# CSE 473: Introduction to Artificial Intelligence
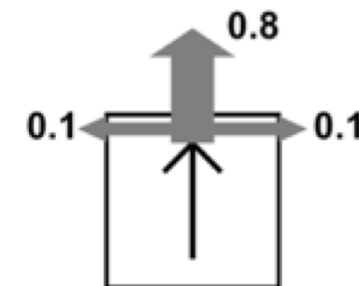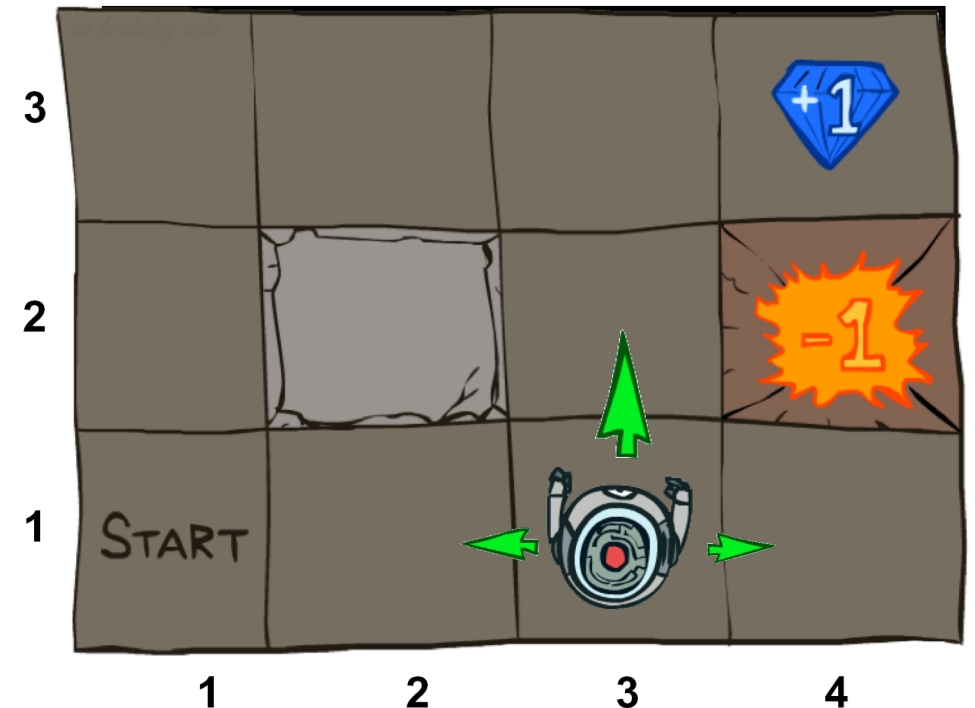
Hanna Hajishirzi

Markov Decision Processes

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

# Example: Grid World

- A maze-like problem
  - The agent lives in a grid
  - Walls block the agent's path

- Noisy movement: actions do not always go as planned
  - 80% of the time, the action North takes the agent North
    (if there is no wall there)
  - 10% of the time, North takes the agent West; 10% East
  - If there is a wall in the direction the agent would have been taken, the agent stays put

- The agent receives rewards each time step
  - Small "living" reward each step (can be negative)
  - Big rewards come at the end (good or bad)

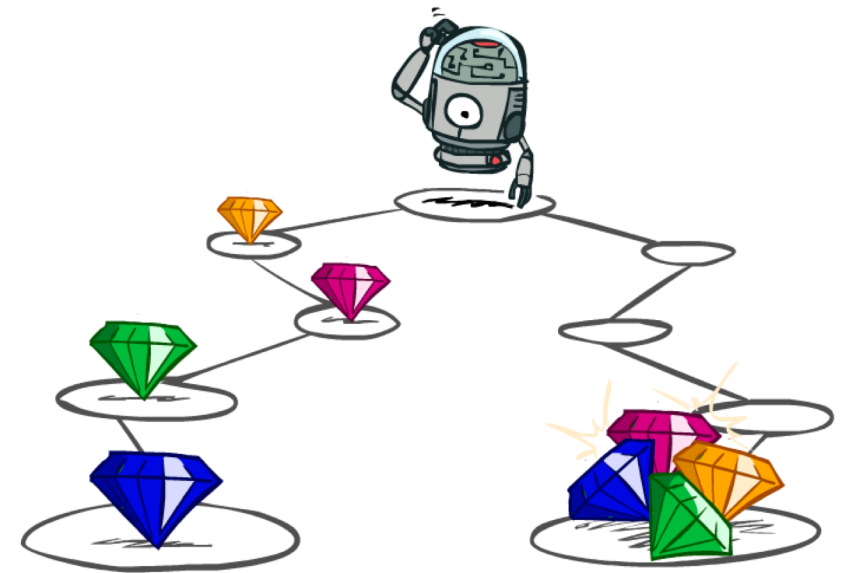- Goal: maximize sum of rewards

# Recap: Defining MDPs

- Markov decision processes:
  - Set of states S
  - Start state $s_0$
  - Set of actions A
  - Transitions $P(s'|s,a)$ (or $T(s,a,s')$)
  - Rewards $R(s,a,s')$ (and discount $\gamma$)

- MDP quantities so far:
  - Policy = Choice of action for each state
  - Utility = sum of (discounted) rewards

# Utilities of Sequences

o What preferences should an agent have over reward sequences?

o More or less?   [1, 2, 2]      or      [2, 3, 4]

o Now or later?  [0, 0, 1]      or      [1, 0, 0]

# Discounting

- It's reasonable to maximize the sum of rewards
- It's also reasonable to prefer rewards now to rewards later
- One solution: values of rewards decay exponentially



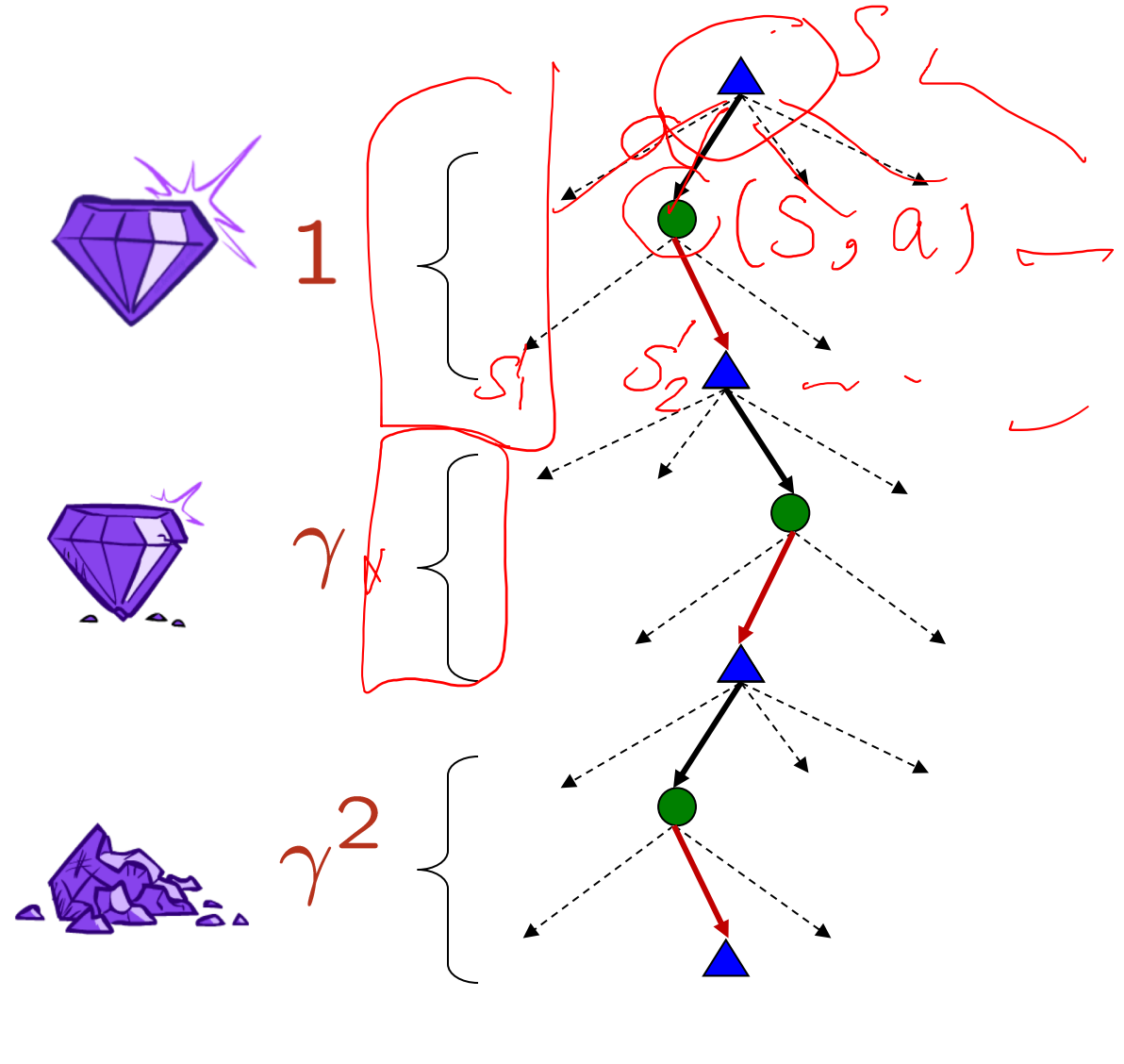| 1 | $\gamma$ | $\gamma^2$ |
|---|---|---|
| Worth Now | Worth Next Step | Worth In Two Steps |

# Discounting

o How to discount?
  o Each time we descend a level,
    we multiply in the discount once

o Why discount?
  o Think of it as a gamma chance of
    ending the process at every step
  o Also helps our algorithms
    converge

o Example: discount of 0.5
  o U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
  o U([1,2,3]) < U([3,2,1])

# Quiz: Discounting

o Given:

| 10 | | | | 1 |
|---|---|---|---|---|
| a | b | c | d | e |

  o Actions: East, West, and Exit (only available in exit states a, e)

  o Transitions: deterministic

o Quiz 1: For γ = 1, what is the optimal policy?

| 10 | <- | <- | <- | 1 |
|---|---|---|---|---|

o Quiz 2: For γ = 0.1, what is the optimal policy?

| 10 | <- | <- | -> | 1 |
|---|---|---|---|---|

o Quiz 3: For which γ are West and East equally good when in state d?

$1\gamma = 10\,\gamma^3$  $\quad \gamma^3 \times 10 = \gamma$
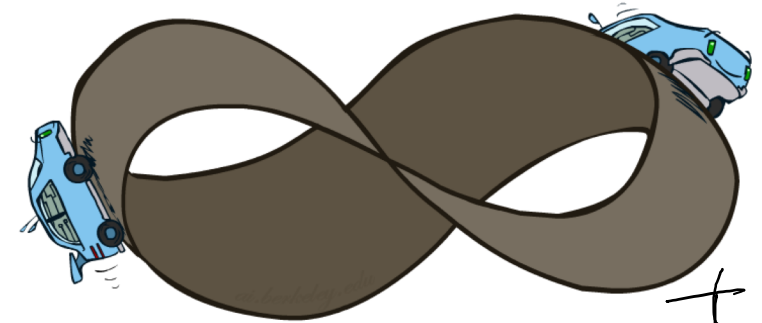
# Infinite Utilities?!

- Problem: What if the game lasts forever?  Do we get infinite rewards?

- Solutions:
  - Finite horizon: (similar to depth-limited search)
    - Terminate episodes after a fixed T steps (e.g. life)
    - Policy π depends on time left
  - Discounting: use $0 < \gamma < 1$

    $$U([r_0, \ldots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{max}/(1-\gamma)$$

    $r_0 + \gamma r_1 + \cdots + \gamma^t r_t$
    $\leq R_{max}(1 + \gamma + \cdots + \gamma^t)$

    - Smaller γ means smaller "horizon" – shorter term focus
  - Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "overheated" for racing)
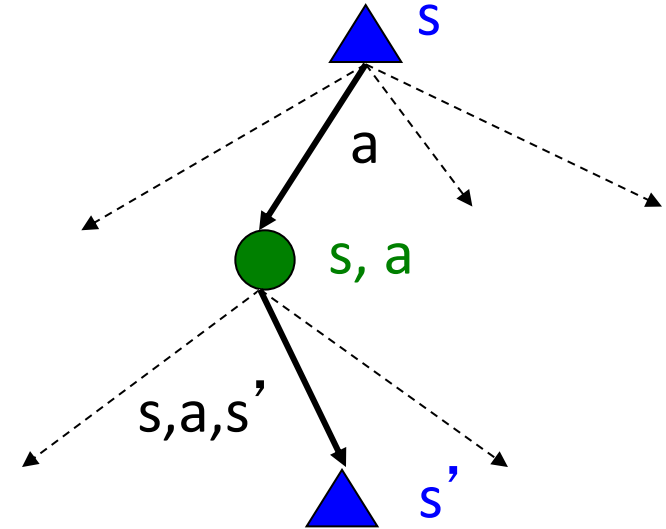
# MDP Search Trees

