# CSE 573:
# Artificial Intelligence

## Hanna Hajishirzi
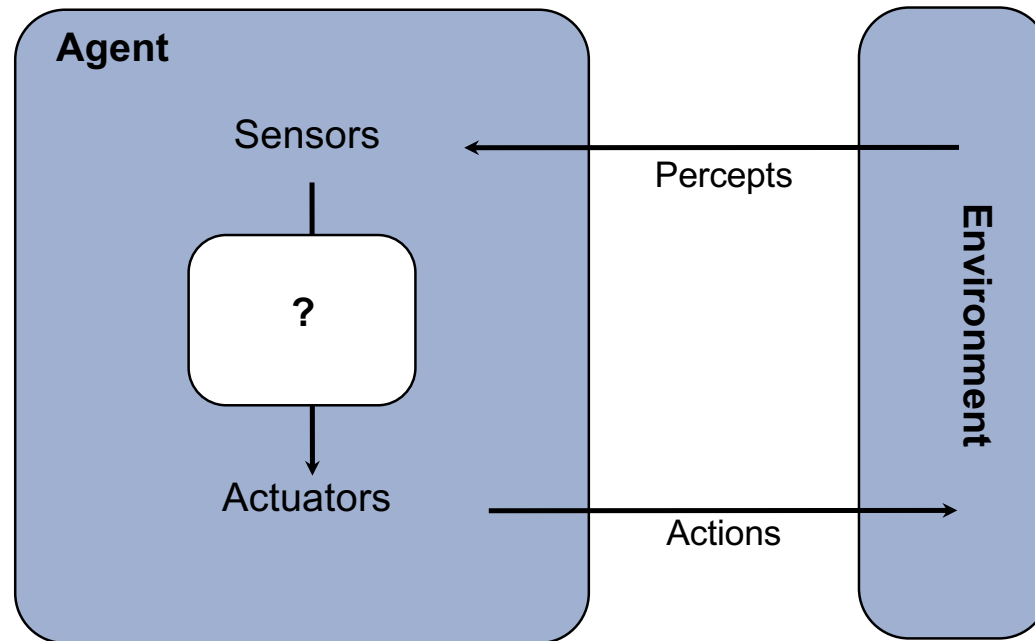
## Markov Decision Processes

# Review and Outline

- **Adversarial Games**
  - Minimax search
  - α-β search
  - Evaluation functions
  - Multi-player, non-0-sum
- **Stochastic Games**
  - Expectimax

  - Markov Decision Processes
  - Reinforcement Learning

# Agents vs. Environment

- An **agent** is an entity that *perceives* and *acts*.

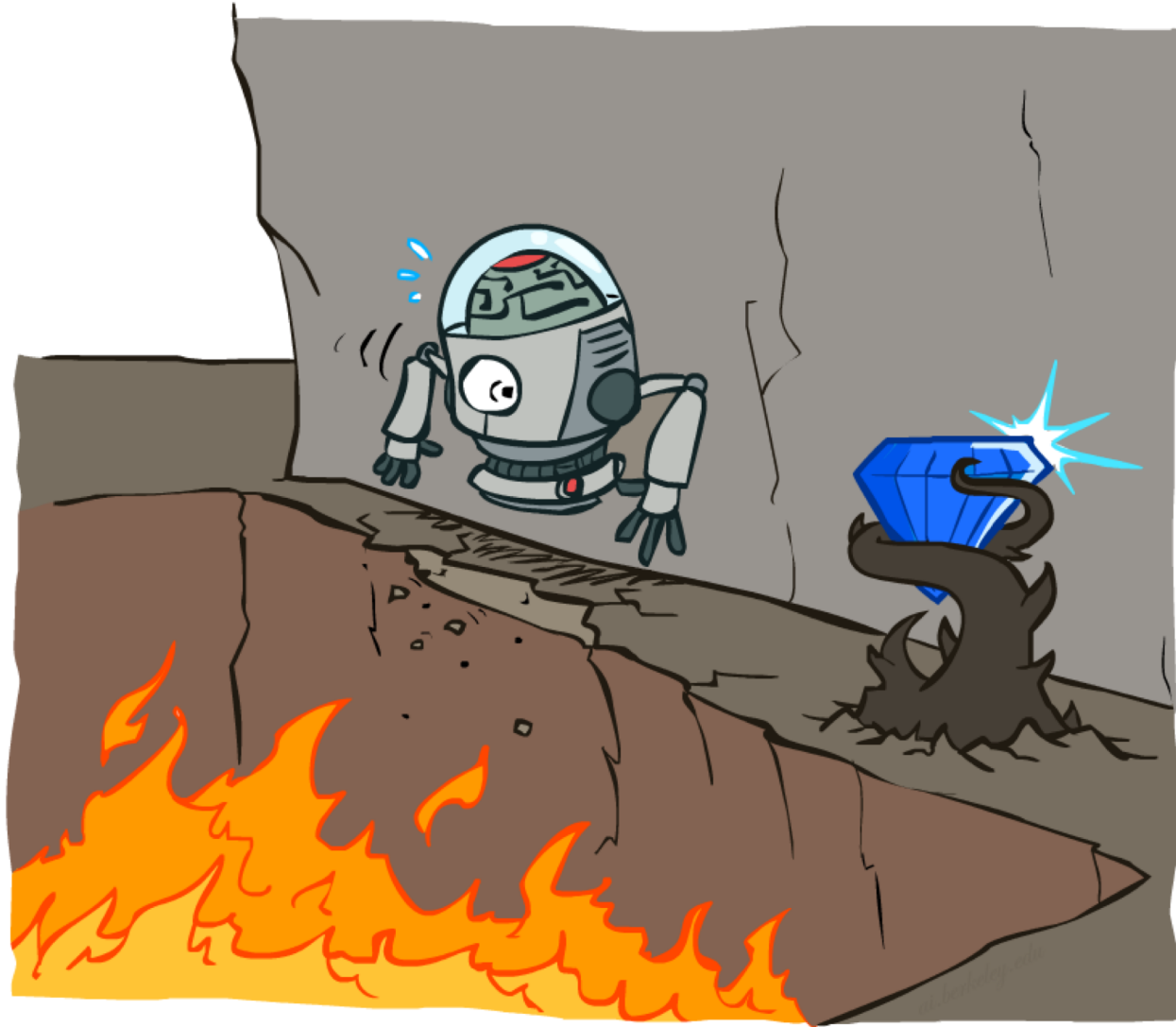- A **rational agent** selects actions that *maximize its utility function.*



