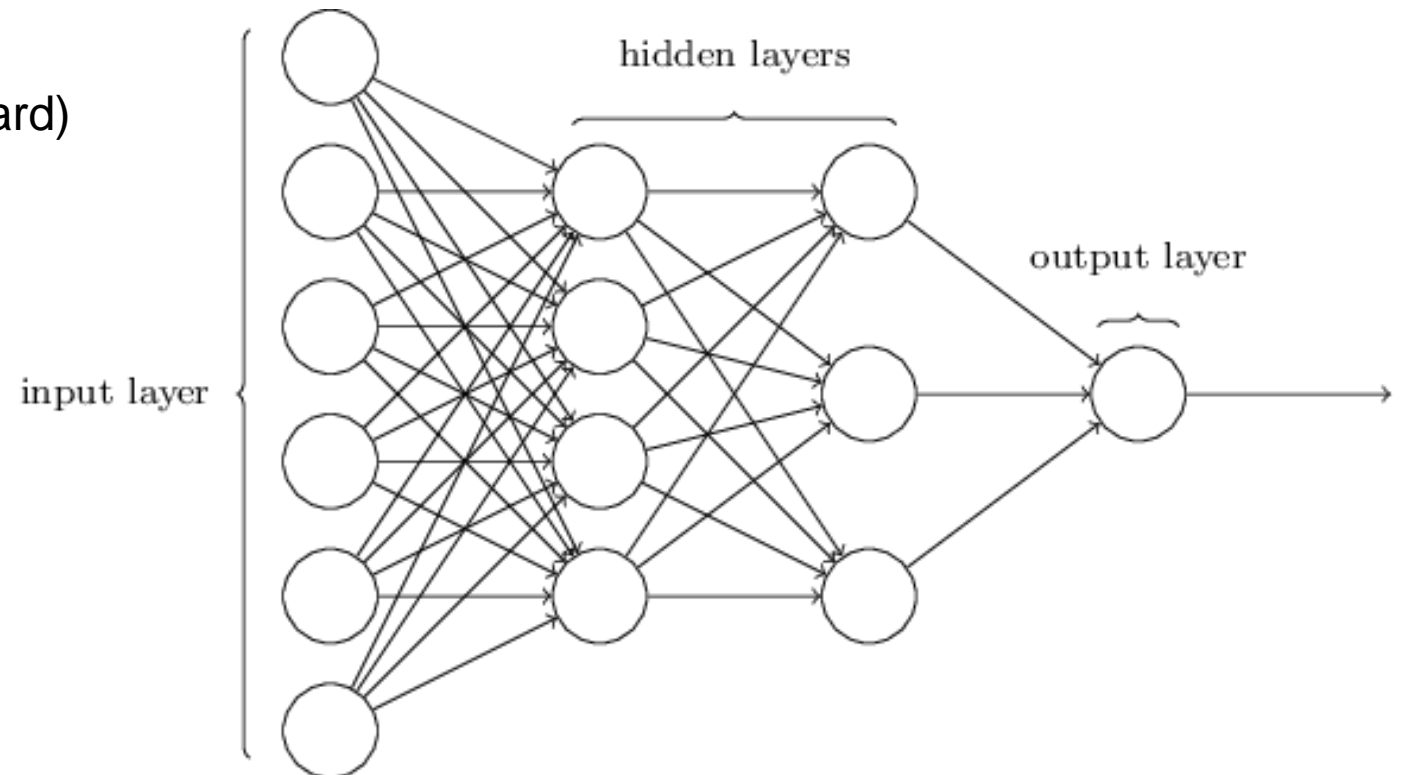
But should you go there? Three factors to decide:

1. Is the weather good?
2. Does your boyfriend/girlfriend want to join?
3. Is it easy to get there?



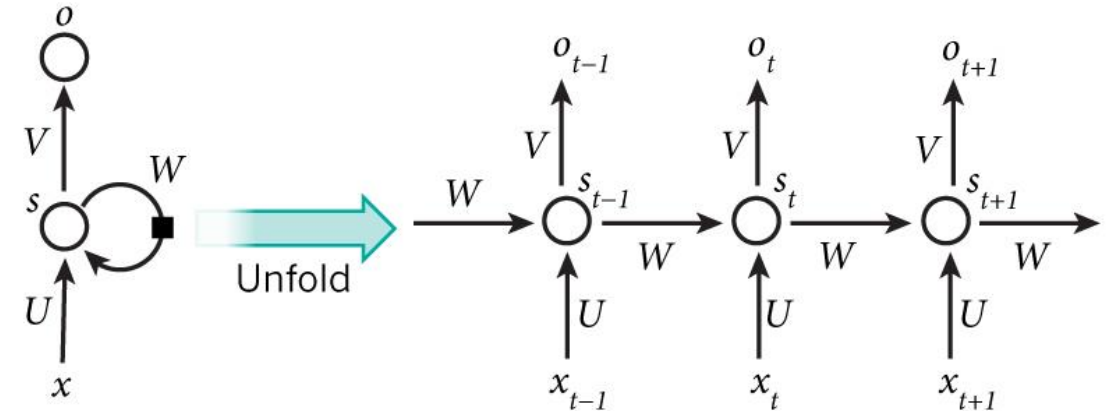
Layered network

- First layer – Input layer
- Last layer – Output layer
- Intermediate layers – Hidden layers
 - For more subtle and sophisticated "decision making"
- Connections usually go forward (feed-forward)

