# CS 3346/3121/9146: Artificial Intelligence I

## Markov Decision Processes
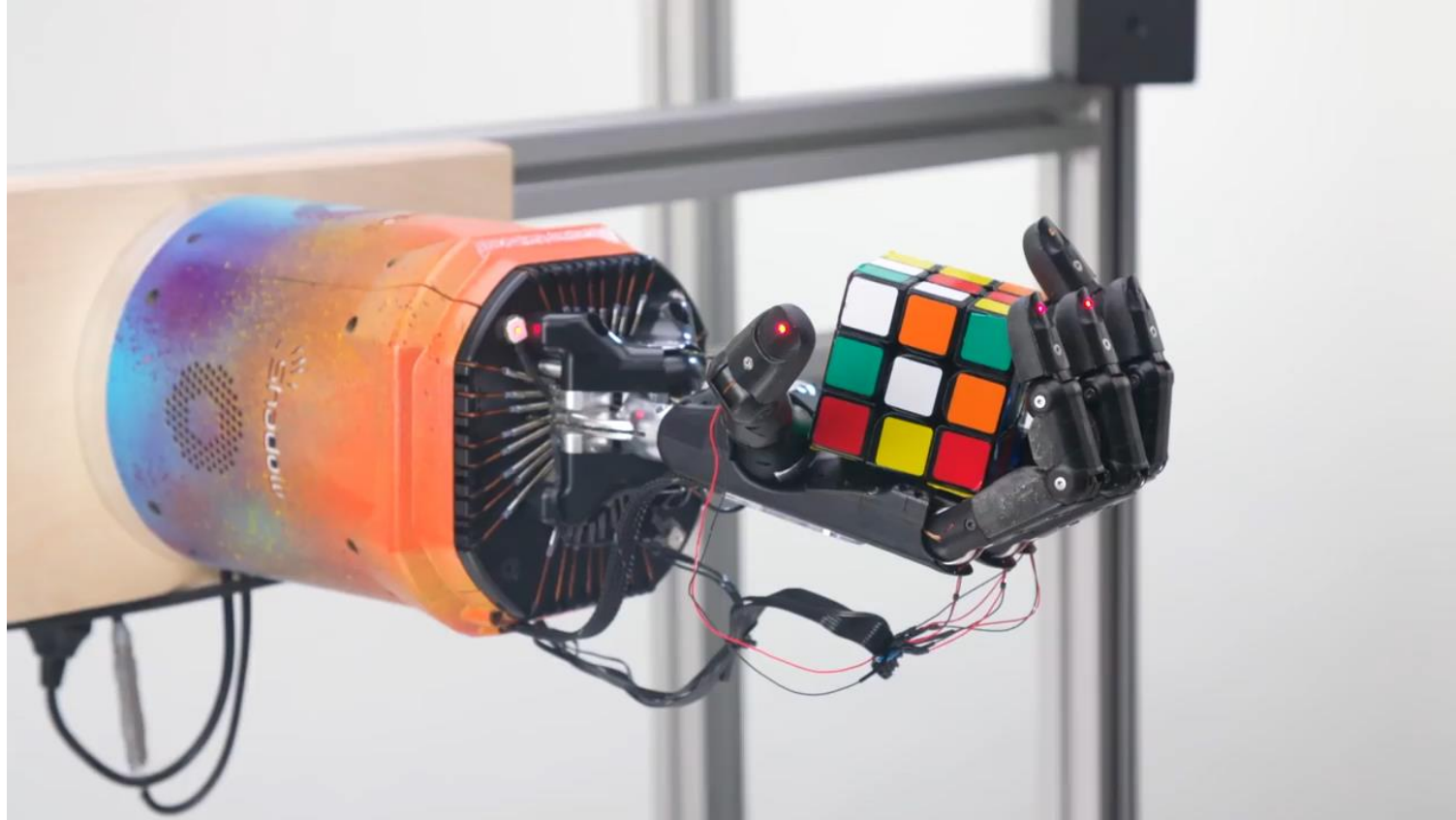


Instructor: Yalda Mohsenzadeh, Western University
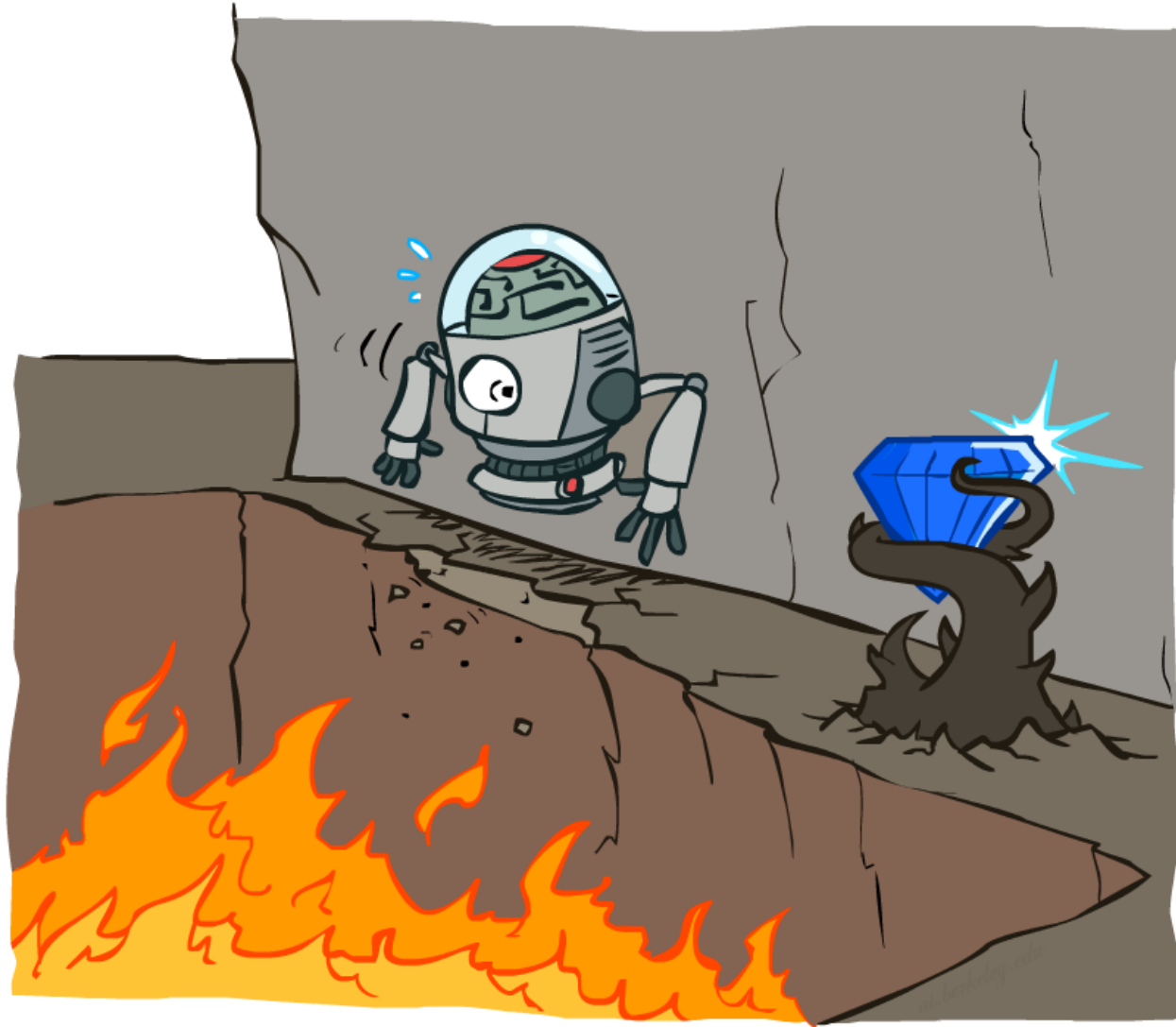
# Deep Reinforcement Learning



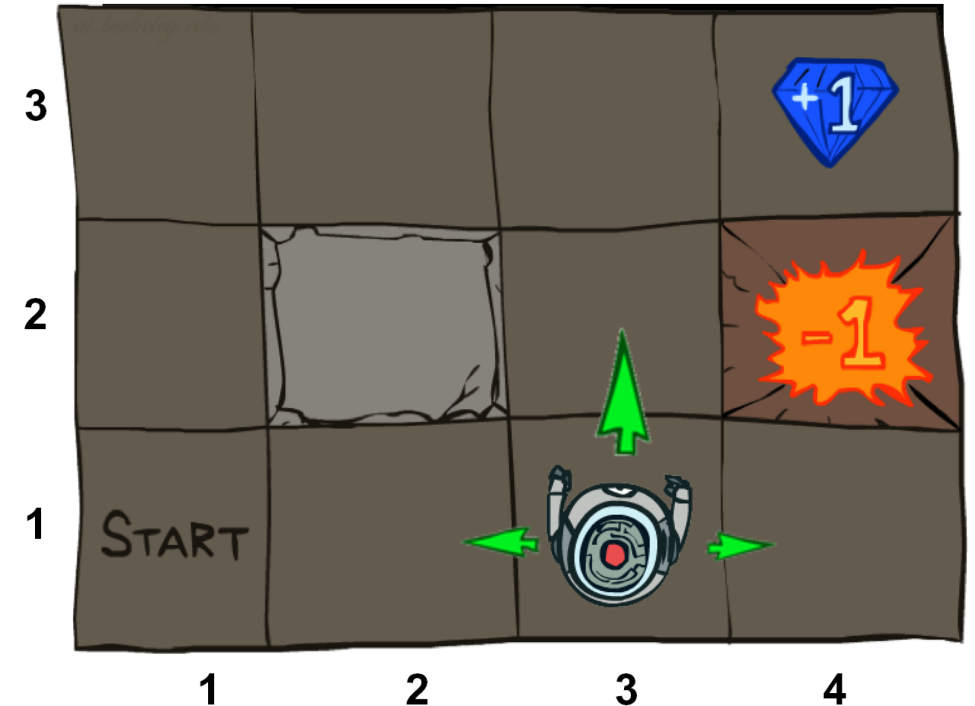**2019**    **Rubik's Cube (PPO+DR)**
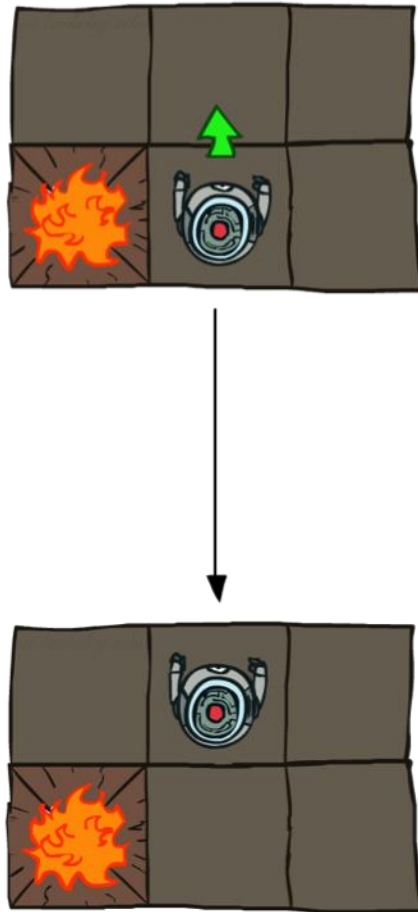      **[OpenAI]**