Deterministic *vs.* **stochastic**

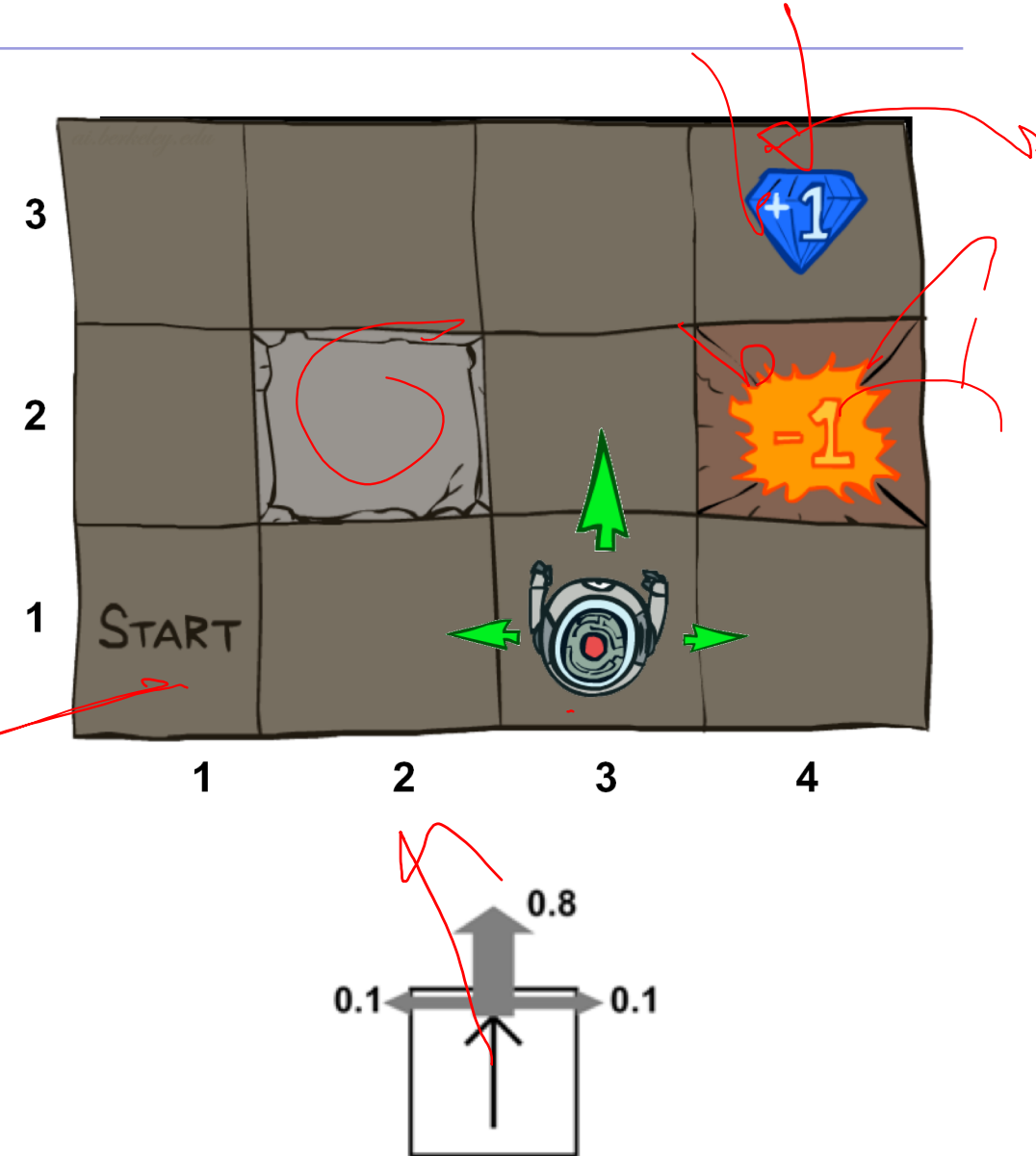**Fully observable** *vs.* partially observable
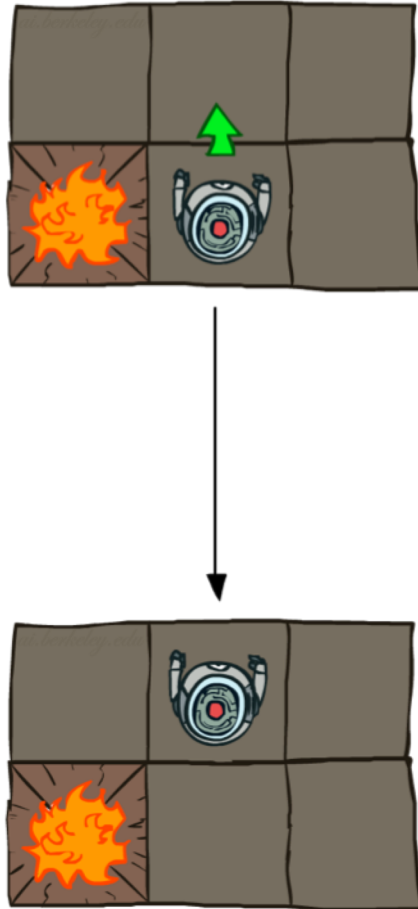
# Non-Deterministic Search

# 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)
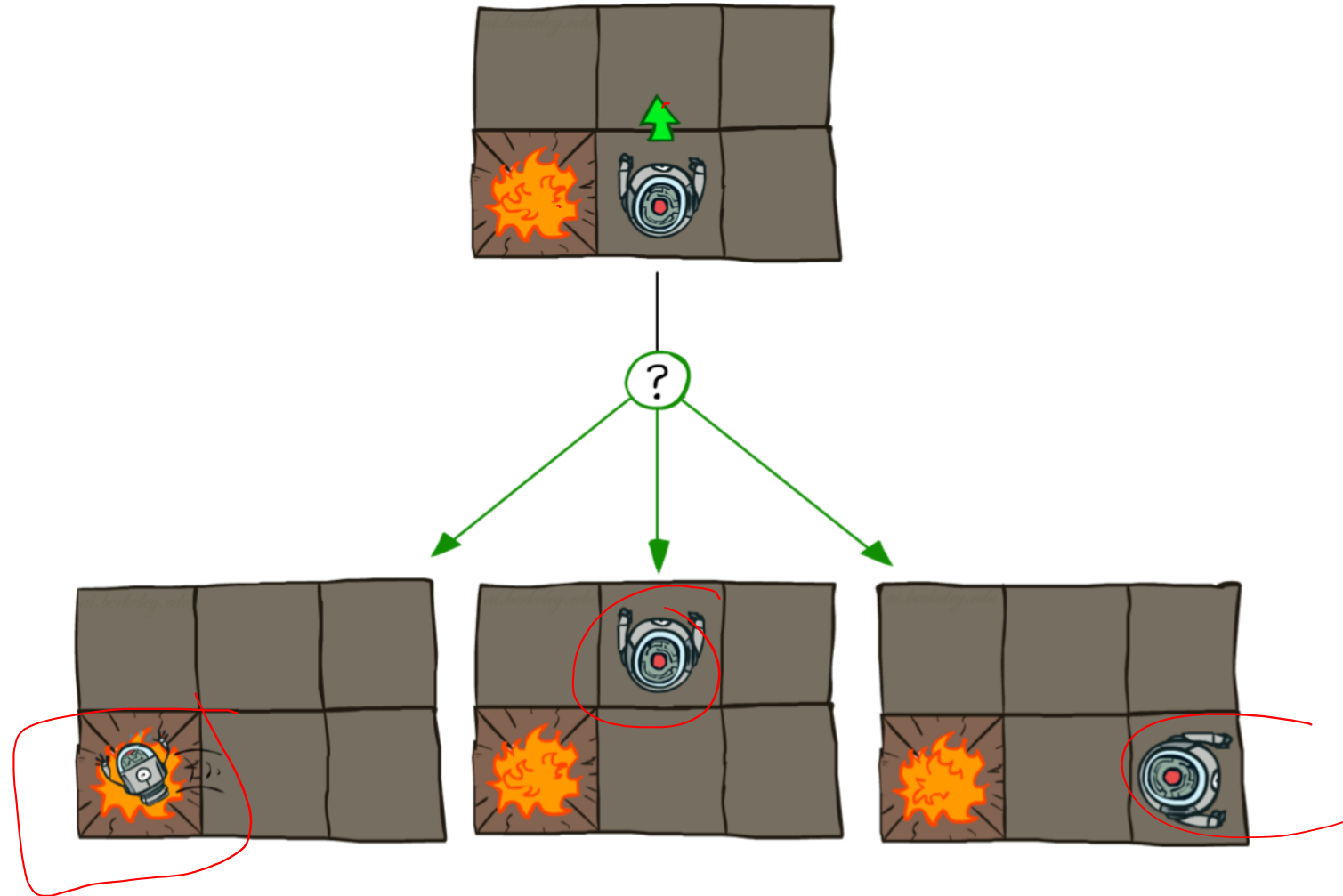
- Goal: maximize sum of rewards

