

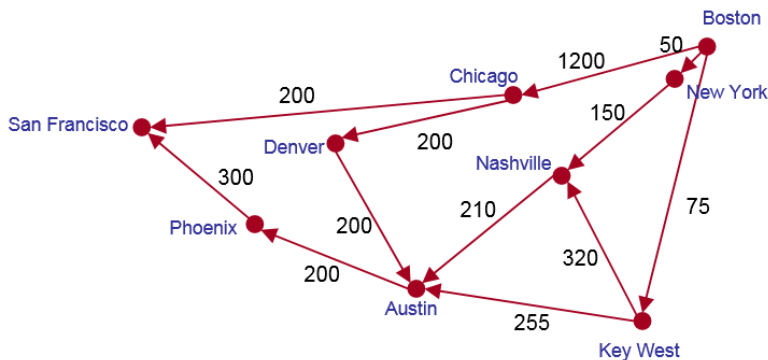
Planning Recap & Theory of Mind Papers

Prof. Brad Hayes

University of Colorado Boulder

Planning Recap

Search as a Problem-Solving Technique



Depth

0

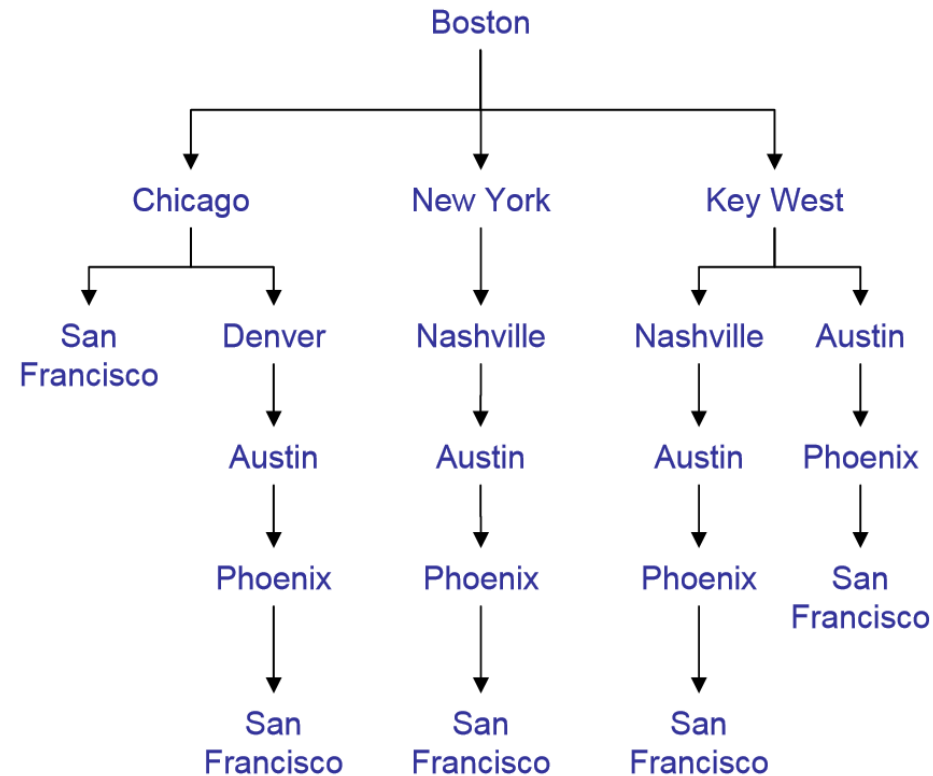
1

2

3

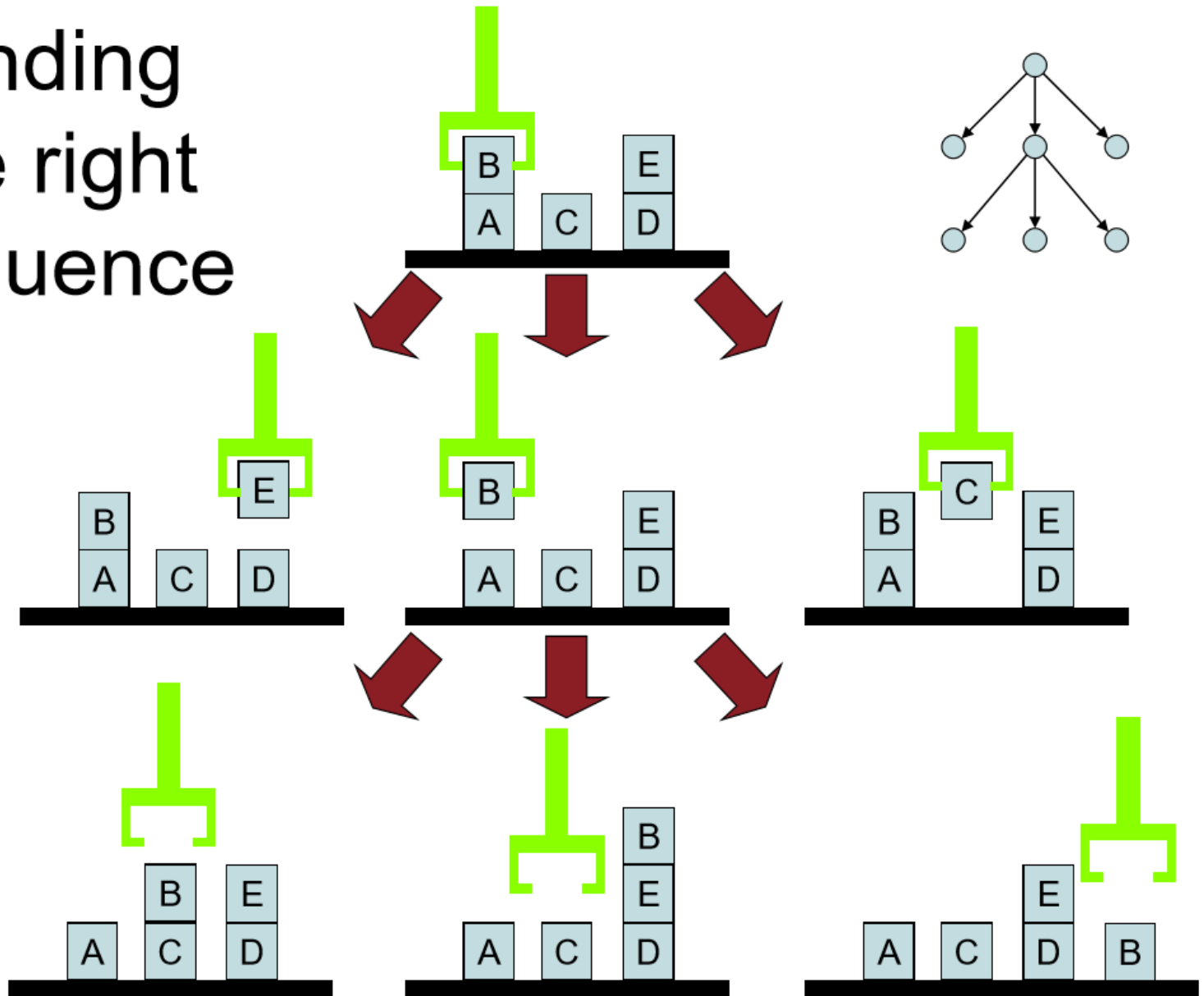
4

5



Branching Factor $b=3$

Finding the right Sequence



Lifted vs. Grounded Planning

Plan in this space

Execute in this space

Predicate Representation

On(a,b)
Moving(d)
Clear(Lhand)

...
...
...
...
...

1.0
1.0
0
0
0
-1.0
-1.0
1.0

Task
Planning

Mapping Function

Compression
Function

World State Vector

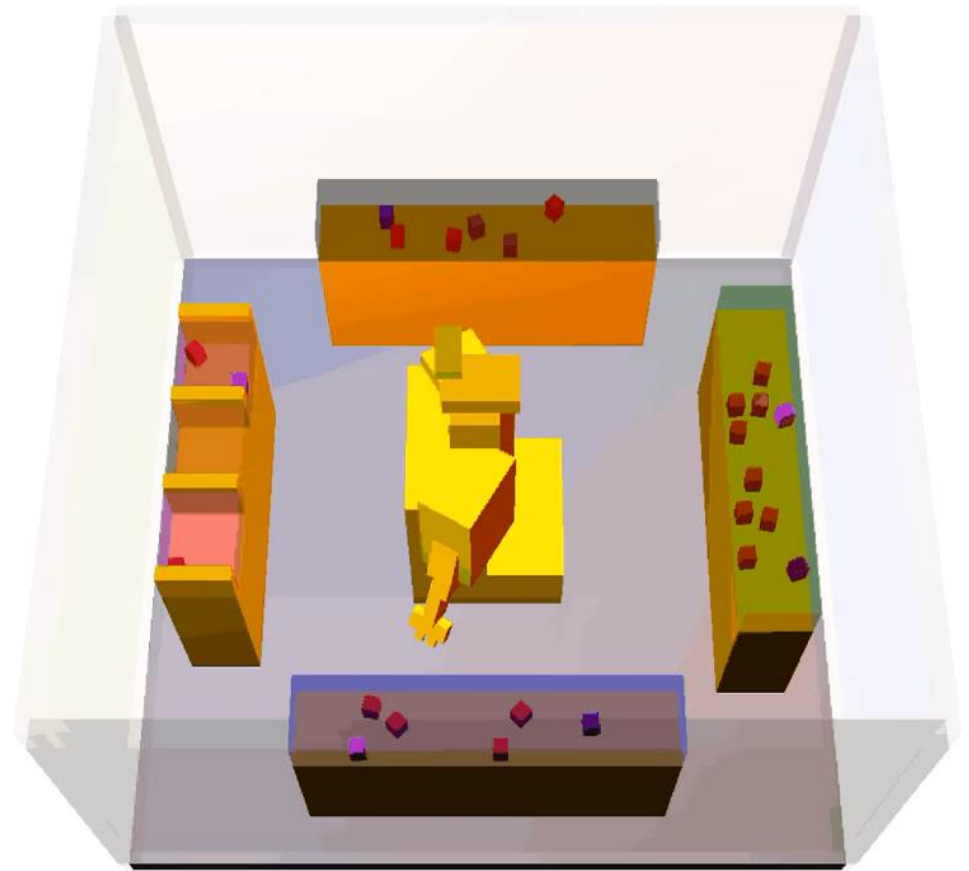
0.124
543.1
491.3
-50.11
194.2

...
...
...
...

Motion
Planning

Strategies for Tough Planning Problems

- Breadth-first search with sampled action parameterizations
- Pre-determined action parameterizations
 - e.g., `Place(object,x,y,z)`
- Plan for abstract 'areas'
 - **Lazy Grounding:**
Sample parameterizations and add that action to plan if a valid parameterization is found
 - **Greedy Grounding:**
Sample parameterizations and add that specific parameterization to the plan once a valid one is found



Planning Heuristics

FastForward Heuristic

Remove deletions from all operators and run planner

Resulting plan gives a heuristic for # steps from current state to goal

FF Plan can provide ordering constraints on actions, even within partially specified plans

Learned Cost Function

Learning Real-Time A* (LRTA*)

Update heuristic dynamically: Value function $V(s)$

Connection back to reinforcement learning: **$V(s)$ and $Q(s,a)$** for search heuristic

Admissible Heuristics and Pruning Strategies

Admissible heuristics from delete-relaxation:

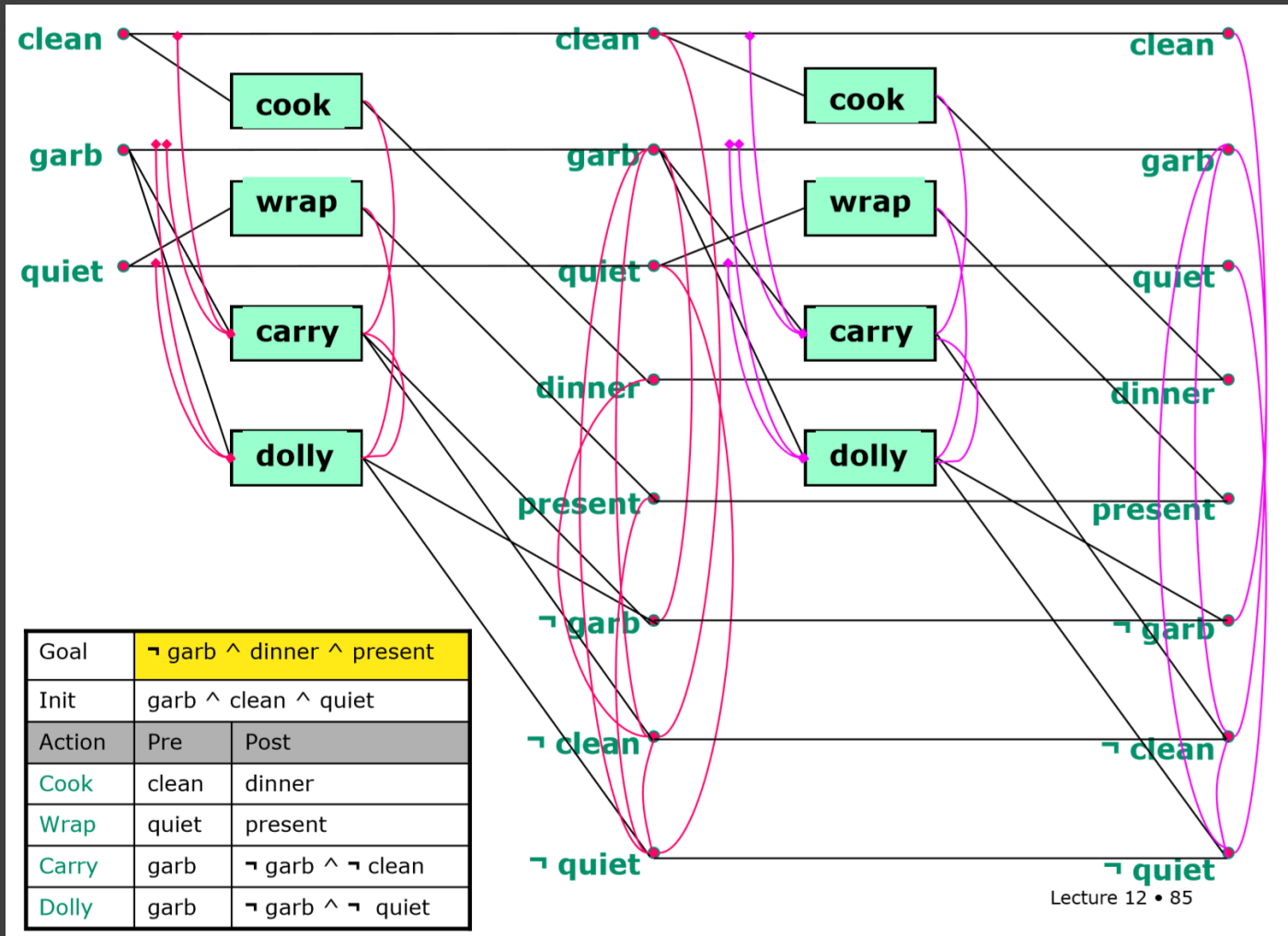
- $h_{max} = i$ iff goal g is reachable in i steps from s
 - Admissible because of superposition (shown in GRAPHPLAN Algorithm)
- $h_{FF} =$ number of **different** actions in $\pi(g)$

Search Pruning Strategies:

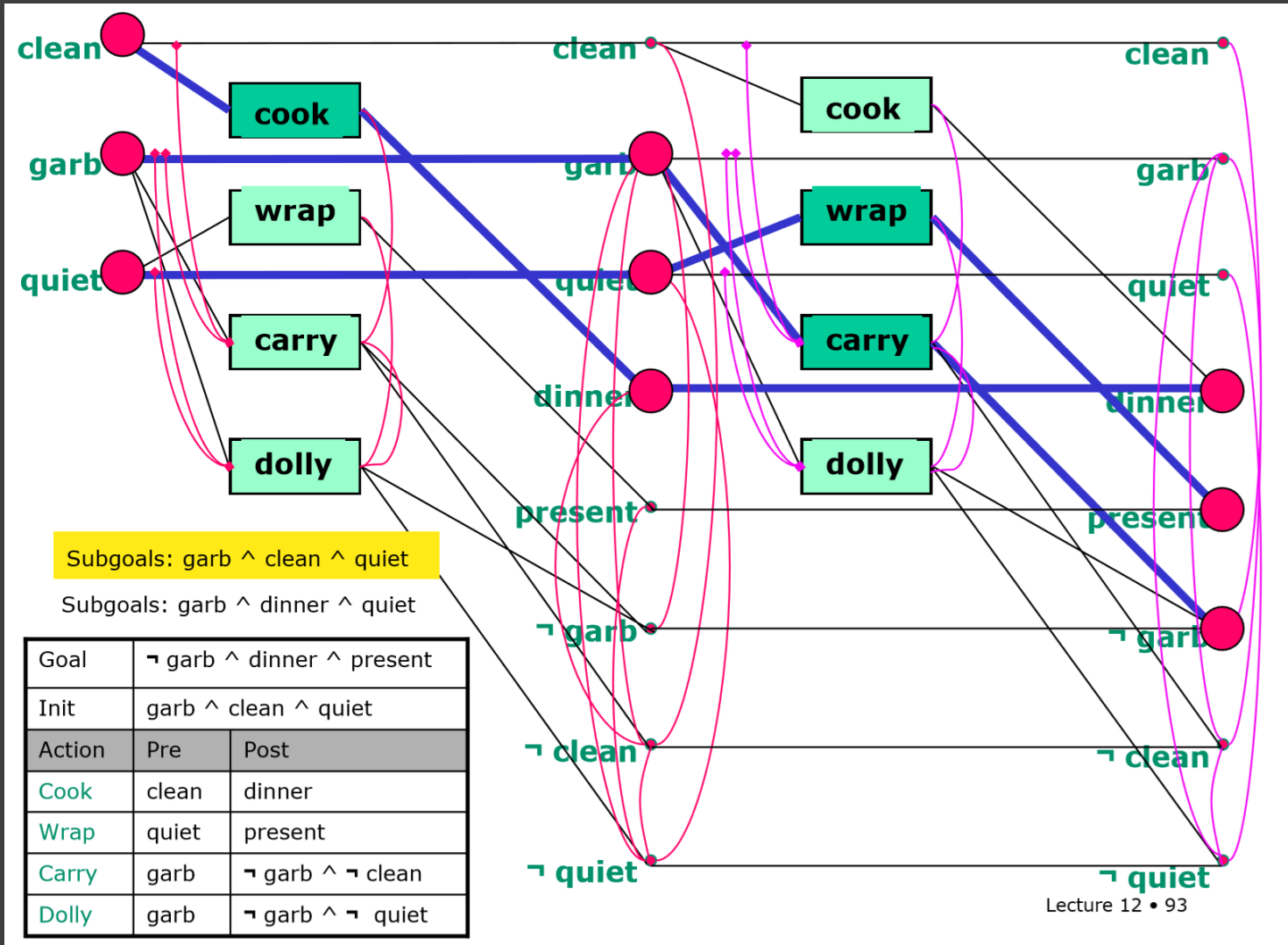
Novelty: Assign each state $novelty(s) = i$ if i predicates take on a value never seen in path from s_0 to s

IW-search (Iterative Width): Ignore any nodes at the search frontier with $novelty(s) > i$.

GRAPHPLAN



GRAPHPLAN Solution



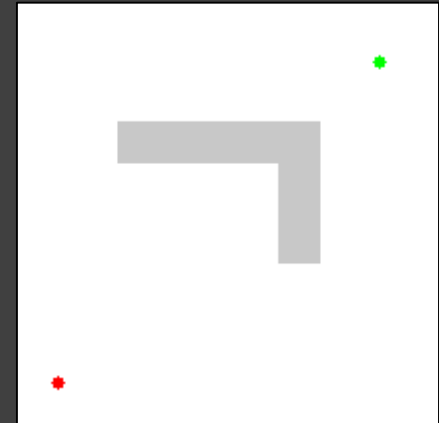
A* Search Algorithm

Adds heuristic to Dijkstra's Algorithm to bias exploration toward optimal path.

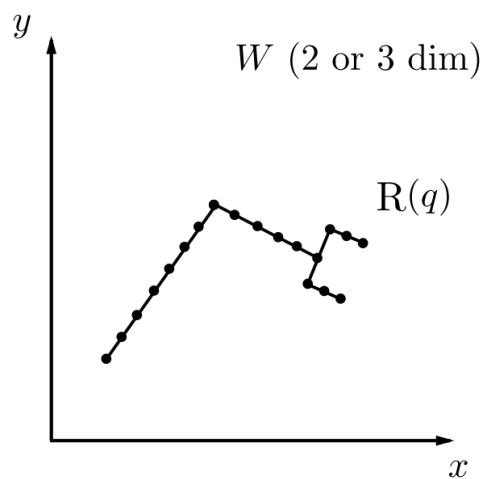
Insight

$$f(n) = g(n) + h(n)$$

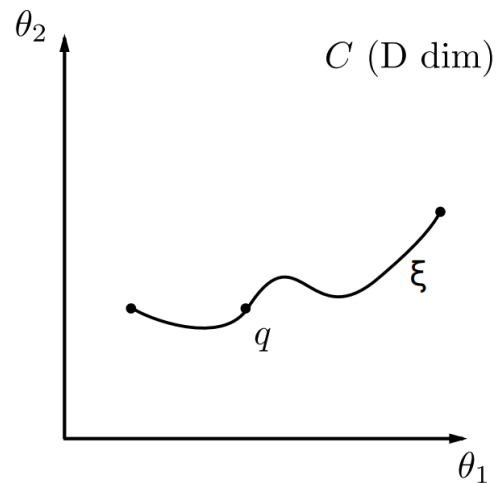
$$f(n) = \text{Cost so far} + \text{Est. cost to go}$$



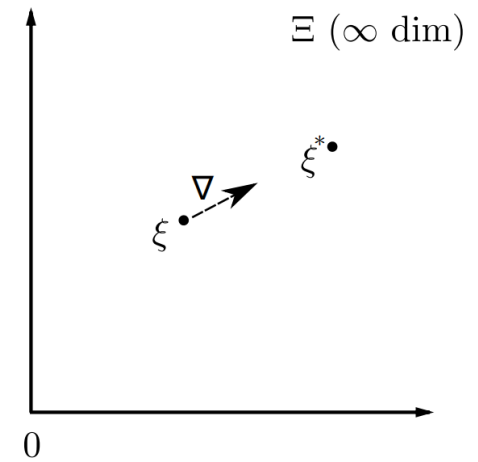
Admissible Heuristic: An underestimate of the shortest possible path from node n to the goal.



World Space $W(\mathbb{R}^3)$



Configuration Space
 $C(\#DoF)$



Trajectory Space
 $\Xi(\infty \text{ dim})$

Robot pose in World
Space (set of points)



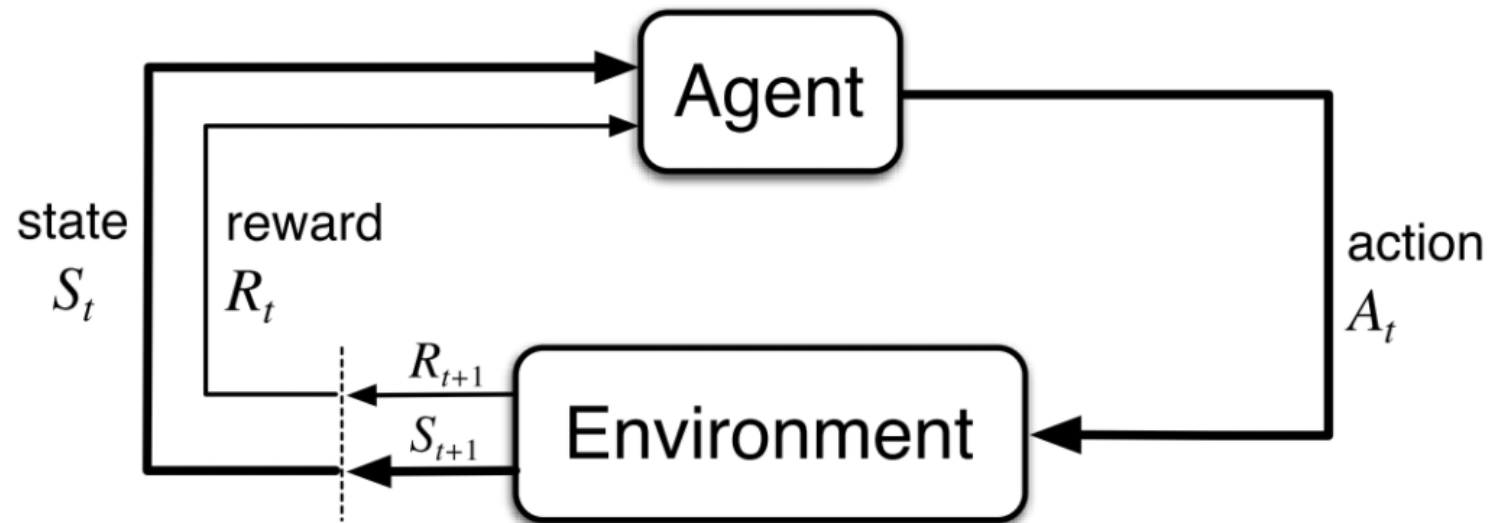
Single point in
Configuration Space

Trajectory through
Configuration Space
(set of points)



Single point in
Trajectory Space

Rewind to January...



Rewind to January...

Reinforcement Learning



\mathbf{o}

Function: $O \rightarrow A$

$\pi_{\theta}(\mathbf{a}|\mathbf{o})$



\mathbf{a}

\mathbf{s}_t – state

\mathbf{o}_t – observation

\mathbf{a}_t – action

$\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ – policy

$\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$ – policy (fully observed)



\mathbf{o}_t – observation



\mathbf{s}_t – state

Reinforcement Learning



o

Function: $O \rightarrow A$

$\pi_{\theta}(\mathbf{a}|\mathbf{o})$



a

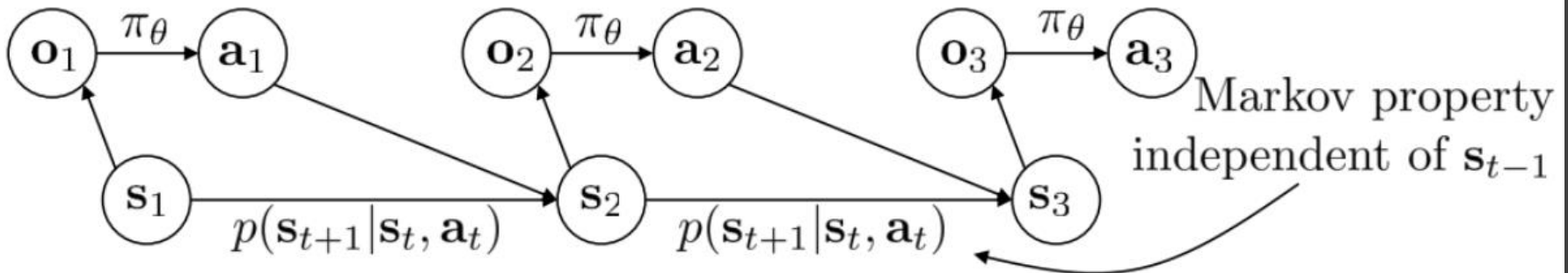
\mathbf{s}_t – state

\mathbf{o}_t – observation

\mathbf{a}_t – action

$\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ – policy

$\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$ – policy (fully observed)



Learning an Action-Value Function (Q-Learning)

- Q-value is the expected utility of taking an action in a state:

$$Q(s,a) = \text{Value of action 'a' in state 's'}$$

- For a state s , choose action 'a' that maximizes $Q(s,a)$
- Q-values do not require an **environment model**:
The transition function isn't necessary, since it will be learned from experience.

$$Q(s, a) = (1 - \alpha) * Q(s, a) + \alpha(r + \gamma * \max_{a'} Q(s', a'))$$

Reliance on
existing vs. new
Knowledge

Current
understanding

reward

Discount
factor

Best we think we can do
in the future from here

Learning Real-Time A*

- Hill-climbing algorithm where the **best child node** is selected and others **discarded**, and the heuristic is updated dynamically.

A* Cost Formula: Dictates order of node expansion

- $f(s) = g(s) + h(s)$
- $g(s)$ remains as cost already paid to get here: $\sum_{i=0}^n c(s_i, a_i)$
- $g(s_{i+1}) = g(s_i) + c(s_i, a_i)$
- Replace $h(s_i)$ with $\operatorname{argmax}_a Q(s_i, a)$

$$f(s) = g(s) + \operatorname{argmax}_a Q(s_{i+1}, a)$$

$$Q(s, a) = (1 - \alpha) * Q(s, a) + \alpha (r + \gamma * \max_{a'} Q(s', a'))$$

Today

Today's Papers

Planning + Predicting Human Behavior

Probabilistically Safe Robot Planning with Confidence-Based Human Predictions

Jaime Fisac et al.

An Implemented Theory of Mind to Improve Human-Robot Shared Plans Execution

Sandra Devin and Rachid Alami

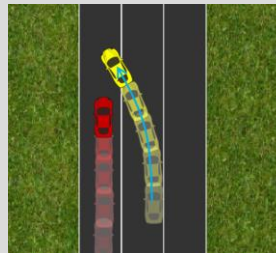
Papers for Next Thursday (2/21)

Modeling Capability and Effect

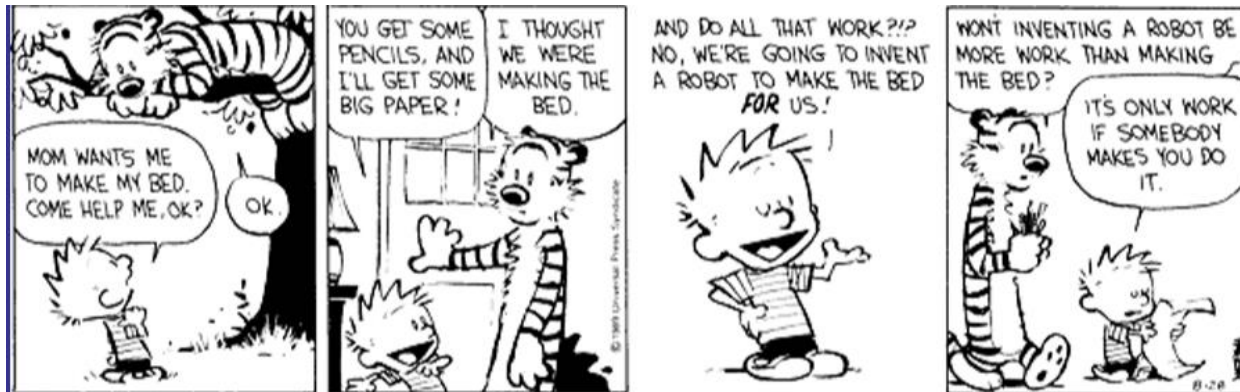
Game-Theoretic Modeling of Human Adaptation in Human-Robot Collaboration
Stefanos Nikolaidis et al.



Planning for Autonomous Cars that Leverage Effects on Human Actions
Dorsa Sadigh et al.



Final Project Outlines



Due Friday (3/1), at 5:00pm via Moodle

- Major deliverables (+ assigned lead)
- Describe technical components (e.g., “face detector”)
- Describe evaluation
- Schedule for design, implementation, evaluation, writing
- Concerns / risks / challenges

Policy Differentiation

$$\theta^* = \arg \max_{\theta} \sum_{t=1}^T \underbrace{E_{(s_t, a_t) \sim p_{\theta}(s_t, a_t)} [r(s_t, a_t)]}_{J(\theta)} \quad \pi_{\theta}(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t | p(s_{t+1} | s_t, a_t))$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [r(\tau)] = \int \pi_{\theta}(\tau) r(\tau) d\tau \quad \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) = \pi_{\theta}(\tau) \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \nabla_{\theta} \pi_{\theta}(\tau)$$

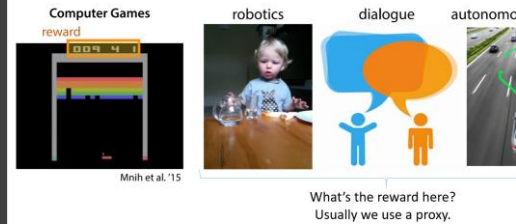
$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} \pi_{\theta}(\tau) r(\tau) d\tau = \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

MaxEnt IRL

1. Initialize ψ (reward function params), gather demonstrations D
2. Solve for optimal policy $\pi(a|s)$ with reward r_{ψ}
3. Solve for state visitation frequencies $p(s|\psi)$
4. Compute gradient: $\nabla_{\psi} L = -\frac{1}{|D|} \sum_{\tau_d \in D} \frac{dr_{\psi}}{d\psi}(\tau_d) - \sum_s p(s|\psi) \frac{dr_{\psi}}{d\psi}(s)$
5. Update ψ with one gradient step using $\nabla_{\psi} L$
6. GOTO 2

Must solve the whole MDP in the inner loop of finding the reward function!

Where does reward come from?



Using the Policy Gradient

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)$$

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(a_t | s_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \right) \left(\sum_t r(s_t^i, a_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Coming Soon

Project Workshops

Stochastic Gradient Descent for Policy Learning (in continuous action spaces)

Inverse Reinforcement Learning