

# CS 188: Artificial Intelligence

## Reinforcement Learning



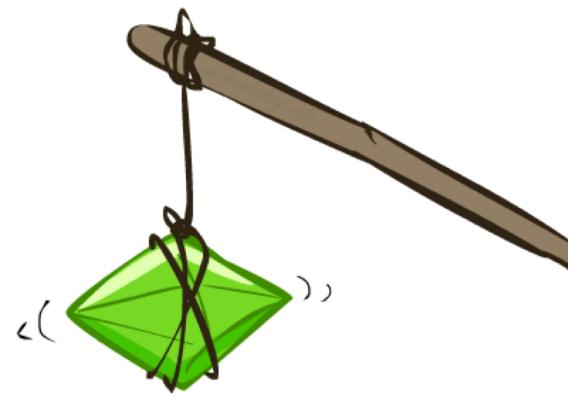
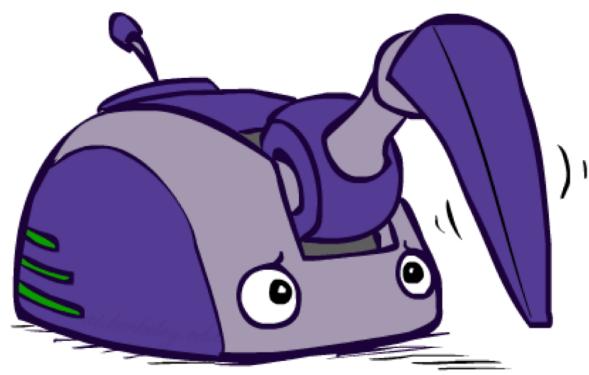
Instructor: Anca Dragan

University of California, Berkeley

[Slides by Dan Klein, Pieter Abbeel, Anca Dragan. [http://ai.berkeley.edu.\]](http://ai.berkeley.edu.)

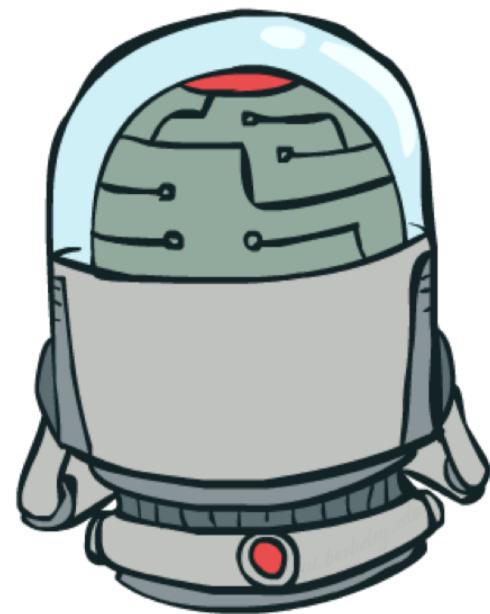
# Reinforcement Learning

---



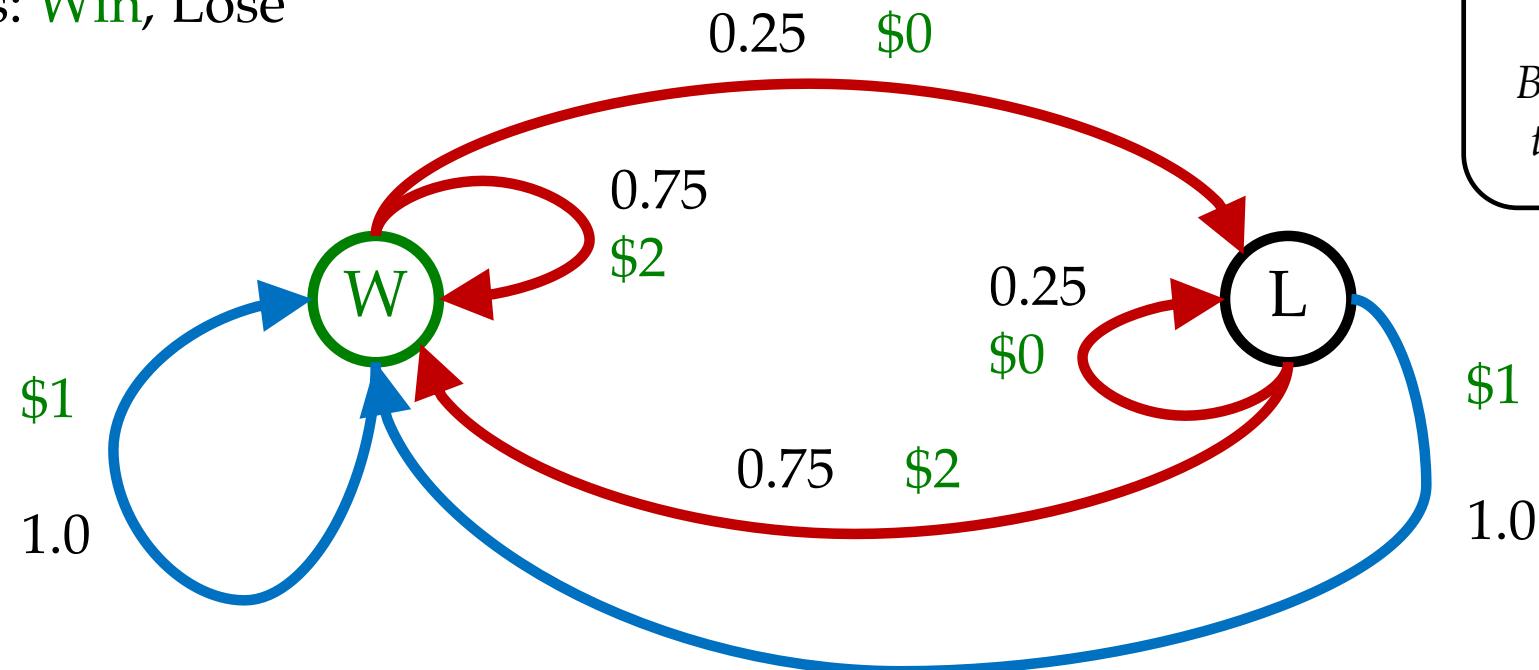
# Double Bandits

---



# Double-Bandit MDP

- Actions: *Blue, Red*
- States: Win, Lose

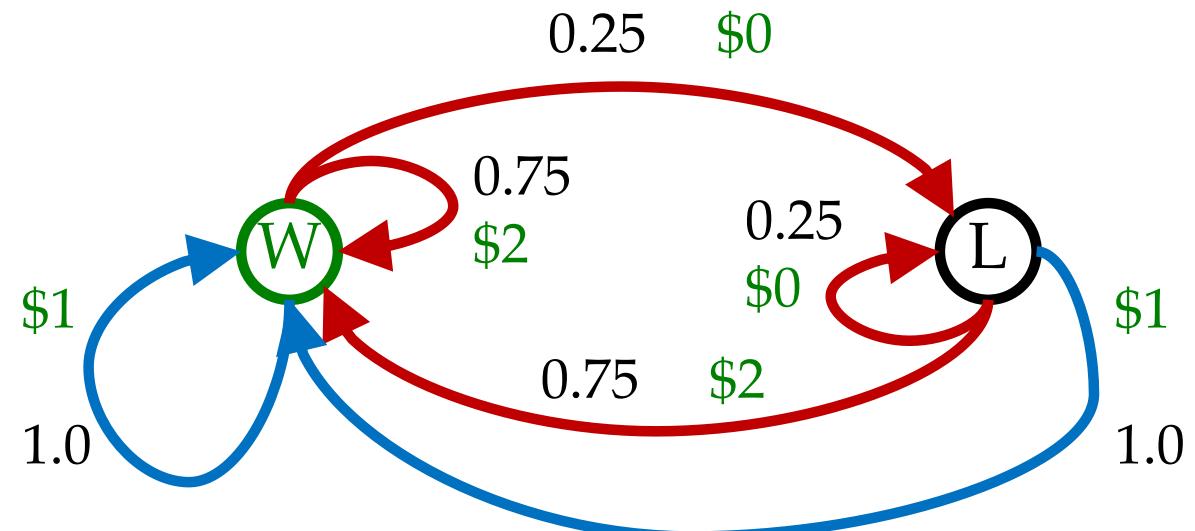


# Offline Planning

- Solving MDPs is offline planning
  - You determine all quantities through computation
  - You need to know the details of the MDP
  - You do not actually play the game!

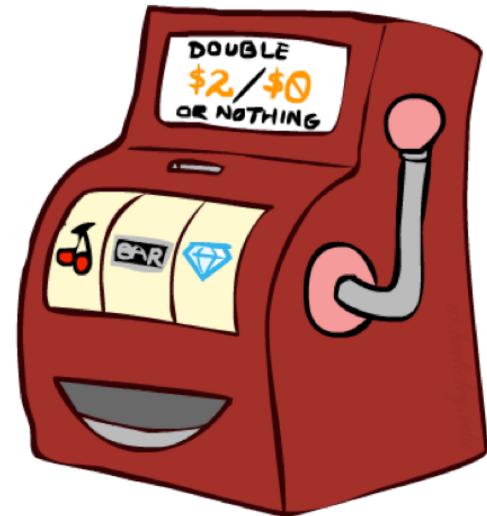
*No discount  
10 time steps  
Both states have  
the same value*

	Value
Play Red	15
Play Blue	10



# Let's Play!

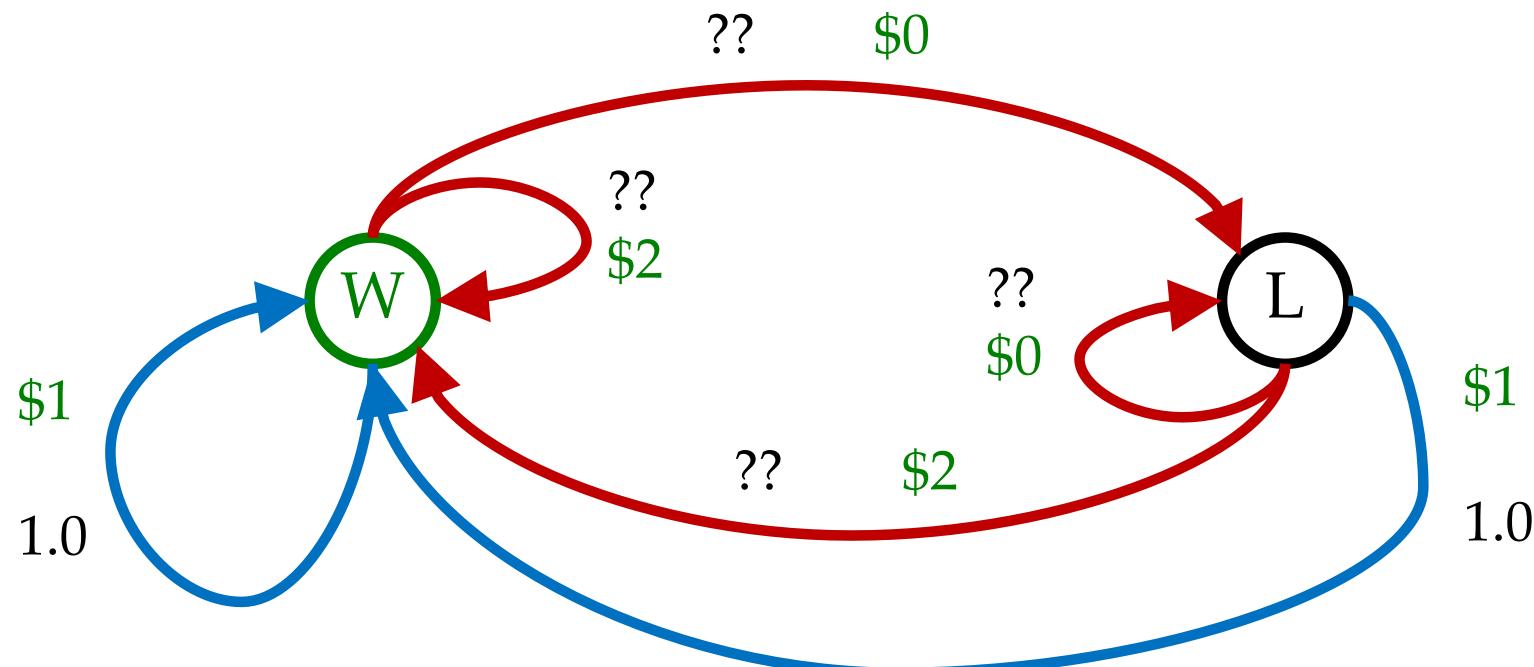
---



\$2	\$2	\$0	\$2	\$2
\$2	\$2	\$0	\$0	\$0

# Online Planning

- Rules changed! Red's win chance is different.



# Let's Play!

---



\$2 \$2 \$2 \$0 \$0  
\$2



\$0 \$0 \$0 \$0

# What Just Happened?

---

- That wasn't planning, it was learning!
  - Specifically, reinforcement learning
  - There was an MDP, but you couldn't solve it with just computation
  - You needed to actually act to figure it out
- Important ideas in reinforcement learning that came up
  - Exploration: you have to try unknown actions to get information
  - Exploitation: eventually, you have to use what you know
  - Regret: even if you learn intelligently, you make mistakes
  - Sampling: because of chance, you have to try things repeatedly
  - Difficulty: learning can be much harder than solving a known MDP



# Reinforcement Learning

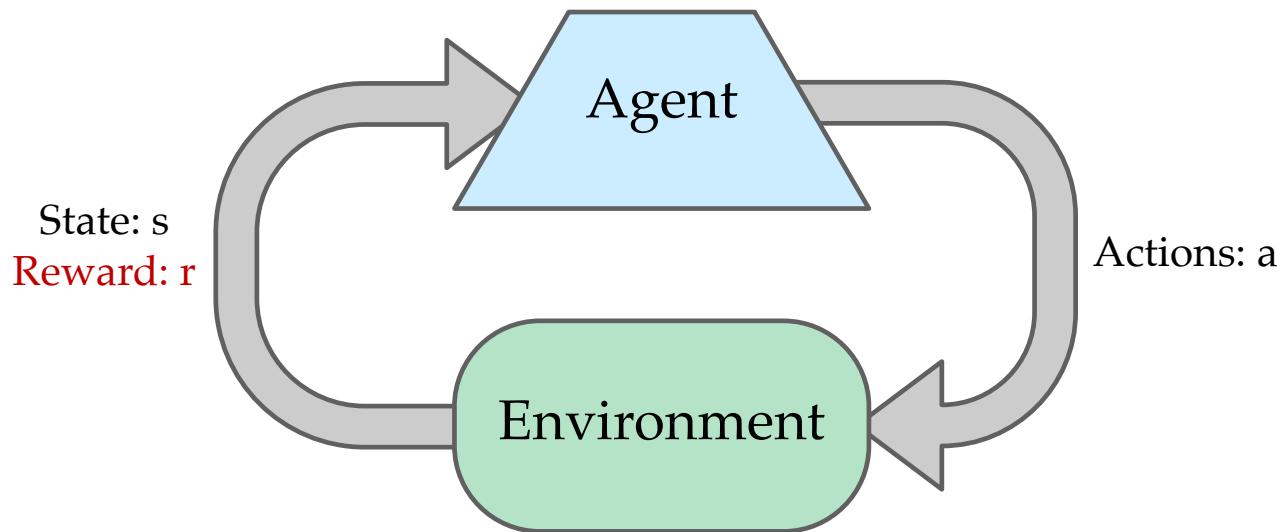
---

- Still assume a Markov decision process (MDP):
  - A set of states  $s \in S$
  - A set of actions (per state)  $A$
  - A model  $T(s,a,s')$
  - A reward function  $R(s,a,s')$
- Still looking for a policy  $\pi(s)$
- New twist: don't know  $T$  or  $R$ 
  - I.e. we don't know which states are good or what the actions do
  - Must actually try actions and states out to learn



# Reinforcement Learning

---



- Basic idea:
  - Receive feedback in the form of **rewards**
  - Agent's utility is defined by the reward function
  - Must (learn to) act so as to **maximize expected rewards**
  - All learning is based on observed samples of outcomes!

# Example: Learning to Walk

---



Initial



A Learning Trial



After Learning [1K Trials]

# Example: Learning to Walk

---



Initial

[Kohl and Stone, ICRA 2004]

[Video: AIBO WALK – initial]

# Example: Learning to Walk

---



Training

[Kohl and Stone, ICRA 2004]

[Video: AIBO WALK – training]

# Example: Learning to Walk

---



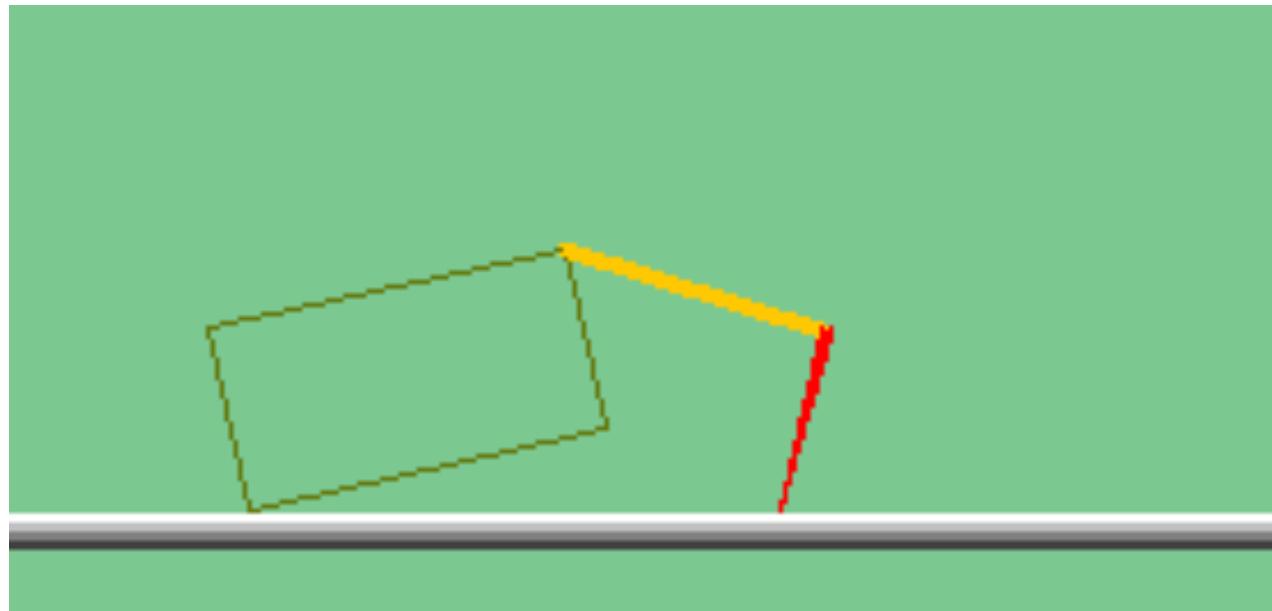
Finished

[Kohl and Stone, ICRA 2004]

[Video: AIBO WALK – finished]

# The Crawler!

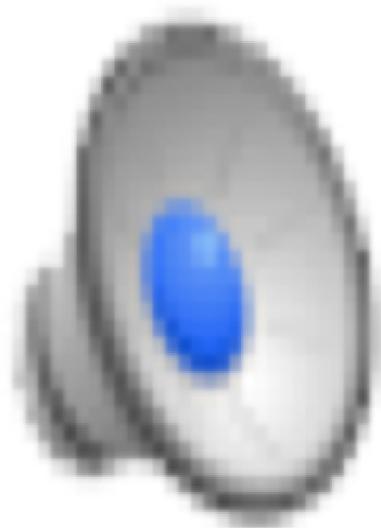
---



[Demo: Crawler Bot (L10D1)] [You, in Project 3]

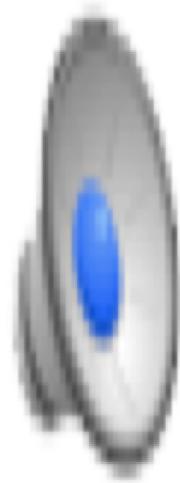
# Video of Demo Crawler Bot

---



# DeepMind Atari (©Two Minute Lectures)

---



# Reinforcement Learning

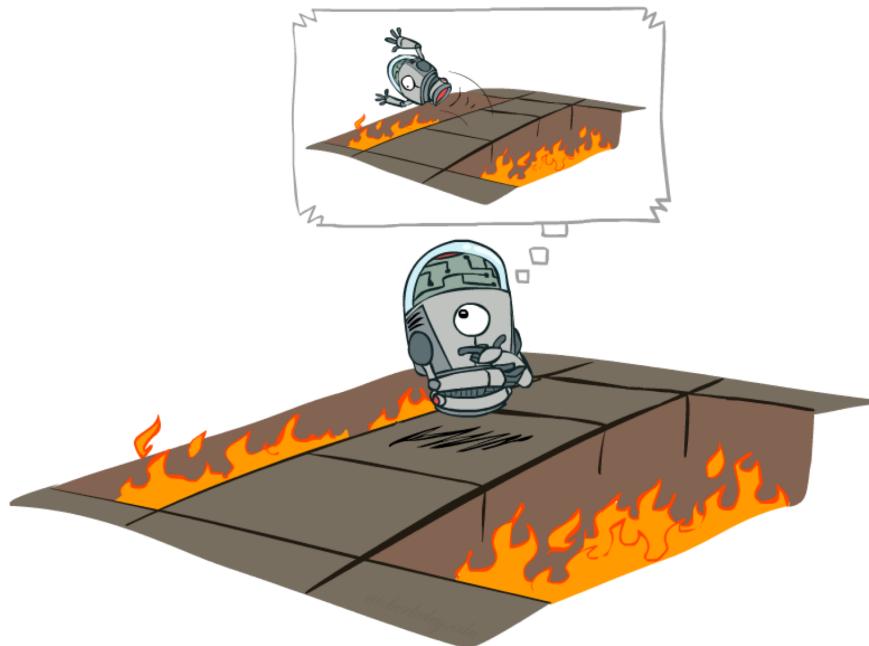
---

- Still assume a Markov decision process (MDP):
  - A set of states  $s \in S$
  - A set of actions (per state)  $A$
  - A model  $T(s,a,s')$
  - A reward function  $R(s,a,s')$
- Still looking for a policy  $\pi(s)$
- New twist: don't know  $T$  or  $R$ 
  - I.e. we don't know which states are good or what the actions do
  - Must actually try actions and states out to learn



# Offline (MDPs) vs. Online (RL)

---



Offline Solution



Online Learning