# Grid World Actions

Deterministic Grid World

Stochastic Grid World

# Markov Decision Processes

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



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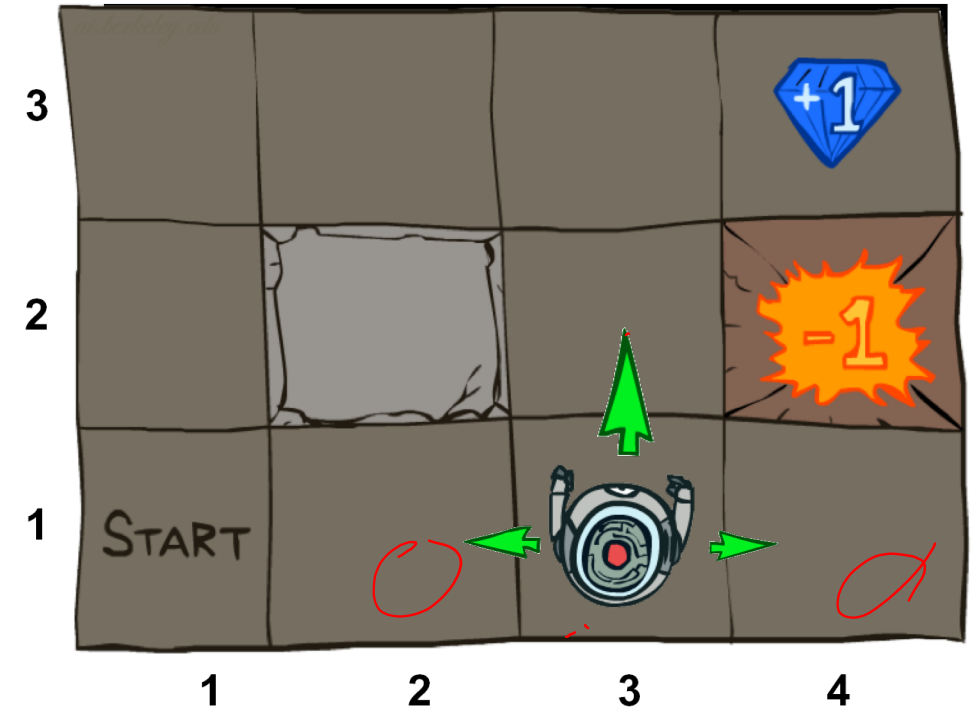$T(s_{31}, N, s_{11}) = 0$
$…$
$T(s_{31}, N, s_{32}) = 0.8$
$T(s_{31}, N, s_{21}) = 0.1$
$T(s_{31}, N, s_{41}) = 0.1$
$…$

*T is a Big Table!*
*11  X 4 x 11 = 484 entries*

**For now, we give this as input to the agent**

# Markov Decision Processes

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

...
R($s_{32}$, N, $s_{33}$) = -0.01    ← *Cost of breathing*
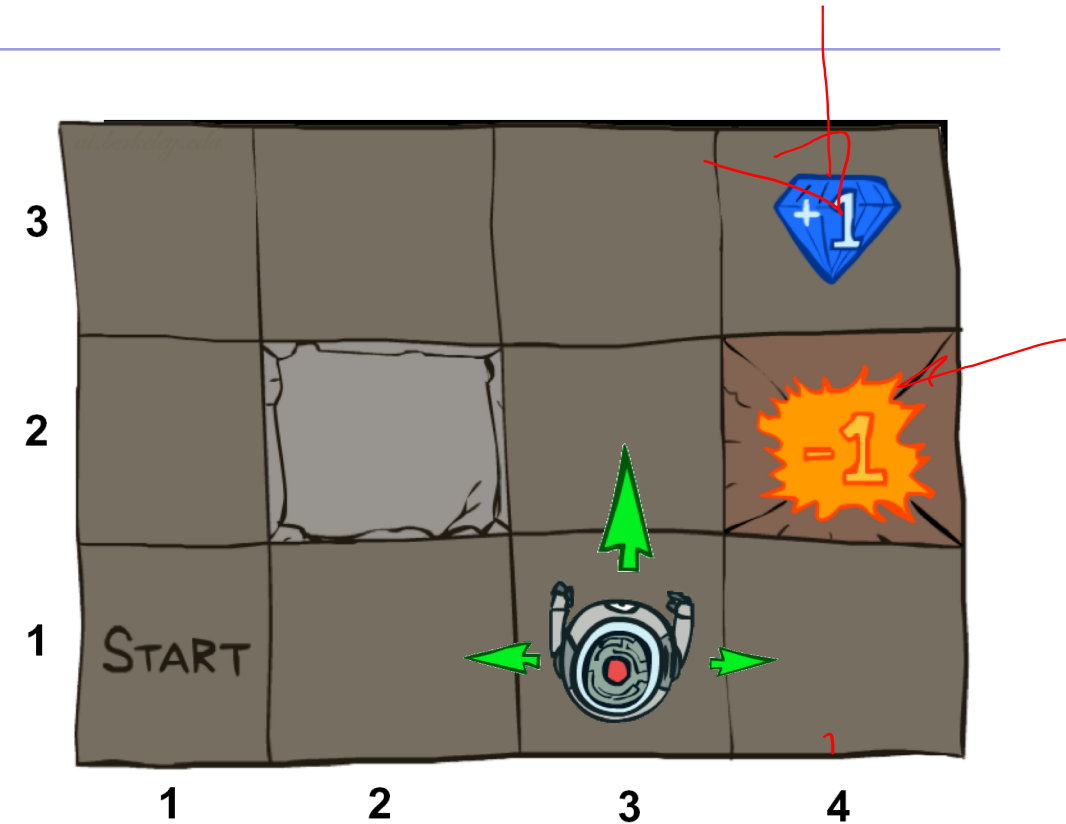...
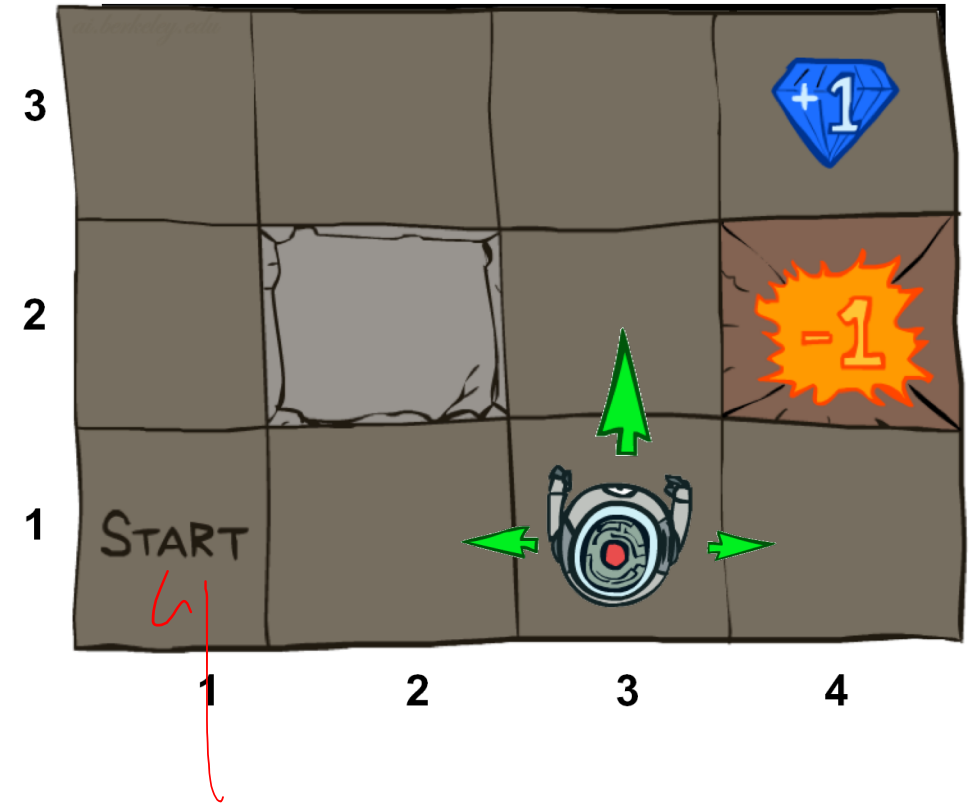R($s_{32}$, N, $s_{42}$) = -1.01
...
R($s_{33}$, E, $s_{43}$) = 0.99

