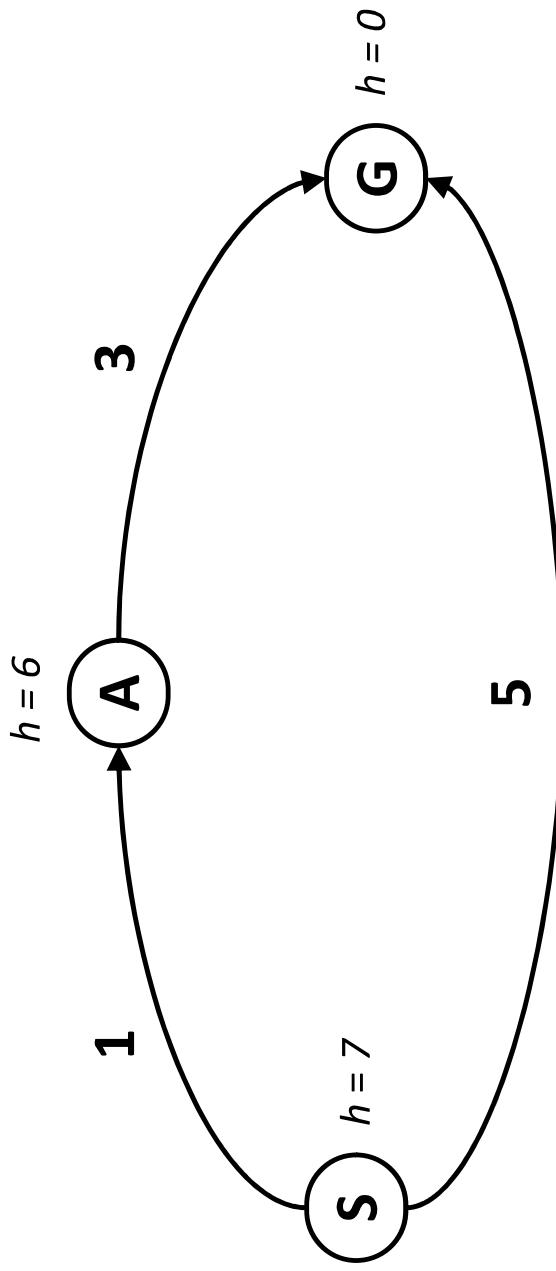
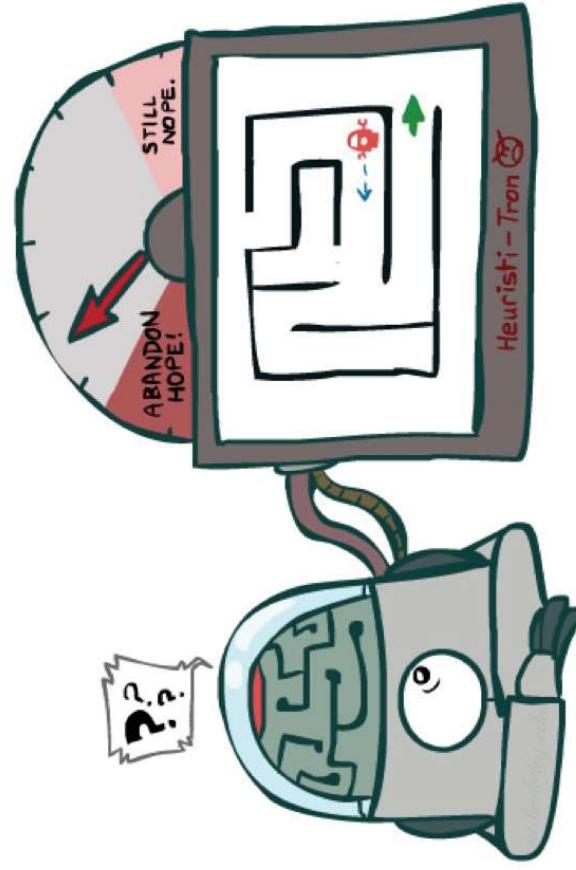


Is A* Optimal?

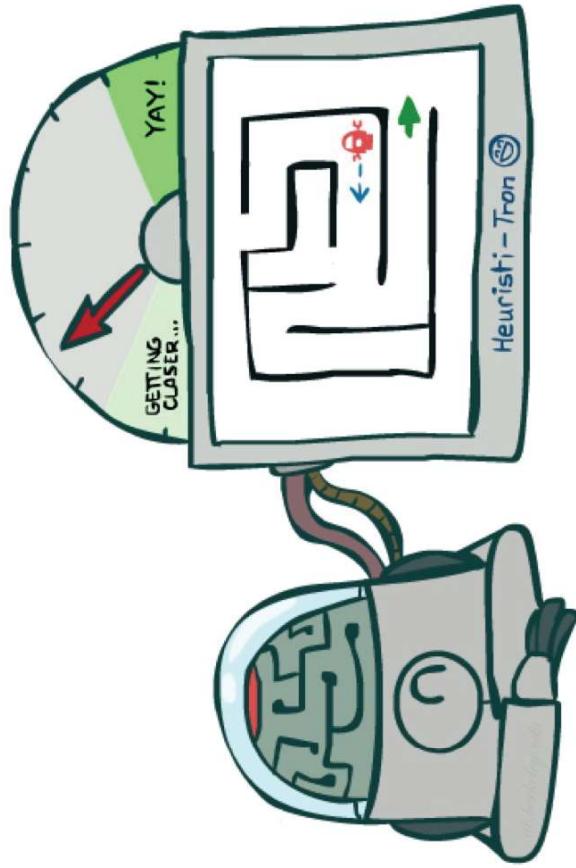


- What went wrong?
 - Actual bad goal cost < estimated good goal cost
 - We need estimates to be less than actual costs!

Idea: Admissibility



Inadmissible (pessimistic) heuristics break optimality by trapping good plans on the fringe



Admissible (optimistic) heuristics never outweigh true costs

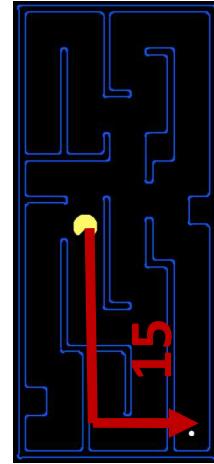
Admissible Heuristics

- A heuristic h is *admissible* (optimistic) if:

$$0 \leq h(n) \leq h^*(n)$$

where $h^*(n)$ is the true cost to a nearest goal

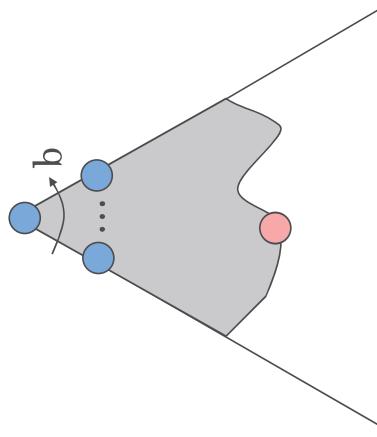
- Examples:



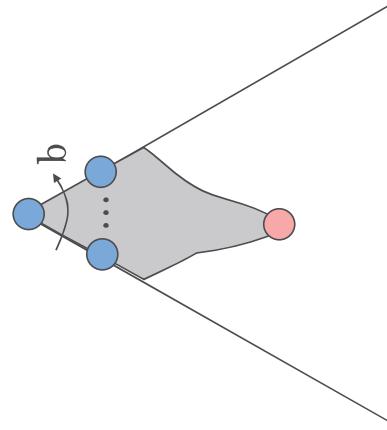
- Coming up with admissible heuristics is most of what's involved in using A* in practice.

Properties of A*

Uniform-Cost

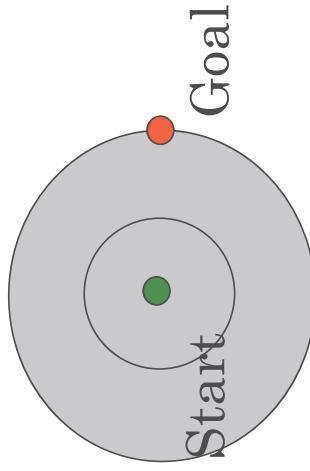


A*

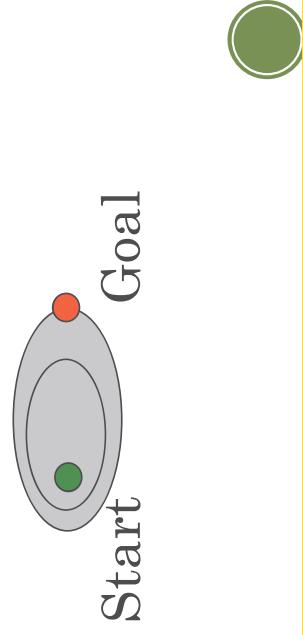


UCS vs A* Contours

- Uniform-cost expands equally in all “directions”



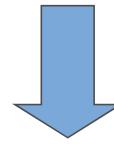
- A* expands mainly toward the goal, but does hedge its bets to ensure optimality



Graph Search

- Very simple fix: never expand a state type twice

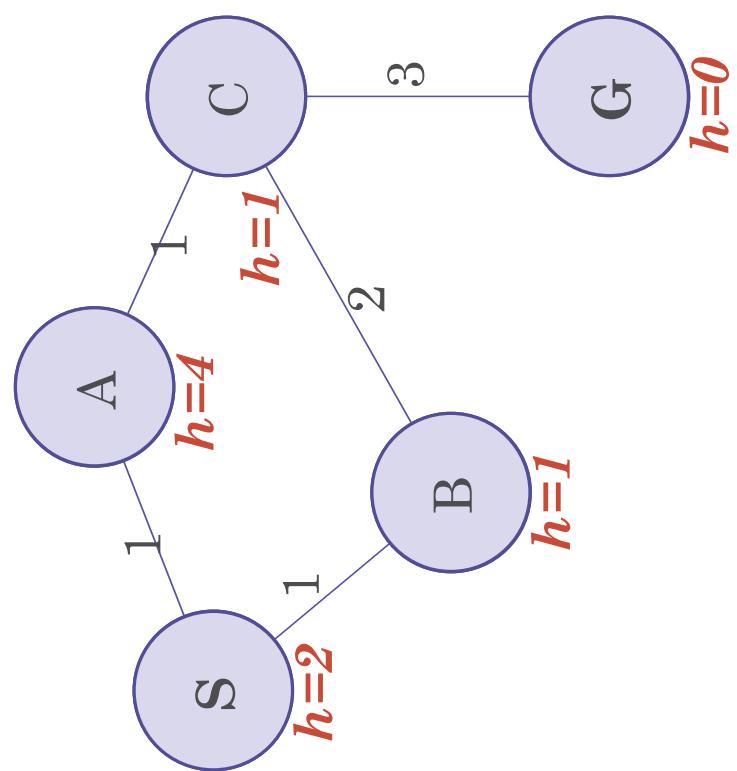
```
function GRAPH-SEARCH(problem, fringe) returns a solution, or failure
  closed  $\leftarrow$  an empty set
  fringe  $\leftarrow$  INSERT(MAKE-NODE(INITIAL-STATE[problem]), fringe)
  loop do
    if fringe is empty then return failure
    node  $\leftarrow$  REMOVE-FRONT(fringe)
    if GOAL-TEST(problem, STATE[node]) then return node
    if STATE[node] is not in closed then
      add STATE[node] to closed
      fringe  $\leftarrow$  INSERTALL(EXPAND(node, problem), fringe)
    end
```



- Can this wreck completeness? Why or why not?
- How about optimality? Why or why not?

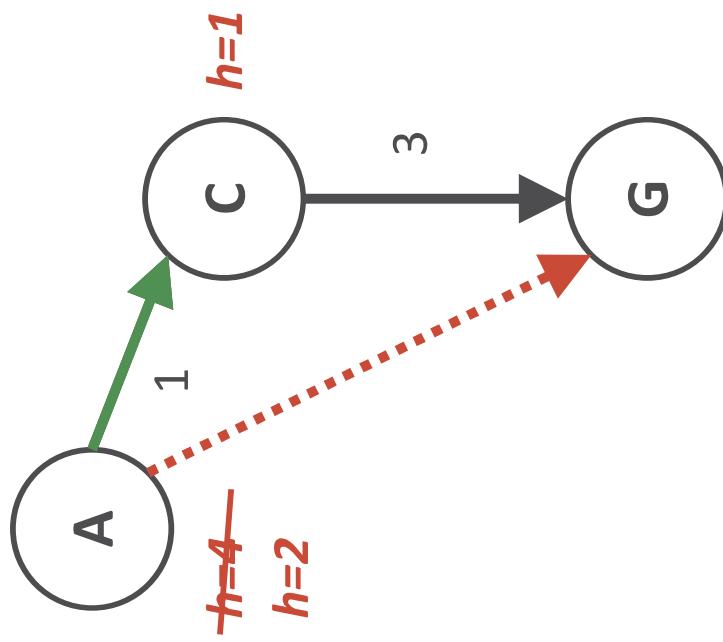
A* Graph Search Gone Wrong

State space graph



Search tree

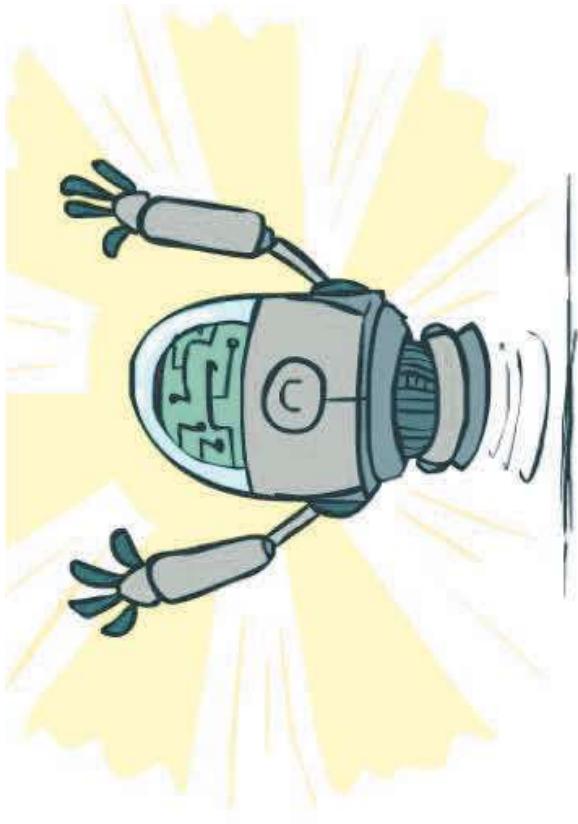
Consistency of Heuristics



- Main idea: estimated heuristic costs \leq actual costs
 - Admissibility: heuristic cost \leq actual cost to goal
$$h(A) \leq \text{actual cost from } A \text{ to } G$$
 - Consistency: heuristic “arc” cost \leq actual cost for each arc
$$h(A) - h(C) \leq \text{cost}(A \text{ to } C)$$
- Consequences of consistency:
 - The f value along a path never decreases
$$h(A) \leq \text{cost}(A \text{ to } C) + h(C)$$
 - $f(A) = g(A) + h(A) \leq g(A) + \text{cost}(A \text{ to } C) + h(C) \leq f(C)$
- A* graph search is optimal

Optimality

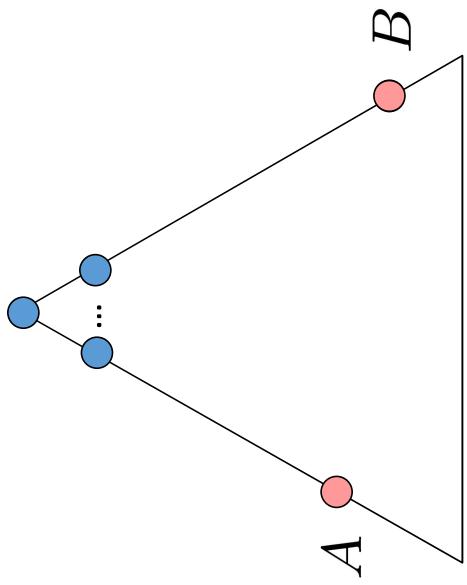
- Tree search:
 - A* is optimal if heuristic is admissible
 - UCS is a special case ($h = 0$)
- Graph search:
 - A* optimal if heuristic is consistent
 - UCS optimal ($h = 0$ is consistent)
- Consistency implies admissibility
- In general, most natural admissible heuristics tend to be consistent, especially if from relaxed problems



Optimality of A* Tree Search

Assume:

- A is an optimal goal node
- B is a suboptimal goal node
- h is admissible



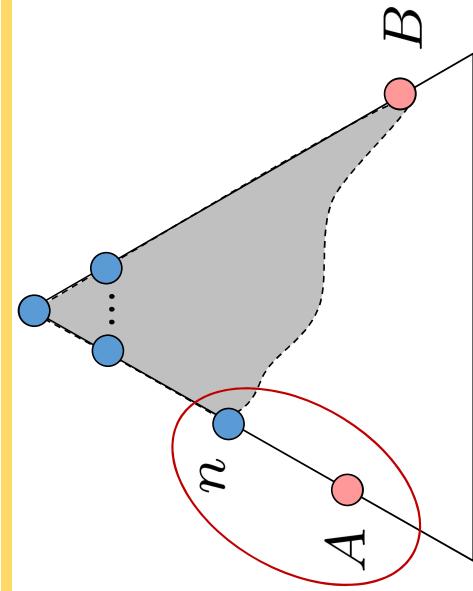
Claim:

- A will exit the fringe before B

Optimality of A* Tree Search: Blocking

Proof:

- Imagine B is on the fringe
- Some ancestor n of A is on the fringe, too (maybe A!)
- Claim: n will be expanded before B
 1. $f(n)$ is less or equal to $f(A)$



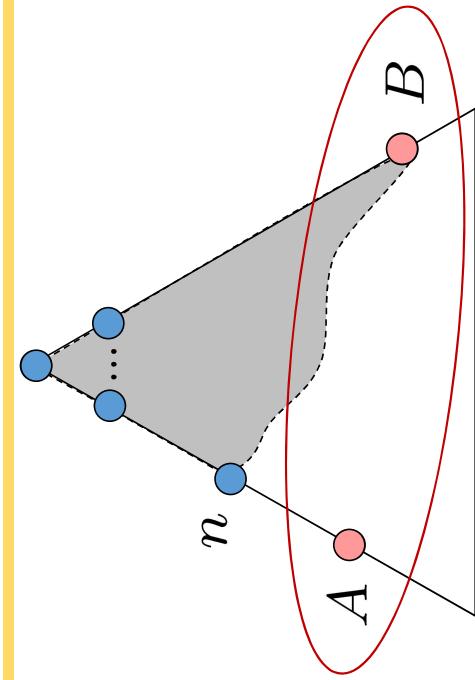
$$\begin{aligned}f(n) &= g(n) + h(n) \\f(n) &\leq g(A) \\g(A) &= f(A)\end{aligned}$$

Definition of f-cost
Admissibility of h
 $h = 0$ at a goal

Optimality of A* Tree Search: Blocking

Proof:

- Imagine B is on the fringe
- Some ancestor n of A is on the fringe, too (maybe A!)
- Claim: n will be expanded before B
 1. $f(n)$ is less or equal to $f(A)$
 2. $f(A)$ is less than $f(B)$



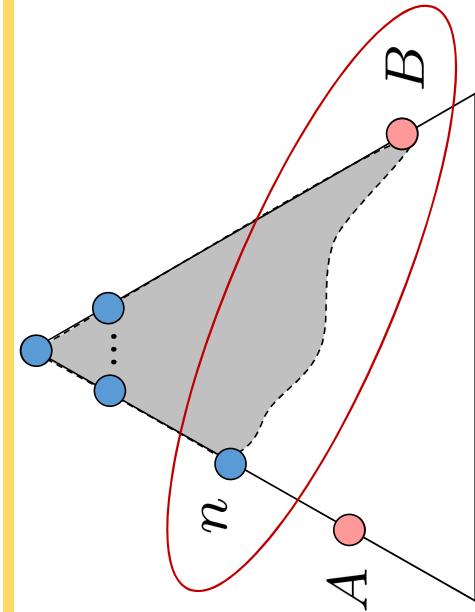
$$\begin{aligned}g(A) &< g(B) \\ f(A) &< f(B)\end{aligned}$$

B is suboptimal
 $h = 0$ at a goal

Optimality of A* Tree Search: Blocking

Proof:

- Imagine B is on the fringe
- Some ancestor n of A is on the fringe, too (maybe A!)
- Claim: n will be expanded before B
 1. $f(n)$ is less or equal to $f(A)$
 2. $f(A)$ is less than $f(B)$
 3. n expands before B
- All ancestors of A expand before B
- A expands before B
- A* search is optimal



$$f(n) \leq f(A) < f(B)$$

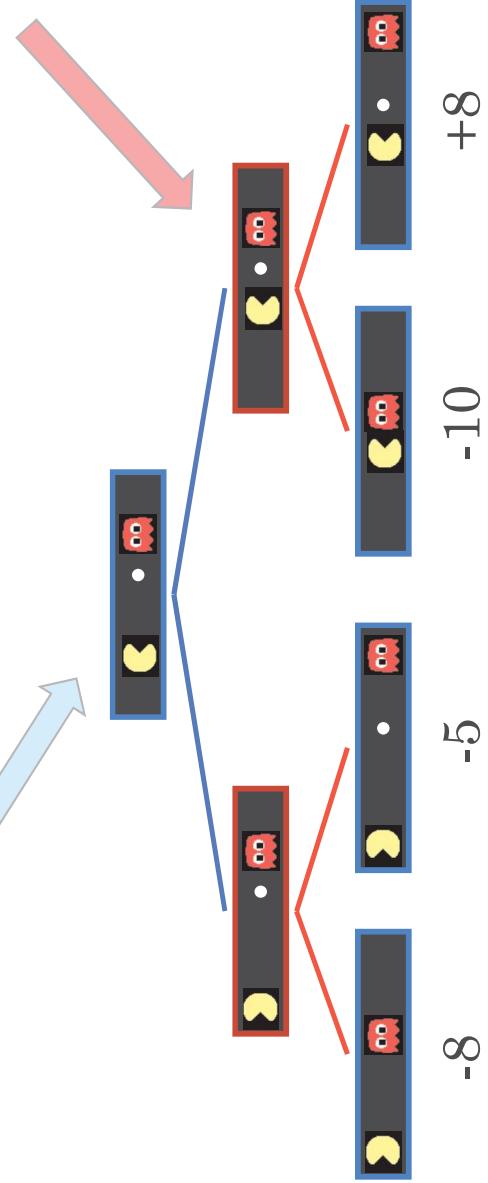
Minimax Values

States Under Agent's Control:

$$V(s) = \max_{s' \in \text{successors}(s)} V(s')$$

States Under Opponent's Control:

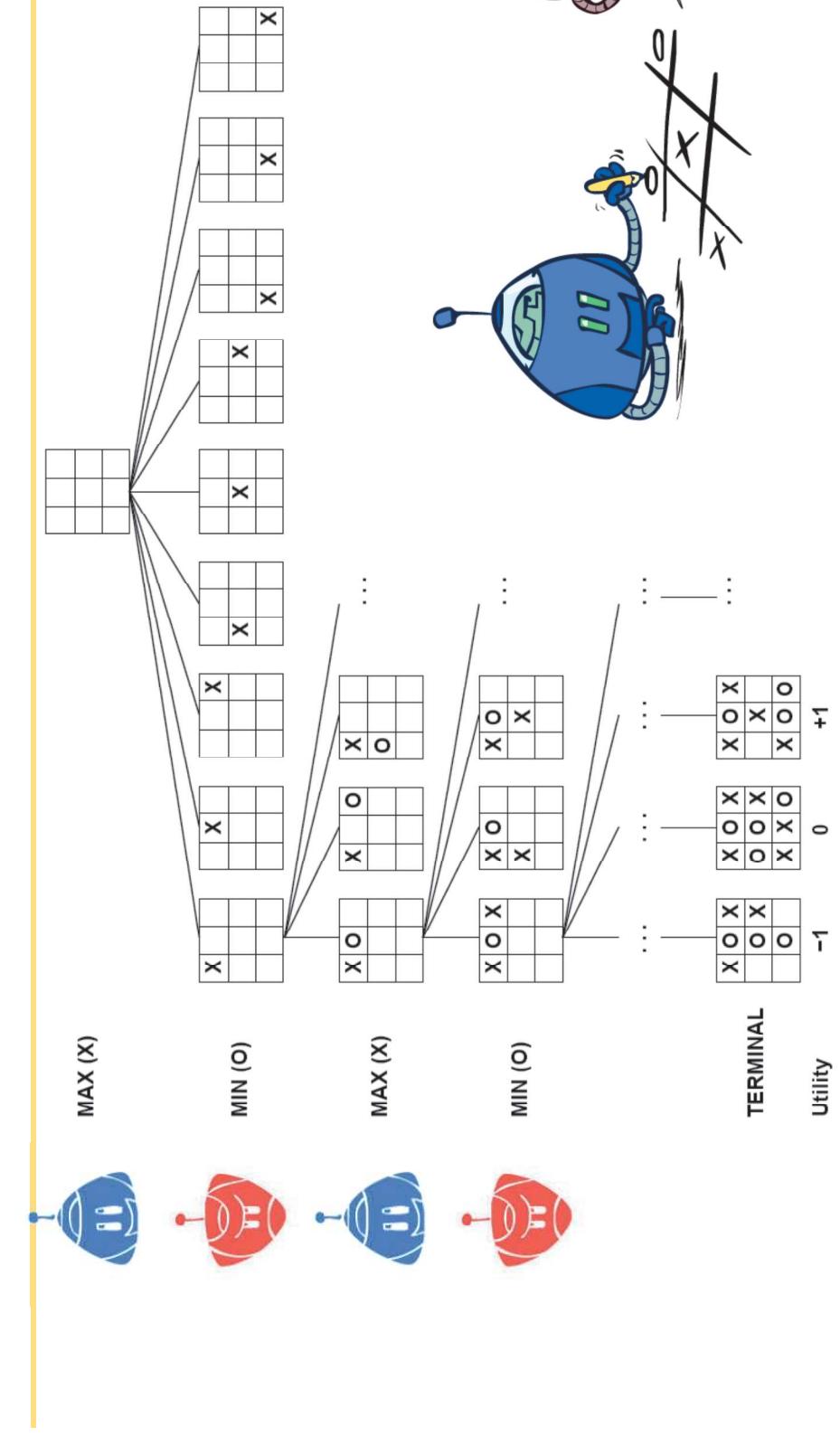
$$V(s') = \min_{s \in \text{successors}(s')} V(s)$$



Terminal States:
 $V(s) = \text{known}$



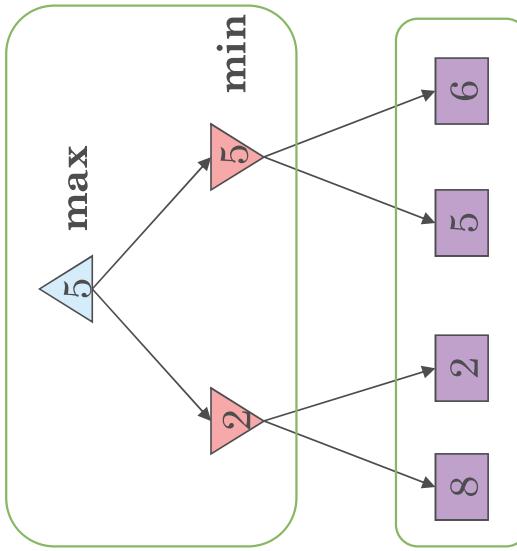
Tic-Tac-Toe Game Tree



Adversarial Search (Minimax)

- Deterministic, zero-sum games:
 - Tic-tac-toe, chess, checkers
 - One player maximizes result
 - The other minimizes result

Minimax values:
computed recursively



- Minimax search:

- A state-space search tree
- Players alternate turns
- Compute each node's **minimax value**:
the best achievable utility against a rational (optimal) adversary

Terminal values:
part of the game

Minimax Implementation

```
def value(state):
    if the state is a terminal state: return the state's utility
    if the next agent is MAX: return max-value(state)
    if the next agent is MIN: return min-value(state)
```

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v
```

```
def min-value(state):
    initialize v = +∞
    for each successor of state:
        v = min(v, value(successor))
    return v
```

Minimax Efficiency

- How efficient is minimax?

- Just like (exhaustive) DFS

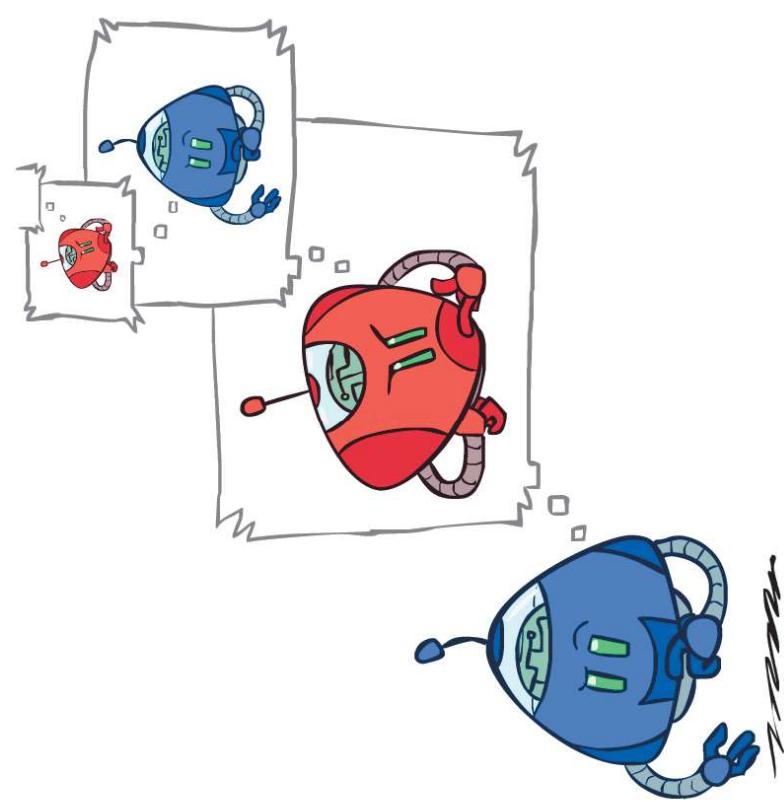
- Time: $O(b^m)$

- Space: $O(bm)$

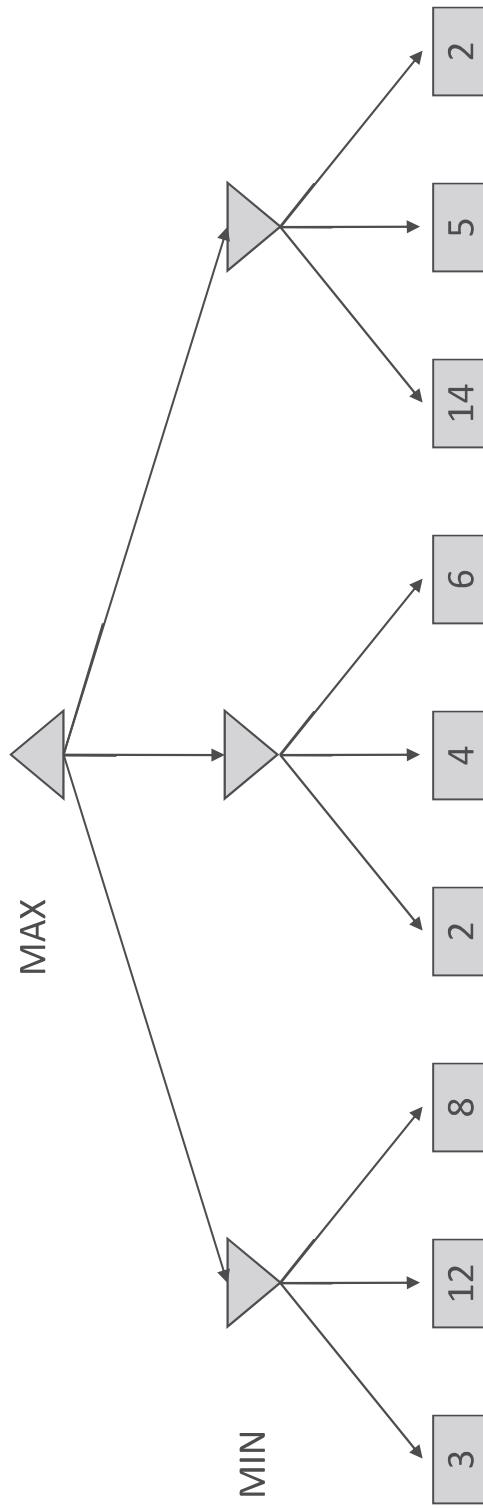
- Example: For chess, $b \approx 35$, $m \approx 100$

- Exact solution is completely infeasible

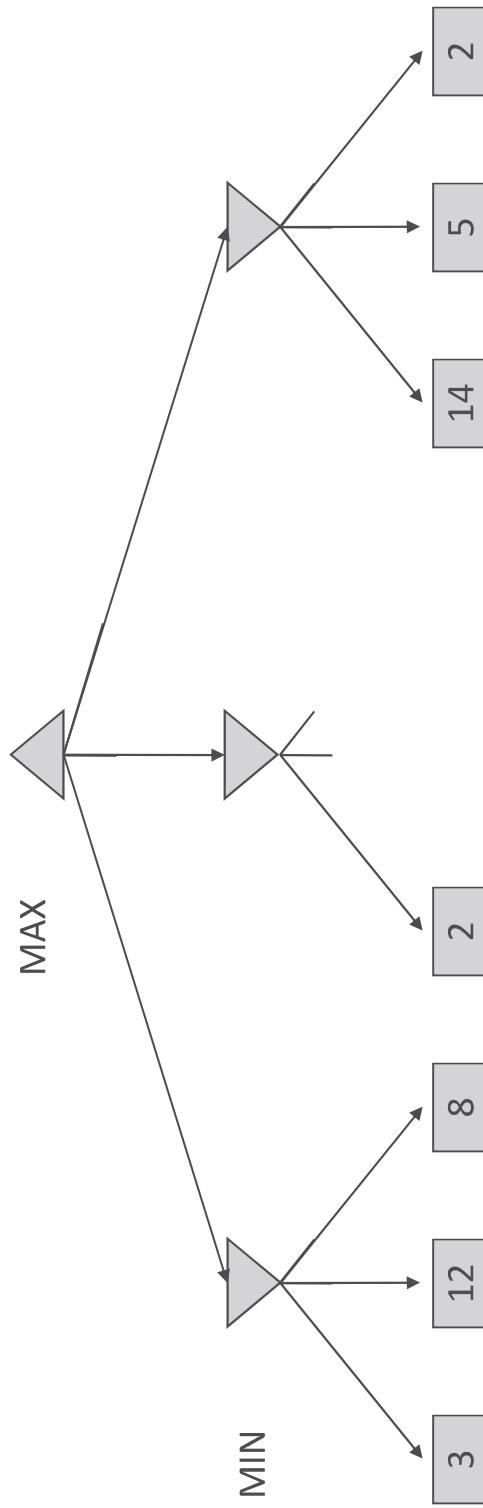
- But, do we need to explore the whole tree?



Minimax Example



Minimax Pruning



Alpha-Beta Pruning

- Alpha α : value of the best choice so far for MAX (lower bound of Max utility)
- Beta β : value of the best choice so far for MIN (upper bound of Min utility)
 - Expanding at MAX node n : update α
 - If a child of n has value greater than β , stop expanding the MAX node n
 - Reason: MIN parent of n would not choose the action which leads to n
 - At MIN node n : update β
 - If a child of n has value less than α , stop expanding the MIN node n
 - Reason: MAX parent of n would not choose the action which leads to n

Alpha-Beta Implementation

```
def value(state, α, β):
```

 if the state is a terminal state: return the state's utility
 if the next agent is MAX: return max-value(state, α, β)
 if the next agent is MIN: return min-value(state, α, β)

```
def max-value(state, α, β):
```

 initialize v = -∞

 for each successor of state:

 v = max(v, value(successor, α, β))

 if v ≥ β return v

 α = max(α, v)

 return v

