

Artificial Neural Networks and Reinforcement Learning:

Lecture 1A

Introduction to the field

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Reading for this Introduction:

Motivational background reading:

Silver et al. 2017, Archive

*Mastering Chess and Shogi by Self-Play with a
General Reinforcement Learning Algorithm*

Goodfellow et al., Ch. 1 of

Deep Learning

Artificial Neural Networks and Reinforcement Learning

Wulfram Gerstner

EPFL, Lausanne, Switzerland

The EPFL logo is displayed in a bold, red, sans-serif font. It consists of the letters 'EPFL' in a stylized, blocky typeface.

Previous slide.

Results with artificial neural networks are discussed in newspaper articles and have inspired people around the world.

These years we experience the third wave of neural networks.

The first wave happened in the 1950s with the first simple computer models of neural networks, with McCulloch and Pitt and Rosenblatt's Perceptron. There was a lot of enthusiasm, and then it died.

The second wave happened in the 1980, around the Hopfield model, the BackPropagation algorithm, and the ideas of 'parallel distributed processing'. It died in the mid-nineties when statistical methods and Support Vector Machines took over.

The third wave started around 2012 with larger neural networks trained on GPUs using data from big image data bases. These neural networks were able to beat the benchmarks of Computer vision and have been called 'deep networks'.

Artificial Neural Networks, how they work, and what they can do, will be in the focus of this lecture series.

Artificial Neural Networks and Reinforcement Learning

Introduction to the field

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. From Biological to Artificial Neurons

Previous slide.

In this first introduction lecture the focus is on a general introduction into the field (with its subparts Reinforcement learning and Supervised Learning for Classification).

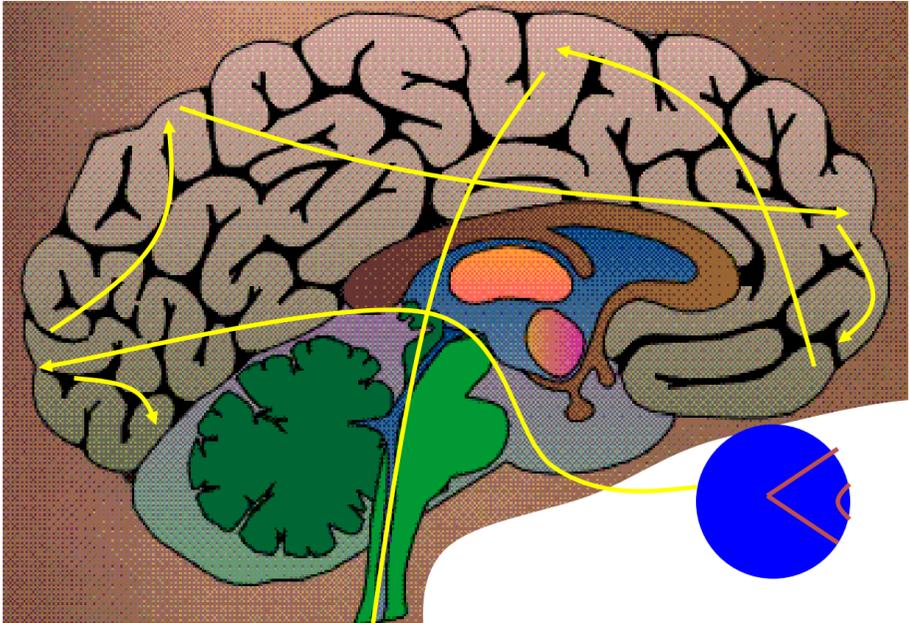
We start with a glimpse of the biological inspirations of the field.

The brain: Cortical Areas

motor cortex

frontal cortex

visual cortex



to muscles

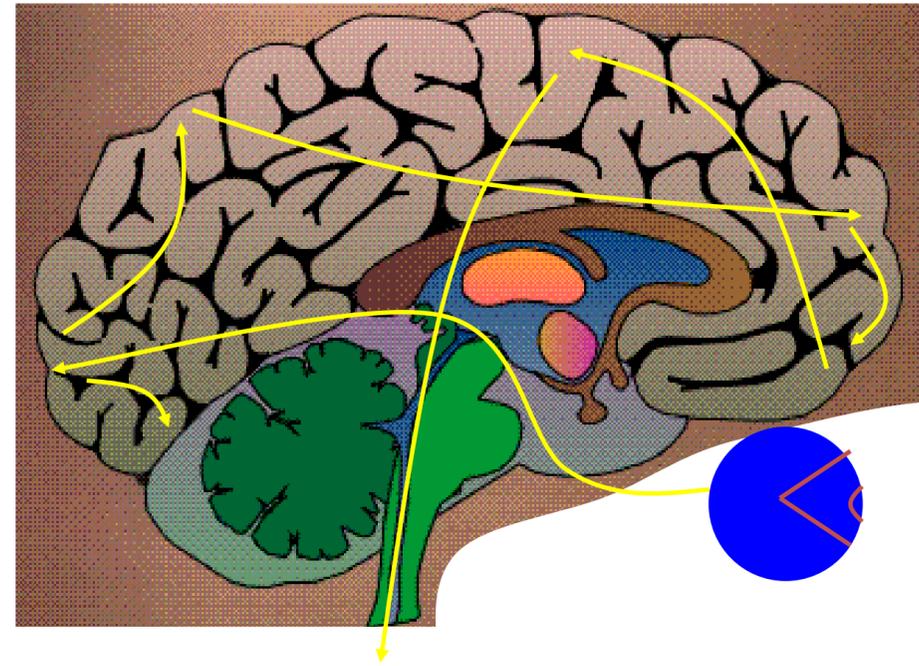


Previous slide.

During all these waves, during 60 years of research, artificial neural networks researchers worked on building intelligent machines that learn, the way humans learn. And for that they took inspiration from the brain.

Suppose you look at an image. Information enters through the eye and then goes to the cortex.

The brain: Cortical Areas



Previous slide.

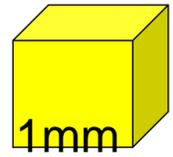
Cortex is divided into different areas:

Information from the eye will first arrive at visual cortex (at the back of the head), and from there it goes on to other areas. Comparison of the input with memory is thought to happen in the frontal area (above the eyes). Movements of the arms are controlled by motor cortex somewhere above your ears.

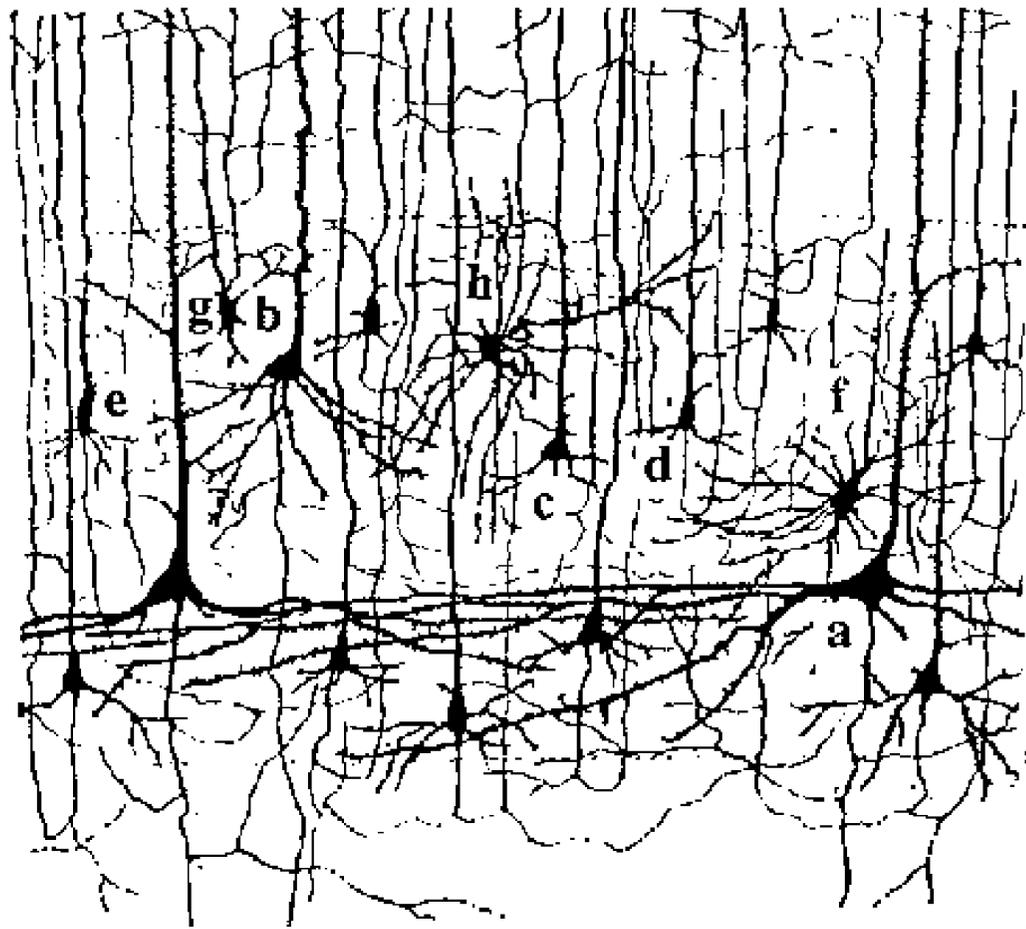
Talking about cortical areas provides a **macroscopic** view of the brain.

The Brain: zooming in

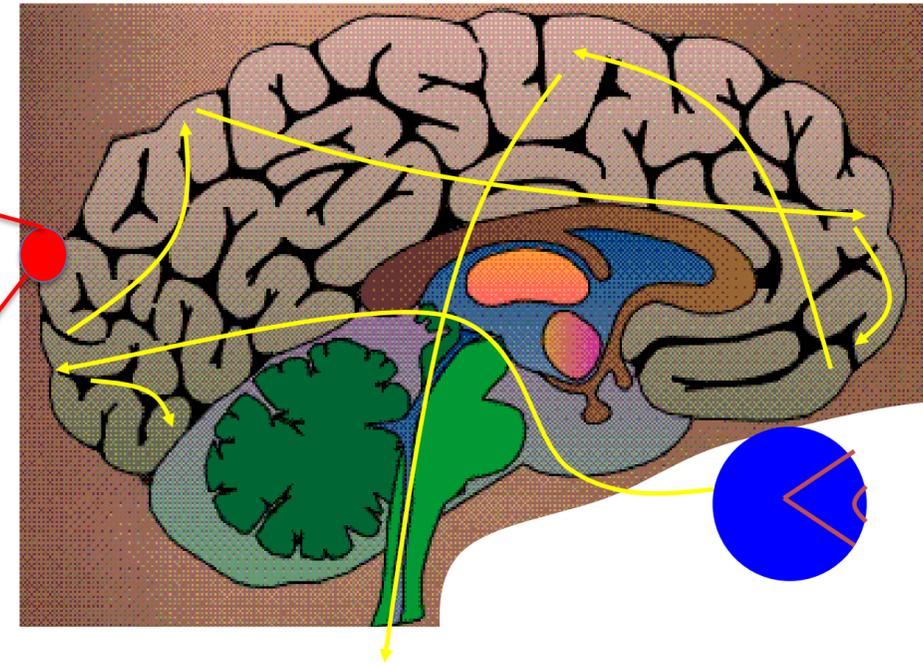
1mm



10 000 neurons
3 km of wire



Ramon y Cajal



Previous slide.

If we zoom in and look at one cubic millimeter of cortical material under the microscope, we see a network of cells.

Each cell has long wire-like extensions.

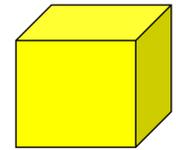
If we counted all the cells in one cubic millimeter, we would get numbers in the range of ten thousand.

Researchers have estimated that, if you put all the wires you find in one cubic millimeter together you would find several kilometers of wire.

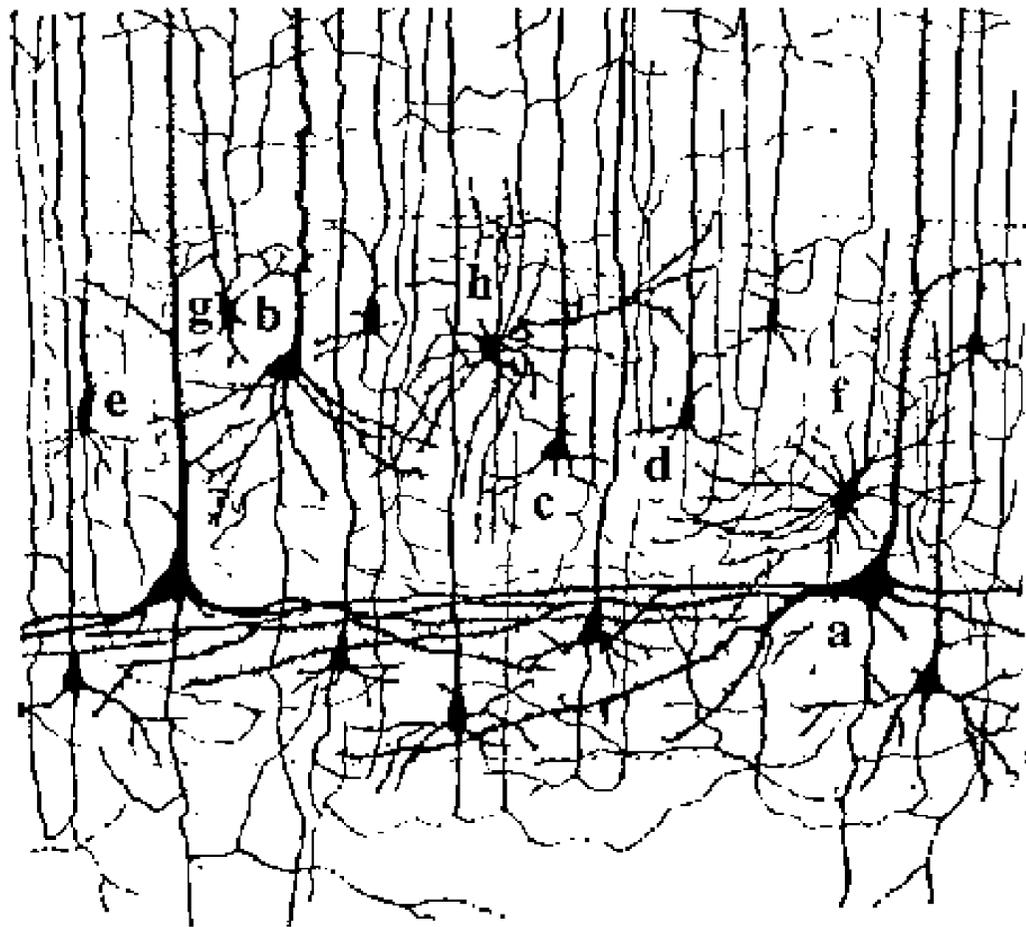
Thus, the neural network of the brain is a densely connected and densely packed network of cells.