R is also a Big Table!

For now, we also give this to the agent

# Markov Decision Processes

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

o MDPs are non-deterministic search problems
  o One way to solve them is with expectimax search
  o 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

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

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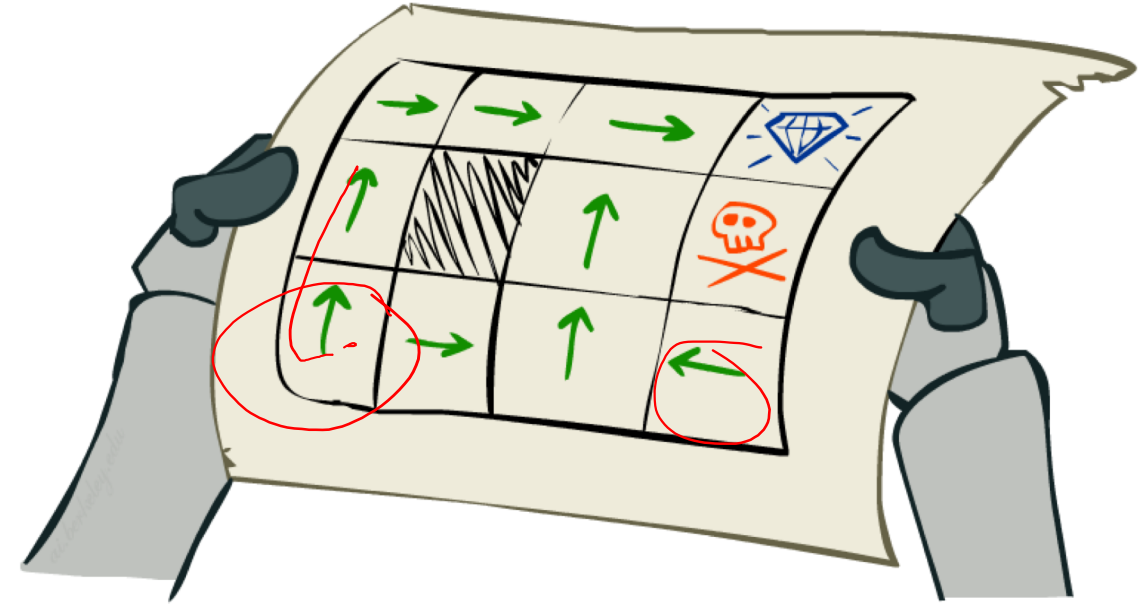
$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

- This is just like search, where the successor function could only depend on the current state (not the history)
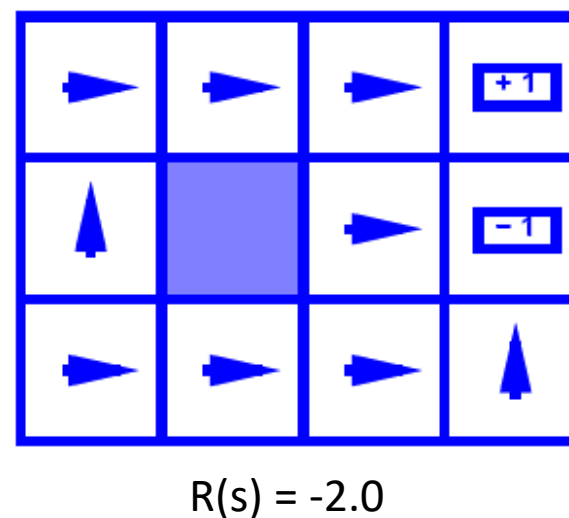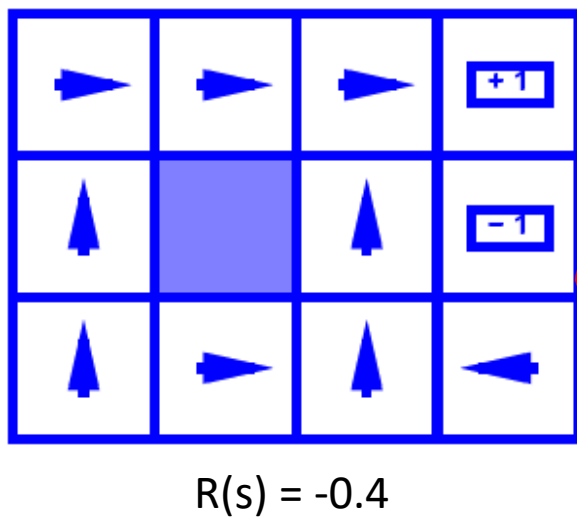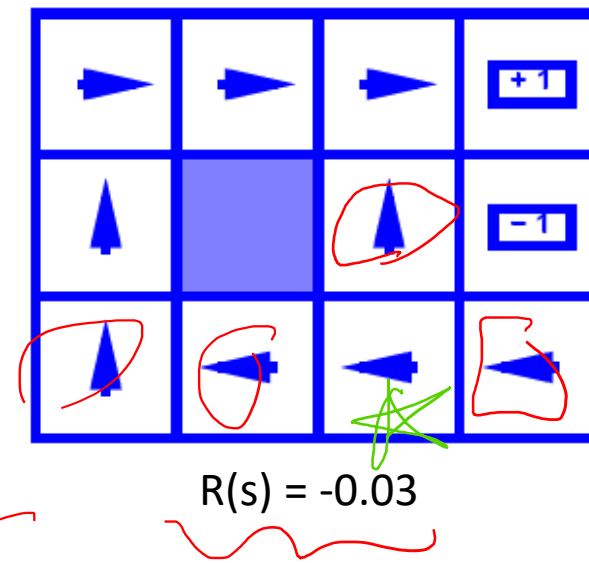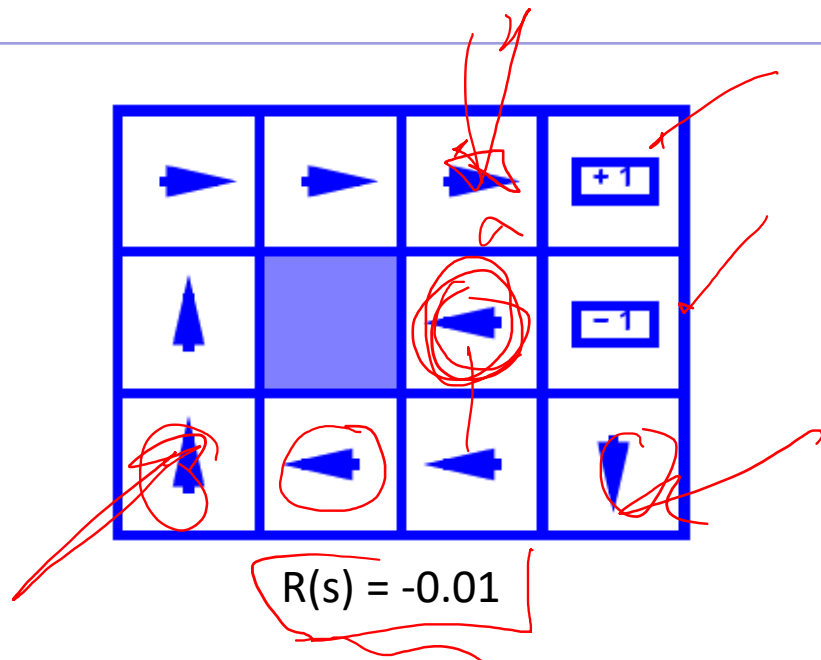
Andrey Markov
(1856-1922)