```
def min-value(state , α, β):
```

 initialize v = +∞

 for each successor of state:

 v = min(v, value(successor, α, β))

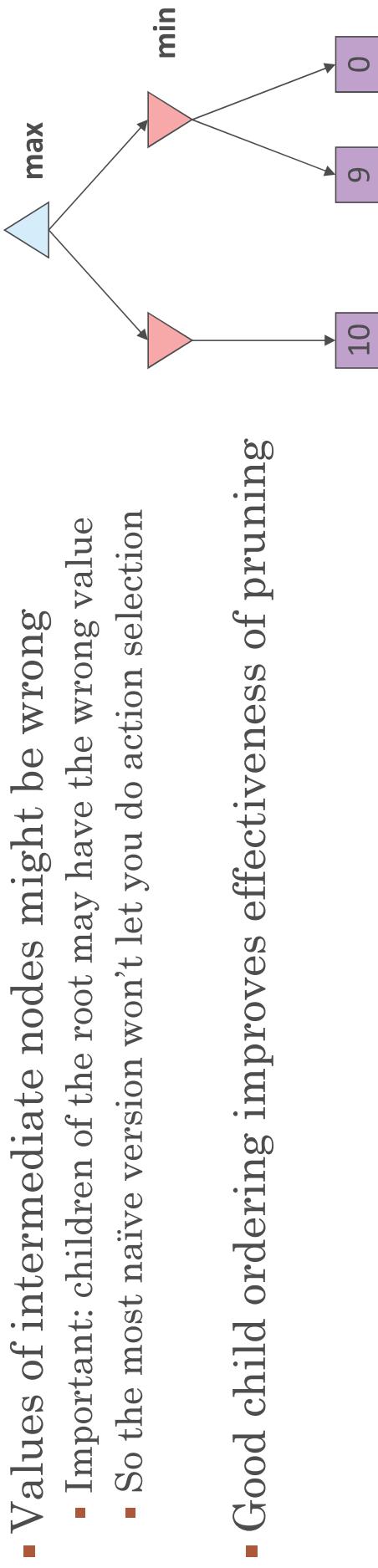
 if v ≤ α return v

 β = min(β, v)

 return v

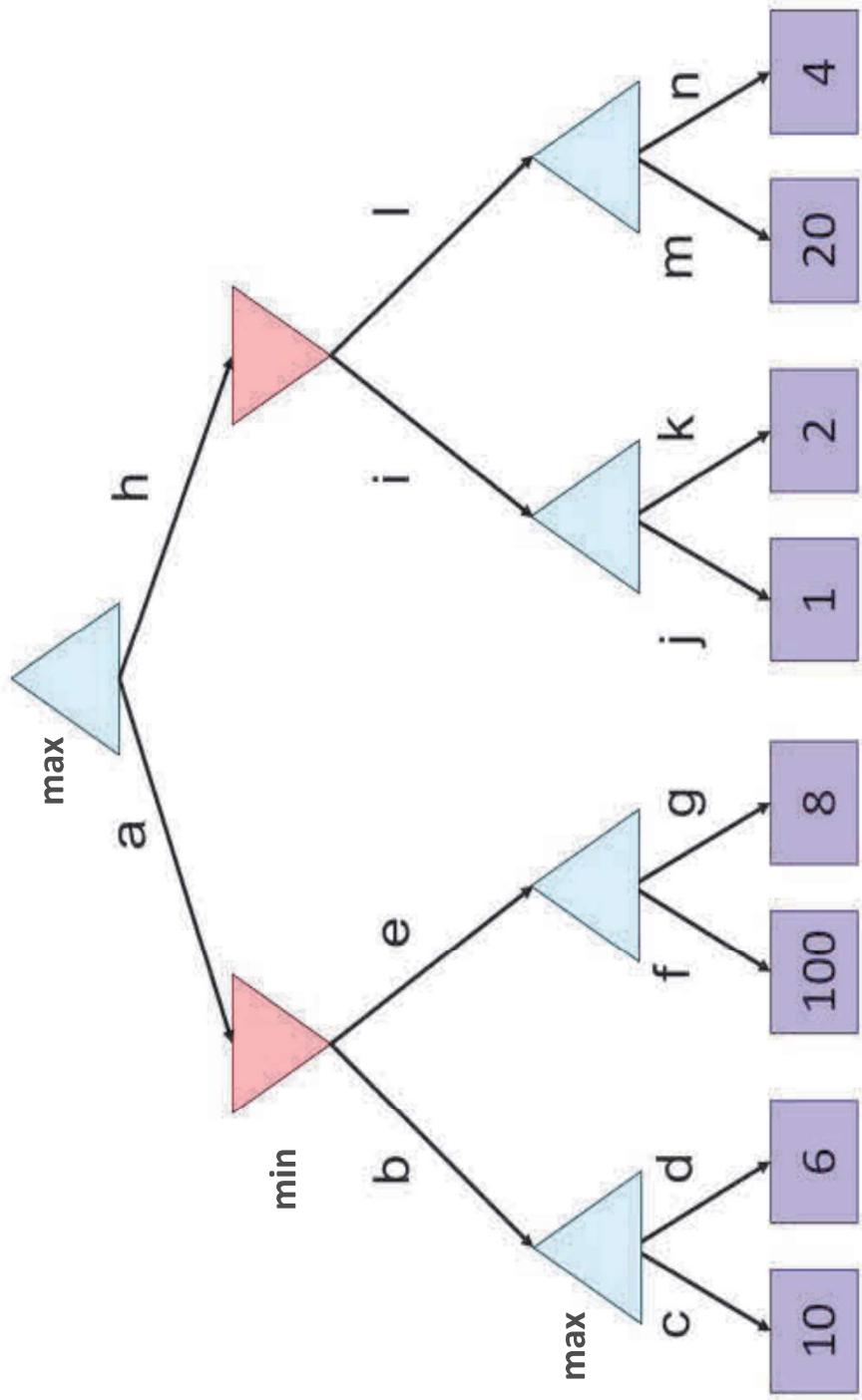
Alpha-Beta Pruning Properties

- This pruning has **no effect** on minimax value computed for the root!

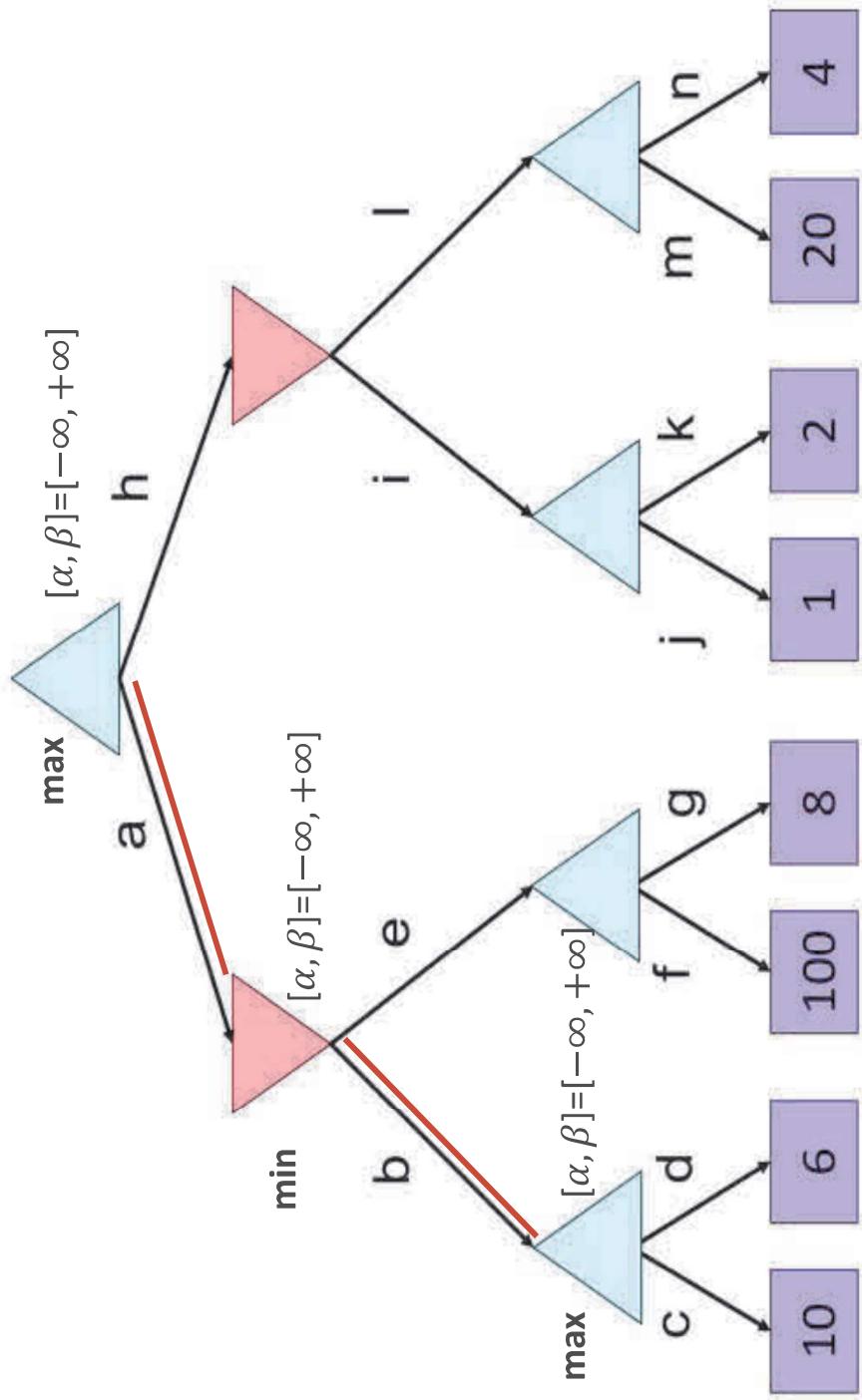


- Values of intermediate nodes might be wrong
 - Important: children of the root may have the wrong value
 - So the most naïve version won't let you do action selection
- Good child ordering improves effectiveness of pruning

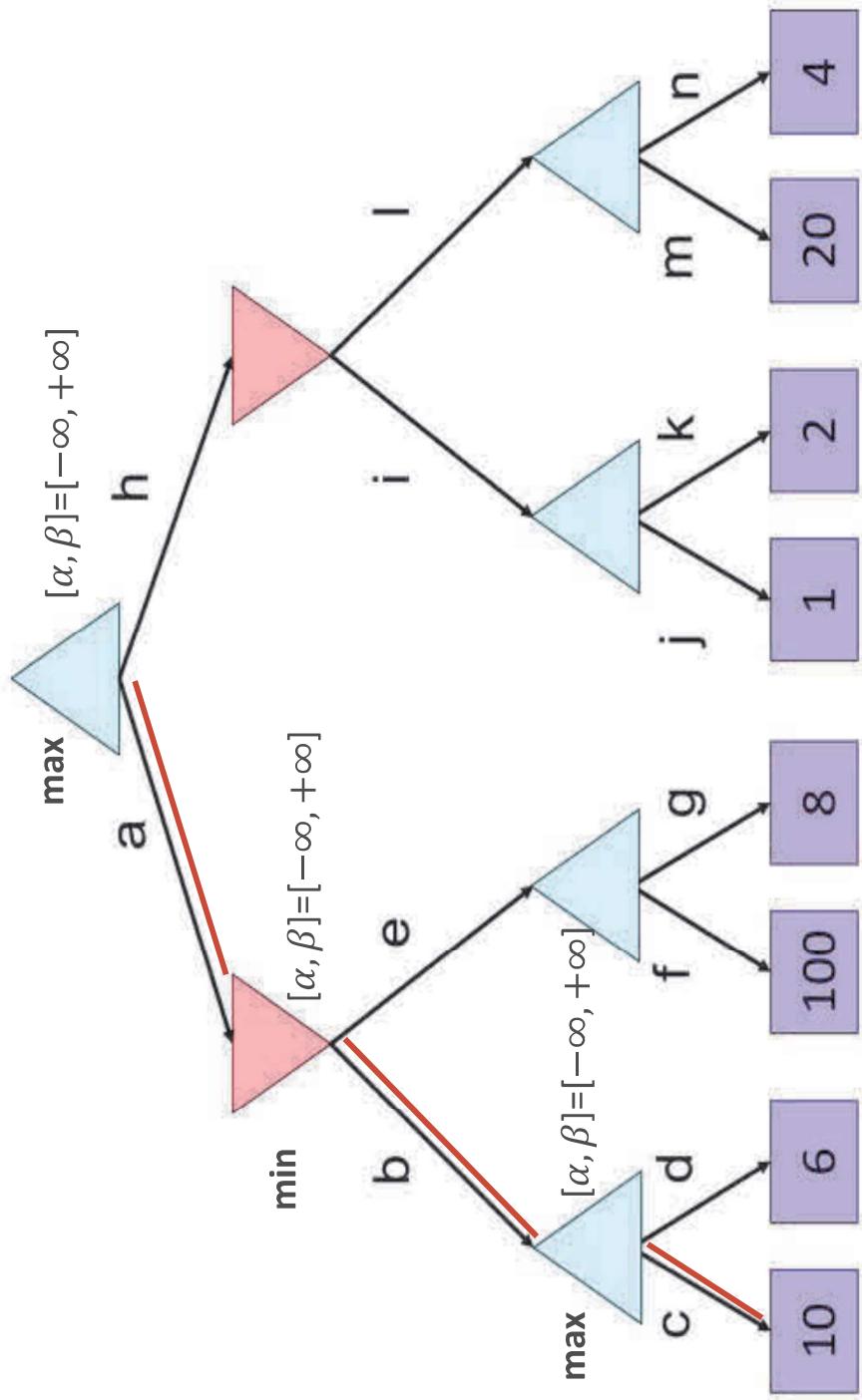
Alpha-Beta Quiz 2



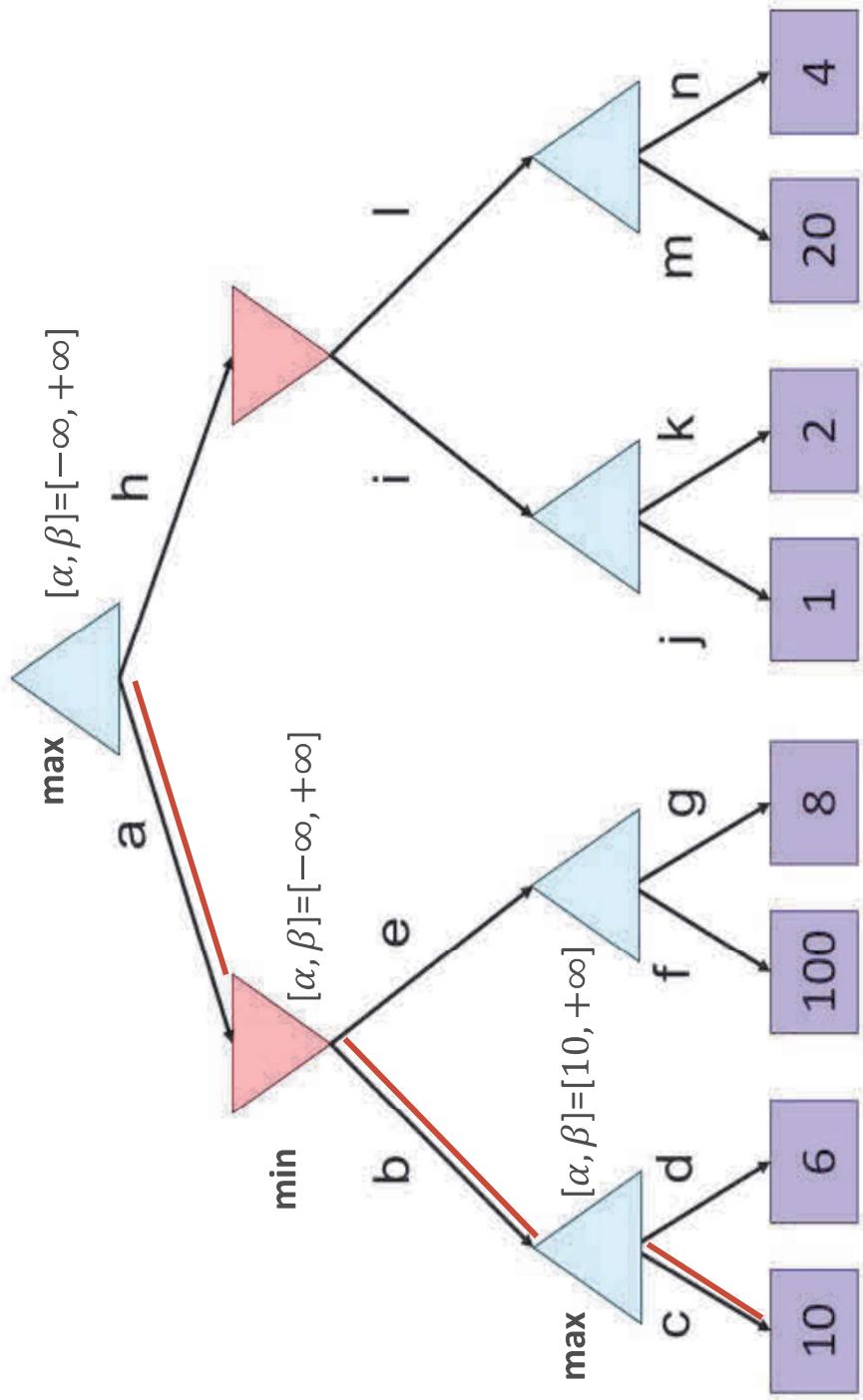
Alpha-Beta Quiz 2



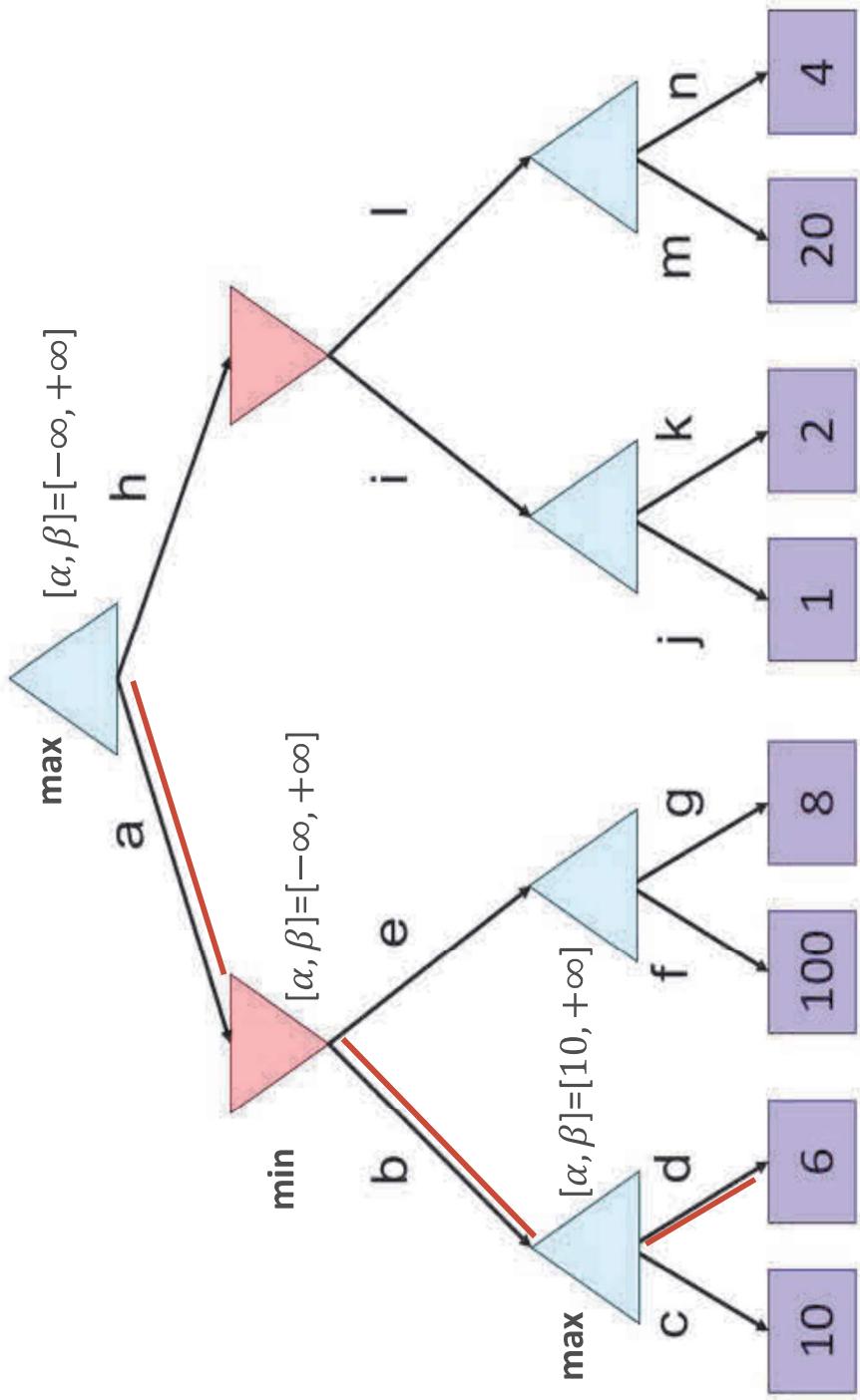
Alpha-Beta Quiz 2



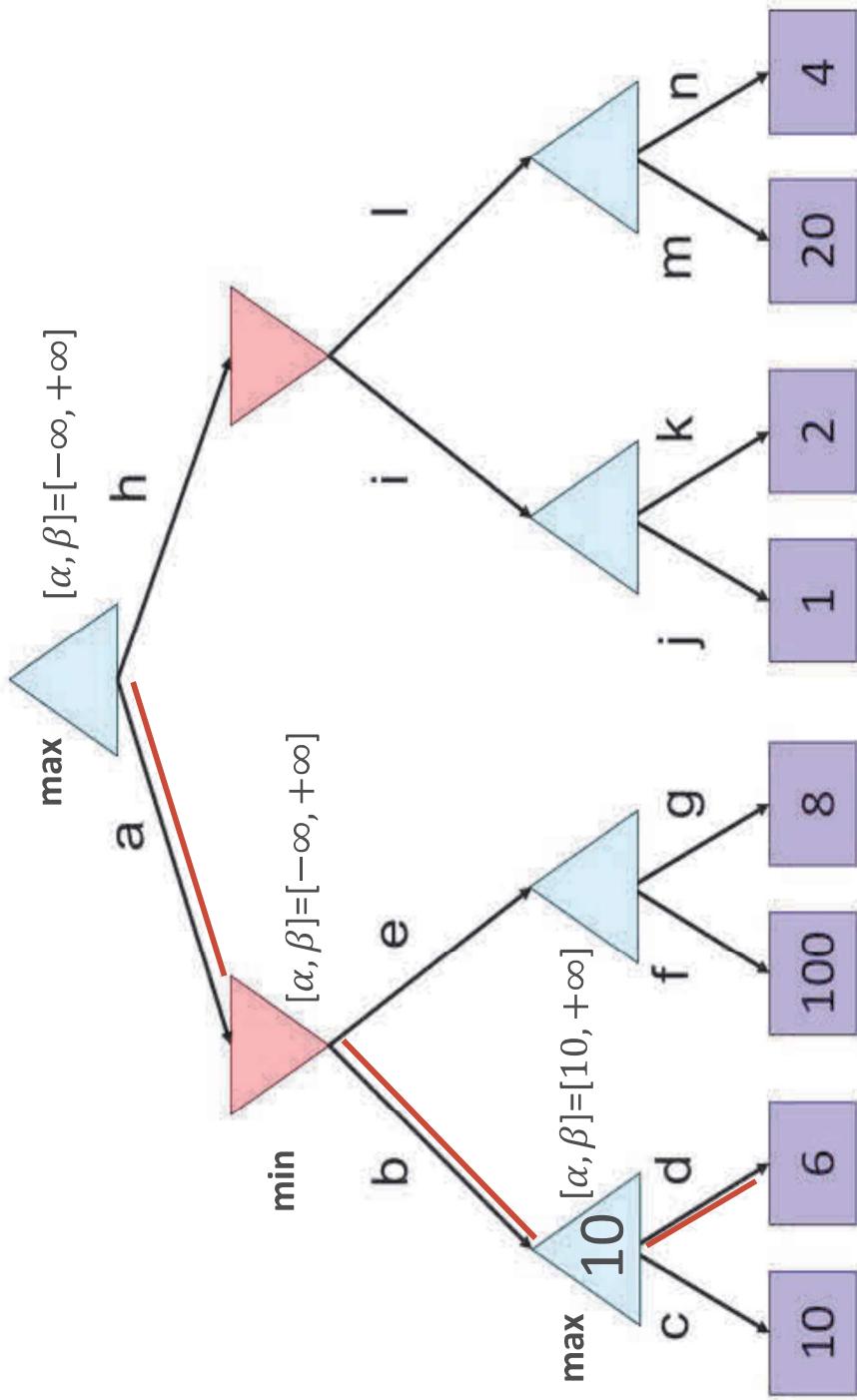
Alpha-Beta Quiz 2



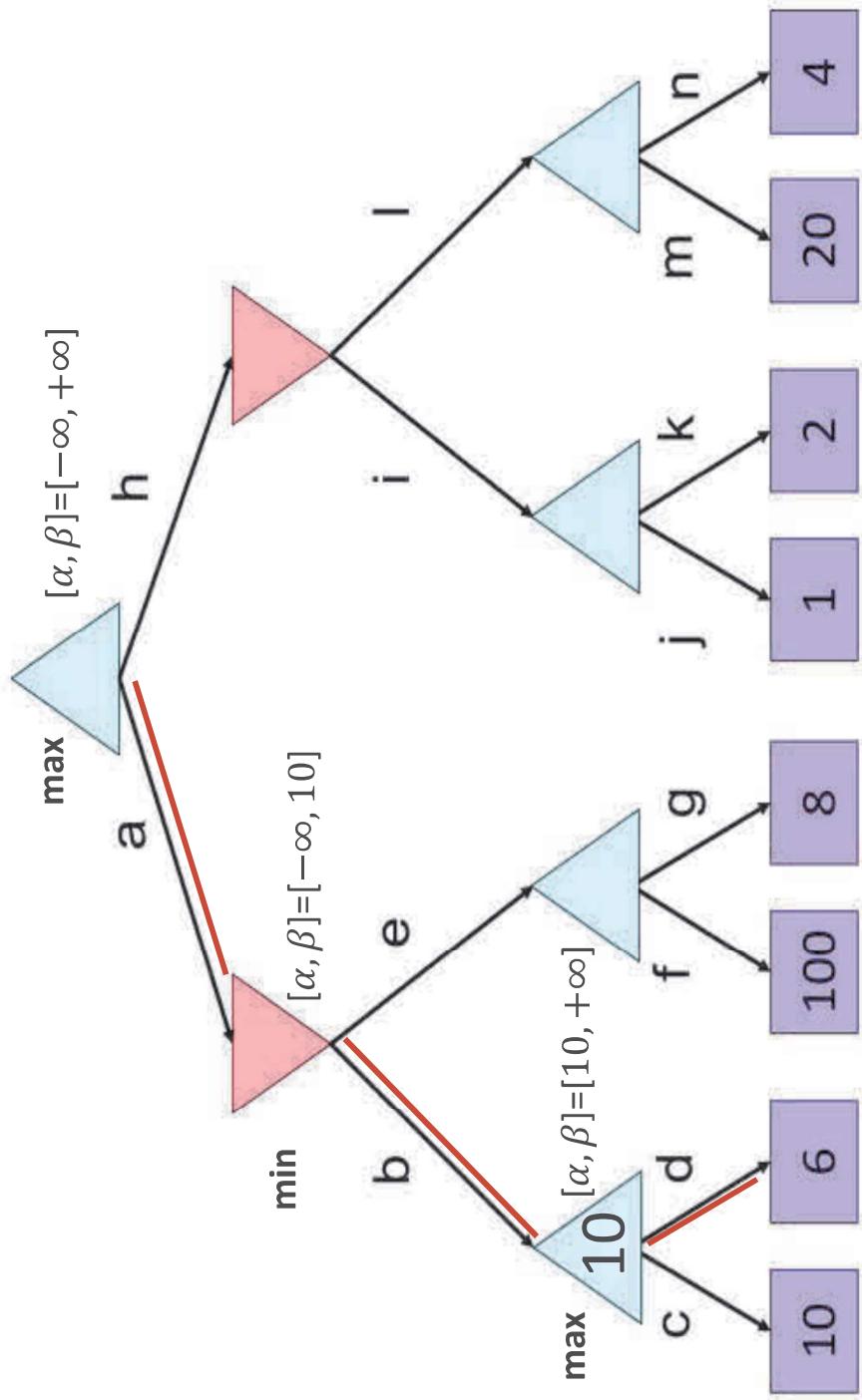
Alpha-Beta Quiz 2



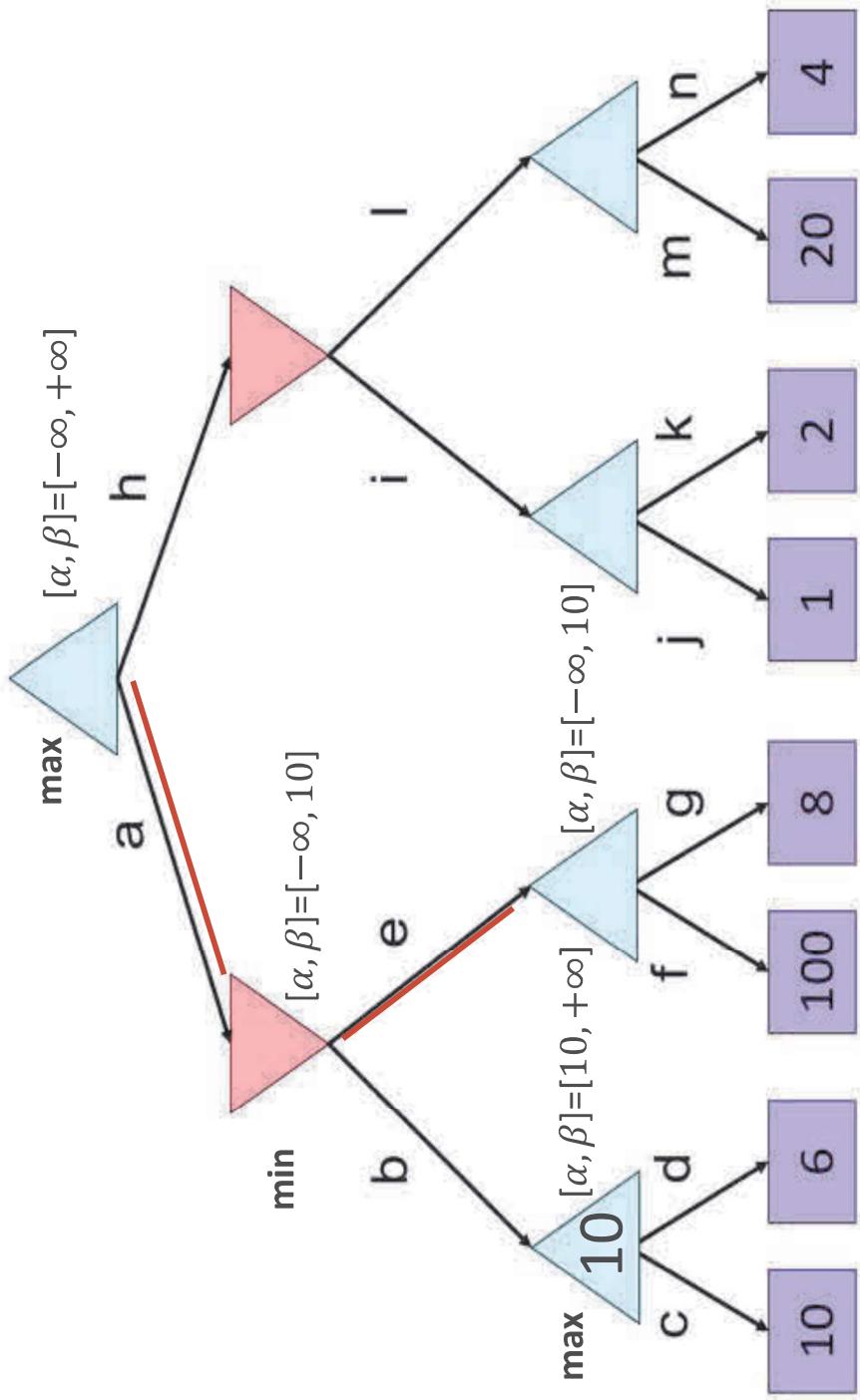
Alpha-Beta Quiz 2



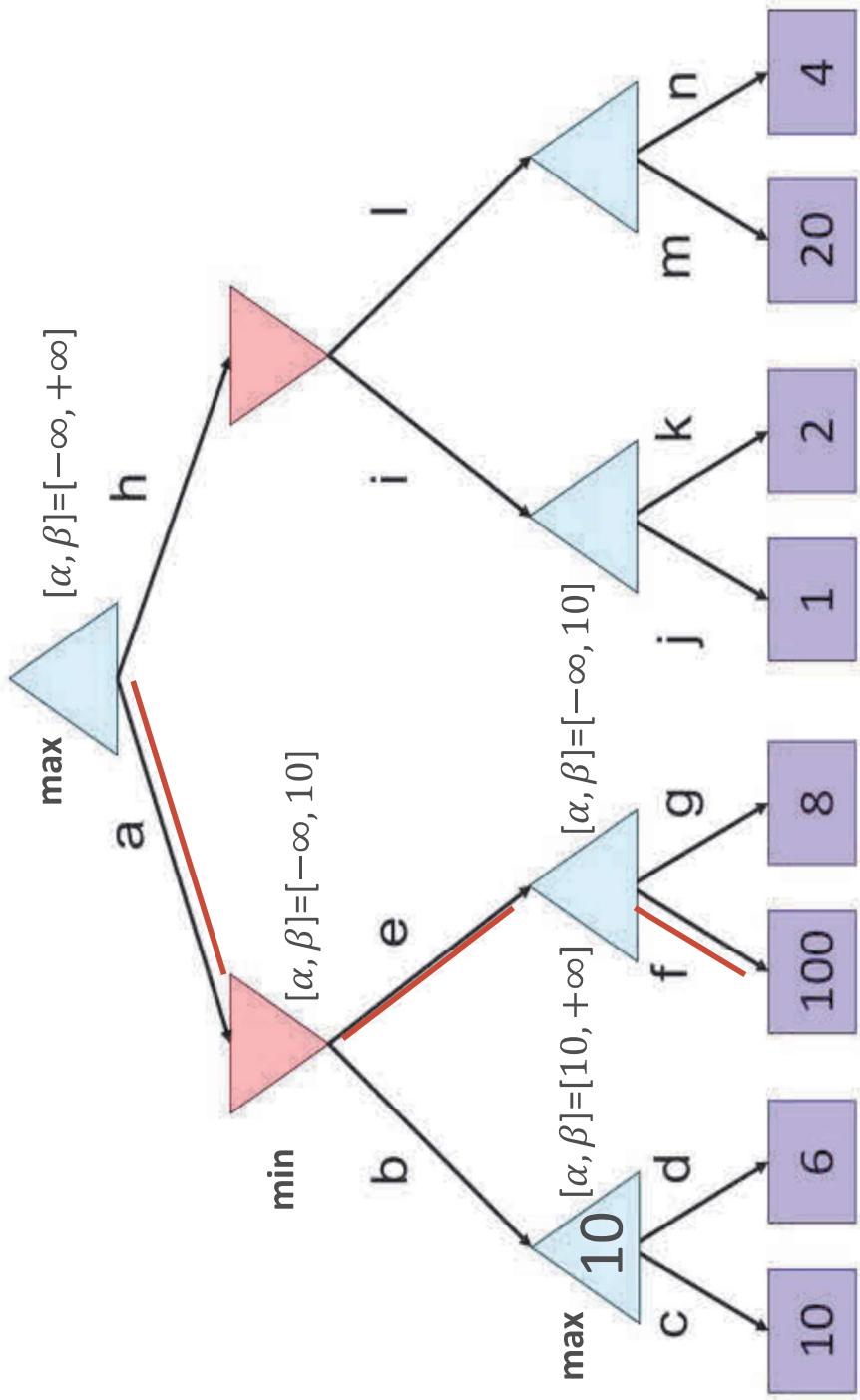
Alpha-Beta Quiz 2



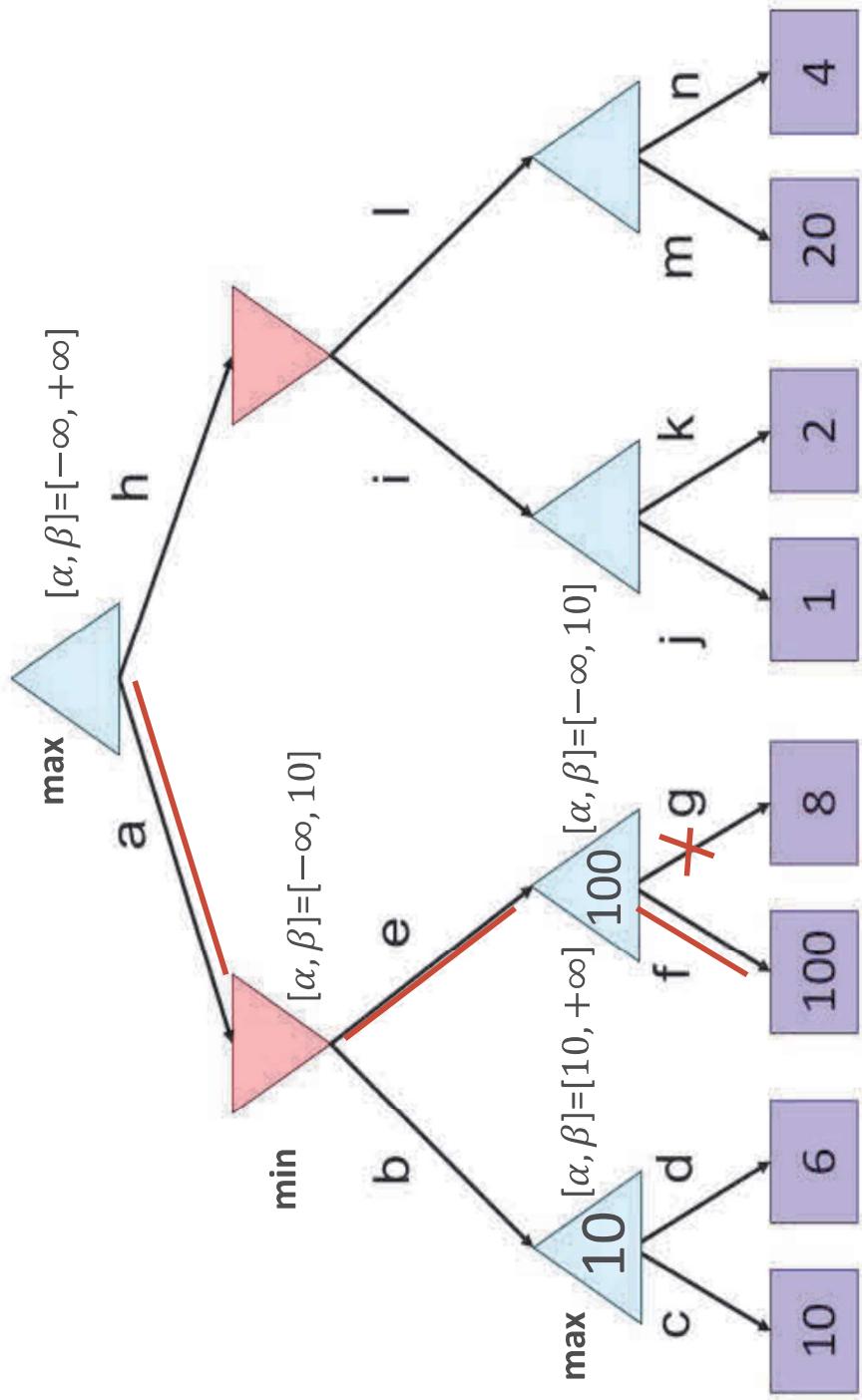
Alpha-Beta Quiz 2



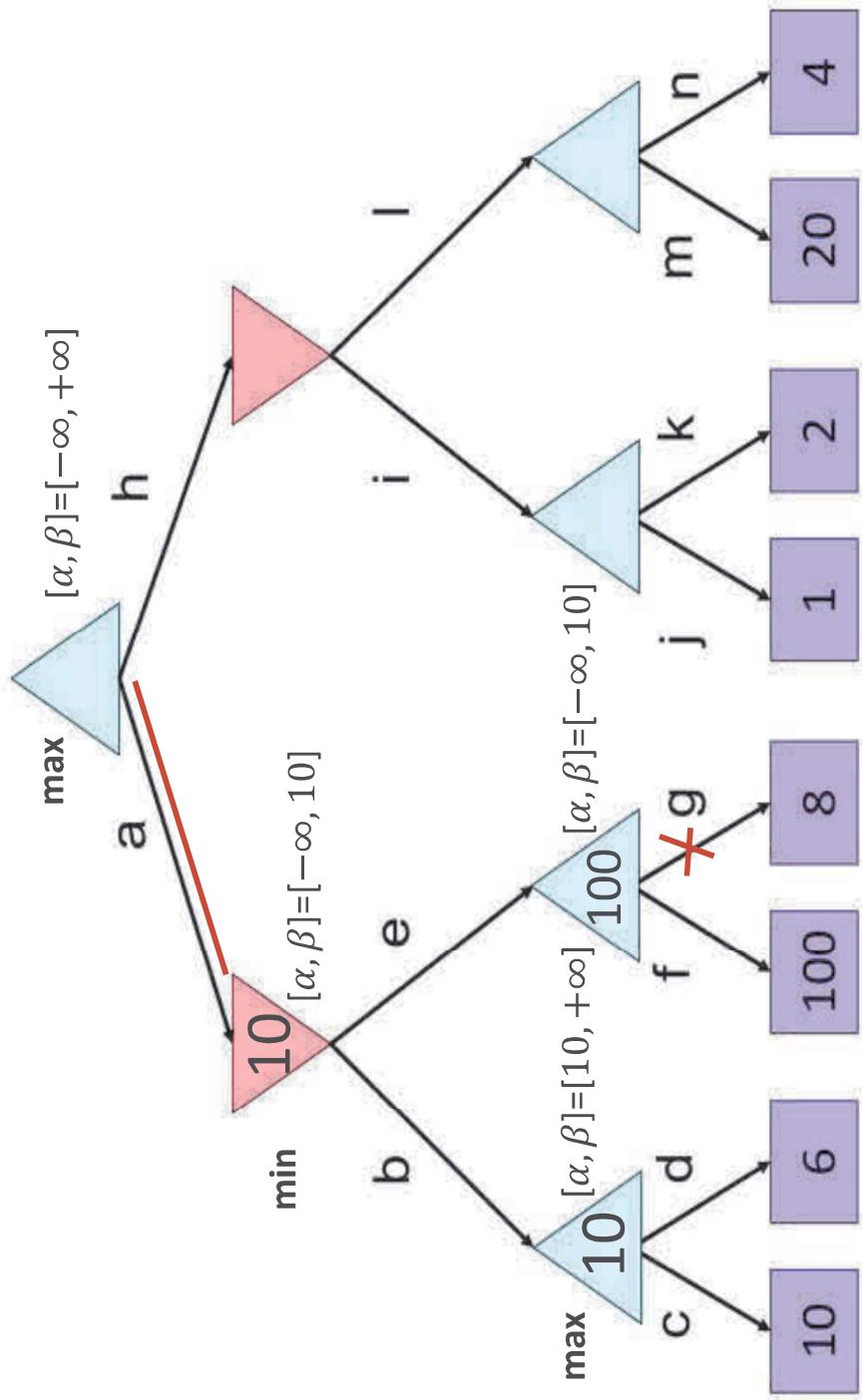
Alpha-Beta Quiz 2



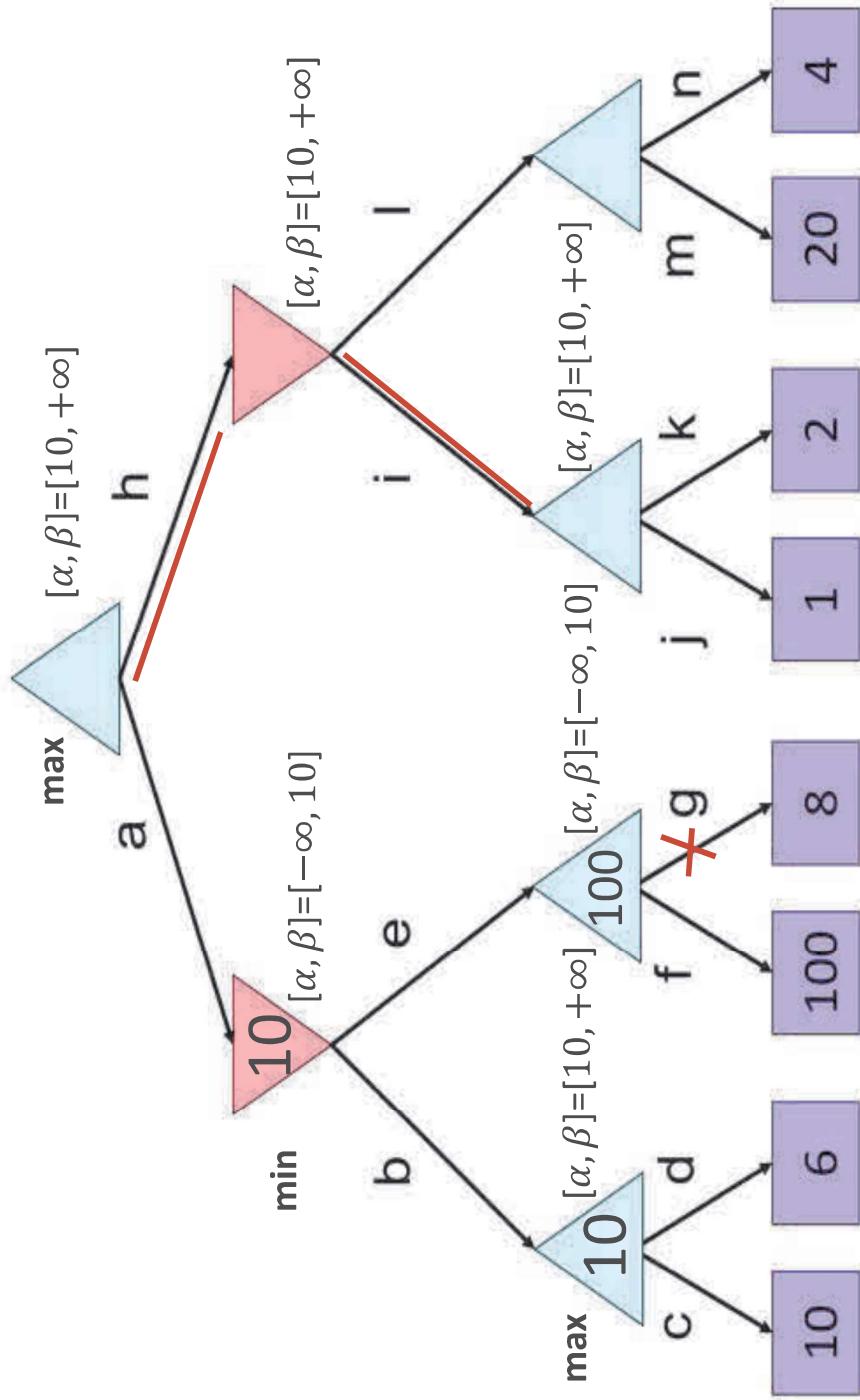
Alpha-Beta Quiz 2



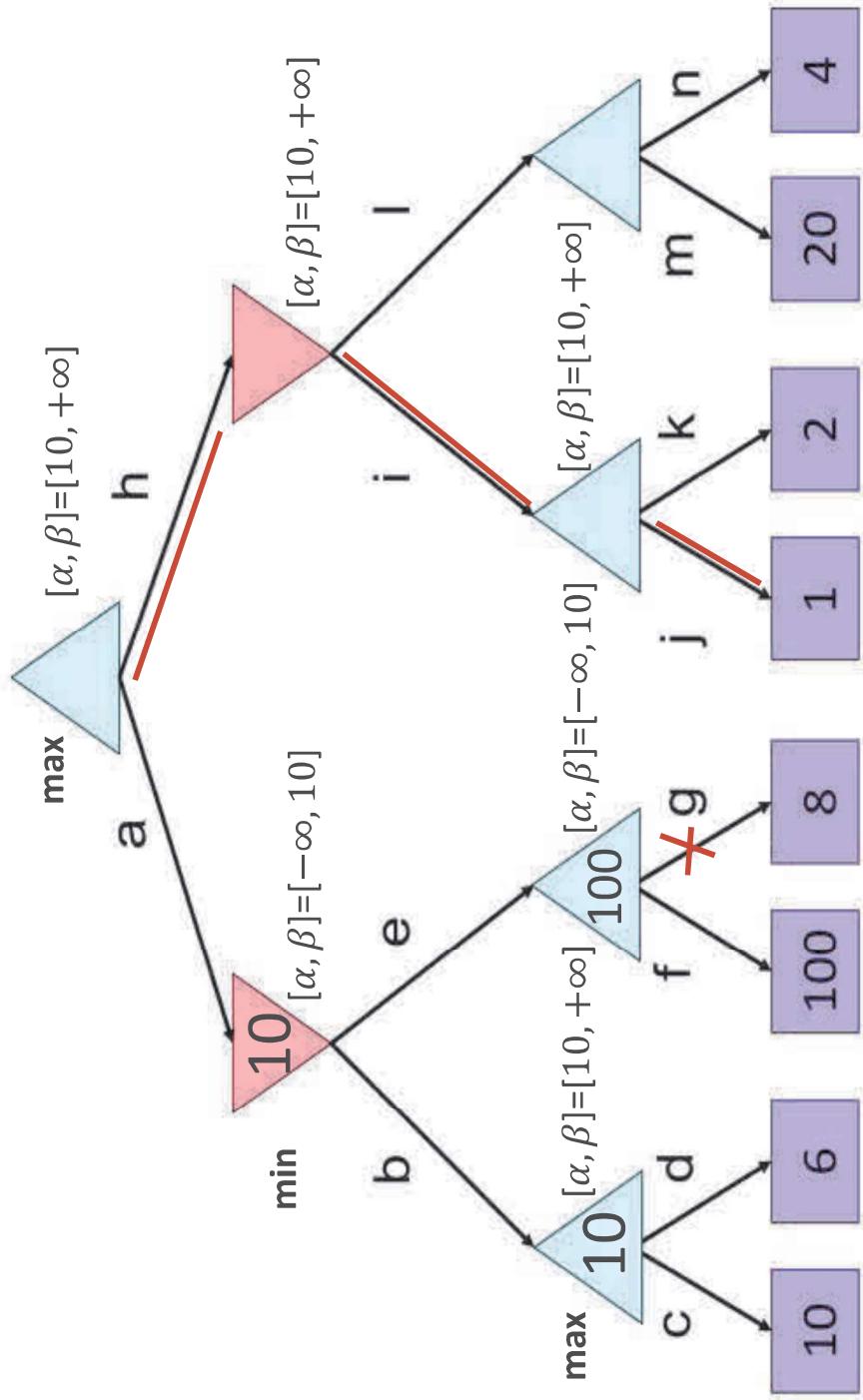
Alpha-Beta Quiz 2



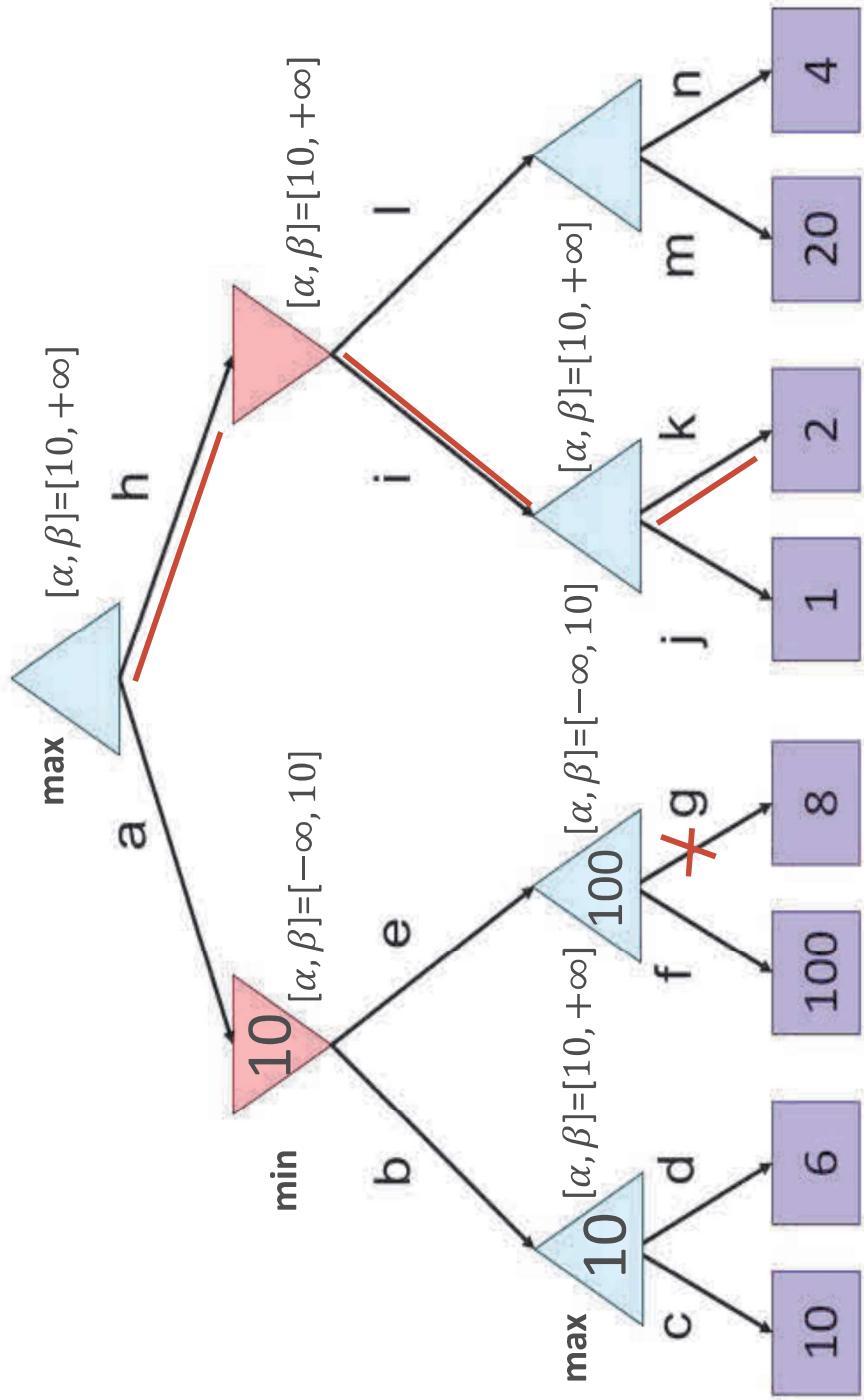
Alpha-Beta Quiz 2



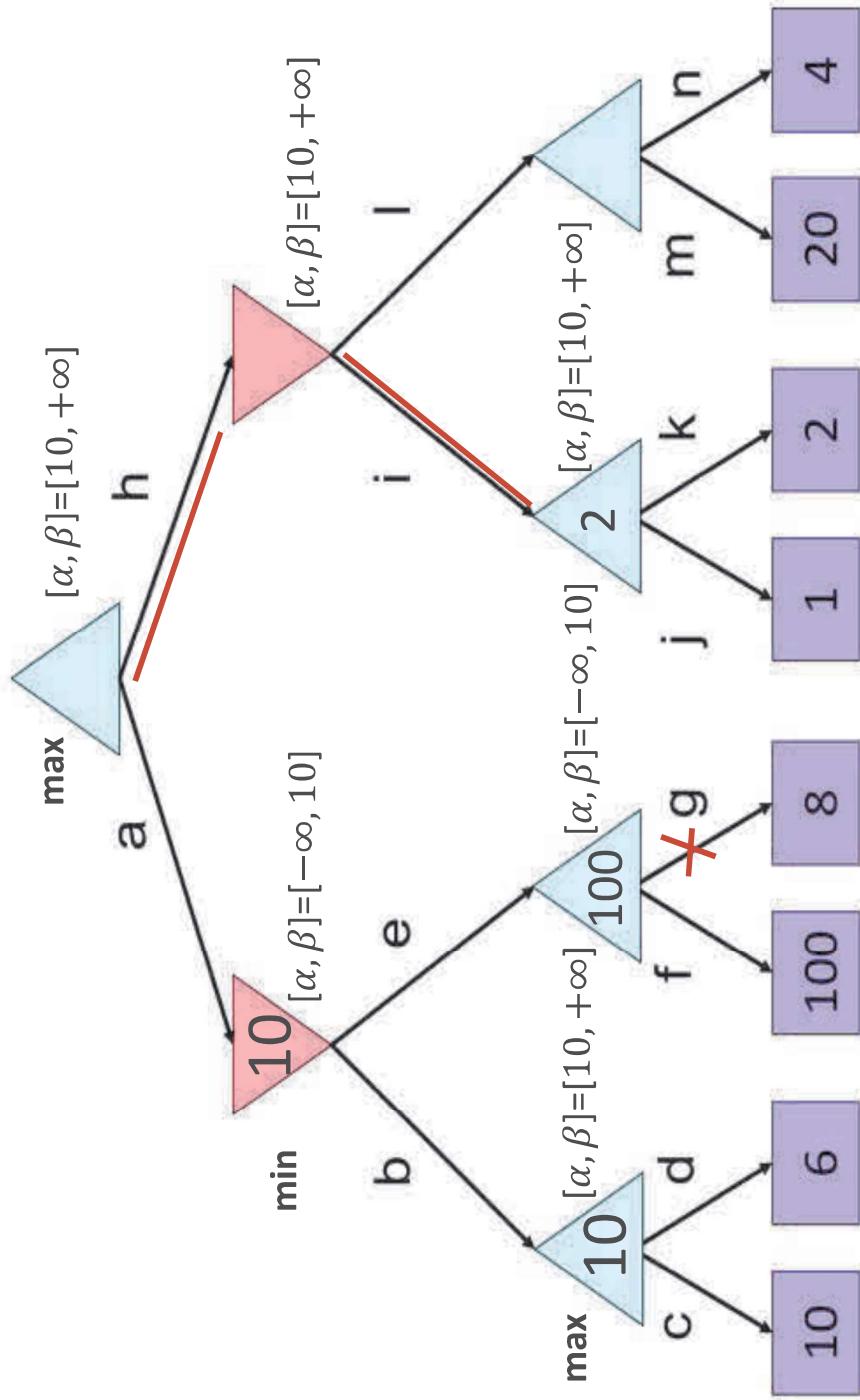
Alpha-Beta Quiz 2



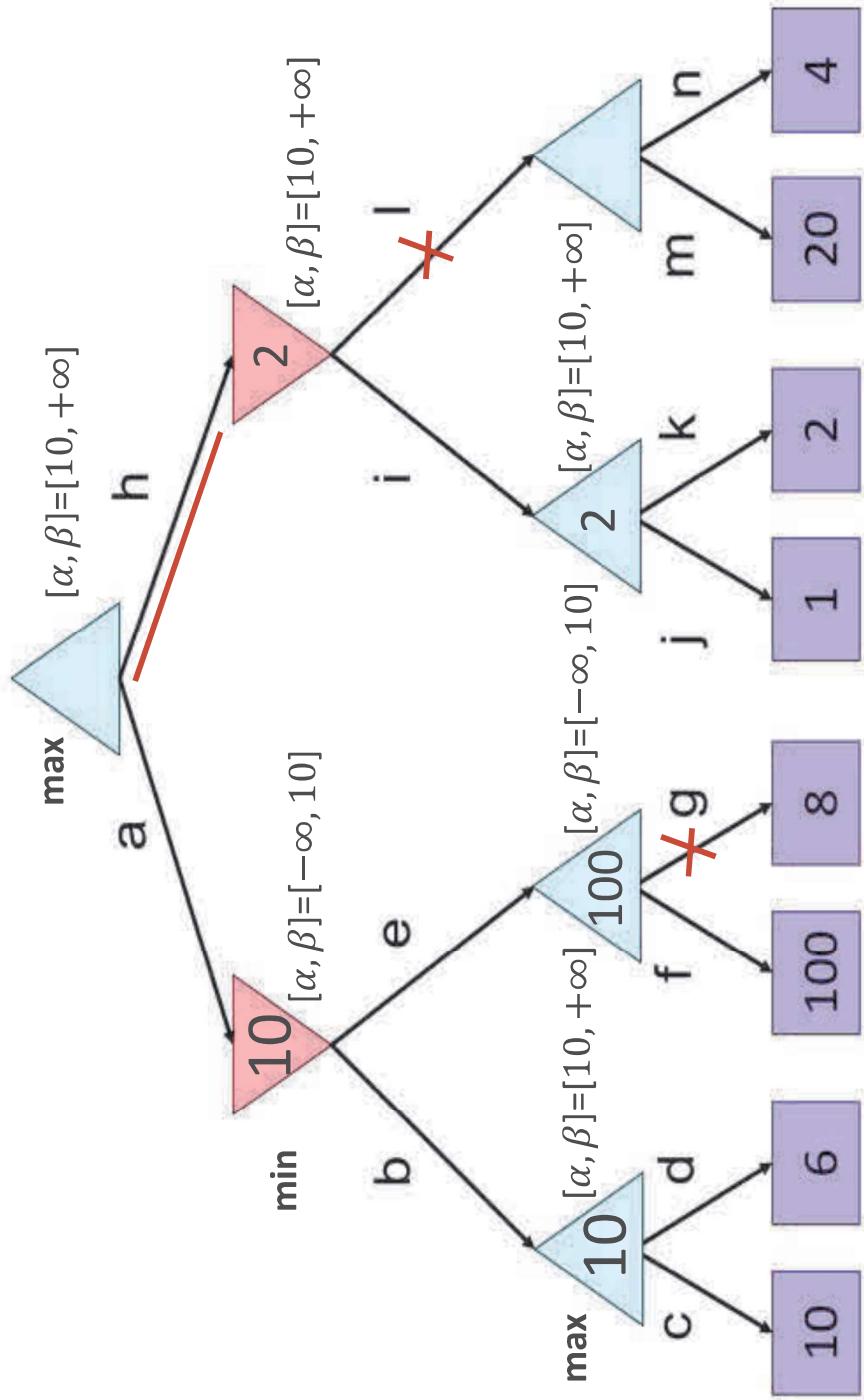
Alpha-Beta Quiz 2



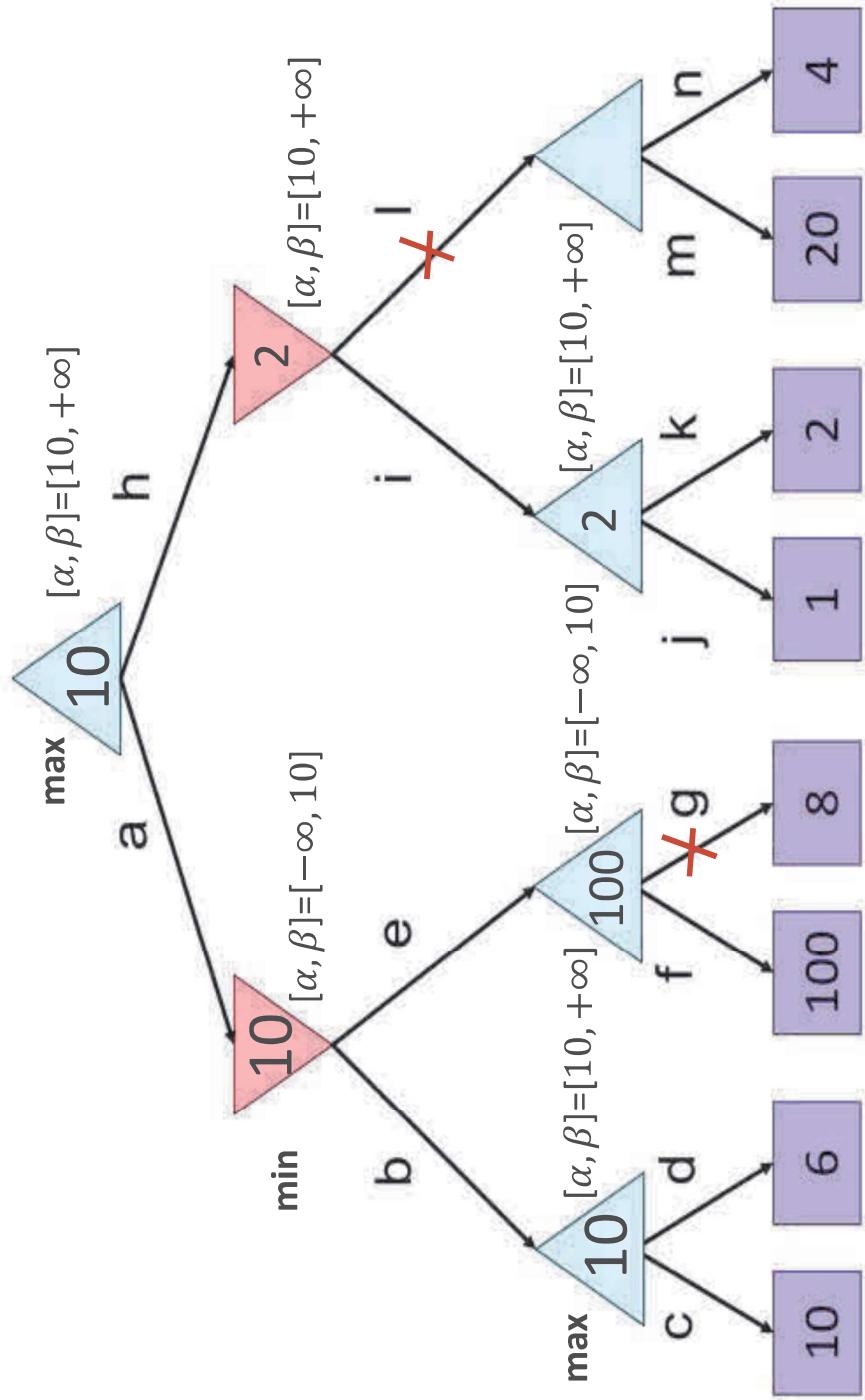
Alpha-Beta Quiz 2



Alpha-Beta Quiz 2

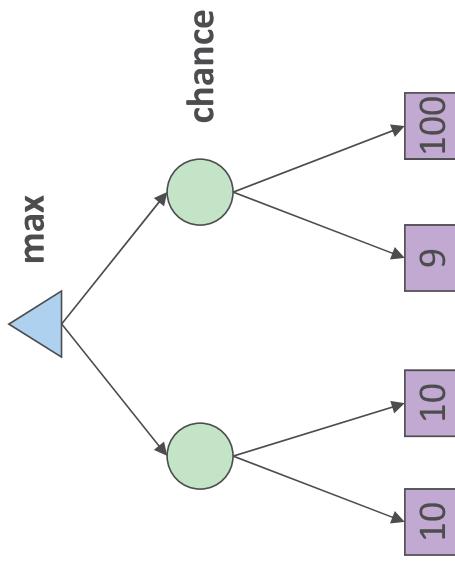


Alpha-Beta Quiz 2



Expectimax Search

- Why wouldn't we know what the result of an action will be?
 - Explicit randomness: rolling dice
 - Unpredictable opponents: the ghosts respond randomly
 - Actions can fail: when moving a robot, wheels might slip
- Values should now reflect average-case (expectimax) outcomes, not worst-case (minimax) outcomes
- **Expectimax search:** compute the average score under optimal play
 - Max nodes as in minimax search
 - Chance nodes are like min nodes but the outcome is uncertain
 - Calculate their **expected utilities**
 - I.e. take weighted average (expectation) of children
- Later, we'll learn how to formalize the underlying uncertain-result problems as **Markov Decision Processes**



[Demo: min vs exp (I-7D1,2)]

Expectimax Pseudocode

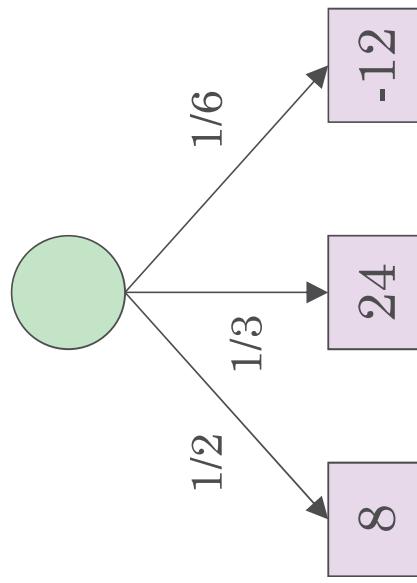
```
def value(state):
    if the state is a terminal state: return the state's utility
    if the next agent is MAX: return max-value(state)
    if the next agent is EXP: return exp-value(state)
```

```
def max-value(state):
    initialize v = -∞
    for each successor of state:
        v = max(v, value(successor))
    return v
```

```
def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v
```

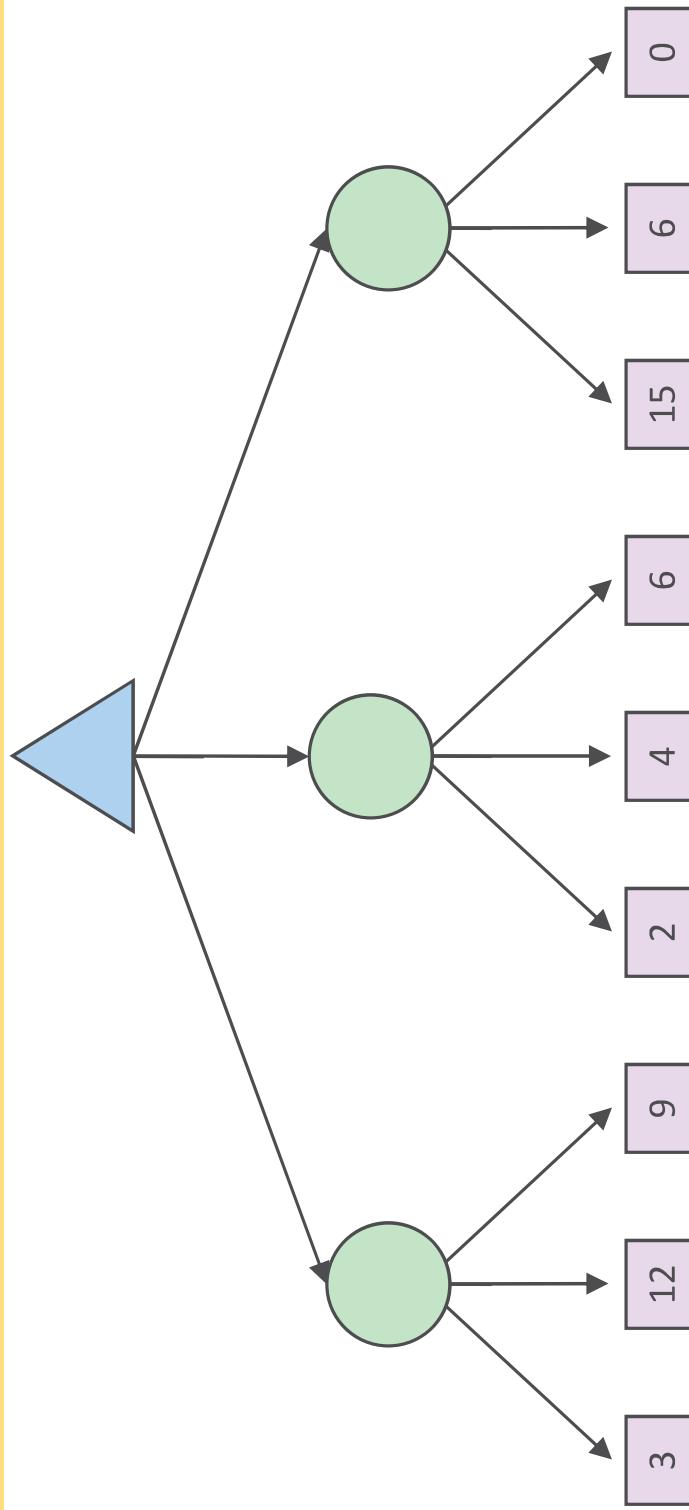
Expectimax Pseudocode

