# CS 573: Artificial Intelligence

## Markov Decision Processes
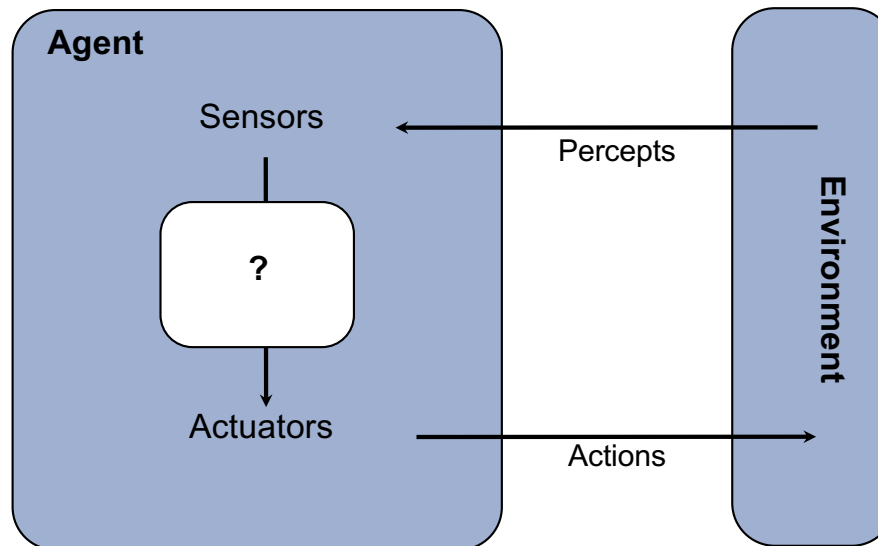
Dan Weld

University of Washington

# Outline

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

  - Markov Decision Processes
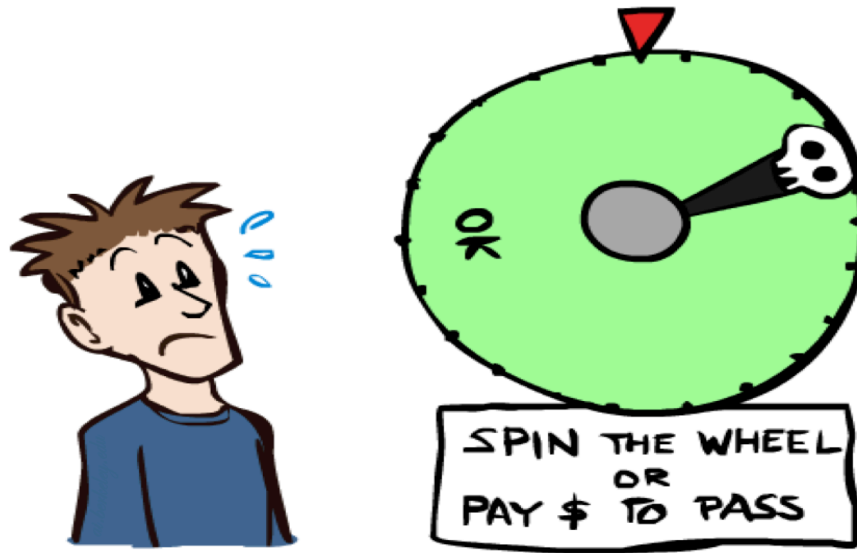  - Reinforcement Learning

# Agent *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

# Human Utilities

# Utility Scales

- WoLoG Normalized utilities: $u_+$ = 1.0, $u_-$ = 0.0

- Micromorts: one-millionth chance of death, useful for paying to reduce product risks, etc.

- QALYs: quality-adjusted life years, useful for medical decisions involving substantial risk

- Note: behavior is invariant under positive *linear* transformation

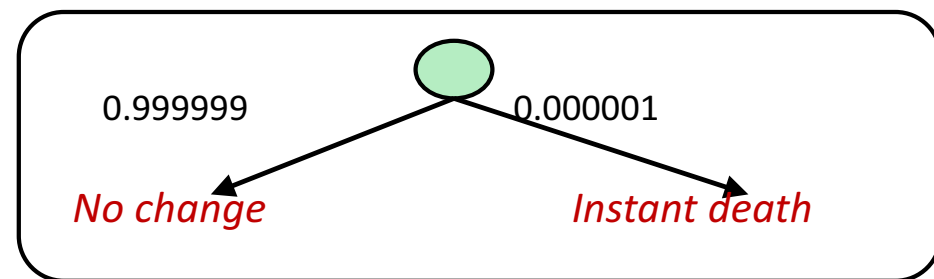$$U'(x) = k_1 U(x) + k_2 \quad \text{where } k_1 > 0$$

# Human Utilities

- Utilities map states to real numbers. Which numbers?
- Standard approach to assessment (elicitation) of human utilities:
  - Compare a prize A to a standard lottery $L_p$ between
    - "best possible prize" $u_+$ with probability p
    - "worst possible catastrophe" $u_-$ with probability 1-p
  - Adjust lottery probability p until indifference: A ~ $L_p$
  - Resulting p is a utility in [0,1]



