# Advanced Topics in AI

## Markov Decision Processes



Instructor: Prof. Dr. techn. Wolfgang Nejdl

Leibniz University Hannover
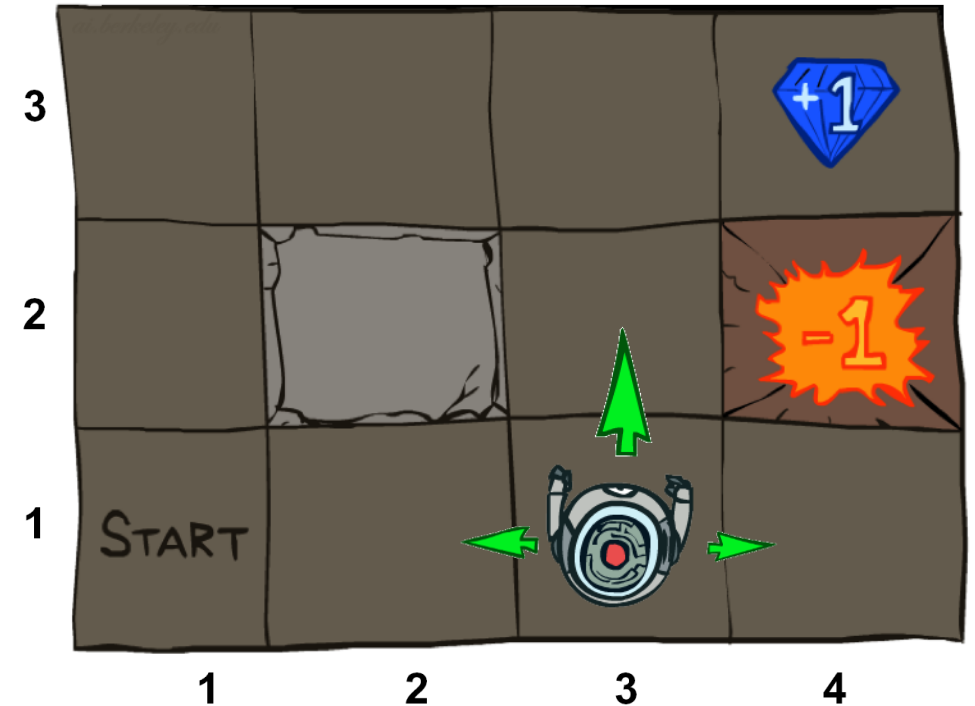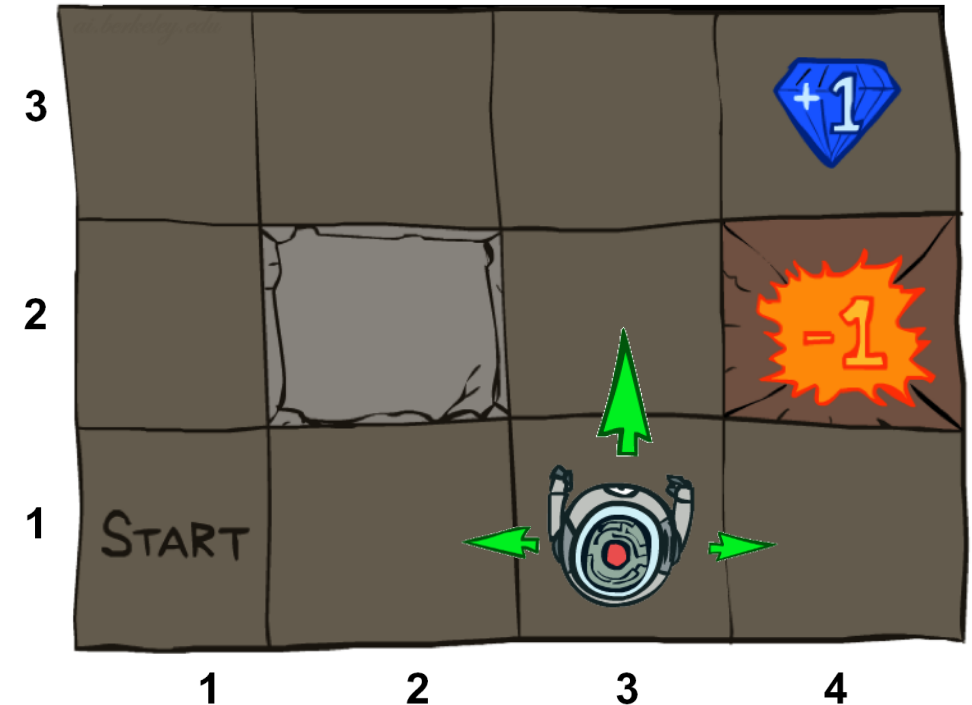
# Markov Decision Processes

