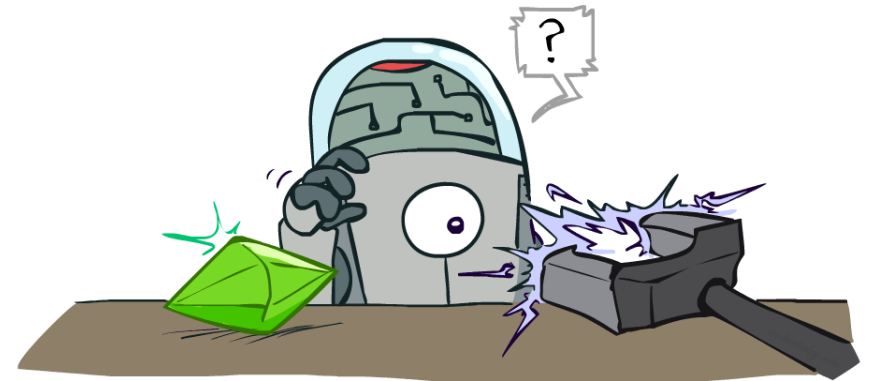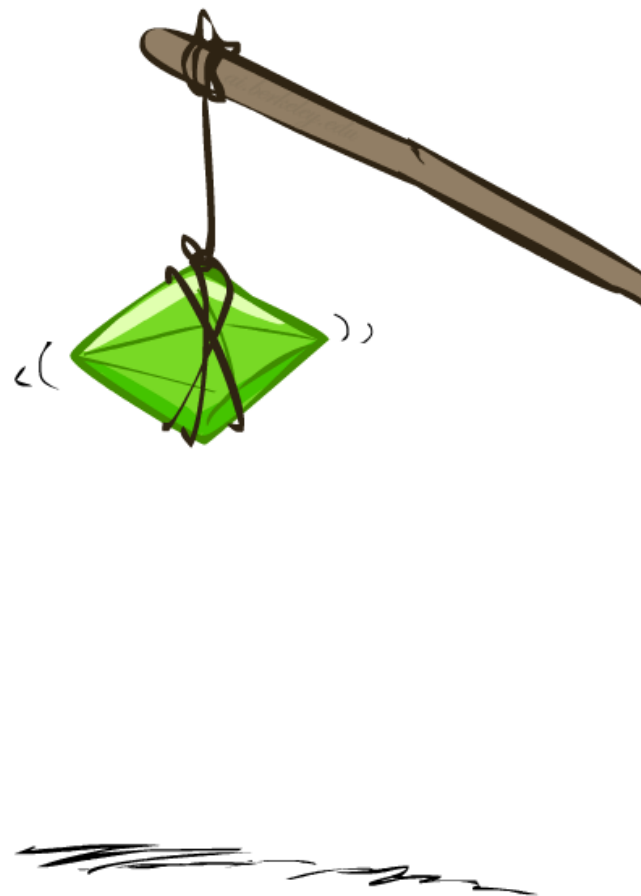# CSE 573 PMP:
# Artificial Intelligence

## Hanna Hajishirzi

## Reinforcement Learning
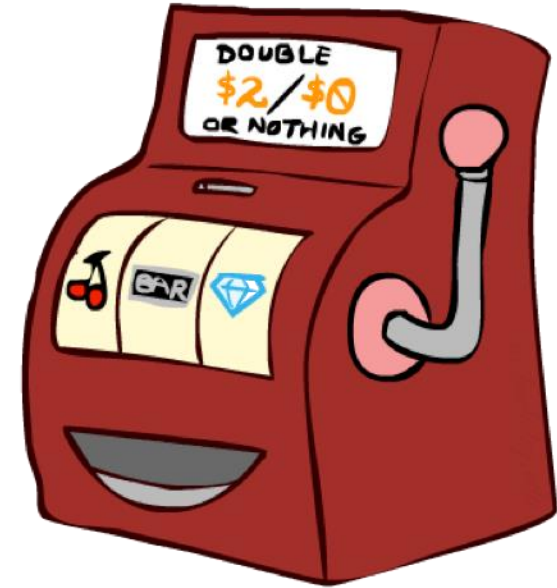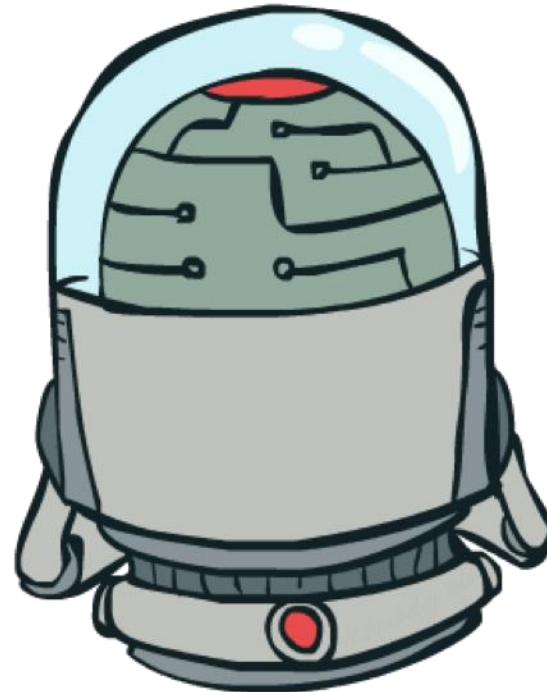
# Reinforcement Learning

# Double Bandits

# Double-Bandit MDP

o Actions: *Blue*, *Red*

o States: Win, Lose



No discount
10 time steps
Both states have
the same value

STRA

0.25    $0

0.75
$2

0.25
$0

$1

1.0

0.75    $2

$1

1.0

Blue : $10
Red : $15

# Offline Planning

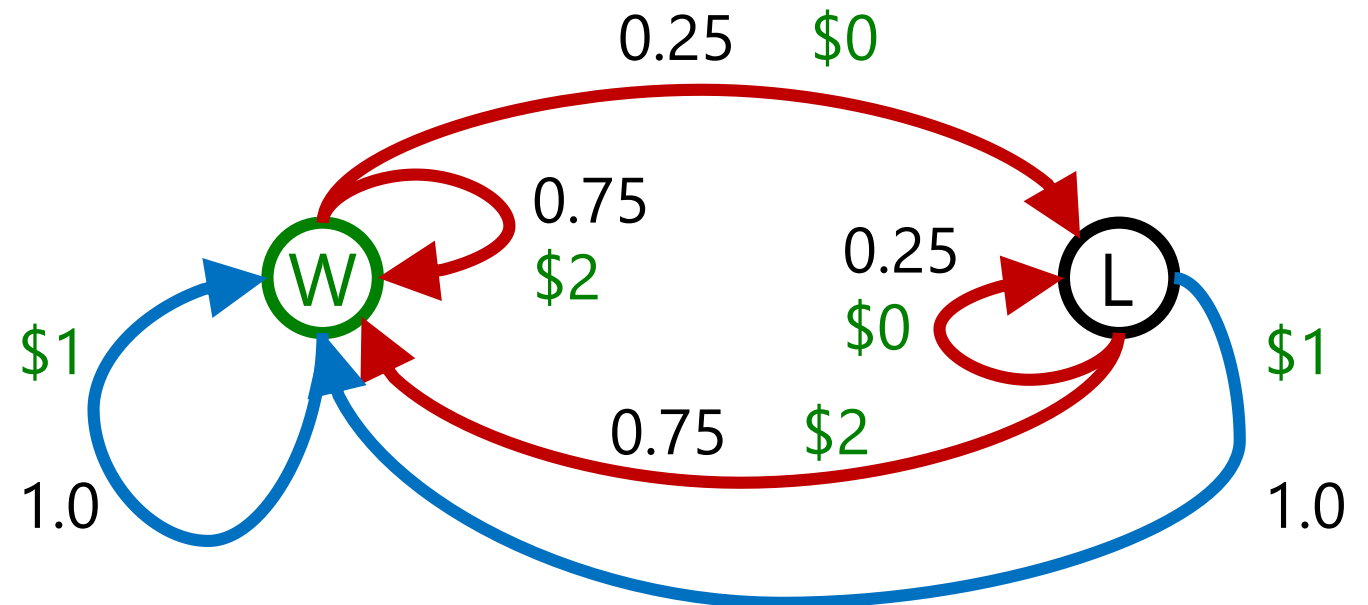o Solving MDPs is offline planning
   o You determine all quantities through computation
   o You need to know the details of the MDP
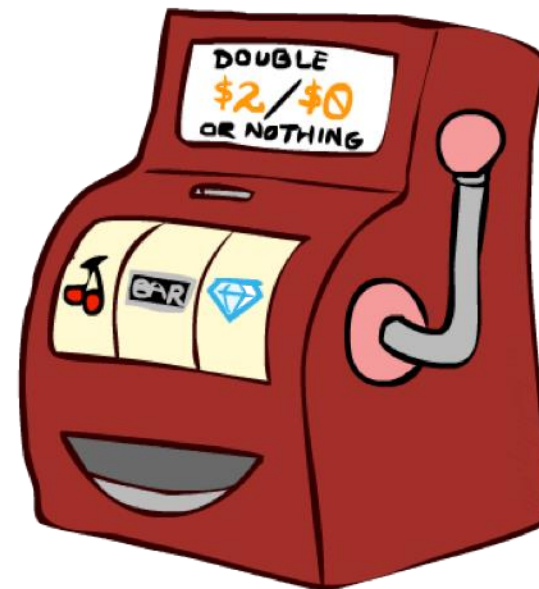   o You do not actually play the game!

*No discount*
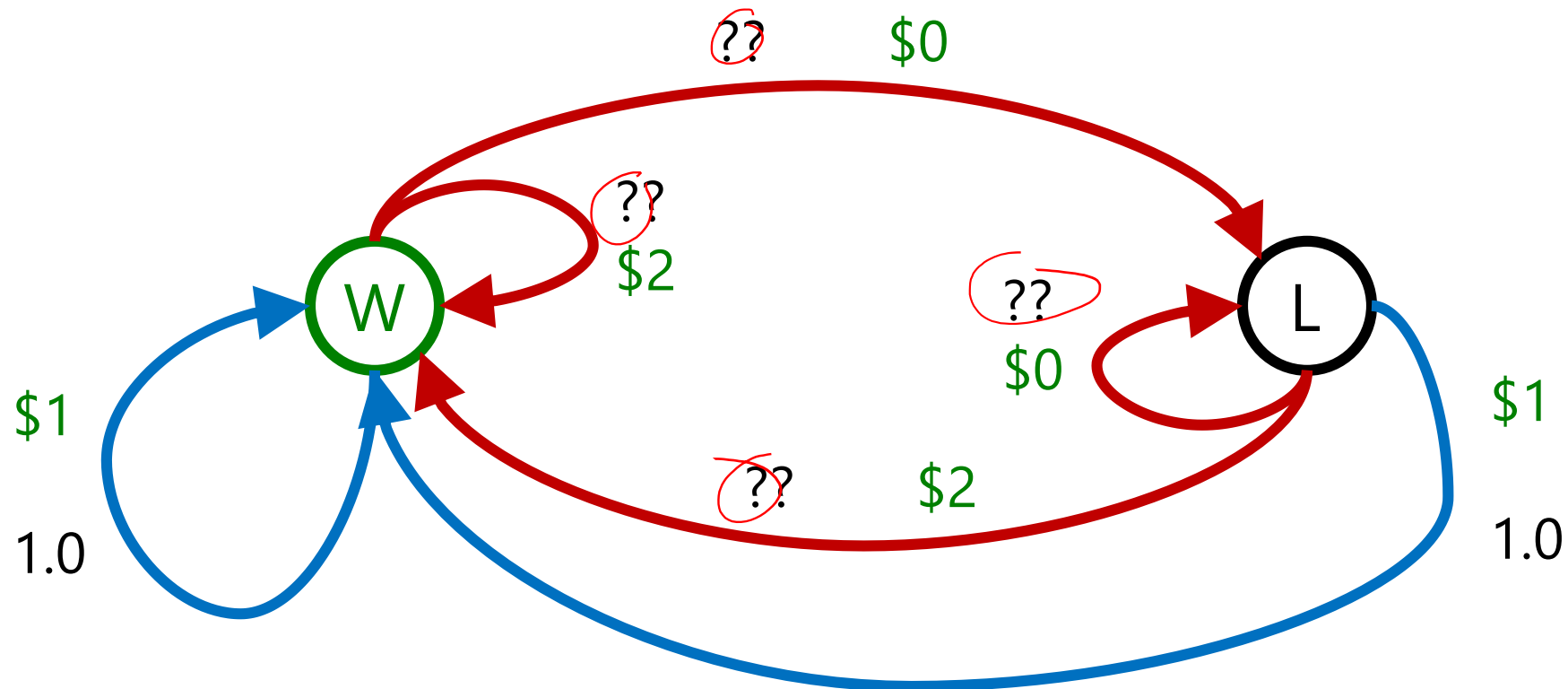
*10 time steps*

# Let's Play!    $15



$2  $2  $0  $2  $2
$2  $2  $0  $0  $0

$112

# Online Planning

○ 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!
  - 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

o 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

o Still assume a Markov decision process (MDP):
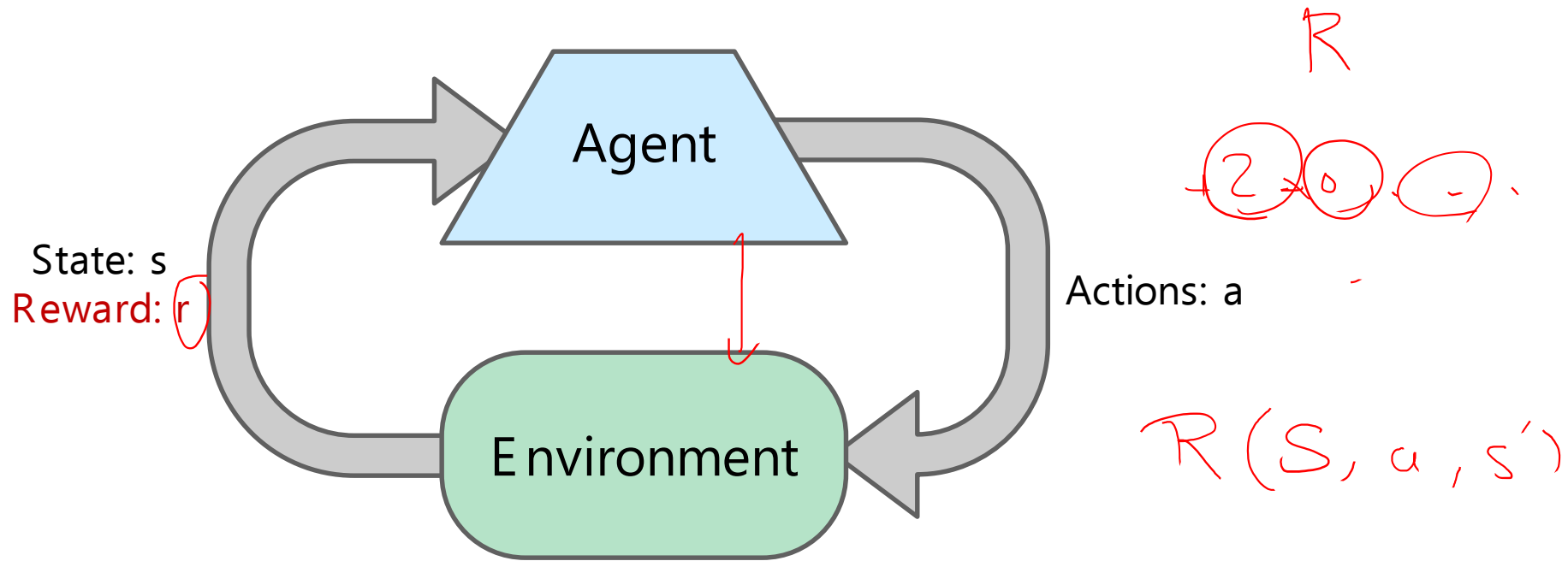  o A set of states s ∈ S
  o A set of actions (per state) A
  o A model T(s,a,s')
  o A reward function R(s,a,s')
o Still looking for a policy π(s)

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

Warm

Cool

Overheated

# Reinforcement Learning



State: s
Reward: r

Agent

Environment

Actions: a

R

+2 0 0 -

R(S, a, s')

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]
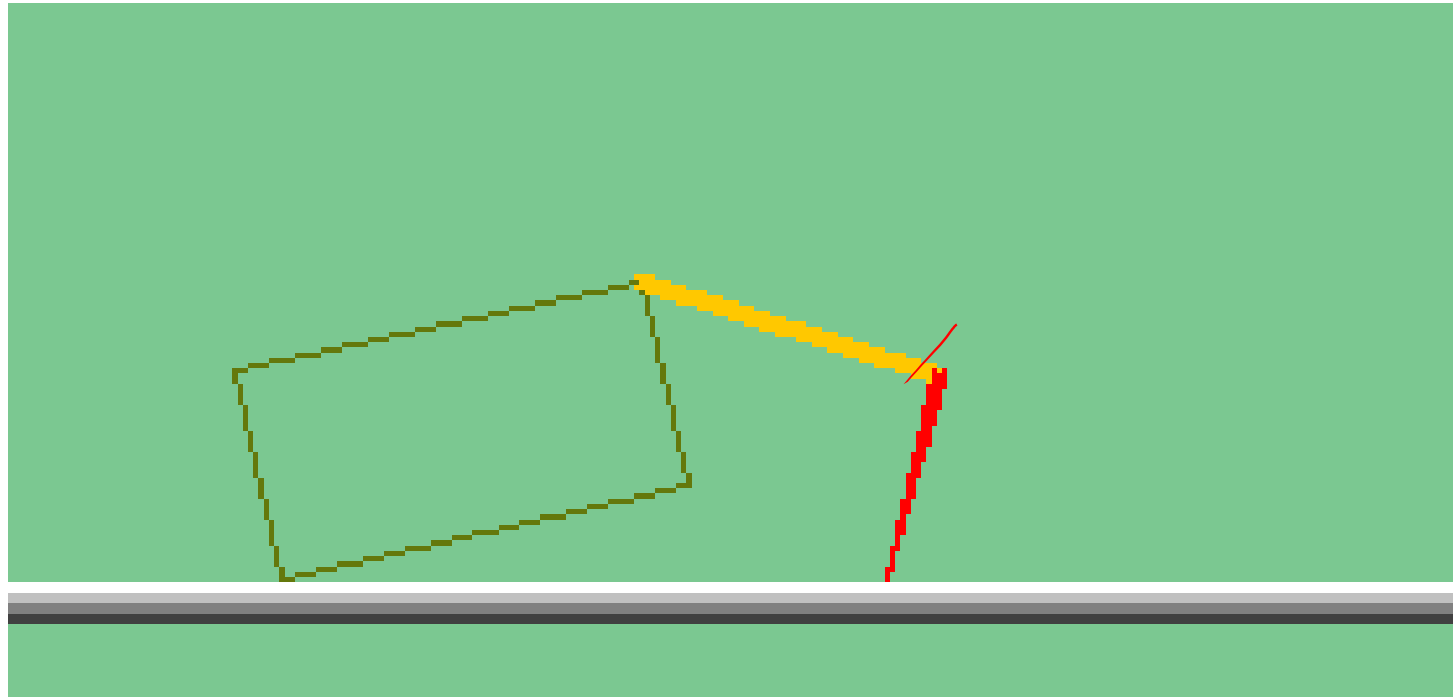
# Example: Toddler Robot



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

# Robotics Rubik Cube
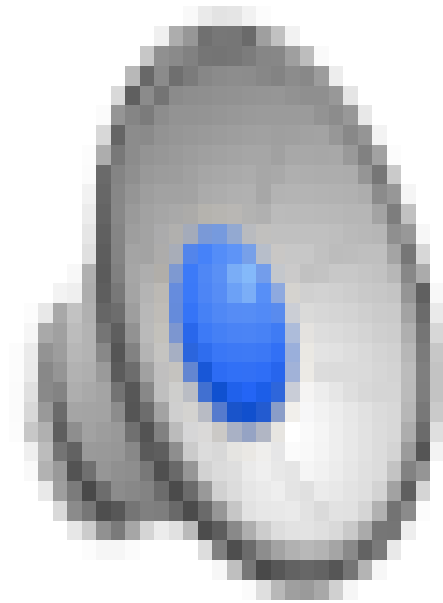
- https://www.youtube.com/watch?v=x4O8pojMF0w

# The Crawler!

# Video of Demo Crawler Bot

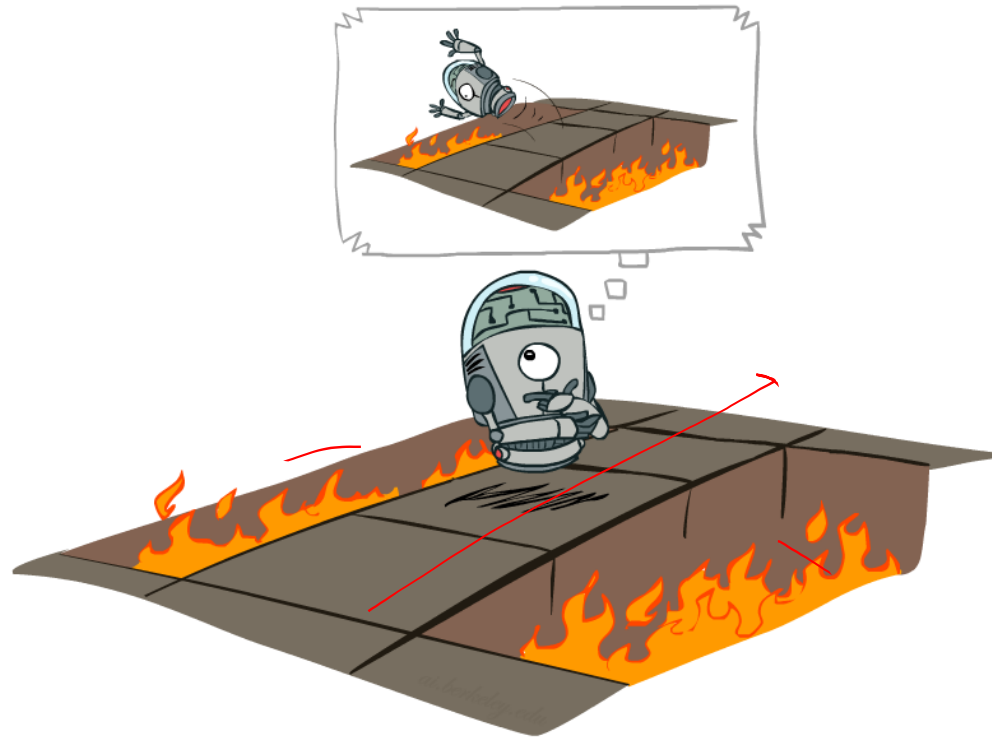# Reinforcement Learning

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



Cool

Warm

Overheated
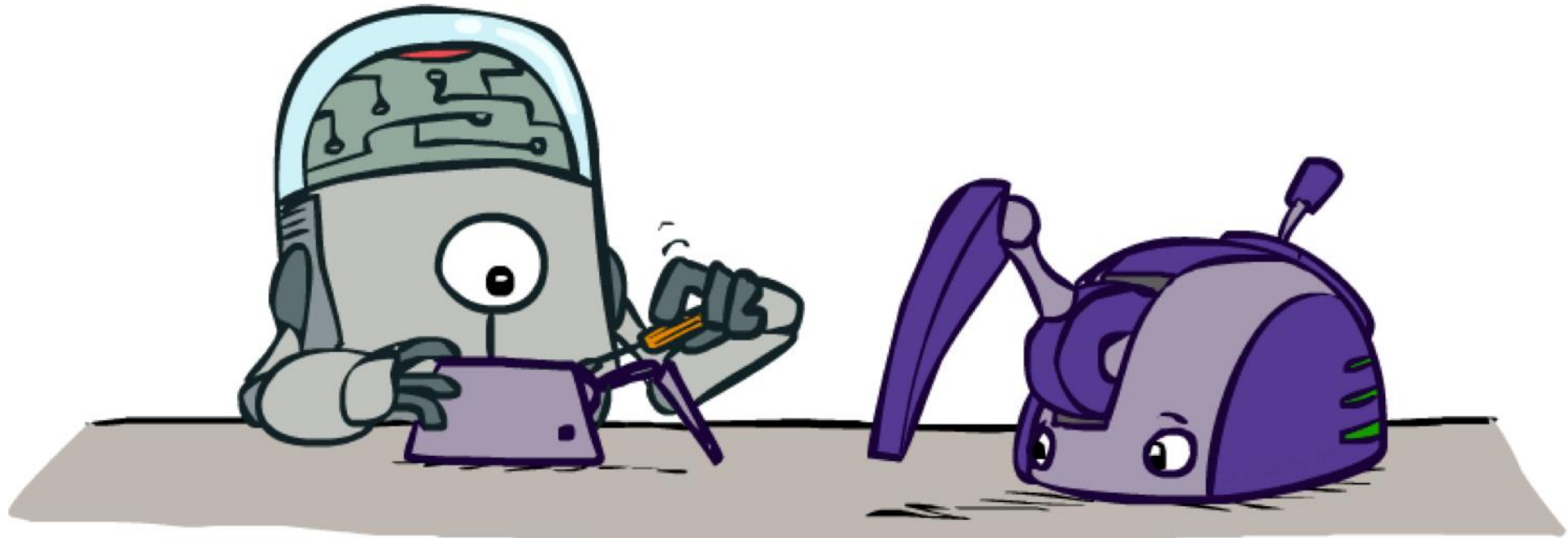
# Offline (MDPs) vs. Online (RL)



Offline Solution

Online Learning

# Model-Based Learning

# Model-Based 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

# Example: Model-Based Learning

## 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

$T(B, east, C) = 100$

$T(C, east, D) = 3|4$

$T(C, east, A) = 1|4$

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

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

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 evaluation

# Example: Direct Evaluation

## Input Policy π

V(B) =



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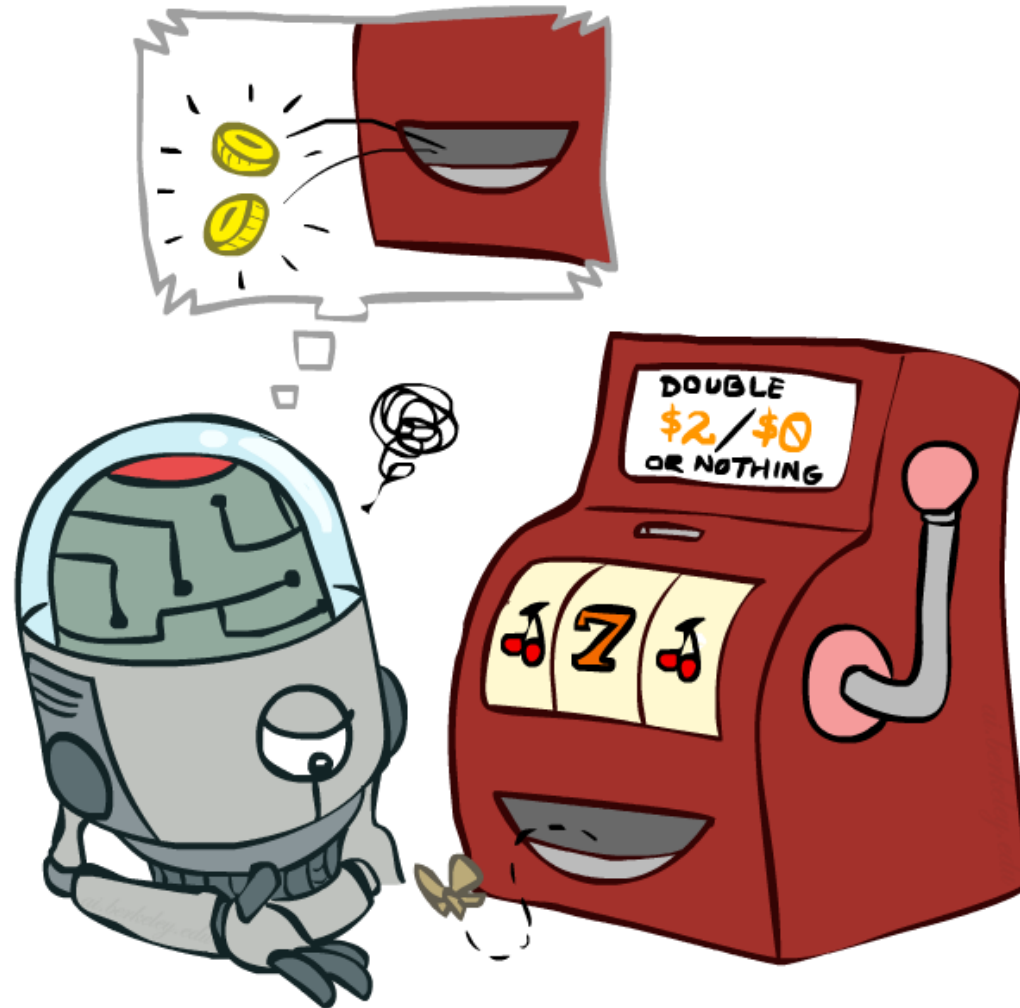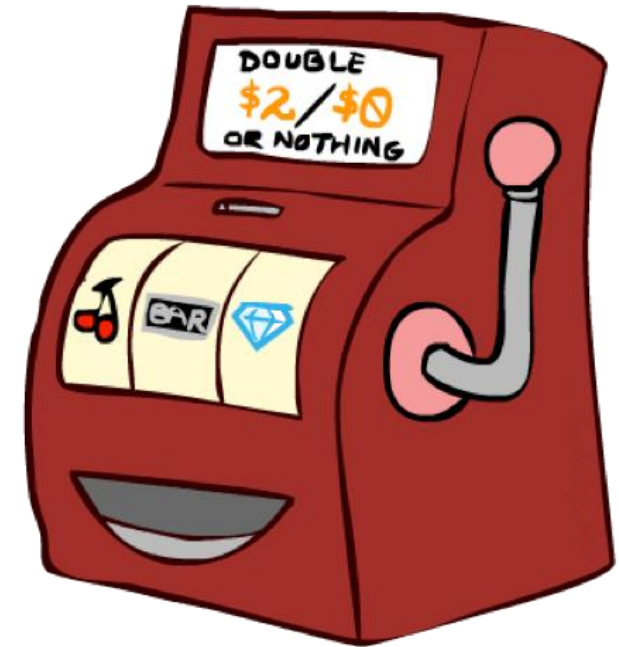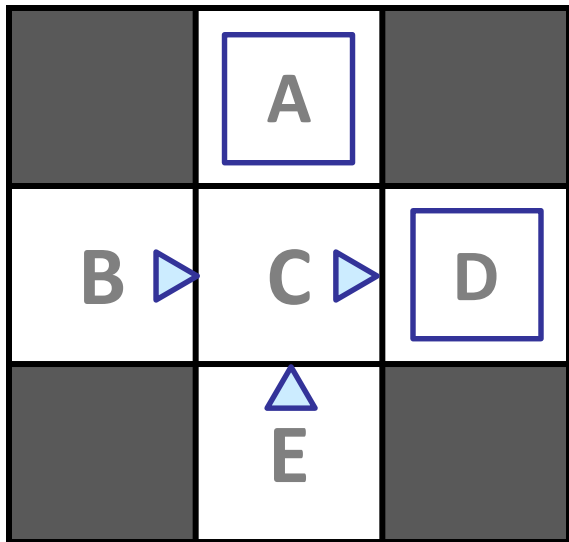
E, north, C, -1
C, east, A, -1
A, exit, x, -10

## Output Values



-10
A
+8  +4  +10
B   C   D
-2
E

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

# Problems with Direct Evaluation

- What's good about direct evaluation?
  - It's easy to understand
  - It doesn't require any knowledge of T, R
  - It eventually computes the correct average values, using just sample transitions

- What bad about it?
  - It wastes information about state connections
  - Each state must be learned separately
  - So, it takes a long time to learn

*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?

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?

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

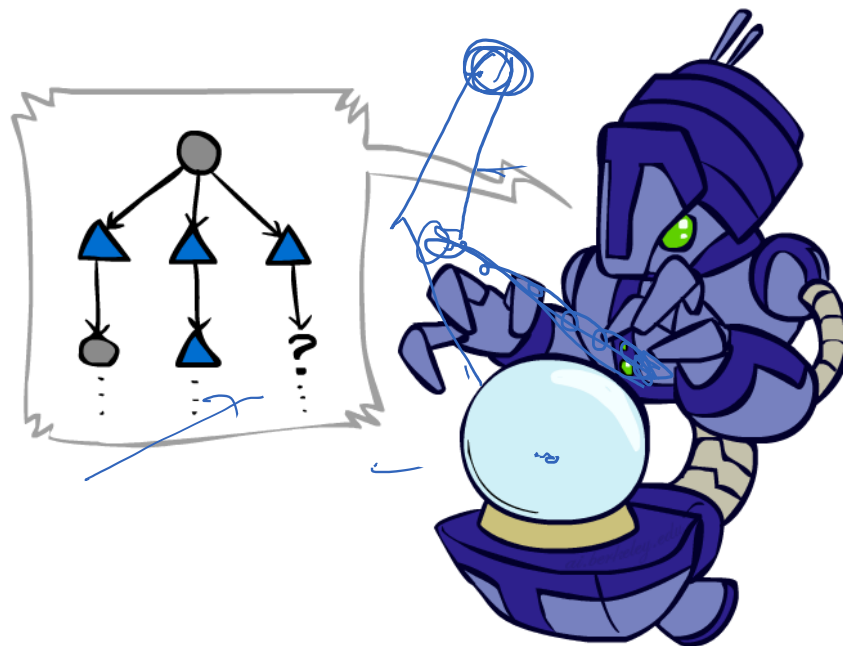$$sample_1 = R(s, \pi(s), s_1') + \gamma V_k^{\pi}(s_1')$$

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

$$\ldots$$

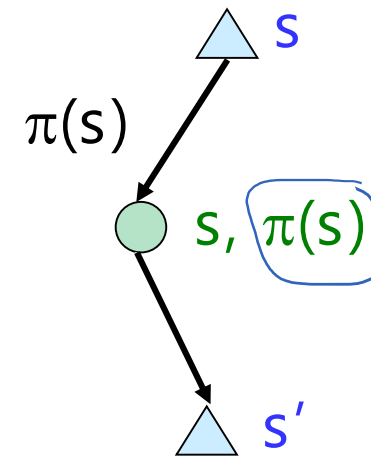$$sample_n = R(s, \pi(s), s_n') + \gamma V_k^{\pi}(s_n')$$

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

from state s.

# Temporal Difference Learning

o Big idea: learn from every experience!
  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

$(1-\alpha)V(s) + (\alpha) sample$

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: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$

  $\left[(1-\alpha)\tilde{x}_{n-2} + \alpha\, x_{n-1}\right] + \alpha \cdot \boxed{x_n}$

  o Makes recent samples more important

  o Forgets about the past (distant past values were wrong anyway)    $+\alpha$

o Decreasing learning rate (alpha) can give converging averages

# Example: Temporal Difference Learning

## States



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

## Observed Transitions

B, east, C, -2

C, east, D, -2



$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha\left[R(s, \pi(s), s') + \gamma V^{\pi}(s')\right]$$
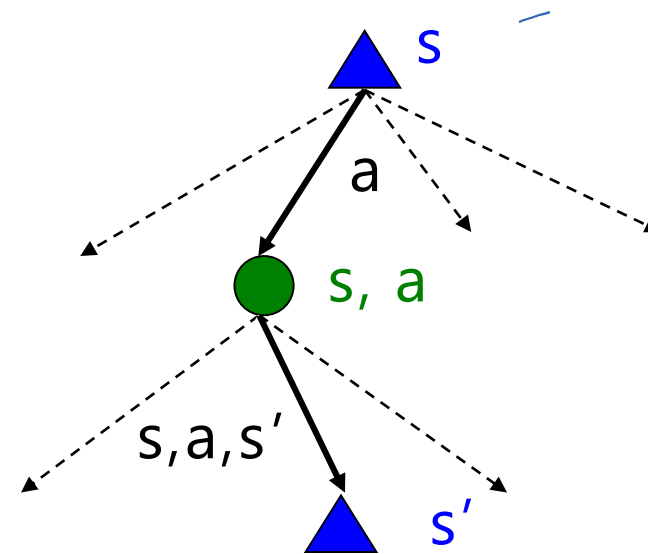
$0$    $+\frac{1}{2}[-2 + 8]$

# 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!

# Announcements

o Project Proposal: Feb 11th

o Paper report: Feb 18th

o PS3: Feb 22nd

o Google Cloud credit is available for you to use.

# Recap: 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
- 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

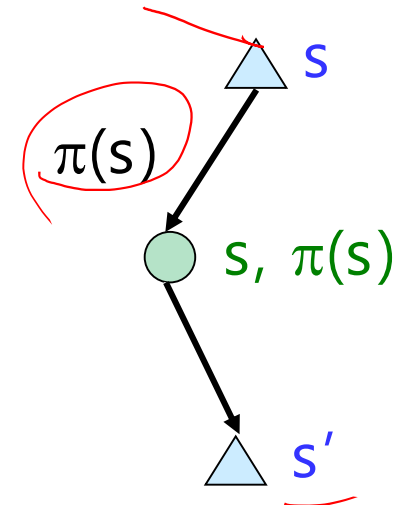| Goal | Technique |
|------|-----------|
| Compute V*, Q*, $\pi$* | Q-learning |
| Evaluate a fixed policy $\pi$ | Value Learning |

# Temporal Difference Learning

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

$\pi(s)$

s

s, $\pi(s)$

s'

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

learning rate

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

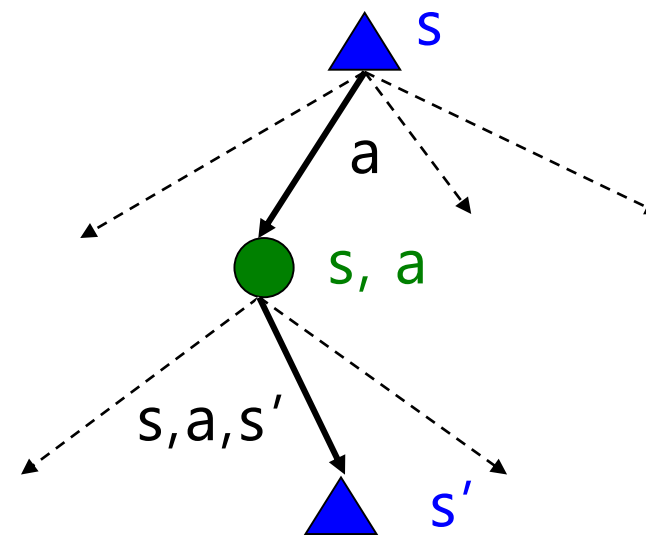Same update: $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$

# 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
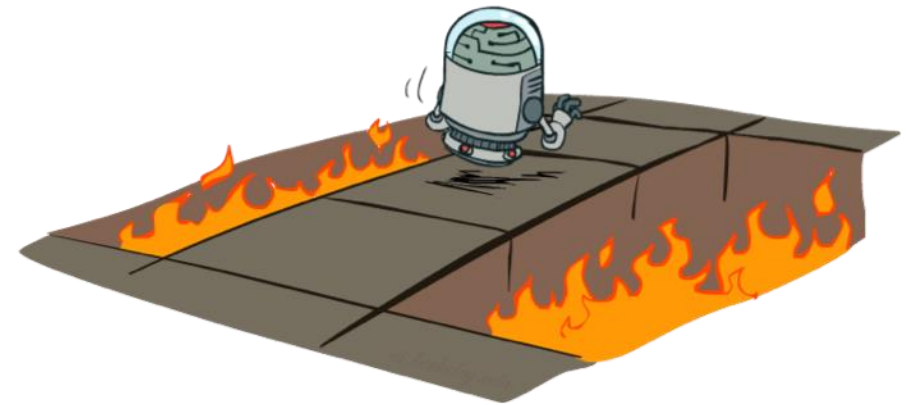- Makes action selection model-free too!
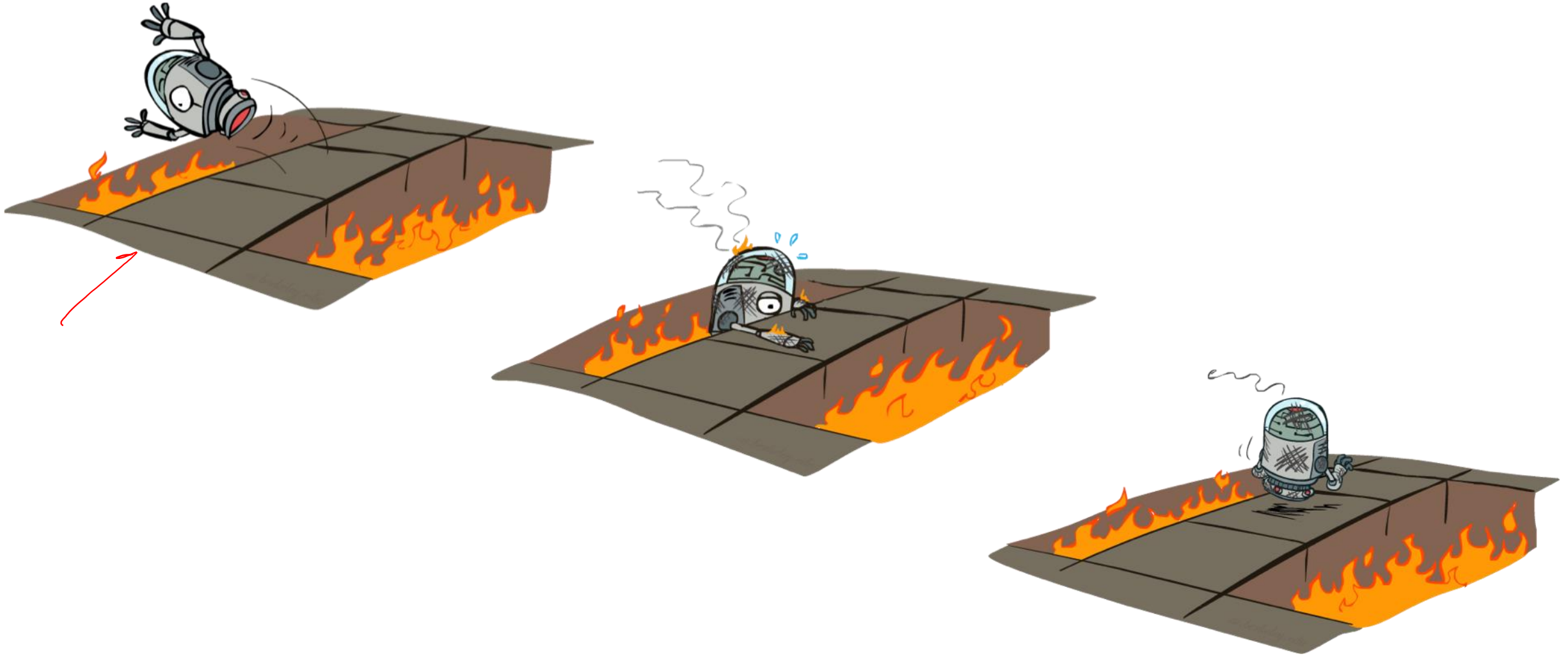
# Discussion: Model-Based vs Model-Free RL

o Model-Based vs. Model Free

TD value

o Active vs. Passive



o Active Reinforcement Learning:
  o act according to current optimal (based on Q-Values)
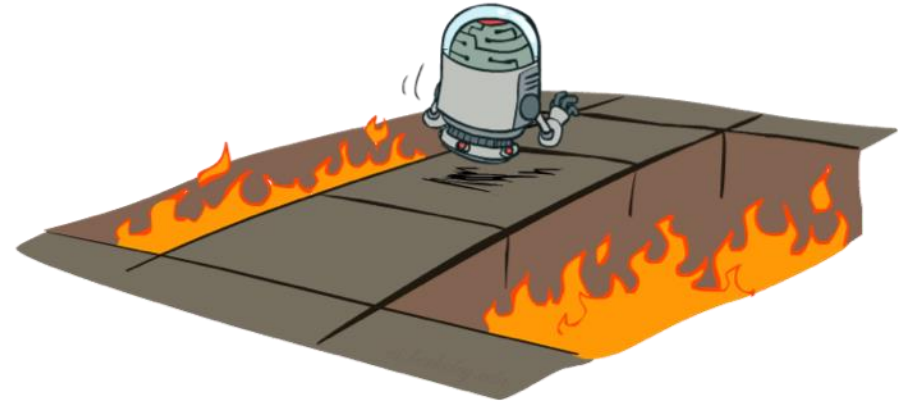  o but also explore…

# Active Reinforcement 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
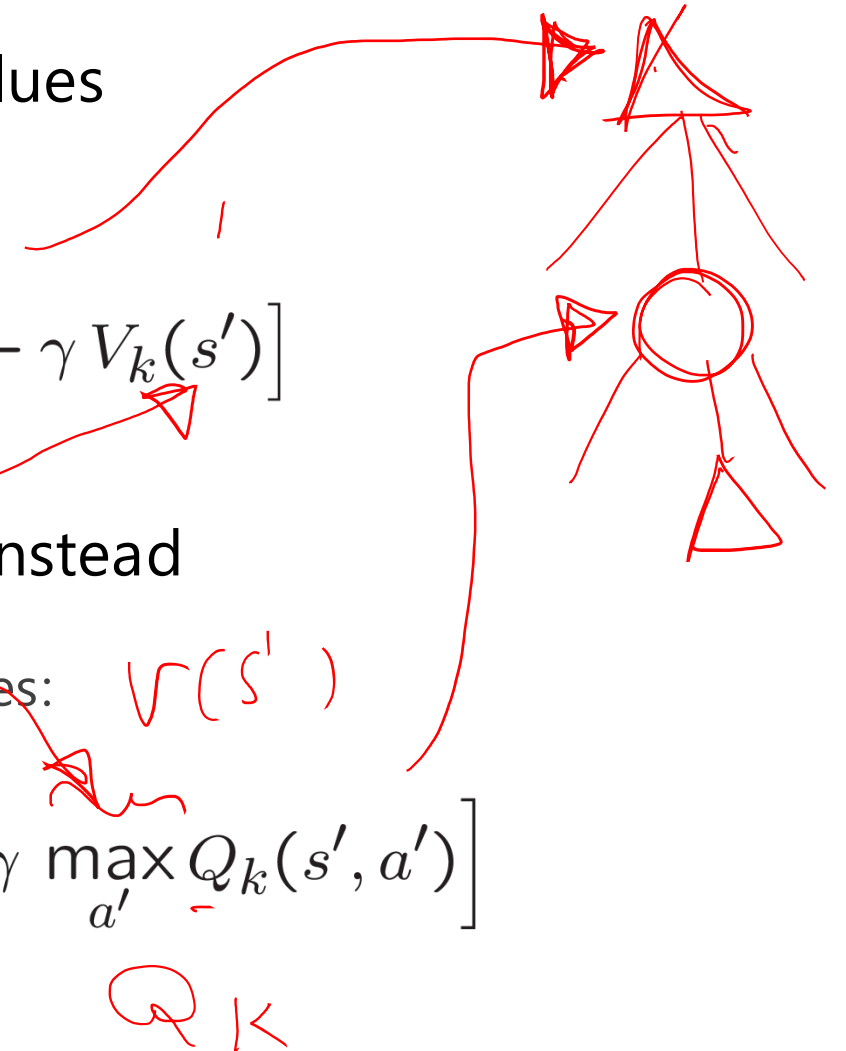- o This is NOT offline planning! You actually take actions in the world and find out what happens...

# Detour: Q-Value Iteration

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

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

- But Q-values are more useful, so compute them instead
  - Start with $Q_0(s,a) = 0$, which we know is right
  - 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

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$
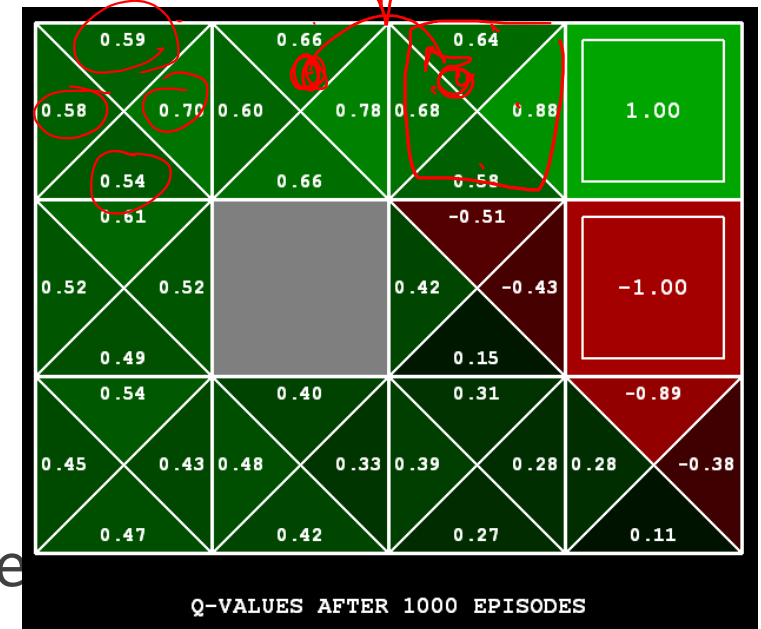
o Learn Q(s,a) values as you go
- o Receive a sample (s,a,s',r)
- o Consider your old estimate $Q(s,a)$
- o Consider your new sample estimate:

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

no longer policy evaluation!
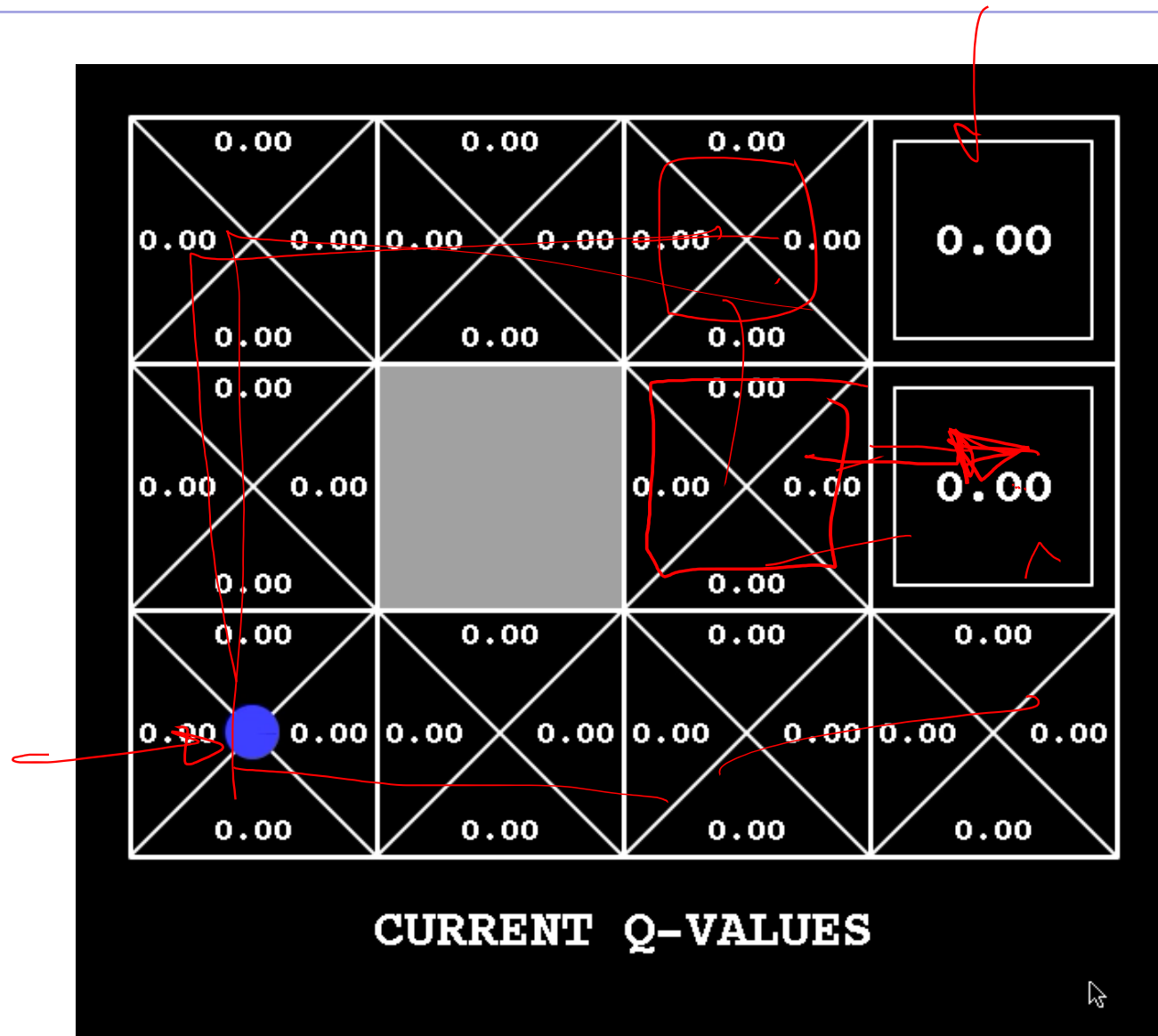
- o Incorporate the new estimate into a running average
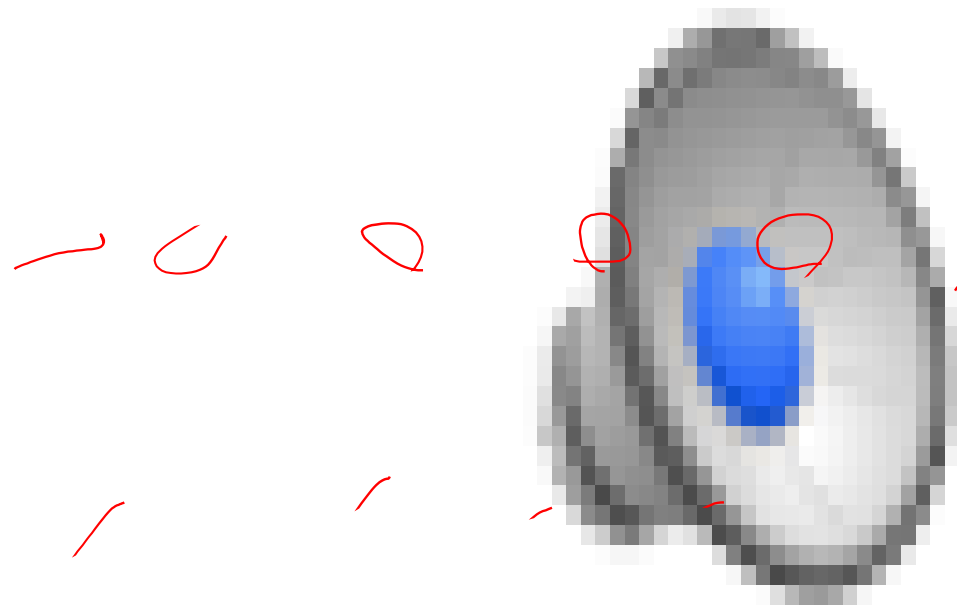
$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + (\alpha)[sample]$$
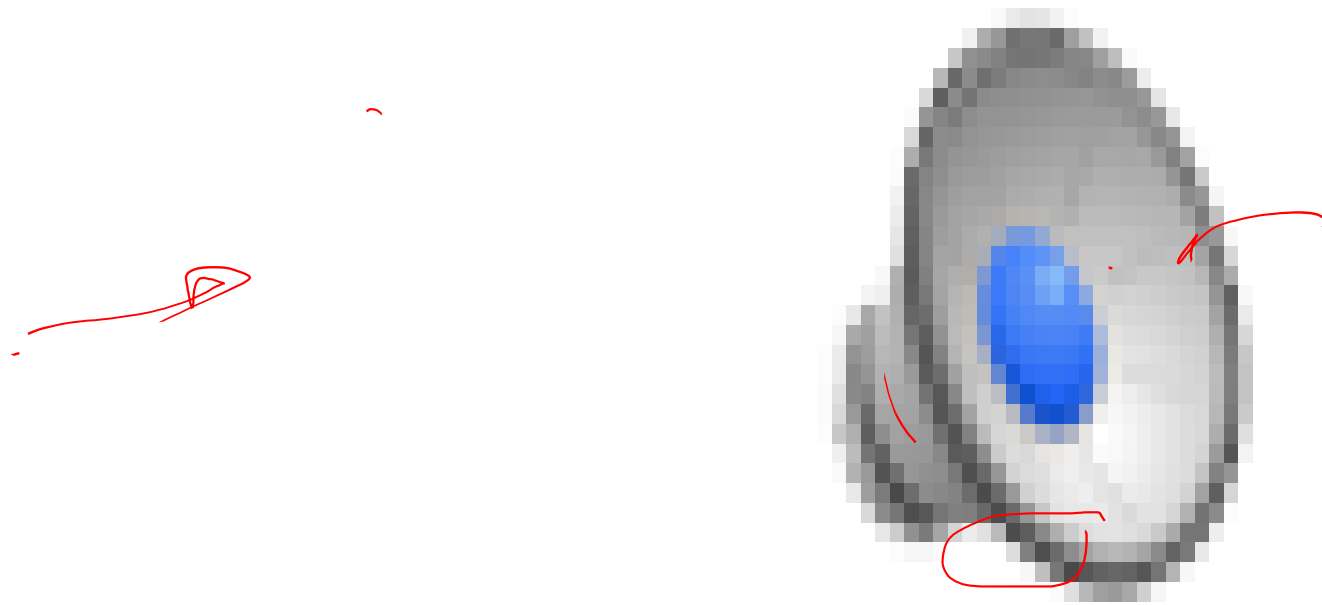


Q-VALUES AFTER 1000 EPISODES

# Q-Learning Demo

# Video of Demo Q-Learning -- Gridworld
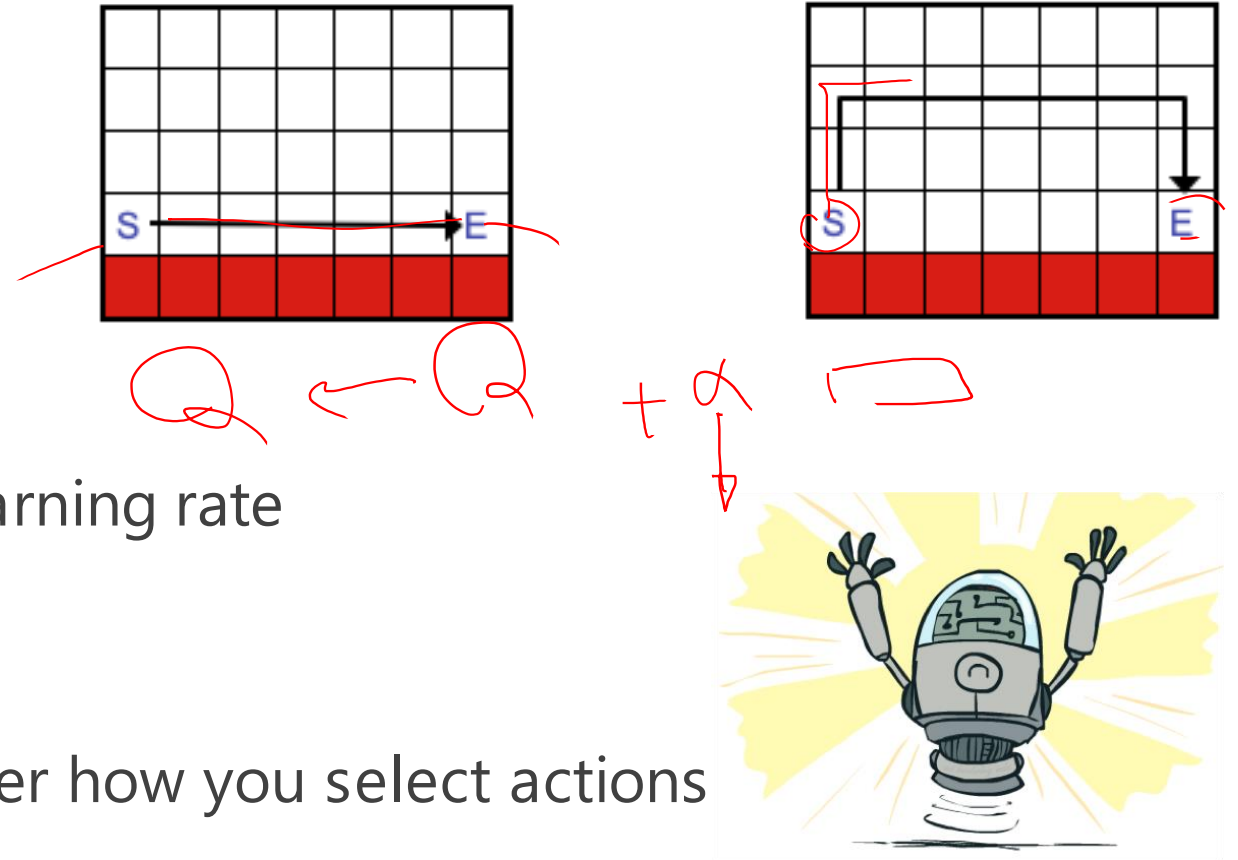
# Video of Demo Q-Learning -- Crawler

# Q-Learning Properties

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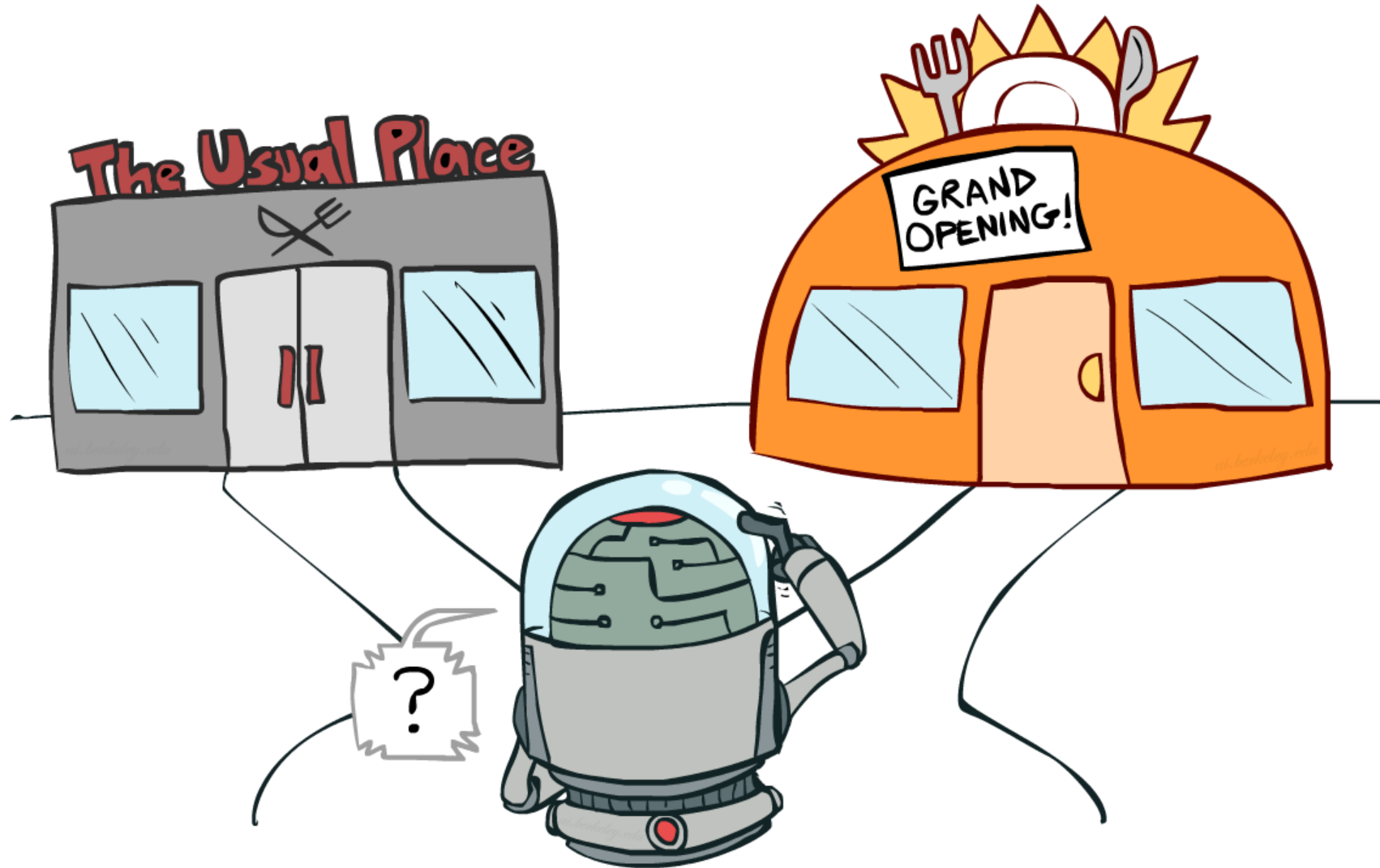even if you're acting suboptimally!

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
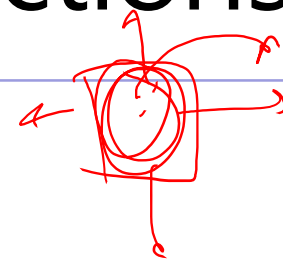
o **Exploration function**

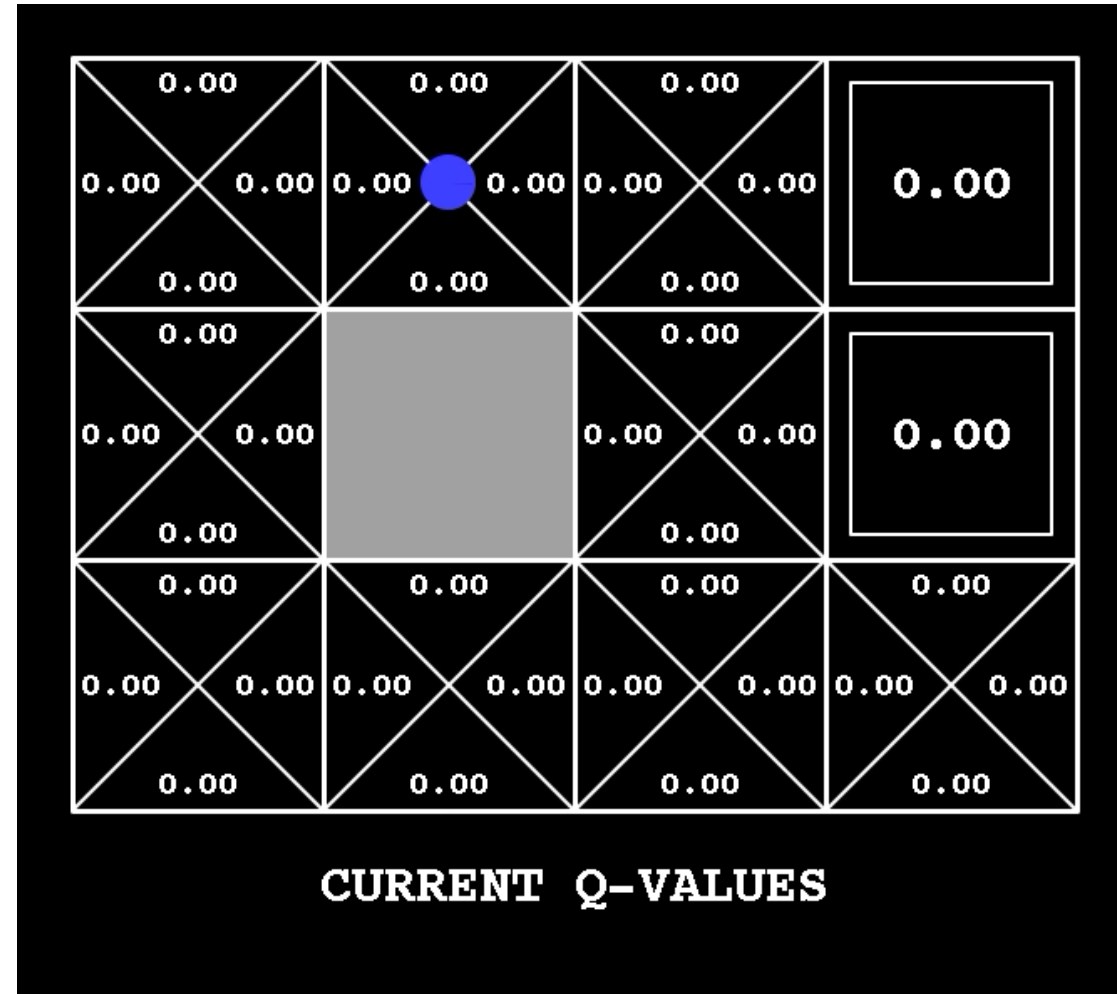   o Takes a value estimate u and a visit count n, and returns an optimistic utility, e.g. $f(u, n) = u + k/n$

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

      Modified Q-Update: $Q(s, a) \leftarrow_\alpha R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$
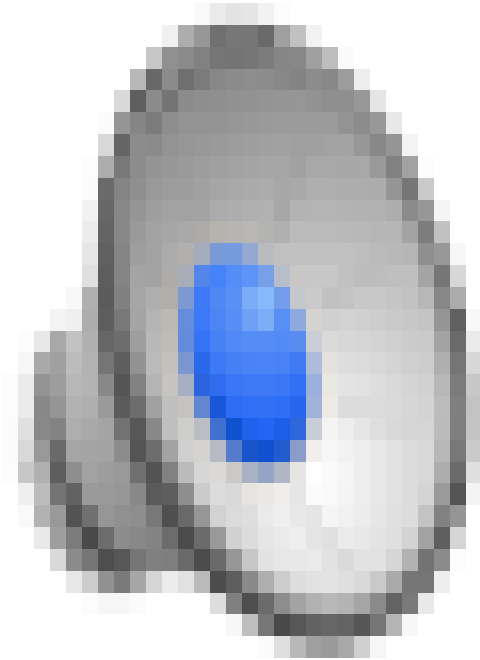
   o Note: this propagates the "bonus" back to states that lead to unknown states as well!
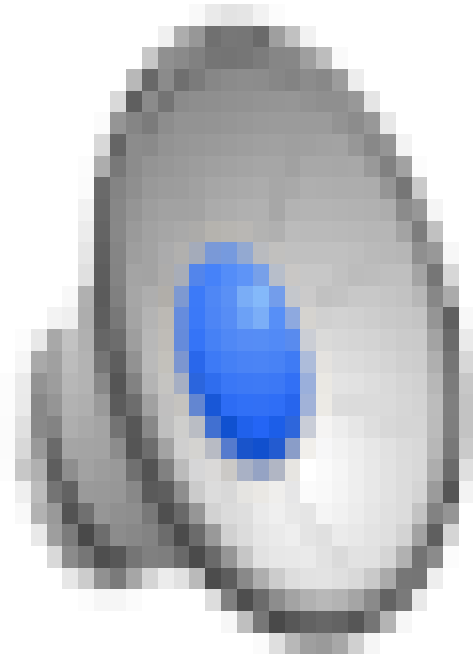
<span style="color:red">[Demo: exploration – Q-learning – crawler – exploration function (L11D4)]</span>

# Q-Learn Epsilon Greedy

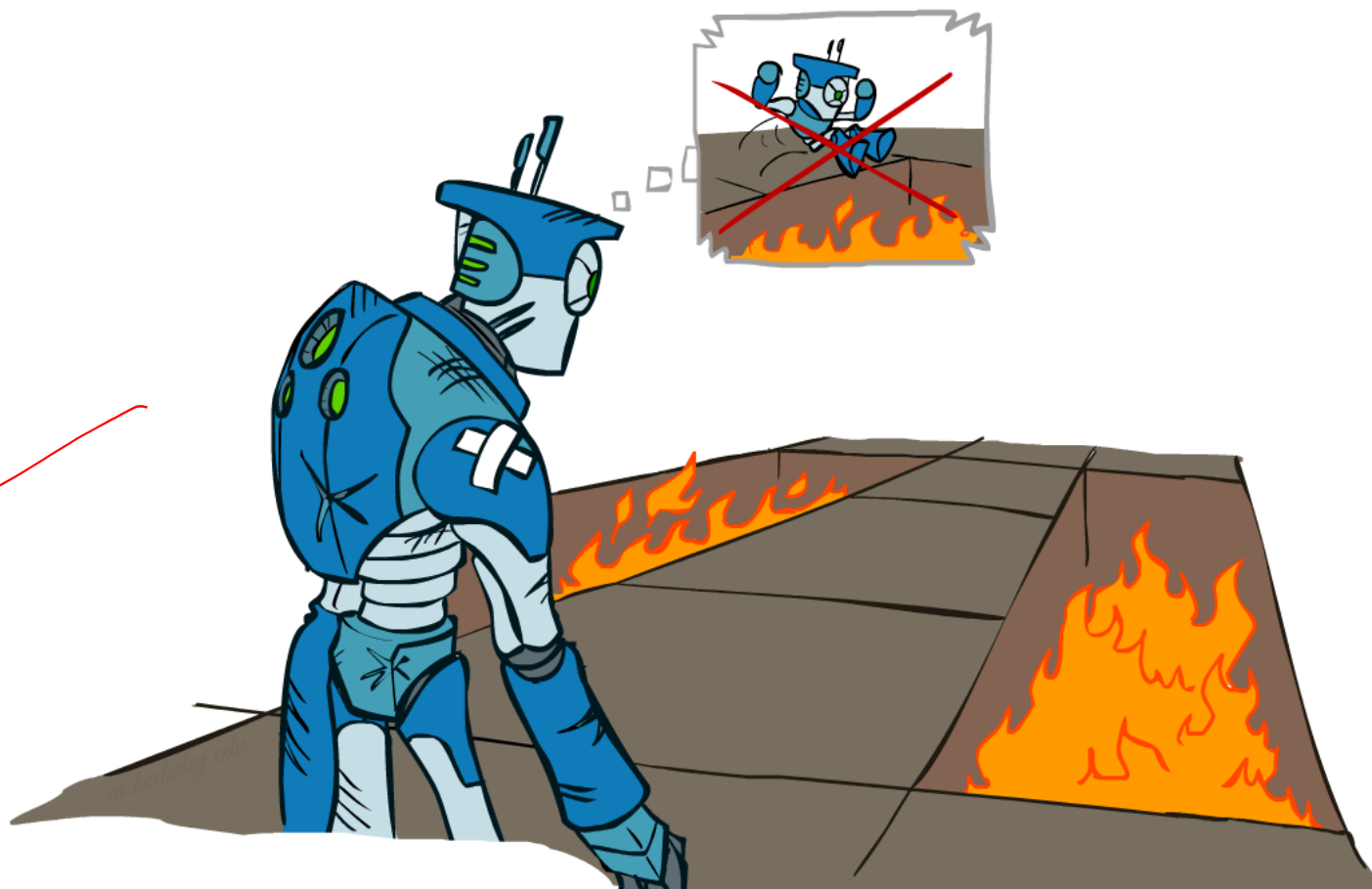# Video of Demo Q-learning – Epsilon-Greedy – Crawler

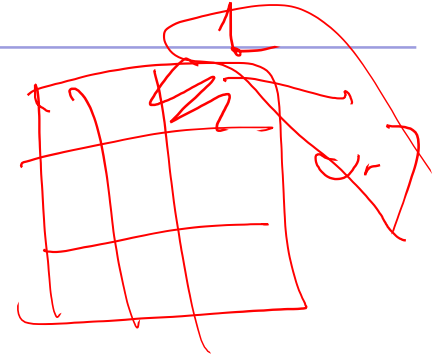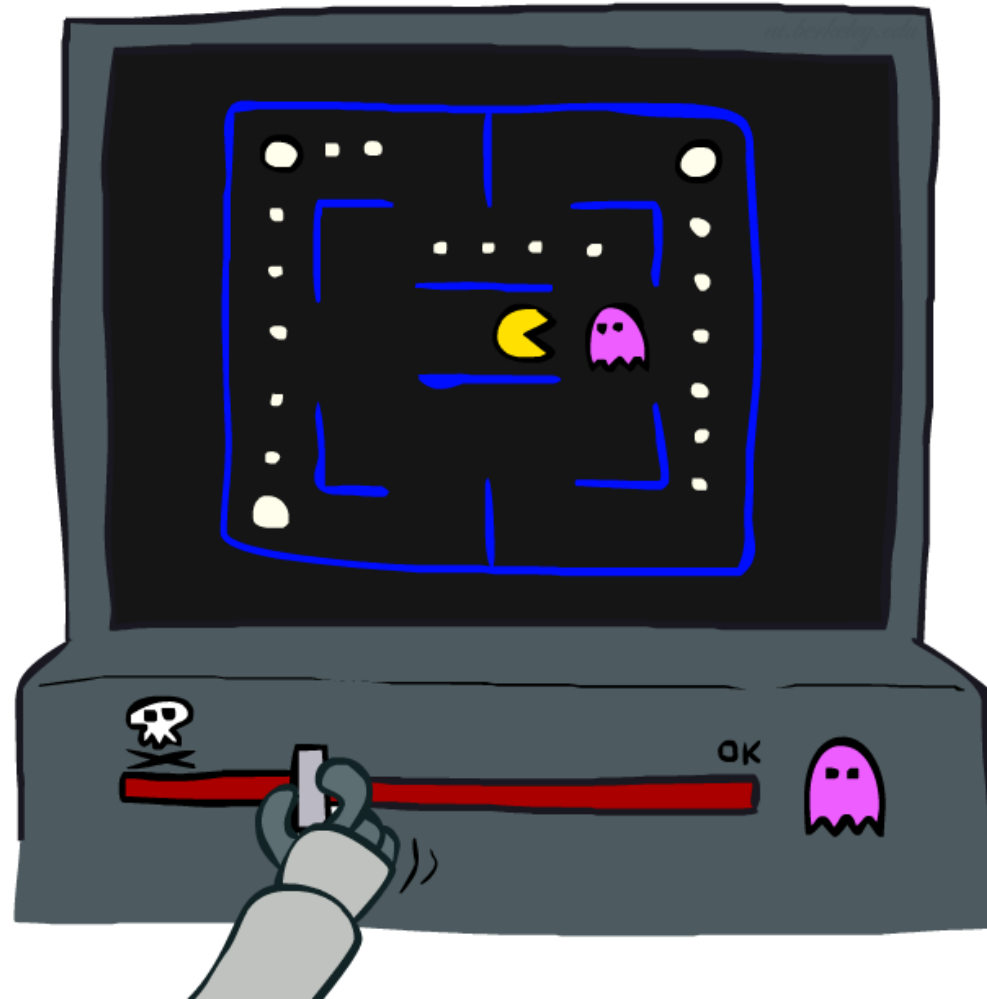# Video of Demo Q-learning – Exploration Function – Crawler

# 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
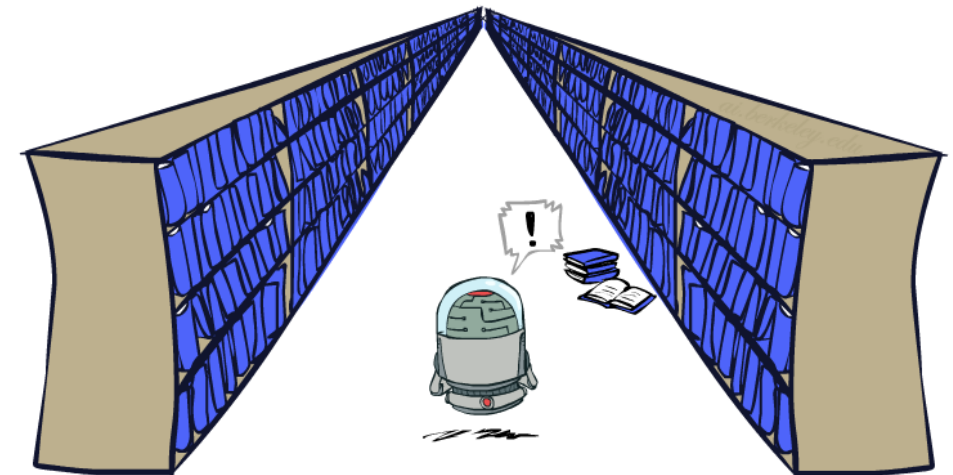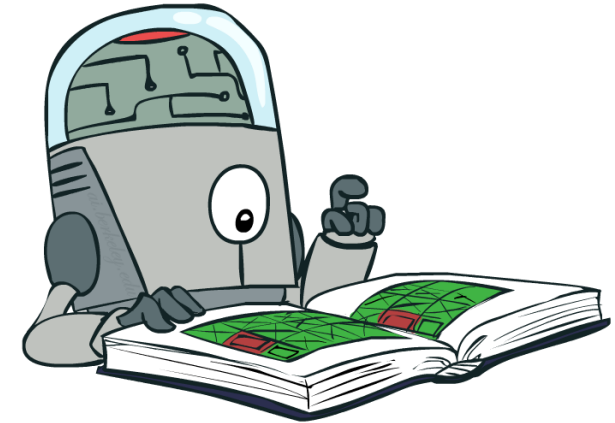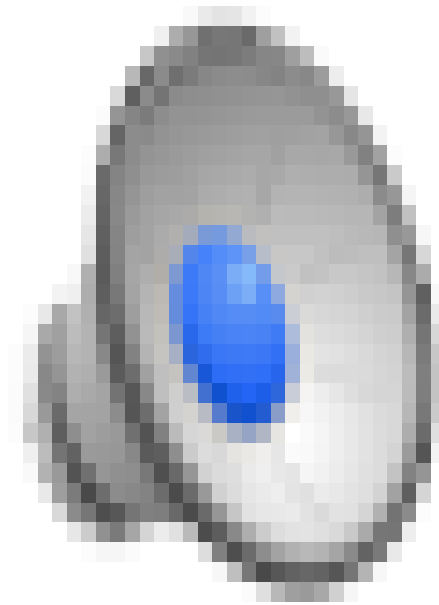
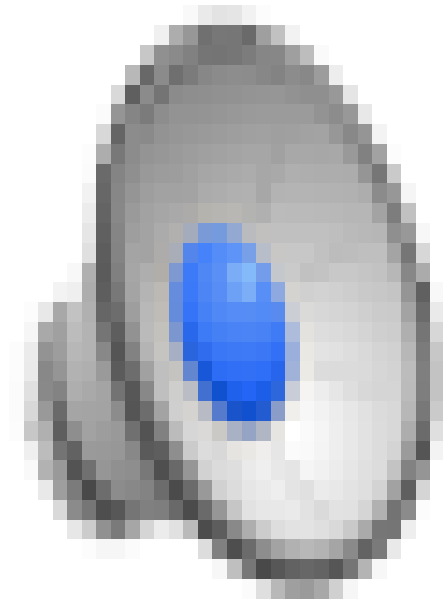# Approximate Q-Learning

# 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
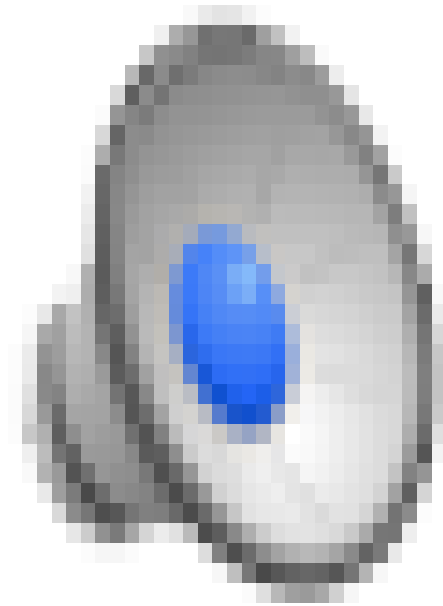
[demo – RL pacman]

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

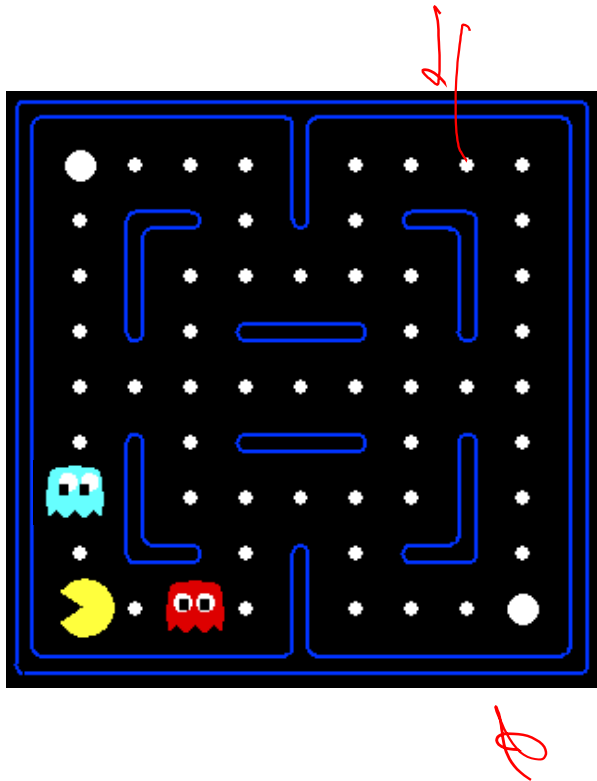# Video of Demo Q-Learning Pacman – Tiny – Silent Train

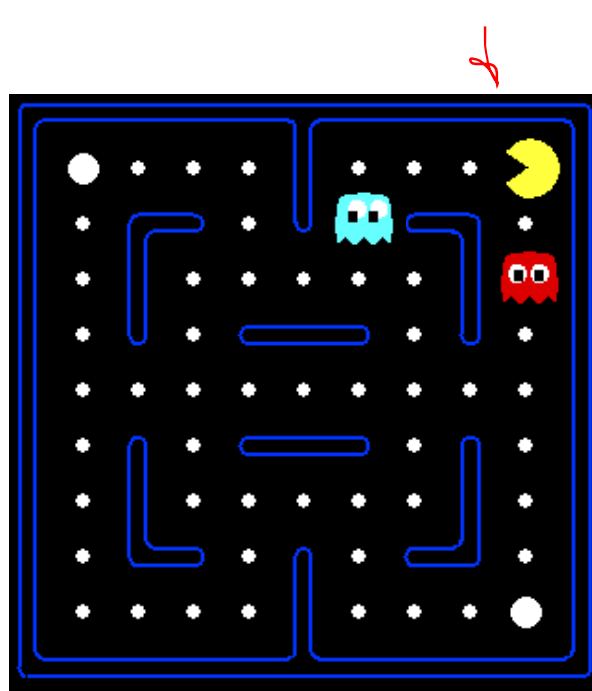# Video of Demo Q-Learning Pacman – Tricky – Watch All
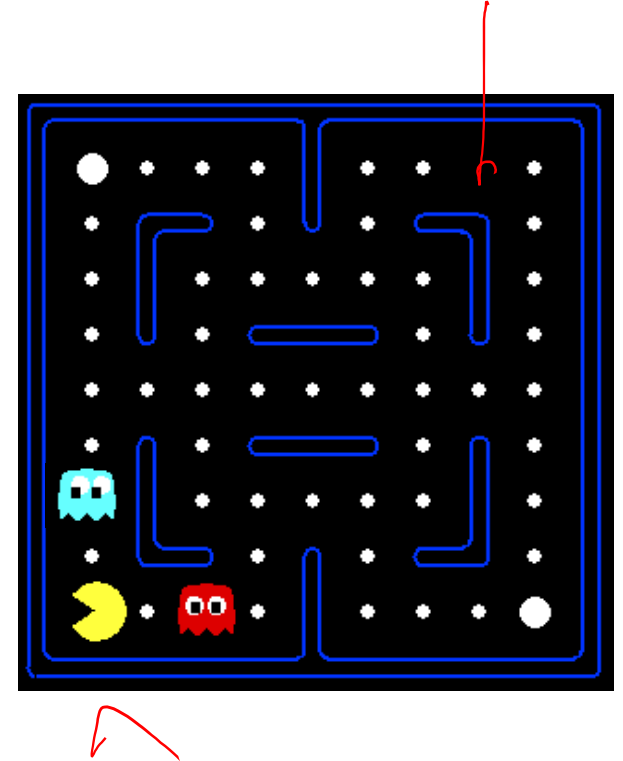
# Example: Pacman

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



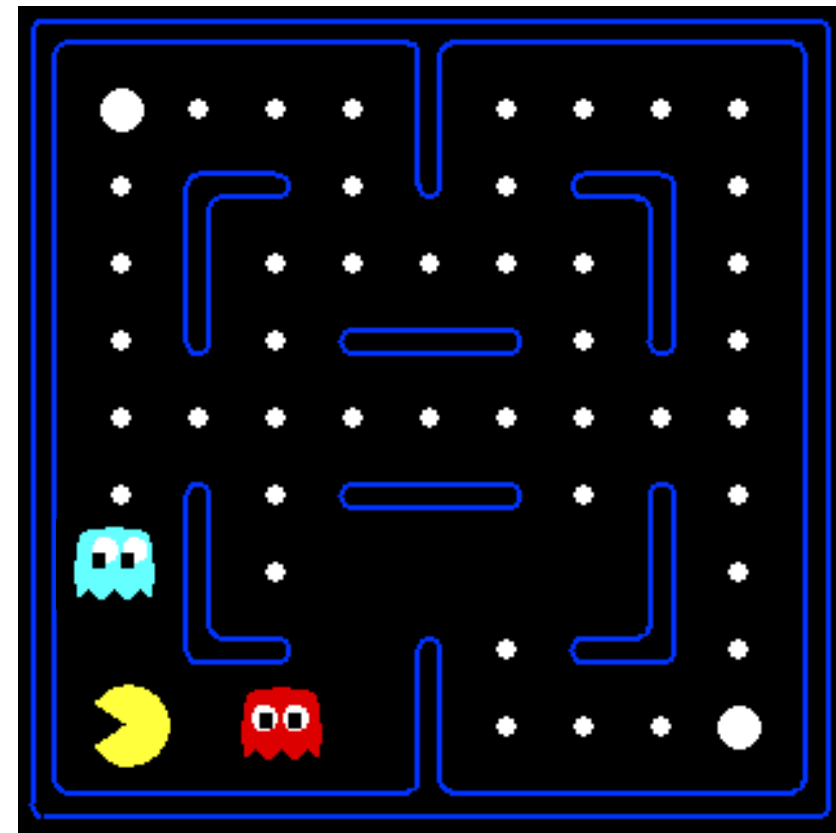In naïve q-learning, we know nothing about this state:



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

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

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

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

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

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

# Approximate Q-Learning

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

- Q-learning with linear Q-functions:

    transition $= (s, a, r, s')$

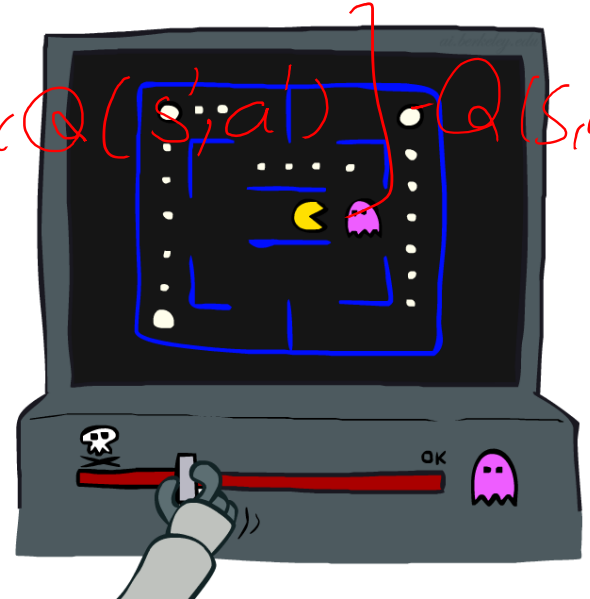    difference $= \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$

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

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

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

    Exact Q's

    Approximate Q's
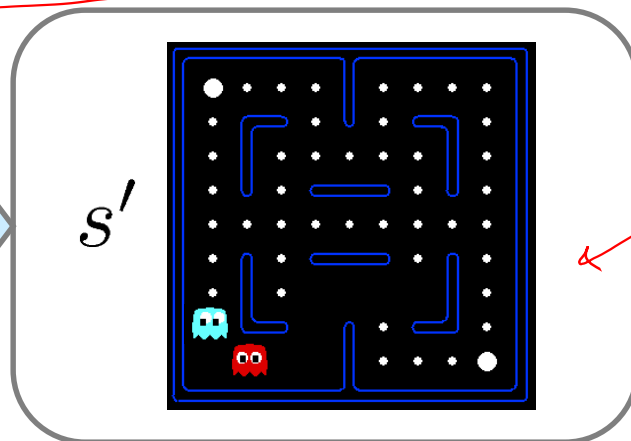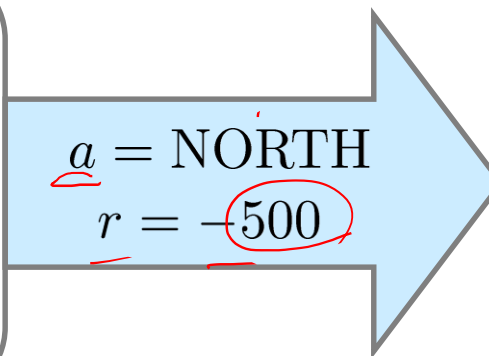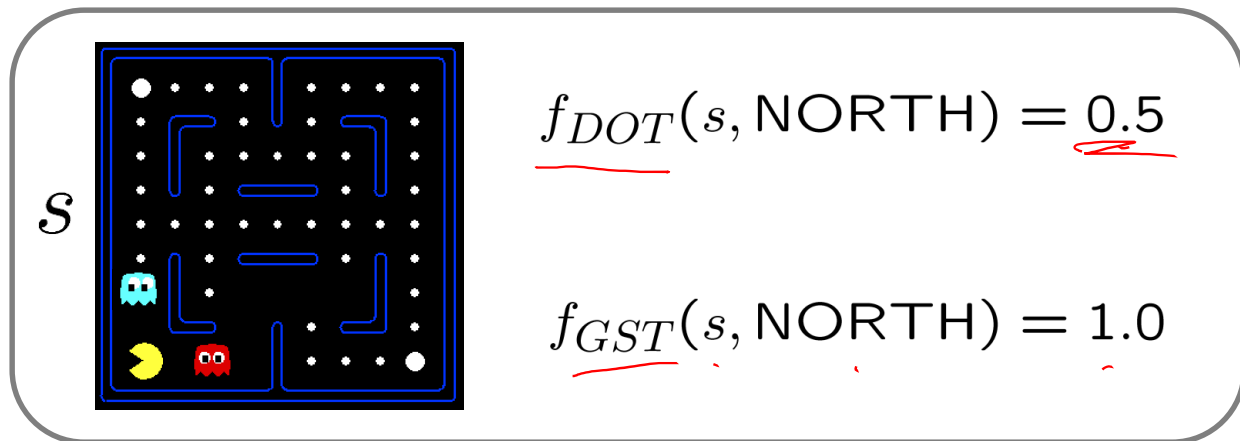
- Intuitive interpretation:
    - Adjust weights of active features
    - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

- Formal justification: online least squares

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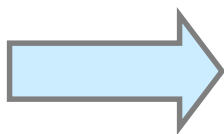
$$Q(s,a) = 4.0 f_{DOT}(s,a) - 1.0 f_{GST}(s,a)$$

$s$

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

$a = \text{NORTH}$
$r = -500$

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

$s'$

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

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

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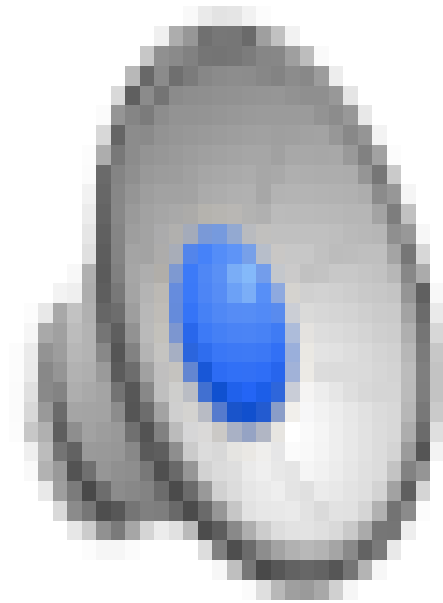
$\text{difference} = -501$

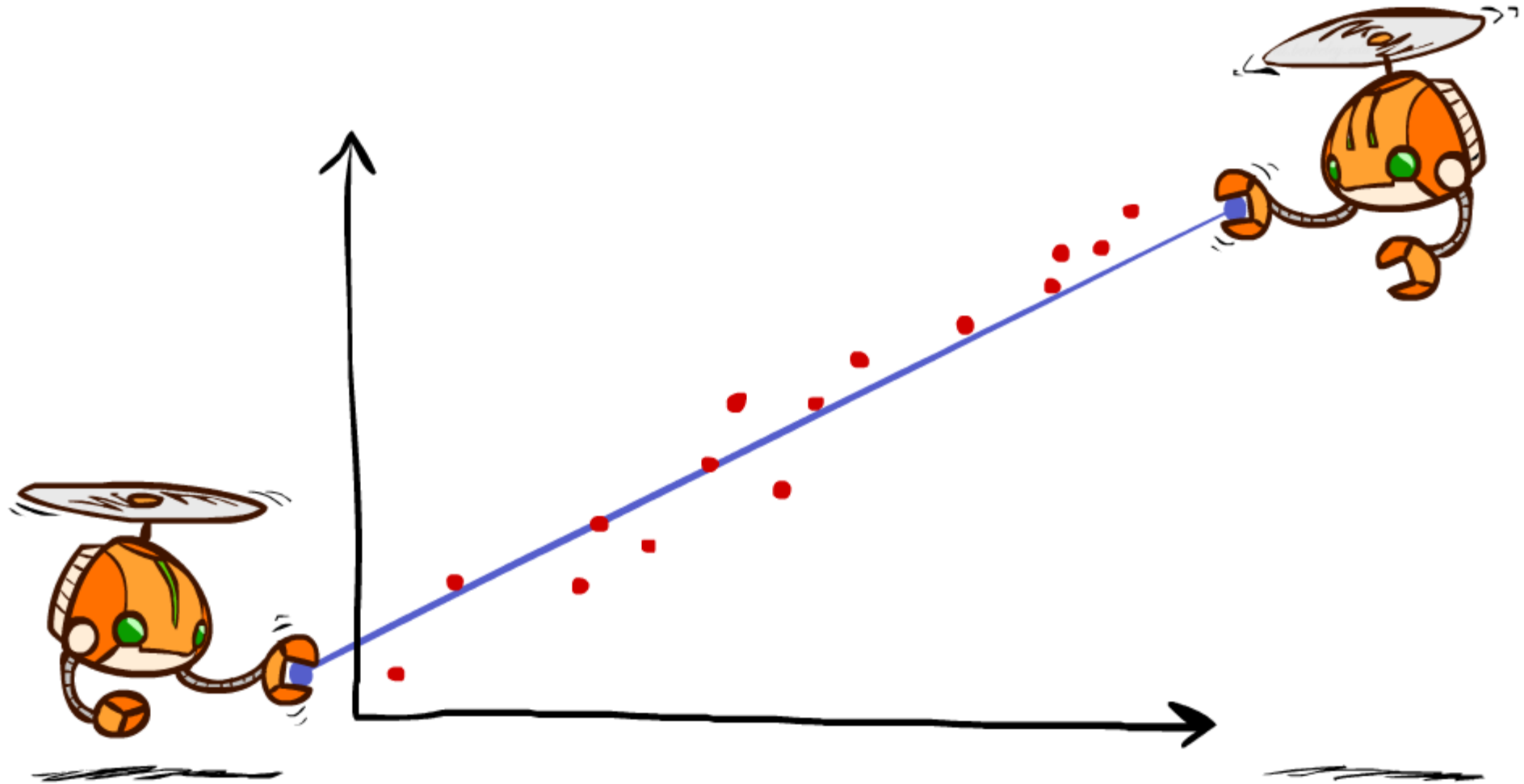$w_{DOT} \leftarrow 4.0 + \alpha [-501] \, 0.5$
$w_{GST} \leftarrow -1.0 + \alpha [-501] \, 1.0$
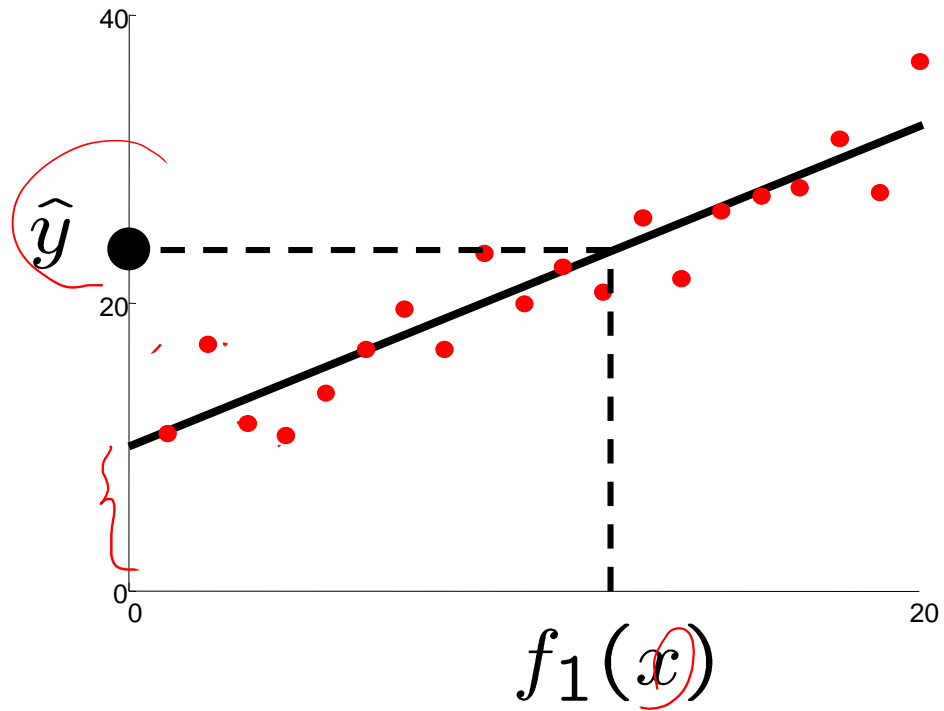
$$Q(s,a) = 3.0 f_{DOT}(s,a) - 3.0 f_{GST}(s,a)$$

# Video of Demo Approximate
# Q-Learning -- Pacman

# Linear Approximation: Regression



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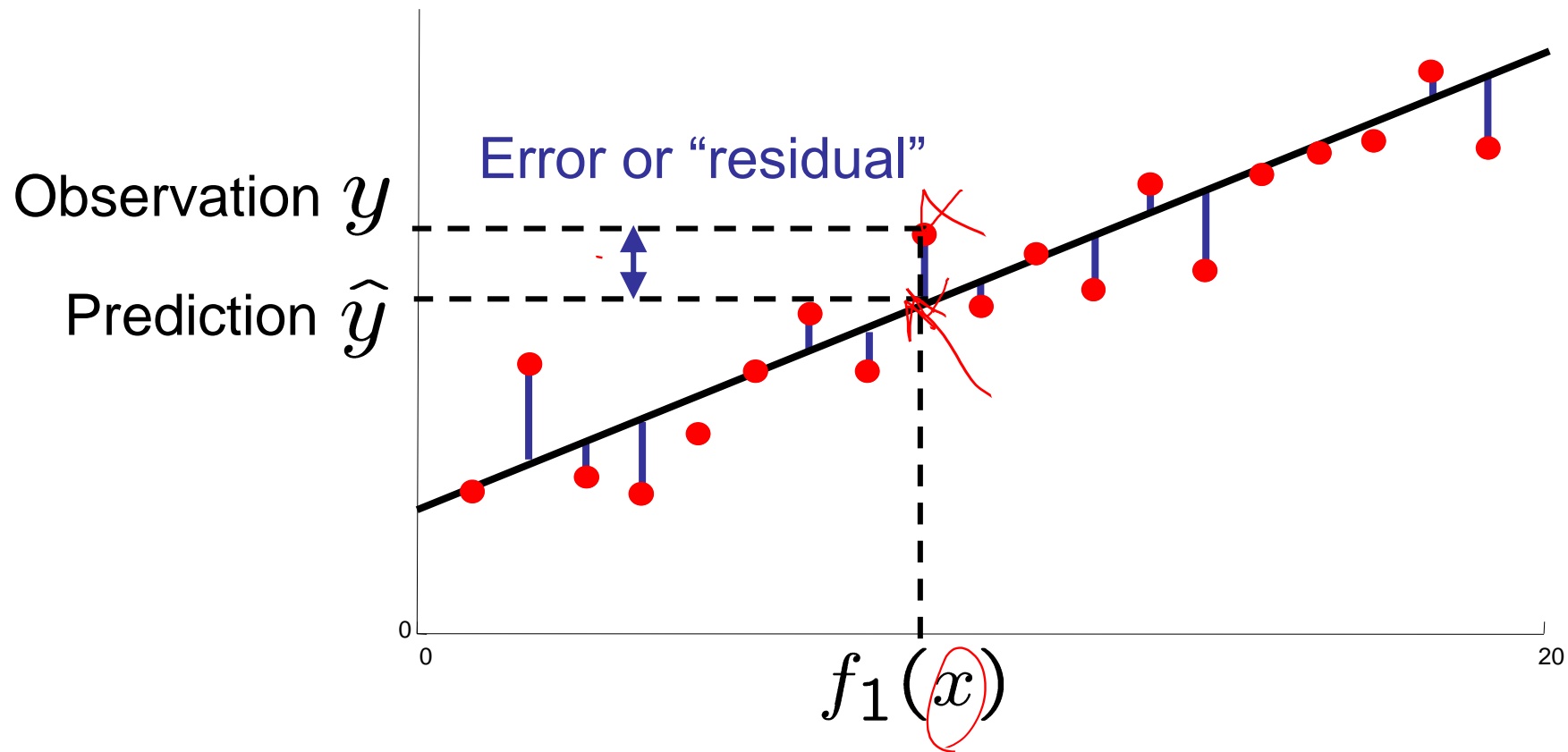
Prediction:
$$\hat{y}_i = w_0 + w_1 f_1(x) + w_2 f_2(x)$$

# Optimization: Least Squares

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



Error or "residual"

Observation $y$

Prediction $\widehat{y}$

$f_1(x)$

# 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"

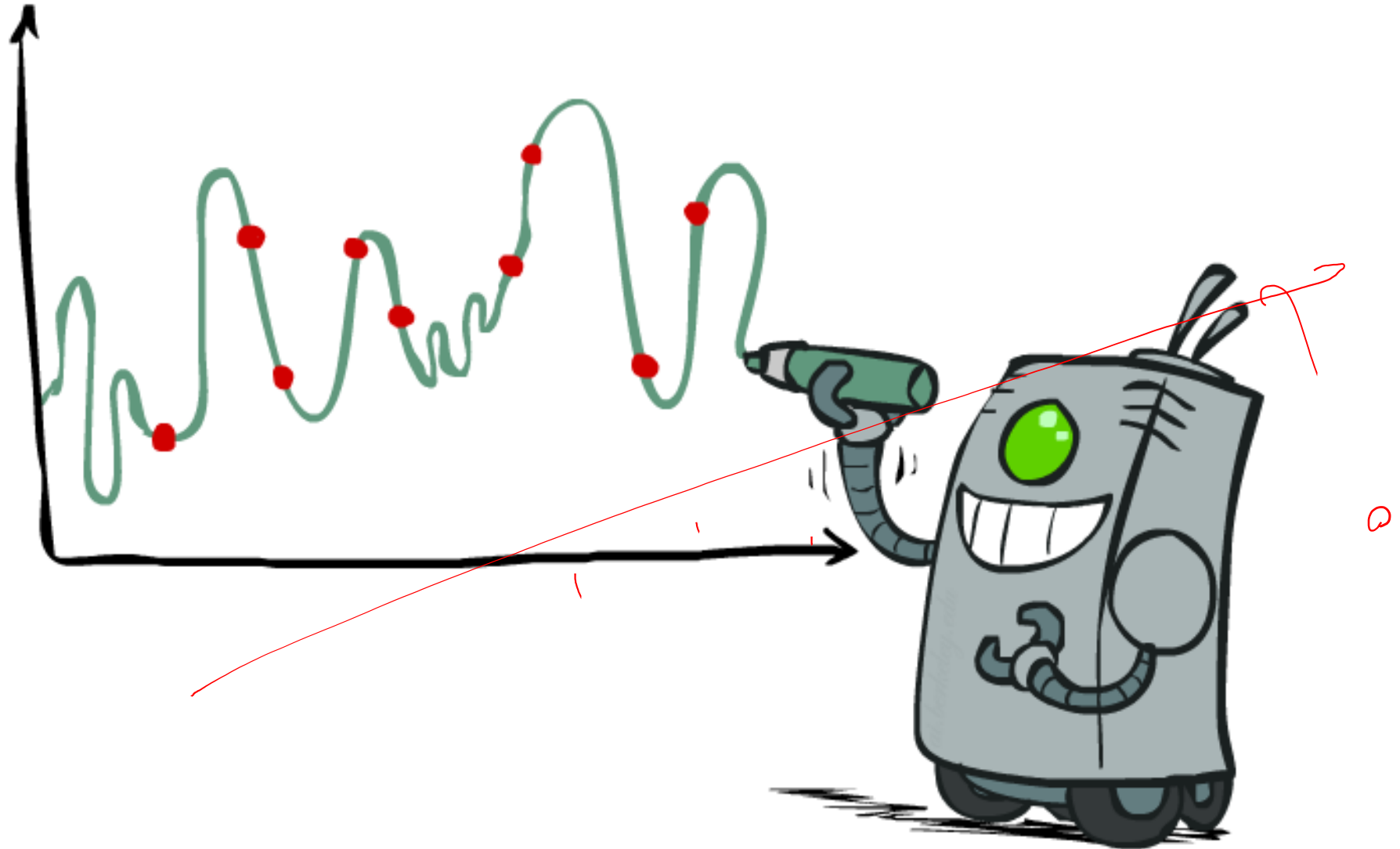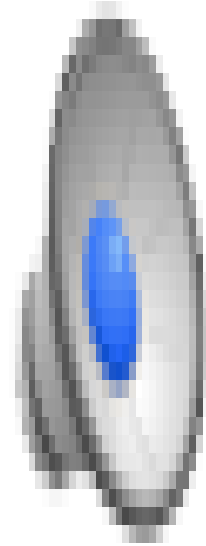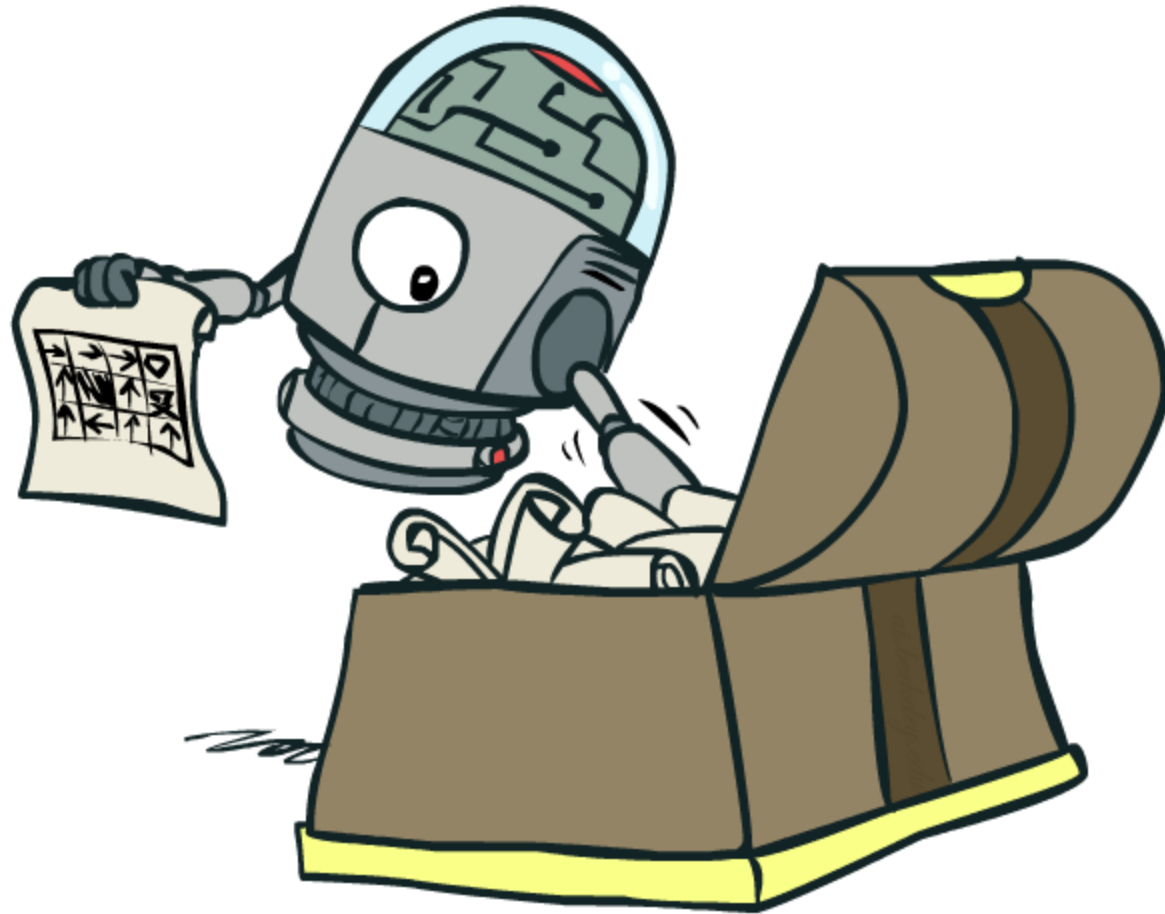# Overfitting: Why Limiting Capacity Can Help

# 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 We'll see this distinction between modeling and prediction again later in the course

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*, $\pi$* | Value / policy iteration |
| Evaluate a fixed policy $\pi$ | Policy evaluation |

## Unknown MDP: Model-Based

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

## Unknown MDP: Model-Free

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

# Conclusion

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

o Next up: Uncertainty and Learning!