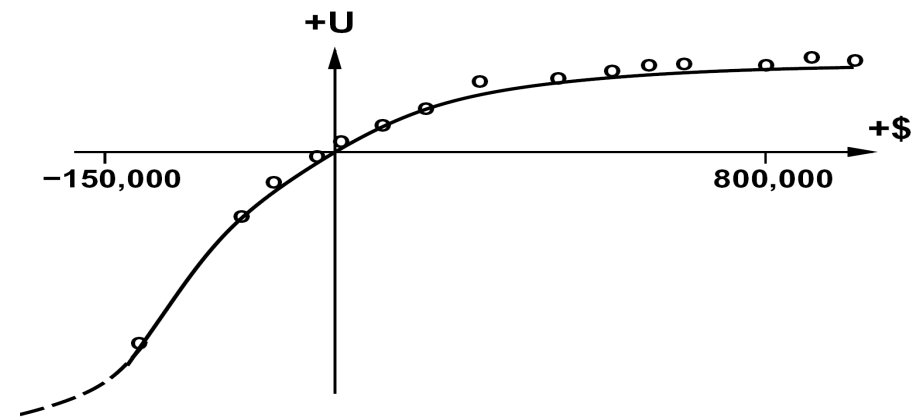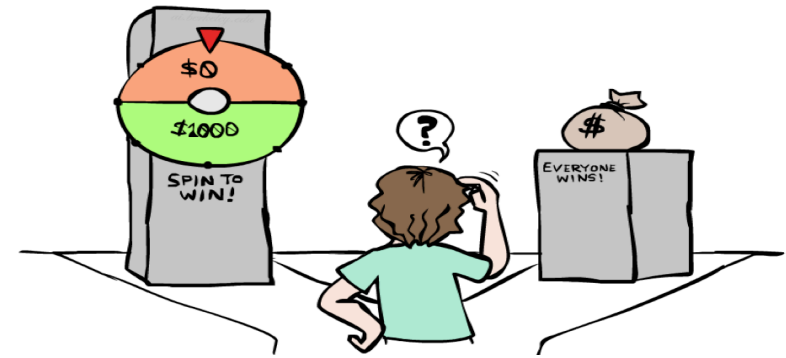*Pay $30*   ~   0.999999   0.000001   *No change*   *Instant death*

# Money

- Money <u>does not</u> behave as a utility function, but we can talk about the utility of having money (or being in debt)
- Given a lottery L = [p, $X; (1-p), $Y]
  - The expected monetary value EMV(L) is p*X + (1-p)*Y
  - U(L) = p*U($X) + (1-p)*U($Y)
  - Typically, U(L) < U( EMV(L) )
  - In this sense, people are risk-averse
  - When deep in debt, people are risk-prone

# Example: Insurance

Consider the lottery [0.5, $1000;  0.5, $0]

- What is its expected monetary value?  ($500)
- What is its certainty equivalent?
    - Monetary value acceptable in lieu of lottery
    - $400 for most people
- Difference of $100 is the insurance premium
    - There's an insurance industry because people will pay to reduce their risk
    - If everyone were risk-neutral, no insurance needed!

- It's win-win: you'd rather have the $400 and the insurance company would rather have the lottery (their utility curve is flat and they have many lotteries)

# Rational Preferences

The Axioms of Rationality



Orderability
$$(A \succ B) \vee (B \succ A) \vee (A \sim B)$$
Transitivity
$$(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$$
Continuity
$$A \succ B \succ C \Rightarrow \exists p \, [p, A; \; 1 - p, C] \sim B$$
Substitutability
$$A \sim B \Rightarrow [p, A; \; 1 - p, C] \sim [p, B; 1 - p, C]$$
Monotonicity
$$A \succ B \Rightarrow$$
$$(p \geq q \Leftrightarrow [p, A; \; 1 - p, B] \succeq [q, A; \; 1 - q, B])$$

Theorem: Rational preferences imply behavior describable as maximization of expected utility

# MEU Principle

- **Theorem [Ramsey, 1931; von Neumann & Morgenstern, 1944]**
  - Given any preferences satisfying these constraints, there exists a real-valued function U such that:

$$U(A) \geq U(B) \iff A \succeq B$$

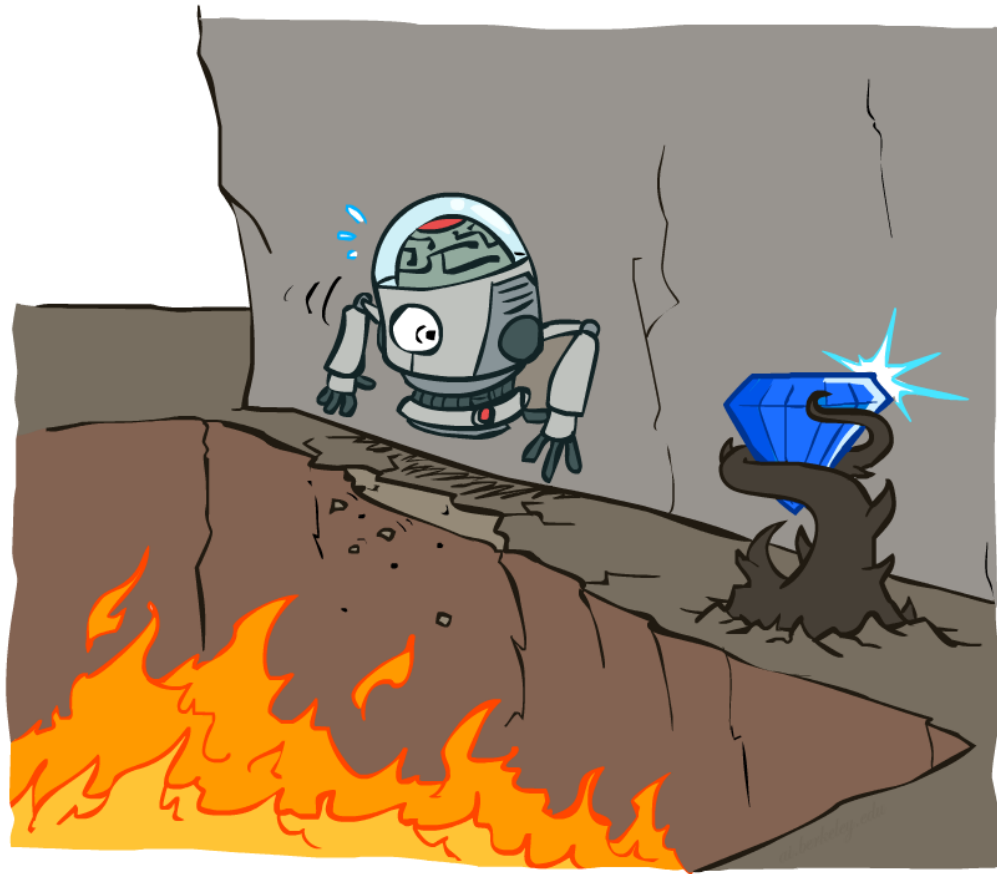$$U([p_1, S_1; \ldots ; p_n, S_n]) = \sum_i p_i U(S_i)$$

  - I.e. values assigned by U preserve preferences of both prizes and lotteries!