# Policies

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

o For MDPs, we want an optimal

   policy $\pi^*$: S $\rightarrow$ A

   o A policy $\pi$ gives an action for each state
   o An optimal policy is one that maximizes expected utility if followed
   o An explicit policy defines a reflex agent



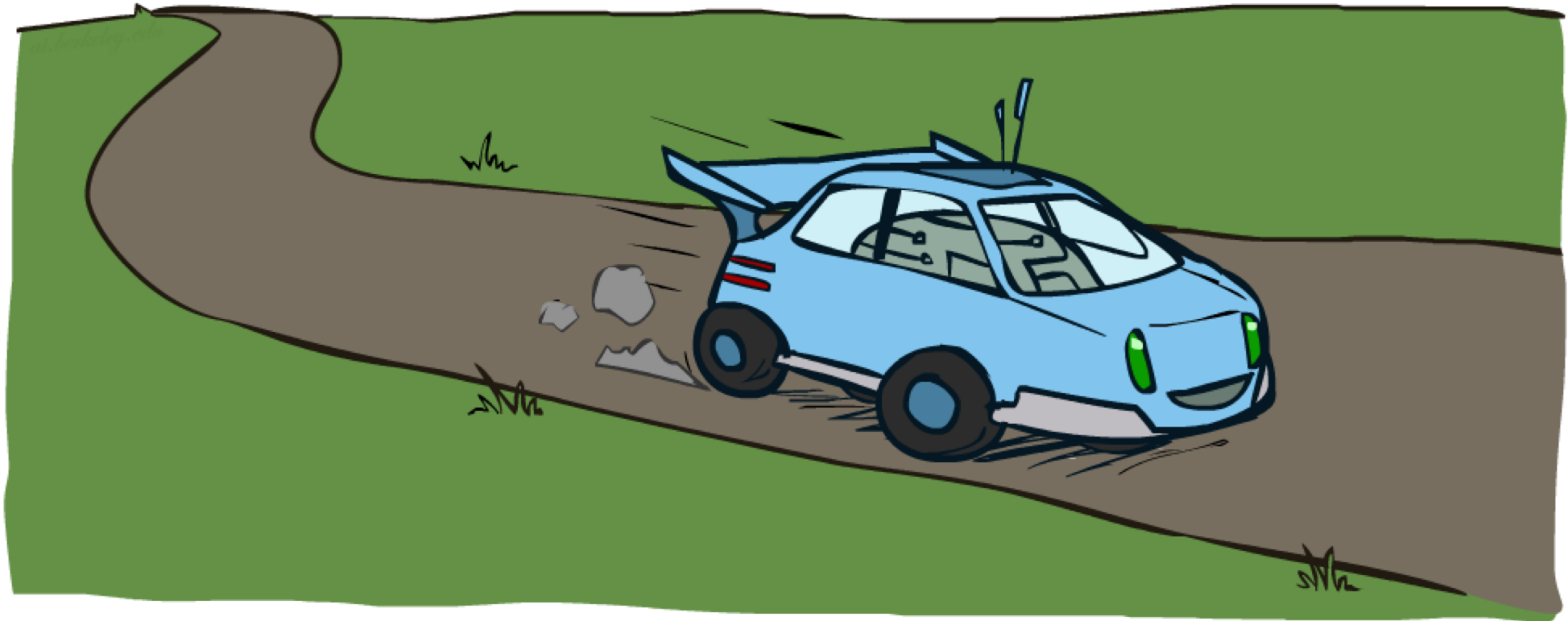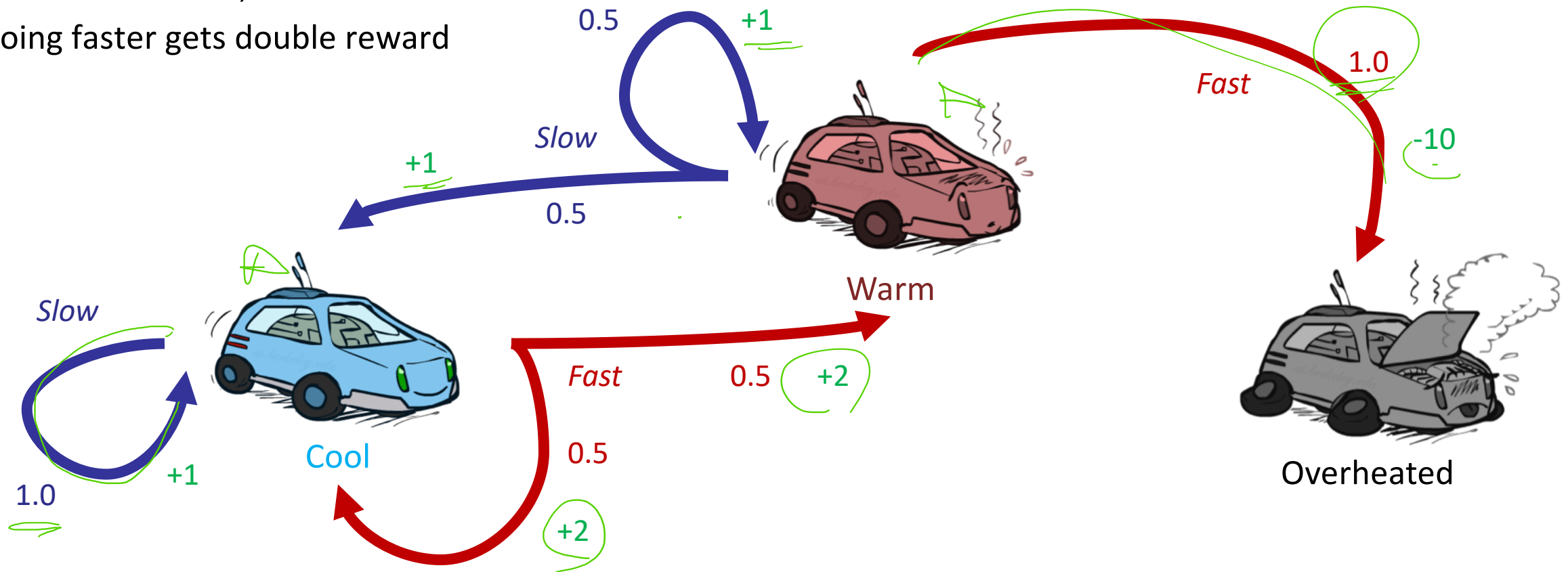Optimal policy when R(s, a, s') = -0.4 for all non-terminals s

