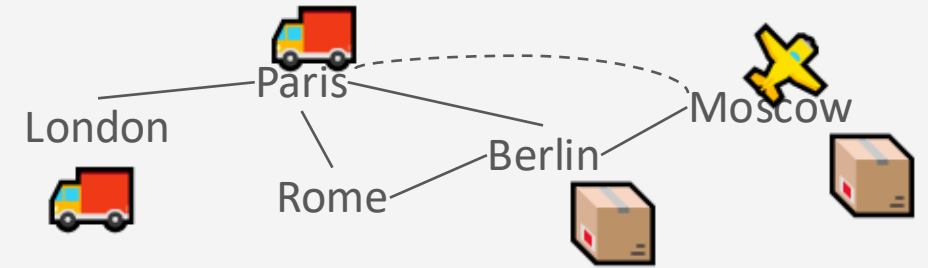


Universität  
Münster

# Automated Planning and Acting

## Decision Making: Structure



## Content: Planning and Acting

1. With **Deterministic** Models
2. With **Temporal** Models
3. With **Nondeterministic** Models
4. With **Probabilistic** Models
5. **By Decision Making**
  - A. *Foundations*
  - B. *Extensions*
  - C. *Structure*
    - Symmetries
    - Relations
6. **Human-aware Planning**

# Outline: Decision Making – Structure

## *Structure by Groups in the Agent Set*

- Agent types
- Partitioned decPOMDPs

## *Structure by Relations in the State Space*

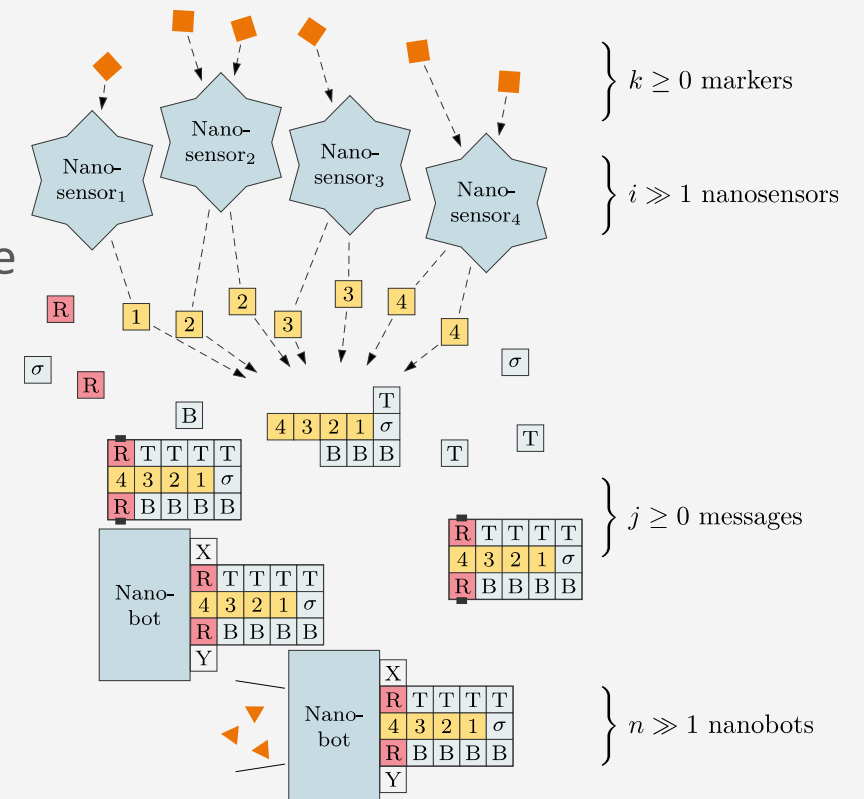
- Situation calculus
- First-order MDPs

## *Structure by Features in the State Space*

- Dynamic Bayesian networks
- Factored MDPs

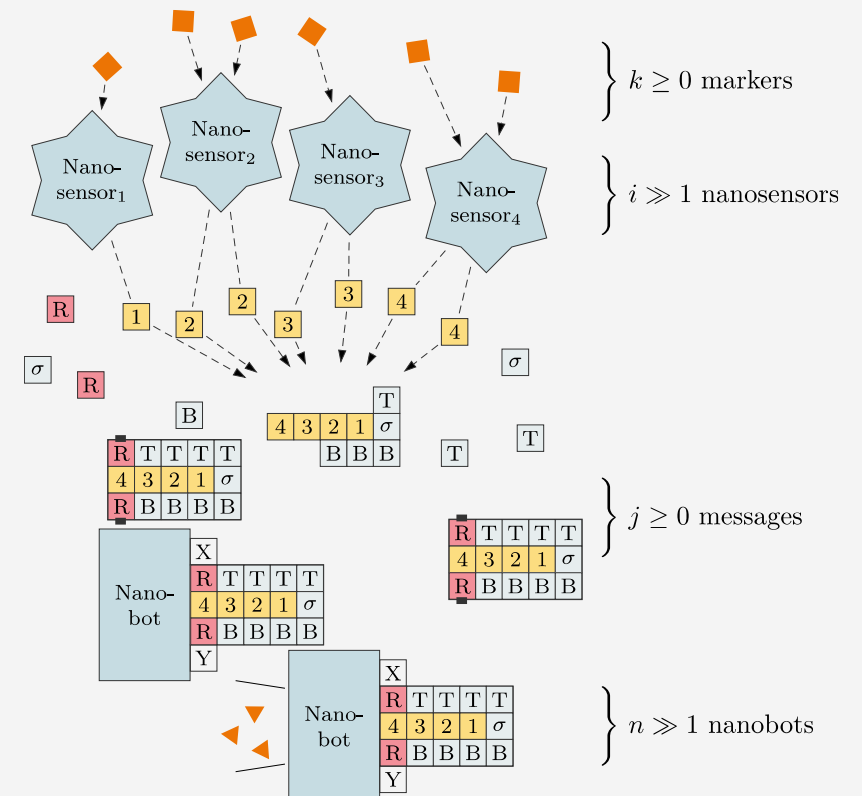
## Example: Medical Nanoscale Systems

- Nanoscale systems regularly consist of  $> 10,000$  nanoagents
  - Different types of agents: nanosensors, nanobots
- Application: DNA-based medical system
  - E.g., for diagnosis (modelled as an AND gate)
    - Nanosensors receptive to individual markers for a specific disease
      - Release individual tiles in presence of their individual markers
    - Tiles assemble themselves to form messages
    - Nanobots receptive to completely formed messages
      - Release markers of their own that signify presense of the disease
- Formal model necessary to argue about
  - Success rates
  - Sizes of agent sets



## Example: Medical Nanoscale Systems as a DecPOMDP

- Set of agents  $I$  consisting of nanosensors, nanobots
- Observations  $O_i$ : markers / messages present (or not)
  - Noisy process  $\rightarrow$  probabilistic behaviour
- Actions  $A_i$ : release of tiles / markers (or not)
  - Noisy process  $\rightarrow$  probabilistic behaviour
- Environment  $\rightarrow$  probabilistic behaviour
  - Presence in general of agents, markers, tiles, messages, or position more specifically  $\rightarrow$  movement over time
- Reward: Qualitative measure
  - Positive diagnosis only in presence of disease



