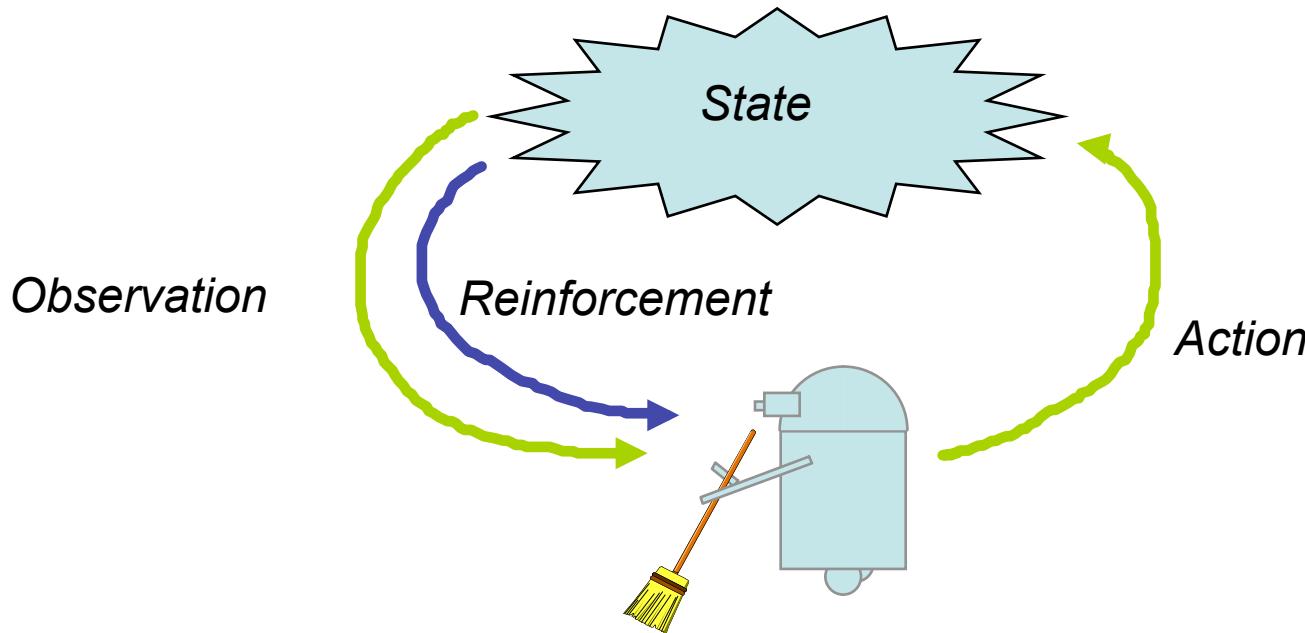
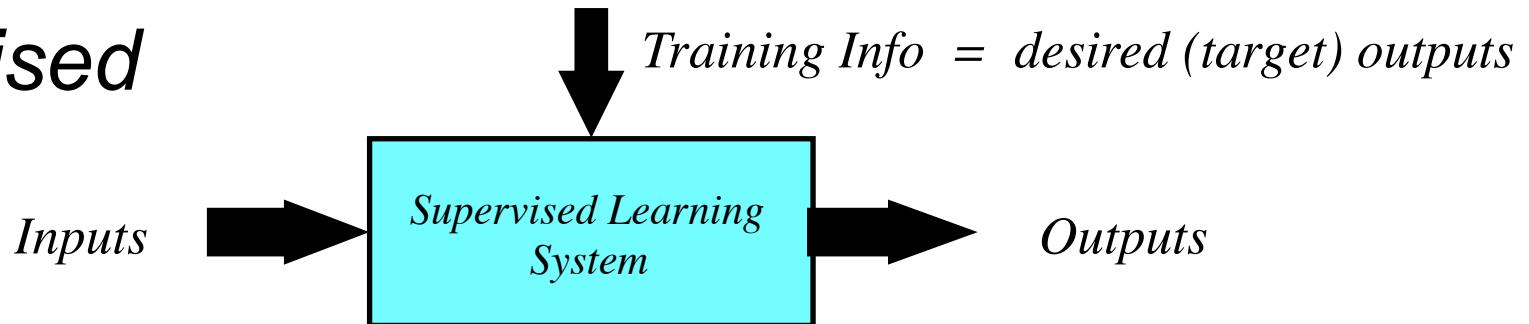


4 Reinforcement Learning (RL)

- Learning by interaction with the environment to accomplish a predefined goal
- Learn what to do – how to map situations with actions – to maximize a numerical reward value

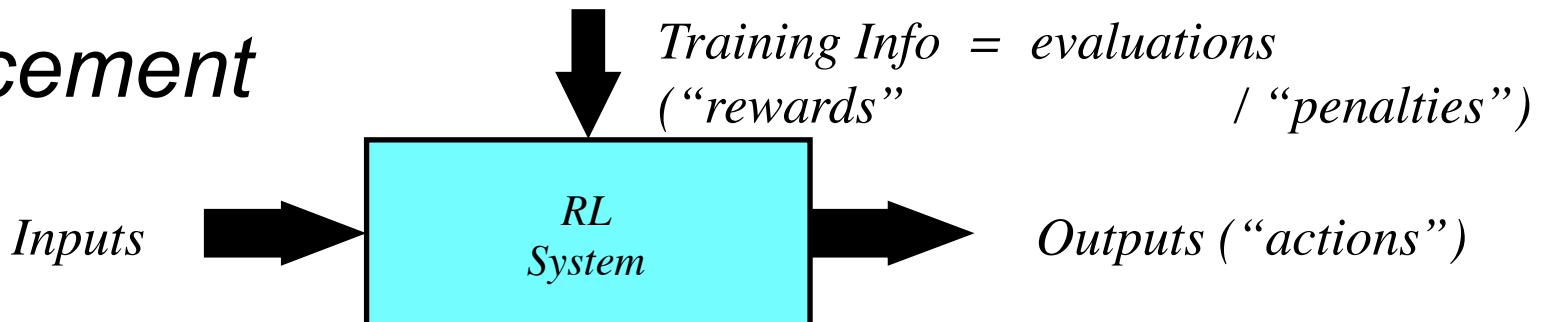


Supervised

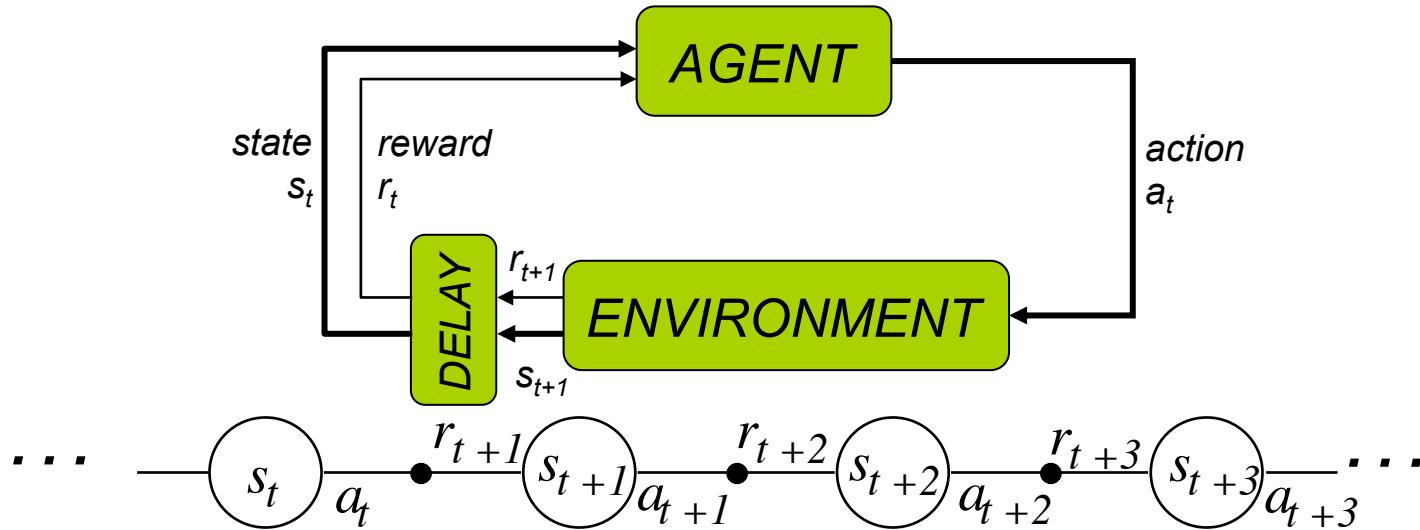


Error = (target output – actual output)

Reinforcement



Objective: get as much reward as possible



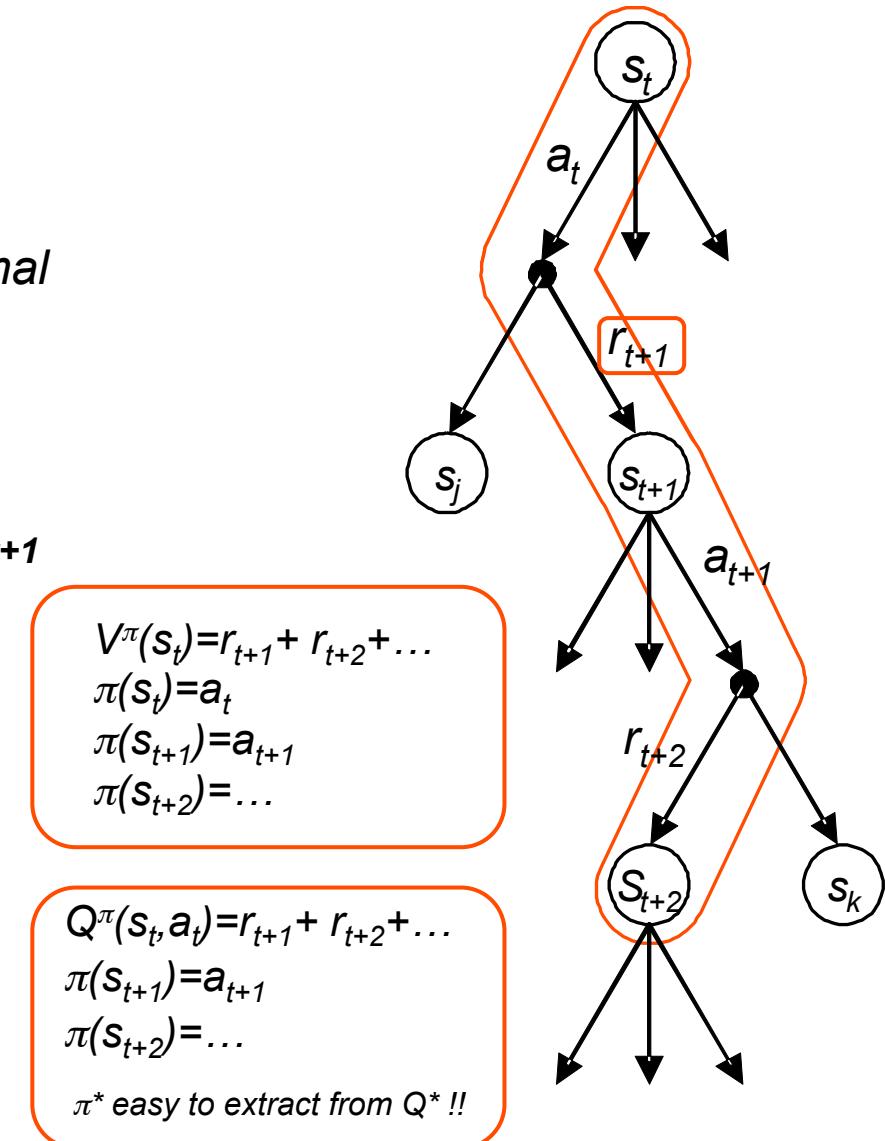
Agent:

- **Learner and decision maker.**
- **Input:** the **state** " s_t " and a **reward** " r_t " (numerical value)
- **Output:** an **action** " a_t "
- **Goal:** **to maximize the amount of future rewards** $\sum r_t$

Environment:

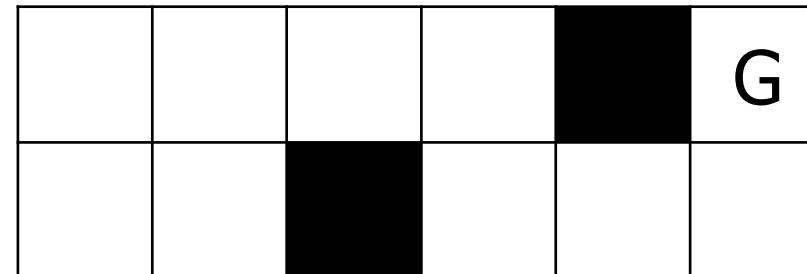
- **Everything outside the agent.**
- **It reacts to actions with a new state.**
- **A Reinforcement Function provides the new reward.**

- **Policy Function**, $\pi(s, a): S \rightarrow A$
 - Learner function*
 - State/action mapping*
 - Kinds: random, greedy, ε -greedy, optimal*
- **Dynamics Function**, $\{s_t, a_t\} \rightarrow s_{t+1}$
 - Environment function*
 - Unknown function*
- **Reinforcement Function**, $r(s_{t+1}, a_t) = r_{t+1}$
 - Environment function*
 - Defines the goal of the learner*
- **Value Function**, $V^\pi(s)$; $Q^\pi(s, a)$
 - Learner function*
 - Expected sum of*
 - Kinds:- State-Value*
 - Action-Value function $Q^\pi(s, a)$*



Example

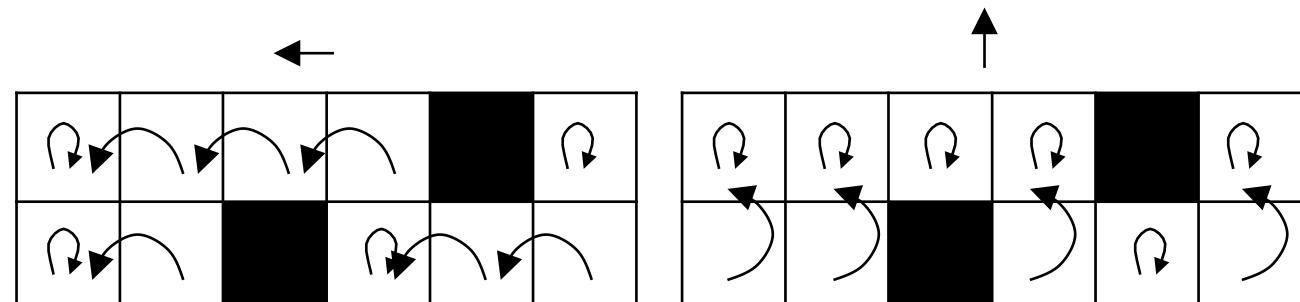
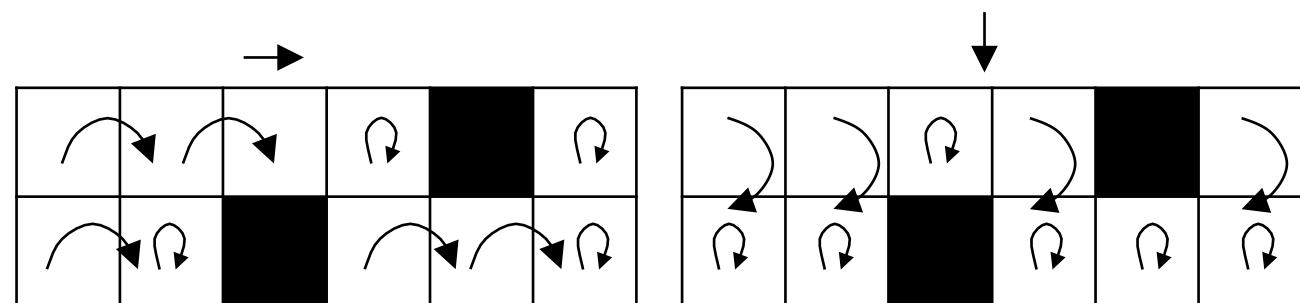
Environment:



Actions:



Dynamics:

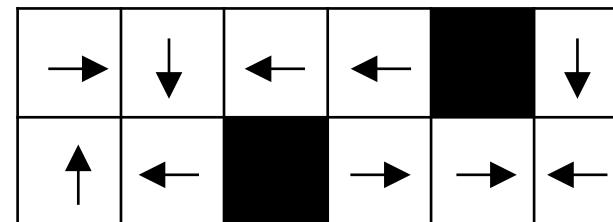


Example

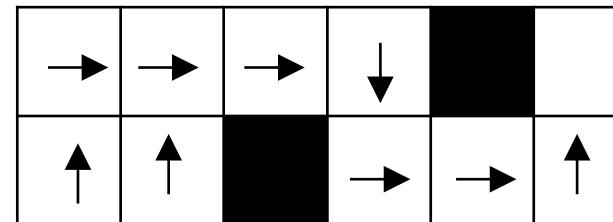
The reinforcement function:

-1	-1	-1	-1		1
-1	-1		-1	-1	-1

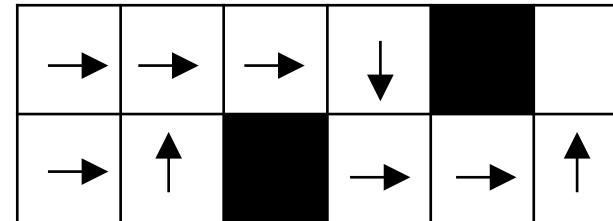
A policy:



*An optimal policy, π^*_1 :*

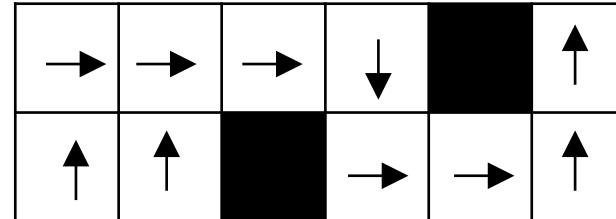


*Another optimal policy, π^*_2 :*



Example

Policy, π_1^ :*



The state value function, $V_1^(s)$:*

-5	-4	-3	-2		
-6	-5		-1	0	1

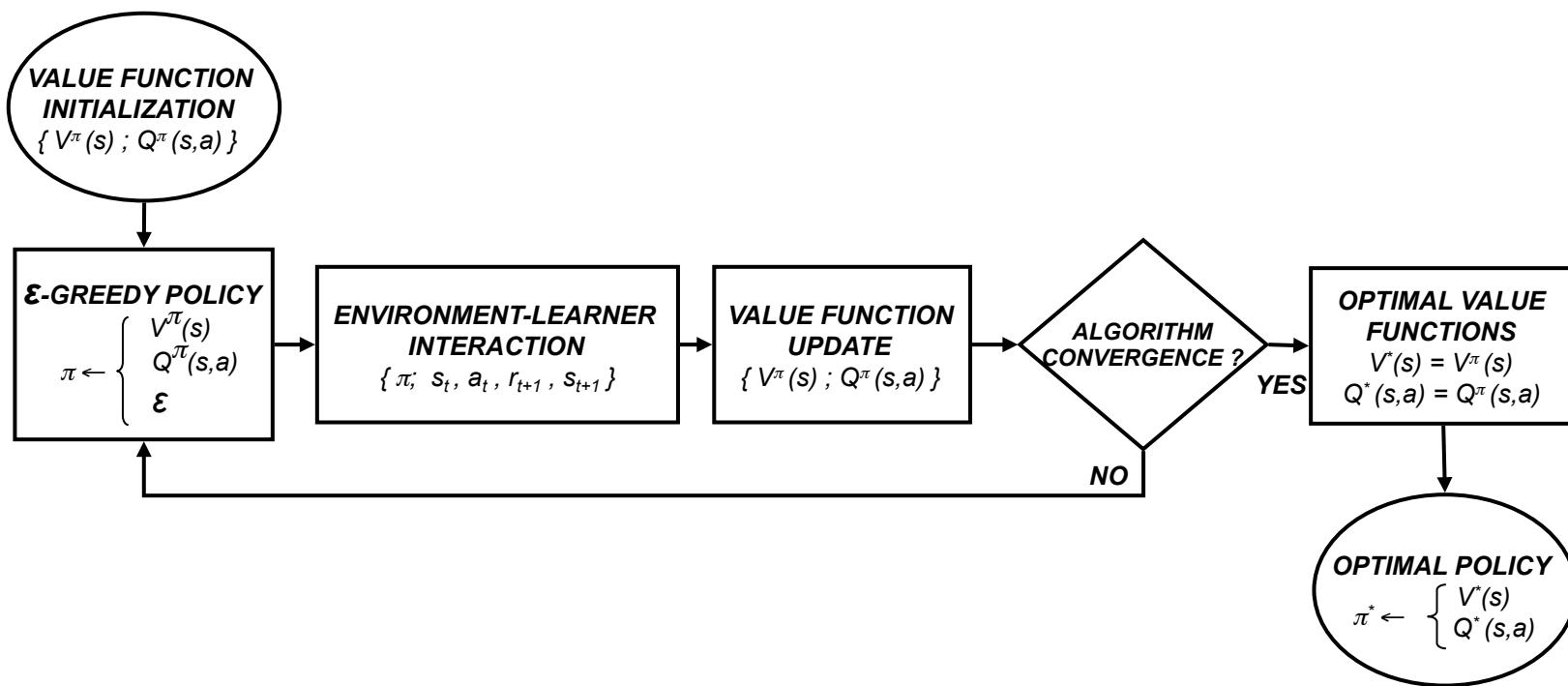
The action value function, $Q_1^(s,a)$:*

-5	-4	-3	-3		
-6	-6		-1	0	0

-7	-6	-4	-2		
-7	-6		-2	-1	0

-6	-6	-5	-4		
-7	-7		-2	-2	-1

-6	-5	-4	-3		
-6	-5		-3	-1	1

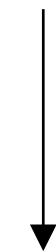




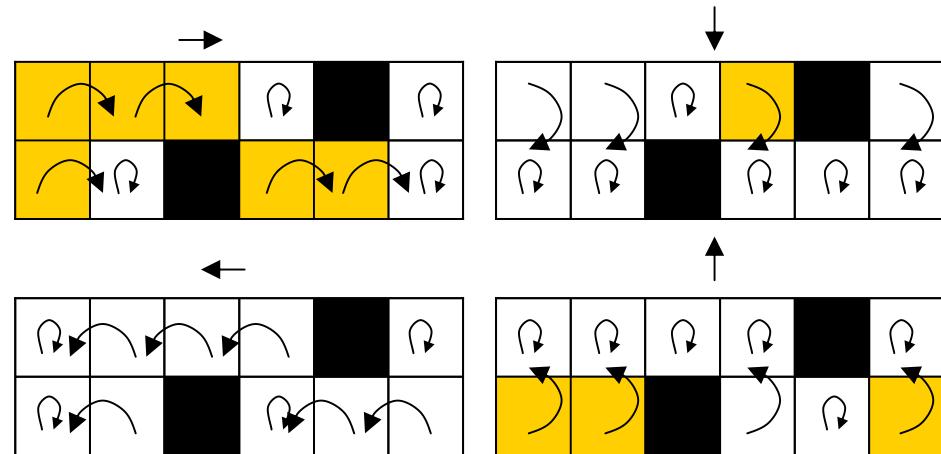
$V^*(s)$:

-5	-4	-3	-2		
-6	-5		-1	0	1

+



Dynamics:



π^* :

→	→	→	↓		
↑→	↑		→	→	↑



$Q^*(s, a)$:

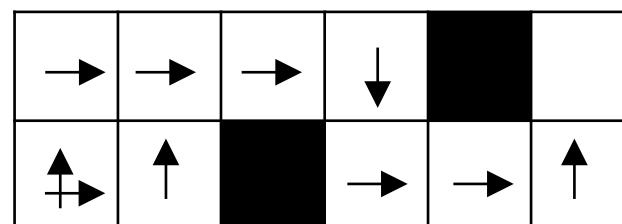
-5	-4	-3	-3		
-6	-6		-1	0	0

-7	-6	-4	-2		
-7	-6		-2	-1	0

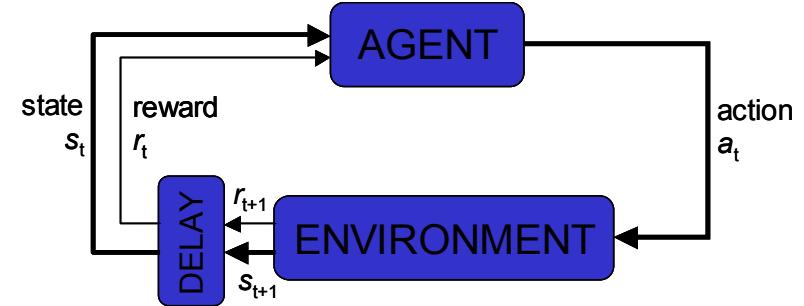
-6	-6	-5	-4		
-7	-7		-2	-2	-1

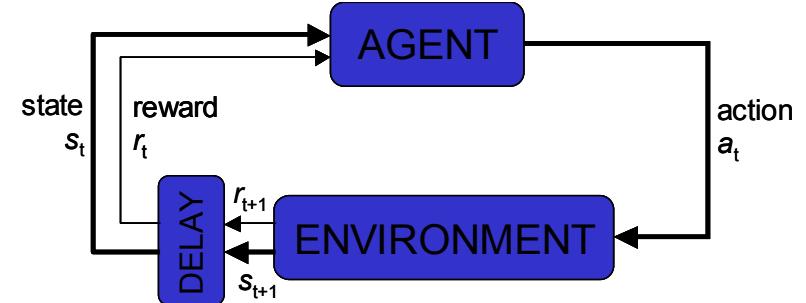
-6	-5	-4	-3		
-6	-5		-3	-1	1

$\pi^*: s_i \rightarrow a_{max} \mid \text{Max } Q^*(s_i, a_{max})$:



- Learning by rewards penalization, without examples
- Trial error search
- Dilemma between exploration / exploitation
 - exploration: to find the optimal actions
 - exploitation: to receive a higher reward
- Delayed rewards
 - ignore immediate reward to obtain better future rewards





- Policy:

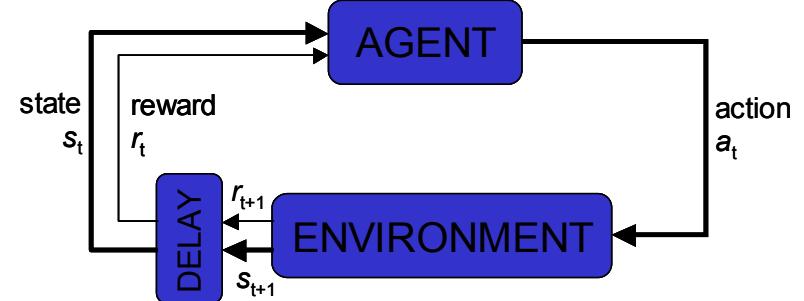
- mapping between states s_t and actions a_t

$$\pi(s, a) : S \rightarrow A$$

- *The goal is to find the **optimal policy** π^**

- Reinforcement function

- Mapping between each state/action pair with a reward
 - Defines the goal to be fulfilled
 - An external knowledge is necessary



- Value functions

- They contain the expected sum of rewards when following a particular policy.
- Kinds:
 - State Value Function: sum of rewards that can be expected when being in a particular state
 $V^\pi(s)$
 - Action Value Function: sum of rewards that can be expected when being in a particular state and a particular action is applied:
 $Q^\pi(s,a)$

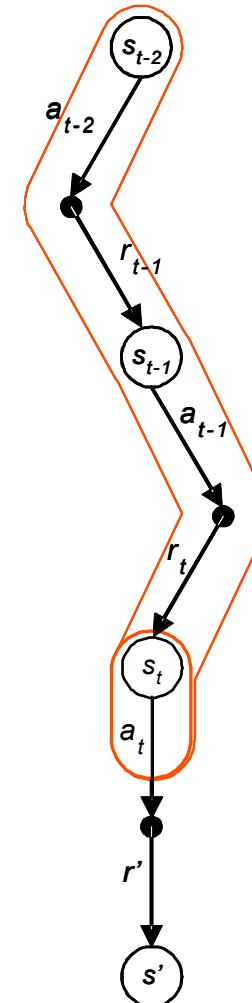
- FMDP: Finite Markov Decision Process
- Mathematical formulation of the RLP
- Constraints
 - Markov Property:

$$\Pr\left\{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t, r_t, s_{t-1}, a_{t-1}, K, r_1, s_0, a_0\right\} =$$

$$\Pr\left\{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t\right\}$$

→ Complete state

- Finite actions and states





- Discrete set of states and actions
- State Transition Probability:

$$P_{ss'}^a = \Pr\{s_{t+1} = s' \mid s_t = s, a_t = a\} \quad \text{for all } s, s' \in S, a \in A(s).$$

- Reinforcements Probability:

$$R_{ss'}^a = E\{r_{t+1} \mid s_t = s, a_t = a, s_{t+1} = s'\} \quad \text{for all } s, s' \in S, a \in A(s).$$

$P_{ss'}^a$, and $R_{ss'}^a$, specify the dynamics of the FMDP

Using the Bellman equation:

- Optimal State Value function

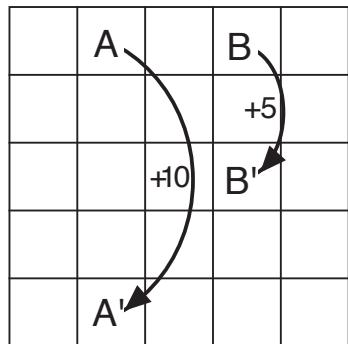
$$V^*(s) = \max_{a \in A(s)} \sum_s P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')]$$

- Optimal Action Value Function

$$Q^*(s, a) = \sum_s P_{ss'}^a [R_{ss'}^a + \gamma \max_{a'} Q^*(s', a')]$$

Obtaining an optimal policy

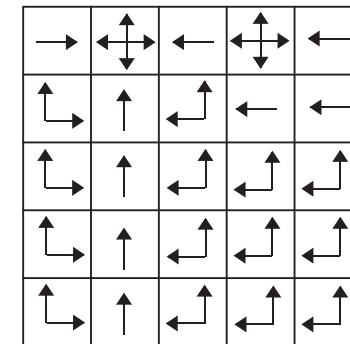
- From V^* : A policy $\pi(s,a)$, for that $\forall s \in S$, the new state s' accomplish that $V^*(s') > V^*(s)$ is an optimal policy



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

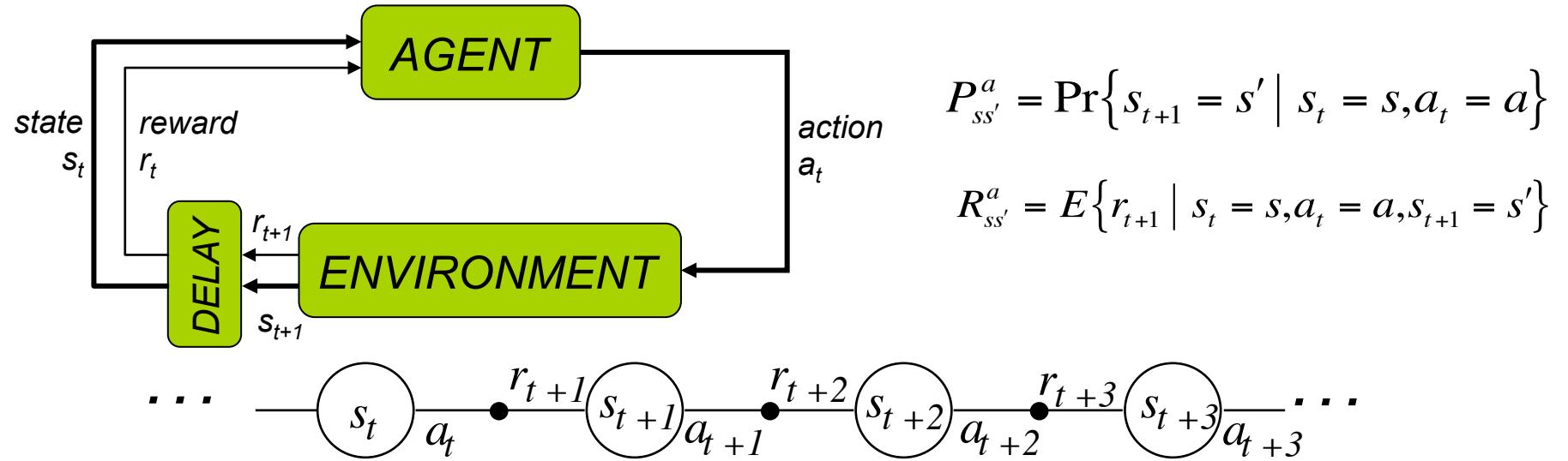
b) V^*



c) π^*

- From Q^* : Given the function $Q^*(s,a)$

$$\pi^*(s) = \arg \max_{a \in A(s)} Q^*(s,a)$$



- Dynamic Programming
- Monte-Carlo Methods
- Temporal Differences

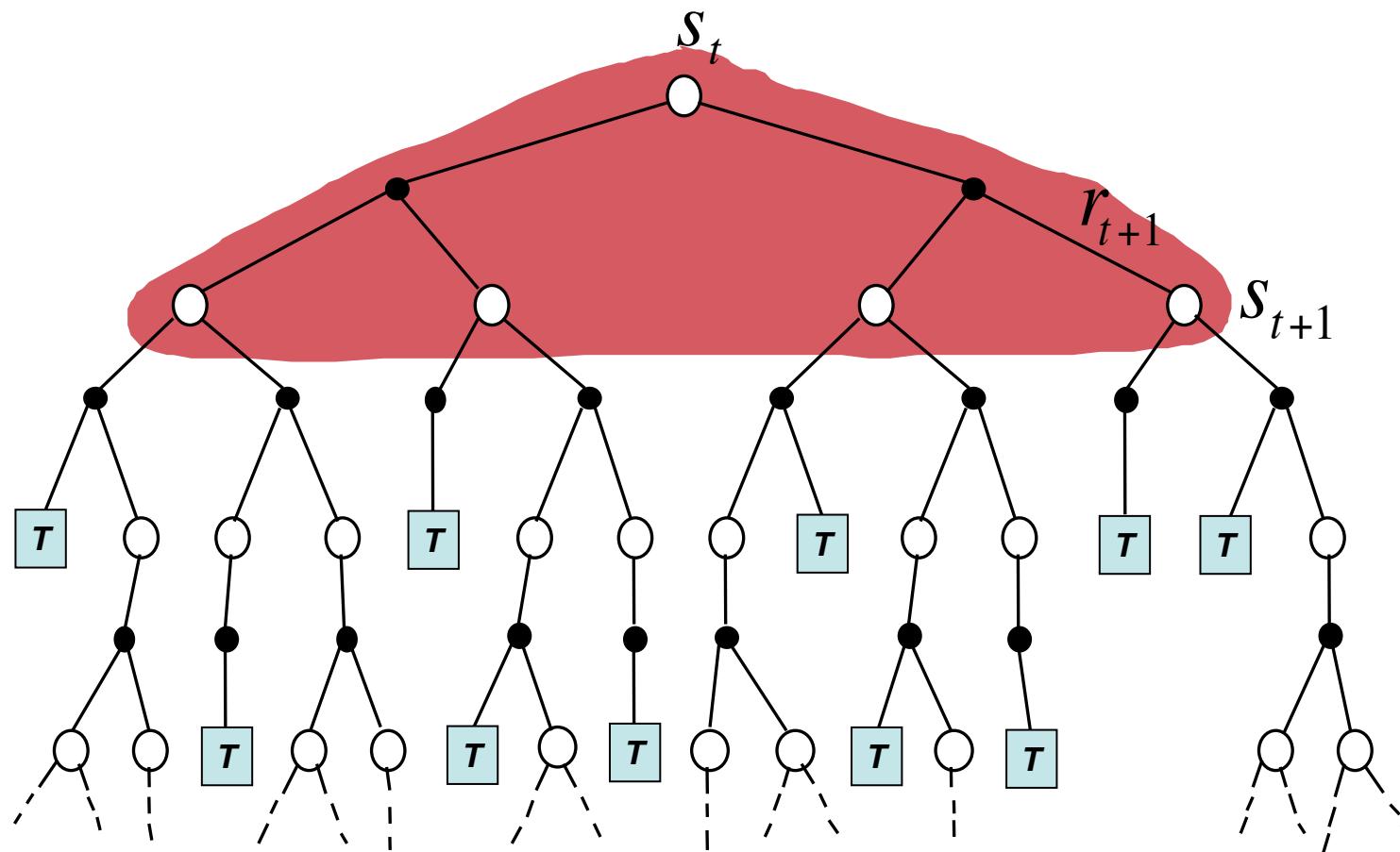


- Algorithms that find the optimal policies from the model of the environment $P_{ss'}^a$, i $R_{ss'}^a$,

$$V^*(s) = \max_{a \in A(s)} \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')]$$

$$Q^*(s, a) = \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma \max_{a'} Q^*(s', a')]$$

$$V(s_t) \leftarrow E_{\pi} \{ r_{t+1} + \gamma V(s_{t+1}) \}$$



- **Do not** use the model of the environment to update the value functions.
- No estimation of the future accumulated reward.
- Value functions are updated at the end of the episode.

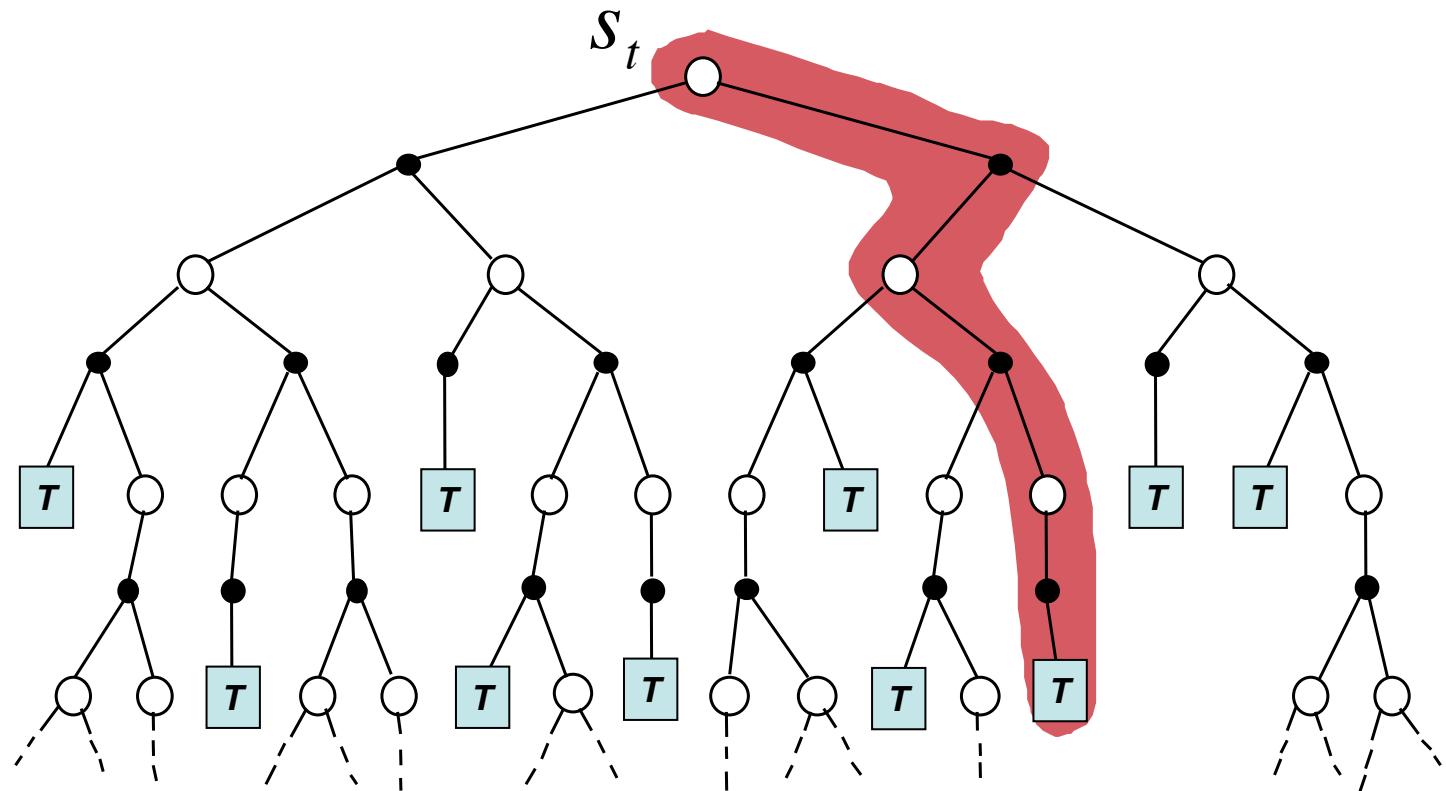
Simple every - visit Monte Carlo method :

$$V(s_t) \leftarrow V(s_t) + \alpha [R_t - V(s_t)]$$


target: the actual return after time t

$$V(s_t) \leftarrow V(s_t) + \alpha [R_t - V(s_t)]$$

where R_t is the actual return following state s_t .



- **Do not** use the model of the environment to update the value functions.
- An estimation of the future accumulated reward is done to update value functions (**online** learning).

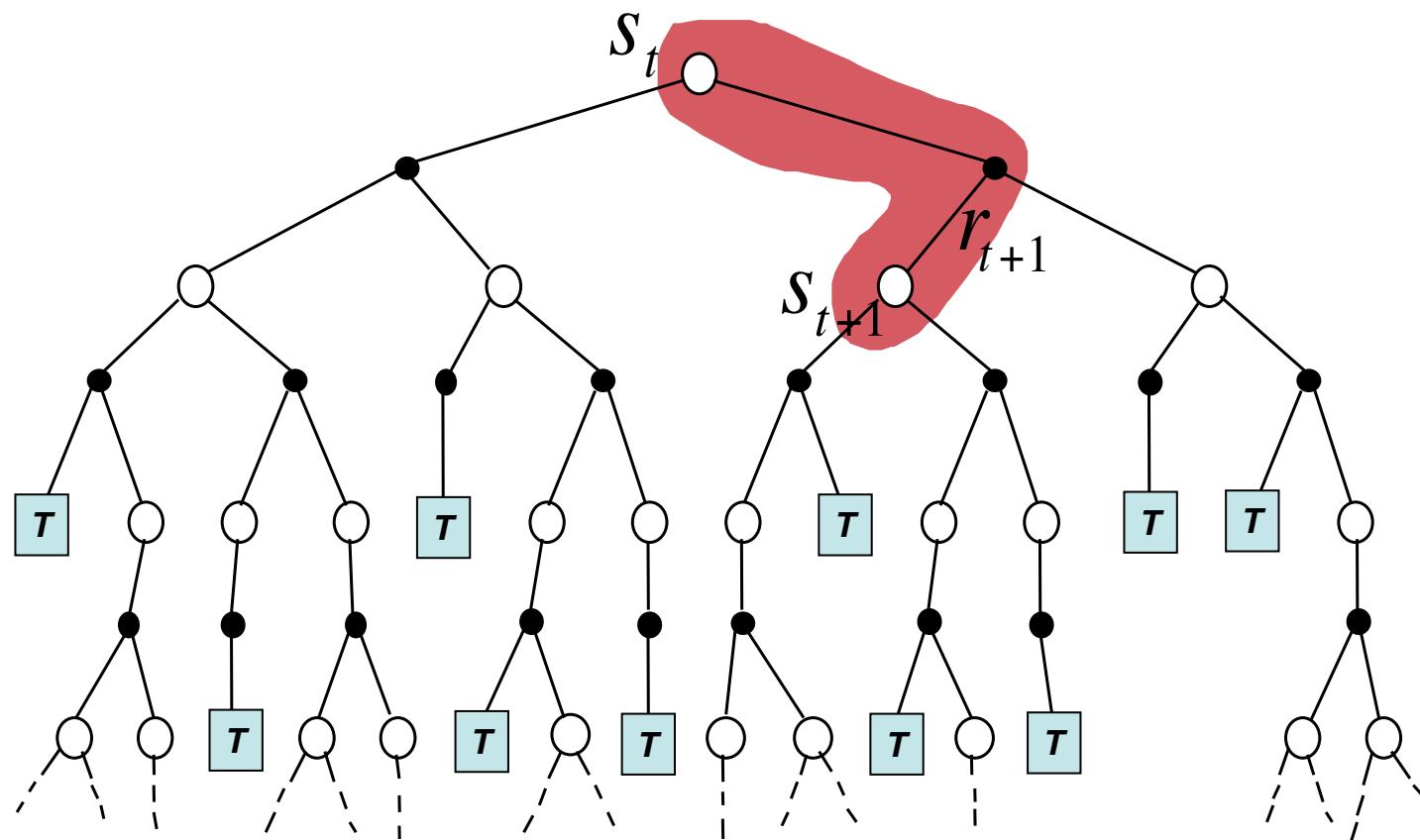
The simplest TD method, TD(0) :

$$V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

target: an estimate of the return

Temporal differences

$$V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$



Methodologies comparison

	ONLINE LEARNING	DYNAMICS FUNCTION
Dynamic Programming	YES	YES
Monte Carlo Methods	NO	NO
Temporal Difference	YES	NO

suitable for robotics

- *Actor Critic methods*
- *Q_learning*
- *SARSA*



- TD method
- Action-Value Function: $Q(s,a)$
- Off-policy learning : any policy can be followed
- Convergence guaranteed as far as all the (s,a) pairs are regularly visited



Initialize $Q(s,a)$ randomly

Initialize s

Repeat

Choose a following ϵ -greedy

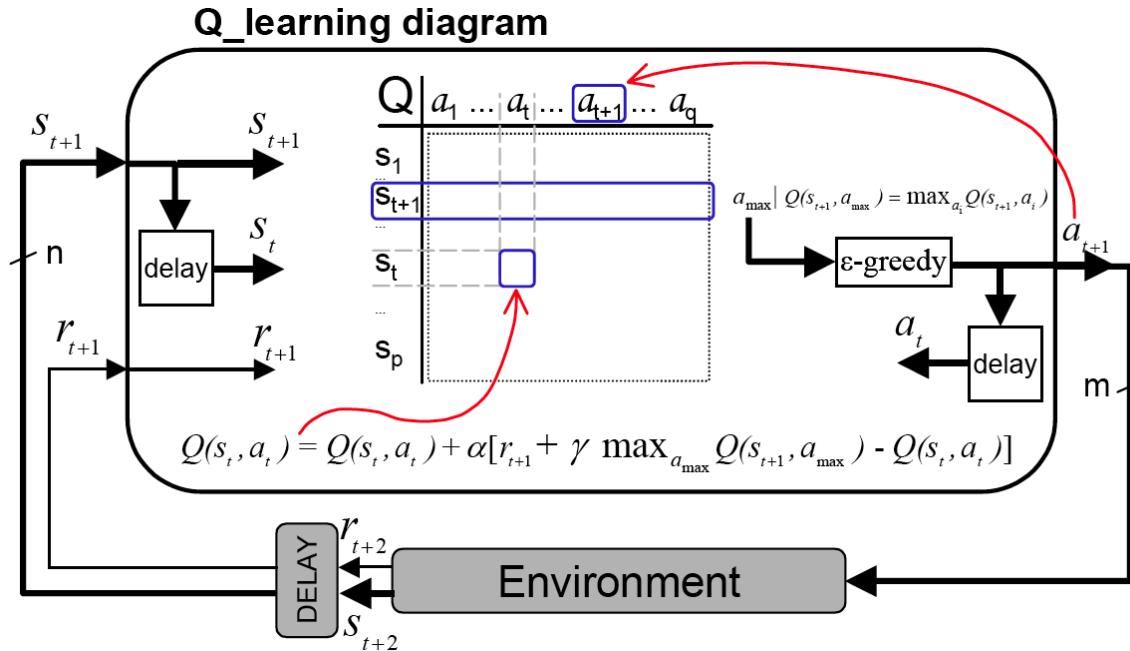
Take action a , observe r, s'

$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \max Q_{a'}(s',a') - Q(s,a)]$

$s \leftarrow s'$

until $Q(s,a)$ is optimal

Q_learning algorithm



$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q_{a'}(s', a') - Q(s, a)]$$

parameters

α : learning rate

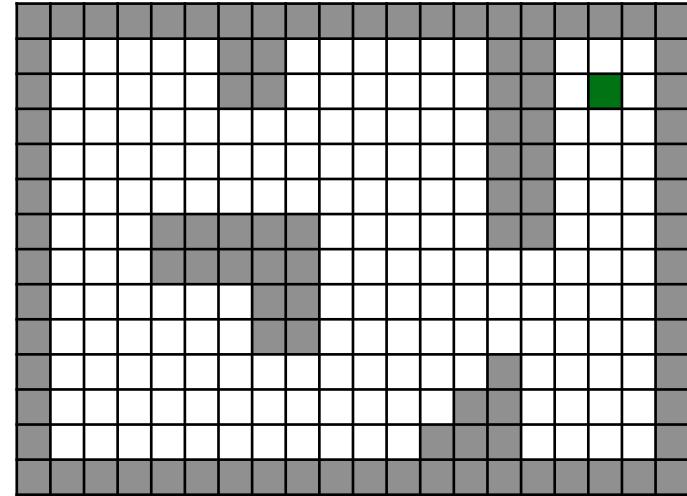
γ : discount factor

ϵ : exploration

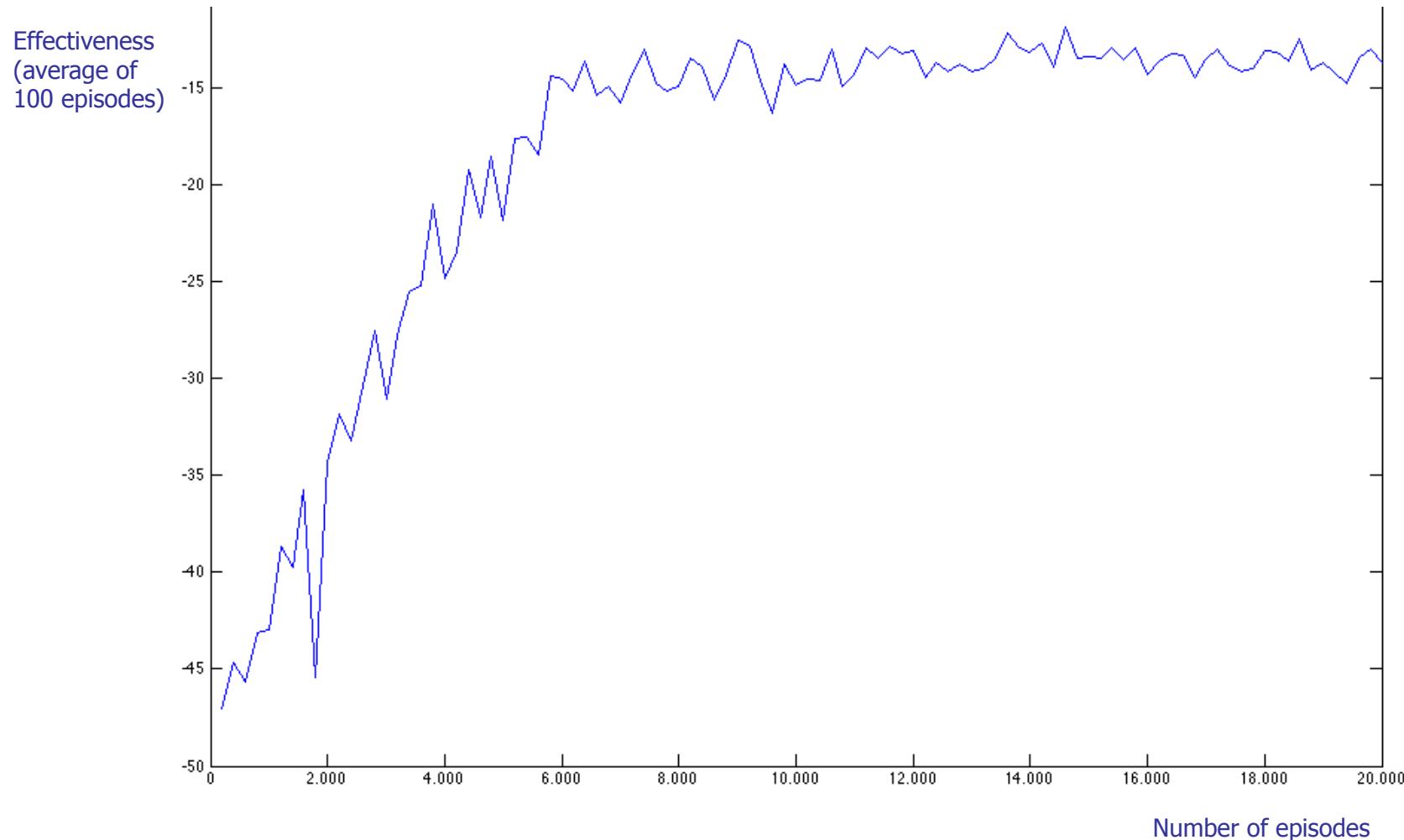
target: an estimate of the return

Example: Path planning

- States: 2D finite environment, 20x14 cells
- Actions: 4-point connectivity
- Reward: "-1" in all cells except the goal "+1"
- Dynamics: 1 step in the action direction
unless obstacle is found
- Q function: $20 \times 14 \times 4 = 1120$ cells
- Random initial position
- Max. iterations per episode = 50
- Number of episodes = 20.000
- Effectiveness : "averaged number of iterations to reach the goal"
- Learning parameters
 - Learning rate: $\alpha = 0.1$
 - Discount factor: $\gamma = 0.9$
 - Exploration: $\epsilon = 0.3$



- Effectiveness evolution



Example: Path planning

- Q function



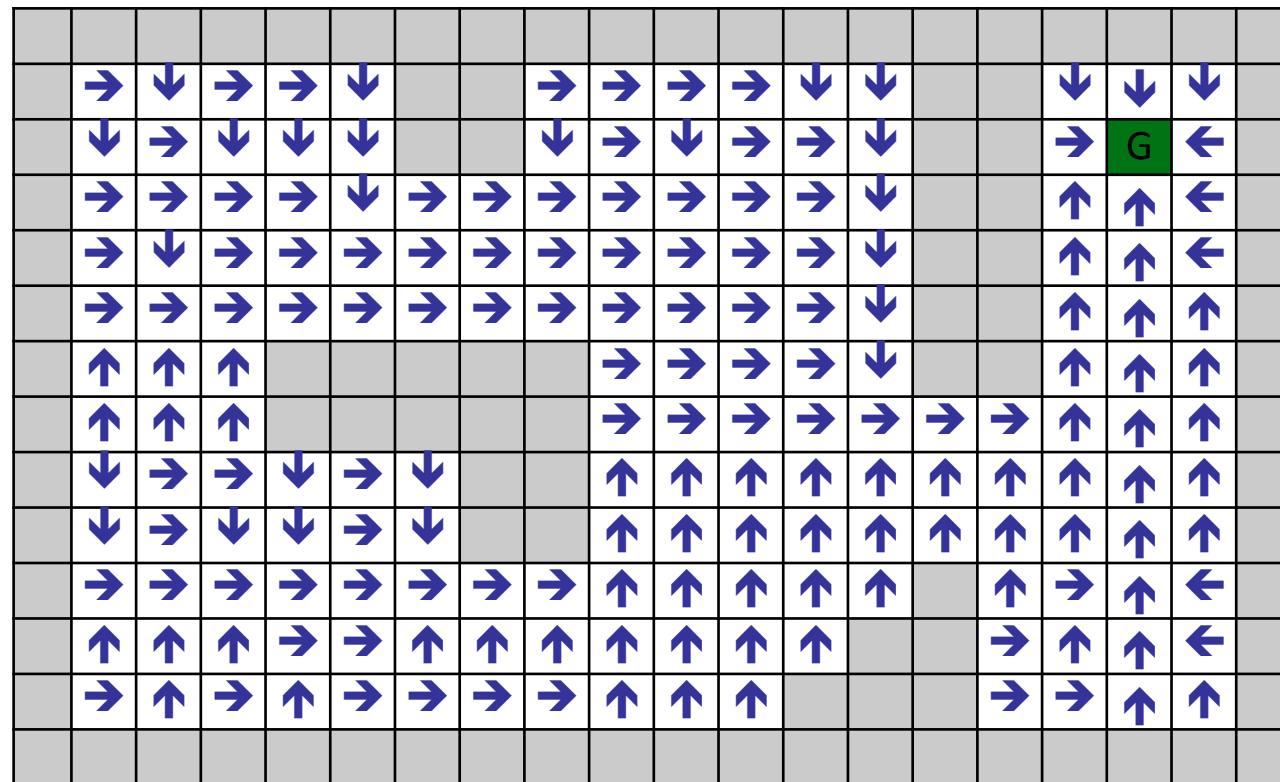
Example: Path planning

Q function



Example: Path planning

- Optimal policy extracted from Q function

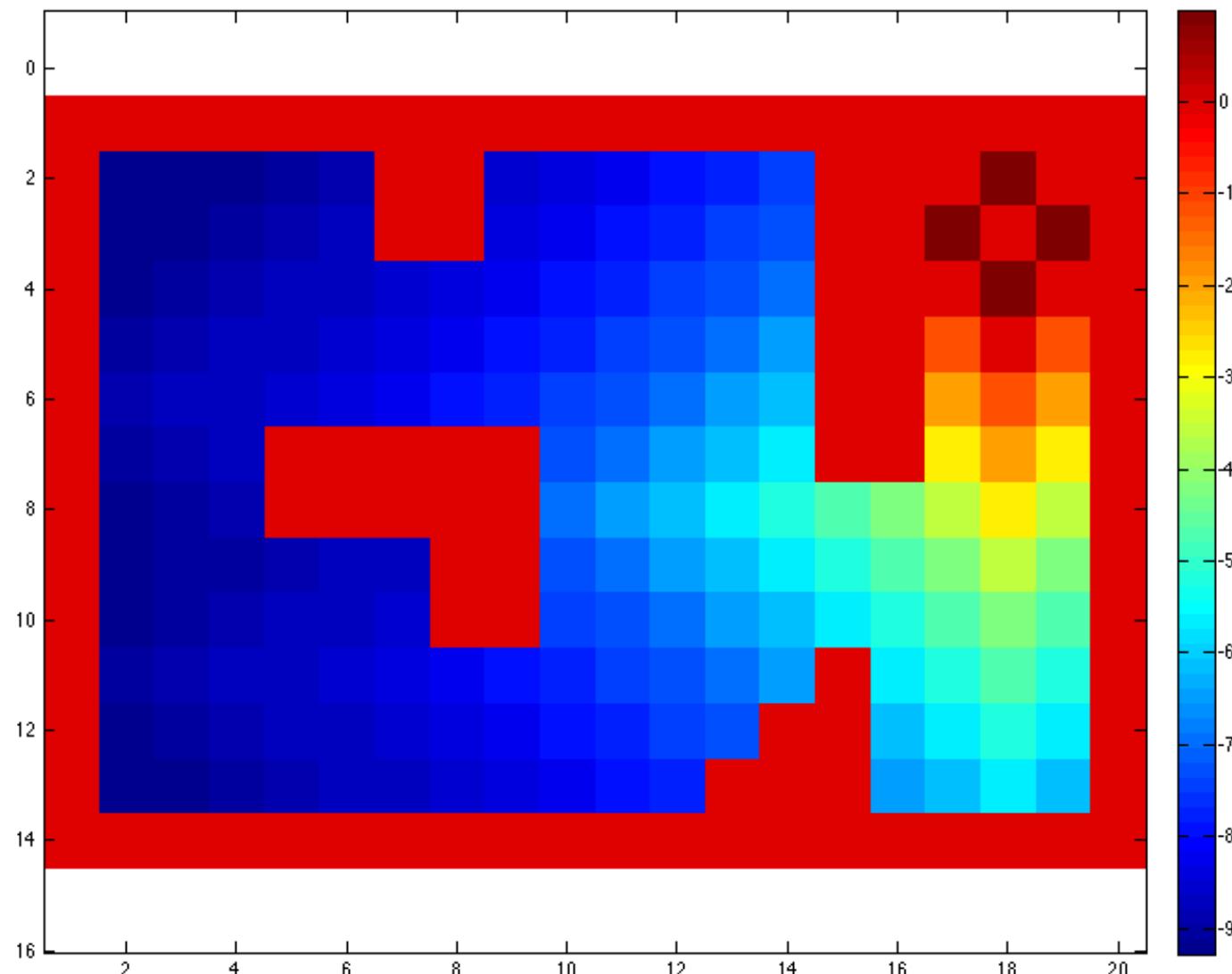


Example: Path planning

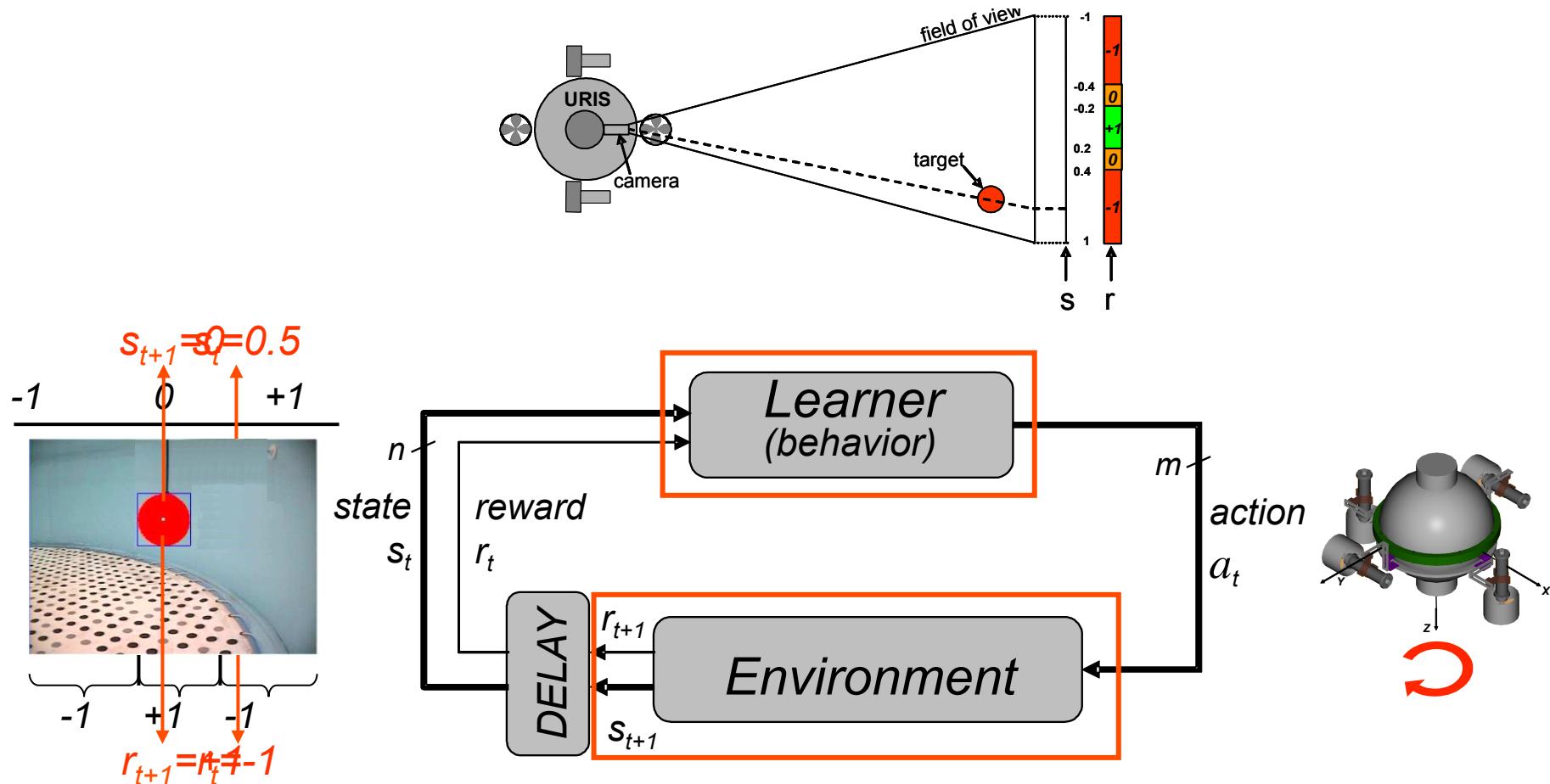
- V function extracted from Q

Example: Path planning

- V function



Case study: behaviour learning





Generalization problem



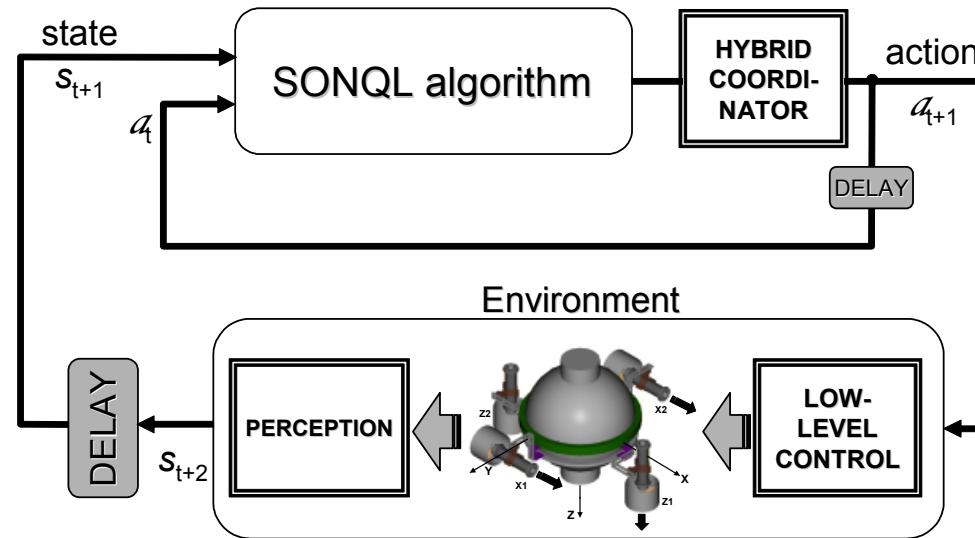
No se puede mostrar la imagen. Puede que su equipo no tenga suficiente memoria para abrir la imagen o que ésta esté dañada. Reinicie el equipo y, a continuación, abra el archivo de nuevo. Si sigue apareciendo la x roja, puede que tenga que borrar la imagen e insertarla de nuevo.



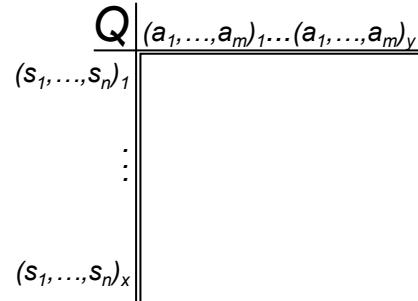
Generalization problem



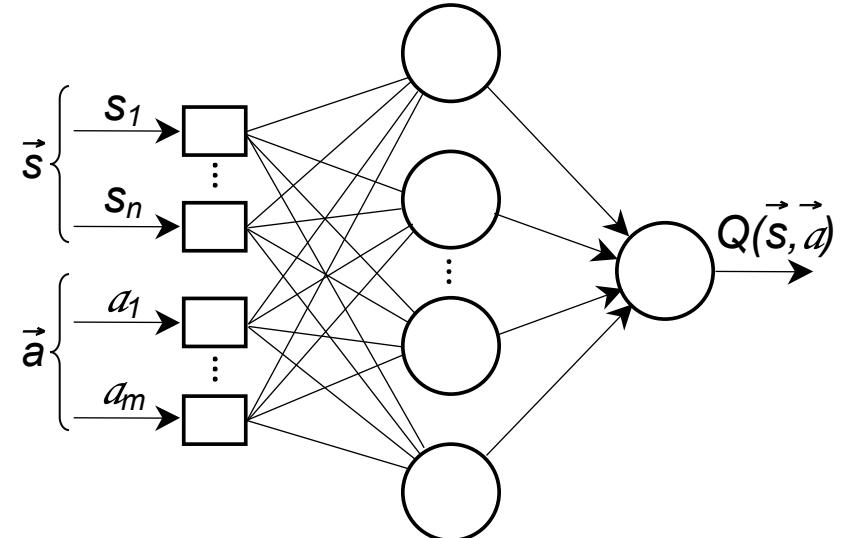
No se puede mostrar la imagen. Puede que su equipo no tenga suficiente memoria para abrir la imagen o que ésta esté dañada. Reinicie el equipo y, a continuación, abra el archivo de nuevo. Si sigue apareciendo la x roja, puede que tenga que borrar la imagen e insertarla de nuevo.



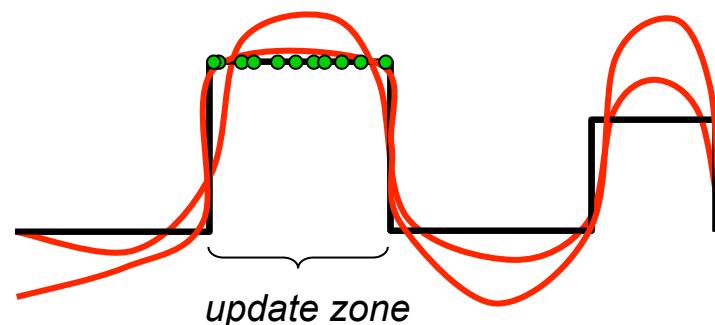
- SONQL goal: **RL algorithm for learning in real-time robotics tasks**
- Components:
 - online, off-policy \rightarrow Q_learning
 - high generalization capability \rightarrow Neural Networks
 - Interference Problem \rightarrow Learning Samples Database



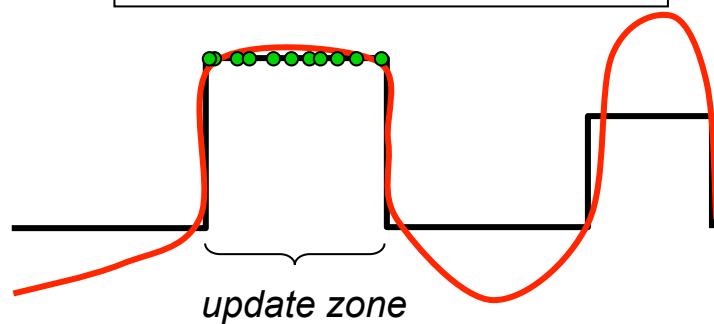
*state: n continuous variables
action: m continuous variable*



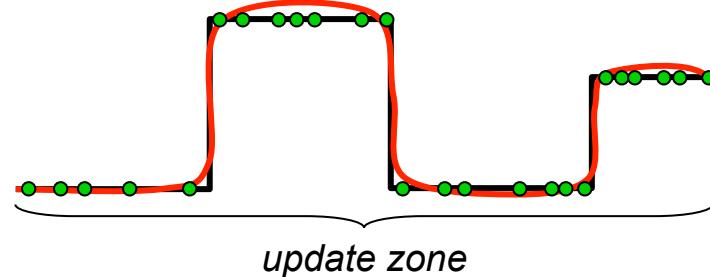
- **NN features:**
 - Multilayer
 - Back-propagation
 - Hidden layers: 1 or 2, hyperbolic tangent function
 - Output layer: linear function
- **Interference problem:**
 - Non-local function approximation



*Updating the NN with
the current samples*



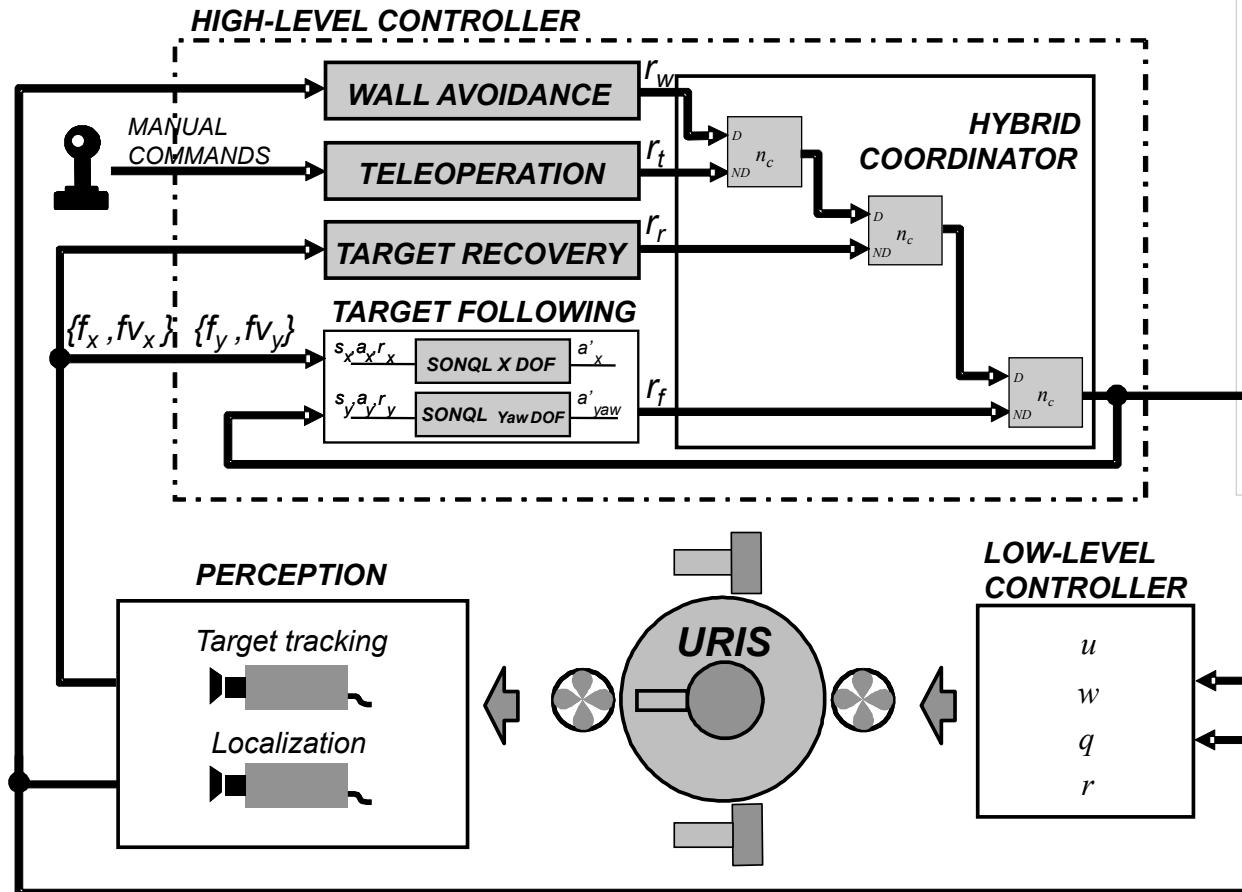
*Updating the NN with a
representative set of samples*



- Interference problem → Uniformly NN updating
- Database of Learning Samples → $\{st, at, rt+1, st+1\}$



Example SONQL: Behaviour learning

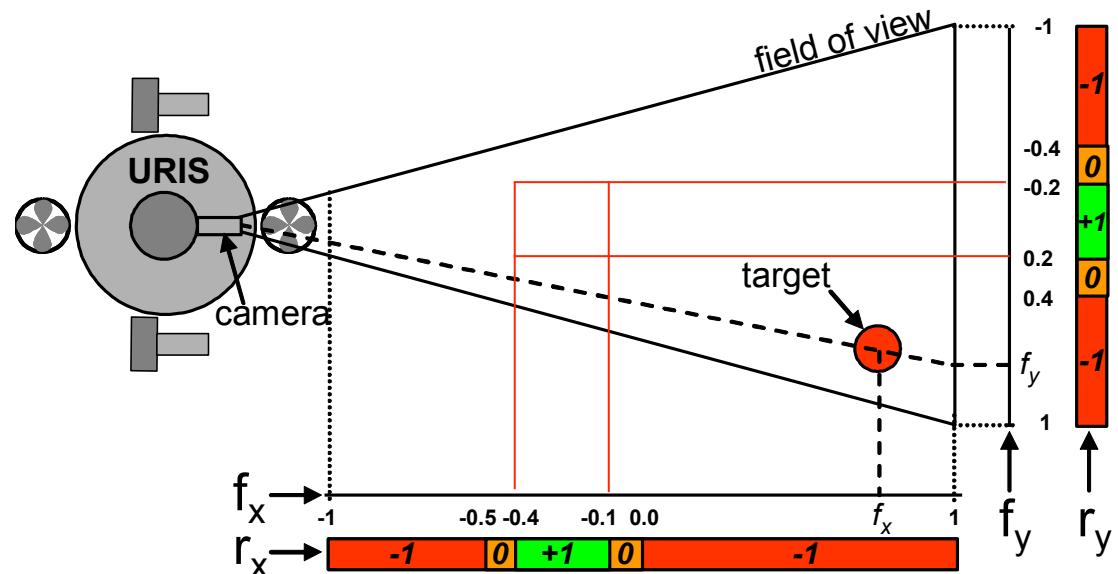


No se puede mostrar la imagen. Puede que su equipo no tenga suficiente memoria para abrir la imagen o que ésta esté dañada. Reinicie el equipo y, a continuación, abra el archivo de nuevo. Si sigue apareciendo la x roja, puede que tenga que borrar la imagen e insertarla de nuevo.

- Task: “**To follow a target avoiding obstacles**”

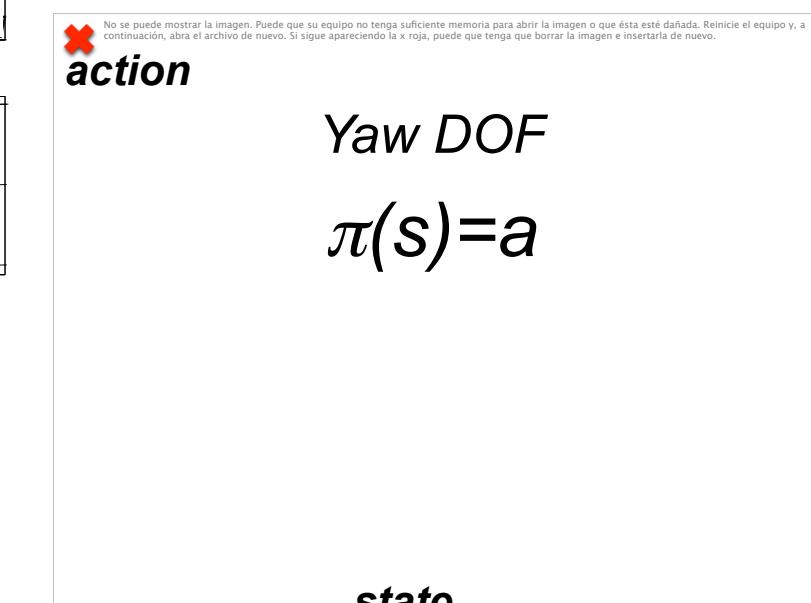
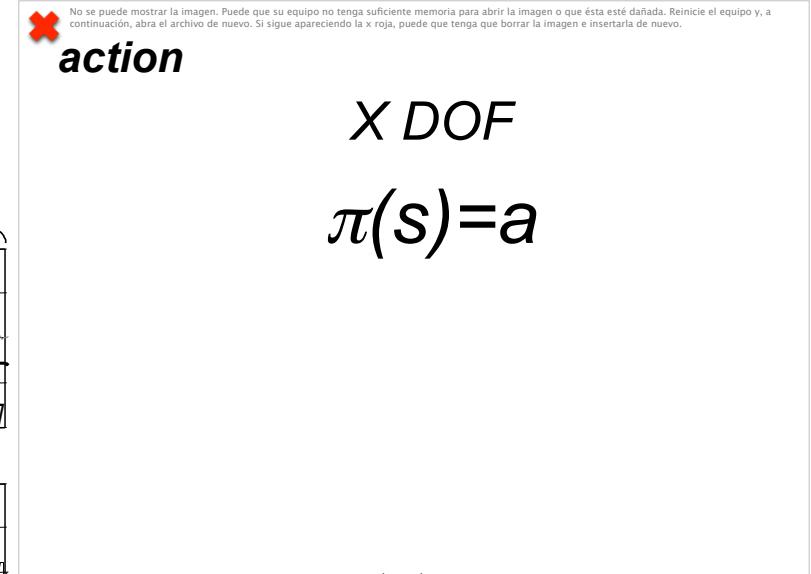
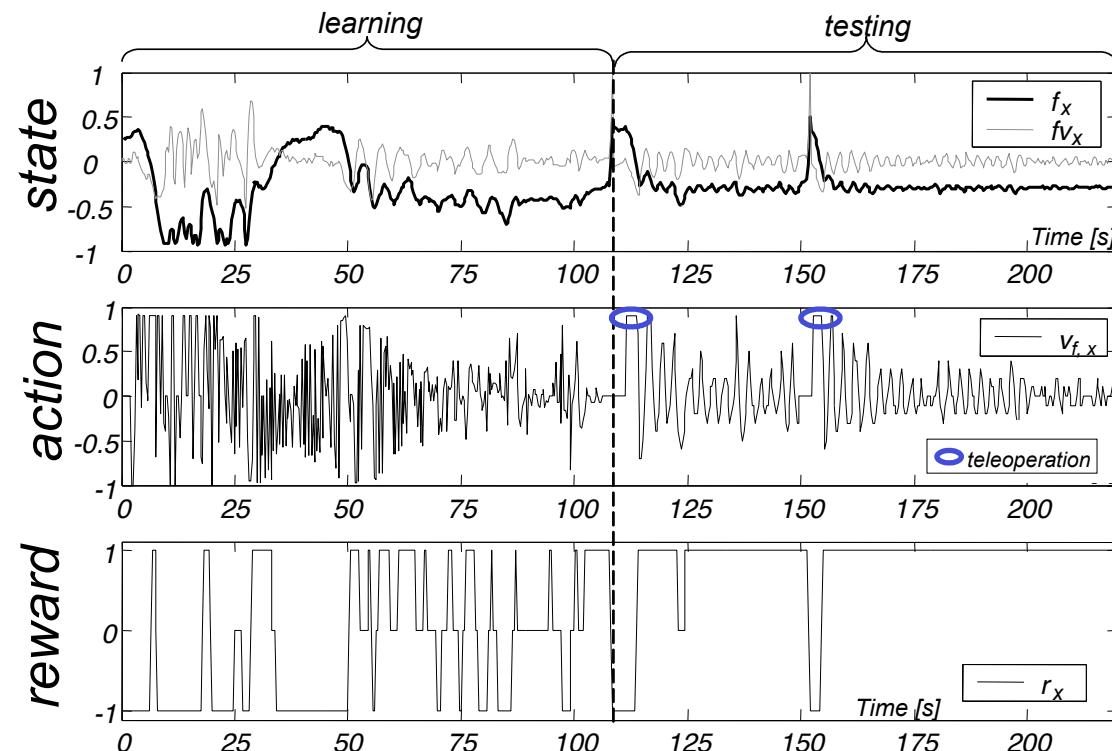
***Target Following
behavior with SONQL***

- *Environment State*
 - X DOF: f_x, fv_x
 - Y DOF: f_y, fv_y
- *Reinforcement Function*
 - X DOF: r_x
 - Y DOF: r_y



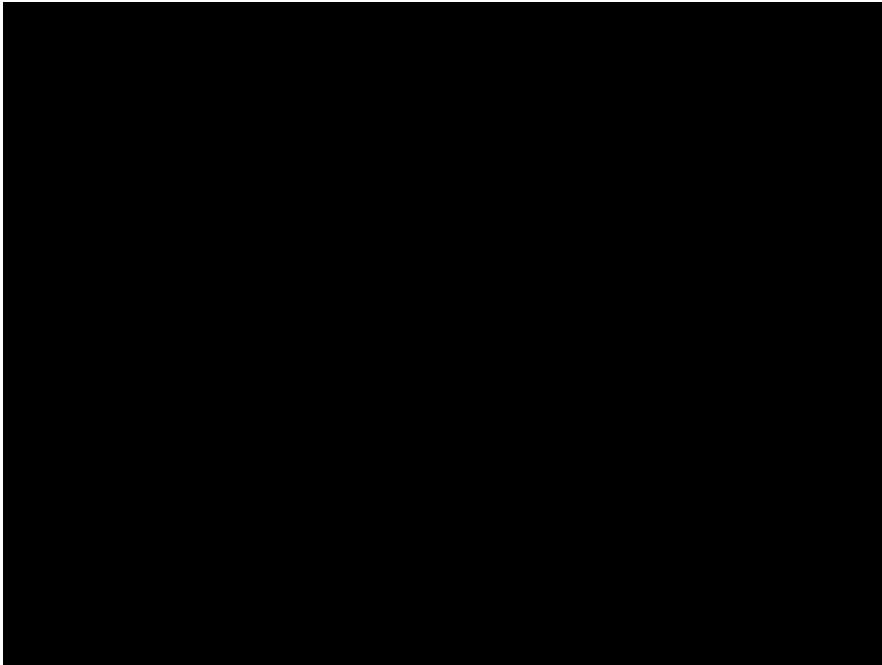
Example SONQL: Behaviour learning

Learning the X DOF

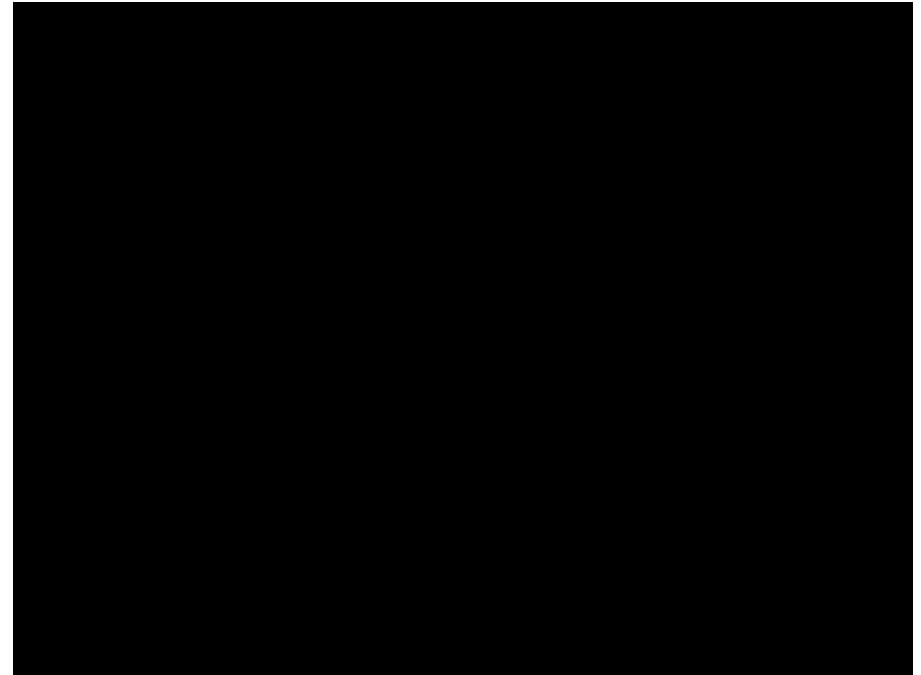




Learning the X DOF

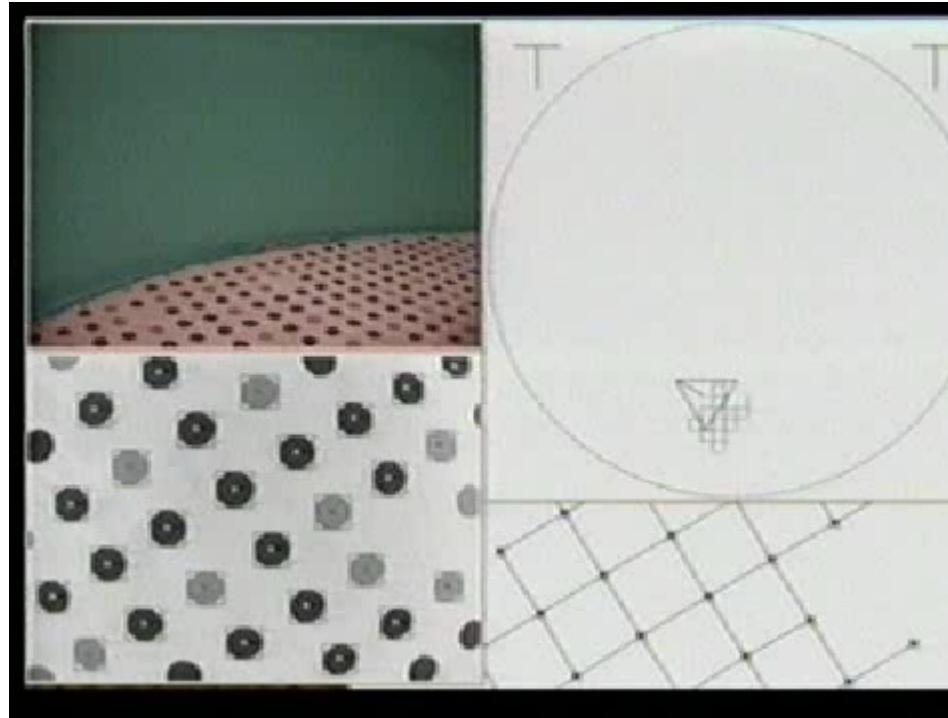


Learning the Yaw DOF





Example SONQL: Behaviour learning

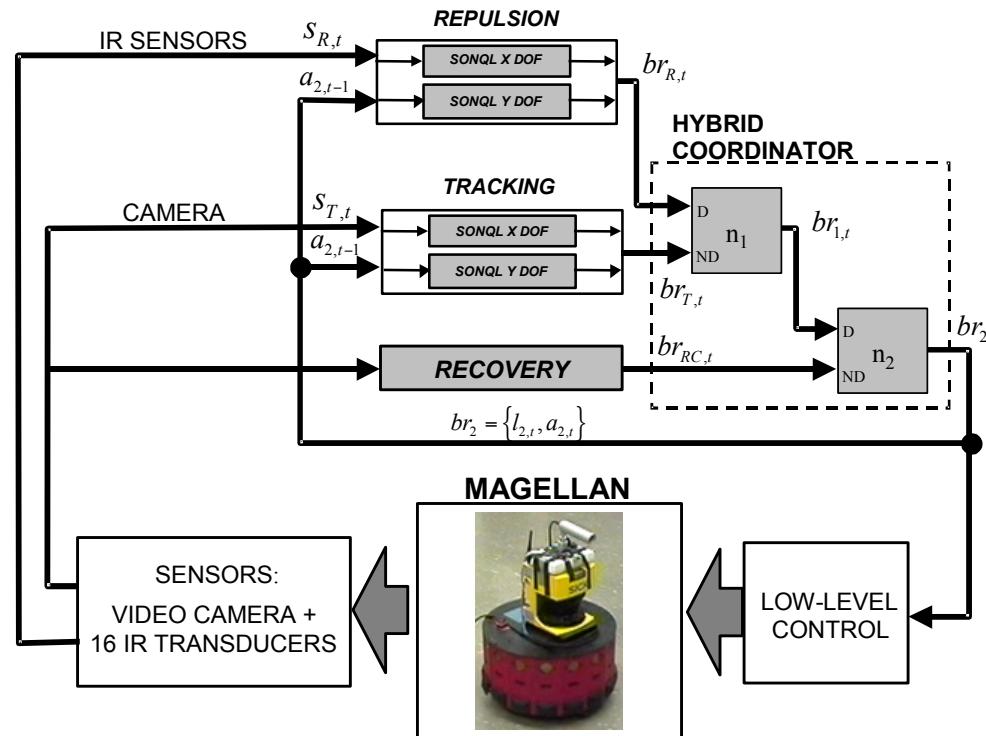
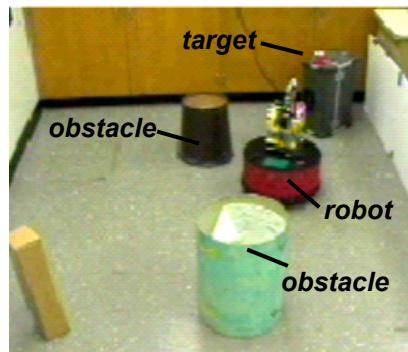
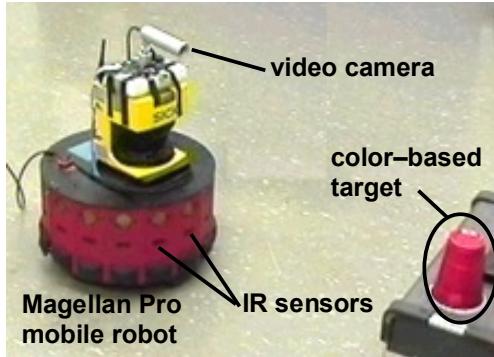




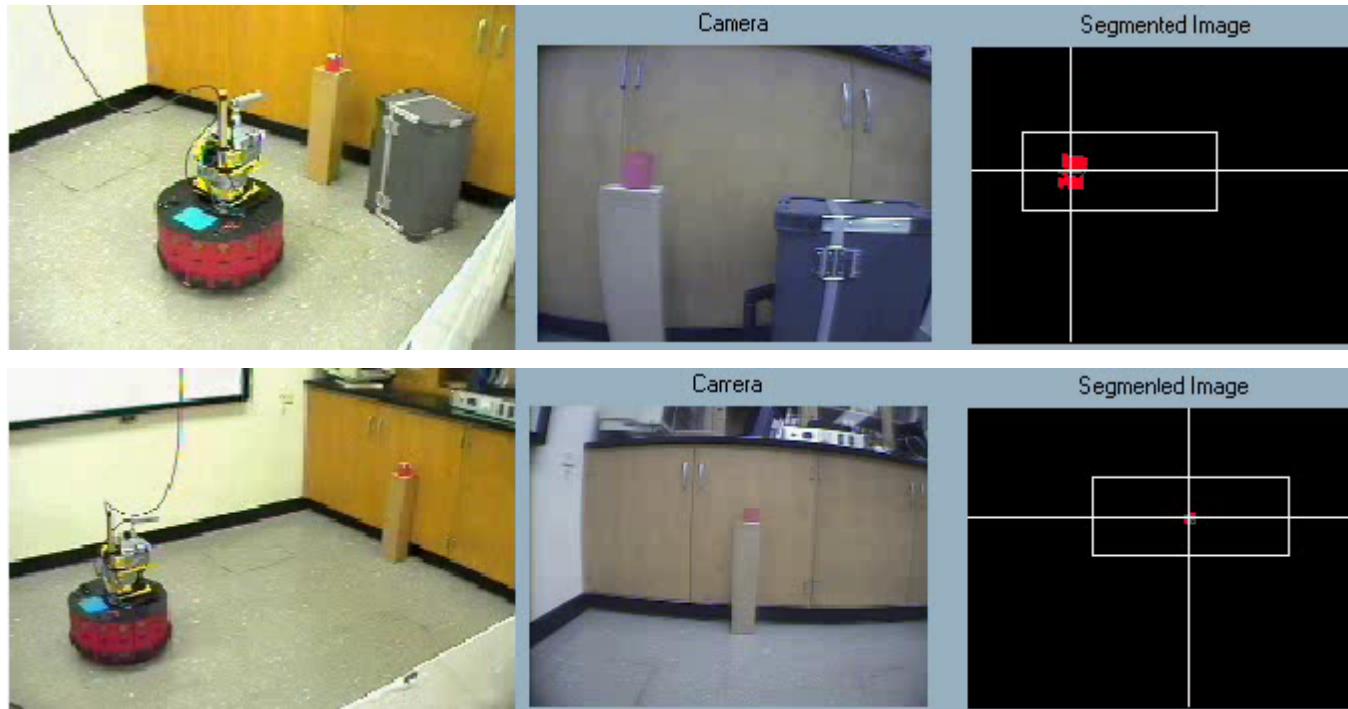
**Efficient Learning of Reactive
Robot Behaviors with a
Neural-Q_learning approach**

**University of Girona
Spain**

IROS 2002, Switzerland



Example SONQL: Behaviour learning 2



Example SONQL: Behaviour learning 2