The brain: a network of neurons

1mm



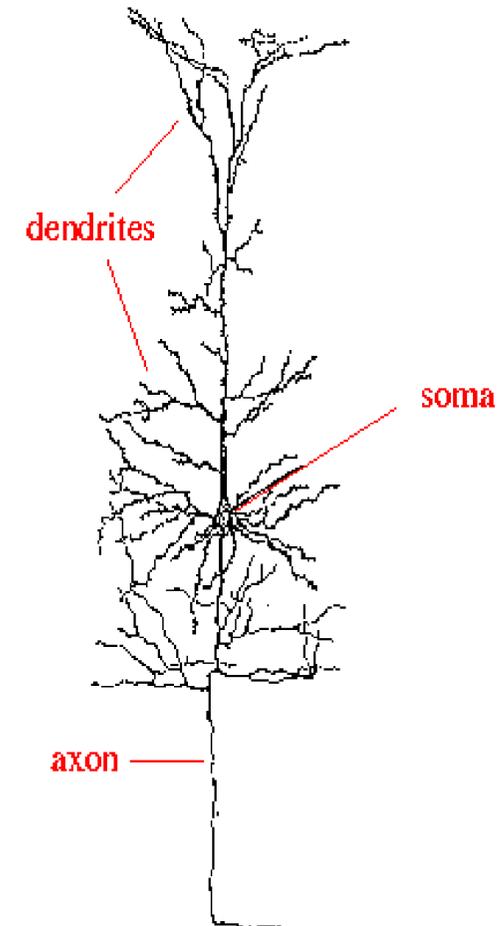
10 000 neurons
3km of wire



Ramon y Cajal

Signal:

Action potential (short pulse)

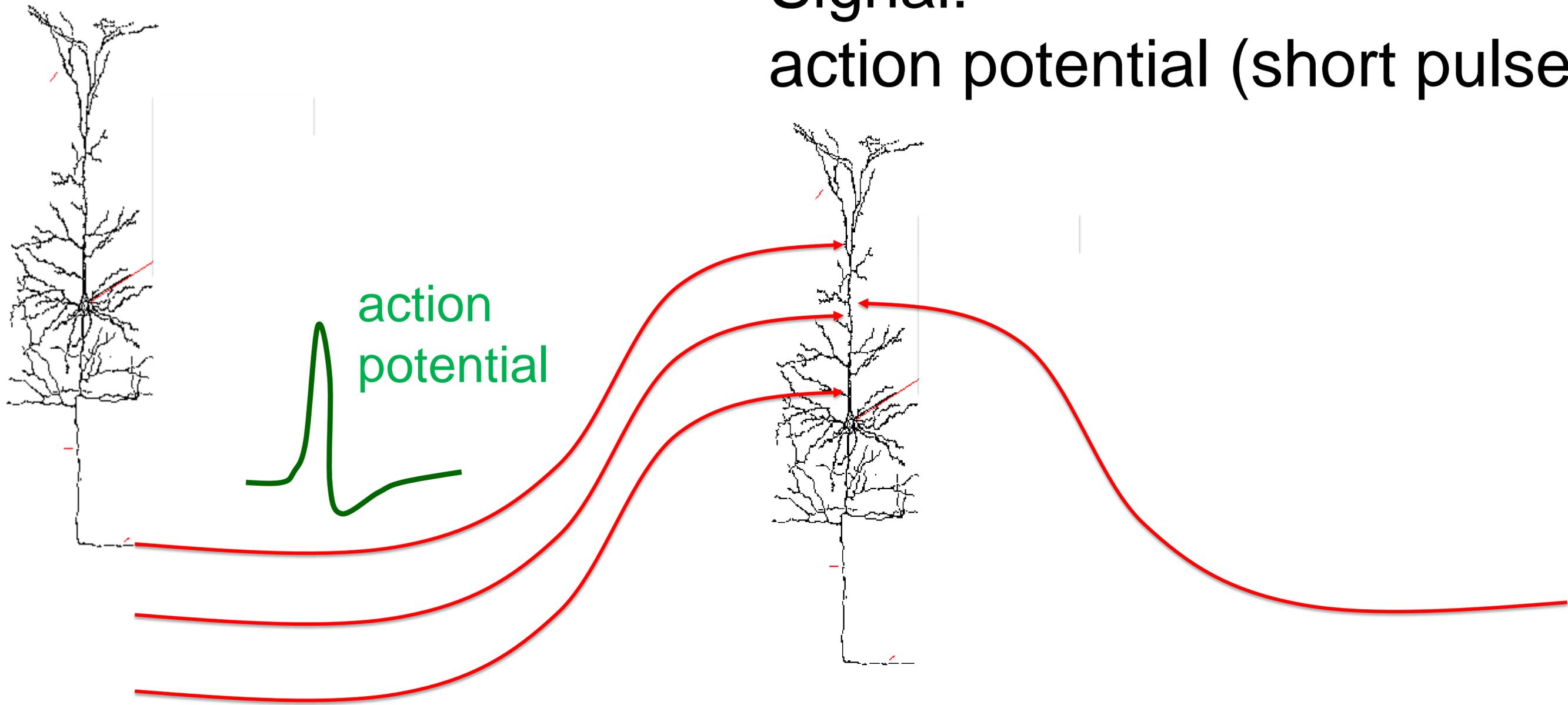


Previous slide.

These cells are called neurons and communicated by short electrical pulses, called action potentials, or 'spikes'.

The brain: signal transmission

Signal:
action potential (short pulse)



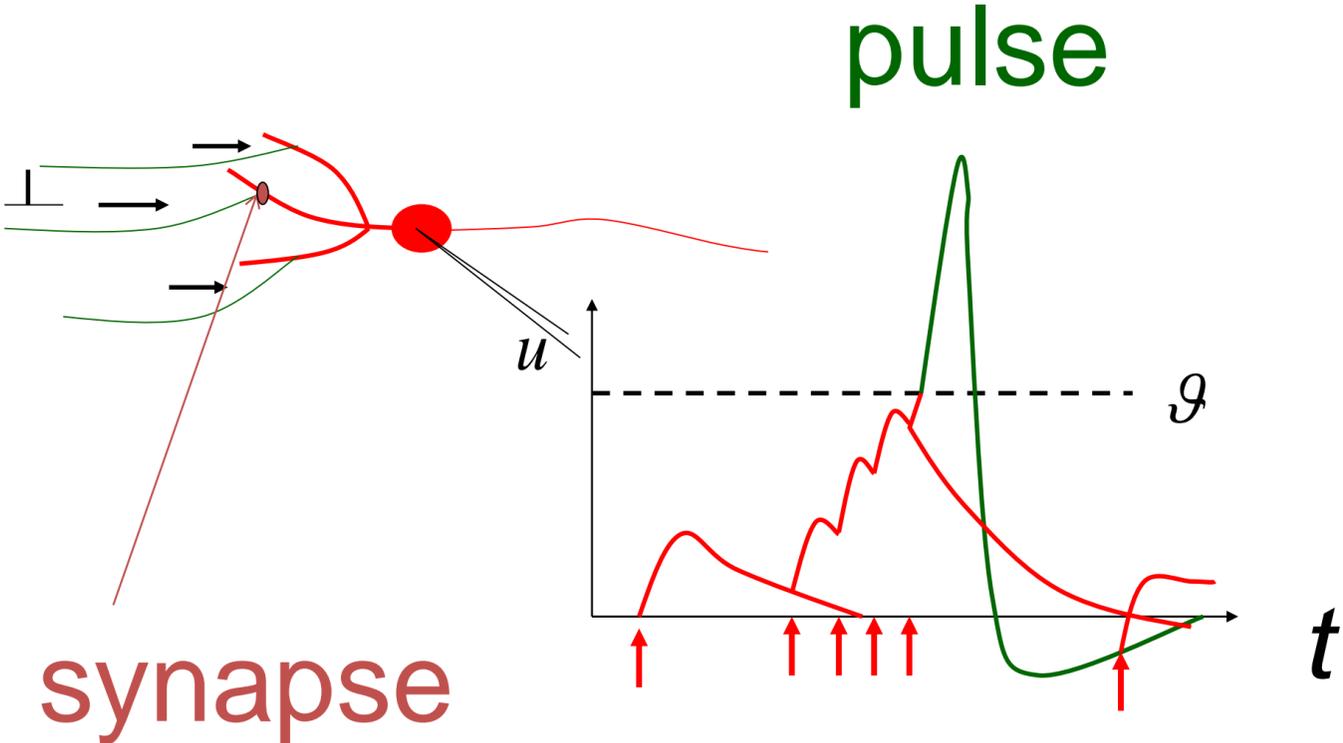
More than 1000 inputs

Previous slide.

Signals are transmitted along the wires (axons). These wires branch out to make contacts with many other neurons.

Each neuron in cortex receives several thousands of wires from other neurons that end in 'synapses' (contact points) on the dendritic tree.

The brain: neurons sum their inputs



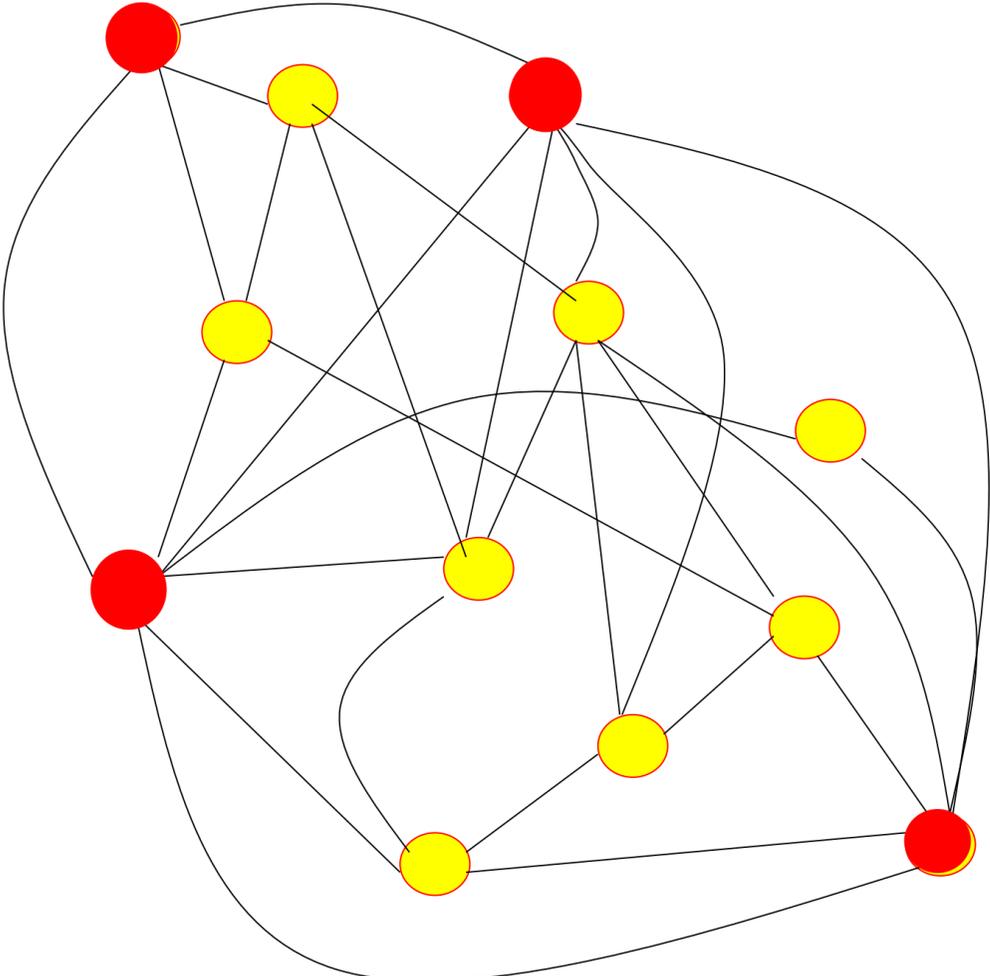
Previous slide.

If a spike arrives at one of the synapses, it causes a measurable response in the receiving neuron.

If several spikes arrive shortly after each other onto the same receiving neuron, the responses add up.

If the summed response reaches a threshold value, this neuron in turn sends out a spike to yet other neurons (and sometimes back to the neurons from which it received a spike).

