

# Artificial Intelligence

## Markov Decision Processes



Instructor: Heni Ben Amor

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

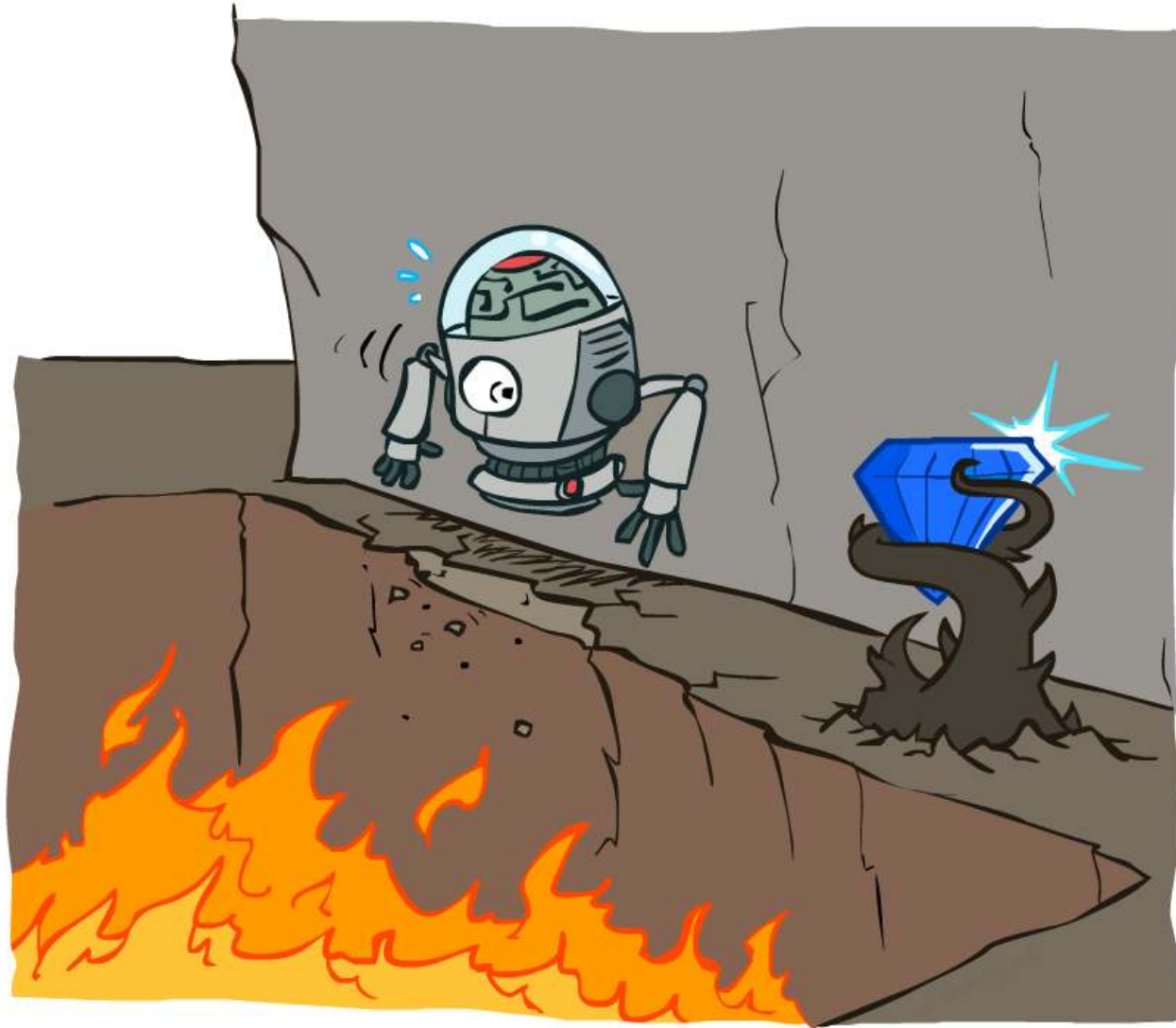
# What we Learned

---

- Trees solve discrete search problems
- When **ordering is not important we have CSP**
- **The challenge is computational complexity**
- Heuristics can help reduce it
- A\* combines **efficiency** and guarantees
- In **adversarial** domain use **alternation** of choices
- Expecti-max allows for **uncertain** events

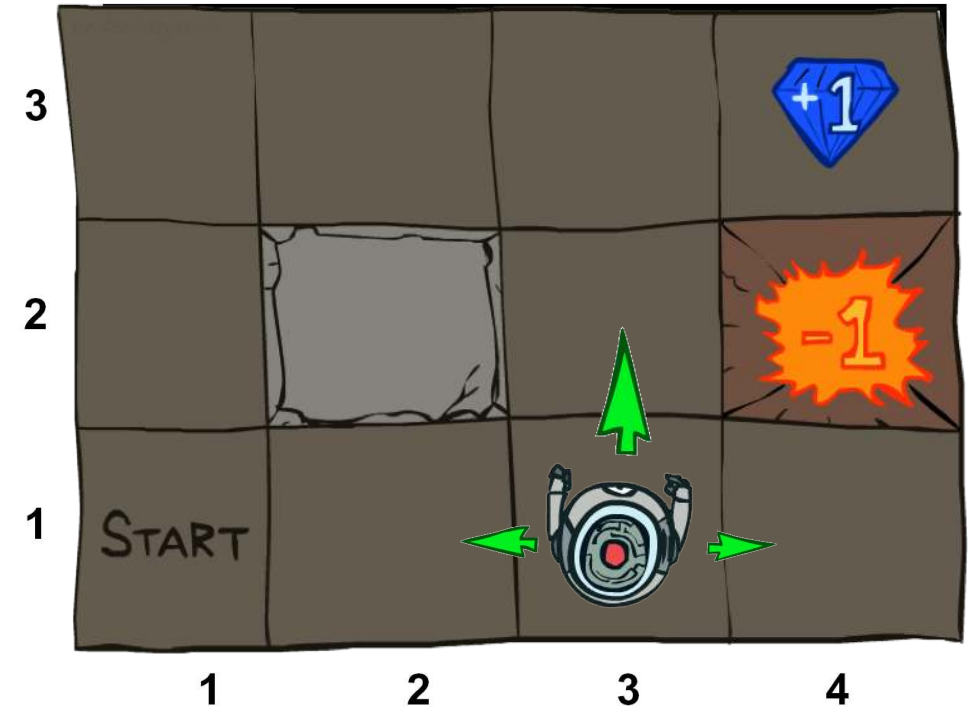
# Non-Deterministic Search

---



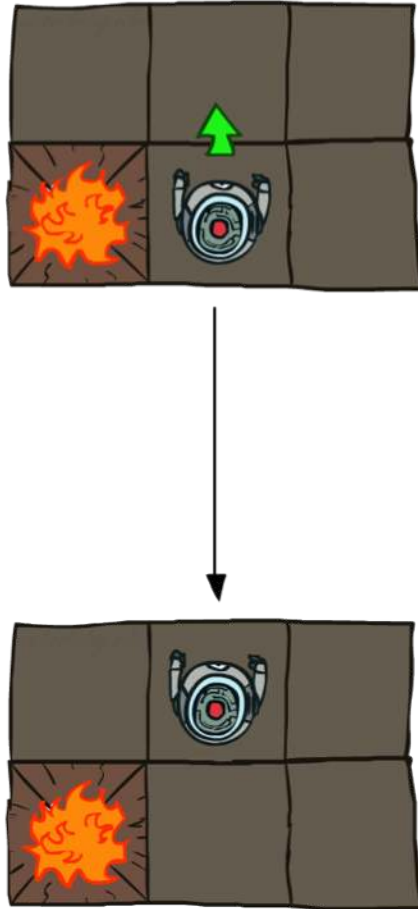
# Example: Grid World

- 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

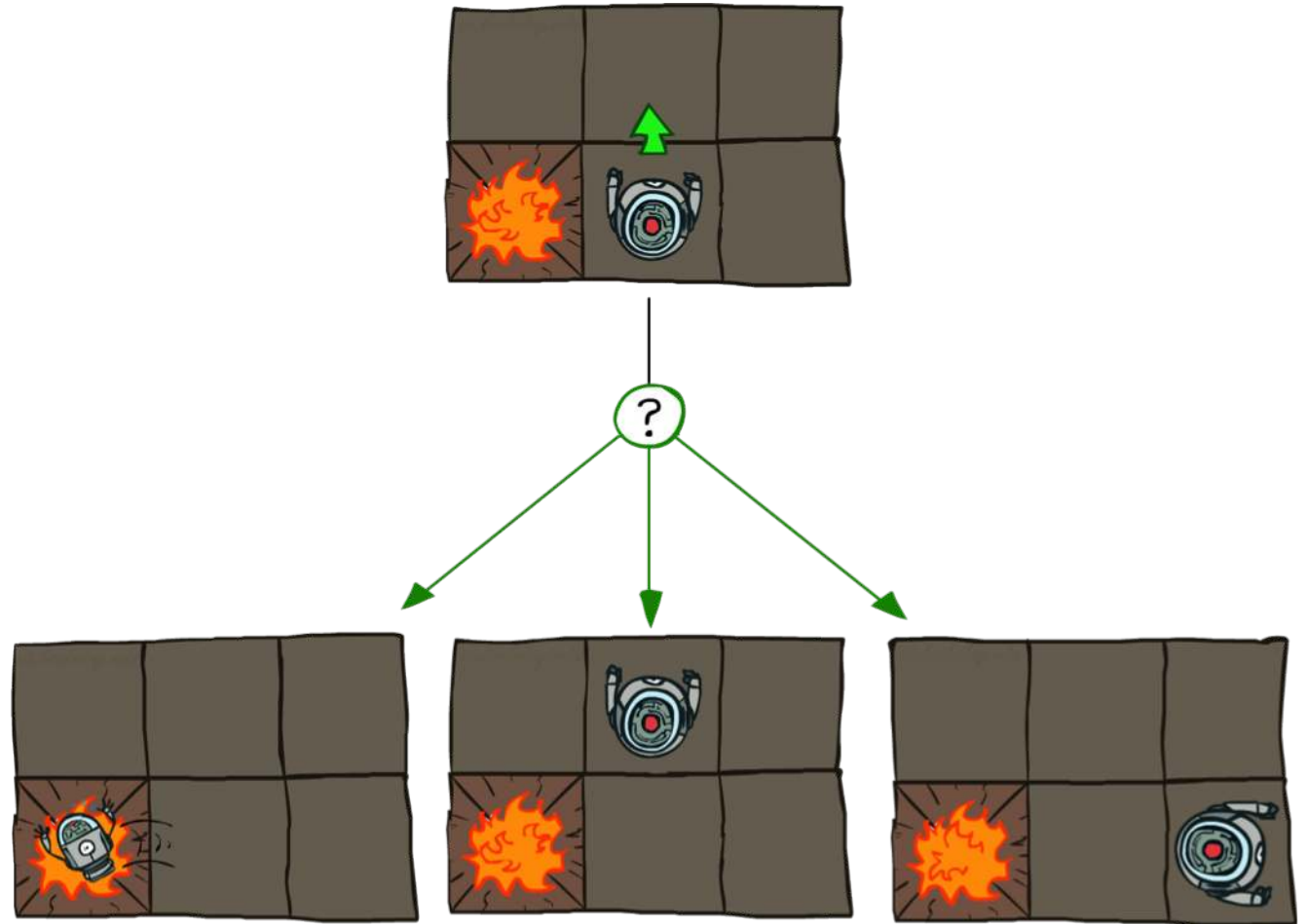


# Grid World Actions

Deterministic Grid World

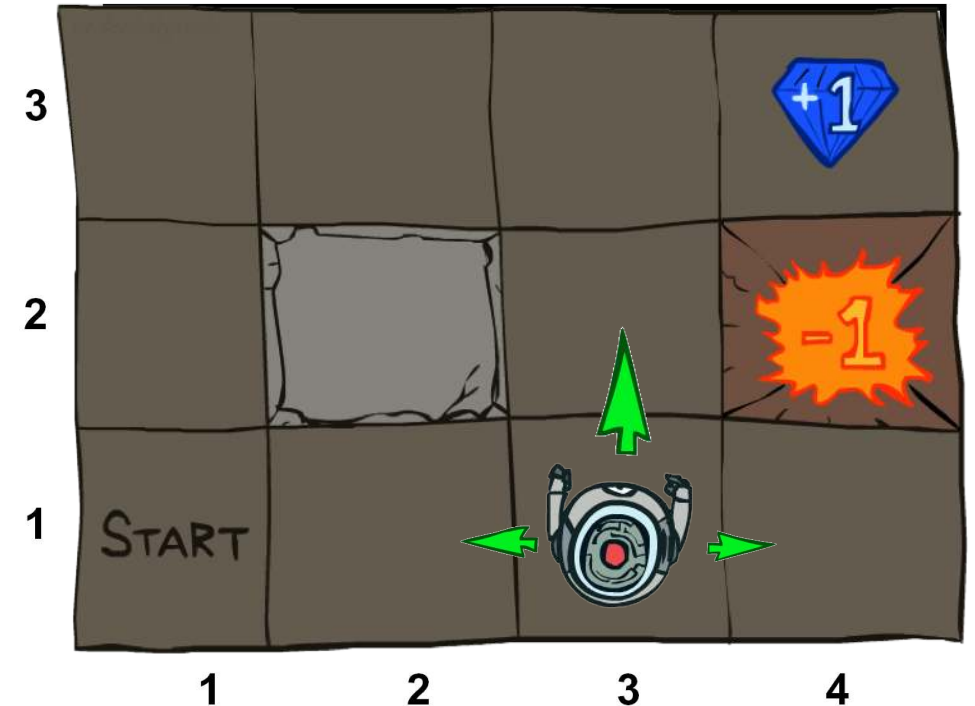


Stochastic Grid World



# Markov Decision Processes

- An MDP is defined by:
  - A set of states  $s \in S$
  - A set of actions  $a \in A$
  - A transition function  $T(s, a, s')$ 
    - Probability that  $a$  from  $s$  leads to  $s'$ , i.e.,  $P(s' | s, a)$   
Also called the model or the dynamics
  - A reward function  $R(s, a, s')$ 
    - Or 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

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

=

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

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

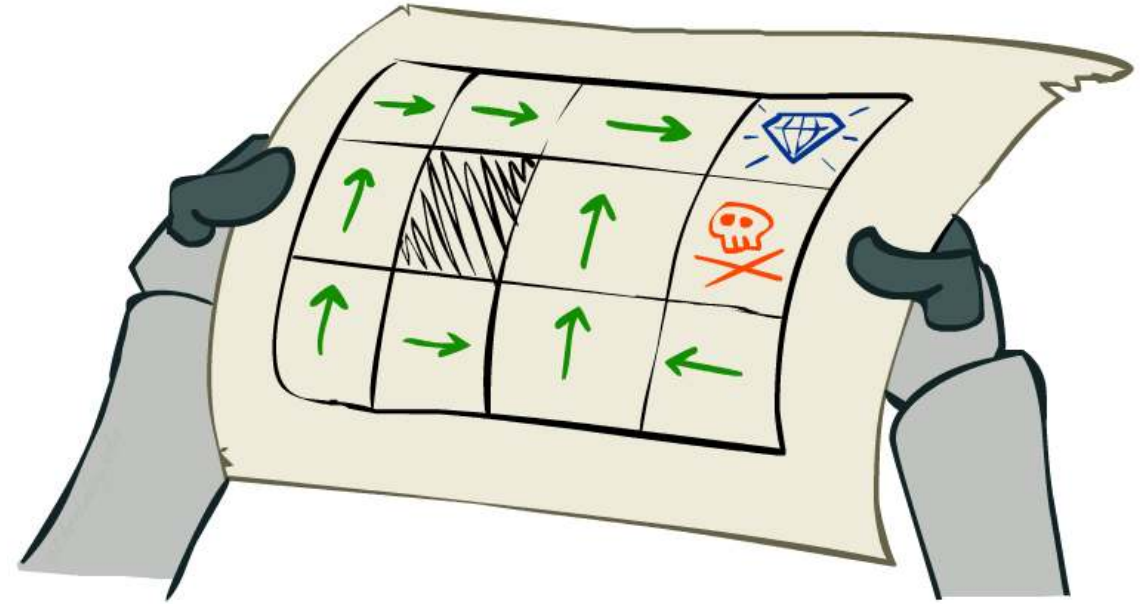


Andrey Markov  
(1856-1922)



# Policies

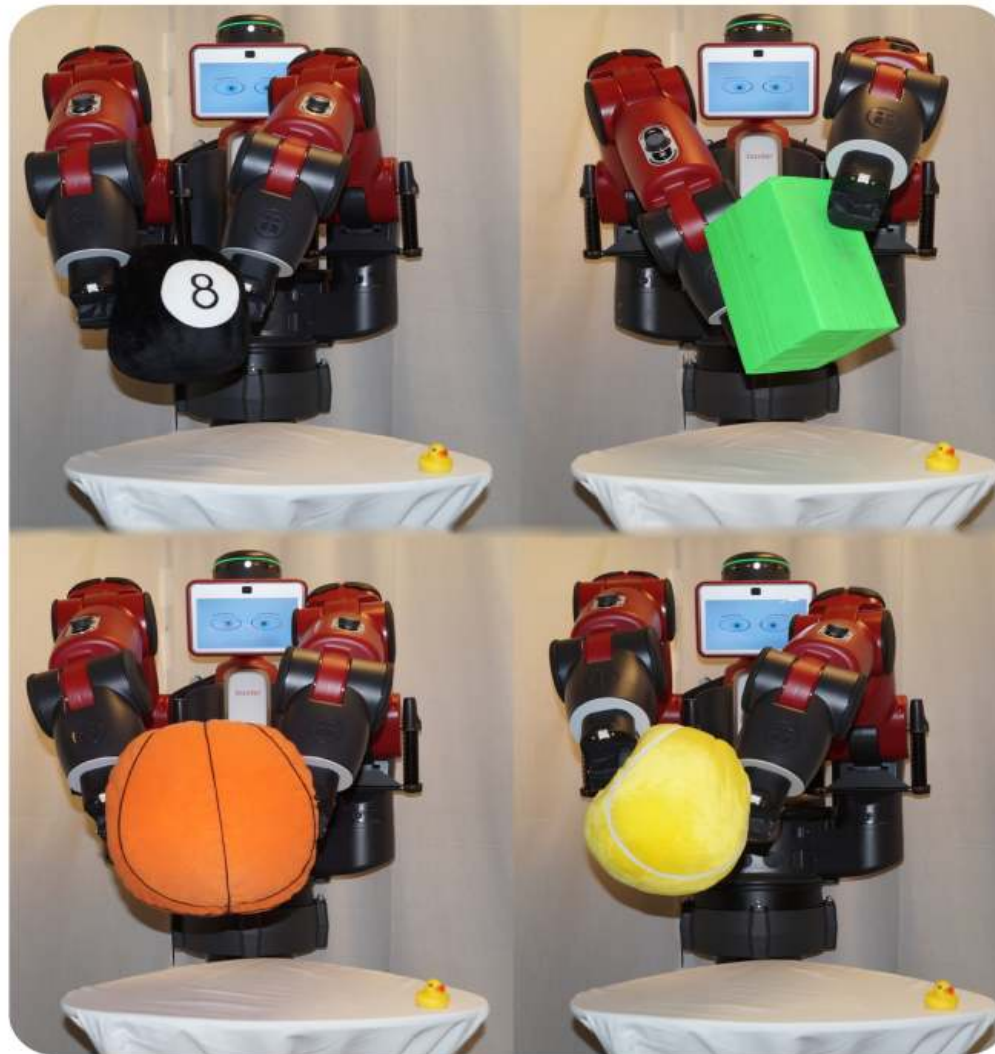
- In deterministic single-agent search problems, we wanted an optimal **plan**, or sequence of actions, from start to a goal
- For MDPs, we want an optimal **policy**  
 $\pi^*: S \rightarrow A$ 
  - A policy  $\pi$  gives an action for each state
  - An optimal policy is one that maximizes expected utility if followed
- Expectimax didn't compute entire policies
  - It computed the action for a single state only
  - MDPs can have **infinite length**!



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



# Example: Learning to Lift



# Example: Learning to Lift

## **Extracting Bimanual Synergies with Reinforcement Learning**

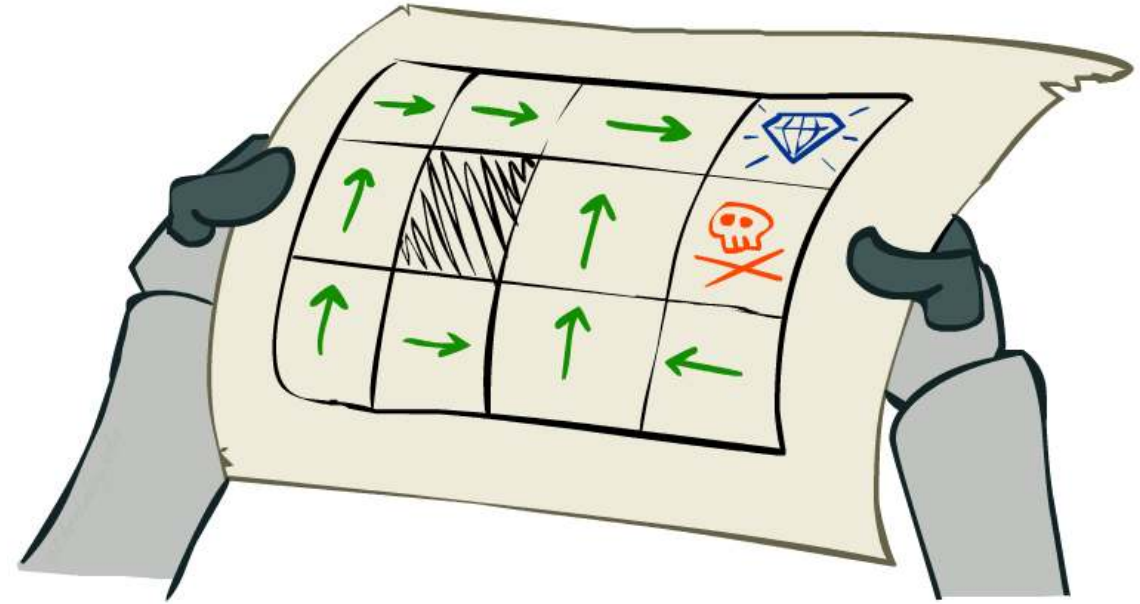
Kevin S. Luck and Heni Ben Amor



To be presented in 2 weeks at **IROS 2017** in Vancouver.

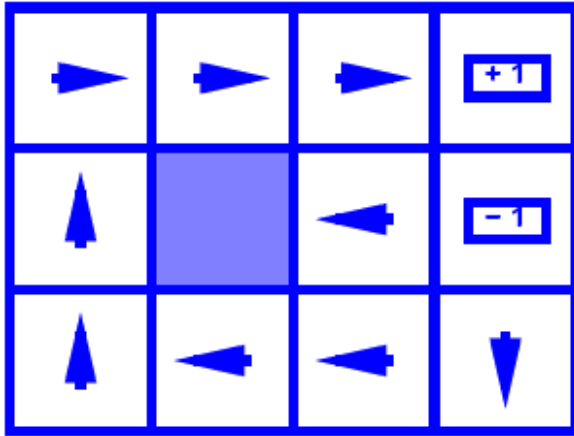
# Policies

- 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
  - MDPs can have **infinite length**!

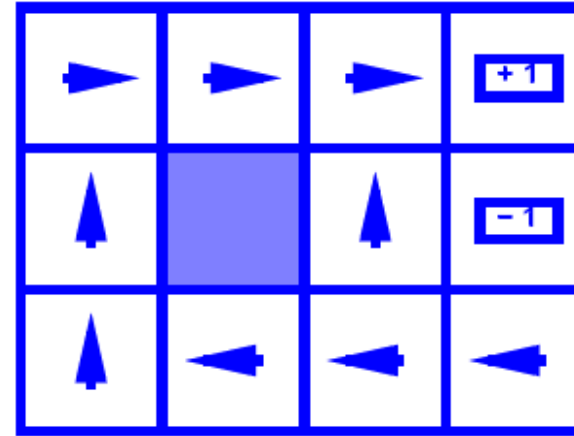


Optimal policy when  $R(s, a, s') = -0.03$  for all non-terminals  $s$

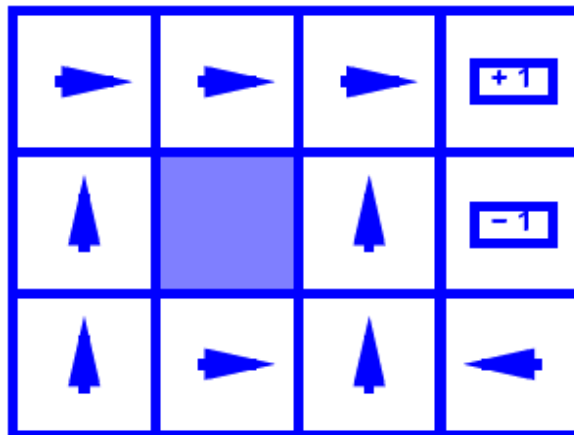
# Optimal Policies



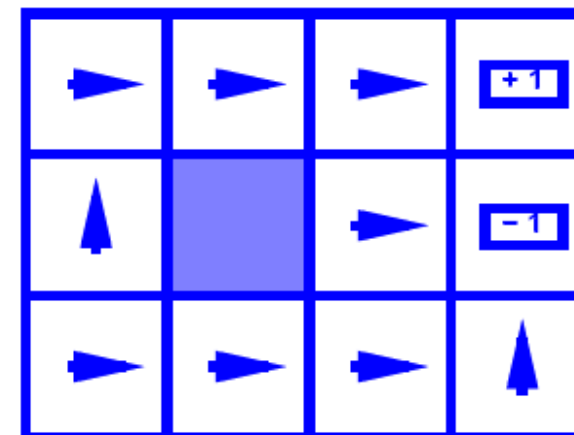
$R(s) =$   
-0.01



$R(s) =$  -0.03



$R(s) =$   
-0.4



$R(s) =$   
-2.0

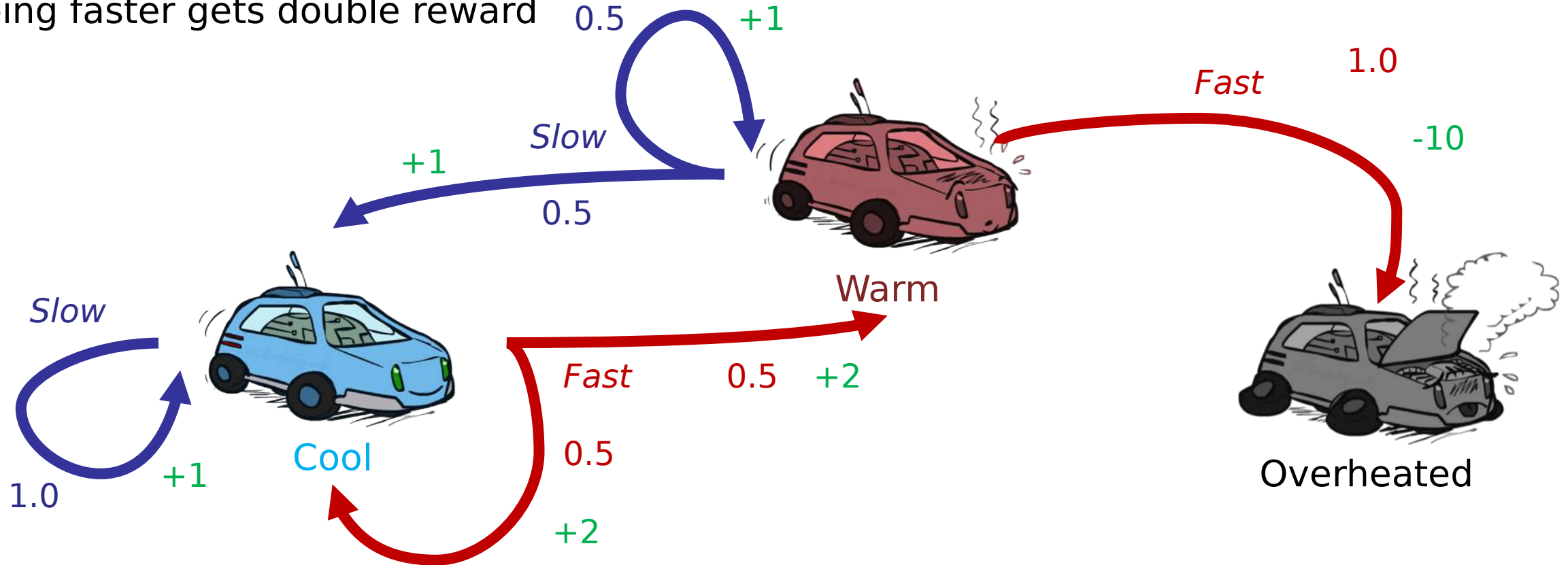
# Example: Racing

---

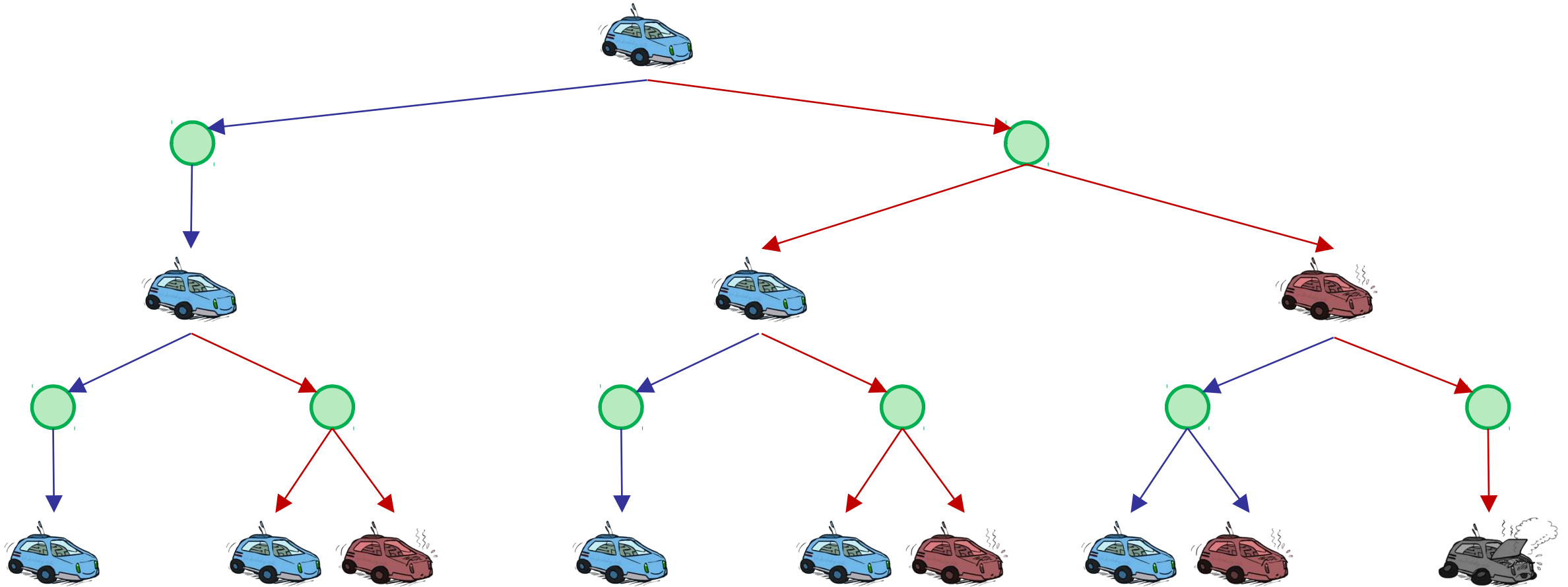


# Example: Racing

- A robot car wants to travel far, quickly
- Three states: **Cool**, **Warm**, Overheated
- Two actions: **Slow**, **Fast**
- Going faster gets double reward



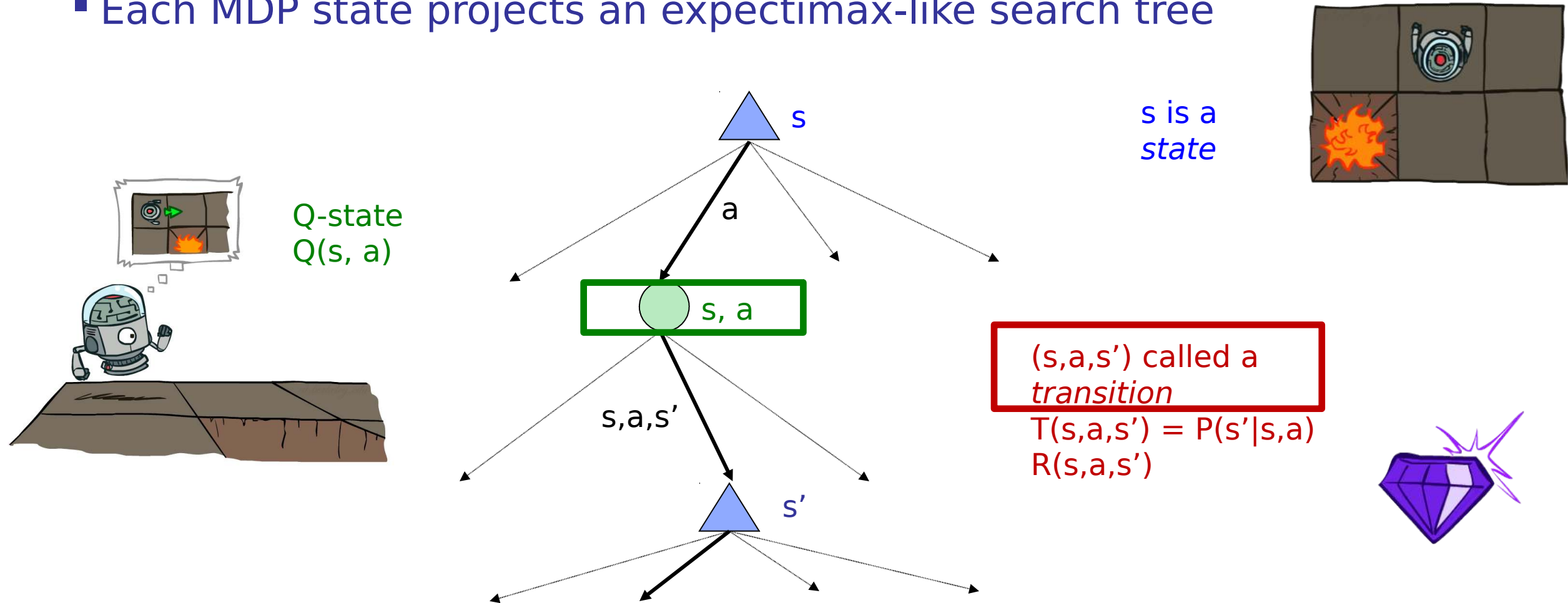
# Racing Search Tree



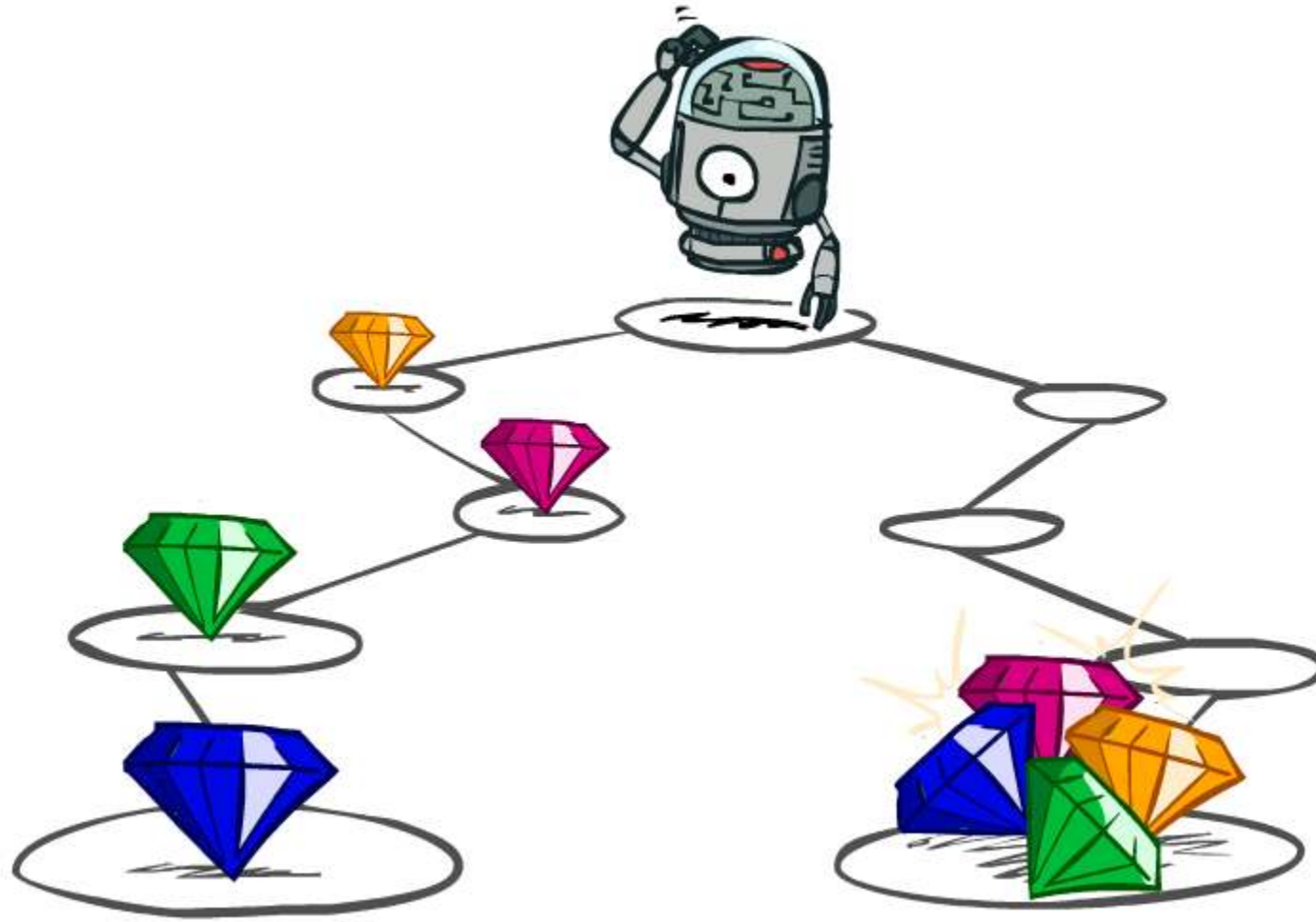


# MDP Search Trees

- Each MDP state projects an expectimax-like search tree

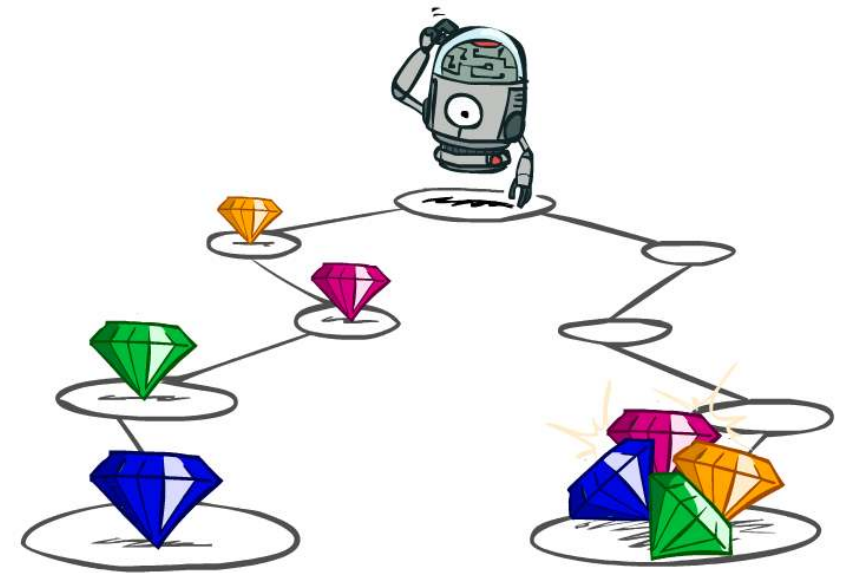


# Utilities of Sequences



# 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  $[1, 0, 0]$



# 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

Worth Now



$\gamma$

Worth Next  
Step



$\gamma^2$

Worth In Two  
Steps

# Discounting

- How to discount?

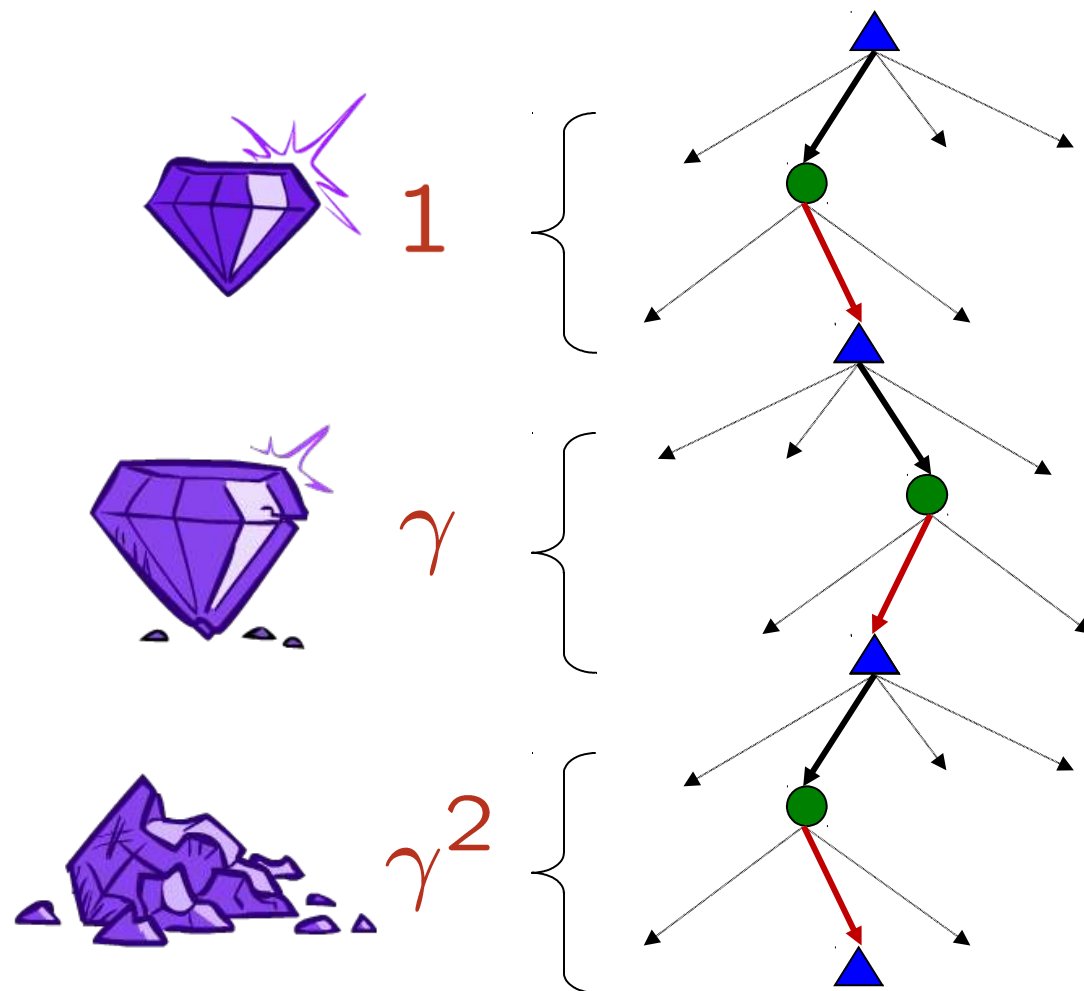
- Each time we descend a level, we multiply in the discount once

- Why discount?

- Sooner rewards probably do have higher utility than later rewards
- Also helps our algorithms converge

- Example: discount of 0.5

- $U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3$
- $U([1,2,3]) < U([3,2,1])$

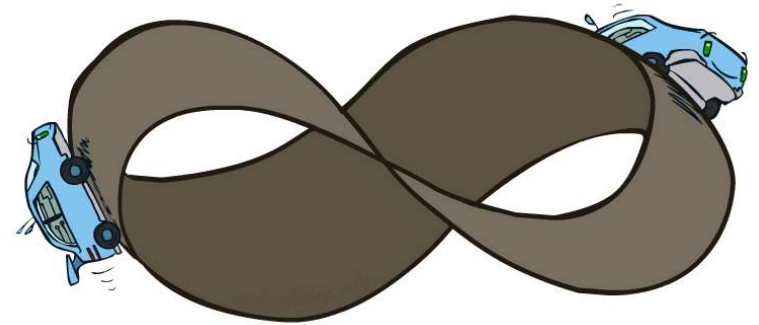


# Infinite Utilities?!

- 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  $0 < \gamma < 1$

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

- 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)

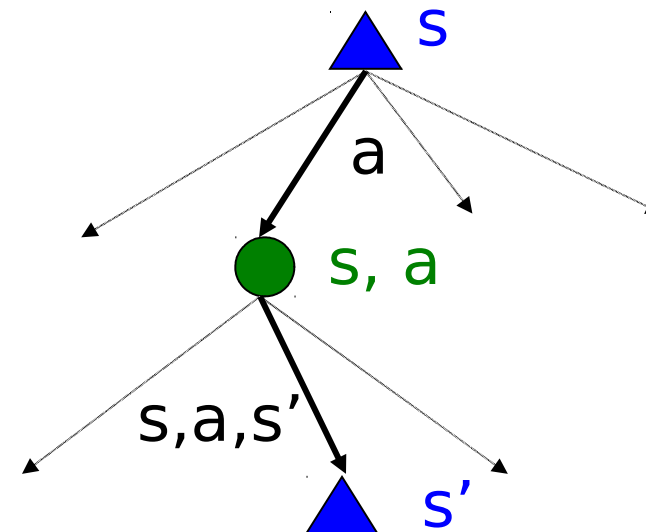
# Recap: Defining MDPs

- 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  $\gamma$ )

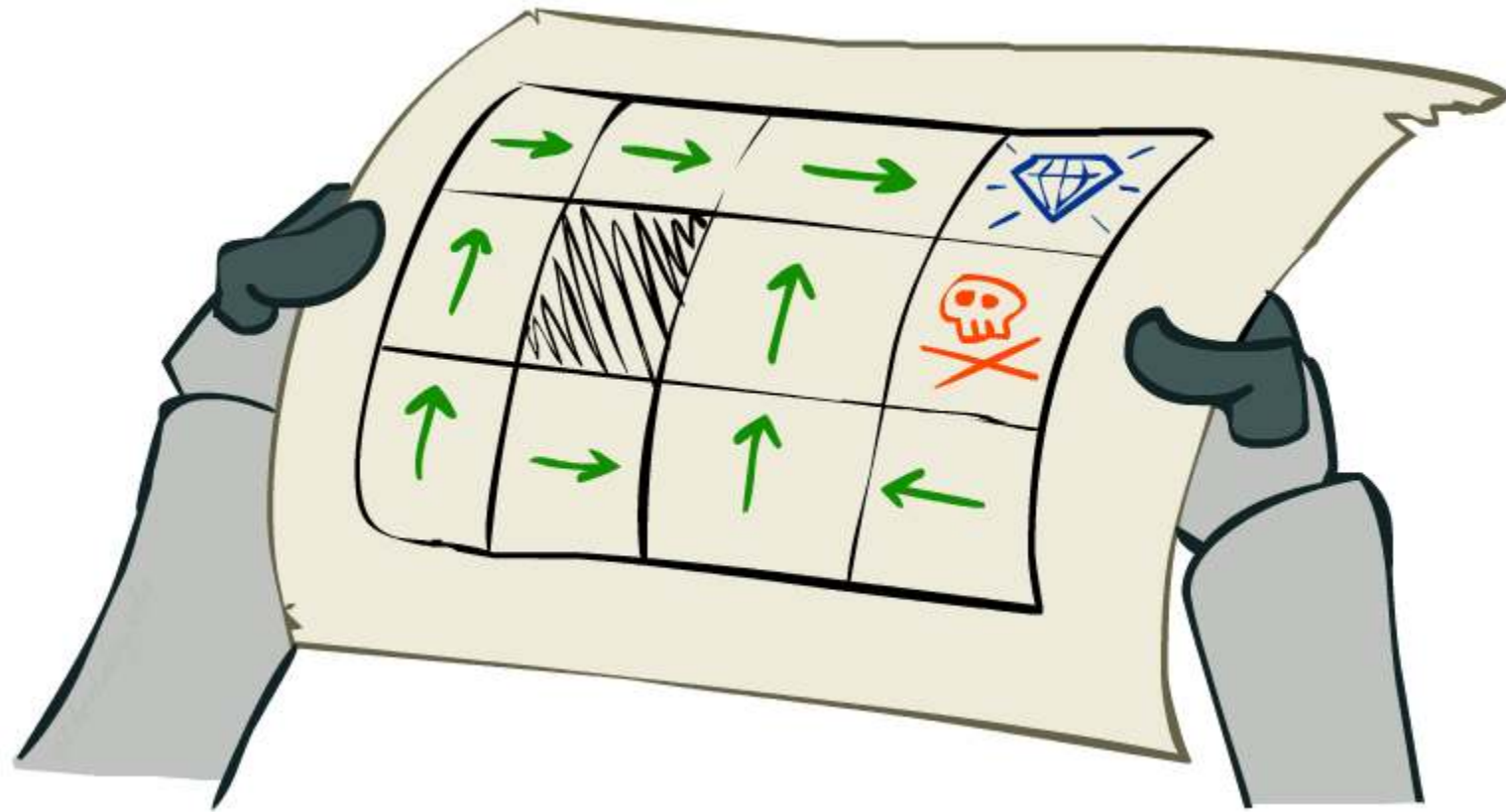
- MDP quantities so far:

- Policy = Choice of action for each state
- Utility = sum of (discounted) rewards



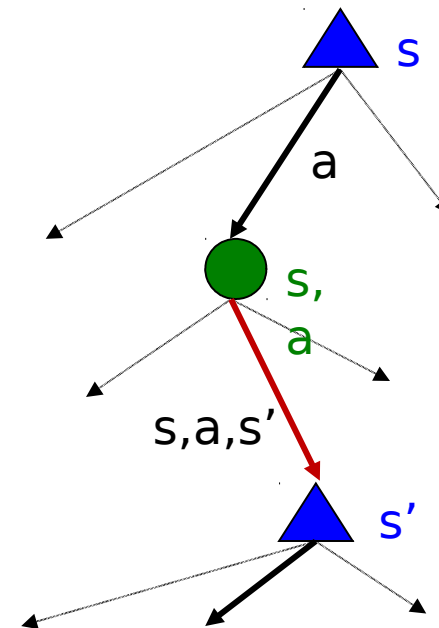


# Solving MDPs



# 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

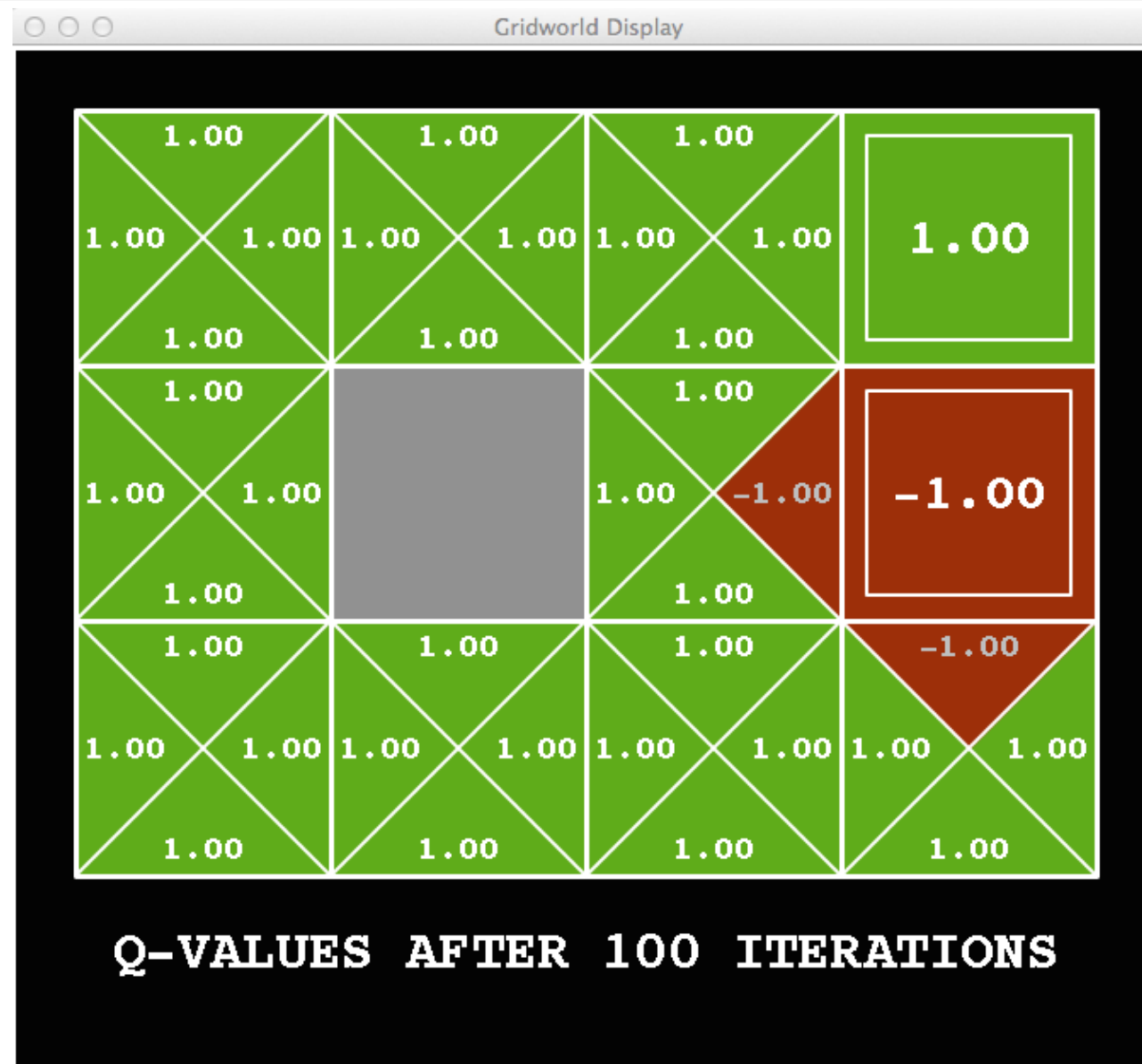
$(s, a, s')$  is a  
transition

# Snapshot of Demo – Gridworld V Values



Noise = 0  
Discount = 1  
Living reward = 0

# Snapshot of Demo – Gridworld Q Values



Noise = 0  
Discount = 1  
Living reward = 0

# Snapshot of Demo – Gridworld V Values



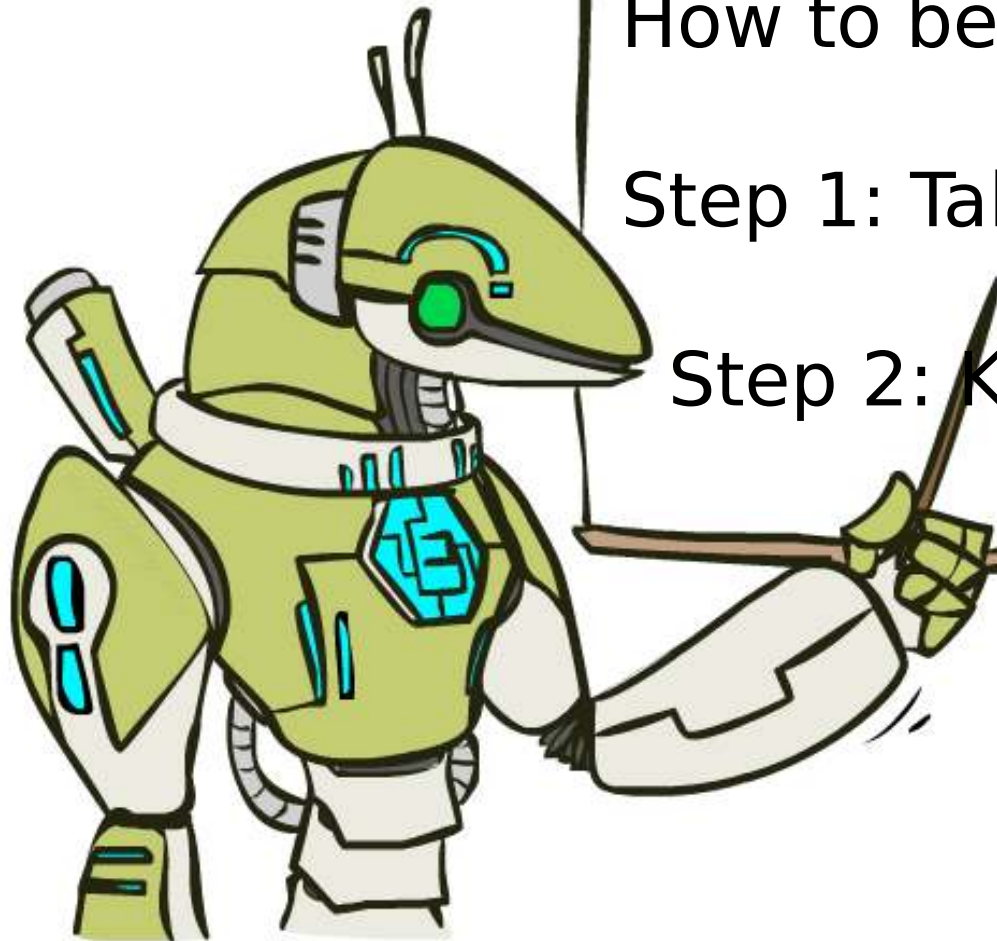
Noise = 0.2  
Discount = 0.9  
Living reward = 0

# Snapshot of Demo – Gridworld Q Values



Noise = 0.2  
Discount = 0.9  
Living reward = 0

# The Bellman Equations



How to be **optimal**:

Step 1: Take correct first action

Step 2: Keep being optimal



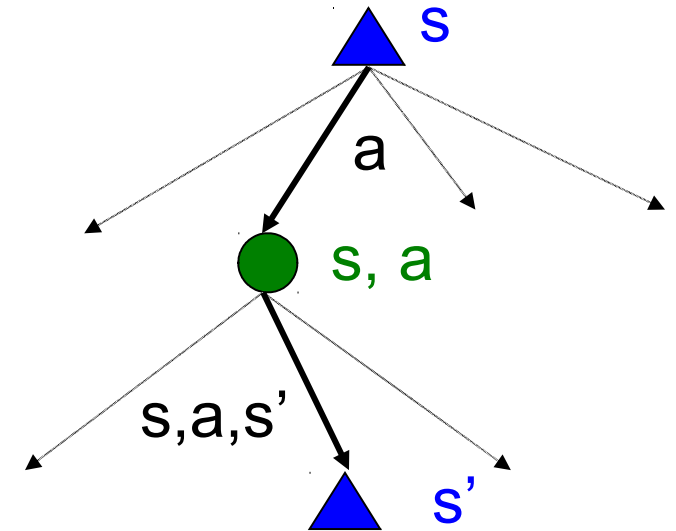
# Bellman Equations

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

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

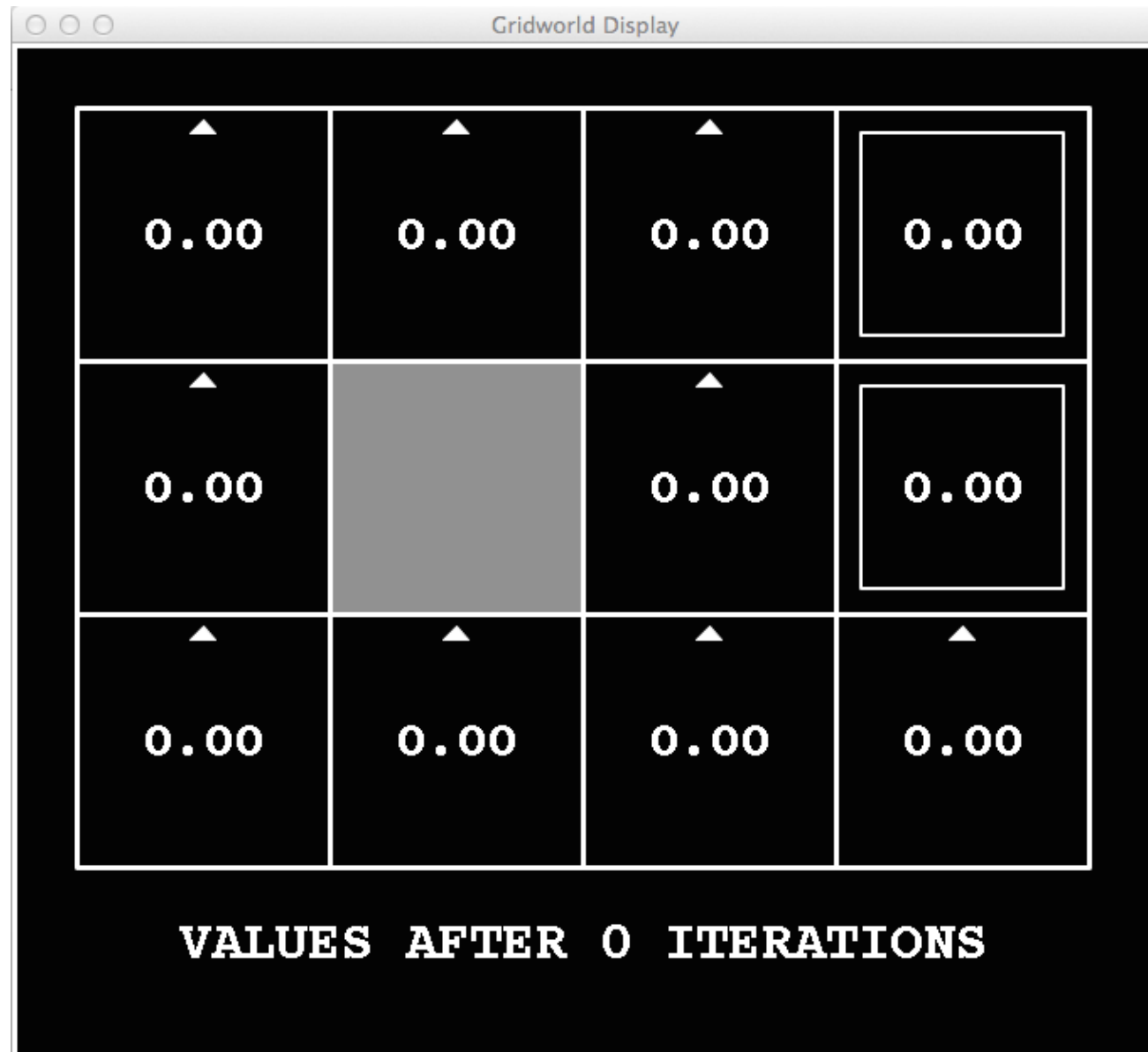
$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

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



- 
- An important step is calculating the **utilities** (Value) of state
  - The Bellman equations define these utilities
  - Missing: how to **efficiently compute** them
  - Iterative algorithm: Value Iteration

# $k=0$



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

# $k=4$



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



# Next Time: Value Iteration

---