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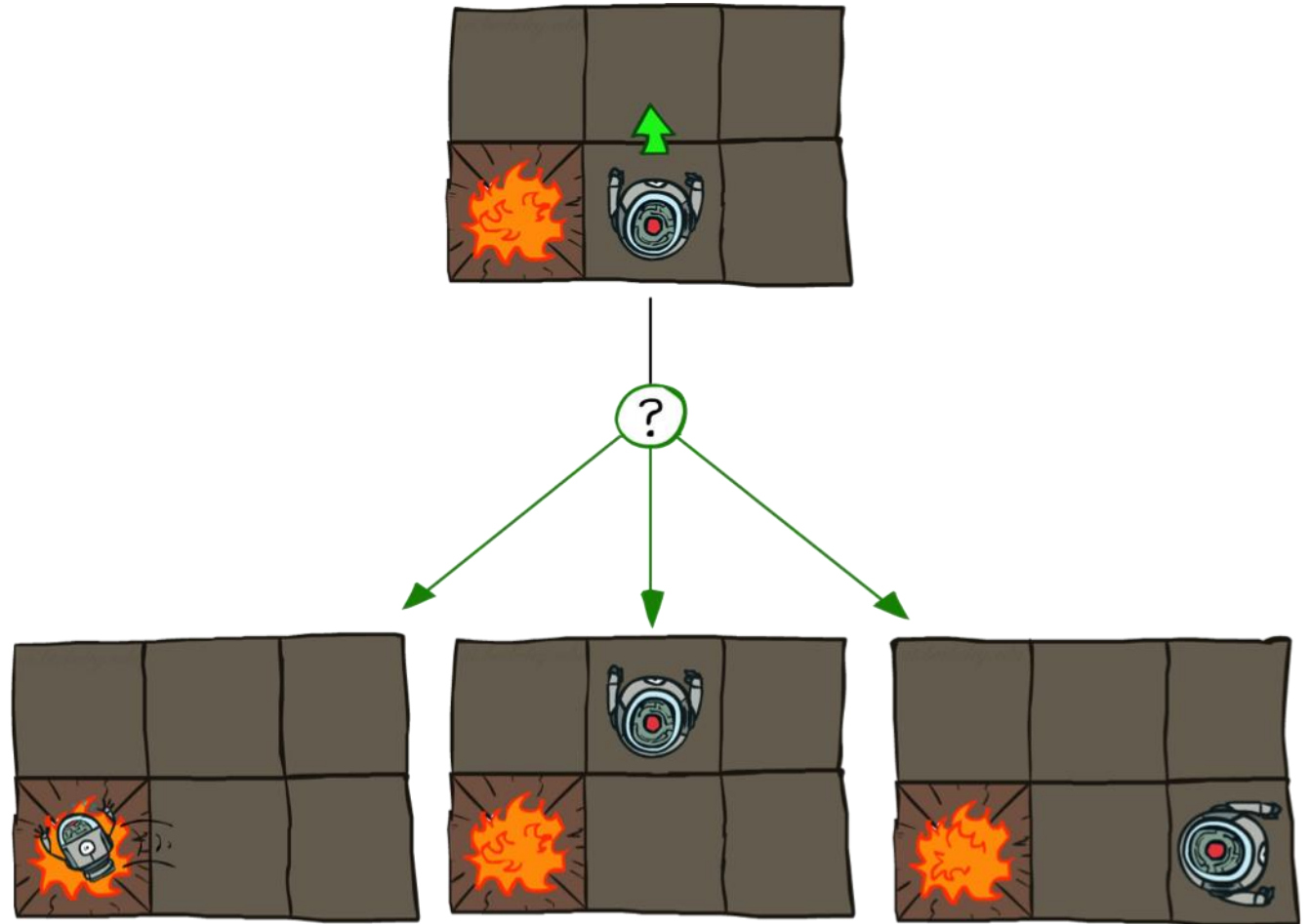  - Big rewards come at the end (good or bad)

- Goal: maximize sum of rewards

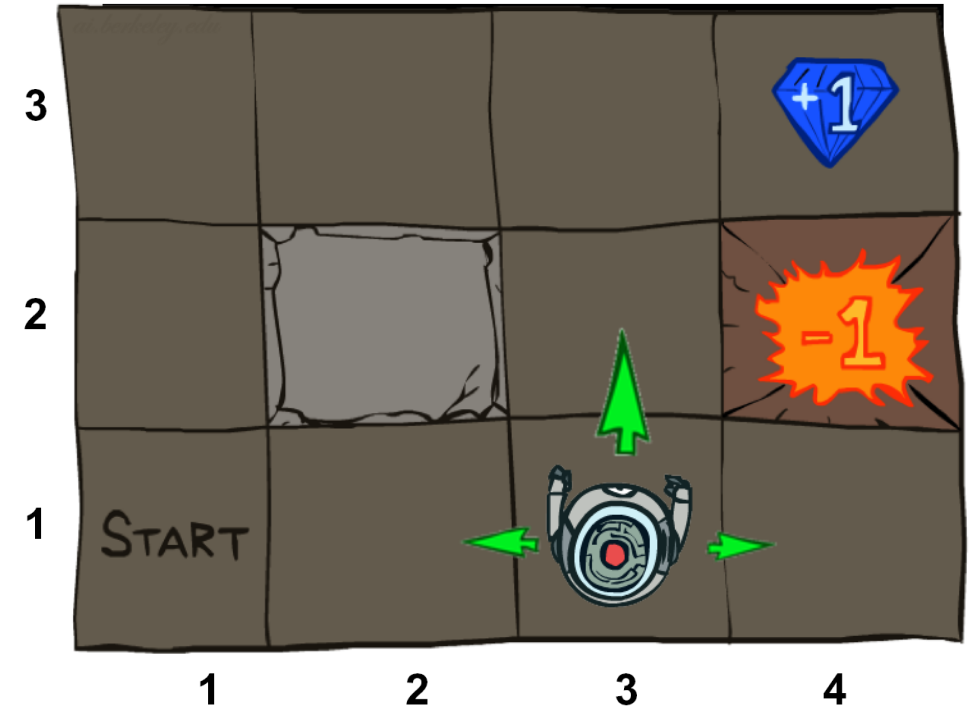# Grid World Actions

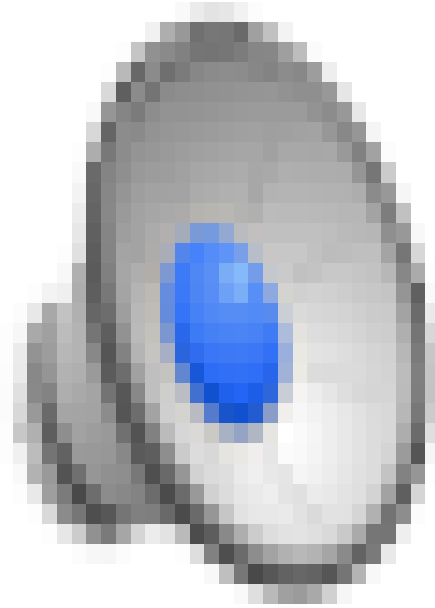Deterministic Grid World

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

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

# Video of Demo Gridworld Manual Intro

# What is Markov about MDPs?

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

o 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)

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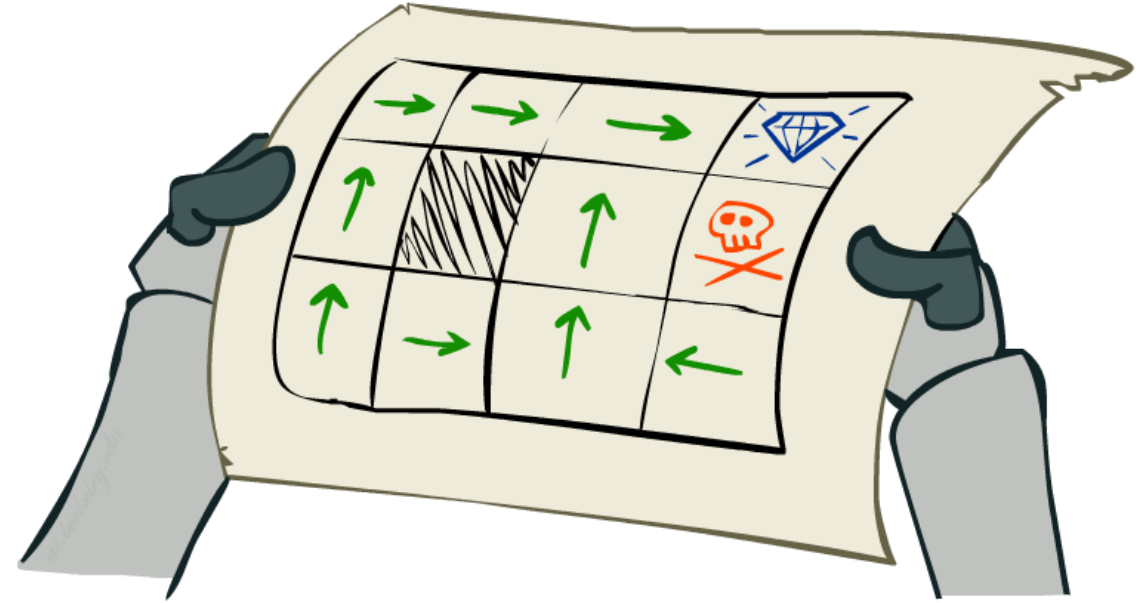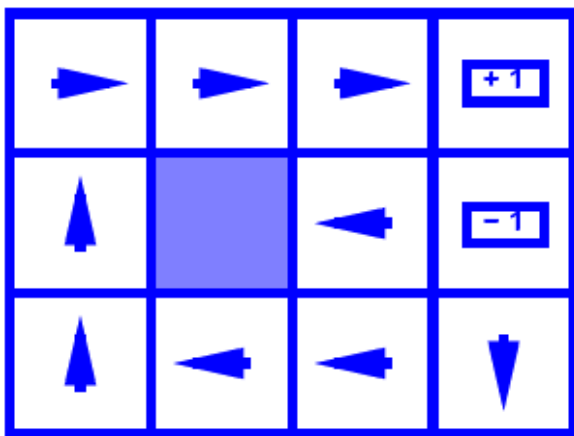
# Policies

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
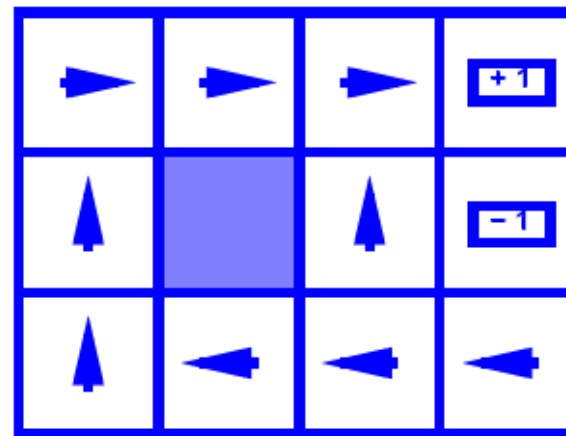   o An explicit policy defines a reflex agent



