# Artificial Intelligence - INFOF311

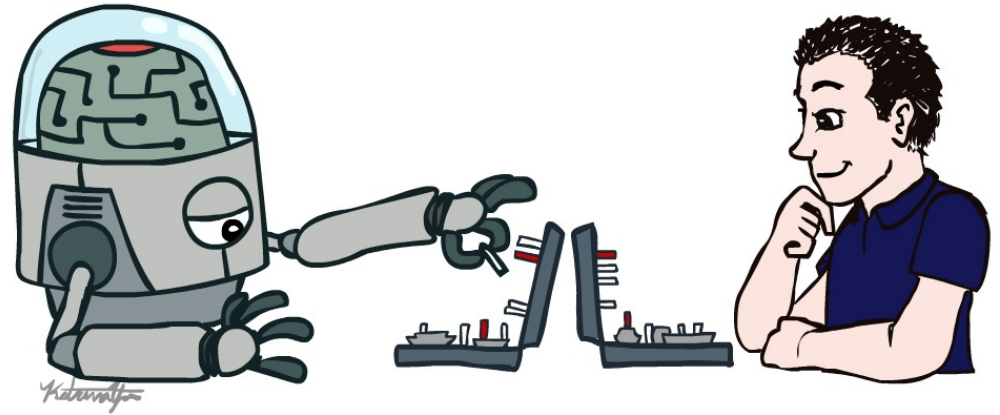**Reinforcement Learning part I**

**Instructor : Tom Lenaerts**

# Acknowledgement

We thank Stuart Russell for his generosity in allowing us to use the slide set of the UC Berkeley Course CS188, Introduction to Artificial Intelligence. These slides were created by Dan Klein, Pieter Abbeel and Anca Dragan for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.
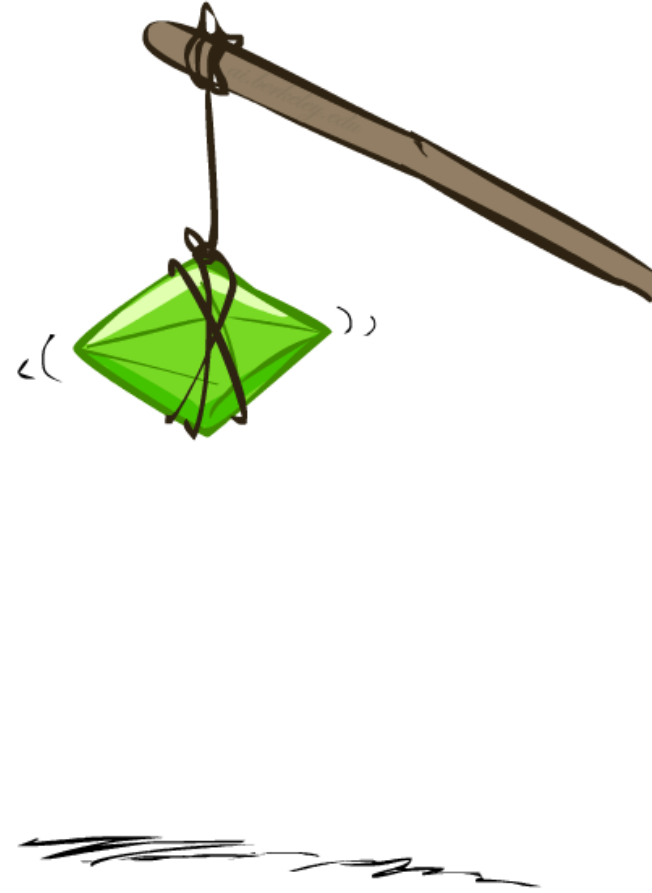
Center for
Human-Compatible
Artificial
Intelligence

The slides for INFOF311 are slightly modified versions of the slides of the spring and summer CS188 sessions in 2021 and 2022
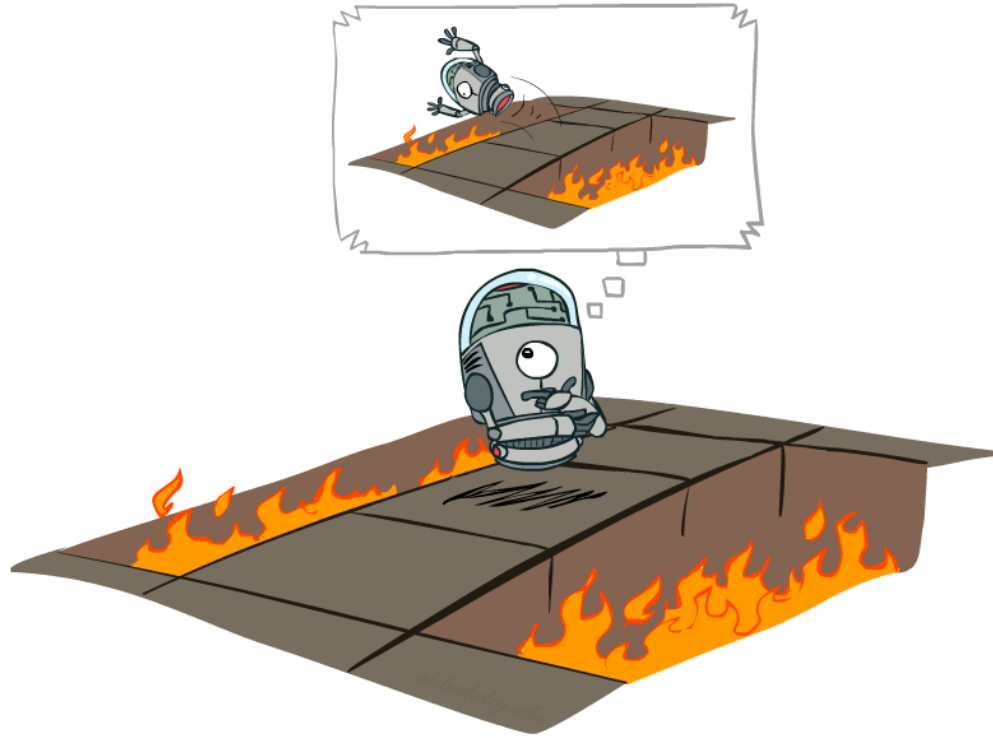
# Reinforcement Learning

# Reinforcement Learning

o Still assume a Markov decision process (MDP):

- o A set of states s ∈ S
- o A set of actions (per state) A
- o A model T(s,a,s')
- o A reward function R(s,a,s')

o Still looking for a policy π(s)

o New twist: don't know T or R

- o I.e. we don't know which states are good or what the actions do
- o Must actually try actions and states out to learn

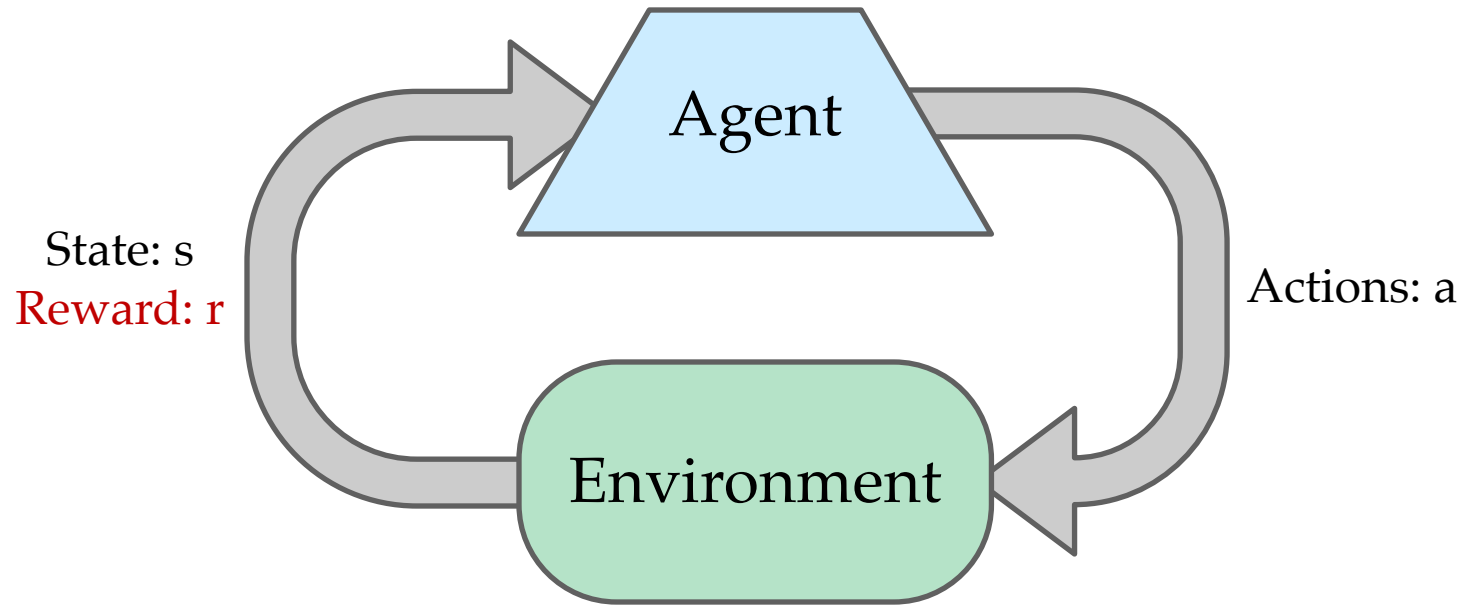# Offline (MDPs) vs. Online (RL)



Offline Solution

Online Learning

# Reinforcement Learning



State: s
Reward: r

Agent

Actions: a

Environment

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

# Example: Samuel's checker player (1956-67)

# Example: Learning to Walk



Initial            A Learning Trial            After Learning [1K Trials]

[Kohl and Stone, ICRA 2004]

# 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]

# DeepMind Atari (©Two Minute Lectures)

# Example: AlphaGo (2016)
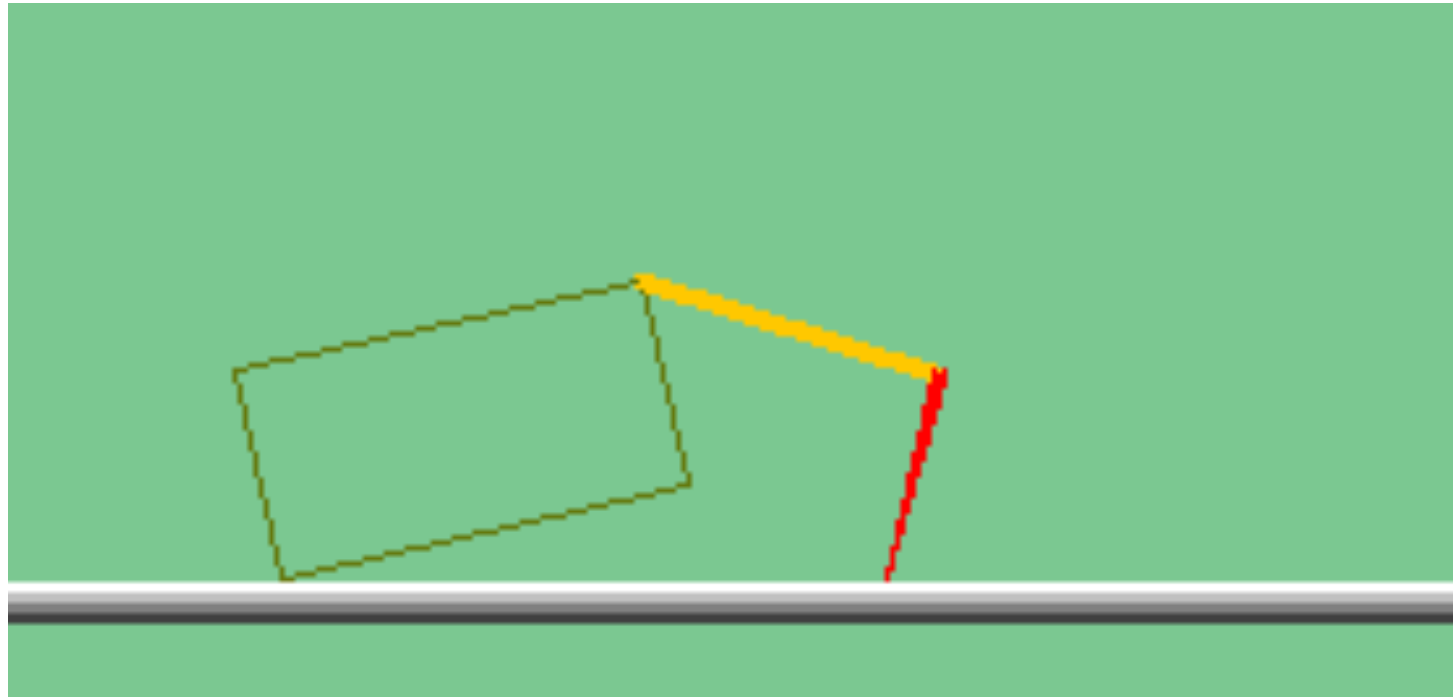
# The Crawler!

# Video of Demo Crawler Bot

# Reinforcement Learning -- Overview

o Passive Reinforcement Learning (= how to learn from experiences)

  o Model-based Passive RL

    o Learn the MDP model from experiences, then solve the MDP

  o Model-free Passive RL

    o Forego learning the MDP model, directly learn V or Q:

      o Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
      o Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)

o Active Reinforcement Learning (= agent also needs to decide how to collect experiences)

  o Key challenges:

    o How to efficiently explore?
    o How to trade off exploration <> exploitation

  o Applies to both model-based and model-free. We'll cover only in context of Q-learning
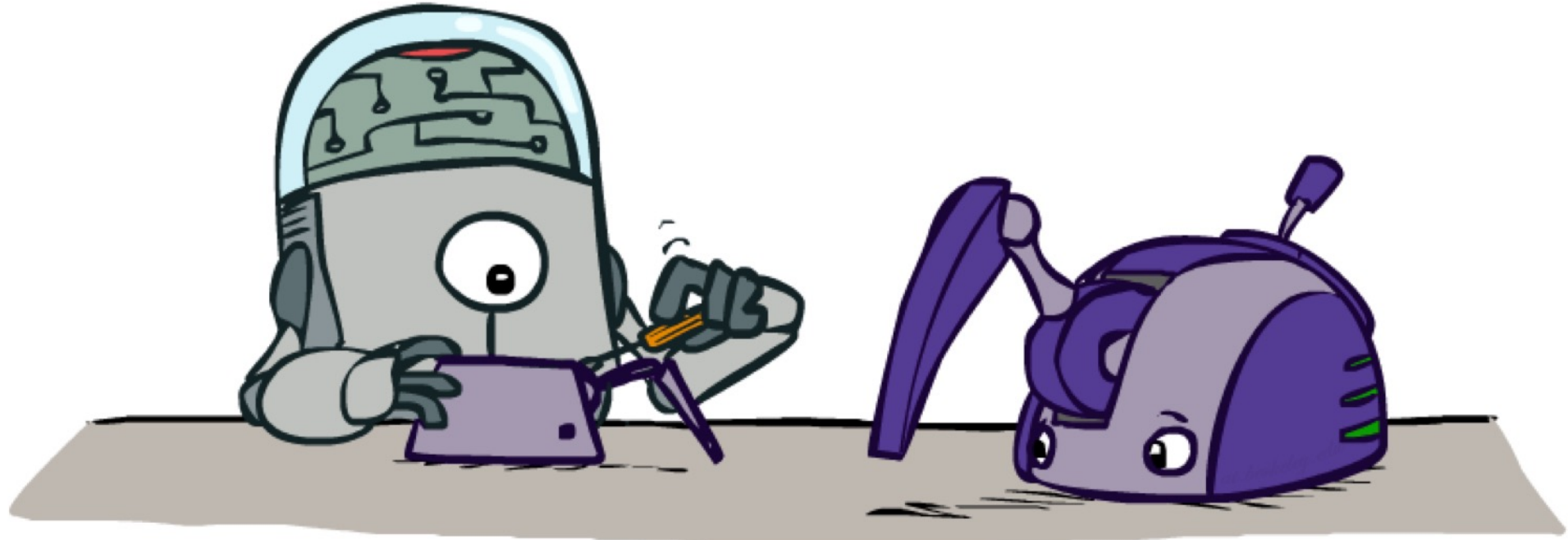
# Reinforcement Learning -- Overview

- **Passive Reinforcement Learning (= how to learn from experiences)**
  - **Model-based Passive RL**
    - Learn the MDP model from experiences, then solve the MDP
  - Model-free Passive RL
    - Forego learning the MDP model, directly learn V or Q:
      - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
      - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
  - Key challenges:
    - How to efficiently explore?
    - How to trade off exploration <> exploitation
  - Applies to both model-based and model-free. We'll cover only in context of Q-learning

# Model-Based Reinforcement Learning
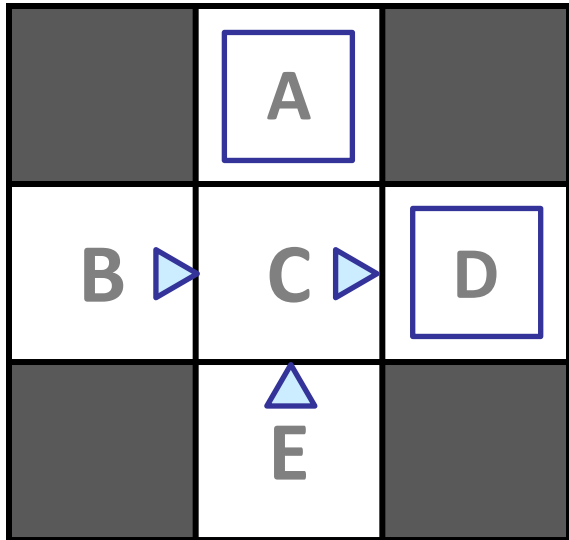
# Model-Based Reinforcement Learning

o Model-Based Idea:
  o Learn an approximate model based on experiences
  o Solve for values as if the learned model were correct

o Step 1: Learn empirical MDP model
  o Count outcomes s' for each s, a
  o Normalize to give an estimate $\hat{T}(s, a, s')$
  o Discover each $\hat{R}(s, a, s')$ when we experience (s, a, s')

o Step 2: Solve the learned MDP
  o For example, use value iteration, as before

(and repeat as needed)

# Example: Model-Based RL

## Input Policy π



*Assume:* γ = 1

## Observed Episodes (Training)

### Episode 1

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 2

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 3

E, north, C, -1
C, east,   D, -1
D, exit,    x, +10

### Episode 4

E, north, C, -1
C, east,   A, -1
A, exit,    x, -10

## Learned Model

$\hat{T}(s, a, s')$

T(B, east, C) = 1.00
T(C, east, D) = 0.75
T(C, east, A) = 0.25
...

$\hat{R}(s, a, s')$

R(B, east, C) = -1
R(C, east, D) = -1
R(D, exit, x) = +10
...

# Analogy: Expected Age

Goal: Compute expected age of cs188 students

**Known P(A)**

$$E[A] = \sum_a P(a) \cdot a \qquad = 0.35 \times 20 + \ldots$$

Without P(A), instead collect samples $[a_1, a_2, \ldots a_N]$

**Unknown P(A): "Model Based"**

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$

$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

Why does this work? Because eventually you learn the right model.

**Unknown P(A): "Model Free"**

$$E[A] \approx \frac{1}{N} \sum_i a_i$$

Why does this work? Because samples appear with the right frequencies.
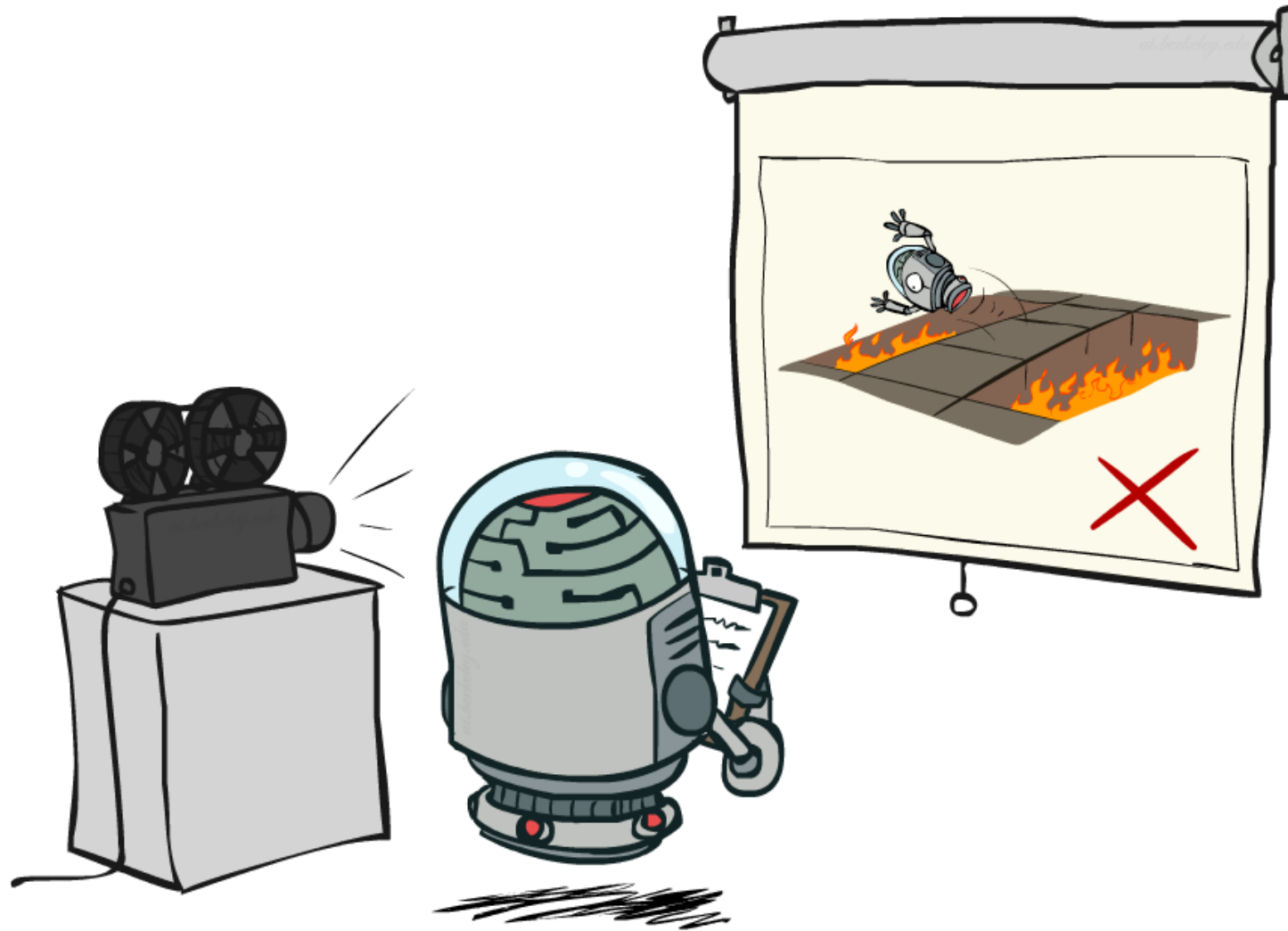
# Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
  - Model-based Passive RL
    - Learn the MDP model from experiences, then solve the MDP
  - **Model-free Passive RL**
    - **Forego learning the MDP model, directly learn V or Q:**
      - **Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning**
      - Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
  - Key challenges:
    - How to efficiently explore?
    - How to trade off exploration <> exploitation
  - Applies to both model-based and model-free. In INFOF311 we'll cover only in context of Q-learning
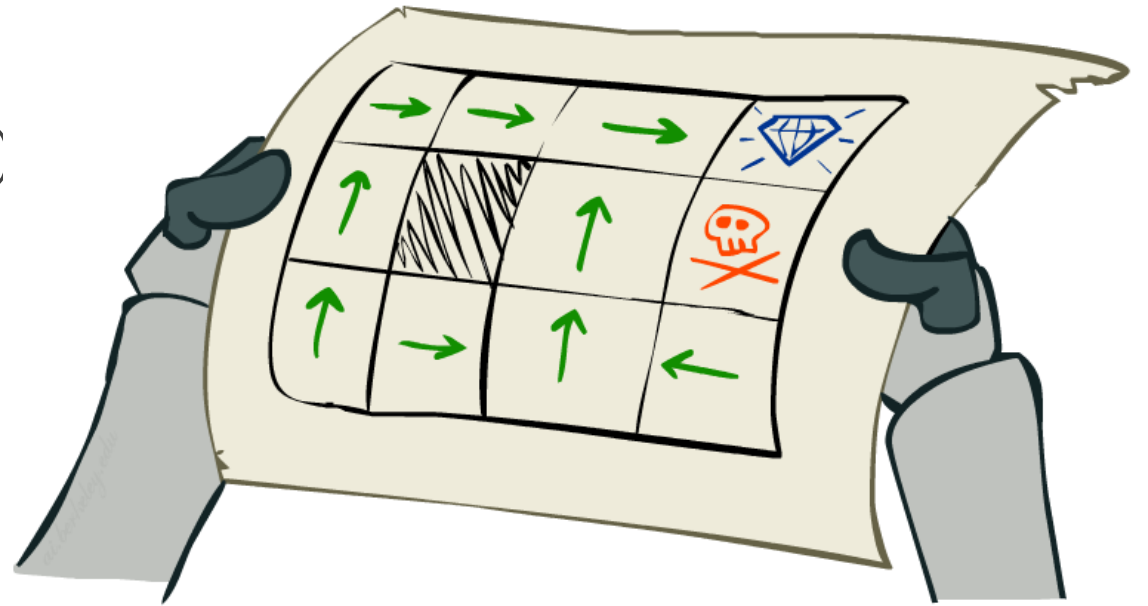
# Passive Model-Free Reinforcement Learning

# Passive Model-Free Reinforcement Learning

o Simplified task: policy evaluation
  o Input: a fixed policy $\pi(s)$
  o You don't know the transitions $T(s,a,s')$
  o You don't know the rewards $R(s,a,s')$
  o Goal: learn the state values

o In this case:
  o Learner is "along for the ride"
  o No choice about what actions to take
  o Just execute the policy and learn from experience
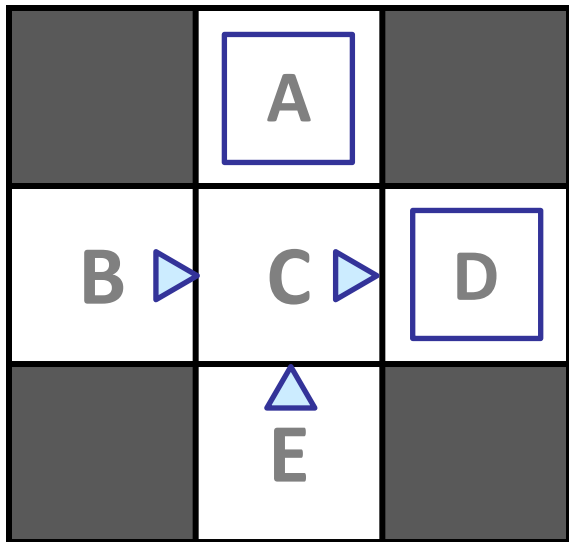  o This is NOT offline planning!  You actually take actions in the world.

# Direct Evaluation

o Goal: Compute values for each state under $\pi$

o Idea: Average together observed sample values

   o Act according to $\pi$

   o Every time you visit a state, write down what the sum of discounted rewards turned out to be

   o Average those samples

o This is called direct (or Monte Carlo) evaluation

# Example: Direct Evaluation

## Input Policy π



*Assume:* $\gamma = 1$

## Observed Episodes (Training)

### Episode 1

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 2

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 3

E, north, C, -1
C, east,   D, -1
D, exit,    x, +10

### Episode 4

E, north, C, -1
C, east,   A, -1
A, exit,    x, -10

## Output Values



*If B and E both go to C under this policy, how can their values be different?*

# Problems with Direct Evaluation

o What's good about direct evaluation?

   o It's easy to understand

   o It doesn't require any knowledge of T, R

   o It eventually computes the correct average values, using just sample transitions

o What bad about it?

   o It wastes information about state connections

   o Each state must be learned separately

   o So, it takes a long time to learn

## Output Values



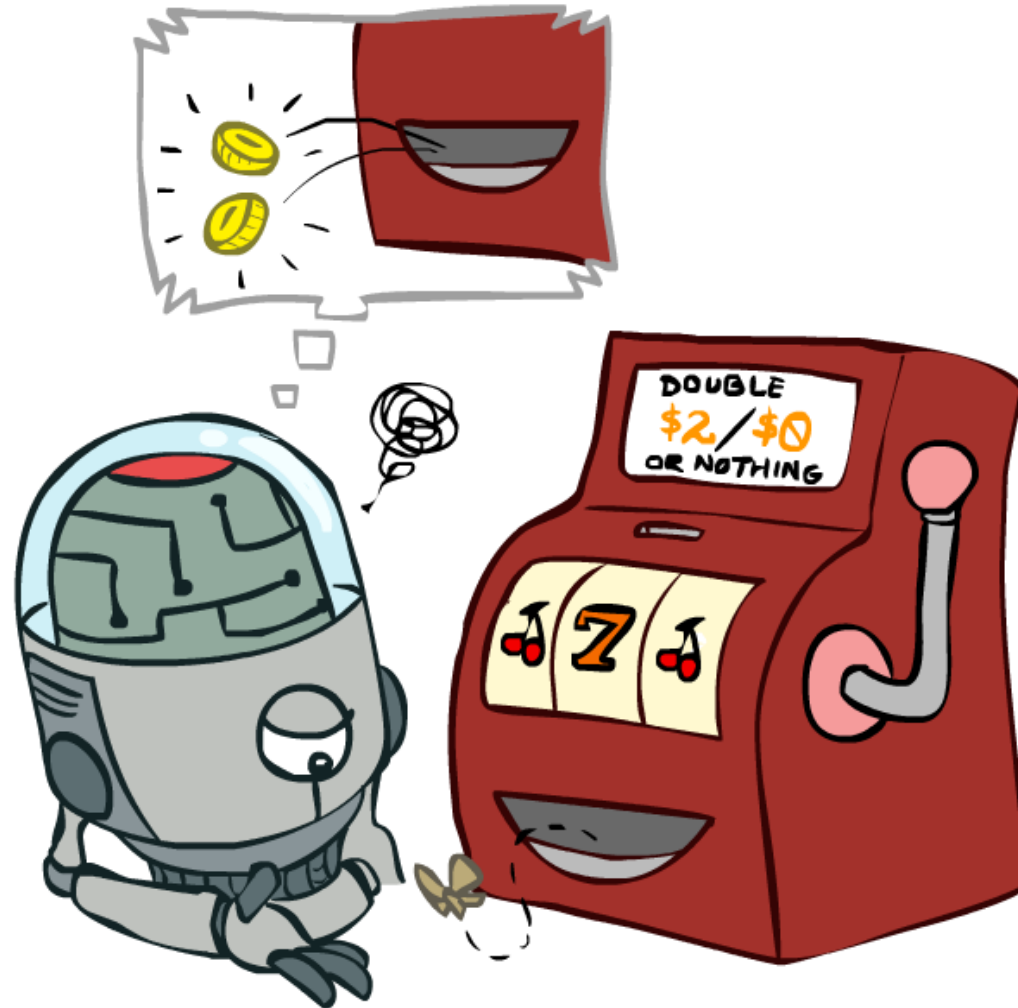*If B and E both go to C under this policy, how can their values be different?*

# Reinforcement Learning -- Overview

o Passive Reinforcement Learning (= how to learn from experiences)
  o Model-based Passive RL
    o Learn the MDP model from experiences, then solve the MDP
  o **Model-free Passive RL**
    o **Forego learning the MDP model, directly learn V or Q:**
      o **Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning**
      o Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)

o Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
  o Key challenges:
    o How to efficiently explore?
    o How to trade off exploration <> exploitation
  o Applies to both model-based and model-free. We'll cover only in context of Q-learning
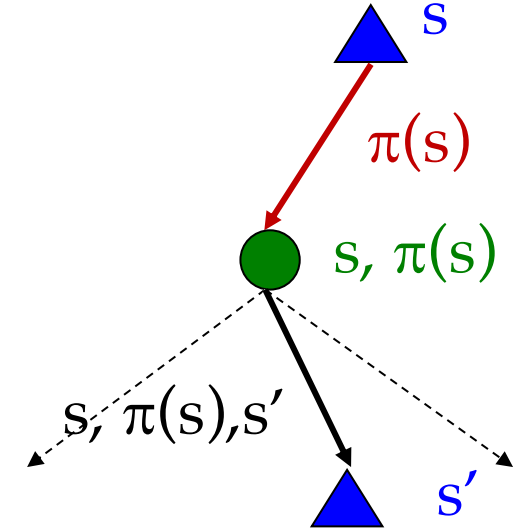
# Temporal Difference Value Learning

# Why Not Use Policy Evaluation?

o **Simplified Bellman updates calculate V for a fixed policy:**
  o Each round, replace V with a one-step-look-ahead layer over V

$$V_0^\pi(s) = 0$$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

  o This approach fully exploited the connections between the states
  o Unfortunately, we need T and R to do it!

o **Key question: how can we do this update to V without knowing T and R?**
  o In other words, how to we take a weighted average without knowing the weights?

s

$\pi(s)$

s, $\pi(s)$

s, $\pi(s)$,s'

s'

# Sample-Based Policy Evaluation?

○ We want to improve our estimate of V by computing these averages:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

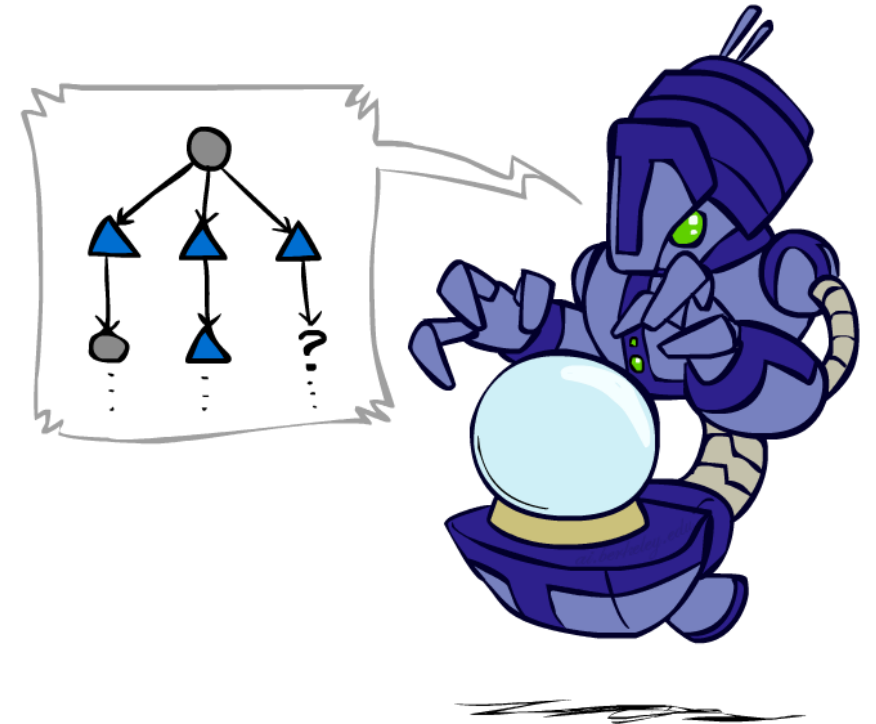○ Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_1 = R(s, \pi(s), s_1') + \gamma V_k^{\pi}(s_1')$$

$$sample_2 = R(s, \pi(s), s_2') + \gamma V_k^{\pi}(s_2')$$
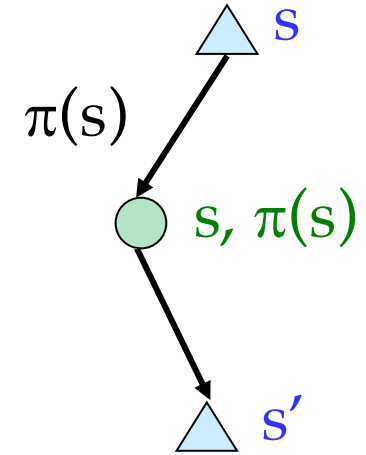
$$\dots$$

$$sample_n = R(s, \pi(s), s_n') + \gamma V_k^{\pi}(s_n')$$

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_{i} sample_i$$

# Temporal Difference Value Learning

o Big idea: learn from every experience as we move forward!

    o Update V(s) each time we experience a transition (s, a, s', r)

    o Likely outcomes s' will contribute updates more often

o Temporal difference learning of values

    o Policy still fixed, still doing evaluation!

    o Move values toward value of whatever successor occurs: running average

Sample of V(s):     $sample = R(s, \pi(s), s') + \gamma V^{\pi}(s')$

Update to V(s):     $V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + (\alpha)sample$

Same update:     $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$

s

$\pi(s)$

s, $\pi(s)$

s'

# Exponential Moving Average

o Exponential moving average

   o The running interpolation update: $\quad \bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$

   o Makes recent samples more important

   o Forgets about the past (distant past values were wrong anyway)

o Decreasing learning rate (alpha) can give converging averages

# Example: Temporal Difference Value Learning

## States



Assume: $\gamma = 1$, $\alpha = 1/2$

## Observed Transitions

B, east, C, -2

C, east, D, -2

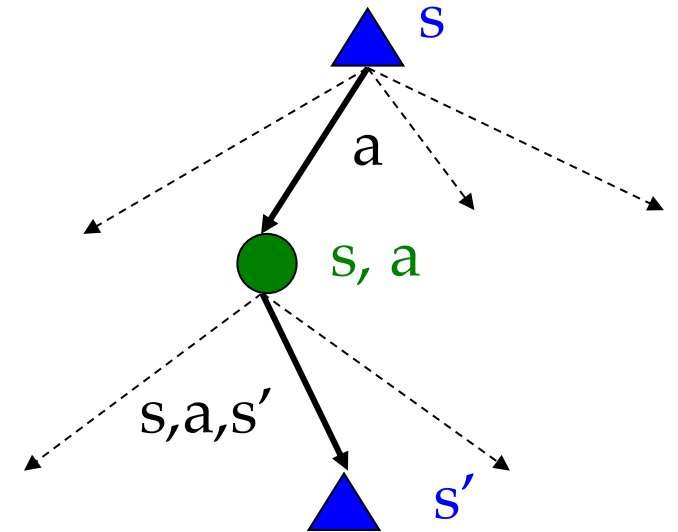$$V^\pi(s) \leftarrow (1-\alpha)V^\pi(s) + \alpha \Big[ R(s, \pi(s), s') + \gamma V^\pi(s') \Big]$$

# Problems with TD Value Learning

o TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages

o However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg\max_a Q(s, a)$$

$$Q(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V(s') \right]$$

o Idea: learn Q-values, not values

o Makes action selection model-free too!

# Reinforcement Learning -- Overview

- Passive Reinforcement Learning (= how to learn from experiences)
  - Model-based Passive RL
    - Learn the MDP model from experiences, then solve the MDP
  - **Model-free Passive RL**
    - Forego learning the MDP model, directly learn V or Q:
      - Value learning – learns value of a fixed policy; 2 approaches: Direct Evaluation & TD Learning
      - **Q learning – learns Q values of the optimal policy (uses a Q version of TD Learning)**
- Active Reinforcement Learning (= agent also needs to decide how to collect experiences)
  - Key challenges:
    - How to efficiently explore?
    - How to trade off exploration <> exploitation
  - Applies to both model-based and model-free. We'll cover only in context of Q-learning

# Q-Value Iteration

o Value iteration: find successive (depth-limited) values
- o Start with $V_0(s) = 0$, which we know is right
- o Given $V_k$, calculate the depth k+1 values for all states:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$

o But Q-values are more useful, so compute them instead
- o Start with $Q_0(s,a) = 0$, which we know is right
- o Given $Q_k$, calculate the depth k+1 q-values for all q-states:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

# Q-Learning

o Q-Learning: sample-based Q-value iteration
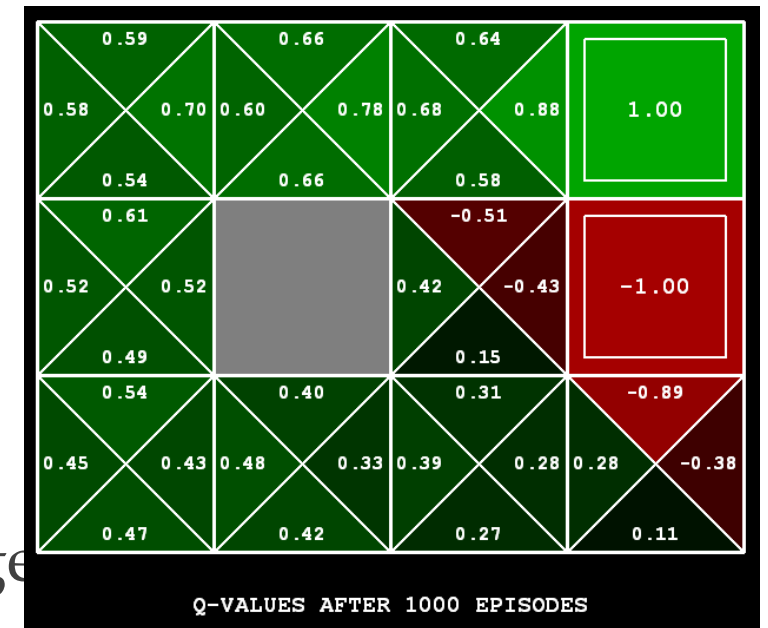
$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

o Learn Q(s,a) values as you go

  o Receive a sample (s,a,s',r)
  o Consider your old estimate: $Q(s,a)$
  o Consider your new sample estimate:

  $$sample = R(s,a,s') + \gamma \max_{a'} Q(s',a')$$

  no longer policy evaluation!

  o Incorporate the new estimate into a running average

  $$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha)[sample]$$
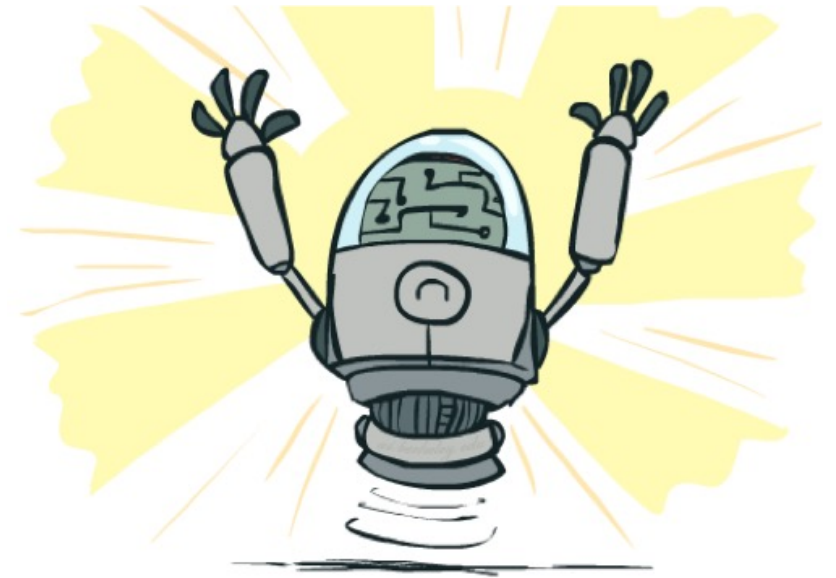


Q-VALUES AFTER 1000 EPISODES

[Demo: Q-learning – gridworld (L10D2)]

[Demo: Q-learning – crawler (L10D3)]

# Q-Learning Properties

o Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!

o This is called off-policy learning

o Caveats:
  o You have to explore enough
  o You have to eventually make the learning rate small enough
  o … but not decrease it too quickly
  o Basically, in the limit, it doesn't matter how you select actions (!)

# Video of Demo Q-Learning -- Gridworld

# Video of Demo Q-Learning -- Crawler