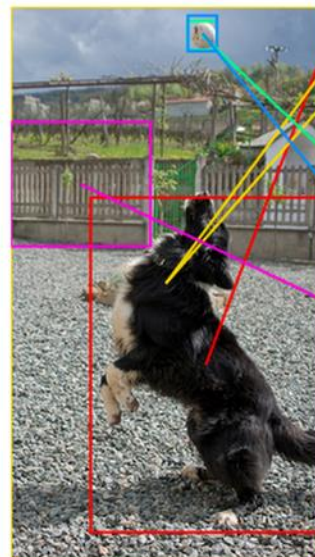


Recurrent neural networks

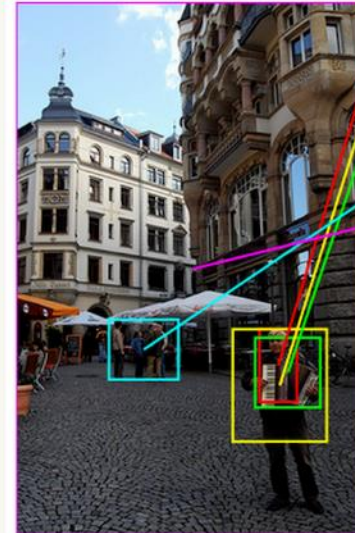
- Allowing for connections within and to previous layers
- The network obtains "memory" of previous states
- Useful for analysing flows of temporally connected data
 - Speech, language, ...
- Best implemented with spiking neurons



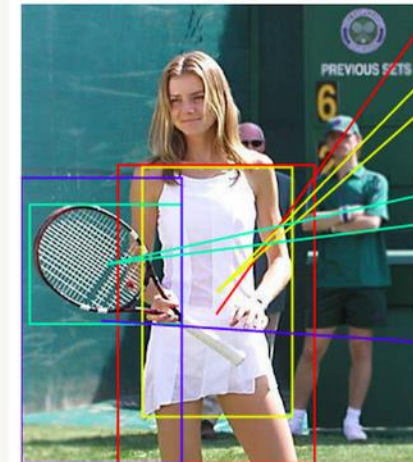
Generation of image caption text



- 1.31 dog
- 0.31 plays
- 0.45 catch
- 0.02 with
- 0.25 white
- 1.62 ball
- 0.10 near
- 0.07 wooden
- 0.22 fence



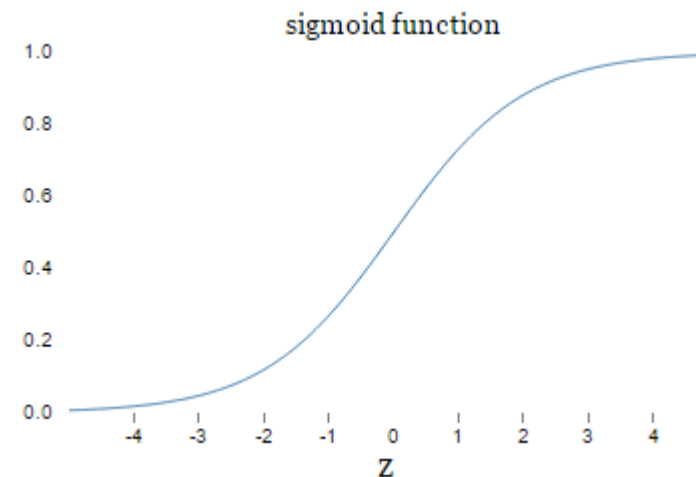
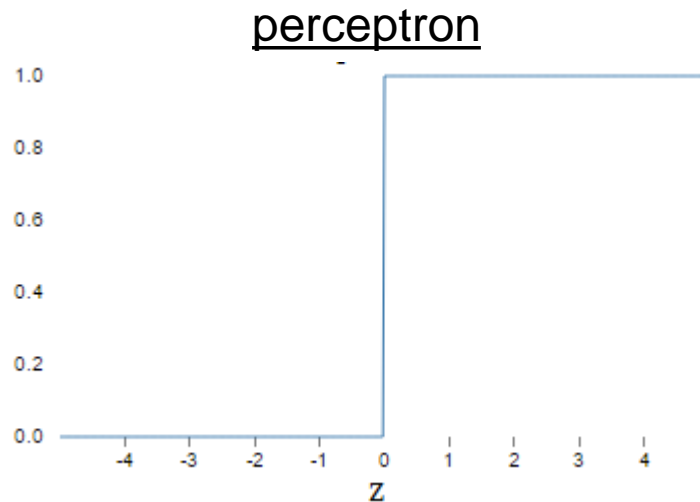
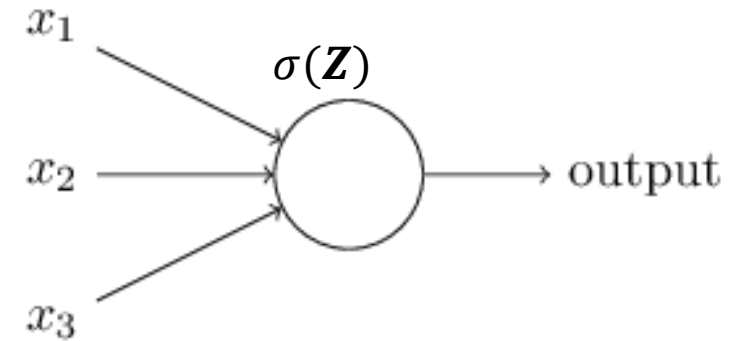
- 0.26 man
- 0.31 playing
- 1.51 accordion
- 0.07 among
- 0.08 in
- 0.42 public
- 0.30 area



- 1.12 woman
- 0.28 in
- 1.23 white
- 1.45 dress
- 0.06 standing
- 0.13 with
- 3.58 tennis
- 1.81 racket
- 0.06 two
- 0.05 people
- 0.14 in
- 0.30 green
- 0.09 behind
- 0.14 her

Sigmoid neuron

- Smooth activation → To facilitate learning algorithms
- Sigmoid function:
- $\sigma(\mathbf{Z}) = \frac{1}{1+\exp(-\mathbf{Z})}$



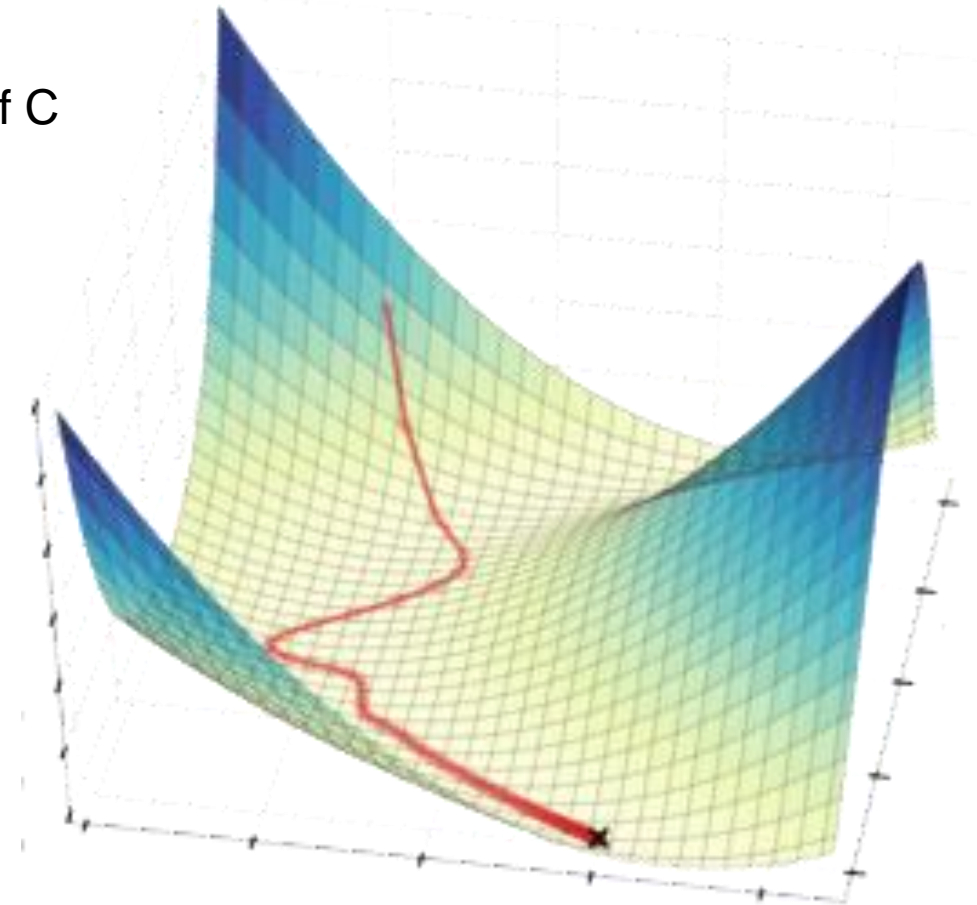
- Idea:
 - Feed network with known examples x that should give output $y(x)$
 - Look at the outcome and define how good was the result
 - The cost function
A smooth function in terms of weights and thresholds to evaluate how well trained the network is
 - Adjust the weights (\mathbf{w}) and biases (\mathbf{b}) to improve
 - Repeat until good enough
- Quadratic cost function
 - $C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2$
 - $y(x)$ – correct output, a – output of network, n – nbr of training inputs

Gradient Descent Algorithm

- A general method to minimize a multivariable function (cost function)
- Change parameters v in the direction of maximum gradient of C
- $\nabla C = \left(\frac{\partial C}{\partial v_1}, \dots, \frac{\partial C}{\partial v_m} \right)^T$
- $\Delta v = -\eta \nabla C$, η is the *learning rate*

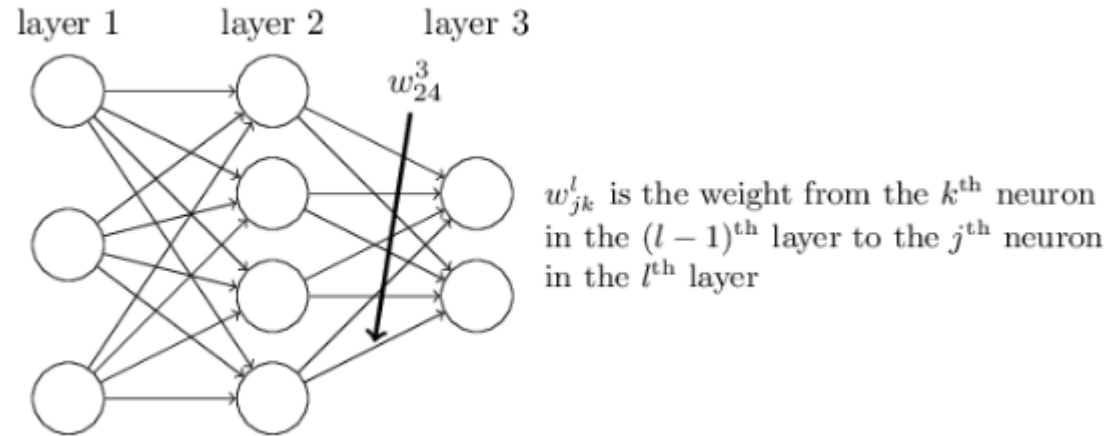
$$w_k \rightarrow w'_k = w_k - \eta \frac{\partial C}{\partial w_k}$$

$$b_l \rightarrow b'_l = b_l - \eta \frac{\partial C}{\partial b_l}$$



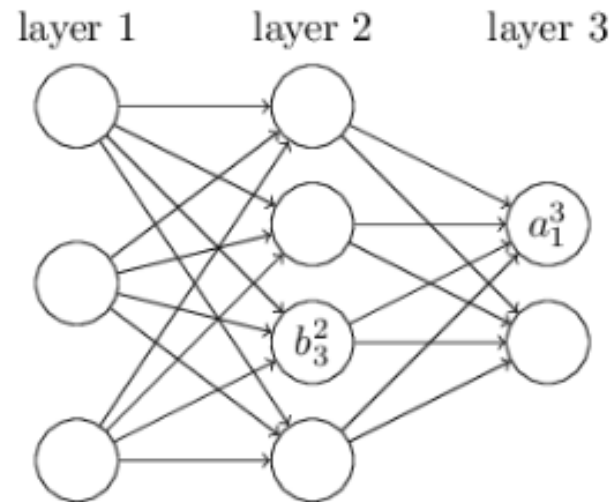
Matrix representation

$$\mathbf{W}^l = \begin{bmatrix} W_{11} & \cdots & W_{1,n} \\ \vdots & \ddots & \vdots \\ W_{m,1} & \cdots & W_{m,n} \end{bmatrix}$$



$$\mathbf{A}^l = \begin{bmatrix} a_1^l \\ \vdots \\ a_m^l \end{bmatrix}$$

$$\mathbf{B}^l = \begin{bmatrix} b_1^l \\ \vdots \\ b_m^l \end{bmatrix}$$



The weighted input of the l^{th} layer:

$$\mathbf{Z}^l = \mathbf{W}^l \mathbf{A}^{l-1} + \mathbf{B}^l$$

The activation of the l^{th} layer is thus:

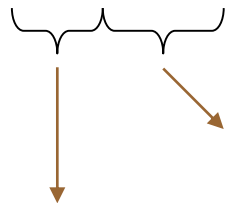
$$\mathbf{A}^l = \sigma(\mathbf{Z}^l)$$

The error function

- A way to quantify the error of the estimation due to a change in the network parameters

- The error in the output layer:

- $\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L) \rightarrow \Delta^L = \nabla_a C .* \sigma'(\mathbf{Z}^L)$

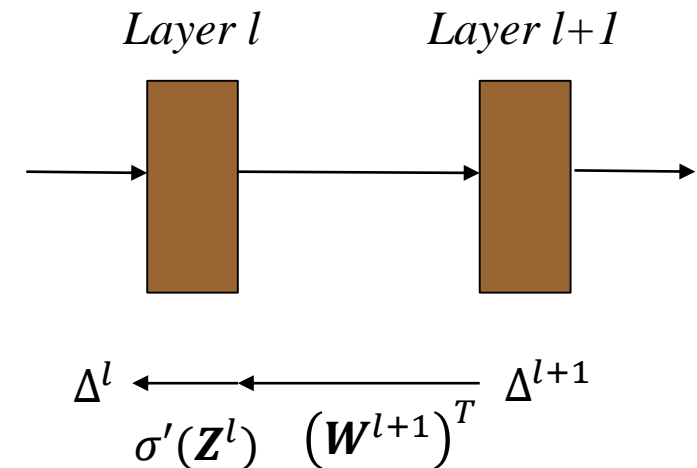


How fast the activation changes at z_j

How fast the cost is changing with change of activation at j

- For the previous layers:

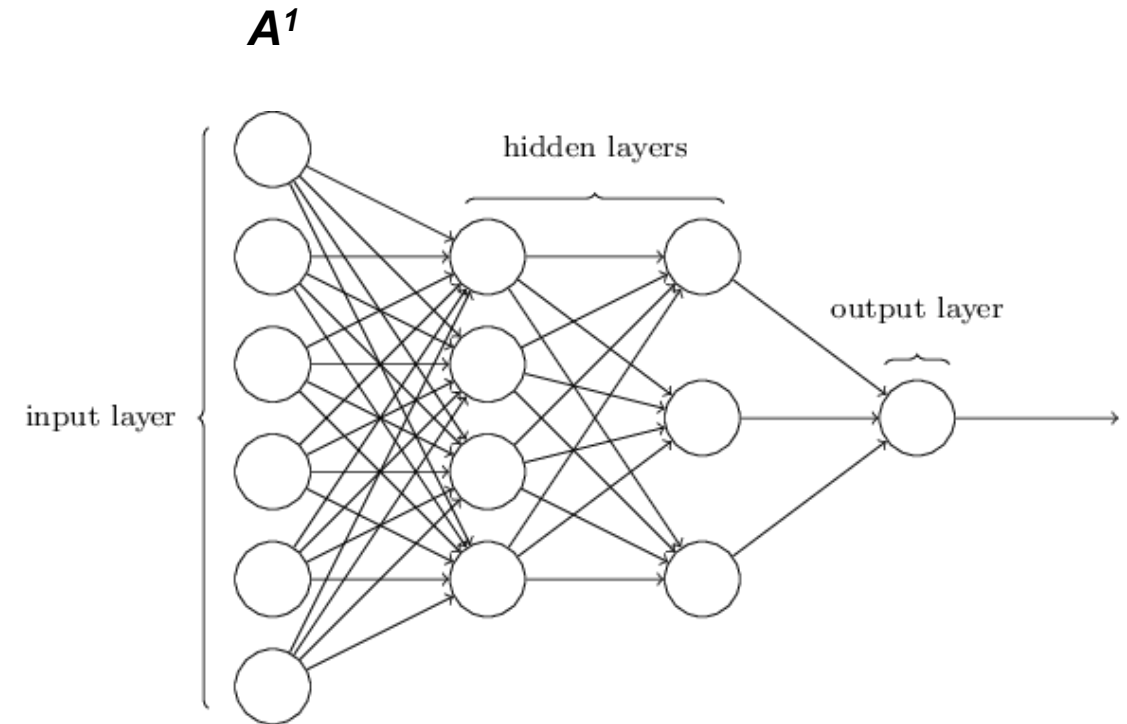
- $\Delta^l = \left[(\mathbf{W}^{l+1})^T \Delta^{l+1} \right] .* \sigma'(\mathbf{Z}^l)$



Backpropagation algorithm

- A way to calculate ∇C
1. **Input x:** Set the activation A^l at the input layer
 2. **Feed-forward:** for each $l = 2, 3, \dots, L$ compute Z^l and A^l
 3. **Output error:** Compute Δ^L of the output layer
 4. **Back-propagate error:** For each layer starting with output, calculate $\Delta^l = \left[(W^{l+1})^T \Delta^{l+1} \right] .* \sigma'(Z^l)$
 5. **Output:**

$$\frac{\partial C}{\partial B^l} = \Delta^l, \quad \frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$



Improving speed at calculating ∇C

- Need to compute the gradient for each training example x and then average

$$- C = \frac{1}{n} \sum_x C_x \rightarrow \nabla C = \frac{1}{n} \sum_x \nabla C_x$$