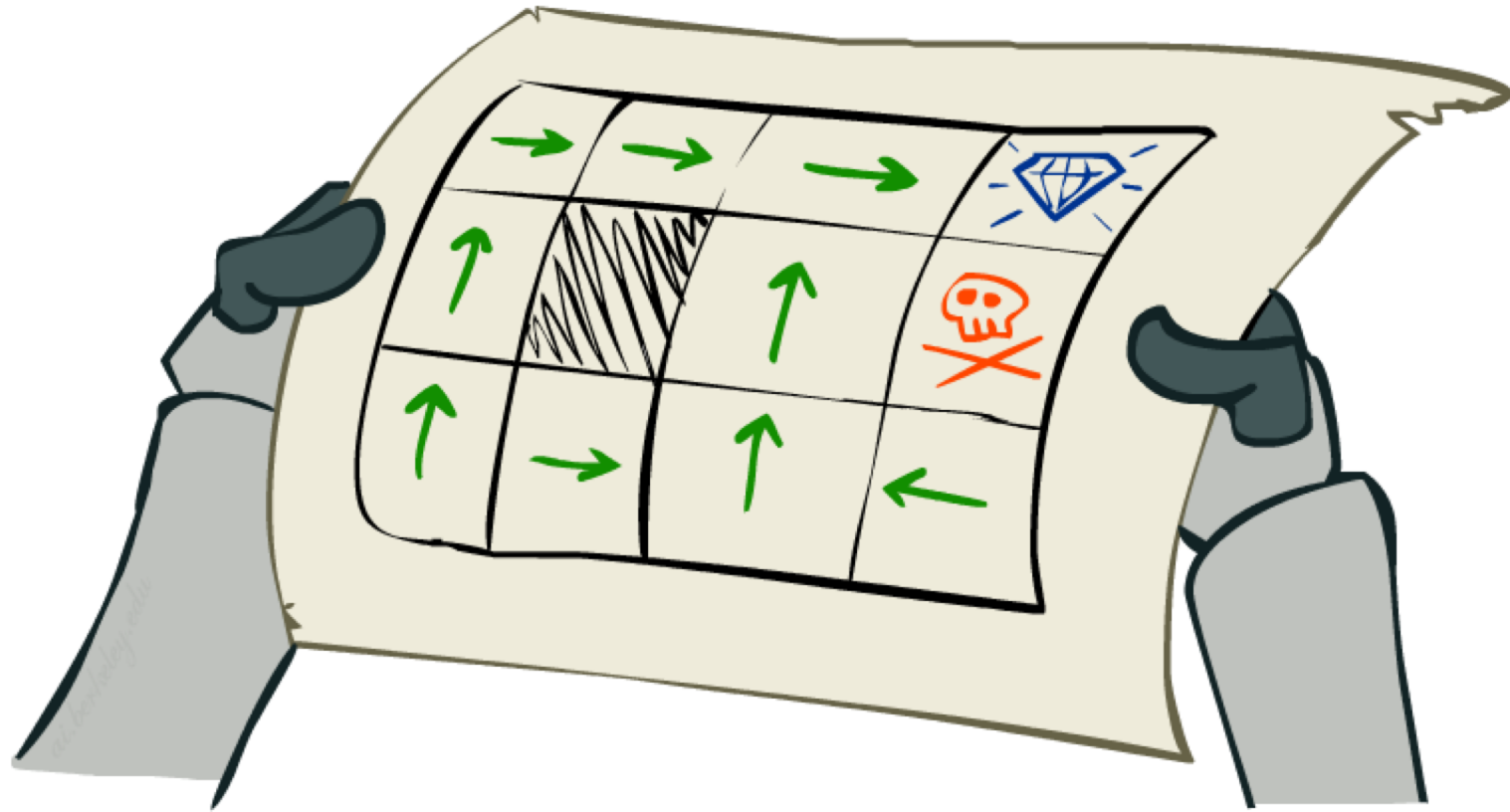o Each MDP state projects an expectimax-like search tree

s is a *state*

s, a

(s,a,s' ) called a *transition*

$T(s,a,s') = P(s'|s,a)$

$R(s,a,s')$

a

s,a,s'

s'

# Recap: Defining MDPs

o Markov decision processes:
  o Set of states S
  o Start state $s_0$
  o Set of actions A
  o Transitions P(s' | s,a) (or T(s,a,s'))
  o Rewards R(s,a,s') (and discount $\gamma$)

o MDP quantities so far:
  o Policy = Choice of action for each state
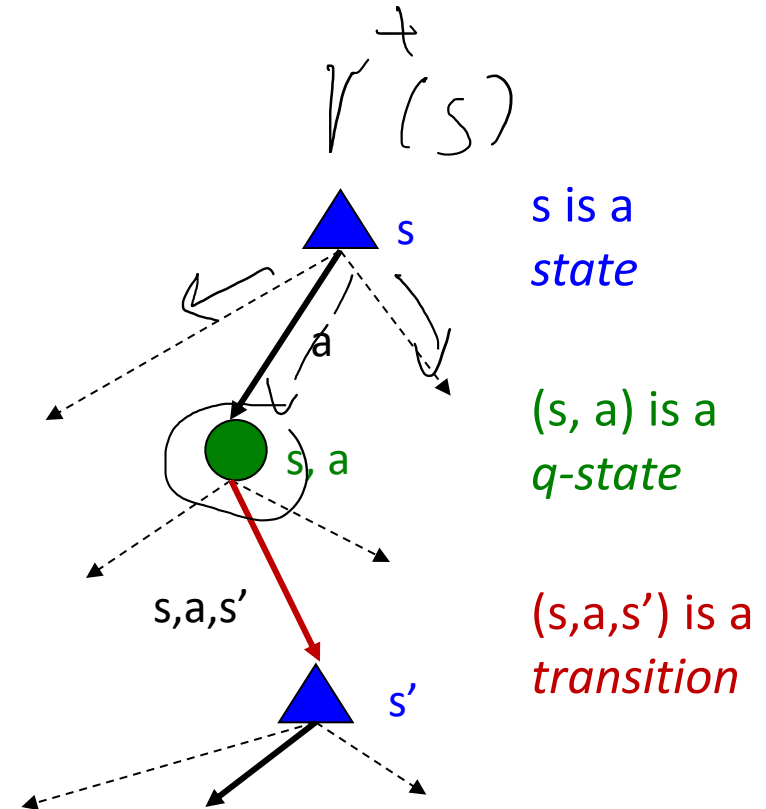  o Utility = sum of (discounted) rewards

# Solving MDPs

# Optimal Quantities

$V(s)$

- **The value (utility) of a state s:** $V^+(s)$

  $V^*(s)$ = expected utility starting in s and acting optimally



s is a *state*

(s, a) is a *q-state*

(s,a,s') is a *transition*

- **The value (utility) of a q-state (s,a):**

  $Q^*(s,a)$ = expected utility starting out having taken action a from state s and (thereafter) acting optimally

- **The optimal policy:**

  $\pi^*(s)$ = optimal action from state s
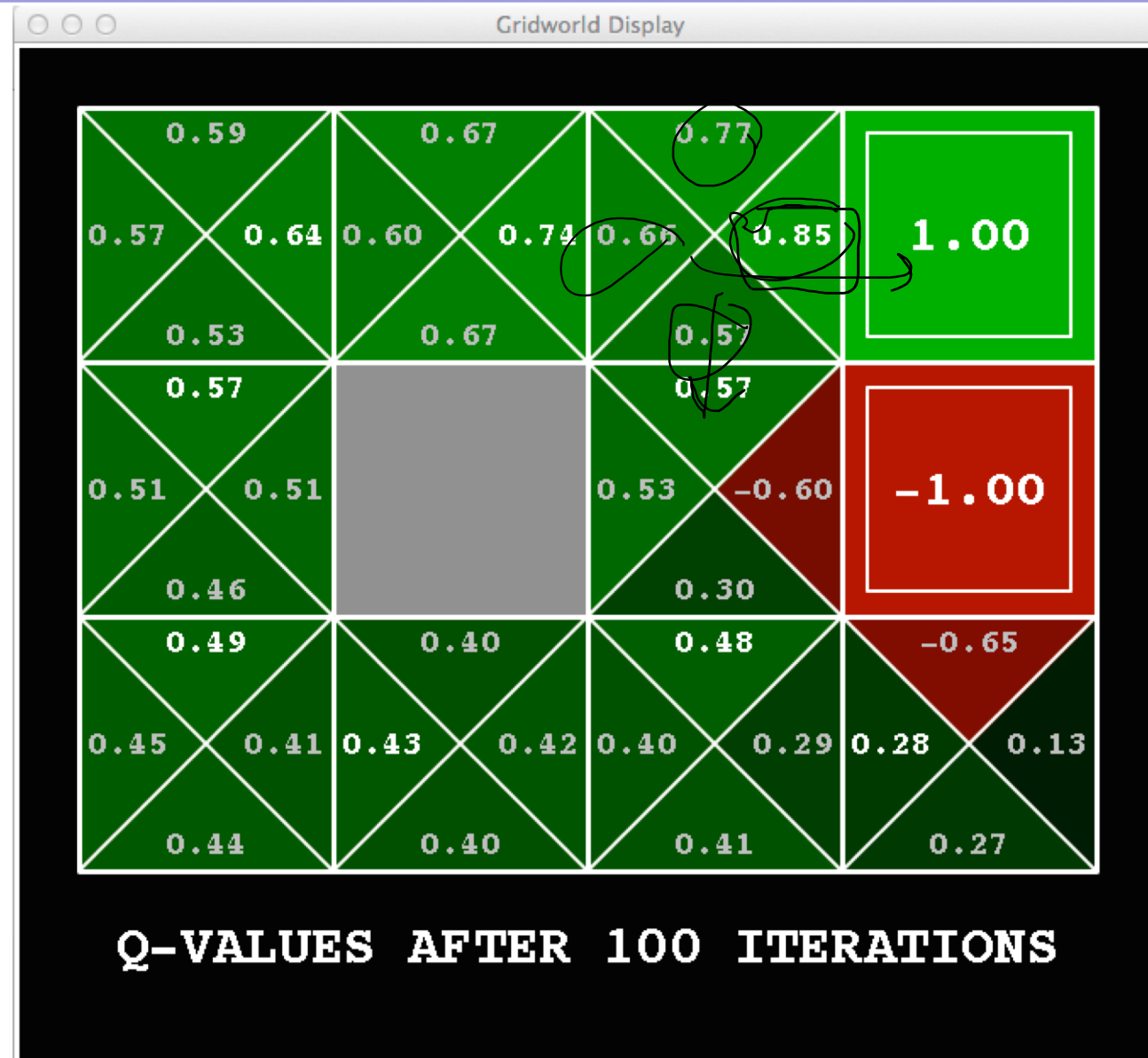
# Snapshot of Demo – Gridworld V Values



Noise = 0.2
Discount = 0.9
Living reward = 0

# Snapshot of Demo – Gridworld Q Values



Q-VALUES AFTER 100 ITERATIONS

Noise = 0.2
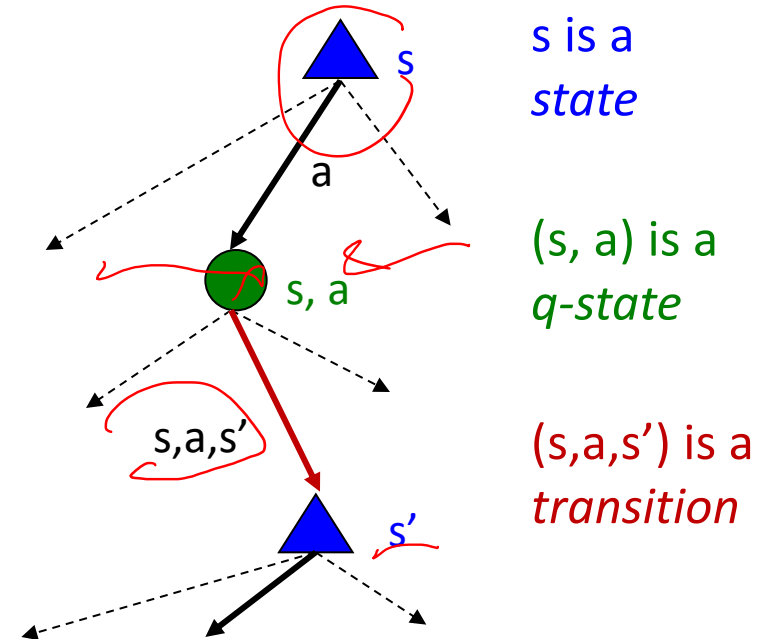Discount = 0.9
Living reward = 0

# Announcements

o Midterm (Search-Games-MDPs)

  o Take home (Nov. 4th-Nov 6th)

  o Midterm Review Session:

    o Respond to the Piazza poll regarding Review session for next week.

  o Additional office hour

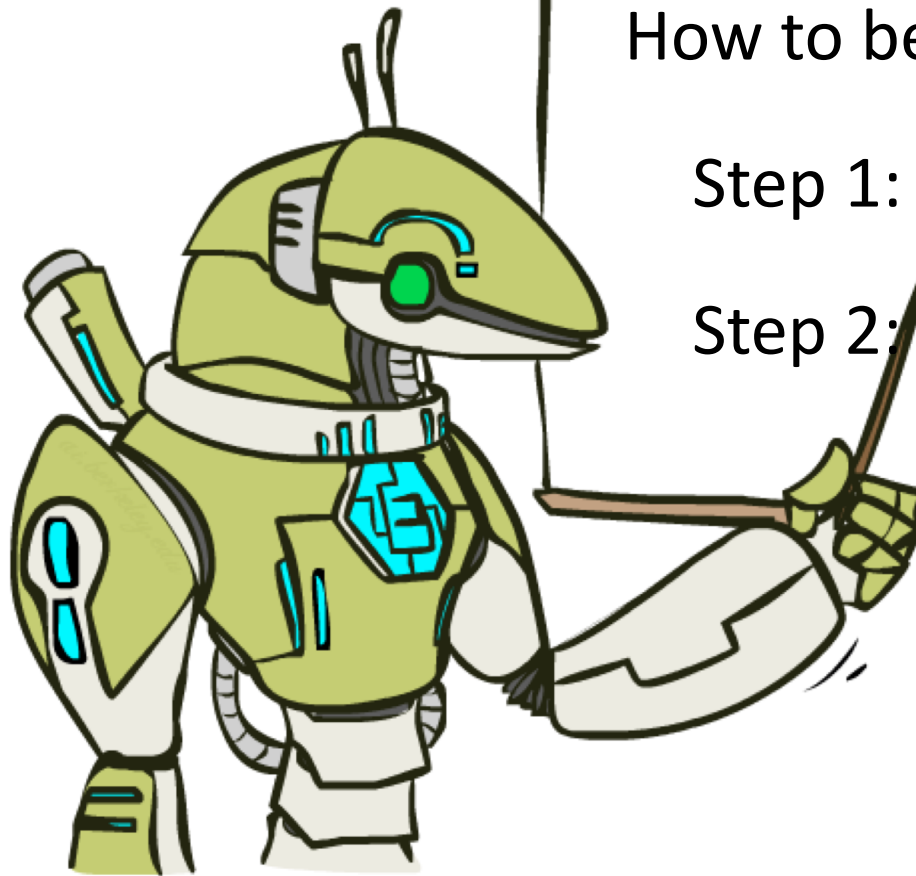o Programming assignment 2 is due Oct. 30th

Mid-Quarter Review

# Recap: MDPs Optimal Quantities

- The value (utility) of a state s:

    $V^*(s)$ = expected utility starting in s and acting optimally

- The value (utility) of a q-state (s,a):

    $Q^*(s,a)$ = expected utility starting out having taken action a from state s and (thereafter) acting optimally

- The optimal policy:

    $\pi^*(s)$ = optimal action from state s

s is a *state*

(s, a) is a *q-state*

(s,a,s') is a *transition*

a

s, a

s,a,s'

s'

# The Bellman Equations

How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal
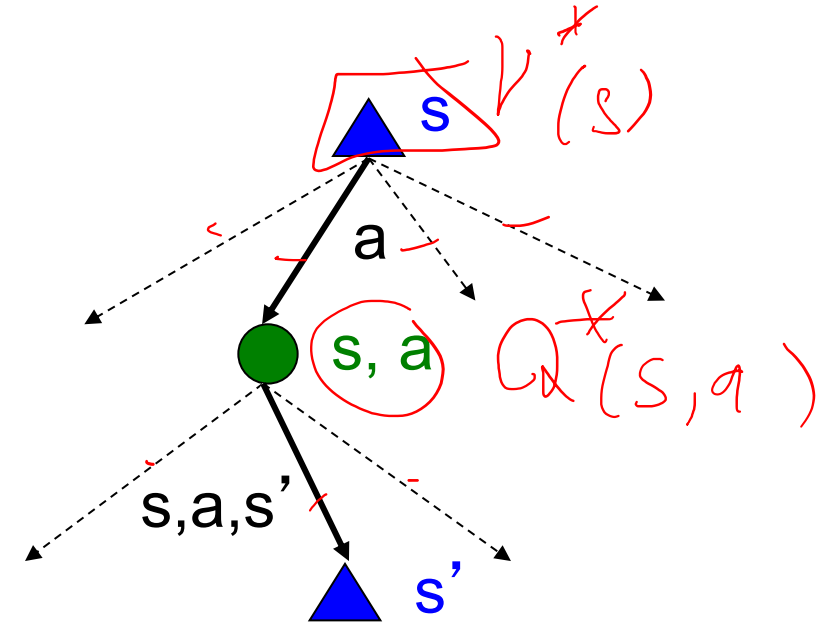
# Values of States (Bellman Equations)

o Fundamental operation: compute the (expectimax) value of a state
  - o Expected utility under optimal action
  - o Average sum of (discounted) rewards
  - o This is just what expectimax computed!
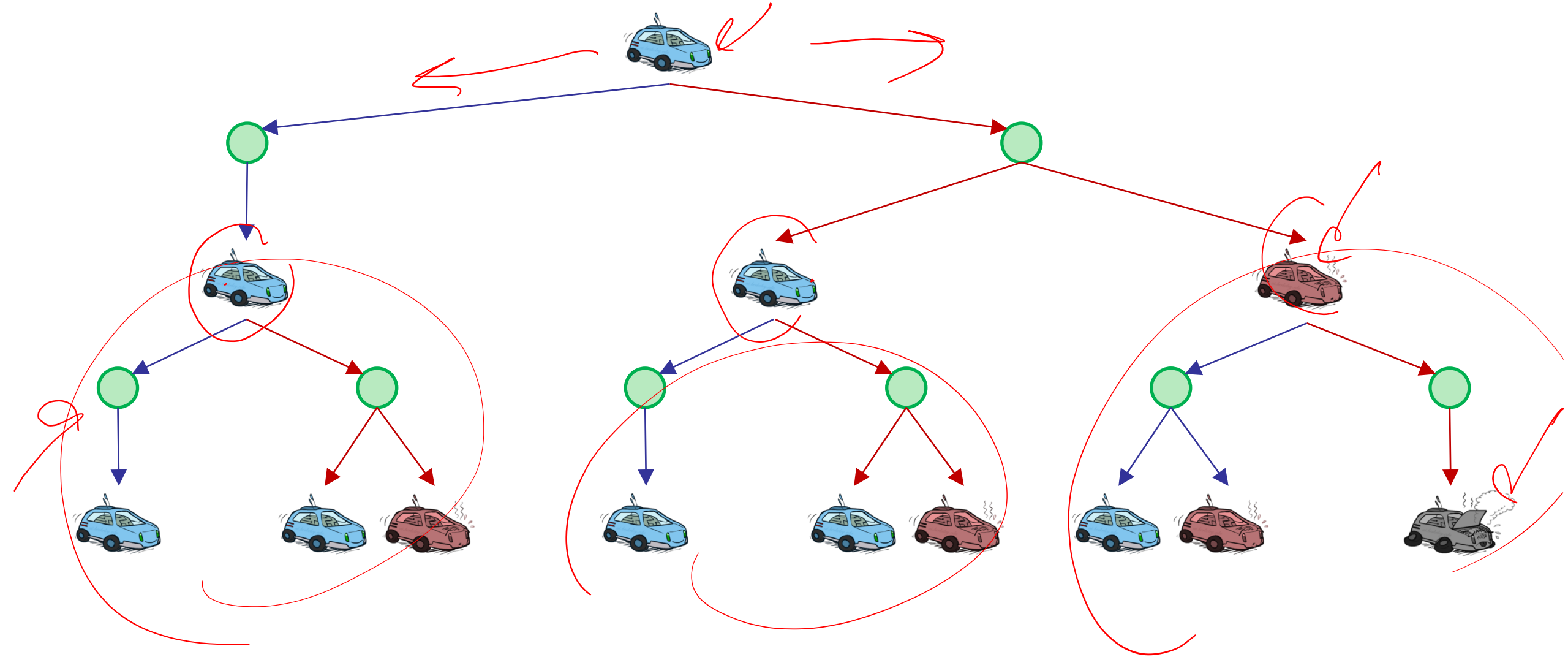
o Recursive definition of value:

$$V^*(s) = \max_a Q^*(s,a)$$

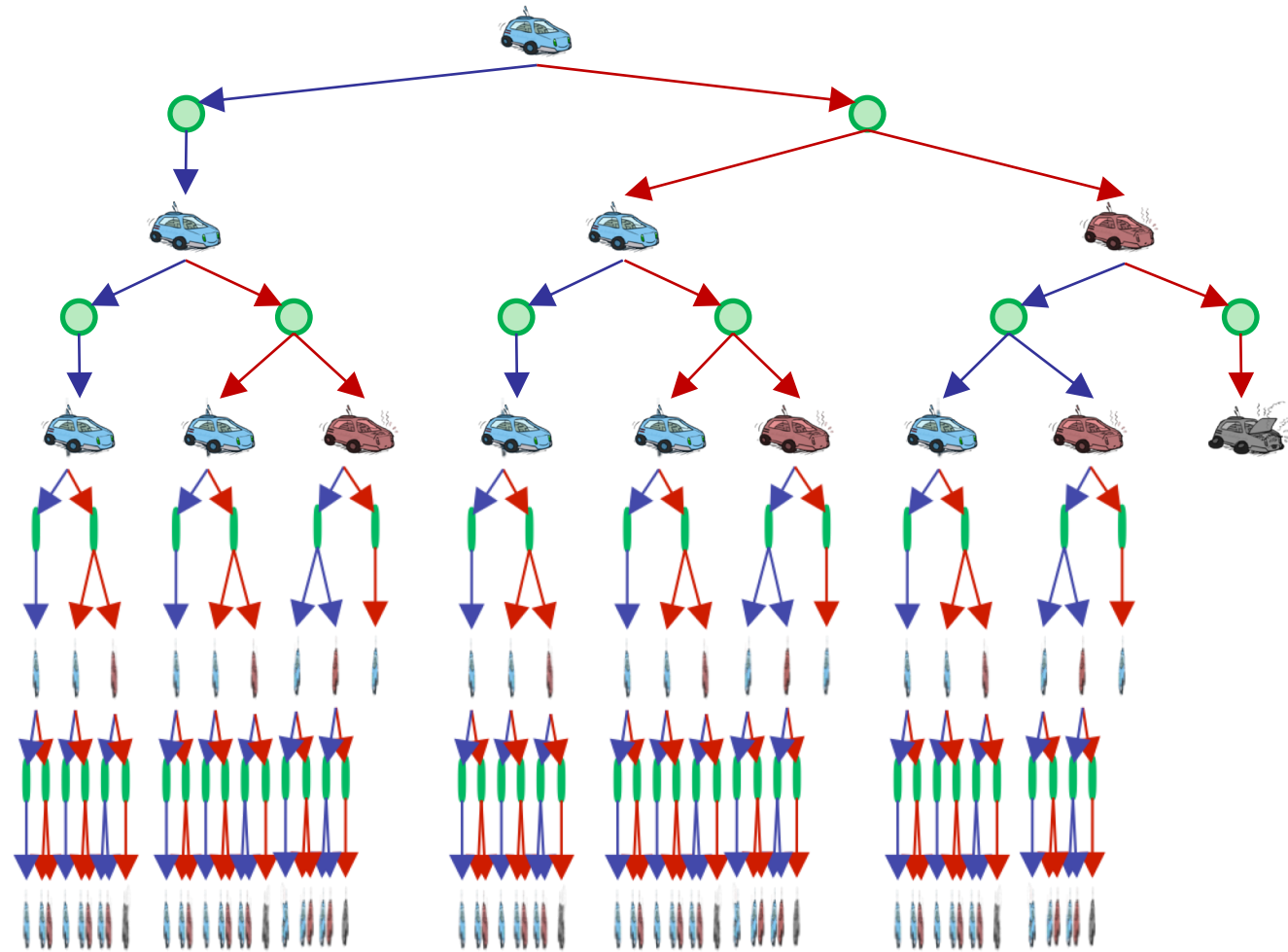$$Q^*(s,a) = \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V^*(s') \right]$$

$$V^*(s) = \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V^*(s') \right]$$
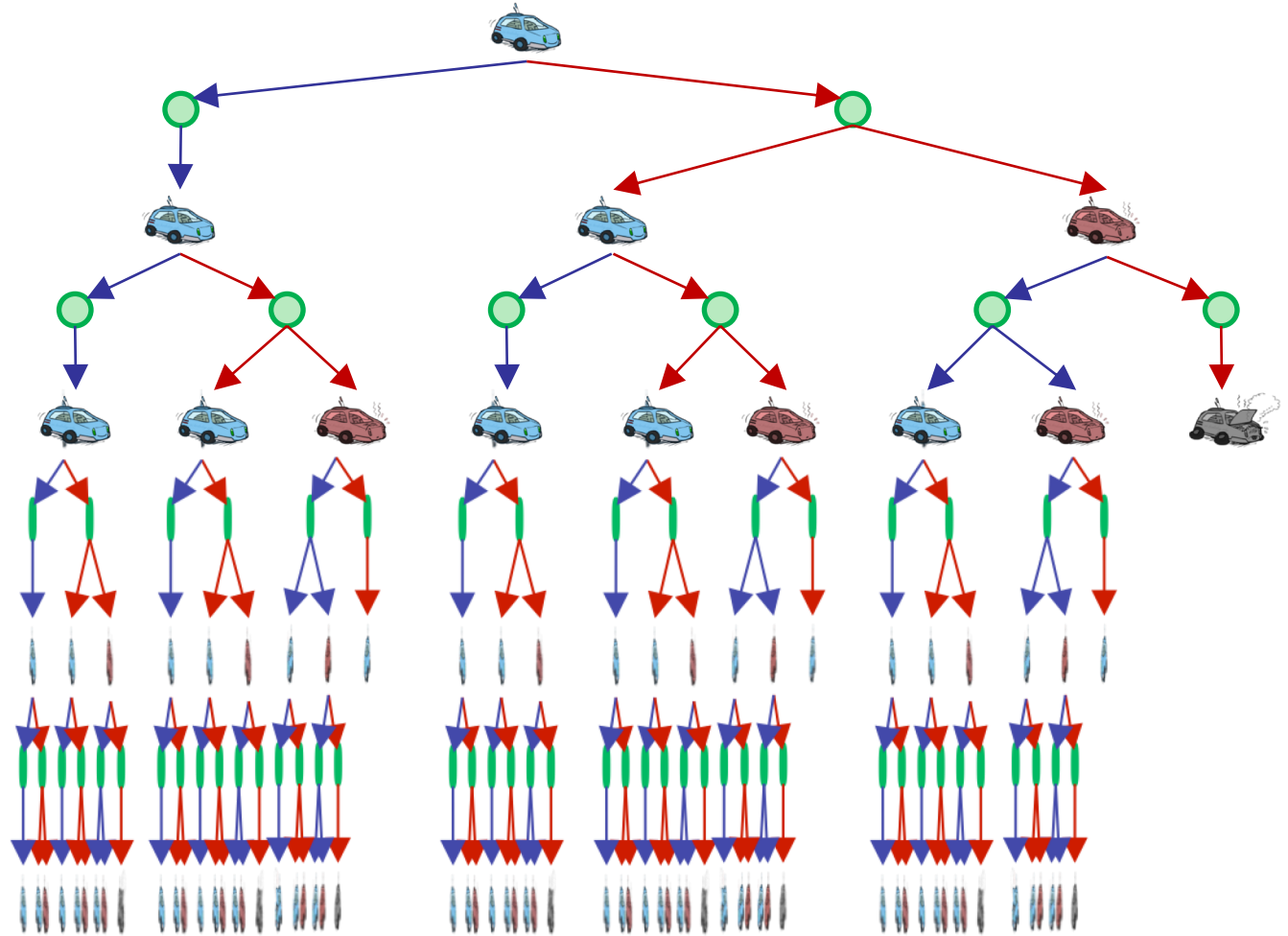
# Racing Search Tree

# Racing Search Tree

# Racing Search Tree

- We're doing way too much work with expectimax!

- Problem: States are repeated
  - Idea: Only compute needed quantities once

- Problem: Tree goes on forever
  - Idea: Do a depth-limited computation, but with increasing depths until change is small
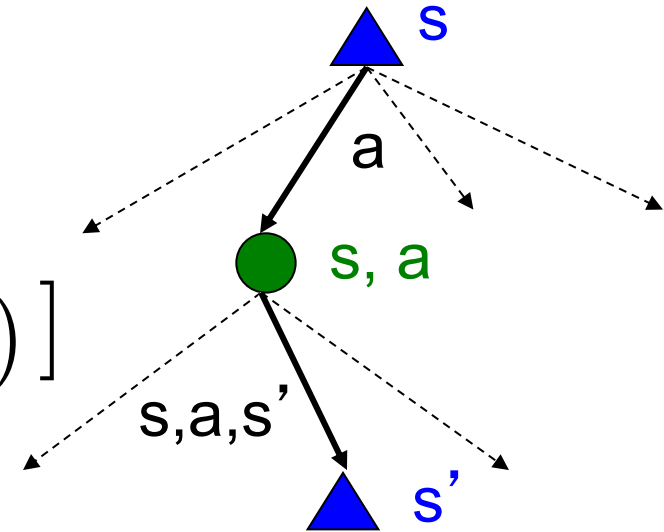  - Note: deep parts of the tree eventually don't matter if $\gamma < 1$

# Values of States

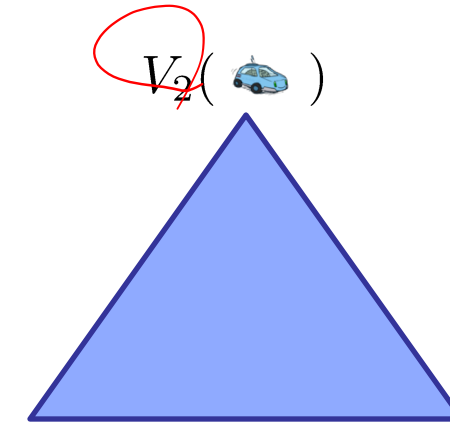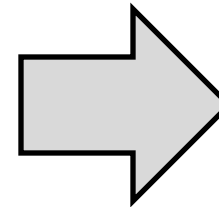o Recursive definition of value:
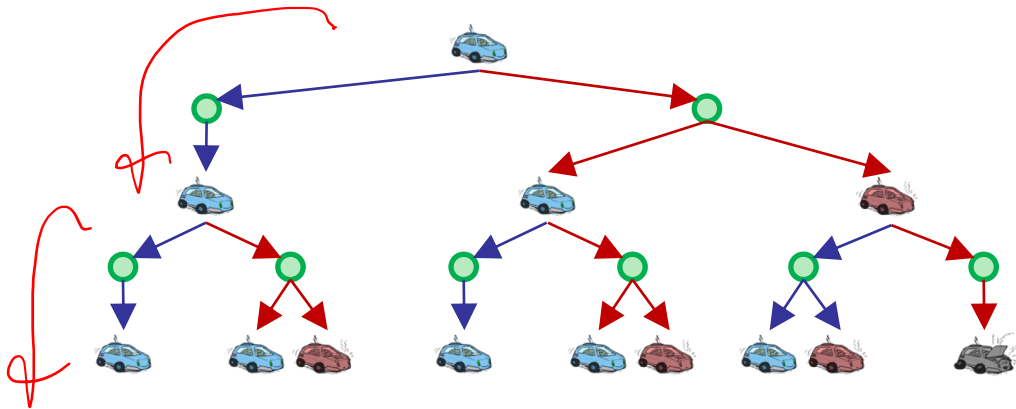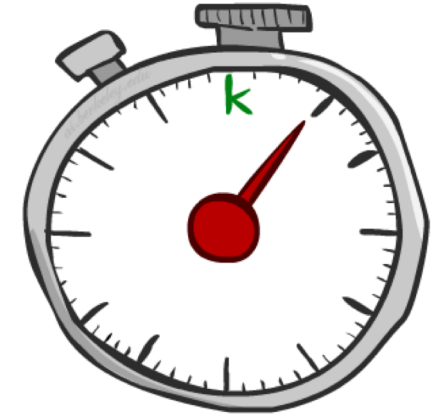
$$V^*(s) = \max_a \ Q^*(s,a)$$

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

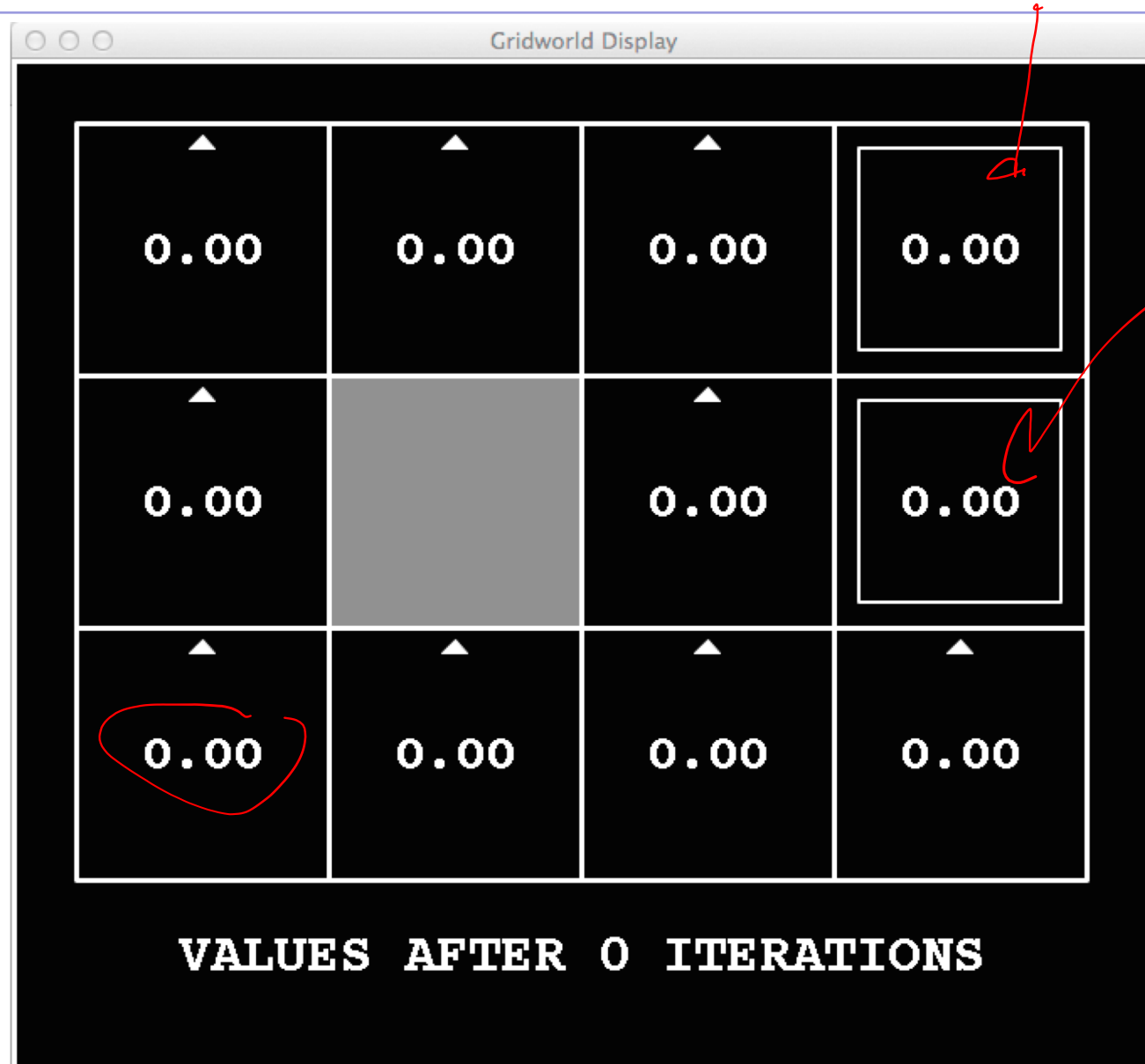$$V^*(s) = \max_a \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma V^*(s')]$$

# Time-Limited Values

o Key idea: time-limited values

$$V^*_k(s)$$

o Define $V_k(s)$ to be the optimal value of s if the game ends in k more time steps

o Equivalently, it's what a depth-k expectimax would give from s



$V_2( \quad )$

# k=0



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3
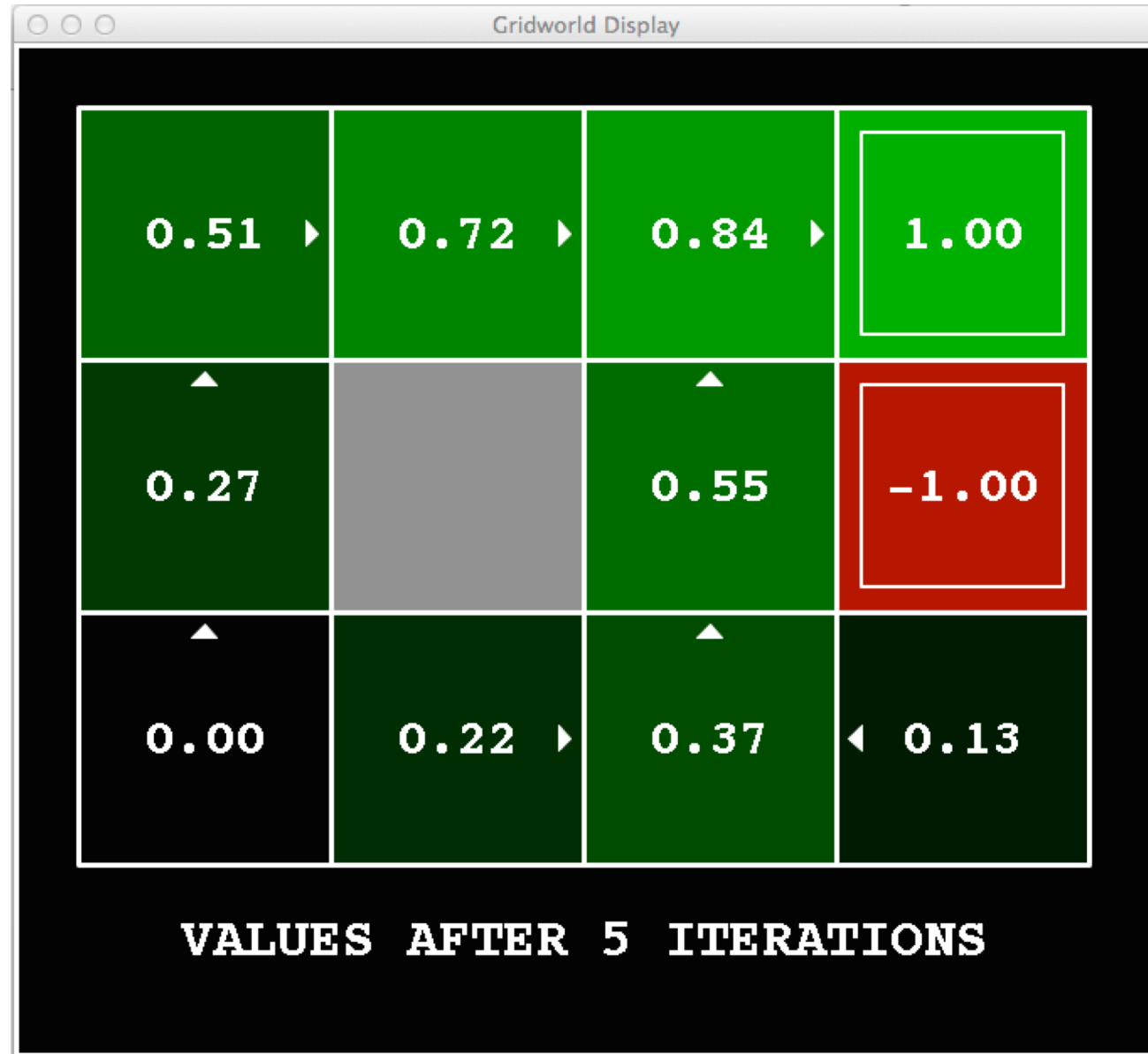


Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=5



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



Noise = 0.2
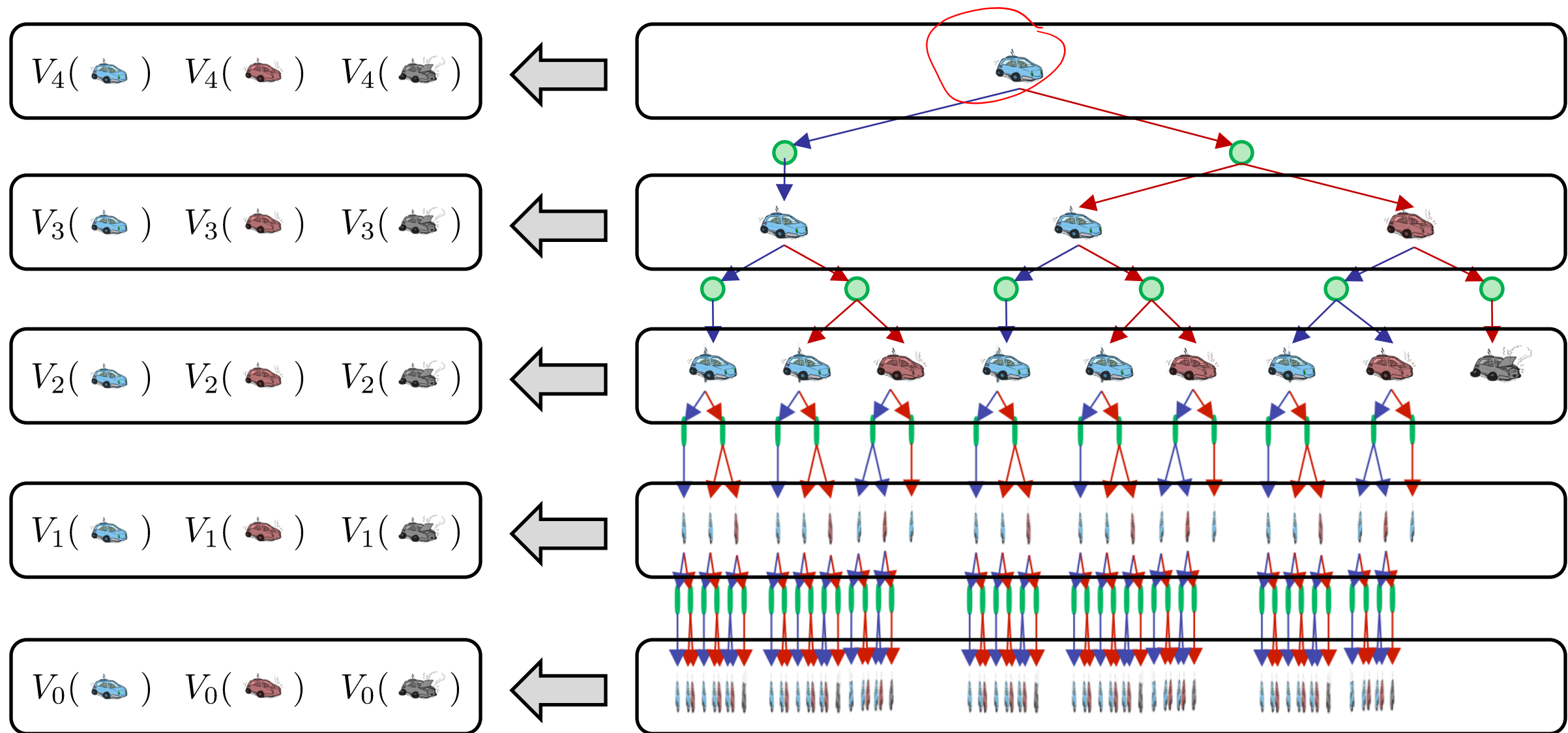Discount = 0.9
Living reward = 0

# k=12



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



Noise = 0.2
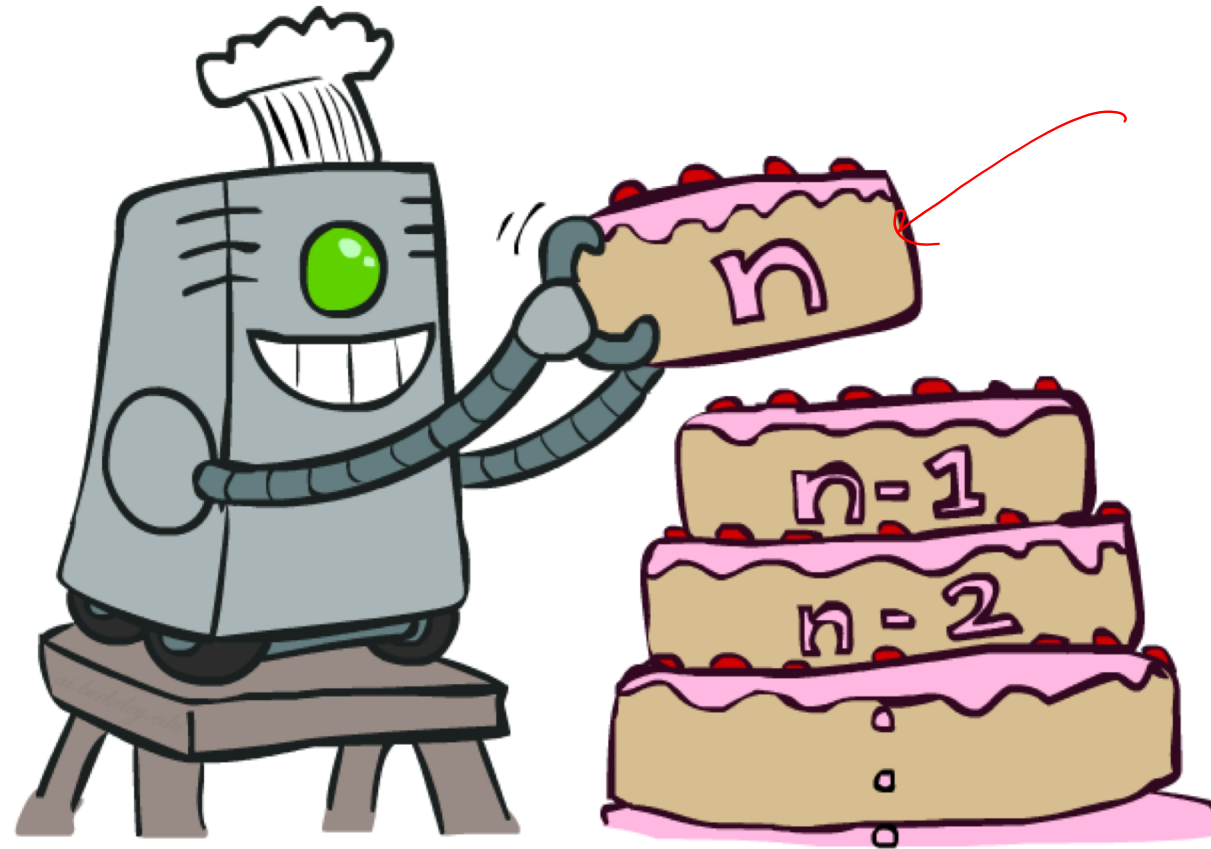Discount = 0.9
Living reward = 0

# Computing Time-Limited Values

$V_4(\text{🚗}) \quad V_4(\text{🚗}) \quad V_4(\text{🚗})$

$V_3(\text{🚗}) \quad V_3(\text{🚗}) \quad V_3(\text{🚗})$

$V_2(\text{🚗}) \quad V_2(\text{🚗}) \quad V_2(\text{🚗})$

$V_1(\text{🚗}) \quad V_1(\text{🚗}) \quad V_1(\text{🚗})$

$V_0(\text{🚗}) \quad V_0(\text{🚗}) \quad V_0(\text{🚗})$

# Announcements

- PS2 is due Oct. 30th

- Midterm: Take home
  - Due: Nov. 6th, will be released: Nov 4th Midterm Review Session:
  - Tue, 3:30-5pm at Allen 403
  - Solving SP 19 midterm  + Open Questions

  - Hanna: Holding extra office hour on Fri 12-1pm

# MDP Value Iteration

# Value Iteration

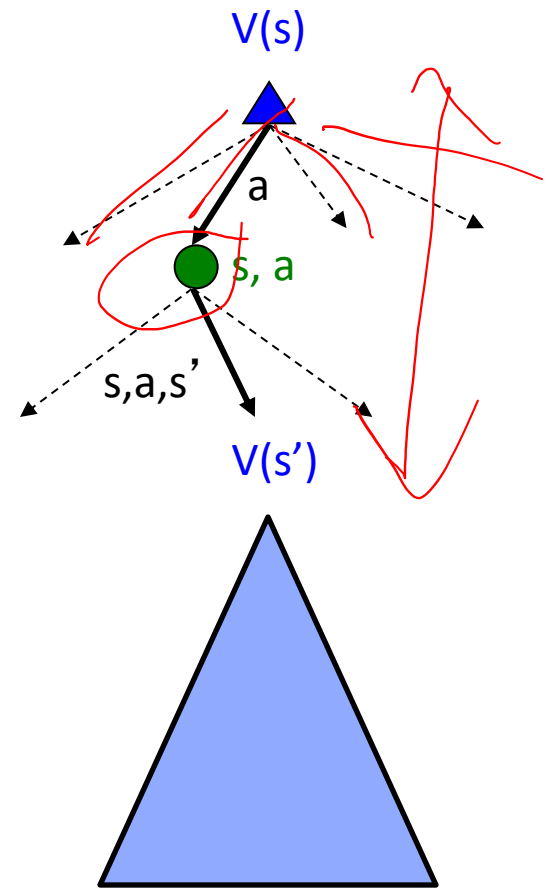- Bellman equations characterize the optimal values:

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

- Value iteration computes them:

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

- Value iteration is just a fixed point solution method
  - ... though the $V_k$ vectors are also interpretable as time-limited values
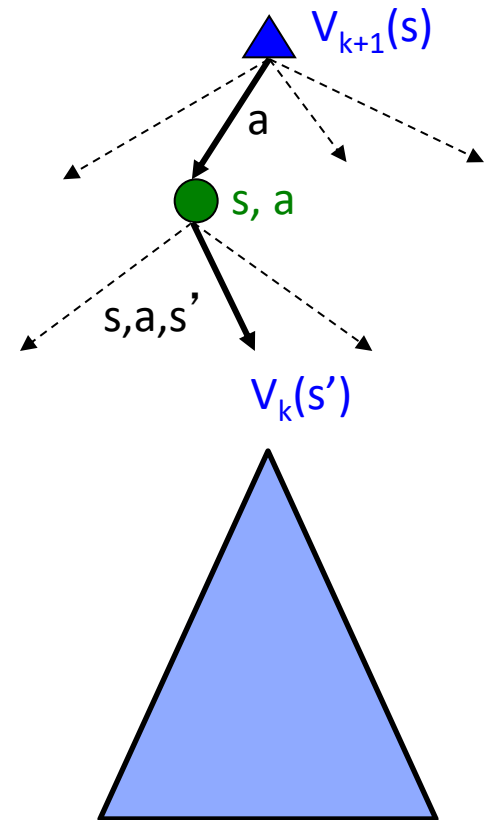
V(s)

a

s, a

s,a,s'

V(s')

# Value Iteration

o Start with $V_0(s) = 0$: no time steps left means an expected reward sum of zero

o Given vector of $V_k(s)$ values, do one ply of expectimax from each state:

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

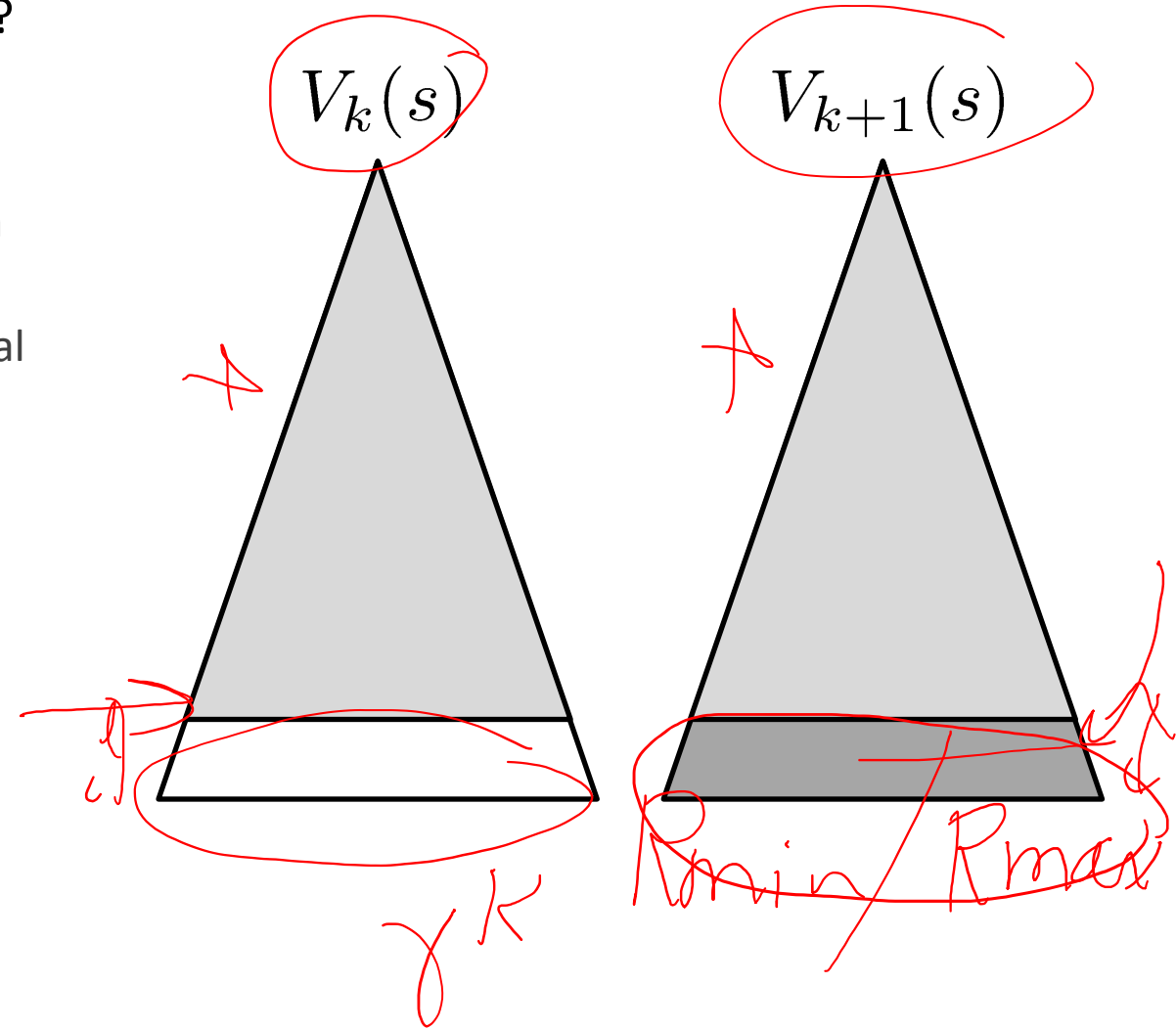o Repeat until convergence

$$|S| \times |A| \times |S|$$

o Complexity of each iteration: $O(S^2 A)$

o Theorem: will converge to unique optimal values
   o Basic idea: approximations get refined towards optimal values
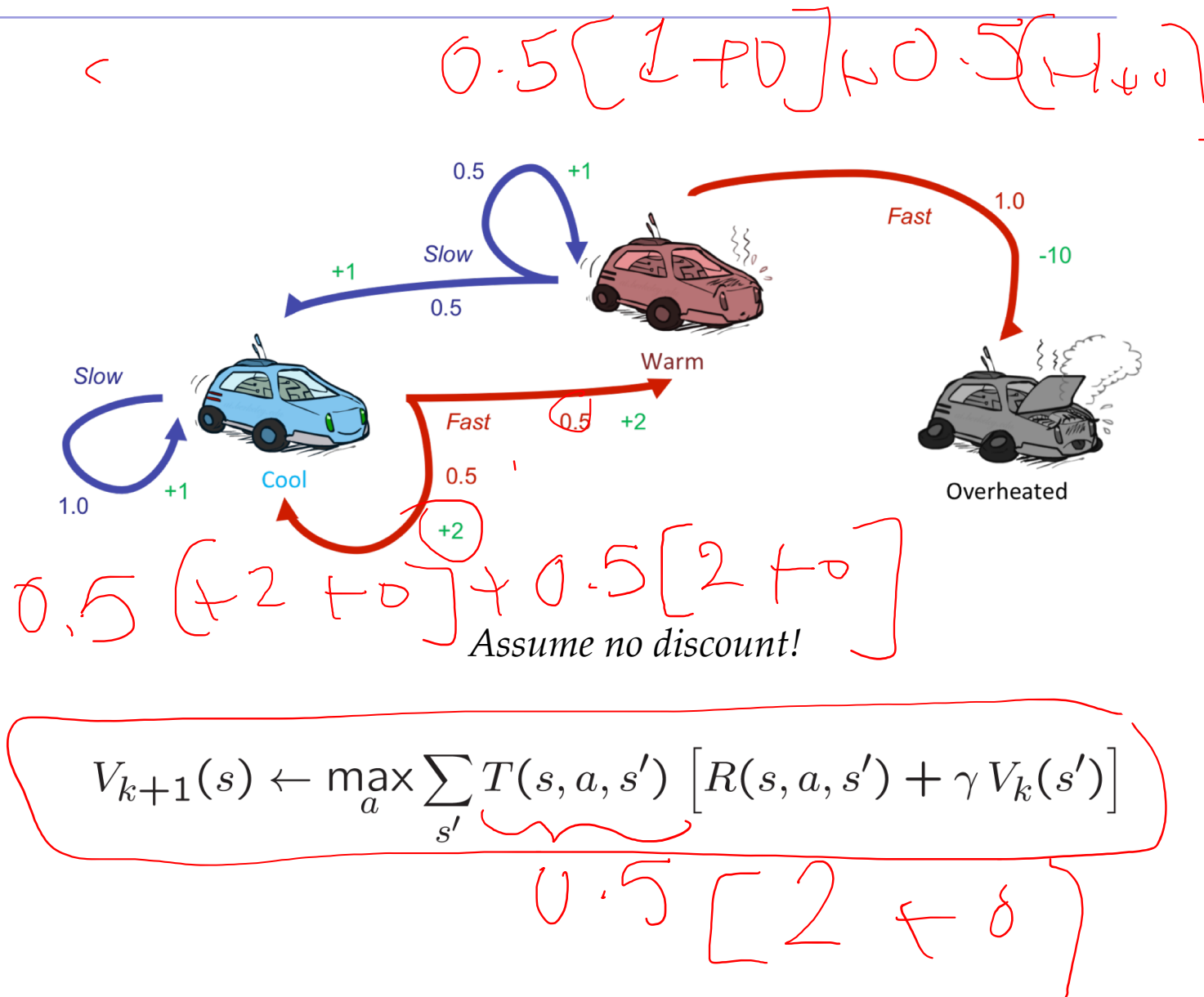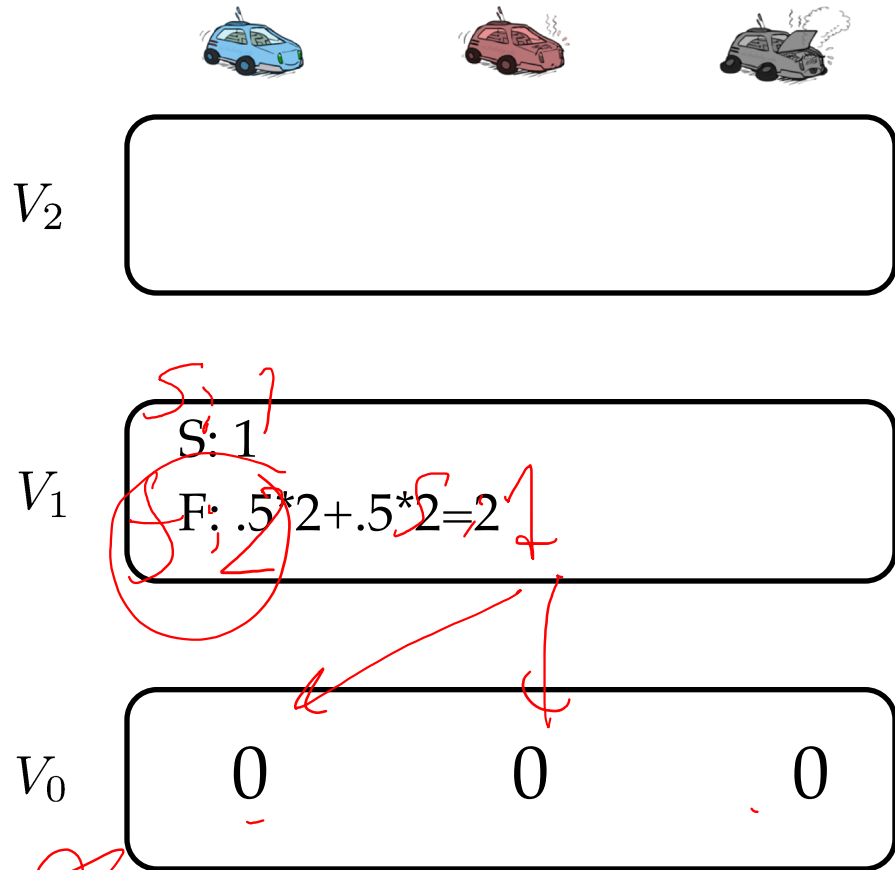   o Policy may converge long before values do

# Convergence*

- How do we know the $V_k$ vectors are going to converge?

- If the discount is less than 1
  - Sketch: For any state $V_k$ and $V_{k+1}$ can be viewed as depth k+1 expectimax results in nearly identical search trees
  - The difference is that on the bottom layer, $V_{k+1}$ has actual rewards while $V_k$ has zeros
  - That last layer is at best all $R_{MAX}$
  - It is at worst $R_{MIN}$
  - But everything is discounted by $\gamma^k$ that far out
  - So $V_k$ and $V_{k+1}$ are at most $\gamma^k \max|R|$ different
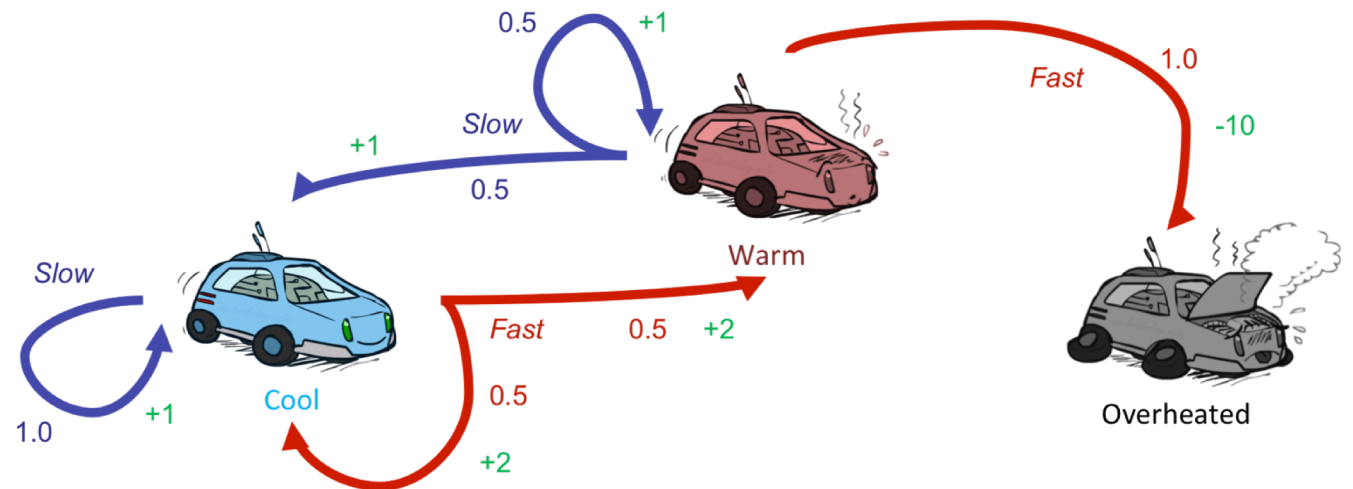  - So as k increases, the values converge

# Example: Value Iteration



$V_2$ [ ]

$V_1$

S: 1

F: .5*2+.5*2=2

$V_0$ [ 0      0      0 ]

Handwritten annotations:

0.5[1+0] & 0.5(H+0)

0.5[+2+0]+0.5[2+0]

*Assume no discount!*

0.5[2+0]

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

Diagram labels:
Slow +1  0.5
Slow +1  0.5
Fast 1.0  -10
Warm
Fast 0.5 +2
0.5 +2
Slow 1.0 +1
Cool
Overheated

# Example: Value Iteration

$V_2$

$V_1$    **2**    S: .5*1+.5*1=1
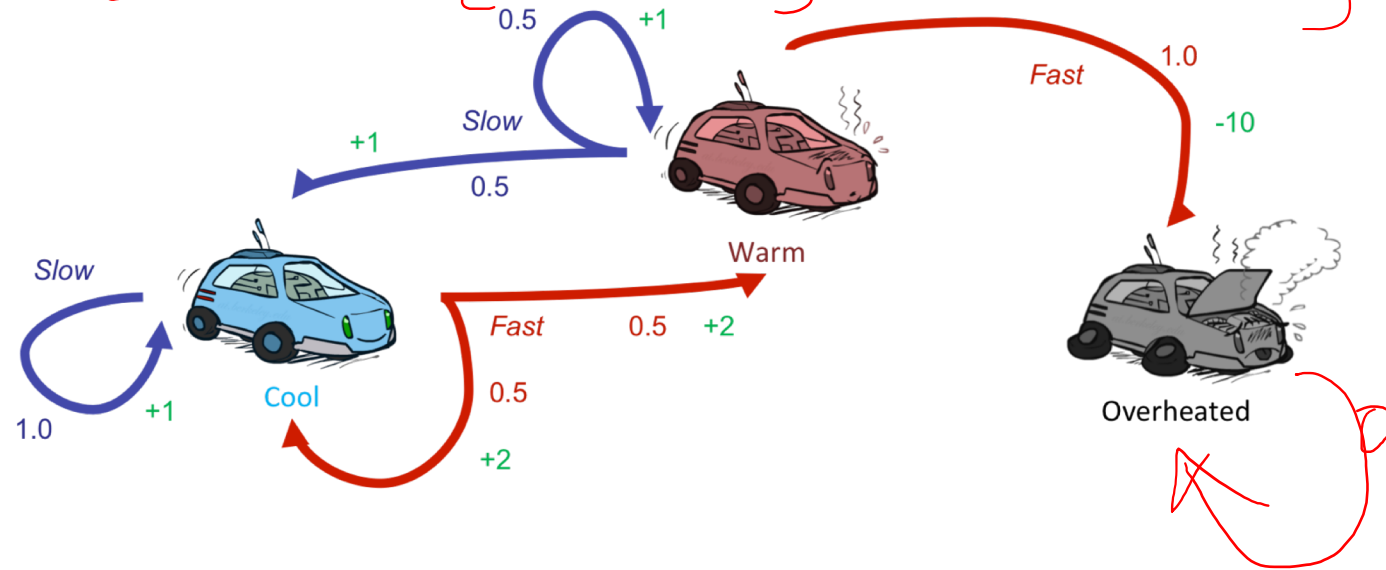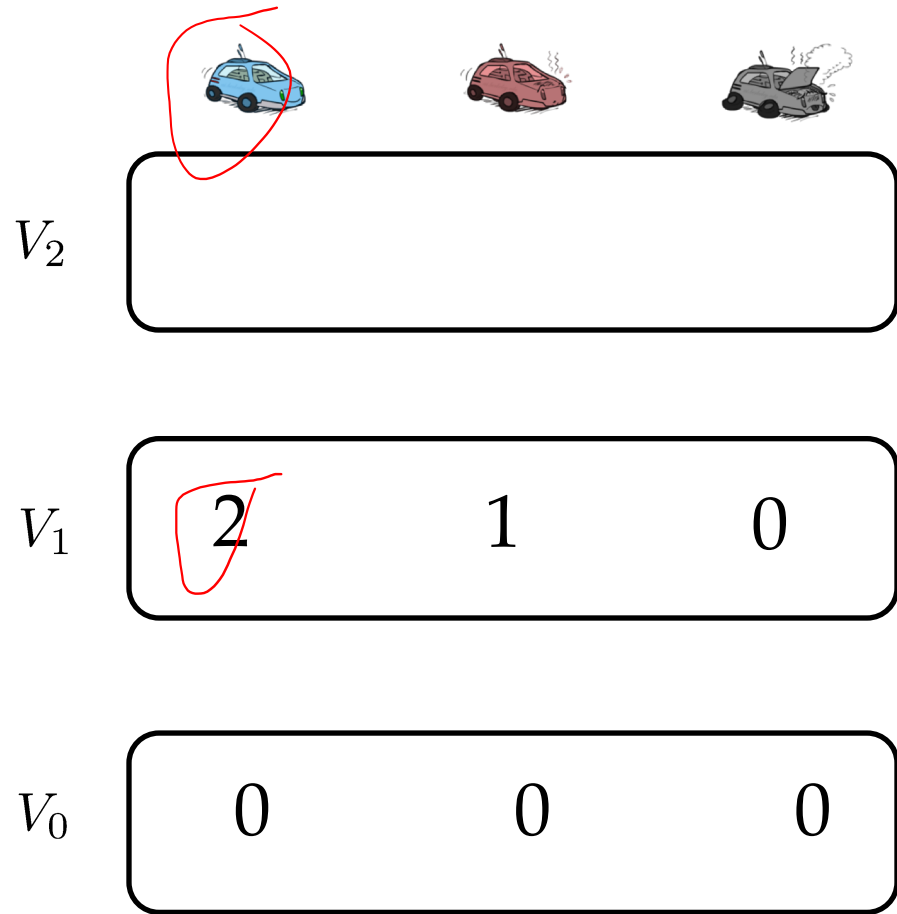
F: -10

$V_0$    0      0      0



*Assume no discount!*

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

# Example: Value Iteration

$S_i + 1 + 2 = 3$

$f, 0.5[+2+2] + 0.5[+2 + 1]$



$V_2$

$V_1$ : 2    1    0

$V_0$ : 0    0    0

*Assume no discount!*

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

# Example: Value Iteration



$V_2$

S: 1+2=3
F:
.5*(2+2)+.5*(2+1)=3.5

$V_1$ | 2 | 1 | 0

$V_0$ | 0 | 0 | 0

0.5    +1
Fast    1.0
Slow
-10
+1
Slow
0.5
Warm

Slow
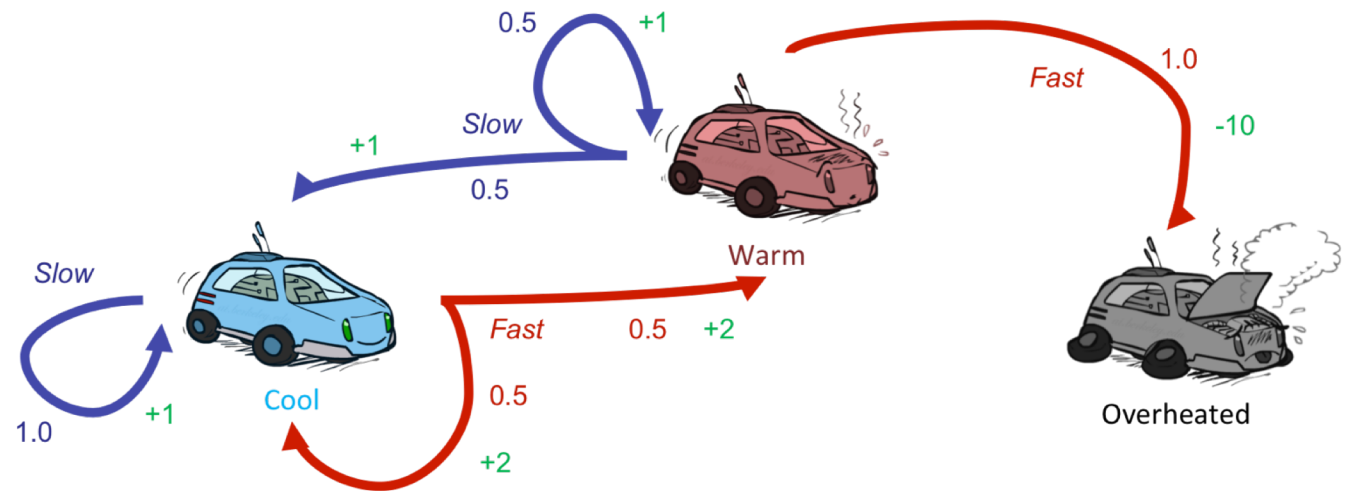Fast    0.5    +2
0.5
Cool    Overheated
1.0    +1
+2

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_k(s') \right]$$
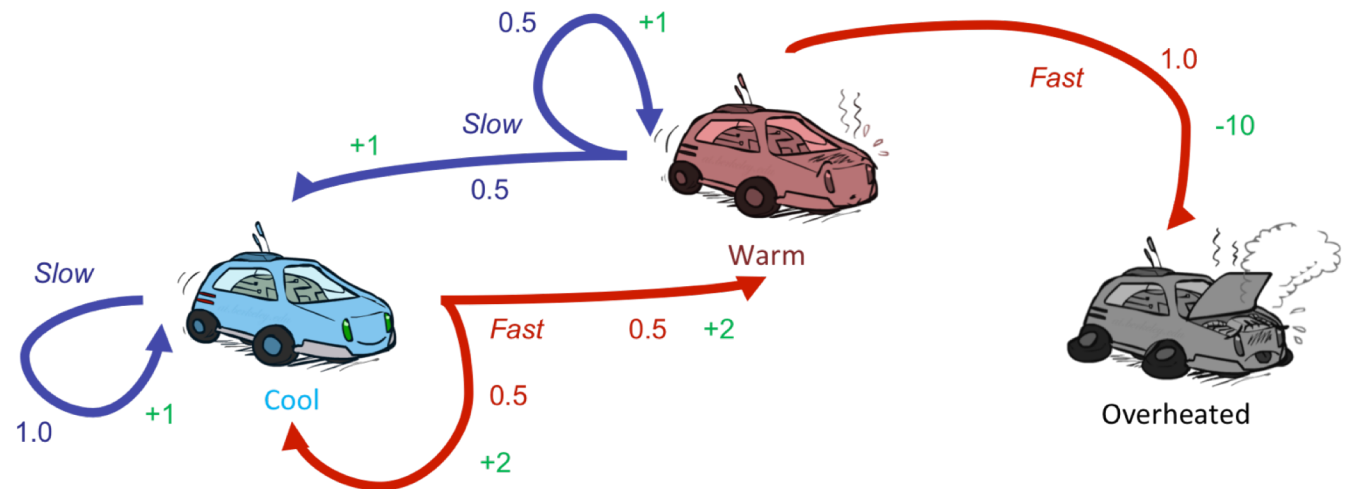
# Example: Value Iteration
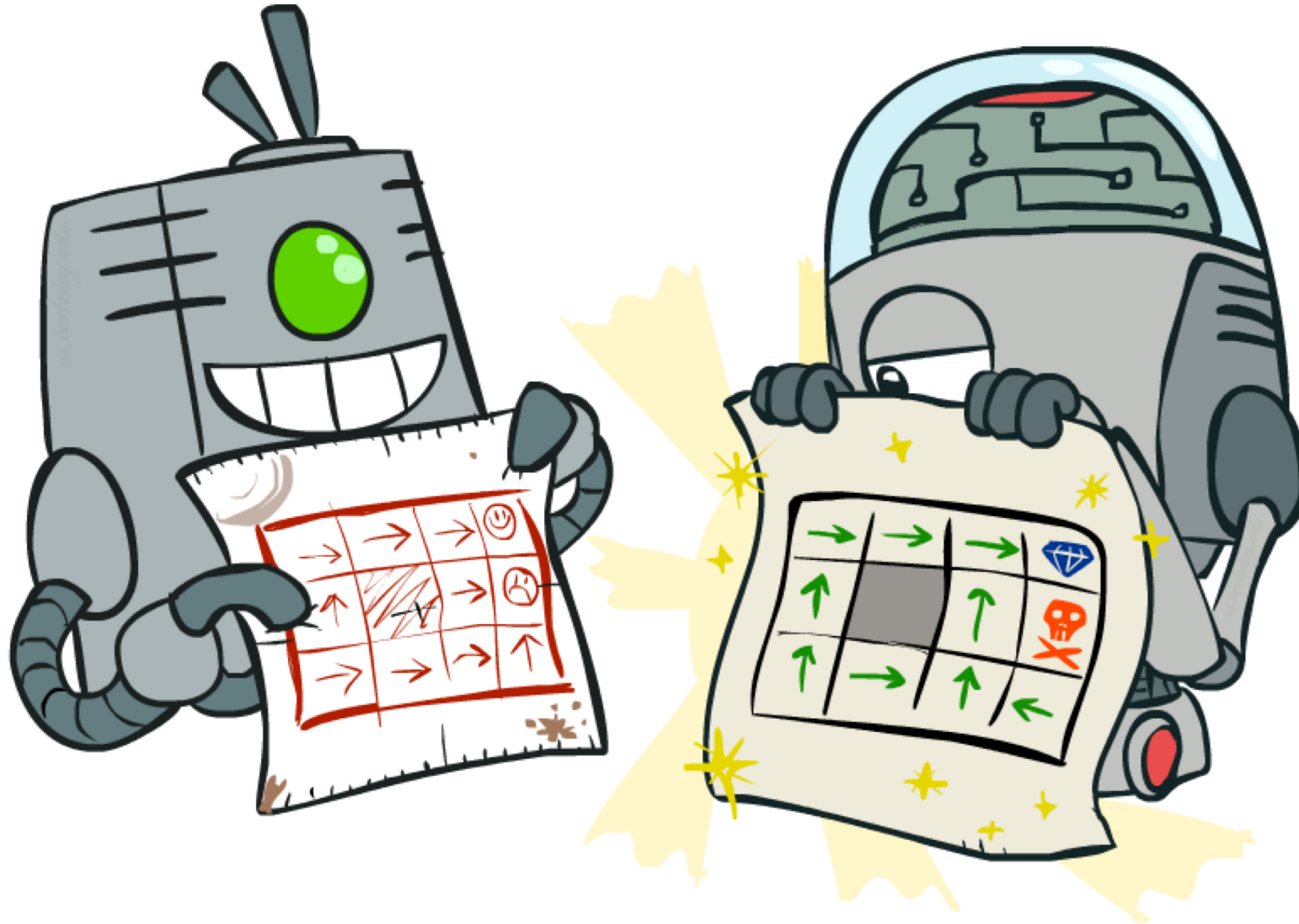
$V_2$

| 3.5 | 2.5 | 0 |

$V_1$

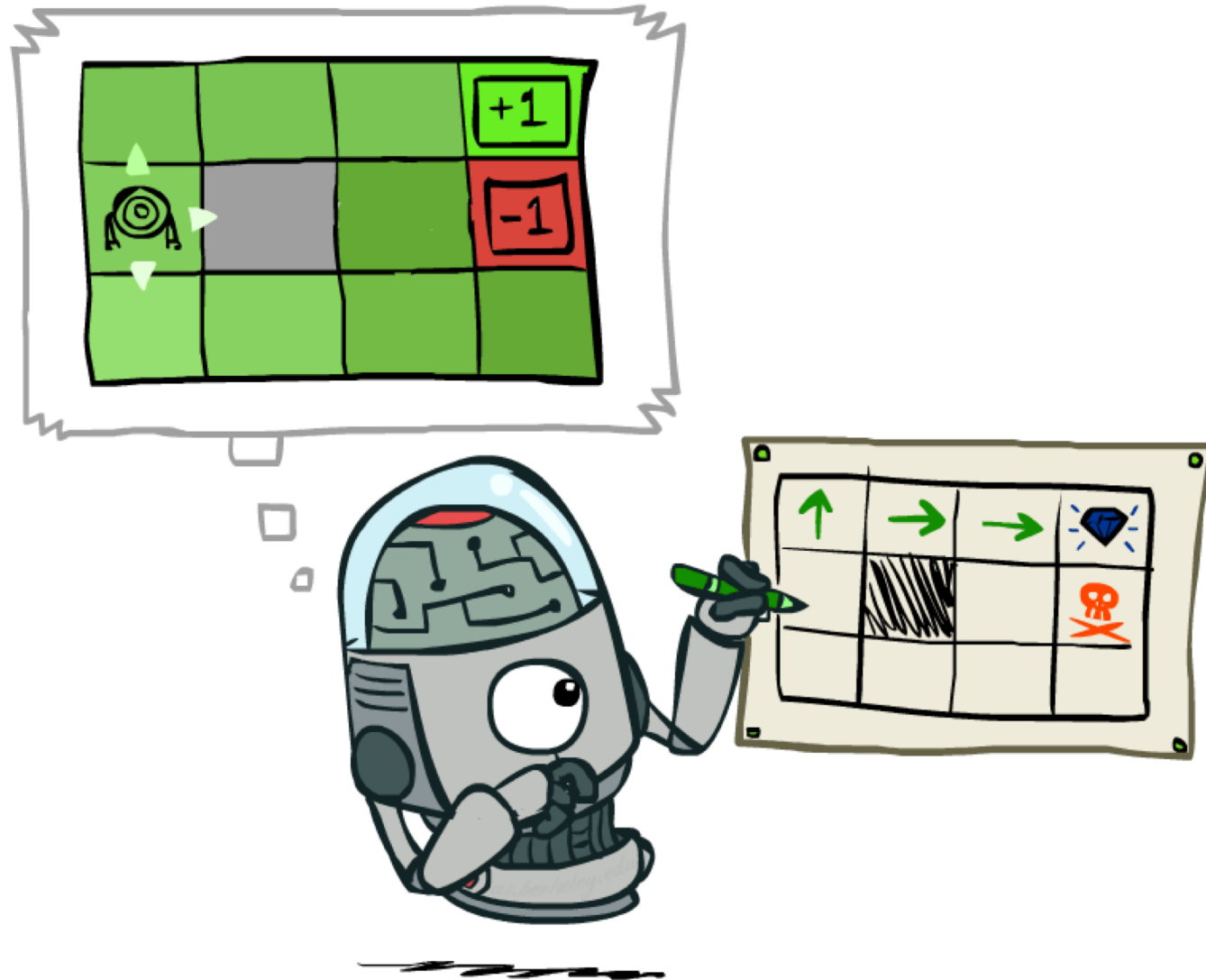| 2 | 1 | 0 |



*Assume no discount!*

$V_0$

| 0 | 0 | 0 |

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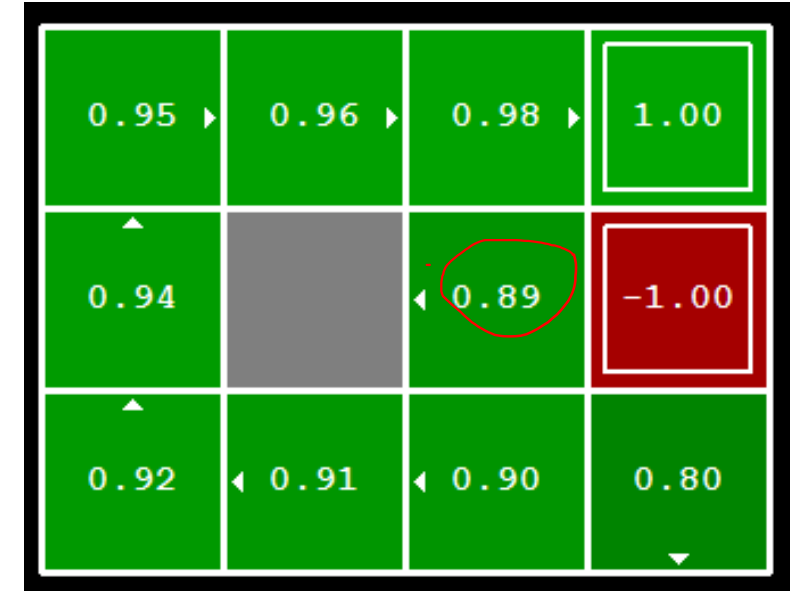
# Policy Methods

# Policy Extraction

# Computing Actions from Values

o Let's imagine we have the optimal values V*(s)



o How should we act?
   o It's not obvious!

o We need to do a mini-expectimax (one step)

$$\pi^*(s) = \arg\max_a \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma V^*(s')]$$

o This is called policy extraction, since it gets the policy implied by the values
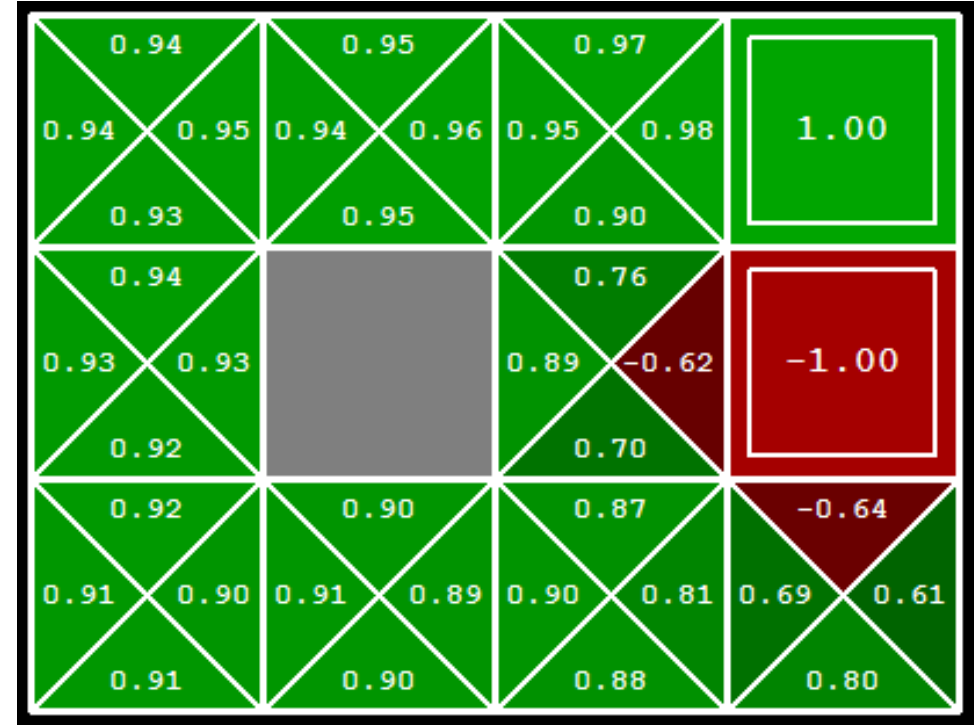
# Computing Actions from Q-Values



o Let's imagine we have the optimal q-values:

$Q_{state}(s, a)$

o How should we act?

  o Completely trivial to decide!

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

o Important lesson: actions are easier to select from q-values than values!

# Policy Evaluation

# Fixed Policies

Do the optimal action

Do what π says to do



o Expectimax trees max over all actions to compute the optimal values

o If we fixed some policy π(s), then the tree would be simpler – only one action per state

   o … though the tree's value would depend on which policy we fixed

# Utilities for a Fixed Policy

- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy
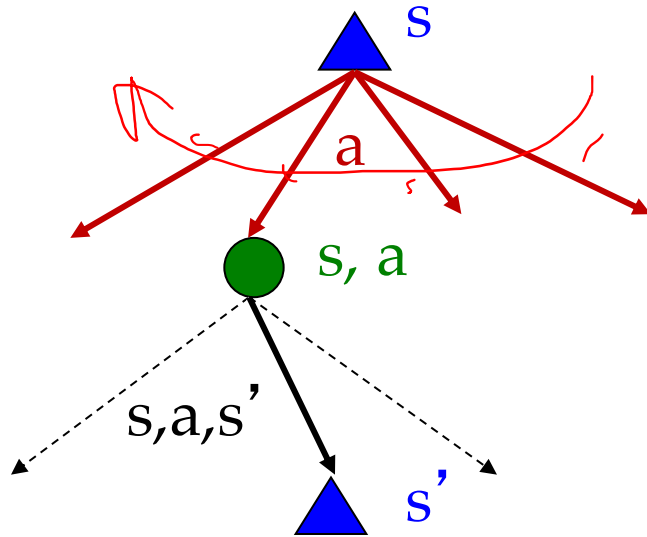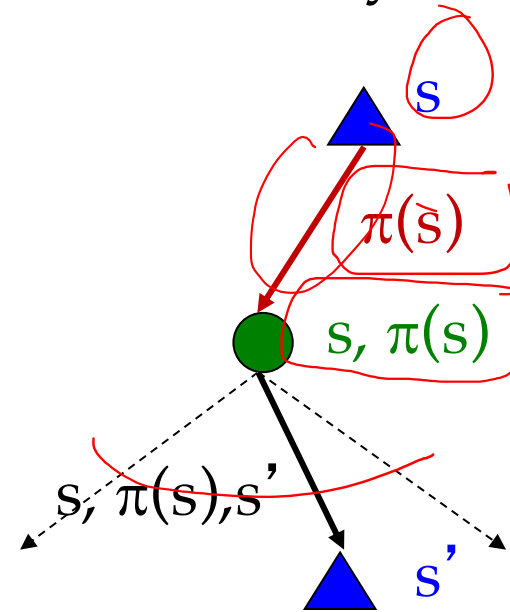
- Define the utility of a state s, under a fixed policy π:

  $V^\pi(s)$ = expected total discounted rewards starting in s and following π

- Recursive relation (one-step look-ahead / Bellman equation):

$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V^\pi(s')]$$

s

π(s)

s, π(s)

s, π(s),s'

s'

# Example: Policy Evaluation

Always Go Right

Always Go Forward

# Example: Policy Evaluation



Always Go Right

Always Go Forward

# Policy Evaluation

o How do we calculate the V's for a fixed policy π?

o Idea 1: Turn recursive Bellman equations into updates
  (like value iteration)

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

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

|S||A||S|

o Efficiency: O(S²) per iteration

o Idea 2: Without the maxes, the Bellman equations are just a linear system
  o Solve with Matlab (or your favorite linear system solver)

s

π(s)

s, π(s)

s, π(s),s'

s'

# Policy Iteration

# Problems with Value Iteration

o Value iteration repeats the Bellman updates:

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

s

a

s, a

s,a,s'

s'

o Problem 1: It's slow – O(S²A) per iteration

o Problem 2: The "max" at each state rarely changes

o Problem 3: The policy often converges long before the values

# k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=5



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# MDPs: Policy Iteration

o Alternative approach for optimal values:

   o **Step 1: Policy evaluation:** calculate utilities for some fixed policy (not optimal utilities!) until convergence

   o **Step 2: Policy improvement:** update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values

   o Repeat steps until policy converges

o This is policy iteration

   o It's still optimal!

   o Can converge (much) faster under some conditions

# Policy Iteration

o Evaluation: For fixed current policy π, find values with policy evaluation:
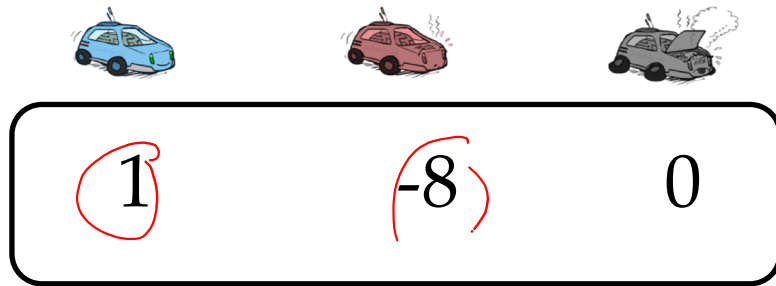
   o Iterate until values converge:

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

o Improvement: For fixed values, get a better policy using policy extraction

   o One-step look-ahead:

$$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$

# Example: Policy Improvement

*Assume: the values for the current policy π*



*Policy π*  *Slow in Cool*
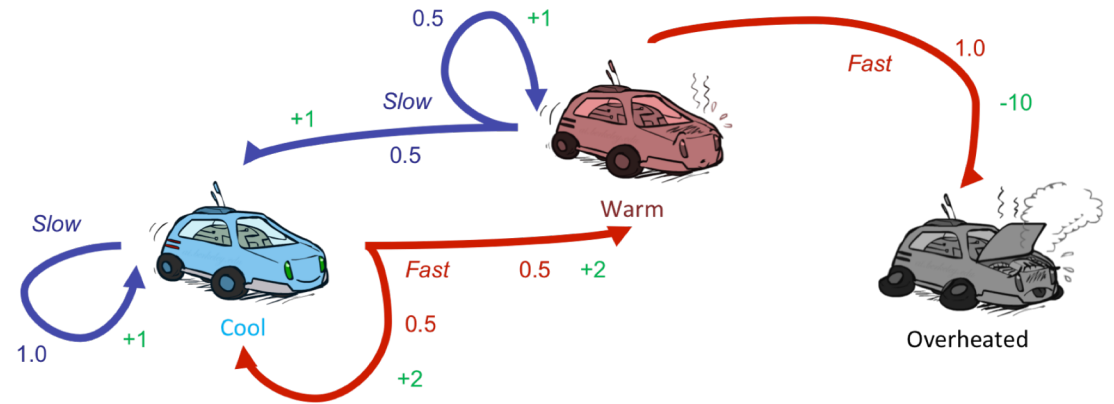*Fast in warm*

$$\boxed{1 \quad\quad -8 \quad\quad 0}$$

*Policy Improvement*

S: .5*(1+ $\gamma$ * 1)+.5*(1-$\gamma$ *8)
F: -10

*Improve policy for warm to: slow*



$$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$
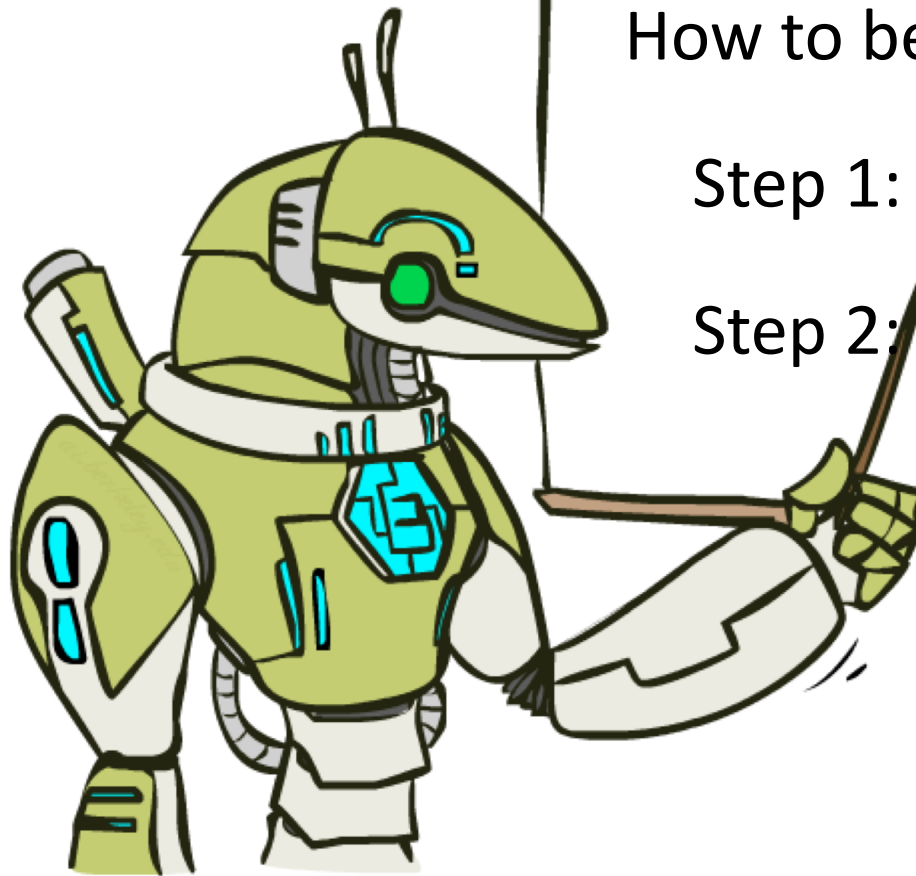
# Comparison

o Both value iteration and policy iteration compute the same thing (all optimal values)

o In value iteration:
  o Every iteration updates both the values and (implicitly) the policy
  o We don't track the policy, but taking the max over actions implicitly recomputes it

o In policy iteration:
  o We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
  o After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
  o The new policy will be better (or we're done)

o Both are dynamic programs for solving MDPs

# Summary: MDP Algorithms

o So you want to….

  o Compute optimal values: use value iteration or policy iteration
  o Compute values for a particular policy: use policy evaluation
  o Turn your values into a policy: use policy extraction (one-step lookahead)

o These all look the same!

  o They basically are – they are all variations of Bellman updates
  o They all use one-step lookahead expectimax fragments
  o They differ only in whether we plug in a fixed policy or max over actions

# The Bellman Equations



How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal