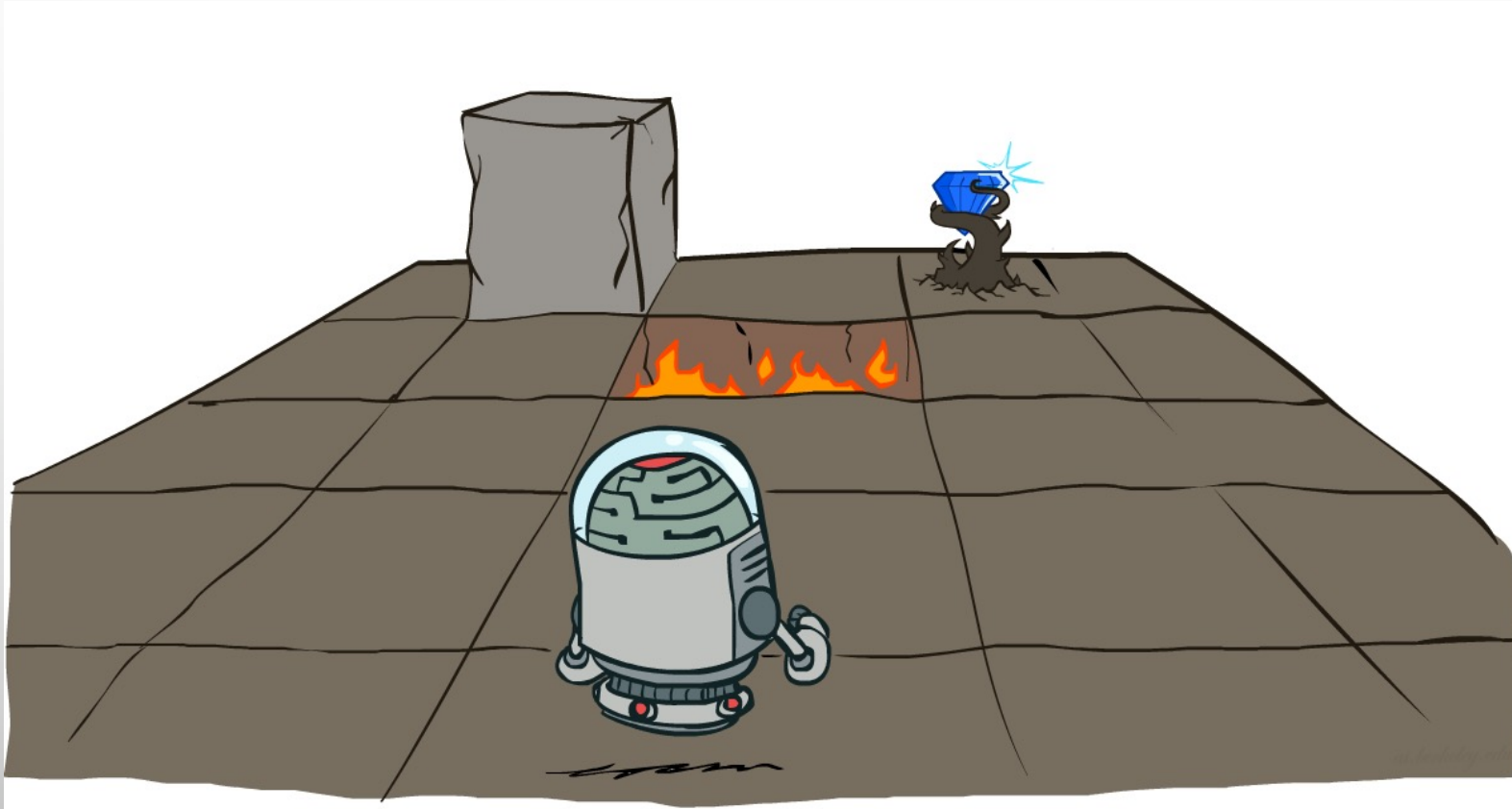# Artificial Intelligence
# CE-417, Group 1
# Computer Eng. Department
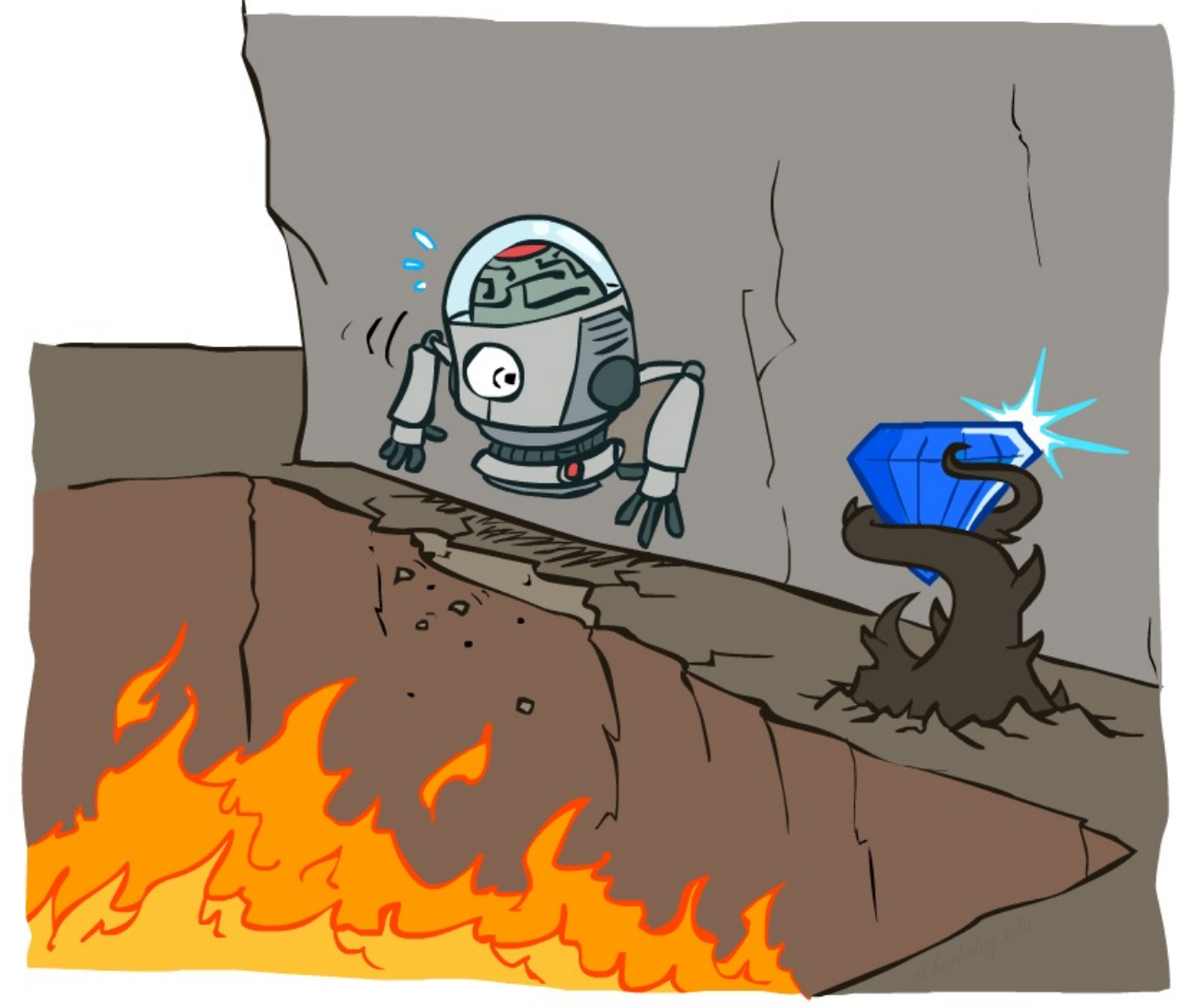# Sharif University of Technology

Fall 2023

By Mohammad Hossein Rohban, Ph.D.

Courtesy: Most slides are adopted from CSE-573 (Washington U.), original slides for the textbook, and CS-188 (UC. Berkeley).
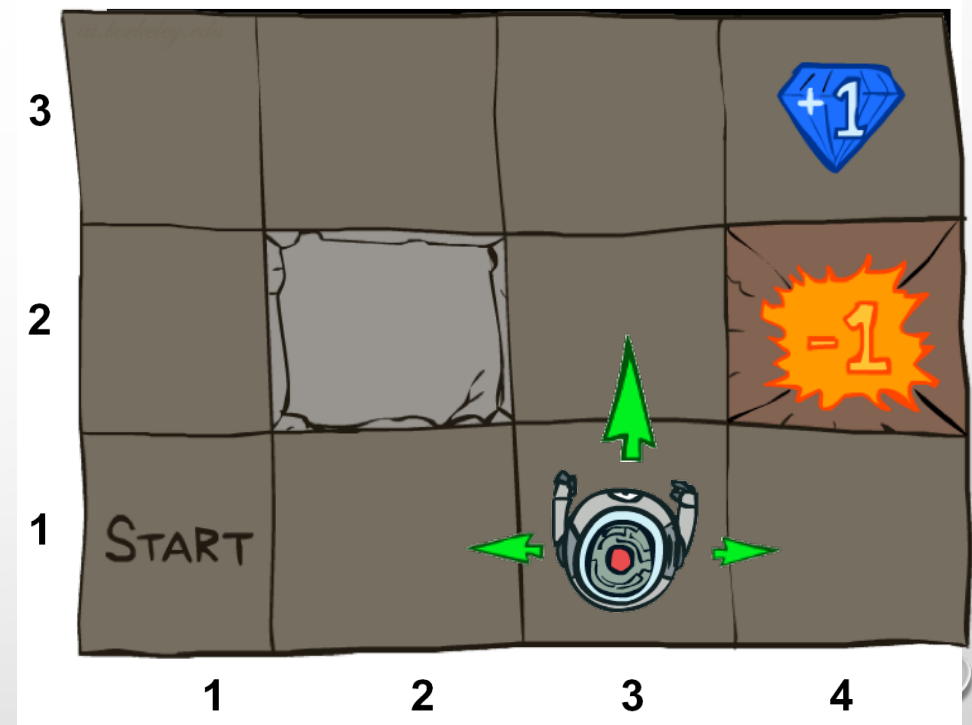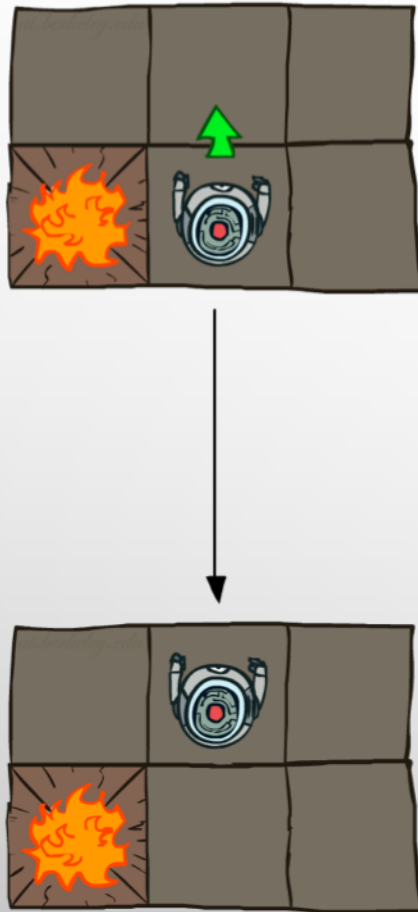
# Markov Decision Processes

# (MDPs)

# 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)
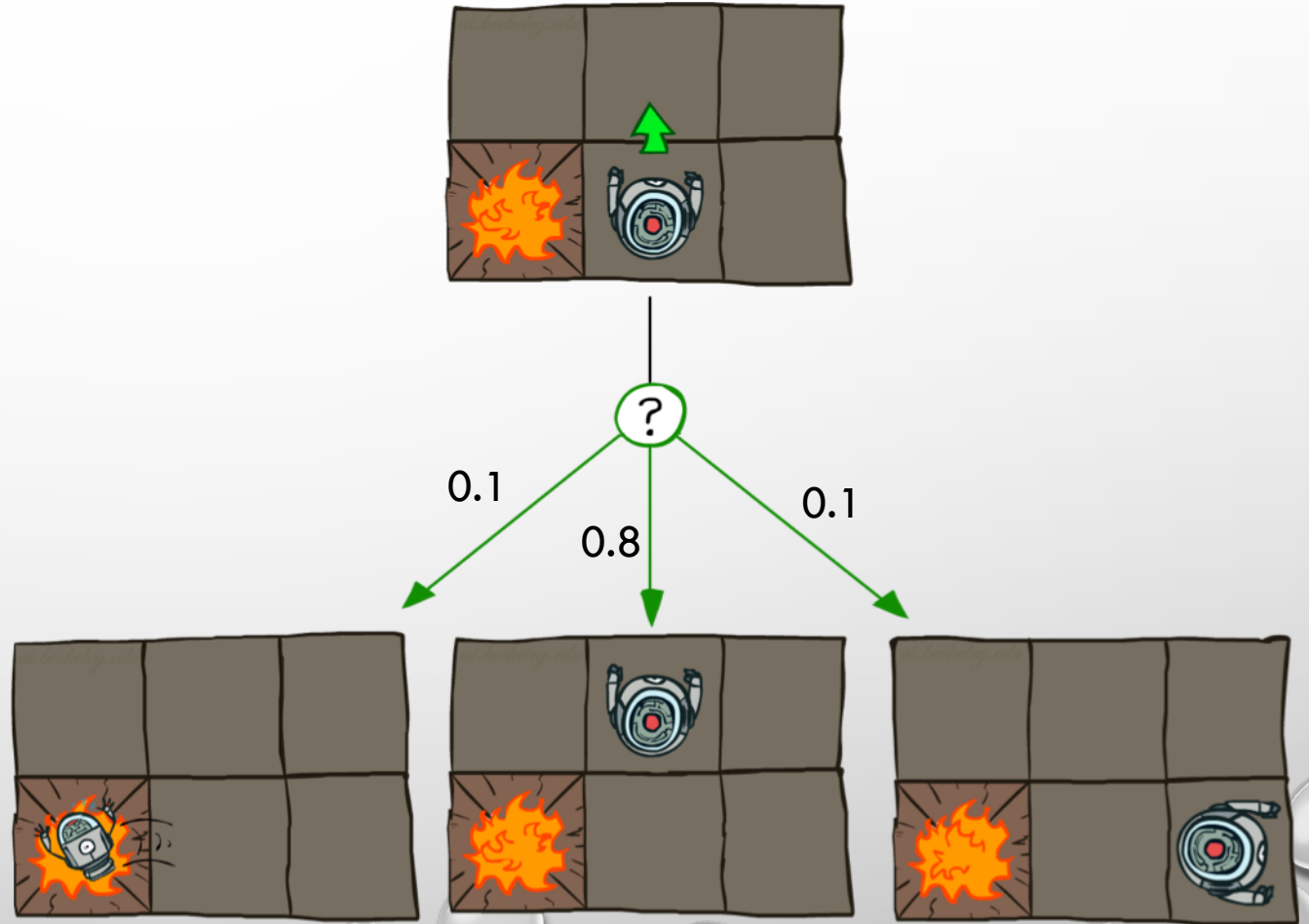
- Goal: maximize sum of rewards

# Grid World Actions

**Deterministic Grid World**

**Stochastic Grid World**

# Markov Decision Processes

- An MDP is defined by:
  - A set of states s ∈ S
  - A set of actions a ∈ A
  - A transition function T(s, a, s')
    - Probability that a from s leads to s', i.e., P(s' | s, a)
    - Also called the model or the dynamics
  - A reward function R(s, a, s')
    - Sometimes just R(s) or R(s')
  - A start state
  - Maybe a terminal state



**T Table**

$$T(s_{11}, E, \cdots$$
$$T(s_{31}, N, s_{11}) = 0$$
$$\vdots$$
$$T(s_{31}, N, s_{32}) = 0.8$$
$$T(s_{31}, N, s_{21}) = 0.1$$
$$T(s_{31}, N, s_{41}) = 0.1$$
$$\vdots$$