- 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 $R(s, a, s')$
    - Sometimes just $R(s)$ or $R(s')$
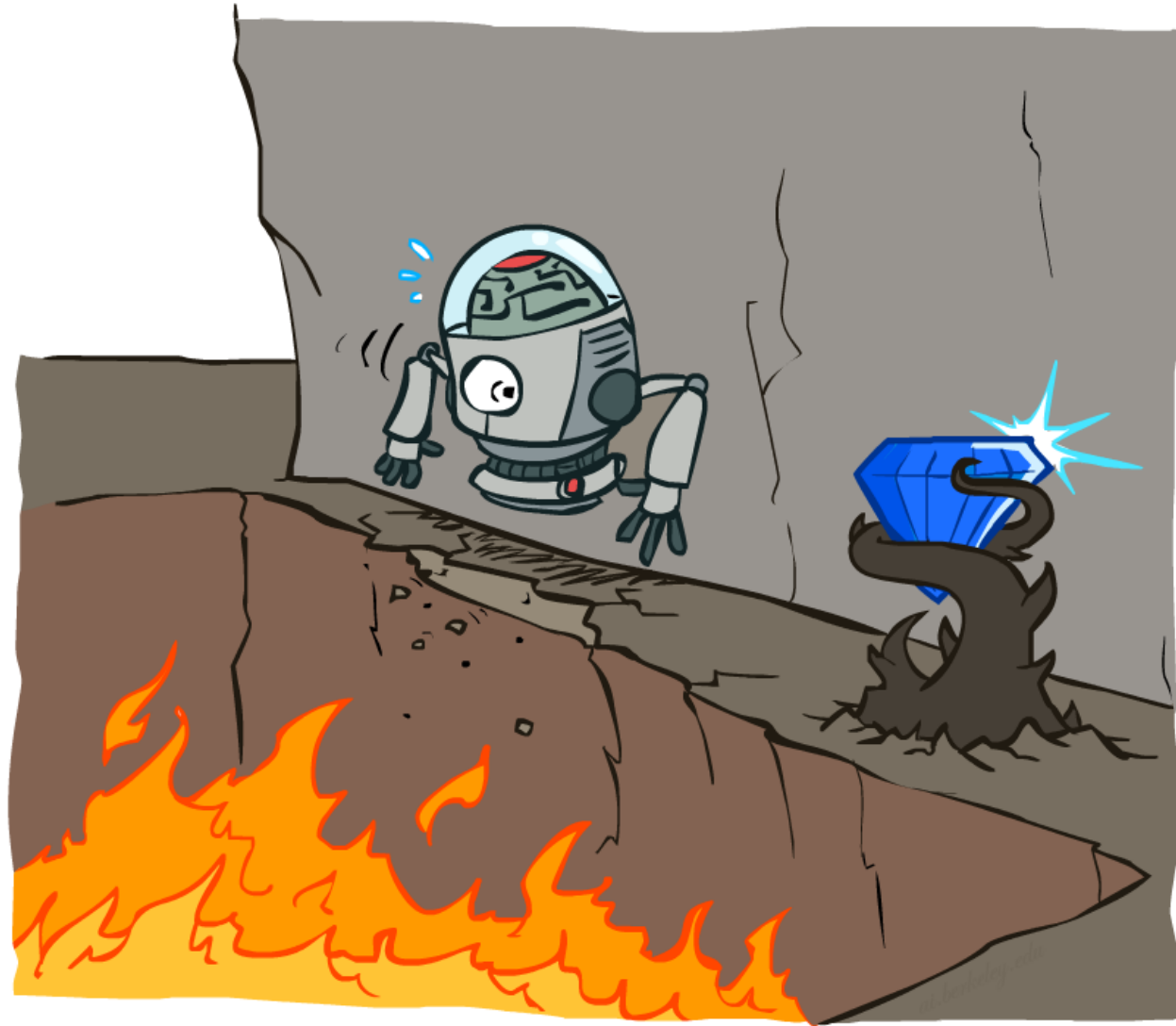  - A start state
  - Maybe a terminal state