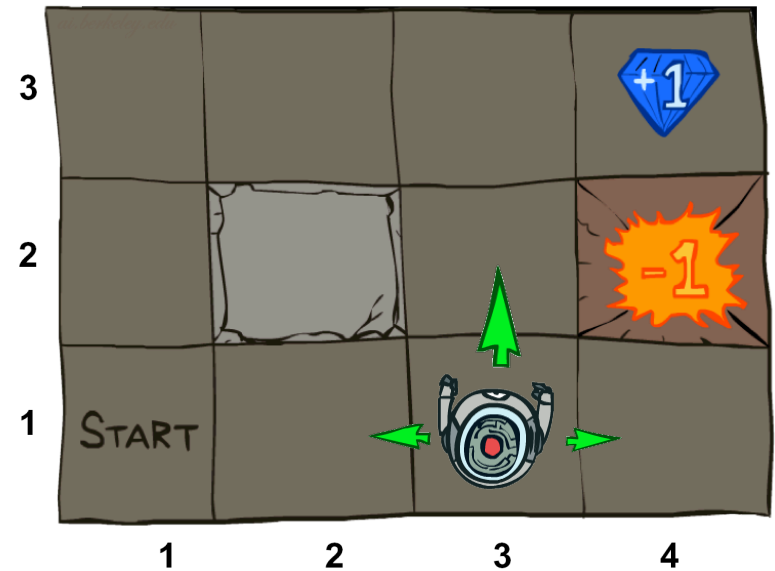- **Maximum expected utility (MEU) principle:**
  - Choose the action that maximizes expected utility
  - Note: an agent can be entirely rational (consistent with MEU) without ever representing or manipulating utilities and probabilities
  - E.g., a lookup table for perfect tic-tac-toe, a reflex vacuum cleaner

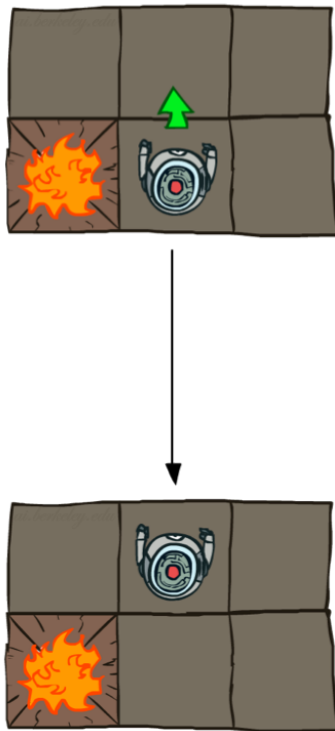# Non-Deterministic Search

# 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)
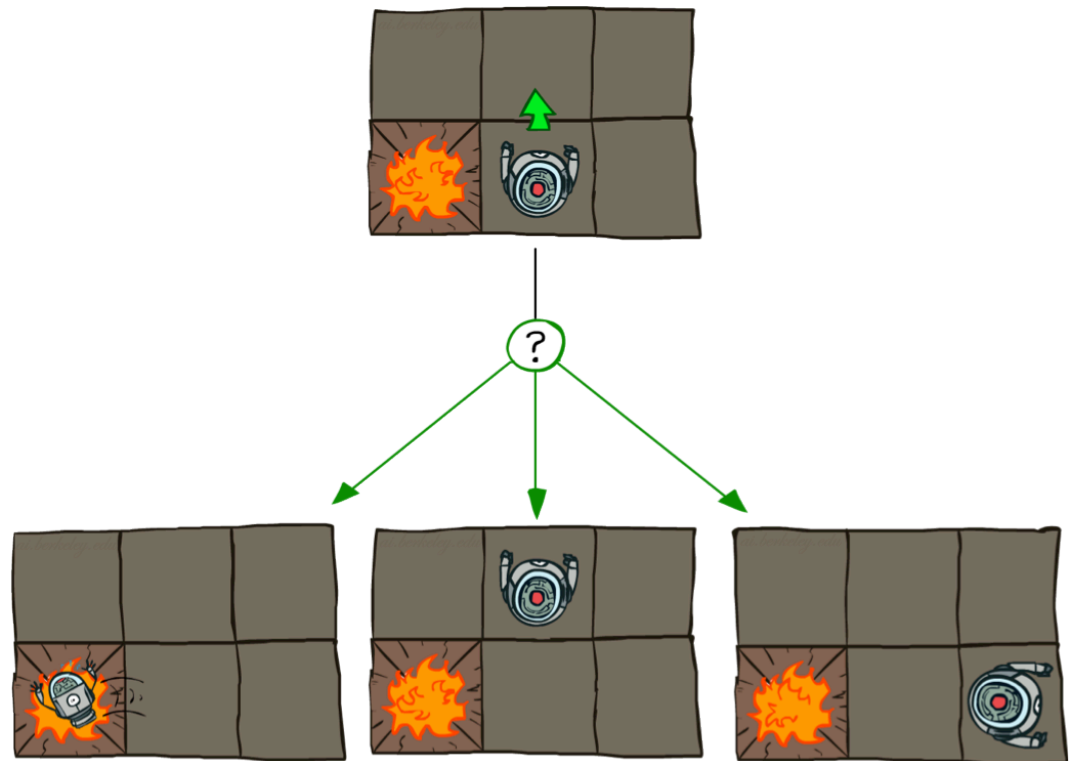
- Goal: maximize sum of rewards