## Reprise: Worst-case Complexity of DecPOMDP

- Space complexity
  - Transition model:  $\mathcal{O}(s \cdot s \cdot a^N)$
  - Sensor model:  $\mathcal{O}(s \cdot o^N)$  or  $\mathcal{O}(s \cdot o^N \cdot a^N)$
  - Reward function:  $\mathcal{O}(s)$  or  $\mathcal{O}(s \cdot a^N)$
- Runtime complexity of brute-force search
  - Evaluation cost of a joint policy:  $\mathcal{O}(s \cdot o^{Nh})$
  - Policy space:  $\mathcal{O}\left(a \frac{o^h - 1}{o - 1}\right)$
- Notations
  - $s = |S|$ 
    - State space size
  - $a = \max_{i \in I} |A_i|$ 
    - Largest individual action space size
  - $o = \max_{i \in I} |O_i|$ 
    - Largest individual action space size
  - $h$ 
    - Horizon

## Agent Types & Partitioned DecPOMDPs

- Types: Agents with the same sets of actions and observations
  - E.g., two nanosensors 1,2 receptive to the same marker and releasing the same tile
    - $A_1 = A_2 = \{0,1\}$ ; 0: do nothing, 1: release tile
    - $O_1 = O_2 = \{0,1\}$ ; 0: marker not present, 1: marker present
- Partitions the set of agents regarding actions, observations
  - Agent set  $I = \{I_1, \dots, I_K\}$  with  $I_1, \dots, I_K$  a partitioning of  $I$  ( $I = \bigcup_k I_k, I_k \cap I_{k'} = \emptyset, I_k \neq \emptyset$ )
  - For each partition  $I_k$ : one set of actions  $A_k$ , one set of observations  $O_k$  for all agents in  $I_k$
  - Expectation that  $K \ll N$
- Additional constraints / assumptions on same behaviour in  $T, R, \Omega$
- Partitions the set of agents completely, enabling more compact encodings

## Additional Assumption 1

- **Symmetric behaviour** in  $T, R, \Omega$ 
  - Idea: It does not matter which agent of a partition does something, only how many
  - Formally, for two agents  $i, j \in I_k$ , exchanging their actions / observations does not have an effect

$$T(s, s', a_i, a_j, a_{-i,-j}) = T(s, s', a_j, a_i, a_{-i,-j})$$

$$R(s, a_i, a_j, a_{-i,-j}) = R(s, a_j, a_i, a_{-i,-j})$$

$$\Omega(s, o_i, o_j, o_{-i,-j}) = \Omega(s, s', o_j, o_i, o_{-i,-j})$$

- Can be generalised to whole partitions, over all possible permutations  $\theta$  of actions  $\vec{a}_k$  / observations  $o_k$ , e.g., for  $T$ :

$$T(s, s', \vec{a}_k, \vec{a}_{-k}) = T(s, s', \theta(\vec{a}_k), \vec{a}_{-k})$$

- Enables counting how many agents do something and not which in particular did something
  - Encode in a histogram  $[\#(a_1), \dots, \#(a_l)]$  how many agents did actions  $A_k = \{a_1, \dots, a_l\}$
  - Number of histograms  $\binom{|I_k|+l-1}{l-1} \leq |I_k|^l$

## Counting in DecPOMDPs

- Also called counting constraint / assumption in  $T, R, \Omega$ 
  - All permutations  $\theta(\vec{a}_k)$  of a partition action  $\vec{a}_k$  map to the same probability
  - Allows for replacing joint actions / observations over  $N$  agents in  $T, R, \Omega$  with histograms for each partition
  - Example with two possible actions and a single partition

$S$	$S'$	$A_1^\#$	$\bar{T}(s, s', a'_1)$ $= P(s' s, a'_1)$
0	0	[0,2]	0.01
0	0	[1,1]	0.02
0	0	[2,0]	0.03
0	1	[0,2]	0.015
0	1	[1,1]	0.012
0	1	[2,0]	0.01
1	0	[0,2]	0.01
		⋮	

$S$	$S'$	$A_1$	$A_2$	$T(s, s', a_1, a_2)$ $= P(s' s, a_1, a_2)$
0	0	0	0	0.01
0	0	0	1	0.02
0	0	1	0	0.02
0	0	1	1	0.03
0	1	0	0	0.015
0	1	0	1	0.012
0	1	1	0	0.012
0	1	1	1	0.01
1	0	0	0	0.01
			⋮	

