

# What is On / Off Policy?

- Q learns how to perform optimally even when we are following a non-optimal policy
- In  $\epsilon$ -greedy,  $\epsilon$  leaves no trace in Q
- SARSA is on-policy
- Learns the best policy given our systematic departures from true optimal
- In  $\epsilon$ -greedy,  $\epsilon$  is reflected within SARSA's Q values

# On Policy vs Off Policy

- Q, an off-policy learner:

$$Q(a, s) \leftarrow Q(a, s) + \alpha \cdot (r_s + \gamma \max_{a'} Q(a', s') - Q(a, s))$$

- SARSA, an on-policy learner:

$$Q(a, s) \leftarrow Q(a, s) + \alpha \cdot (r_s + \gamma \cdot Q(a', s') - Q(a, s))$$

- How do they differ?
  - Remember the persistent need for exploration
  - Consider cliff walking

# Can Q use more Information?

- Efficient = use all information
- Indirect learners can use the approx. model
- No information until rewards
- Information is propagated one step back
- Can we do better?

## Eligibility traces

Another parameter:  $\lambda$

Between 0 and 1 (close to 1)

TD( $\lambda$ ), Q( $\lambda$ ), SARSA( $\lambda$ )

# Deficiencies with RL

(and statistical learning generally)

- Constrained expressiveness
- Generally limited to propositional representations
- Stems from formalizing with Random Variables
- What is Propositional?
- What are the alternatives?
- Complexity of change
- Structure among world states

# A Relevant Example

Robot R delivers gifts from  
location J to locations  
G1 and G2

Thief T may steal the gifts

Actions: Left, Right, Up, Down

	X	G1		X
	R			
X			T	
			X	G2
		J		

MP2 World

Propositional state representation

Effects of actions “Right:  $x \leftarrow x+1$ ”

Structure among world states

Goals / subgoals

# The Curse of Statistical Learning

- Information required for approx.  $\pi^*$
- Looser approximations require less information
  - How close to optimal?
  - 5x5 grid vs. 50x50
  - (how many places can *you* be)
- Empirical data makes up the difference
- Impoverished expressiveness
- Everything we know but do not tell, inflates the empirical data requirement
- Consider learning brain surgery by example

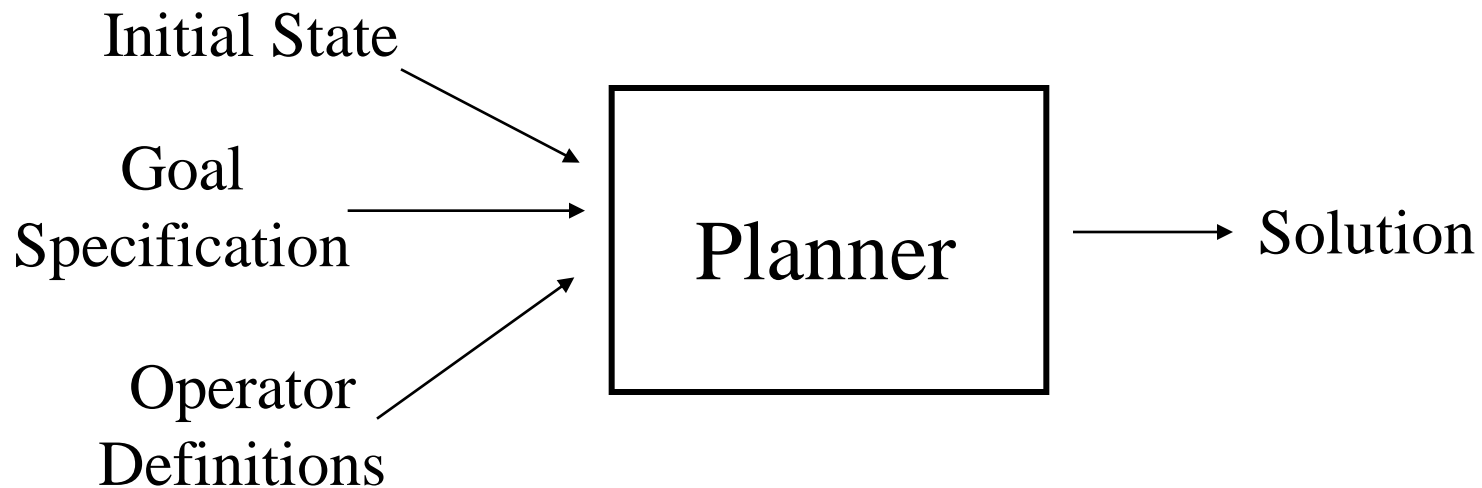
# Reinforcement Learning (vs. Classical Planning)

- Robust
  - Fewer *a priori* assumptions (esp. actions)
  - Empirical model
  - Fit (via parameter adjustment) to the *observed* world
- Scaling difficulties
  - Propositional expressiveness
  - Complexity space / time (?)
  - States / Features (e.g., block positions)
- Markov assumption
  - Our world?
  - Implications for sensors & convergence
  - Discretizing may not respect Markov

# Classical Planning

## Domain Independent Planning

Using inference to find a sequence of operator instances (actions) that transform an initial state into a state in which the goal is satisfied.



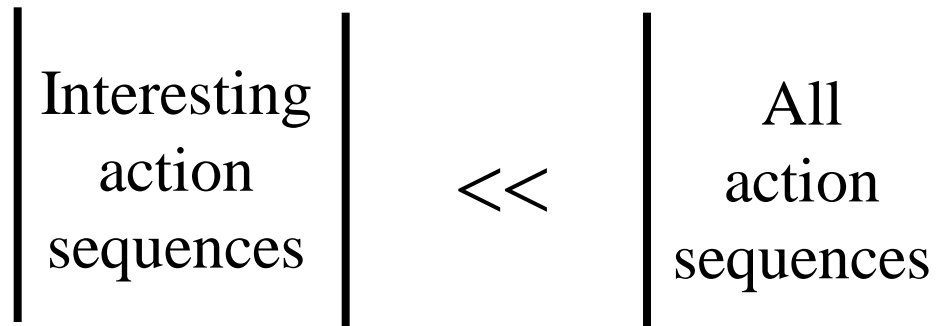
Real World Applications:

Scheduling, Semantic web support, Computer gaming, ...

(limited by assumption of an analytic model)



# Recall from Search



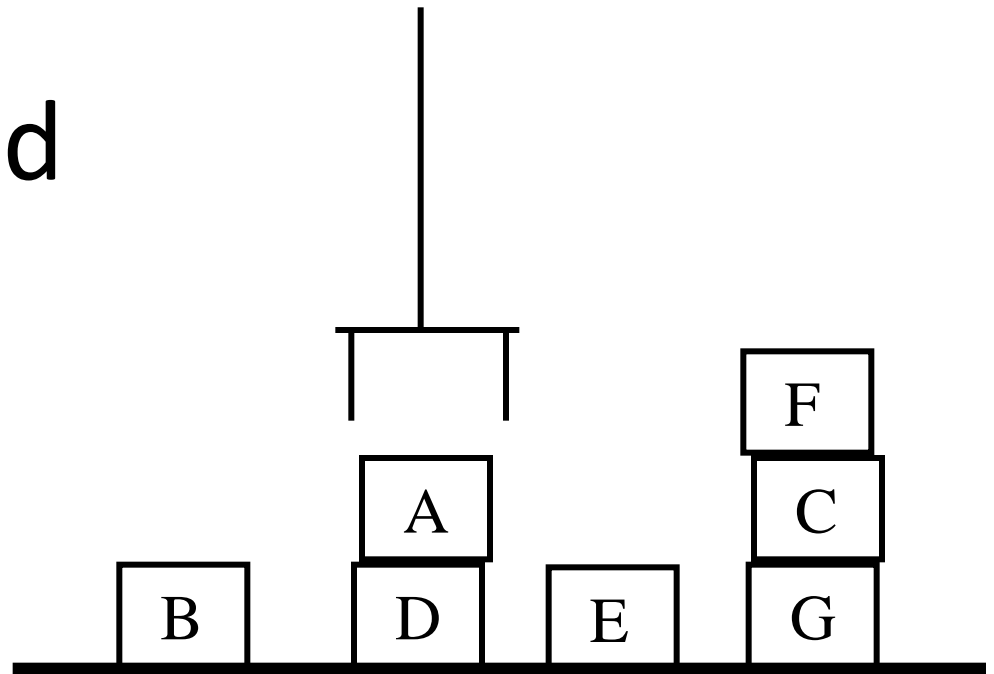
RL actions are “inferentially opaque”

World state is a collection of relevant *objects* and their *relationships*

Features impose a structure among states

Planning allows reasoning about state features

# Blocks World



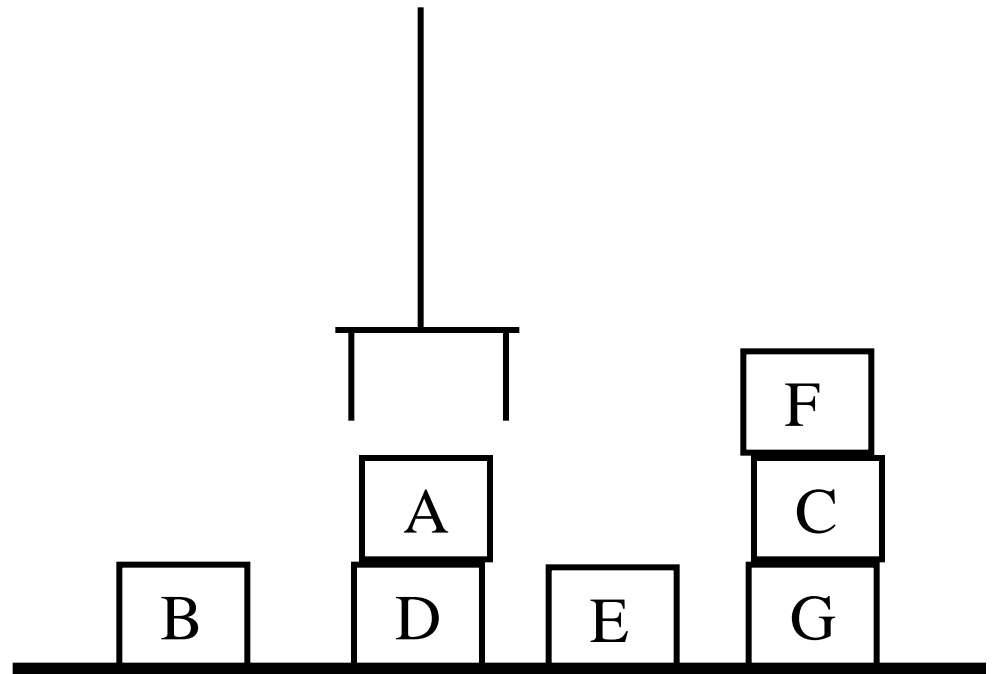
Several ontologies possible (ways to conceptualize the world and its changes)

Operator - General knowledge of one kind of change

Action - Ground instance of an operator (RL has only these)

Silly domain but concisely illustrates many GENERAL planning issues

# Traditional Blocks World



Only keep track of support relationships: On, Clr

A block can support at most one other block

The table can support any number of blocks

Generalized block movement – no gripper

# Alternative Ontologies

change a block's position differently

Move-Block

Move-Gripper

Grasp-Block

Move-Gripper

UnGrasp-Block

Move-Gripper

Open-Gripper

Move-Gripper

Close-Gripper

...

Motor1-Velocity

Motor2-Velocity

...

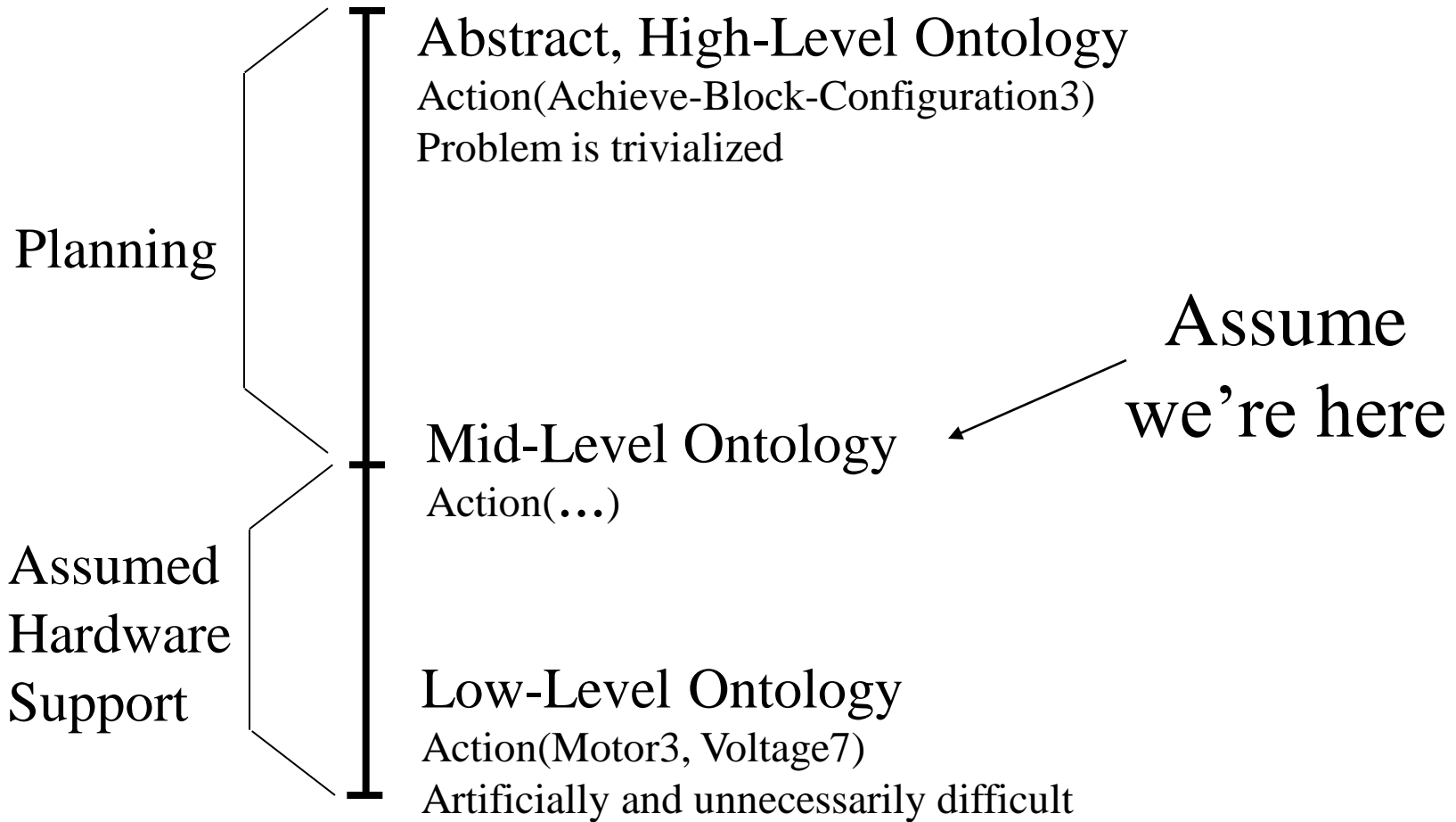
Motor1-Voltage (Current, Duty Cycle)

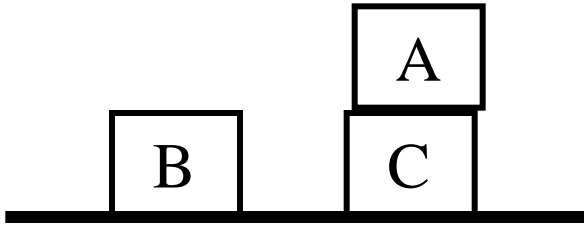
Motor2-Voltage

...

...

# Levels of Ontological Commitment





Initial State

On(A, C)

On(C, Tbl)

On(B, Tbl)

Blk(A)

Blk(B)

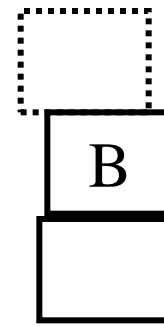
Blk(C)

Table(Tbl)

Clr(A)

Clr(B)

Clr(Tbl)



Goal

On(B, ?x)

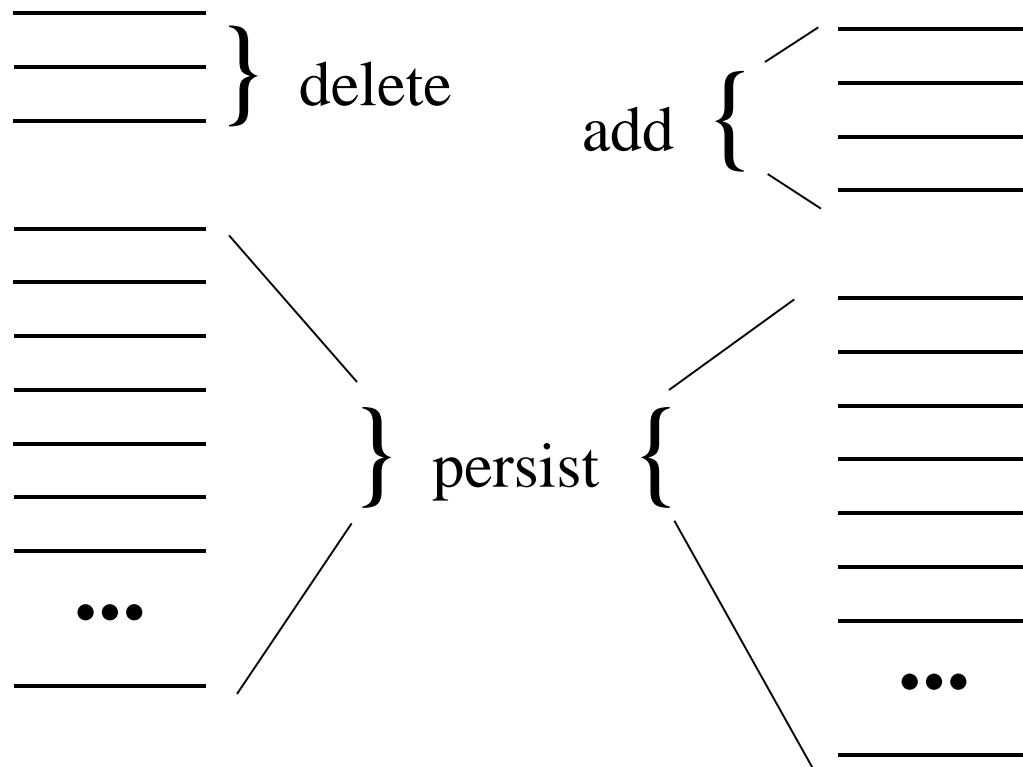
Blk(?x)

# World Changes

Operator must fully define resulting world state  
from any legal previous state

$S_i$

Result (Action,  $S_i$ )

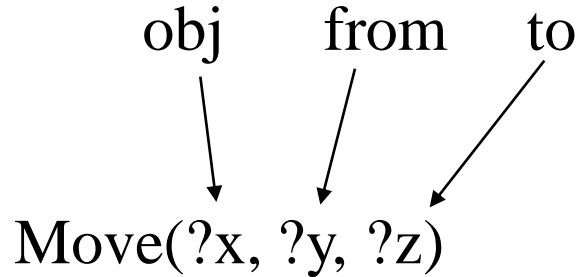


# Strips Operators

- Preconditions: list of features that must hold in the world to apply the operator
- Effects
  - Delete list: features ceasing to hold
  - Add list: features asserted by the action



# the Move operator



What holds in the resulting situation?

?x is on ?z

?y is clear

...

What must hold now?

?x is on ?y

?z is clear

?x is a block

?x is clear

...

# Two Strips Operators for Move

MoveToBlock (x, y, z):

PC: Clr (x), Clr (z), On (x, y), Blk (x),  
Blk (z), Diff (x, z), Diff (y, z)

Delete: On (x, y), Clr (z),  
Add On (x, z), Clr (y)

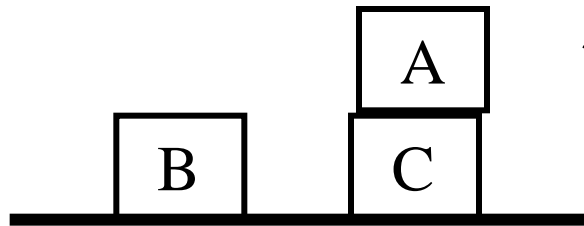
MoveToTable (x, y):

PC: Clr (x), On (x, y), Blk (x), Diff (y, Tbl)

Delete: On (x, y)  
Add: On (x, Tbl), Clr (y)

# World Change

Initial State:  $S_i$



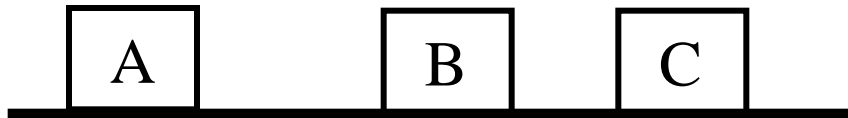
...

$\text{On}(A, C)$

...

$\text{MoveToTable}(A, C)$

State after  $\text{MoveToTable}(A, C)$  in  $S_i$



...

$\text{On}(A, \text{Tbl})$

...

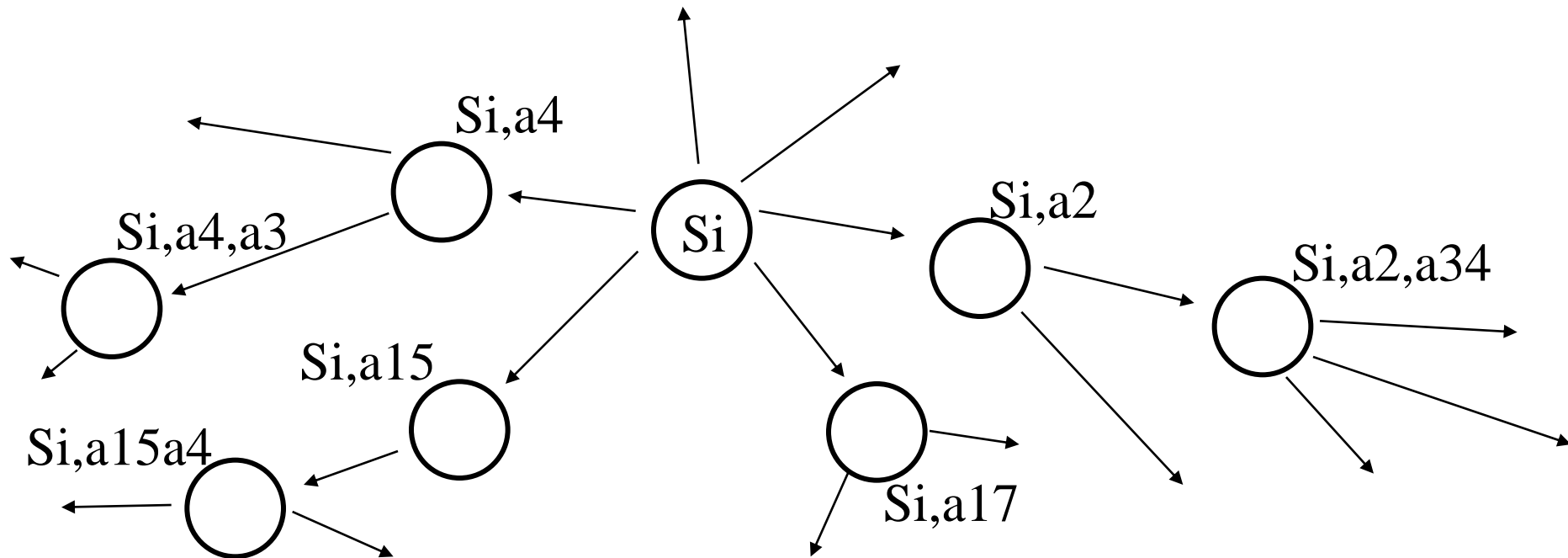
# All Reachable Situations are Defined

Given: 1) the Initial State

2) Axioms of World Change (operator definitions)

Can be realized as predicate calculus theorem proving

$\Delta \equiv \text{Initial State} \cup \text{Operator Definitions}$



# Planners

- State space vs Plan space
- Linear vs. Nonlinear
- Goal-stack planning
- POP
- Difficulties with classical planning

Learning search heuristics

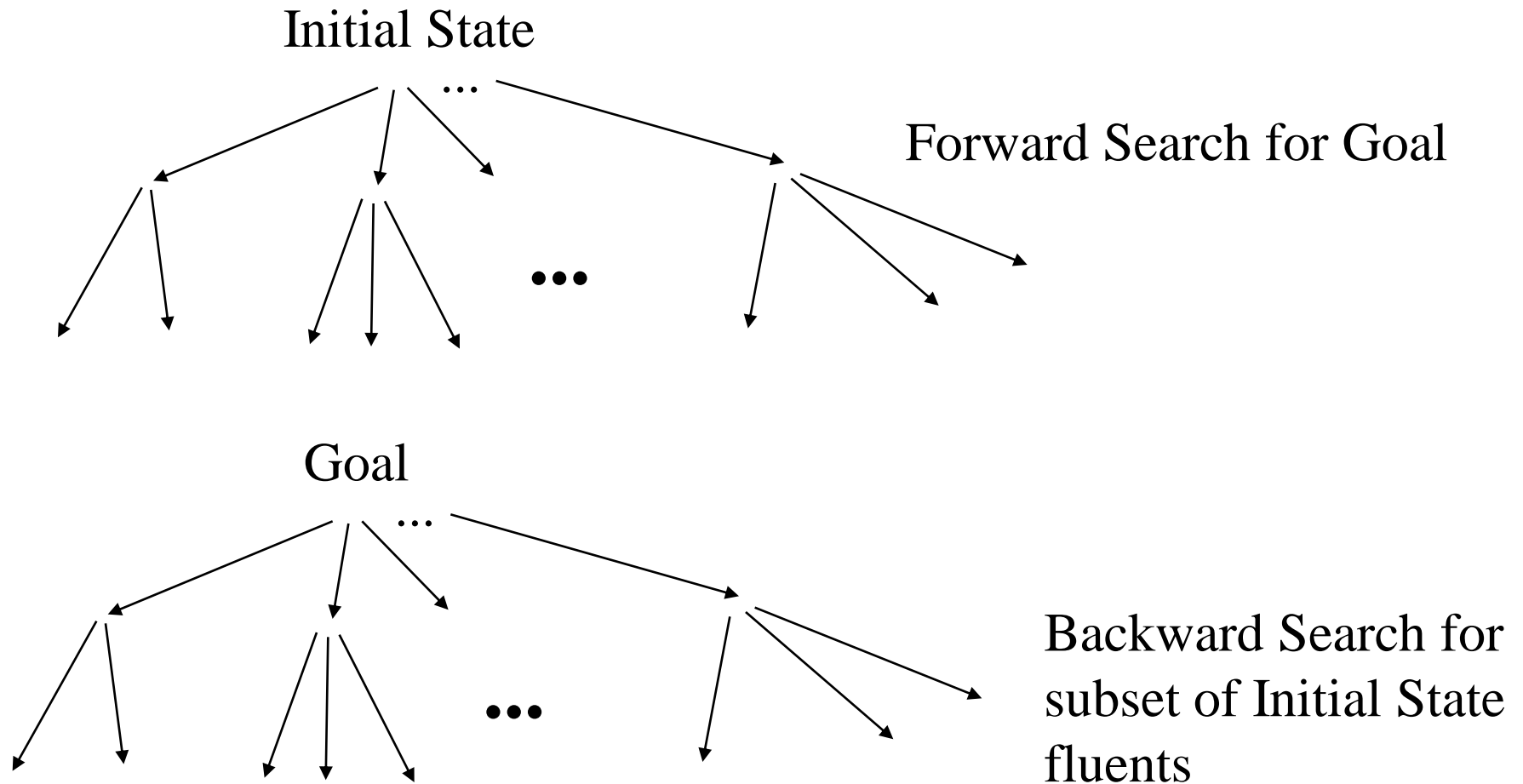
ADL – action description language

HTN – hierarchical task network planning

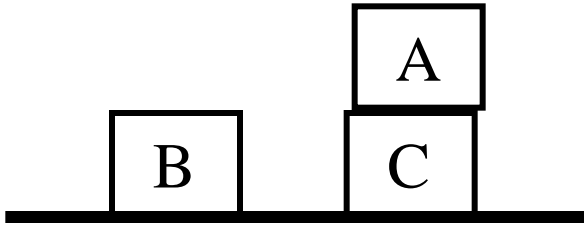
GraphPlan

SatPlan

# State Space Planner



One reason why planning beats searching



Initial State  $S_i$

$On(A, C, S_i)$

$On(C, Tbl, S_i)$

$On(B, Tbl, S_i)$

$Blk(A)$

$Blk(B)$

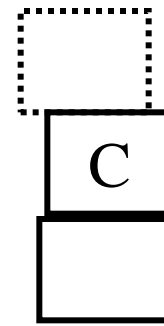
$Blk(C)$

$Table(Tbl)$

$Clr(A, S_i)$

$Clr(B, S_i)$

$Clr(Tbl, S_i)$



Goal  $?s$

Find an  $?x$  and  $?s$  s.t.:

$On(C, ?x, ?s)$

$Blk(?x)$

Negate Goal, add to axioms w/ Answer literal

Answer ( $?s$ ) should yield something like

Answer (Result (Move (C, Tbl, B),

Result (Move (A, C, Tbl),  $S_i$ ))