# 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
  - Small "living" reward each step (can be negative)
  - Big rewards come at the end (good or bad)
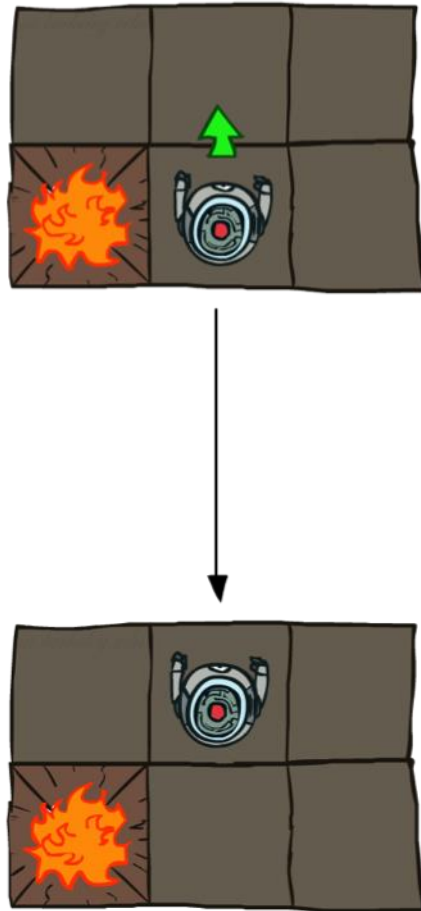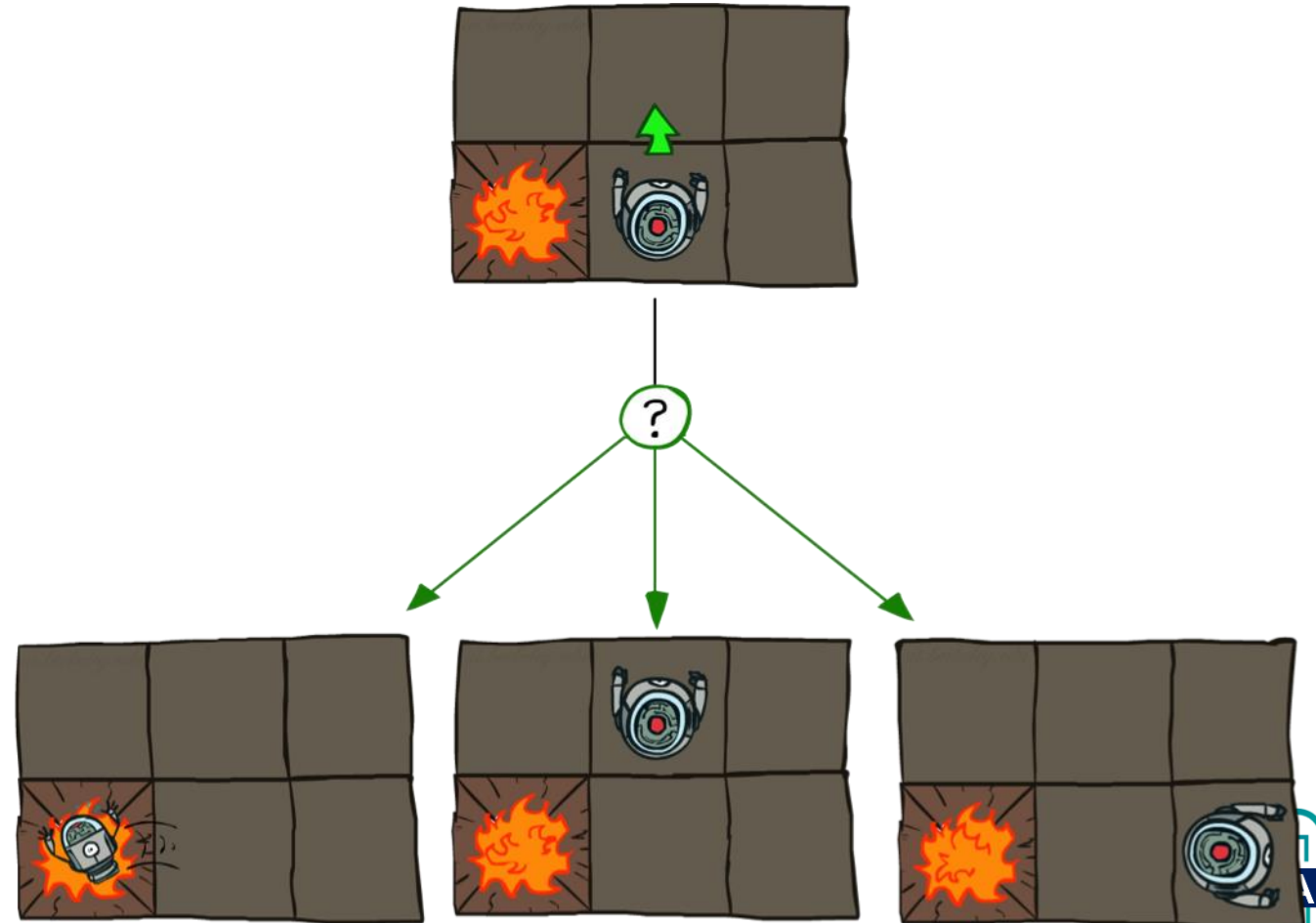- Goal: maximize sum of rewards

# Non-Deterministic Search

# Grid World Actions

## Deterministic Grid World



## Stochastic Grid World

# 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

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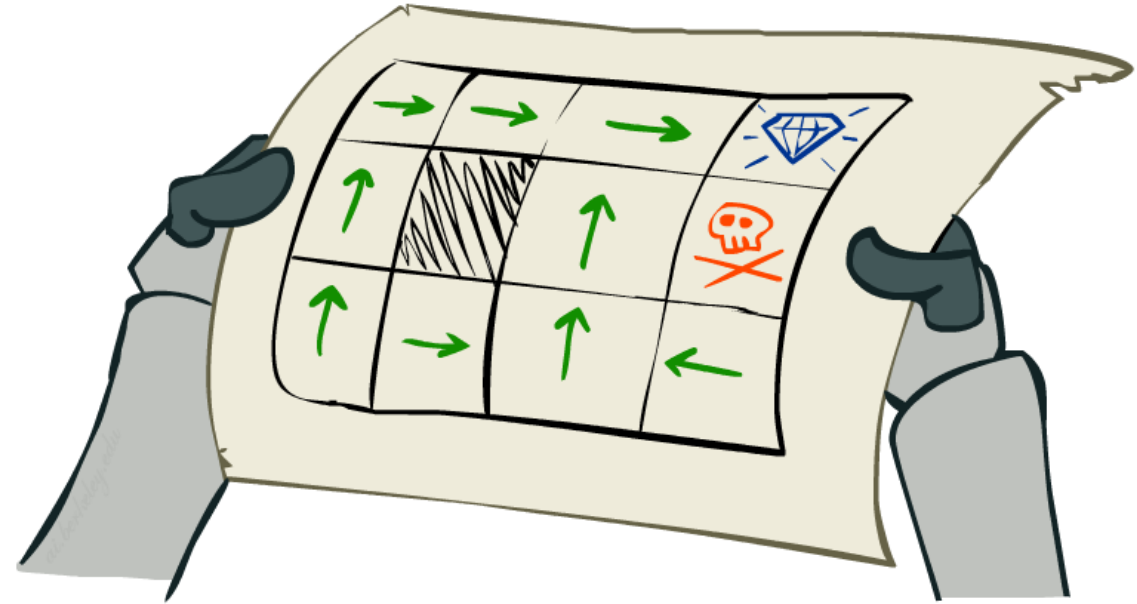
$$=$$

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

Andrey Markov
(1856-1922)

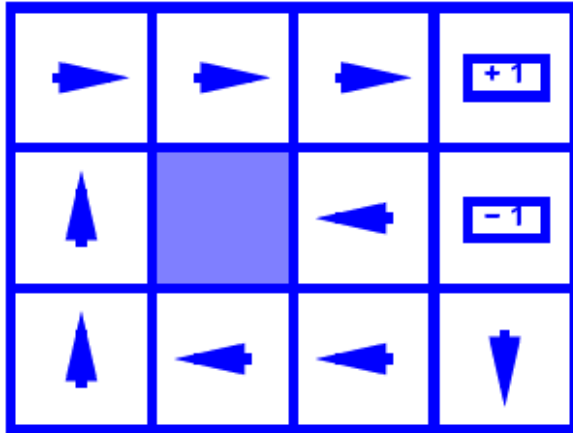- This is just like search, where the successor function could only depend on the current state (not the history)

# 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
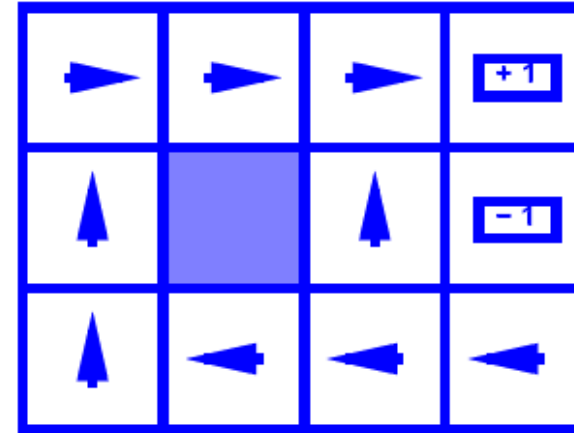    - An explicit policy defines a reflex agent



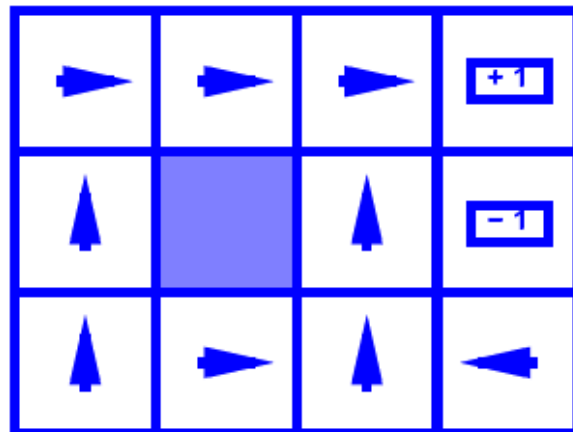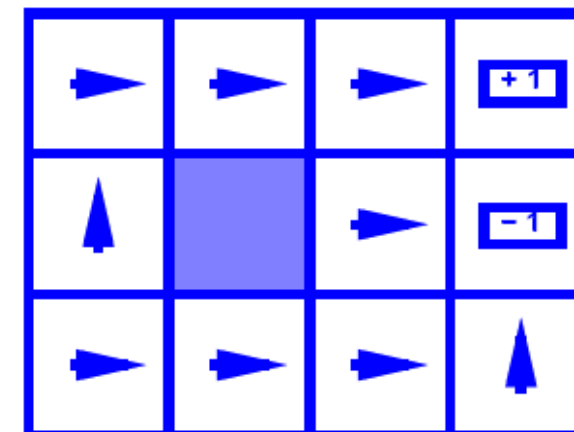Optimal policy when R(s, a, s') = -0.03 for all non-terminals s

# Optimal Policies



R(s) = -0.01



R(s) = -0.03



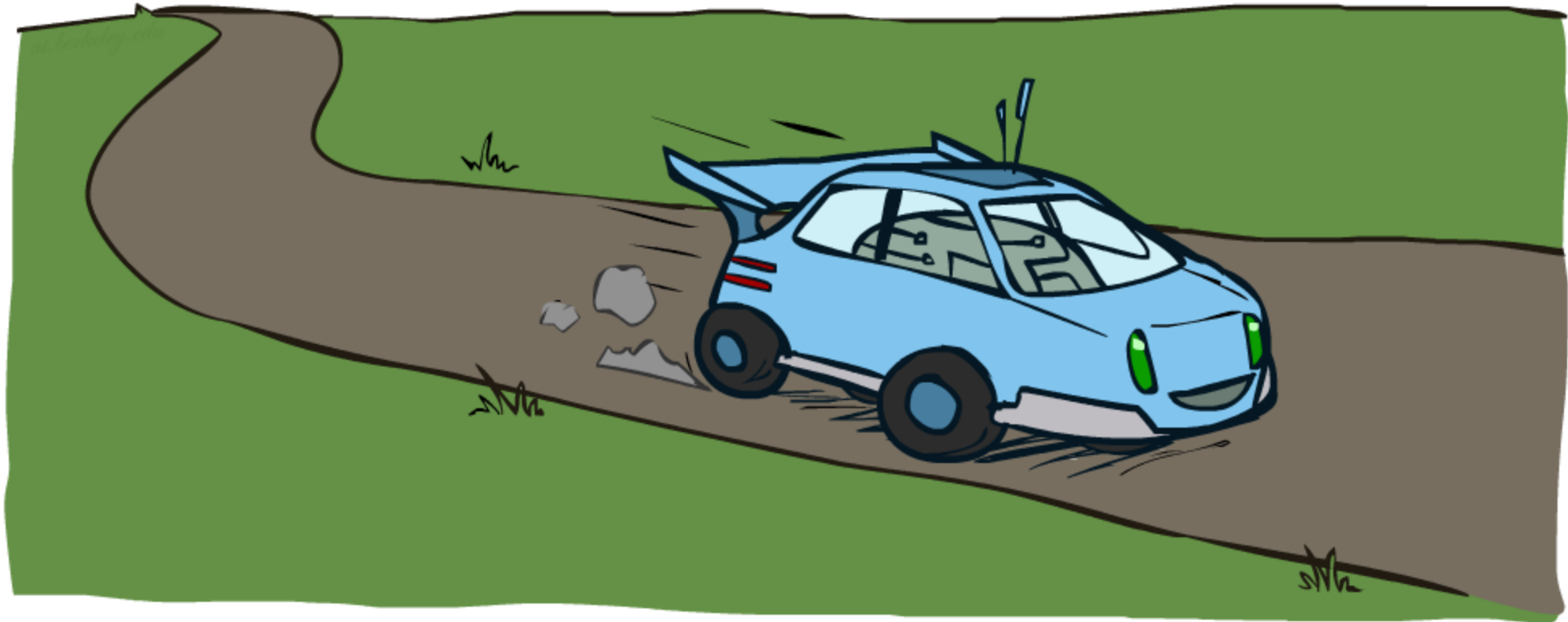