- With $n > 10000$ this is extremely slow

→ Stochastic gradient descent

- Randomly choose sample m of x (mini-batch)

$$- \nabla C \approx \frac{1}{m} \sum_{j=1}^m \nabla C_{x_j}$$

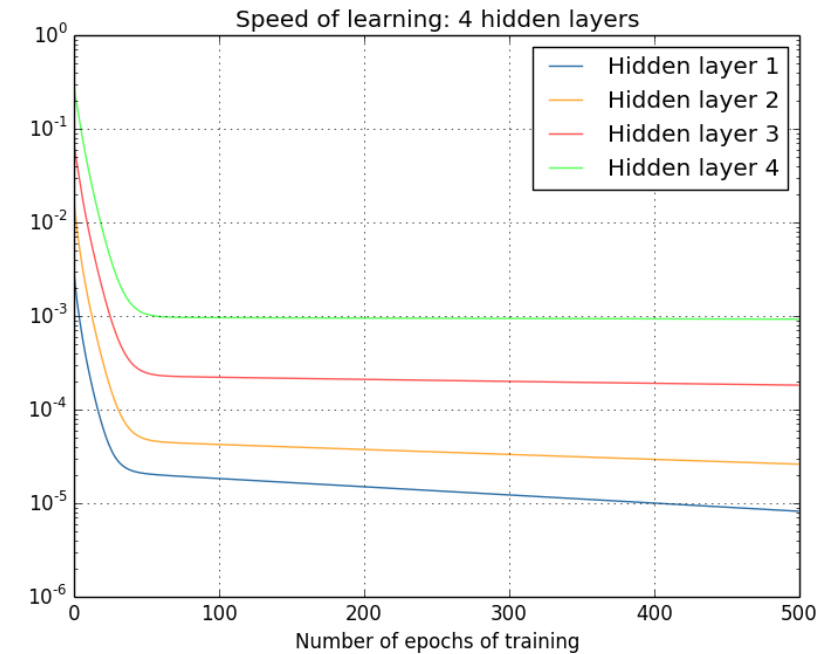
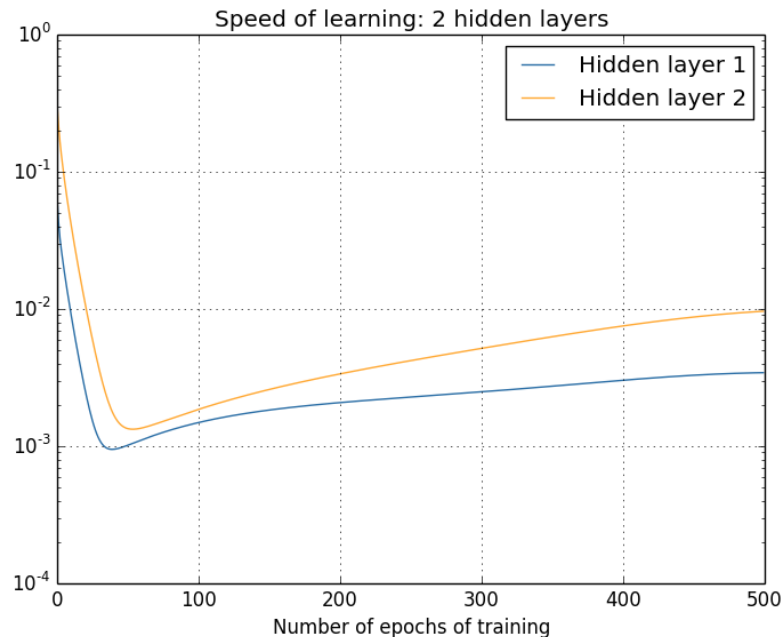
- Train on this mini-batch, then choose a new batch until having trained with all (1 epoch)
- Then start over and redo until finished

$$w_k \rightarrow w'_k = w_k - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial w_k}$$

$$b_l \rightarrow b'_l = b_l - \frac{\eta}{m} \sum_j \frac{\partial C_{x_j}}{\partial b_l},$$

Deep networks

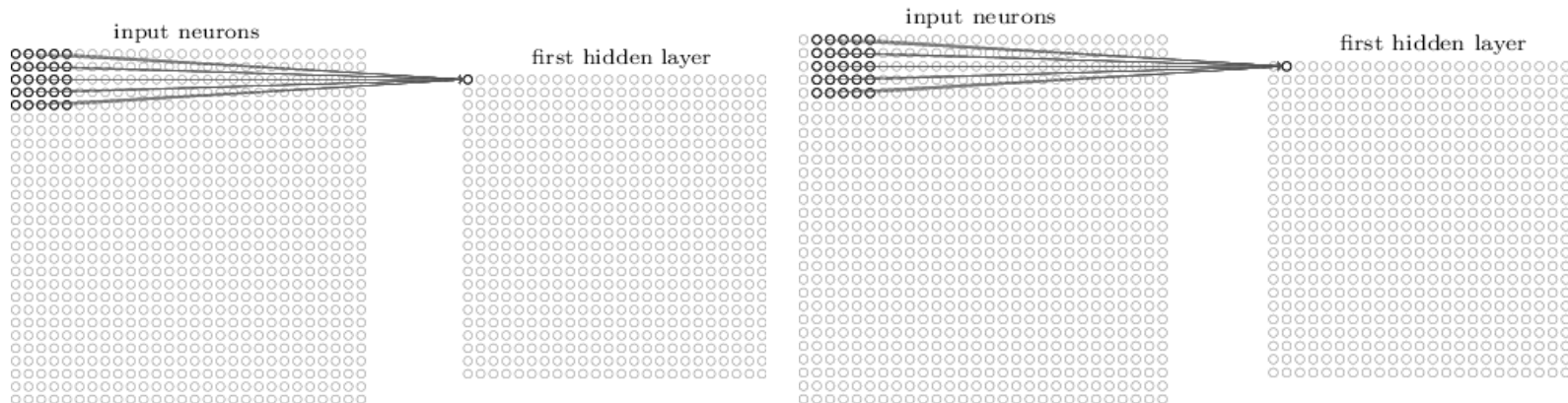
- More than one hidden layer → much improved abstraction power
- Different layers learn at vastly different speed → backprop and gradient descent works poorly
 - Vanishing gradient problem: Built-in instability with gradient descent techniques.



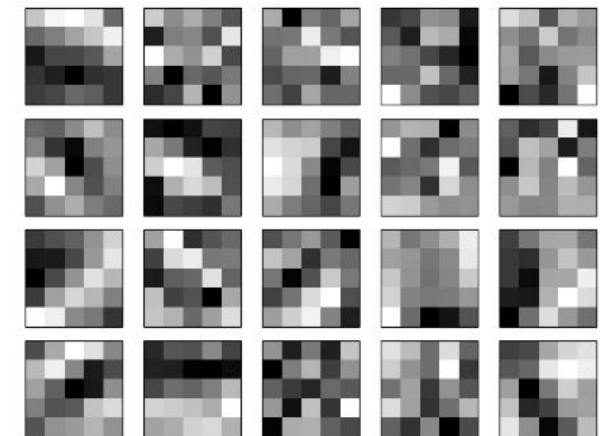
Convolutional deep neural networks

- Utilizing the fact that neighboring inputs are related to each other → great for image analysis
- Create hidden layer: Each neuron connects to a subset of neurons in the previous layer
 - A *local receptive field*
- All neurons in first hidden layer use same *shared bias* and *shared weights*!
 - i.e. This layer can detect the same feature **anywhere** in the image
 - Drastically reduces number of parameters: $5 \times 5 + 1 = 26$
- A convolutional network uses many parallel feature maps to build up an understanding of images

$$\sigma \left(b + \sum_{l=0}^4 \sum_{m=0}^4 w_{l,m} a_{j+l,k+m} \right)$$

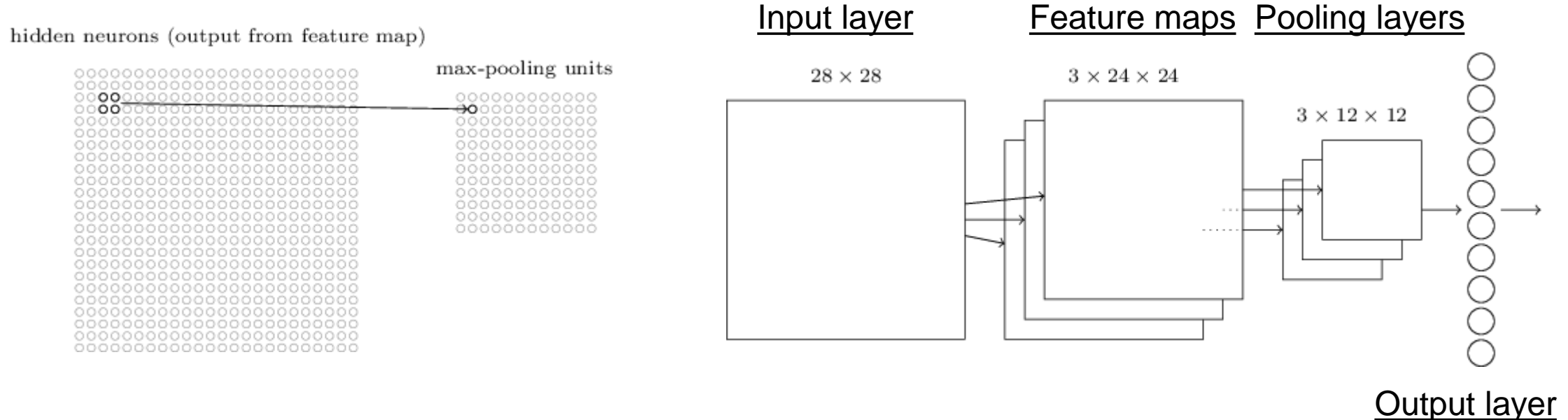


20 feature maps used for recognizing handwritten numbers



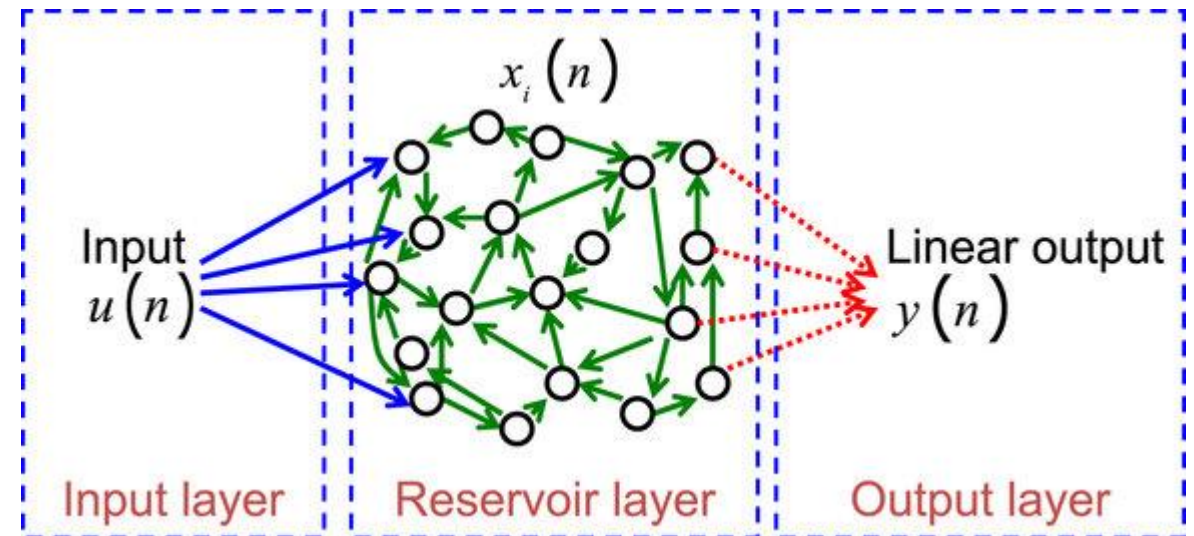
Complete Convolutional network

- One pooling layer per feature map
 - To condense information from feature maps
 - Max pooling: activation = the max of the neurons in the connected neurons in feature map
 - L2 pooling: root of the sum of squares
- The Output layer is a normal completely connected layer to all neurons in the pooling layers
- Avoids learning issues by reducing parameters...