# Optimal Policies



R(s) = -0.01
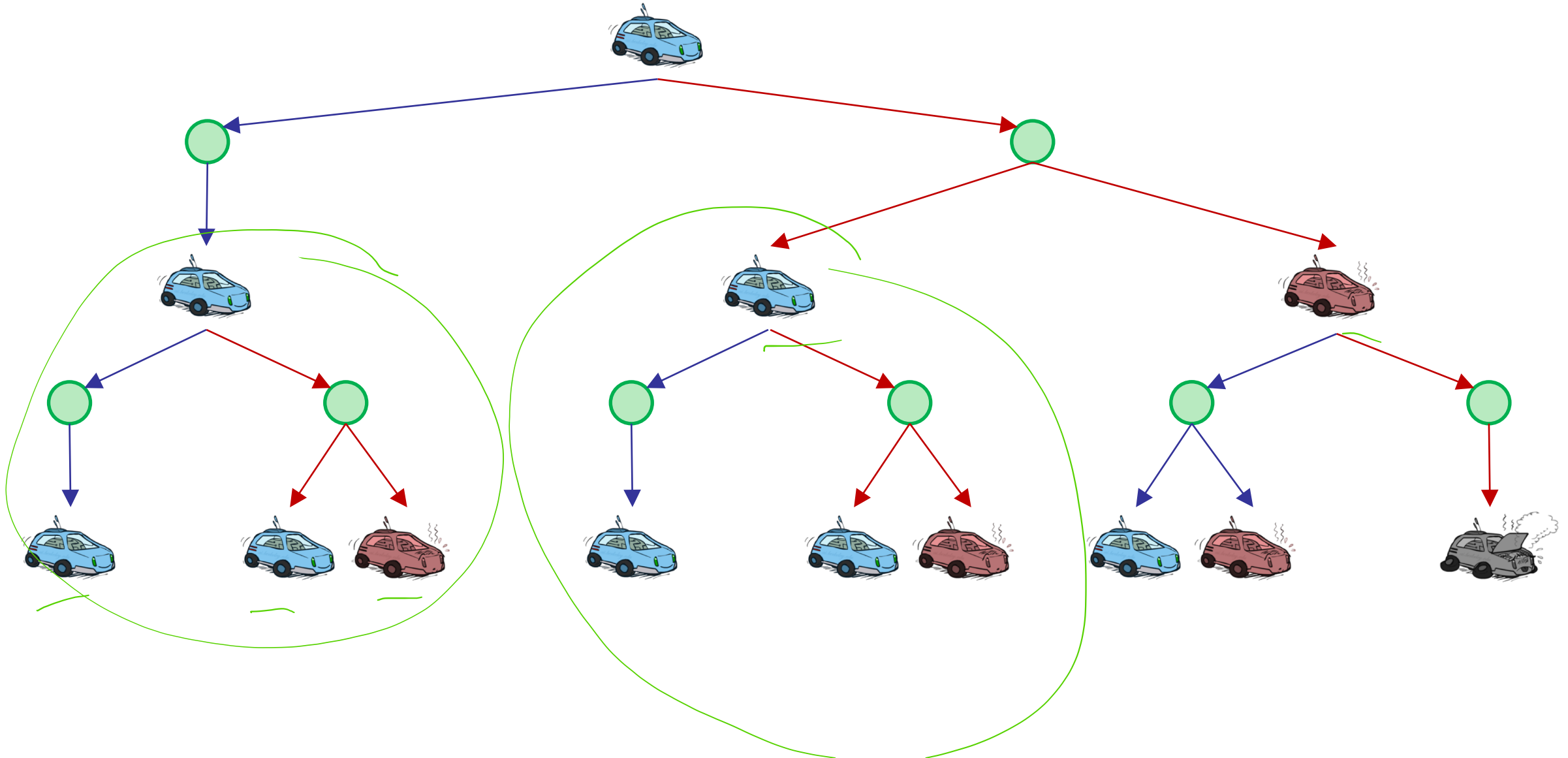
R(s) = -0.03

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

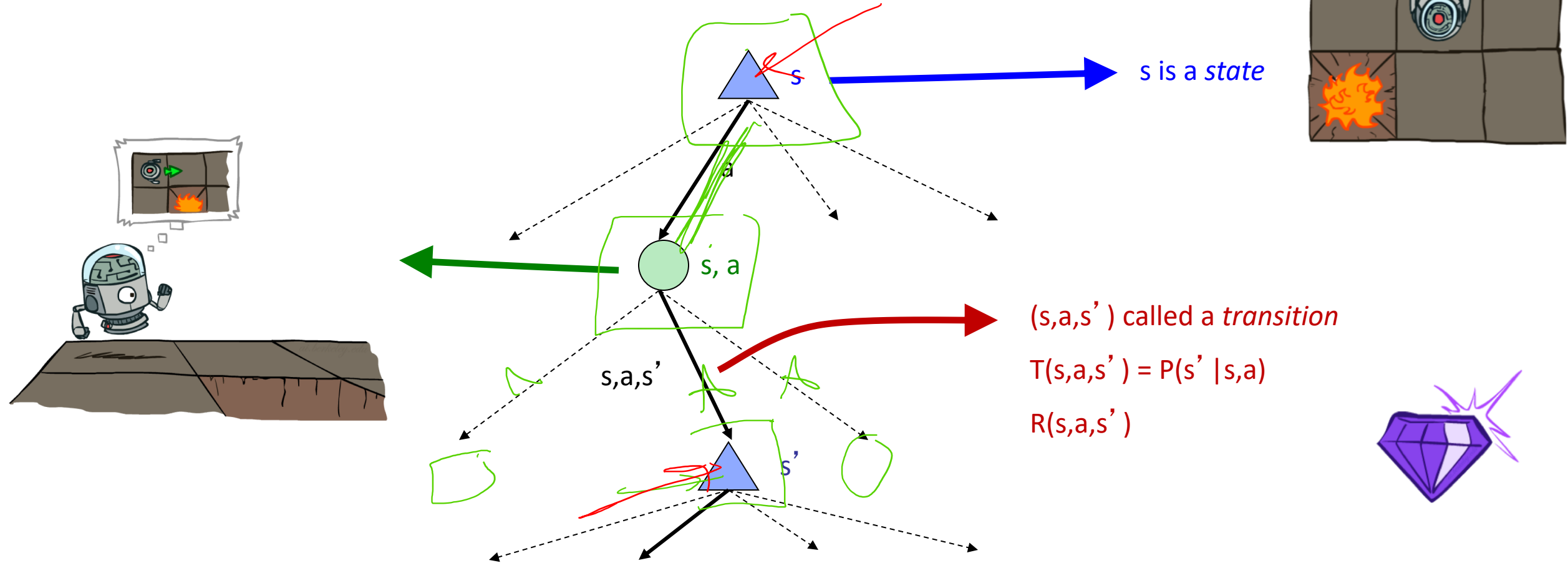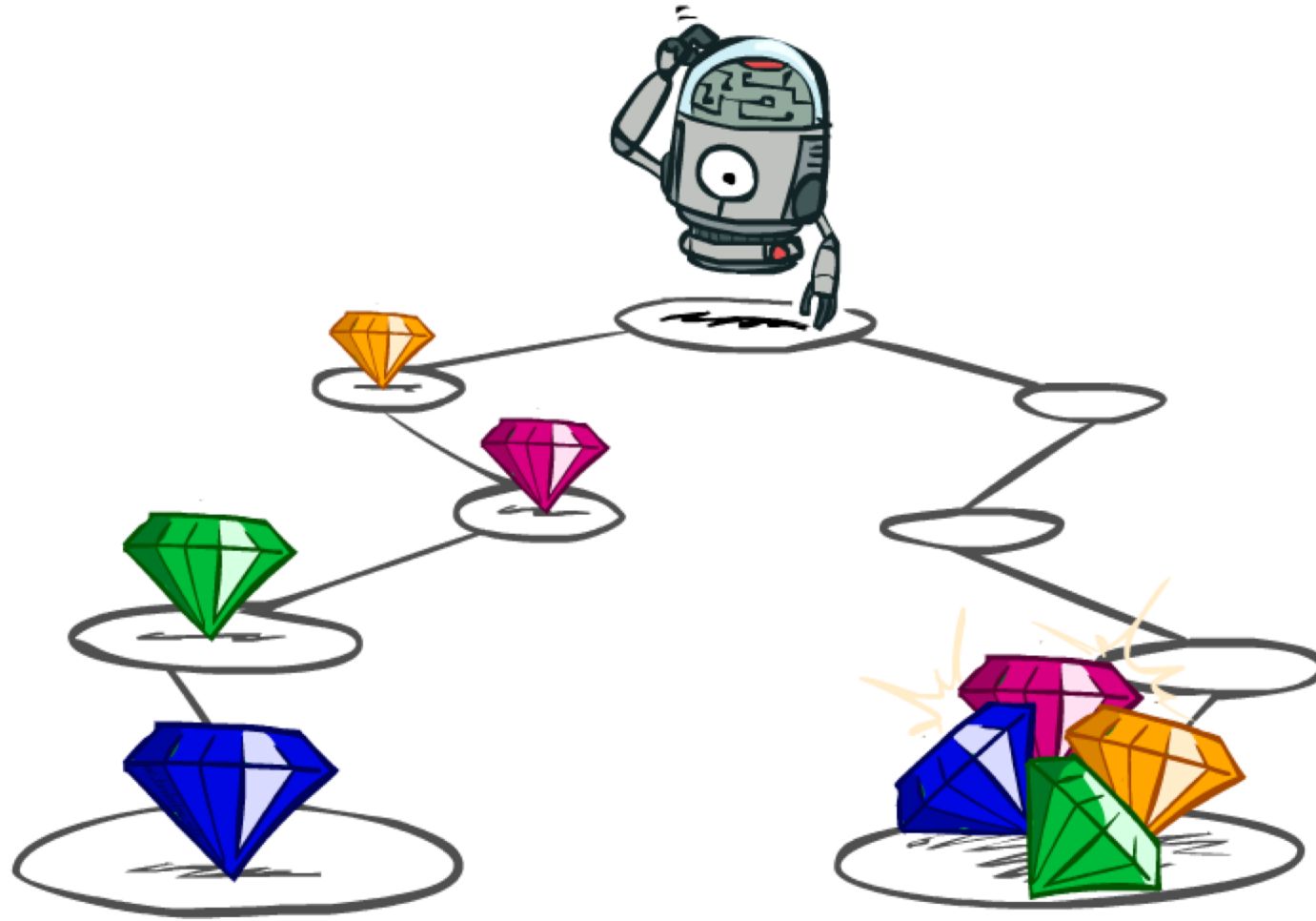# Racing Search Tree