# 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



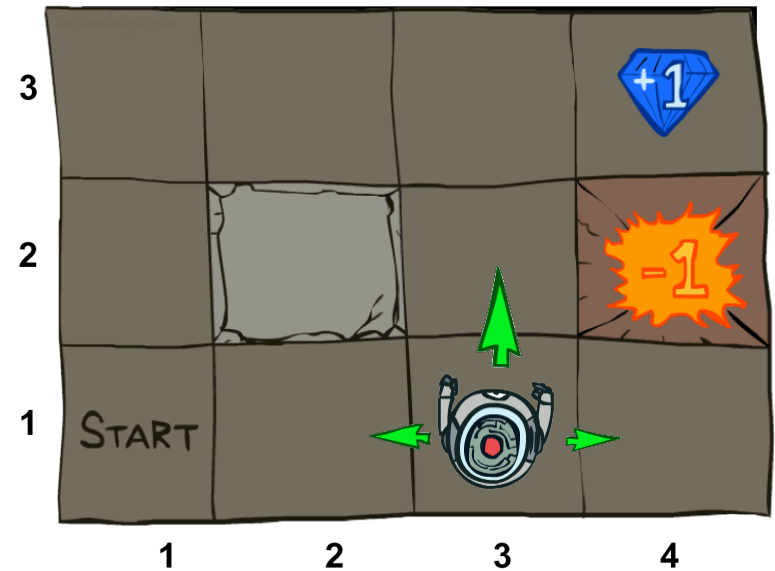| | | | |
|---|---|---|---|
| **3** | | | |
| **2** | | | |
| **1** | START | | |
| | **1** | **2** | **3** | **4** |

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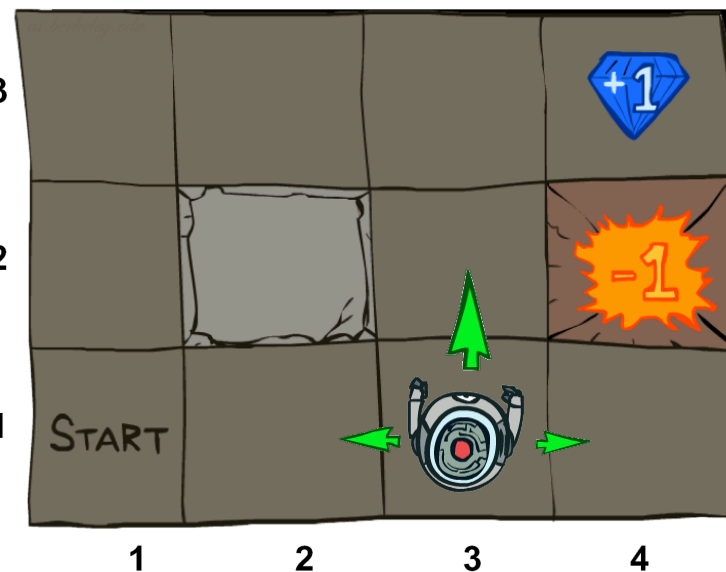$T(s_{31}, N, s_{41}) = 0.1$
$\ldots$

*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

- 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')



$$R(s_{32}, \ddot{N}, s_{33}) = -0.01$$
...
$$R(s_{32}, \ddot{N}, s_{42}) = -1.01$$
...
$$R(s_{33}, \ddot{E}, s_{43}) = 0.99$$

***Cost of breathing***
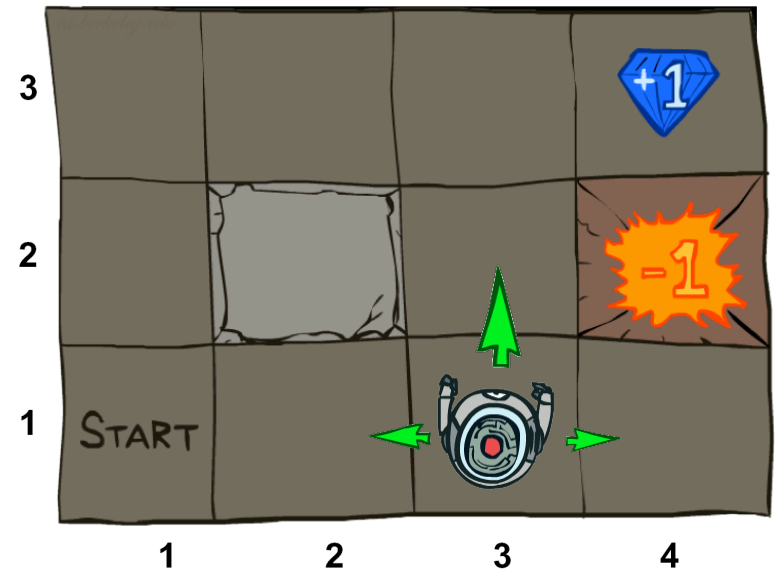
R is also a Big Table!

For now, we also give this to the agent

# 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 $R(s, a, s')$
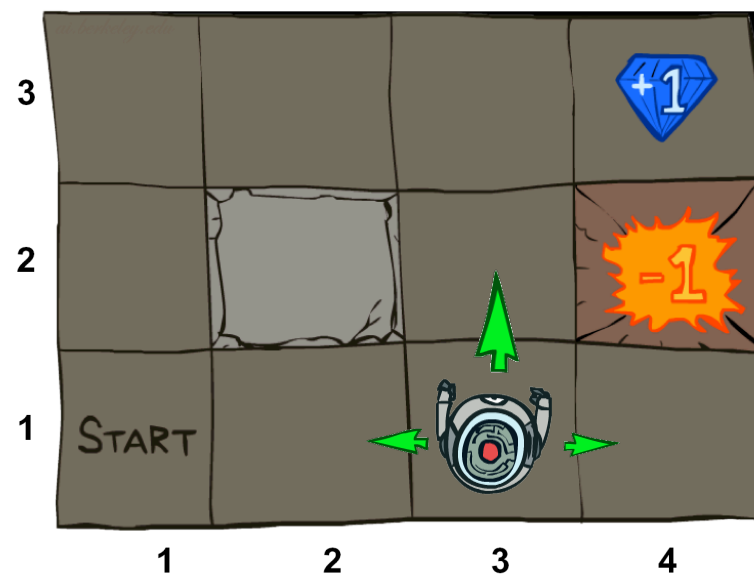    - Sometimes just $R(s)$ or $R(s')$

...