Summary: the brain is a large recurrent network of neurons



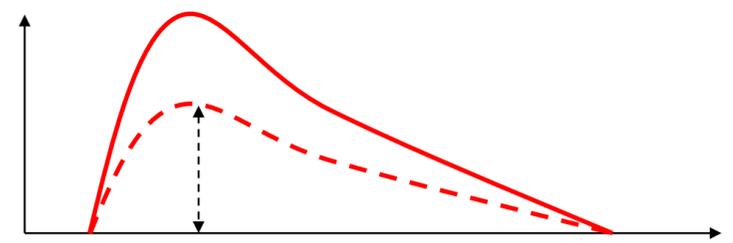
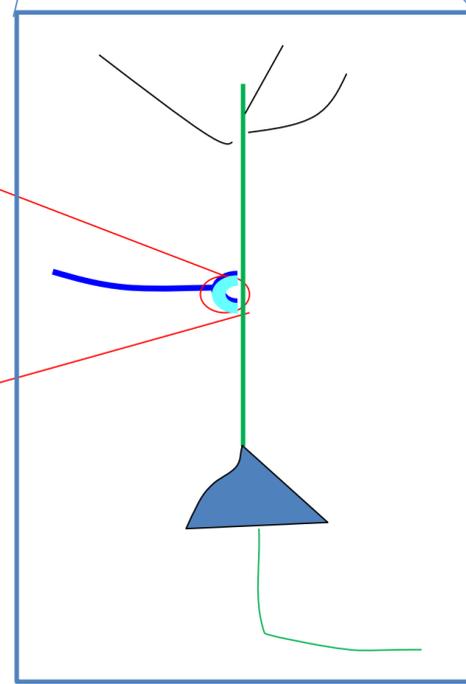
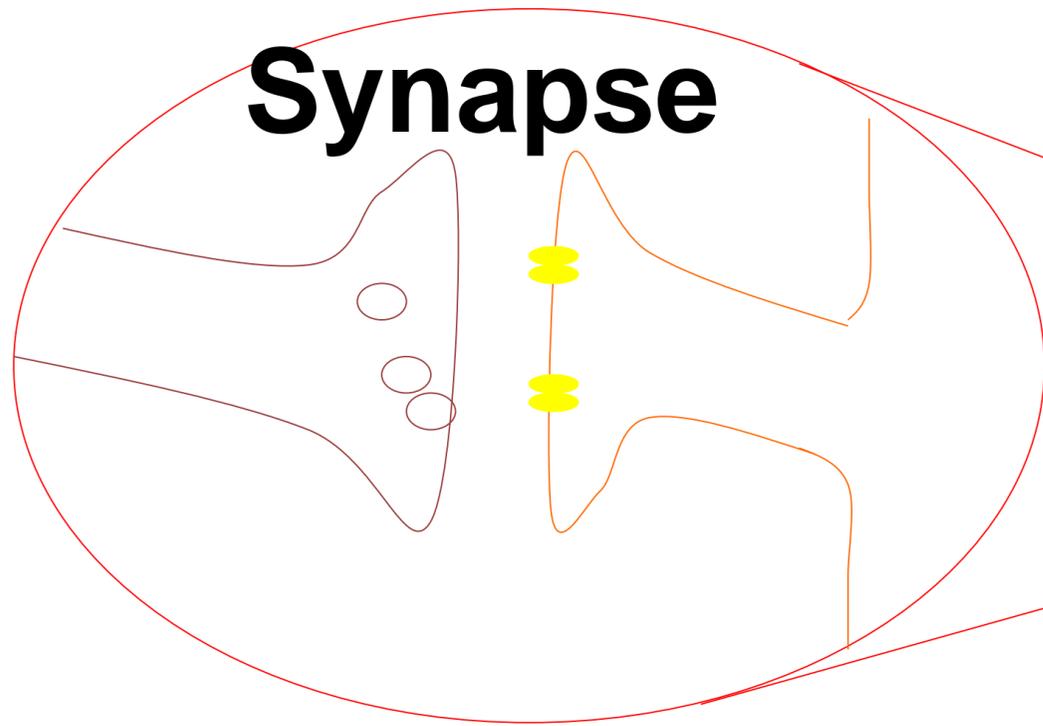
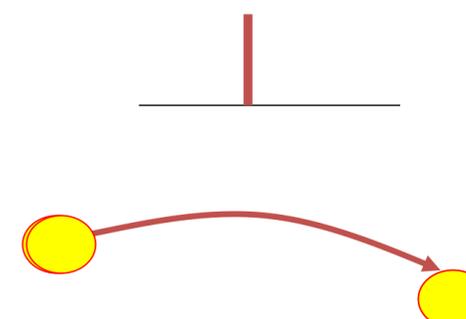
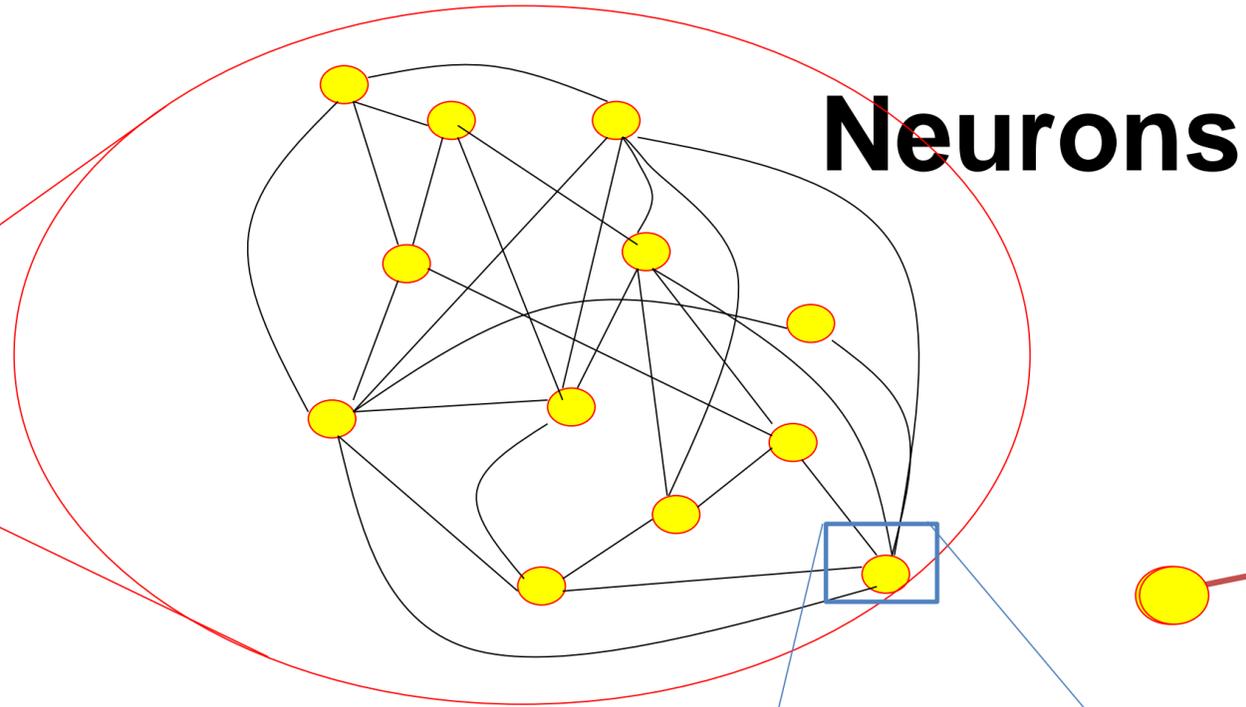
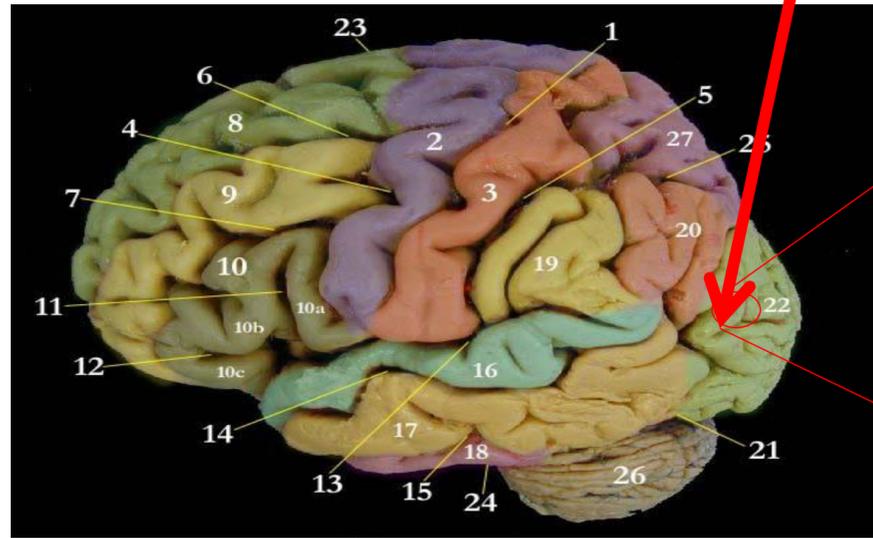
● Active neuron

Previous slide.

Thus, signals travel along the connections in a densely connected network of neurons.

Sometimes I draw an active neuron (that is a neuron that currently sends out a spike) with a filled red circle, and an inactive one with a filled yellow circle.

Learning in the brain: changes between connections



learning = change of connection

Previous slide.

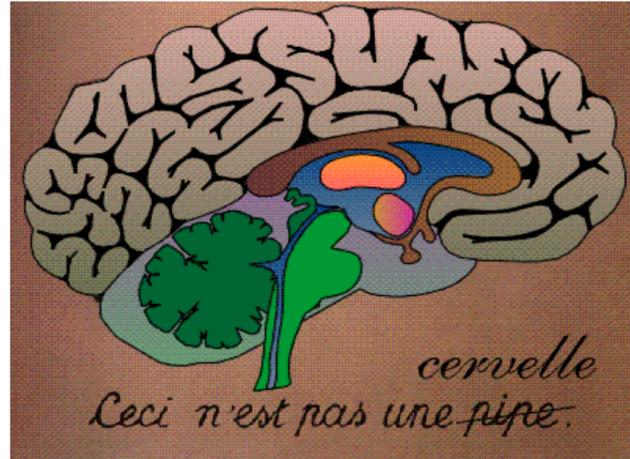
Synapses are not just simple contact points between neurons, but they are crucial for learning.

Any change in the behavior of an animal (or a human, or an artificial neural network) is thought to be linked to a change in one or several synapses.

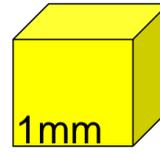
Synapses have a 'weight'. Spike arrival at a synapse with a large weight causes a strong response; while the same spike arriving at a synapses with a small weight would cause a low-amplitude response.

All Learning corresponds to a change of synaptic weights. For example, forming new memories corresponds to a change of weights. Learning new skills such as table tennis corresponds to a change of weights.

Neurons and Synapses form a big network



Brain



1mm

1mm

10 000 neurons

3 km of wire

10 billions neurons

10 000 connexions/neurons

memory in the connections

Distributed Architecture

**No separation of
processing and memory**

Previous slide.

Even though we are not going to work with the Hebb rule during this class, the above example still shows that

- Memory is located in the connections
- Memory is largely distributed
- Memory is not separated from processing

(as opposed to classical computing architectures such as the van Neumann architecture or the Turing machine)

Artificial Neural Networks

Wulfram Gerstner

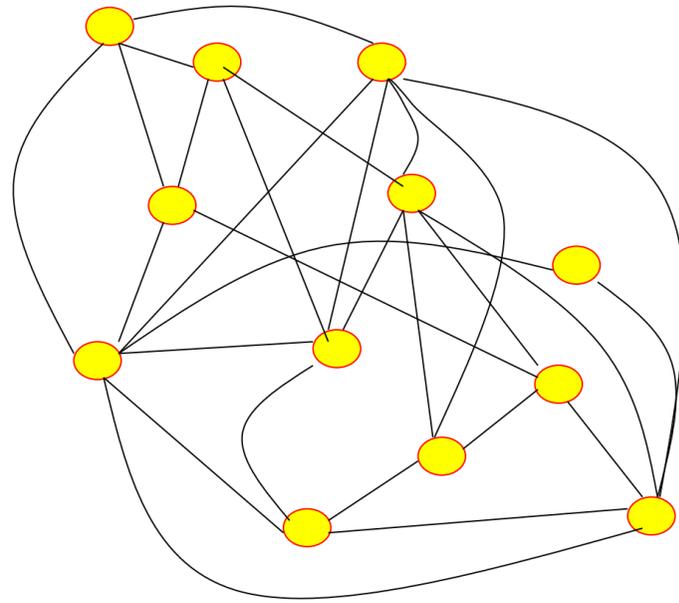
EPFL, Lausanne, Switzerland

From biological neurons to artificial neurons

Previous slide.

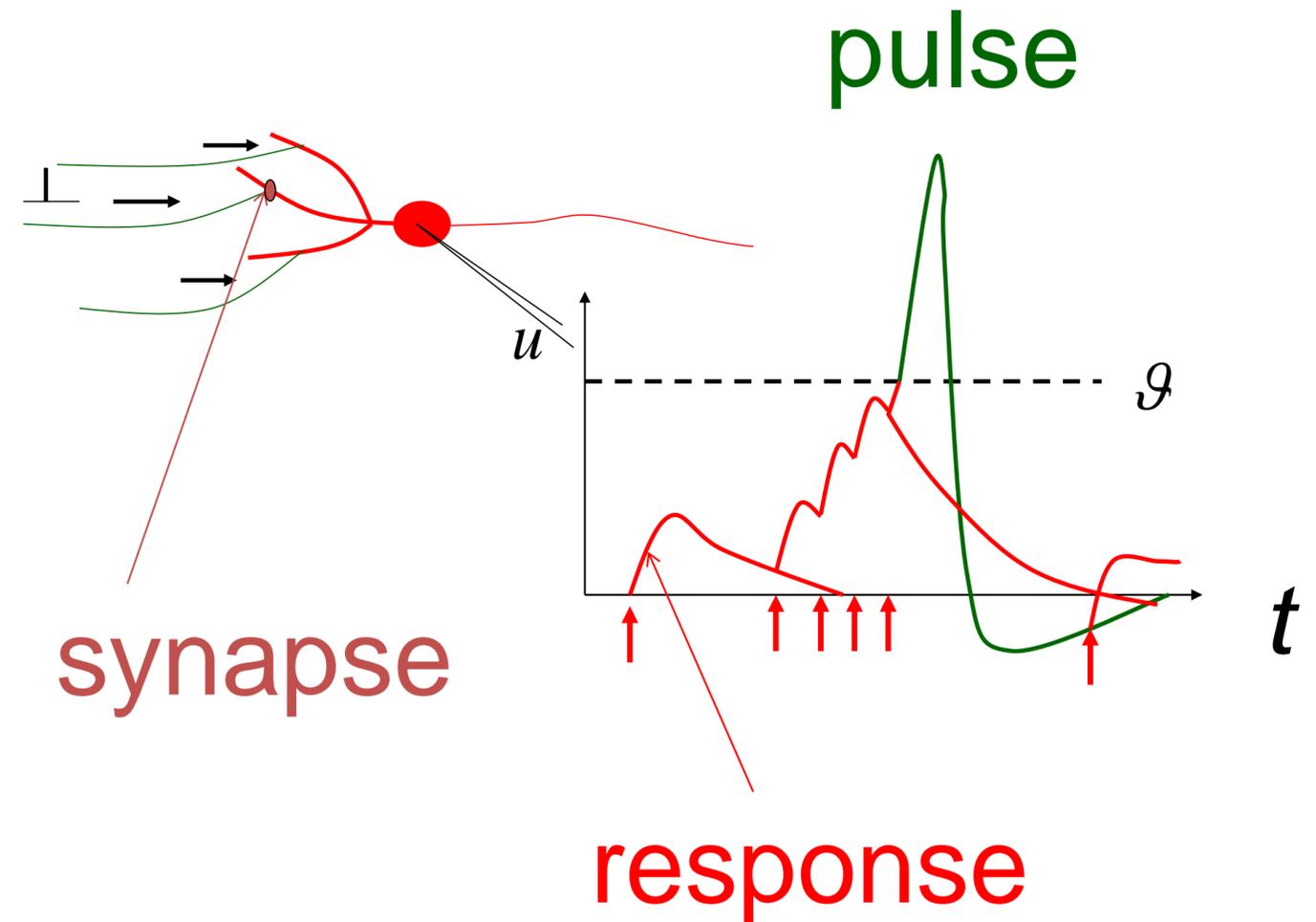
After this super-short overview of the brain, we now turn to artificial neural networks: highly simplified models of neurons and synapses.

Modeling: artificial neurons



- responses are added
- pulses created at threshold
- transmitted to other

→ Mathematical description



Previous slide.

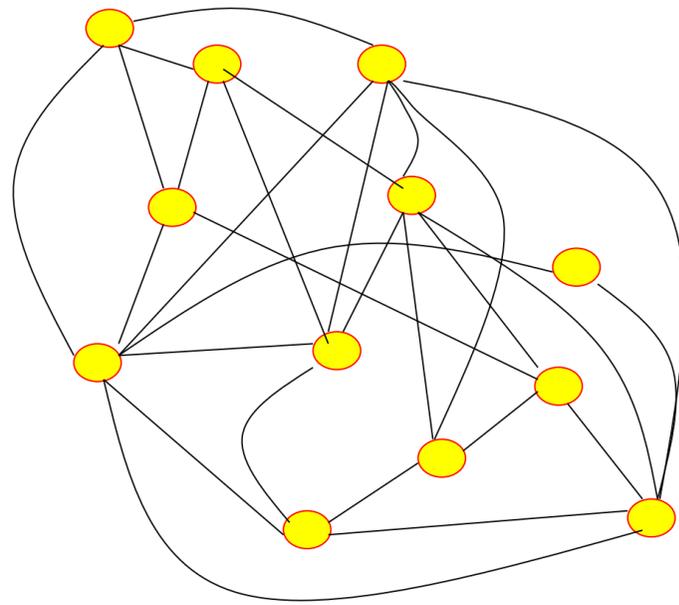
In the previous part we have seen that response are added and compared with a threshold.

This is the essential ideal that we keep for the abstract mathematical model in the following.

We drop the notion of pulses or spikes and just talk of neurons as active or inactive.

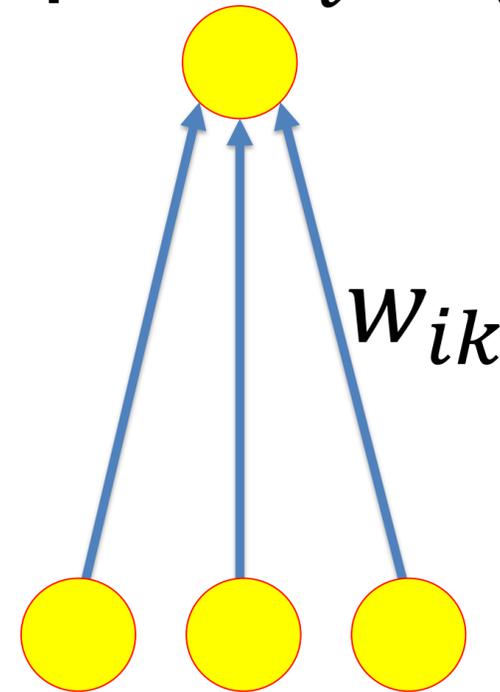
Modeling: artificial neurons

forget spikes: continuous activity x
forget time: discrete updates



activity of output $x_i = g \left(\sum_k w_{ik} x_k \right)$

nonlinearity/threshold



activity of inputs x_k

weights =
adaptive
parameters

Previous slide.

The activity of inputs (or input neurons) is denoted by x_k

The weight of a synapse is denoted by w_{ik}

The nonlinearity (or threshold function) is denoted by g

The output of the receiving neuron is given by

$$x_i = g \left(\sum_k w_{ik} x_k \right)$$

Quiz: biological neural networks

- Neurons in the brain have a threshold.
- Learning means a change of the connection weights
- The total input to a neuron is the weighted sum of individual inputs
- The neuronal network in the brain is feedforward: it has no recurrent connections

Previous slide. Your notes

Artificial Neural Networks and Reinforcement Learning

Introduction to the field

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. From Biological to Artificial Neurons
2. **Artificial Neural Networks for Classification**
- layered feedforward networks

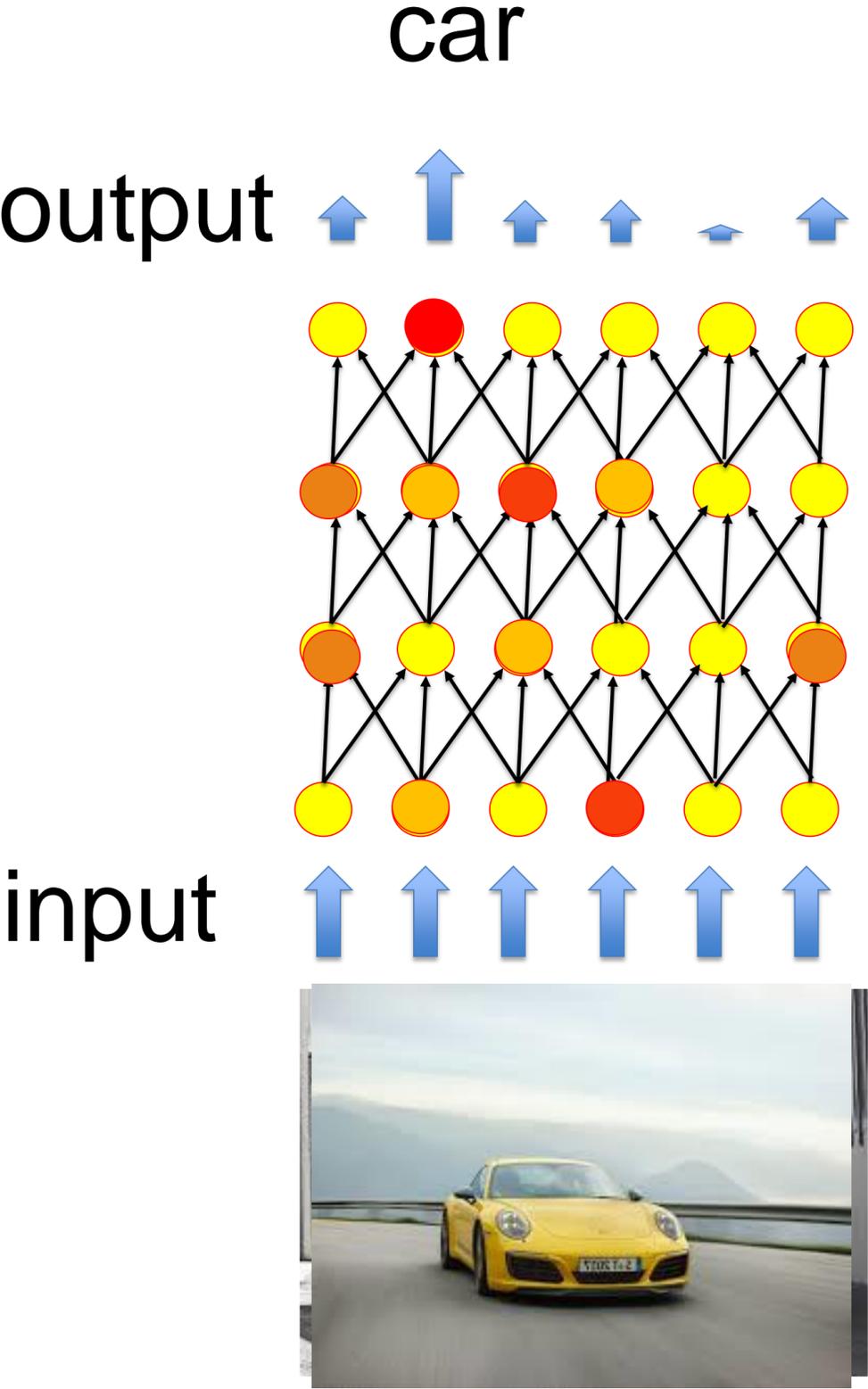
Previous slide.

Now that we know about artificial neurons and synaptic weights, let us construct a useful network.

The first task we study is classification

Artificial Neural Networks for classification

feedforward network



Previous slide.

An input is presented at the bottom of the network.

It passes through several layers of neurons.

All connections are directed from the bottom to the next layer further up: this architecture is called a feedforward network.

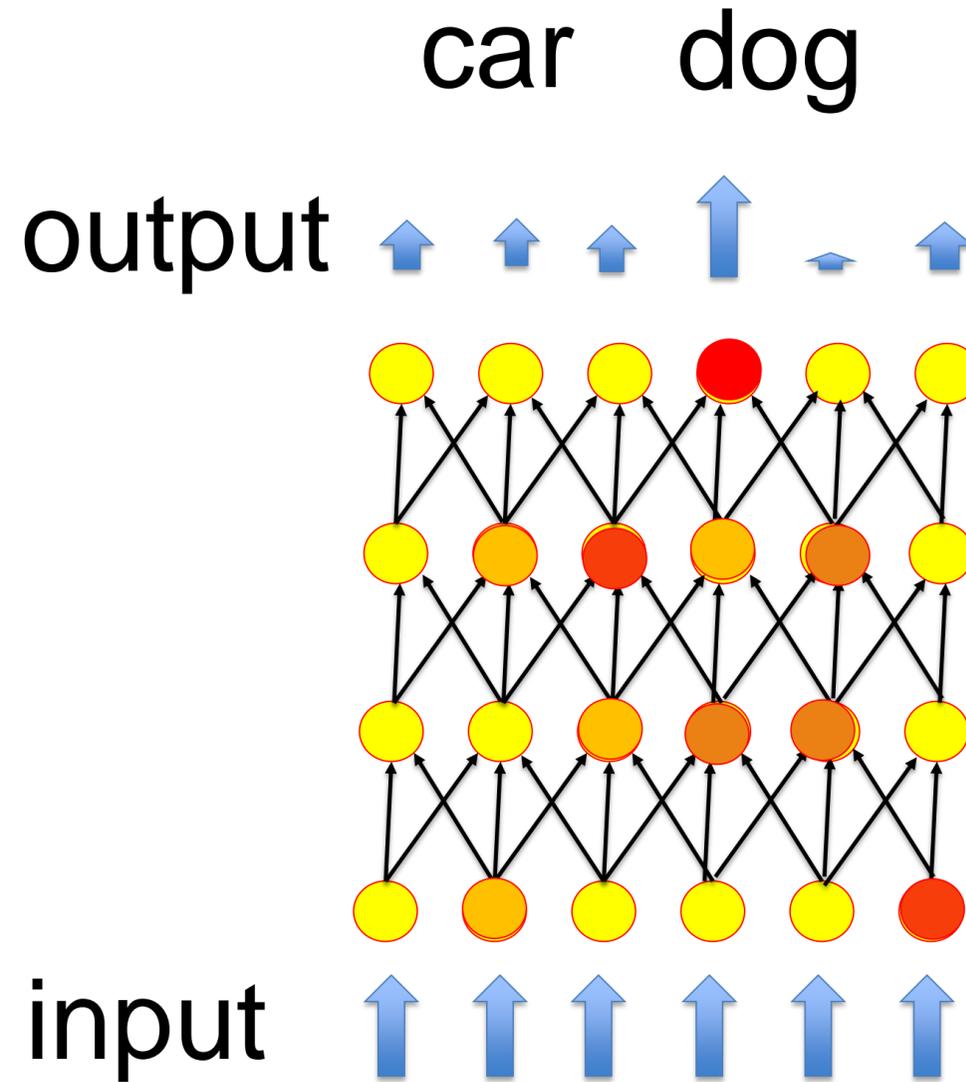
The output is a set of neurons that correspond to different 'classes'.