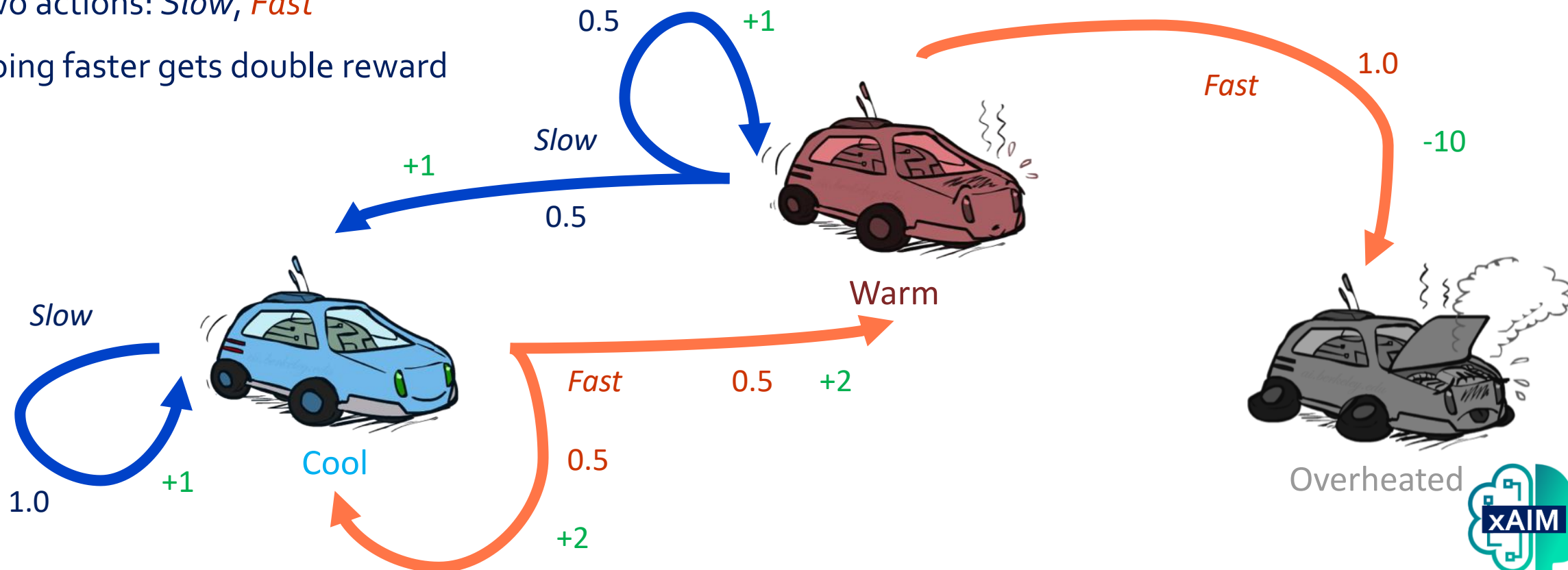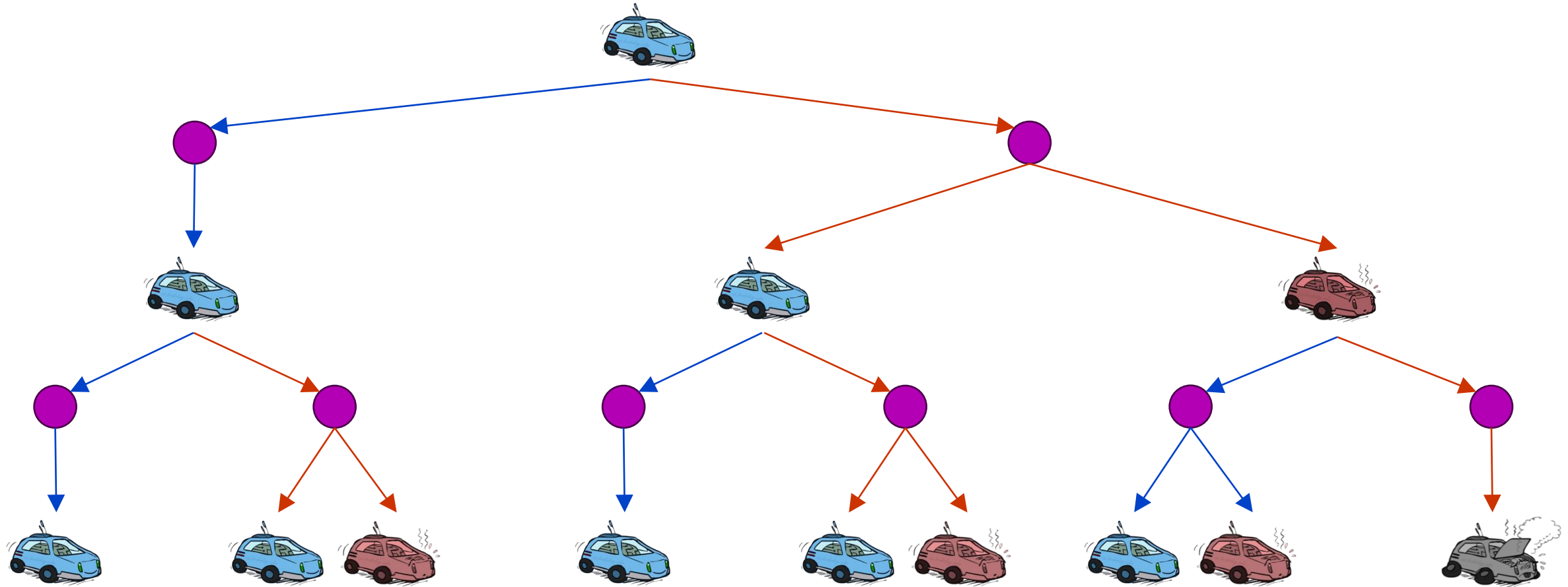R(s) = -0.4



R(s) = -2.0

# Example: Racing

# Example: Racing

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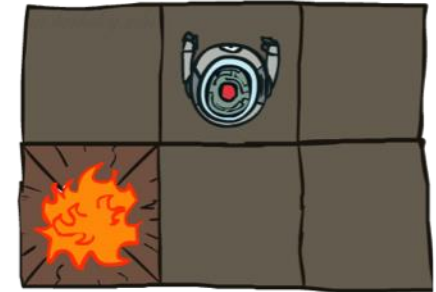- Two actions: *Slow*, *Fast*
- Going faster gets double reward

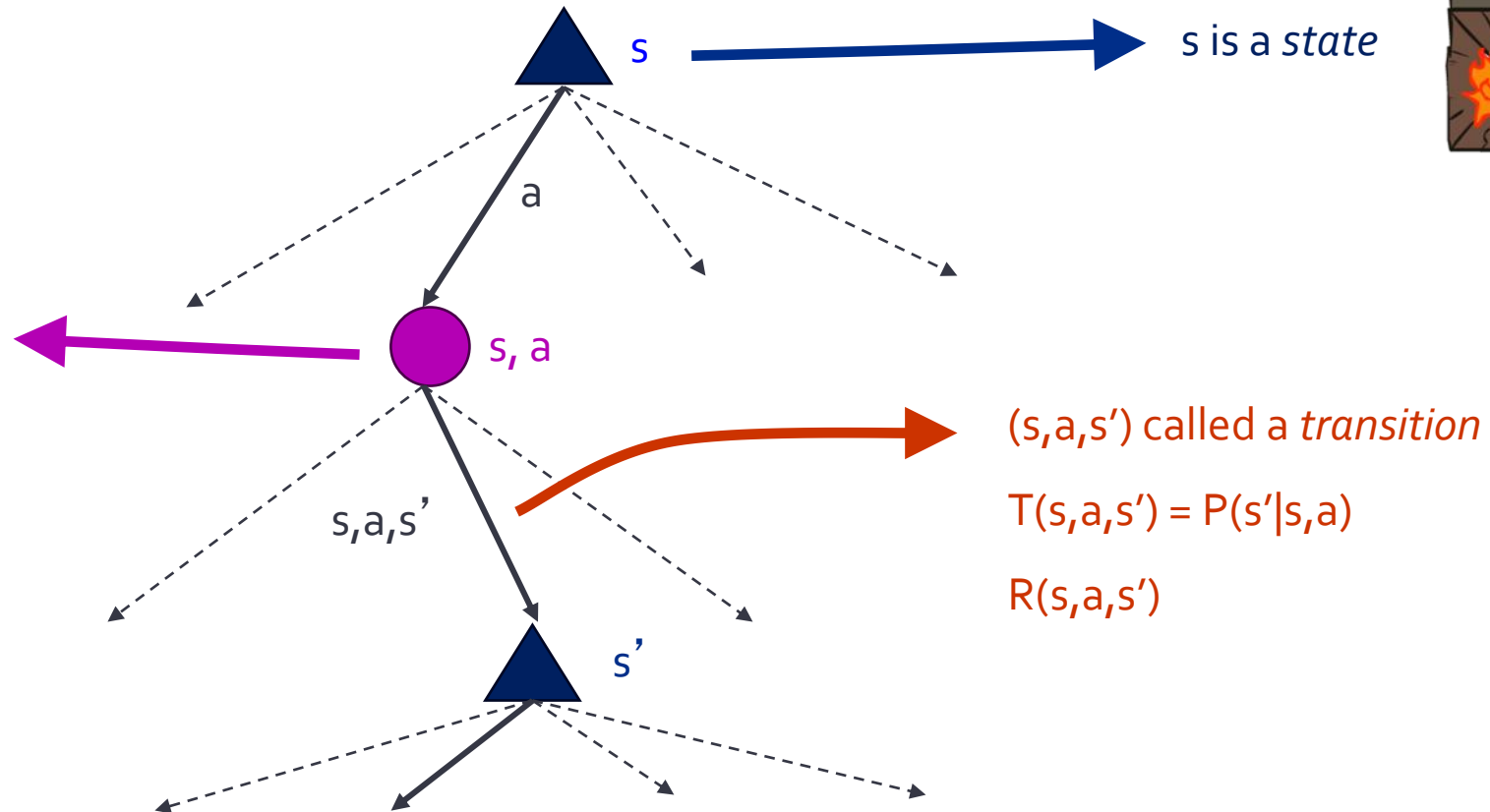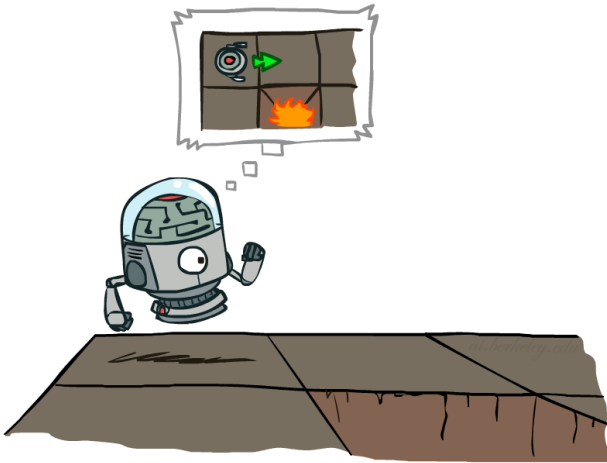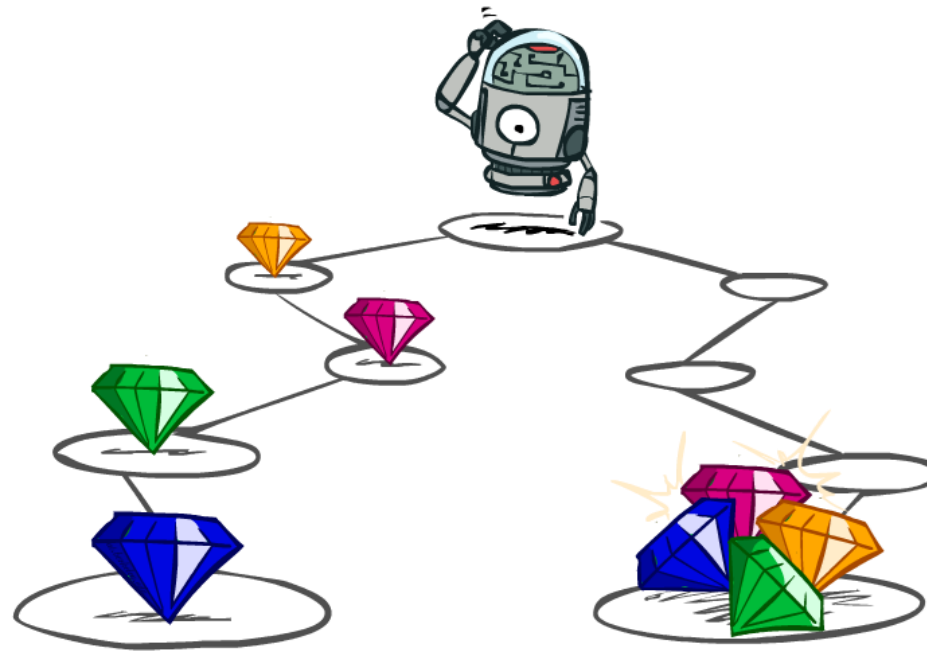# Racing Search Tree

# MDP Search Trees

- Each MDP state projects an expectimax-like search tree

s — 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'

# Advanced Topics in AI

## Next: Finite Horizons and Discounting



Instructor: Prof. Dr. techn. Wolfgang Nejdl

Leibniz University Hannover