$R(s_{33}) = -0.01$

$R(s_{42}) = -1.01$

$R(s_{43}) = 0.99$

# 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'), e.g. in R&N
  - 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

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

$$=$$

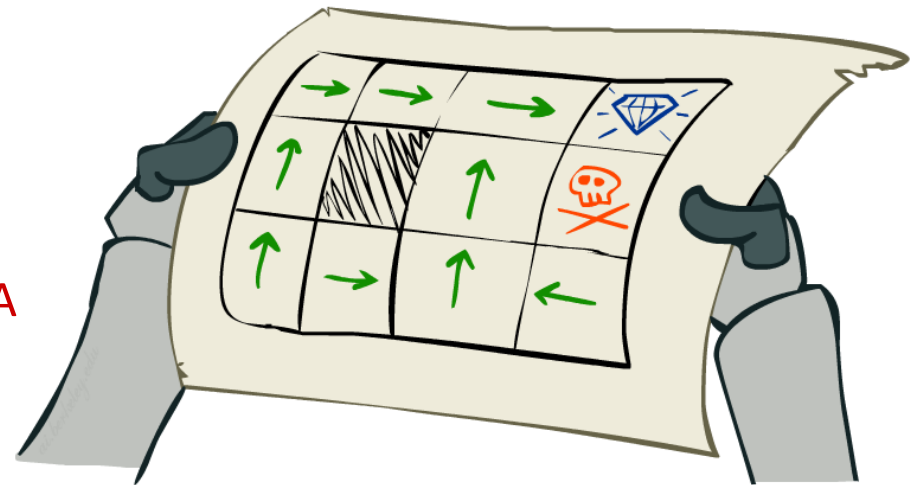$$P(S_{t+1} = s' | S_t = s_t, A_t = a_t)$$

Andrey Markov
(1856-1922)

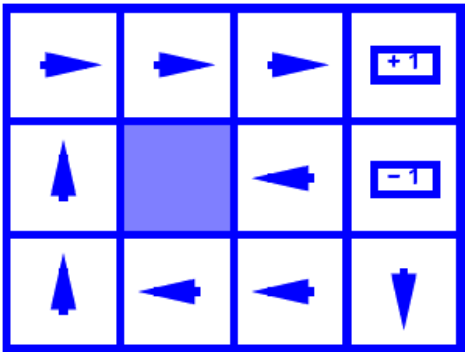- This is just like search, where the successor function can only depend on the current state (not the history)

# 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
  - An explicit policy defines a reflex agent

- Expectimax didn't output an entire policy
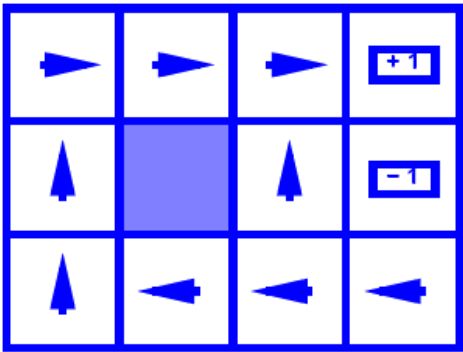  - It computed the action for a single state only



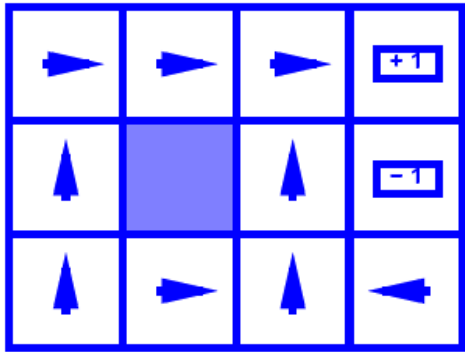Optimal policy when R(s, a, s') = -0.03 for all non-terminals s

# Optimal Policies



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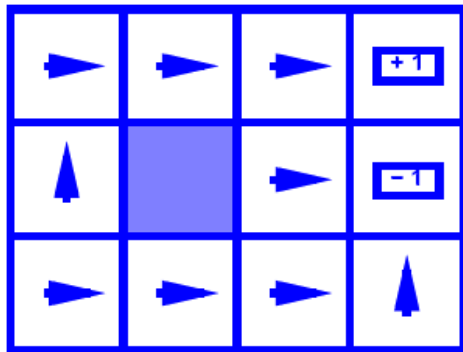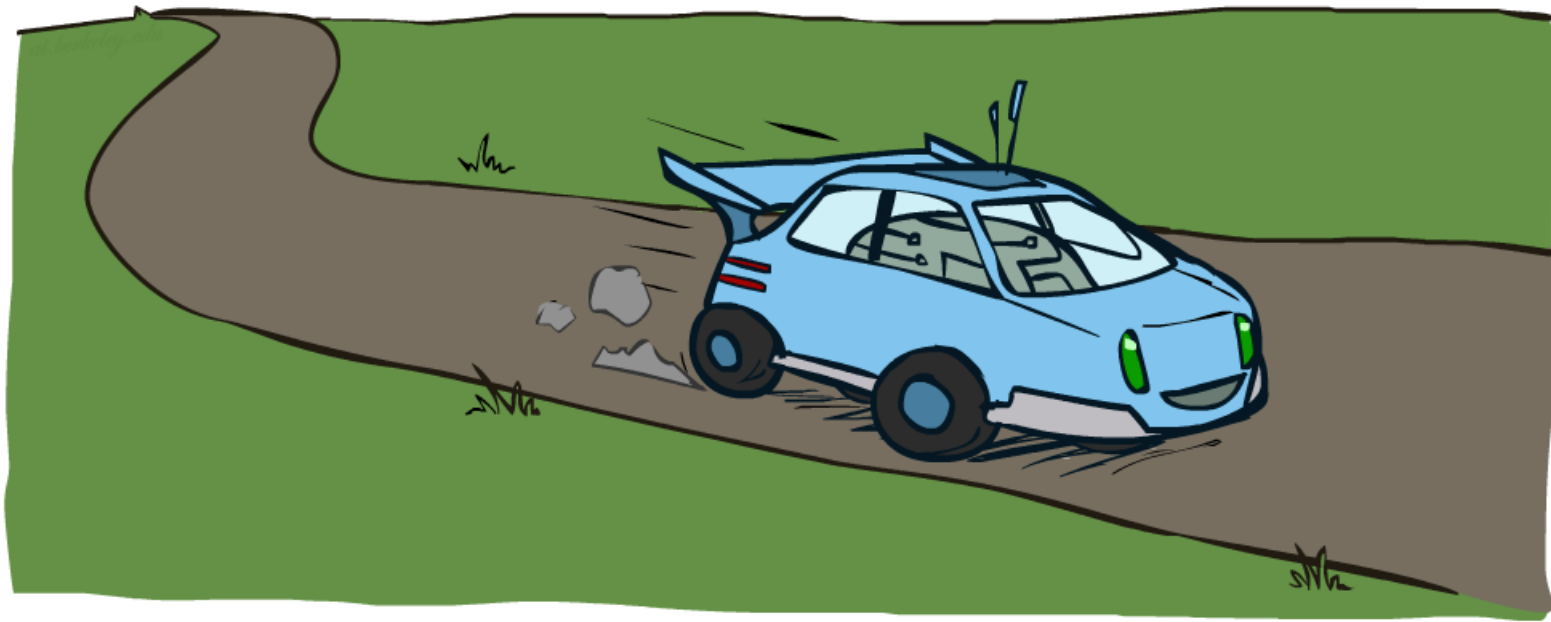
R(s) = -0.03