**R Table**

$$\cdots$$
$$R(s_{32}, N, s_{33}) = -0.01 \quad \text{(Breathing cost)}$$
$$\cdots$$
$$R(s_{32}, N, s_{33}) = -1.01$$
$$\cdots$$
$$R(s_{32}, N, s_{33}) = 0.99$$
$$\cdots$$

# Markov Decision Processes
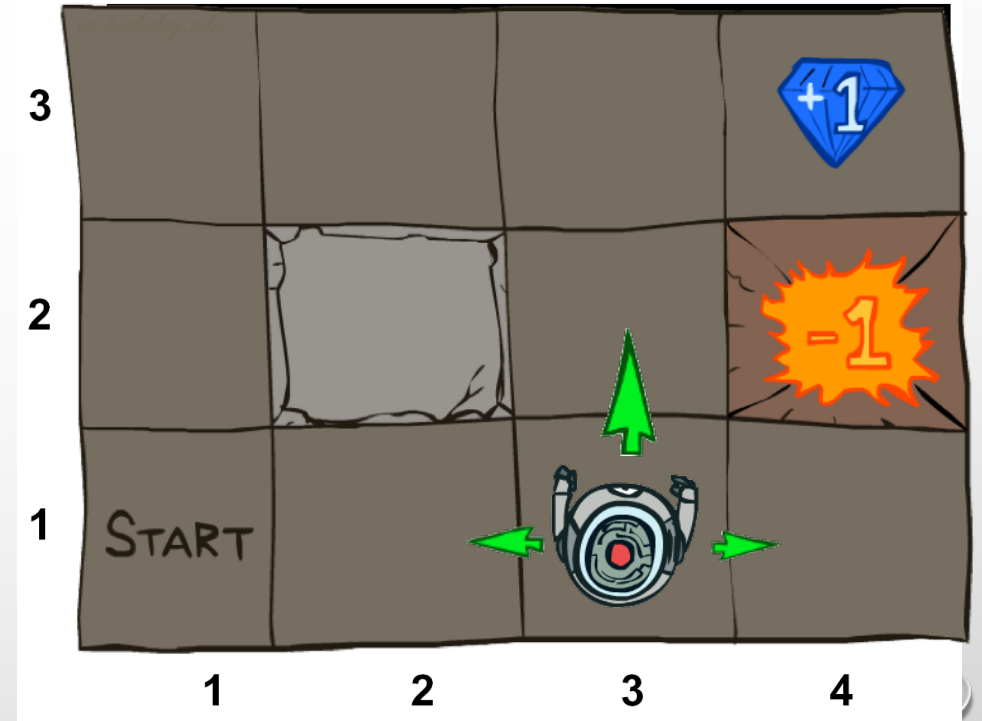
- An MDP is defined by:
  - A set of states $s \in S$
  - A set of actions $a \in A$
  - A transition function $T(s, a, s')$
    - Probability that a from s leads to s', i.e., $P(s' \mid s, a)$
    - Also called the model or the dynamics
  - A reward function $R(s, a, s')$
    - Sometimes just $R(s)$ or $R(s')$
  - A start state
  - Maybe a terminal state

- MDPs are non-deterministic search problems
  - One way to solve them is with expectimax search
  - We'll have a new tool soon

# What is Markov about MDPs?

- "Markov" generally means that given the present state, the future and the past are independent

- For Markov decision processes, "Markov" means action outcomes depend only on the current state

$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \ldots S_0 = s_0)$$

$$=$$
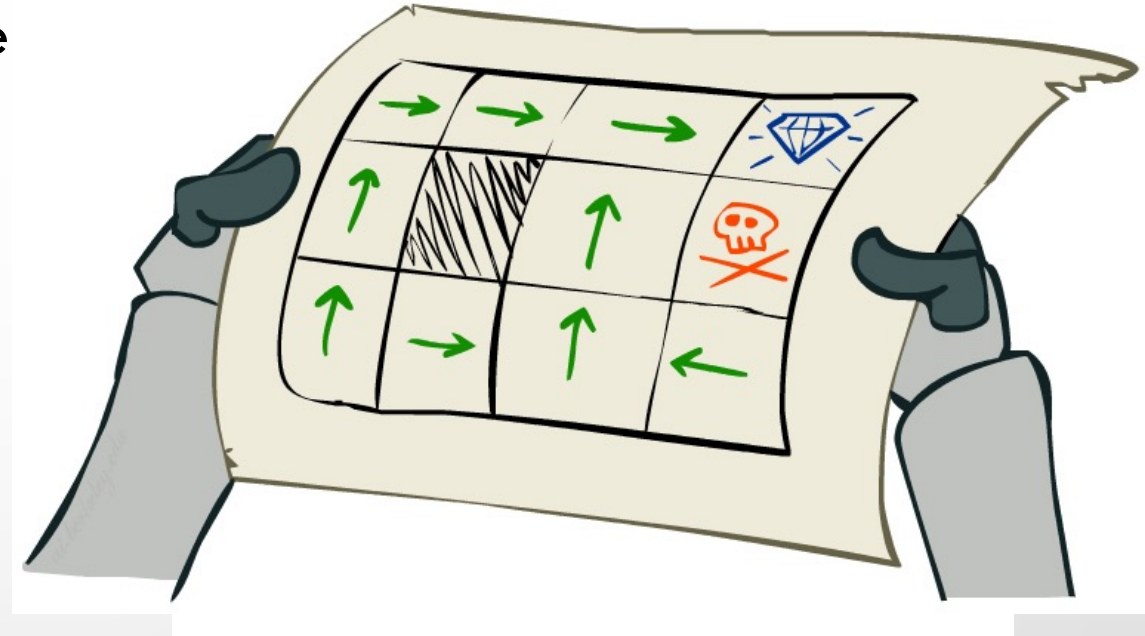
$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

Andrey Markov
(1856-1922)

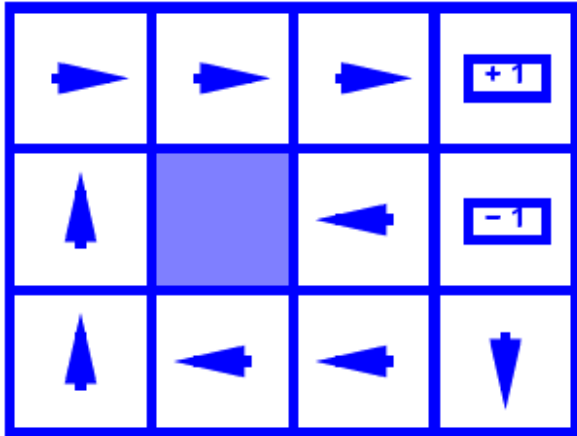- This is just like search, where the successor function could only depend on the current state (not the history)

# Policies

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal

- For MDPs, we want an optimal policy $\pi^*$: $S \rightarrow A$

  - A policy $\pi$ gives an action for each state

  - An optimal policy is one that maximizes expected utility if followed

- Expectimax didn't compute entire policies

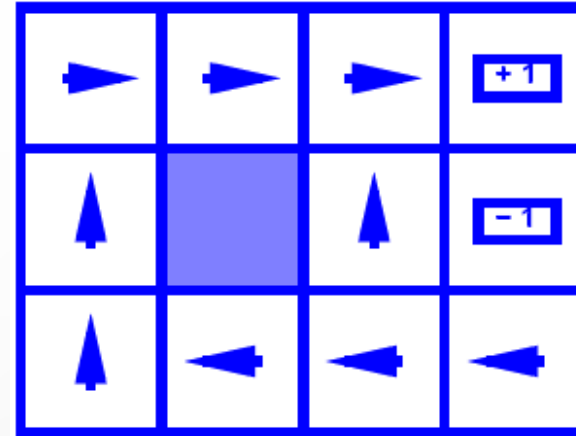  - It computed the action for a single state only



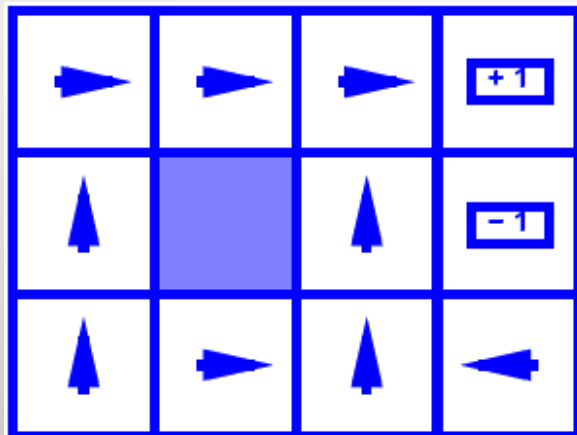Optimal policy when R(s, a, s') = -0.03
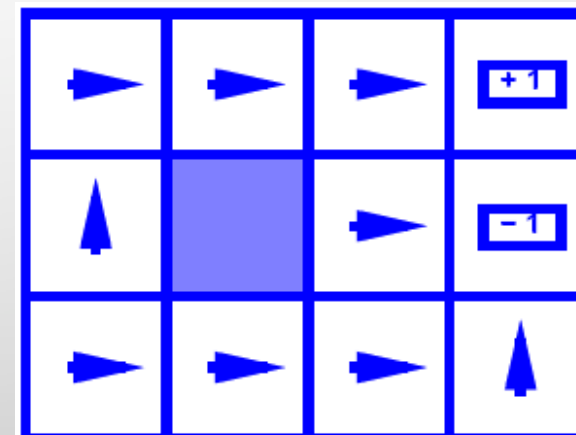for all non-terminals s

# Optimal Policies



R(s) = -0.01

R(s) = -0.03
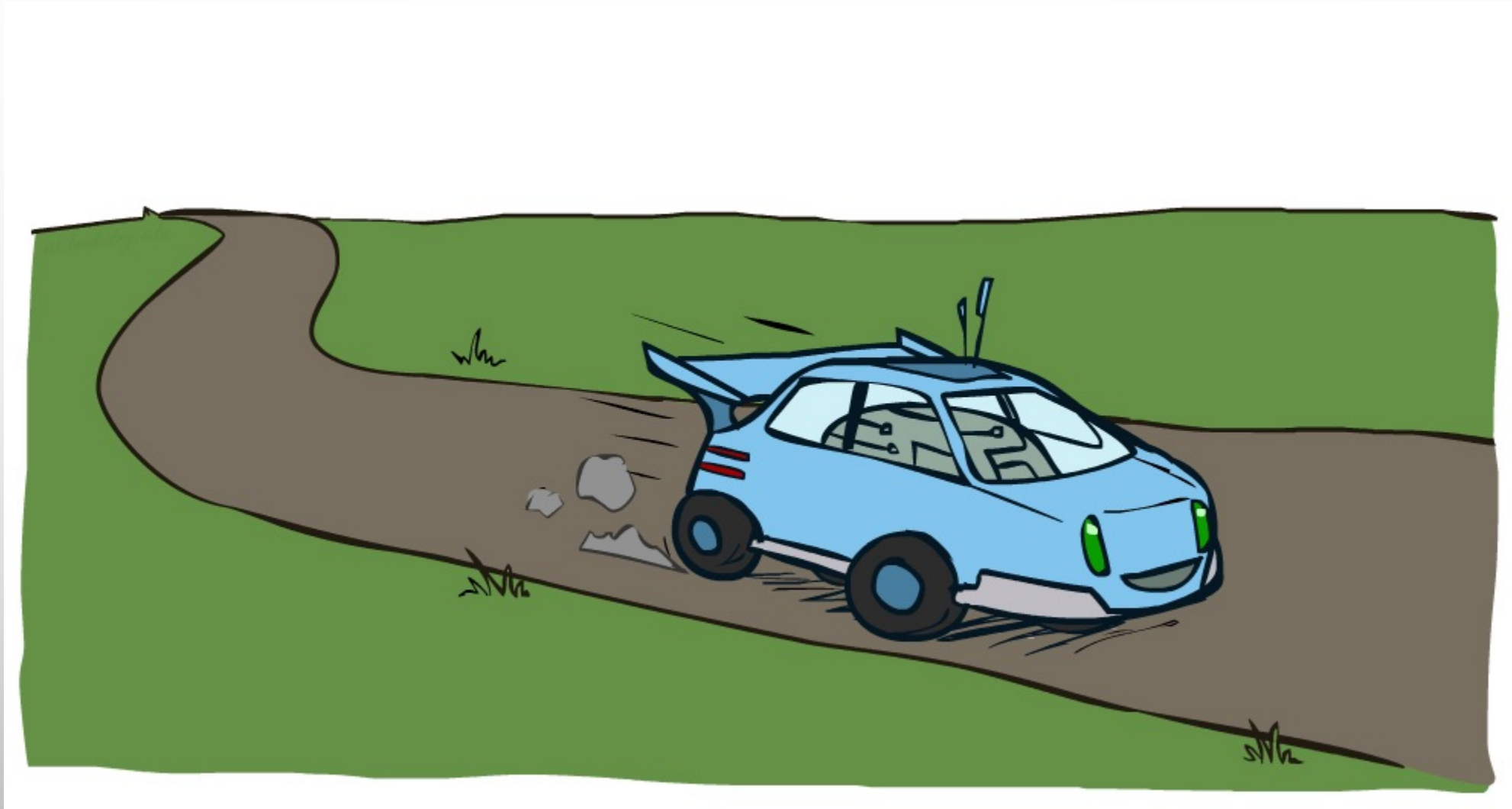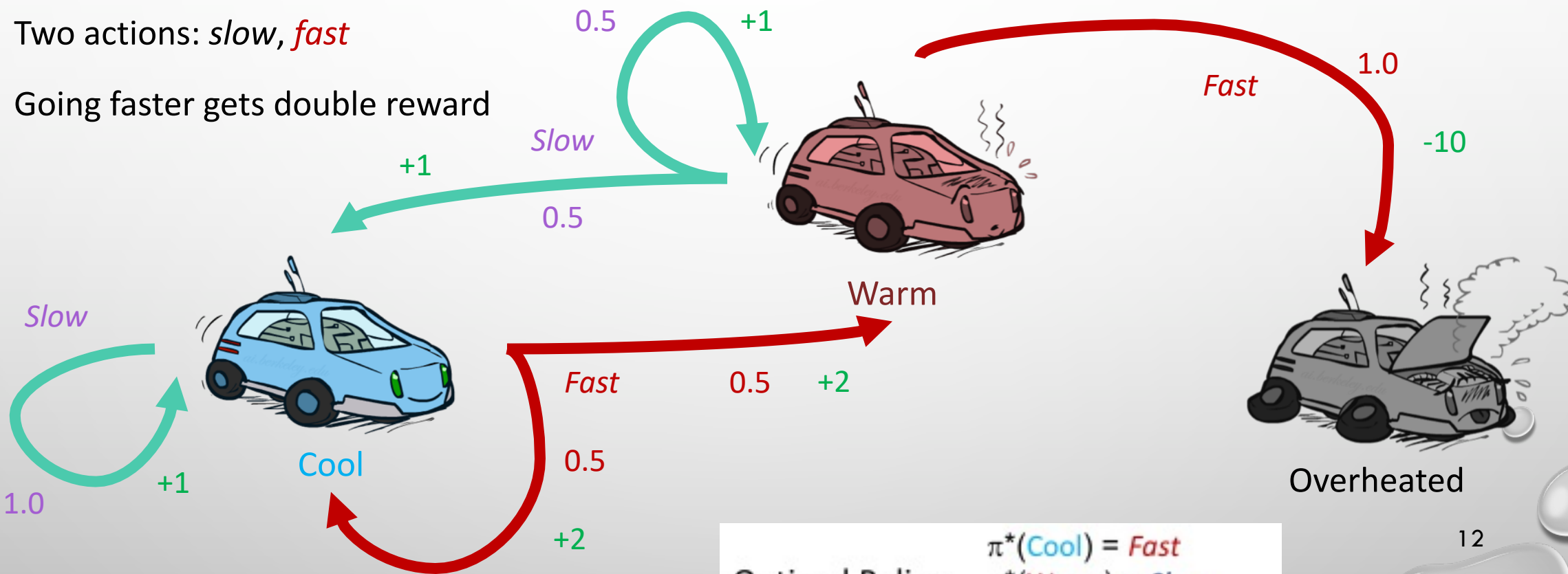
R(s) = -0.4

R(s) = -2.0

# Example: Racing

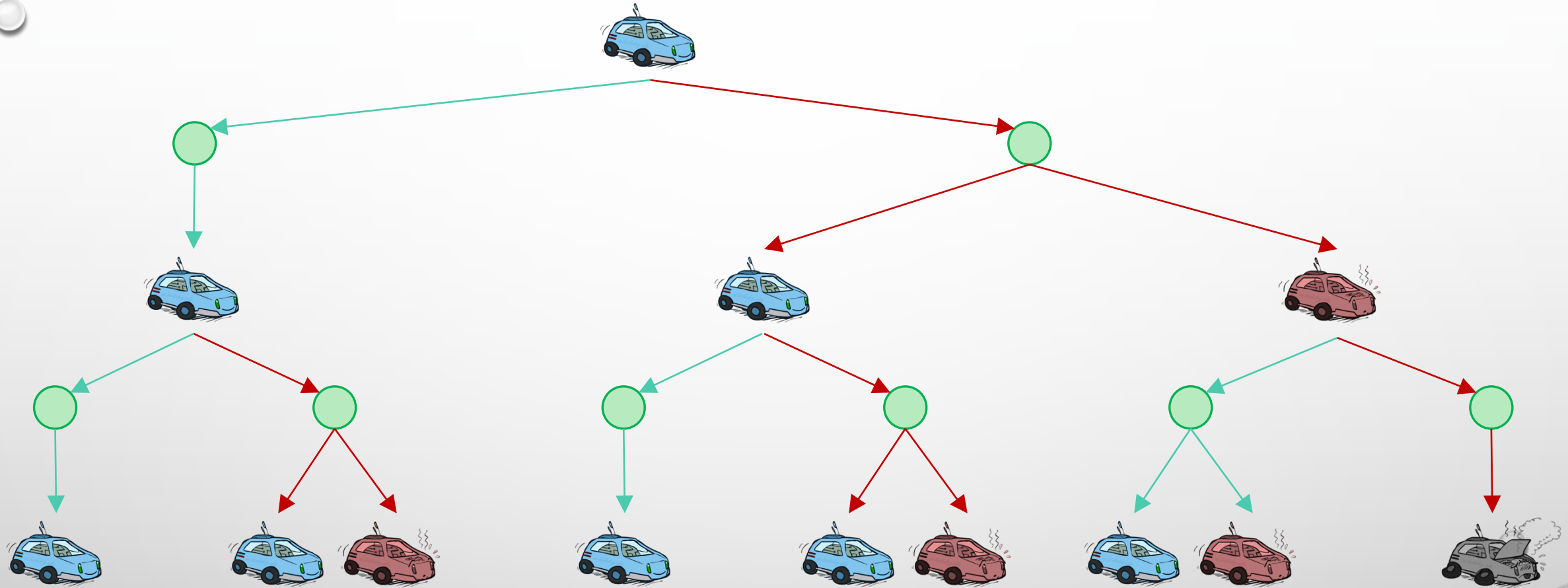# Example: Racing

- A robot car wants to travel far, quickly

- Three states: cool, warm, overheated

- Two actions: *slow*, *fast*

- Going faster gets double reward



0.5   +1

*Fast*   1.0

*Slow*   -10

+1

0.5

*Slow*

+1

1.0

Cool

Warm

Overheated

*Fast*   0.5   +2

0.5

+2

Optimal Policy:  $\pi^*(Cool) = Fast$
$\pi^*(Warm) = Slow$
$\pi^*(Overheated) = end$

# Racing Search Tree

# MDP Search Trees

- Each MDP state projects an expectimax-like search tree

s **s is a** *state*

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

s, a

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

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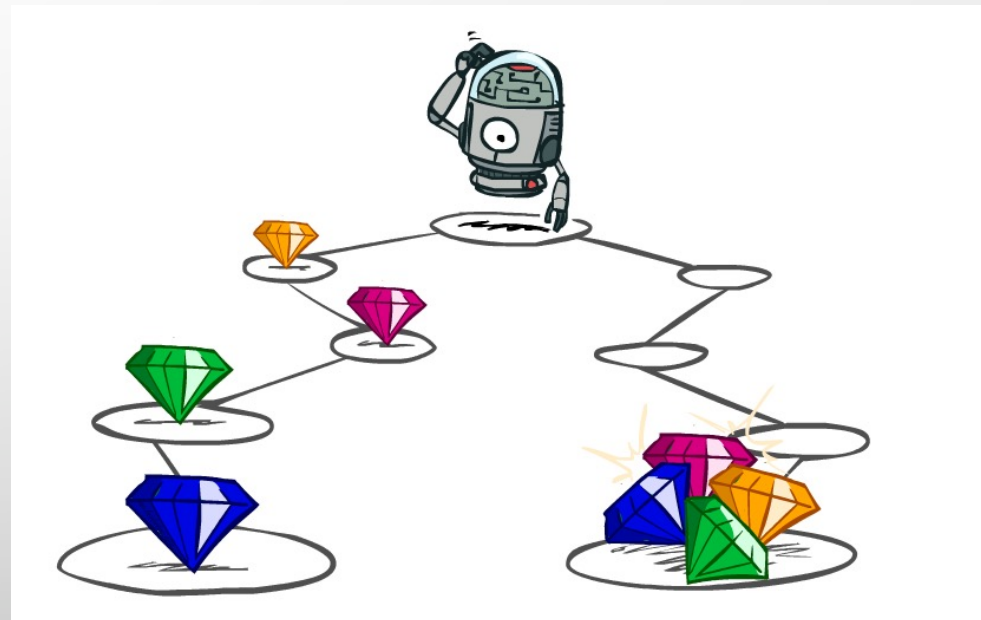
R(s,a,s' )

s,a,s'

s'

# Utilities of Sequences

# Utilities of Sequences
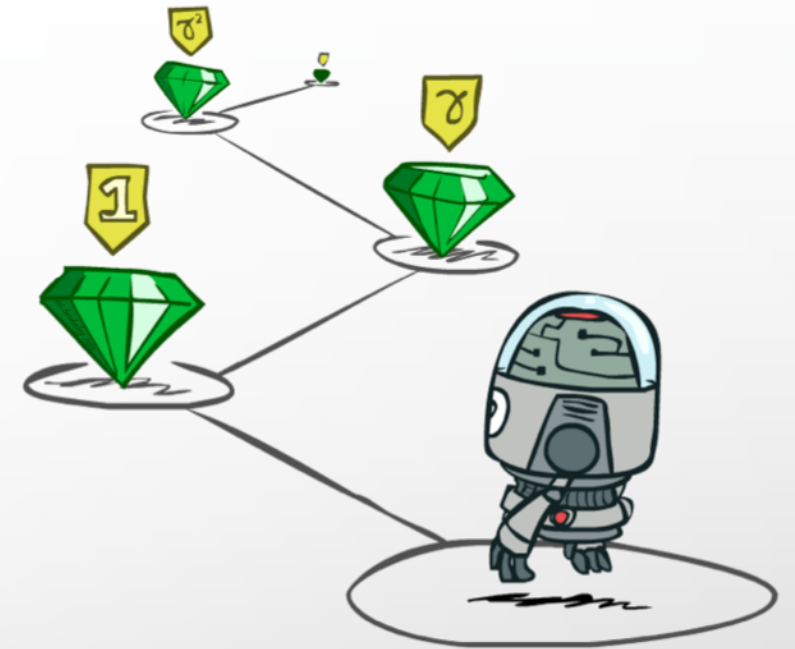
- What preferences should an agent have over reward sequences?

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

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

# Stationary Preferences

- In order to formalize optimality of a policy, we need assumption about preferences remaining the same independent of time.

- <span style="color:red">If you prefer one future to another starting tomorrow, then you should still prefer that future if it were to start today:</span>

$$[r, r_0, r_1, r_2, \ldots] \succ [r, r'_0, r'_1, r'_2, \ldots]$$
$$\Leftrightarrow$$
$$[r_0, r_1, r_2, \ldots] \succ [r'_0, r'_1, r'_2, \ldots]$$

- Given stationary preferences, there are two ways to assign utilities to sequences:

  - Additive utility: $U([r_0, r_1, r_2, \ldots]) = r_0 + r_1 + r_2 + \cdots$

  - Discounted utility: $U([r_0, r_1, r_2, \ldots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \cdots$

# 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

$$1$$

$$\gamma$$

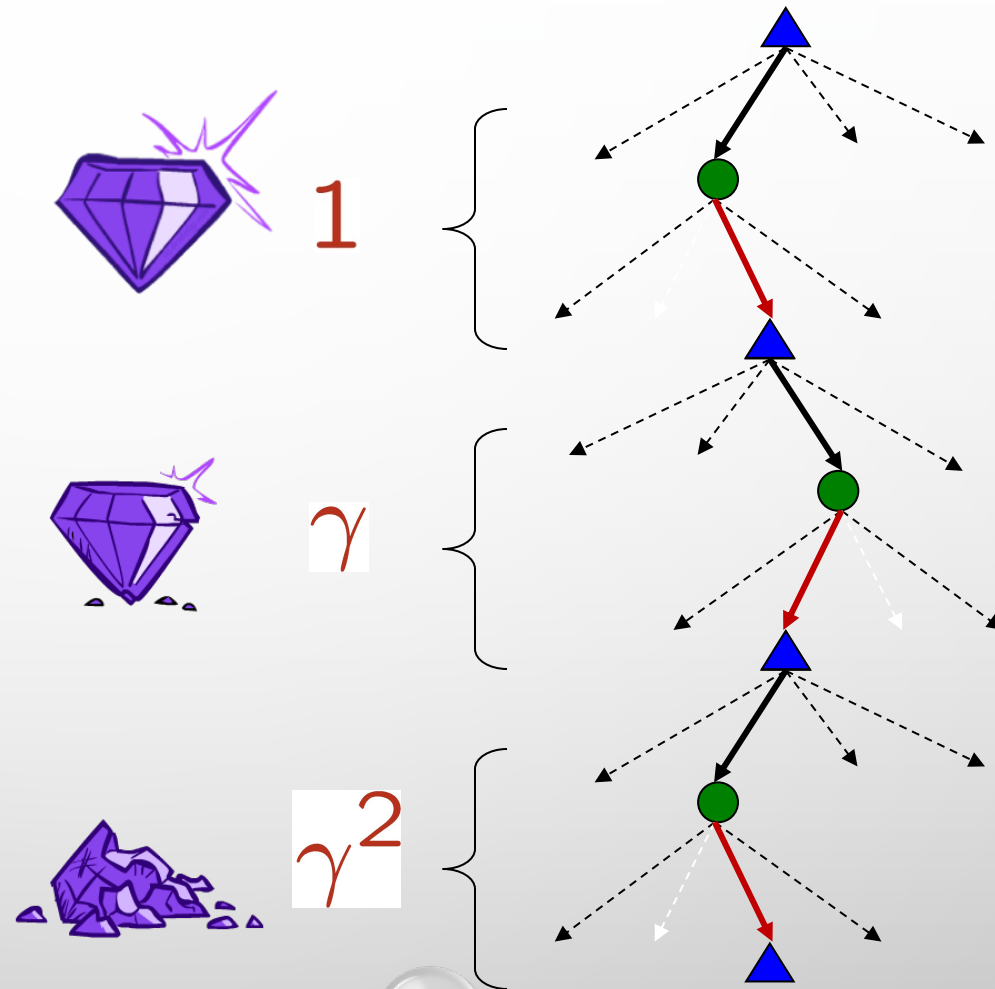$$\gamma^2$$
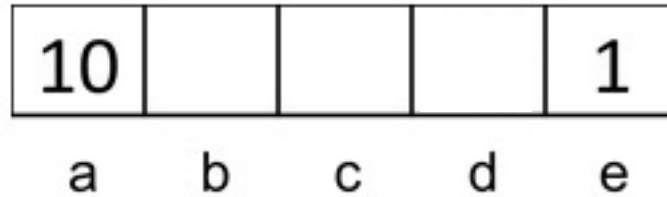
Worth Now          Worth Next Step          Worth In Two Steps

# Discounting

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

- Why discount?
  - Sooner rewards probably do have higher utility than later rewards
  - Think of it as a gamma chance of ending the process at every step(chance of death!)
  - Also helps our algorithms converge

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

$1$

$\gamma$

$\gamma^2$

# Quiz: Discounting

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

- Given:

  - Actions: east, west, and exit (only available in exit states a, e)
  - Transitions: deterministic

- Quiz 1: for $\gamma = 1$, what is the optimal policy?

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

- Quiz 2: for $\gamma = 0.1$, what is the optimal policy?

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

- Quiz 3: for which $\gamma$ are west and east equally good when in state d?

# Infinite Utilities?!

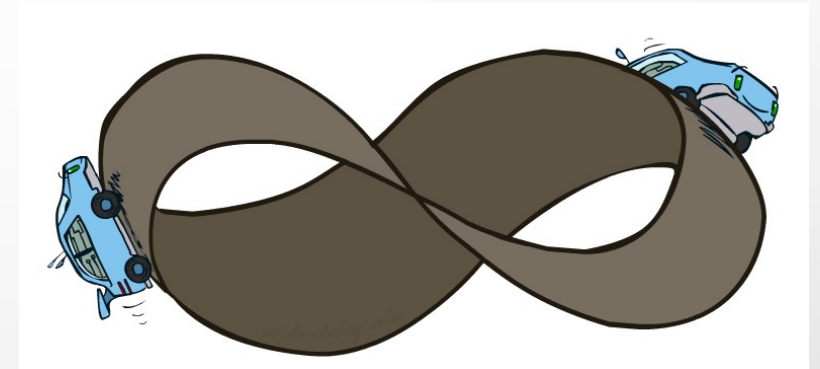- Problem: what if the game lasts forever?  Do we get infinite rewards?

- Solutions:
  - Finite horizon: (similar to depth-limited search)
    - Terminate episodes after a fixed T steps (e.g. Life)
    - Gives nonstationary policies ($\pi$ depends on time left)

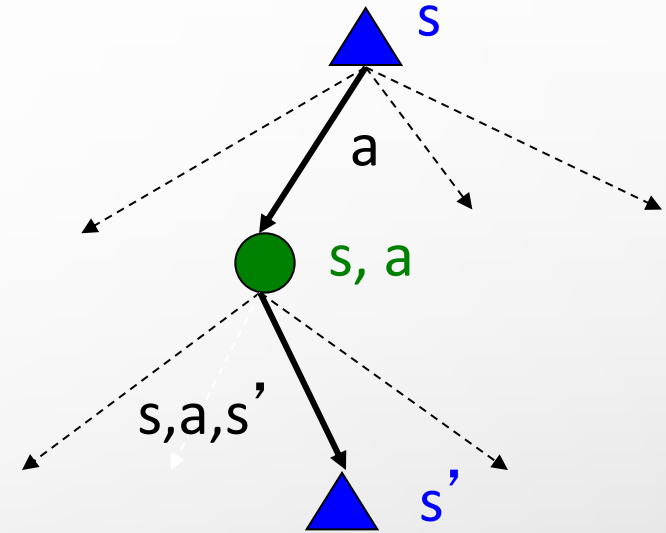  - Discounting: use $0 < \gamma < 1$

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

    - Smaller $\gamma$ means smaller "horizon" – shorter term focus

  - Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "overheated" for racing)
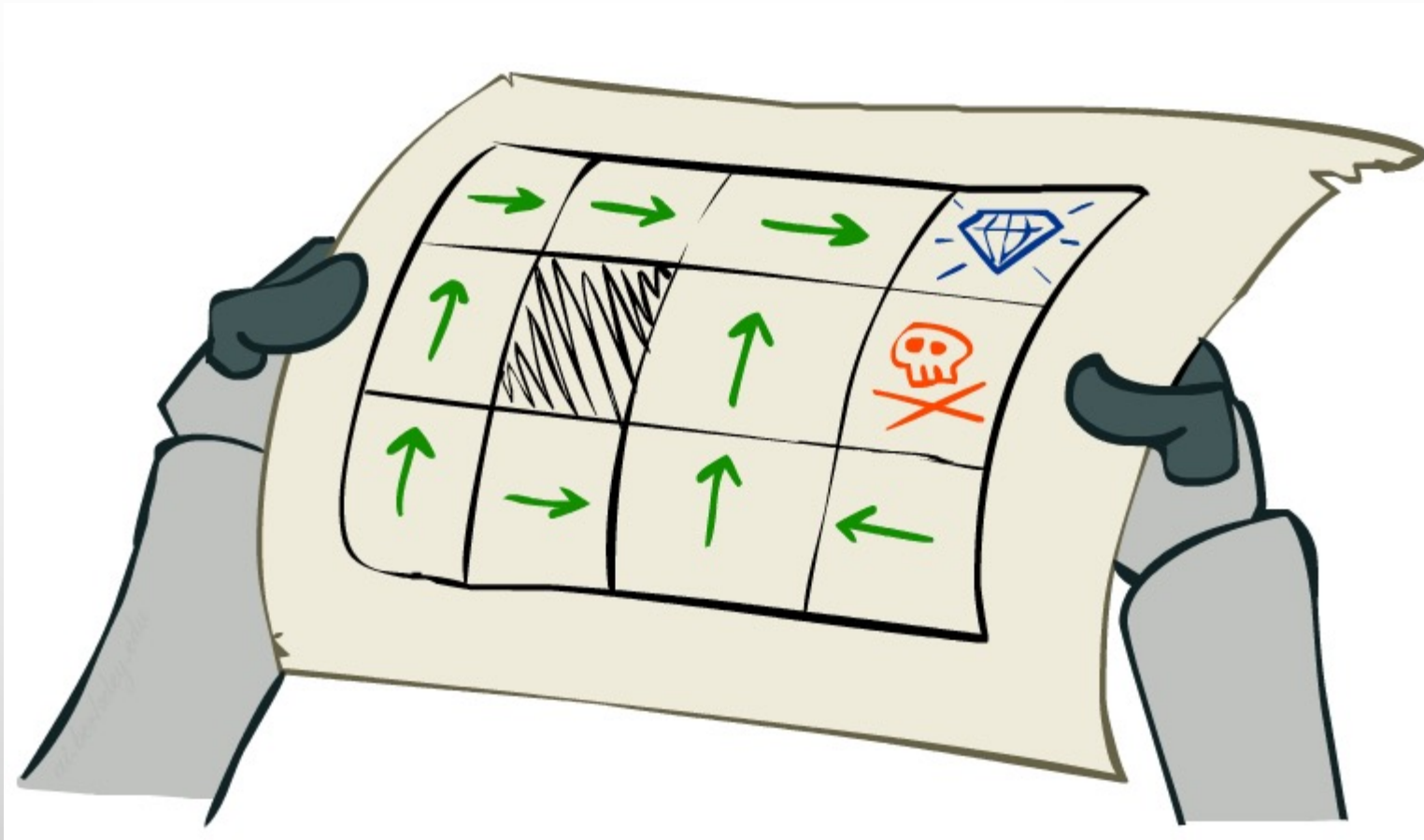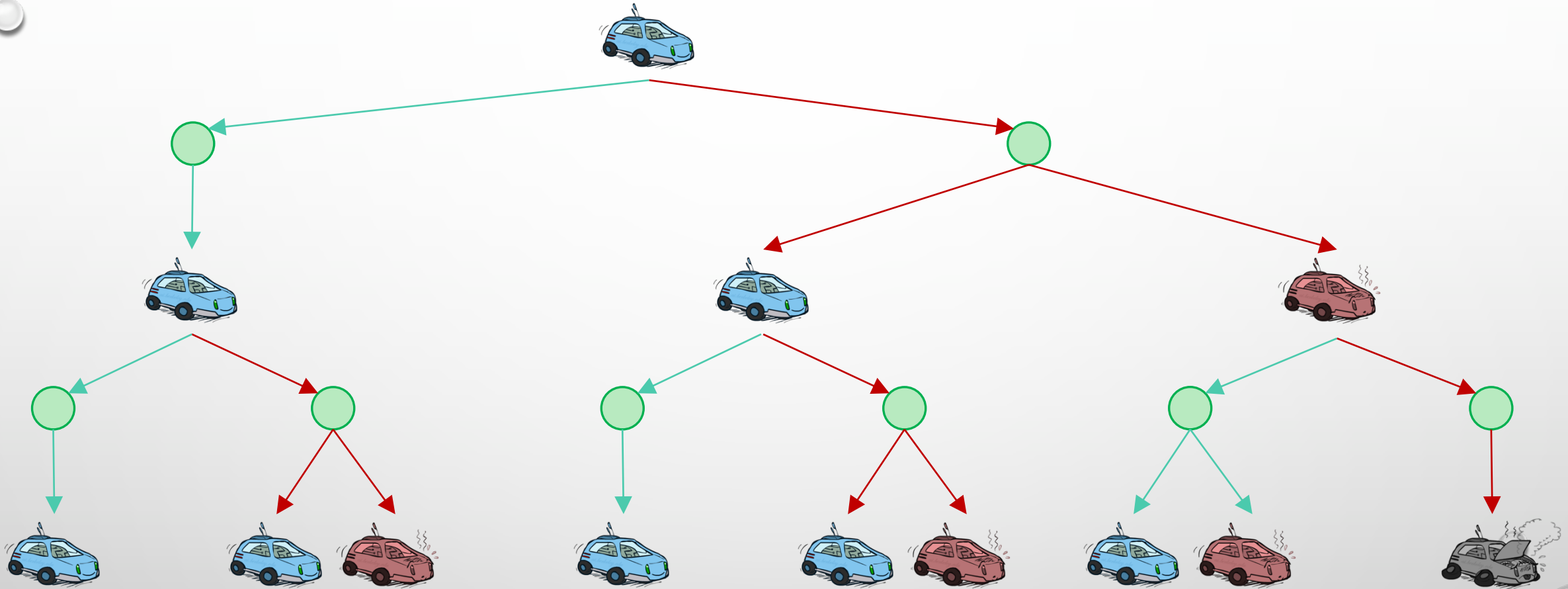
# Recap: Defining MDPs

- 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
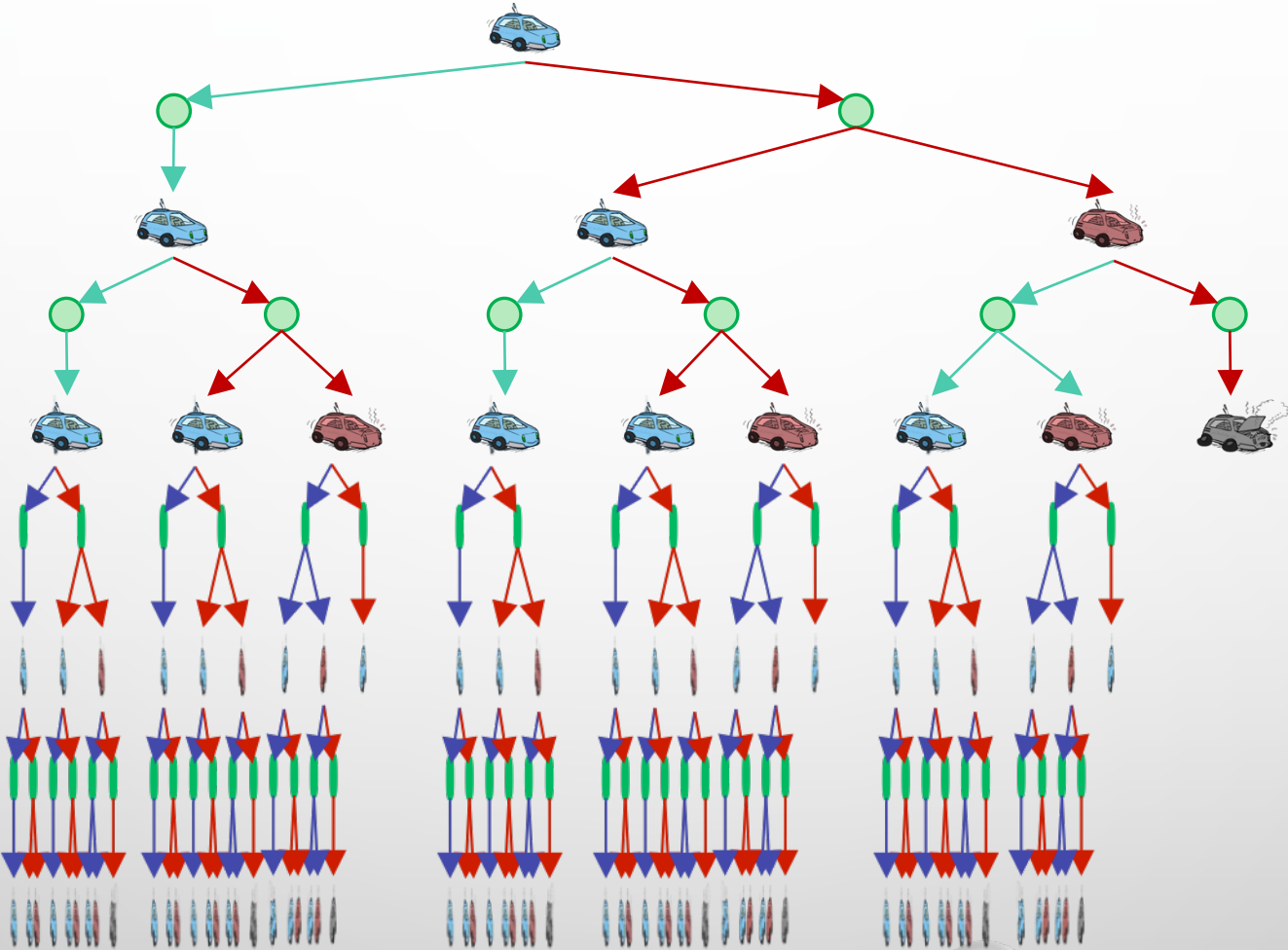
s

a

s, a

s,a,s'

s'
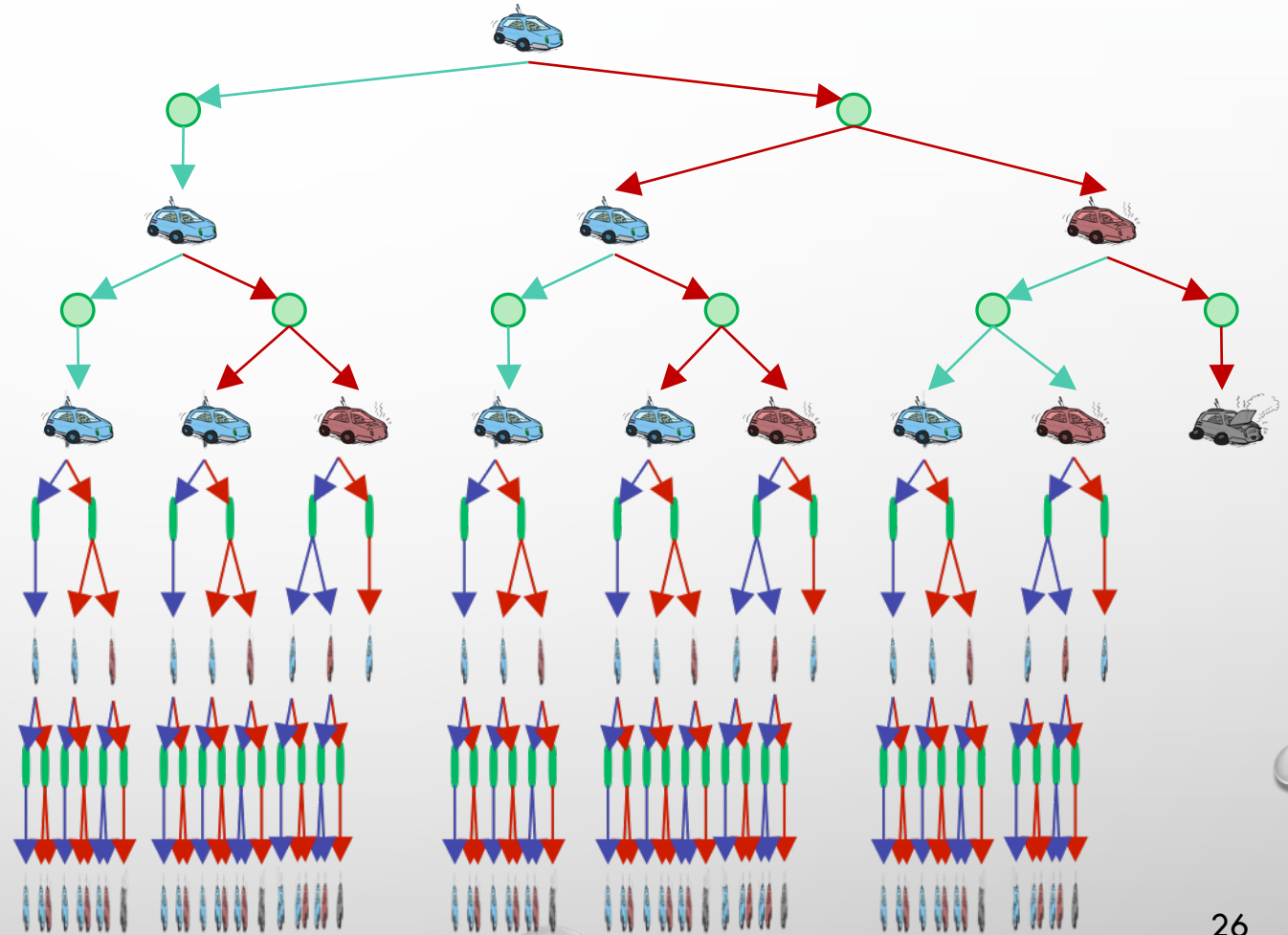
# Solving MDPs

# Recall: Racing Search Tree

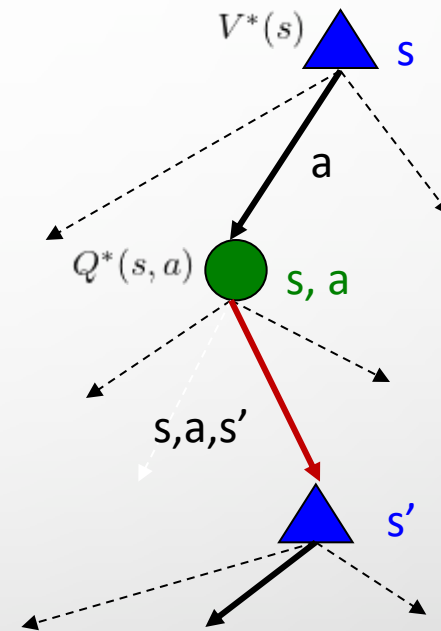# Racing Search Tree

# Racing Search Tree

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

- Problem: states are repeated
  - Idea: Only compute needed quantities once, cache the rest in a lookup table

- 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 $\gamma < 1$

# 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

$V^*(s)$  ▲ s

a

$Q^*(s,a)$  ● s, a

s,a,s'

▲ s'

s is a *state*

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

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

# Snapshot of Demo – Gridworld V Values



Gridworld Display

| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.57 | | ▲ 0.57 | -1.00 |
| ▲ 0.49 | ◀ 0.43 | ▲ 0.48 | ◀ 0.28 |

