



Lee Sedol playing Go  
against AlphaGo

# AN INTRODUCTION TO REINFORCEMENT LEARNING

Guest lecture by Francis wyffels

 @fwyffels - [francis.wyffels@UGent.be](mailto:francis.wyffels@UGent.be)

*airplane*



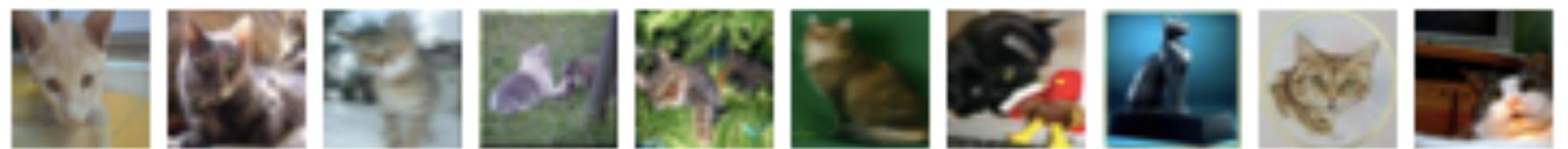
*automobile*



*bird*



*cat*



*deer*



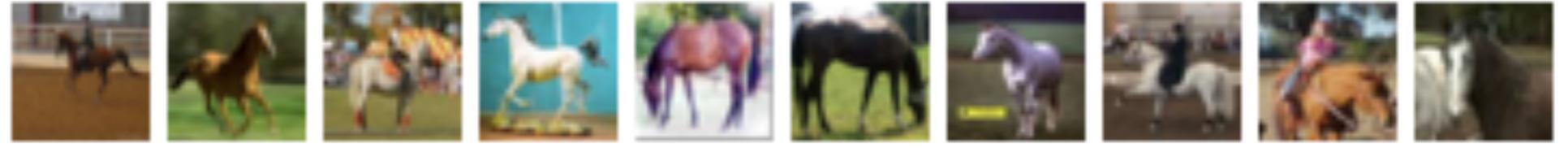
*dog*



*frog*



*horse*



*ship*



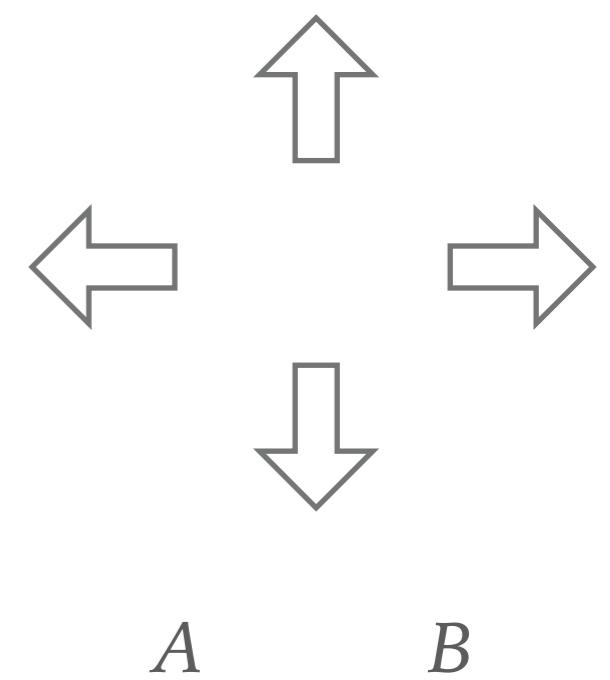
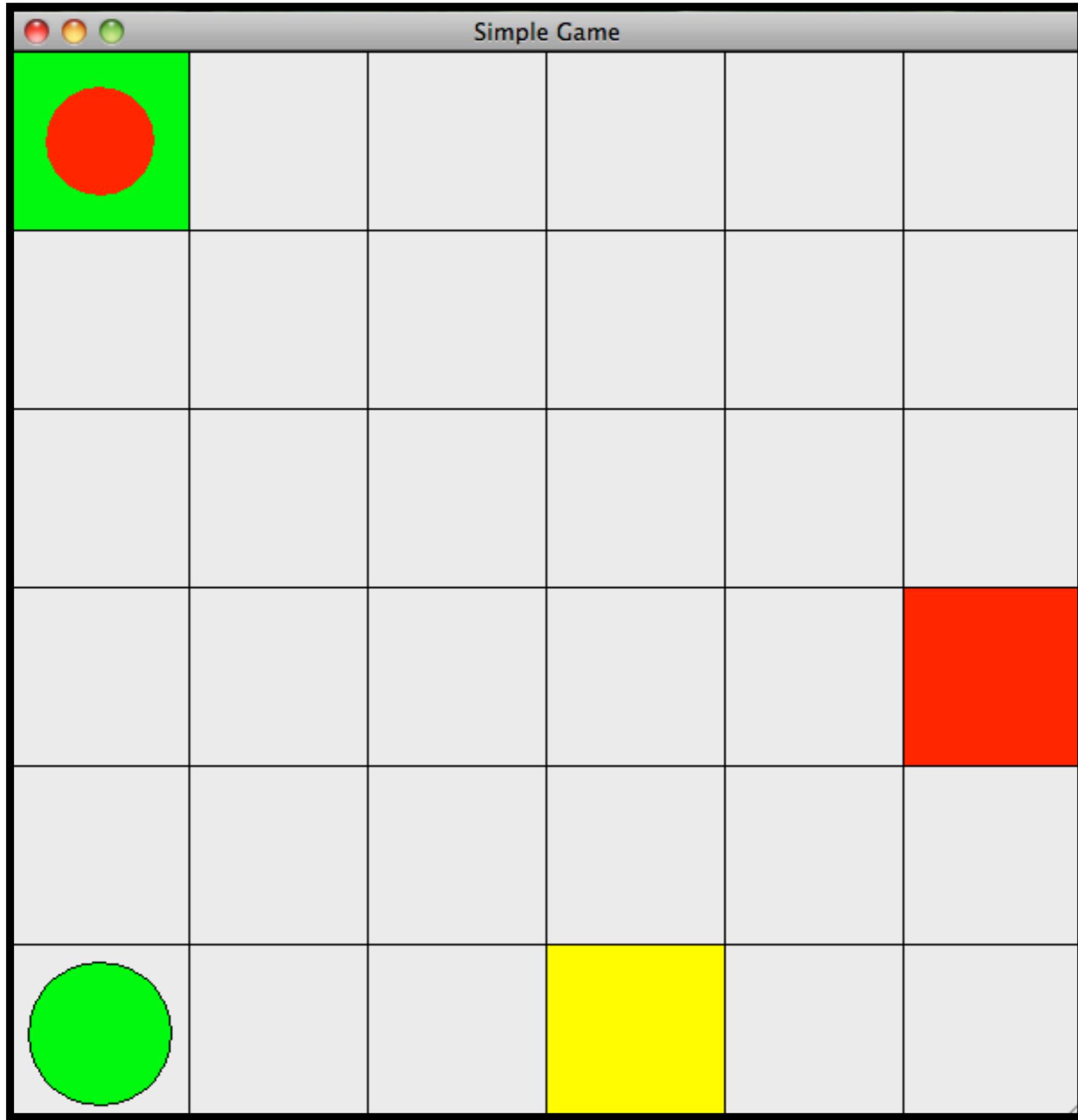
*truck*





[Stanford University - AI]

**LET'S PLAY A GAME**



“

Reinforcement learning is an interesting solution for a class of problems where *no explicit solution* is available and for which a *control policy* has to be found in order to maximise the *cumulative reward*.

*Discrete state space &  
discrete action space*

# PART I

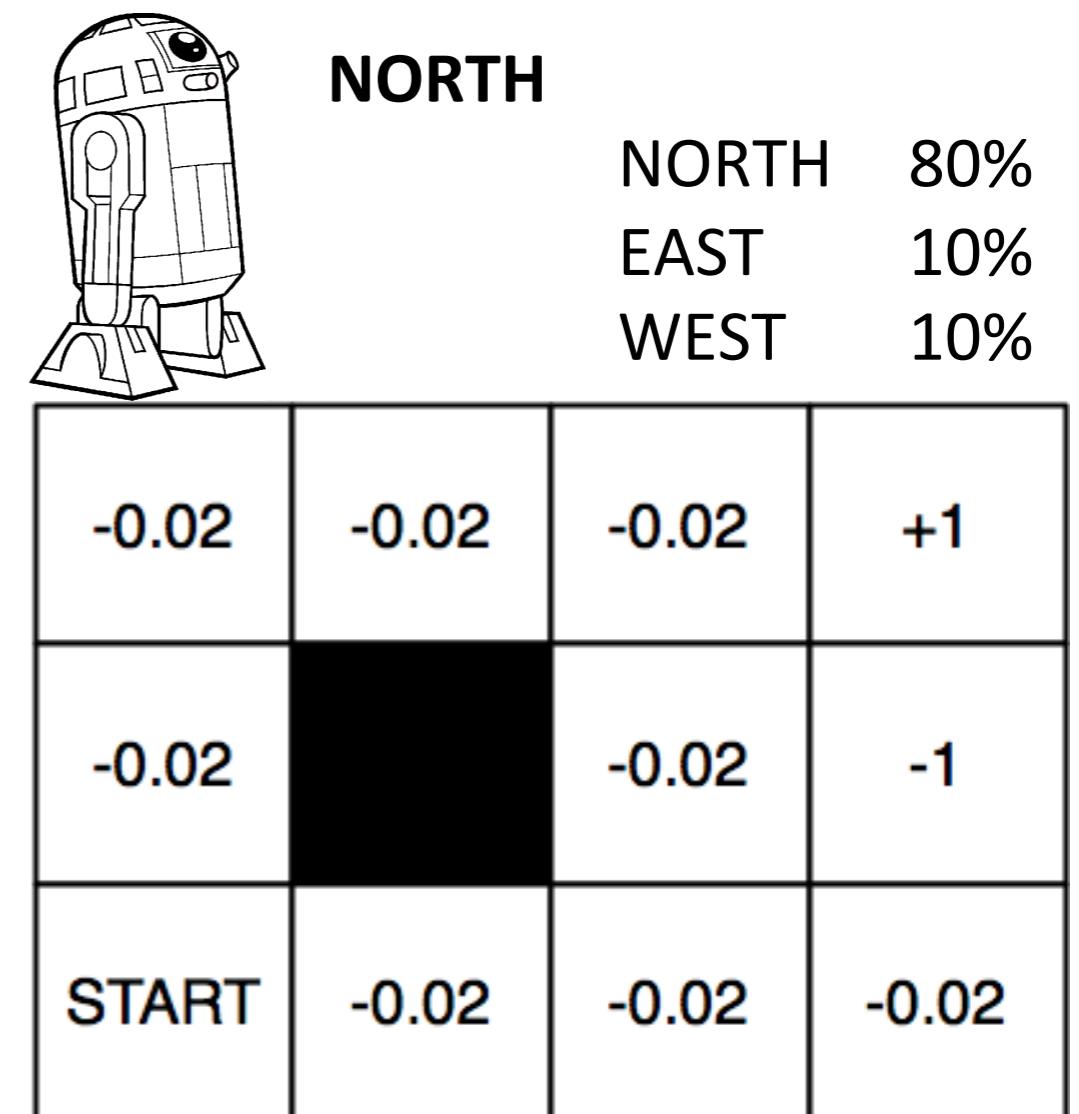
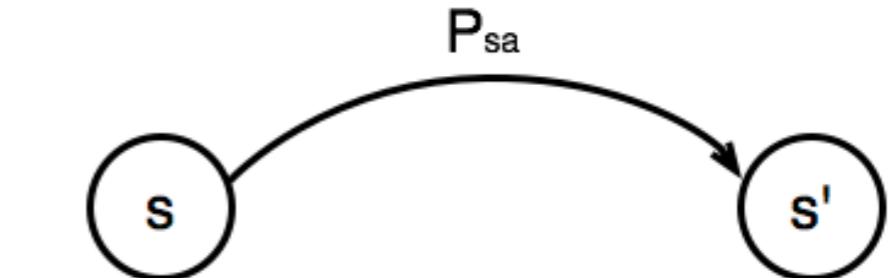


# MARKOV DECISION PROCESSES

- Also known as MDP's
- Defined by  $(S, A, \{P_{SA}\}, \gamma, R)$

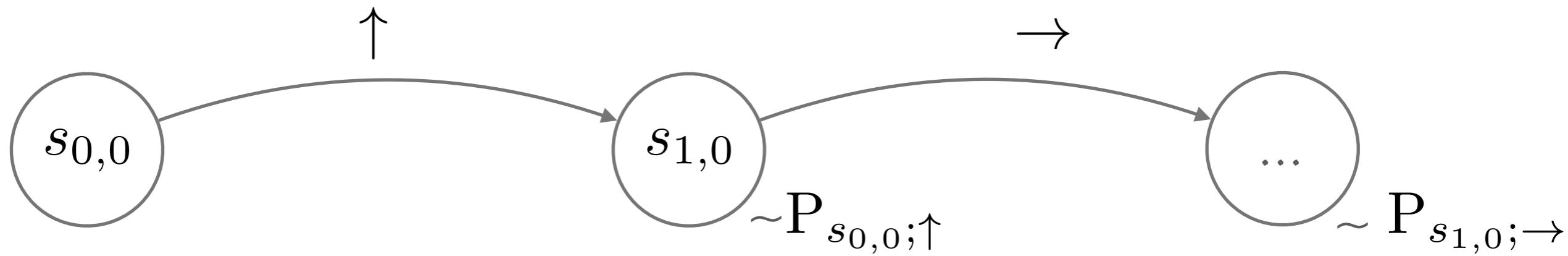
- ✓  $S$ : state
- ✓  $A$ : action
- ✓  $\{P_{SA}\}$ : state transition probability distribution
- ✓  $\gamma$ : discount factor
- ✓  $R$ : reward

$$\sum_{\forall s'} P_{sa}(s') = 1$$
$$P_{sa}(s') \geq 0$$



Possible actions: NORTH, EAST, SOUTH, WEST

# MARKOV DECISION PROCESSES: HOW IT WORKS



$$P_{s_{0,0};\uparrow}(s_{1,0}) = 0.8$$

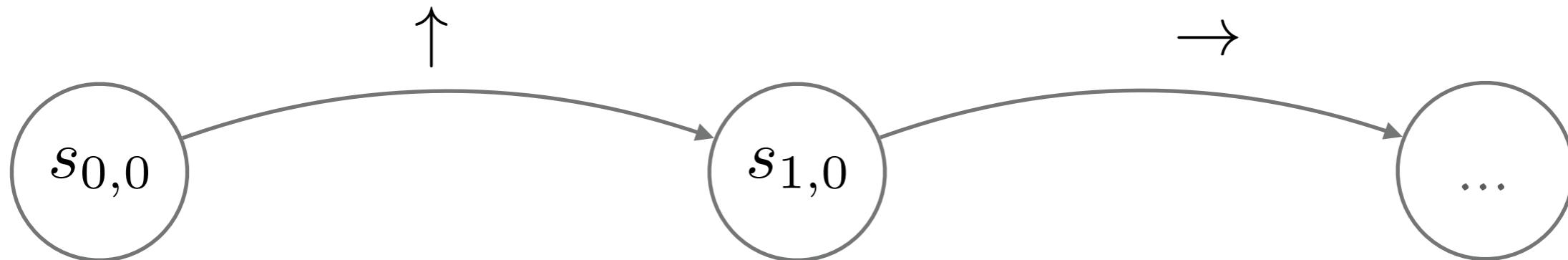
$$P_{s_{0,0};\uparrow}(s_{0,1}) = 0.1$$

$$P_{s_{0,0};\uparrow}(s_{0,0}) = 0.1$$

-0.02	-0.02	-0.02	+1
-0.02		-0.02	-1
START	-0.02	-0.02	-0.02

Possible actions: NORTH, EAST, SOUTH, WEST

# MARKOV DECISION PROCESSES: MAXIMISE THE TOTAL PAYOFF



$$R(s_{0,0}) + \gamma R(s_{1,0}) + \gamma^2 R(s_{\dots}) + \dots$$

*Goal: choose actions to maximise the expected value of the total payoff!*

$$E(R(s_{0,0}) + \gamma R(s_{1,0}) + \gamma^2 R(s_{\dots}) + \dots)$$

*For this we need to find a policy:*

$$\Pi : S \rightarrow A$$

-0.02	-0.02	-0.02	+1
-0.02		-0.02	-1
START	-0.02	-0.02	-0.02

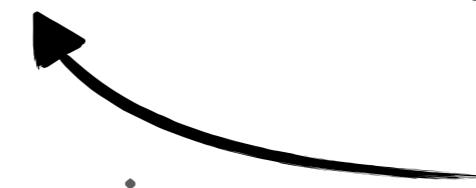
Possible actions: NORTH, EAST, SOUTH, WEST

# THE BELLMAN EQUATIONS

---

- Value function for a specific policy  $\pi$

$$V^\pi(s) = R(s) + \gamma \sum_{\forall s'} P_{s\pi(s)}(s') V^\pi(s')$$



- Optimal value function

$$V^*(s) = R(s) + \max_a \gamma \sum_{\forall s'} P_{sa}(s') V^*(s')$$

- Optimal policy

$$\pi^*(s) = \arg \max_a \sum_{\forall s'} P_{sa}(s') V^*(s')$$

# VALUE ITERATION

---

1. Initialise  $V(s) = 0 \quad \forall s$

2. For every  $s$  update

$$V(s) = R(s) + \max_a \gamma \sum_{\forall s'} P_{sa}(s') V(s')$$

3. Repeat until there are no changes

$$V(s) \rightarrow V^*(s)$$

Afterwards, get the optimal policy:

$$\pi^*(s) = \arg \max_a \sum_{\forall s'} P_{sa}(s') V^*(s')$$

# POLICY ITERATION

---

1. Initialise policy  $\pi$  arbitrarily

2. Solve  $V := V^\pi$

3. Let  $\pi(s) := \arg \max_a \sum_{\forall s'} P_{sa}(s')V(s')$

4. Repeat until there are no changes

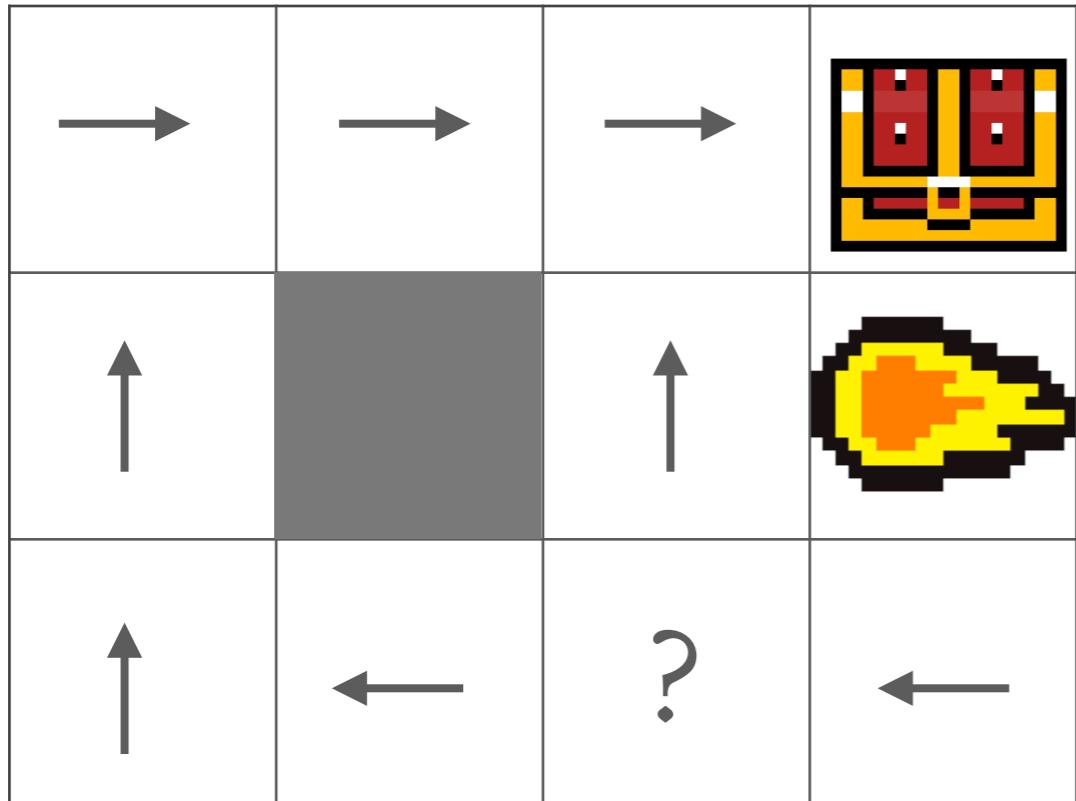
$$V(s) \rightarrow V^*(s)$$

$$\pi(s) \rightarrow \pi^*(s)$$

# SOME EXAMPLES

# WHAT IS THE OPTIMAL ACTION?

---

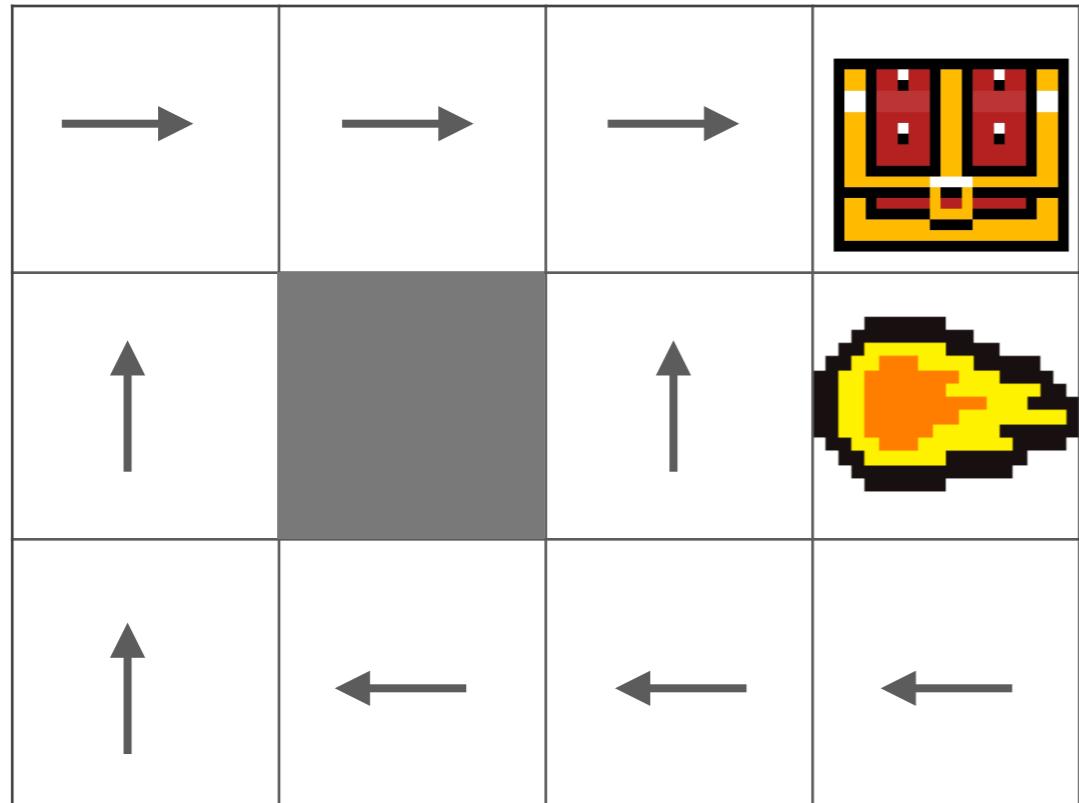


0.86	0.90	0.93	+1
0.82		0.69	-1
0.78	0.75	0.71	0.49

- $P_{SA} \sim \{0.8, 0.1, 0.1\}$
- $\gamma = 0.99$

# GOING WEST IS OPTIMAL

---



0.86	0.90	0.93	+1
0.82		0.69	-1
0.78	0.75	0.71	0.49

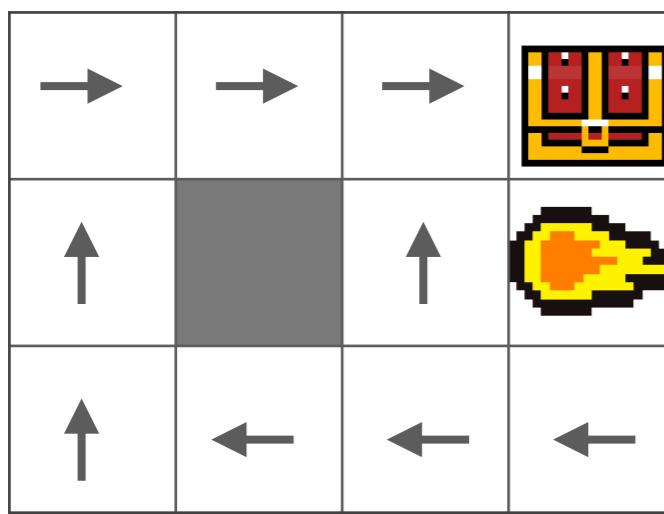
►  $P_{SA} \sim \{0.8, 0.1, 0.1\}$

►  $\gamma = 0.99$

$$0.8 * 0.69 + 0.1 * 0.75 + 0.1 * 0.49 = 0.676$$

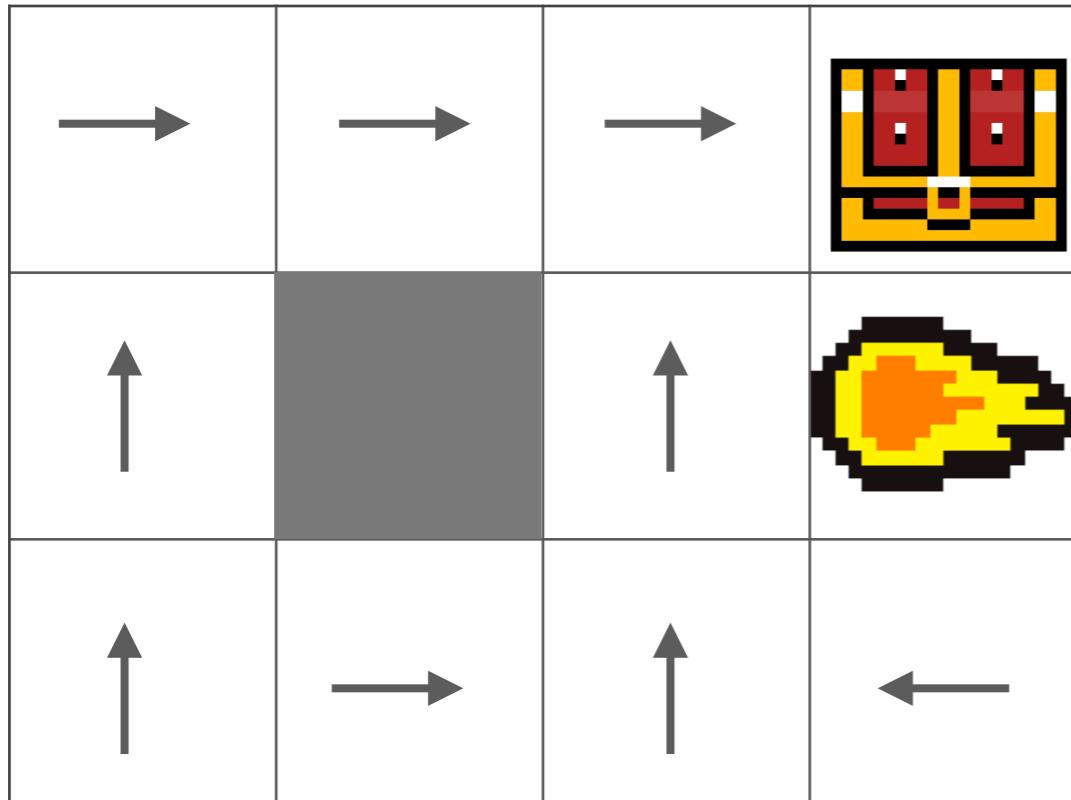


$$0.8 * 0.75 + 0.1 * 0.71 + 0.1 * 0.69 = 0.74$$



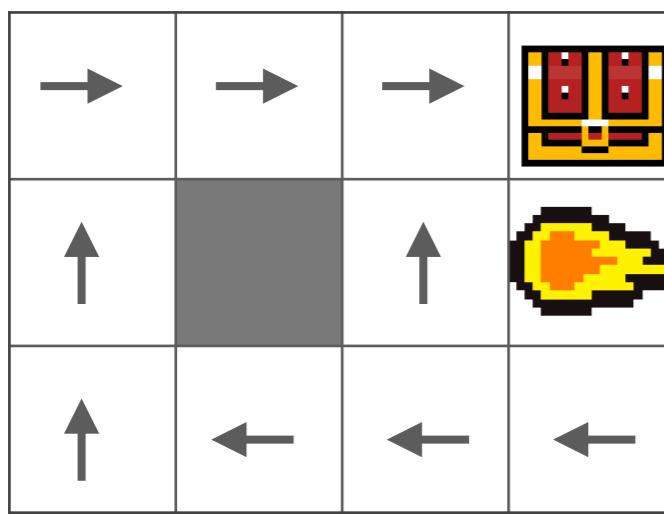
0.86	0.90	0.93	+1
0.82		0.69	-1
0.78	0.75	0.71	0.49

- $P_{SA} \sim \{0.8, 0.1, 0.1\}$
- $\gamma = 0.99$



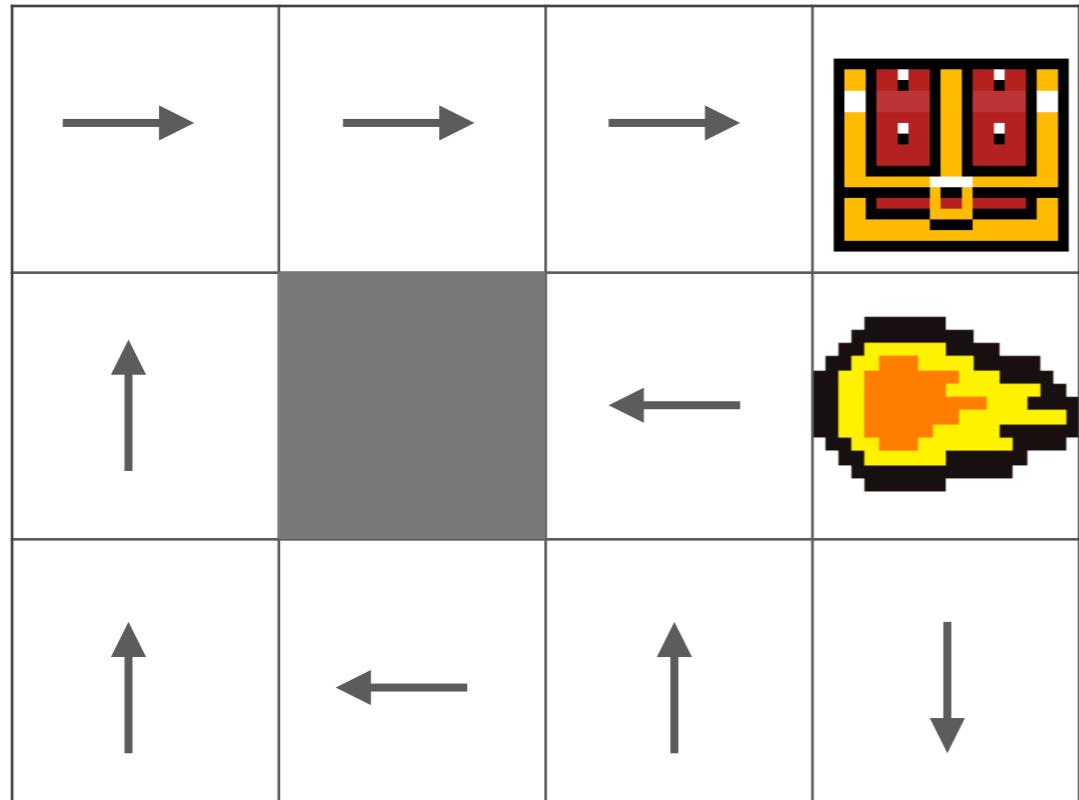
0.91	0.94	0.97	+1
0.88		0.93	-1
0.85	0.87	0.90	0.86

- $P_{SA} \sim \{0.99, 0.005, 0.005\}$
  - $\gamma = 0.99$
- the robot is very accurate*



0.86	0.90	0.93	+1
0.82		0.69	-1
0.78	0.75	0.71	0.49

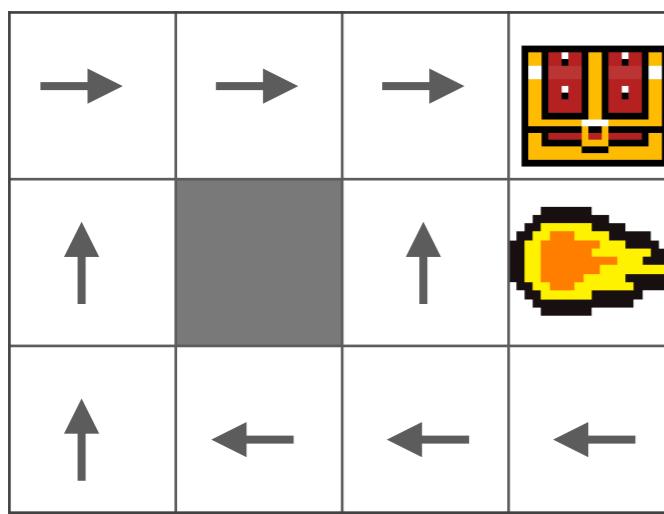
- $P_{SA} \sim \{0.8, 0.1, 0.1\}$
- $\gamma = 0.99$



0.70	0.78	0.83	+1
0.65		0.69	-1
0.57	0.52	0.53	0.43

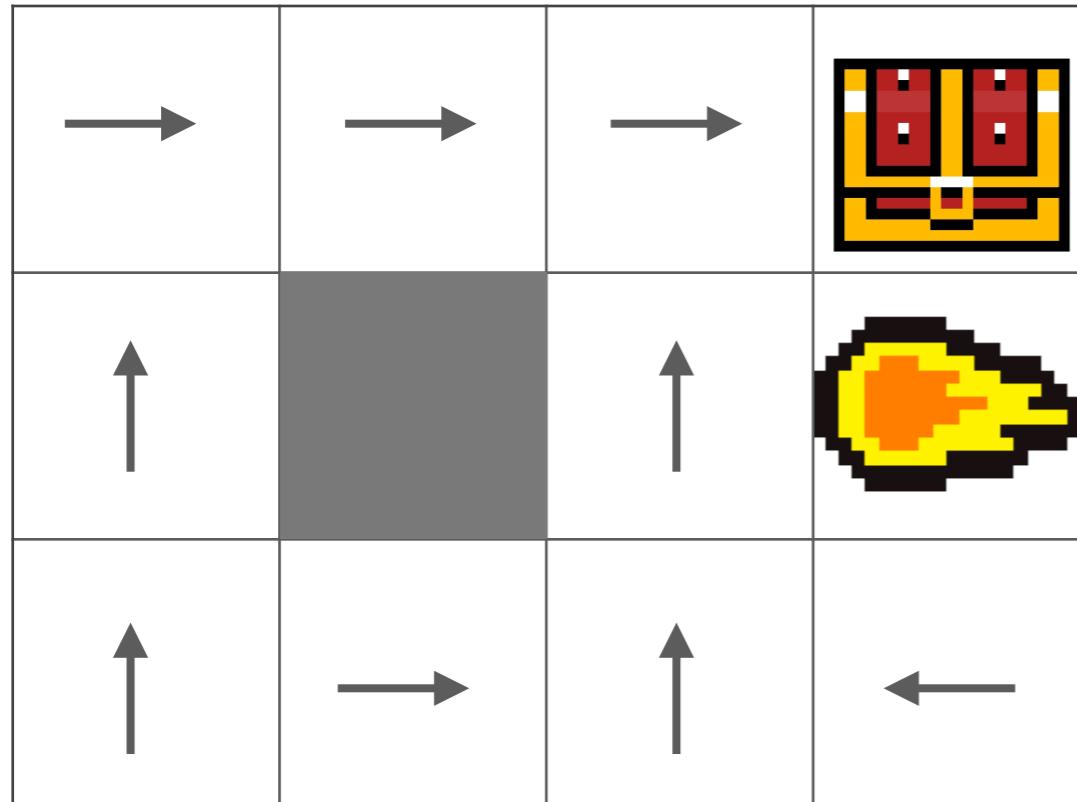
- $P_{SA} \sim \{0.5, 0.25, 0.25\}$
- $\gamma = 0.99$

the robot is very sloppy



0.86	0.90	0.93	+1
0.82		0.69	-1
0.78	0.75	0.71	0.49

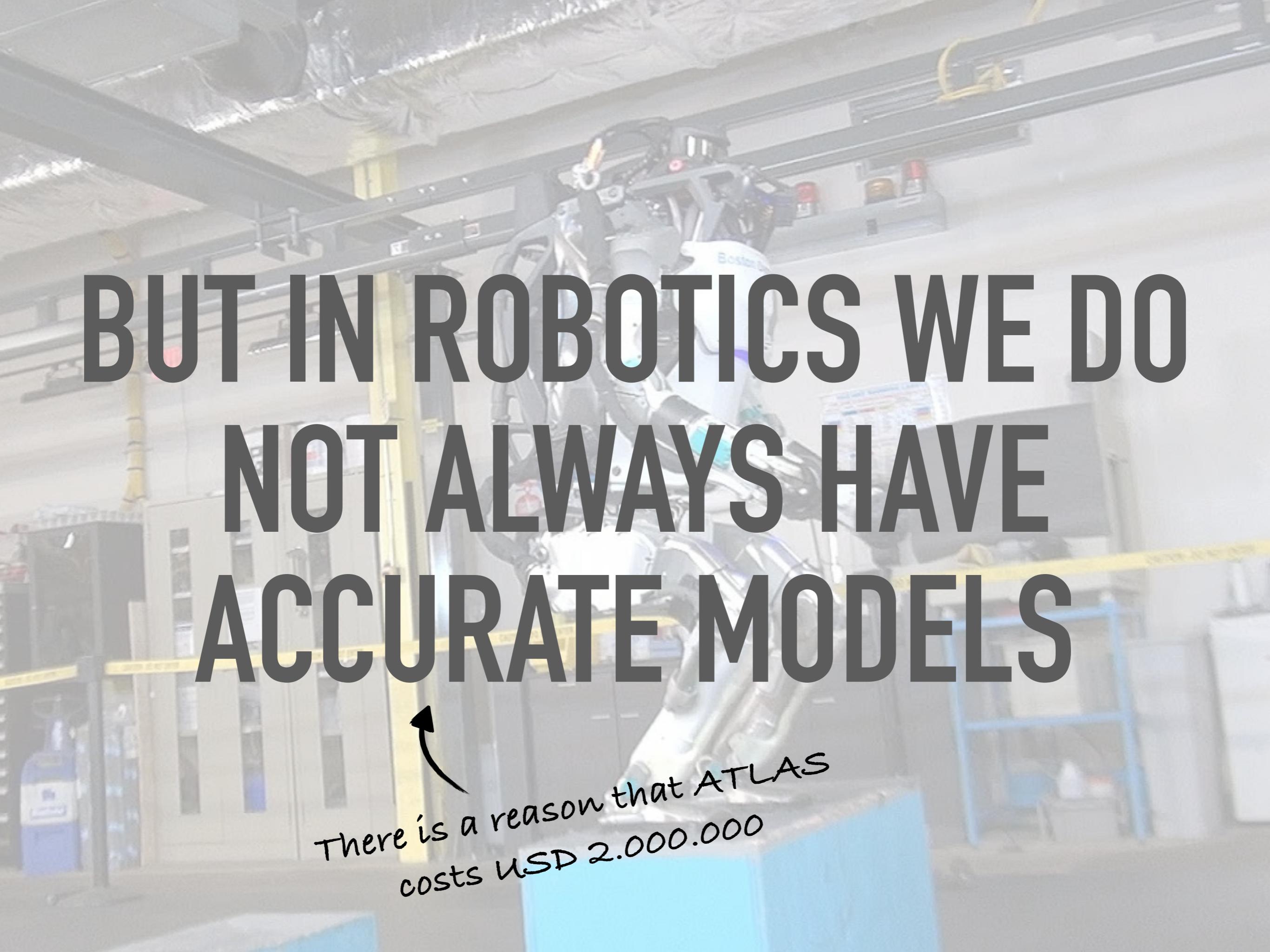
- $P_{SA} \sim \{0.8, 0.1, 0.1\}$
- $\gamma = 0.99$



0.58	0.70	0.82	+1
0.48		0.52	-1
0.39	0.34	0.41	0.20

- $P_{SA} \sim \{0.8, 0.1, 0.1\}$
- $\gamma = 0.9$

slight decrease in  
discount, dramatic effects!



**BUT IN ROBOTICS WE DO  
NOT ALWAYS HAVE  
ACCURATE MODELS**

There is a reason that ATLAS  
costs USD 2.000.000

# WHAT IF THERE IS NO MODEL AVAILABLE OF THE ENVIRONMENT

---

- $P_{sa}(s')$  is not available
  - value iteration & policy iteration won't work
- One solution is to estimate  $P_{sa}(s')$  based on interactions within the environment

$$P_{sa}(s') = \frac{\# \text{ } a \text{ taken in } s \text{ and got to } s'}{\# \text{ } a \text{ taken in } s}$$

- Use Temporal Difference (TD) learning

# TEMPORAL DIFFERENCE LEARNING: TD(0)

---

1. Initialise  $V(s)$  arbitrarily, evaluate policy  $\pi$
2. Repeat until convergence
  1. Initialise  $s$
  2. Repeat (until  $s$  is terminal)
    - $a$  is an action given by  $\pi$
    - do  $a$ , observe reward  $r$  and next state
$$V(s) = V(s) + \alpha(r + \gamma V(s') - V(s))$$
    - update  $s$  to  $s'$

*No need for a model of the environment,  
implemented in an online, incremental fashion*

# TEMPORAL DIFFERENCE LEARNING: TD(0)

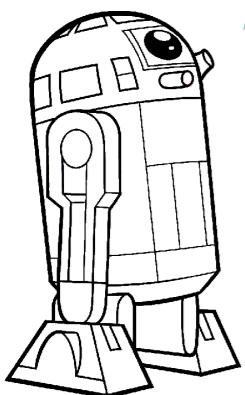
---

1. Initialise  $V(s)$  arbitrarily, evaluate policy  $\pi$
2. Repeat until convergence:

What about finding the optimal control policy?

$$\text{New } V(s) = V(s) + \alpha(r + \gamma V(s') - V(s))$$

- update  $s$  to  $s'$



*No need for a model of the environment,  
implemented in an online, incremental fashion*

# POLICY IMPROVEMENT

---

- It's all about the question whether action  $a$  improves the policy  $\pi$  or not!
- Define the quality (the payoff) of a state-action combination:

$$Q : S \times A \rightarrow \mathbb{R}$$

$$Q^\pi(s, a) = \mathbb{E}[r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s, a_t = a]$$

# SOLUTION: Q-LEARNING

---

- Q-learning: initialise  $Q$  randomly and update the  $Q$  function every time step  $k$  using the following equation

$$Q(s_k, a_k) := (1 - \alpha)Q(s_k, a_k) + \alpha[r_{k+1} + \gamma \max_a Q(s_{k+1}, a)]$$

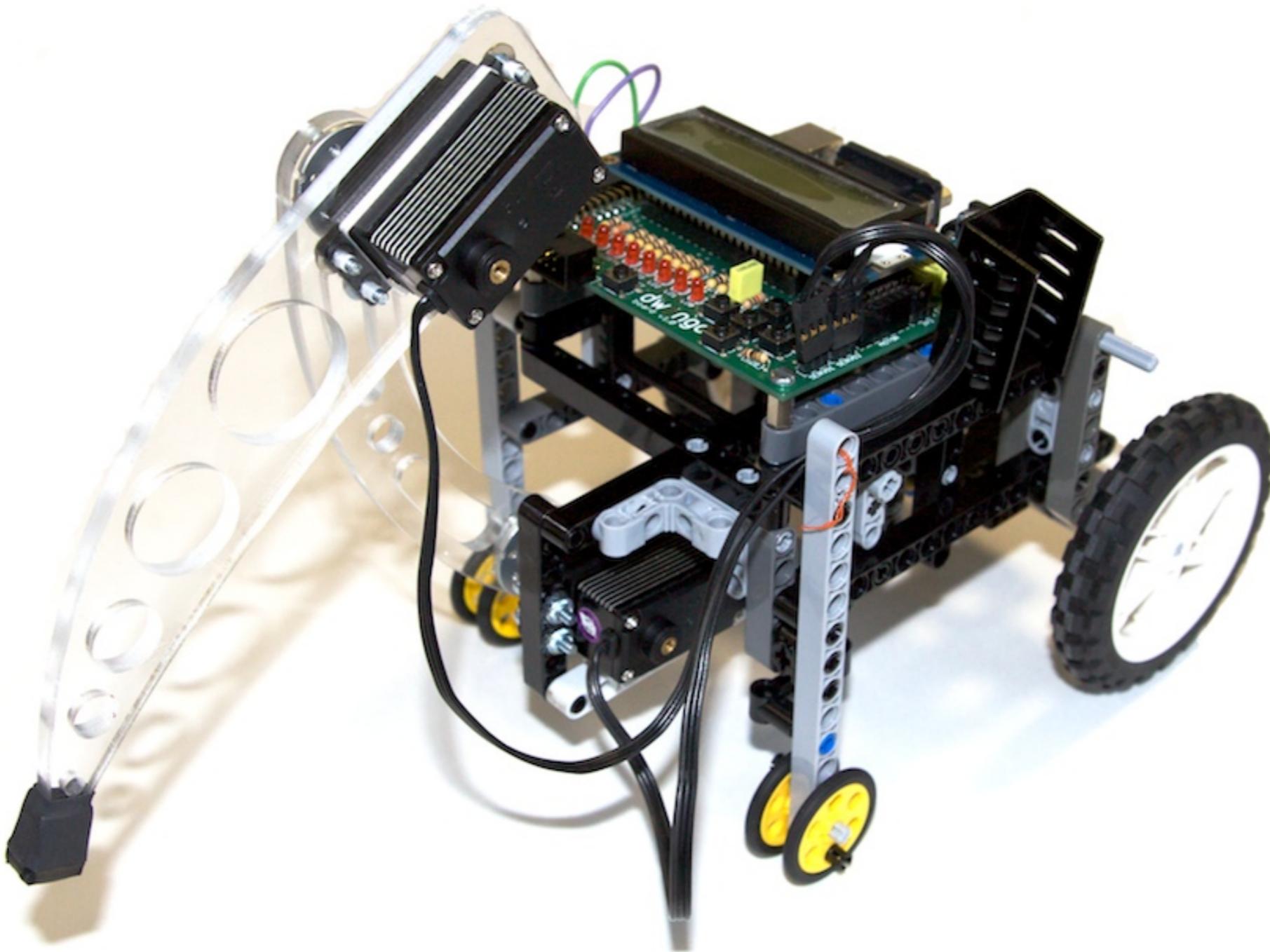

- Selection of the action can be done in an  $\varepsilon$ -greedy way
- This involves an exploration-exploitation trade-off

↑  
avoid that the agent gets stuck!

↑  
off-policy learner

# A CRAWLING ROBOT EXAMPLE

---



[<https://www.youtube.com/watch?v=2iNrJx6IDEo>]

# A CRAWLING ROBOT EXAMPLE

---

1st UP	1st	2nd UP	2nd
-1.2	0.6	0.9	1.2
0.2	-0.8	-0.5	-1.6
-1.4	0.5	1	-1.1
0.9	-0.9	0.6	0.6
-0.1	1.1	-0.1	0.2
0.2	1	0.9	-1
1.1	-0.1	-0.7	1.9
-0.8	-0.8	-1.1	0.5
0.5	-0.3	1	-0.3
-0.6	1.2	-0.6	-0.9
-1	1.3	0.2	-0.1
-0.1	-0.5	-0.9	1
0.9	0.2	1.1	-0.6
0.8	-0.1	0.3	1.3
-0.1	-1.4	1	-0.4
-1	-0.1	-1.1	0.8

# A CRAWLING ROBOT EXAMPLE

1st UP	1st	2nd UP	2nd
-1.2	0.6	0.9	1.2
0.2	-0.8	-0.5	-1.6
-1.4	0.5	1	-1.1
0.9	-0.9	0.6	0.6
-0.1	1.1	-0.1	0.2
0.2	1	0.9	-1
1.1	-0.1	-0.7	1.9
-0.8	-0.8	-1.1	0.5
0.5	-0.3	1	-0.3
-0.6	1.2	-0.6	-0.9
-1	1.3	0.2	-0.1
-0.1	-0.5	-0.9	1
0.9	0.2	1.1	-0.6
0.8	-0.1	0.3	1.3
-0.1	-1.4	1	-0.4
-1	-0.1	-1.1	0.8

The diagram shows a crawling robot's path on a grid. The path starts at the bottom-left cell (-2, -6) and ends at the top-right cell (112, 1). The path moves right to (35, -6), up-right to (37, -5), up to (1, -5), up-right to (17, -5), up to (17, 1), up-right to (112, 1), and finally up-left to (112, 124). The path is highlighted with thick grey arrows.

1st UP	1st	2nd UP	2nd
5	-15	-29	10
12	16	3	28
41	19	19	24
6	11	-8	-18
-4	5	36	11
0	18	13	-16
-2	35	17	112
124	37	6	5
-35	-80	17	-32
-16	-10	-25	-9
-3	93	1	0
-59	-6	22	-40
-4	38	13	1
8	-4	17	14
16	5	15	23
-35	-34	-39	-12

## PART II

---

*Continuous state space &  
discrete action space*



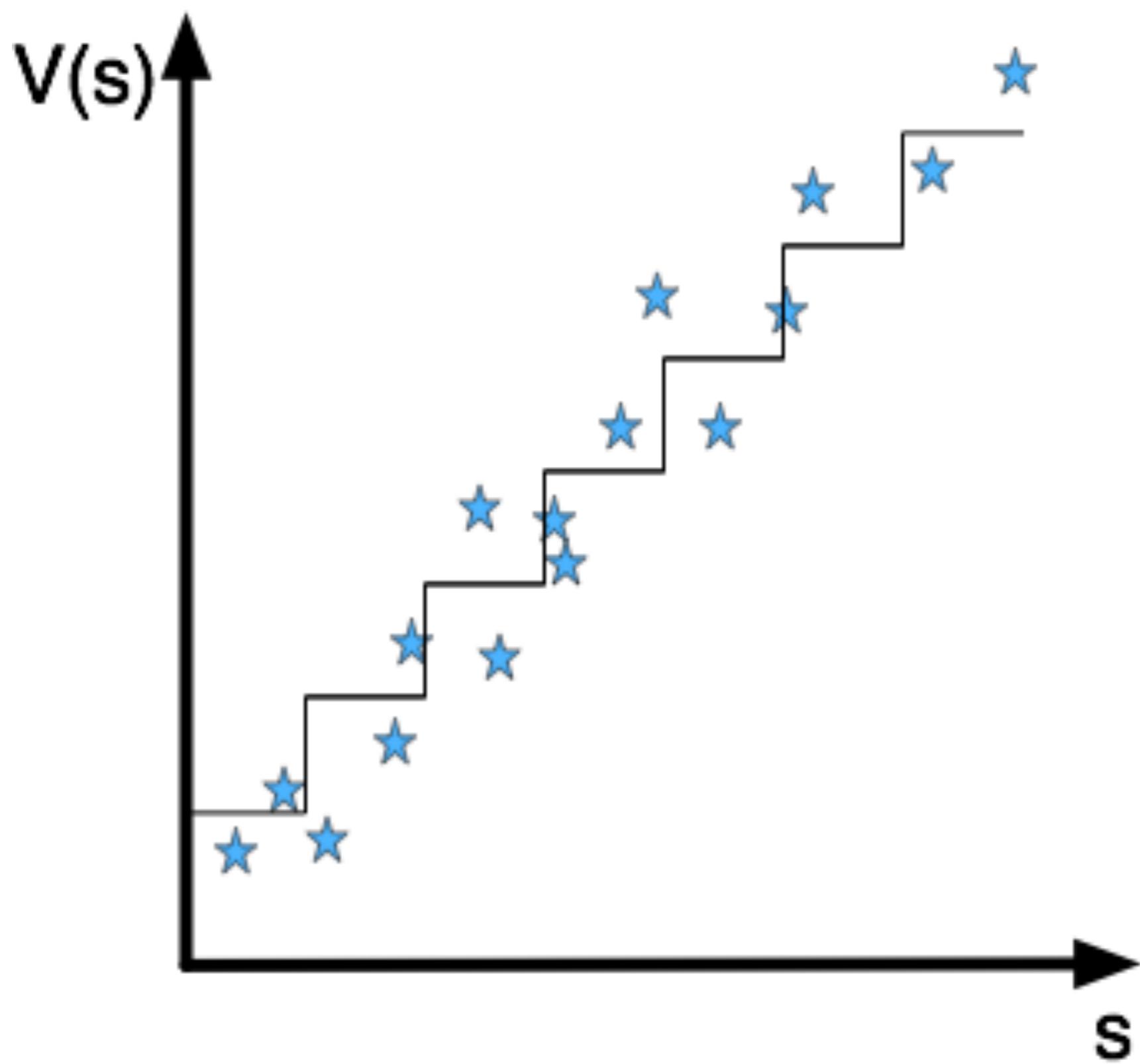
# CONTINUOUS STATE MDP'S

---


$$x, y, z, \alpha, \beta, \gamma, \dot{x}, \dot{y}, \dot{z}, \dot{\alpha}, \dot{\beta}, \dot{\gamma}$$

# DISCRETISATION LEADS TO POOR REPRESENTATIONS OF THE DATA

---





## DISCRETISATION: A CURSE

---

- the curse of dimensionality
  
- Let's say  $s \in \mathbb{R}^n$
  
- If  $s$  can be divided into  $k$  grids, then  $k^n$  possibilities
  
- For chess is this already  $64^{\wedge} 32$  (minus some illegal positions...)

*Tip: see [https://en.wikipedia.org/wiki/Game\\_complexity](https://en.wikipedia.org/wiki/Game_complexity) for game complexities*

# IN SUMMARY: CONTINUOUS STATE SPACE

---

- Too many states to visit
- Too many states to memorise
- Solution: generalise!

# THE BACKGAMMON CASE

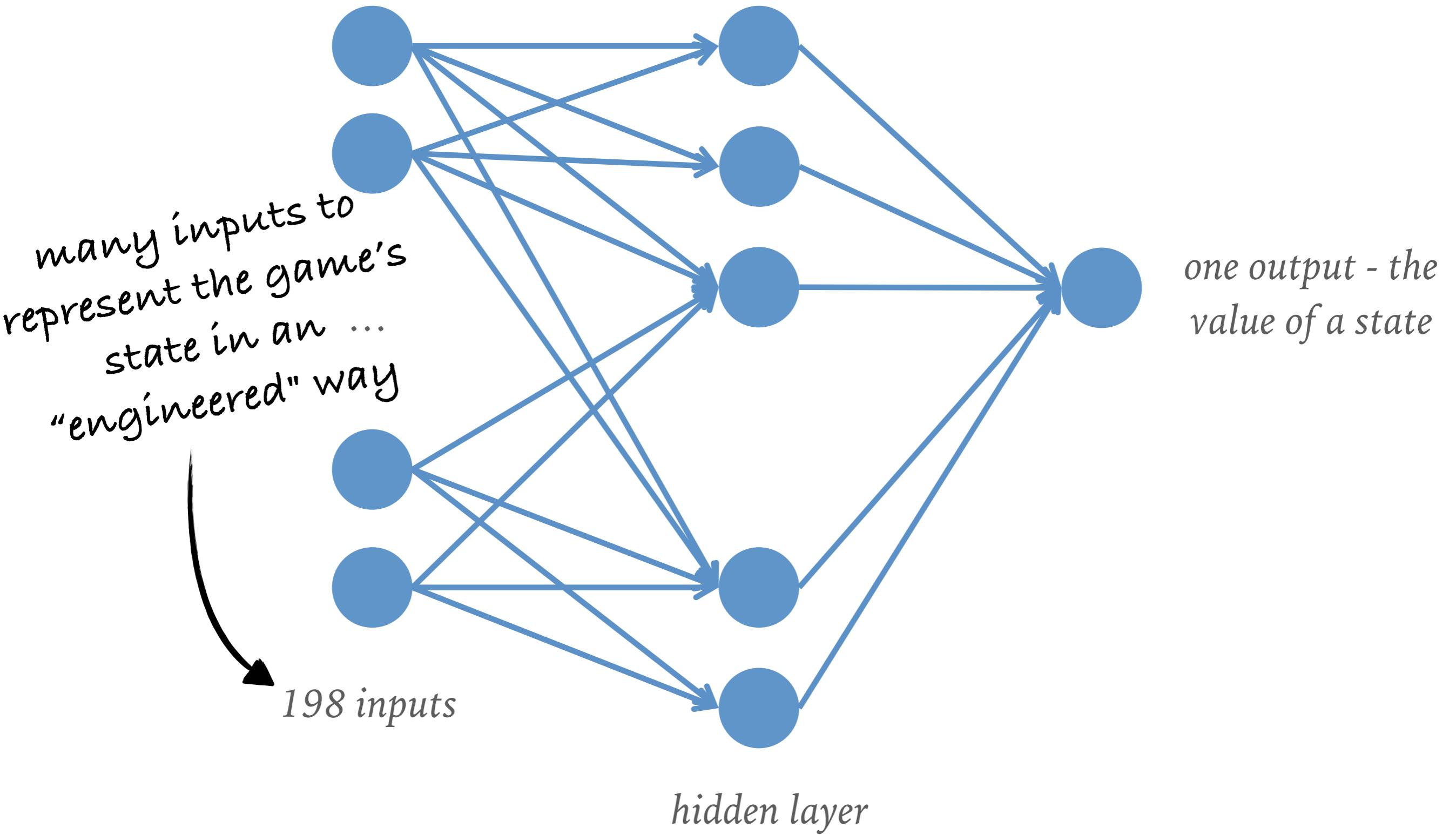
# THE TD-GAMMON CASE BY TESAURO, 1995

---



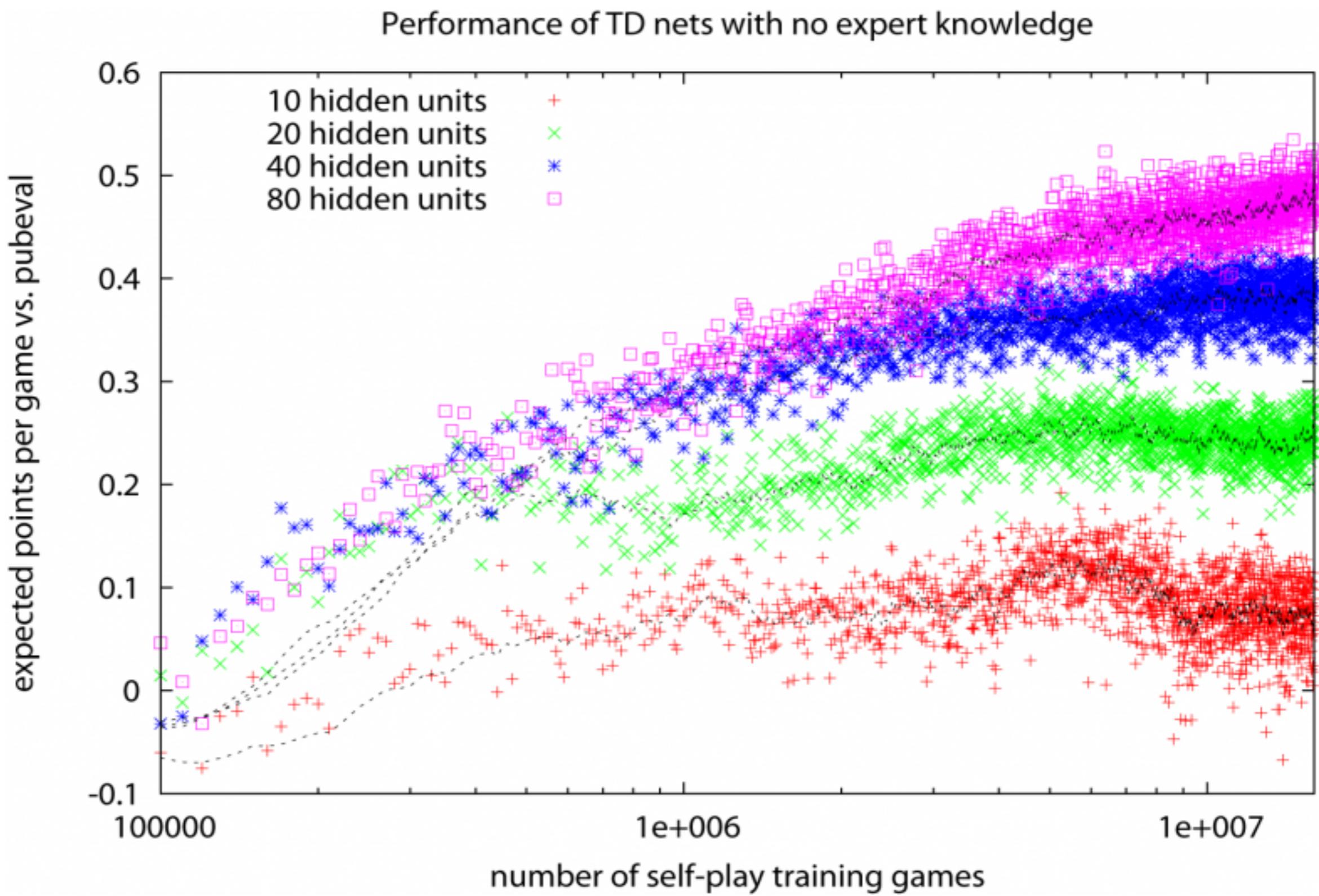
# APPROXIMATION OF THE VALUE-FUNCTION WITH A NEURAL NETWORK

---



# TD-GAMMON - PERFORMANCE

---

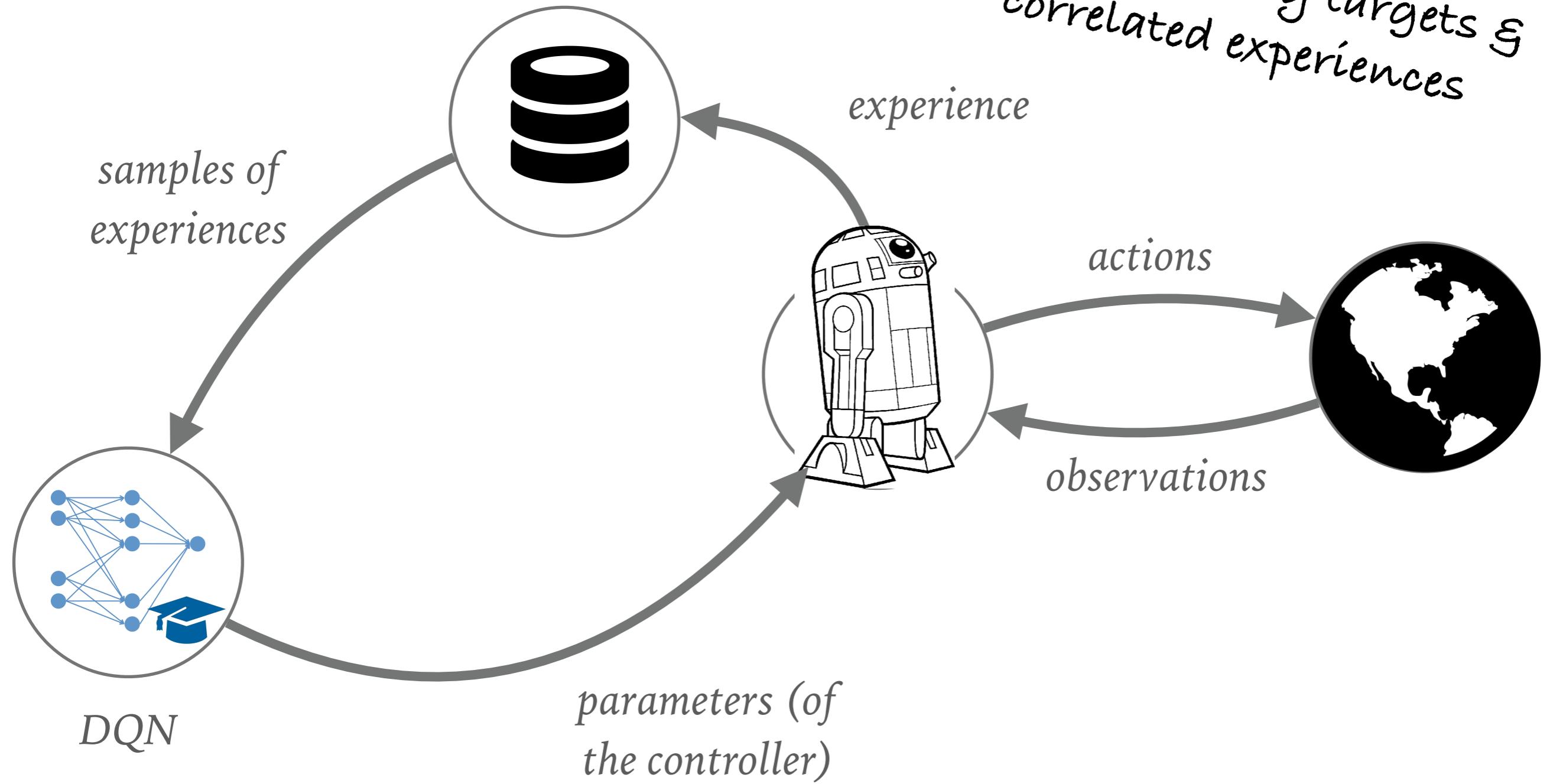


# TD-GAMMON GREATLY INFLUENCED HOW HUMAN PLAYERS PLAY TODAY

version	hidden neurons	trained games	human opponent	result
1	80	300000	Robertie	-13 after 51 games
2	40	800000	Variousmasters	-7 after 38 games
2.1	80	1500000	Robertie	-1 after 40 games
3	80	1500000	Kazaros	6 after 20 games

similar to the  
achievements of AlphaGo  
Zero today

# (NEURAL) FITTED Q-ITERATION



# FITTED Q-ITERATION

random sampling  
experiences in order to  
avoid correlations

- Approximate  $Q(s, a)$  instead of discretisation
- Update is done with a set of historical data (interactions)  
also known as the replay memory  $H = \{\langle s_i, a_i, r_i, s'_i \rangle\}$
- Apply a well known regression method (e.g. neural network)
- Use bootstrapping to incorporate already well known information

we keep the network fixed to enforce  
some stationarity

$$\mathcal{L}_i(\theta_i) = \mathbb{E}_{s,a,r,s' \sim H} \left( r + \gamma \max_{a'} Q(s', a', \theta_i^-) - Q(s, a; \theta_i) \right)^2$$

# Learning to Swing-Up and Balance from Scratch

"Neuroinformatics Group"  
University of Osnabrück  
Prof. Dr. Martin Riedmiller

**THAT IS SO NINETIES,  
NOW EVERYTHING IS  
DEEP . . .**

# DEEP Q-LEARNING: THE ARCADE GAME CASE

---

- Approximate the Q-function with a deep neural network
- State = image
- Output = Q-value for each possible action

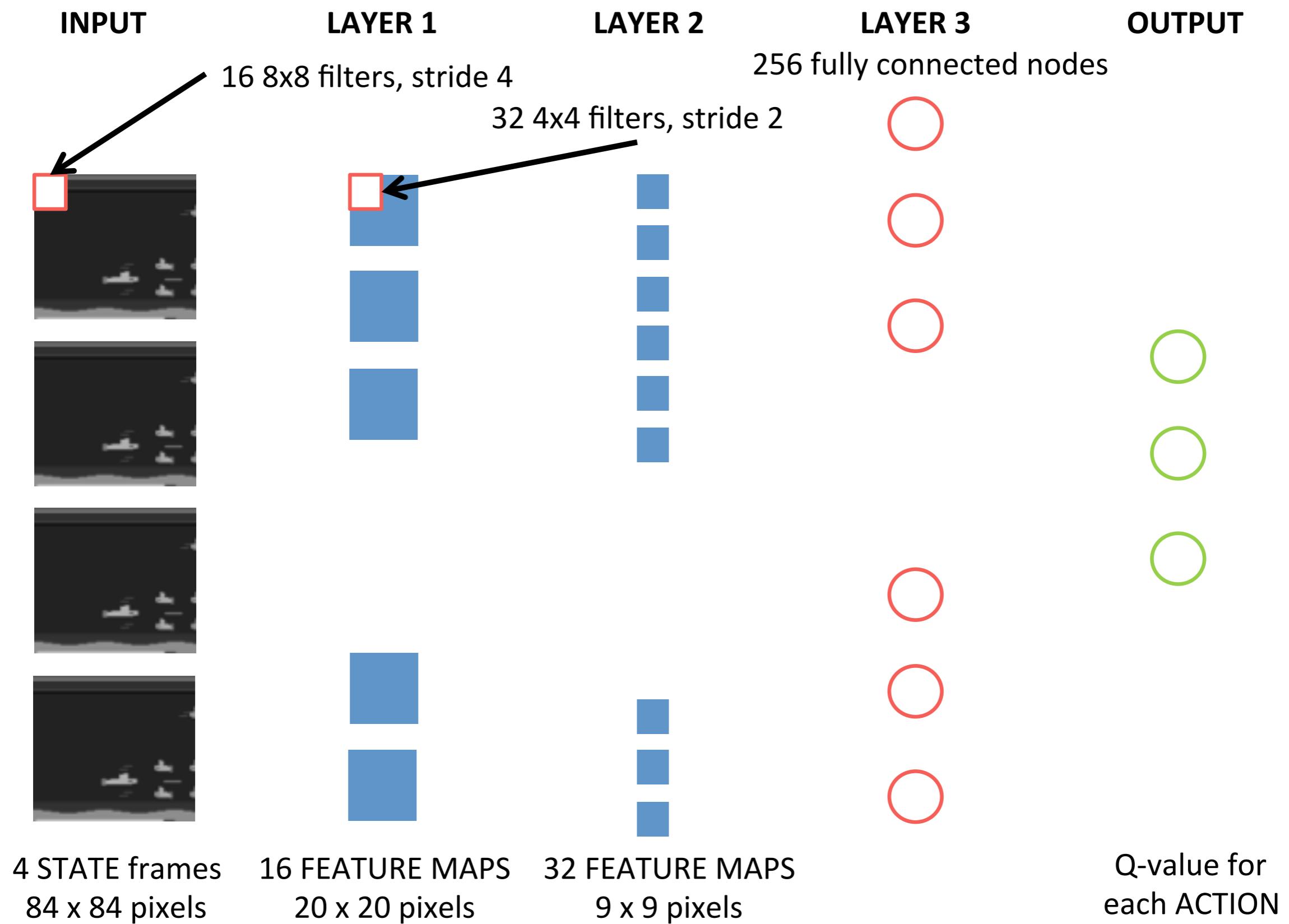


*210 x 160 pixels - 128 colours*

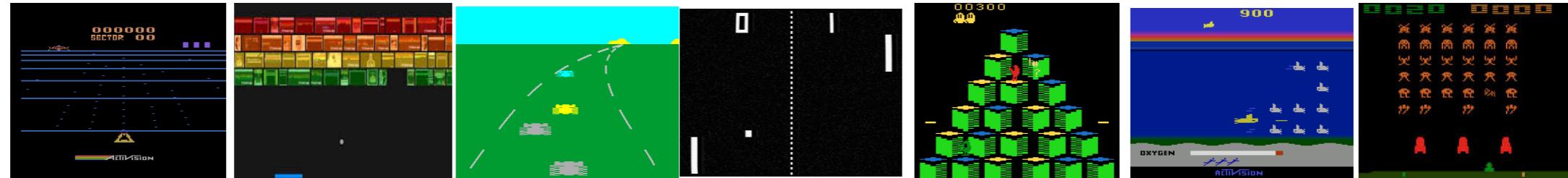


*84 x 84 pixels - 4 colours*

# DEEP Q-LEARNING: THE ARCADE GAME CASE



# DEEP Q-LEARNING: THE ARCADE GAME CASE

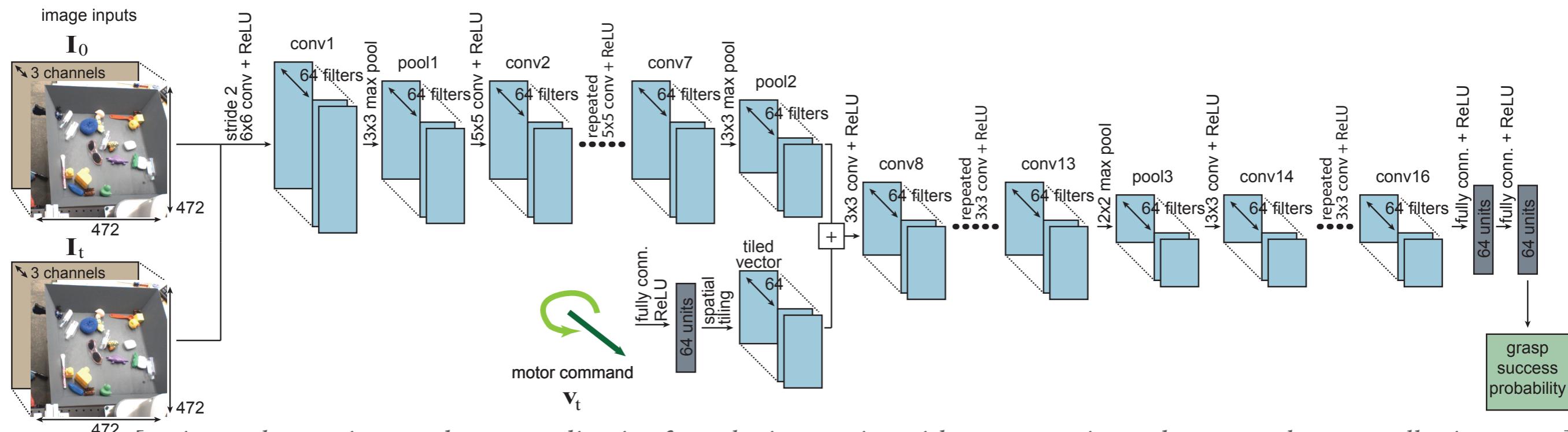


	Score	Frames	Deaths	Score	Frames	Deaths	Score
Random	354	1.2	0	-20.4	157	110	179
SARSA	996	5.2	129	-19	614	665	271
Contingency	1743	6	159	-17	960	723	268
DQN	4092	<b>168</b>	<b>470</b>	<b>20</b>	1952	1705	581
HUMAN	<b>7456</b>	31	368	-3	<b>18900</b>	<b>28010</b>	<b>3690</b>

[Mnih et al, Playing Atari with Deep Reinforcement Learning, 2014]



# THE GRASPING ROBOTS CASE: 14 ROBOTS, 2 MONTHS, > 800.000 ATTEMPTS



[Levine et al, Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection, 2016]

# HOE LEREN ROBOTS VAN MENSEN?

## PART III

*What if you don't know  
the reward function?*



[picture by Wouter Van Vooren]

**LET'S EVALUATE  
FOLDING SKILLS**



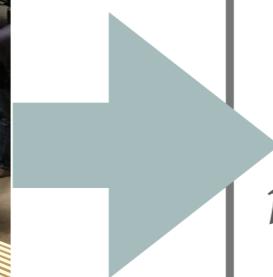
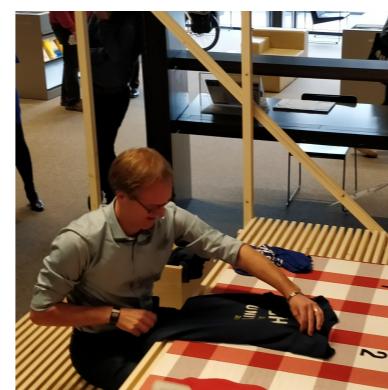
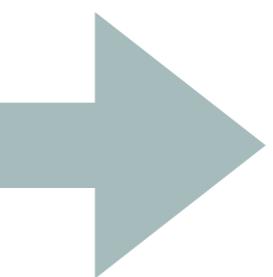
[<https://www.youtube.com/watch?v=VAhTgKT86W4>]



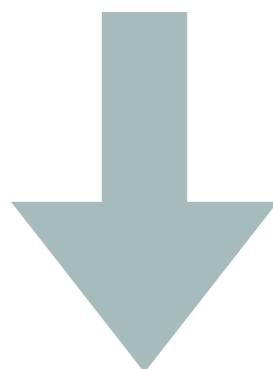
[<https://www.youtube.com/watch?v=VAhTgKT86W4>]



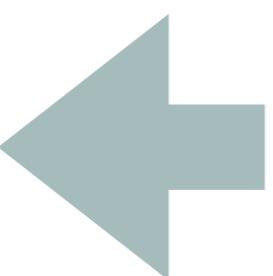
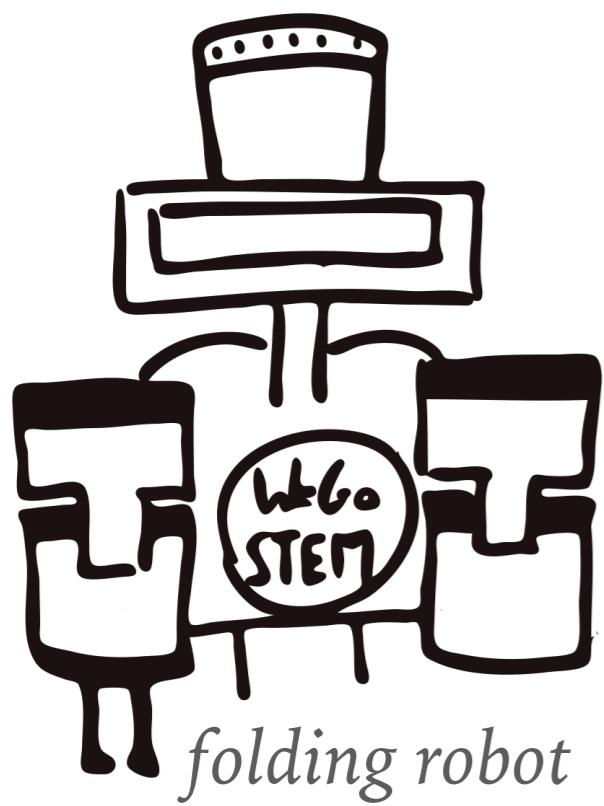
*data collection*



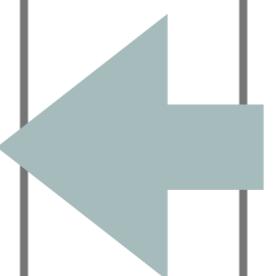
*segmentation of the most important steps*



*clustering*



*reinforcement learning*



*inverse reinforcement learning*

# INVERSE REINFORCEMENT LEARNING

---

- We try to find the reward function  $R$  that explains the expert behaviour
- Assume that the expert is indeed optimal
- Assume we can enumerate all policies -> computational challenge!
- Ongoing research topic

**AVAILABLE FOR  
EVERYONE?**

# DEEP Q-LEARNING: PATENTED BY GOOGLE

---

## Methods and apparatus for reinforcement learning

US 20150100530 A1

### ABSTRACT

We describe a method of reinforcement learning for a subject system having multiple states and actions to move from one state to the next. Training data is generated by operating on the system with a succession of actions and used to train a second neural network. Target values for training the second neural network are derived from a first neural network which is generated by copying weights of the second neural network at intervals.

<b>Publication number</b>	US20150100530 A1
<b>Publication type</b>	Application
<b>Application number</b>	US 14/097,862
<b>Publication date</b>	Apr 9, 2015
<b>Filing date</b>	Dec 5, 2013
<b>Priority date</b> ⓘ	Oct 8, 2013
<b>Also published as</b>	<a href="#">CN105637540A</a> , <a href="#">EP3055813A1</a> , <a href="#">WO2015054264A1</a>
<b>Inventors</b>	<a href="#">Volodymyr Mnih</a> , <a href="#">Koray Kavukcuoglu</a>
<b>Original Assignee</b>	<a href="#">Google Inc.</a>
<b>Export Citation</b>	<a href="#">BiBTeX</a> , <a href="#">EndNote</a> , <a href="#">RefMan</a>
<b>Patent Citations</b> (10), <b>Non-Patent Citations</b> (6), <b>Referenced by</b> (2), <b>Classifications</b> (6), <b>Legal Events</b> (2)	

**External Links:** [USPTO](#), [USPTO Assignment](#), [Espacenet](#)

### IMAGES (11)

“

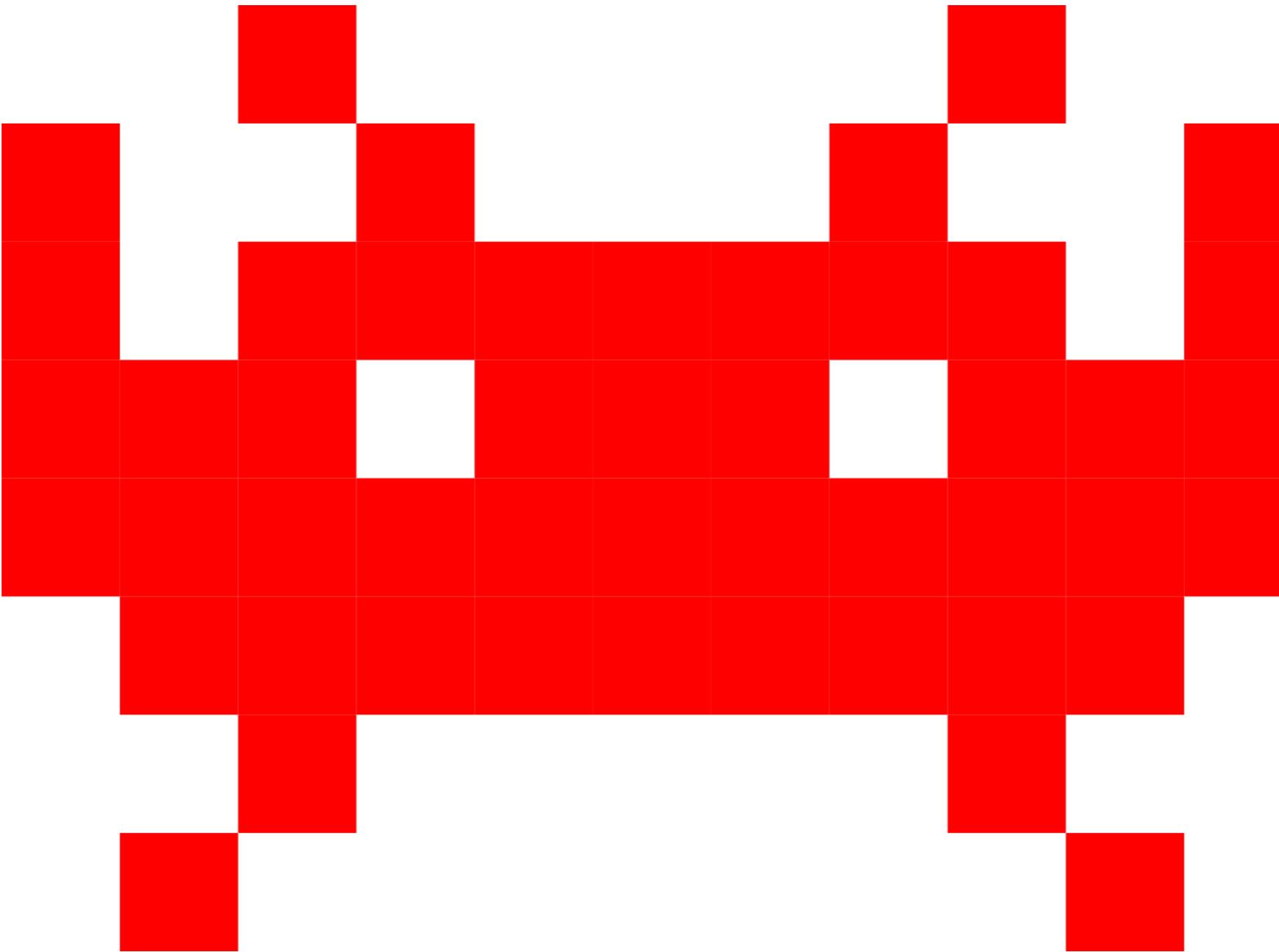
There has, nonetheless, been an attempt to use a multilayer perceptron to approximate a Q-value function, in “Neural fitted Q iteration—first experiences with a data efficient neural reinforcement learning method”, Machine Learning: ECML 2005, Springer 2005, pages 317-328, Riedmiller, Martin. The technique described there is based on the principle of storing and reusing transition experiences but has some significant practical disadvantages: broadly speaking a neural network is trained based on the stored experience, and when the experience is updated with a new (initial state—action—resulting state) triple the previous neural network is discarded and an entirely new neural network is trained on the updated experience.

*-Excerpt from US 20150100530 A1*

# TAKE HOME MESSAGES

---

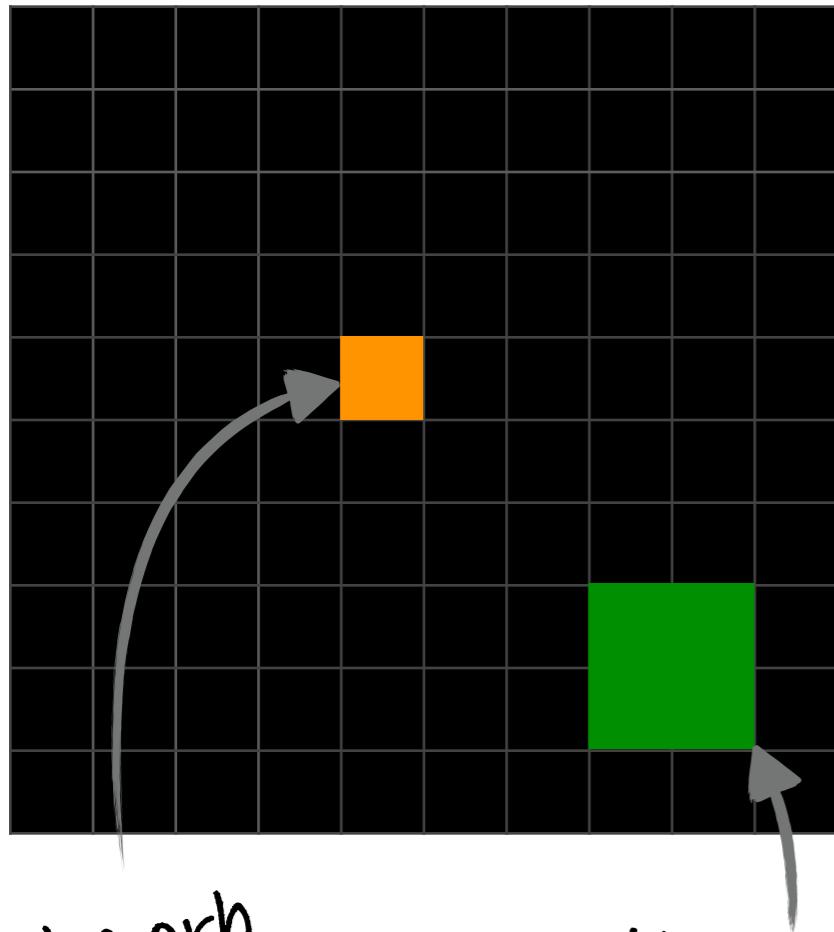
1. Use RL in problems for which you don't have explicit solutions on the condition that you can define some reward ~ learning with a critic
2. Is there a model of the world
  - ✓ YES -> PI or VI
  - ✓ NO -> TD-learning
3. Continuous state space? Try to approximate the value function (or Q function)
4. There are many interesting challenges out there!



# PROJECT ASSIGNMENT

# ORB CATCHING GAME

---



the orb

the agent

- Hands-on experience with the principles of reinforcement learning
- Starts with fixed tasks
  - ✓ reward function
  - ✓ preprocessing
  - ✓ DQN architecture
  - ✓ policy
  - ✓ debugging
- Possible questions on your exam!

# ORB CATCHING GAME: ADDING HIGH LEVEL GOALS & INCREASING SIZE

---