Cost of breathing

R(s) = -0.4
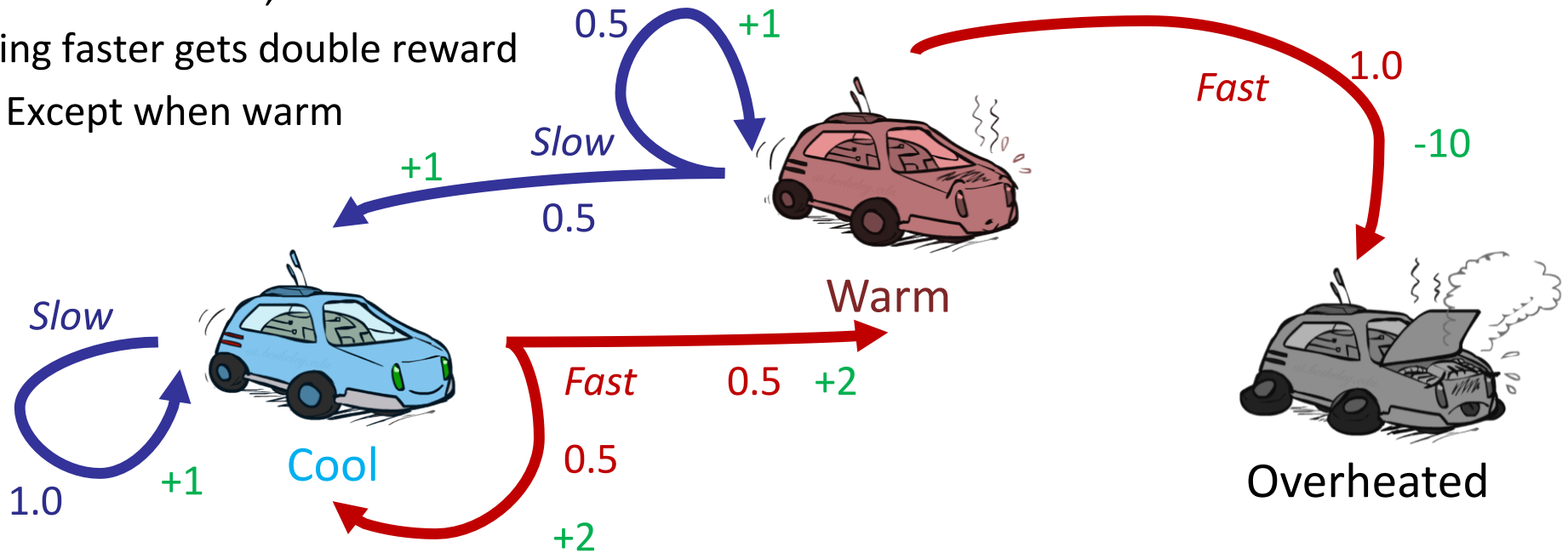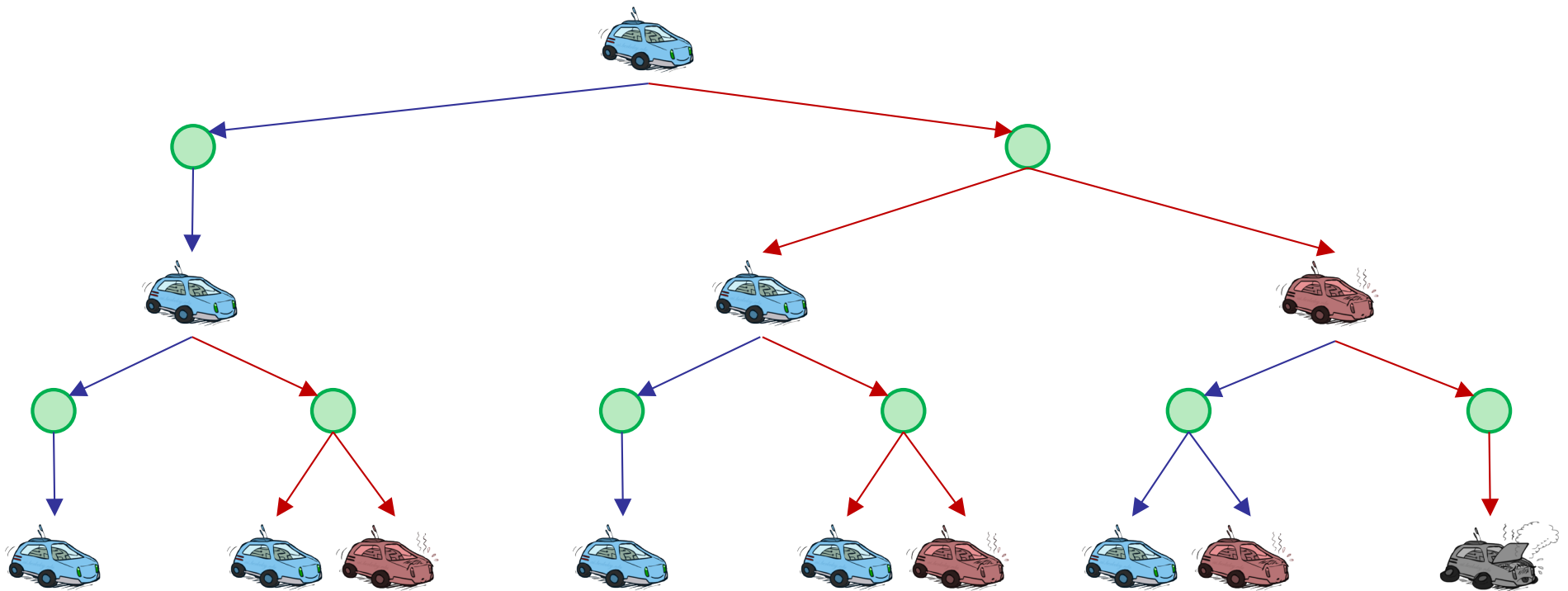
R(s) = -2.0

# 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
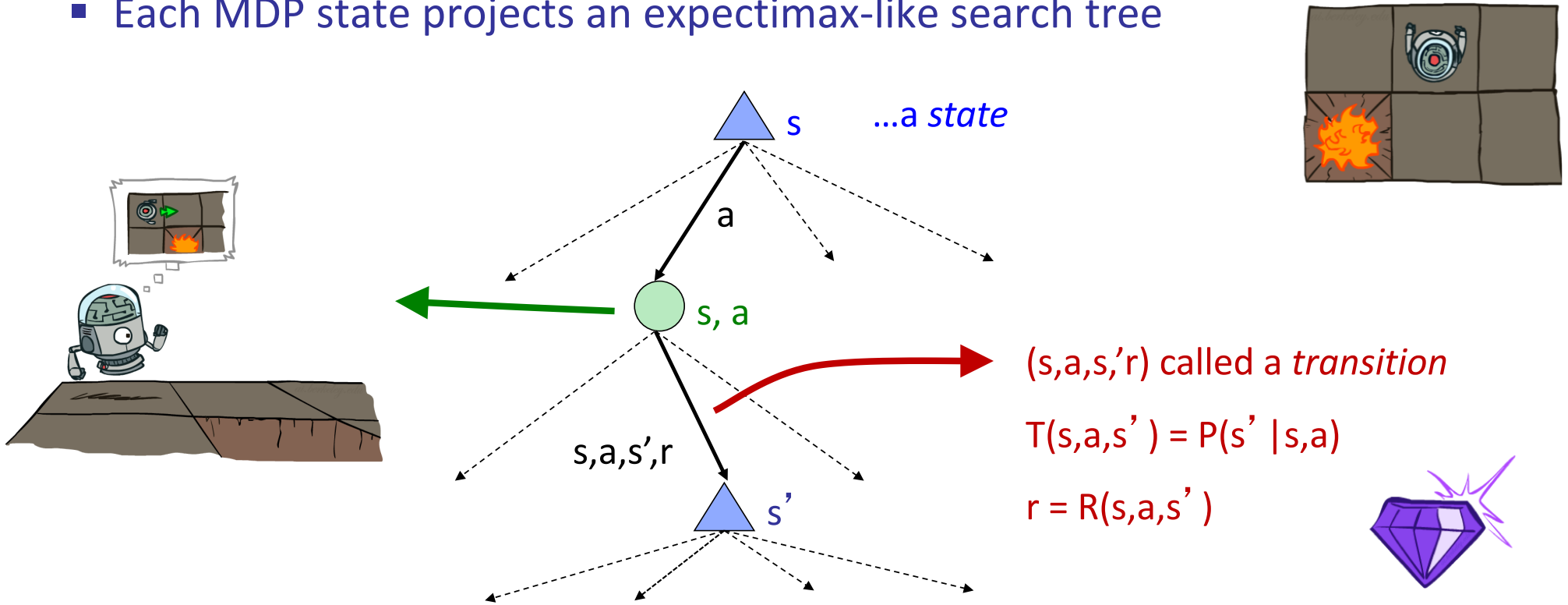  - Except when warm

# Racing: Search Tree



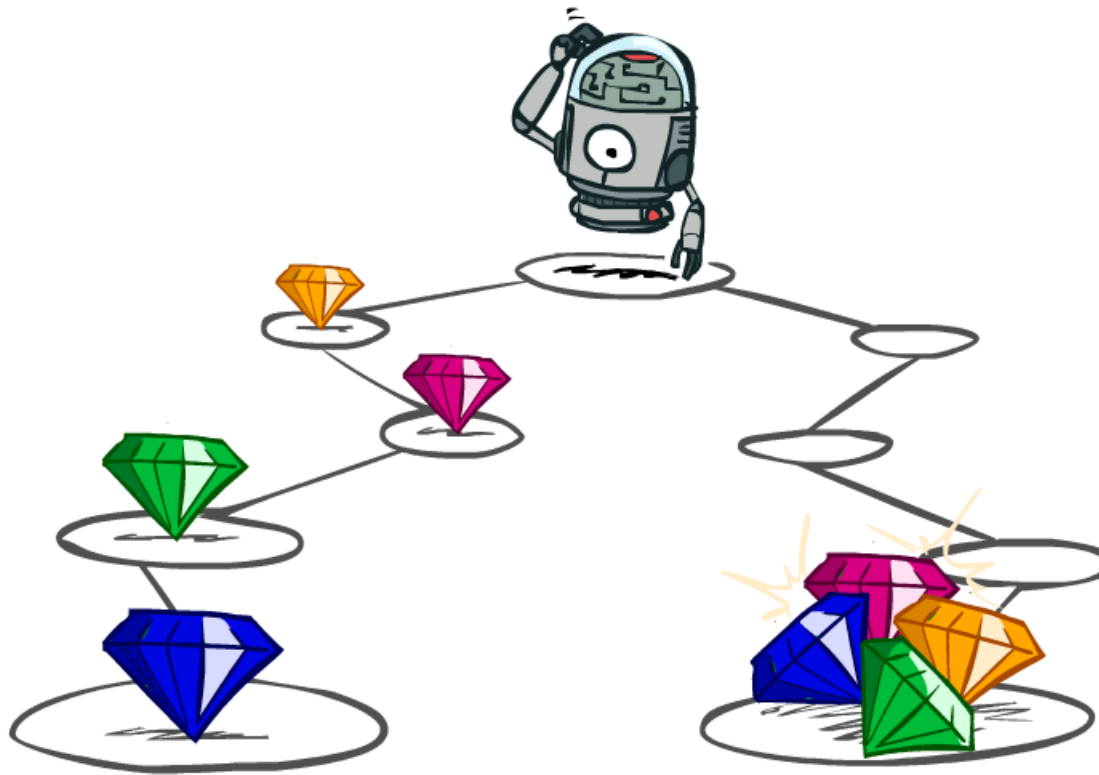Might be generated with ExpectiMax, but …?