Reservoir computing

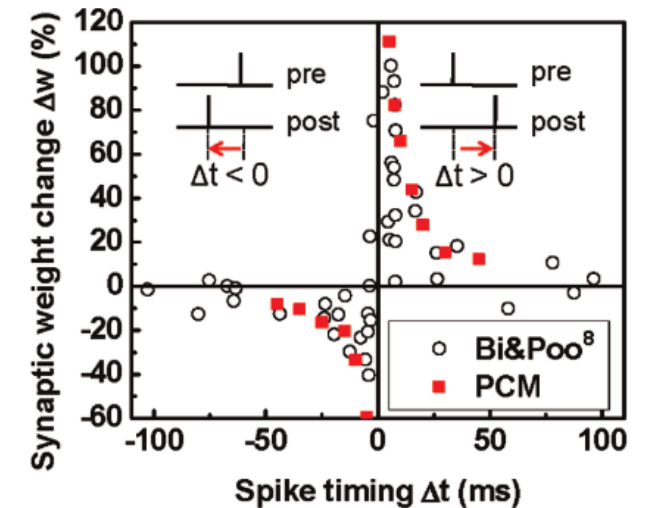
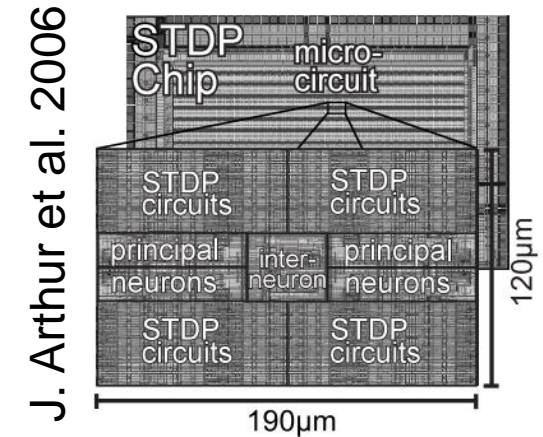
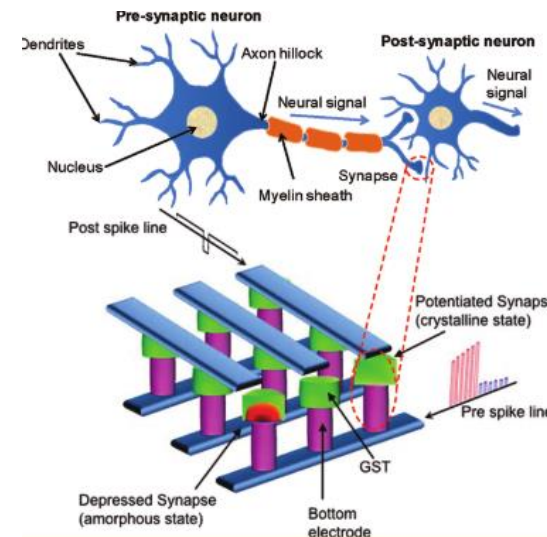
- Viewed as an extension of neural networks
- Non-linear randomly connected nodes form a reservoir
 - Reservoir nodes and connections are usually constant
- The dynamic behavior of the reservoir creates the logic
- Read out by linear combination of the reservoir output
 - Training by linear regression of the output to known inputs



Hardware for Neural networks

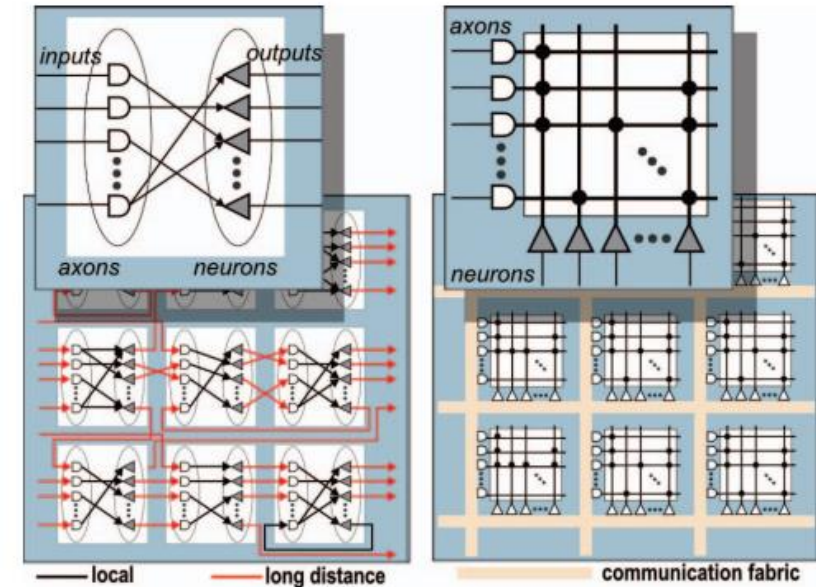
- The brain is very energy efficient 10 W
 - Simulating 5s of brain activity with supercomputer → 500s and 1.4 MW
- Specialized hardware needed to improve energy efficiency
- Using CMOS as synapse → ~10 transistors/synapse
- PCMs as synapses → dense and energy efficient
- GPU instead of CPU → Faster at matrix operations
 - Nvidia leads this market

Nvidia Tesla P100
\$2 Billion in R&D
15 billion transistors



Kuzum et al Nano Lett. 2011

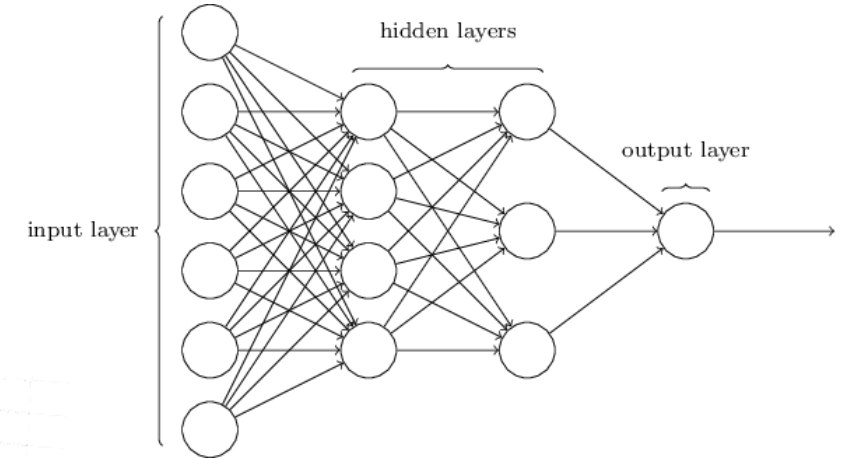
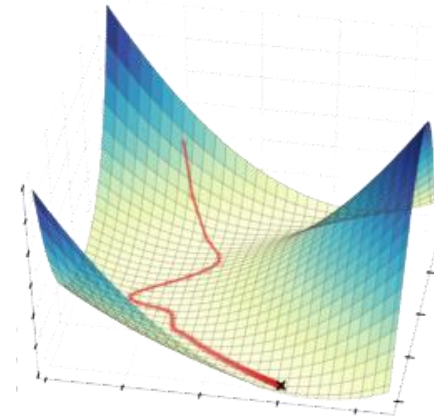
- The first (?) commercial neuromorphic processing chip
- Neurosynaptic cores of 256 output neurons and 256 input axons, connected by 256x256 synapses
- Cores are connected by network-on-chip → large neural network
- 4096 cores as 64x64 array
- $>1 \times 10^6$ neurons, $>256 \times 10^6$ synapses
- 5.4 billion transistors (!!)
- Uses spiking neurons
- Each TrueNorth consumes ~ 0.23 W in active mode
 - A one chip system: 3.5 W
- Currently scaled up to 4x4 TrueNorth system (NS16e)



NS16e

Summary – Cognitive computing

- A new revolution in computing
 - Data-centered computing
 - Big business → industry is leaping at it
- Based on the neural network of the brain
- Simplified neuron models
 - Sigmoid neuron
 - Spiking neurons
- Learning by:
 - gradient descent
 - Back-propagation
- Convolutional deep networks
- Reservoir computing
- GPU instead of CPU
- Specialized hardware for energy efficiency
 - PCM...or something else?



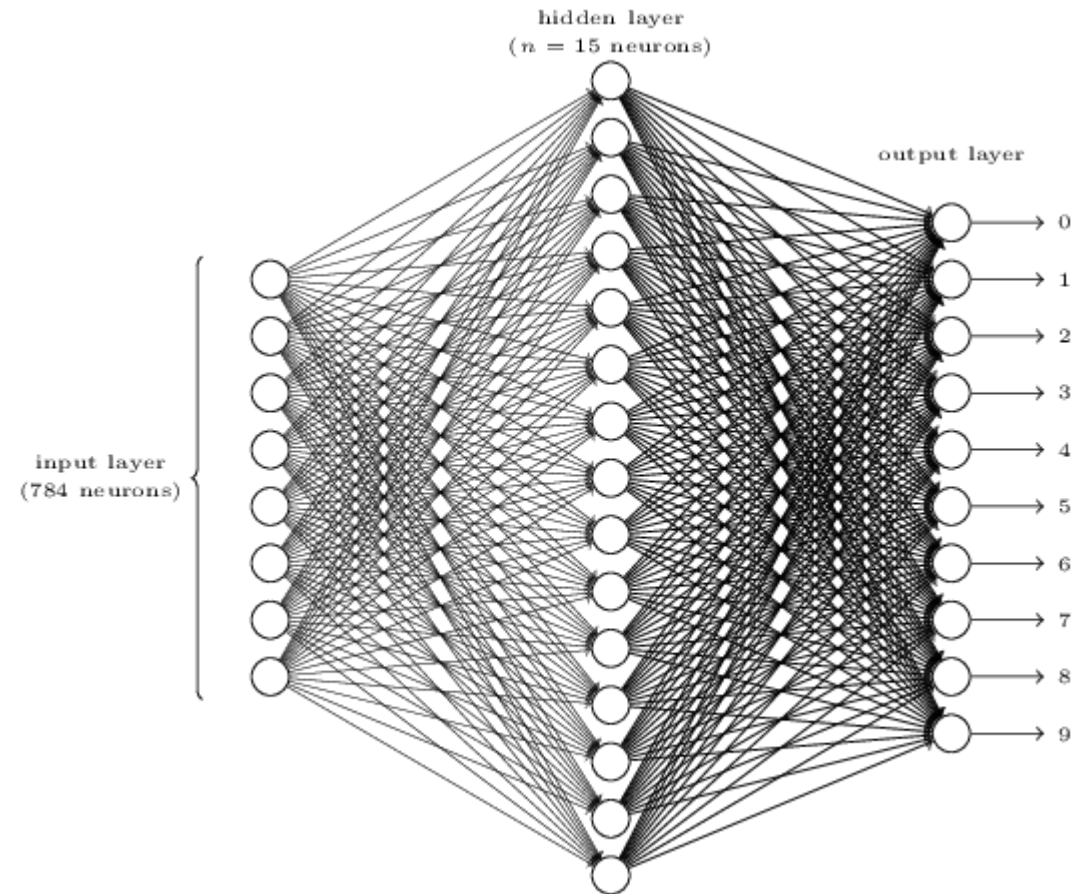
Hand-in assignment

- Task: "Design and implement a neural network that can identify handwritten numbers 0-9"
- MNIST data base (<http://yann.lecun.com/exdb/mnist/>): 60 000 examples of handwritten numbers
 - 28x28 pixels grey-scale
 - Target: > 95% accuracy (but who will be the best?)
- 1-2 persons per project
- Written report: 5-10 pages
 - What kind of network architecture was used and why?
 - How was learning done and why did you do it that way?
 - Include code as Appendix
- Deadline: **12 March**
- Use literature list as an aid
- Detailed instructions come tomorrow...



Some tips

- Start with designing your network
 - # neurons in input layer: one/pixel
 - # neurons in output layer: 10
 - # hidden layers: experiment
 - # neurons in hidden layer: experiment
 - Use matrix representation!
 - Feed-forward network
 - Use gradient descent and back-propagation
- Important to consider:
- Choice of cost function
- Size of mini-batch
- Learning step size



Oral exams

- 7-10 March
 - 30 min discussion 1 on 1
 - Based on content from lectures and literature list
 - 1. Be able to describe devices and phenomena that were brought up on the lectures
 - 2. Be able to describe the impact on society/industry of emerging technologies
-
- Sign up on Doodle to decide on exact times.