# MDP Search Trees

o Each MDP state projects an expectimax-like search tree



s is a *state*

s, a

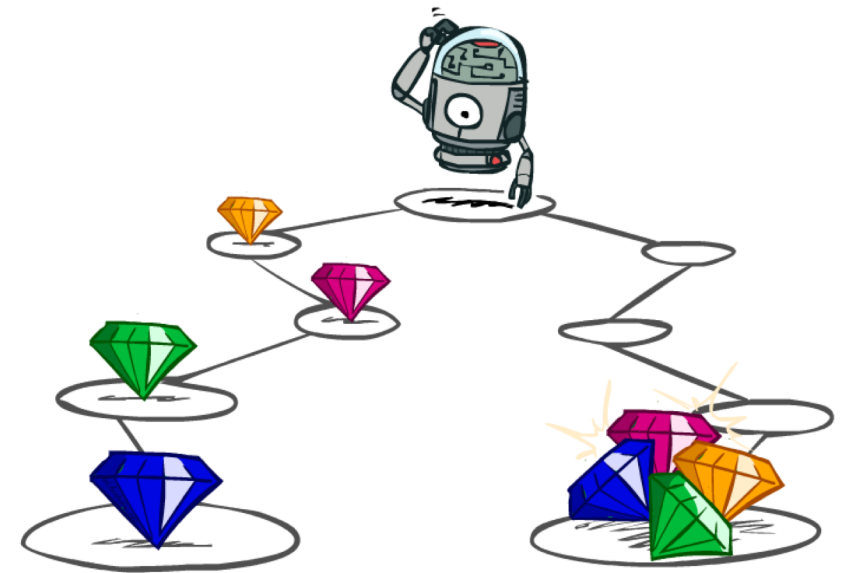(s,a,s') called a *transition*

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

$R(s,a,s')$

s,a,s'

s'

# Utilities of Sequences

# 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

o It's reasonable to maximize the sum of rewards

o It's also reasonable to prefer rewards now to rewards later

o One solution: values of rewards decay exponentially



$1$

Worth Now

$\gamma$

Worth Next Step
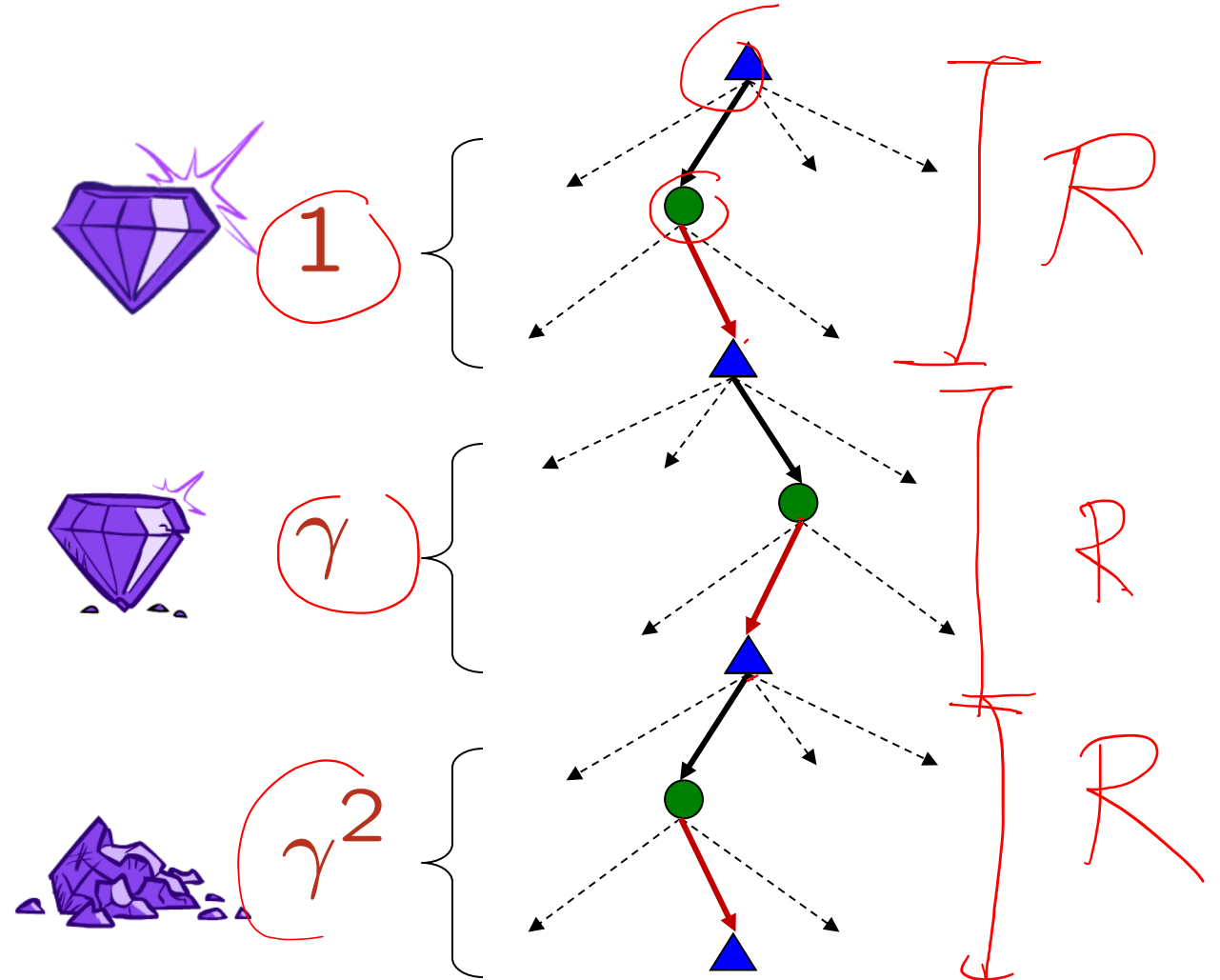
$\gamma^2$

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
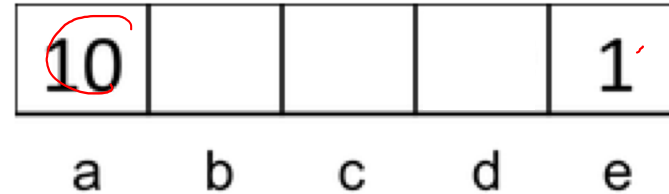  - o Also helps our algorithms converge

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
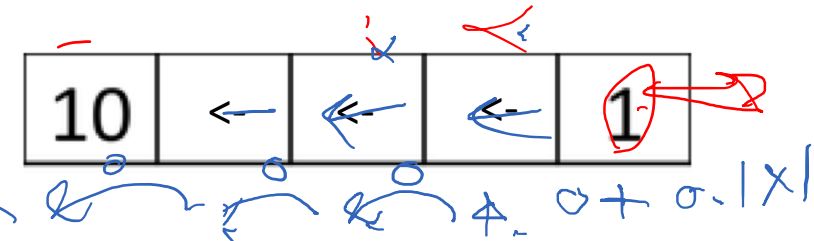
o Given:

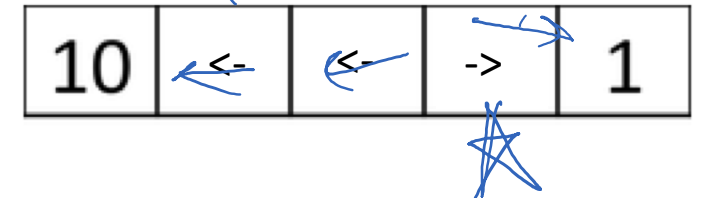| 10 | | | | 1 |
|----|---|---|---|---|
| a | b | c | d | e |

　　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?

1γ=10 γ³

$$\gamma = 10\gamma^3$$

# Infinite Utilities?!

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

- Solutions:
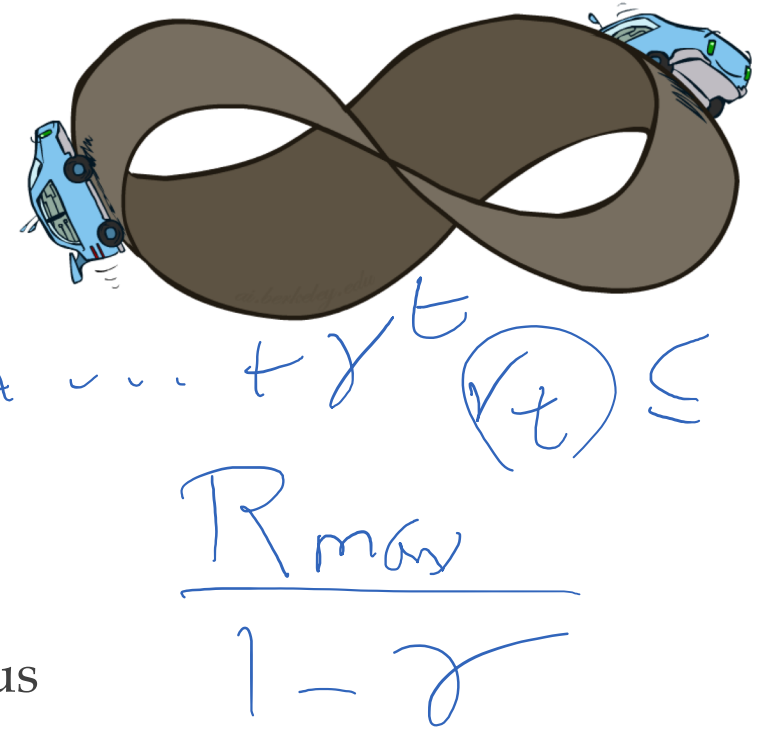  - Finite horizon: (similar to depth-limited search)
    - Terminate episodes after a fixed T steps (e.g. life)
    - Policy π depends on time left
  - Discounting: use $0 < \gamma < 1$
    $$U([r_0, \ldots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \le R_{\max}/(1 - \gamma)$$
    - Smaller γ means smaller "horizon" – shorter term focus

  $r_0 + \gamma(r_1) + \cdots + \gamma^t (r_t) \le$

  $\dfrac{R_{max}}{1 - \gamma}$

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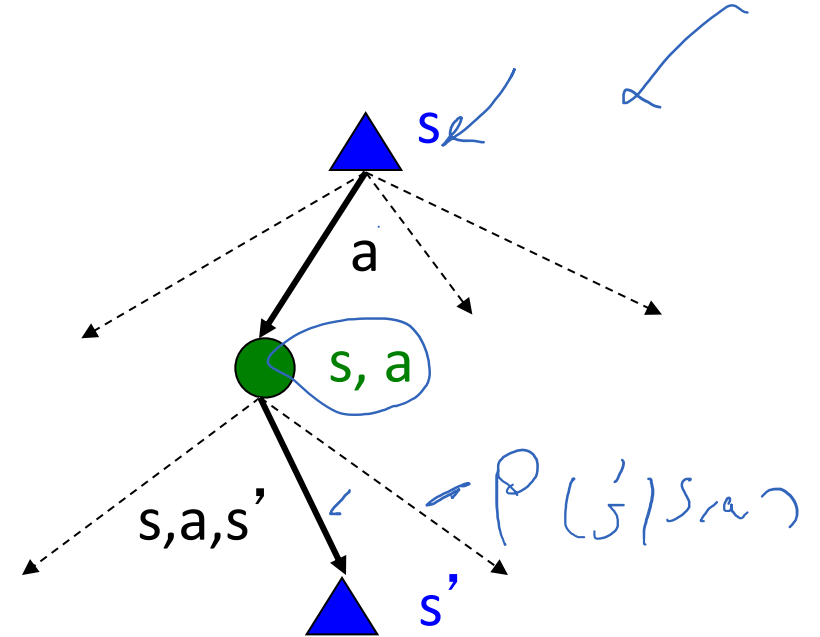
# Recap: Defining MDPs

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