# todo

- Add rewards into previous slide
- Next slide seems weirdly placed – totally unnecessary here

# MDP Search Trees

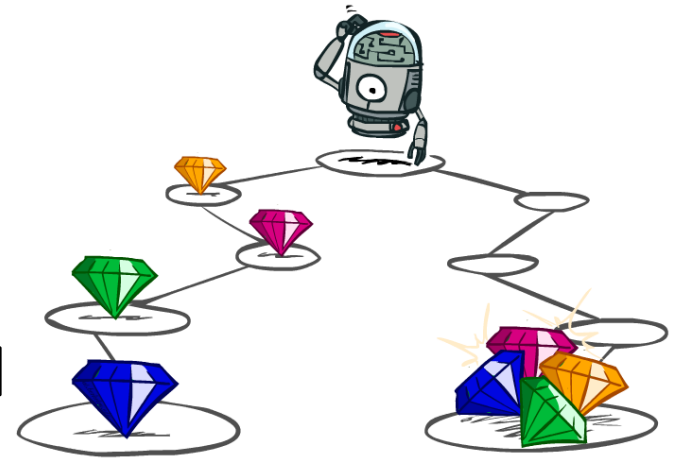- Each MDP state projects an expectimax-like search tree

$s$  ...a *state*

$a$

$s, a$

$s,a,s',r$

$s'$

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

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

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

# 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]

- Harder…             [1, 2, 3]    or    [3, 1, 1]

- Infinite sequences?  [1, 2, 1, …]  or    [2, 1, 2, …]

# 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
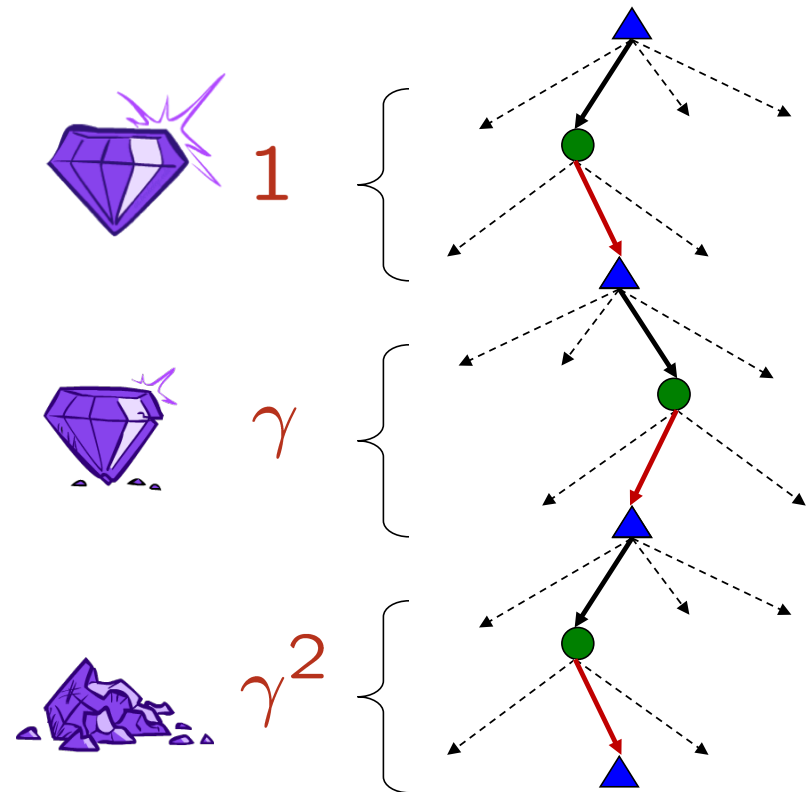
# Discounting

- **How to discount?**
  - Each time we descend a level, we multiply by the discount

- **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 = 2.75
  - U([3,1,1]) = 1*3 + 0.5*1 + 0.25*1 = 3.75
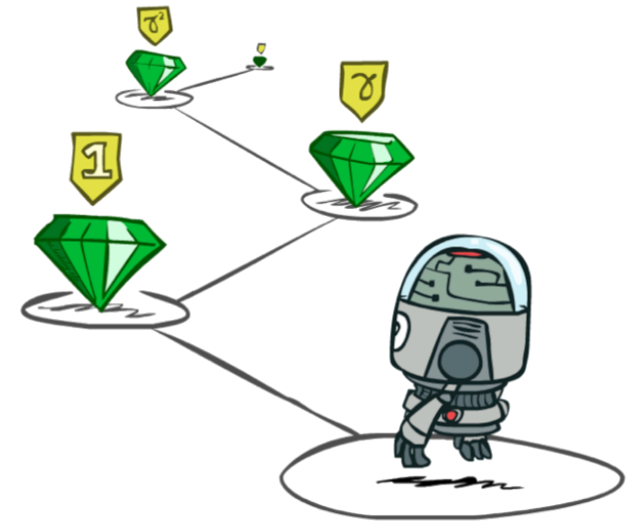  - U([1,2,3]) < U([3,1,1])

$1$

$\gamma$

$\gamma^2$

# Stationary Preferences

- Theorem: if we assume stationary preferences:

$$[a_1, a_2, \ldots] \succ [b_1, b_2, \ldots]$$

$$\updownarrow$$

$$[r, a_1, a_2, \ldots] \succ [r, b_1, b_2, \ldots]$$

- Then: there are **only two ways** to define utilities

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

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