```
def exp-value(state):
    initialize v = 0
    for each successor of state:
        p = probability(successor)
        v += p * value(successor)
    return v
```

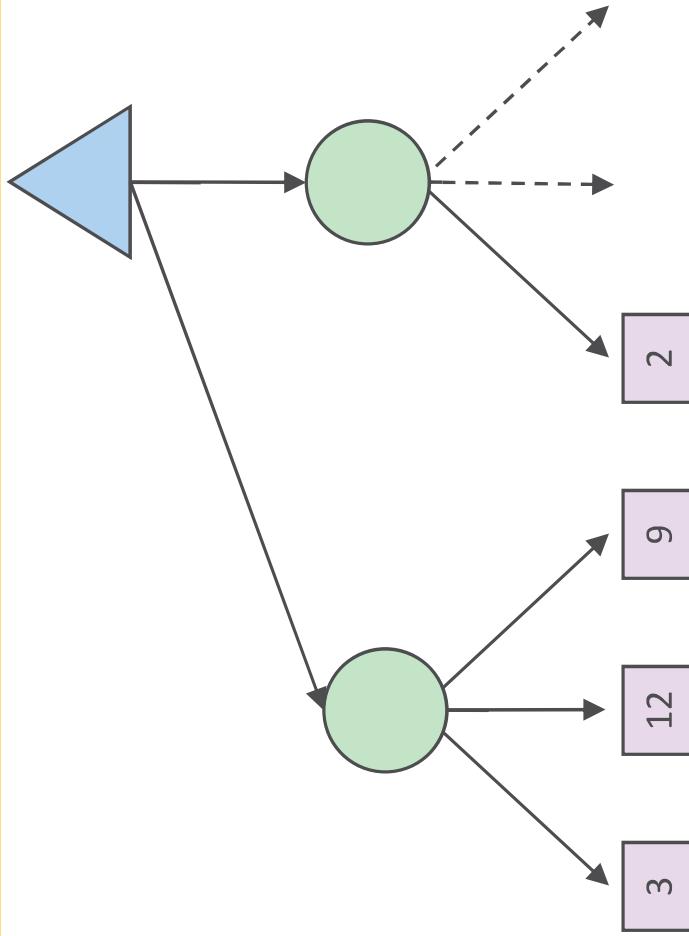


$$v = (1/2)(8) + (1/3)(24) + (1/6)(-12) = 10$$

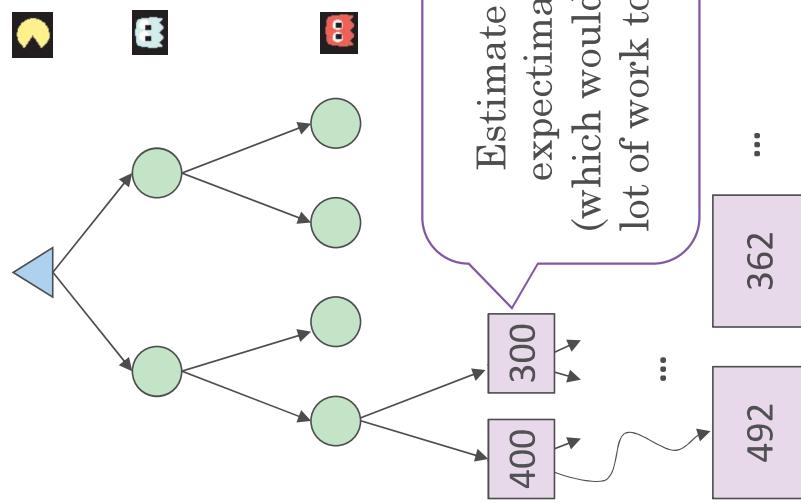
Expectimax Example



Expectimax Pruning?

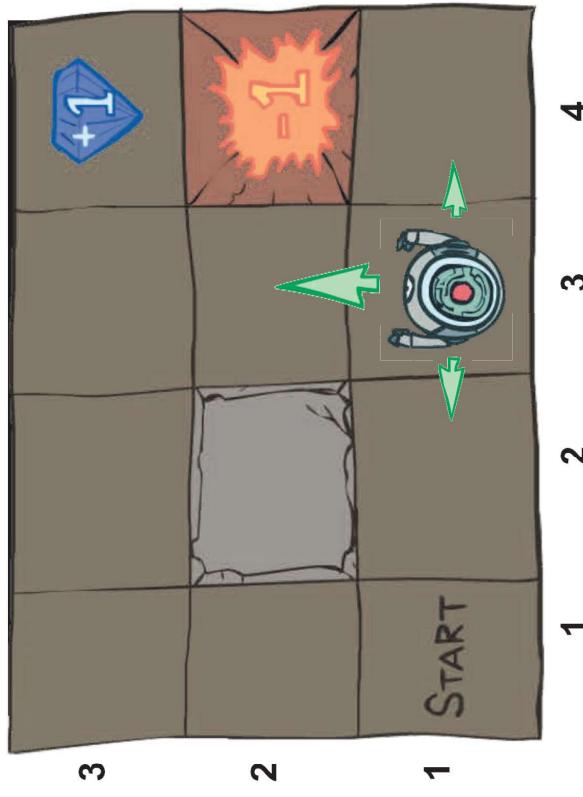


Depth-Limited Expectimax



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**
- 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

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

=

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

Andrey Markov
(1856-1922)

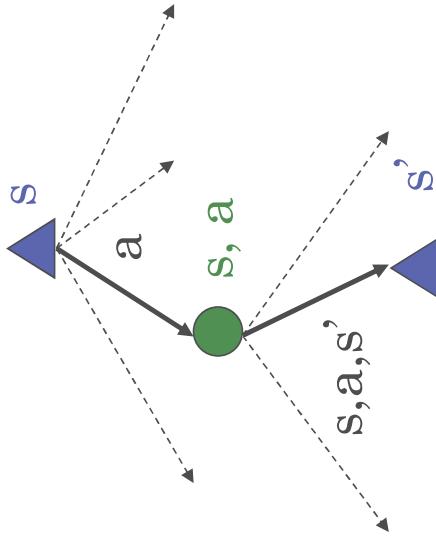
- This is just like search, where the successor function could only depend on the current state (not the history)



Recap: Defining MDPs

- Markov decision processes:

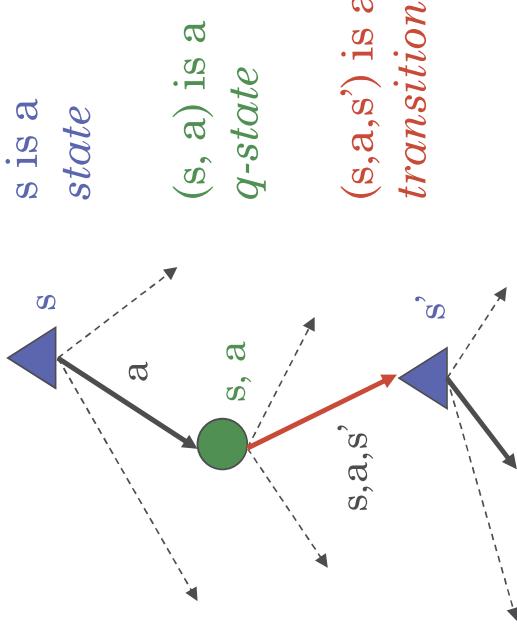
- Set of states S
- Start state s_0
- Set of actions A
- Transitions $P(s' | s, a)$ (or $T(s, a, s')$)
- Rewards $R(s, a, s')$ (and discount γ)



- MDP quantities so far:
 - Policy = Choice of action for each state
 - Utility = sum of (discounted) rewards

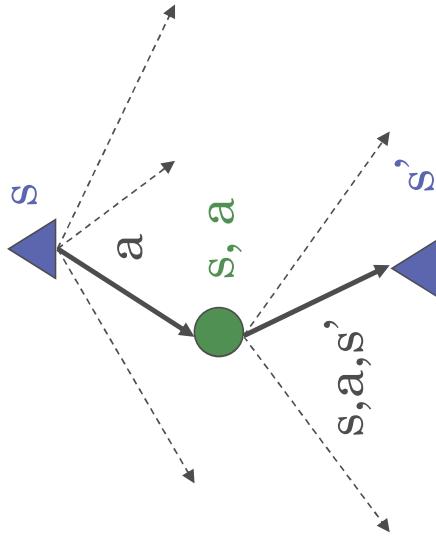
Optimal Quantities

- **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) :**
 $Q^*(s,a) = \text{expected utility starting out having taken action } a \text{ from state } s \text{ and (thereafter) acting optimally}$
- **The optimal policy:**
 $\pi^*(s)$ = optimal action from state s



Values of States

- 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:

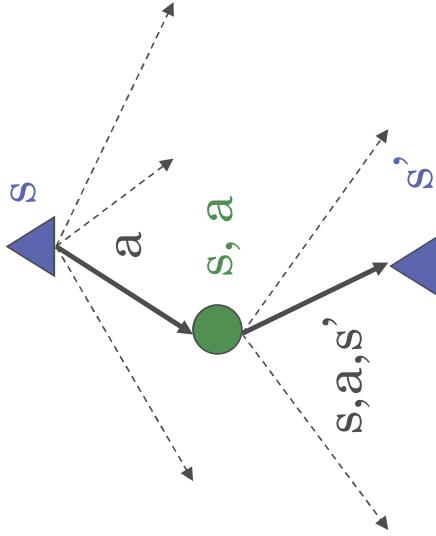
$$V^*(s) = \max_a Q^*(s, a)$$

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

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

Recap: MDPs

- Markov decision processes:
 - States S
 - Actions A
 - Transitions $P(s' | s, a)$ (or $T(s, a, s')$)
 - Rewards $R(s, a, s')$ (and discount γ)
 - Start state s_0



- Quantities:
 - Policy = map of states to actions
 - Utility = sum of discounted rewards
 - Values = expected future utility from a state (max node)
 - Q-Values = expected future utility from a q-state (chance node)

CS 471/571 (Fall 2023): Introduction to Artificial Intelligence

Lecture 11: MDPs (Part 2)

Thanh H. Nguyen

Source: <http://ai.berkeley.edu/home.html>

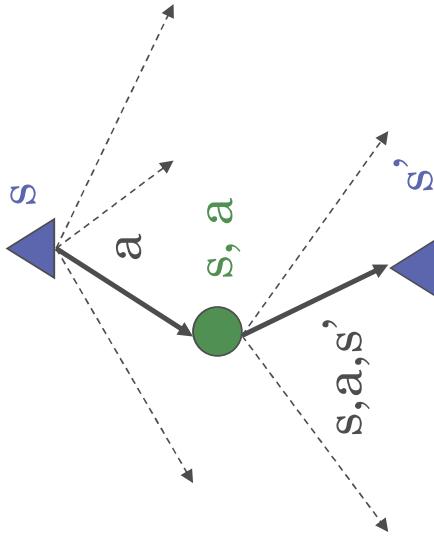
Announcement

- Project 2: Multi-agent Search

- Deadline: Nov 03, 2023

Recap: MDPs

- Markov decision processes:
 - States S
 - Actions A
 - Transitions $P(s' | s, a)$ (or $T(s, a, s')$)
 - Rewards $R(s, a, s')$ (and discount γ)
 - Start state s_0



- Quantities:
 - Policy = map of states to actions
 - Utility = sum of discounted rewards
 - Values = expected future utility from a state (max node)
 - Q-Values = expected future utility from a q-state (chance node)

Optimal Quantities

- The value (utility) of a state s :

$V^*(s)$ = expected utility starting in s and acting optimally

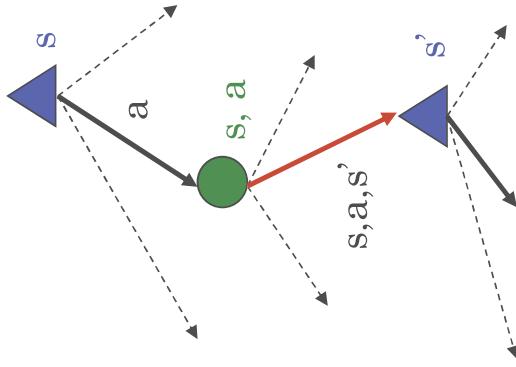
s is a state

- 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

(s, a) is a q-state

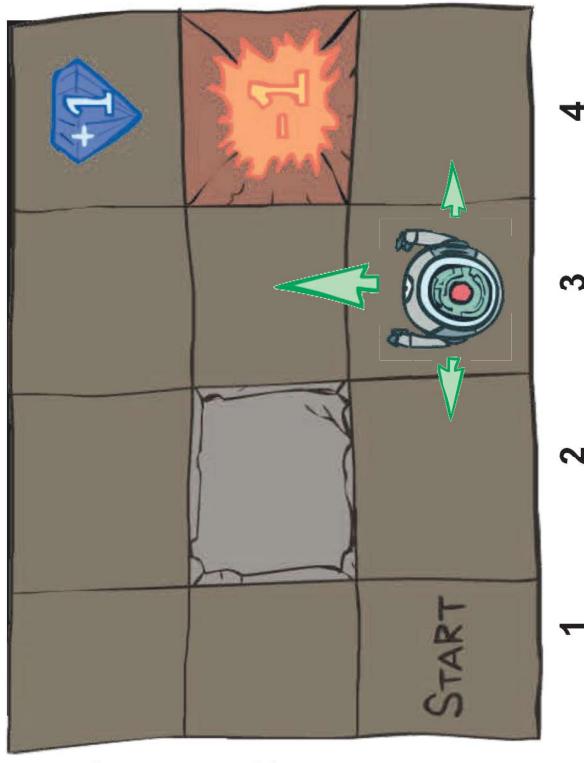
(s,a,s') is a transition



- The optimal policy:
- $\pi^*(s)$ = optimal action from state s

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
 - 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 (discounted) rewards

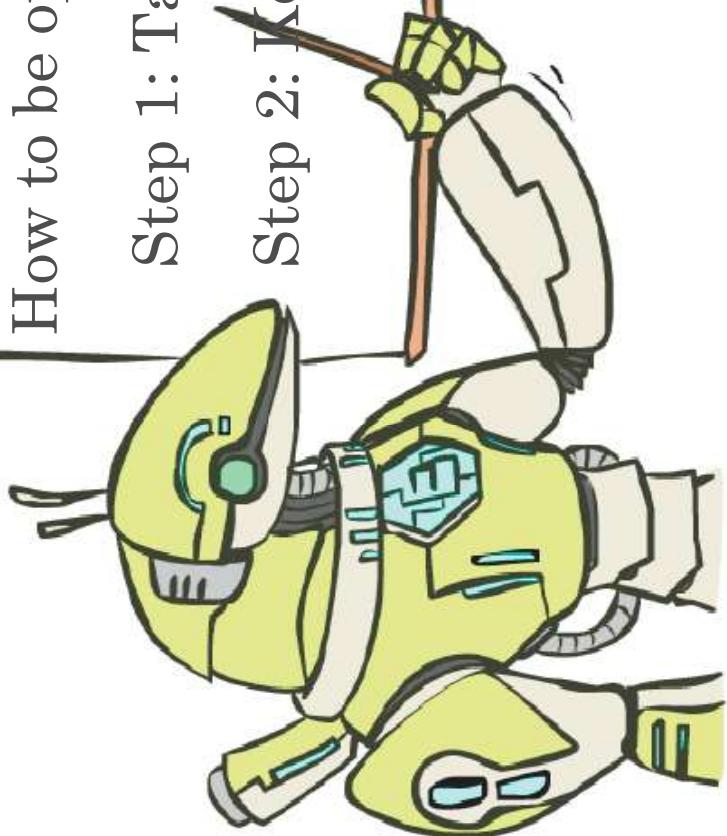


The Bellman Equations

How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal



The Bellman Equations

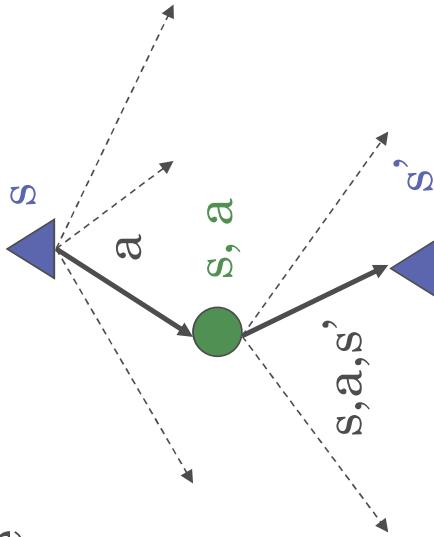
- Definition of “optimal utility” via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

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

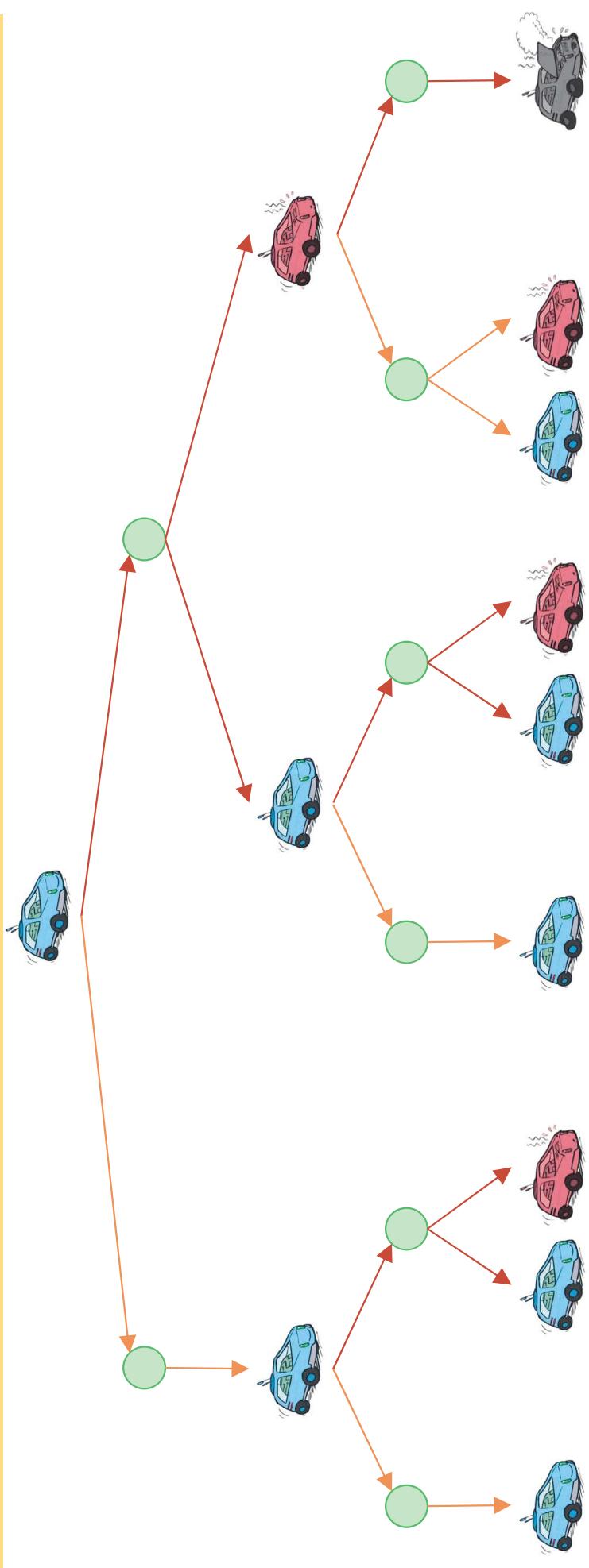
$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

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

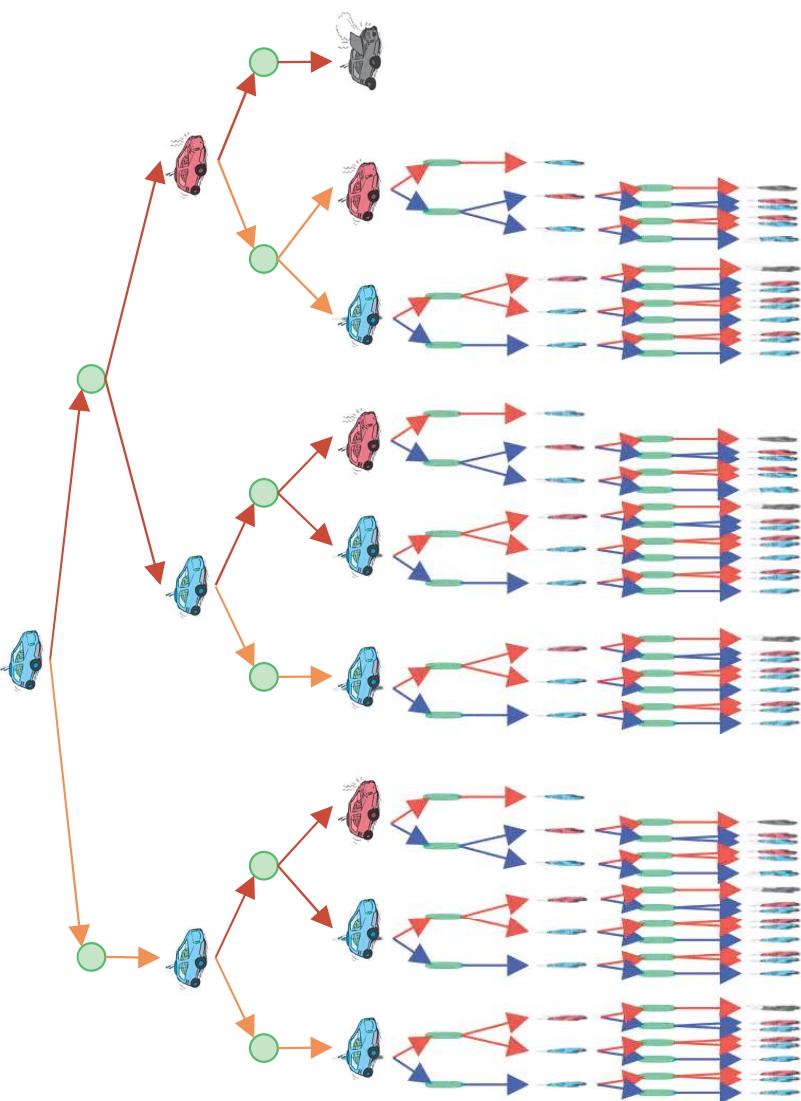
- These are the Bellman equations, and they characterize optimal values in a way we'll use over and over



Racing Search Tree



Racing Search Tree



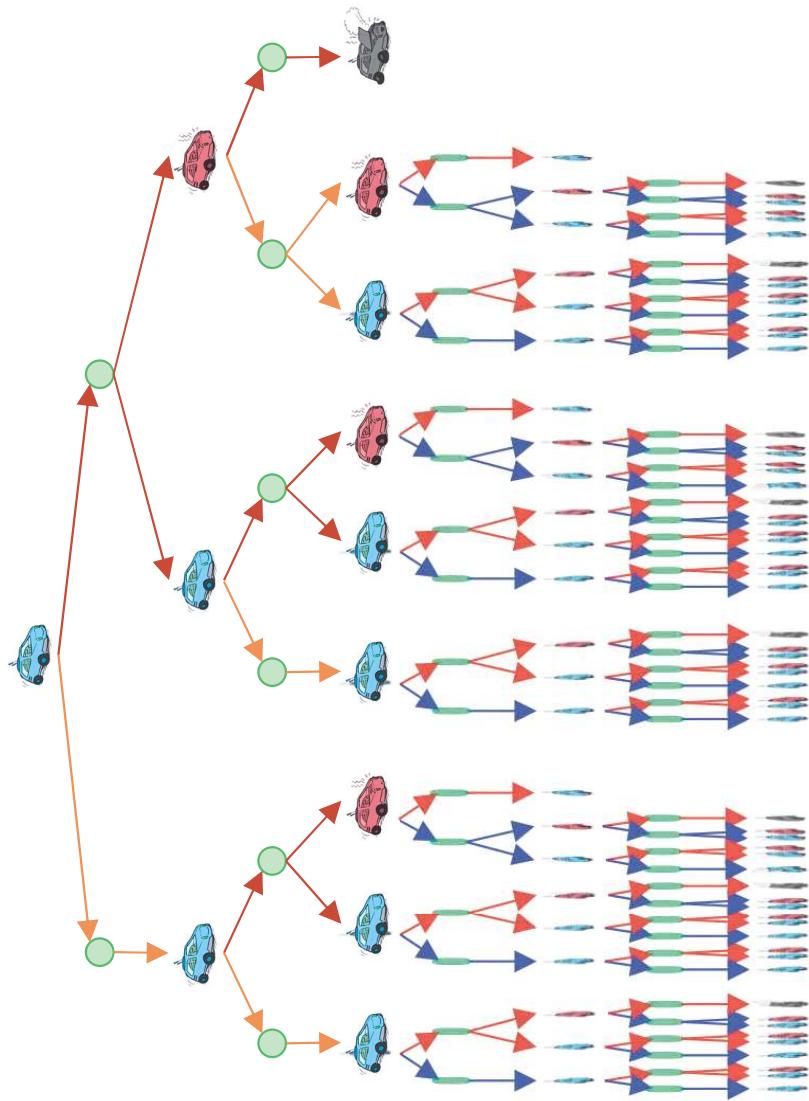
Racing Search Tree

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

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

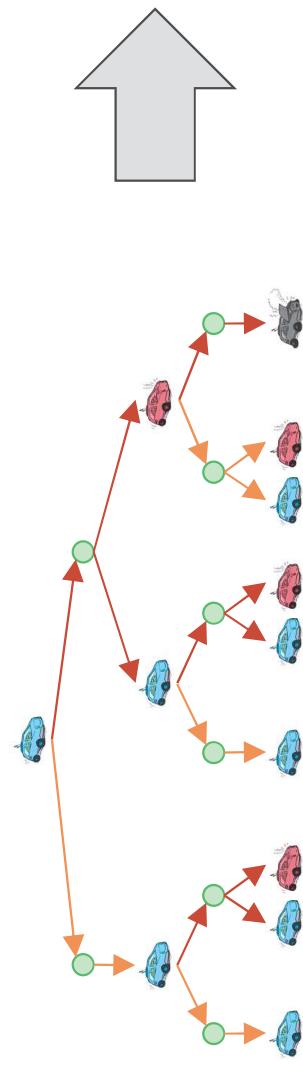
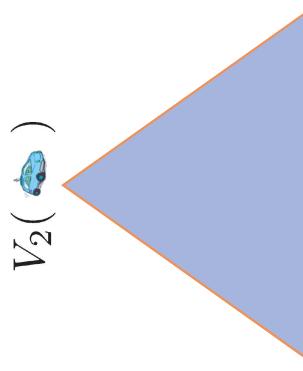
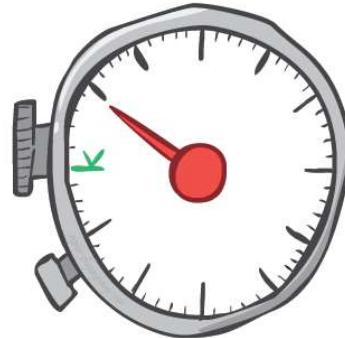
- Problem: Tree goes on forever
 - Idea: Do a depth-limited computation, but with increasing depths until change is small

- Note: deep parts of the tree eventually don't matter if $\gamma < 1$



Time-Limited Values

- Key idea: time-limited values
- Define $V_k(s)$ to be the optimal value of s if the game ends in k more time steps
 - Equivalently, it's what a depth- k expectimax would give from s



$k=0$

Gridworld Display

VALUES AFTER 0 ITERATIONS			
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00

Noise = 0.2

Discount = 0.9

Living reward = 0

$k=1$

Gridworld Display

		0.00	0.00	1.00	-1.00	0.00	0.00	0.00
		0.00	0.00	0.00	0.00	0.00	0.00	0.00
		0.00	0.00	0.00	0.00	0.00	0.00	0.00
		0.00	0.00	0.00	0.00	0.00	0.00	0.00
		0.00	0.00	0.00	0.00	0.00	0.00	0.00

VALUES AFTER 1 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

k=2

Gridworld Display

0.00	0.00	0.72	1.00	-1.00	0.00	0.00	0.00
0.00	0.00	0.72	1.00	-1.00	0.00	0.00	0.00
0.00	0.00	0.72	1.00	-1.00	0.00	0.00	0.00
0.00	0.00	0.72	1.00	-1.00	0.00	0.00	0.00
0.00	0.00	0.72	1.00	-1.00	0.00	0.00	0.00

VALUES AFTER 2 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

$k=3$

Gridworld Display

OOO

Orange line

0.00	0.52	0.78	1.00	-1.00	0.43	0.00	0.00	0.00
0.00					0.43			
0.00						0.00		
0.00							0.00	
0.00								0.00

VALUES AFTER 3 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0



Green line

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	0.00
0.00	0.00	0.00	0.00

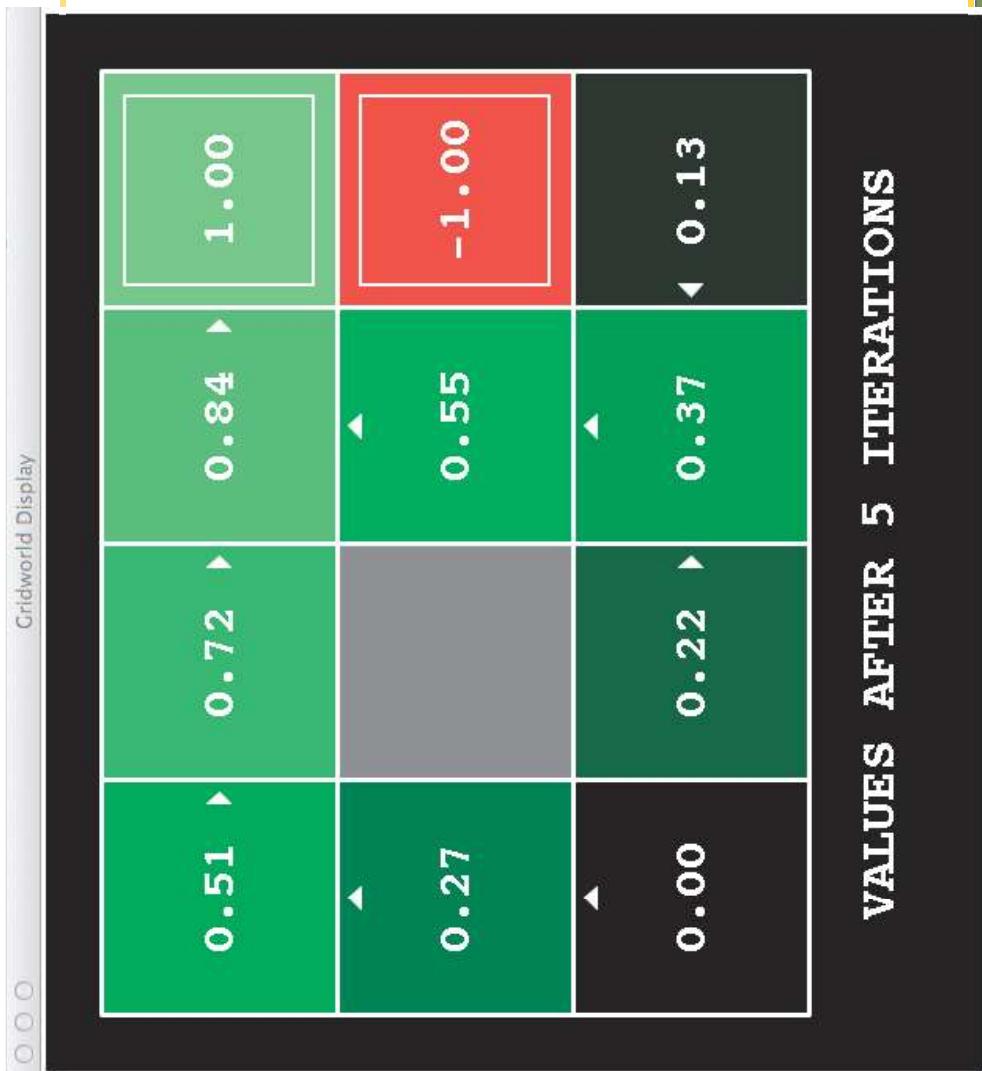
VALUES AFTER 4 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

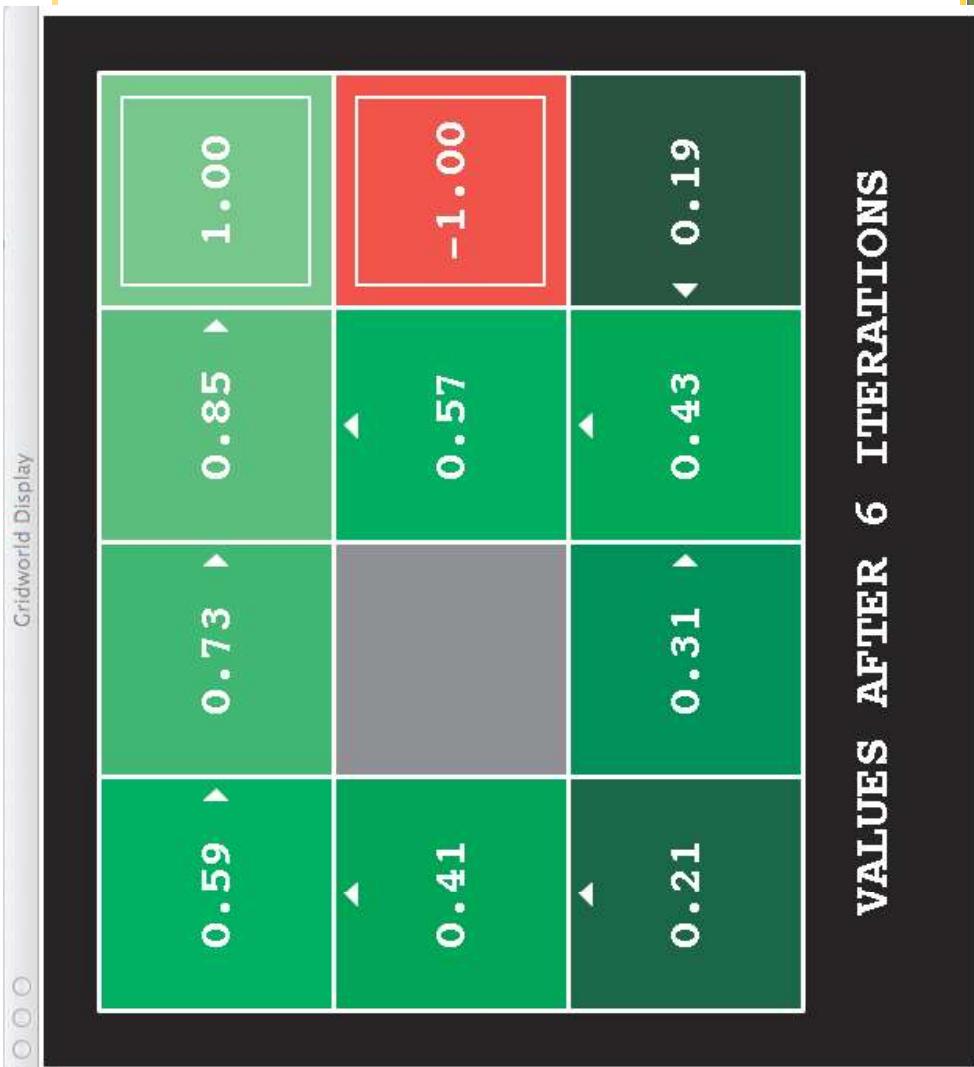
$k=5$



Noise = 0.2
Discount = 0.9
Living reward = 0

$k=6$

Gridworld Display

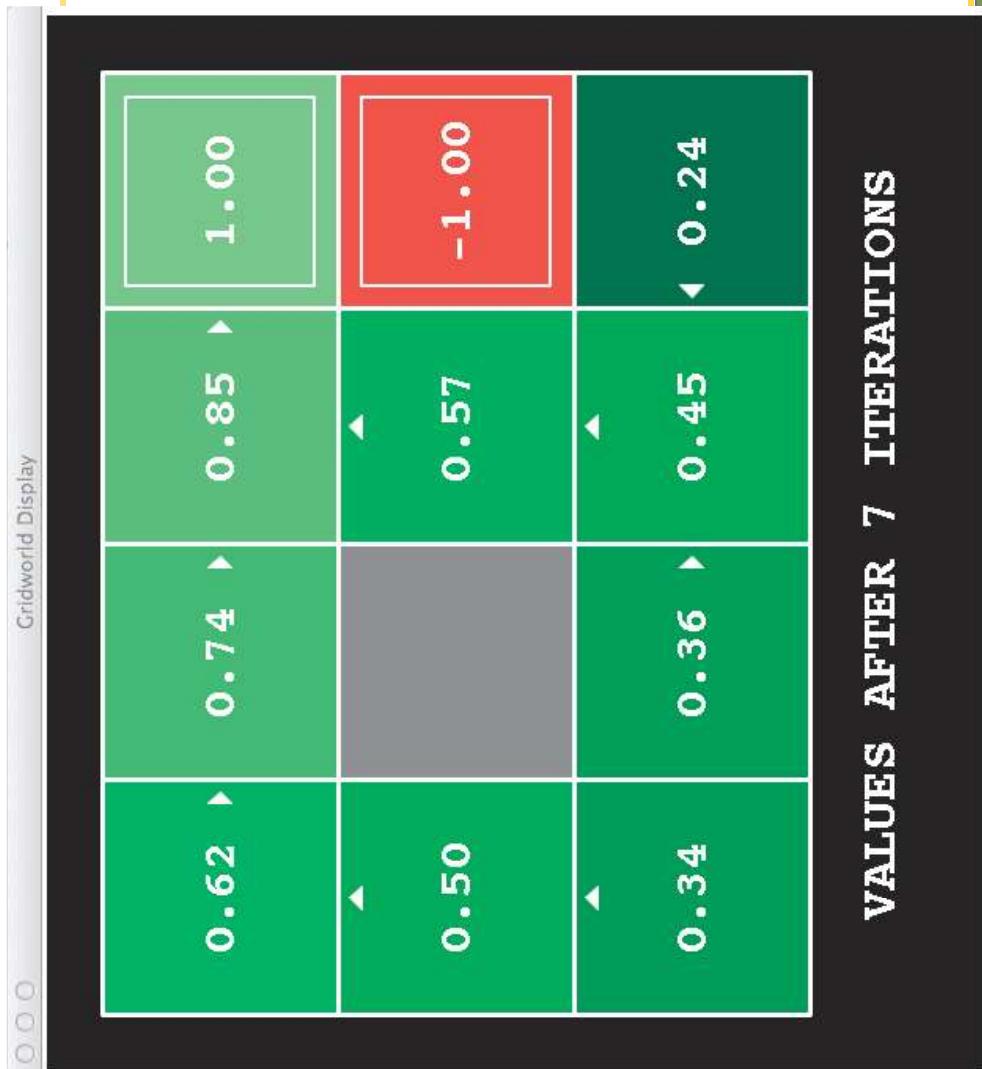


Noise = 0.2

Discount = 0.9

Living reward = 0

$k=7$



Noise = 0.2

Discount = 0.9

Living reward = 0

$k=8$

Gridworld Display

○○○

0.63 ▶	0.74 ▶	0.85 ▶	1.00
▲		0.57	-1.00
0.53		▲	
▲	0.39 ▶	0.46 ▶	0.26
0.42			

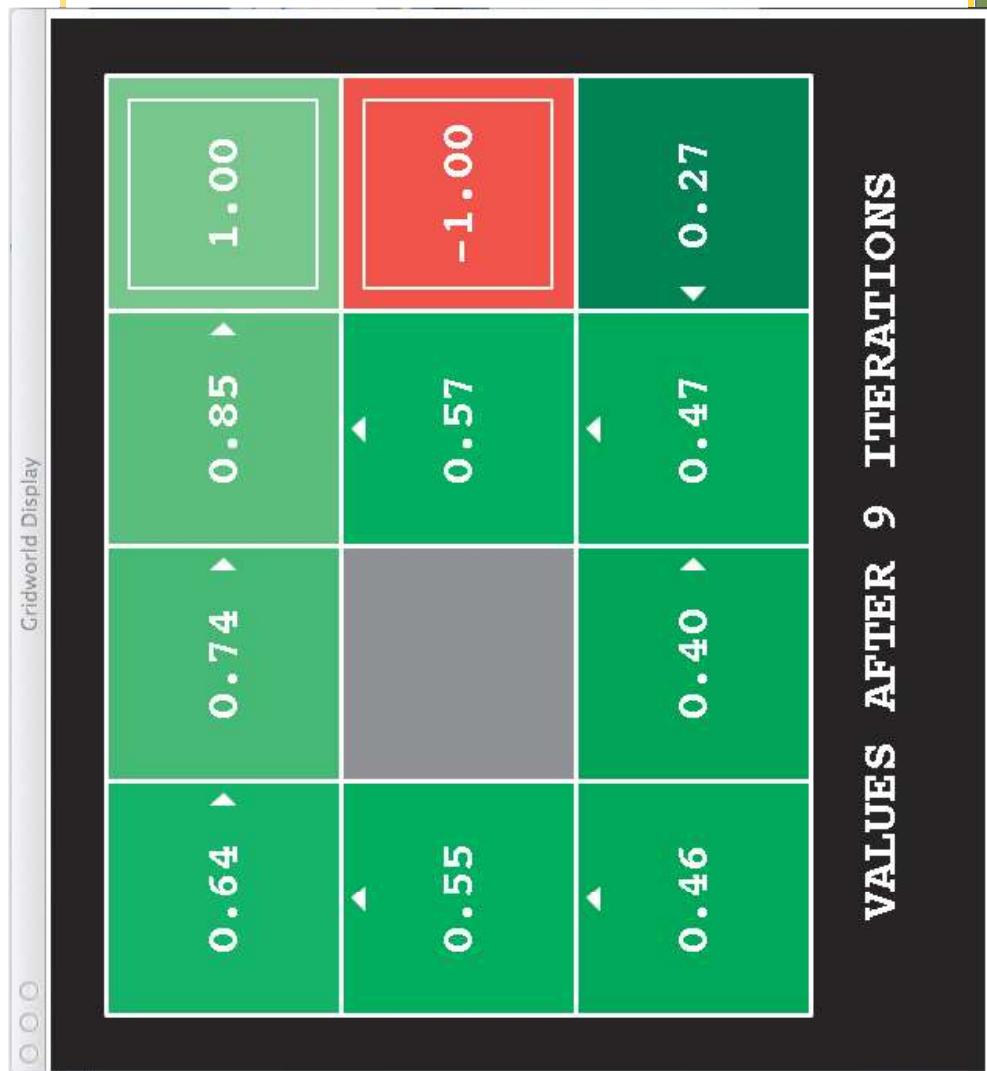
VALUES AFTER 8 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

k=9



Noise = 0.2
Discount = 0.9
Living reward = 0

k=10

Gridworld Display

0.64 ▶	0.74 ▶	0.85 ▶	1.00
◀ 0.56		0.57	-1.00
◀ 0.48		0.47	◀ 0.27

VALUES AFTER 10 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

k=11

Gridworld Display

0.64	0.74	0.85	1.00
0.56	0.57	-1.00	
0.48	0.42	0.47	0.27

VALUES AFTER 11 ITERATIONS

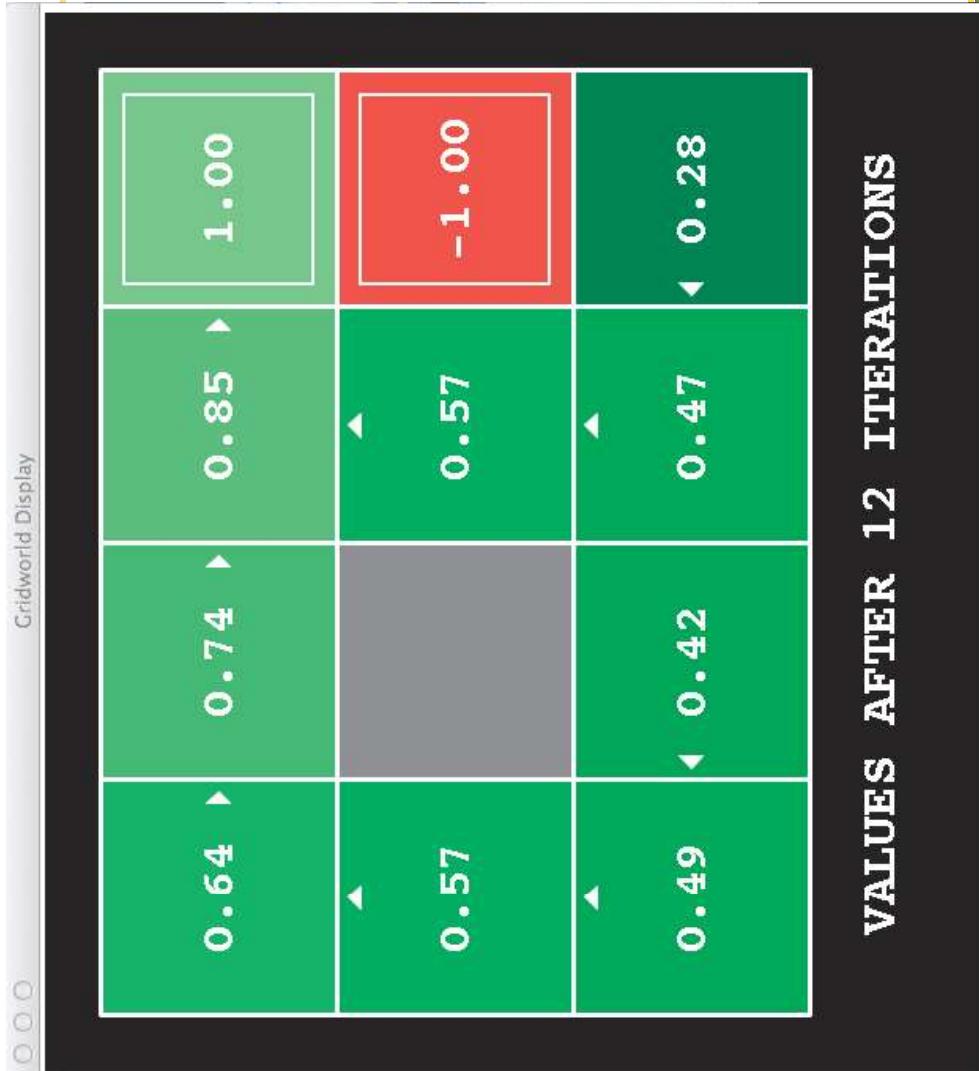
Noise = 0.2

Discount = 0.9

Living reward = 0

k=12

Gridworld Display

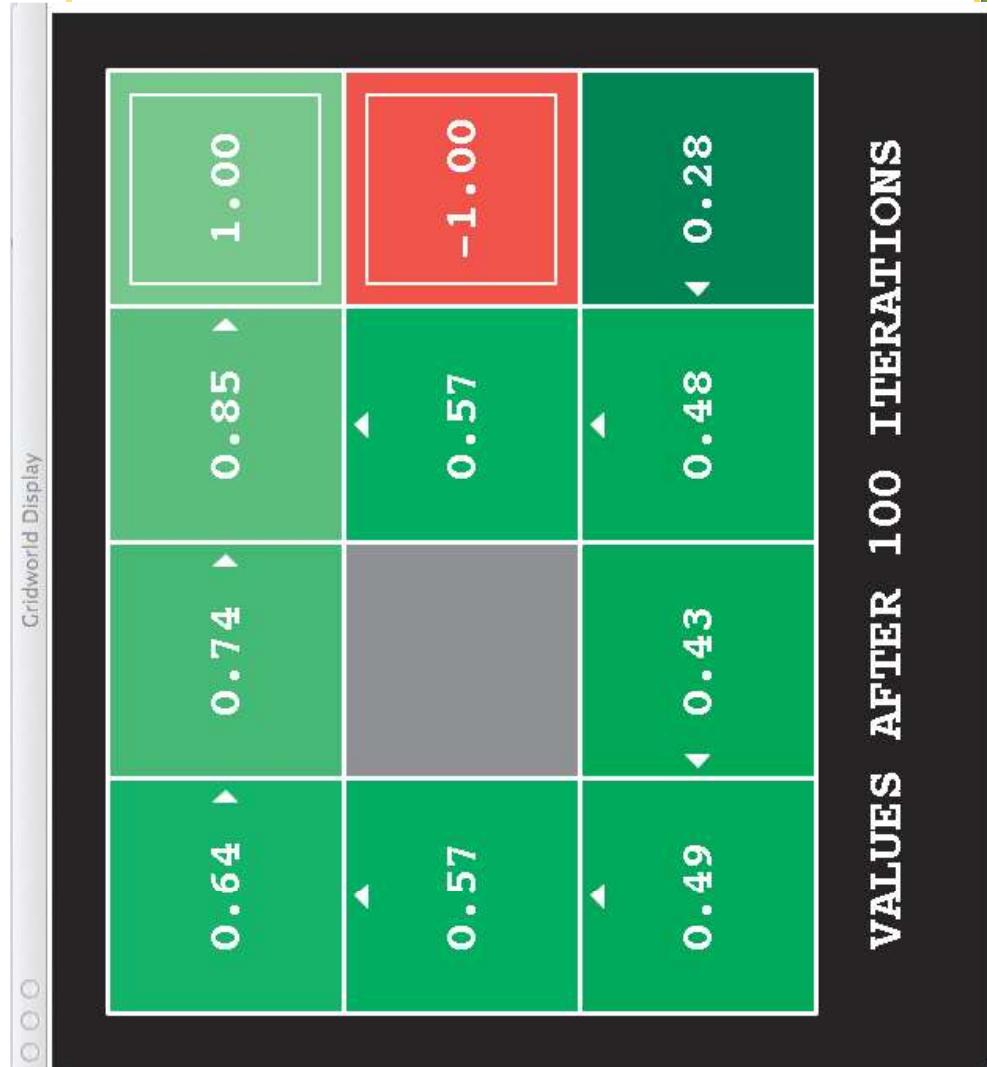


Noise = 0.2

Discount = 0.9

Living reward = 0

$k=100$

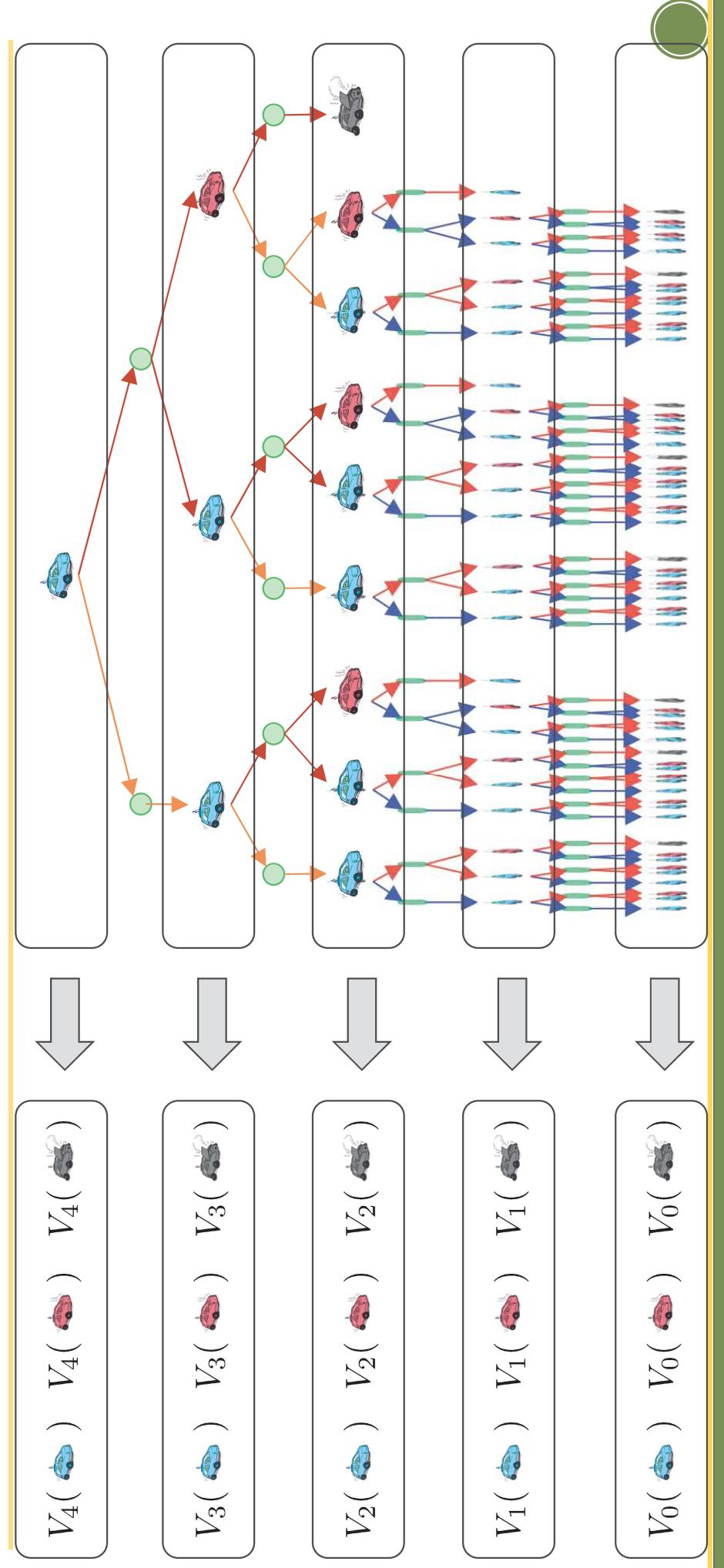


Noise = 0.2

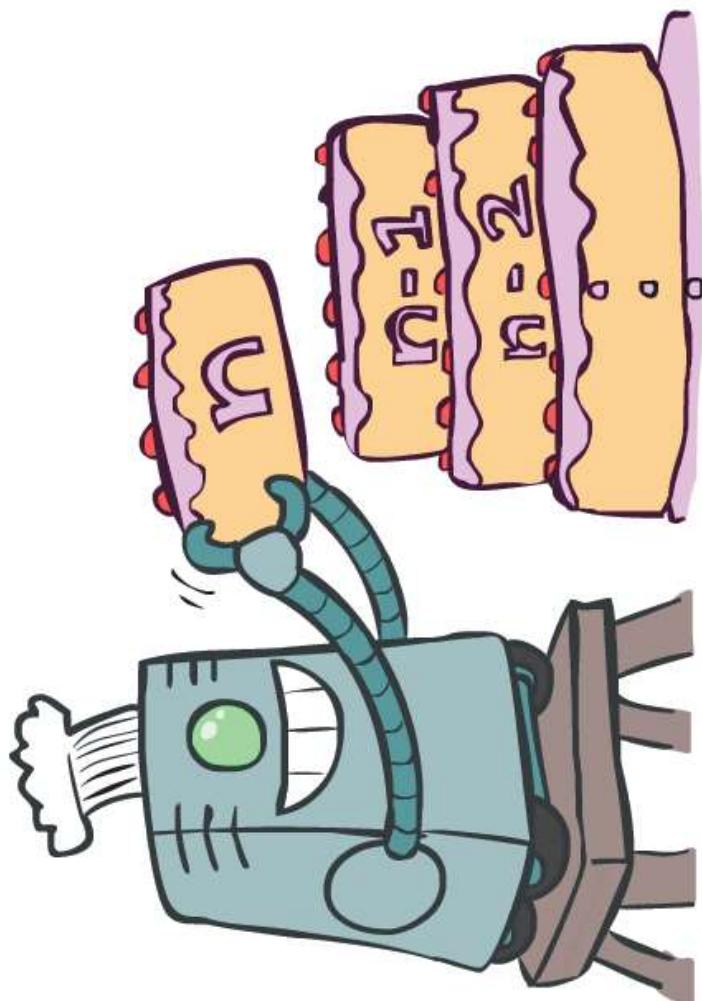
Discount = 0.9

Living reward = 0

Computing Time-Limited Values



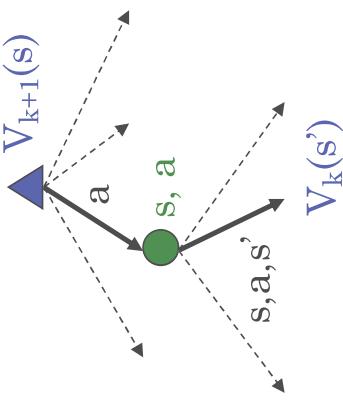
Value Iteration



Value Iteration

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

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

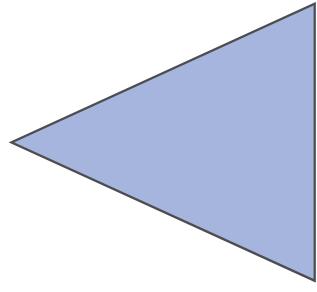


$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

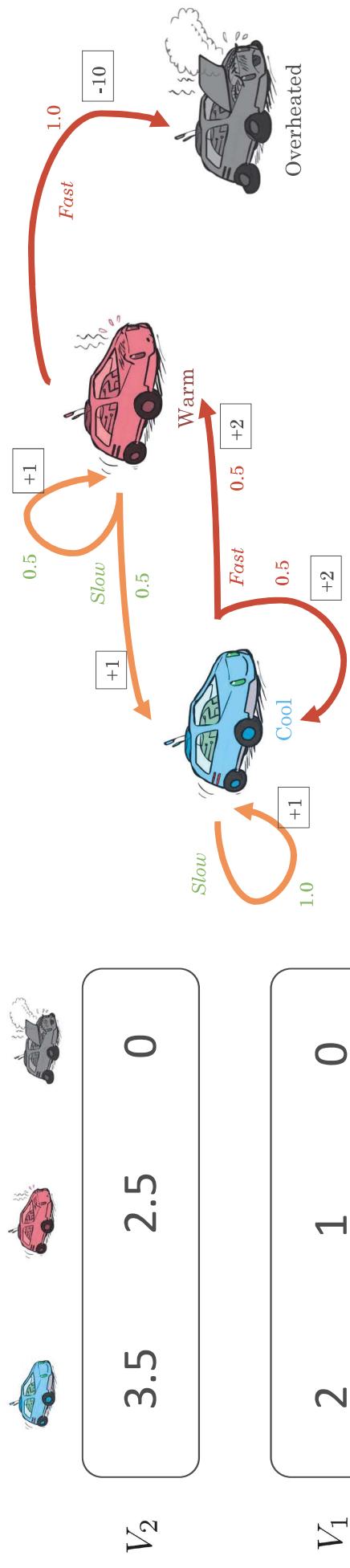
- Repeat until convergence

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

- Theorem: will converge to unique optimal values
 - Basic idea: approximations get refined towards optimal values
 - Policy may converge long before values do



Example: Value Iteration



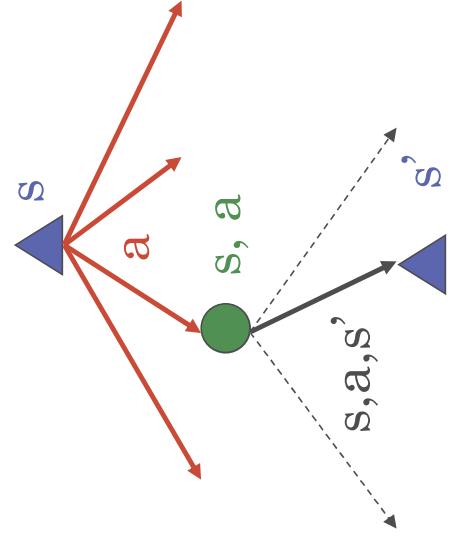
Assume no discount!

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

Problems with Value Iteration

- Value iteration repeats the Bellman updates:

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



- Problem 1: It's slow – $O(S^2A)$ per iteration
- Problem 2: The “max” at each state rarely changes

- Problem 3: The policy often converges long before the values

$k=0$

Gridworld Display

VALUES AFTER 0 ITERATIONS			
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00

Noise = 0.2

Discount = 0.9



Living reward = 0

$k=1$

Gridworld Display

		0.00	1.00
	0.00	0.00	-1.00
0.00		0.00	
0.00		0.00	
0.00		0.00	

VALUES AFTER 1 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

$$k=2$$

Gridworld Display

VALUES AFTER 2 ITERATIONS

Noise = 0.2

Discount = 0.9

Discount = 0.9

$k=3$

Gridworld Display

0 0 0

Orange line

0.00	0.52	0.78	1.00	-1.00	0.43	0.00	0.00	0.00
0.00					0.43			
0.00						0.00		
0.00							0.00	
0.00								0.00

VALUES AFTER 3 ITERATIONS

Noise = 0.2

Discount = 0.9



Living reward = 0

Green line

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 ◀	0.00 ◀
0.00 ◀	0.00 ◀	0.00 ◀	0.00 ◀

VALUES AFTER 4 ITERATIONS

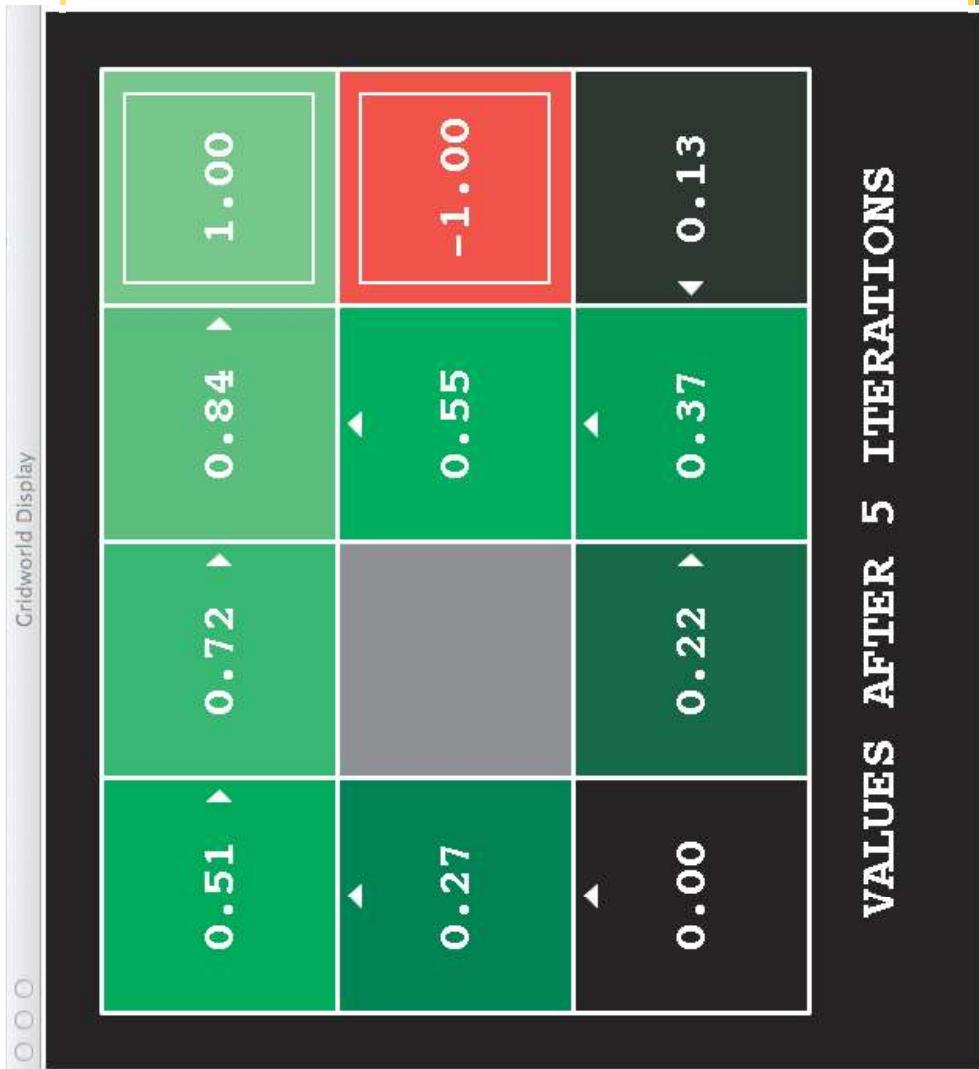
Noise = 0.2

Discount = 0.9

Living reward = 0

$k=5$

Gridworld Display



VALUES AFTER 5 ITERATIONS

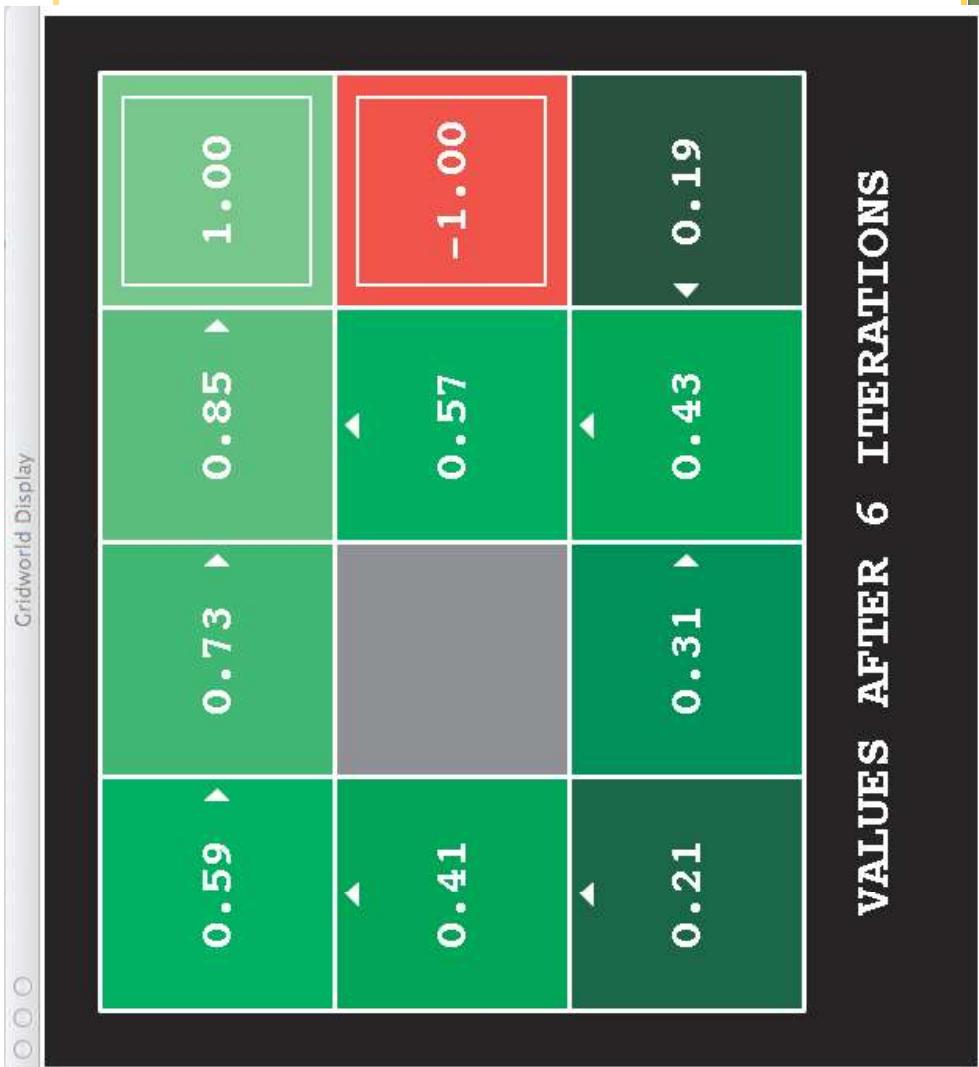
Noise = 0.2

Discount = 0.9

Living reward = 0

$k=6$

Gridworld Display

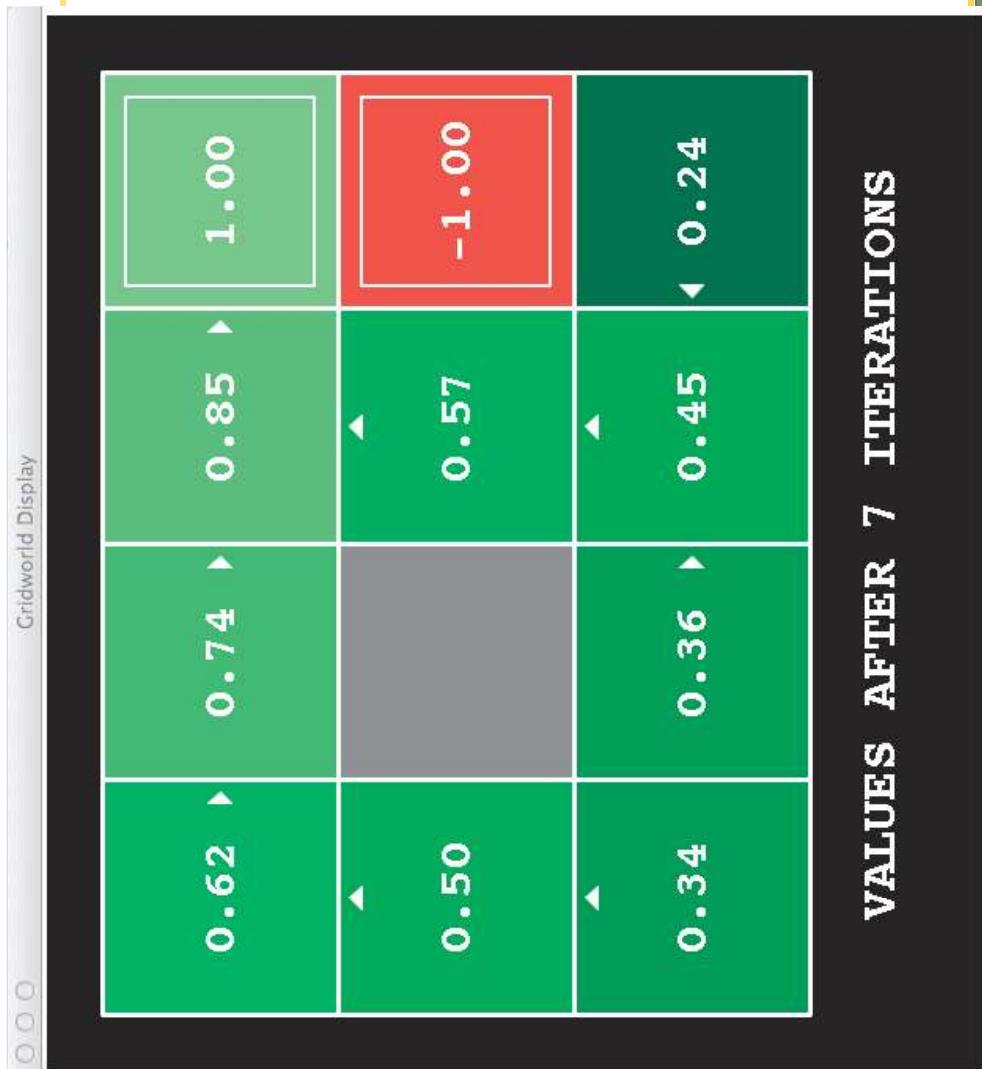


Noise = 0.2

Discount = 0.9

Living reward = 0

$k=7$



Noise = 0.2

Discount = 0.9

Living reward = 0

$k=8$

Gridworld Display

○○○



0.63 ▶	0.74 ▶	0.85 ▶	1.00
▲		0.57	-1.00
0.53		▲	
▲	0.39 ▶	0.46 ▶	0.26
0.42			

VALUES AFTER 8 ITERATIONS

Noise = 0.2

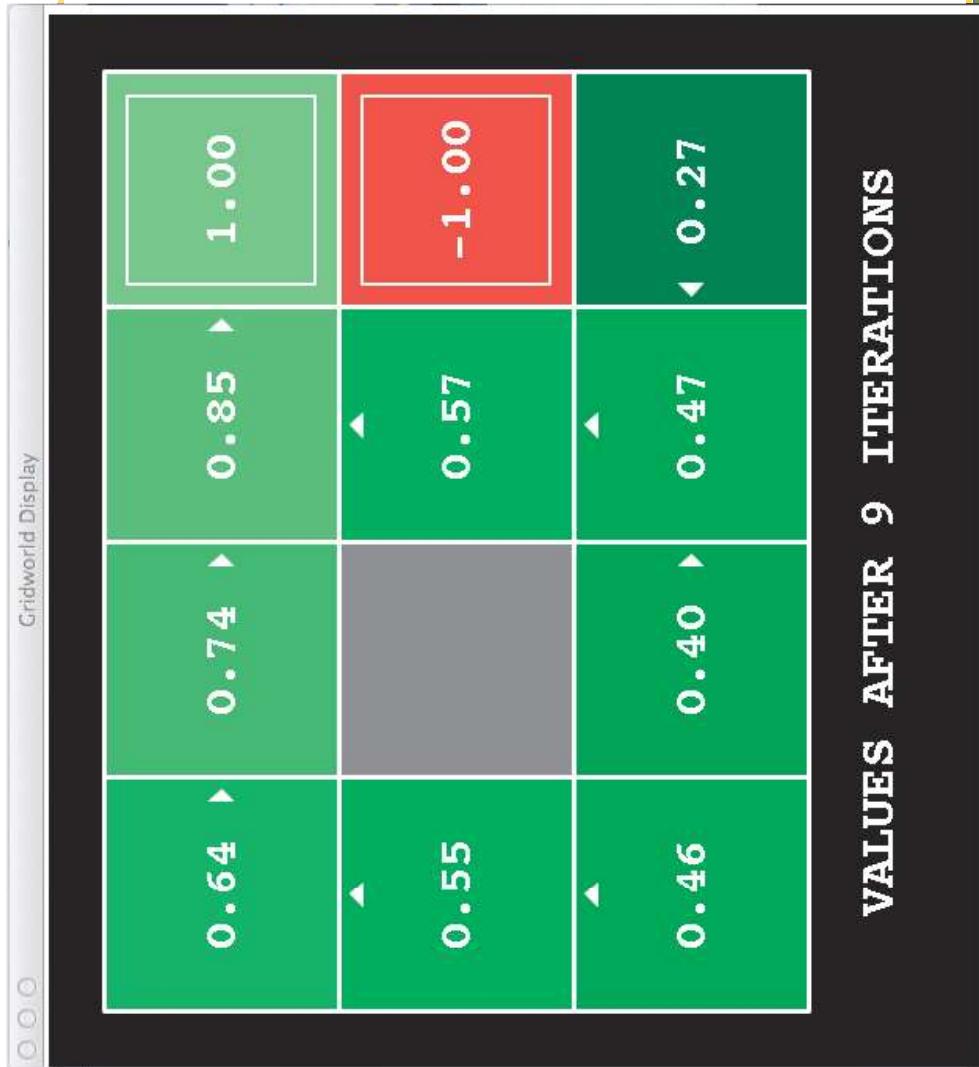
Discount = 0.9

Living reward = 0



$k=9$

Gridworld Display



VALUES AFTER 9 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

$k=10$

Gridworld Display

0.64 ▶	0.74 ▶	0.85 ▶	1.00
◀ 0.56		0.57	-1.00
◀ 0.48		0.47	◀ 0.27

VALUES AFTER 10 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

k=11

Gridworld Display

0.64	0.74	0.85	1.00
0.56	0.57	-1.00	
0.48	0.42	0.47	0.27

VALUES AFTER 11 ITERATIONS

Noise = 0.2

Discount = 0.9

Living reward = 0

$k=12$

Gridworld Display

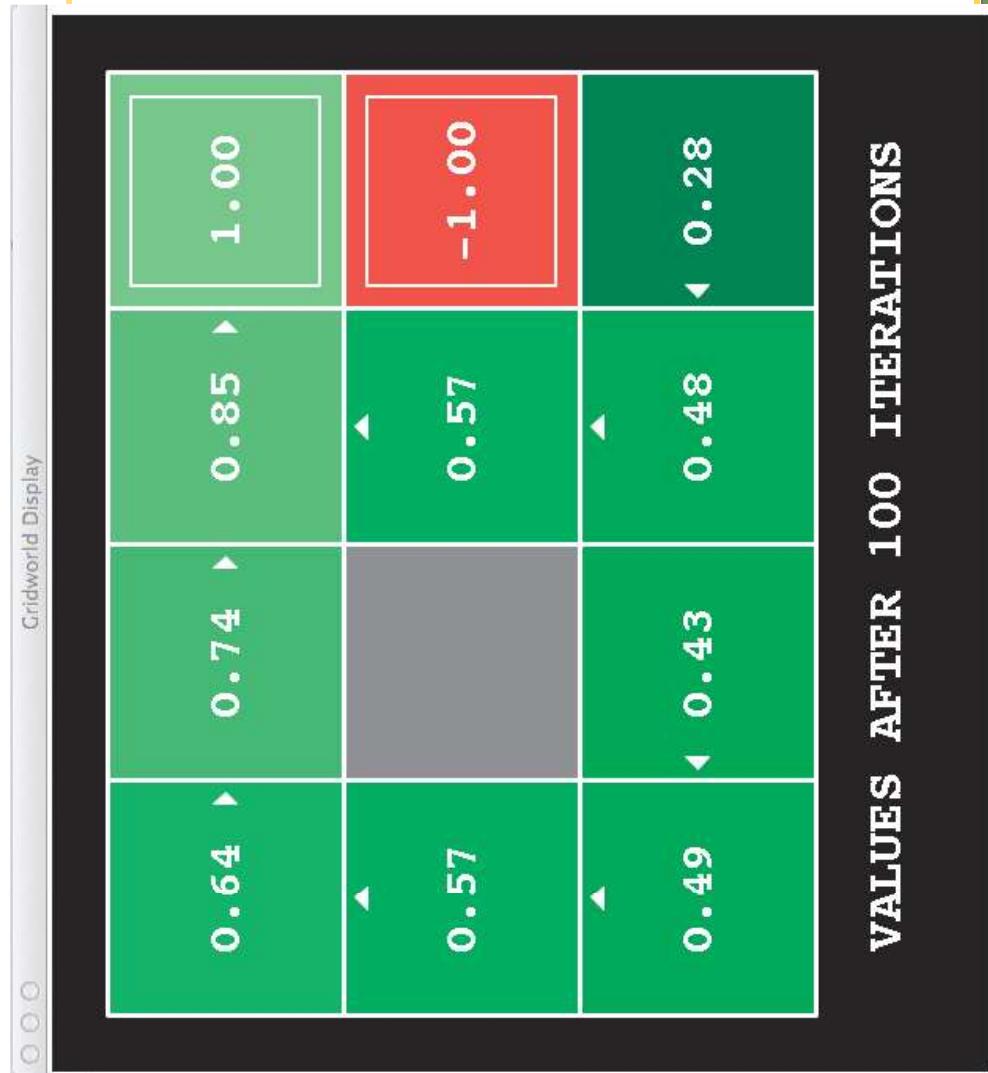
VALUES AFTER 12 ITERATIONS		
0.64 ▶	0.74 ▶	0.85 ▶
◀ 0.57	0.57	◀ 0.47
0.49 ▶	0.42 ▶	0.28 ▶
1.00	-1.00	

Noise = 0.2

Discount = 0.9

Living reward = 0

$k=100$



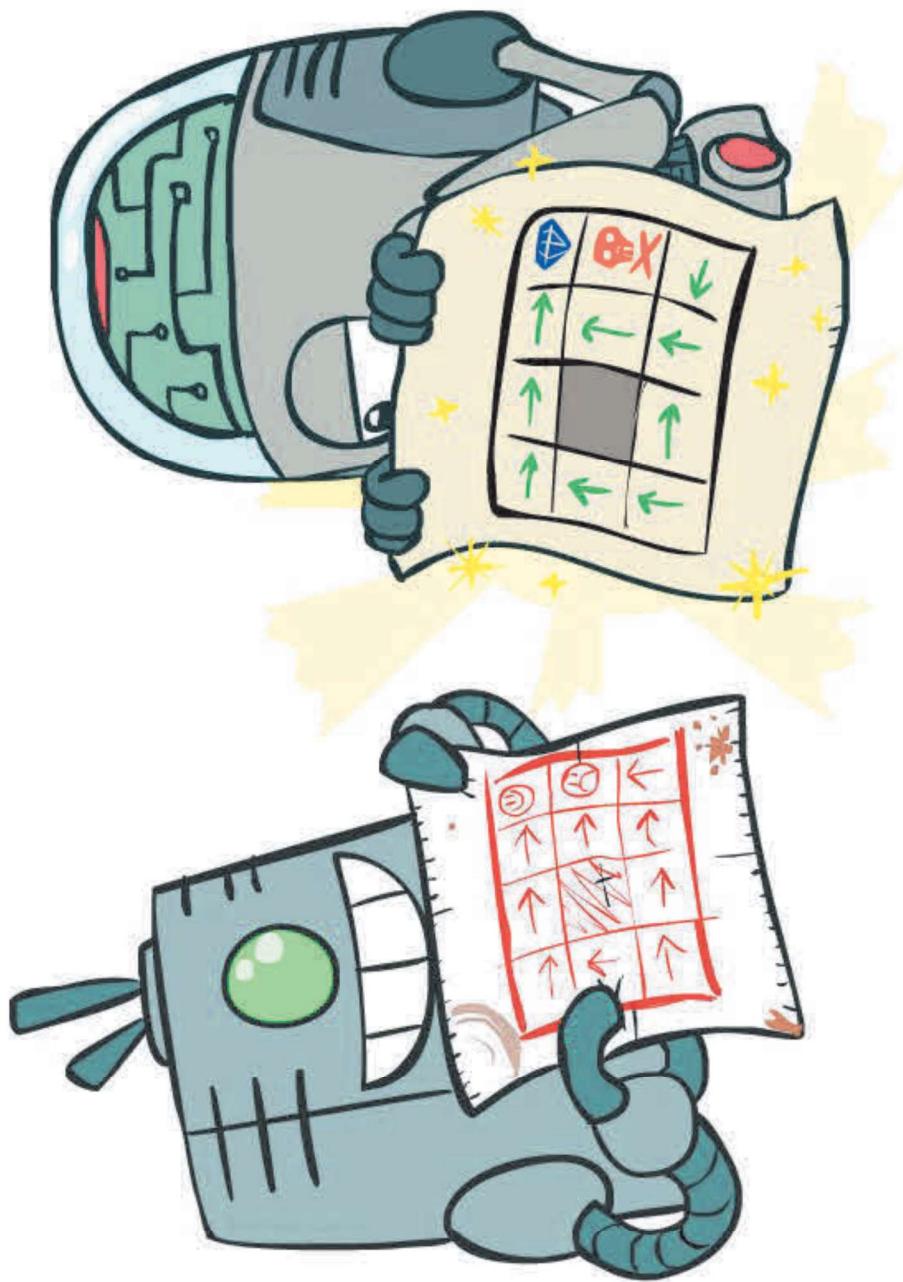
Noise = 0.2

Discount = 0.9

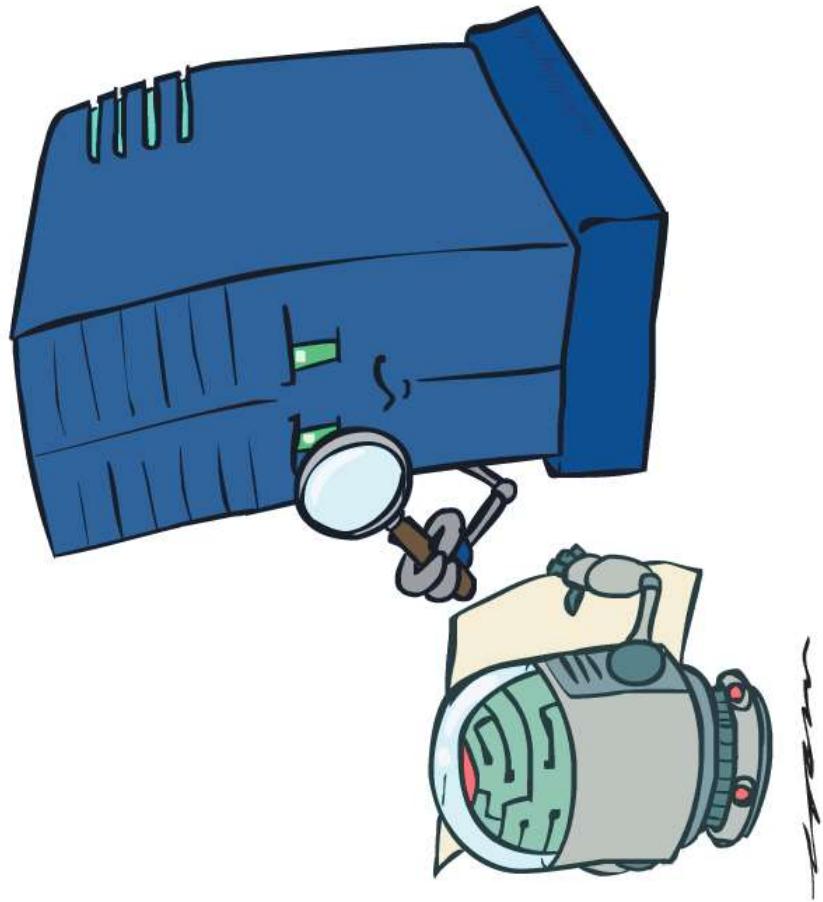


Living reward = 0

Policy Methods



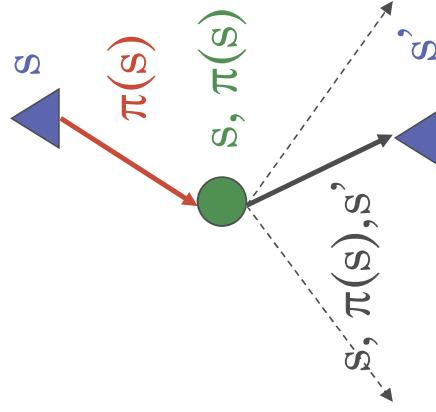
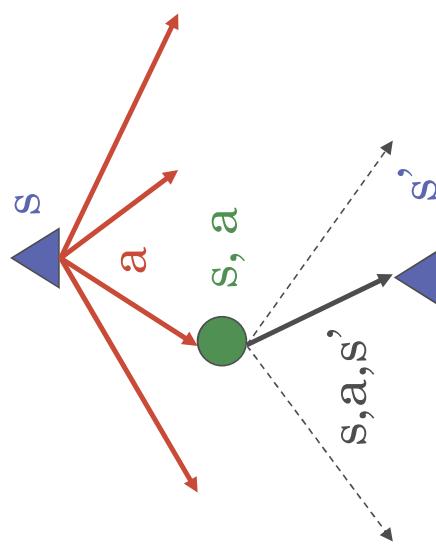
Policy Evaluation



Fixed Policies

Do the optimal action

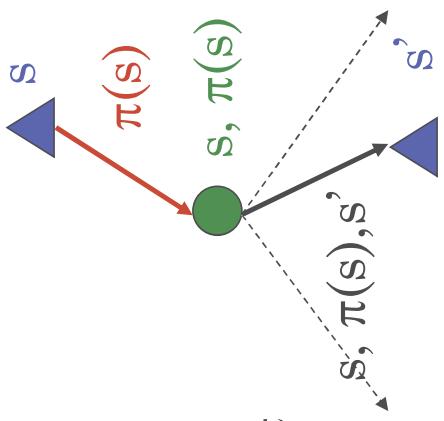
Do what π says to do



- Expectimax trees max over all actions to compute the optimal values
- If we fixed some policy $\pi(s)$, then the tree would be simpler – only one action per state
 - ... though the tree's value would depend on which policy we fixed

Utilities for a Fixed Policy

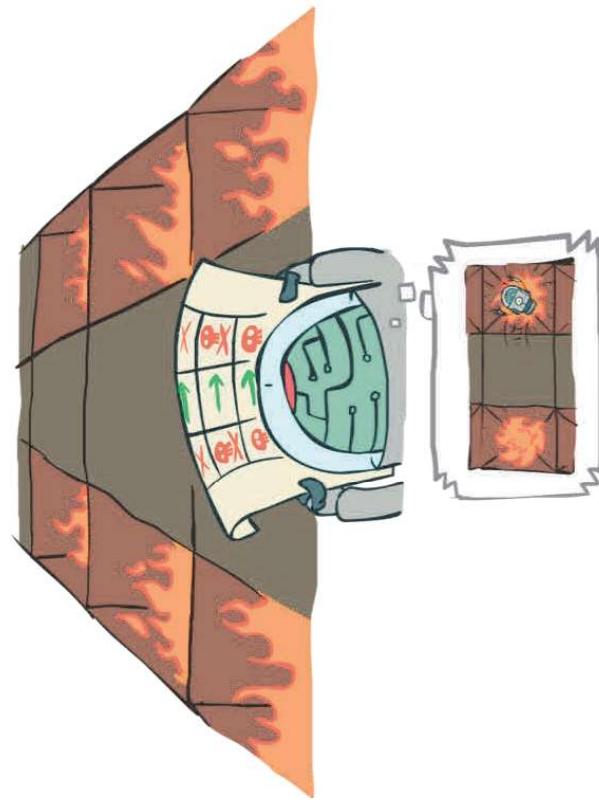
- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy



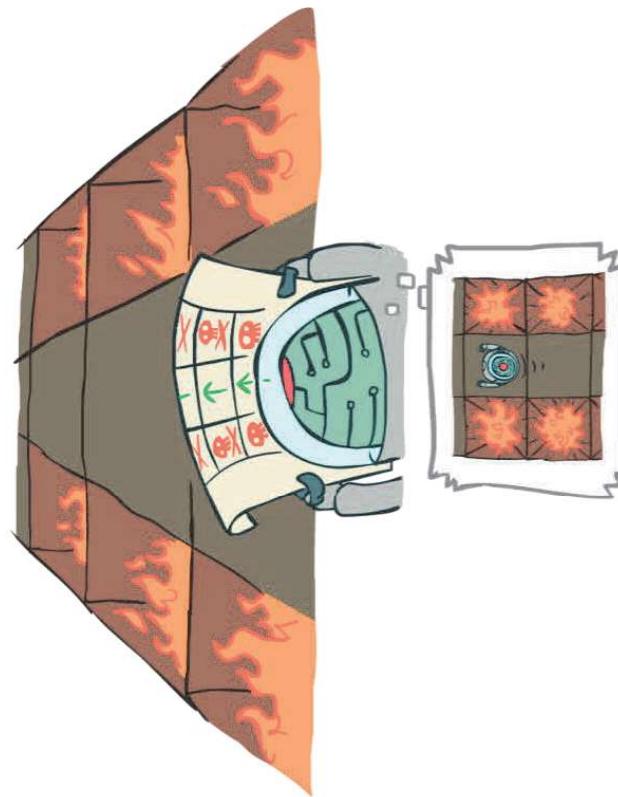
- Define the utility of a state s , under a fixed policy π :
 $V^\pi(s)$ = expected total discounted rewards starting in s and following π
- Recursive relation (one-step look-ahead / Bellman equation):
$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Example: Policy Evaluation

Always Go Right



Always Go Forward



Example: Policy Evaluation

Always Go Right

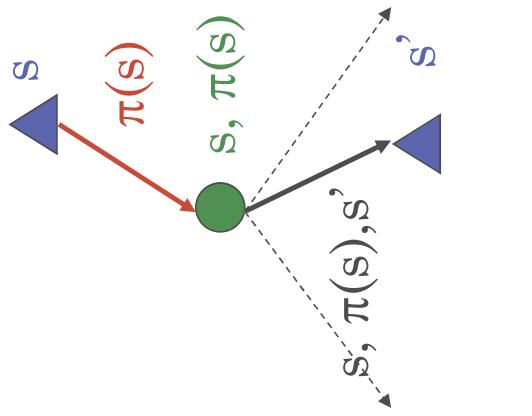
-10.00	100.00	1.09	-10.00	-10.00	-7.88	-8.69	-10.00
-10.00			-10.00		-10.00		-10.00
-10.00				-10.00		-10.00	
-10.00					-10.00		-10.00
-10.00						-10.00	-10.00

Always Go Forward

-10.00	100.00	70.20	48.74	33.30
-10.00		-10.00	-10.00	-10.00
-10.00			-10.00	-10.00
-10.00				-10.00
-10.00				

Policy Evaluation

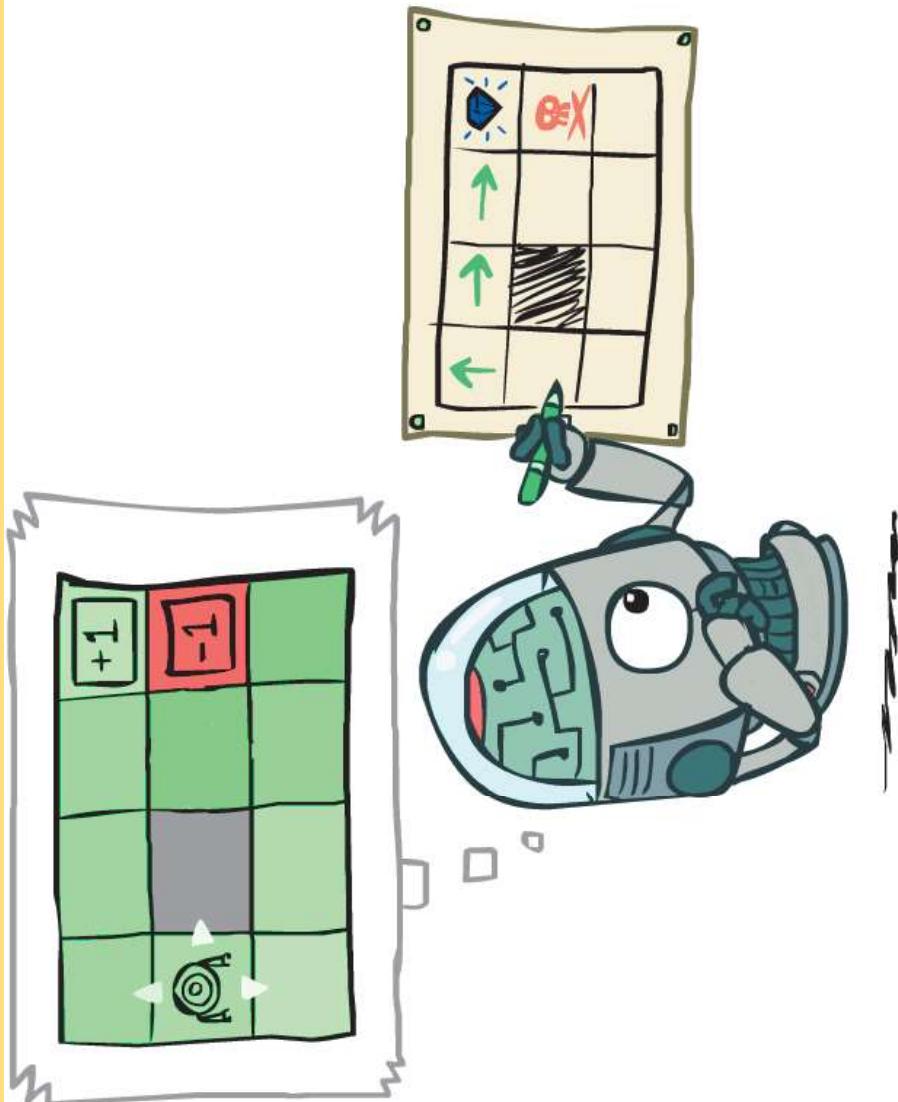
- How do we calculate the V's for a fixed policy π ?
- Idea 1: Turn recursive Bellman equations into updates (like value iteration)



$$V_0^\pi(s) = 0$$
$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

- Efficiency: $O(S^2)$ per iteration
- Idea 2: Without the maxes, the Bellman equations are just a linear system
 - Solve with Matlab (or your favorite linear system solver)

Policy Extraction



Computing Actions from Values

- Let's imagine we have the optimal values $V^*(s)$

0.95	0.96	0.98	1.00
0.94		0.89	-1.00
0.92	0.91	0.90	0.80

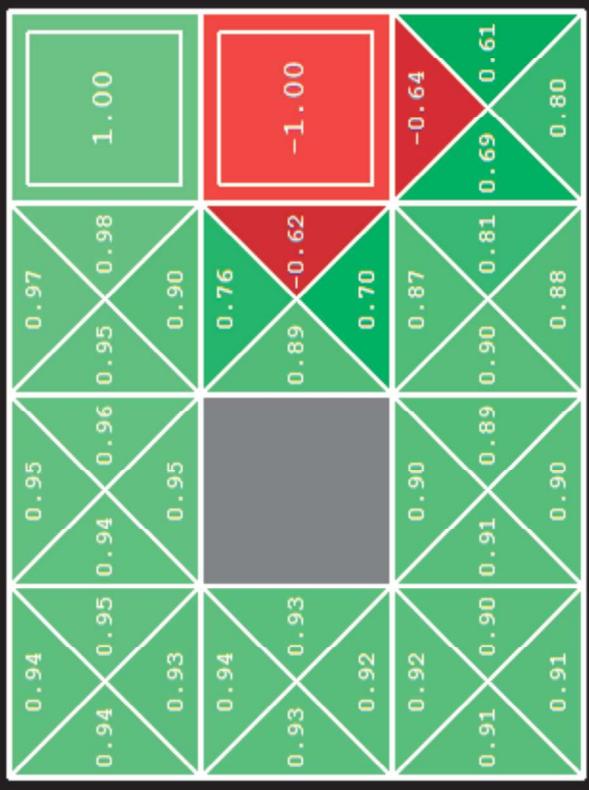
- How should we act?
 - It's not obvious!
- We need to do a mini-expectimax (one step)

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

- This is called **policy extraction**, since it gets the policy implied by the values

Computing Actions from Q-Values

- Let's imagine we have the optimal q-values:

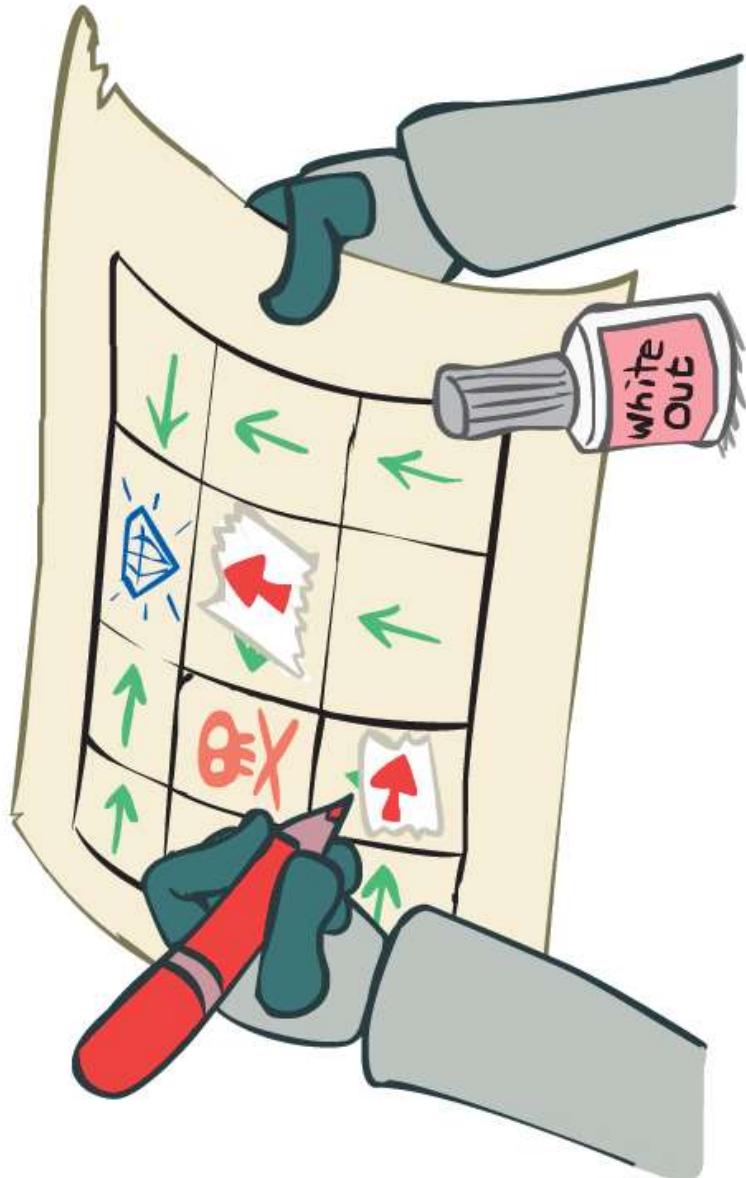


- How should we act?
 - Completely trivial to decide!

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

- Important lesson: actions are easier to select from q-values than values!

Policy Iteration



Policy Iteration

- Alternative approach for optimal values:
 - **Step 1: Policy evaluation:** calculate utilities for some fixed policy (not optimal utilities!) until convergence
 - **Step 2: Policy improvement:** update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
 - Repeat steps until policy converges
- This is **policy iteration**
 - It's still optimal!
 - Can converge (much) faster under some conditions

Policy Iteration

- Evaluation: For fixed current policy π , find values with policy evaluation:

- Iterate until values converge:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

- Improvement: For fixed values, get a better policy using policy extraction

- One-step look-ahead:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

Comparison

- Both value iteration and policy iteration compute the same thing (all optimal values)
- In value iteration:
 - Every iteration updates both the values and (implicitly) the policy
 - We don't track the policy, but taking the max over actions implicitly recomputes it
- In policy iteration:
 - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
 - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
 - The new policy will be better (or we're done)
- Both are dynamic programs for solving MDPs

Summary: MDP Algorithms

- So you want to....
 - Compute optimal values: use value iteration or policy iteration
 - Compute values for a particular policy: use policy evaluation
 - Turn your values into a policy: use policy extraction (one-step lookahead)
- These all look the same!
 - They basically are – they are all variations of Bellman updates
 - They all use one-step look-ahead expectimax fragments
 - They differ only in whether we plug in a fixed policy or max over actions

Example: Racing

- Discount: $\gamma = 0.1$

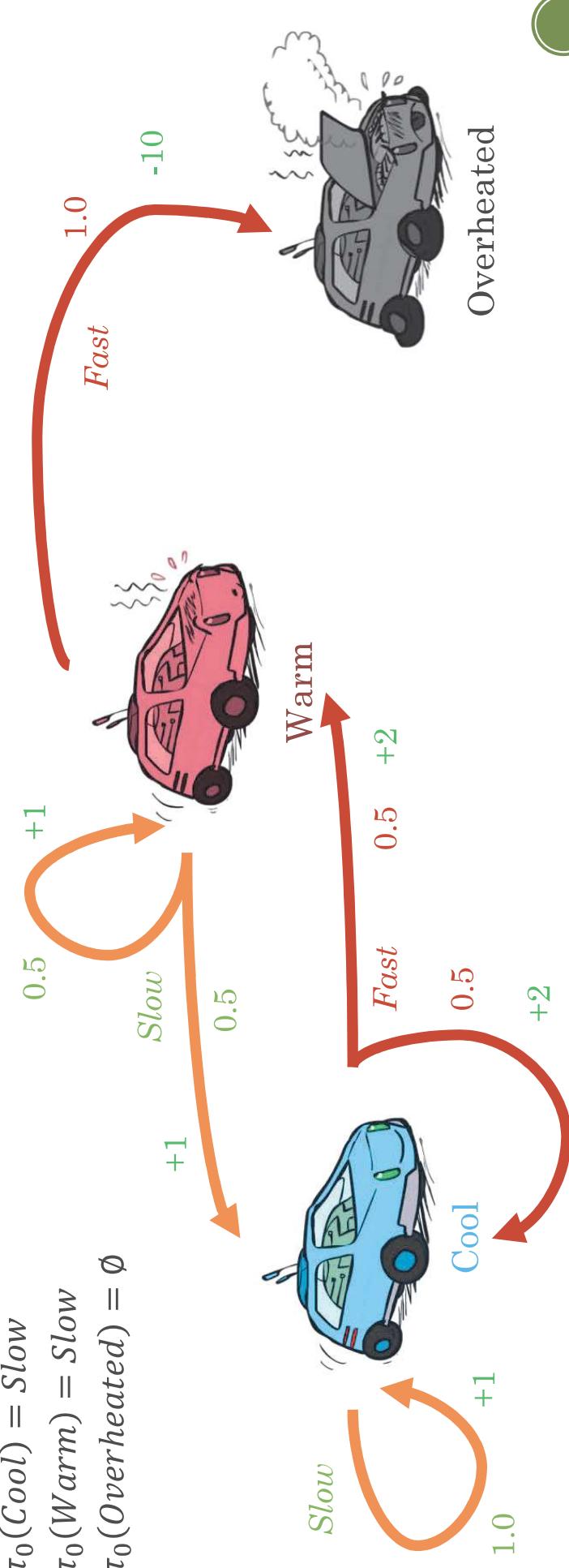
- Initial policy

- $\pi_0(Cool) = Slow$

▪ $\pi_s(Warm) \equiv Slow$

■ $\pi_1(Overheated) = \emptyset$

- Discount: $\gamma = 0.1$
- Initial policy
 - $\pi_0(Cool) = Slow$
 - $\pi_0(Warm) = Slow$
 - $\pi_0(Overheated) = Cool$



Example: Racing

- Discount: $\gamma = 0.1$

- Initial policy

- $\pi_0(Cool) = Slow$
- $\pi_0(Warm) = Slow$
- $\pi_0(Overheated) = \emptyset$

+1

+1

Slow

0.5

0.5

+1



Slow

+1

Cool

0.5

+2



Warm

0.5

+2

+2

Fast

0.5

+2

Fast

-1



Overheated

The Bellman Equations

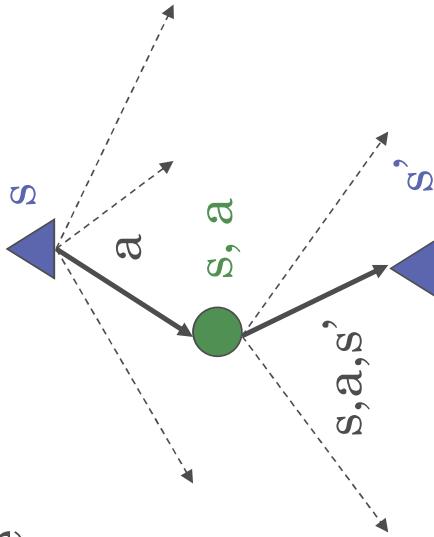
- Definition of “optimal utility” via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

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

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

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

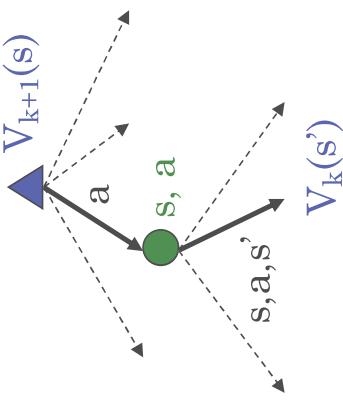
- These are the Bellman equations, and they characterize optimal values in a way we'll use over and over



Value Iteration

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

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

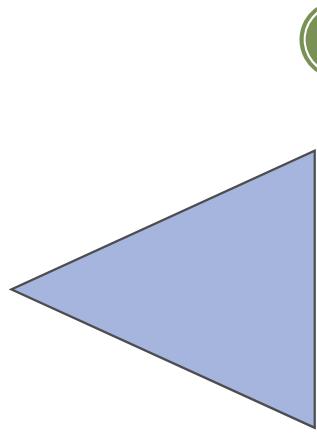


$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

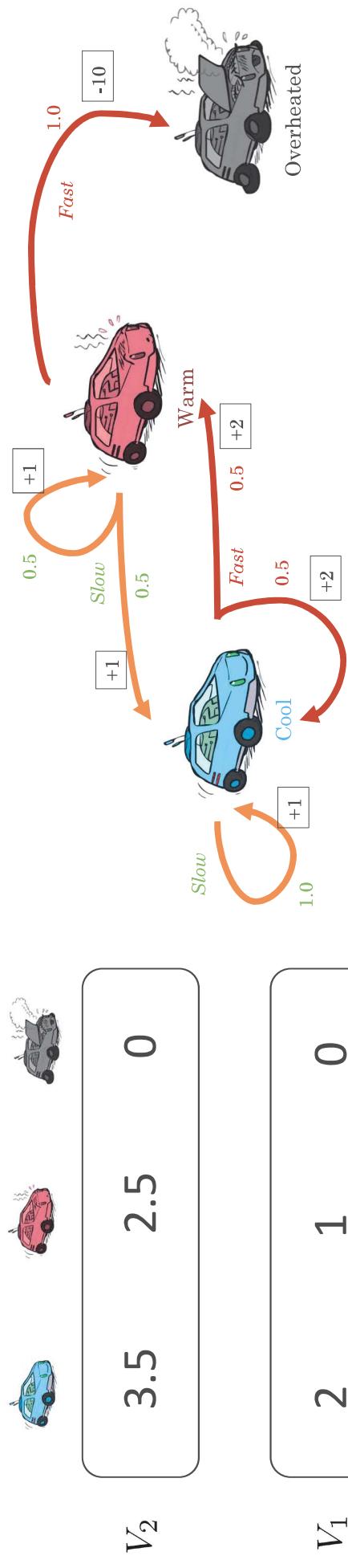
- Repeat until convergence

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

- Theorem: will converge to unique optimal values
 - Basic idea: approximations get refined towards optimal values
 - Policy may converge long before values do



Example: Value Iteration



Assume no discount!

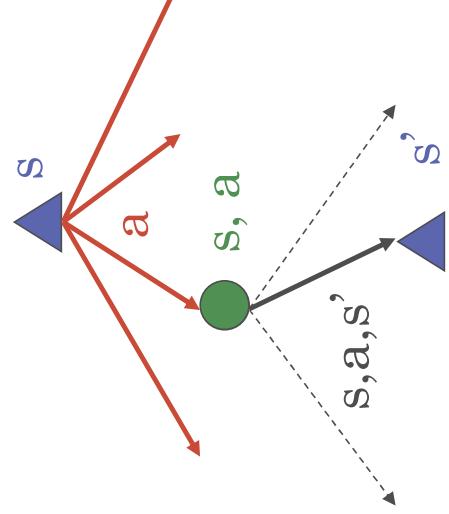
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

V_0	0	0	0
V_1	2	1	0
V_2	3.5	2.5	0

Problems with Value Iteration

- Value iteration repeats the Bellman updates:

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



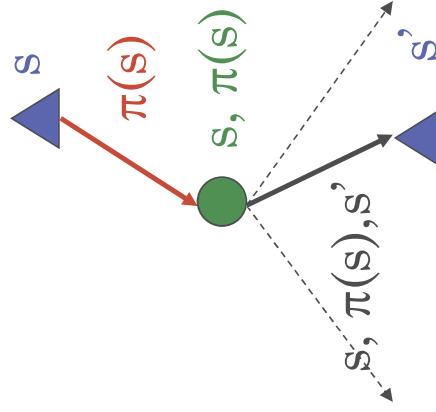
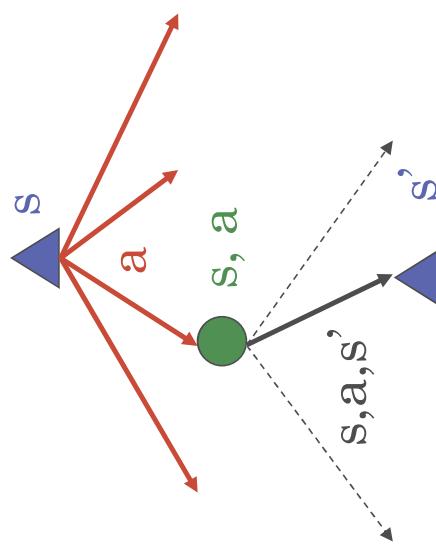
- Problem 1: It's slow – $O(S^2A)$ per iteration
- Problem 2: The “max” at each state rarely changes

- Problem 3: The policy often converges long before the values

Fixed Policies

Do the optimal action

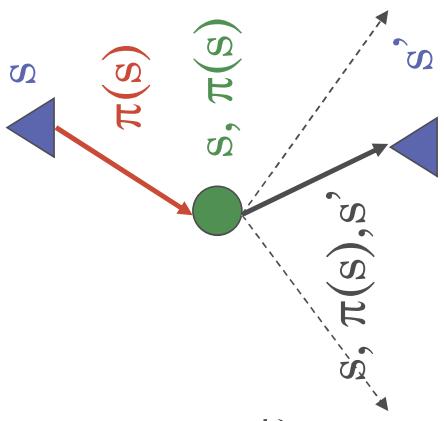
Do what π says to do



- Expectimax trees max over all actions to compute the optimal values
- If we fixed some policy $\pi(s)$, then the tree would be simpler – only one action per state
 - ... though the tree's value would depend on which policy we fixed

Utilities for a Fixed Policy

- Another basic operation: compute the utility of a state s under a fixed (generally non-optimal) policy



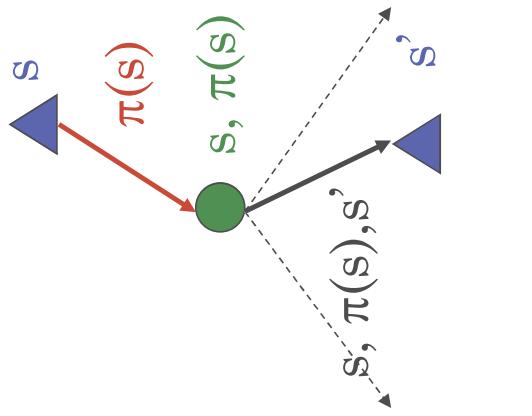
- Define the utility of a state s , under a fixed policy π :
 $V^\pi(s)$ = expected total discounted rewards starting in s and following π

- Recursive relation (one-step look-ahead / Bellman equation):

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

Policy Evaluation

- How do we calculate the V's for a fixed policy π ?
- Idea 1: Turn recursive Bellman equations into updates (like value iteration)

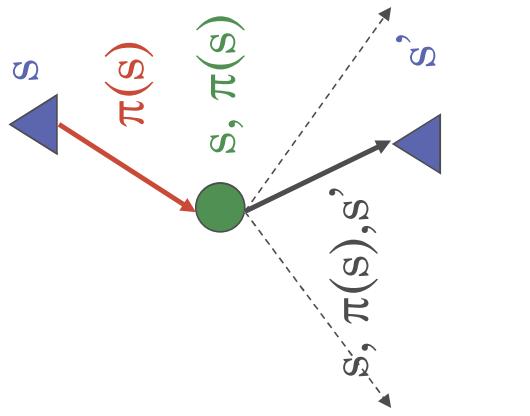


$$V_0^\pi(s) = 0$$
$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

- Efficiency: $O(S^2)$ per iteration
- Idea 2: Without the maxes, the Bellman equations are just a linear system
 - Solve with Matlab (or your favorite linear system solver)

Policy Evaluation

- How do we calculate the V's for a fixed policy π ?
- Idea 1: Turn recursive Bellman equations into updates (like value iteration)



$$\begin{aligned}V_0^\pi(s) &= 0 \\V_{k+1}^\pi(s) &\leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]\end{aligned}$$

- Efficiency: $O(S^2)$ per iteration
- Idea 2: Without the maxes, the Bellman equations are just a linear system
 - Solve with Matlab (or your favorite linear system solver)

Computing Actions from Values

- Let's imagine we have the optimal values $V^*(s)$

0.95	0.96	0.98	1.00
0.94		0.89	-1.00
0.92	0.91	0.90	0.80

- How should we act?
 - It's not obvious!
- We need to do a mini-expectimax (one step)

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

- This is called **policy extraction**, since it gets the policy implied by the values

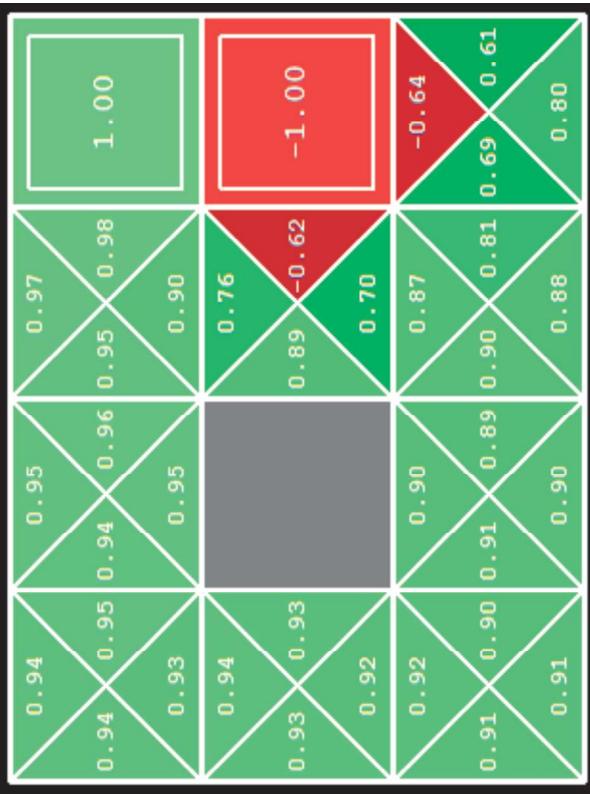
Computing Actions from Q-Values

- Let's imagine we have the optimal q-values:

- How should we act?

- Completely trivial to decide!

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



- Important lesson: actions are easier to select from q-values than values!

Policy Iteration

- Alternative approach for optimal values:
 - **Step 1: Policy evaluation:** calculate utilities for some fixed policy (not optimal utilities!) until convergence
 - **Step 2: Policy improvement:** update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
 - Repeat steps until policy converges
- This is **policy iteration**
 - It's still optimal!
 - Can converge (much) faster under some conditions

Policy Iteration

- Evaluation: For fixed current policy π , find values with policy evaluation:

- Iterate until values converge:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

- Improvement: For fixed values, get a better policy using policy extraction

- One-step look-ahead:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

Comparison

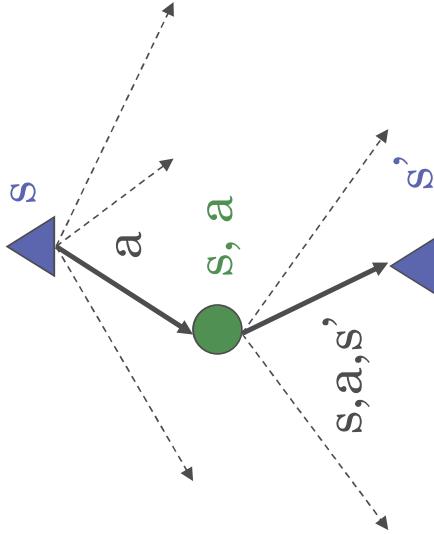
- Both value iteration and policy iteration compute the same thing (all optimal values)
- In value iteration:
 - Every iteration updates both the values and (implicitly) the policy
 - We don't track the policy, but taking the max over actions implicitly recomputes it
- In policy iteration:
 - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
 - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
 - The new policy will be better (or we're done)
- Both are dynamic programs for solving MDPs

Summary: MDP Algorithms

- So you want to....
 - Compute optimal values: use value iteration or policy iteration
 - Compute values for a particular policy: use policy evaluation
 - Turn your values into a policy: use policy extraction (one-step lookahead)
- These all look the same!
 - They basically are – they are all variations of Bellman updates
 - They all use one-step look-ahead expectimax fragments
 - They differ only in whether we plug in a fixed policy or max over actions

Recap: MDPs

- Markov decision processes:
 - States S
 - Actions A
 - Transitions $P(s' | s, a)$ (or $T(s, a, s')$)
 - Rewards $R(s, a, s')$ (and discount γ)
 - Start state s_0

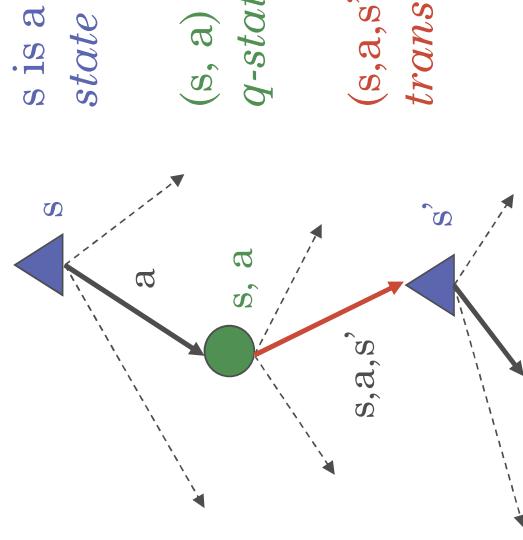


- Quantities:
 - Policy = map of states to actions
 - Utility = sum of discounted rewards
 - Values = expected future utility from a state (max node)
 - Q-Values = expected future utility from a q-state (chance node)

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

Policy Iteration

- Alternative approach for optimal values:
 - **Step 1: Policy evaluation:** calculate utilities for some fixed policy (not optimal utilities!) until convergence
 - **Step 2: Policy improvement:** update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
 - Repeat steps until policy converges
- This is **policy iteration**
 - It's still optimal!
 - Can converge (much) faster under some conditions

Policy Iteration

- Evaluation: For fixed current policy π , find values with policy evaluation:

- Iterate until values converge:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

- Improvement: For fixed values, get a better policy using policy extraction

- One-step look-ahead:

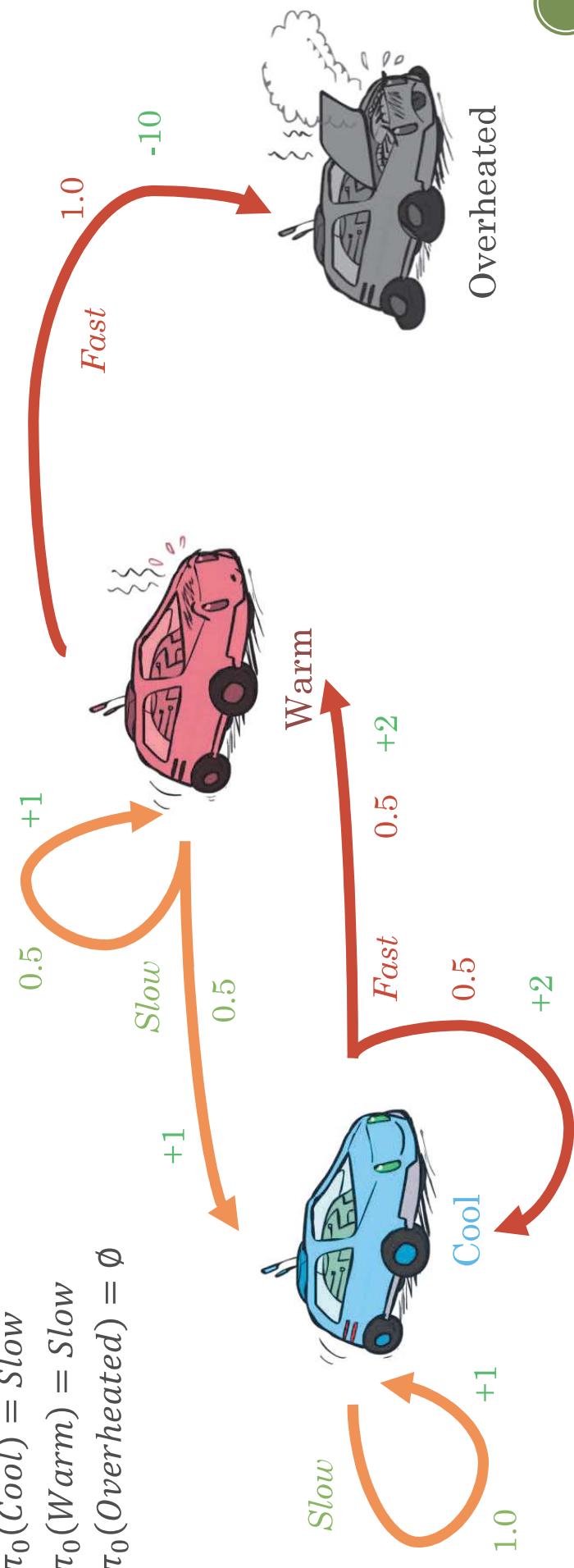
$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

Example: Racing

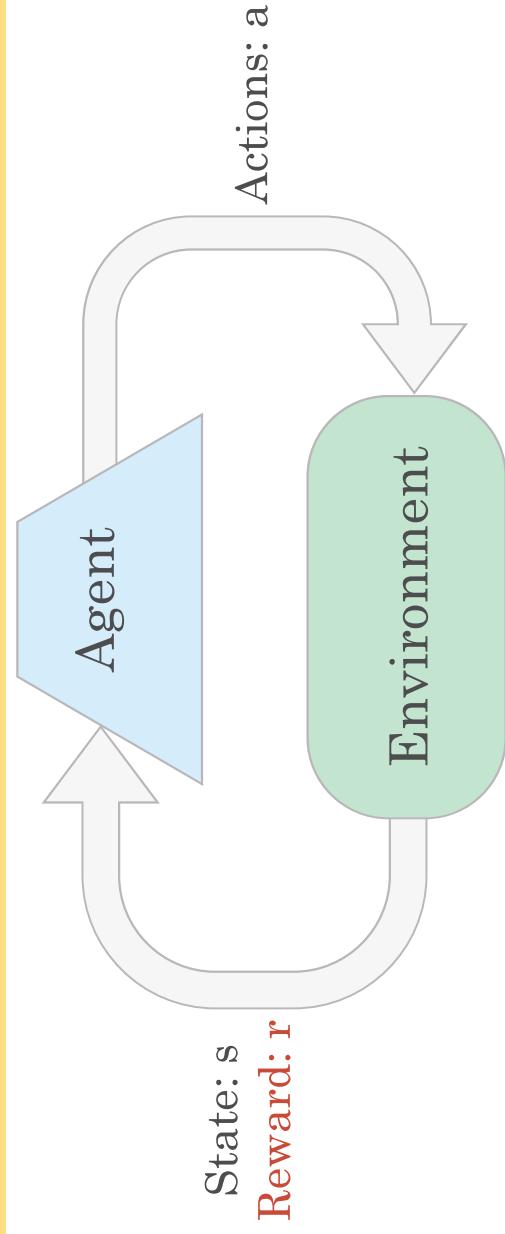
- Discount: $\gamma = 0.1$

- Initial policy

- $\pi_0(Cool) = Slow$
- $\pi_0(Warm) = Slow$
- $\pi_0(Overheated) = \emptyset$



Reinforcement Learning



- Basic idea:
 - Receive feedback in the form of **rewards**
 - Agent's utility is defined by the reward function
 - Must (learn to) act so as to **maximize expected rewards**
 - All learning is based on observed samples of outcomes!

Reinforcement Learning

- Still assume a Markov decision process (MDP):

- A set of states $s \in S$
- A set of actions (per state) A
- A model $T(s, a, s')$
- A reward function $R(s, a, s')$

- Still looking for a policy $\pi(s)$

- New twist: **don't know T or R**

- I.e. we don't know which states are good or what the actions do
- Must actually try out actions and states to learn



The Story So Far: MDPs and RL

Known MDP: Offline Solution

Goal

Compute V^* , Q^* , π^*

Technique

Value / policy iteration

Evaluate a fixed policy π

Policy evaluation

Unknown MDP: Model-Free

Goal

Technique

Compute V^* , Q^* , π^*

Q-learning

Evaluate a fixed policy π

Value Learning

Unknown MDP: Model-Based

Goal

Technique

Compute V^* , Q^* , π^*

VI/PI on approx. MDP

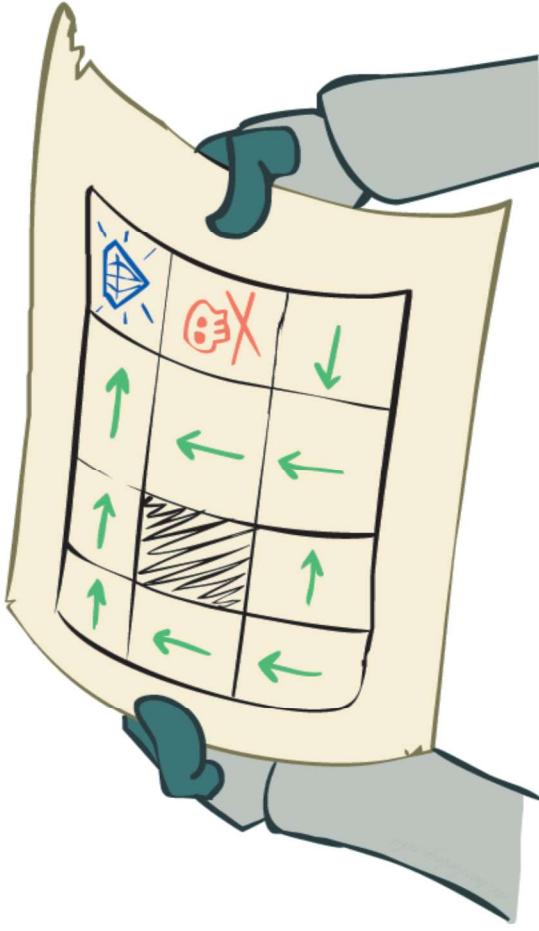
Evaluate a fixed policy π

PE on approx. MDP

Passive Reinforcement Learning

- Simplified task: policy evaluation

- Input: a fixed policy $\pi(s)$
- You don't know the transitions $T(s,a,s')$
- You don't know the rewards $R(s,a,s')$
- Goal: learn the state values

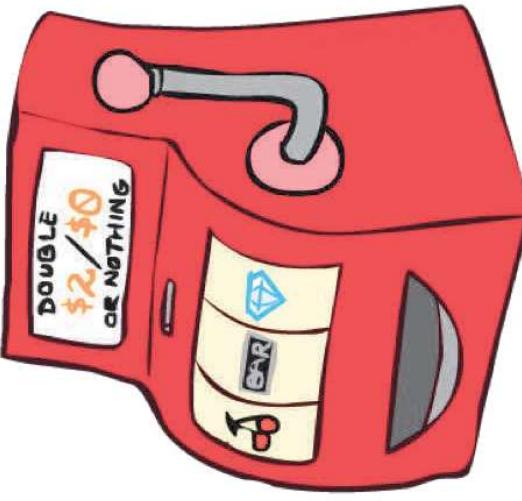


- In this case:

- Learner is “along for the ride”
- No choice about what actions to take
- Just execute the policy and learn from experience
- This is NOT offline planning! You actually take actions in the world.

Direct Evaluation

- Goal: Compute values for each state under π
- Idea: Average together observed sample values
 - Act according to π
 - Every time you visit a state, write down what the sum of discounted rewards turned out to be
 - Average those samples
- This is called direct evaluation



Example: Direct Evaluation

Input Policy π	Observed Episodes (Training)	Output Values
	Episode 1	Episode 2
	B, east, C, -1 C, east, D, -1 D, exit, x, +10	B, east, C, -1 C, east, D, -1 D, exit, x, +10
	Episode 3	Episode 4
	E, north, C, -1 C, east, D, -1 D, exit, x, +10	E, north, C, -1 C, east, A, -1 A, exit, x, -10

Assume: $\gamma = 1$

A	D	
B	C	E

	-10	
A		
B	+8	+4
C		
D		+10
E		-2

Problems with Direct Evaluation

- What's good about direct evaluation?

- It's easy to understand
- It doesn't require any knowledge of T, R
- It eventually computes the correct average values, using just sample transitions

Output Values

		-10 A		
	+8 B		+4 C	+10 D
				-2 E

- What bad about it?

- It wastes information about state connections
- Each state must be learned separately
- So, it takes a long time to learn

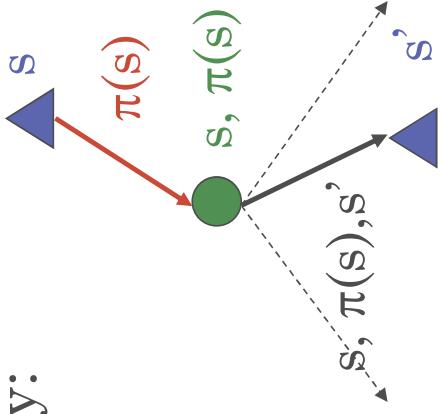
If B and E both go to C under this policy, how can their values be different?

Why Not Use Policy Evaluation?

- Simplified Bellman updates calculate V for a fixed policy:
 - Each round, replace V with a one-step-look-ahead layer over V

$$V_0^\pi(s) = 0$$

$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$



- This approach fully exploited the connections between the states
 - Unfortunately, we need T and R to do it!

- Key question: how can we do this update to V without knowing T and R?
 - In other words, how to we take a weighted average without knowing the weights?

Sample-Based Policy Evaluation?

- We want to improve our estimate of V by computing these averages:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

- Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

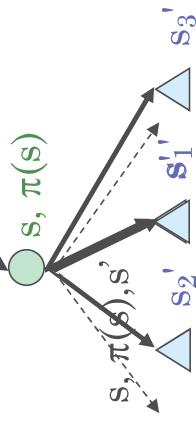
$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

...

$$sample_n = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$

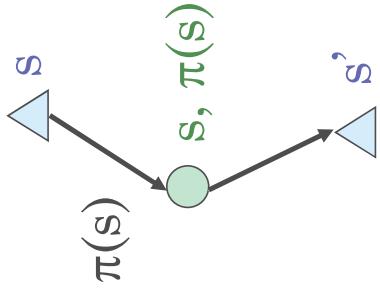
Almost! But we can't
rewind time to get
sample after sample
from state s .

$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_i sample_i$$



Temporal Difference Learning

- Big idea: learn from every experience!
 - Update $V(s)$ each time we experience a transition (s, a, s', r)
 - Likely outcomes s' will contribute updates more often
- Temporal difference learning of values
 - Policy still fixed, still doing evaluation!
 - Move values toward value of whatever successor occurs: running average



Sample of $V(s)$: $sample = R(s, \pi(s), s') + \gamma V^\pi(s')$

Update to $V(s)$: $V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$

Same update: $V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$

Exponential Moving Average

- Exponential moving average
 - The running interpolation update: $\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$
 - Makes recent samples more important:

$$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- Forgets about the past (distant past values were wrong anyway)
- Decreasing learning rate (alpha) can give converging averages

Example: Temporal Difference Learning

States

Observed Transitions

B, east, C, -2

C, east, D, -2

A			D
B	C		
		E	

0	0	8	
0	0	8	
0	0	0	
0	0	0	

0	0	8	
-1	0	8	
-1	0	0	
0	0	0	

0			8
-1	3	8	
-1	3	0	
0	0	0	

Assume: $\gamma = 1$, $\alpha = 1/2$

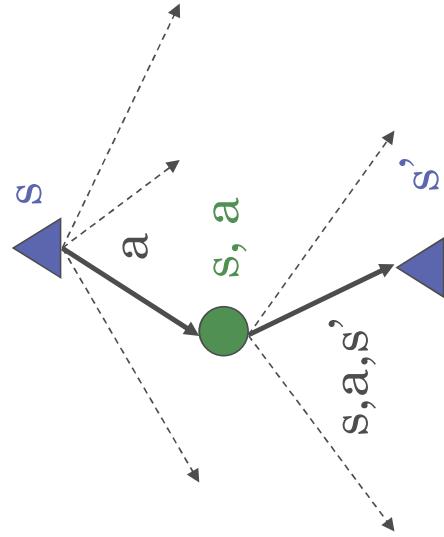
$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Problems with TD Value Learning

- TD value learning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages
- However, if we want to turn values into a (new) policy, we're sunk:

$$\pi(s) = \arg \max_a Q(s, a)$$

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



- Idea: learn Q-values, not values
- Makes action selection model-free too!

Detour: Q-Value Iteration

- Value iteration: find successive (depth-limited) values
 - Start with $V_0(s) = 0$, which we know is right
 - Given V_k , calculate the depth $k+1$ values for all states:

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

- But Q-values are more useful, so compute them instead
 - Start with $Q_0(s, a) = 0$, which we know is right
 - Given Q_k , calculate the depth $k+1$ q-values for all q-states:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

Q-Learning

- We'd like to do Q-value updates to each Q-state:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

- But can't compute this update without knowing T, R

- Instead, compute average as we go

- Receive a sample transition (s, a, r, s')
- This sample suggests

$$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

- But we want to average over results from (s, a) (Why?)
- So keep a running average

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$$

Example

- Two states: A, B
- Two actions: Up, Down
- Discount factor: $\gamma = 0.5$
- Learning rate: $\alpha = 0.5$
- $Q(A, \text{Down}) = ?$
- $Q(B, \text{Up}) = ?$

t	s_t	a_t	s_{t+1}	r_t
0	A	Down	B	2
1	B	Down	B	-4
2	B	Up	B	0
3	B	Up	A	3
4	A	Up	A	-1

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$$

CS 471/571 (Fall 2023): Introduction to Artificial Intelligence

Lecture 14: Reinforcement Learning (Part 3)

Thanh H. Nguyen

Source: <http://ai.berkeley.edu/home.html>

Reminder

- Written assignment 3: MDPs and Reinforcement Learning
 - Deadline: Nov 08th, 2023

Reinforcement Learning

- We still assume an MDP:

- A set of states $s \in S$
- A set of actions (per state) A
- A model $T(s,a,s')$
- A reward function $R(s,a,s')$
- Still looking for a policy $\pi(s)$



- New twist: don't know T or R , so must try out actions

- Big idea: Compute all averages over T using sample outcomes

Model-Free Learning

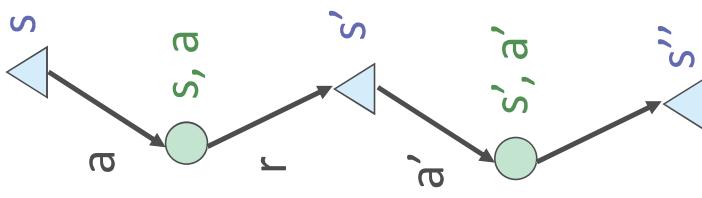
- Model-free (temporal difference) learning

- Experience world through episodes

$$(s, a, r, s', a', r', s'', a'', r'', s''', \dots)$$

- Update estimates each transition (s, a, r, s')

- Over time, updates will mimic Bellman updates



Q-Learning

- We'd like to do Q-value updates to each Q-state:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma \max_{a'} Q_k(s', a') \right]$$

- But can't compute this update without knowing T, R

- Instead, compute average as we go

- Receive a sample transition (s, a, r, s')
- This sample suggests

$$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

- But we want to average over results from (s, a) (Why?)
- So keep a running average

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$$

Example

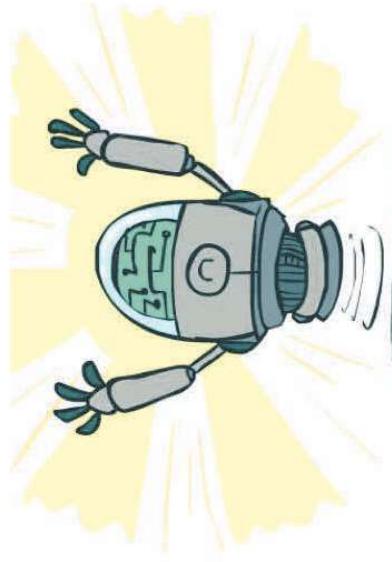
- Two states: A, B
- Two actions: Up, Down
- Discount factor: $\gamma = 0.5$
- Learning rate: $\alpha = 0.5$
- $Q(A, \text{Down}) = ?$
- $Q(B, \text{Up}) = ?$

t	s_t	a_t	s_{t+1}	r_t
0	A	Down	B	2
1	B	Down	B	-4
2	B	Up	B	0
3	B	Up	A	3
4	A	Up	A	-1

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$$

Q-Learning Properties

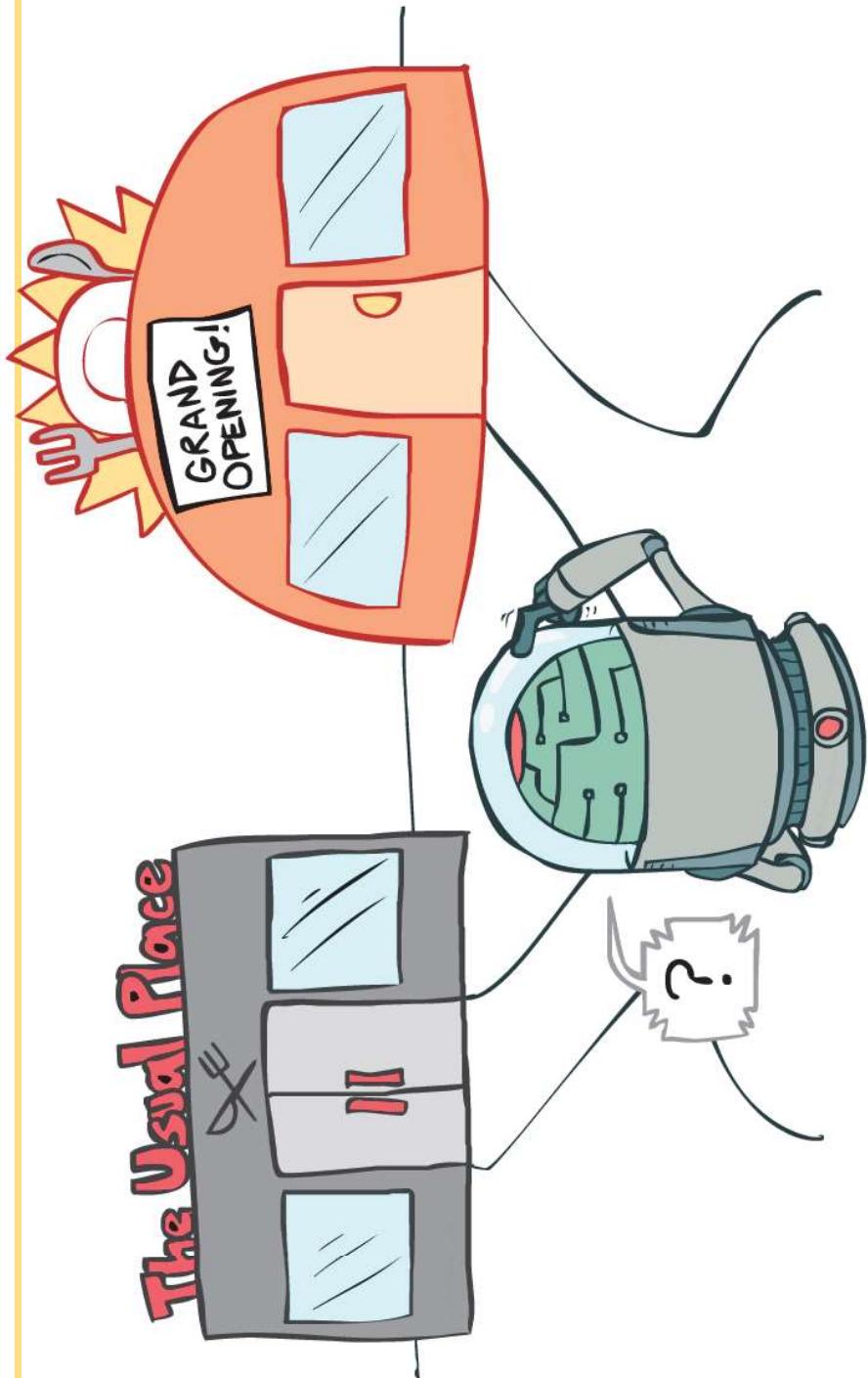
- Amazing result: Q-learning converges to optimal policy
 - even if you're acting suboptimally!



- Caveats:

- You have to explore enough
- You have to eventually make the learning rate small enough
 - ... but not decrease it too quickly
- Basically, in the limit, it doesn't matter how you select actions (!)

Exploration vs. Exploitation



How to Explore?

- Several schemes for forcing exploration
 - Simplest: random actions (ϵ -greedy)
 - Every time step, flip a coin
 - With (small) probability ϵ , act randomly
 - With (large) probability $1-\epsilon$, act on current policy
 - Problems with random actions?
 - You do eventually explore the space, but keep thrashing around once learning is done
 - One solution: lower ϵ over time
 - Another solution: exploration functions



Exploration Functions

- When to explore?
 - Random actions: explore a fixed amount
 - Better idea: explore areas whose badness is not (yet) established, eventually stop exploring



- Exploration function
 - Takes a value estimate u and a visit count n , and returns an optimistic utility, e.g.

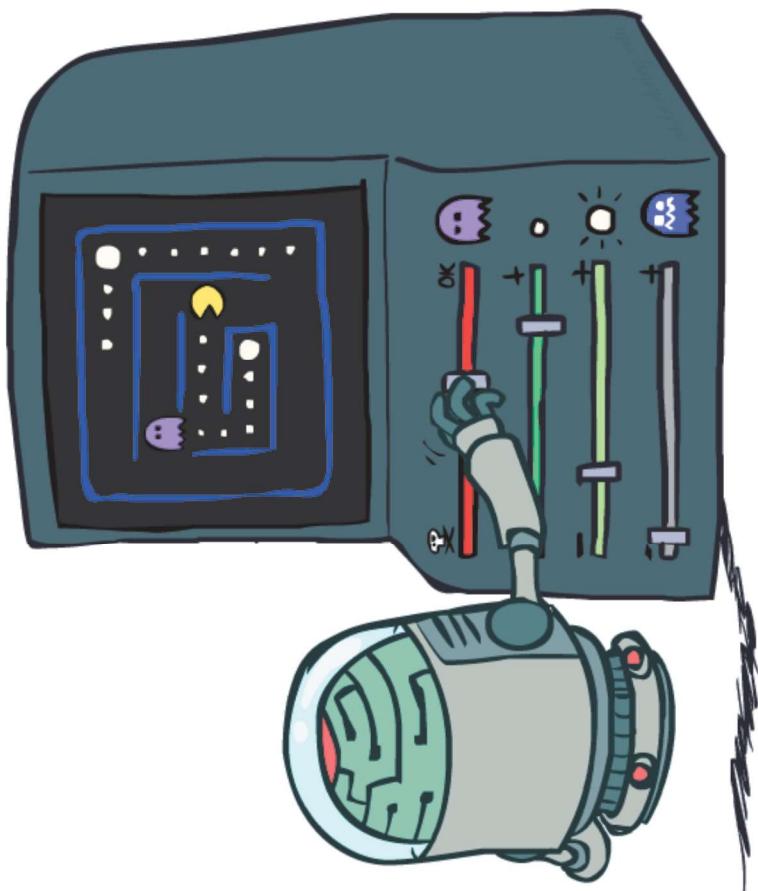
$$f(u, n) = u + k/n$$

Regular Q-Update: $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} Q(s', a')$

Modified Q-Update: $Q(s, a) \leftarrow_{\alpha} R(s, a, s') + \gamma \max_{a'} f(Q(s', a'), N(s', a'))$

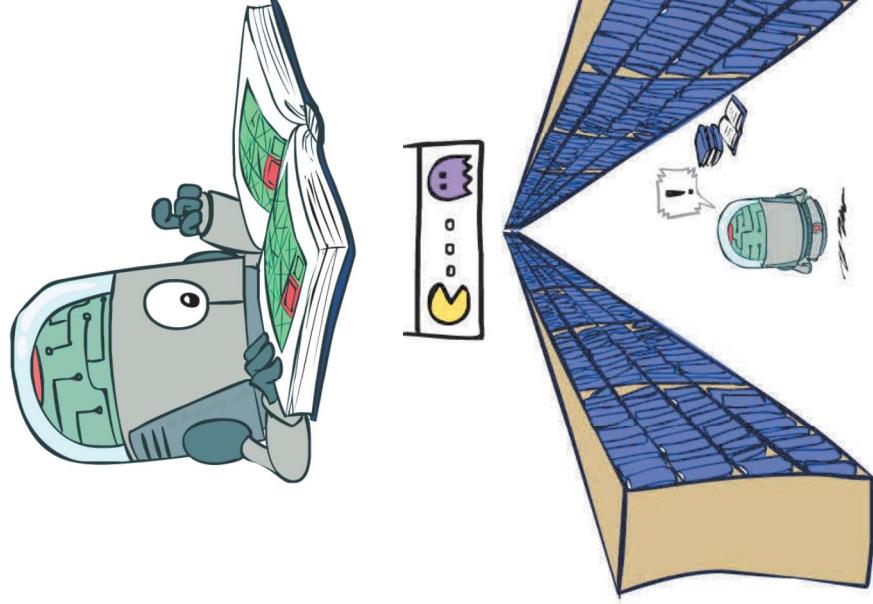
- Note: this propagates the “bonus” back to states that lead to unknown states as well!

Approximate Q-Learning



Generalizing Across States

- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state!
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - Generalize that experience to new, similar situations
 - This is a fundamental idea in machine learning, and we'll see it over and over again

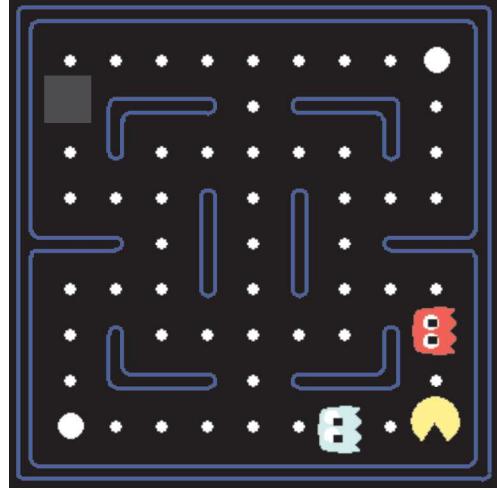
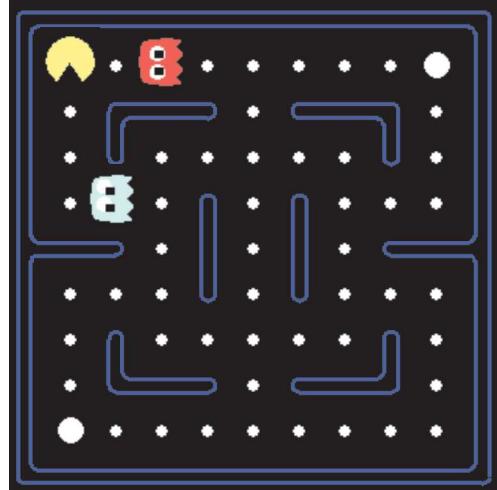
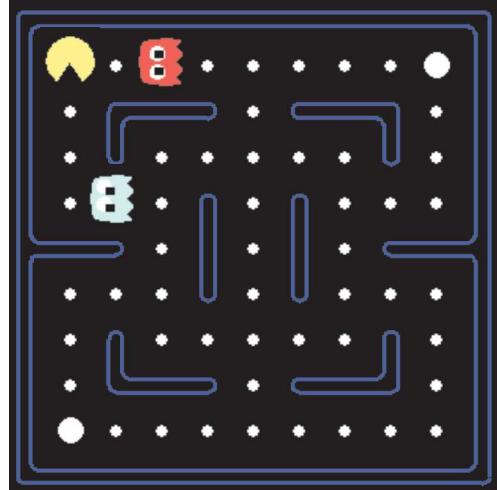
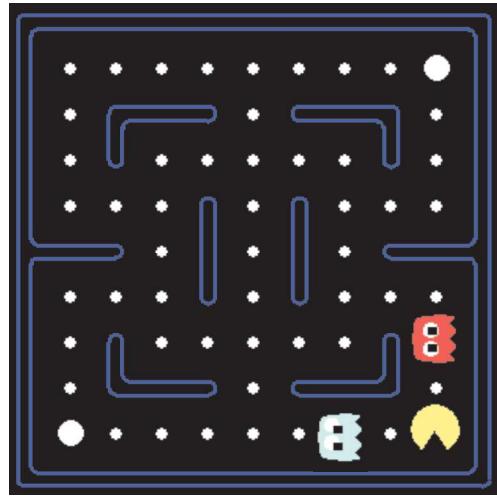


Example: Pacman

Let's say we discover through experience that this state is bad:

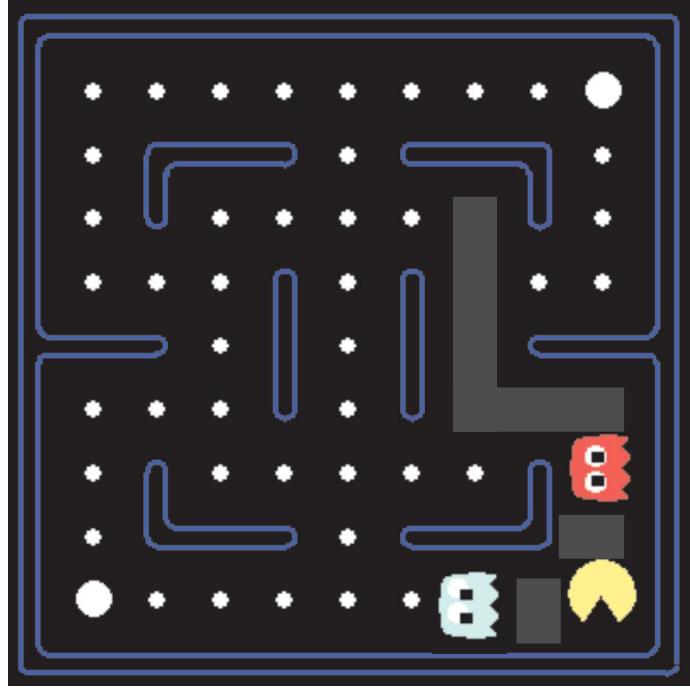
In naïve q-learning, we know nothing about this state:

Or even this one!



Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
 - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
 - Example features:
 - Distance to closest ghost
 - Distance to closest dot
 - Number of ghosts
 - $1 / (\text{dist to dot})^2$
 - Is Pacman in a tunnel? (0/1)
 - etc.
 - Is it the exact state on this slide?
- Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Advantage: our experience is summed up in a few powerful numbers

- Disadvantage: states may share features but actually be very different in value!

Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Q-learning with linear Q-functions:

transition = (s, a, r, s')

$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

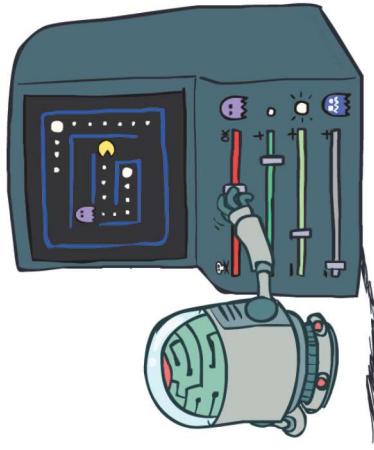
$$Q(s, a) \leftarrow Q(s, a) + \alpha \text{ [difference]}$$

$$w_i \leftarrow w_i + \alpha \text{ [difference]} f_i(s, a)$$

- Intuitive interpretation:

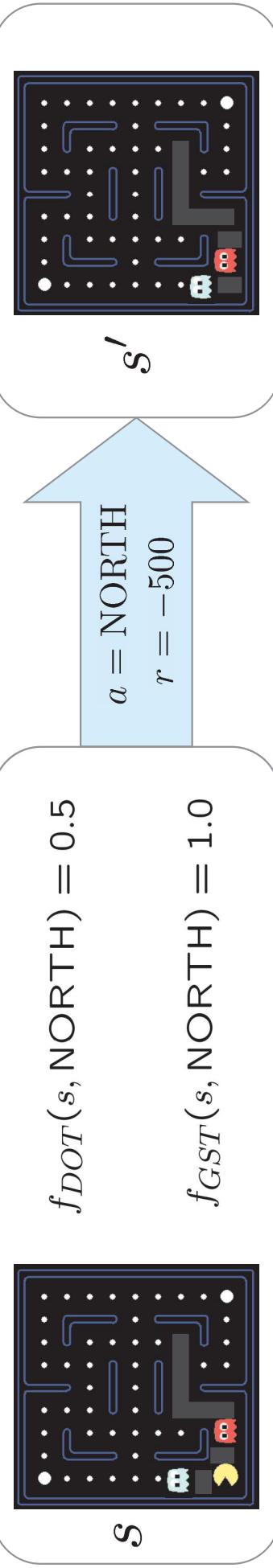
- Adjust weights of active features
 - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

- Formal justification: online least squares



Example: Q-Pacman

$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$



$$Q(s, \text{NORTH}) = +1$$
$$r + \gamma \max_{a'} Q(s', a') = -500 + 0$$
$$Q(s', \cdot) = 0$$

$$\begin{aligned} \text{difference} &= -501 & \xrightarrow{\quad \text{blue arrow} \quad} \\ w_{DOT} &\leftarrow 4.0 + \alpha [-501] 0.5 \\ w_{GST} &\leftarrow -1.0 + \alpha [-501] 1.0 \end{aligned}$$

$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

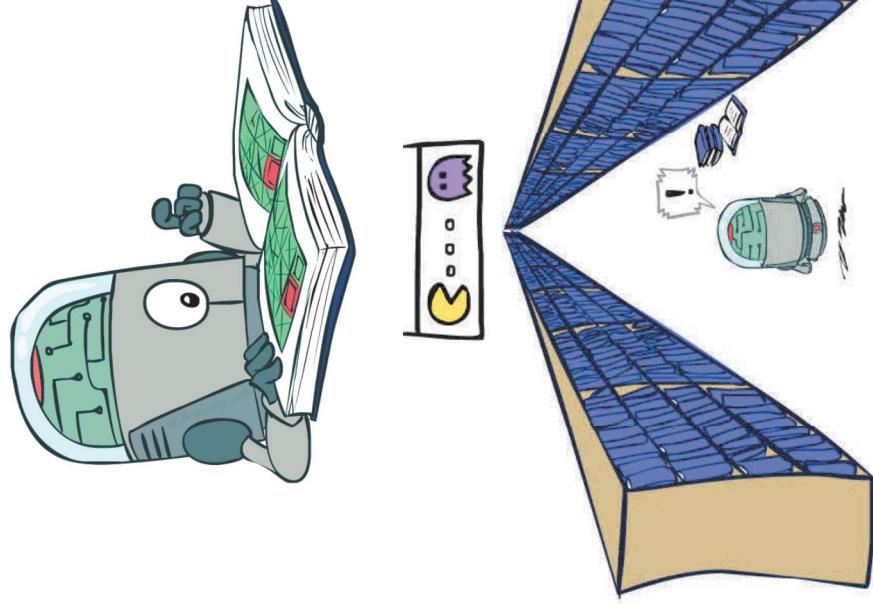
Q-learning with Linear Approximation

Algorithm 4: Q-learning with linear approximation.

```
1 Initialize q-value function  $Q$  with random weights  $w$ :  $Q(s, a; w) = \sum_m w_m f_m(s, a)$ ;
2 for  $episode = 1 \rightarrow M$  do
3   Get initial state  $s_0$ ;
4   for  $t = 1 \rightarrow T$  do
5     With prob.  $\epsilon$ , select a random action  $a_t$ ;
6     With prob.  $1 - \epsilon$ , select  $a_t \in \operatorname{argmax}_a Q(s_t, a; w)$ ;
7     Execute selected action  $a_t$  and observe reward  $r_t$  and next state  $s_{t+1}$ ;
8     Set target  $y_t = \begin{cases} r_t & \text{if episode terminates at step } t + 1 \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$ 
9     Perform a gradient descent step to update  $w$ :  $w_m \leftarrow w_m + \alpha [y_t - Q(s_t, a_t; w)] f_m(s, a)$ ;
```

Generalizing Across States

- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state!
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - Generalize that experience to new, similar situations
 - This is a fundamental idea in machine learning, and we'll see it over and over again

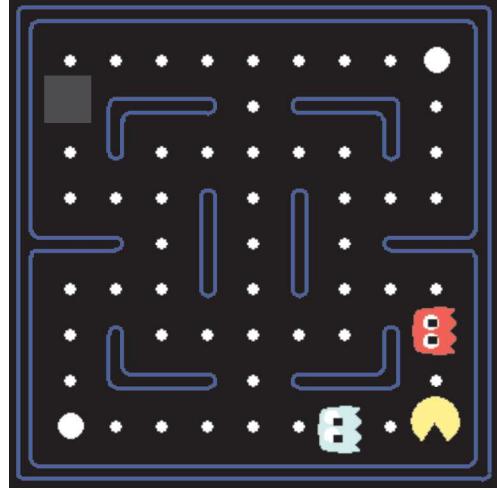
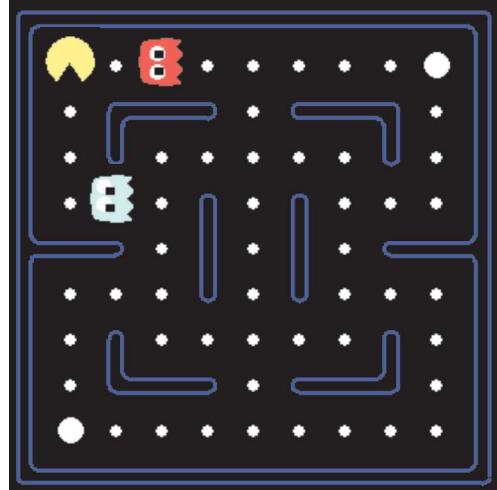
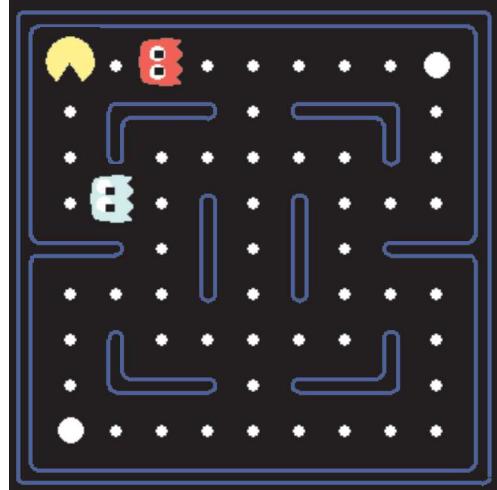
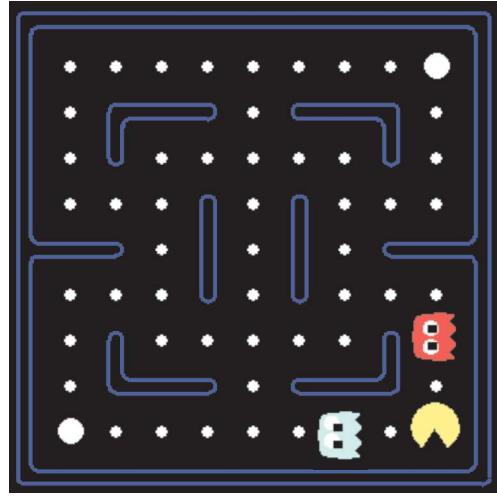


Example: Pacman

Let's say we discover through experience that this state is bad:

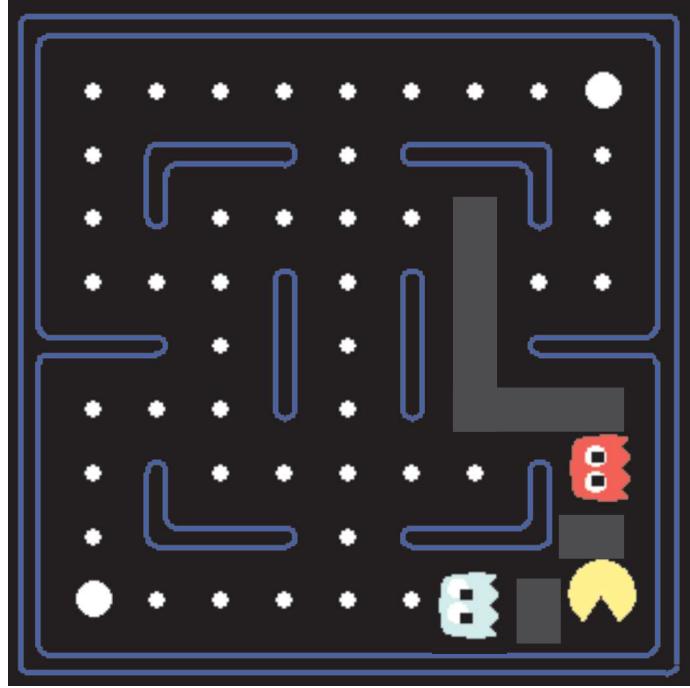
In naïve q-learning, we know nothing about this state:

Or even this one!



Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
 - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
 - Example features:
 - Distance to closest ghost
 - Distance to closest dot
 - Number of ghosts
 - $1 / (\text{dist to dot})^2$
 - Is Pacman in a tunnel? (0/1)
 - etc.
 - Is it the exact state on this slide?
- Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Advantage: our experience is summed up in a few powerful numbers

- Disadvantage: states may share features but actually be very different in value!

Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Q-learning with linear Q-functions:

transition = (s, a, r, s')

$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

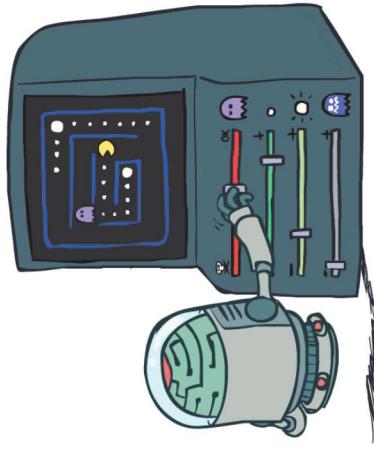
$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$$

- Intuitive interpretation:

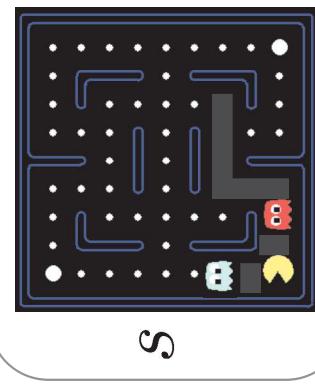
- Adjust weights of active features
 - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

- Formal justification: online least squares



Example: Q-Pacman

$$Q(s, a) = 4.0 f_{DOT}(s, a) - 1.0 f_{GST}(s, a)$$



$$f_{DOT}(s, \text{NORTH}) = 0.5$$

$$f_{GST}(s, \text{NORTH}) = 1.0$$

s'

$$a = \text{NORTH}$$

$$r = -500$$

$$Q(s, \text{NORTH}) = +1$$

$$r + \gamma \max_{a'} Q(s', a') = -500 + 0$$

$$Q(s', \cdot) = 0$$

$$\begin{aligned} w_{DOT} &\leftarrow 4.0 + \alpha [-501] 0.5 \\ w_{GST} &\leftarrow -1.0 + \alpha [-501] 1.0 \end{aligned}$$

$$Q(s, a) = 3.0 f_{DOT}(s, a) - 3.0 f_{GST}(s, a)$$

Q-learning with Linear Approximation

Algorithm 4: Q-learning with linear approximation.

```
1 Initialize q-value function  $Q$  with random weights  $w$ :  $Q(s, a; w) = \sum_m w_m f_m(s, a)$ ;
2 for  $episode = 1 \rightarrow M$  do
3   Get initial state  $s_0$ ;
4   for  $t = 1 \rightarrow T$  do
5     With prob.  $\epsilon$ , select a random action  $a_t$ ;
6     With prob.  $1 - \epsilon$ , select  $a_t \in \operatorname{argmax}_a Q(s_t, a; w)$ ;
7     Execute selected action  $a_t$  and observe reward  $r_t$  and next state  $s_{t+1}$ ;
8     Set target  $y_t = \begin{cases} r_t & \text{if episode terminates at step } t + 1 \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$ 
9     Perform a gradient descent step to update  $w$ :  $w_m \leftarrow w_m + \alpha [y_t - Q(s_t, a_t; w)] f_m(s, a)$ ;
```

Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

- Q-learning with linear Q-functions:

transition = (s, a, r, s')

$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

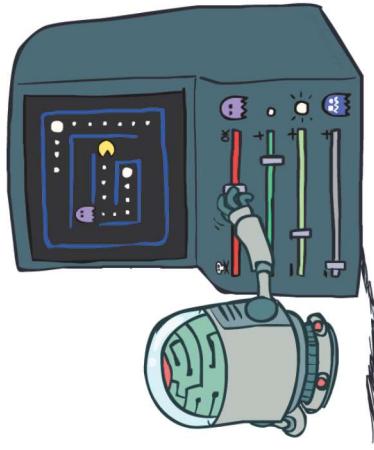
$$Q(s, a) \leftarrow Q(s, a) + \alpha \text{ [difference]}$$

$$w_i \leftarrow w_i + \alpha \text{ [difference]} f_i(s, a)$$

- Intuitive interpretation:

- Adjust weights of active features
 - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features

- Formal justification: online least squares



Q-learning with Linear Approximation

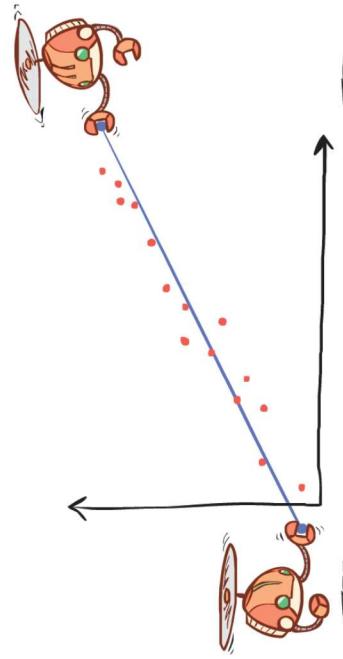
Algorithm 4: Q-learning with linear approximation.

```
1 Initialize q-value function  $Q$  with random weights  $w$ :  $Q(s, a; w) = \sum_m w_m f_m(s, a)$ ;
2 for  $episode = 1 \rightarrow M$  do
3   Get initial state  $s_0$ ;
4   for  $t = 1 \rightarrow T$  do
5     With prob.  $\epsilon$ , select a random action  $a_t$ ;
6     With prob.  $1 - \epsilon$ , select  $a_t \in \operatorname{argmax}_a Q(s_t, a; w)$ ;
7     Execute selected action  $a_t$  and observe reward  $r_t$  and next state  $s_{t+1}$ ;
8     Set target  $y_t = \begin{cases} r_t & \text{if episode terminates at step } t + 1 \\ r_t + \gamma \max_{a'} Q(s_{t+1}, a'; w) & \text{otherwise} \end{cases}$ 
9     Perform a gradient descent step to update  $w$ :  $w_m \leftarrow w_m + \alpha [y_t - Q(s_t, a_t; w)] f_m(s, a)$ ;
```

Minimizing Error*

Imagine we had only one point \mathbf{x} , with features $f(\mathbf{x})$, target value y , and weights w :

$$\text{error}(w) = \frac{1}{2} \left(y - \sum_k w_k f_k(x) \right)^2$$
$$\frac{\partial \text{error}(w)}{\partial w_m} = - \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$
$$w_m \leftarrow w_m + \alpha \left(y - \sum_k w_k f_k(x) \right) f_m(x)$$



Approximate q update explained:

$$w_m \leftarrow w_m + \alpha [r + \gamma \max_a Q(s', a') - Q(s, a)] f_m(s, a)$$

“target”

“prediction”

Marginal Distributions

- Marginal distributions are sub-tables which eliminate variables
- Marginalization (summing out): Combine collapsed rows by adding

$$P(T, W)$$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

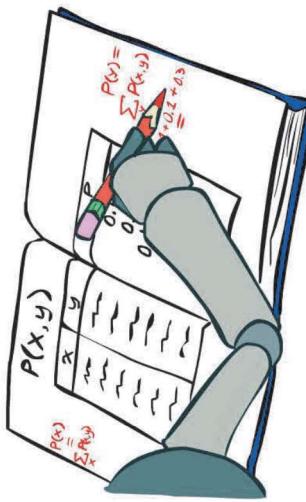
$$P(t) = \sum_s P(t, s)$$

T	P
hot	0.5
cold	0.5

$$P(W)$$

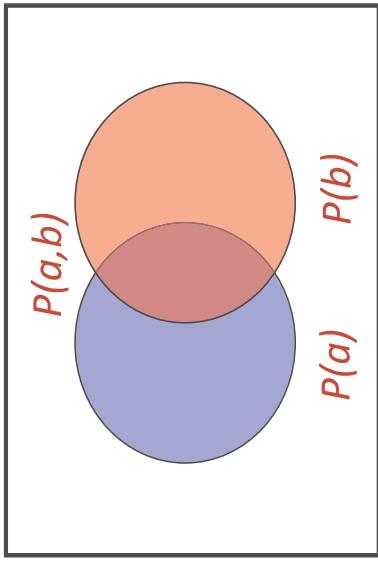
W	P
sun	0.6
rain	0.4

$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2)$$



Conditional Probabilities

- A simple relation between joint and marginal probabilities
 - In fact, this is taken as the *definition* of a **conditional probability**



$$P(a|b) = \frac{P(a,b)}{P(b)}$$

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

$$\begin{aligned} P(W = s | T = c) &= \frac{P(W = s, T = c)}{P(T = c)} = \frac{0.2}{0.5} = 0.4 \\ &= P(W = s, T = c) + P(W = r, T = c) \\ &= 0.2 + 0.3 = 0.5 \end{aligned}$$



Conditional Distributions

- Conditional distributions are probability distributions over some variables given fixed values of others

Conditional Distributions

$$P(W|T = \text{hot})$$

W	P
sun	0.8
rain	0.2

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

$$P(W|T = \text{cold})$$

W	P
sun	0.4
rain	0.6

$P(W|T)$

Normalization Trick

$$\begin{aligned} P(W = s | T = c) &= \frac{P(W = s, T = c)}{P(T = c)} \\ &= \frac{P(W = s, T = c) + P(W = r, T = c)}{P(W = s, T = c) + P(W = r, T = c)} \\ &= \frac{0.2}{0.2 + 0.3} = 0.4 \\ P(W | T = c) &= \frac{P(W = r, T = c)}{P(T = c)} \\ &= \frac{P(W = r, T = c) + P(W = s, T = c)}{P(W = r, T = c) + P(W = s, T = c)} \\ &= \frac{0.3}{0.2 + 0.3} = 0.6 \end{aligned}$$

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

W	P
sun	0.4
rain	0.6

Normalization Trick

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

SELECT the joint
probabilities
matching the
evidence

$P(c, W)$ (make it sum to one)

T	W	P
cold	sun	0.2
cold	rain	0.3



$P(W|T = c)$

W	P
sun	0.4
rain	0.6



$$\begin{aligned}
 P(W = s|T = c) &= \frac{P(W = s, T = c)}{P(T = c)} \\
 &= \frac{P(W = s, T = c) + P(W = r, T = c)}{0.2 + 0.3} \\
 &= \frac{0.3}{0.2 + 0.3} = 0.6
 \end{aligned}$$



Normalization Trick

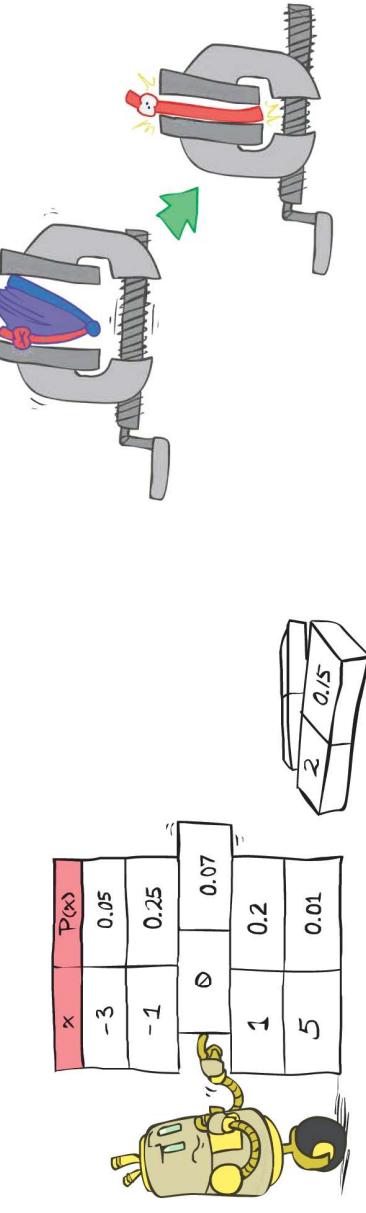
$P(T, W)$			SELECT the joint probabilities matching the evidence			NORMALIZE the selection (make it sum to one)		$P(W T = c)$
T	W	P	T	W	P			
hot	sun	0.4						
hot	rain	0.1						
cold	sun	0.2						
cold	rain	0.3						

- Why does this work? Sum of selection is $P(\text{evidence})!$ ($P(T=c)$, here)

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\sum_{x_1} P(x_1, x_2)}$$

Inference by Enumeration

- General case:
 - Evidence variables: $E_1 \dots E_k = e_1 \dots e_k$
 - Query variable: Q
 - Hidden variables: $H_1 \dots H_r$
- Step 1: Select the entries consistent with the evidence
- Step 2: Sum out H to get joint of Query and evidence
- Step 3: Normalize



$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(Q, h_1 \dots h_r, e_1 \dots e_k) / Z$$

$$Z = \sum_q P(Q, e_1 \dots e_k)$$

$$P(Q | e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

* Works fine with multiple query variables, too

$$\times \frac{1}{Z}$$

Inference by Enumeration

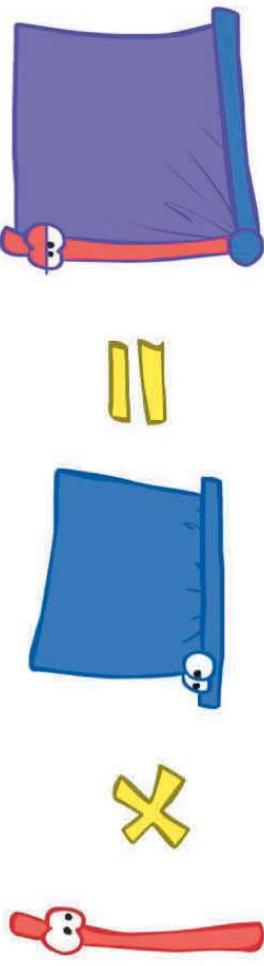
- **Obvious problems:**

- Worst-case time complexity $O(d^n)$
- Space complexity $O(d^n)$ to store the joint distribution

The Product Rule

- Sometimes have conditional distributions but want the joint

$$P(y) P(x|y) = P(x, y) \quad \leftrightarrow \quad P(x|y) = \frac{P(x, y)}{P(y)}$$



The Product Rule

$$P(y)P(x|y) = P(x,y)$$

- Example:

$$P(D|W)$$

D	W	P
wet	sun	0.1
dry	sun	0.9
wet	rain	0.7
dry	rain	0.3

$$P(D,W)$$

D	W	P
wet	sun	
dry	sun	
wet	rain	
dry	rain	



R	P
sun	0.8
rain	0.2

The Chain Rule

- More generally, can always write any joint distribution as an incremental product of conditional distributions

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i|x_1 \dots x_{i-1})$$

- Why is this always true?

Bayes Rule

- Two ways to factor a joint distribution over two variables:

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

- Dividing, we get:

$$P(x|y) = \frac{P(y|x)}{P(y)} P(x)$$

- Why is this at all helpful?

- Lets us build one conditional from its reverse
 - Often one conditional is tricky but the other one is simple

- In the running for most important AI equation!



Quiz

■ Given:

$$P(D|W)$$

D	W	P
wet	sun	0.1
dry	sun	0.9

$$P(W)$$

R	P
sun	0.8
rain	0.2

■ What is $P(W \mid \text{dry})$?

Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

$$P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$$

- Example:

- M: meningitis, S: stiff neck

$$\left. \begin{array}{l} P(+m) = 0.0001 \\ P(+s|m) = 0.8 \\ P(+s|-m) = 0.01 \end{array} \right\} \text{Example givens}$$

$$P(+m|+s) = \frac{P(+s|m)P(+m)}{P(+s)} = \frac{P(+s|m)P(+m)}{P(+s|m)P(+m) + P(+s|-m)P(-m)} = \frac{0.8 \times 0.0001}{0.8 \times 0.0001 + 0.01 \times 0.999} = \frac{0.8 \times 0.0001}{0.8 \times 0.0001 + 0.01 \times 0.999}$$

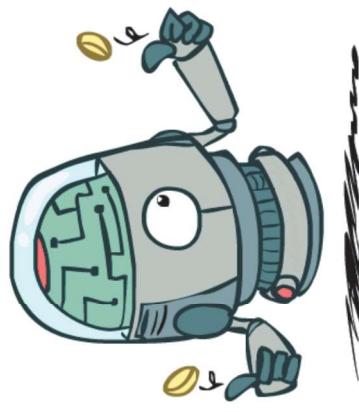
- Note: posterior probability of meningitis still very small
 - Note: you should still get stiff necks checked out! Why?

Independence

- Two variables are *independent* if:

$$\forall x, y : P(x, y) = P(x)P(y)$$

- This says that their joint distribution *factors* into a product two simpler distributions
- Another form: $\forall x, y : P(x|y) = P(x)$
- We write: $X \perp\!\!\!\perp Y$
- Independence is a simplifying *modeling assumption*
- What could we assume for {Weather, Traffic, Cavity, Toothache}?



Conditional Independence

- Unconditional (absolute) independence very rare
- *Conditional independence* is our most basic and robust form of knowledge about uncertain environments.

- X is conditionally independent of Y given Z

$$X \perp\!\!\!\perp Y | Z$$

if and only if:

$$\forall x, y, z : P(x, y | z) = P(x | z)P(y | z)$$

or, equivalently, if and only if

$$\forall x, y, z : P(x | z, y) = P(x | z)$$

Bayes' Net Semantics



Bayes' Net Semantics

- A set of nodes, one per variable X

- A directed, acyclic graph

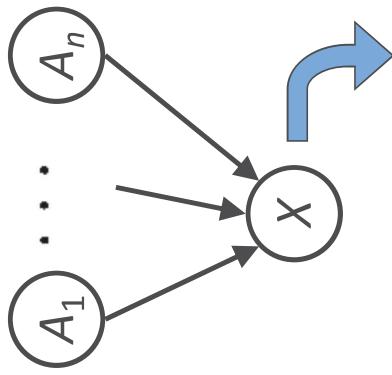
- A conditional distribution for each node

- A collection of distributions over X, one for each combination of parents' values

$$P(X|a_1 \dots a_n)$$

- CPT: conditional probability table

- Description of a noisy “causal” process



$$P(X|A_1 \dots A_n)$$

A Bayes net = Topology (graph) + Local Conditional Probabilities

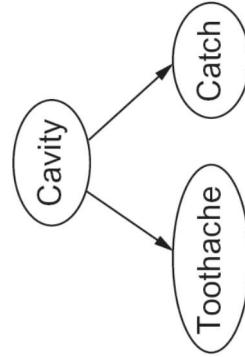
Probabilities in BNs

- Bayes' nets **implicitly** encode joint distributions

- As a product of local conditional distributions

- To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



- Example:

$$P(+\text{cavity}, +\text{catch}, -\text{toothache})$$

Probabilities in BNs

- Why are we guaranteed that setting

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

results in a proper joint distribution?

- Chain rule (valid for all distributions):

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1 \dots x_{i-1})$$

- Assume conditional independences:

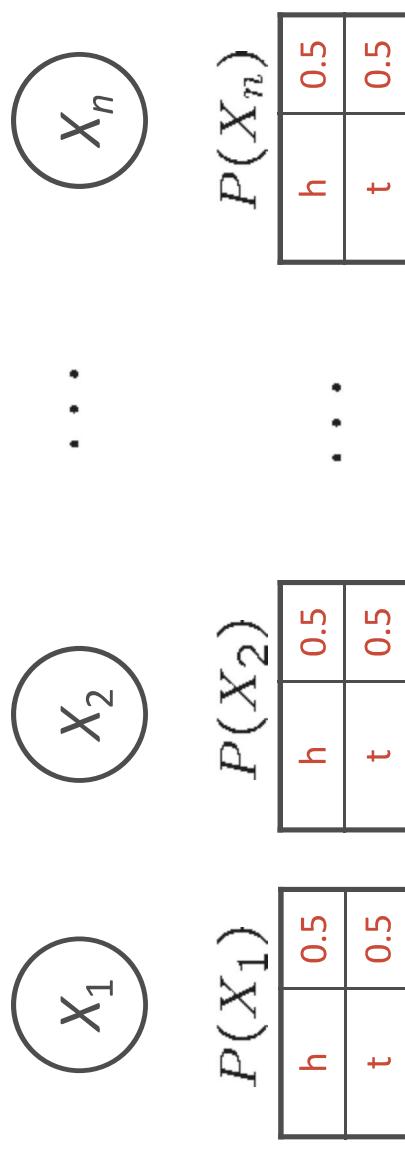
$$P(x_i | x_1, \dots, x_{i-1}) = P(x_i | \text{parents}(X_i))$$

$$\rightarrow \text{Consequence: } P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

- Not every BN can represent every joint distribution

- The topology enforces certain conditional independencies

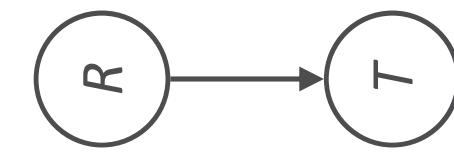
Example: Coin Flips



$$P(h, h, t, h) =$$

Only distributions whose variables are absolutely independent can be represented by a Bayes' net with no arcs.

Example: Traffic



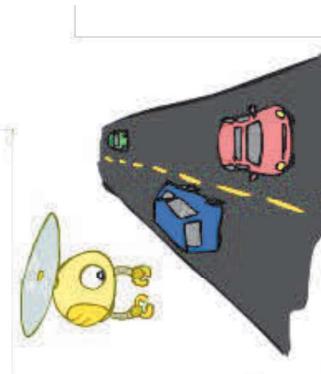
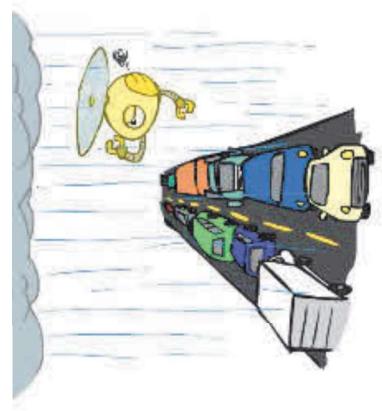
$$P(R)$$

+r	1/4
-r	3/4

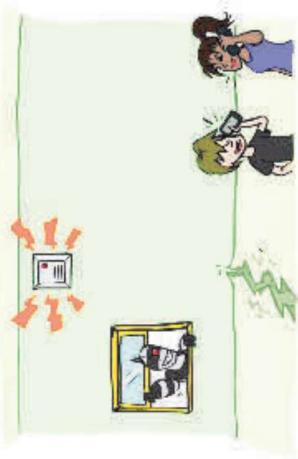
$$P(+r, -t) =$$

$$P(T|R)$$

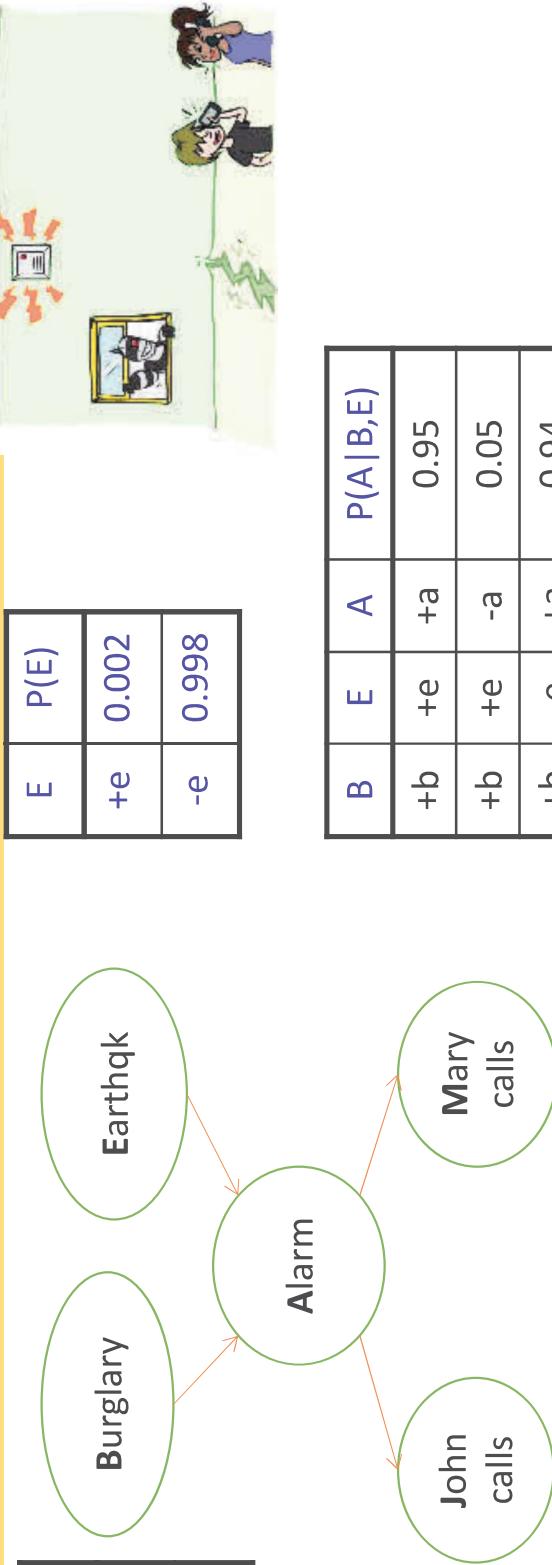
+t	3/4
-t	1/4
+t	1/2
-t	1/2



Example: Alarm Network



B	P(B)
+b	0.001
-b	0.999



B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

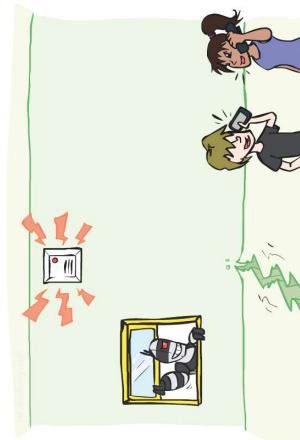
Example: Alarm Network

B	P(B)	P(J A)
+b	0.001	0.9
-b	0.999	

E	P(E)	P(A B,E)
+e	0.002	
-e	0.998	

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

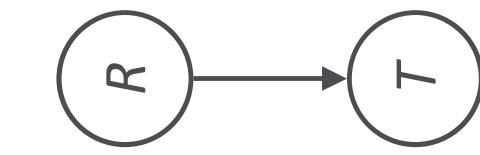


$$P(+b, -e, +a, -j, +m) =$$

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
-b	-e	+a	0.94
-b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

Example: Traffic

- Causal direction



$$P(R)$$

+r	1/4
-r	3/4

$$P(T|R)$$

+r	+t	3/4
-r	-t	1/4
-r	+t	1/2
-r	-t	1/2

$$P(T, R)$$

+r	+t	3/16
+r	-t	1/16
-r	+t	6/16
-r	-t	6/16

Example: Reverse Traffic

- Reverse causality?



$$P(T)$$

+t	9/16
-t	7/16

$$P(R|T)$$

+t	+r	1/3
-t	-r	2/3
-t	+r	1/7
-t	-r	6/7

$$P(T, R)$$

+r	+t	3/16
+r	-t	1/16
-r	+t	6/16
-r	-t	6/16

Size of a Bayes' Net

- How big is a joint distribution over N Boolean variables?
 2^N
- Both give you the power to calculate
 $P(X_1, X_2, \dots, X_n)$
- BNs: Huge space savings!
- How big is an N-node net if nodes have up to k parents?
 $O(N * 2^{k+1})$



Causality?

- When Bayes' nets reflect the true causal patterns:

- Often simpler (nodes have fewer parents)
- Often easier to think about
- Often easier to elicit from experts

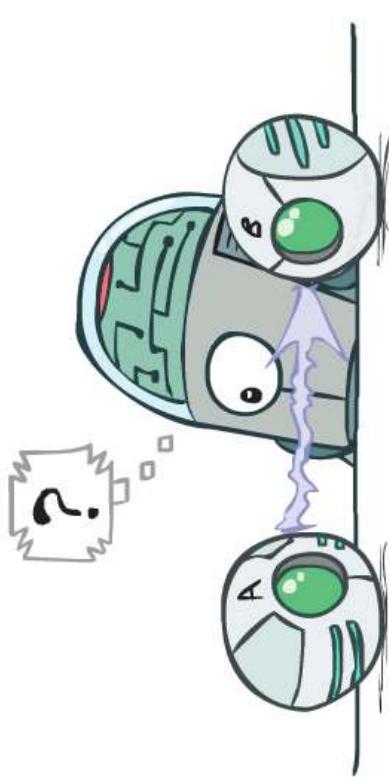
- BNs need not actually be causal

- Sometimes no causal net exists over the domain (especially if variables are missing)
- E.g. consider the variables *Traffic* and *Drips*
- End up with arrows that reflect correlation, not causation

- What do the arrows really mean?

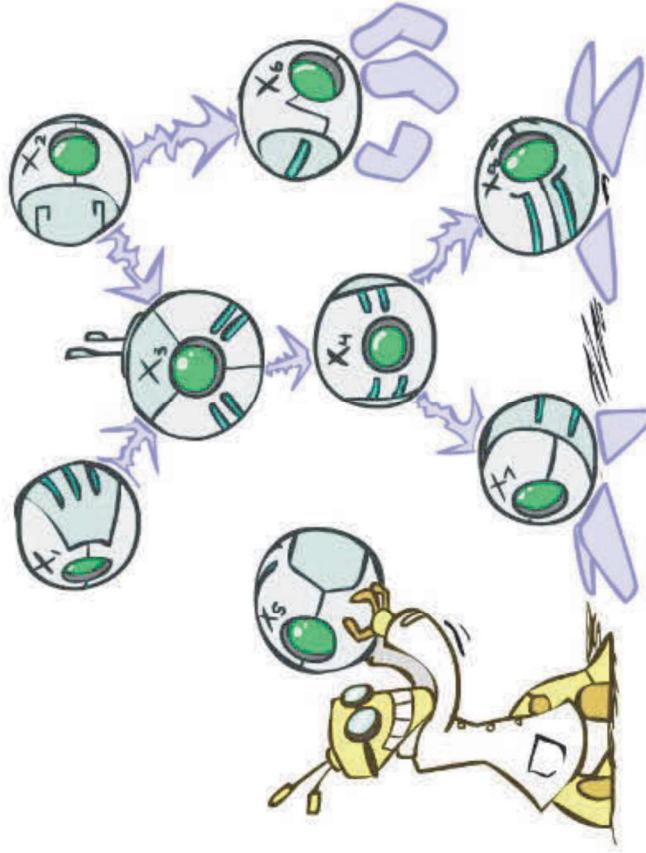
- Topology may happen to encode causal structure
- **Topology really encodes conditional independence**

$$P(x_i | x_1, \dots, x_{i-1}) = P(x_i | \text{parents}(X_i))$$



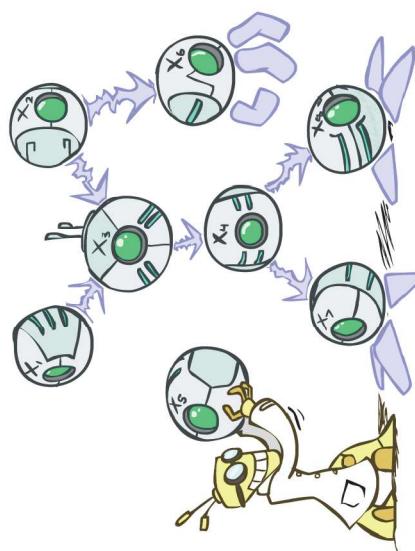
Bayes' Nets

- So far: how a Bayes' net encodes a joint distribution
- Next: how to answer queries about that distribution
- Today:
 - First assembled BNs using an intuitive notion of conditional independence as causality
 - Then saw that key property is conditional independence
 - Main goal: answer queries about conditional independence and influence
- After that: how to answer numerical queries (inference)



Bayes' Net Semantics

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
 - A collection of distributions over X_i , one for each combination of parents' values
$$P(X|a_1 \dots a_n)$$
 - Bayes' nets implicitly encode joint distributions



- As a product of local conditional distributions
- To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

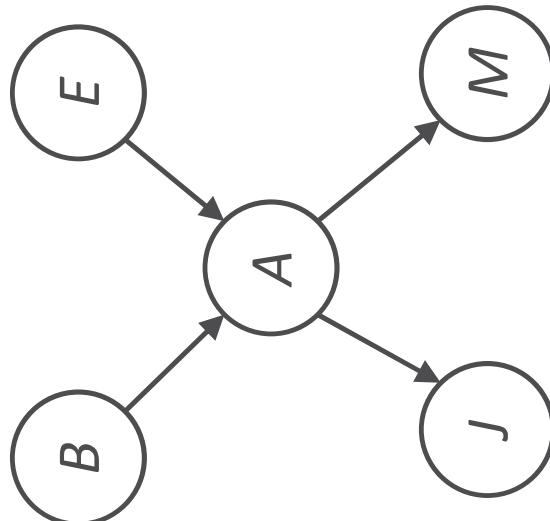
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

Example: Alarm Network

B	P(B)	P(E B)
+b	0.001	0.002
-b	0.999	0.998

E	P(E)
+e	0.002
-e	0.998

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95



A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
-b	-e	+a	0.94
-b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

$$\begin{aligned}
 P(+b, -e, +a, -j, +m) &= \\
 P(+b)P(-e)P(+a|+b, -e)P(-j|+a)P(+m|+a) &= \\
 0.001 \times 0.998 \times 0.94 \times 0.1 \times 0.7 &
 \end{aligned}$$

Conditional Independence

- X and Y are **independent** if

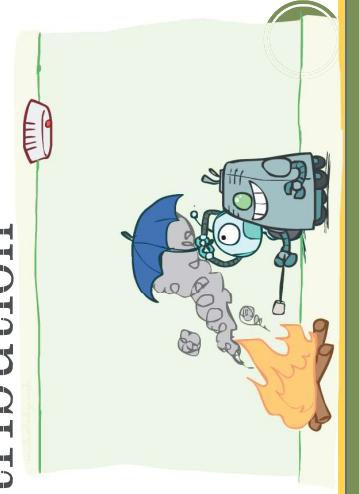
$$\forall x, y \quad P(x, y) = P(x)P(y) \dashrightarrow X \perp\!\!\!\perp Y$$

- X and Y are **conditionally independent** given Z

$$\forall x, y, z \quad P(x, y|z) = P(x|z)P(y|z) \dashrightarrow X \perp\!\!\!\perp Y|Z$$

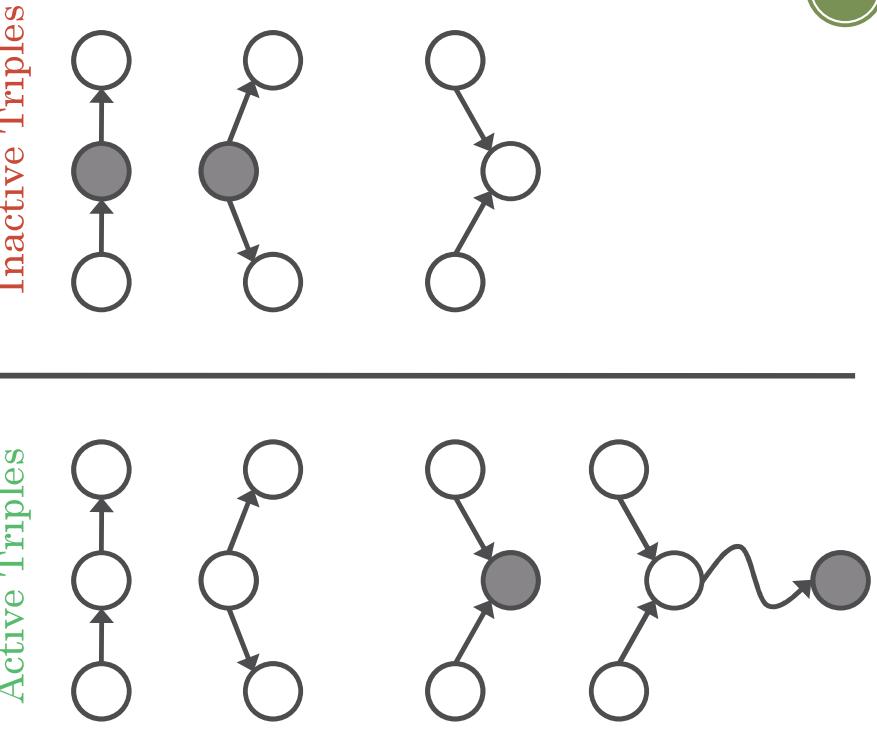
- (Conditional) independence is a property of a distribution

- Example:
 $Alarm \perp\!\!\!\perp Fire|Smoke$



Active / Inactive Paths

- Question: Are X and Y conditionally independent given evidence variables $\{Z\}$?
 - Yes, if X and Y “d-separated” by Z
 - Consider all (undirected) paths from X to Y
 - No active paths = independence!



- A path is active if each triple is active:
 - Causal chain $A \rightarrow B \rightarrow C$ where B is unobserved (either direction)
 - Common cause $A \leftarrow B \rightarrow C$ where B is unobserved
 - Common effect (aka v-structure)
 $A \rightarrow B \leftarrow C$ where B or one of its descendants is observed
- All it takes to block a path is a single inactive segment

D-Separation

▪ Query: $X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$?

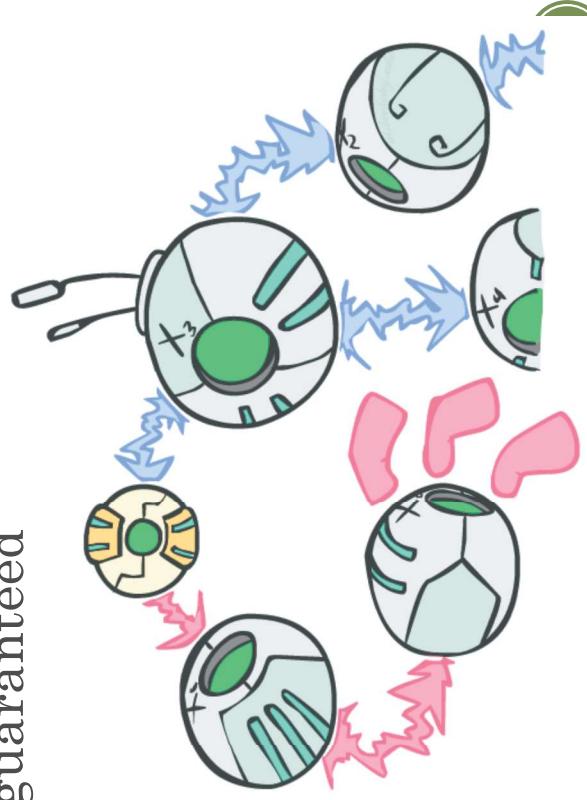
▪ Check all (undirected!) paths between X_i and X_j

▪ If one or more active, then independence not guaranteed

$$X_i \not\perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$

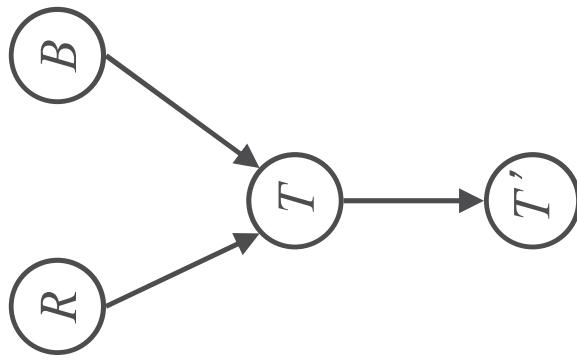
▪ Otherwise (i.e. if all paths are inactive),
then independence is guaranteed

$$X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$

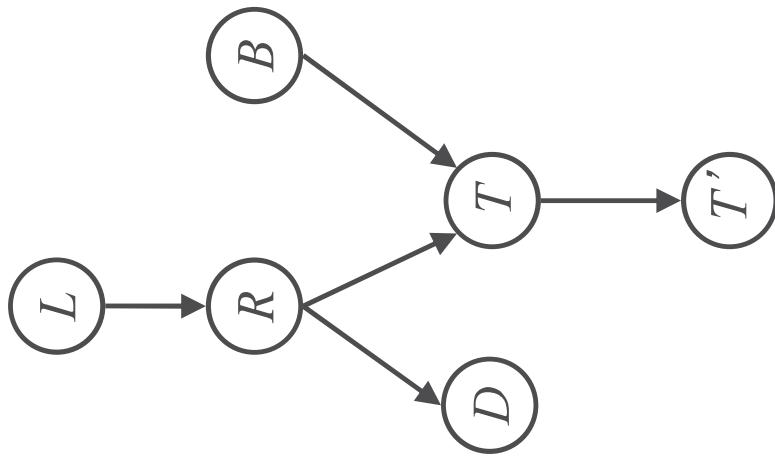


Example

$$\begin{aligned} R \perp\!\!\!\perp B \\ R \perp\!\!\!\perp B | T \\ R \perp\!\!\!\perp B | T' \end{aligned}$$



Example



$L \perp\!\!\!\perp T'|T$

$L \perp\!\!\!\perp B$

$L \perp\!\!\!\perp B|T$

$L \perp\!\!\!\perp B|T'$

$L \perp\!\!\!\perp B|T, R$

Yes

Yes

Yes

Example

- Variables:

- R: Raining
- T: Traffic
- D: Roof drips
- S: I'm sad

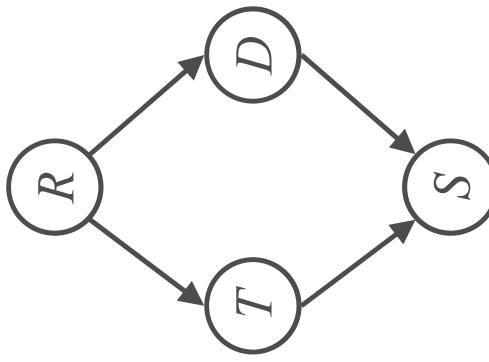
- Questions:

$$T \perp\!\!\!\perp D$$

$$T \perp\!\!\!\perp D|R$$

$$T \perp\!\!\!\perp D|R, S$$

Yes



Inference

- Inference: calculating some useful quantity from a joint probability distribution

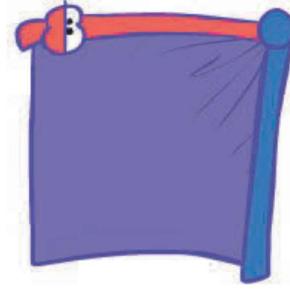
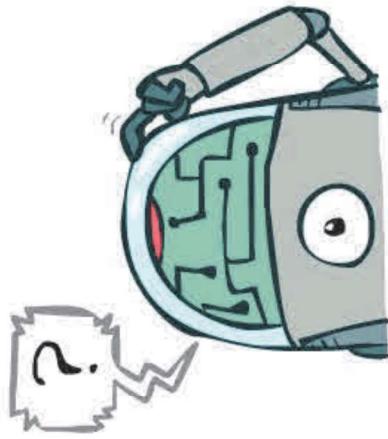
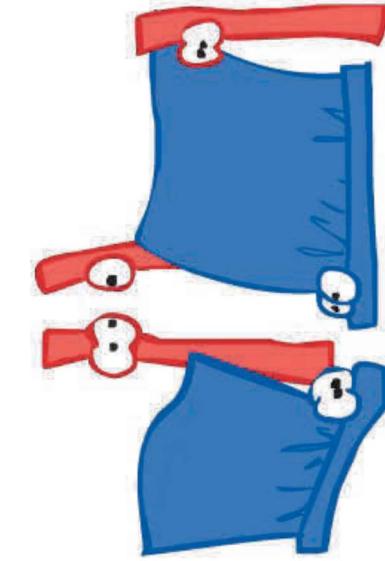
- Examples:

- Posterior probability

$$P(Q|E_1 = e_1, \dots, E_k = e_k)$$

- Most likely explanation:

$$\operatorname{argmax}_q P(Q = q | E_1 = e_1, \dots)$$



Inference by Enumeration

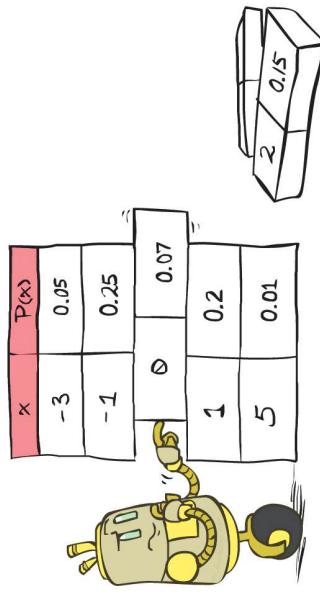
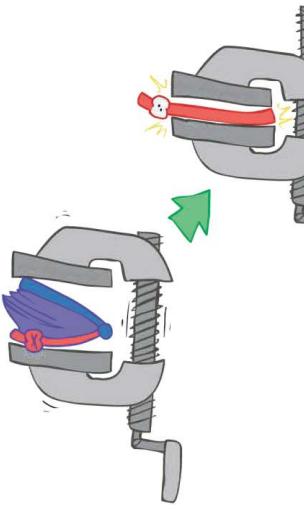
- General case:
 - Evidence variables: $E_1 \dots E_k = e_1 \dots e_k$
 - Query variable: Q
 - Hidden variables: $H_1 \dots H_r$

$$P(Q | e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

Z = \sum_q P(Q, e_1 \dots e_k)

- Step 1: Select the entries consistent with the evidence

- Step 2: Sum out H to get joint of Query and evidence

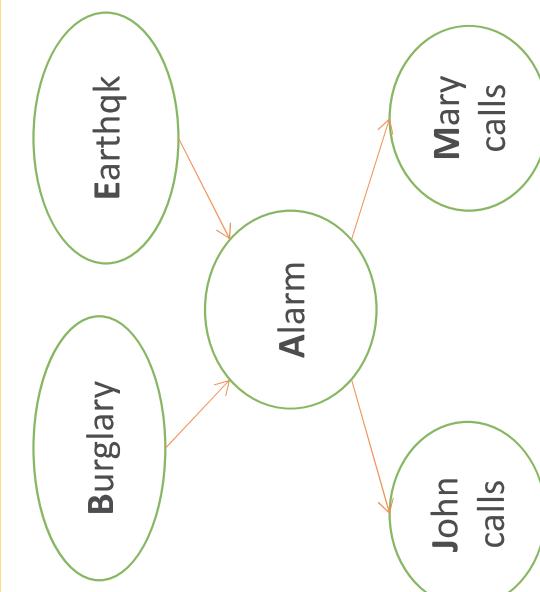


- We want: $P(Q | e_1 \dots e_k)$
- Step 3: Normalize

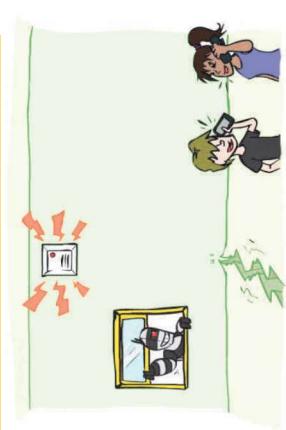
$$P(Q | e_1 \dots e_k) = \frac{1}{Z} \sum_q P(Q, e_1 \dots e_k)$$

Example: Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998



B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06

A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

Inference by Enumeration in Bayes' Net

- Given unlimited time, inference in BNs is easy

- Reminder of inference by enumeration by example:

$$P(B \mid +j, +m) \propto_B P(B, +j, +m)$$

$$= \sum_{e,a} P(B, e, a, +j, +m)$$

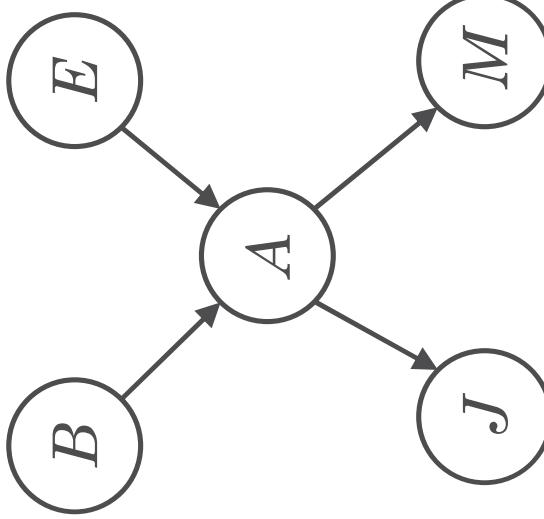
Sum-out hidden variables

$$= \sum_{e,a} P(B)P(e)P(a|B, e)P(+j|a)P(+m|a)$$

Select entries consistent with evidences

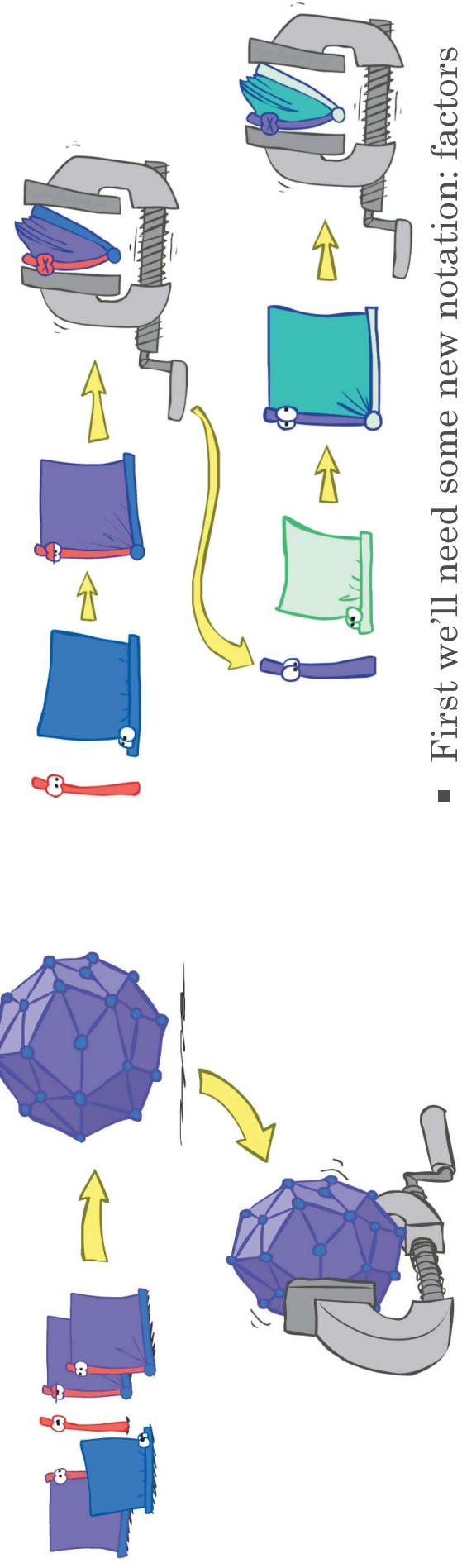
$$= P(B)P(+e)P(+a|B, +e)P(+j|+a)P(+m|+a) + P(B)P(+e)P(-a|B, +e)P(+j|-a)P(+m|-a)$$

$$P(B)P(-e)P(+a|B, -e)P(+j|+a)P(+m|+a) + P(B)P(-e)P(-a|B, -e)P(+j|-a)P(+m|-a)$$



Inference by Enumeration VS. Variable Elimination

- Why is inference by enumeration so slow?
 - You join up the whole joint distribution before you sum out the hidden variables
- Idea: interleave joining and marginalizing!
 - Called “Variable Elimination”
 - Still NP-hard, but usually much faster than inference by enumeration



- First we'll need some new notation: factors

Factor Zoo I

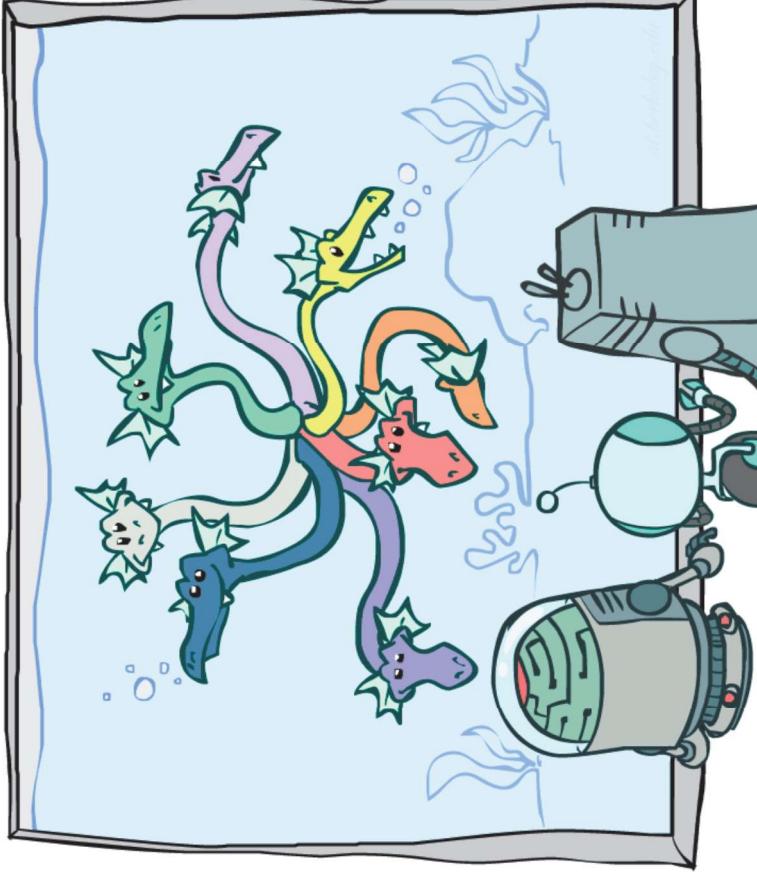
- Joint distribution: $P(X, Y)$
 - Entries $P(x,y)$ for all x, y
 - Sums to 1

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

- Selected joint: $P(x, Y)$
 - A slice of the joint distribution
 - Entries $P(x,y)$ for fixed x , all y
 - Sums to $P(x)$

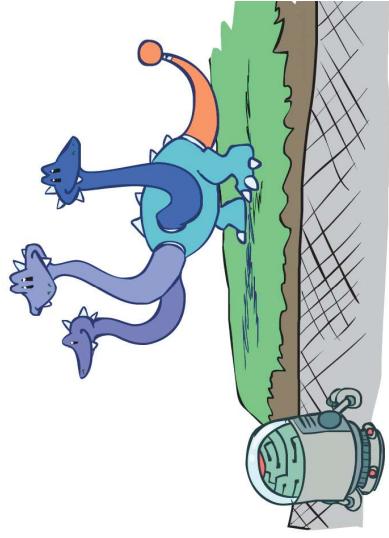
T	W	P
cold	sun	0.2
cold	rain	0.3

- Number of capitals = dimensionality of the table



Factor Zoo III

- Single conditional: $P(Y \mid x)$
- Entries $P(y \mid x)$ for fixed x , all y
- Sums to 1

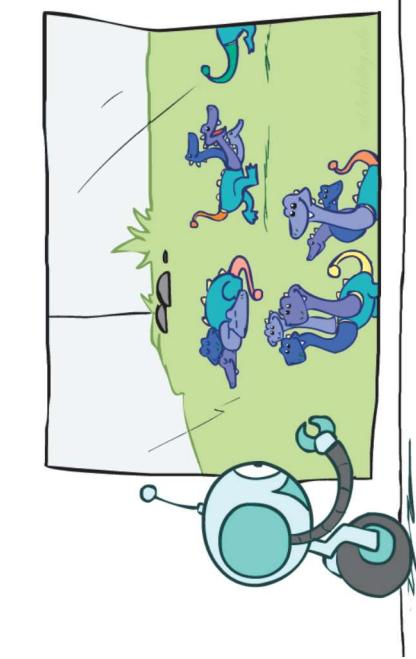


$P(W|cold)$

T	W	P
cold	sun	0.4
cold	rain	0.6

$P(W|T)$

T	W	P
hot	sun	0.8
hot	rain	0.2
cold	sun	0.4
cold	rain	0.6



- Family of conditionals:

$P(Y \mid X)$

- Multiple conditionals
- Entries $P(y \mid x)$ for all x, y
- Sums to $|X|$

$P(W|cold)$

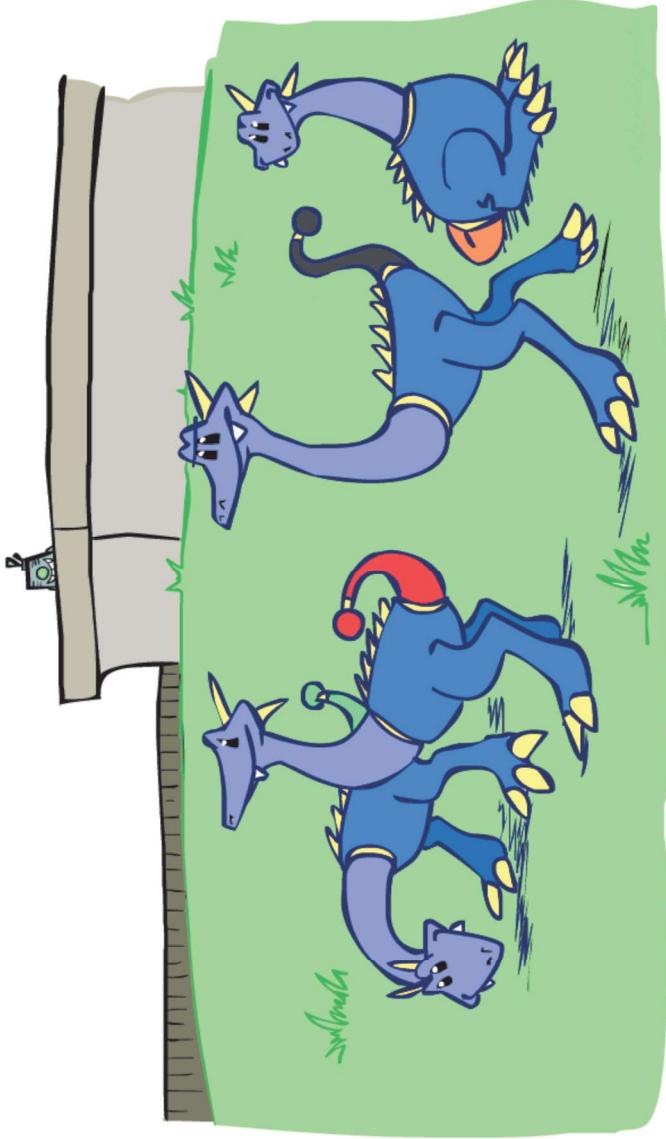
Factor Zoo III

- Specified family: $P(y \mid X)$
 - Entries $P(y \mid x)$ for fixed y ,
but for all x
 - Sums to ... who knows!

$P(\text{rain} \mid T)$

T	W	P
hot	rain	0.2
cold	rain	0.6

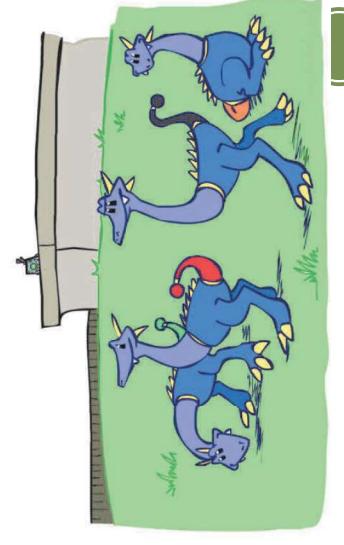
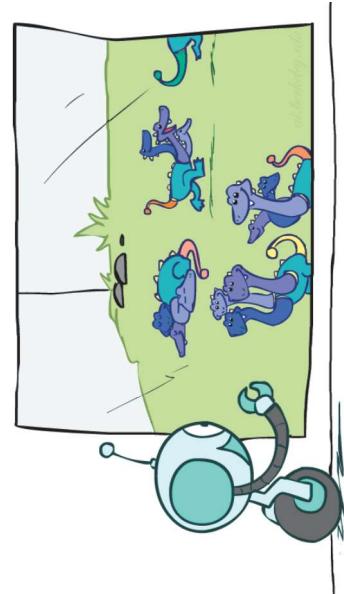
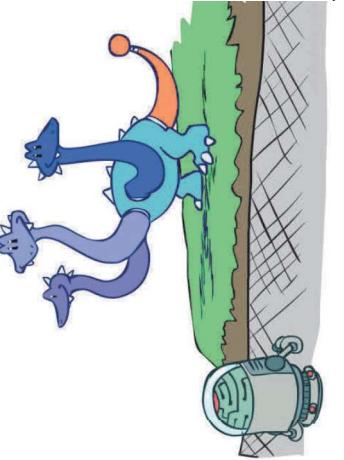
$$\left. \begin{array}{l} P(\text{rain} \mid \text{hot}) \\ P(\text{rain} \mid \text{cold}) \end{array} \right\}$$



Factor Zoo Summary

- In general, when we write $P(Y_1 \dots Y_N | X_1 \dots X_M)$

- It is a “factor,” a multi-dimensional array
- Its values are $P(y_1 \dots y_N | x_1 \dots x_M)$
- Any assigned (=lower-case) X or Y is a dimension missing (selected) from the array



Example: Traffic Domain

- Random Variables

- R: Raining

- T: Traffic

- L: Late for class!



$P(R)$

+r	0.1
-r	0.9

$P(T|R)$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$P(L|T)$

+t	+	0.3
+t	-	0.7
-t	+	0.1
-t	-	0.9

$$P(L) = ?$$

$$= \sum_{r,t} P(r, t, L)$$

$$= \sum_{r,t} P(r) P(t|r) P(L|t)$$

Inference by Enumeration: Procedural Outline

- Track objects called factors

- Initial factors are local CPTs (one per node)

$$P(R) \quad P(T|R)$$

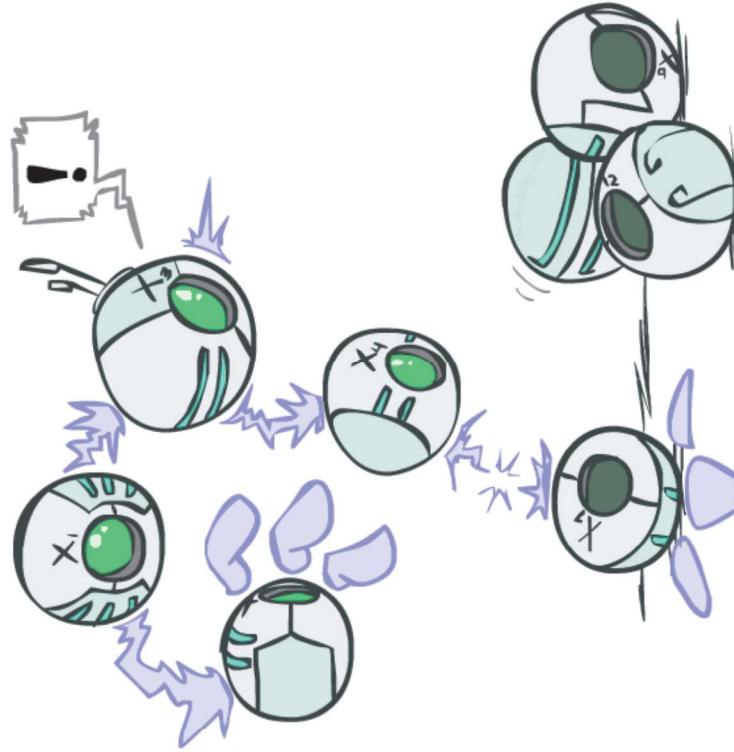
+r	0.1	+r	0.8	+t	+l	0.3
-r	0.9	+r	0.2	+t	-l	0.7
		-r	0.1	-t	+l	0.1
		-r	0.9	-t	-l	0.9

- Any known values are selected

- E.g. if we know $L = +\ell$ the initial factors are

$$P(R) \quad P(T|R)$$

+r	0.1	+r	0.8	+t	+l	0.3
-r	0.9	+r	0.2	-t	+l	0.1
		-r	0.1	+t	0.1	
		-r	0.9	-t	-l	0.9



- Procedure: Join all factors, eliminate all hidden variables, normalize

Operation 1: Join Factors

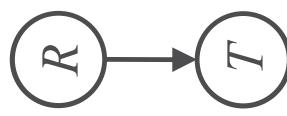
- First basic operation: joining factors

- Combining factors:

- **Just like a database join**

- Get all factors over the joining variable
- Build a new factor over the union of the variables involved

- Example: Join on R



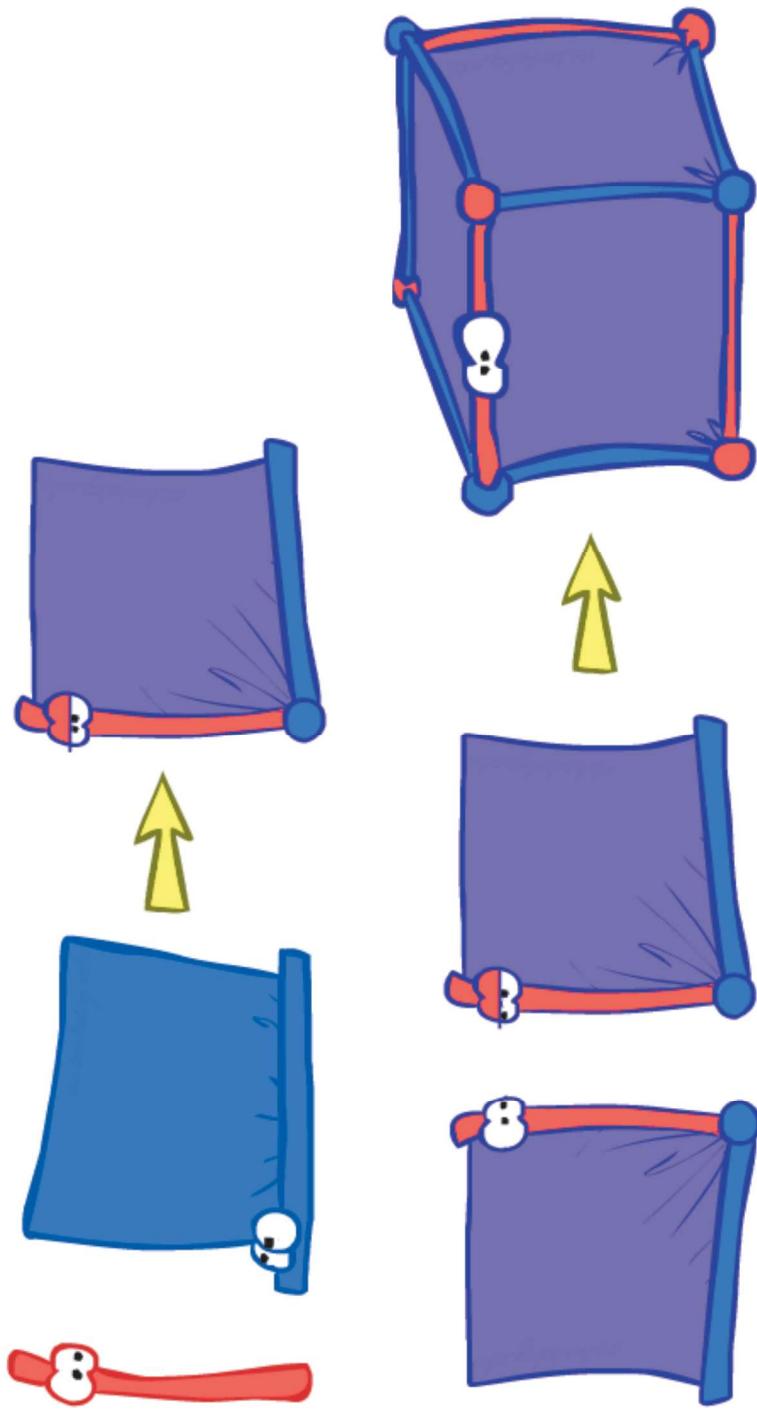
$$P(R) \times P(T|R)$$

	+r	-r	
+t	0.1	0.9	
-t	0.9	0.1	

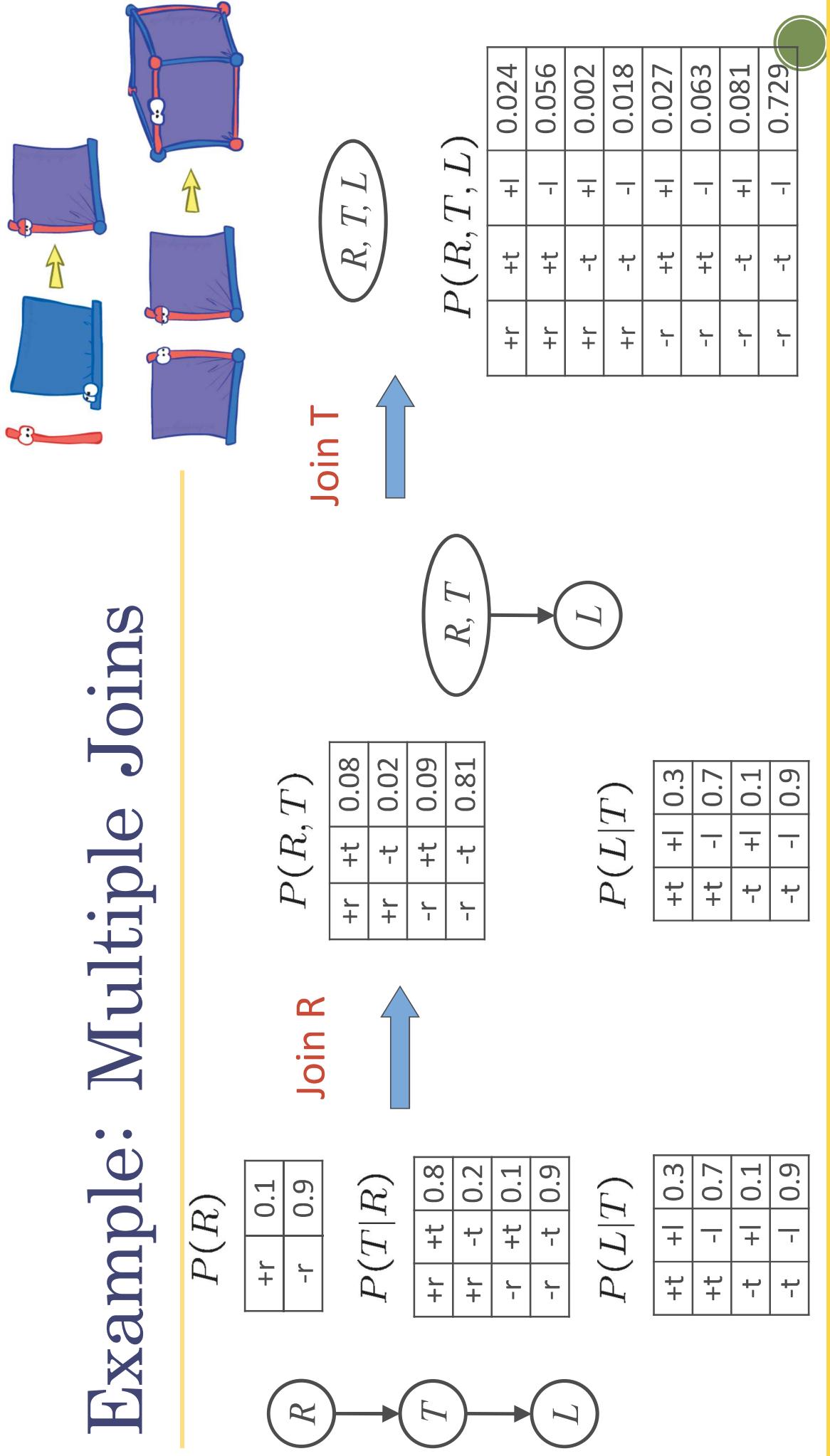
	+r	-r	
+t	0.08	0.02	
-t	0.02	0.09	

- Computation for each entry: pointwise products $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$

Example: Multiple Joins

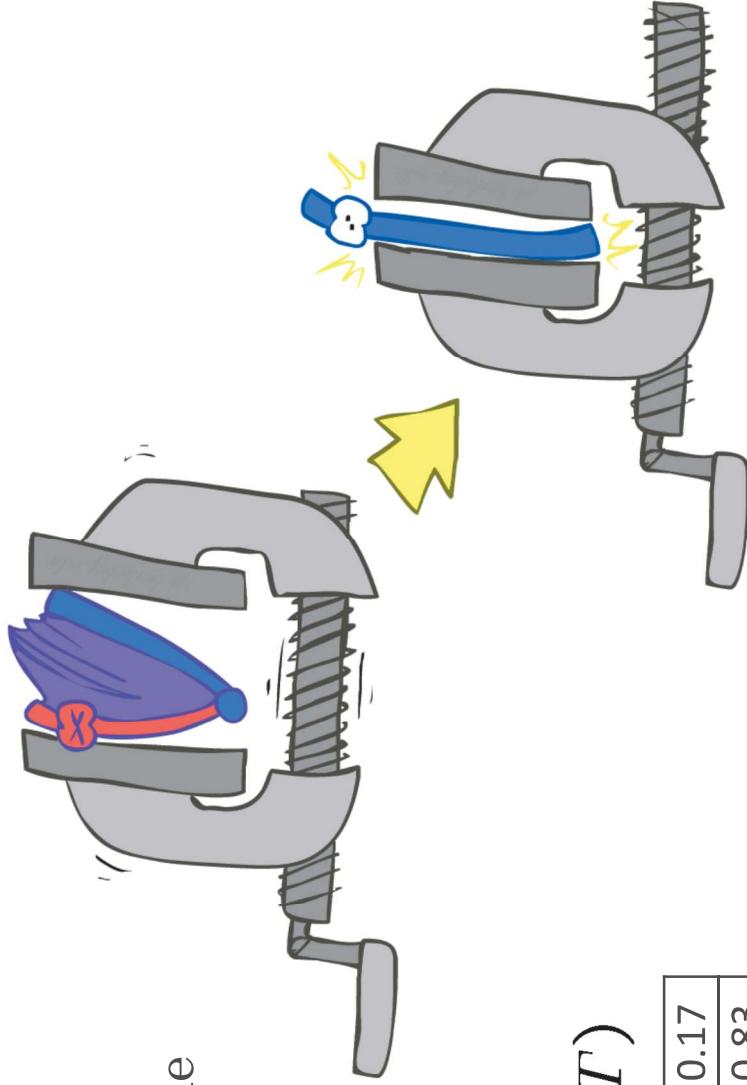


Example: Multiple Joins



Operation 2: Eliminate

- Second basic operation:
marginalization



- Take a factor and sum out a variable
 - Shrinks a factor to a smaller one
 - A projection operation

- Example:

$$P(R, T)$$

sum R

$$P(T)$$

+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

+t	0.17
-t	0.83

Multiple Elimination

$P(R, T, L)$	$+r$	$+t$	$+l$	0.024
	$+r$	$+t$	$-l$	0.056
	$+r$	$-t$	$+l$	0.002
	$+r$	$-t$	$-l$	0.018
	$+r$	$+t$	$+l$	0.027
	$-r$	$+t$	$-l$	0.063
	$-r$	$+t$	$+l$	0.081
	$-r$	$-t$	$-l$	0.729

Sum
out R

$P(T, L)$

Sum
out T

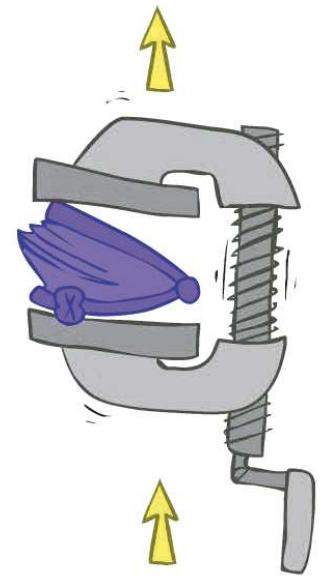
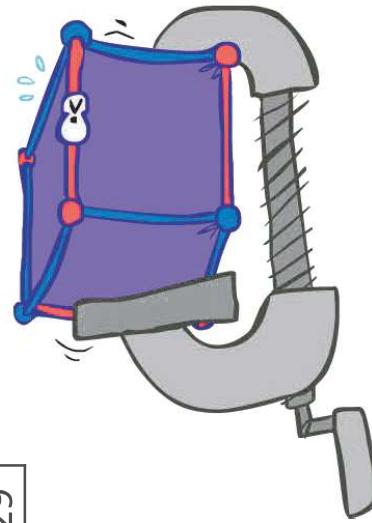
$+l$	0.134
$-l$	0.886

$P(L)$

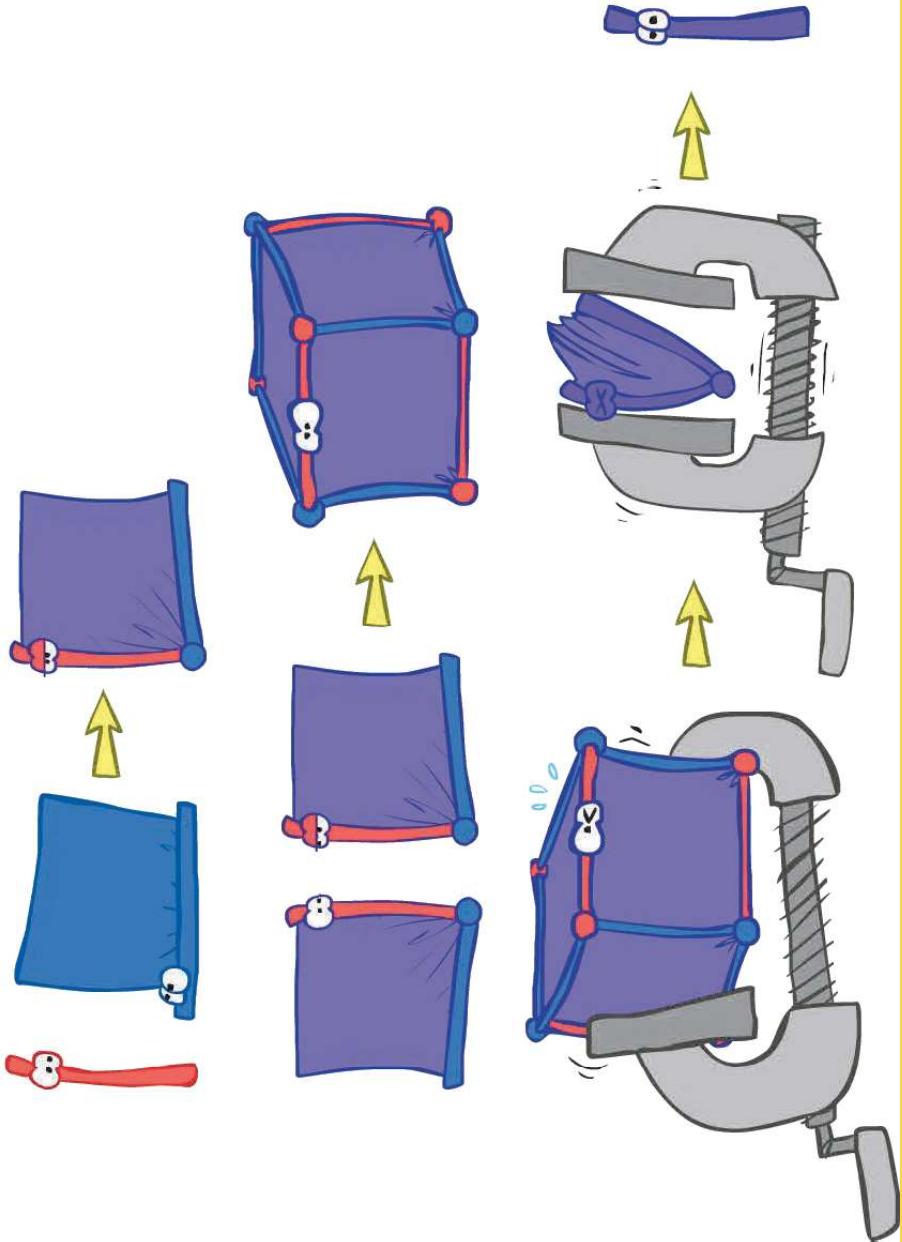
L

T, L

R, T, L



Thus Far: Multiple Join, Multiple Eliminate (= Inference by Enumeration)



Inference by Enumeration

- General case:
 - Evidence variables:
 - Query* variable:
 - Hidden variables:

$$\left. \begin{array}{l} E_1 \dots E_k = e_1 \dots e_k \\ Q \\ H_1 \dots H_r \end{array} \right\} X_1, X_2, \dots X_n$$

All variables

- * Works fine with multiple query variables, too
- We want:
 $P(Q|e_1 \dots e_k)$

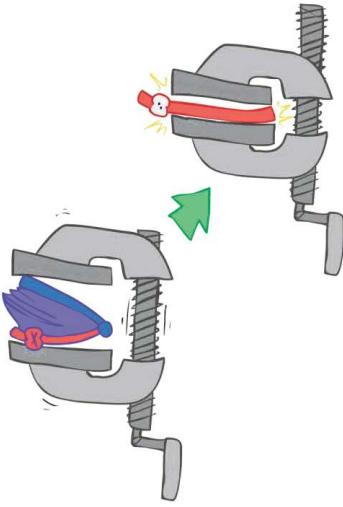
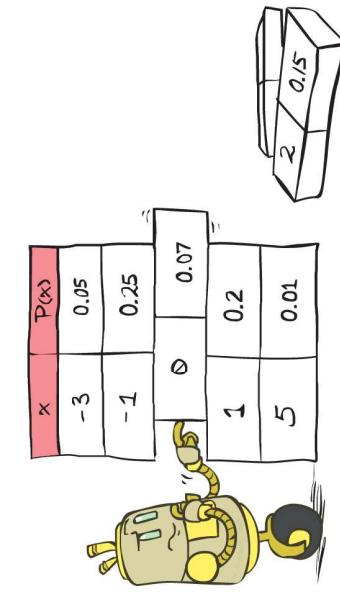
- Step 1: Select the entries consistent with the evidence

- Step 2: Sum out H to get joint of Query and evidence

- Step 3: Normalize

$$P(Q|e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

$$Z = \sum_q P(Q, e_1 \dots e_k)$$



Inference by Enumeration: Procedural Outline

- Track objects called factors

- Initial factors are local CPTs (one per node)

$$P(R) \quad P(T|R)$$

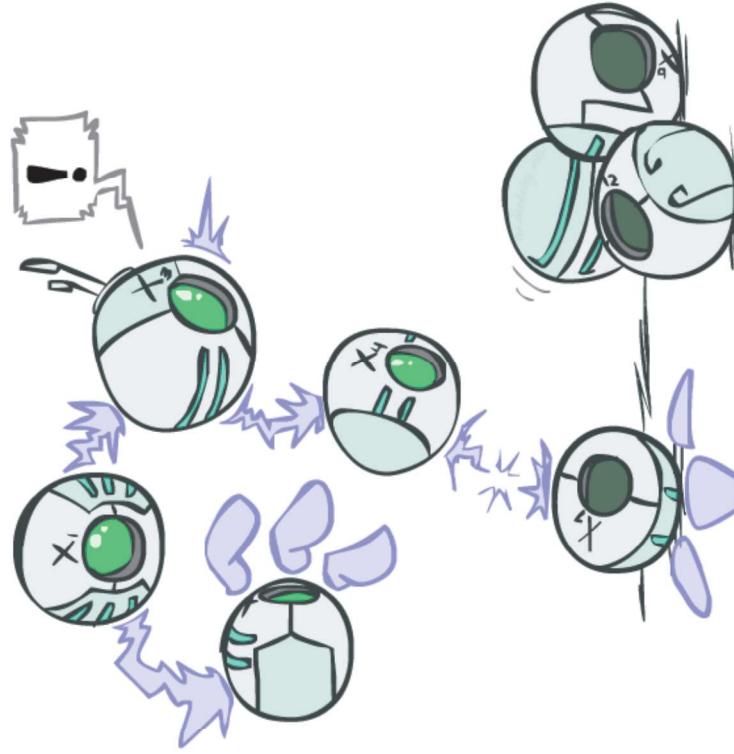
+r	0.1	+r	0.8	+t	+l	0.3
-r	0.9	+r	0.2	+t	-l	0.7
		-r	0.1	-t	+l	0.1
		-r	0.9	-t	-l	0.9

- Any known values are selected

- E.g. if we know $L = +\ell$ the initial factors are

$$P(R) \quad P(T|R)$$

+r	0.1	+r	0.8	+t	+l	0.3
-r	0.9	+r	0.2	-t	+l	0.1
		-r	0.1	+t	0.1	
		-r	0.9	-t	-l	0.9



- Procedure: Join all factors, eliminate all hidden variables, normalize

Operation 1 : Join Factors

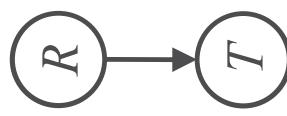
- First basic operation: joining factors

- Combining factors:

- **Just like a database join**

- Get all factors over the joining variable
- Build a new factor over the union of the variables involved

- Example: Join on R



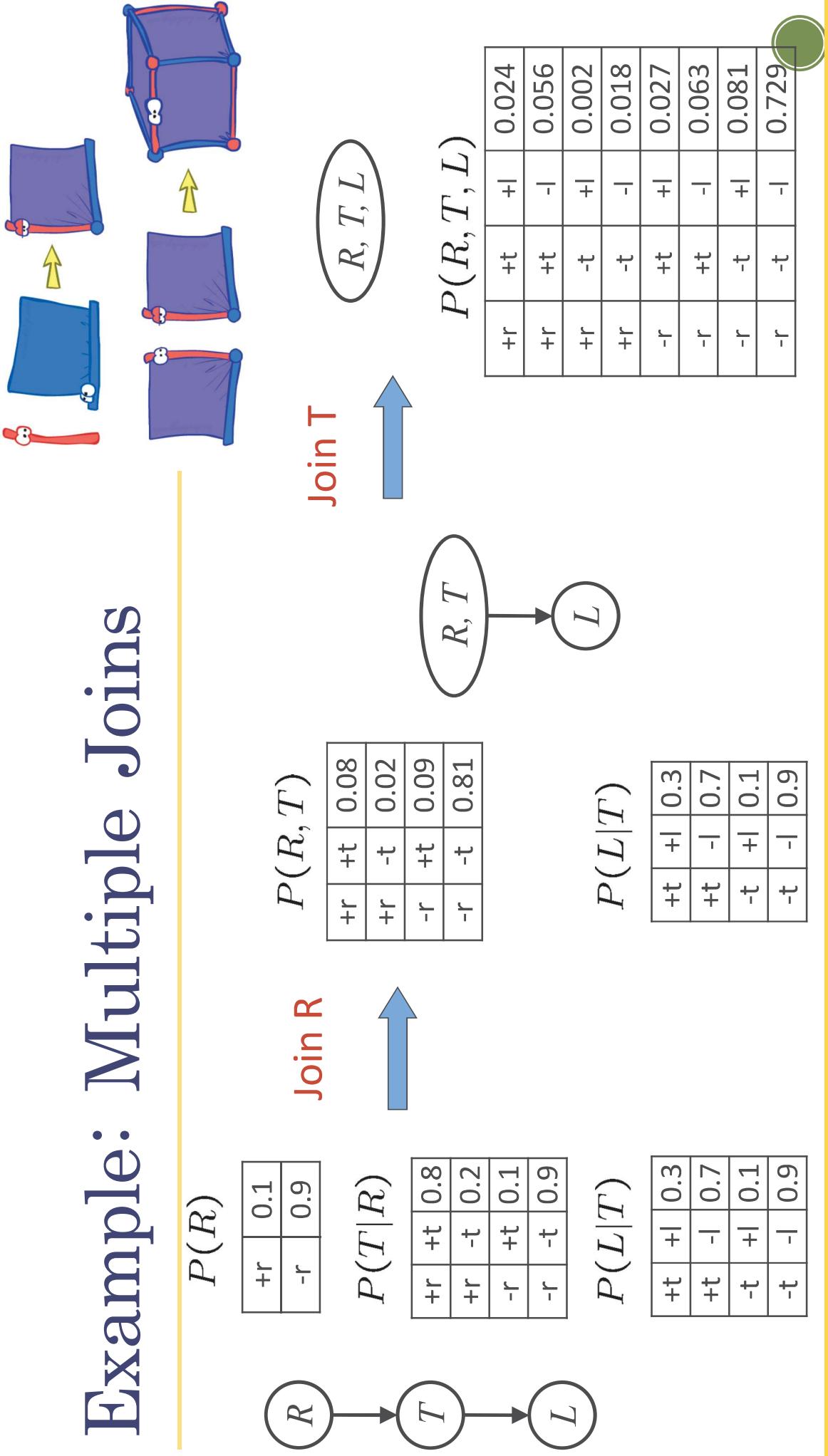
$$P(R) \times P(T|R)$$

	+r	-r	
+t	0.1	0.9	
-t	0.9	0.1	

	+r	-r	
+t	0.08	0.02	
-t	0.02	0.09	

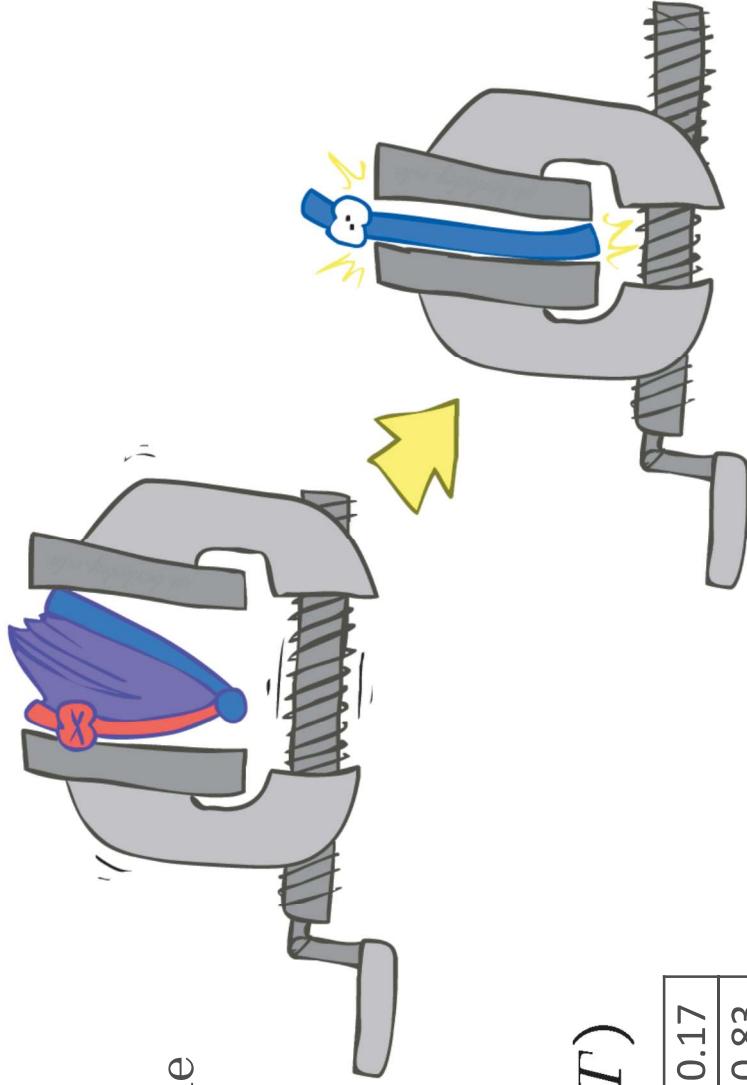
- Computation for each entry: pointwise products $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$

Example: Multiple Joins



Operation 2: Eliminate

- Second basic operation:
marginalization



- Take a factor and sum out a variable
 - Shrinks a factor to a smaller one
 - A projection operation

- Example:

$$P(R, T)$$

sum R

$$P(T)$$

+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

+t	0.17
-t	0.83

Multiple Elimination

$P(R, T, L)$	$+r$	$+t$	$+l$	0.024
	$+r$	$+t$	$-l$	0.056
	$+r$	$-t$	$+l$	0.002
	$+r$	$-t$	$-l$	0.018
	$+r$	$+t$	$+l$	0.027
	$-r$	$+t$	$-l$	0.063
	$-r$	$+t$	$+l$	0.081
	$-r$	$-t$	$-l$	0.729

Sum
out R

$P(T, L)$

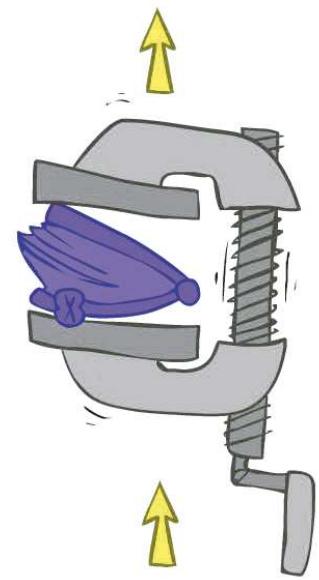
Sum
out T

$P(L)$

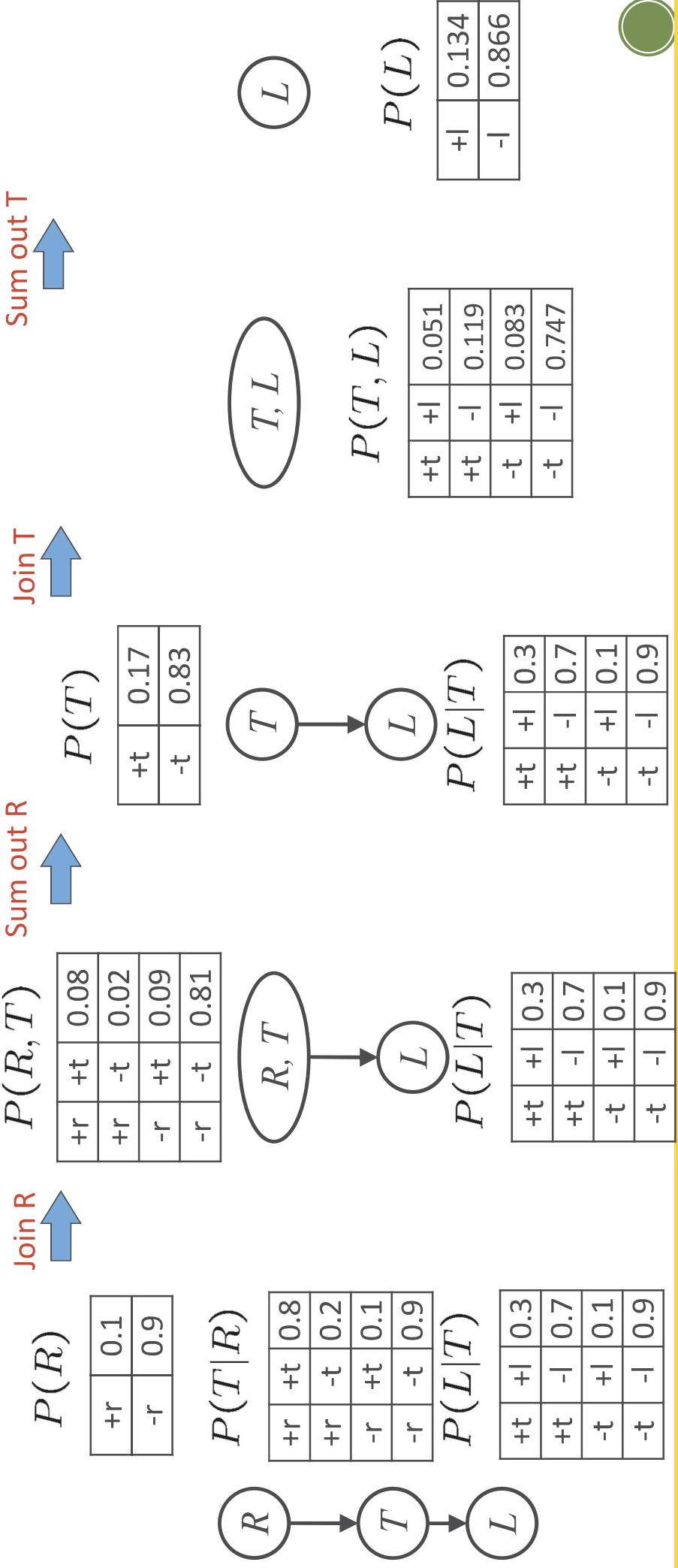
R, T, L

T, L

L



Marginalizing Early! (aka VE)



Evidence

- If evidence, start with factors that select that evidence
 - No evidence uses these initial factors:

$$P(R) \quad P(T|R)$$

	+r	0.1
+r		0.8
-r		0.2

	+t	0.8
+r		0.2
-r		0.1

$$P(L|T) \quad P(T|R)$$

	+t	0.3
+t		0.7
-t		0.1

- Computing $P(L|+r)$ the initial factors become:

$$P(+r) \quad P(T|+r)$$

	+r	0.1
+r		0.8
-r		0.2



- We eliminate all vars other than query + evidence

Evidence II

- Result will be a selected joint of query and evidence
 - E.g. for $P(L \mid +r)$, we would end up with:

$$P(+r, L)$$

+r	+l	0.026
+r	-l	0.074

Normalize
↑

$$P(L \mid +r)$$

+l	0.26
-l	0.74

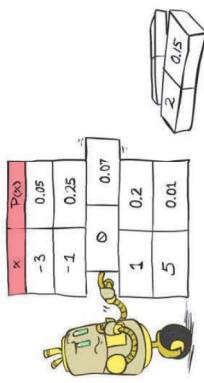
- To get our answer, just normalize this!

- That's it!



General Variable Elimination

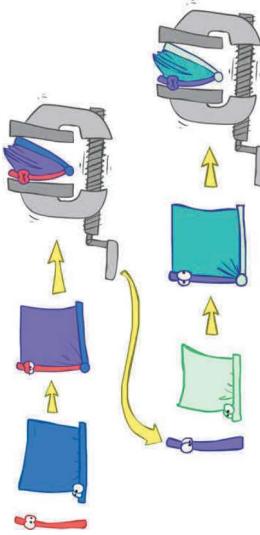
▪ Query: $P(Q|E_1 = e_1, \dots, E_k = e_k)$



X	P(X)
-3	0.05
-1	0.25
0	0.07
1	0.2
5	0.01

- Start with initial factors:
 - Local CPTs (but instantiated by evidence)

- While there are still hidden variables (not Q or evidence):
 - Pick a hidden variable H
 - Join all factors mentioning H
 - Eliminate (sum out) H

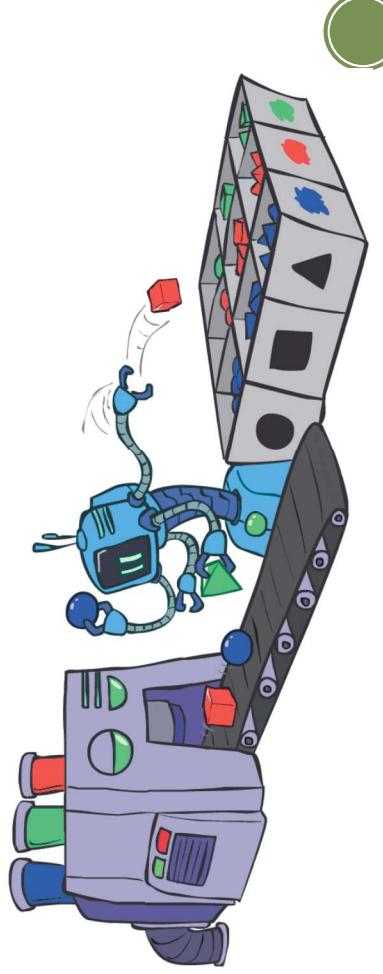


- Join all remaining factors and normalize


$$\frac{1}{Z} \times \begin{array}{|c|c|}\hline X & P(X) \\ \hline \text{Blue} & 0.05 \\ \text{Green} & 0.25 \\ \text{Purple} & 0.2 \\ \text{Yellow} & 0.01 \\ \hline \end{array}$$

Sampling

- Sampling is a lot like repeated simulation
 - Why sample?
- Predicting the weather, basketball games, ...
- Basic idea
- Draw N samples from a sampling distribution S
- Compute an approximate posterior probability
- Show this converges to the true probability P



Sampling

- Sampling from given distribution

- Step 1: Get sample u from uniform distribution over $[0, 1]$
 - E.g. `random()` in python
- Step 2: Convert this sample u into an outcome for the given distribution
 - Each target outcome is associated with a sub-interval of $[0,1]$
 - Sub-interval size is equal to probability of the outcome.

- Example

C	P(C)
red	0.6
green	0.1
blue	0.3

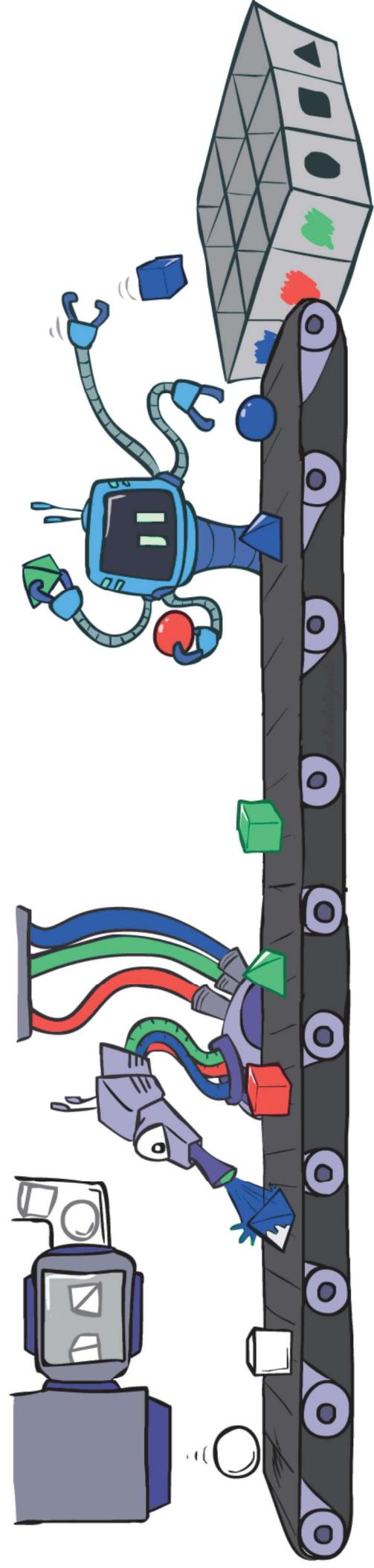
- If `random()` returns $u = 0.83$, then our sample is $C = \text{blue}$
- E.g, after sampling 8 times:



Sampling in Bayes' Nets

- Prior Sampling
- Rejection Sampling
- Likelihood Weighting
- Gibbs Sampling

Prior Sampling

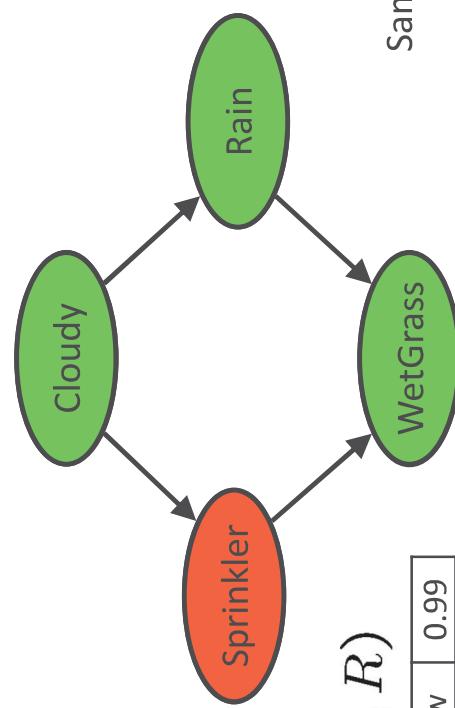


Prior Sampling

$$P(C) = \begin{array}{|c|c|} \hline +c & 0.5 \\ \hline -c & 0.5 \\ \hline \end{array}$$

$$P(S|C) = \begin{array}{|c|c|c|} \hline +c & +s & 0.1 \\ \hline & -s & 0.9 \\ \hline -c & +s & 0.5 \\ \hline & -s & 0.5 \\ \hline \end{array}$$

$$P(R|C) = \begin{array}{|c|c|c|} \hline +c & +r & 0.8 \\ \hline & -r & 0.2 \\ \hline -c & +r & 0.2 \\ \hline & -r & 0.8 \\ \hline \end{array}$$



$$P(W|S, R) = \begin{array}{|c|c|c|c|} \hline +s & +r & +w & 0.99 \\ \hline & -r & +w & 0.01 \\ \hline -s & +r & +w & 0.90 \\ \hline & -r & -w & 0.10 \\ \hline \end{array}$$

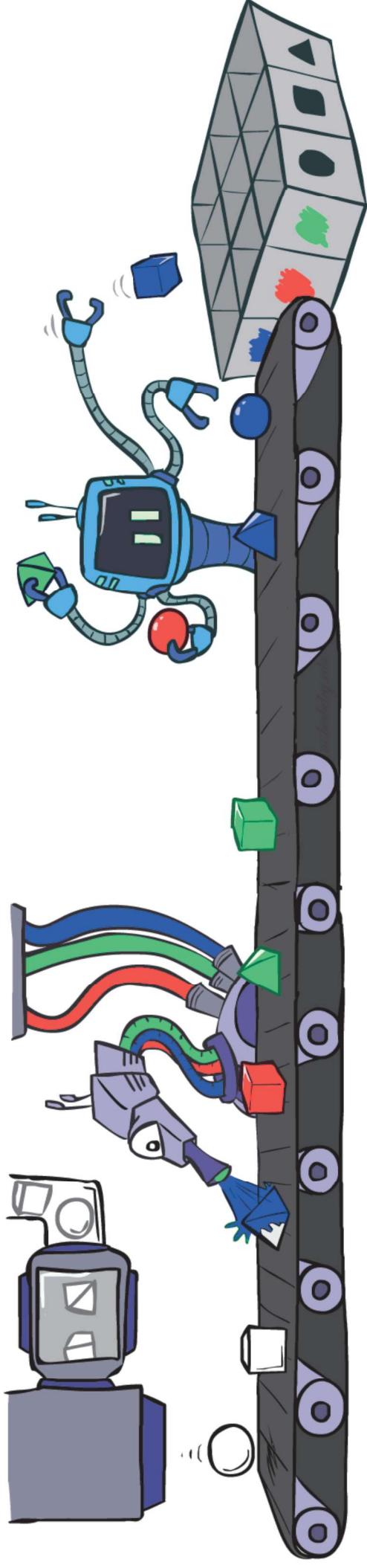
Samples:

+c, -s, +r, +w
-c, +s, -r, +w
...



Prior Sampling

- For $i = 1, 2, \dots, n$
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
 - Return (x_1, x_2, \dots, x_n)



Prior Sampling

- This process generates samples with probability:

$$S_{PS}(x_1 \dots x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i)) = P(x_1 \dots x_n)$$

...i.e. the BN's joint probability

- Let the number of samples of an event $\bar{N}_{PS}(x_1 \dots x_n)$

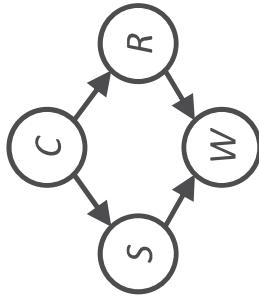
$$\begin{aligned}\text{■ Then } \lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N \\ &= S_{PS}(x_1, \dots, x_n) \\ &= P(x_1 \dots x_n)\end{aligned}$$

- I.e., the sampling procedure is **consistent**

Example

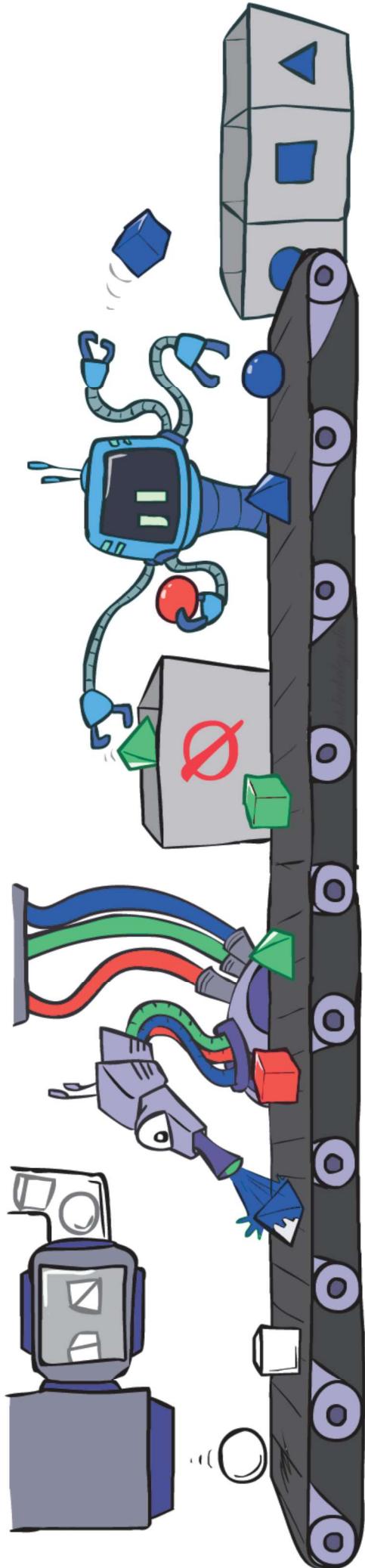
- We'll get a bunch of samples from the BN:

+c, -s, +r, +w
+c, +s, +r, +w
-c, +s, +r, -w
+c, -s, +r, +w
-c, -s, -r, +w



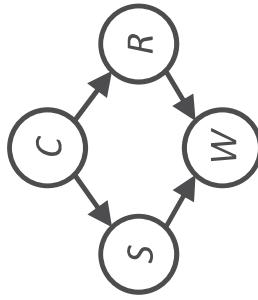
- If we want to know $P(W)$
 - We have counts $<+w:4, -w:1>$
 - Normalize to get $P(W) = <+w:0.8, -w:0.2>$
 - This will get closer to the true distribution with more samples
 - Can estimate anything else, too
 - What about $P(C \mid +w)$? $P(C \mid +r, +w)$? $P(C \mid -r, -w)$?
 - Fast: can use fewer samples if less time (what's the drawback?)

Rejection Sampling



Rejection Sampling

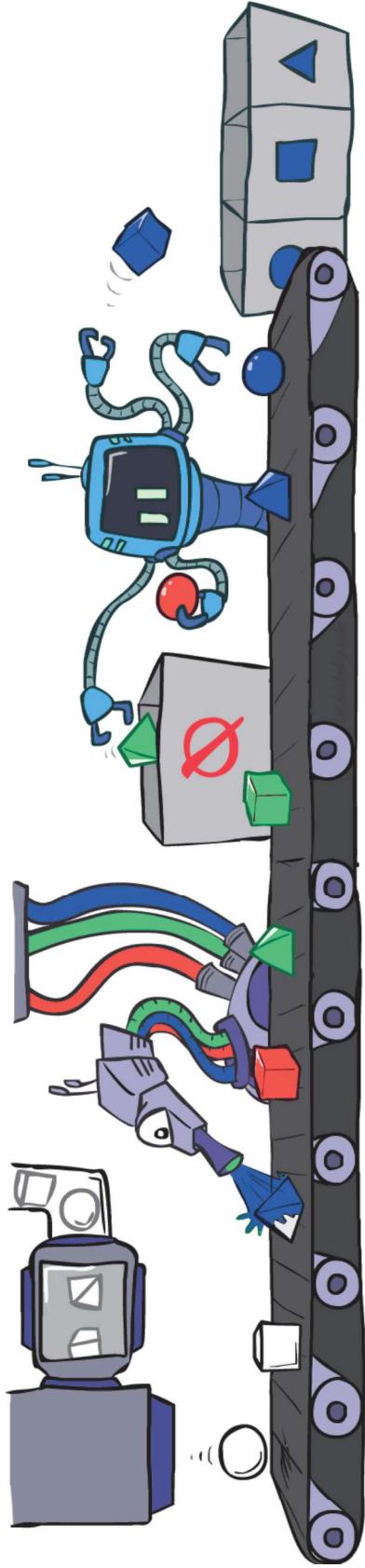
- Let's say we want $P(C)$
 - No point keeping all samples around
 - Just tally counts of C as we go
- Let's say we want $P(C \mid +S)$
 - Same thing: tally C outcomes, but ignore (reject) samples which don't have $S=+S$
 - This is called rejection sampling
 - It is also consistent for conditional probabilities (i.e., correct in the limit)



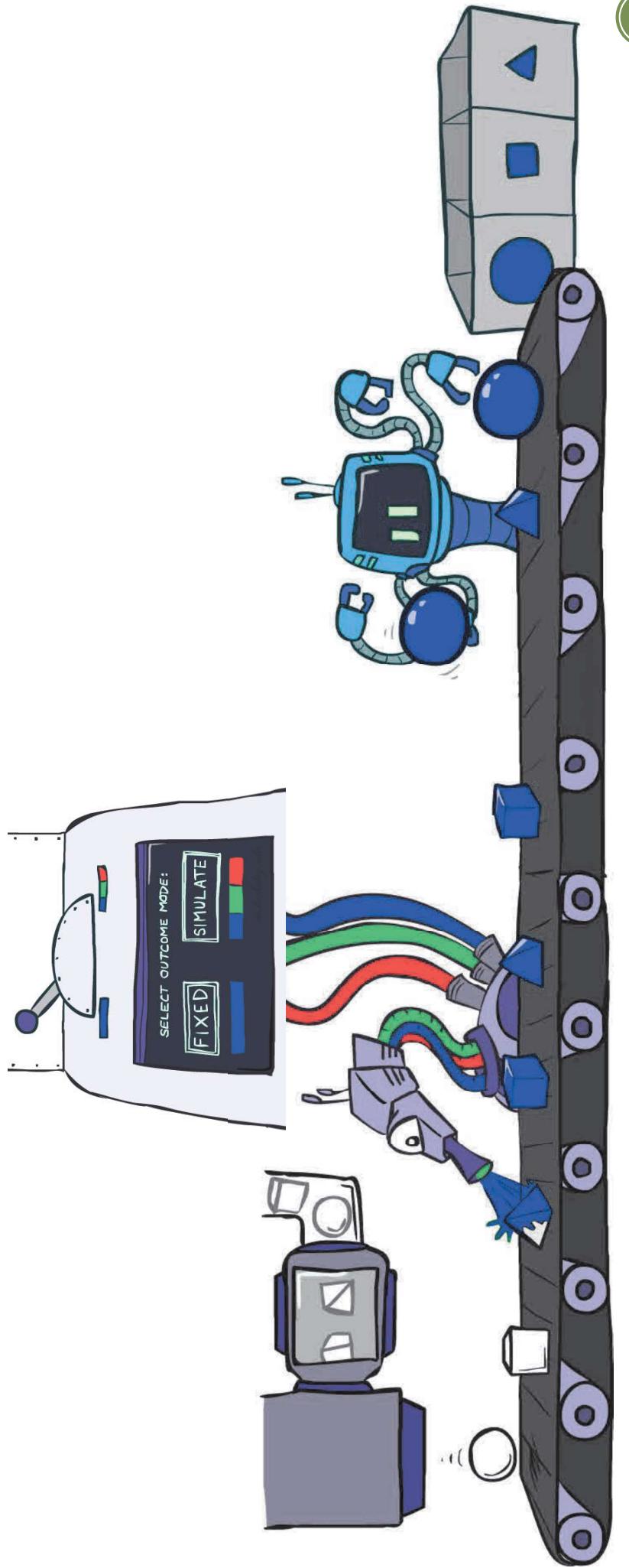
+C, -S, +R, +W
+C, +S, +R, +W
-C, +S, +R, -W
+C, -S, +R, +W
-C, -S, -R, +W

Rejection Sampling

- Input: evidence instantiation
- For $i = 1, 2, \dots, n$
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
 - If x_i not consistent with evidence
 - Reject: return – no sample is generated in this cycle
 - Return (x_1, x_2, \dots, x_n)

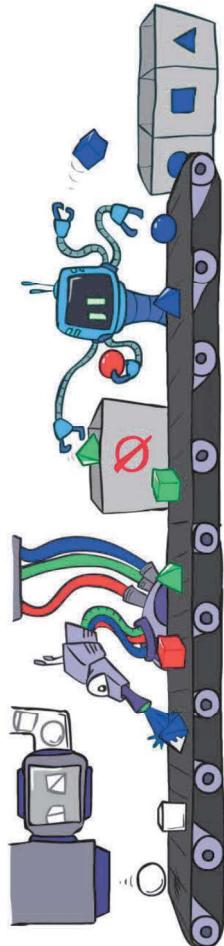
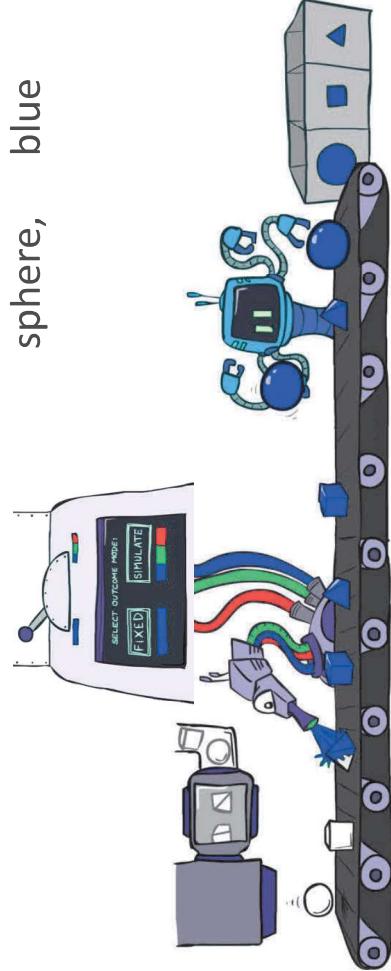


Likelihood Weighting



Likelihood Weighting

- Problem with rejection sampling:
 - If evidence is unlikely, rejects lots of samples
 - Evidence not exploited as you sample
 - Consider $P(\text{Shape} \mid \text{blue})$
 - Idea: fix evidence variables and sample the rest!
 - Problem: sample distribution not consistent!
 - Solution: weight by probability of evidence given parents
- Color
- Shape
- pyramid, green
- pyramid, red
- sphere, blue
- cube, red
- sphere, green
- pyramid, blue
- pyramid, blue
- sphere, blue
- cube, blue
- sphere, blue
- Color
- Shape
- pyramid, green
- pyramid, red
- sphere, blue
- cube, red
- sphere, green



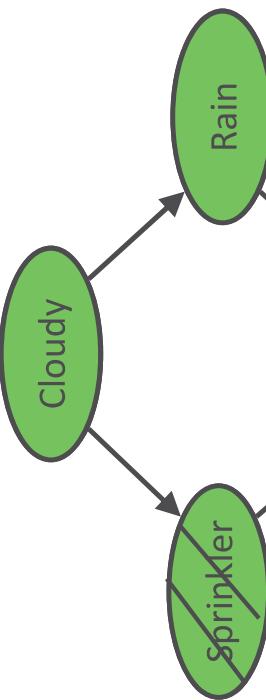
Likelihood Weighting

$$P(S|C)$$

+c	+s	0.1
-c	-s	0.9
+c	+s	0.5
-c	-s	0.5

$$P(R|C)$$

+c	+r	0.8
-c	-r	0.2
-c	+r	0.2
+c	-r	0.8



$$P(W|S, R)$$

+s	+r	+w	0.99
-s	-r	-w	0.01
+s	-r	+w	0.90
-s	+r	-w	0.10
+s	-r	-w	0.90
-s	+r	+w	0.10
+s	-r	+w	0.01
-s	+r	-w	0.99

Samples:

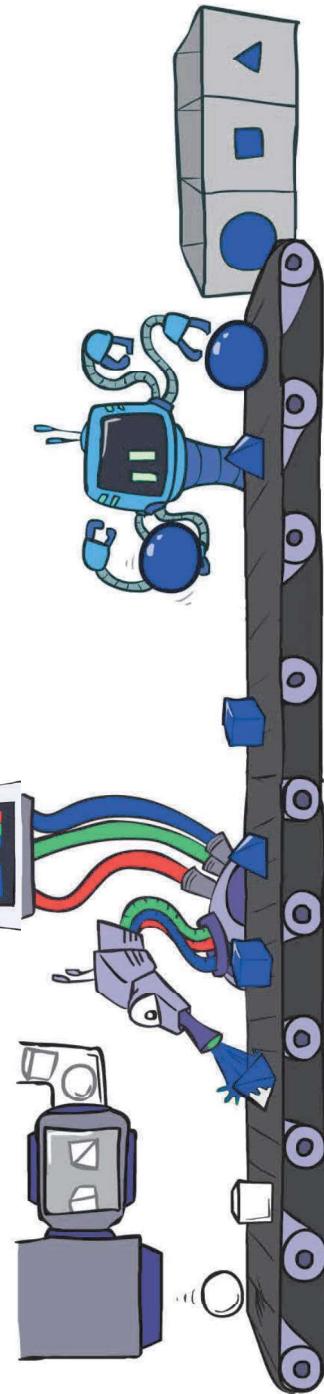
+c, +s, +r, +w

...

$$w = 1.0 \times 0.1 \times 0.99$$

Likelihood Weighting

- Input: evidence instantiation
- $w = 1.0$
- for $i = 1, 2, \dots, n$
 - if X_i is an evidence variable
 - $X_i = \text{observation } x_i \text{ for } X_i$
 - Set $w = w * P(x_i \mid \text{Parents}(X_i))$
 - else
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
- return $(x_1, x_2, \dots, x_n), w$



Likelihood Weighting

- Sampling distribution if \mathbf{z} sampled and \mathbf{e} fixed evidence

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^l P(z_i | \text{Parents}(Z_i))$$

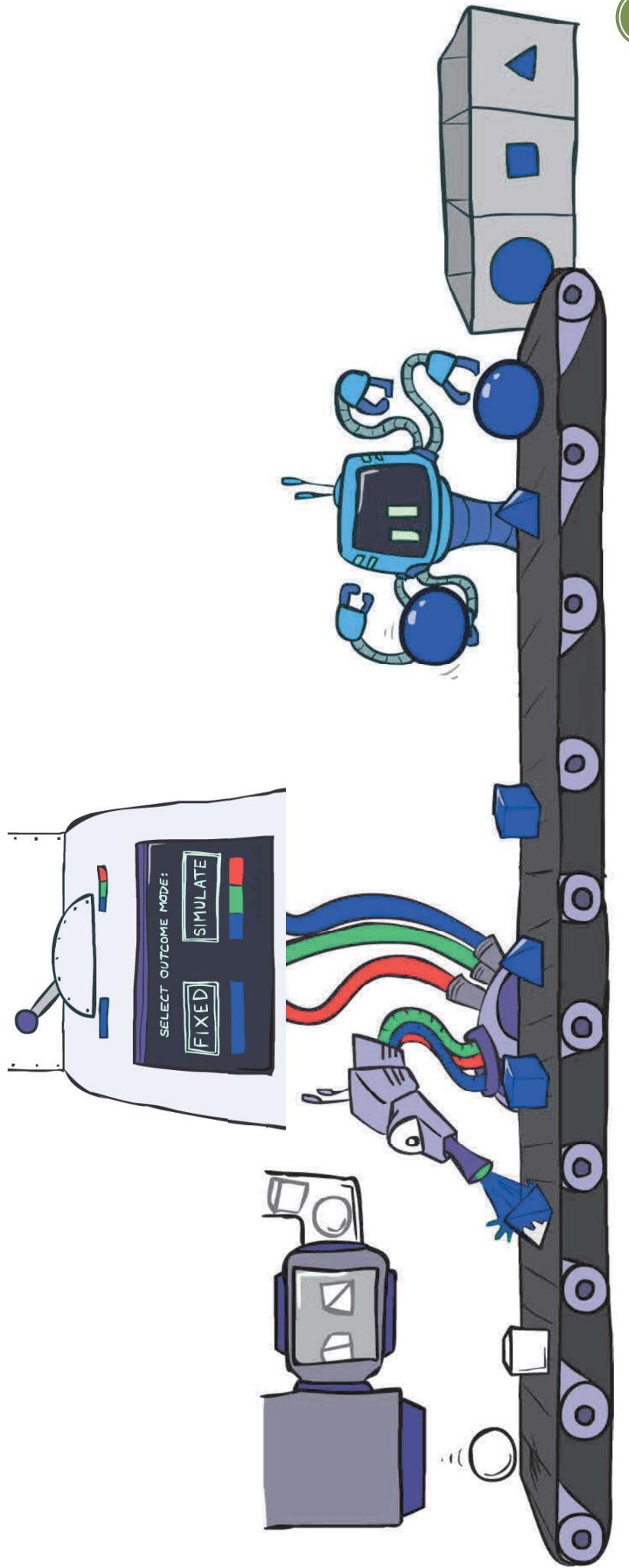
- Now, samples have weights

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^m P(e_i | \text{Parents}(E_i))$$

- Together, weighted sampling distribution is consistent

$$\begin{aligned} S_{WS}(z, e) \cdot w(z, e) &= \prod_{i=1}^l P(z_i | \text{Parents}(z_i)) \prod_{i=1}^m P(e_i | \text{Parents}(e_i)) \\ &= P(\mathbf{z}, \mathbf{e}) \end{aligned}$$

Likelihood Weighting



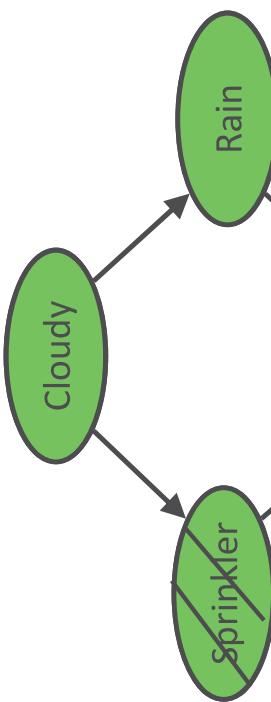
Likelihood Weighting

$$P(S|C)$$

+c	+s	0.1
-c	-s	0.9
+c	+s	0.5
-c	-s	0.5

$$P(R|C)$$

+c	+r	0.8
-c	-r	0.2
-c	+r	0.2
+c	-r	0.8



$$P(W|S, R)$$

+s	+r	+w	0.99
-s	-r	-w	0.01
+s	-r	+w	0.90
-s	+r	-w	0.10
+s	-r	-w	0.90
-s	+r	+w	0.10
+s	-r	+w	0.01
-s	+r	-w	0.99

Samples:

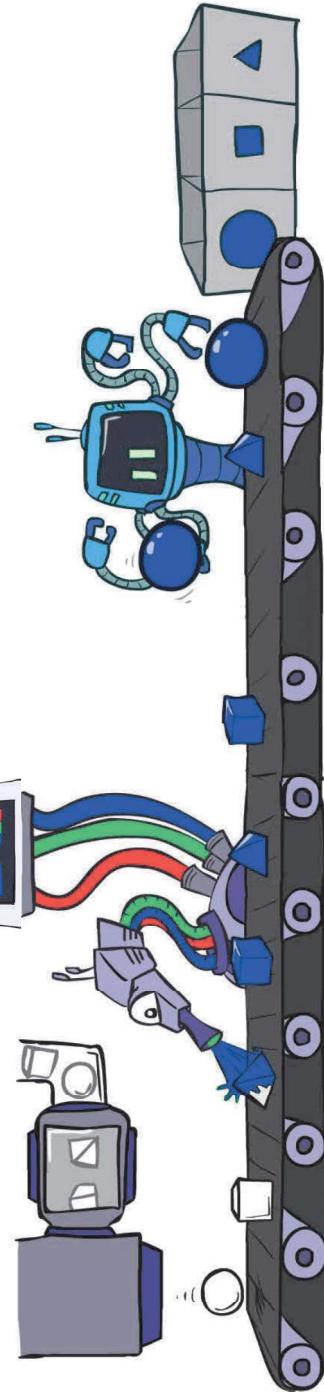
+c, +s, +r, +w

...

$$w = 1.0 \times 0.1 \times 0.99$$

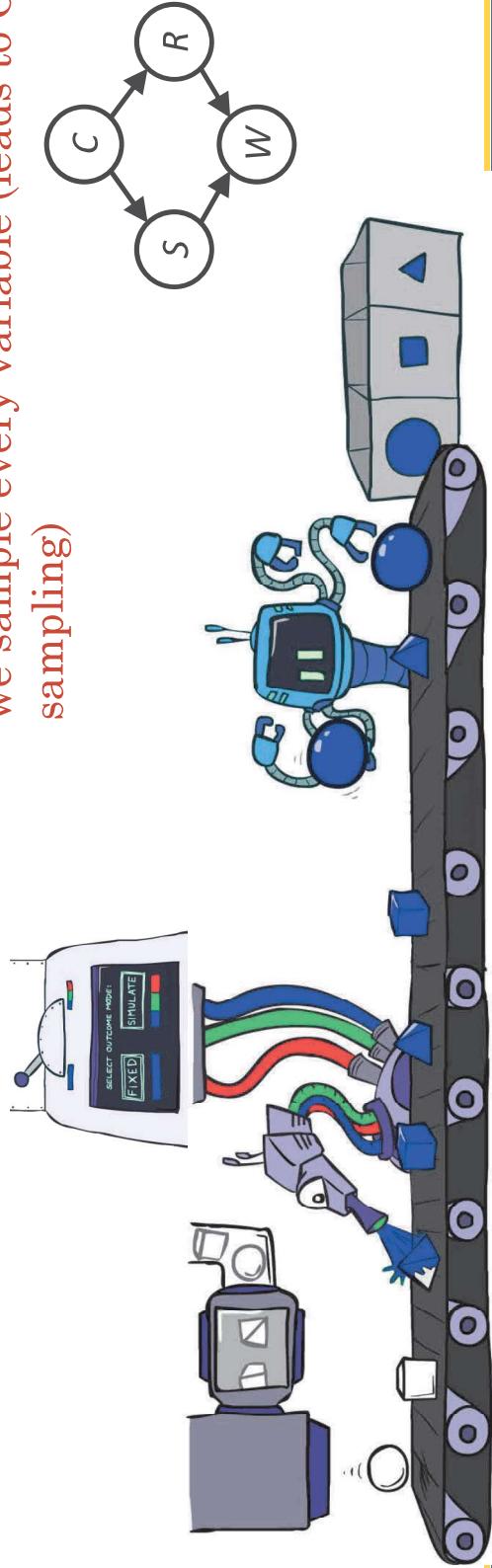
Likelihood Weighting

- Input: evidence instantiation
- $w = 1.0$
- for $i = 1, 2, \dots, n$
 - if X_i is an evidence variable
 - $X_i = \text{observation } x_i \text{ for } X_i$
 - Set $w = w * P(x_i \mid \text{Parents}(X_i))$
 - else
 - Sample x_i from $P(X_i \mid \text{Parents}(X_i))$
- return $(x_1, x_2, \dots, x_n), w$

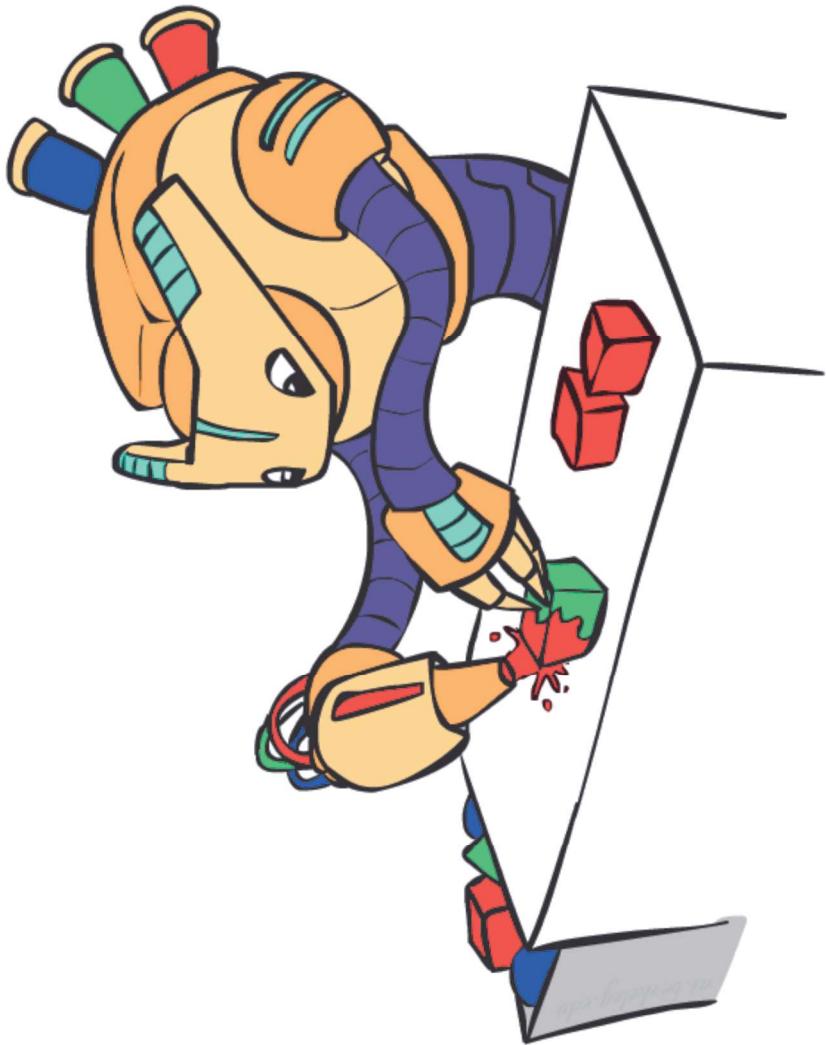


Likelihood Weighting

- Likelihood weighting is good
 - We have taken evidence into account as we generate the sample
 - E.g. here, W's value will get picked based on the evidence values of S, R
 - More of our samples will reflect the state of the world suggested by the evidence
- Likelihood weighting doesn't solve all our problems
 - Evidence influences the choice of downstream variables, but not upstream ones (C isn't more likely to get a value matching the evidence)
 - We would like to consider evidence when we sample every variable (leads to Gibbs sampling)



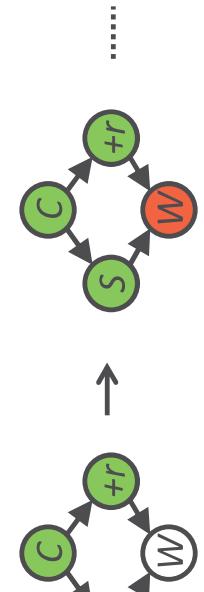
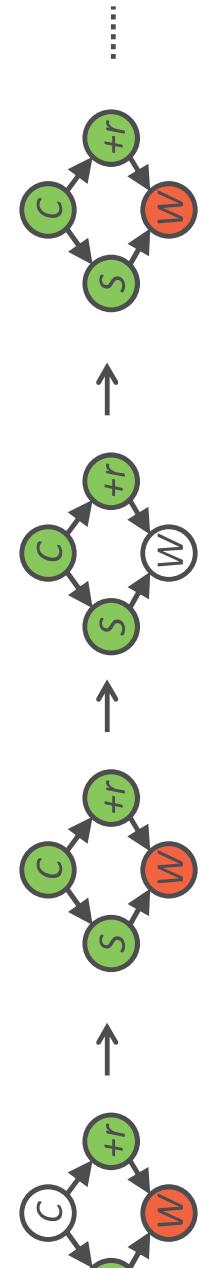
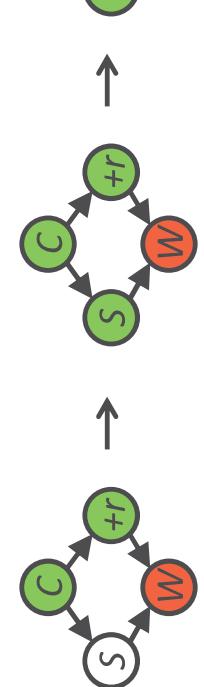
Gibbs Sampling



Gibbs Sampling

- *Procedure:* keep track of a full instantiation x_1, x_2, \dots, x_n . Start with an arbitrary instantiation consistent with the evidence. Sample one variable at a time, conditioned on all the rest, but keep evidence fixed. Keep repeating this for a long time.
- *Property:* in the limit of repeating this infinitely many times the resulting samples come from the correct distribution (i.e. conditioned on evidence).
- *Rationale:* both upstream and downstream variables condition on evidence.
- In contrast: likelihood weighting only conditions on upstream evidence, and hence weights obtained in likelihood weighting can sometimes be very small. Sum of weights over all samples is indicative of how many “effective” samples were obtained, so we want high weight.

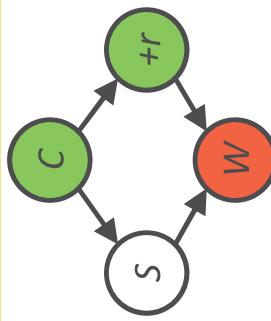
Gibbs Sampling Example: $P(S \mid +r)$

- Step 1: Fix evidence
 - $R = +r$
 - Step 2: Initialize other variables
 - Randomly
 - Steps 3: Repeat
 - Choose a non-evidence variable X
 - Resample X from $P(X \mid \text{all other variables})$
- Sample from $P(S \mid +c, -w, +r)$ Sample from $P(C \mid +s, -w, +r)$ Sample from $P(W \mid +s, -w, +r)$ Sample from $P(W \mid +s, +c, +r)$
- 
- 
- 
- 

Efficient Resampling of One Variable

- Sample from $P(S \mid +c, +r, -w)$

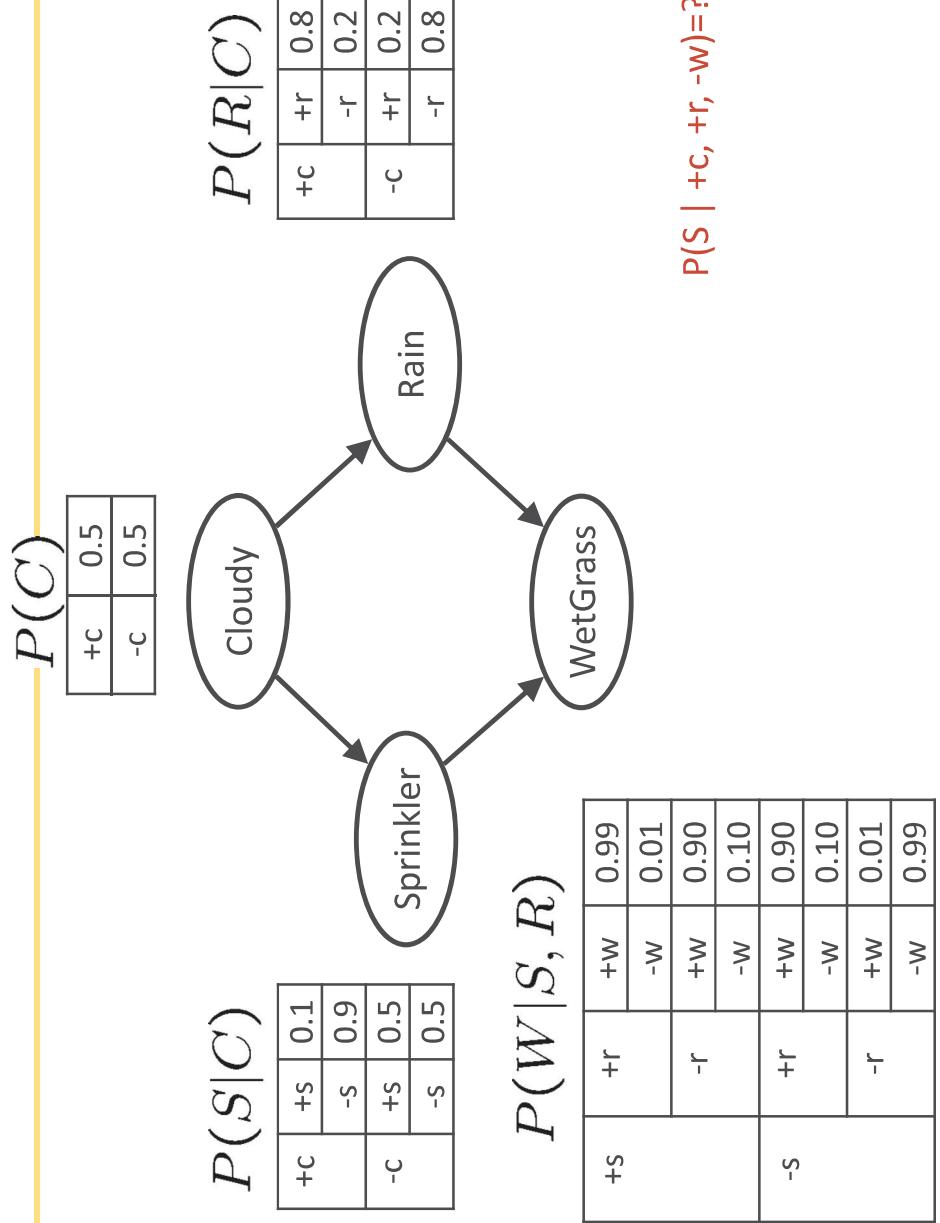
$$\begin{aligned} P(S \mid +c, +r, -w) &= \frac{P(S, +c, +r, -w)}{P(+c, +r, -w)} \\ &= \frac{P(S, +c, +r, -w)}{\sum_s P(s, +c, +r, -w)} \\ &= \frac{P(+c)P(S \mid +c)P(+r \mid +c)P(-w \mid S, +r)}{\sum_s P(+c)P(s \mid +c)P(+r \mid +c)P(-w \mid s, +r)} \\ &= \frac{P(+c)P(S \mid +c)P(+r \mid +c)P(-w \mid S, +r)}{P(+c)P(+r \mid +c) \sum_s P(s \mid +c)P(-w \mid s, +r)} \\ &= \frac{P(S \mid +c)P(-w \mid S, +r)}{\sum_s P(s \mid +c)P(-w \mid s, +r)} \end{aligned}$$



- Many things cancel out – only CPTs with S remain!

- More generally: only CPTs that have resampled variable need to be considered, and joined together

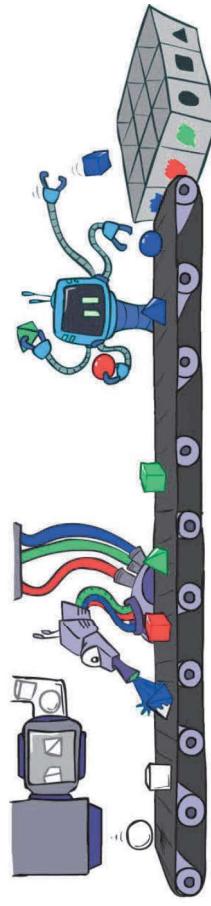
Example



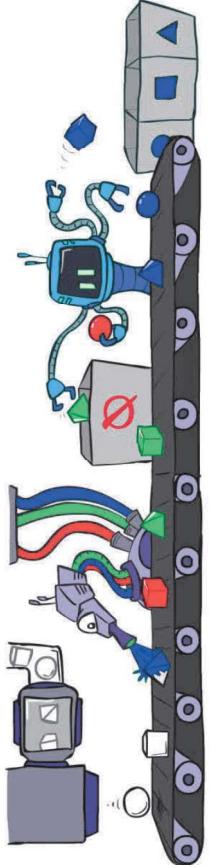
$$P(S \mid +c, +r, -w) = ?$$

Bayes' Net Sampling Summary

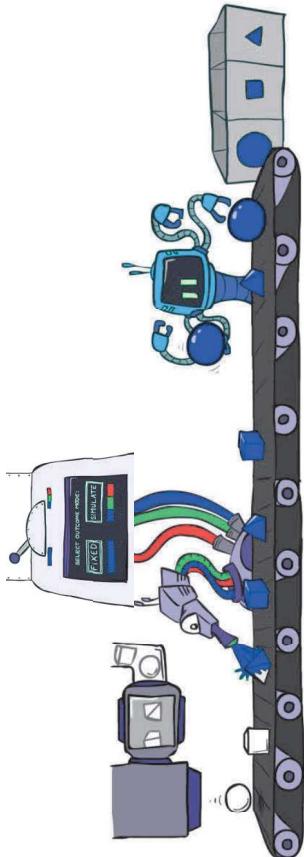
- Prior Sampling $P(Q)$



- Rejection Sampling $P(Q | e)$



- Likelihood Weighting $P(Q | e)$



- Gibbs Sampling $P(Q | e)$



Further Reading on Gibbs Sampling*

- Gibbs sampling produces sample from the query distribution $P(Q | e)$ in limit of re-sampling infinitely often
- Gibbs sampling is a special case of more general methods called Markov chain Monte Carlo (MCMC) methods
- Metropolis-Hastings is one of the more famous MCMC methods (in fact, Gibbs sampling is a special case of Metropolis-Hastings)
- You may read about Monte Carlo methods – they're just sampling