## Additional Assumption 2

- **Conditional independence** among agents in  $T, R, \Omega$ 
  - Idea: Given a state and the other partitions, two agents are conditionally independent
  - Formally, for two agents  $i, j \in I_k$ ,
 
$$T(s, s', a_i, a_j, a_{-i,-j}) = P(s'|s, a_i, a_j, a_{-i,-j}) = P_i(s'|s, a_i, a_{-i,-j}) \cdot P_j(s'|s, a_j, a_{-i,-j})$$
    - Analogous for  $R$  (with addition) and  $\Omega$
    - Can also be generalised to whole partitions, e.g., for  $T$ 

$$T(s, s', \vec{a}_k, \vec{a}_{-k}) = P(s'|s, \vec{a}_k, \vec{a}_{-k}) = \prod_{i \in I_k} P_i(s'|s, a_i, \vec{a}_{-k})$$
- Combining Assumption 1 + 2 yields that  $P_i = P_j$  for all  $i, j \in I_k$  for all  $k \in \{1, \dots, K\}$ 
  - Instead of having one large function  $T$  over  $S, S', A_1, \dots, A_N$
  - We have  $K$  smaller functions  $T_k$  over  $S, S', A_1, \dots, A_K$

## Isomorphism in DecPOMDPs

- Assumption 1 only used with Assumption 2 (not on its own)
- Also called isomorphic constraint / assumption in  $T, R, \Omega$ :  
Conditional independence between agents of a partition given joint state  
→ Enables factorisation of  $T, R, \Omega$

$$\bullet \text{ E.g., } T(s, s', a_1, a_2) = \underbrace{T_1(s, s', a_1)}_{T_1} \cdot \underbrace{T_2(s, s', a_2)}_{T_2} = \prod_{i \in I_K} T'(s, s', a_i)$$

$$T_1 = T_2 = T'$$

- Allows for replacing  $T, R, \Omega$  with joint actions / observations over  $N$  agents with  $K$  functions each, using single representative actions / observations

$S$	$S'$	$A_i$	$T'(s, s', a_i)$ $= P(s' s, a_i)$
0	0	0	0.01
0	0	1	0.03
0	1	0	0.015
0	1	1	0.01
1	0	0	0.01
			$\vdots$



## Remember: Multi-agent Dynamic Programming

- Iteratively eliminate weakly-dominated policy trees using a dynamic programming
  - A policy  $\pi_{i,j}^t$  of agent  $i$  is **weakly dominated** if there always exists another policy  $\pi_{i,j'}^t$  with higher value over all possible states  $s$  and all possible policies of the other agents  $\boldsymbol{\pi}_{-i}^t$  :
$$\forall s \in S, \boldsymbol{\pi}_{-i}^t \in \boldsymbol{\Pi}_{-i}^t \exists \pi_{i,j'}^t \in \Pi_i^t : V_i(s, \pi_{i,j}^t, \boldsymbol{\pi}_{-i}^t) \leq V_i(s, \pi_{i,j'}^t, \boldsymbol{\pi}_{-i}^t)$$
  - Solvable by a linear programme

### *Dynamic Programming Operator:*

**Input:** sets of policies and corresponding value vectors for each agent.

#### **Process:**

Perform exhaustive backup using policies for each agent.

Calculate new value vectors.

Prune weakly dominated policies per agent.

**Output:** Updated sets of policies and value vectors for each agent.

## Multi-agent Dynamic Programming for Isomorphic DecPOMDPs

- Isomorphic DecPOMDP: tuple  $(I = \{I_1, \dots, I_K\}, S, \{A_k\}_{k=1}^K, \{O_k\}_{k=1}^K, \bar{T}, \bar{R}, \bar{\Omega})$ 
  - Transition function  $\bar{T} = \{T_k(s', s, a_1, \dots, a_K)\}_{k=1}^K$
  - Reward function  $\bar{R} = \{R_k(s, a_1, \dots, a_K)\}_{k=1}^K$
  - Sensor model (observation function)  $\bar{\Omega} = \{\Omega_k(s, o_1, \dots, o_K)\}_{k=1}^K$
- Dynamic programming:
  - Exhaustive backup: once per partition using  $A_k, O_k$
  - Value computation: once per partition (value function needs to be updated accordingly)
  - Pruning: Check for weakly dominated policies per partition
    - A weakly dominated policy of  $i \in I_k$  is also a weakly dominated policy of  $j \in I_k$
    - Only need to consider every possible state and representative policy of each partition
  - Complexity: polynomial in the number of agents (due to value function computation)

## Counting DecPOMDPs

- Counting DecPOMDP: tuple  $(I = \{I_1, \dots, I_K\}, S, \{\#A_k\}_{k=1}^K, \{\#O_k\}_{k=1}^K, \bar{T}, \bar{R}, \bar{\Omega})$ 
  - Transition function  $\bar{T} = T(s, s', h_1^a, \dots, h_K^a)$
  - Reward function  $\bar{R} = R(s, h_1^a, \dots, h_K^a)$
  - Sensor model (observation function)  $\bar{\Omega} = \Omega(s, h_1^o, \dots, h_K^o)$
- **Problem:**
  - Builds policies out of histograms
  - Actually leads to a blow-up in the policy space:  $\mathcal{O}\left(\frac{aK(n^{ho}-1)}{n^{o-1}}\right)$
- **Solution:**
  - Use  $A_k, O_k$  in partitions
  - Count how many agents of each partitions follow which representative policies instead

## Multi-agent Dynamic Programming for Policy-counted DecPOMDPs

- Policy-counted DecPOMDP: tuple  $(I = \{I_1, \dots, I_K\}, S, \{A_k\}_{k=1}^K, \{O_k\}_{k=1}^K, \bar{T}, \bar{R}, \bar{\Omega})$ 
  - Transition function  $\bar{T} = T(s, s', h_1^a, \dots, h_K^a)$
  - Reward function  $\bar{R} = R(s, h_1^a, \dots, h_K^a)$
  - Sensor model (observation function)  $\bar{\Omega} = \Omega(s, h_1^o, \dots, h_K^o)$
- Dynamic programming: Difference to isomorphic case
  - As before: Backup, value computation per partition (accordingly updated value function)
  - But, cannot check pruning for a single representative policy for each partition
    - Cannot just consider the representative policies of other partitions
  - Thus, exploit structure in policy space by counting how many agents follow which policy
    - Consider all possible counted policies (all histograms of representative policies)
  - Complexity: polynomial in the number of agents

## Partitioned POSGs

- Remember: **Partially observable stochastic game (POSG)**
  - Dec-POMDP  $(I, S, \{A_i\}_{i \in I}, \{O_i\}_{i \in I}, T, \mathbf{R}, \Omega)$  but with individual reward functions  $\{R_i\}_{i \in I}$
  - Reward function  $R_i$  for each agent  $i \in I$
- Analogous assumptions for POSGs
  - Additionally,  $R_i = R_j$  for all  $i, j \in I_k$  for each  $k$
  - Same mathematical transformations possible
  - Allow for defining isomorphic and policy-counted (and counting) POSGs
- Multi-agent dynamic programming also a solution method for POSGs, therefore results carry over
  - Slightly different value computation, rest stays the same
- Complexities: polynomial in the number of agents

## Interim Summary: Structure by Groups in the Agent Set

- Types of agents with identical action and observation space
- Assumptions: symmetric behaviour, conditional independence
  - Counting
    - Permutations of actions of agents of the same partition map to the same probability / reward
    - Count occurrences → encode in histograms
  - Model complexity polynomial in the number of agents
- Isomorphic DecPOMDP: Use single representative agent per partition for policies
- