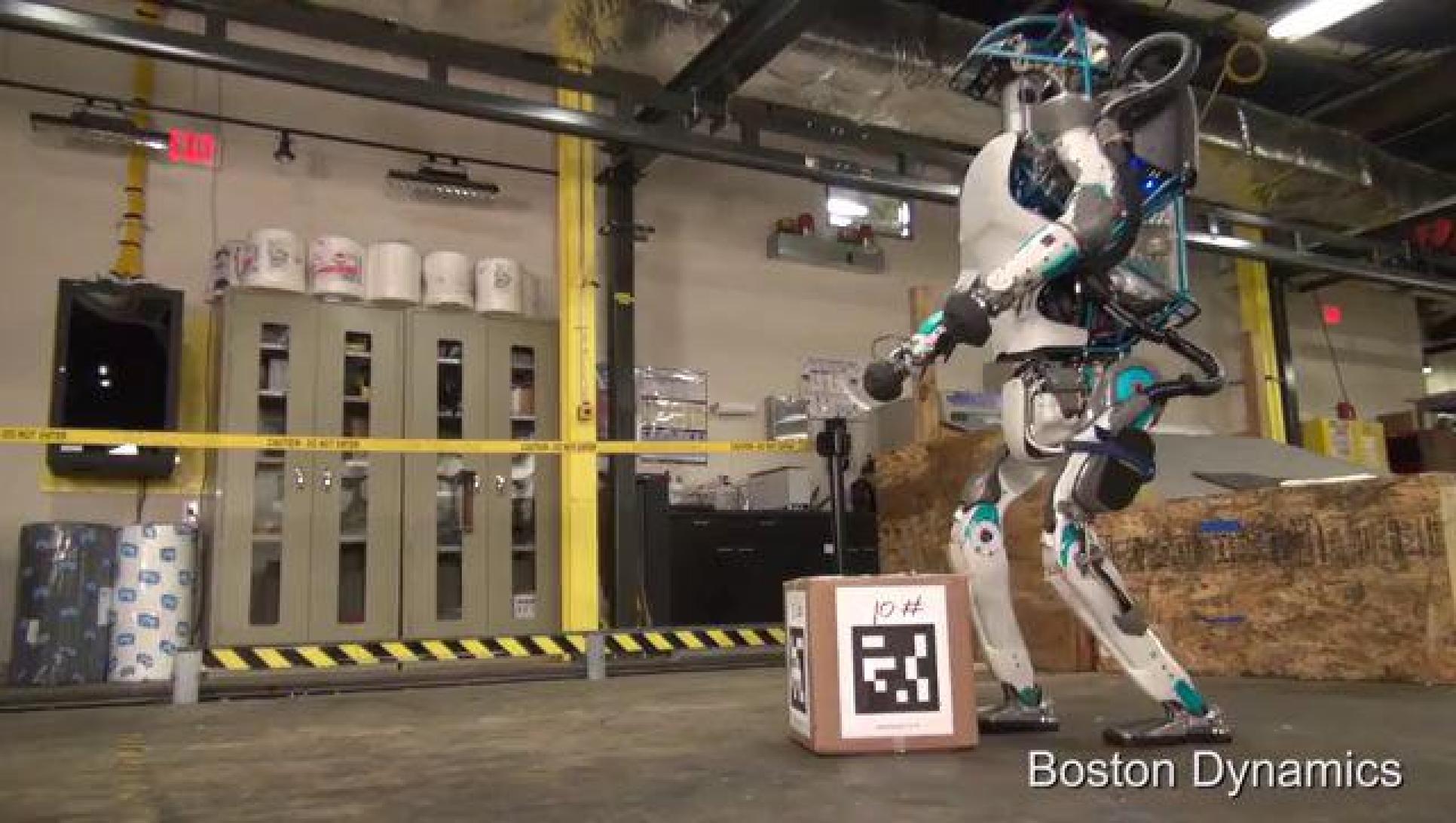


Learning to Do

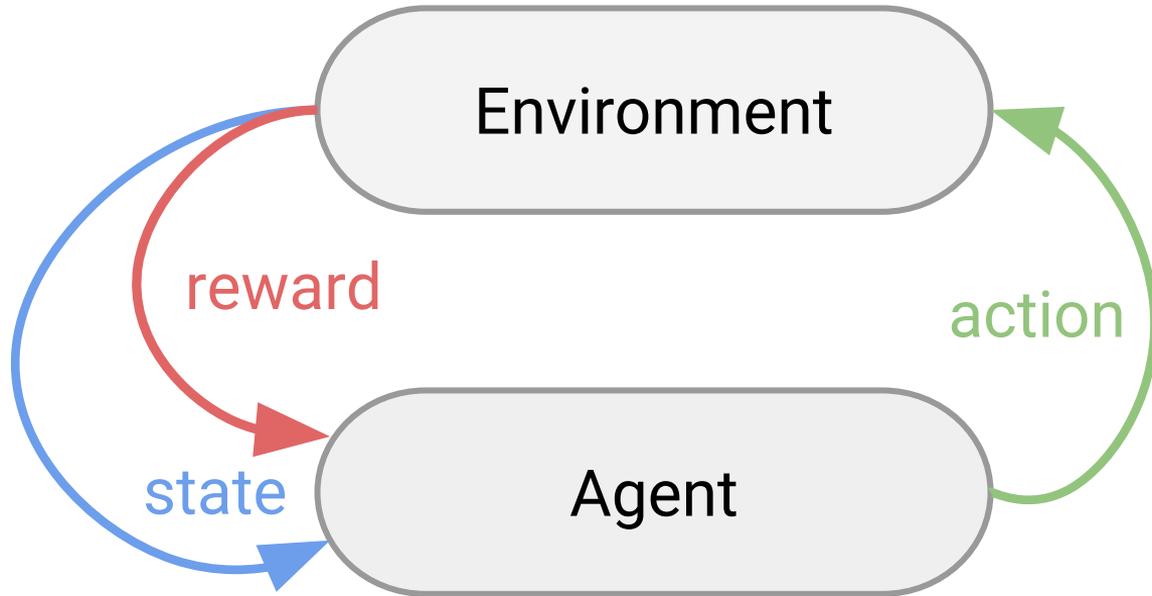
Supervised Learning: given **data**,
predict **labels**

Unsupervised Learning: given **data**,
learn about that **data**

Reinforcement Learning: given **data**,
choose **action** to maximize expected
long-term reward



Boston Dynamics

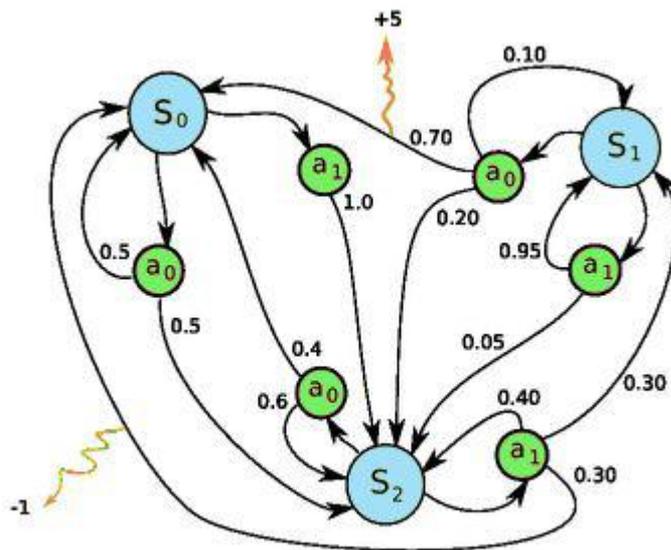


Episode: sequence of states and actions

$s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T, r_T$

Transition function:

$$P(s_{t+1}, r_t \mid s_t, a_t)$$





We want to find a policy $\pi(s) = p(a|s)$ which maximizes

$$\cancel{r_0 + r_1 + r_2 + \dots + r_T} \quad \gamma OLO$$

$$\cancel{r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^T r_T}$$

$$E \left[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^T r_T \right]$$

$$\max_{\pi} E \left[\sum_{i=0}^T \gamma^i r_i \right]$$

Policy Learning

Find $\pi(s)$

$$a \sim \pi(s)$$

Value Learning

Find $Q(s, a)$

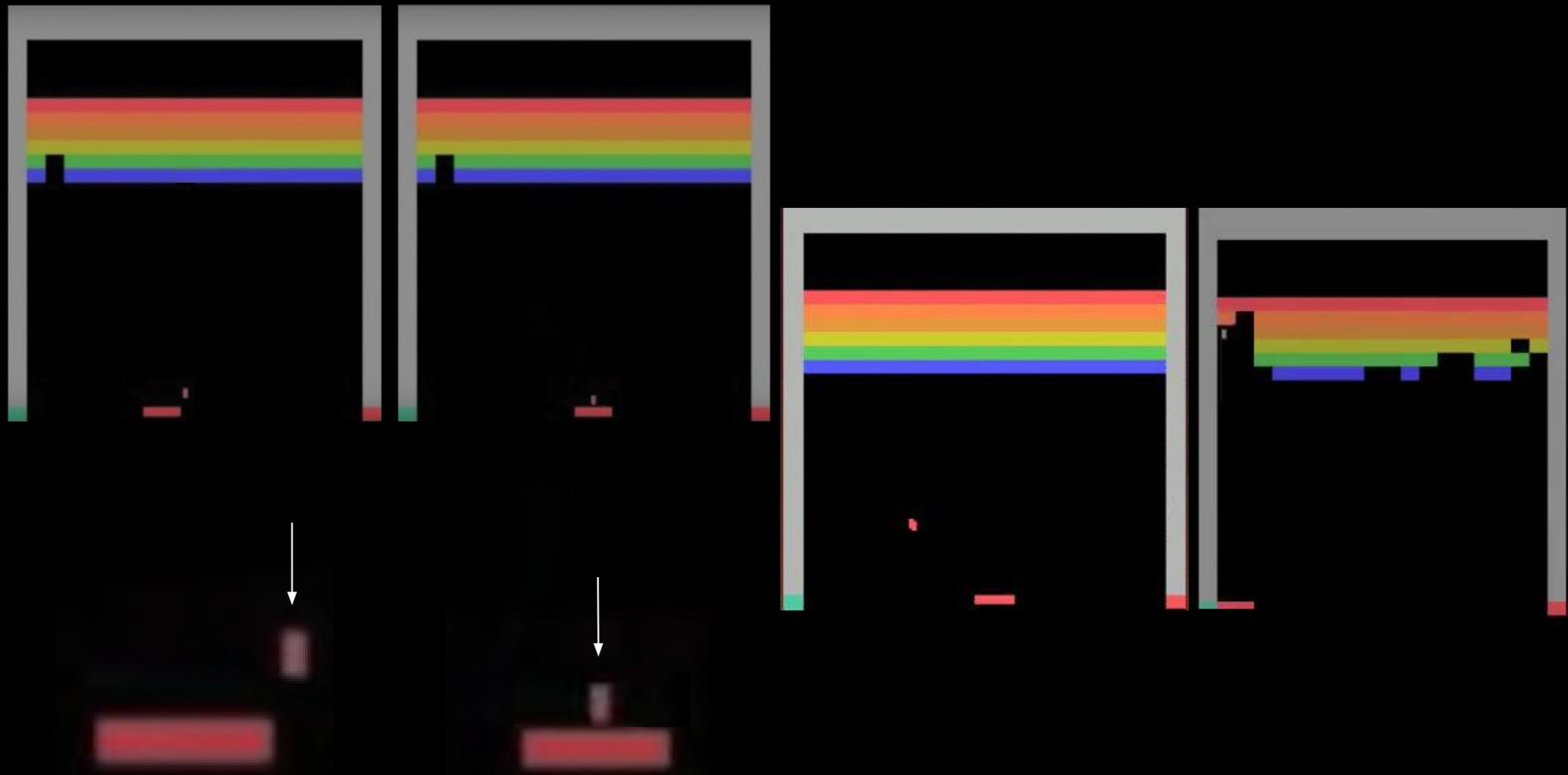
$$a = \arg \max_{a'} Q(s, a')$$

$$Q^*(s_t, a_t) = \max_{\pi} E \left[\sum_{i=t}^T \gamma^i r_i \right]$$

maximum expected future rewards starting at **state s_t** ,
choosing **action a_t** , and then following an **optimal policy π^***

001 5 1





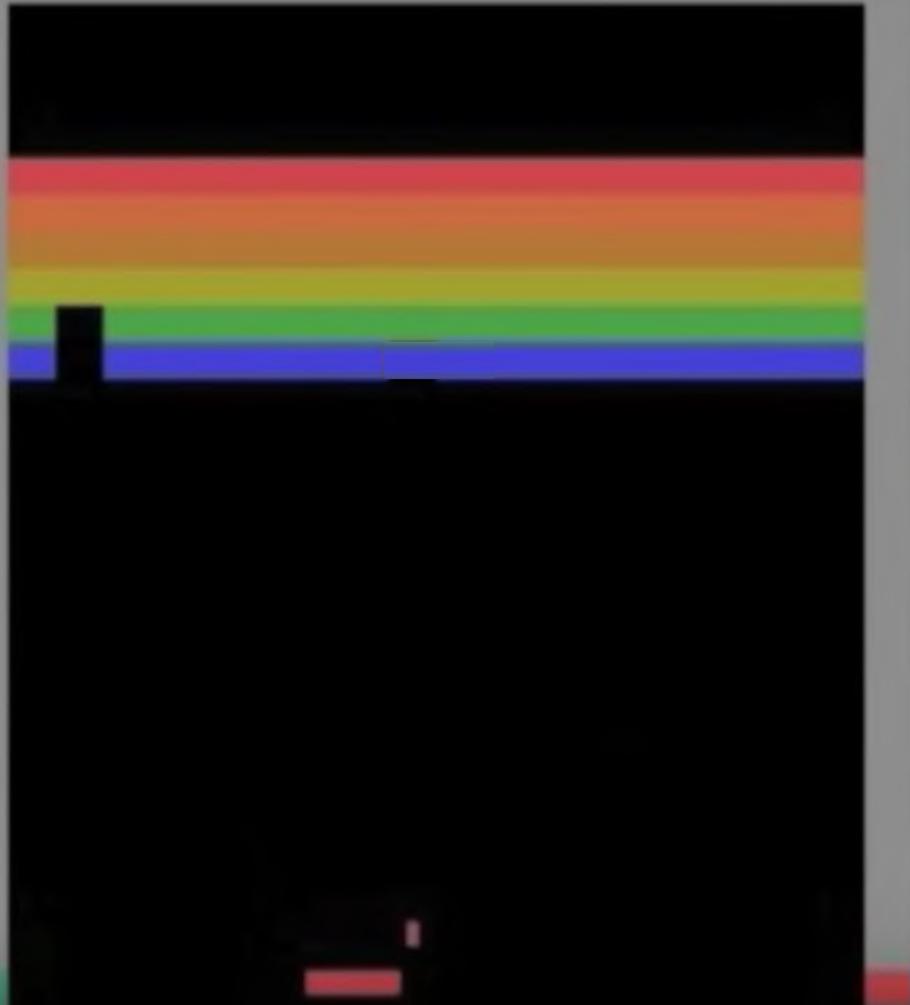
$$Q^*(s_t, a_t) = E \left[r_t + \gamma \max_{a'} Q^*(s_{t+1}, a') \right]$$

The **max future reward** for taking **action a_t** is the **current reward** plus the next step's **max future reward** from taking the best next **action a'**

$$\widehat{Q}_{j+1}(s_t, a_t) \leftarrow E \left[r_t + \gamma \max_{a'} \widehat{Q}_j(s_{t+1}, a') \right]$$

$$\widehat{Q}_j \rightarrow \widehat{Q}_{j+1} \rightarrow \widehat{Q}_{j+2} \rightarrow \dots \rightarrow Q^*$$

But... how large is $Q(\cdot)$?



states: $\sim 2^{96} \cdot 60 \cdot 60$

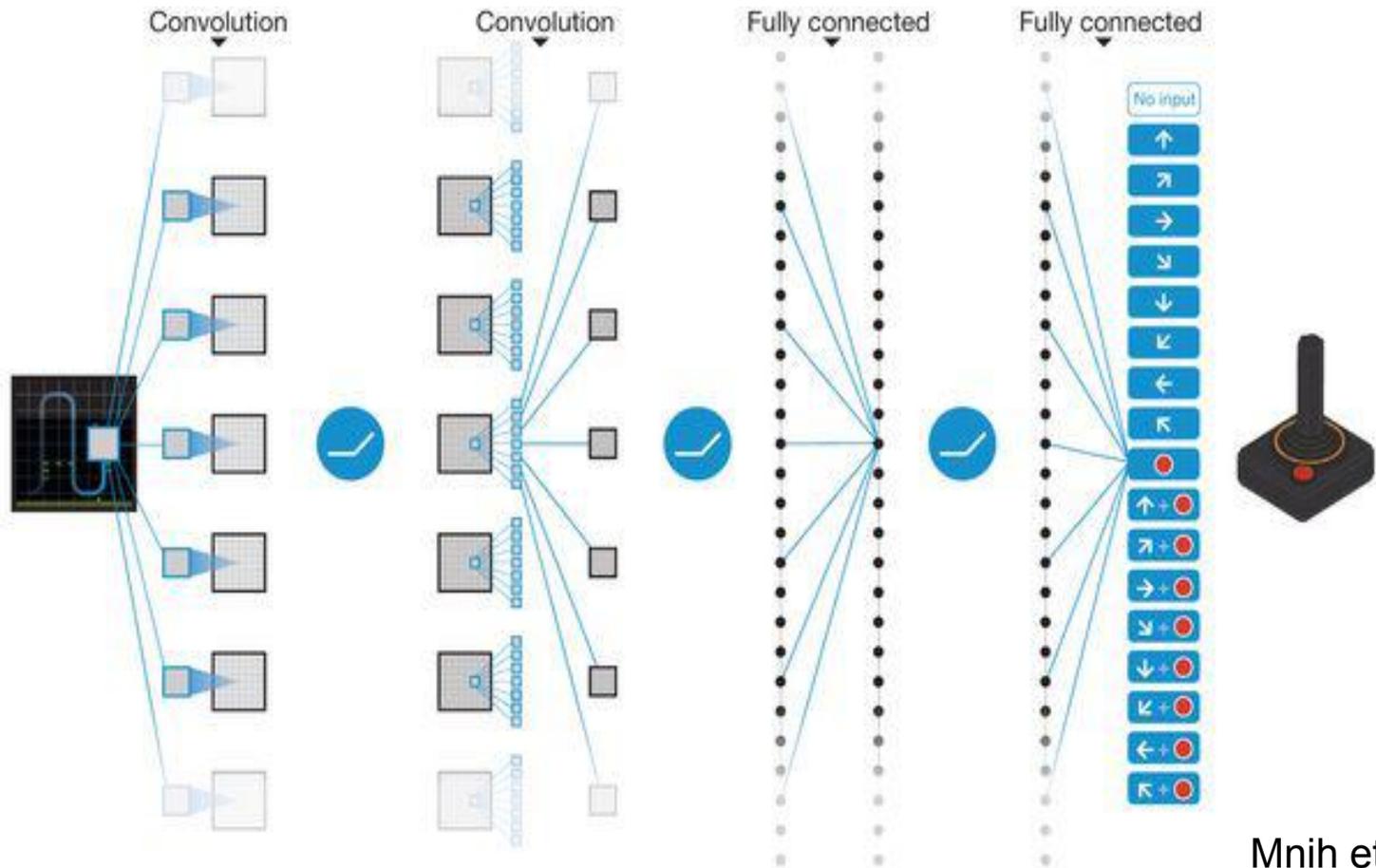
actions: 3

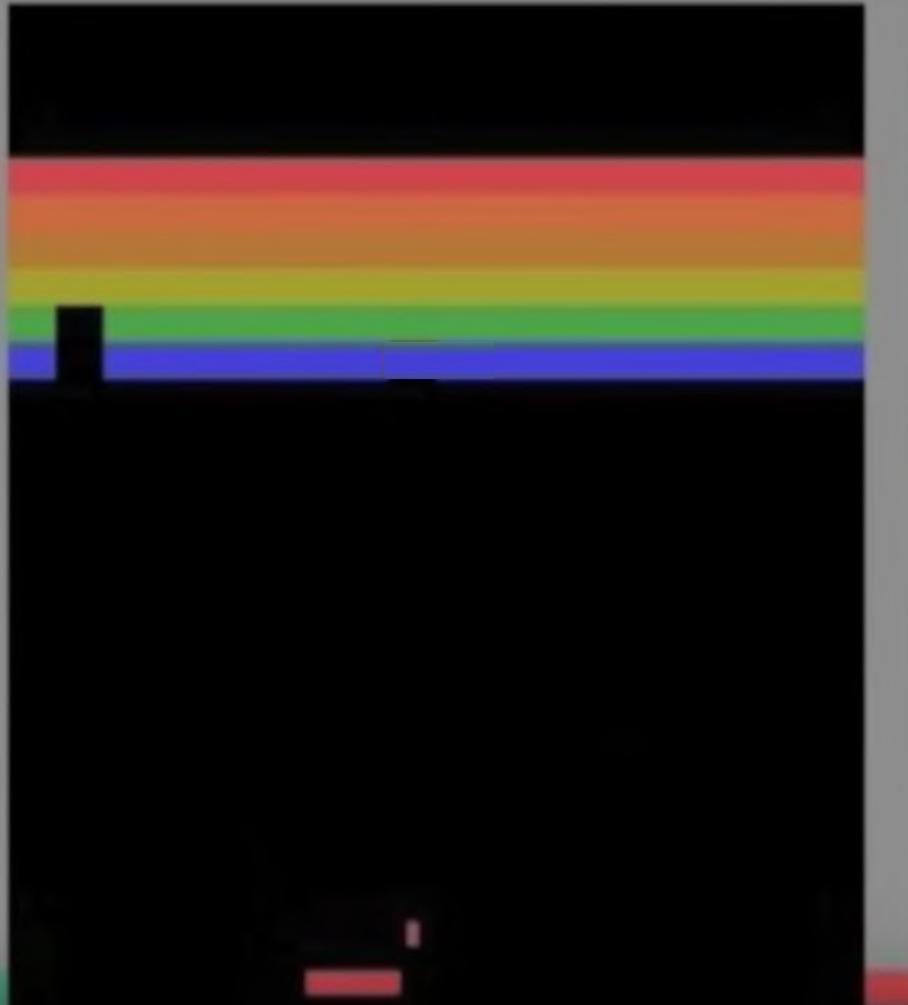
Q values: $\sim 2^{111}$

1957 - 2013

:(

ENTER THE DEEP

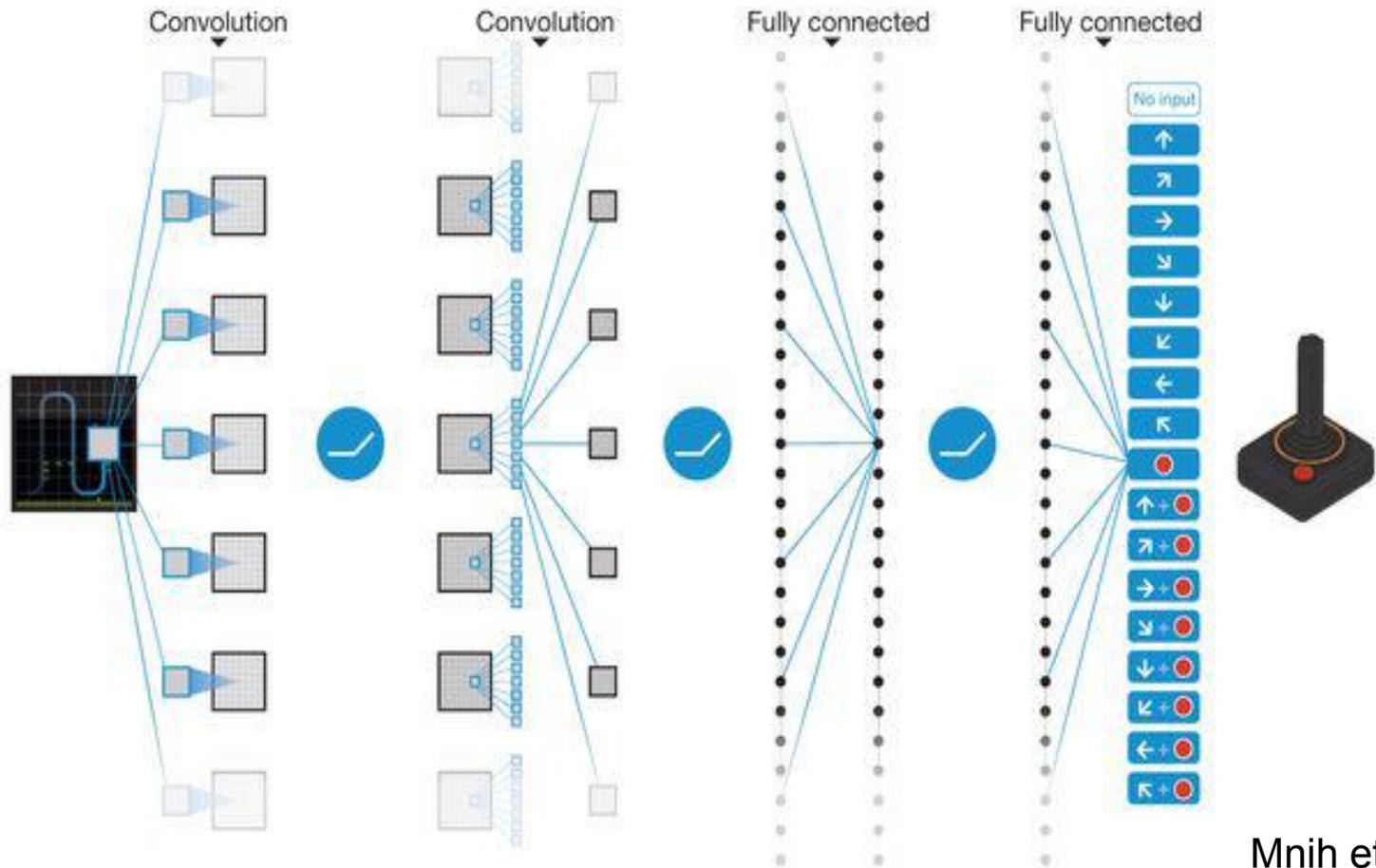




Features for Estimating Q

- Whether the paddle can reach the ball
- # remaining blocks
- How are the blocks spatially arranged?

ENTER THE DEEP



Define approximate Q^* function

$$\widehat{Q}_\theta(s, a|\theta) \sim Q^*(s, a)$$

and choose θ to minimize

$$\min_{\theta} \sum_{e \in E} \sum_{t=0}^T \left\| \widehat{Q}(s_t, a_t|\theta) - \left(r_t + \gamma \max_{a'} \widehat{Q}(s_{t+1}, a'|\theta) \right) \right\|$$

1: **function** Q-LEARNING

2: Initialize θ

3: $s = s_0$

4: **while** not bored yet **do**

5: Choose a from some policy $\pi(s)$, and store results r, s_{new}

6: Compute $\nabla_{\theta} E_Q = \nabla_{\theta} \left\| \hat{Q}(s, a|\theta) - \left(r + \gamma \max_{a'} \hat{Q}(s, a'|\theta_{old}) \right) \right\|$

7: $\theta = \theta - \eta \nabla_{\theta} E_Q$

8: $s = s_{new}$ (or s_0 if episode ended)

9: $\theta_{old} = \theta$

1: **function** Q-LEARNING

2: Initialize θ

3: $s = s_0$

4: **while** not bored **yet do**

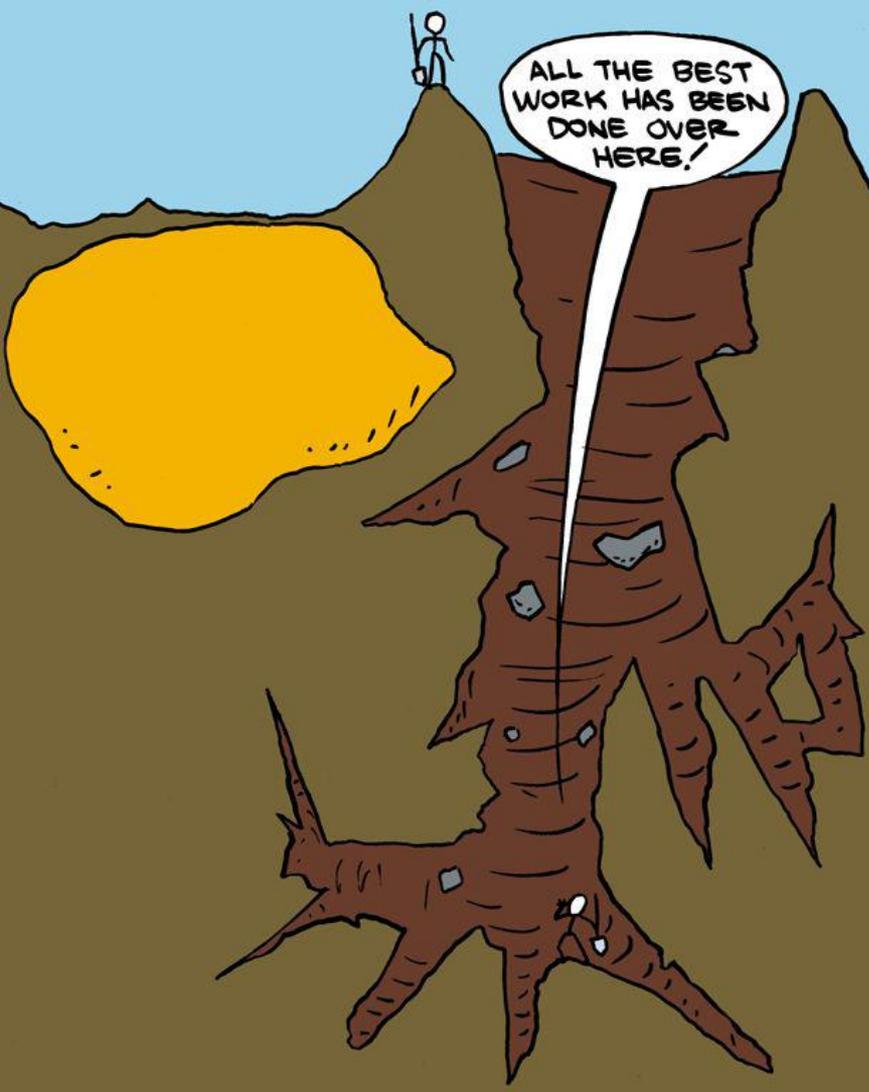
5: **Choose** a from some policy $\pi(s)$, and store results r, s_{new}

6: Compute $\nabla_{\theta} E_Q = \nabla_{\theta} \left\| \hat{Q}(s, a | \theta) - \left(r + \gamma \max_{a'} \hat{Q}(s, a' | \theta_{old}) \right) \right\|$

7: $\theta = \theta - \eta \nabla_{\theta} E_Q$

8: $s = s_{new}$ (or s_0 if episode ended)

9: $\theta_{old} = \theta$



We need to balance
exploration and **exploitation**

ϵ -greedy exploration

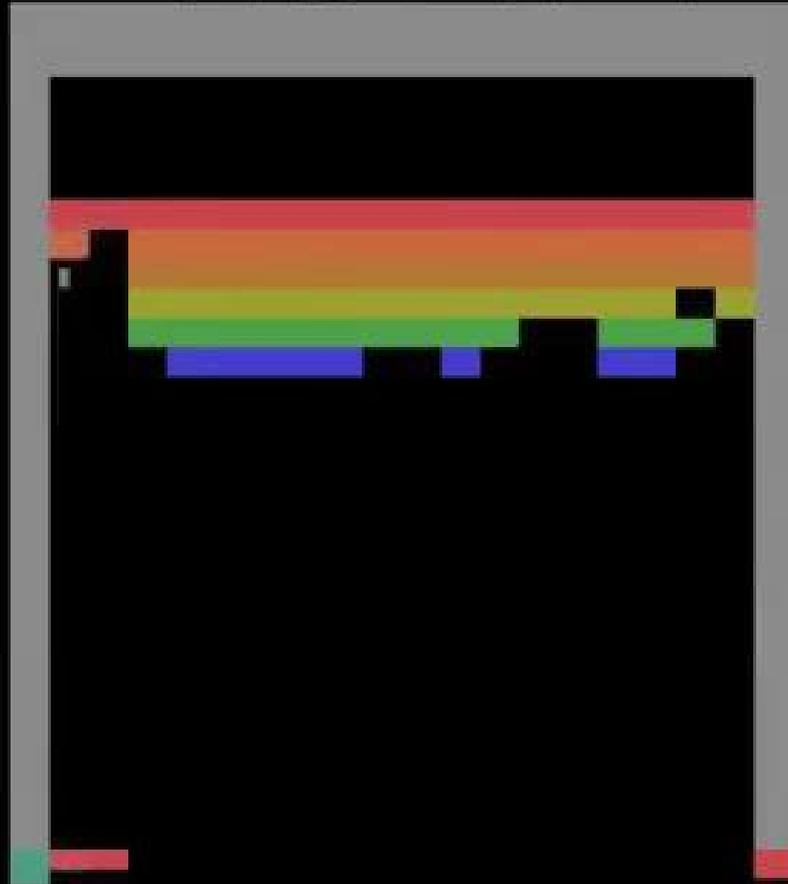
With probability $1 - \epsilon$:

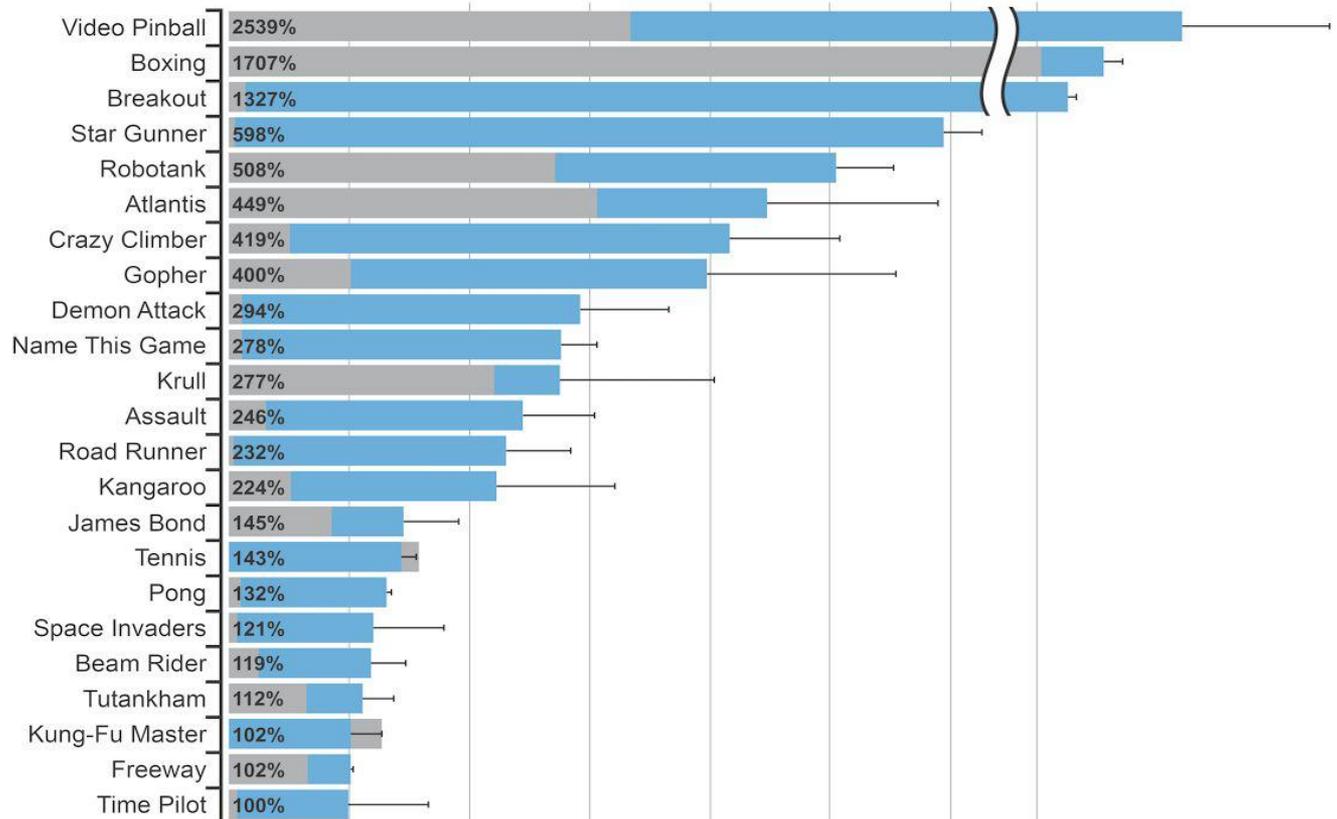
Pick $a_{t+1} \sim \text{soft max}_{a'} \hat{Q}(s_{t+1}, a')$

With probability ϵ :

Pick a_{t+1} at random

042 2 1





% improvement over professional player

Mnih et al, 2013



Try it out!

<http://selfdrivingcars.mit.edu/deeptraffic/>

(Kudos to Lex Fridman & the 6.S094 team!)

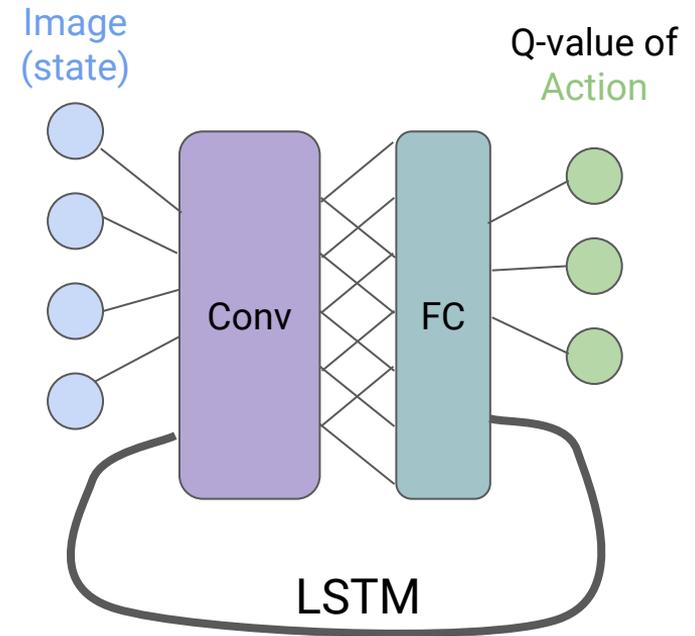
Road Overlay:

The problems with Q-learning

- Restrictive Assumptions
- Handles long horizons poorly
- Requires a simulator

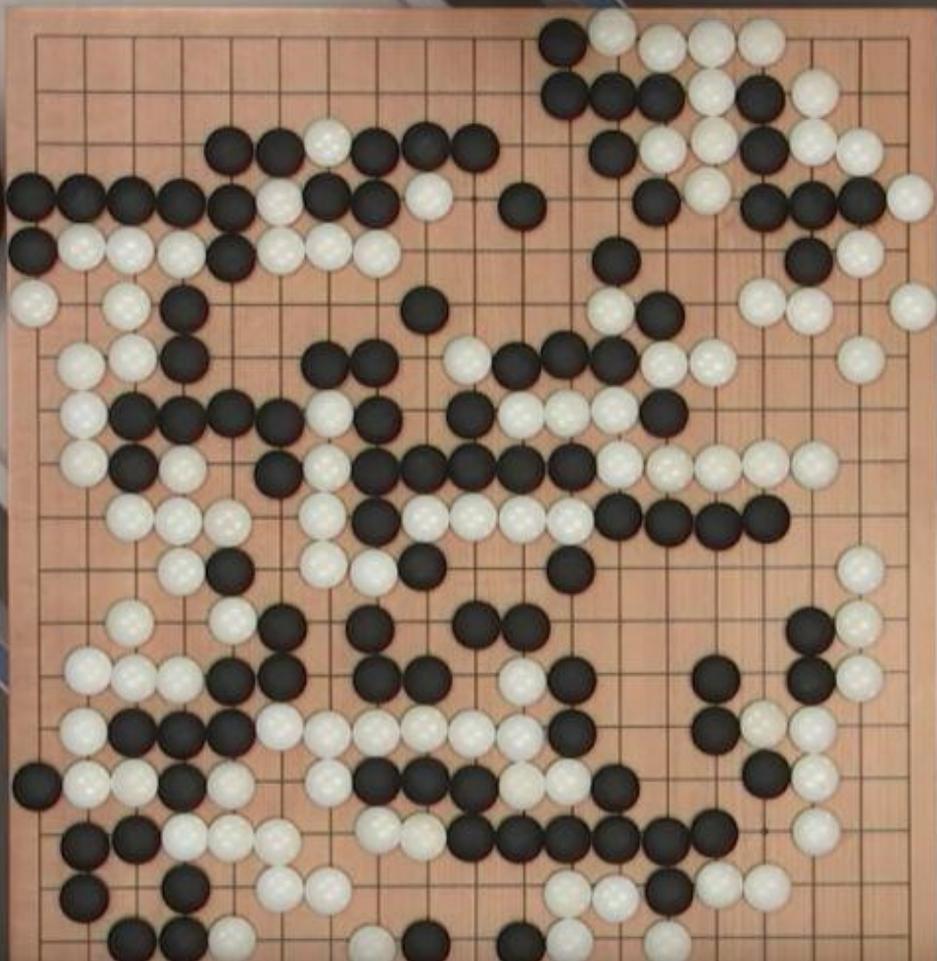


LSTM RNNs!



The problems with Q-learning

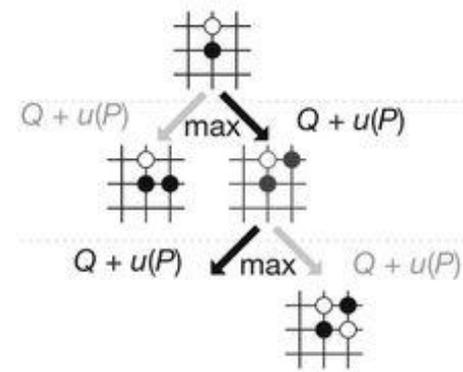
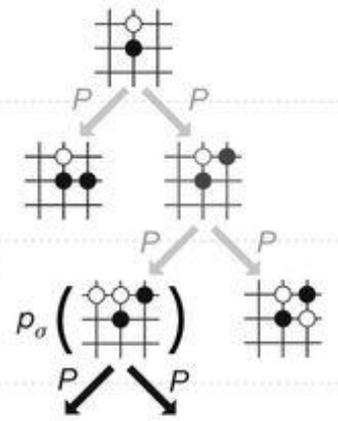
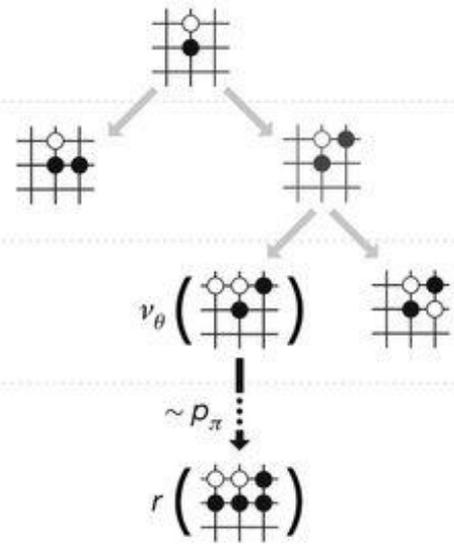
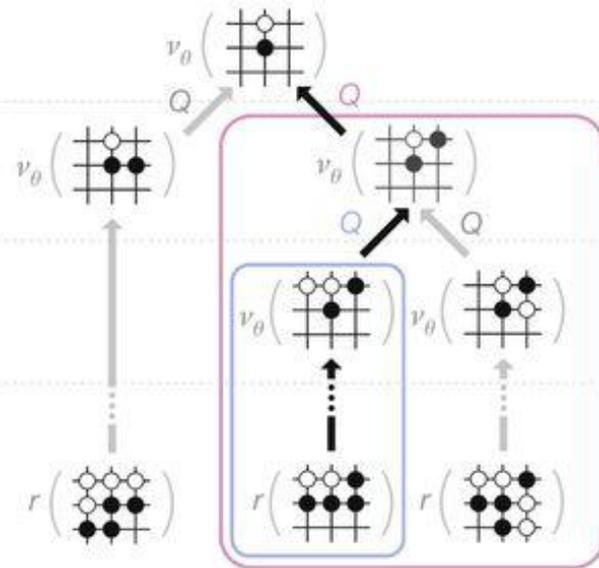
- Restrictive Assumptions
- Handles long horizons poorly
- Requires a simulator



ALPHAGO
00:00:48

LEE SEDOL
00:01:00



a Selection**b** Expansion**c** Evaluation**d** Backup

The problems with Q-learning

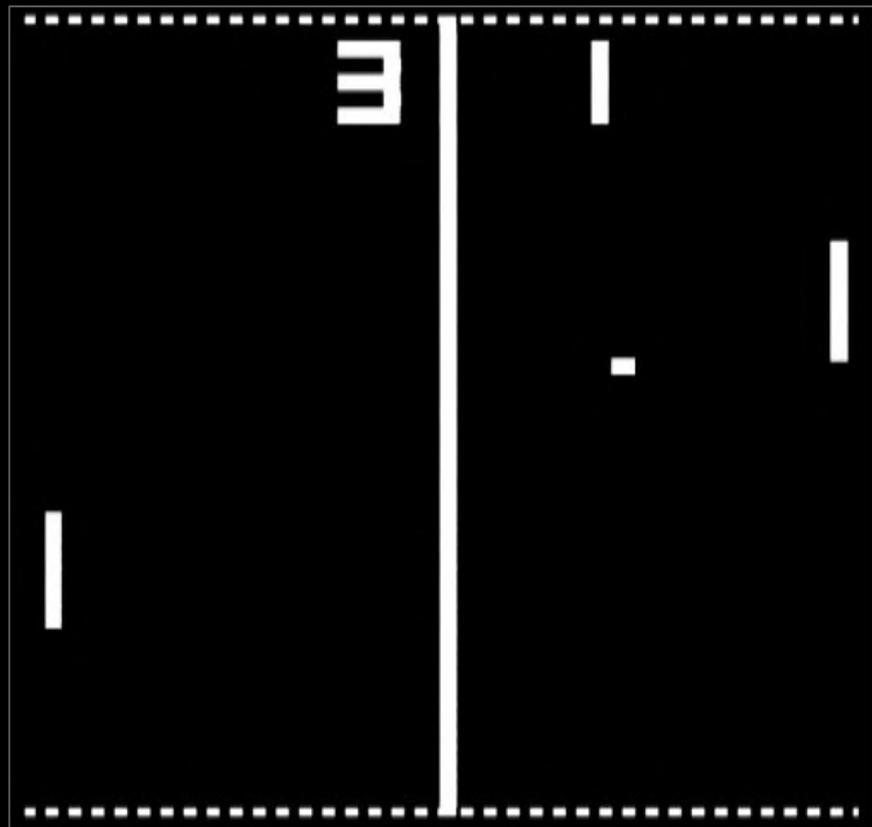
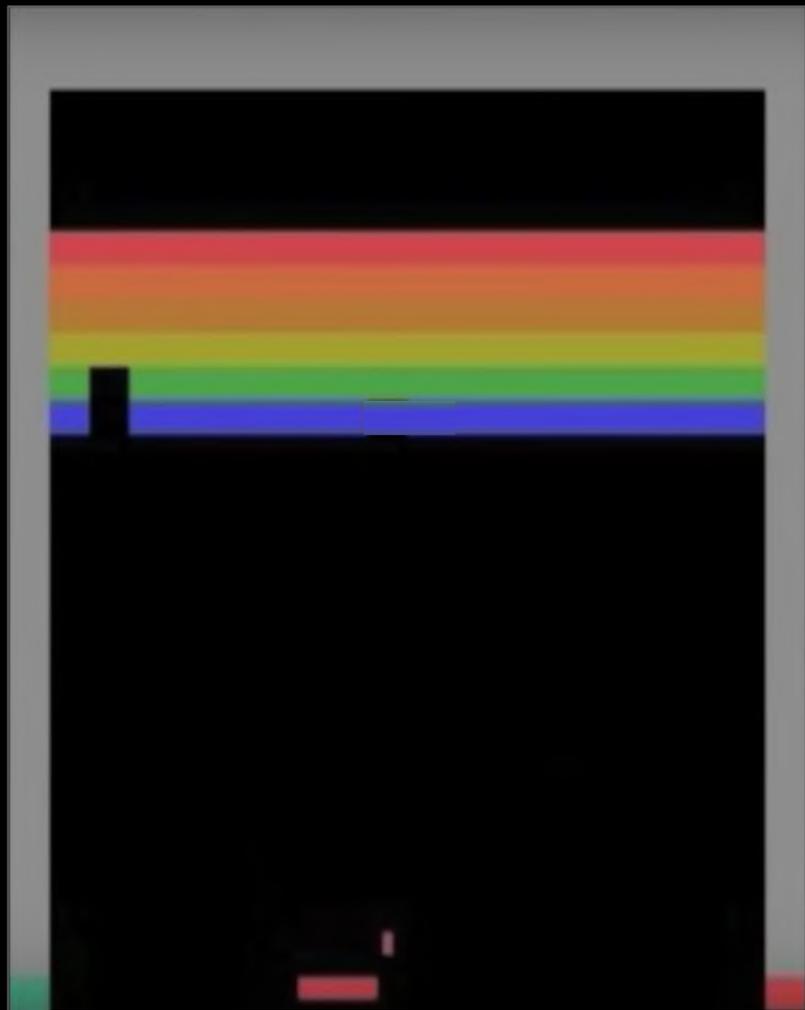
- Restrictive Assumptions
- Handles long horizons poorly
- Requires a simulator





Levine et al, 2016

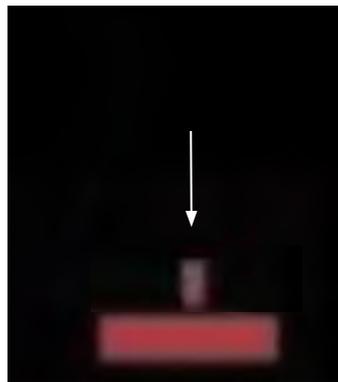
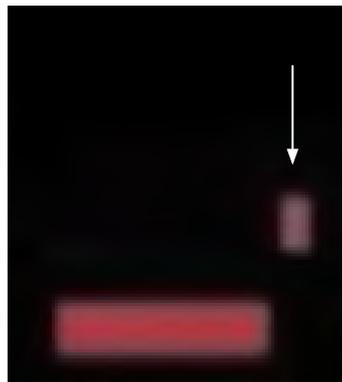
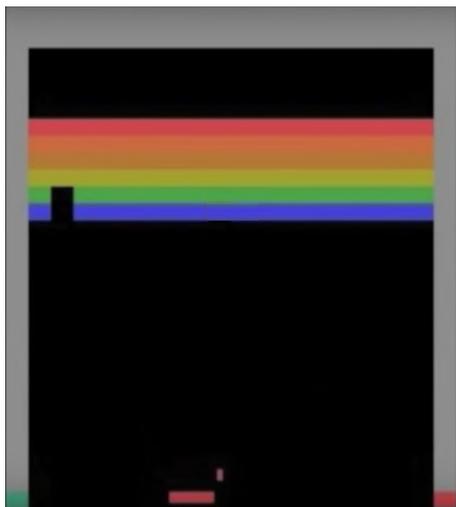








Boston Dynamics



```
1: function  $\epsilon$ -GREEDY-Q-LEARNING
2:   Initialize  $\theta$ 
3:    $\epsilon =$  some tiny number
4:   while not bored yet do
5:      $p = \text{randf}(0,1)$ 
6:     if  $p < \epsilon$  then
7:       Choose random action  $a$  and store results  $r, s_{new}$ 
8:     else
9:       Choose  $a = \arg \max_{a'} Q(s, a')$ , and store results  $r, s_{new}$ 
10:    Compute  $\nabla_{\theta} E_Q = \nabla_{\theta} \left\| \widehat{Q}(s, a|\theta) - \left( r + \gamma \max_{a'} \widehat{Q}(s, a'|\theta_{old}) \right) \right\|$ 
11:     $\theta = \theta - \eta \nabla_{\theta} E_Q$ 
12:     $s = s_{new}$  (or  $s_0$  if episode ended)
13:     $\theta_{old} = \theta$ 
```

