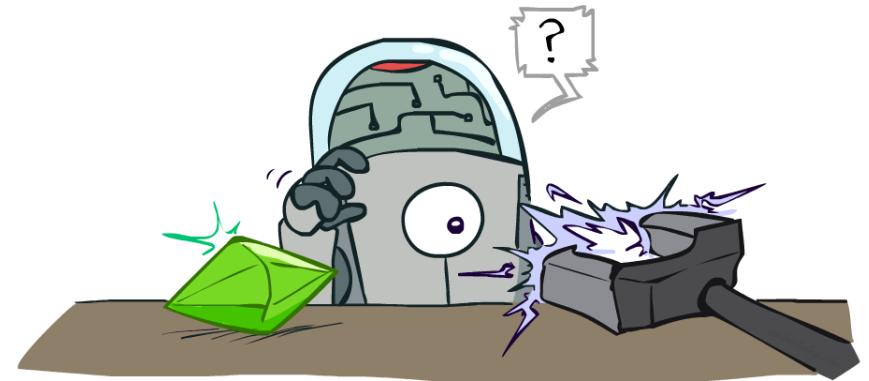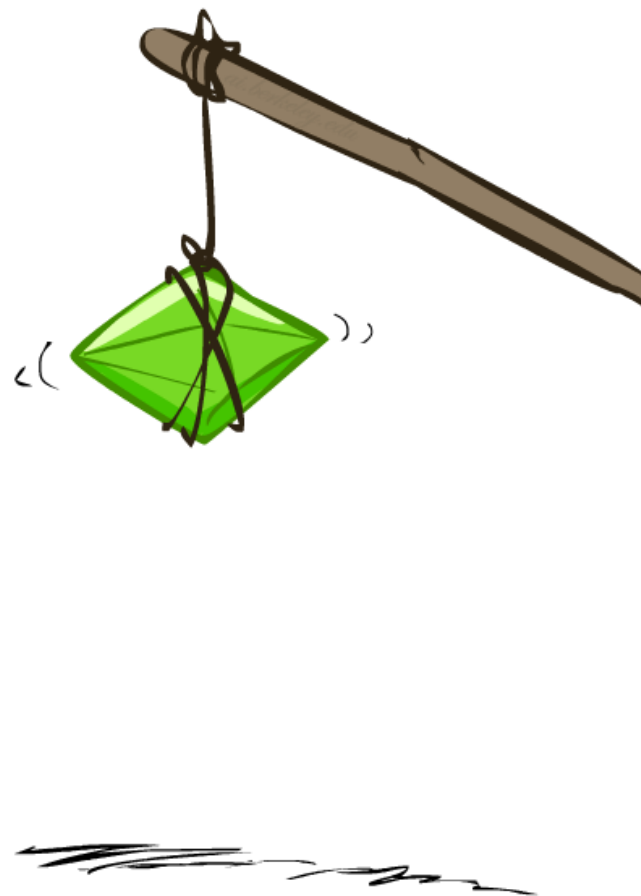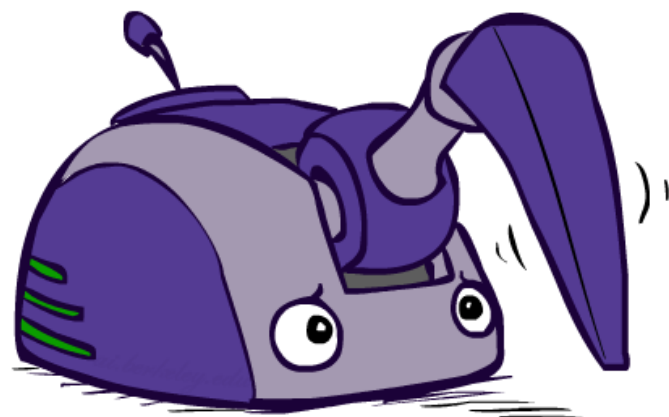# CSE 573:
# Artificial Intelligence
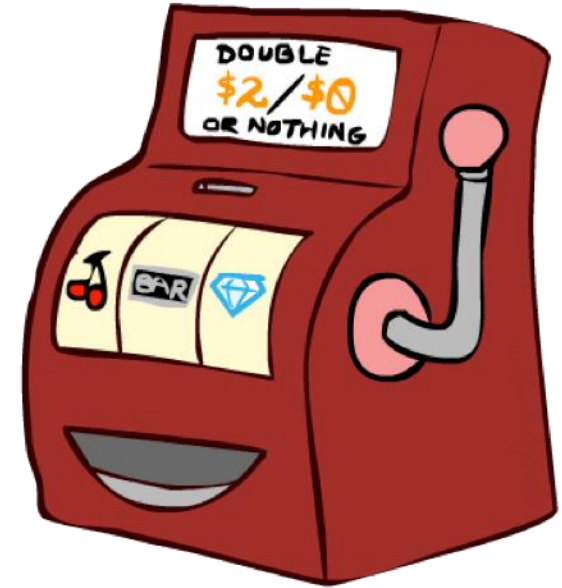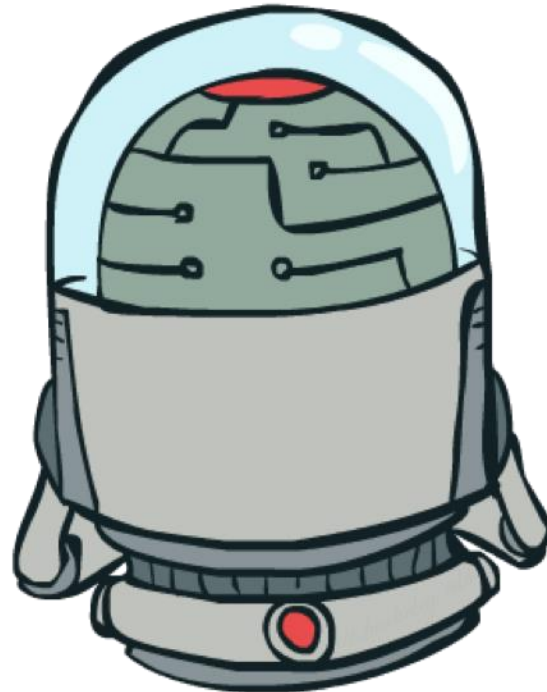
## Hanna Hajishirzi

## Reinforcement Learning

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer

# Reinforcement Learning

# Double Bandits

# Double-Bandit MDP

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

No discount
10 time steps
Both states
have the same
value

# Offline Planning

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

*No discount*

*10 time steps*

| | Value |
|---|---|
| Play Red | 15 |
| Play Blue | 10 |

# Let's Play!



$2  $2  $0  $2  $2

$2  $2  $0  $0  $0

# Online Planning

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

# Let's Play!

$0  $0  $2  $0

$0  $2  $2  $0  $0

$0

# What Just Happened?

o That wasn't planning, it was learning!
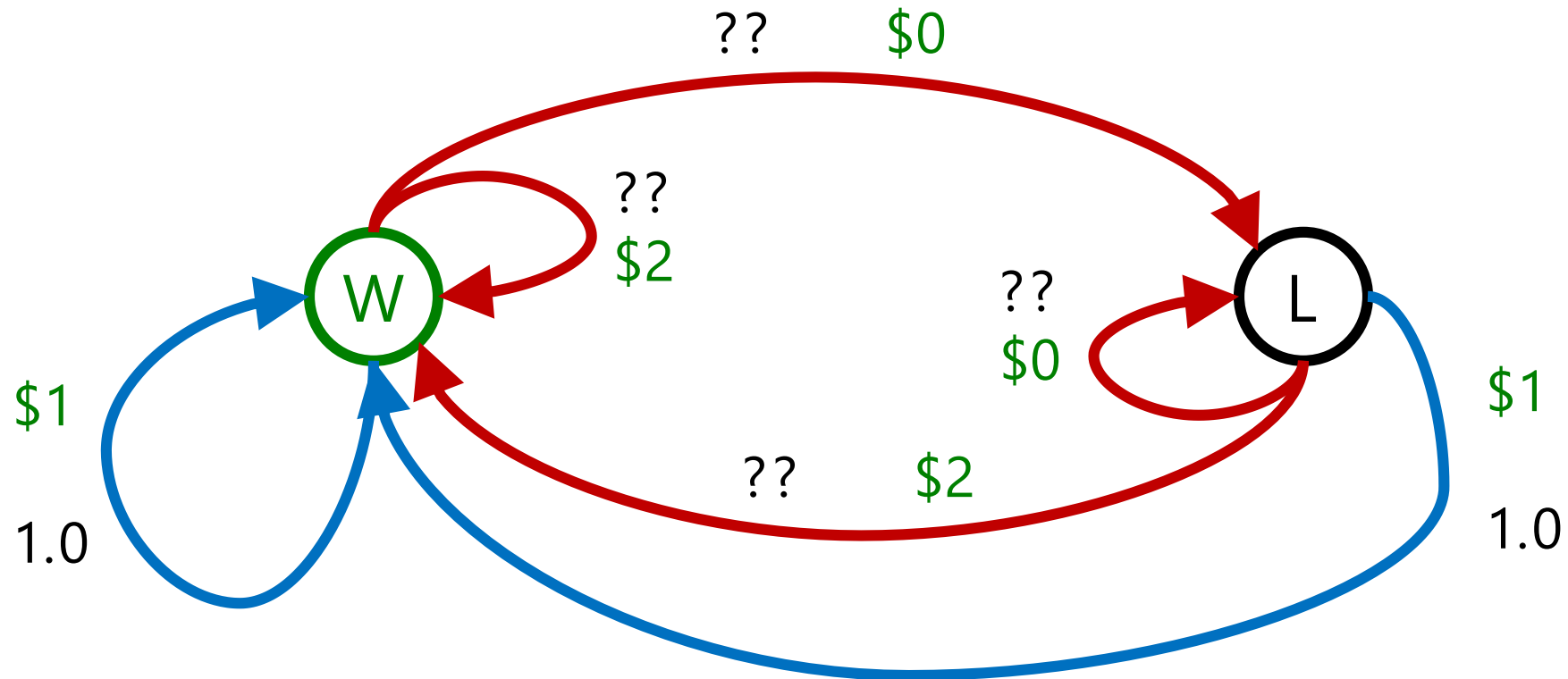  o Specifically, reinforcement learning
  o There was an MDP, but you couldn't solve it with just computation
  o You needed to actually act to figure it out

o Important ideas in reinforcement learning that came up
  o Exploration: you have to try unknown actions to get information
  o Exploitation: eventually, you have to use what you know
  o Regret: even if you learn intelligently, you make mistakes
  o Sampling: because of chance, you have to try things repeatedly
  o Difficulty: learning can be much harder than solving a known MDP

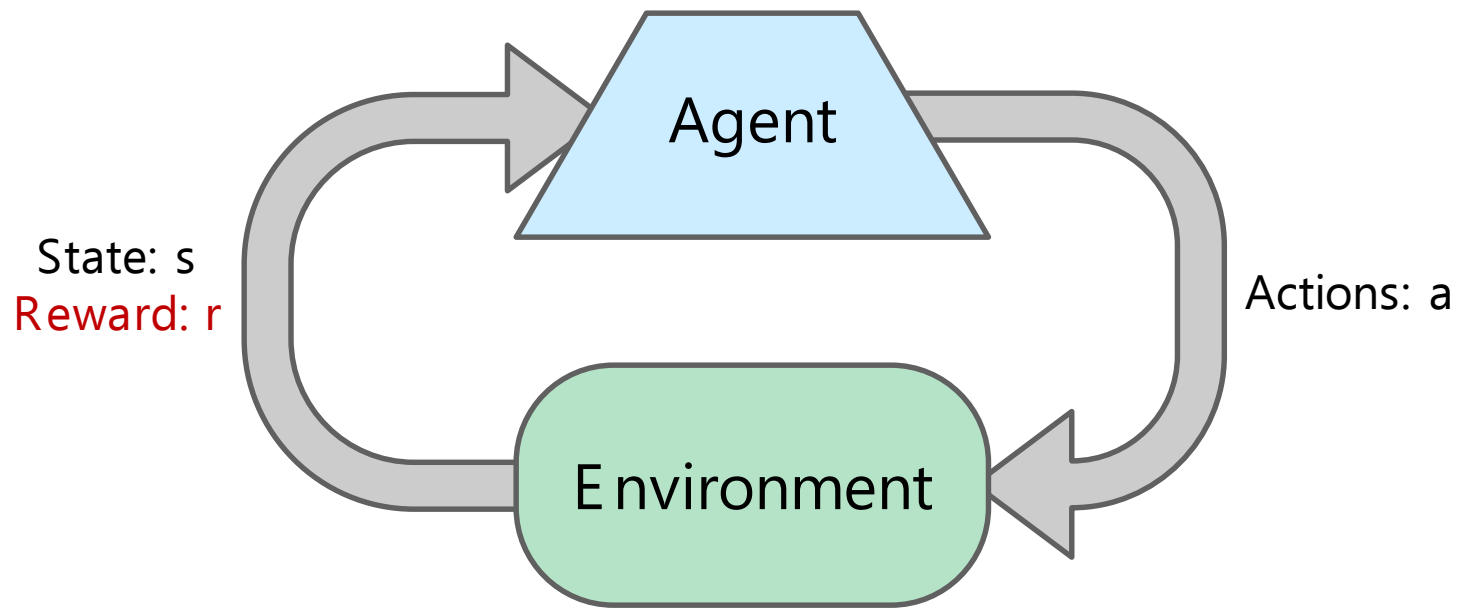# Reinforcement Learning

- 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

Cool

Warm

Overheated

# Reinforcement Learning



State: s
Reward: r

Actions: a

Agent

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: Learning to Walk



Initial           A Learning Trial           After Learning [1K Trials]

[Kohl and Stone, ICRA 2004]

# Example: Toddler Robot



[Tedrake, Zhang and Seung, 2005]                    [Video: TODDLER – 40s]

# Robotics Rubik Cube

○ https://www.youtube.com/watch?v=x4O8pojMF0w



Solving Rubik's Cube with a Robot Hand

# CSE 573:
# Artificial Intelligence

## Hanna Hajishirzi

## Reinforcement Learning

# Announcements

o PS2: April 29

o Project proposals: May 6th

o Paper review: May 13

# Project Proposal

o Project proposals: May 6th

  o Pick projects close to you interest, or select from here: [list of potential projects](). Your final project can also be a re-implementation of one of the recent papers from AI/ML/NLP/Computer vision conferences.

o The project proposal is a 1-page summary of the project topic, motivation, definition, dataset, and resources. It should also include the milestones, detailed experiment plan, and the timeline to complete each milestone.
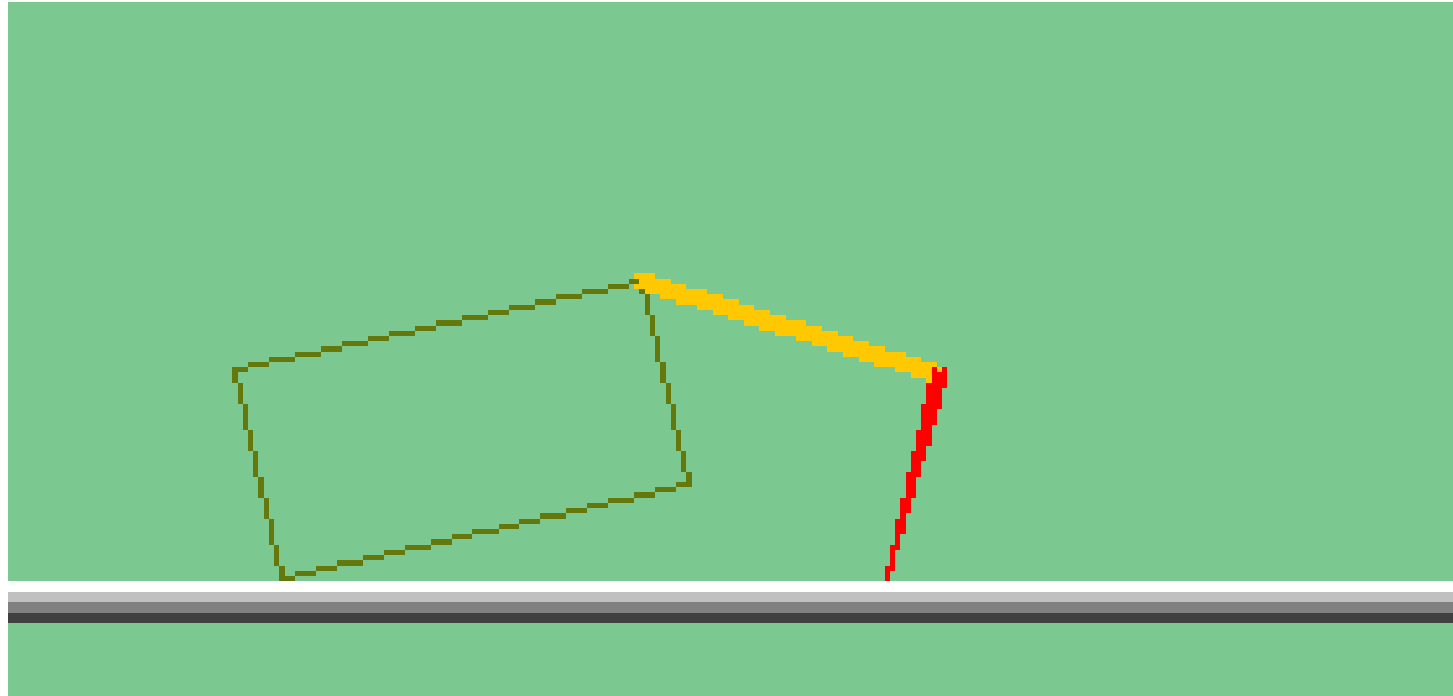
# Paper Review

o Paper review:

  o 1. Describe what problem or question this paper addresses, and the main contributions that it makes towards a solution or answer. a. Problem/Question:

  o b. Solution/approach:

  o c. Contributions (list at least two):

o 2. Evaluate the paper in terms of novelty, significance, and empirical results. 3. Describe the main strengths you see in the paper. 4. Describe critiques and weaknesses you see in the paper.
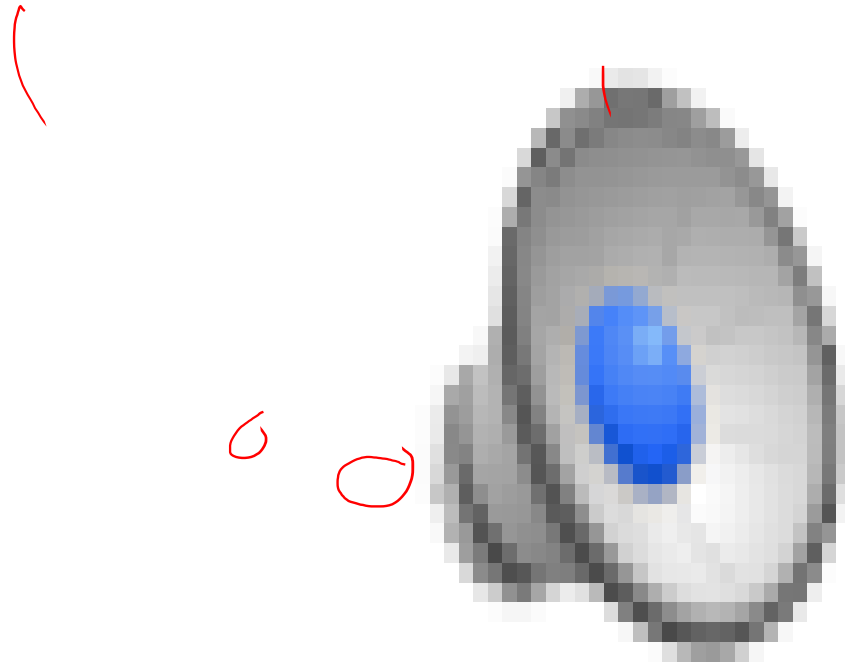
# Reinforcement Learning

o Still assume a Markov decision process (MDP):
   o A set of states s $\in$ 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 $\pi$(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
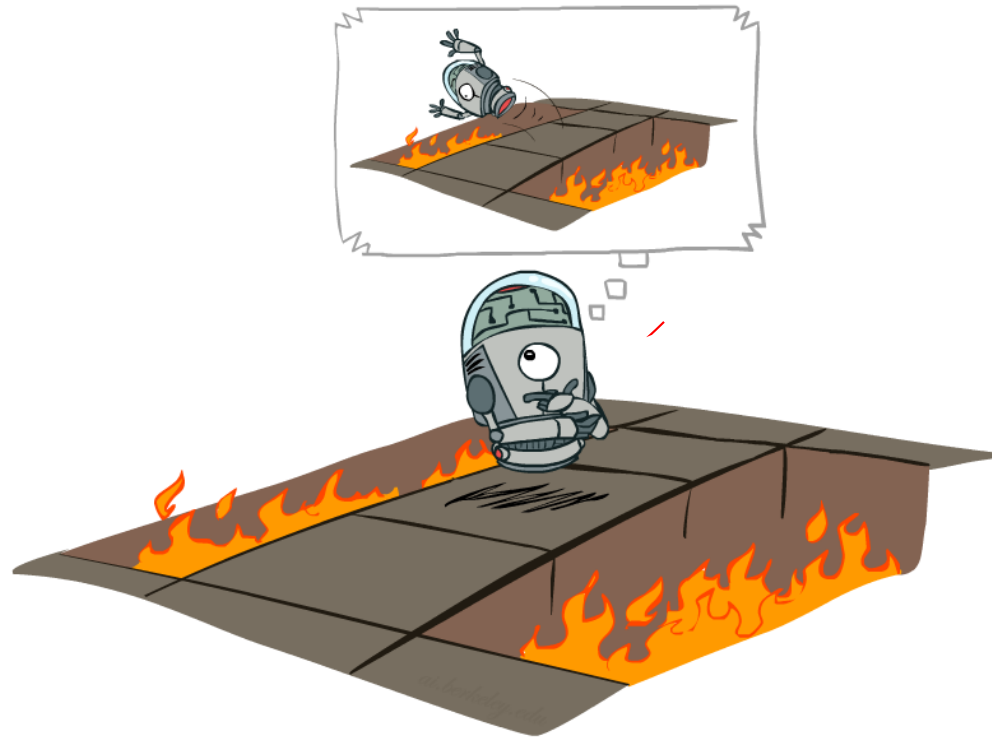   o Must actually try actions and states out to learn

Cool

Warm

Overheated

# The Crawler!

# Video of Demo Crawler Bot
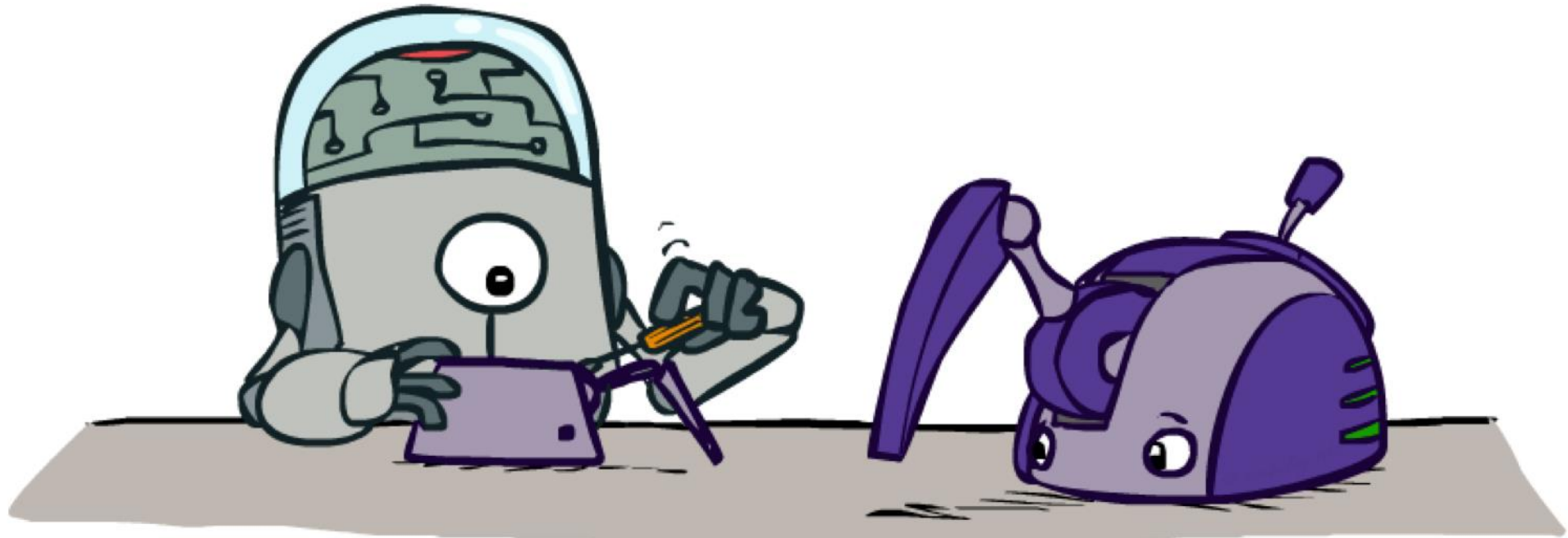
# Offline (MDPs) vs. Online (RL)



Offline Solution

Online Learning

# Model-Based Learning

# Model-Based Learning

$T_j \beta$

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

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

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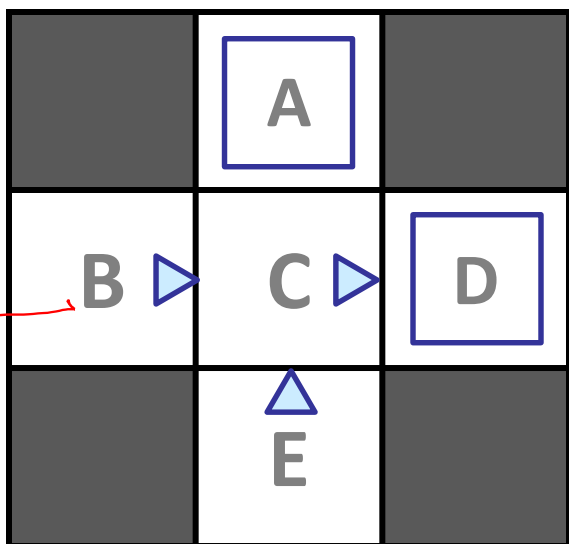  - For example, use value iteration, as before

$R(s, a, s')$

# Example: Model-Based Learning

$T(B, east, -)$
$(C, east, -)$  3/4 D

$T(s, a, s')$

## Input Policy π



*Assume: γ = 1*

## Observed Episodes (Training)

1/4 A

### 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

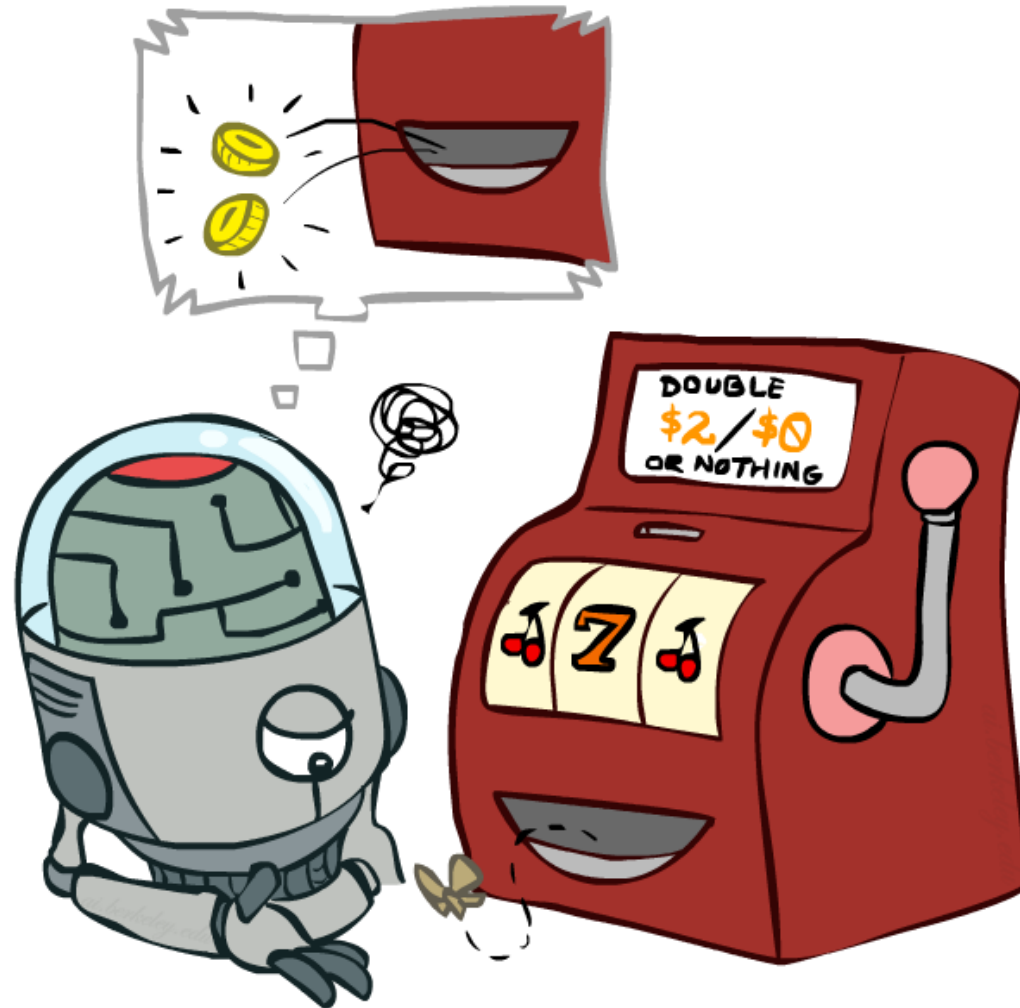E, north, C, -1
C, east,    A, -1
A, exit,    x, -10

## R Learned Model

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

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

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

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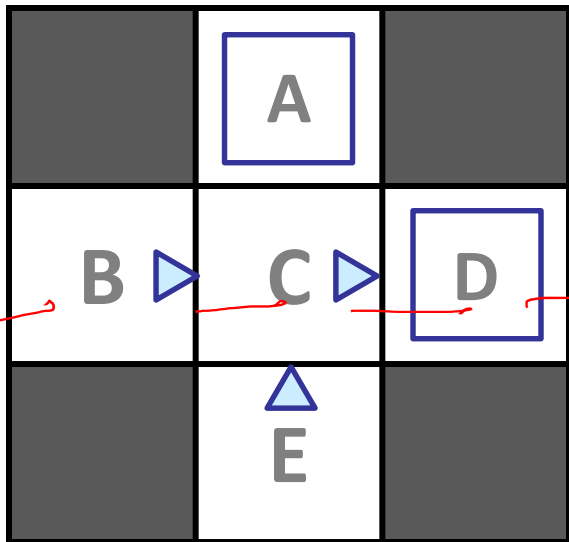R(D, exit, x) = +10
...

# Model-Free Learning

# Direct Evaluation

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

○ Idea: Average together observed sample values

    ○ Act according to $\pi$

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

    ○ Average those samples

○ This is called direct evaluation
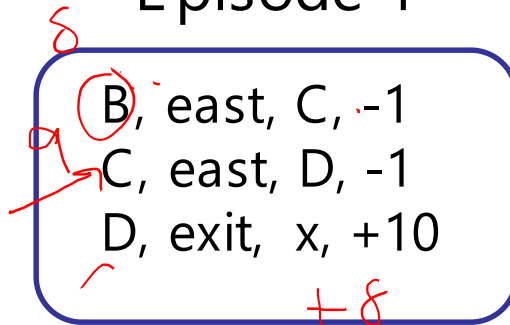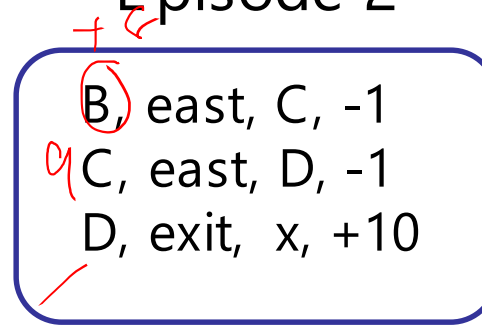
# Example: Direct Evaluation

## Input Policy π



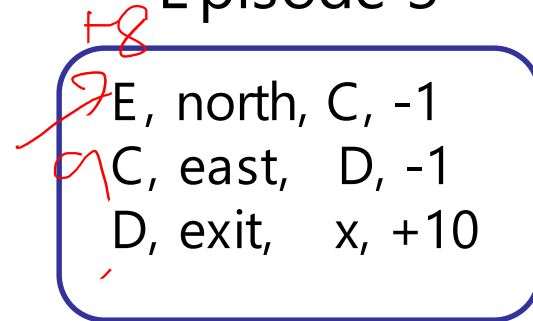Assume: γ = 1

## 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

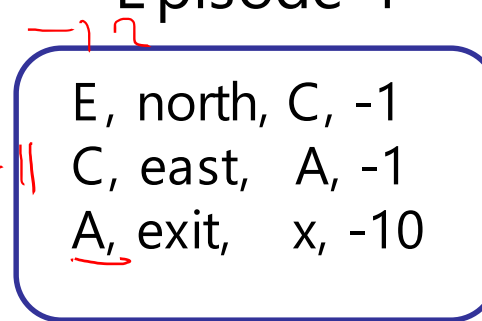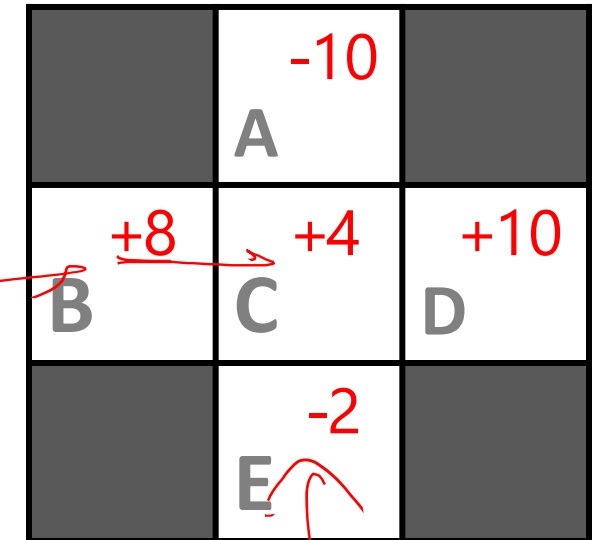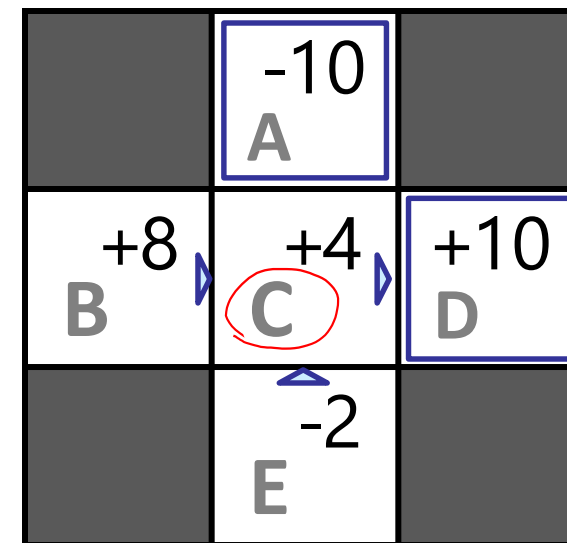## Output Values



*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
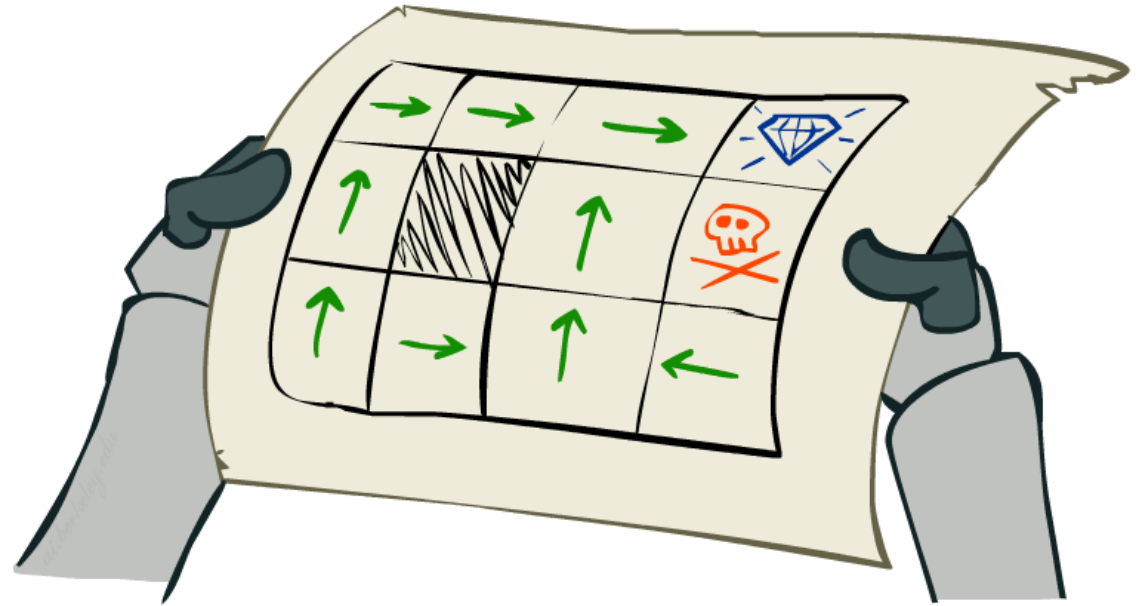
  o So, it takes a long time to learn

Output Values



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

# Passive 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.
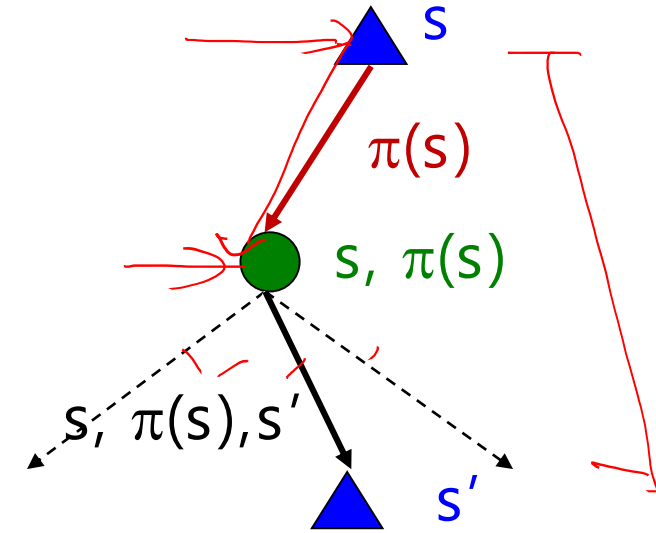
# Why Not Use Policy Evaluation?

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

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

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

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

○ Key question: how can we do this update to V without knowing T and R?
  ○ 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?

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

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

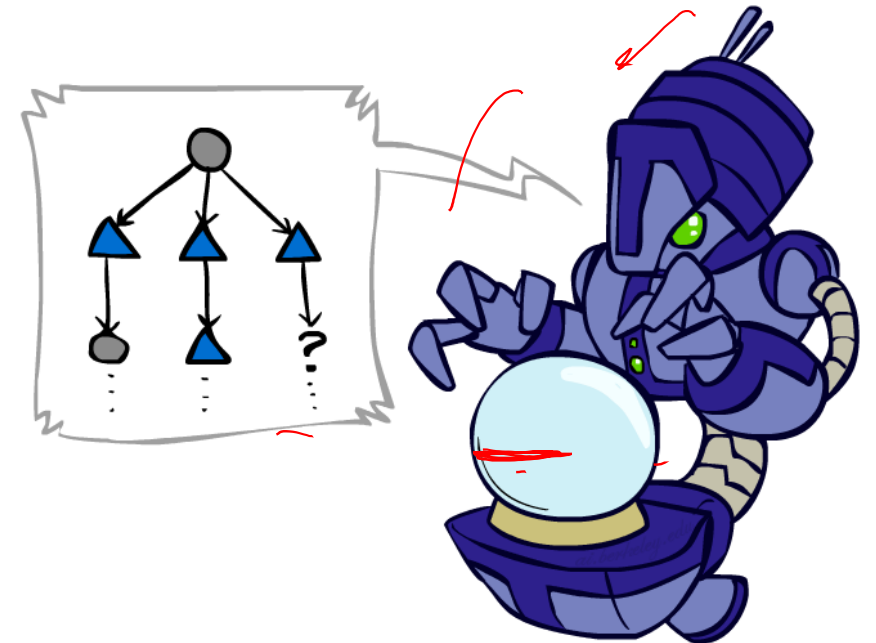o Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

$$\ldots$$

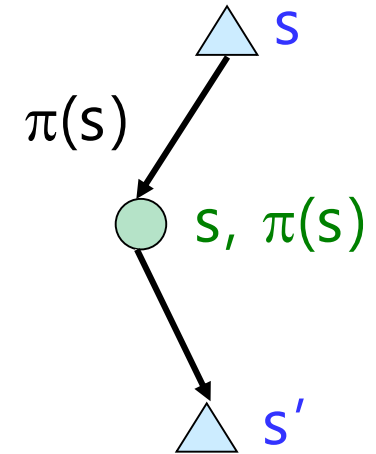$$sample_n = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$

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

*from state s.*

# Temporal Difference Learning

○ Big idea: learn from every experience!
  ○ Update V(s) each time we experience a transition (s, a, s', r)
  ○ Likely outcomes s' will contribute updates more often

○ Temporal difference learning of values
  ○ Policy still fixed, still doing evaluation!
  ○ 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:   $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

$\pi(s)$

s

s, $\pi(s)$

s'

# Exponential Moving Average

o Exponential moving average

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

   o Makes recent samples more important

$$\bar{x}_n = (1 - \alpha) \left[ (1 - \alpha) x_{n-2} + \alpha \cdot x_{n-1} \right] + \alpha \cdot x_n$$

$$\alpha^2$$

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

o Decreasing learning rate (alpha) can give converging averages

# Example: Temporal Difference Learning

## States

## Observed Transitions

| B, east, C, -2 | | C, east, D, -2 |



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

$$V^\pi(s) \leftarrow (1-\alpha)V^\pi(s) + \alpha\left[R(s, \pi(s), s') + \gamma V^\pi(s')\right]$$
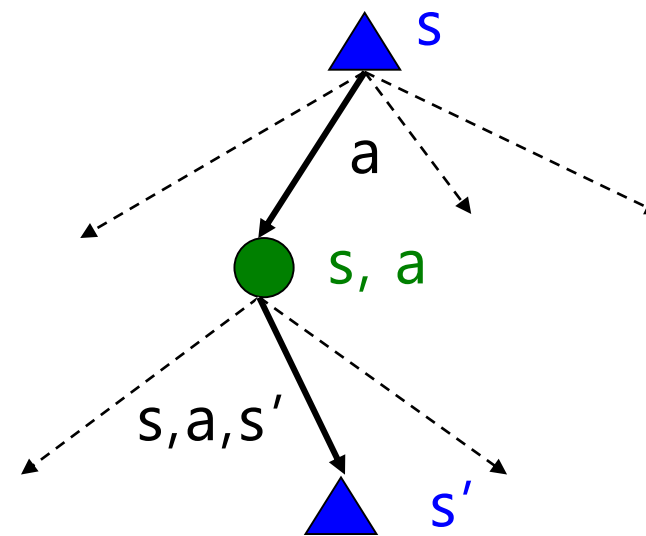
-2

$0 \leftarrow 1/2\left[-2 + 8\right]$

# Problems with TD Value Learning

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

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

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

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



- Idea: learn Q-values, not values
- Makes action selection model-free too!

# Recap: Reinforcement Learning

o Still assume a Markov decision process (MDP):
  - o A set of states $s \in S$
  - 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 $\pi$(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
o Big Idea: Compute all averages over T using sample outcomes

Warm

Cool

Overheated

# The Story So Far: MDPs and RL

## Known MDP: Offline Solution

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Value / policy iteration |
| Evaluate a fixed policy $\pi$ | Policy evaluation |

## Unknown MDP: Model-Based

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | VI/PI on approx. MDP |
| Evaluate a fixed policy $\pi$ | PE on approx. MDP |

## Unknown MDP: Model-Free

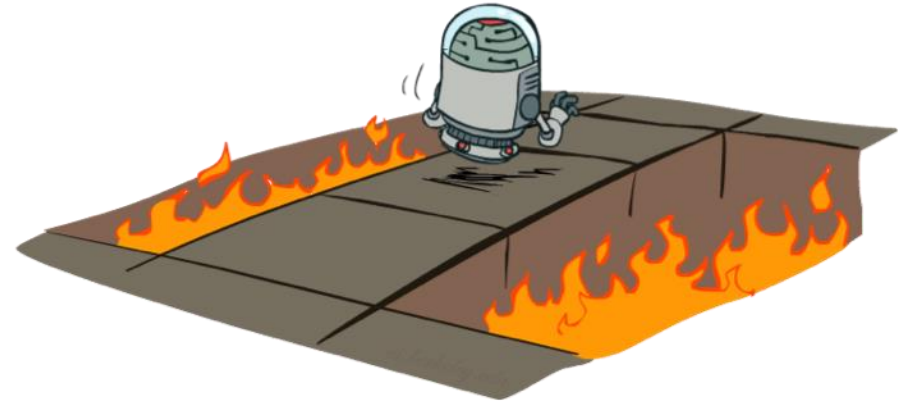| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Q-learning |
| Evaluate a fixed policy $\pi$ | Value Learning |

# Active Reinforcement Learning

o Full reinforcement learning: optimal policies (like value iteration)

    o You don't know the transitions $T(s,a,s')$

    o You don't know the rewards $R(s,a,s')$

    o You choose the actions now

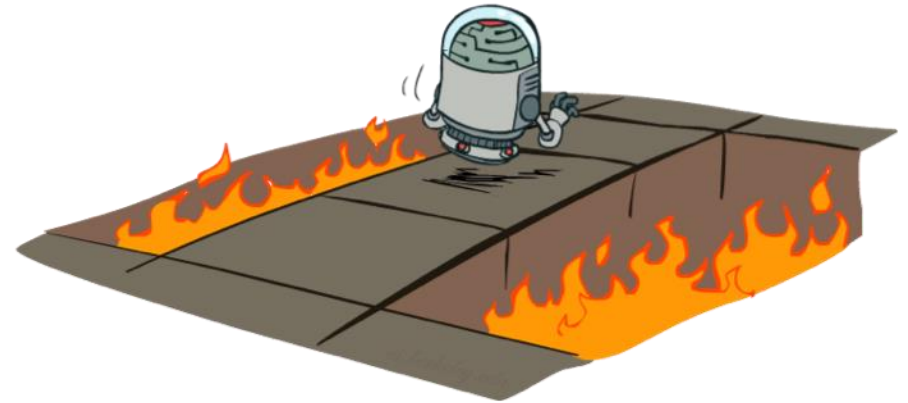    o Goal: learn the optimal policy / values

o In this case:

    o Learner makes choices!

    o Fundamental tradeoff: exploration vs. exploitation

    o This is NOT offline planning!  You actually take actions in the world and find out what happens…

# Model-Free Learning

○ act according to current optimal (based on Q-Values)
○ but also explore…

# Detour: 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

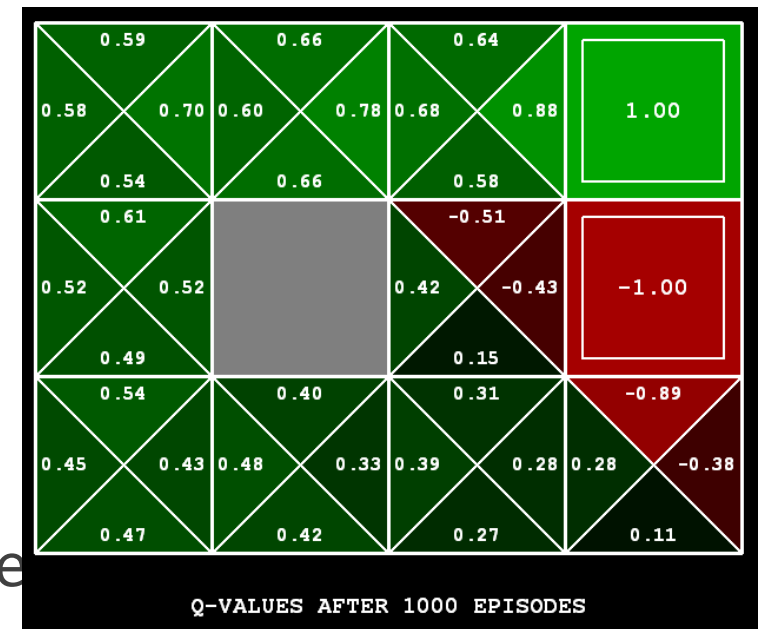$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

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



Q-VALUES AFTER 1000 EPISODES

  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!
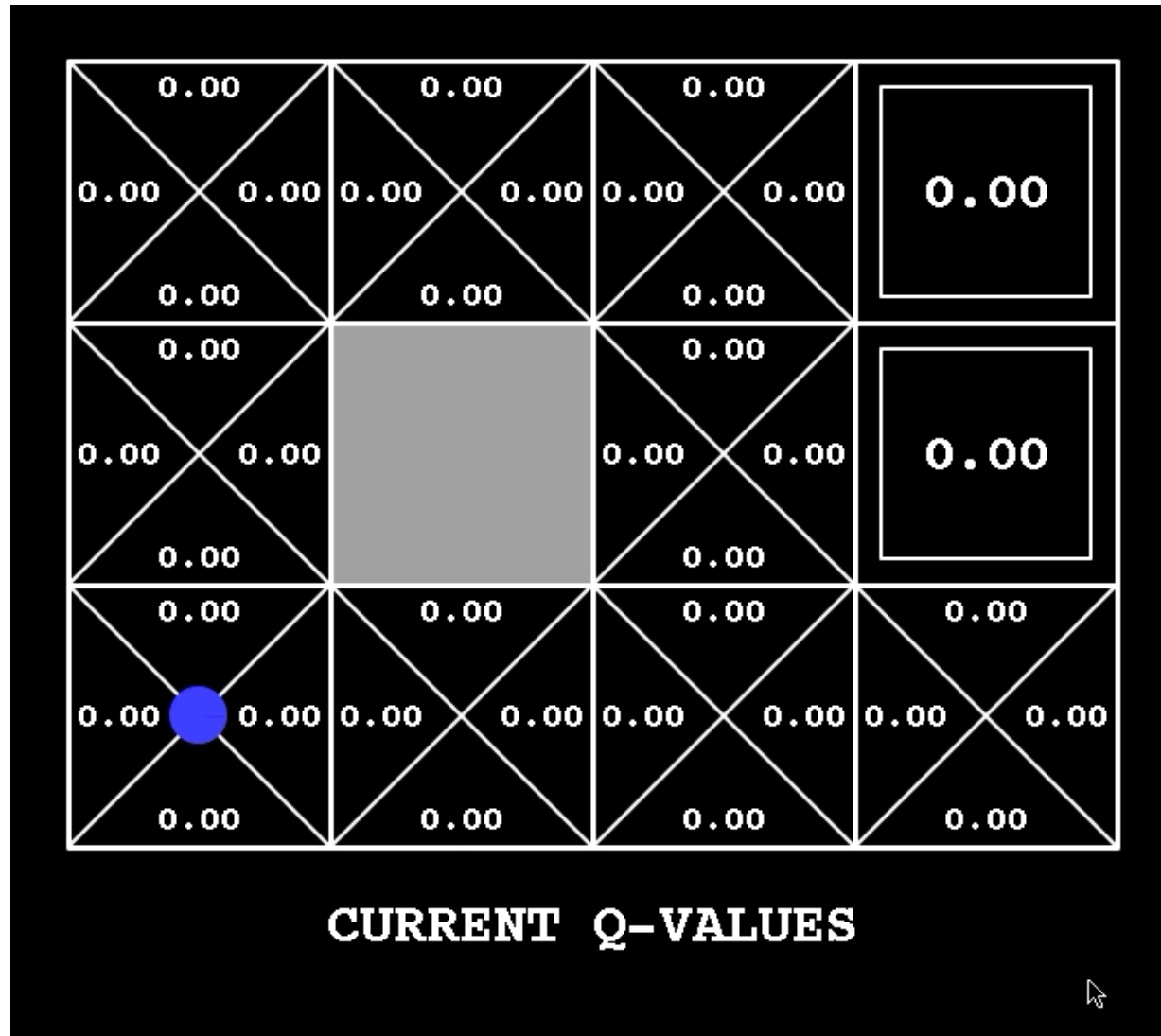
o Incorporate the new estimate into a running average

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

[Demo: Q-learning – gridworld (L10D2)]
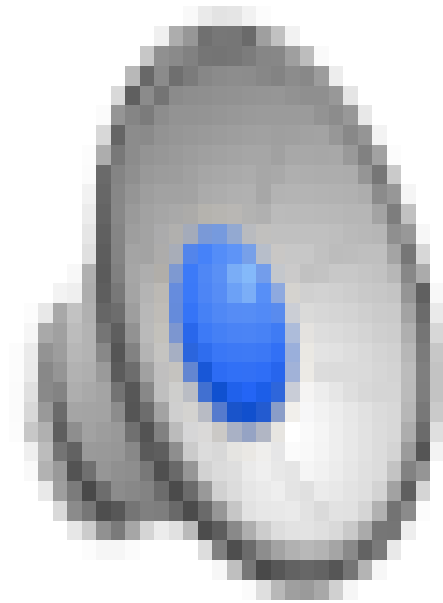[Demo: Q-learning – crawler (L10D3)]

# Q-Learning Demo

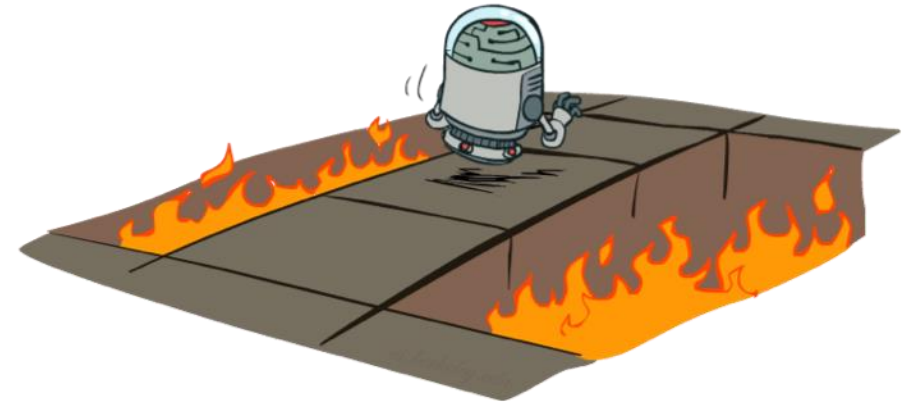# Video of Demo Q-Learning -- Gridworld

# Video of Demo Q-Learning -- Crawler

# Q-Learning:
# act according to current optimal (and also explore...)
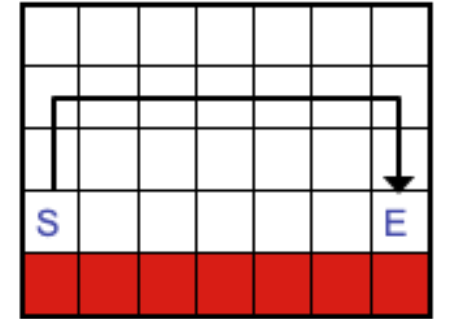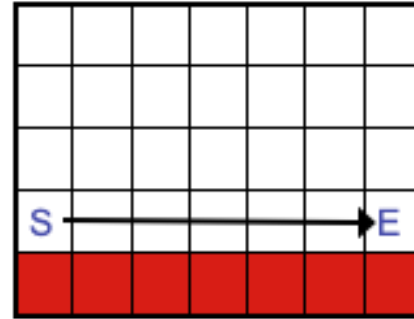
o Full reinforcement learning: optimal policies (like value iteration)
  - o You don't know the transitions T(s,a,s')
  - o You don't know the rewards R(s,a,s')
  - o You choose the actions now
  - o Goal: learn the optimal policy / values

o In this case:
  - o Learner makes choices!
  - o Fundamental tradeoff: exploration vs. exploitation
  - o This is NOT offline planning! You actually take actions in the world and find out what happens...
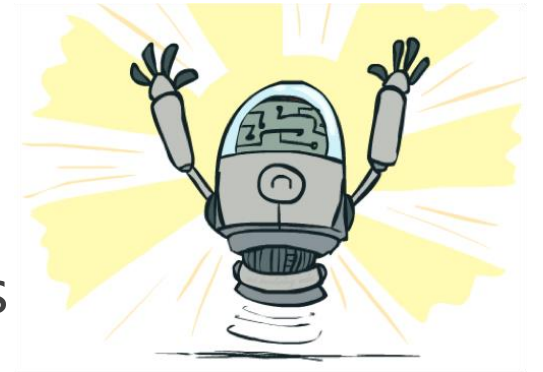
# Q-Learning Properties

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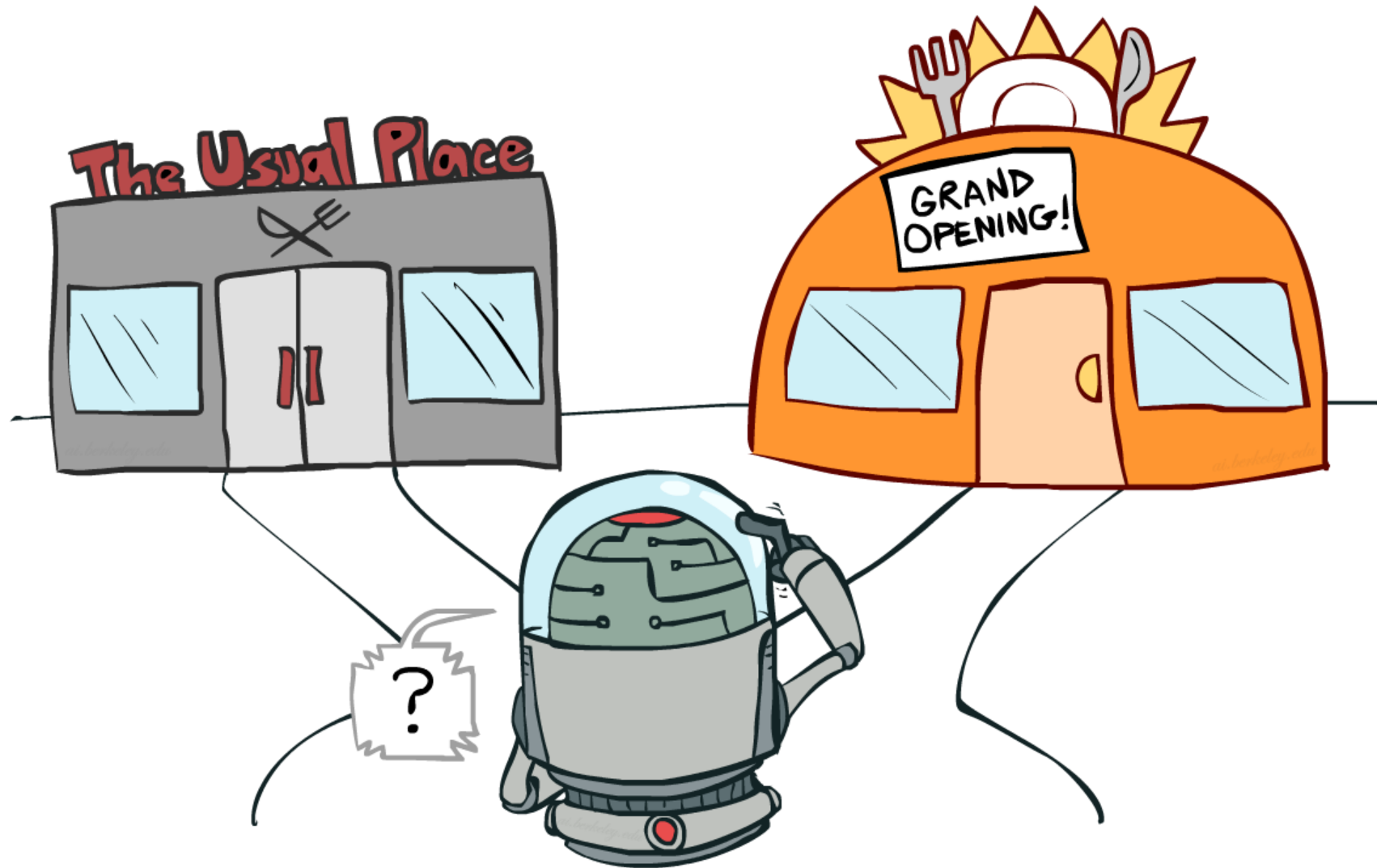even if you're acting suboptimally!

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

# Exploration vs. Exploitation

# How to Explore?

o **Several schemes for forcing exploration**
  - o Simplest: random actions ($\varepsilon$-greedy)
    - o Every time step, flip a coin
    - o With (small) probability $\varepsilon$, act randomly
    - o With (large) probability $1-\varepsilon$, act on current policy

  - o Problems with random actions?
    - o You do eventually explore the space, but keep thrashing around once learning is done
    - o One solution: lower $\varepsilon$ over time
    - o Another solution: exploration functions

# Exploration Functions

o When to explore?

   o Random actions: explore a fixed amount

   o Better idea: explore areas whose badness is not (yet) established, eventually stop exploring
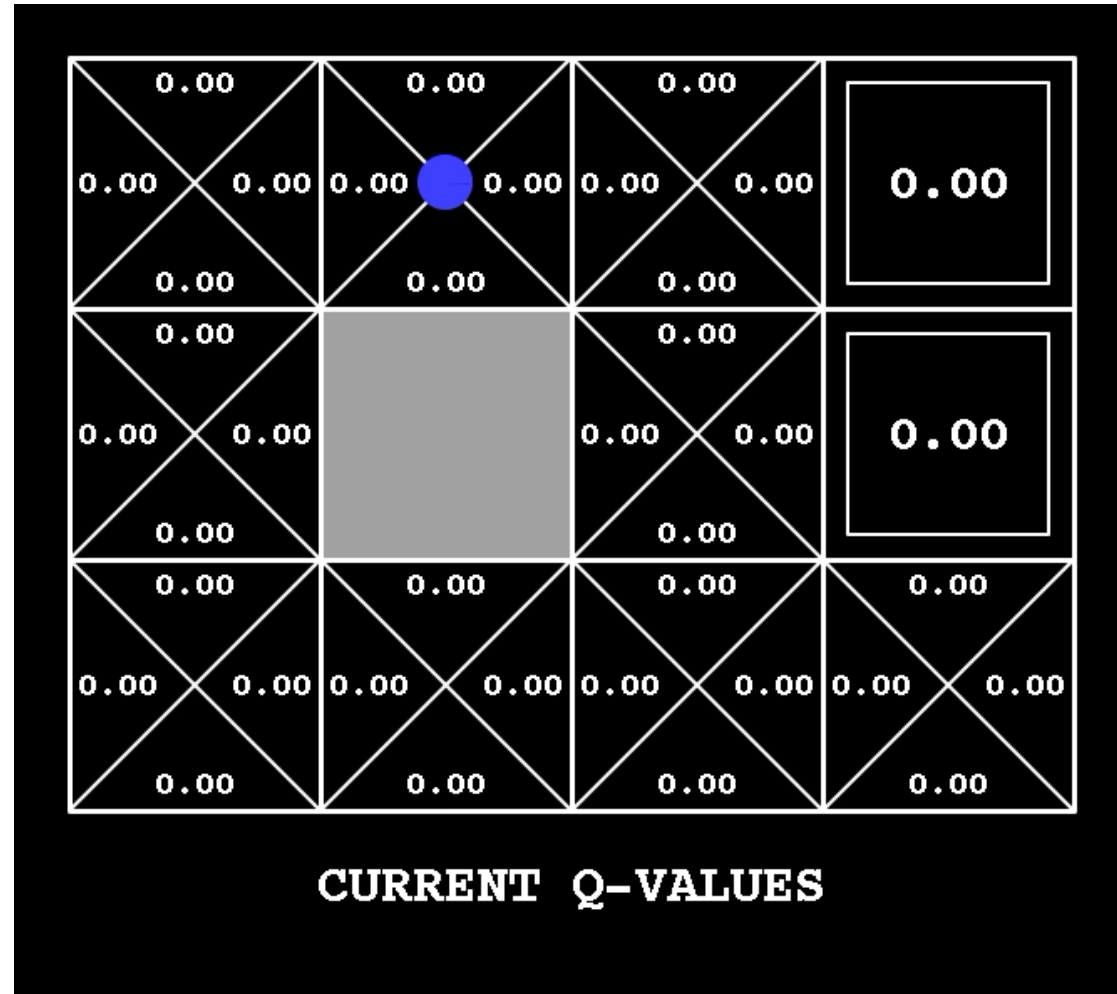
o Exploration function

   o Takes a value estimate u and a visit count n, and returns an optimistic utility, e.g. $f(u,n) = u + k/n$

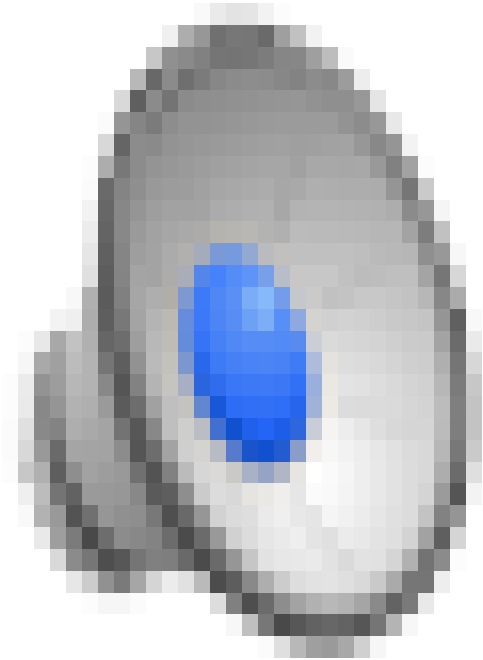     Regular Q-Update: $\quad Q(s,a) \leftarrow_\alpha R(s,a,s') + \gamma \max_{a'} Q(s',a')$

     Modified Q-Update: $Q(s,a) \leftarrow_\alpha R(s,a,s') + \gamma \max_{a'} f(Q(s',a'), N(s',a'))$

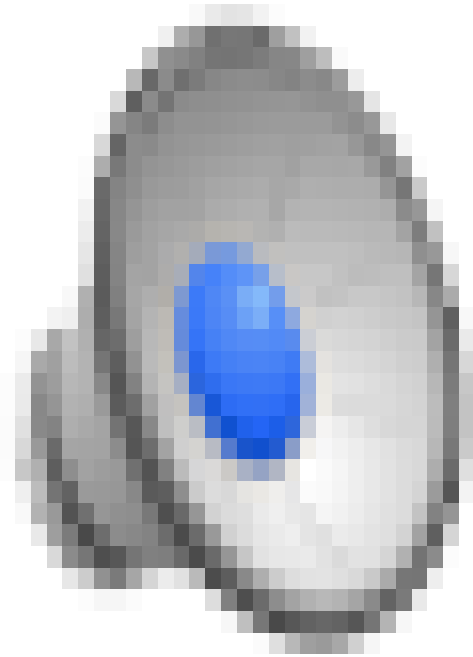   o Note: this propagates the "bonus" back to states that lead to unknown states as well!

<span style="color:red">[Demo: exploration – Q-learning – crawler – exploration function (L11D4)]</span>

# Q-Learn Epsilon Greedy

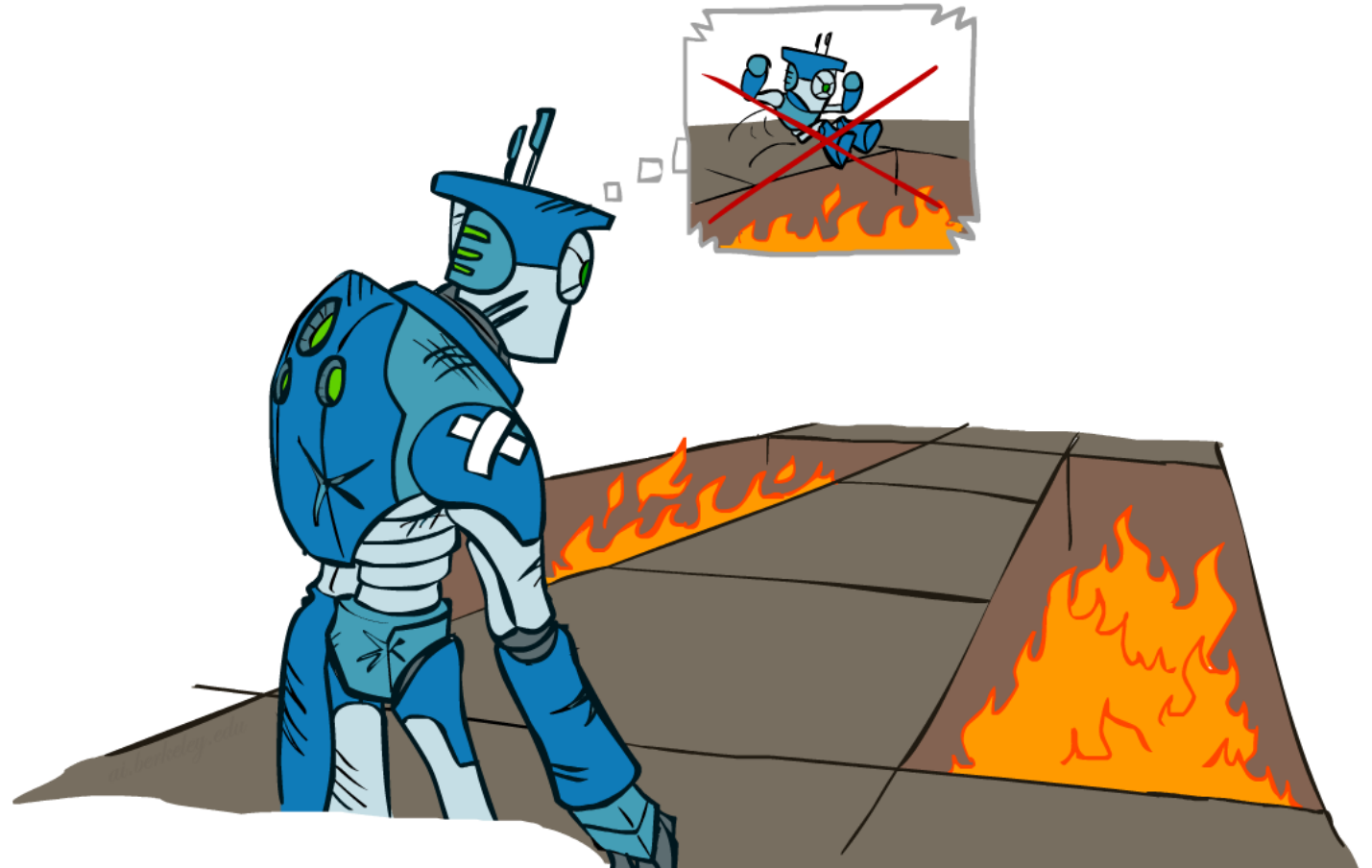# Video of Demo Q-learning – Epsilon-Greedy – Crawler

# Video of Demo Q-learning – Exploration Function – Crawler
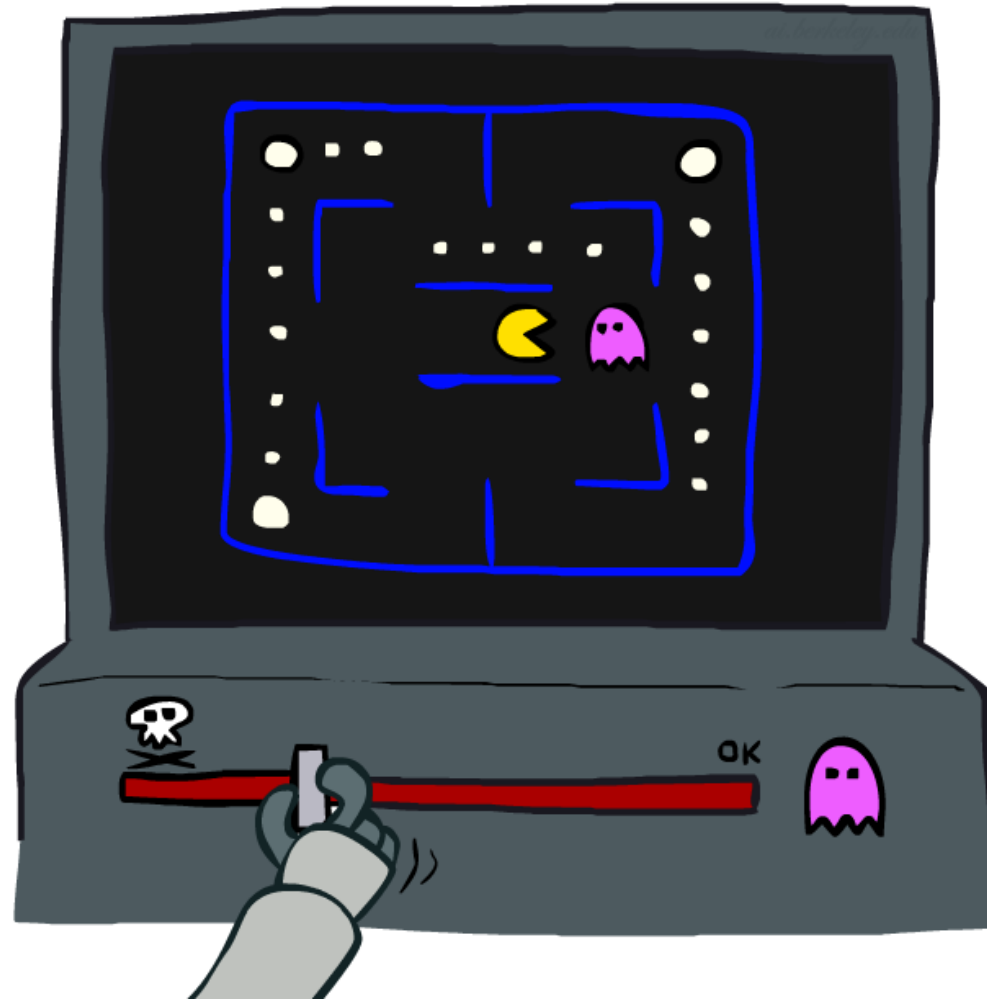
# Regret

o Even if you learn the optimal policy you still make mistakes along the way!

o Regret is a measure of your total mistake cost: the difference between your (expected) rewards and optimal (expected) rewards

o Minimizing regret goes beyond learning to be optimal – it requires optimally learning to be optimal

o Example: random exploration and exploration functions both end up optimal, but random exploration has higher regret
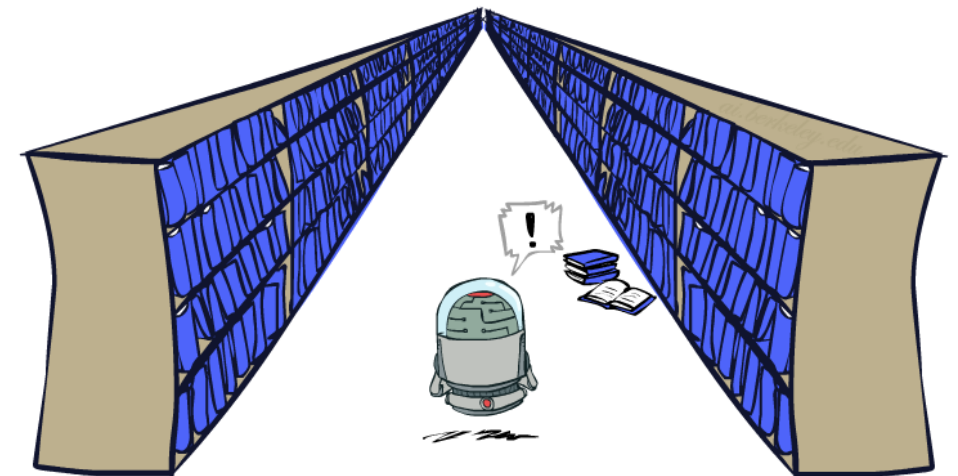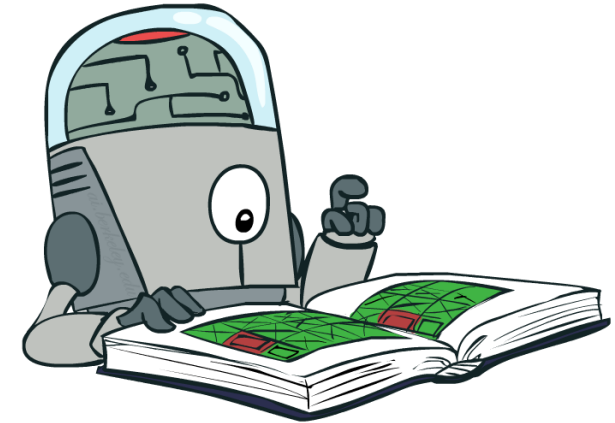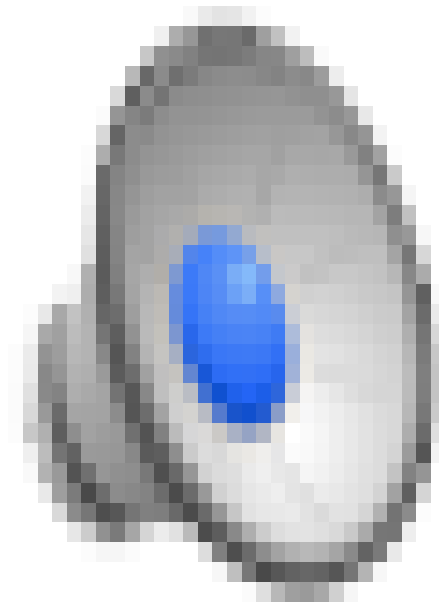
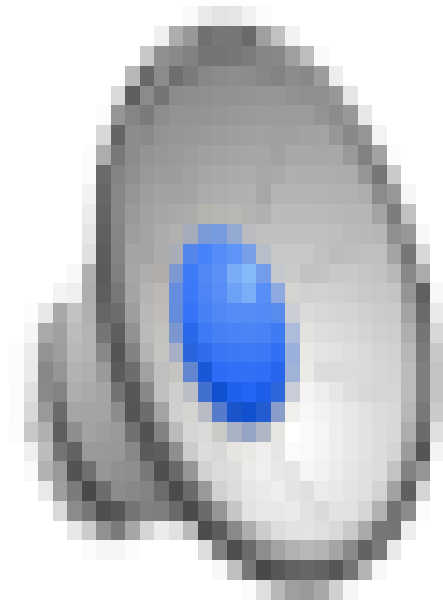# Approximate Q-Learning

# Generalizing Across States

o Basic Q-Learning keeps a table of all q-values

o In realistic situations, we cannot possibly learn about every single state!
   o Too many states to visit them all in training
   o Too many states to hold the q-tables in memory

o Instead, we want to generalize:
   o Learn about some small number of training states from experience
   o Generalize that experience to new, similar situations
   o This is a fundamental idea in machine learning, and we'll see it over and over again
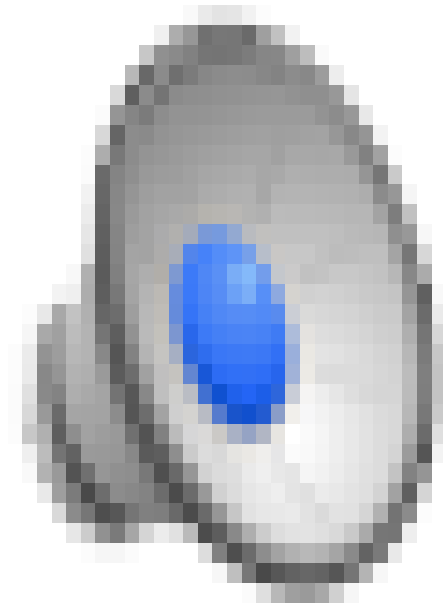
[demo – RL pacman]

# Video of Demo Q-Learning Pacman – Tiny – Watch All

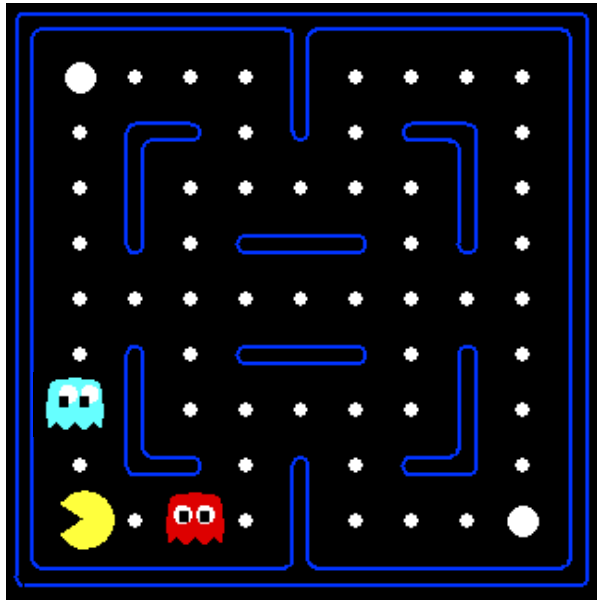# Video of Demo Q-Learning Pacman – Tiny – Silent Train

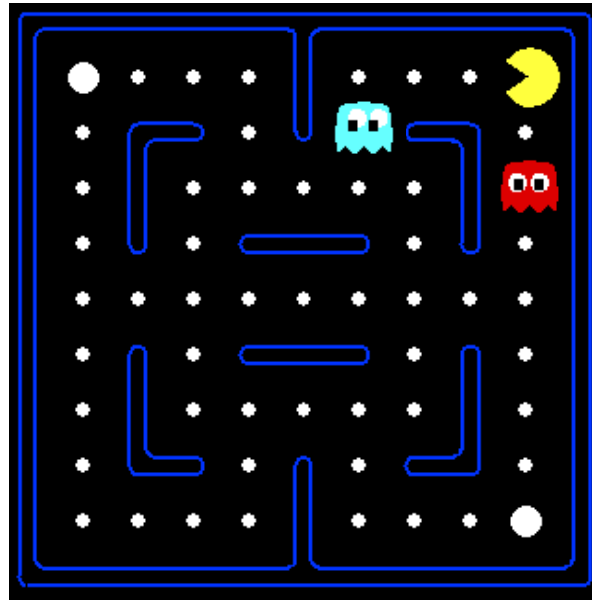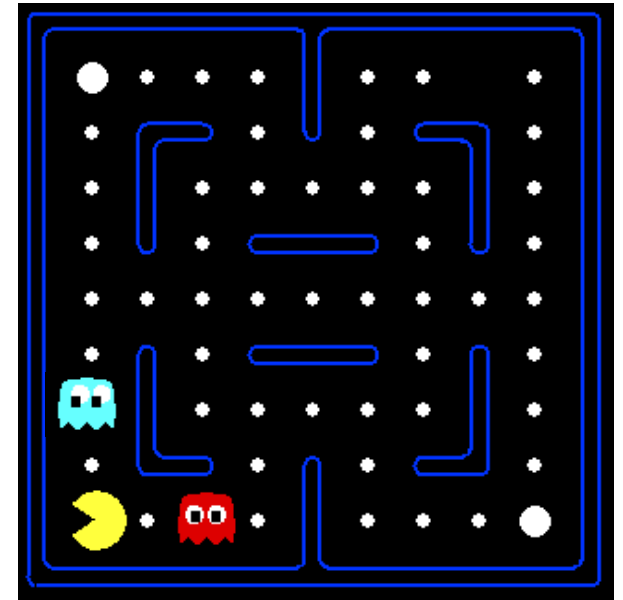# Video of Demo Q-Learning Pacman – Tricky – Watch All

# Example: Pacman

Let's say we discover
through experience
that this state is bad:

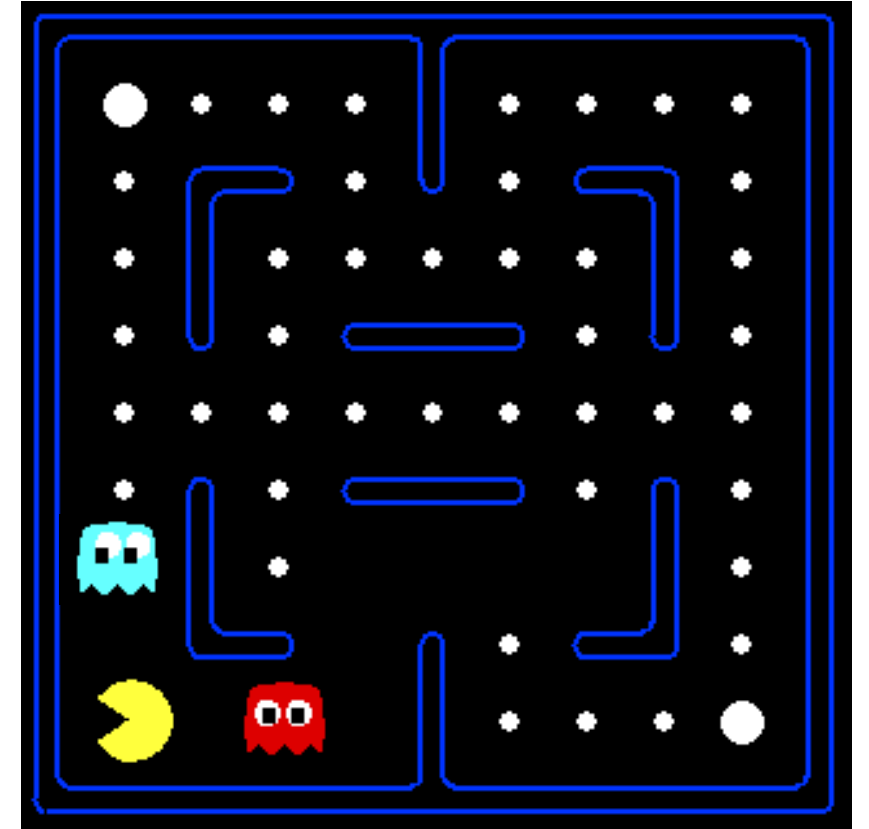In naïve q-learning,
we know nothing
about this state:

Or even this one!

# Feature-Based Representations

o Solution: describe a state using a vector of features (properties)
  - o Features are functions from states to real numbers (often 0/1) that capture important properties of the state
  - o Example features:
    - o Distance to closest ghost
    - o Distance to closest dot
    - o Number of ghosts
    - o 1 / (dist to dot)$^2$
    - o Is Pacman in a tunnel? (0/1)
    - o ... ...  etc.
    - o Is it the exact state on this slide?
  - o Can also describe a q-state (s, a) with features (e.g. action moves closer to food)

# Linear Value Functions

○ Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \ldots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \ldots + w_n f_n(s, a)$$

○ Advantage: our experience is summed up in a few powerful numbers

○ Disadvantage: states may share features but actually be very different in value!

# CSE 573:
# Artificial Intelligence

## Hanna Hajishirzi

## Reinforcement Learning

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer

# Approximate Q-Learning

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

o Q-learning with linear Q-functions:

$$\text{transition } = (s, a, r, s')$$

$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s,a)$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \, [\text{difference}] \qquad \text{Exact Q's}$$

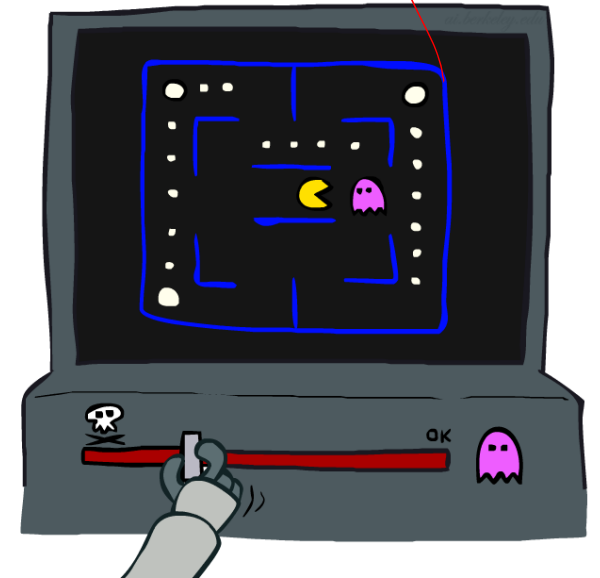$$w_i \leftarrow w_i + \alpha \, [\text{difference}] \, f_i(s,a) \qquad \text{Approximate Q's}$$

o Intuitive interpretation:
  o Adjust weights of active features
  o E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
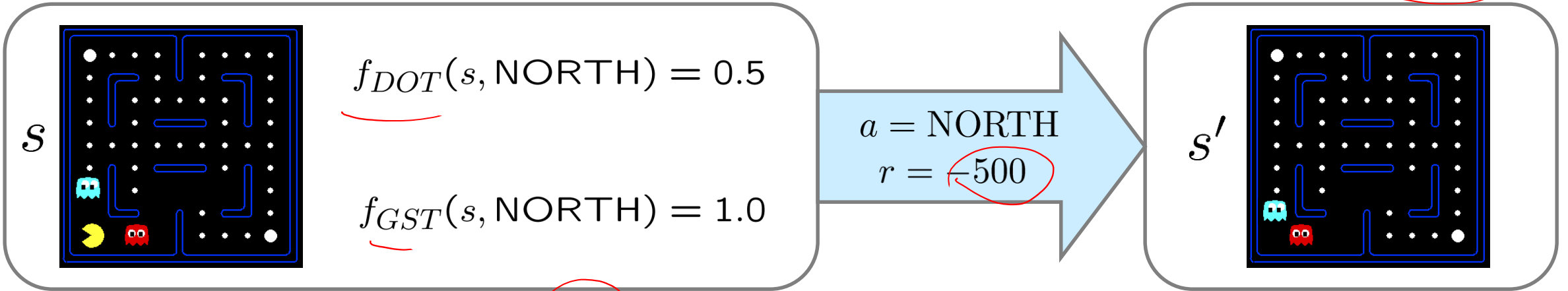
o Formal justification: online least squares

# Example: Q-Pacman

$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$

$\alpha = 2/501$



$s$

$f_{DOT}(s, \text{NORTH}) = 0.5$

$f_{GST}(s, \text{NORTH}) = 1.0$

$a = \text{NORTH}$

$r = -500$

$s'$

$Q(s, \text{NORTH}) = +1$

$Q(s', \cdot) = 0$

$r + \gamma \max_{a'} Q(s', a') = -500 + 0$
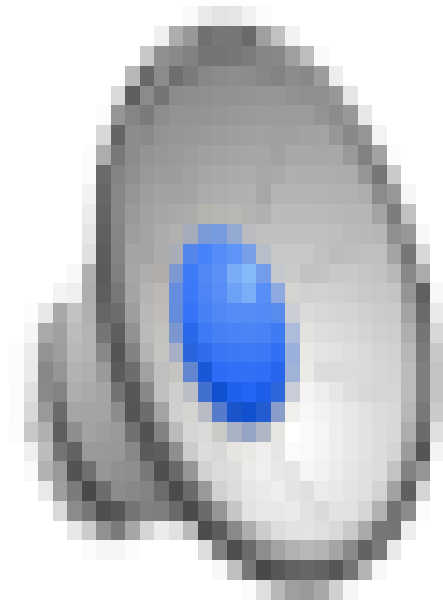
difference $= -501$

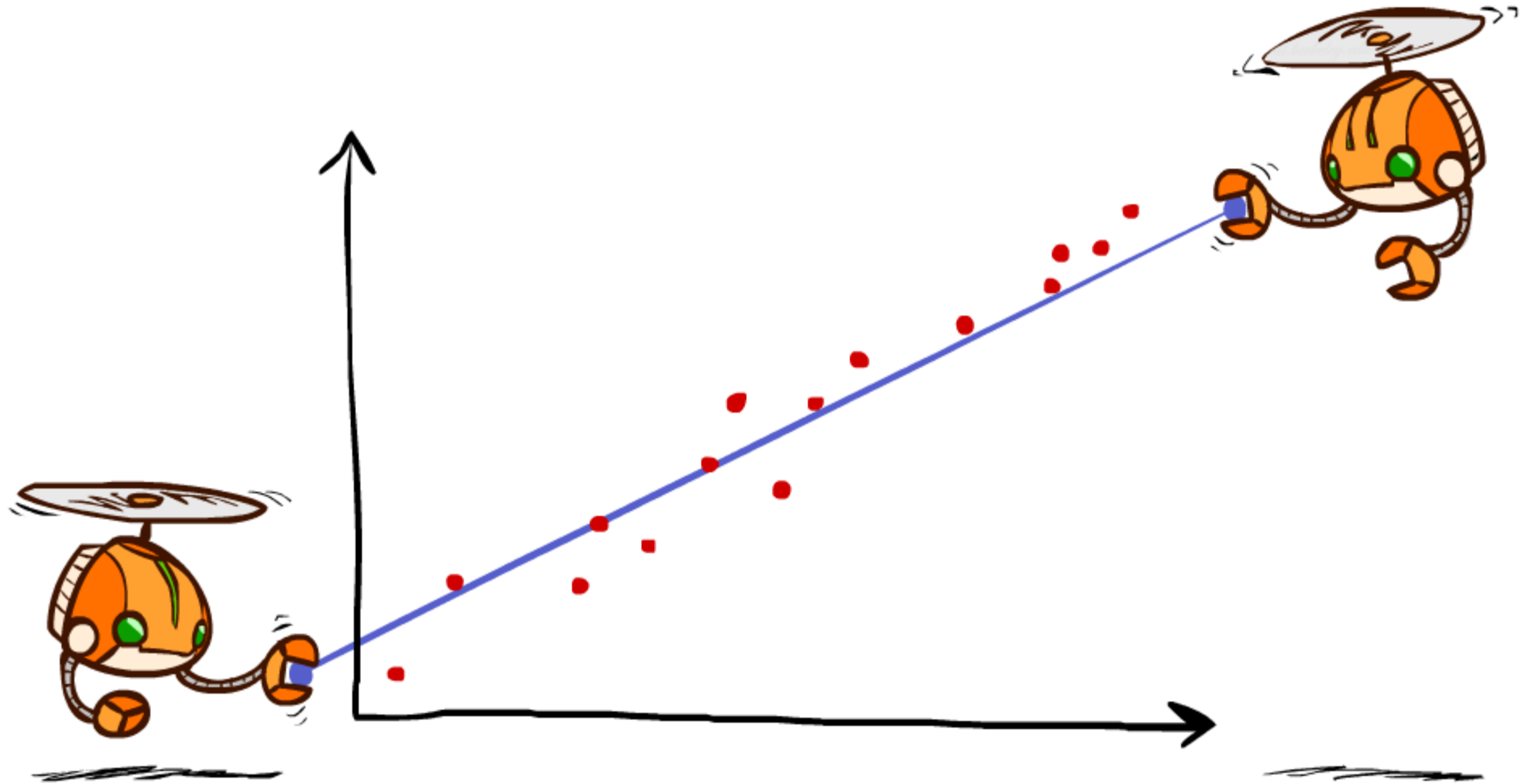$$w_{DOT} \leftarrow 4.0 + \alpha[-501]\,0.5$$
$$w_{GST} \leftarrow -1.0 + \alpha[-501]\,1.0$$

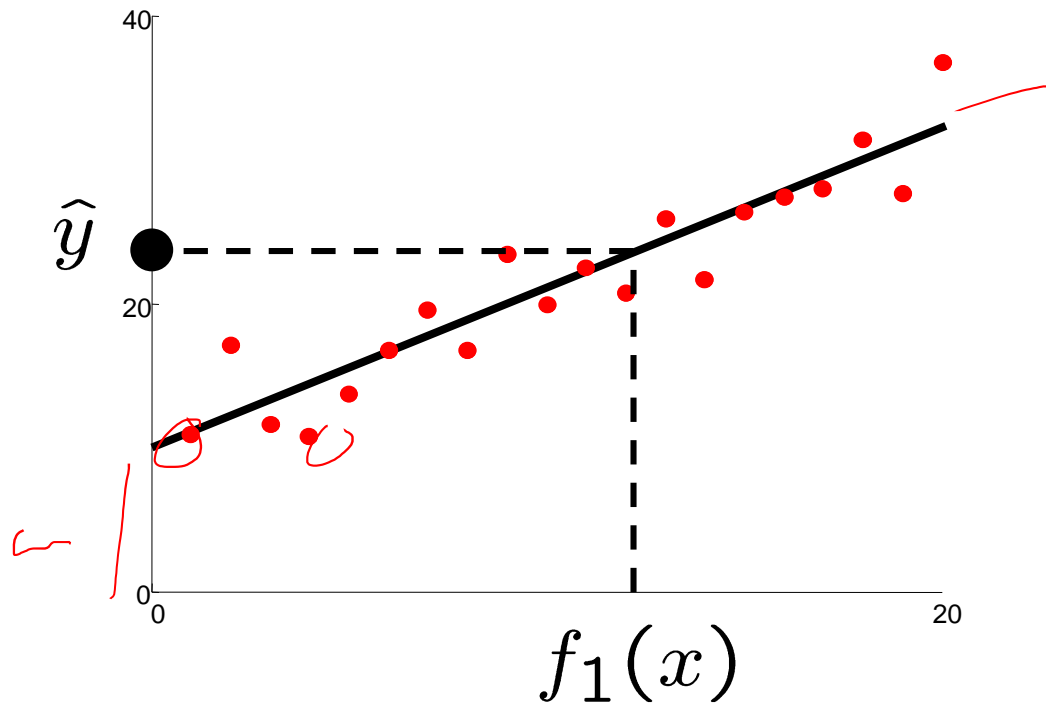$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

# Video of Demo Approximate Q-Learning -- Pacman
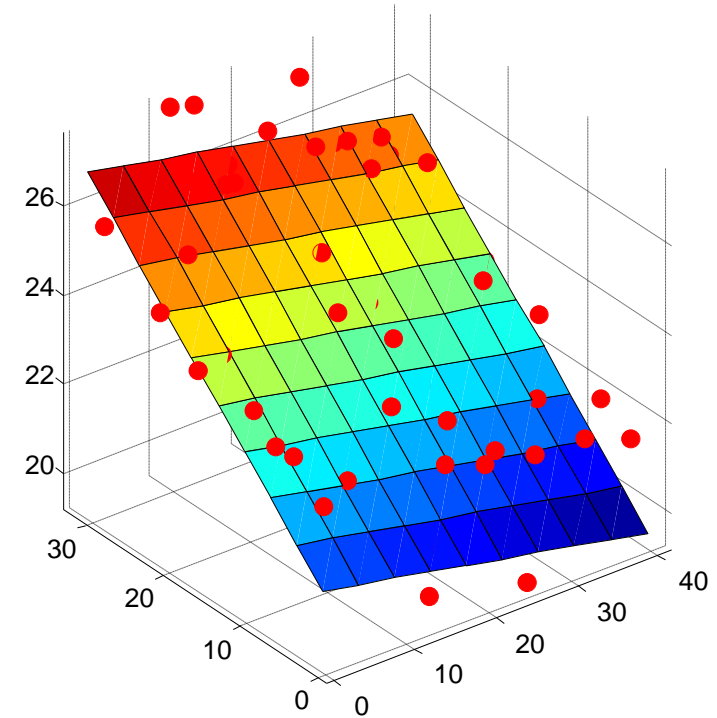
# Q-Learning and Least Squares

# Linear Approximation: Regression



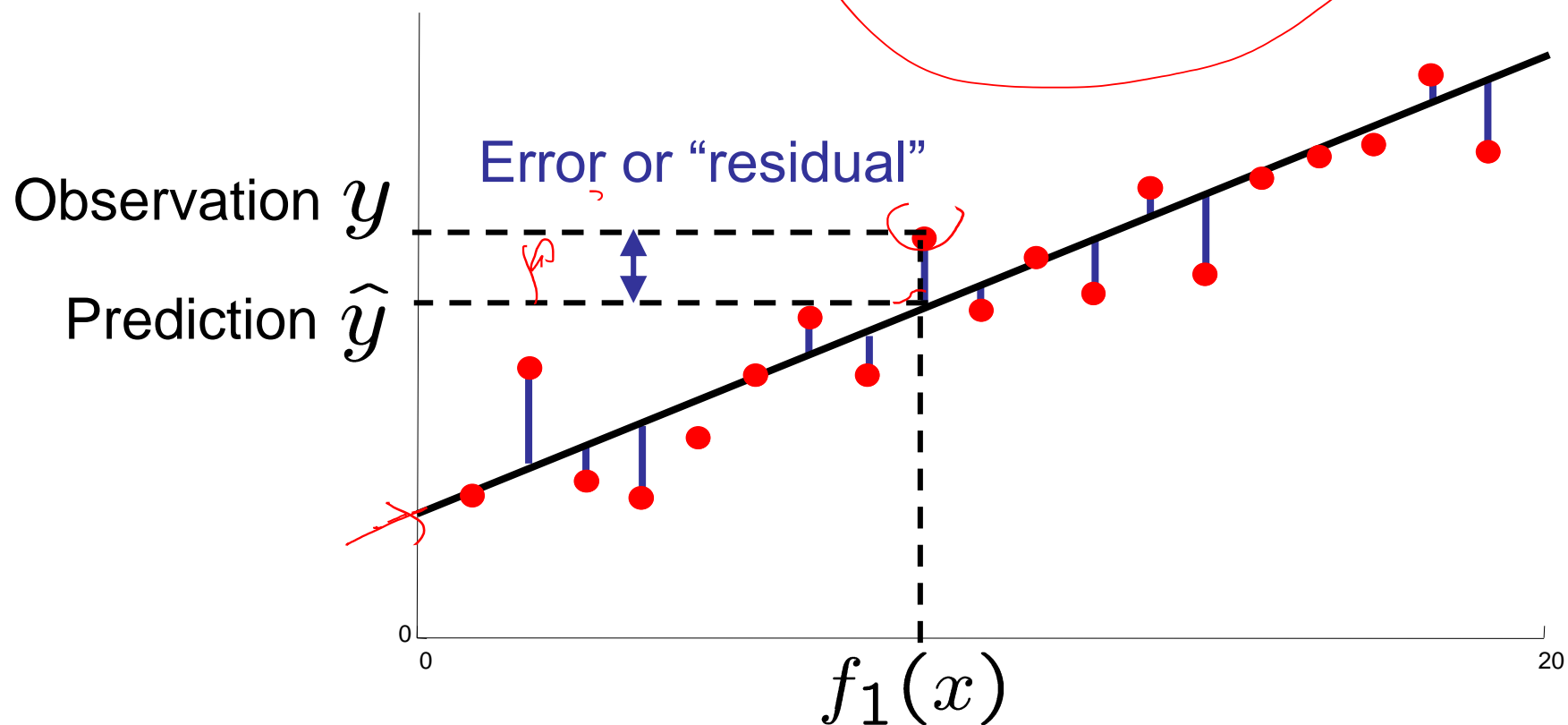Prediction:
$$\hat{y} = w_0 + w_1 f_1(x)$$

Prediction:
$$\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$$
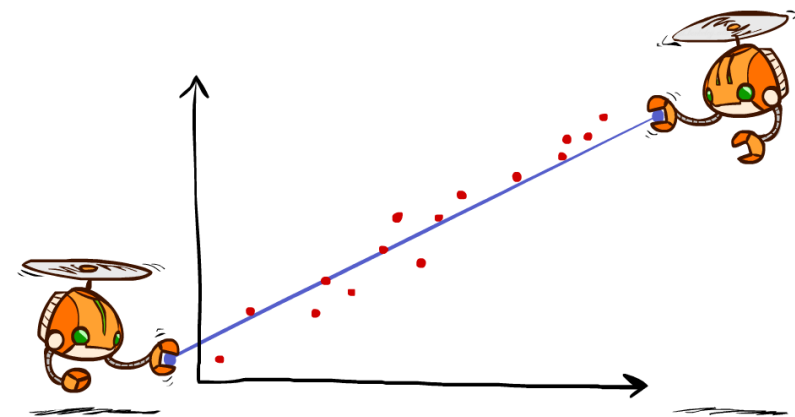
# Optimization: Least Squares

$$\text{total error} = \sum_i (y_i - \widehat{y}_i)^2 = \sum_i \left( y_i - \sum_k w_k f_k(x_i) \right)^2$$



Observation $y$

Prediction $\widehat{y}$

Error or "residual"

$f_1(x)$

# Minimizing Error

Imagine we had only one point x, with features f(x), target value y, and weights w:

$$\text{error}(w) = \frac{1}{2}\left(y - \sum_k w_k f_k(x)\right)^2$$

$$\frac{\partial \text{ error}(w)}{\partial w_m} = -\left(y - \sum_k w_k f_k(x)\right) f_m(x)$$

$$w_m \leftarrow w_m + \alpha \left(y - \sum_k w_k f_k(x)\right) f_m(x)$$

Approximate q update explained:

$$w_m \leftarrow w_m + \alpha \left[ r + \gamma \max_a Q(s', a') - Q(s, a) \right] f_m(s, a)$$
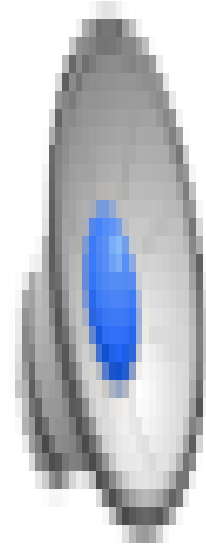
"target"          "prediction"

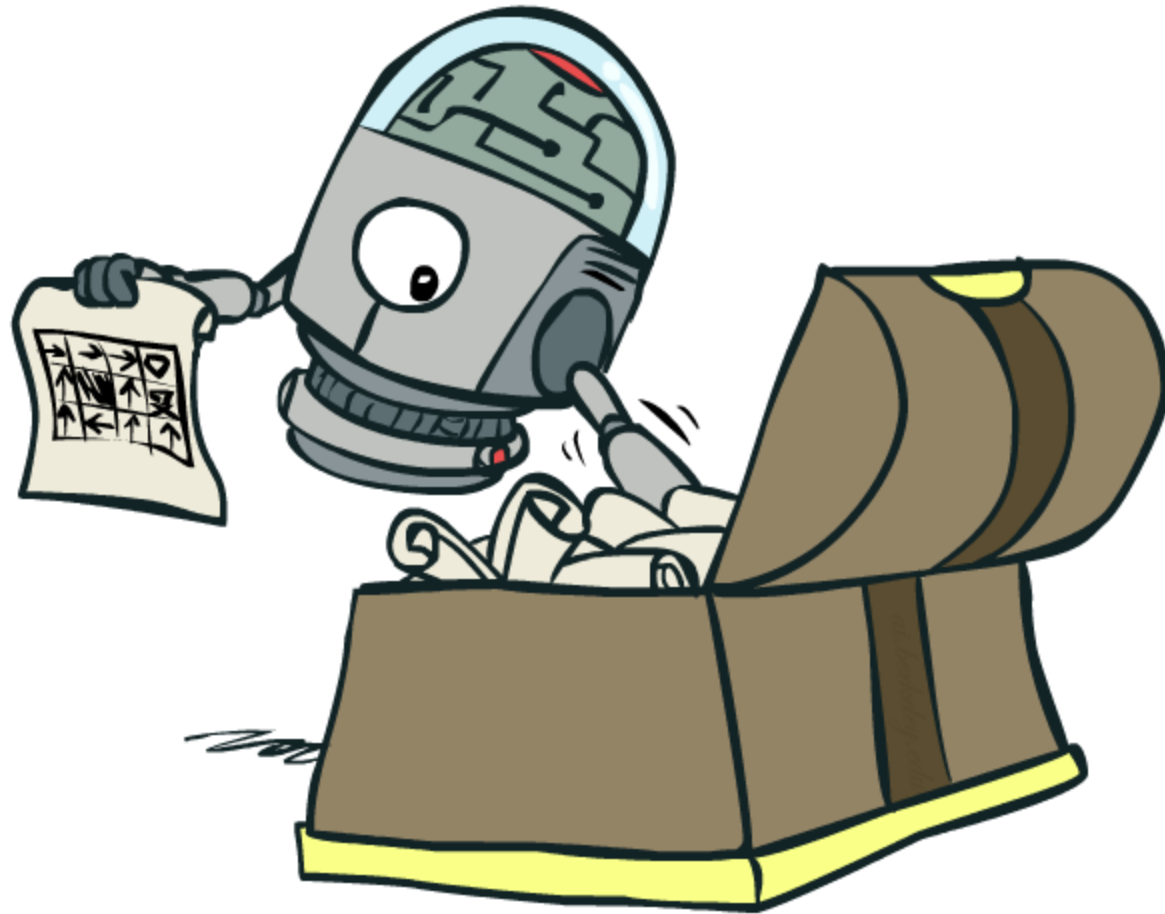# Overfitting: Why Limiting Capacity Can Help

# New in Model-Free RL
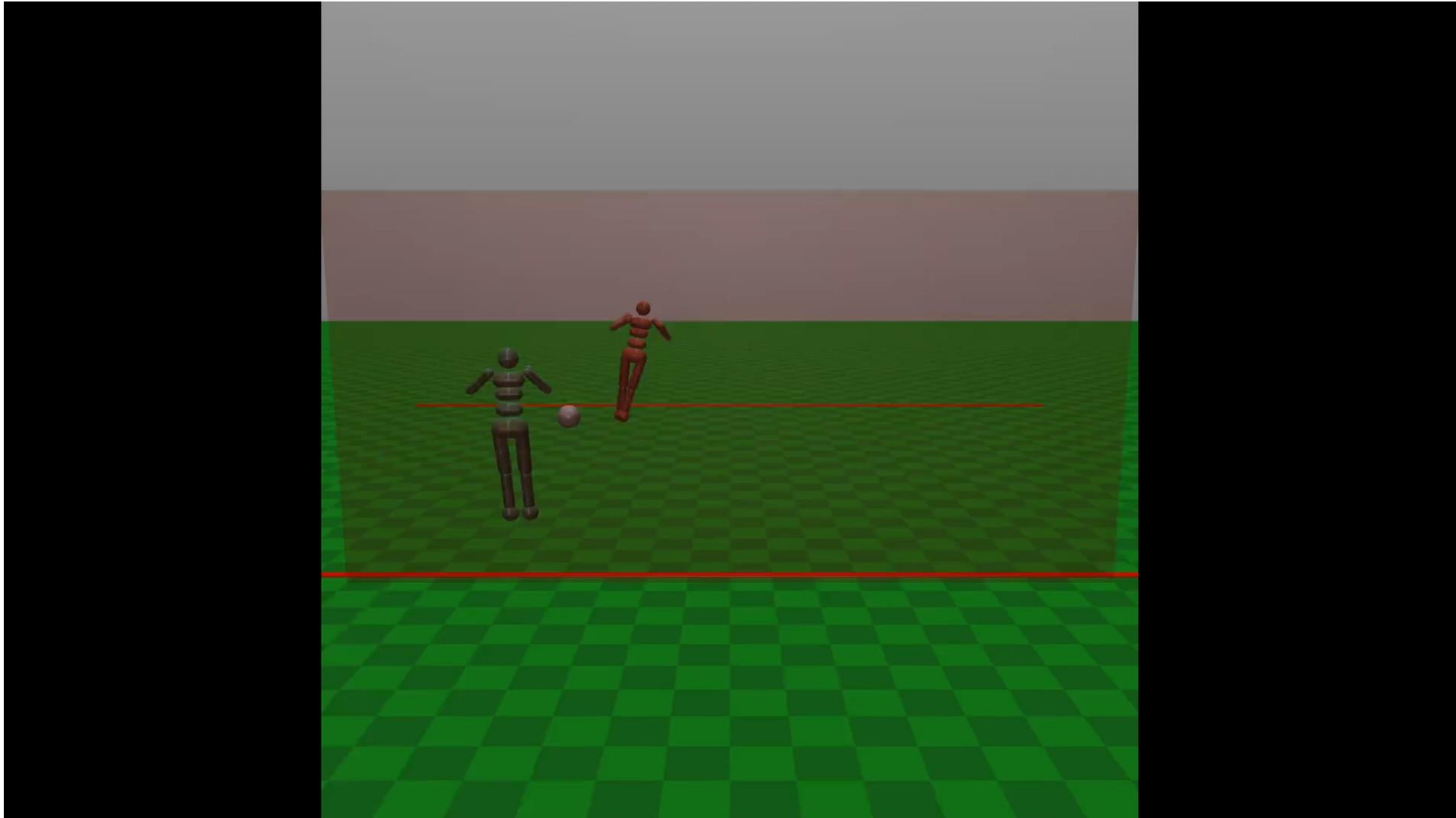# Playing Atari Games

# Policy Search

# Policy Search

o Problem: often the feature-based policies that work well (win games, maximize utilities) aren't the ones that approximate V / Q best

    o E.g. your value functions from project 2 were probably horrible estimates of future rewards, but they still produced good decisions

    o Q-learning's priority: get Q-values close (modeling)

    o Action selection priority: get ordering of Q-values right (prediction)

o Solution: learn policies that maximize rewards, not the values that predict them

o Policy search: start with an ok solution (e.g. Q-learning) then fine-tune by hill climbing on feature weights

# Policy Search

o Simplest policy search:

    o Start with an initial linear value function or Q-function

    o Nudge each feature weight up and down and see if your policy is better than before

o Problems:

    o How do we tell the policy got better?

    o Need to run many sample episodes!

    o If there are a lot of features, this can be impractical

o Better methods exploit lookahead structure, sample wisely, change multiple parameters…

# Summary: MDPs and RL

## Known MDP: Offline Solution

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, π* | Value / policy iteration |
| Evaluate a fixed policy π | Policy evaluation |

## Unknown MDP: Model-Based

*use features to generalize

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, π* | VI/PI on approx. MDP |
| Evaluate a fixed policy π | PE on approx. MDP |

## Unknown MDP: Model-Free

*use features to generalize

| Goal | Technique |
|------|-----------|
| Compute V*, Q*, π* | Q-learning |
| Evaluate a fixed policy π | Value Learning |

# Conclusion

- We've seen how AI methods can solve problems in:
  - Search
  - Games
  - Markov Decision Problems
  - Reinforcement Learning

- Next up: Uncertainty and Learning!