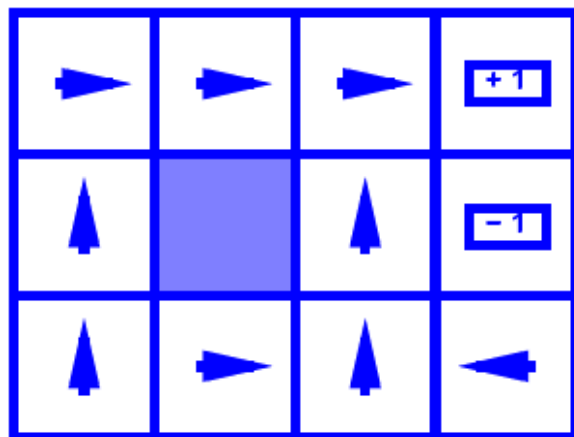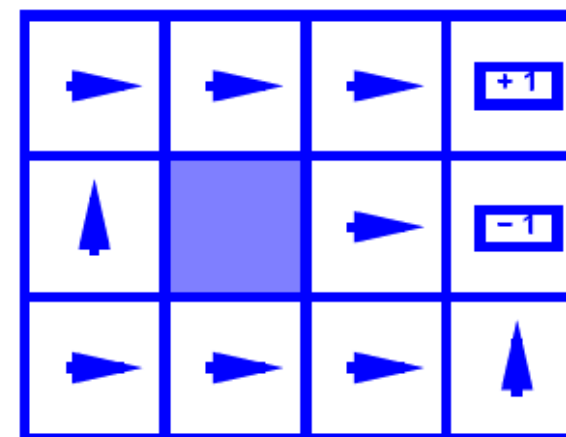Optimal policy when R(s, a, s') = -0.03 for all non-terminals s
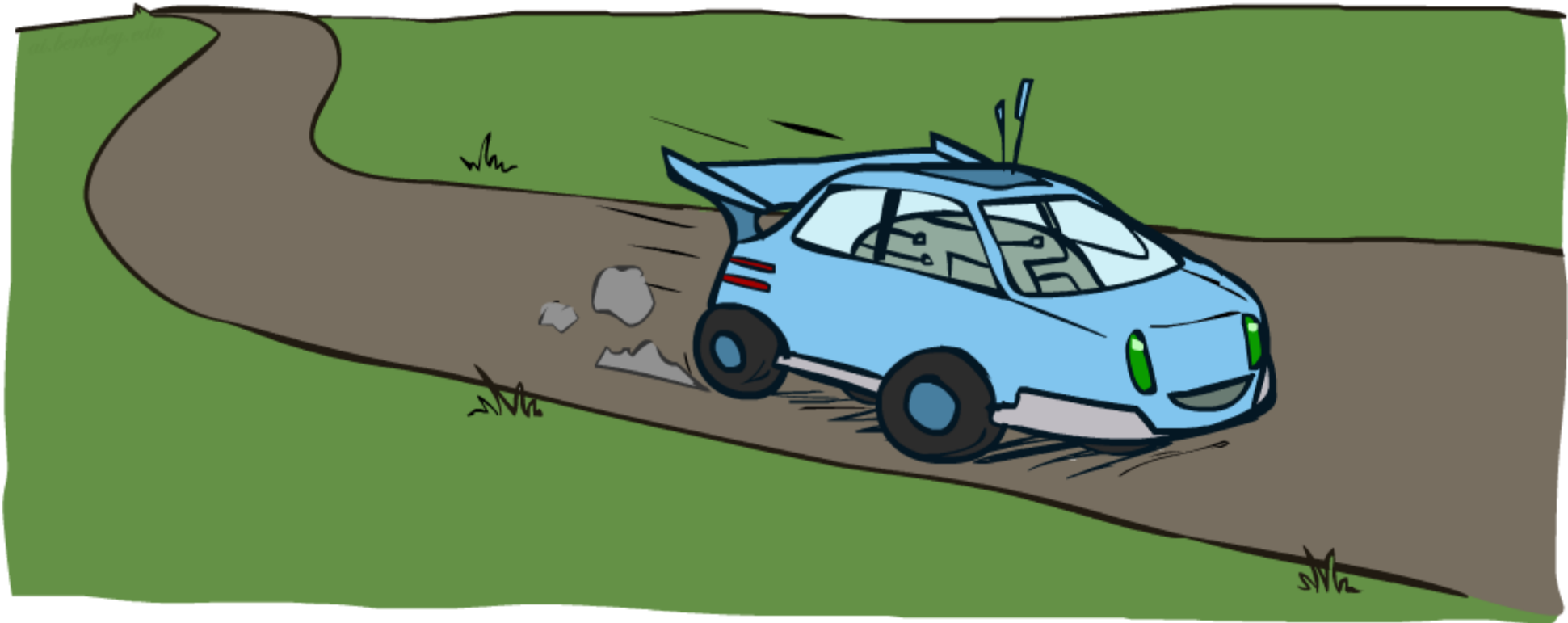
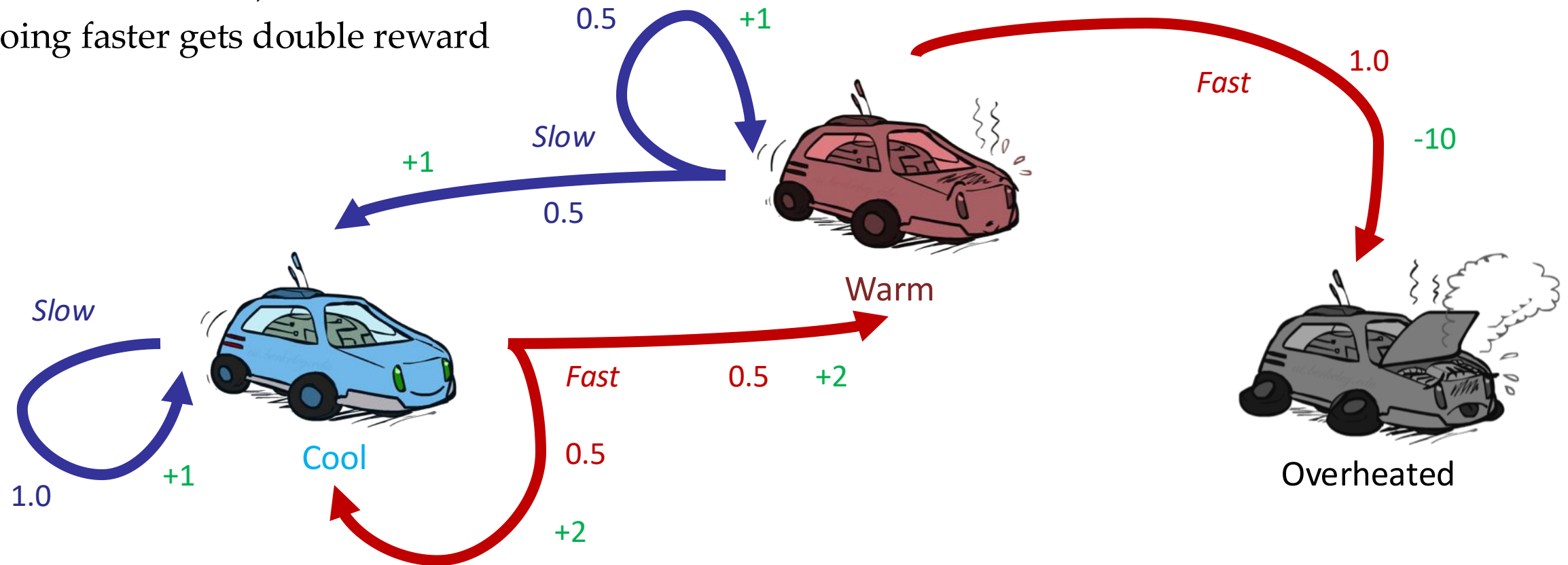# Optimal Policies



R(s) = -0.01

R(s) = -0.03

R(s) = -0.4

R(s) = -2.0

# Example: Racing
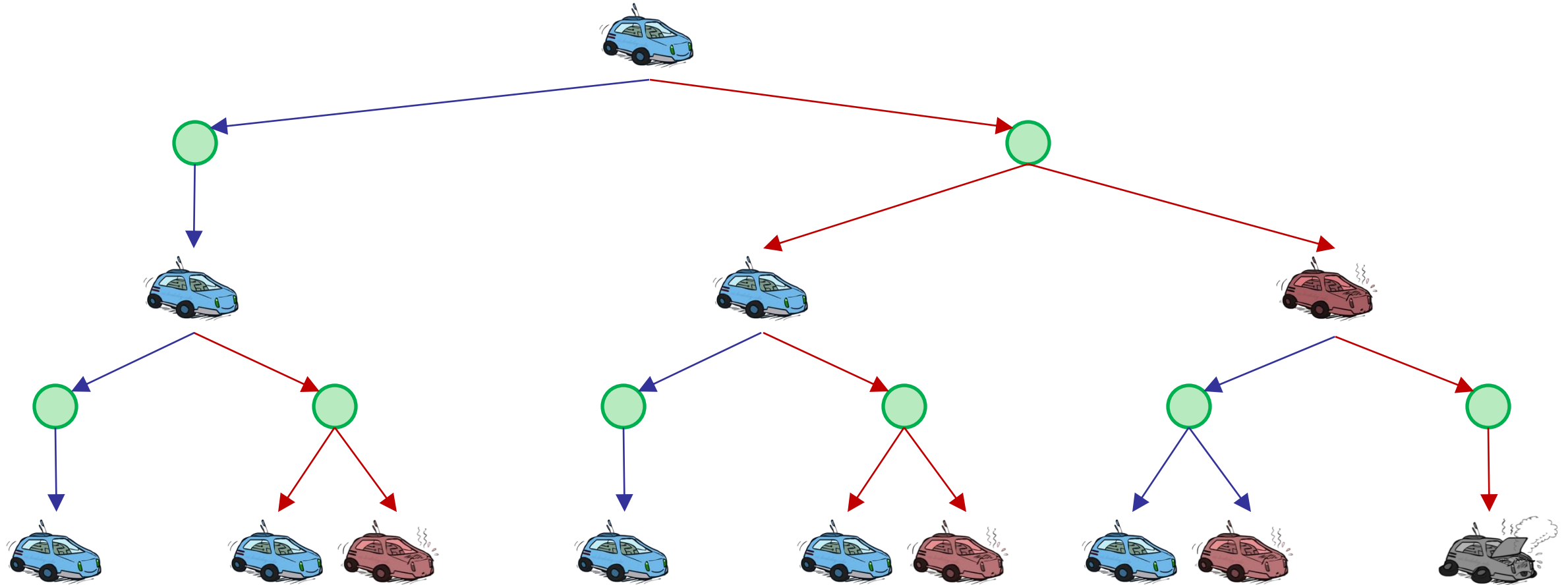
# Example: Racing

o A robot car wants to travel far, quickly
o Three states: Cool, Warm, Overheated
o Two actions: *Slow*, *Fast*
o 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

○ 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

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

o $\gamma \in (0,1)$



$1$

Worth Now

$\gamma$

Worth Next Step

$\gamma^2$

Worth In Two Steps

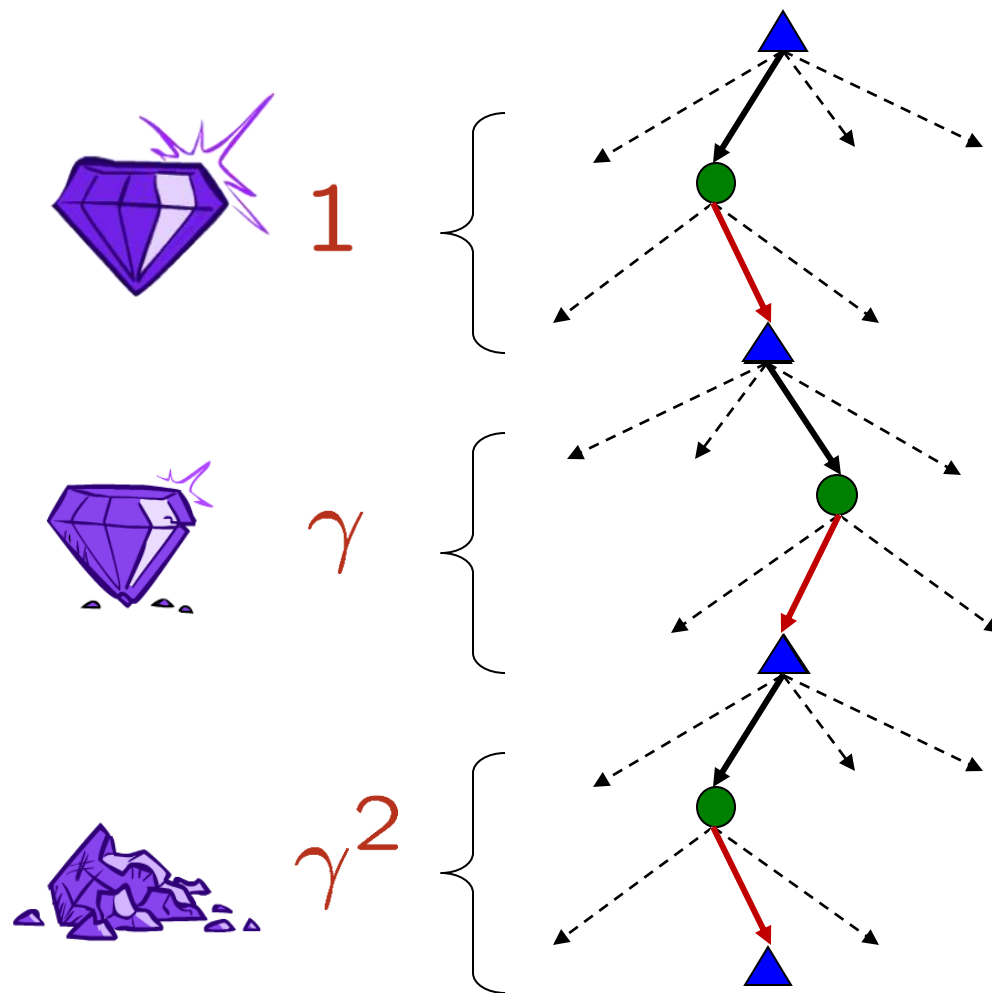# Discounting

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

- Why discount?
  - Reward now is better than later
  - Can also think of it as a 1-gamma chance of ending the process at every step
  - Also helps our algorithms converge

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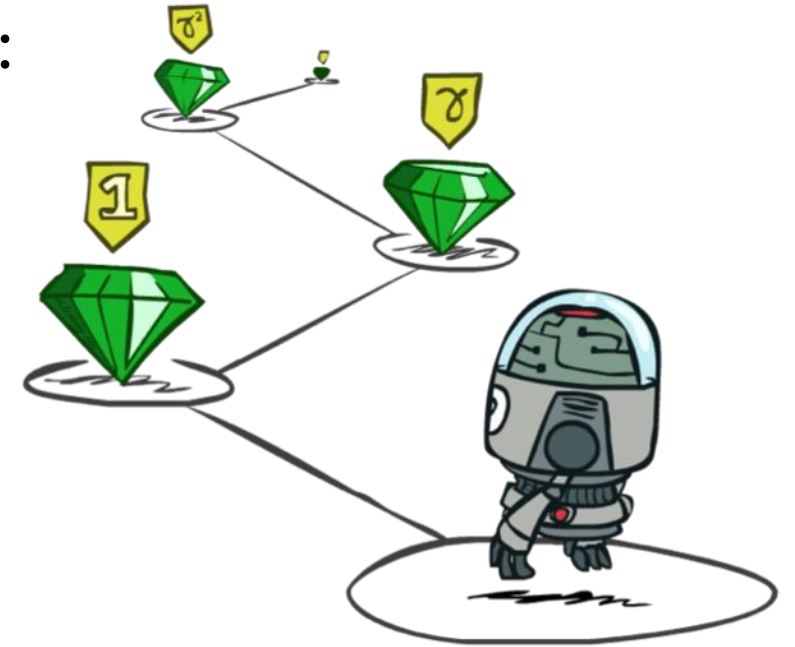  - U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
  - U([1,2,3]) < U([3,2,1])

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 is only ways to define utilities

  - Additive discounted utility:

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

# 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?
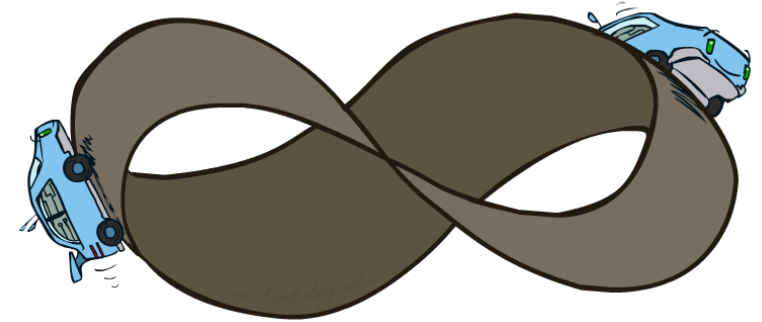
$1\gamma = 10 \ \gamma^3$

# 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)
    - Gives nonstationary policies ($\pi$ depends on time left)
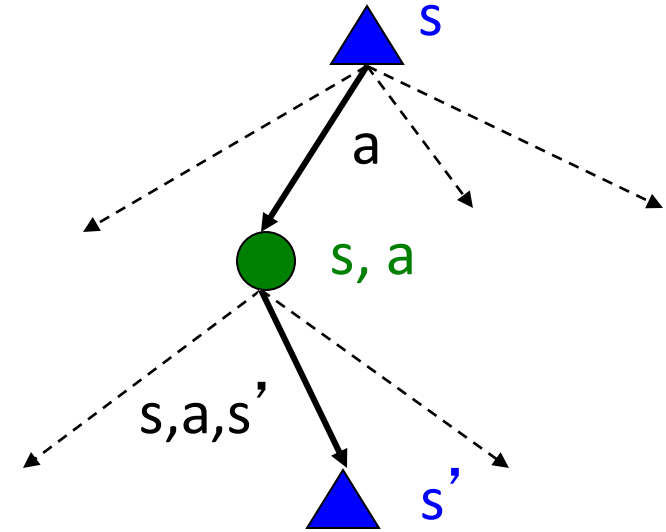
  - Discounting: use $0 < \gamma < 1$
    $$U([r_0, \ldots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\mathsf{max}}/(1 - \gamma)$$

    - 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)
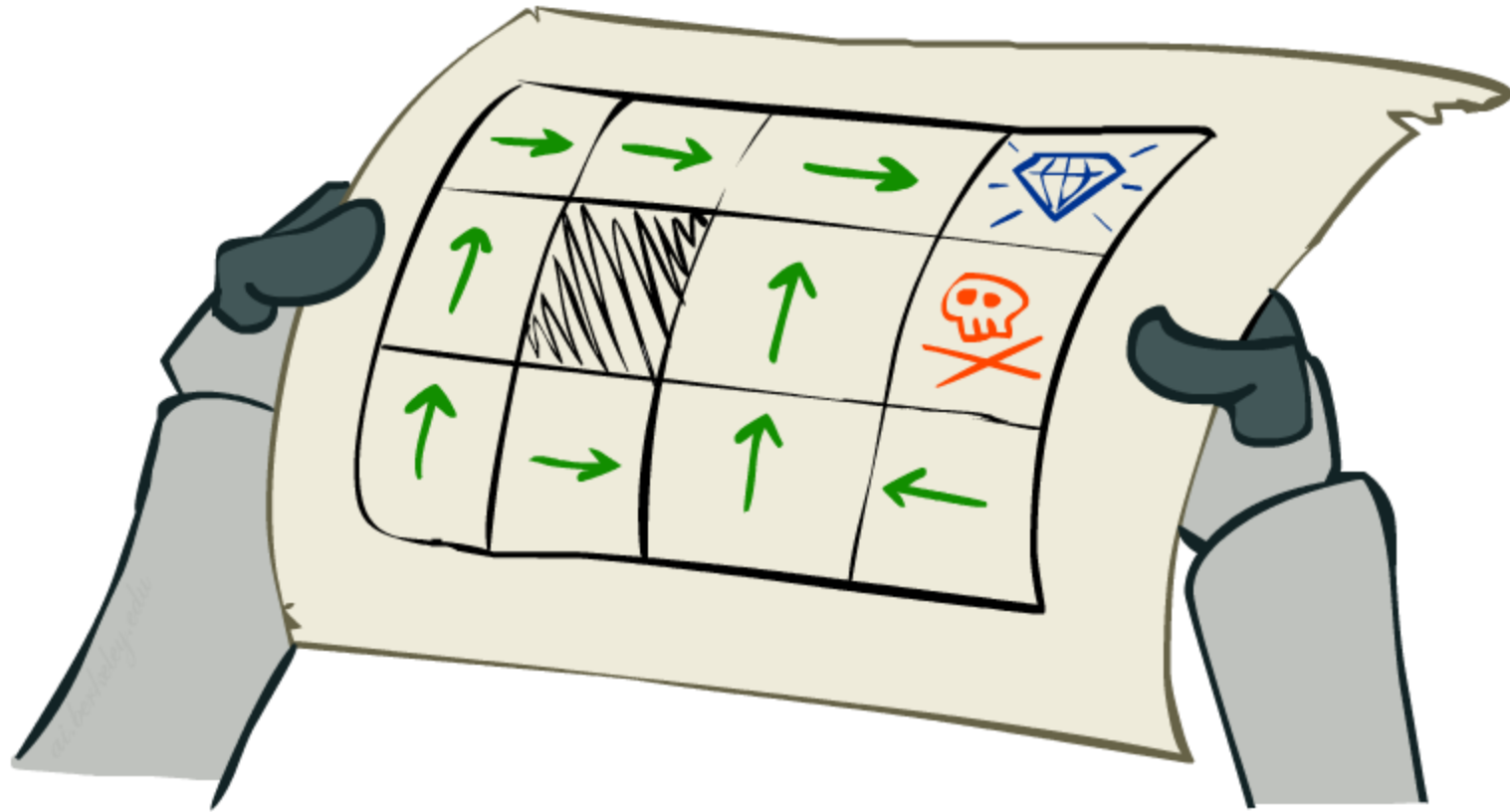
# Recap: Defining MDPs

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
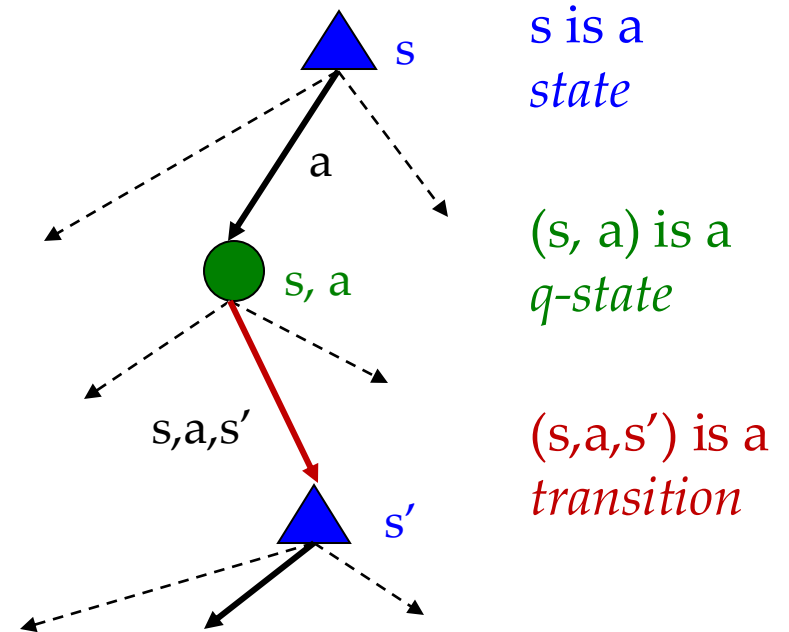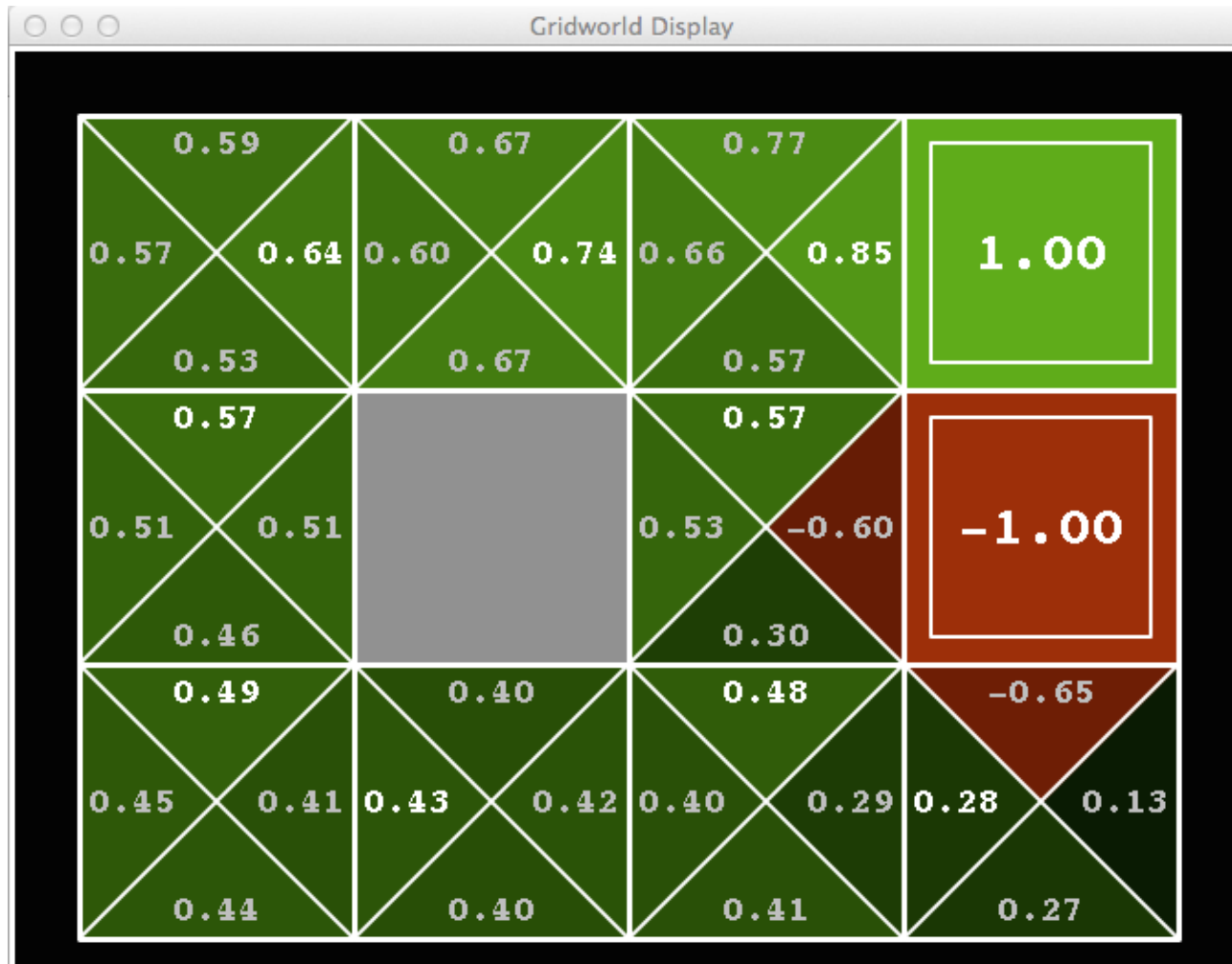
# Solving MDPs

# 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*

s

a

s, a

s,a,s'

s'

# Gridworld V* Values



Noise = 0.2
Discount = 0.9
Living reward = 0

# Gridworld Q* Values



Noise = 0.2
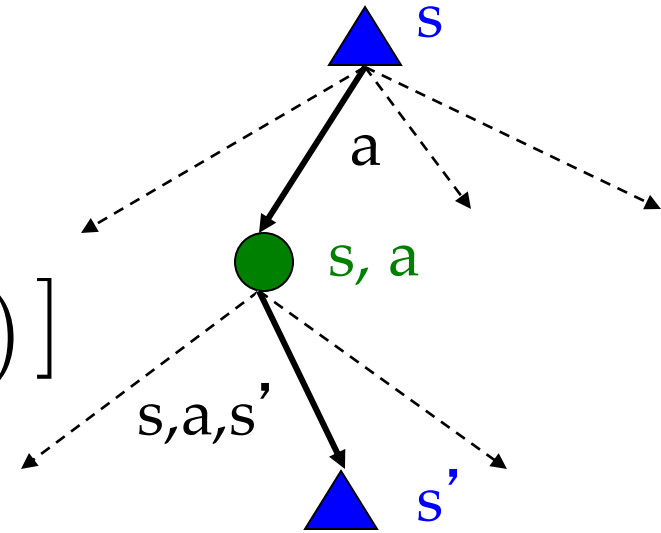Discount = 0.9
Living reward = 0

# Values of States: Bellman Equation

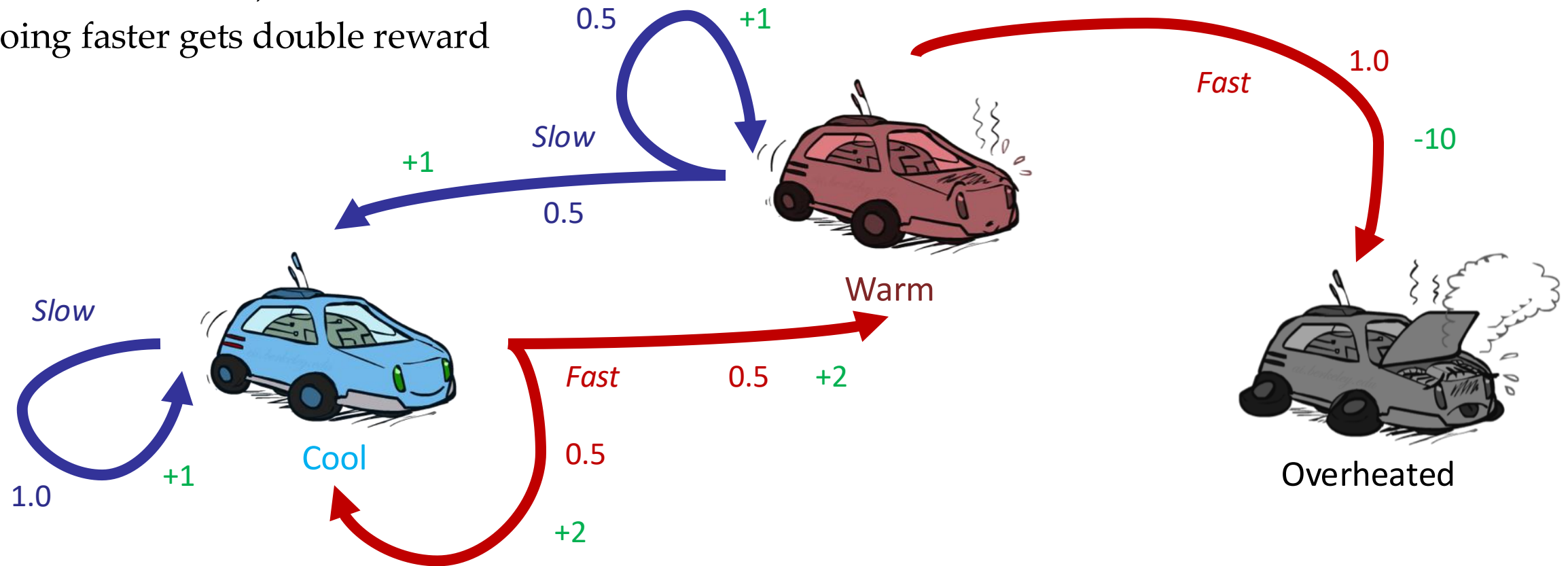○ Recursive definition of value:

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

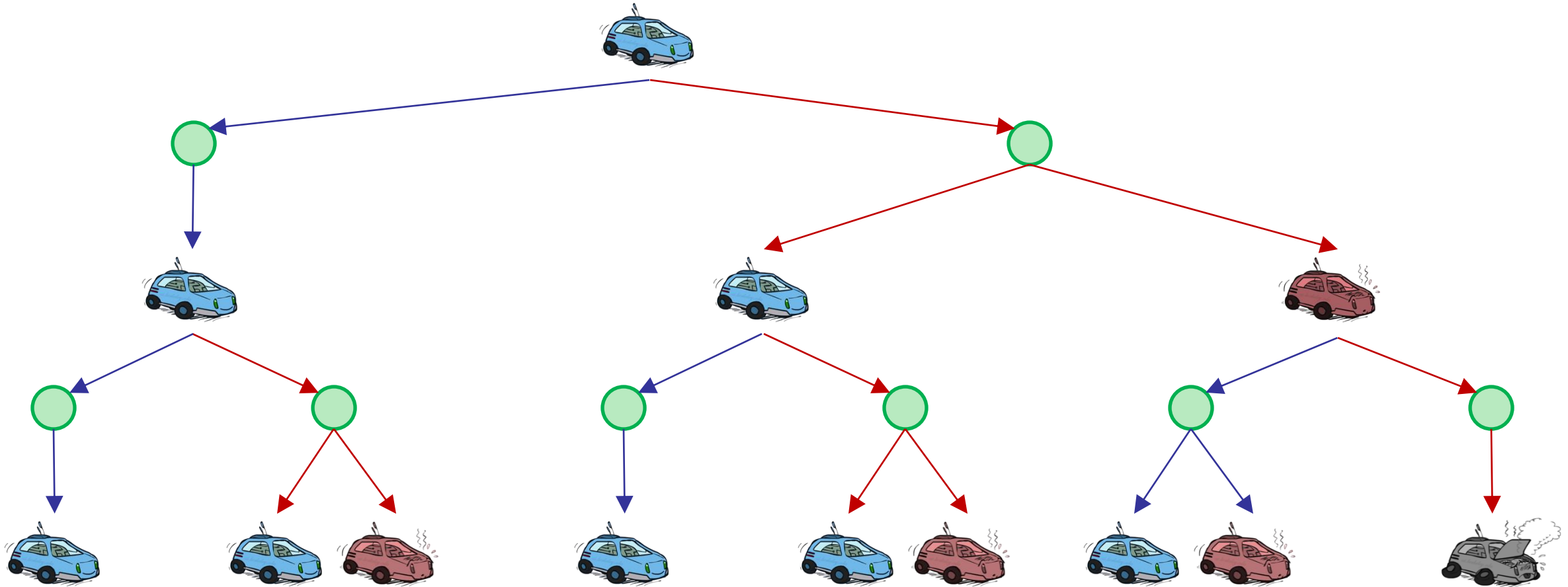$$Q^*(s, a) = \sum_{s'} T(s, a, s')\left[R(s, a, s') + \gamma V^*(s')\right]$$

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

s

a

s, a

s,a,s'

s'

# Recall: Racing MDP

o   A robot car wants to travel far, quickly

o   Three states: Cool, Warm, Overheated

o   Two actions: *Slow*, *Fast*

o   Going faster gets double reward

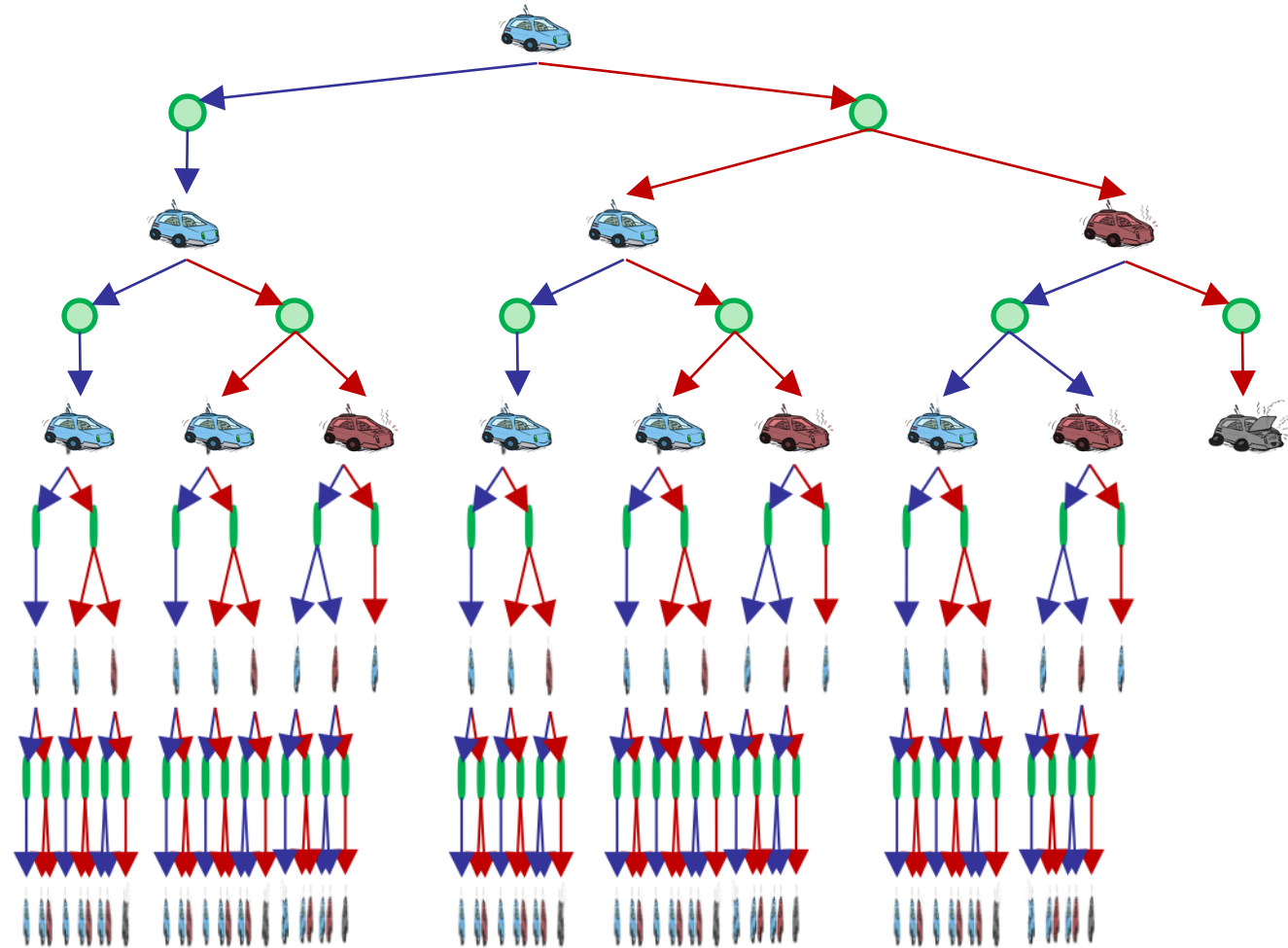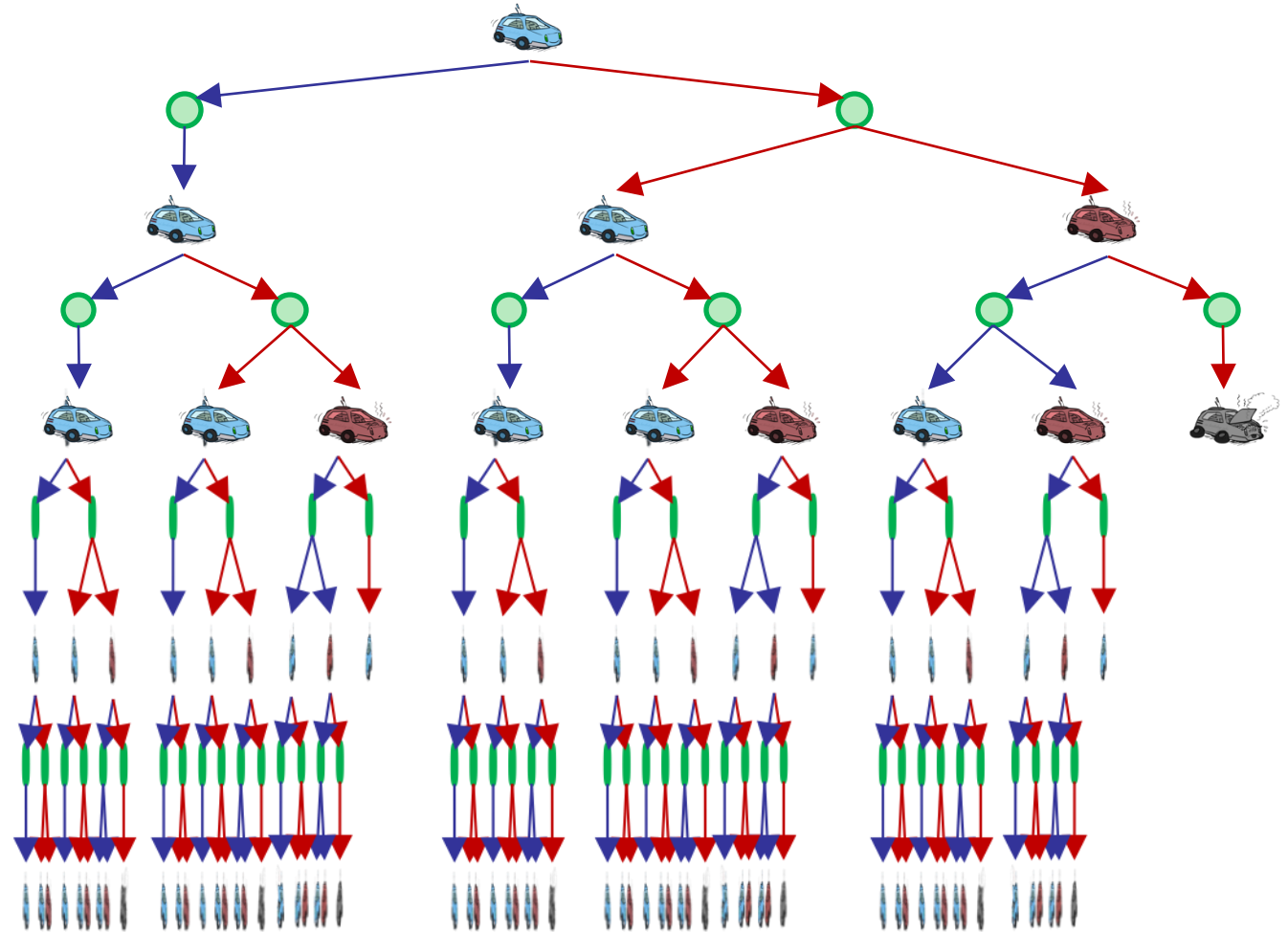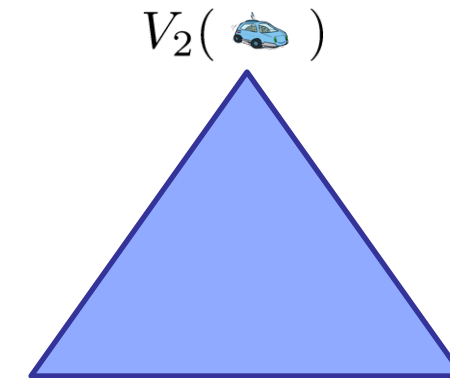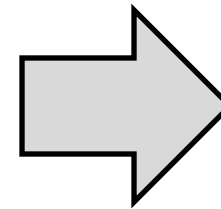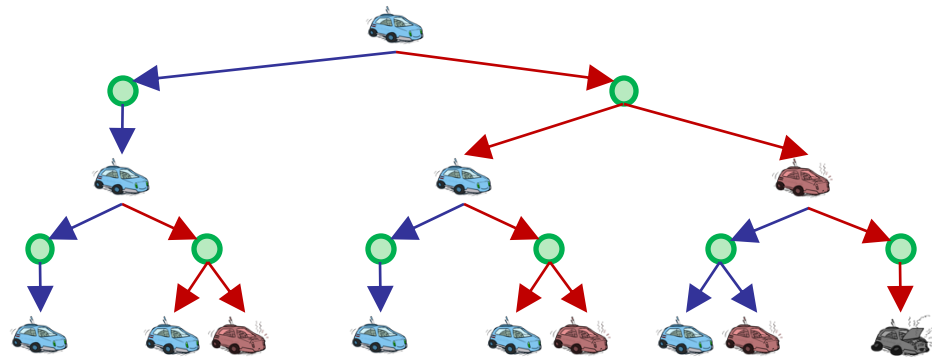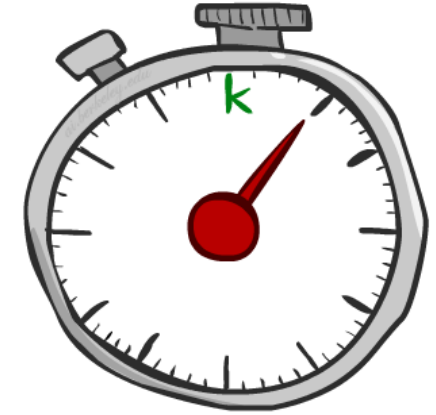# Racing Search Tree

# Racing Search Tree

# Racing Search Tree

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

o Problem: States are repeated
  o Idea: Only compute needed quantities once

o Problem: Tree goes on forever
  o Idea: Do a depth-limited computation, but with increasing depths until change is small
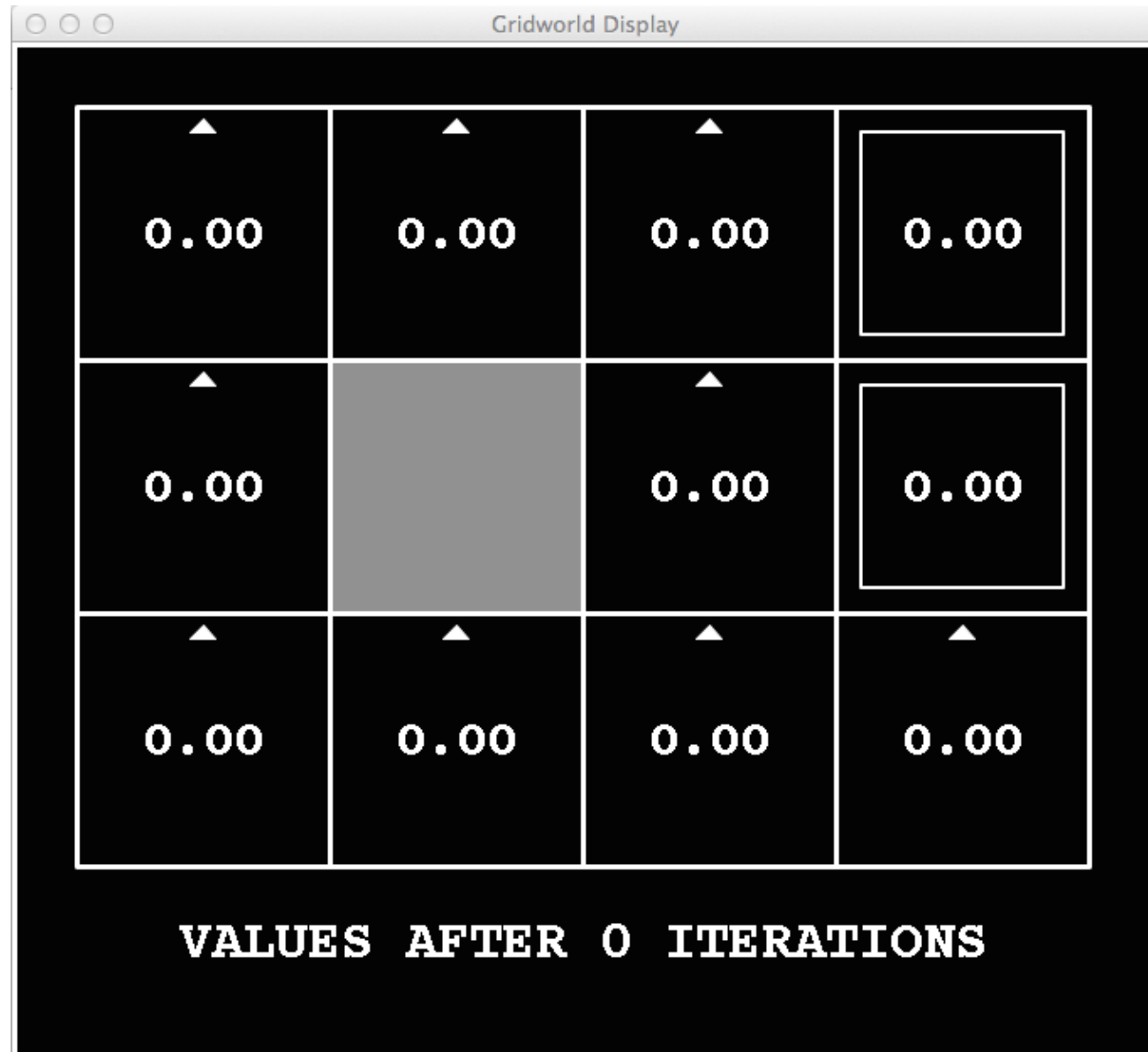  o Note: deep parts of the tree eventually don't matter if $\gamma < 1$

# Time-Limited Values

o Key idea: time-limited values

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

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

$V_2(\ \text{🚗}\ )$

# 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



Gridworld Display

| | | | |
|---|---|---|---|
| 0.37 ▶ | 0.66 ▶ | 0.83 ▶ | 1.00 |
| ▲ 0.00 | | ▲ 0.51 | -1.00 |
| ▲ 0.00 | 0.00 ▶ | ▲ 0.31 | ◀ 0.00 |

**VALUES AFTER 4 ITERATIONS**
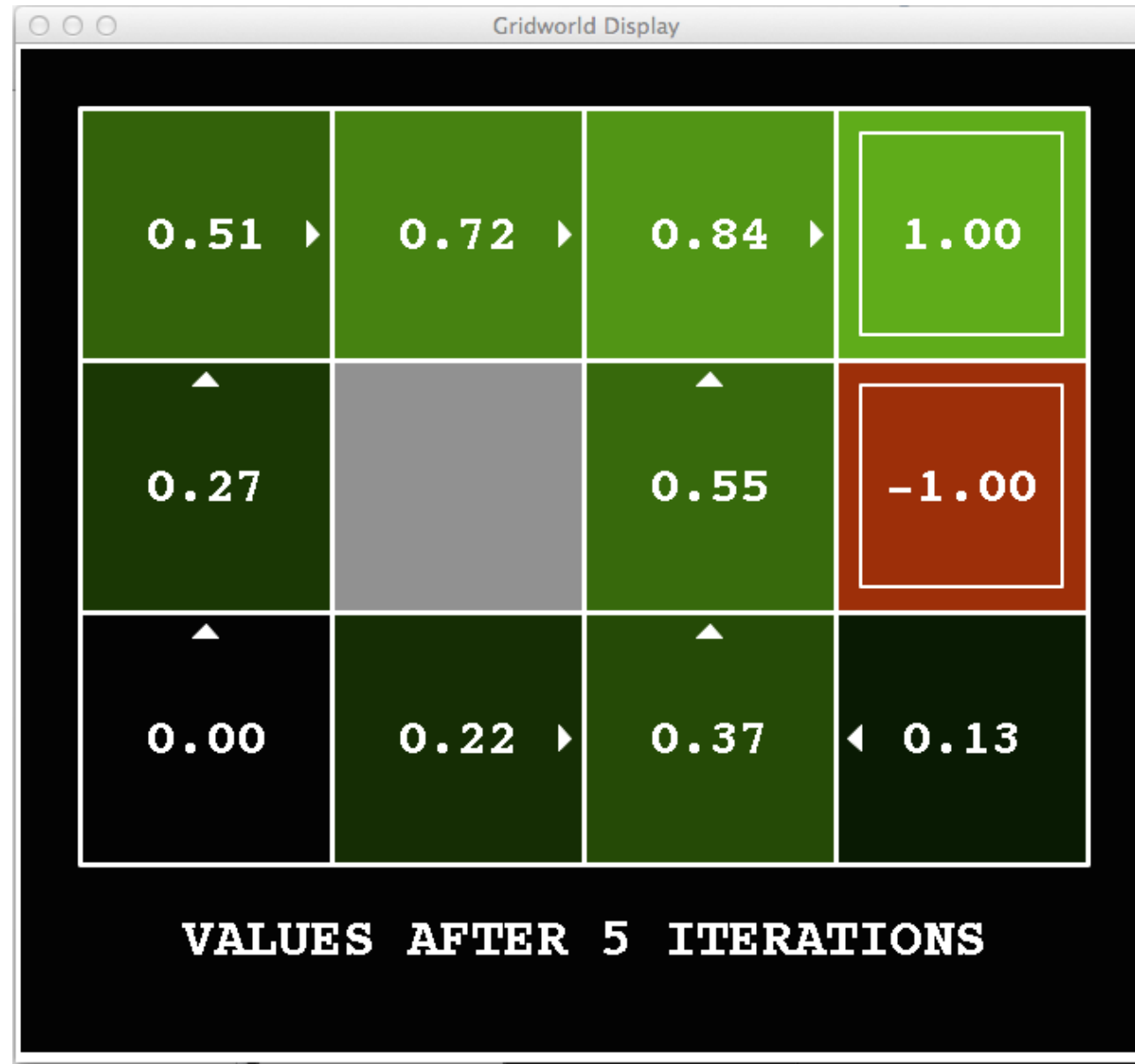
Noise = 0.2
Discount = 0.9
Living reward = 0

# 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
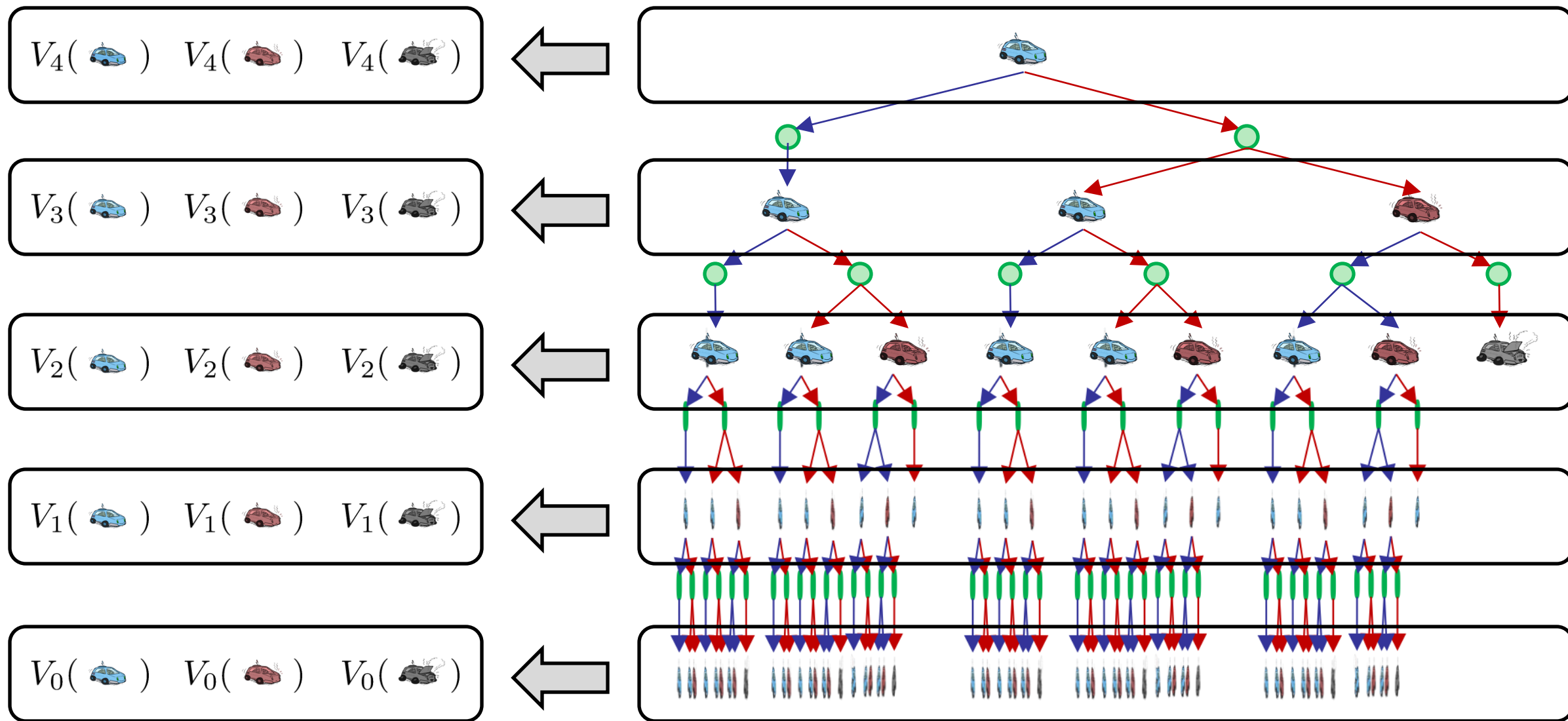


Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



Noise = 0.2
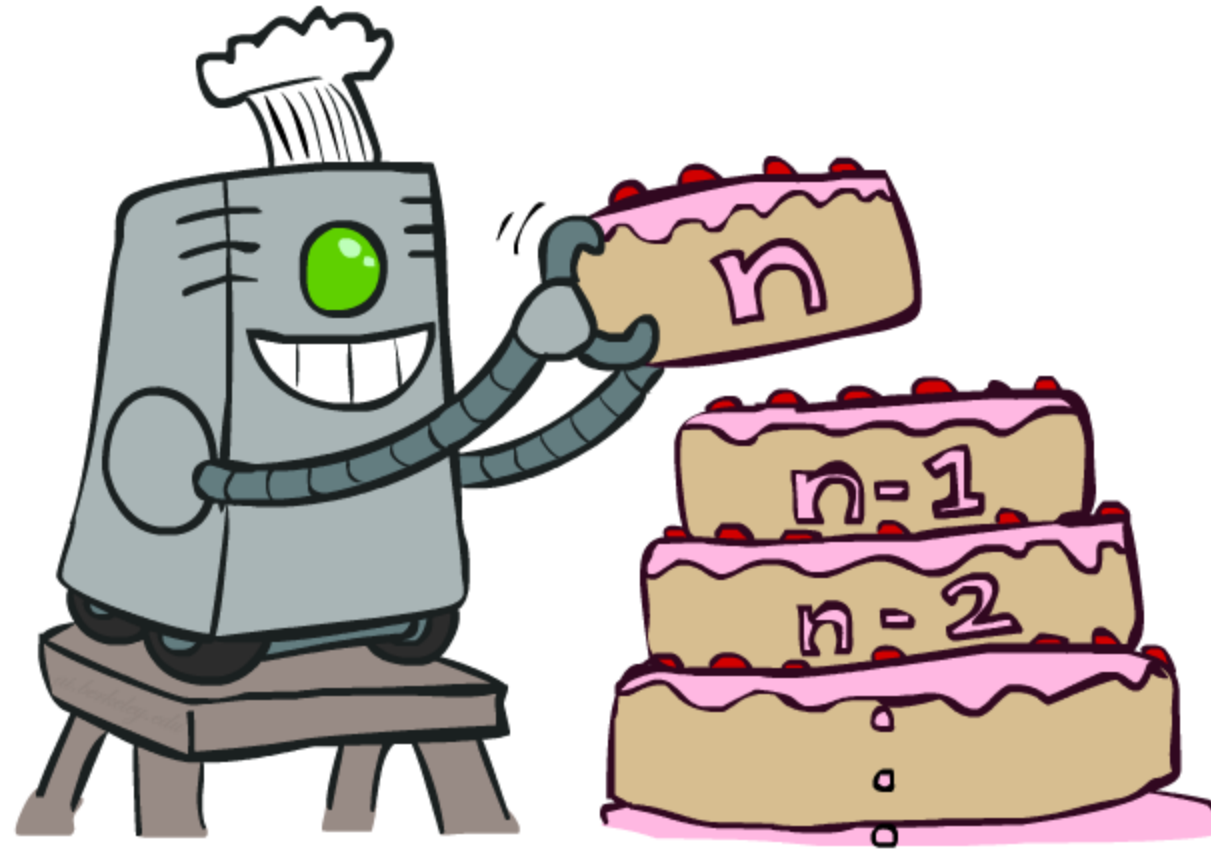Discount = 0.9
Living reward = 0

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# Computing Time-Limited Values



$V_4(\text{🚗}) \quad V_4(\text{🚗}) \quad V_4(\text{🚗})$

$V_3(\text{🚗}) \quad V_3(\text{🚗}) \quad V_3(\text{🚗})$

$V_2(\text{🚗}) \quad V_2(\text{🚗}) \quad V_2(\text{🚗})$

$V_1(\text{🚗}) \quad V_1(\text{🚗}) \quad V_1(\text{🚗})$

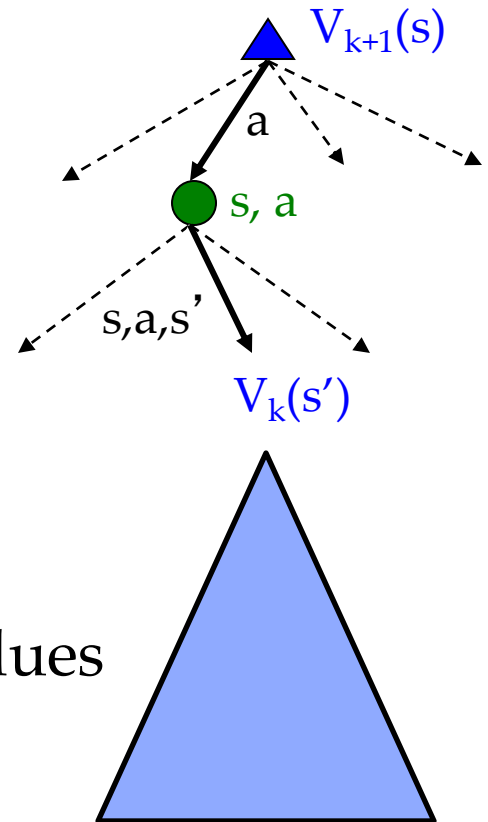$V_0(\text{🚗}) \quad V_0(\text{🚗}) \quad V_0(\text{🚗})$

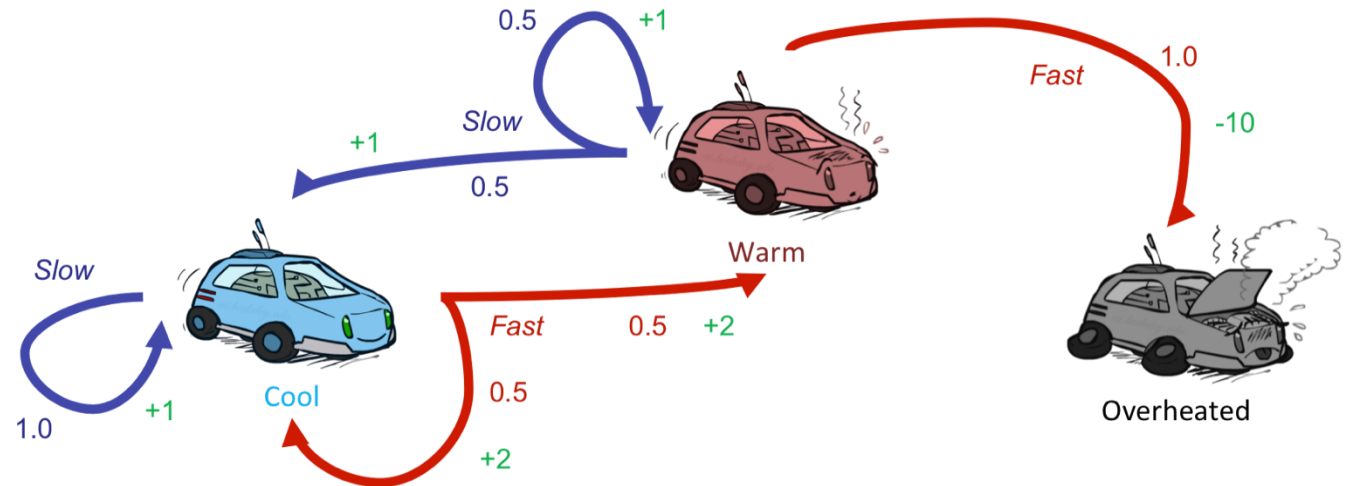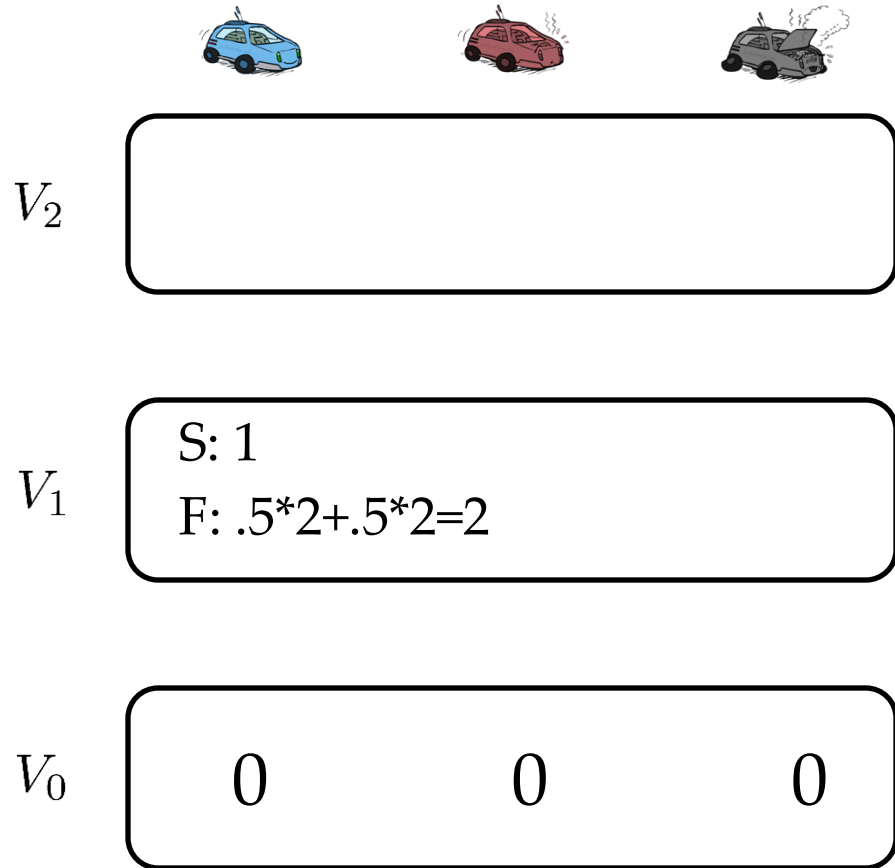# Value Iteration

# Value Iteration: Dynamic Programming

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

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

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



$V_{k+1}(s)$

s, a

s,a,s'

$V_k(s')$

   o V = B(V) Where B is the Bellman update operator

o Repeat until convergence, which yields V*

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

o Theorem: Value Iteration will converge to unique optimal values
   o Basic idea: approximations get refined towards optimal values
   o Policy may converge long before values do
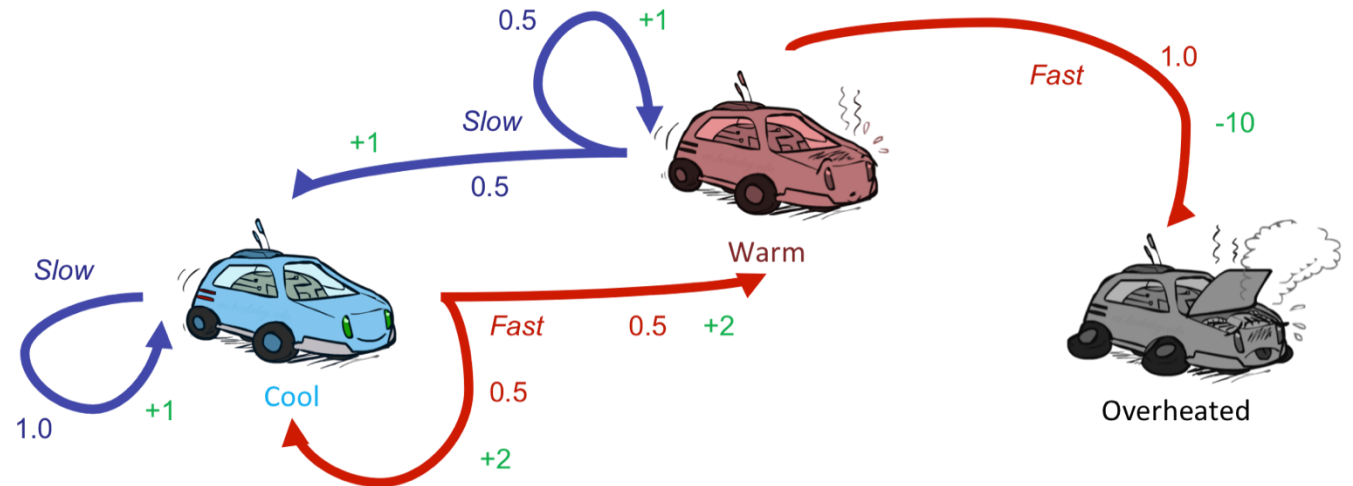
# Example: Value Iteration

$V_2$

$V_1$
S: 1
F: .5*2+.5*2=2

$V_0$
0        0        0



0.5    +1

Fast    1.0

Slow

+1    -10

Slow

Warm

Slow    Fast    0.5    +2

+1    Cool    0.5

1.0    +2    Overheated

*Assume no discount!*

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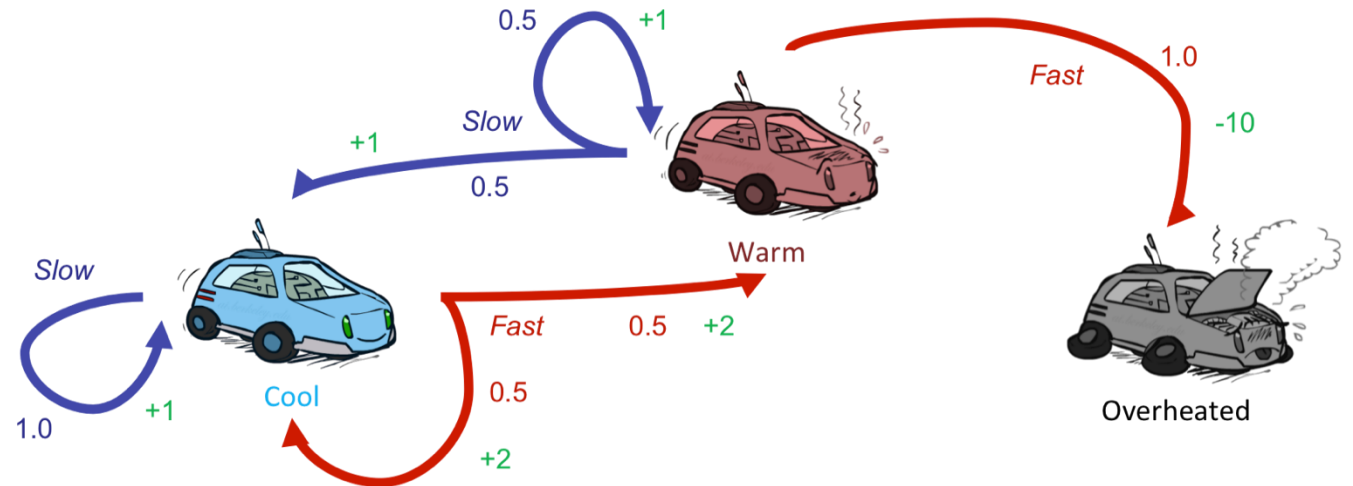
# Example: Value Iteration



$V_2$

$V_1$ | 2 | S: .5*1+.5*1=1
F: -10

$V_0$ | 0 | 0 | 0

*Assume no discount!*

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

Labels in figure: 
- 0.5  +1  Slow  0.5
- Fast  1.0  -10
- Slow  Warm
- Cool  Fast  0.5  +2
- 1.0  +1  0.5  +2
- Overheated

# Example: Value Iteration



$V_2$

$V_1$  |  2  |  1  |  0

$V_0$  |  0  |  0  |  0

*Assume no discount!*

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

# Example: Value Iteration
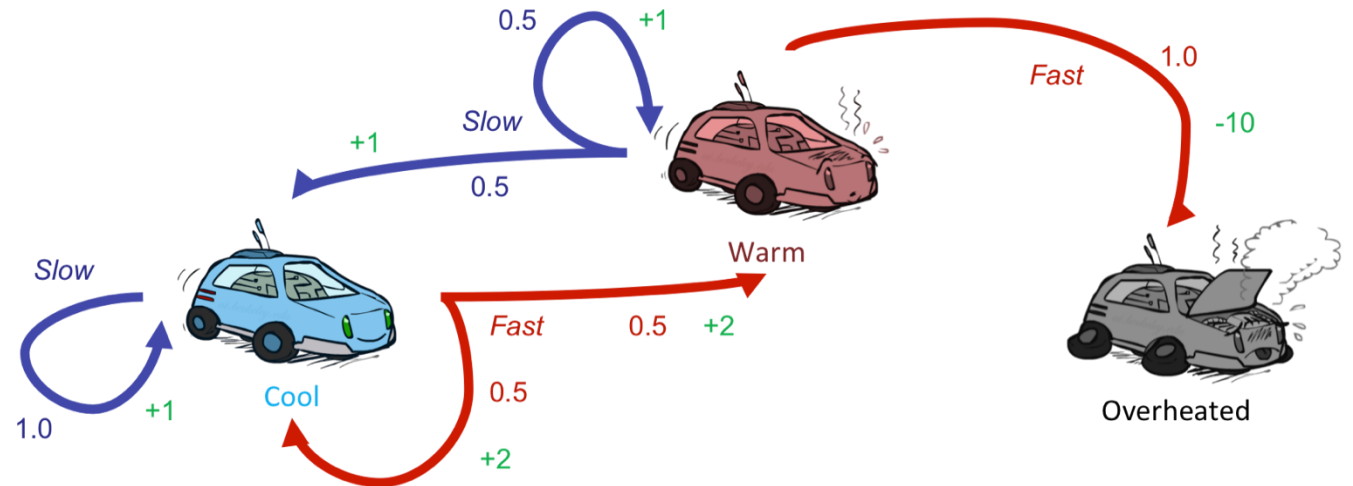


$V_2$

S: 1+2=3
F: .5*(2+2)+.5*(2+1)=3.5

$V_1$

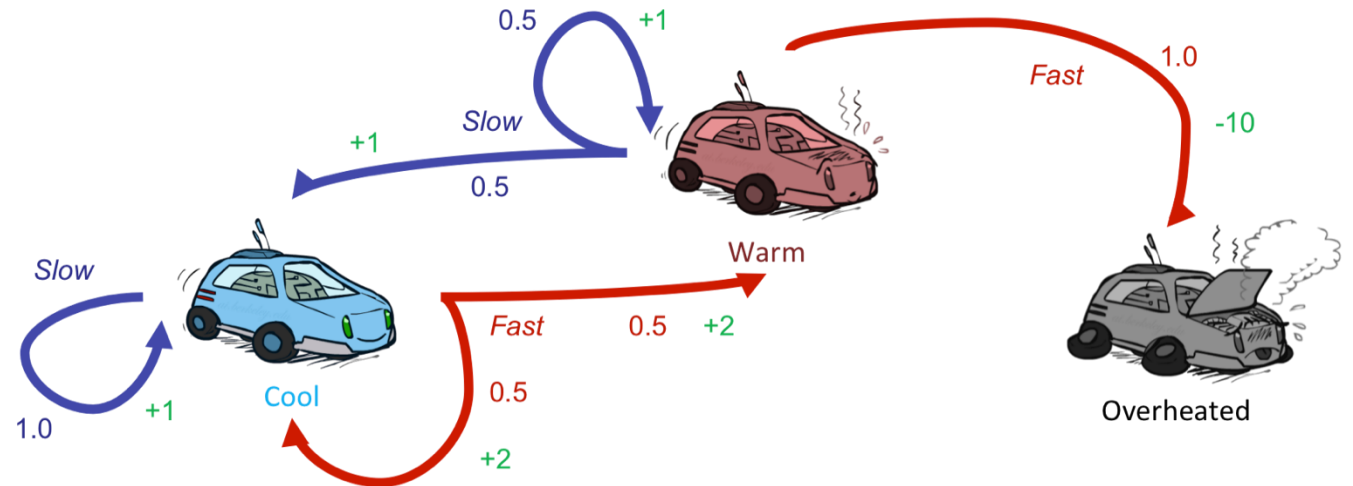| 2 | 1 | 0 |

$V_0$

| 0 | 0 | 0 |

*Assume no discount!*

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

# Example: Value Iteration

$V_2$ | 3.5 | 2.5 | 0

$V_1$ | 2 | 1 | 0

$V_0$ | 0 | 0 | 0



*Assume no discount!*

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