An ideal network should respond with activating the neuron corresponding to 'car', if the input image shows a car.

Artificial Neural Networks for classification

Aim of learning:

Adjust connections such that output class is correct (for each future input)



Previous slide.

The aim of learning is to adjust the connection weights such that, for each future input, the output class is correct.

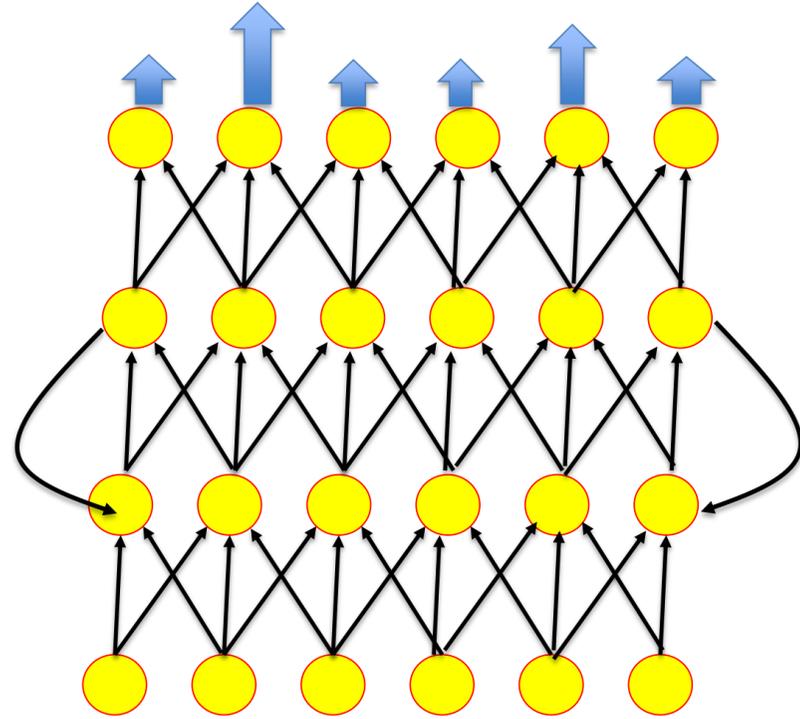
If the input is a dog, the 'dog'-neuron should respond.

If the input is a car, the 'car'-neuron should respond.

In this class, we do not cover the task of building and training artificial neural networks for classification.

Deep networks with recurrent connections

'a man sitting on a couch with a dog'



Network describes the image with the words:

'a man sitting on a couch with a dog'

(Fang et al. 2015)

Previous slide.

An amazing example of sequence production with a recurrent neural network is this network which looks at a static image and outputs the spoken sentence: 'A man is sitting on a couch with a dog'.

A sentence is a temporal sequence of words – and importantly, a sequence that follows grammatical rules.

Sequence learning requires recurrent connections (feedback connections), in contrast to the feedforward architecture that we have seen so far.

And yes, recurrent neural networks can implicitly pick up the statistical rules for sentence formation if they are trained on a sufficiently large data set containing millions of examples.

Summary

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Artificial Neural Networks for Classification

- layered feedforward networks**
(- sometimes also recurrent networks)
- weights are used as adjustable parameters
to learn a stationary classification task or sequence task
- Has been seen in class of ML by Jaggi-Flammariion

Previous slide.

Synaptic weights are the adjustable parameters of artificial neural networks. Artificial neural networks are mostly used in a layered feedforward structure, but recurrent neural networks are often used for sequence learning.

Artificial Neural Networks and Reinforcement Learning

Introduction to the field

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. From Biological to Artificial Neurons
2. **Artificial Neural Networks for Classification**
 - layered feedforward networks
 - recurrent networks
3. **Artificial Neural Networks for action learning**

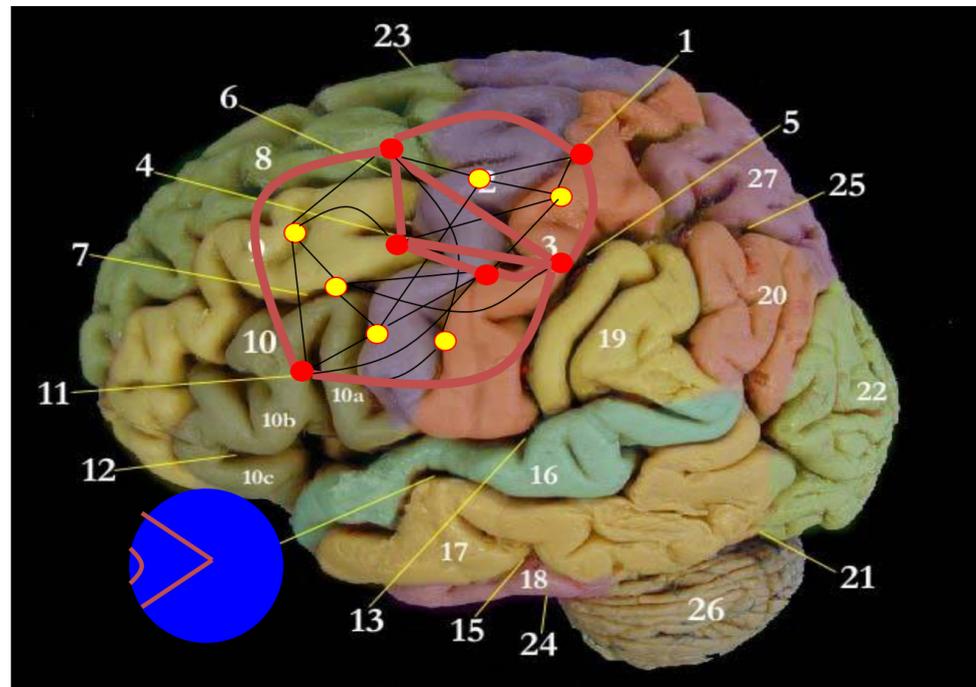
Previous slide.

However, classification is not the only task we are interested in.

Artificial Neural Networks for action learning



Learning from mistakes,
And from (rare!) successes



Missing:

Value of action

- 'goodie' for dog
- 'success'
- 'compliment'



Reinforcement learning = learning based on reward

Previous slide.

Let us go back for a moment to the brain, and how humans or animals learn.

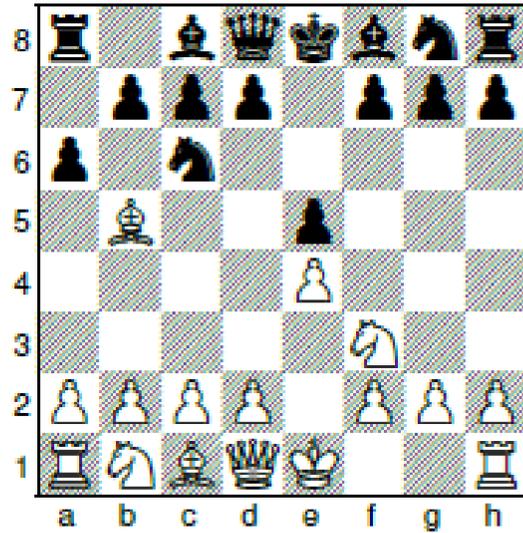
We learn actions by trial and error exploiting rather general feedback: reward or praise on one side, pleasure and pain on the other side.

Important is the notion of value of an action.

Learning actions or sequences of actions based on 'reward' is very different from classification: it falls in the field of 'Reinforcement Learning'.

Deep reinforcement learning

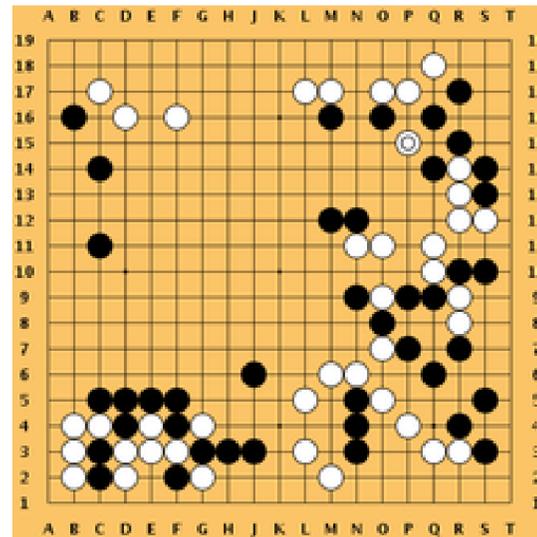
Chess



Artificial neural network
(*AlphaZero*) discovers different
strategies by playing against itself.

In Go, it beats Lee Sedol

Go



Previous slide.

The same kind of ideas (learning by reward) have also been implemented in artificial neural networks that are trained by reinforcement learning.

In a game such a Chess or Go, the reward signal is only given once at the very end of the game: positive reward if the game is won, and negative reward if it is lost.

This very sparse reward information is sufficient to train an artificial neural network to a level where it can win against grand masters in chess or Go.

To improve performance, each network plays against a copy of itself. By doing so it discovers good strategies (such as openings in chess).

Deep reinforcement learning

Network for choosing action

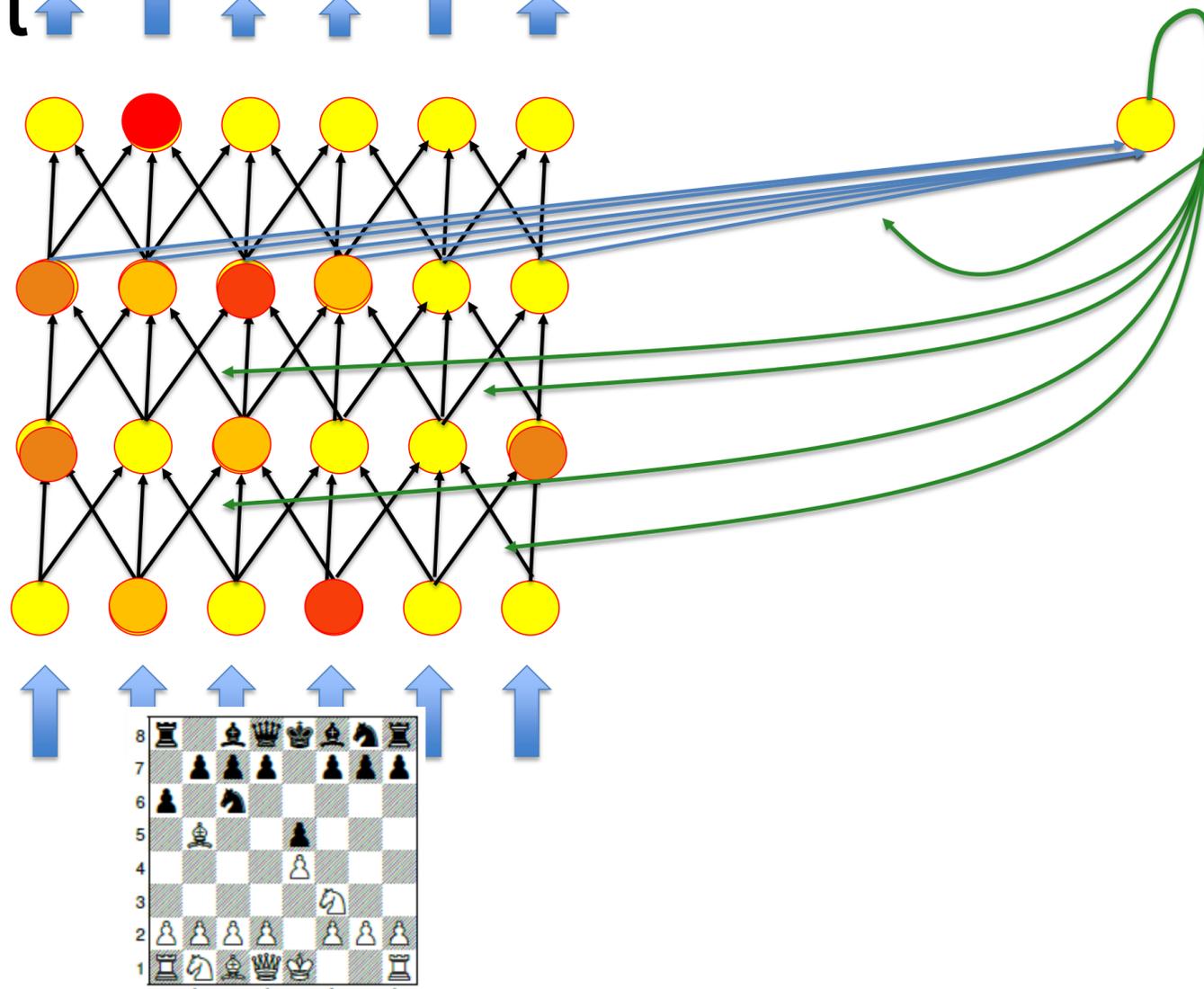
action:

Advance king

2nd output for **value** of action:

probability to win

output



learning:

- change connections

aim:

- Predict value of position

- Choose next action to win

Previous slide.

Schematically, the artificial neural network takes the position of chess as input.

There are two types of outputs:

- The main outputs are the actions such as 'move king to the right'
- An auxiliary output predicts the 'value' of each state. It can be used to explore possible next positions so as to pick the one with the highest value.
- The value can be interpreted as the probability to win (given the position)

In the theory of reinforcement learning, positions are also called 'states'.

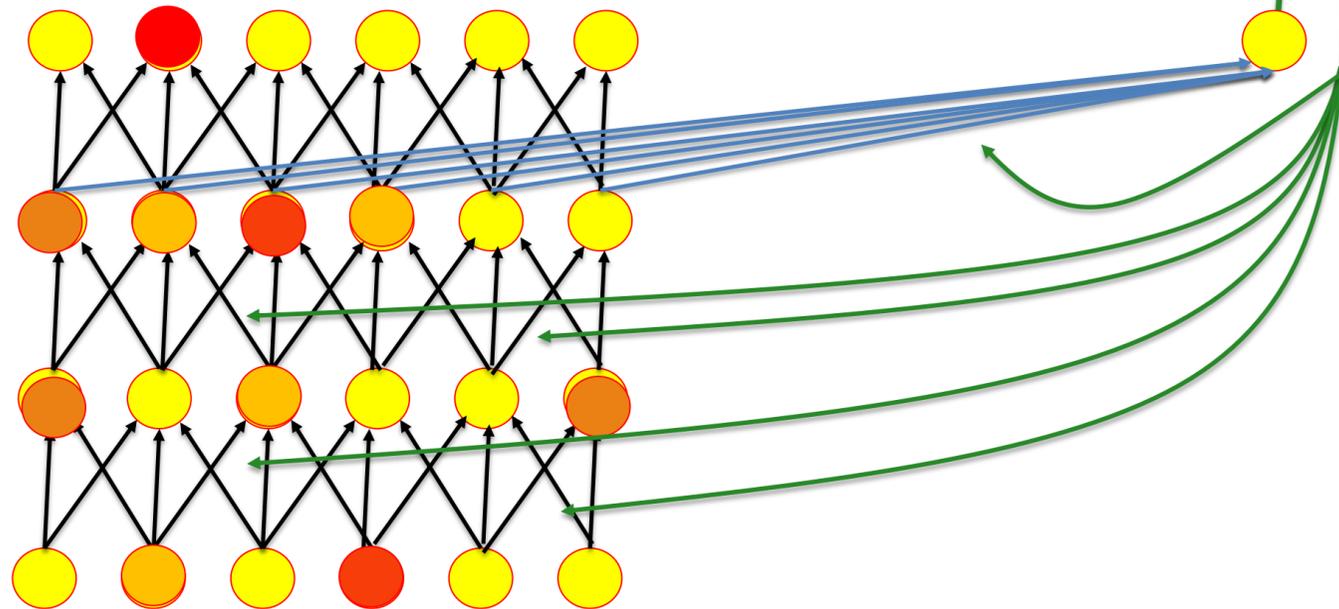
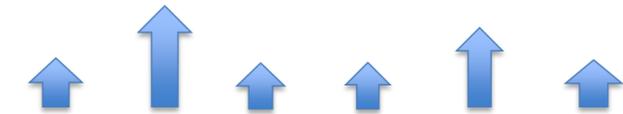
Deep reinforcement learning (alpha zero)

Silver et al. (2017) , Deep Mind

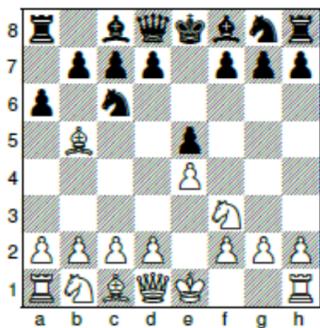
output: 4672 actions

Training 44Mio
games (9 hours)

advance king



Planning:
potential sequences
(during 1s before playing
next action)



input: 64x6x2x8 neurons
(about 10 000)

Previous slide.

Since there are many different positions, the number of input neurons is in the range of ten thousand:

On each of 64 positions there can be one of 6 different 'figures' (king, knight, bishop) of 2 different colors.

To avoid repetitions the 8 last time steps are used as input.

Training is done by playing against itself in 44 million games.

The allotted computer time for planning the next action is 1s.

Deep reinforcement learning (alpha zero)

Silver et al. (2017)

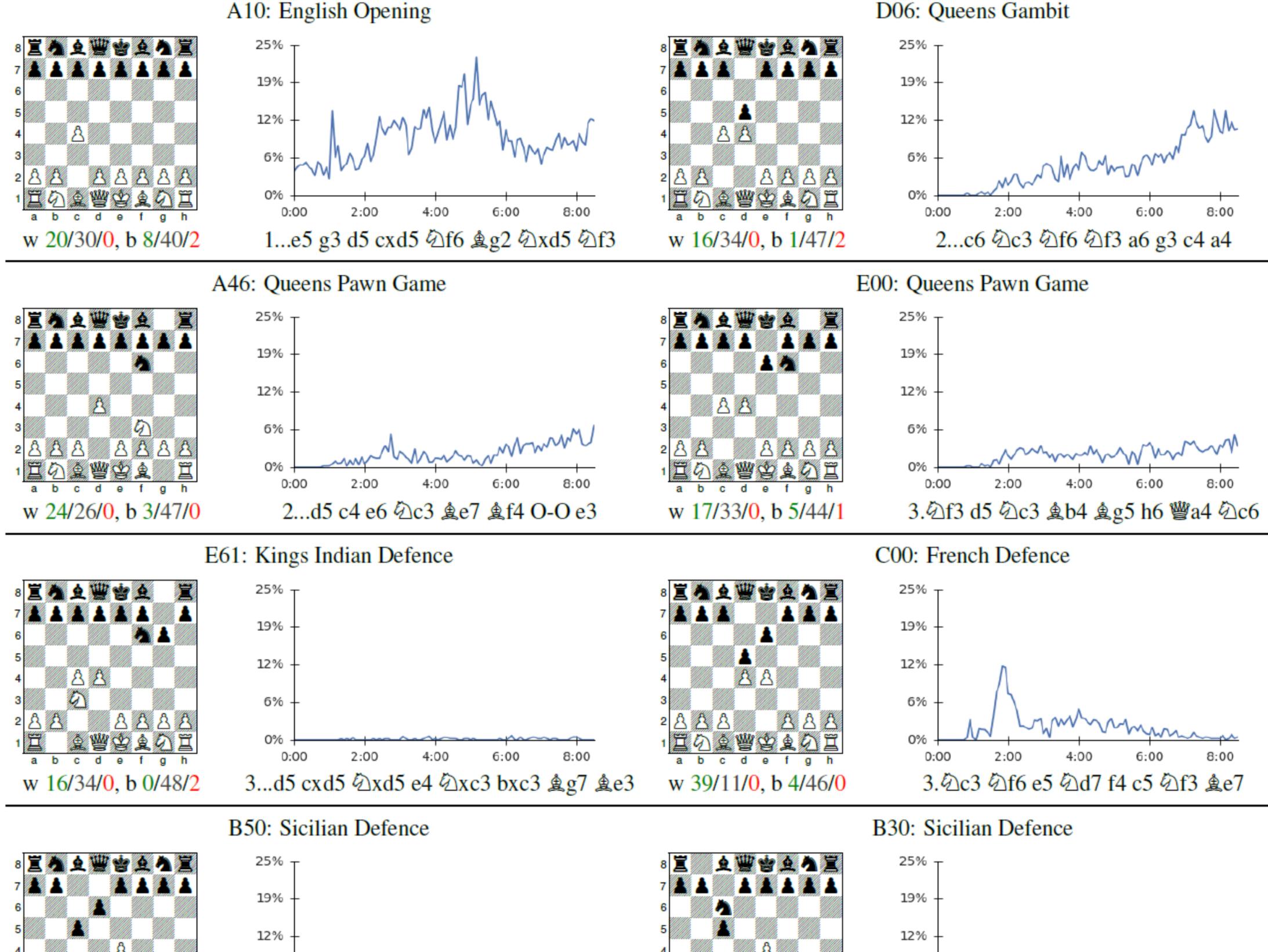
Chess:

- Trained for
44 Mio games

- discovers classic
openings

- beats best human
players

- beats best classic
AI algorithms



Previous slide.

After training for 44 Million self-play games, the algorithm matches or beats classical AI algorithms for chess.

Interestingly, it 'discovers' well-known strategies for openings, corresponding closely to well known openings in textbooks on chess.

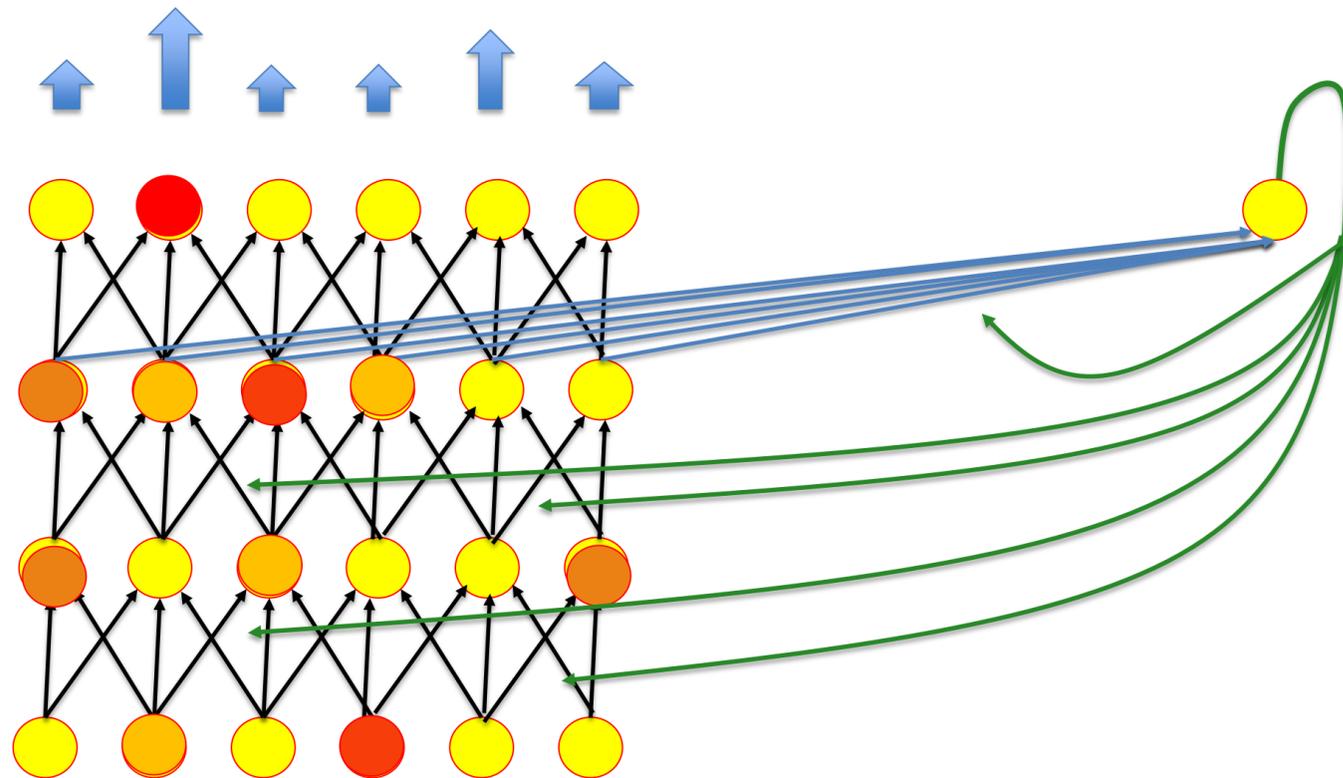
When trained on Go it beats the world champions.

Self-driving cars

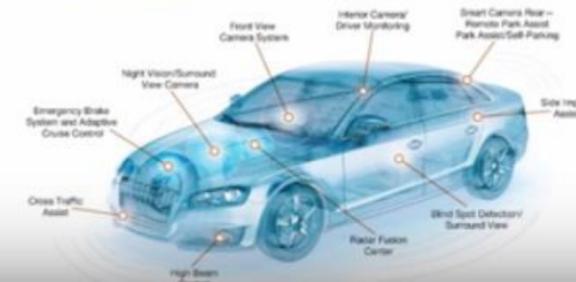
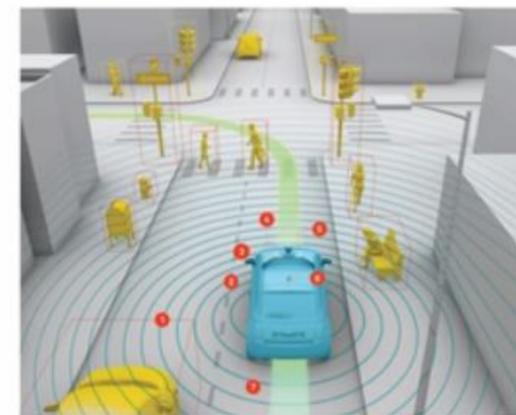
<https://selfdrivingcars.mit.edu/>

Lex Friedman, MIT

advance and accerate



Value: security,
duration of travel

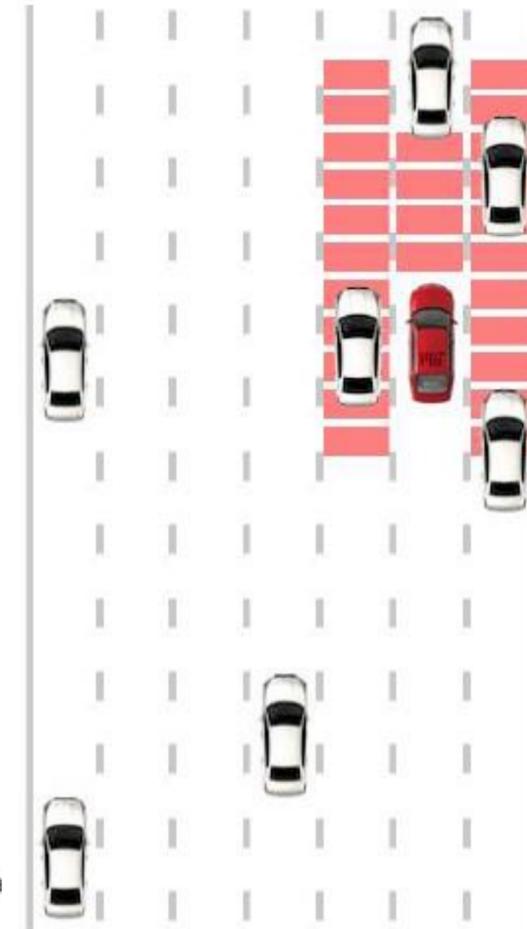


External

1. Radar
2. Visible-light camera
3. LIDAR
4. Infrared camera
5. Stereo vision
6. GPS/IMU
7. CAN
8. Audio

Internal

1. Visible-light camera
2. Infrared camera
3. Audio



Road Overlay:

Safety System ⇅

Previous slide.

Similar reinforcement learning algorithms are also used to train selfdriving cars.

There is a nice series of video lectures by Lex Friedman on the WEB.

Inputs are video images as well as distance sensors.

The value is security (top priority) combined with duration of travel.

The main focus of the class ANN is on reinforcement learning.

Large Language Models/ ChatGPT

- A large phase of supervised learning
(immense text corpora)
- Fine tuning with Reinforcement Learning
(Humans give Feedback: good/bad)

RLHF

Previous slide.

ChatGPT has revolutionized our conception of what language models can do. After a lot of expensive training on large data collections, additional fine tuning is necessary to adapt the system to make it an agreeable companion:
remove biases, remove mistakes etc

This is done by **Reinforcement Learning from Human Feedback (RLHF)**

https://en.wikipedia.org/wiki/Reinforcement_learning_from_human_feedback

Supervised Classification vs. Reinforcement Learning

Machine learning class of Jaggi-Flammarion covers
→ Supervised learning and Deep Networks

This class focuses on

Reinforcement learning/Reward-based learning

- Q-learning, SARSA, TD-learning, function approximation
- policy gradient, actor-critic, eligibility traces
- Deep Reinforcement learning, Applications
- Interdisciplinary Reinforcement Learning

Previous slide.

This class has two main parts:

Classification by Supervised learning.

- simple perceptron, geometry of classification (not covered, recap in Exerc.)
- deep learning in multilayer networks (not covered, recap in Exerc.)
- convolutional networks for image classification (not covered, recap in Exerc.)

Reinforcement learning/action learning driven by sparse rewards

- Q-learning, SARSA, TD-learning, function approximation (week2-5)
- Policy gradient methods and actor-critic (6+7)
- Deep Reinforcement learning
- Games and model-based reinforcement learning

Previous slide.

If you put a rat into an environment it will wander around. Suppose that, at some place, it discovers a food source hidden below the sand of the surface. After a couple of trials it will go straight to the location of the food source which implies that it has learned the appropriate sequence of actions in the environment to find the food source from arbitrary starting positions!

Quiz: Classification versus Reinforcement Learning

- Classification aims at predicting the correct category such as 'car' or 'dog'
- Classification by supervised learning is based on rewards
- Reinforcement learning is based on rewards
- Reinforcement learning aims at optimal action choices

Your notes.

Reinforcement Learning Lecture 1

Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Part 2: Elements of Reinforcement Learning

- Examples of Reward-based Learning
- **Elements of Reinforcement Learning**

Previous slide.

We now start with the formalization of reinforcement learning

Elements of Reinforcement Learning:

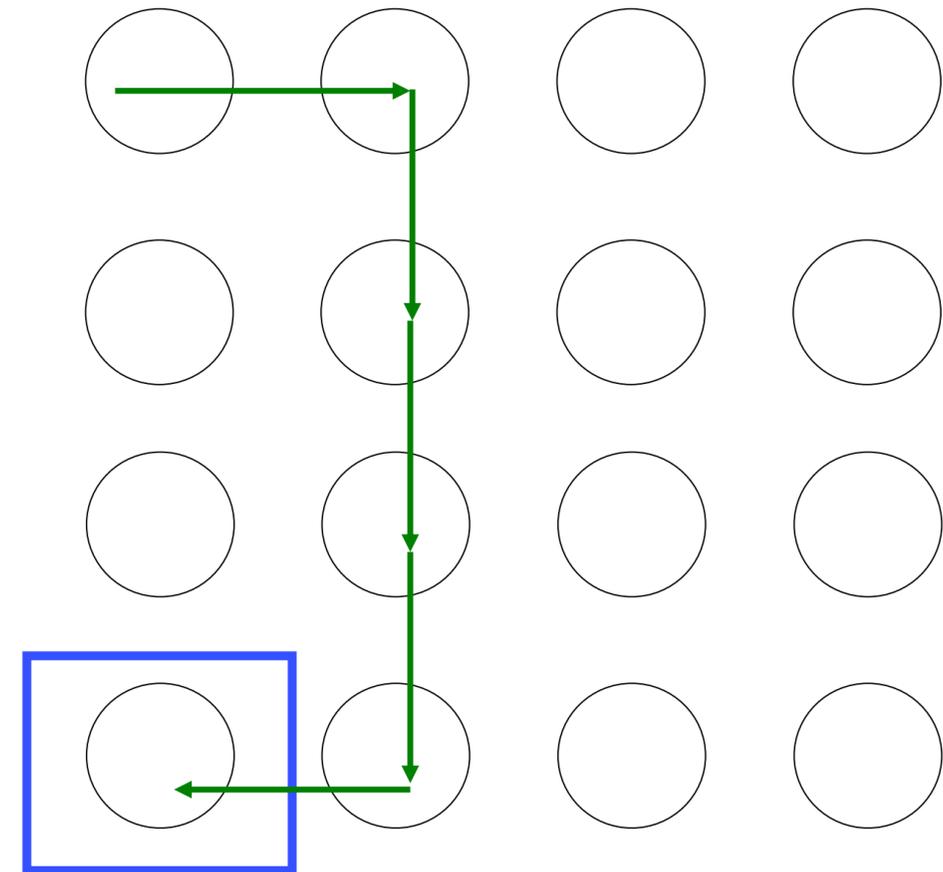
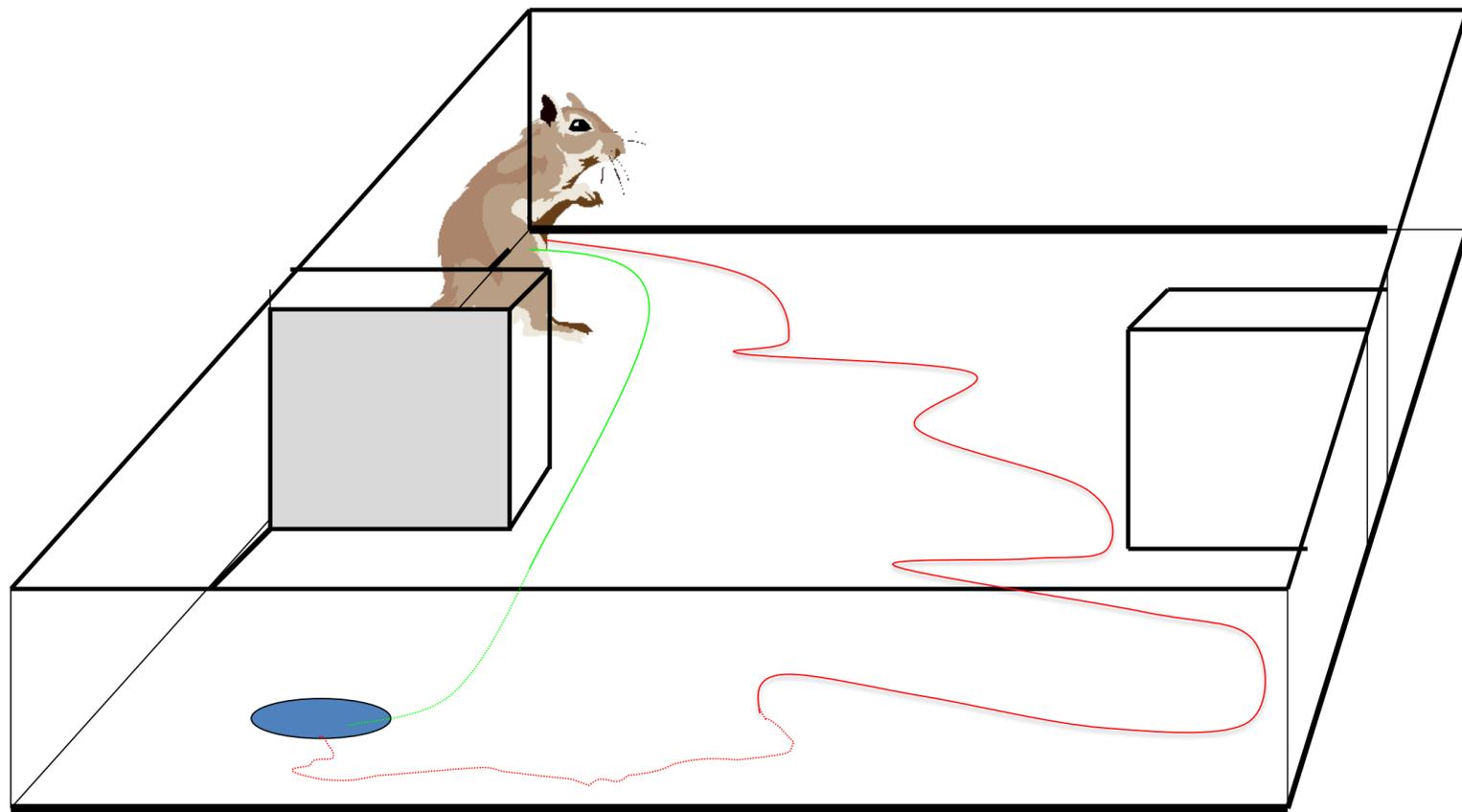
- states
- actions
- rewards

Previous slide.

Reinforcement learning needs states, actions, and rewards.

Elements of Reinforcement Learning:

- discrete states
- discrete actions
- sparse rewards



Previous slide.

Note that, for standard formulations of Reinforcement Learning Theories this (normally) implies discretizing space and actions.

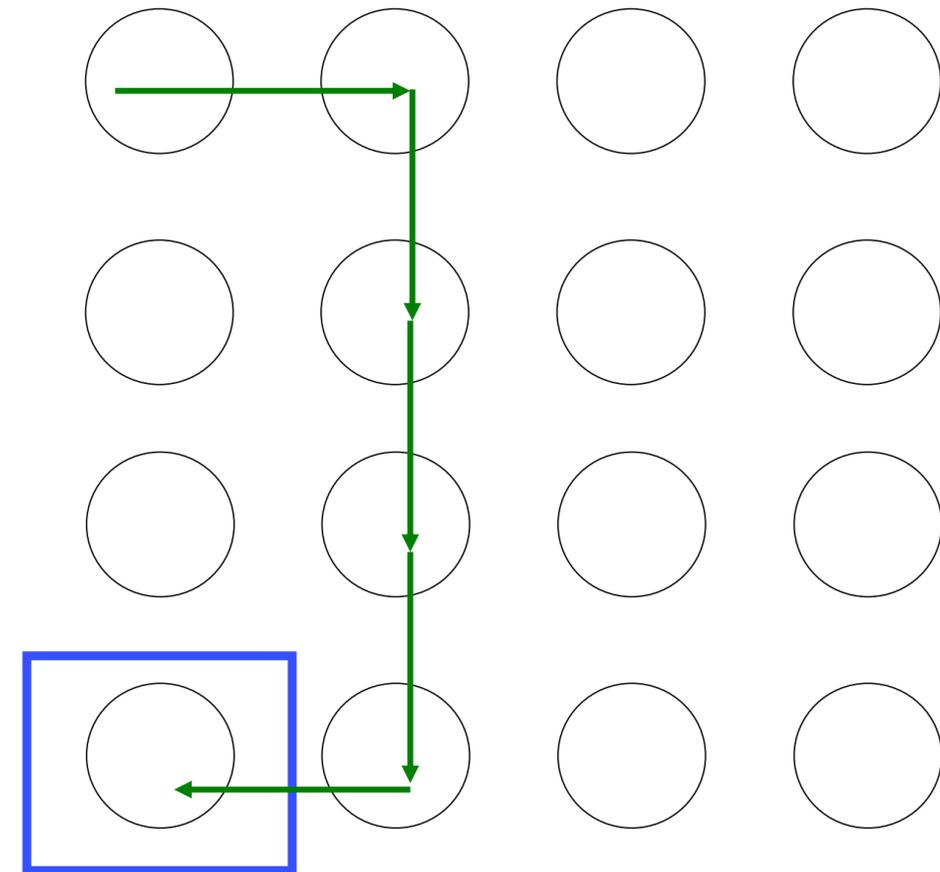
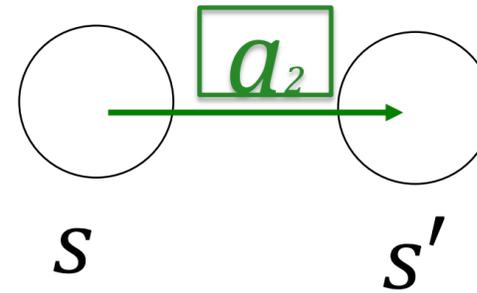
We will study continuous-space formulations only next week.

Elements of Reinforcement Learning:

- discrete states:
 - old state s
 - new state s'
- current state: s_t
- discrete actions: $a_1, a_2 \dots a_A$
- current action: a_t
- current reward: r_t
- Mean rewards for transitions:

$$R_{s \rightarrow s'}^a$$

often most transitions have zero reward



Previous slide.

The elementary step is:

The agent starts in state s .

It takes action a

It arrives in a new state s'

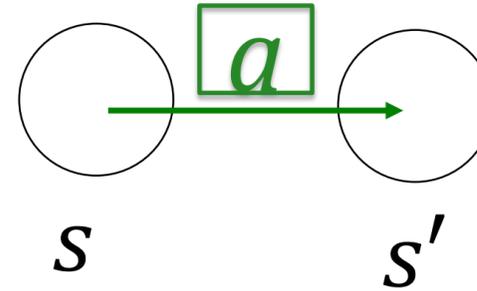
Potentially receiving reward r (during the transition or upon arrival at s').

Since rewards are stochastic we have to distinguish the mean reward at the transition (capital R with indices identifying the transition) from the actual reward (lower-case r with index t) that is received at time t on a transition.

Note that in many practical situations most transitions or states have zero rewards, except a single 'goal' state at the end.

States in Reinforcement Learning:

- discrete states:
 - starting state s
 - arrival state s'
- current state: s_t



state = current configuration/well-defined situation
= generalized 'location' of actor in environment

Previous slide.

What are these discrete states?

Loosely speaking a state is the current configuration that **uniquely** describes the momentary situation. We can think of the generalized 'location' of the actor in the environment

To get acquainted with this, let us look at an example.

Reinforcement Learning: Example Acrobot

3 actions: a_1 = no torque,
 a_2 = torque +1 at elbow,
States? a_3 = torque -1 at elbow

reward if tip above line

→ discretize!

**Suppose 5 states per dimension,
How many states in total?**

- [] 5
- [] 25
- [] 125
- [] 625

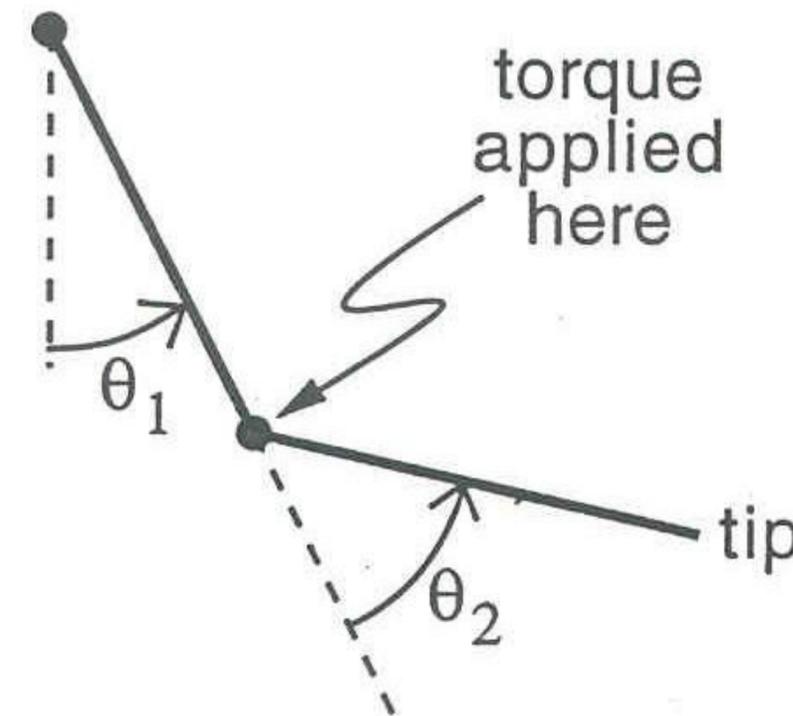


Figure 11.4 The acrobot.

*From Book:
Sutton and Barto*

Previous slide.

The aim of the acrobat is to move the tip above the blue line. To achieve this torque can be applied at the 'elbow' link. The second link is the 'shoulder'.

There are three possible actions.

But what are the states? How many states do we have?

Reinforcement Learning Lecture 1

Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Part 3: One-step horizon (bandit problems)

- Examples of Reward-based Learning
- Elements of Reinforcement Learning
- **One-step horizon (bandit problems)**

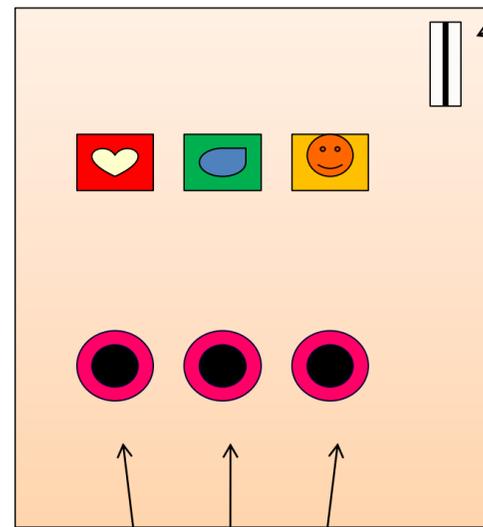
Previous slide.

We start with the simplest discrete example: the game is over and reward is given after a single step.

One-step horizon games (bandit)

action=button press

coins



Slot Machine
3-armed bandit

buttons

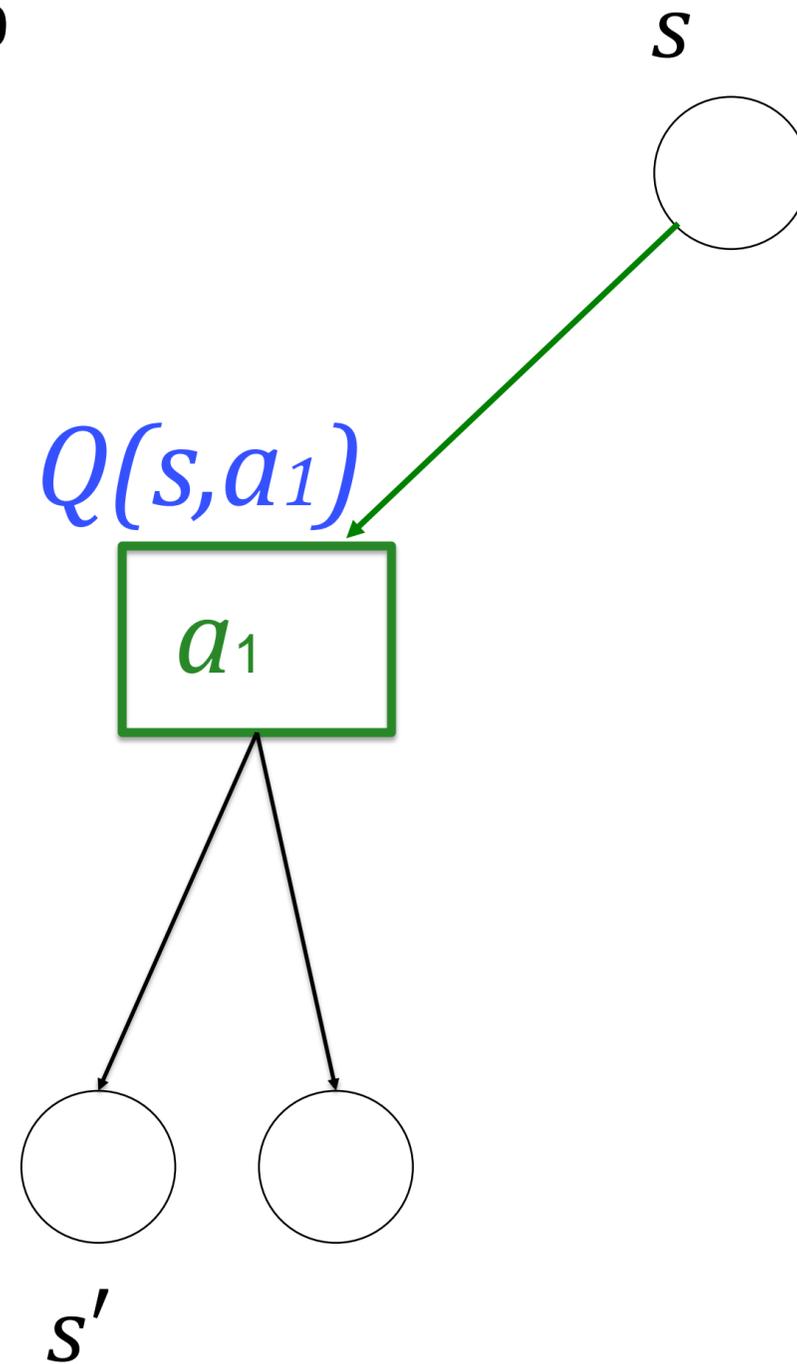
Previous slide.

The standard example is a multi-armed bandit, or slot machine: you have to choose between a few actions, and once you have pressed the button you can just wait and see whether you get reward or not.

One-step horizon games

Q-value: $Q(s,a)$

Expected reward for
action a starting from s



Blackboard1:
Q-values

Previous slide.

One of the most central notion in reinforcement learning is the Q-value.

$Q(s,a)$ has two indices: you start in state s and take action a .

The Q-value $Q(s,a)$ is (an estimate of) the mean expected reward that you will get if you take action a starting from state s .

One-step horizon games

Blackboard1:
Q-values

Your notes.

One-step horizon games: Q-value

Q-value $Q(s,a)$

Expected reward for action a starting from s

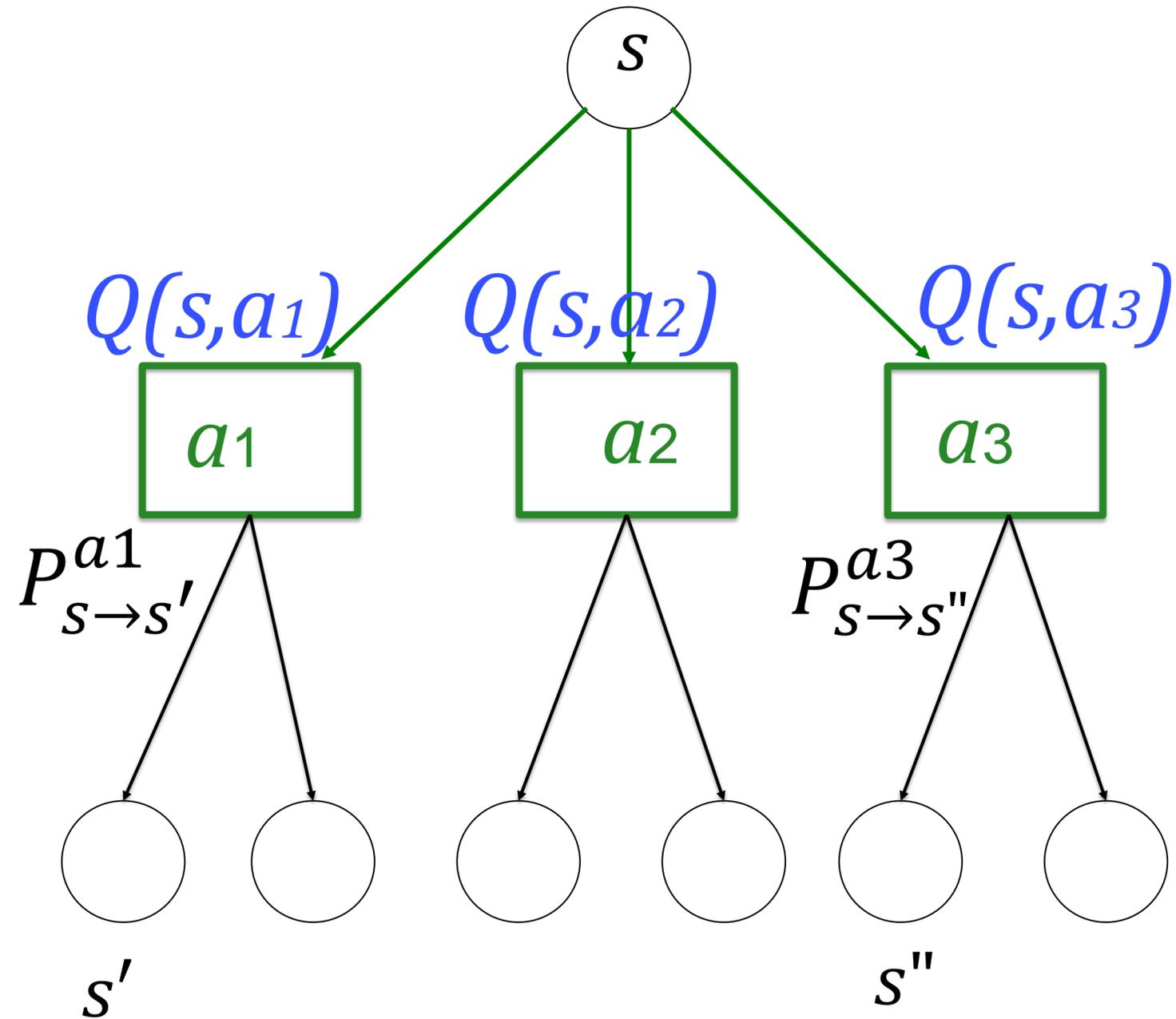
$$Q(s,a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a$$

Reminder:

$$R_{s \rightarrow s'}^a = E(r | s', a, s)$$

Similarly:

$$Q(s,a) = E(r | s, a)$$



Previous slide.

$P_{s \rightarrow s'}^{a1}$ is the probability that you end up in a specific state s' if you take action $a1$ in state s .

We refer to this sometimes as the 'branching ratio' below the 'actions'.

$Q(s,a)$ is attached to the branches linking the state s with the actions.

actions are indicated by green boxes; states are indicated by black circles.

The mean reward $R_{s \rightarrow s'}^a$ is defined as the expected reward given that you start in state s with action a and end up in state s' (see Blackboard 1).

Given the branching ratio and the mean rewards, it is easy to calculate the Q-values (Blackboard 1).