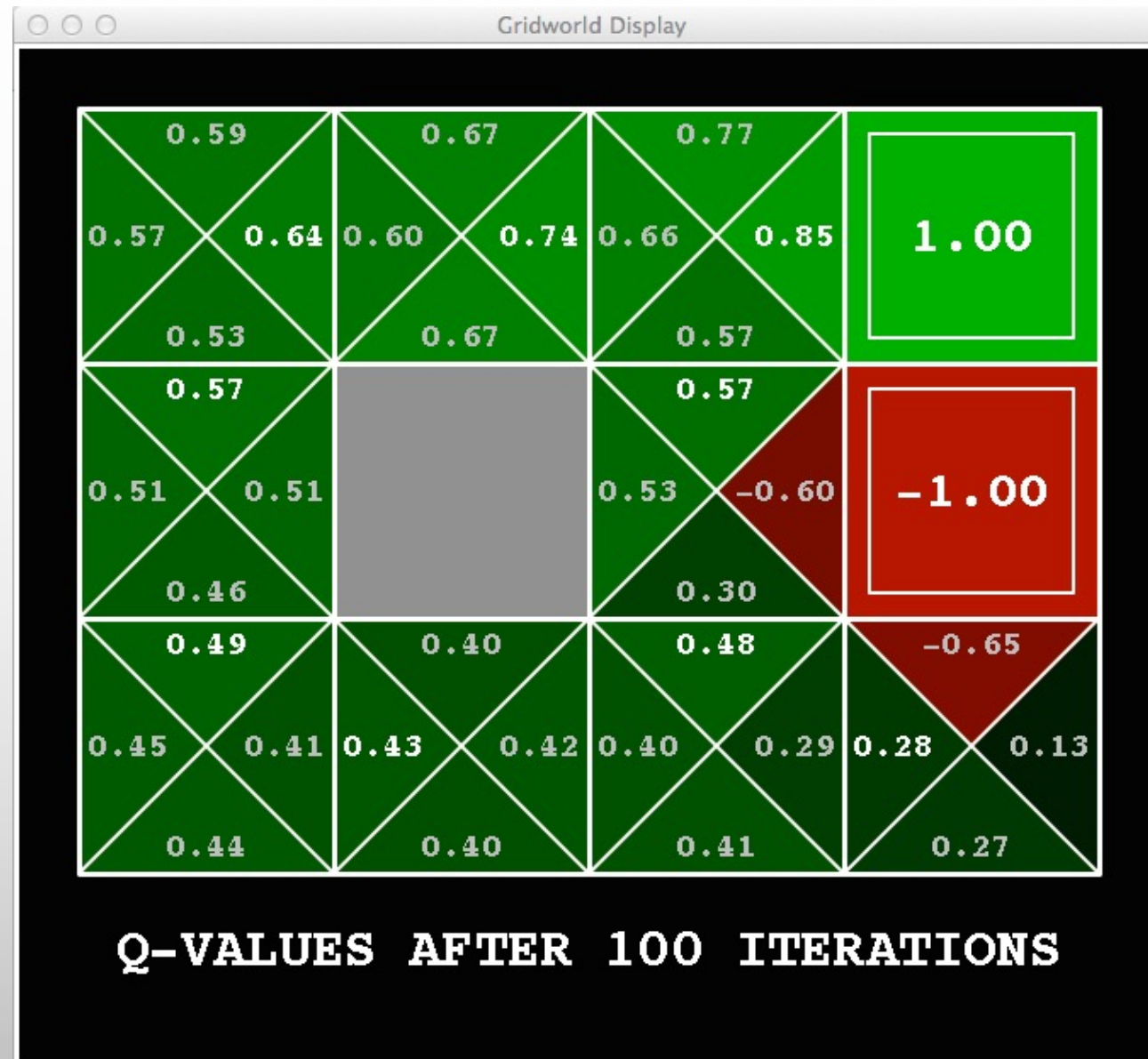VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# Snapshot of Demo – Gridworld Q Values



Noise = 0.2
Discount = 0.9
Living reward = 0
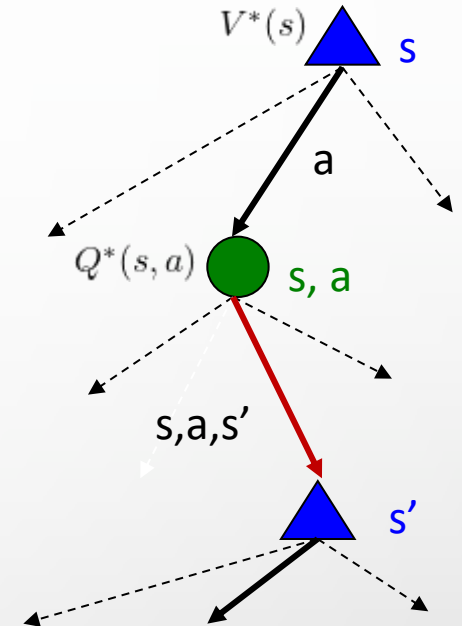
# Values of States (The Bellman Equations)

- Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values
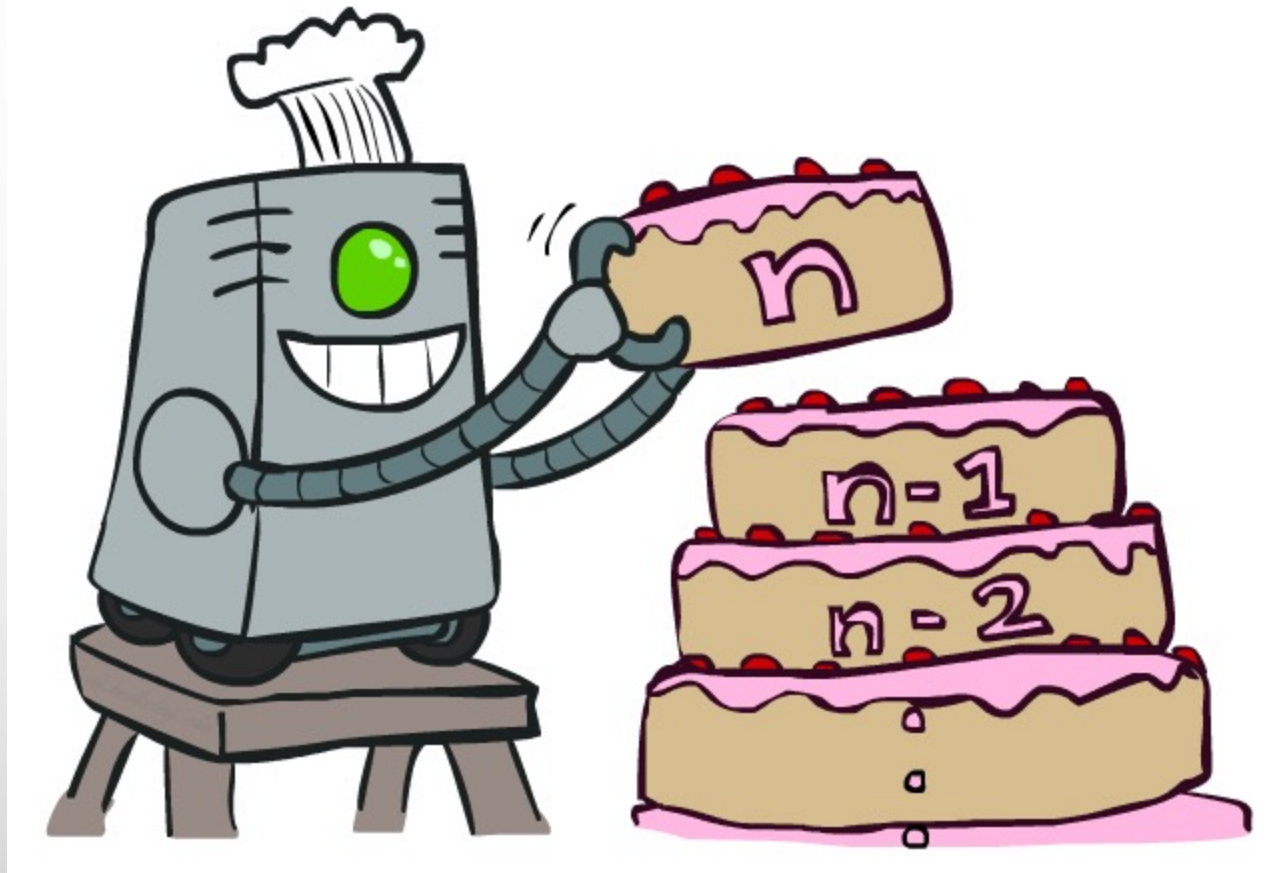
$$V^*(s) = \max_a Q^*(s,a)$$

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

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

- These are the bellman equations, and they characterize optimal values in a way we'll use over and over
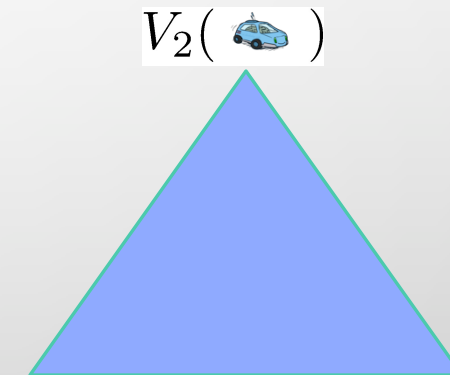
- But how do we solve these equations?

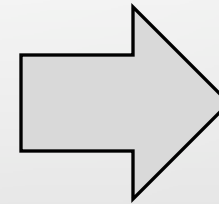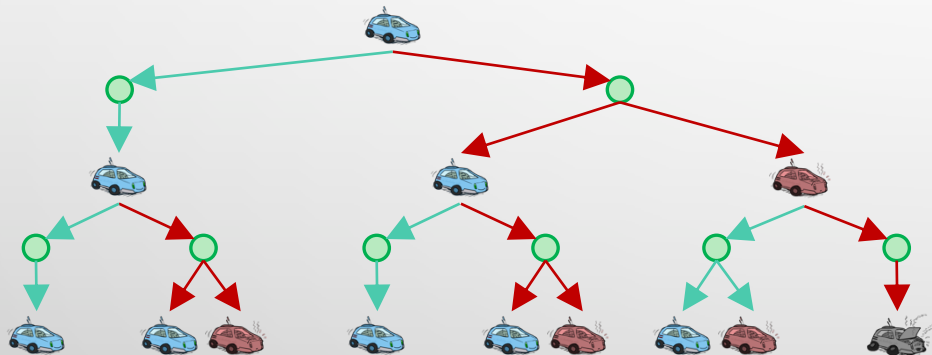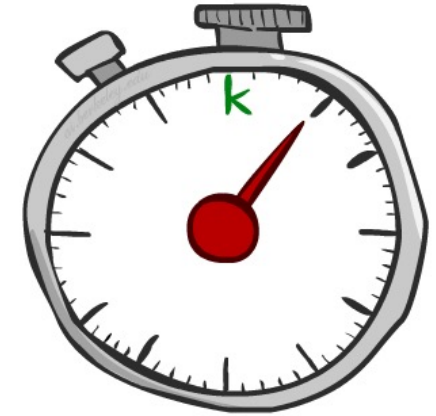# Value Iteration

# Another View: Time-Limited Values

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

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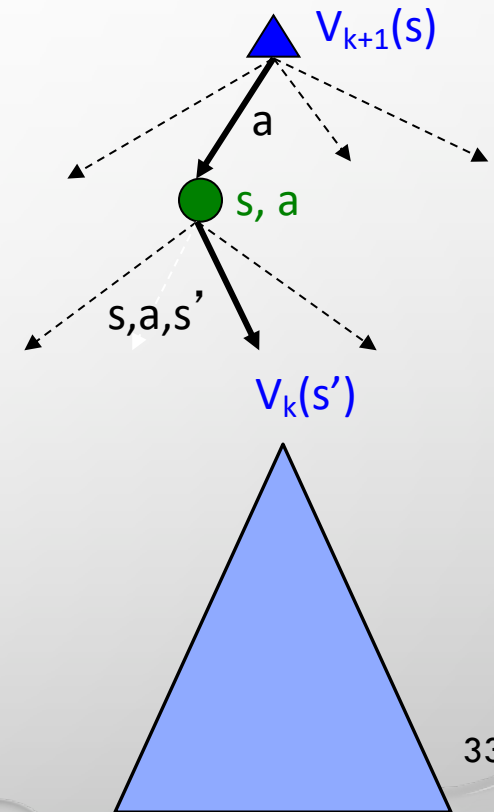


$$V_2(\text{🚙})$$

# Value Iteration

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

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

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

- Repeat until convergence

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

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

$V_{k+1}(s)$

a

s, a

s,a,s'

$V_k(s')$

# k=0

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=1



Noise = 0.2
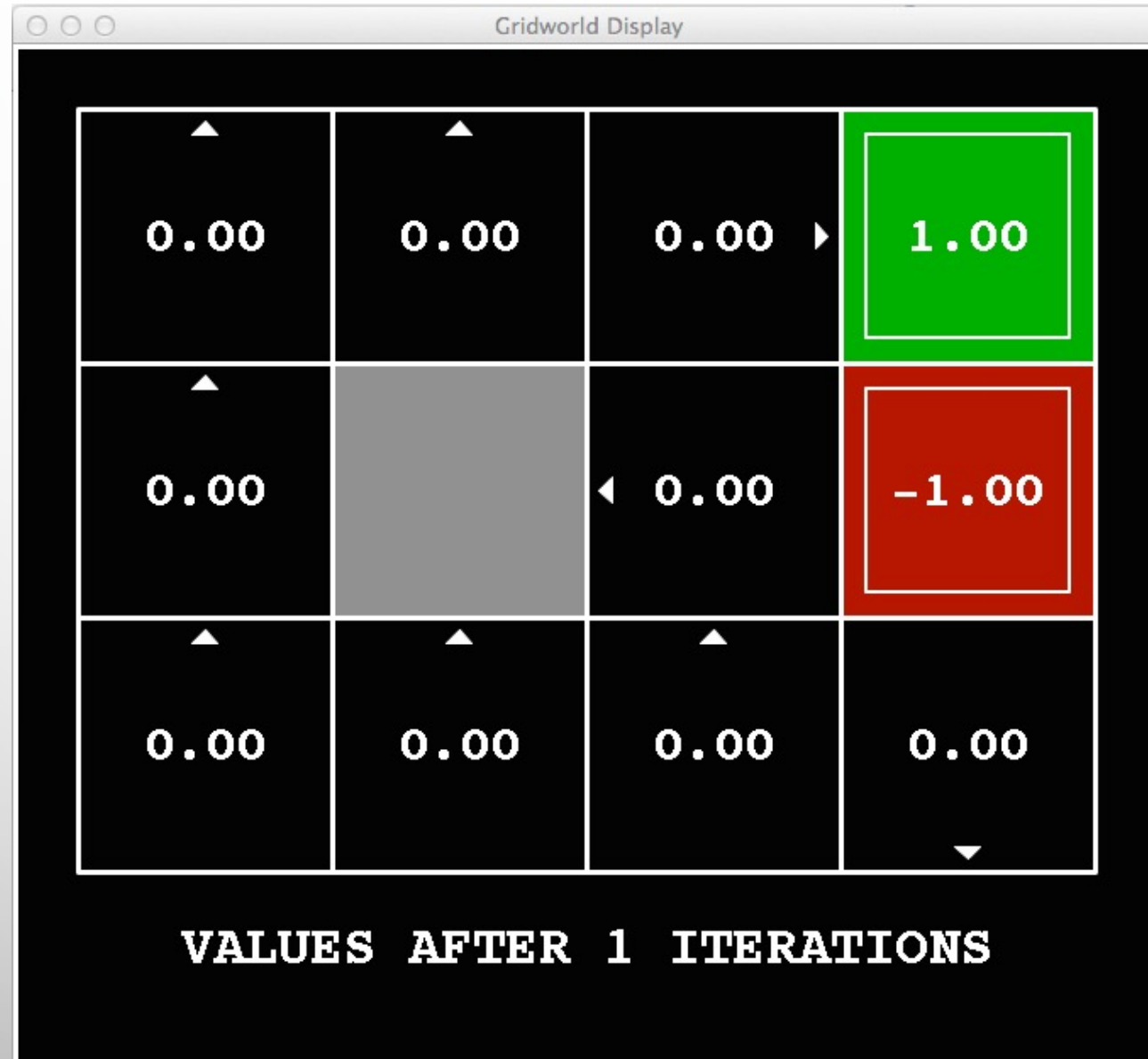Discount = 0.9
Living reward = 0

# k=2

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

38

# k=5



VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

42

# k=9



VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

43

# k=10



Noise = 0.2
Discount = 0.9
Living reward = 0

44

# k=11



VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
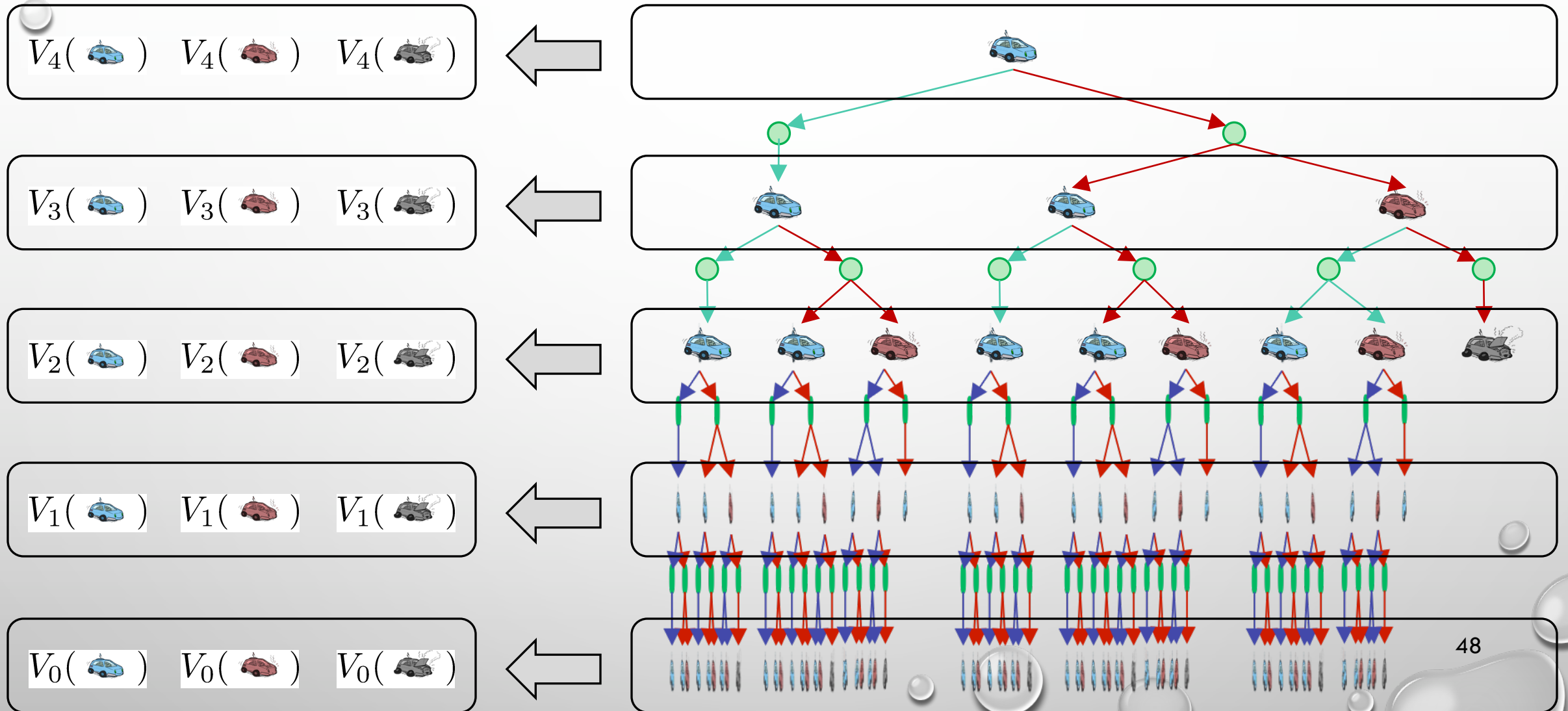Living reward = 0

45

# k=12



VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

46

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# Computing Time-Limited Values



$V_4(\ )$  $V_4(\ )$  $V_4(\ )$

$V_3(\ )$  $V_3(\ )$  $V_3(\ )$

$V_2(\ )$  $V_2(\ )$  $V_2(\ )$

$V_1(\ )$  $V_1(\ )$  $V_1(\ )$

$V_0(\ )$  $V_0(\ )$  $V_0(\ )$

48

# Example: Value Iteration
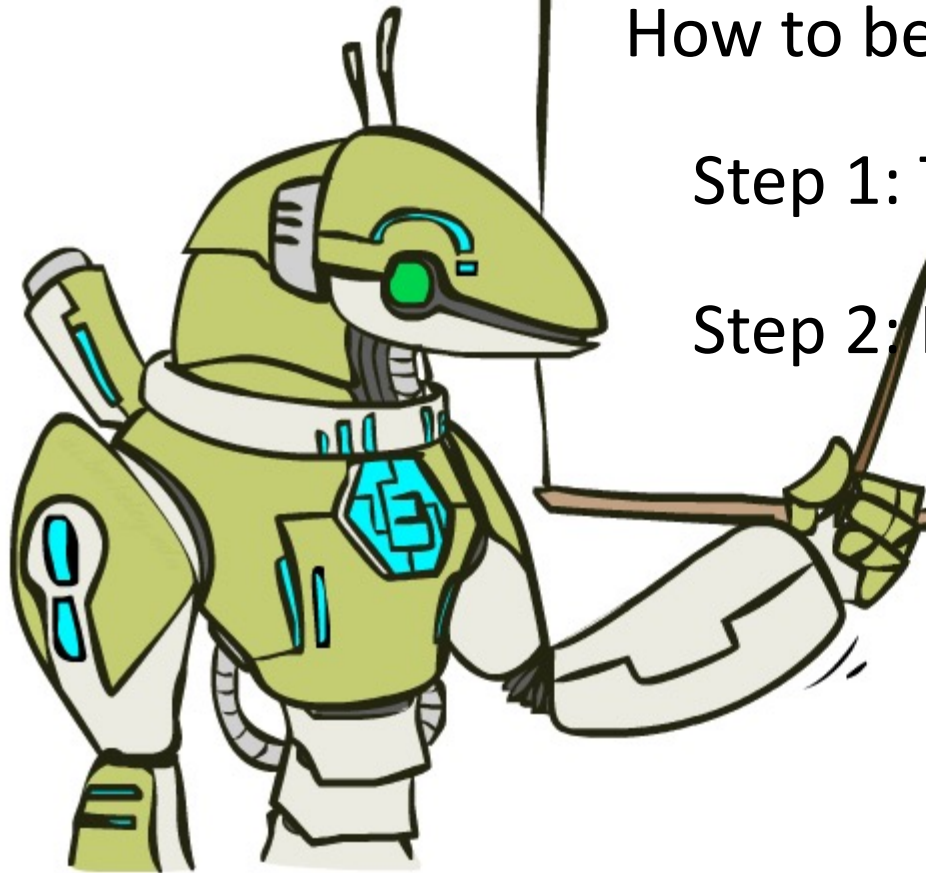


$V_2$ : 3.5    2.5    0

$V_1$ : 2    1    0

$V_0$ : 0    0    0

*Assume no discount!*

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

# Recap: The Bellman Equations

How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal
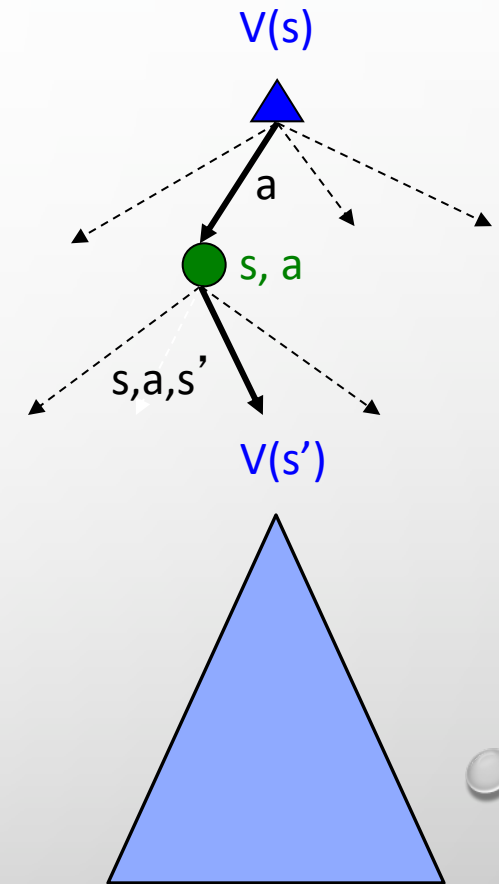
# Value Iteration? The Bellman Equations?

- Bellman equations <span style="color:red">characterize</span> the optimal values:

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

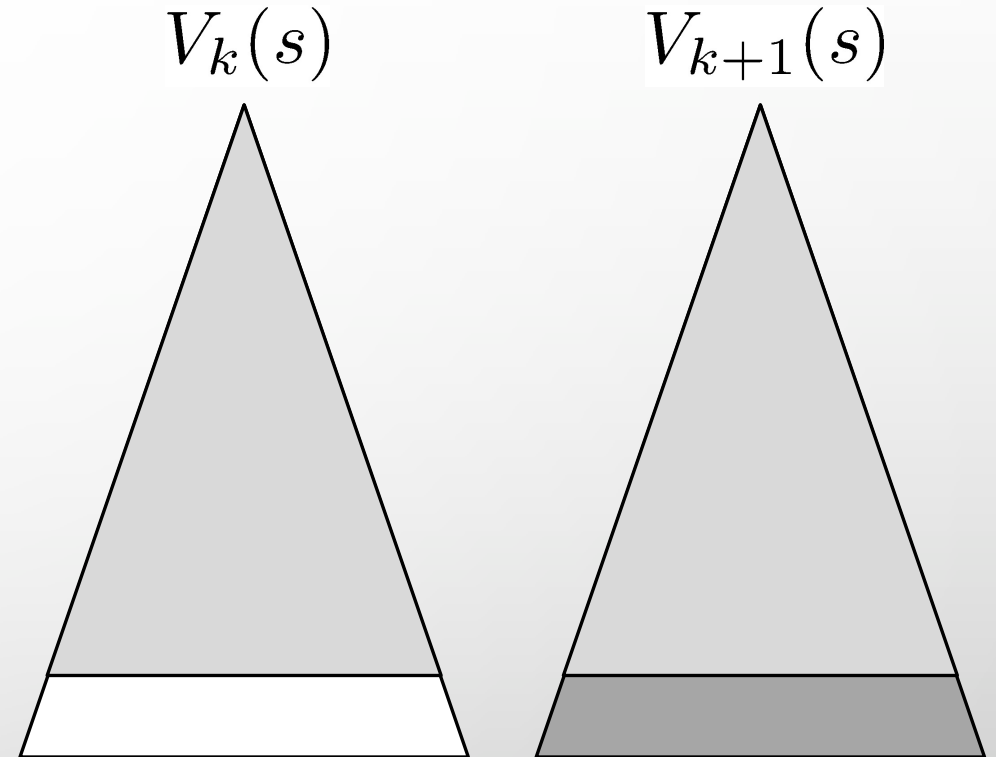- Value iteration <span style="color:red">computes</span> them:

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

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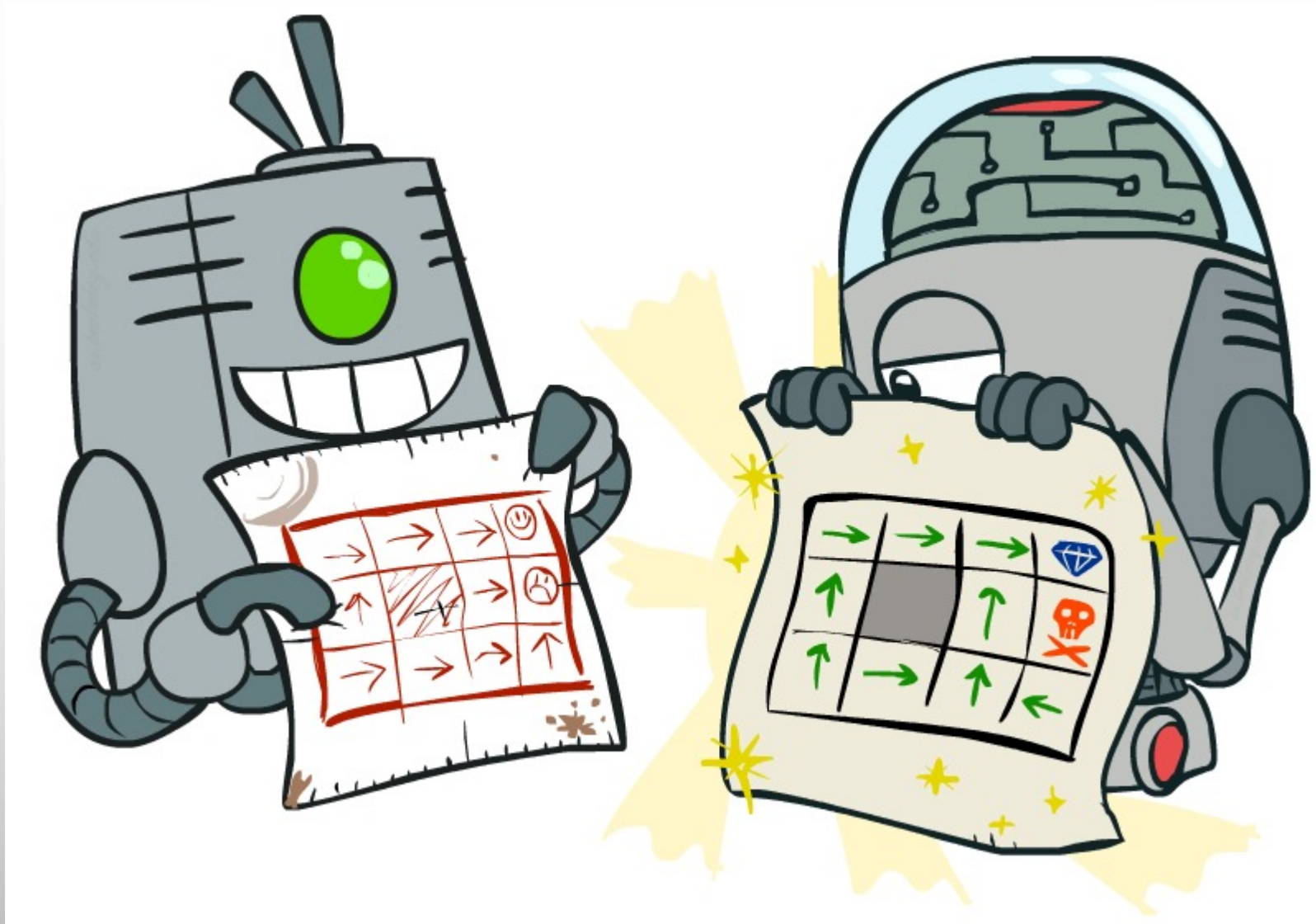
V(s)

a

s, a

s,a,s'

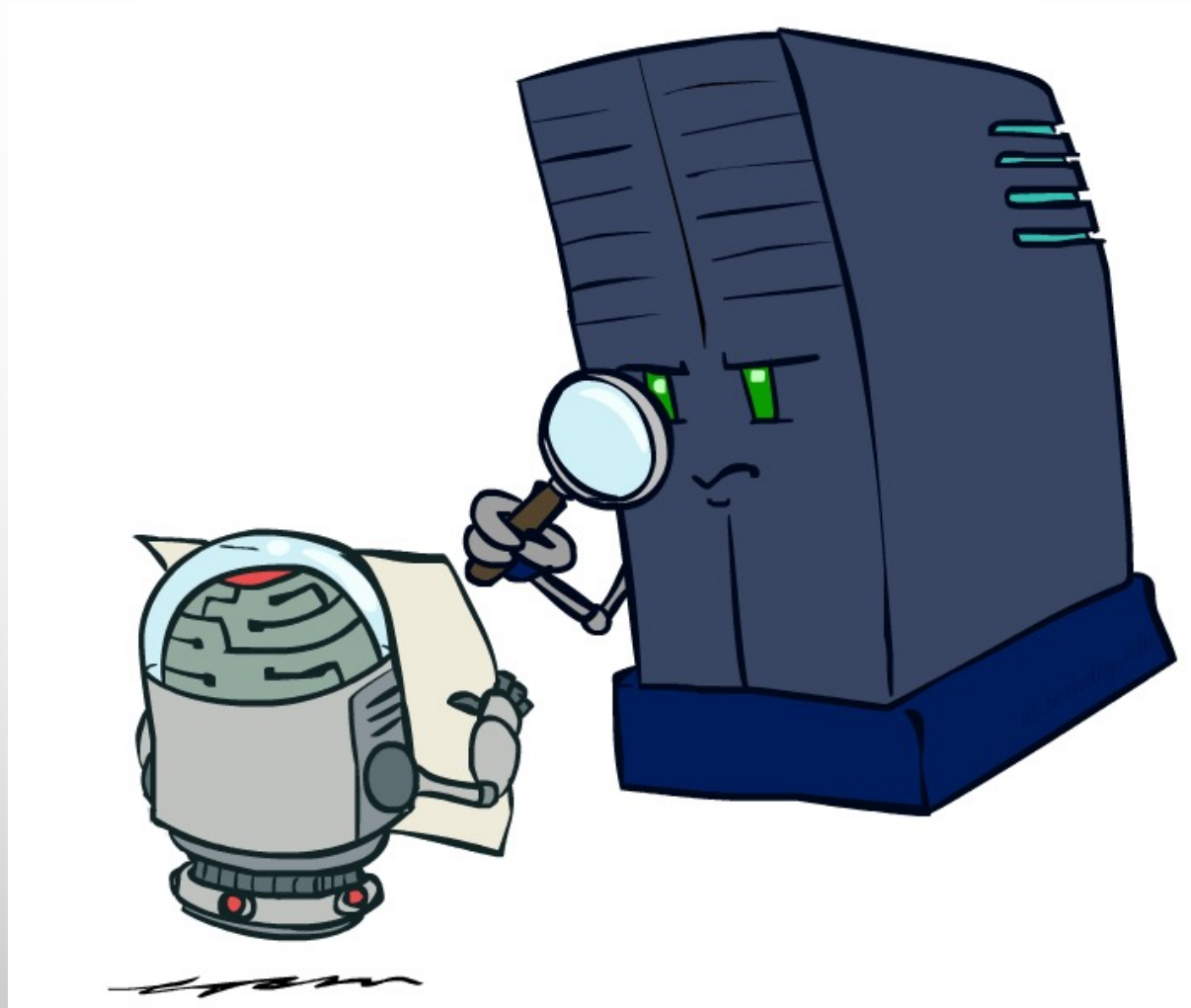V(s')

51

# Value Iteration Convergence

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

- Case 1: if the tree has maximum depth M, then $V_M$ holds the actual untruncated values

- Case 2: 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

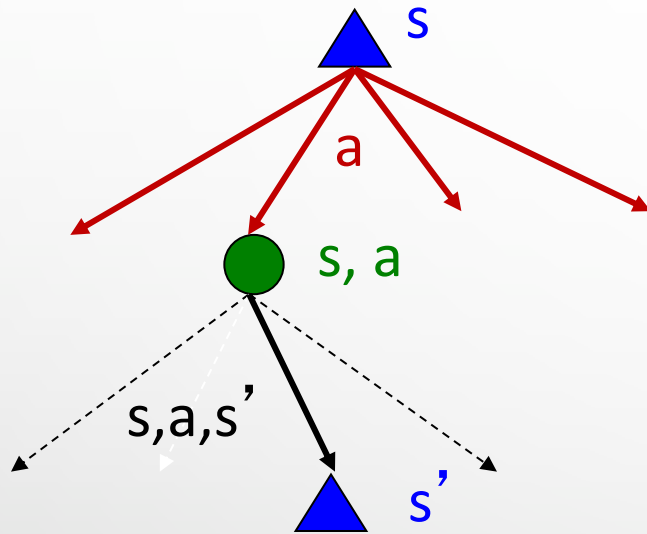$$V_k(s) \qquad V_{k+1}(s)$$

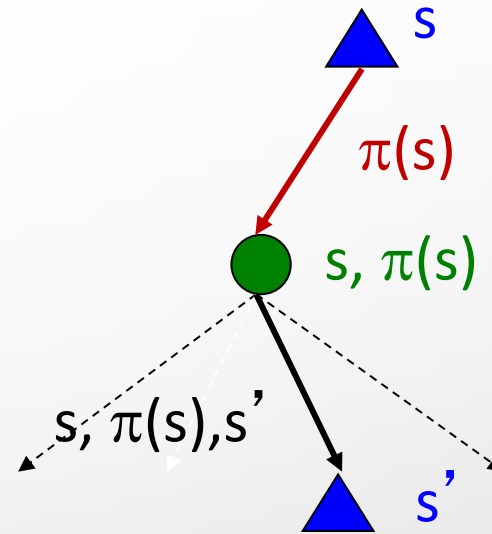# Policy Methods

# Policy Evaluation

# Fixed Policies

Do the optimal action

Do what $\pi$ says to do



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

- If we fixed some policy $\pi(s)$, then the tree would be simpler – only one action per state

  - … 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 $\pi$:

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

- 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

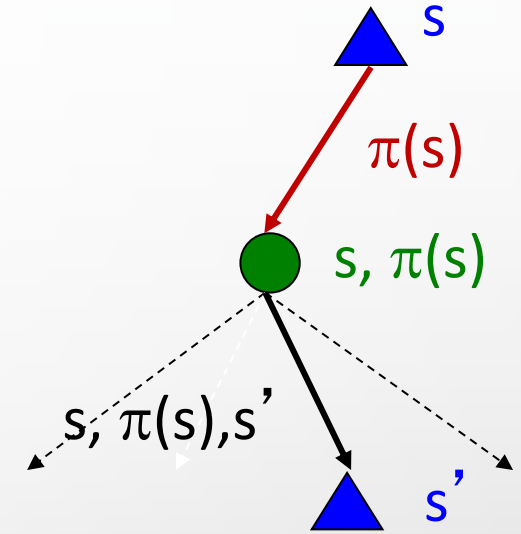$\pi$(s)

s, $\pi$(s)

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

s'

# Policy Evaluation

- How do we calculate the V's for a fixed policy $\pi$?

- Idea 1: turn recursive bellman equations into updates (Like value iteration)

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

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

- Efficiency: $O(S^2)$ per iteration

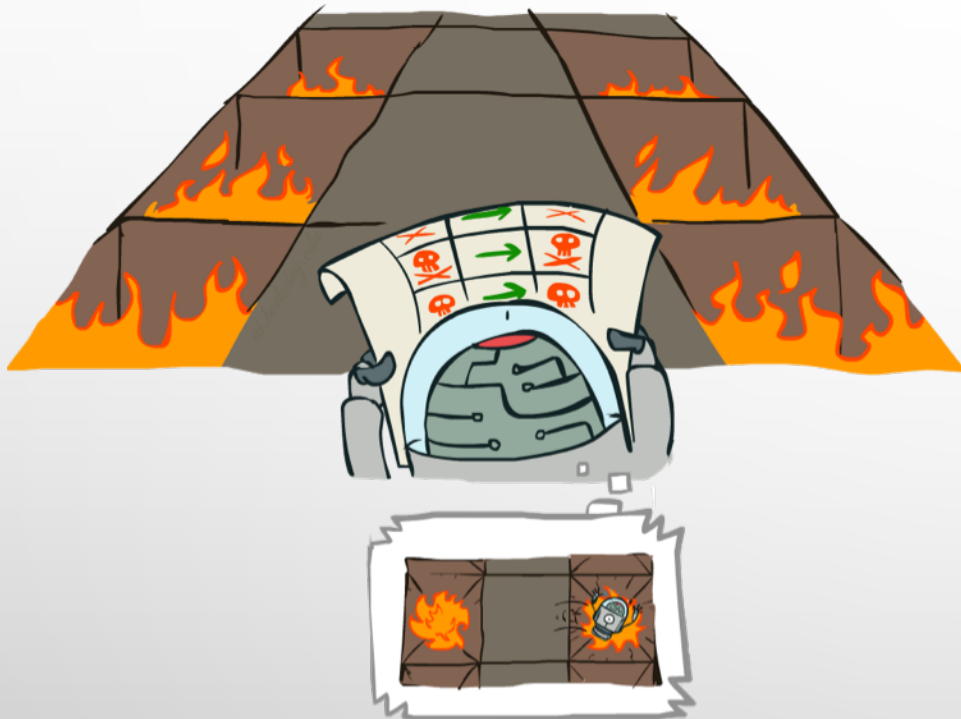- Idea 2: without the maxes, the bellman equations are just a linear system
  - Solve with your favorite linear system solver
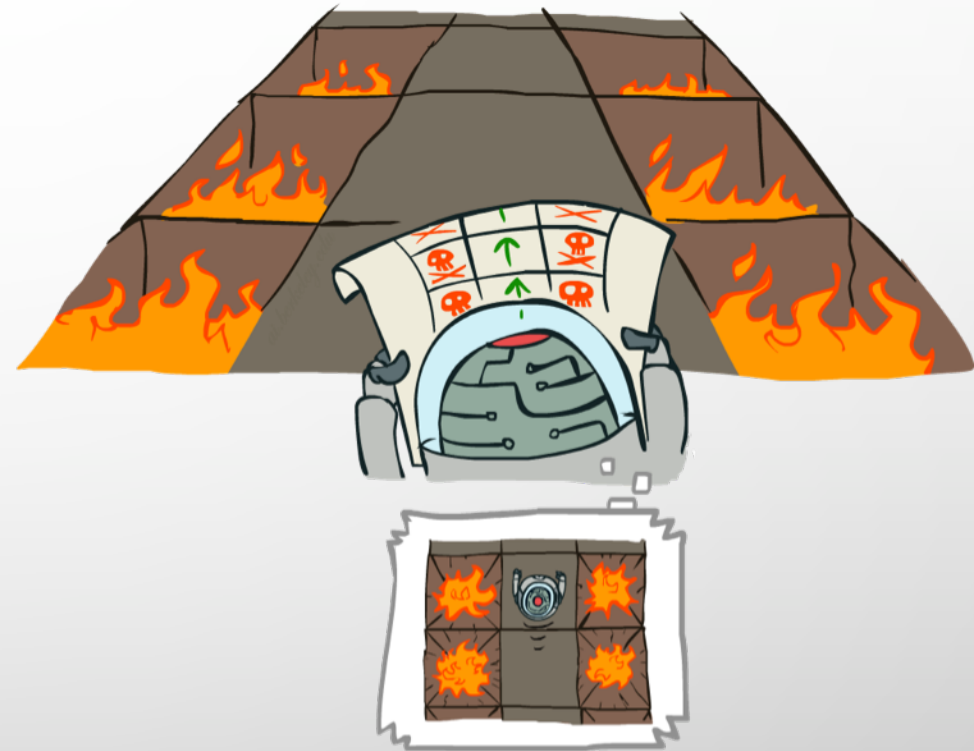
s

$\pi(s)$

s, $\pi(s)$

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

s'

$$\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \times \begin{bmatrix} V^\pi(s_1) \\ V^\pi(s_2) \\ \dots \end{bmatrix} = \begin{bmatrix} \\ \\ \end{bmatrix}$$

# Example: Policy Evaluation

Always Go Right

Always Go Forward

# Example: Policy Evaluation
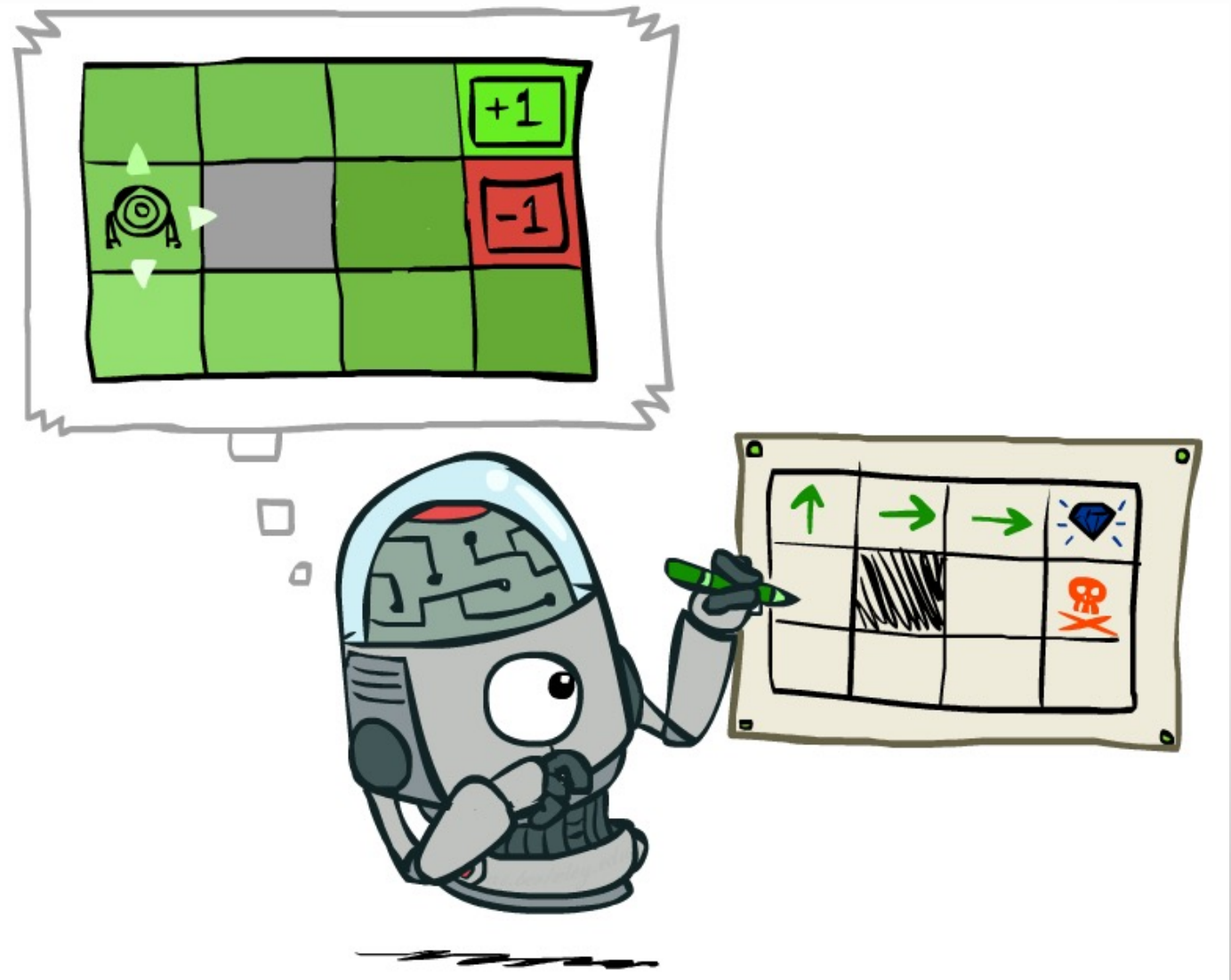
Always Go Right

Always Go Forward

# Policy Extraction

# Computing Actions from Values

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

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

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

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

- This is called policy extraction, since it gets the policy implied by the values

# Computing Actions from Q-Values



- Let's imagine we have the optimal Q-values:

- How should we act?
  - Completely trivial to decide!

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

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

# Policy Iteration

# Problems with Value Iteration

- Value iteration repeats the Bellman updates:

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

- Problem 1: it's slow – O(S²A) per iteration

- Problem 2: the "max" at each state rarely changes

- Problem 3: the policy often converges long before the values

s

a

s, a

s,a,s'

s'

# k=0



VALUES AFTER 0 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

65

# k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

69

# k=5



VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



VALUES AFTER 7 ITERATIONS

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



VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



VALUES AFTER 10 ITERATIONS
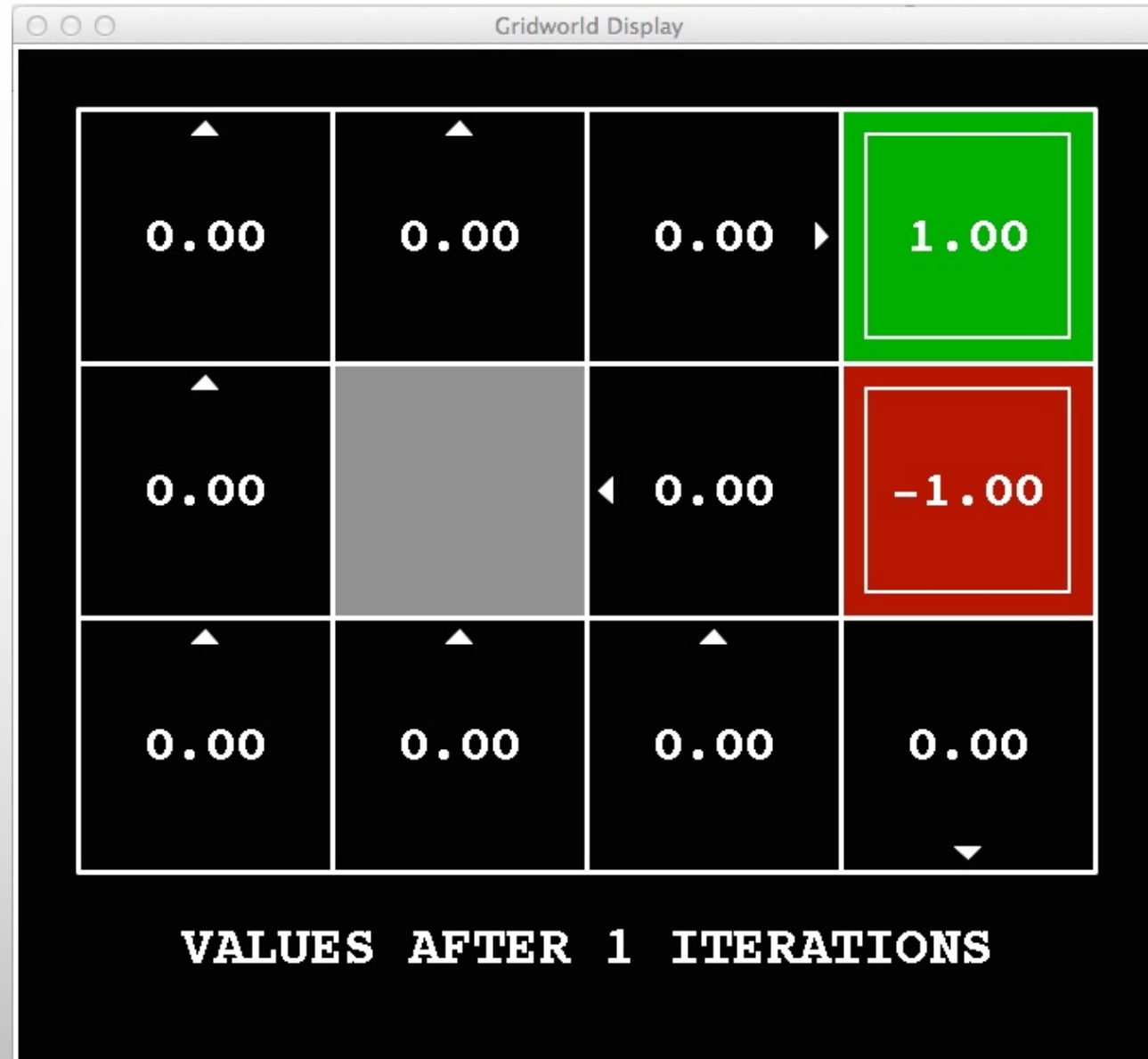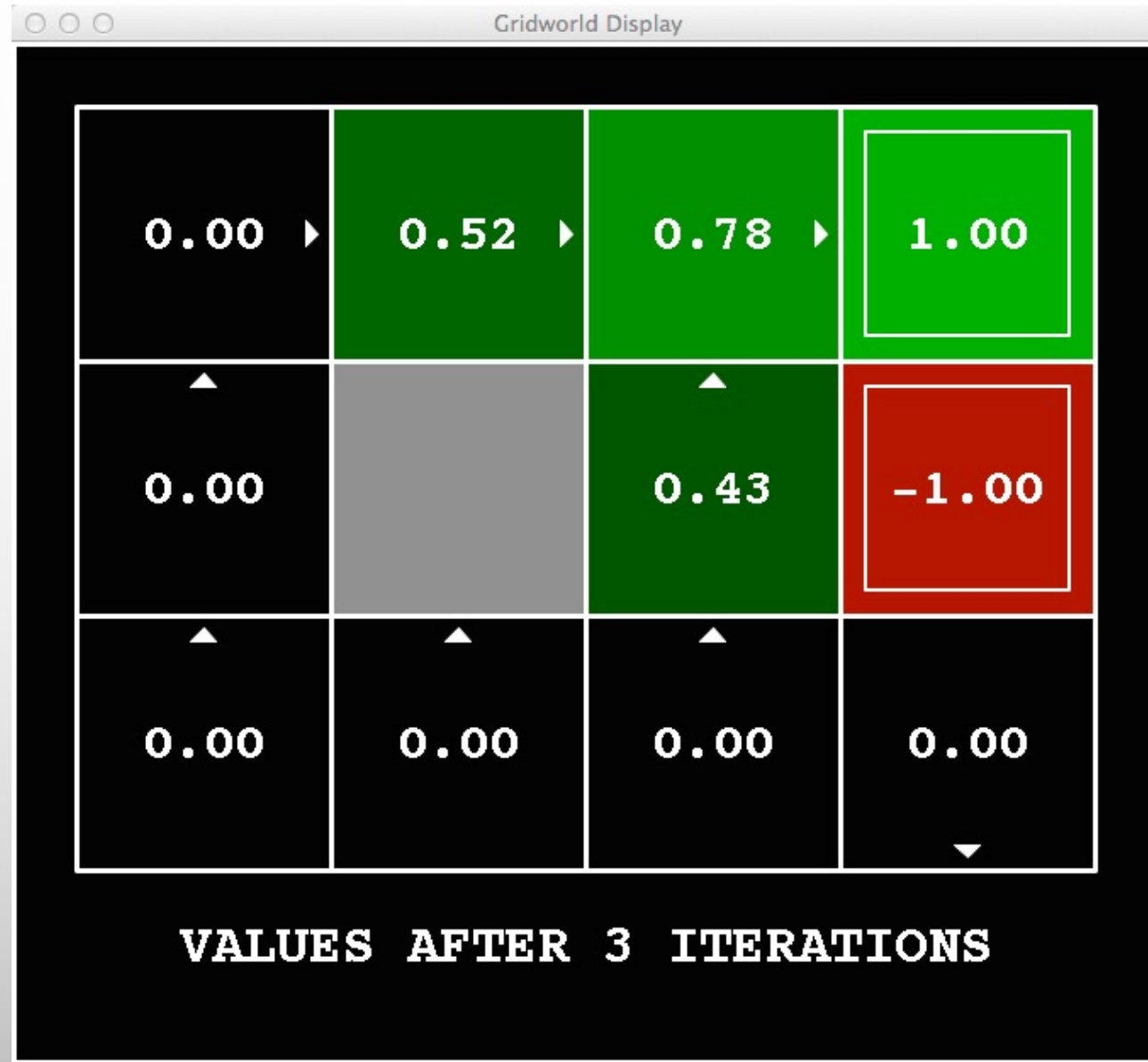
Noise = 0.2
Discount = 0.9
Living reward = 0

# 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

# Policy Iteration

- Alternative approach for optimal values:

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

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

  - Repeat steps until policy converges


- This is policy iteration

  - It's still optimal!

  - Can converge (much) faster under some conditions

# Policy Iteration

- Evaluation: for fixed current policy $\pi$, find values with policy evaluation:

  - Iterate until values converge:

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

- Improvement: for fixed values, get a better policy using policy extraction

  - One-step look-ahead:

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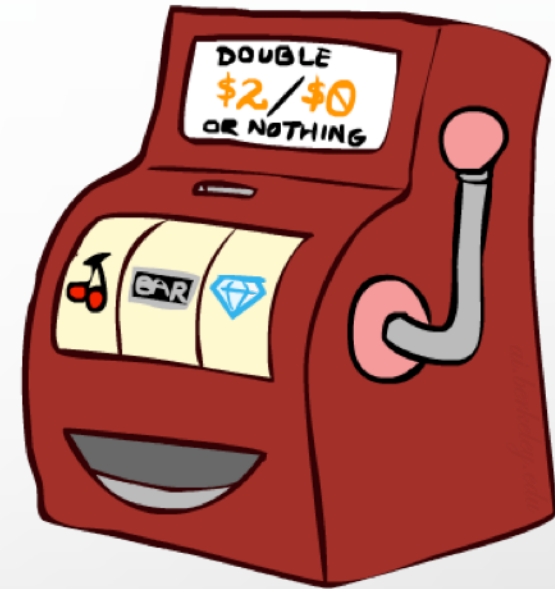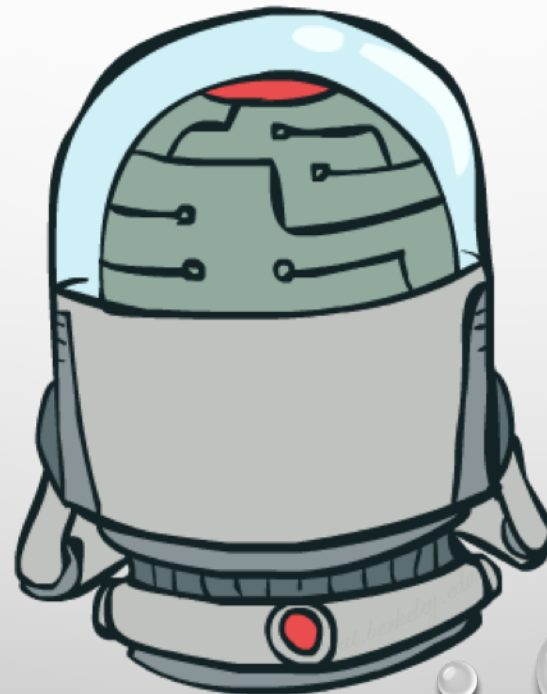
# Comparison

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

- In value iteration:

  - Every iteration updates both the values and (implicitly) the policy

  - We don't track the policy, but taking the max over actions implicitly recomputes it

- In policy iteration:

  - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)

  - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)

  - The new policy will be better (or we're done)

- Both are dynamic programs for solving MDPs
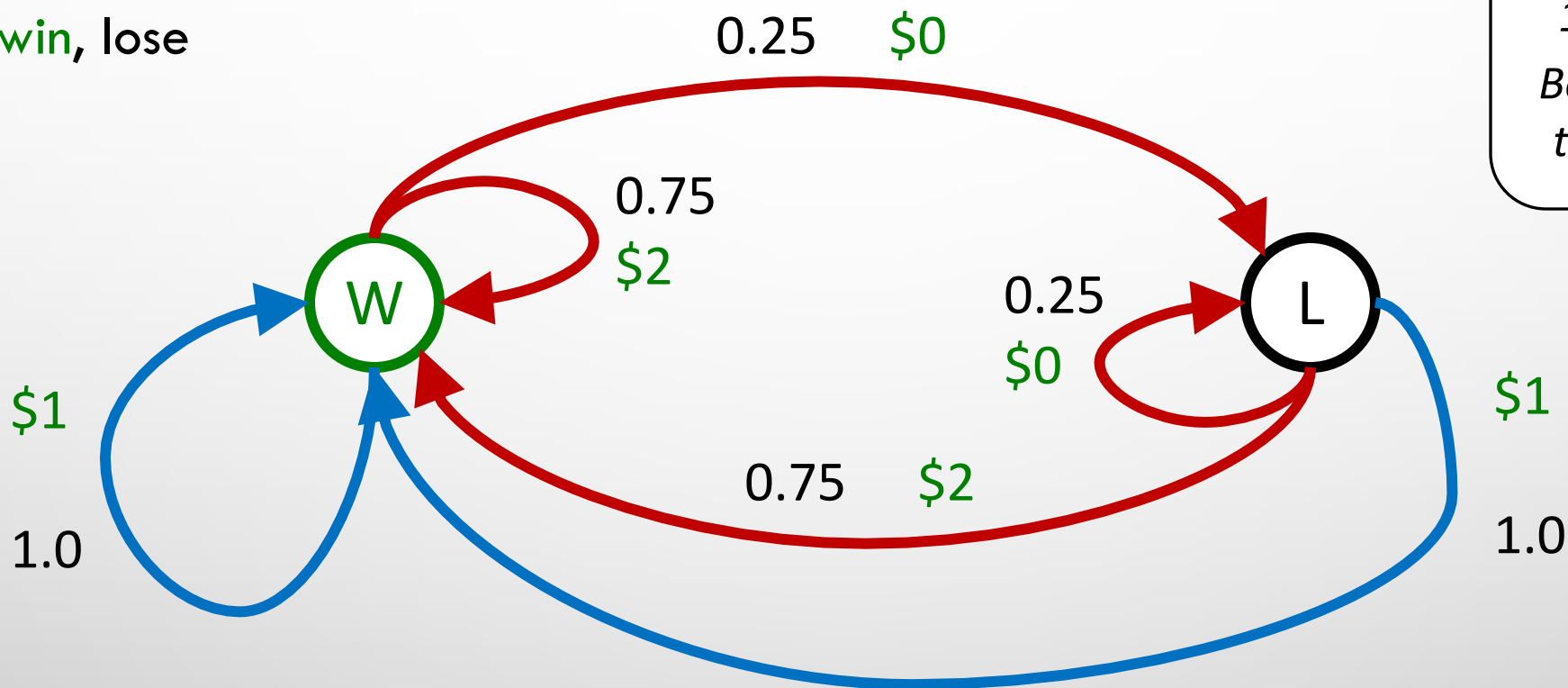
# Summary: MDP Algorithms

- So you want to….
  - Compute optimal values: use value iteration or policy iteration
  - Compute values for a particular policy: use policy evaluation
  - Turn your values into a policy: use policy extraction (one-step lookahead)

- These all look the same!
  - They basically are – they are all variations of bellman updates
  - They all use one-step lookahead expectimax fragments
  - They differ only in whether we plug in a fixed policy or max over actions

# Double Bandits

# Double-Bandit MDP

- Actions: *blue*, *red*

- States: win, lose

No discount
100 time steps
Both states have the same value
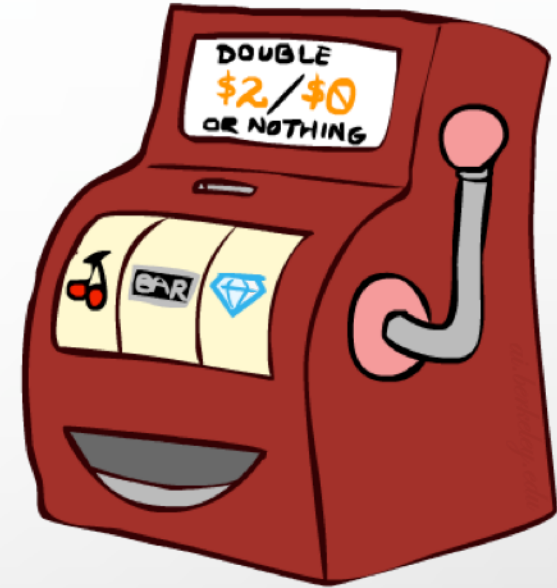


84

# Offline Planning

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

*No discount*

*100 time steps*

*Both states have the same value*

| | Value |
|---|---|
| Play Red | 150 |
| Play Blue | 100 |



0.25  $0

0.75
$2

0.25
$0

$1

$1

0.75  $2

1.0

1.0

# Let's Play!



$2  $2  $0  $2  $2

$2  $2  $0  $0  $0

# Online Planning

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

# Let's Play!

$0  $0  $0  $2  $0

$2  $0  $0  $0  $0

# What Just Happened?

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

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

89

# Next Time: Reinforcement Learning!