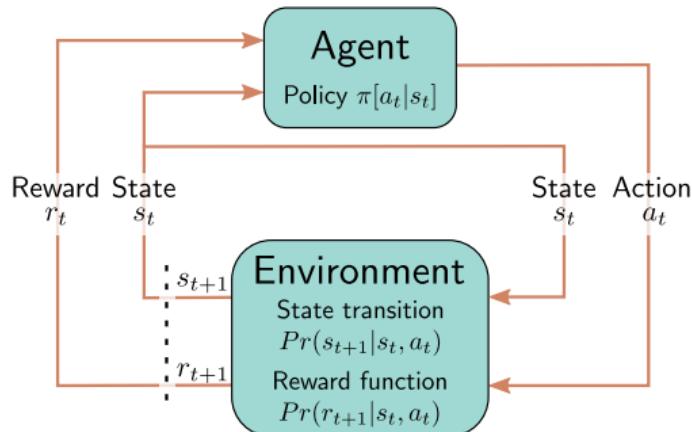


Deep Reinforcement Learning

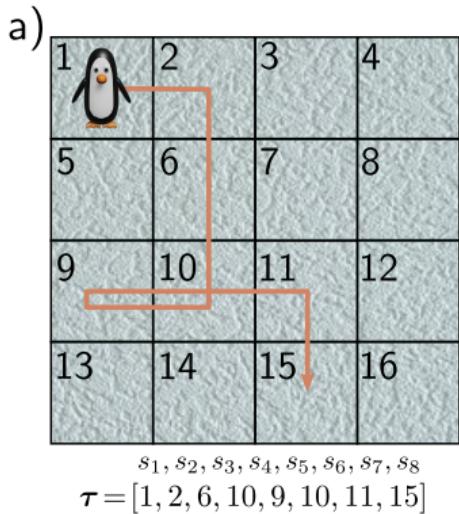
CS 4277 Deep Learning

Kennesaw State University

Reinforcement Learning



Markov Process



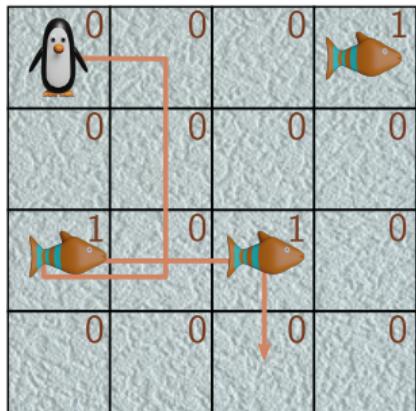
b)

s_{t+1}	s_t															
0	0.33	0	0	0.33	0	0	0	0	0	0	0	0	0	0	0	0
0.5	0	0.33	0	0	0.25	0	0	0	0	0	0	0	0	0	0	0
0	0.33	0	0.5	0	0	0.25	0	0	0	0	0	0	0	0	0	0
0	0	0.33	0	0	0	0	0.33	0	0	0	0	0	0	0	0	0
0.5	0	0	0	0	0.25	0	0	0.33	0	0	0	0	0	0	0	0
0	0.33	0	0	0.33	0	0.25	0	0	0.25	0	0	0	0	0	0	0
0	0	0.33	0	0	0.25	0	0.33	0	0	0.25	0	0	0	0	0	0
0	0	0	0.5	0	0	0.25	0	0	0	0	0.33	0	0	0	0	0
0	0	0	0	0.33	0	0	0	0	0.25	0	0	0.5	0	0	0	0
0	0	0	0	0	0.25	0	0	0.33	0	0.25	0	0	0.33	0	0	0
0	0	0	0	0	0	0.25	0	0	0.25	0	0.33	0	0	0.33	0	0
0	0	0	0	0	0	0.33	0	0	0	0.25	0	0	0.33	0	0.5	0
0	0	0	0	0	0	0	0.33	0	0	0.25	0	0	0.33	0	0	0
0	0	0	0	0	0	0	0	0.33	0	0	0	0.33	0	0	0.33	0
0	0	0	0	0	0	0	0	0	0.25	0	0	0.5	0	0.33	0	0
0	0	0	0	0	0	0	0	0	0.25	0	0	0.33	0	0	0.33	0
0	0	0	0	0	0	0	0	0	0	0.25	0	0	0.33	0	0.5	0
0	0	0	0	0	0	0	0	0	0	0	0.33	0	0	0	0.33	0

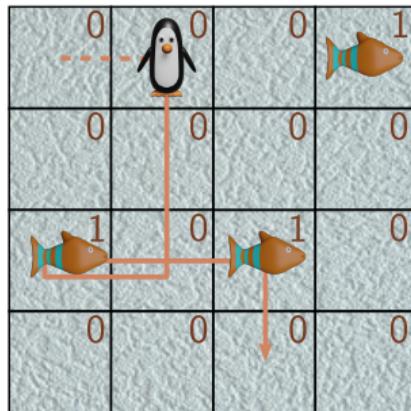
$Pr(s_{t+1}|s_t)$

Markov Reward Process

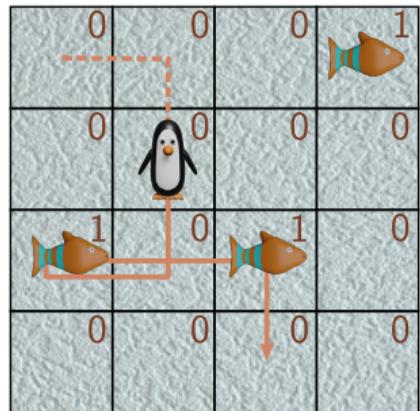
a) $G_1 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 0$
 $+ \gamma^4 \cdot 1 + \gamma^5 \cdot 0 + \gamma^6 \cdot 1 + \gamma^7 \cdot 0 = 1.19$



b) $G_2 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 0 + \gamma^3 \cdot 1$
 $+ \gamma^4 \cdot 0 + \gamma^5 \cdot 1 + \gamma^6 \cdot 0 = 1.31$

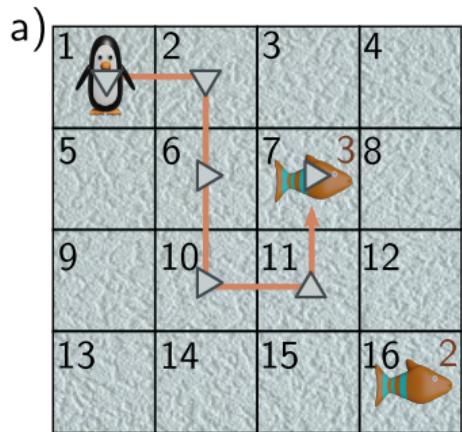
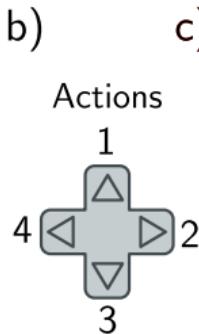


c) $G_3 = 0 + \gamma \cdot 0 + \gamma^2 \cdot 1 + \gamma^3 \cdot 0$
 $+ \gamma^4 \cdot 1 + \gamma^5 \cdot 0 = 1.47$



$s_1 \ r_2 \ s_2 \ r_3 \ s_3 \ r_4 \ s_4 \ r_5 \ s_5 \ r_6 \ s_6 \ r_7 \ s_7 \ r_8 \ s_8 \ r_9$
 $\tau = [1, 0, 2, 0, 6, 0, 10, 0, 9, 1, 10, 0, 11, 1, 15, 0]$

Markov Decision Process


 $\tau = [s_1 \ a_1 \ r_2 \ s_2 \ a_2 \ r_3 \ s_3 \ a_3 \ r_4 \ s_4 \ a_4 \ r_5 \ s_5 \ a_5 \ r_6 \ s_6 \ a_6 \ r_7]$
 $\tau = [1, 3, 0, 2, 3, 0, 6, 2, 0, 10, 2, 0, 11, 1, 0, 7, 2, 3]$


$$Pr(s_{t+1}|s_t=6, a_t=1)$$

$$\begin{bmatrix} 0 \\ 0.5 \\ 0 \\ 0 \\ 0.17 \\ 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$Pr(s_{t+1}|s_t=6, a_t=2)$$

$$\begin{bmatrix} 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

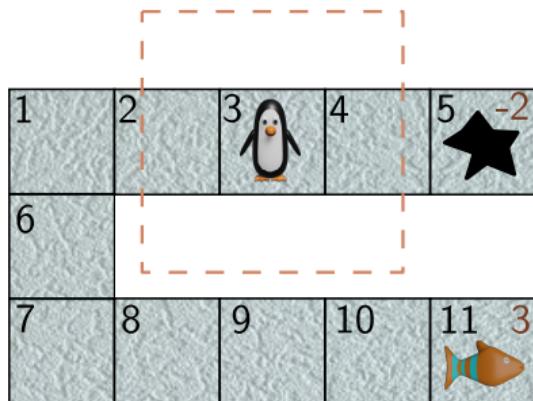
$$Pr(s_{t+1}|s_t=6, a_t=3)$$

$$\begin{bmatrix} 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$Pr(s_{t+1}|s_t=6, a_t=4)$$

$$\begin{bmatrix} 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0.17 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

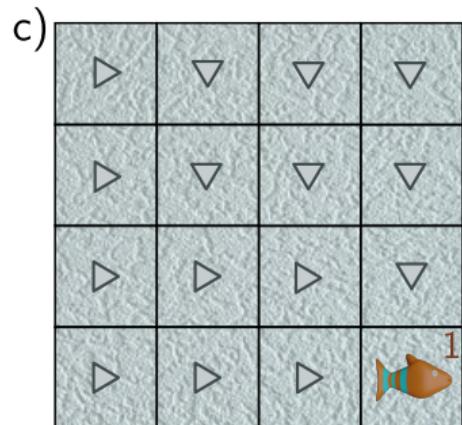
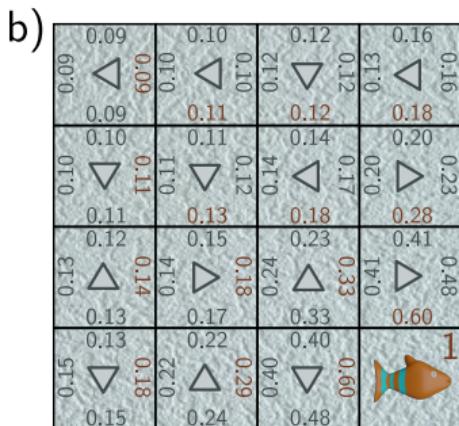
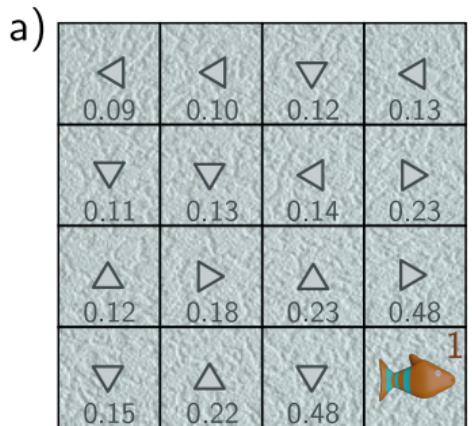
Partially-Observable MDP (POMDP)



Policy

$$\pi(a|s)$$

State and Action Values

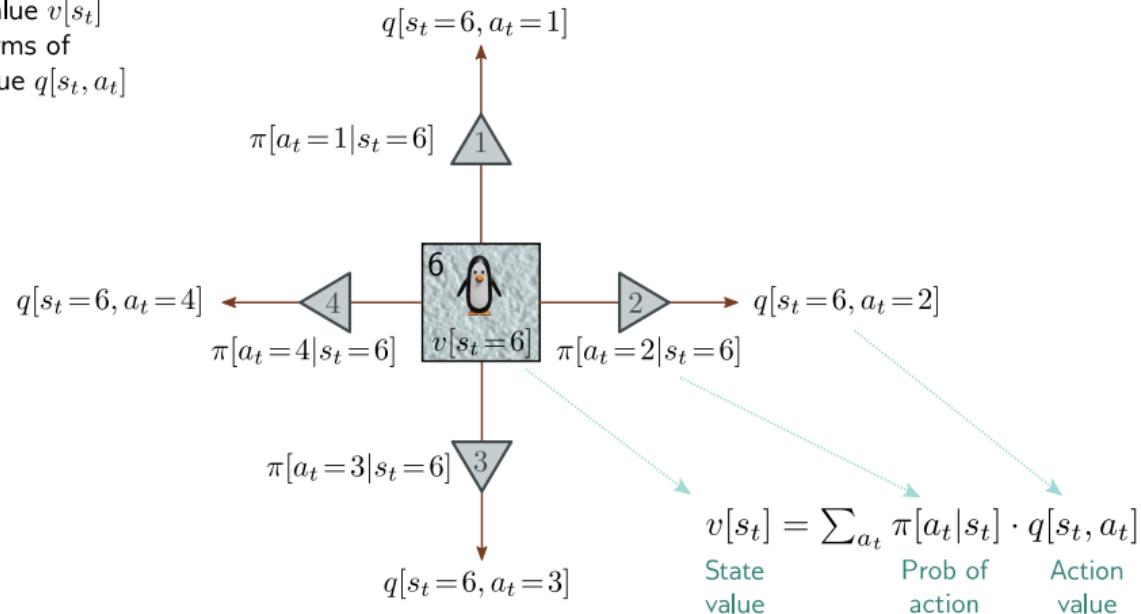


Optimal Policy

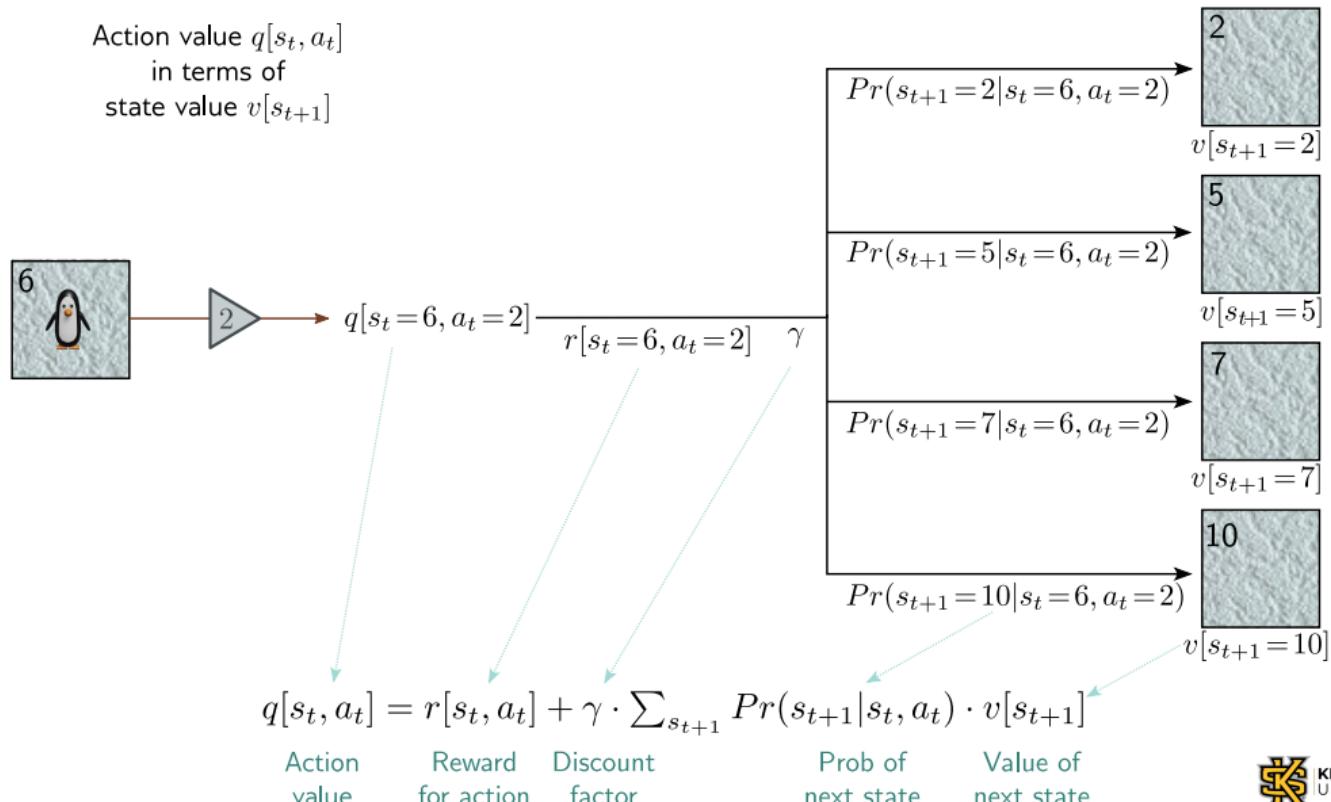
$$\pi(a|s)$$

Bellman Equations

State value $v[s_t]$
in terms of
action value $q[s_t, a_t]$



Bellman Equations 2



Tabular RL

Dynamic Programming

a)

0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00

b)

0.00	0.00	0.00	0.00
0.00	-0.29	-0.29	0.00
-0.45	-2.29	-2.59	0.90
0.00	-0.9	-1.10	3.0

c)

-0.23	-0.17	0.05	0.31
-0.5	-0.62	-0.22	-0.44
-1.16	-3.03	-2.30	0.93
-1.13	-1.53	-1.51	3.00

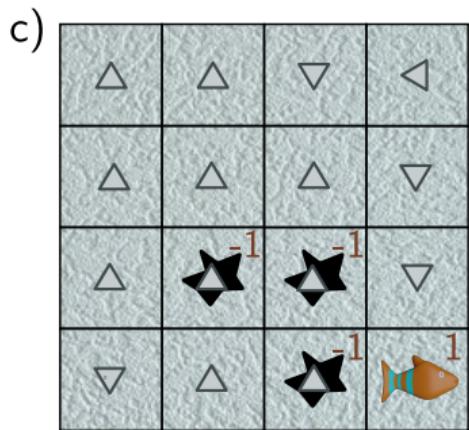
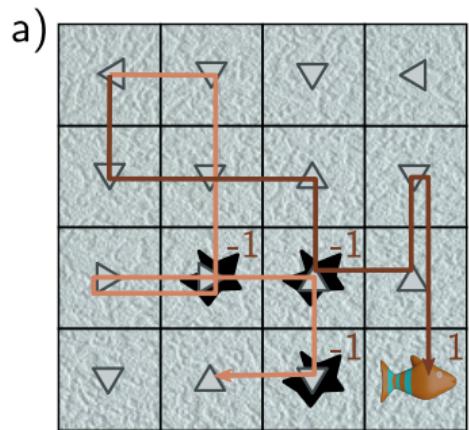
Policy Iteration

$$\pi(a|s)$$

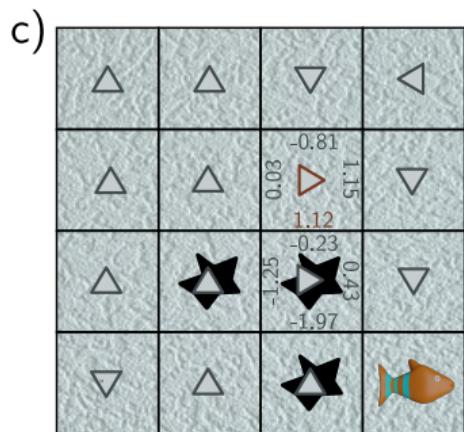
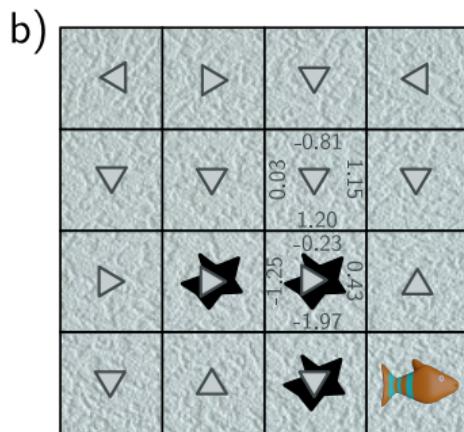
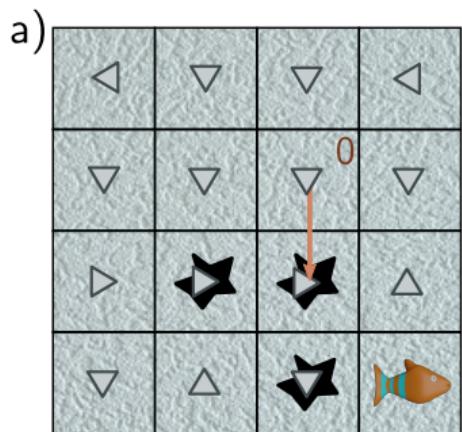
Value Iteration

$$\pi(a|s)$$

Monte Carlo Methods



Temporal Difference Methods

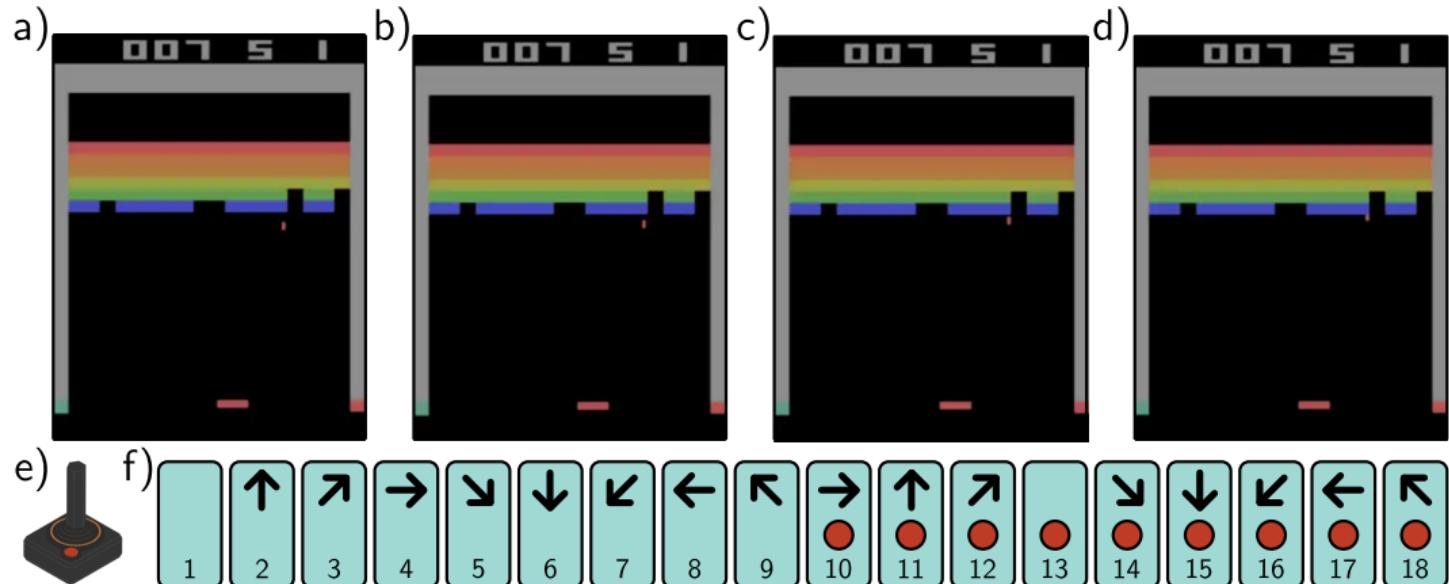


$$q[s_t, a_t] \leftarrow q[s_t, a_t] + \alpha \left(r[s_t, a_t] + \gamma \cdot \max_a [q[s_{t+1}, a]] - q[s_t, a_t] \right)$$
$$1.12 \leftarrow 1.20 + 0.1 \left(0.0 + 0.9 \cdot \max [-0.23, 0.43, -1.97, -1.25] - 1.20 \right)$$

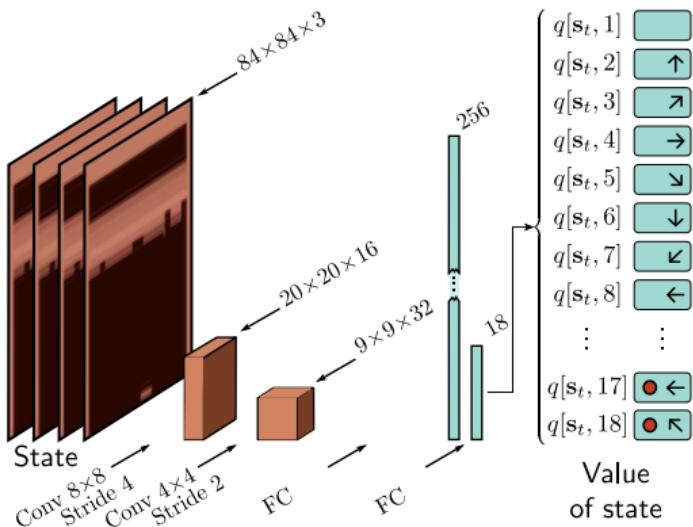
Fitted Q-Learning

$$\pi(a|s)$$

Deep Q-Networks for Playing Atari Games



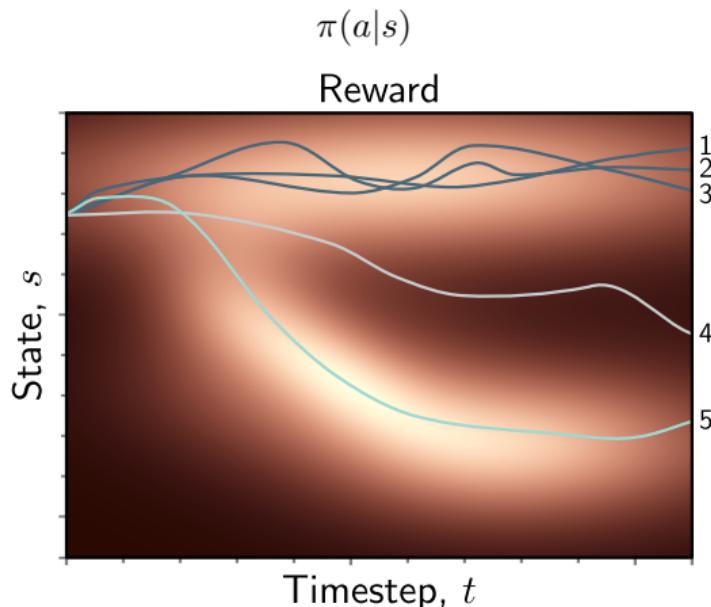
DQN Architecture



Double Q-Learning and Double Deep Q-Networks

$$\pi(a|s)$$

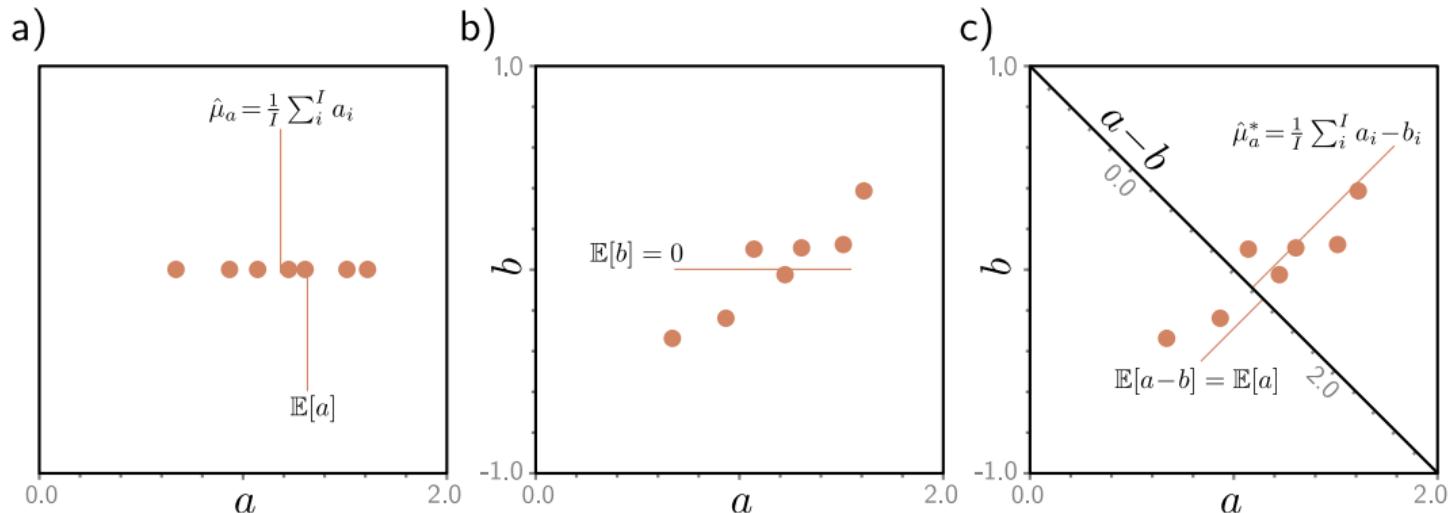
Policy Gradient Methods



REINFORCE Algorithm

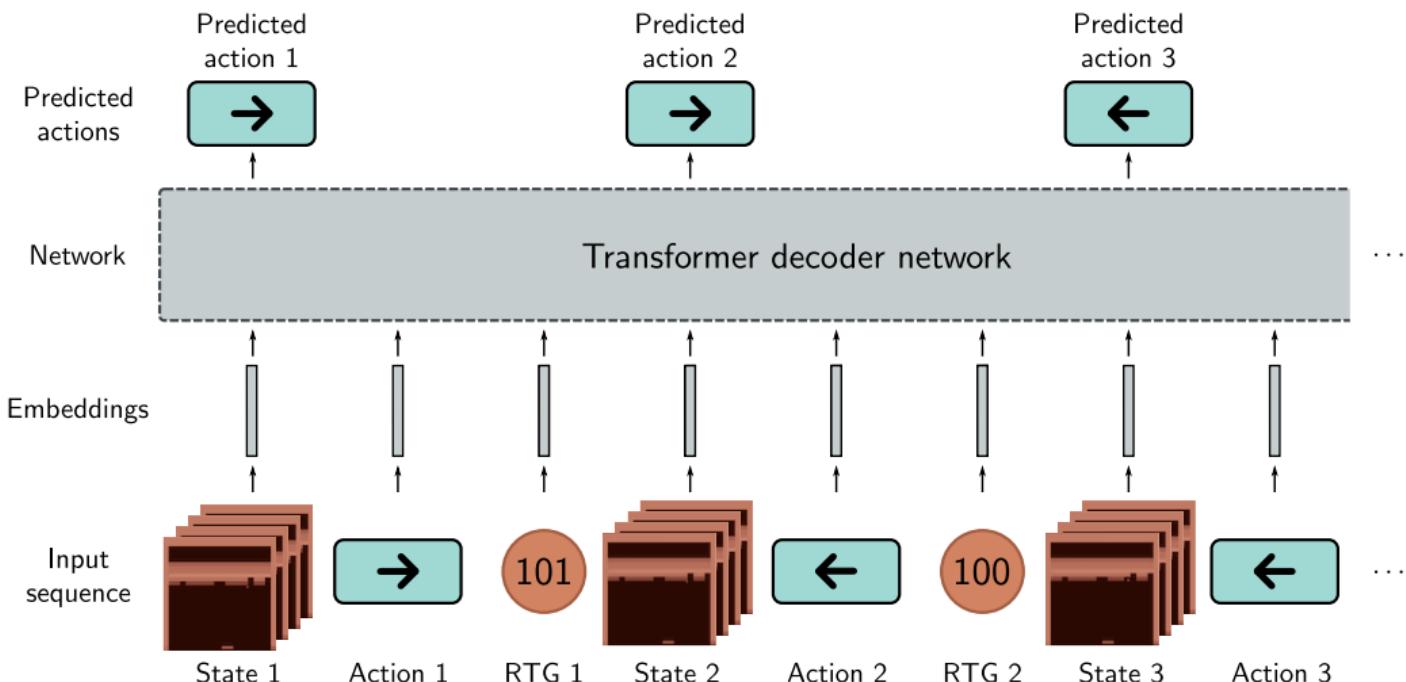
$$\pi(a|s)$$

Baselines



Actor-Critic Methods

Offline Reinforcement Learning



Closing Thoughts

Boom!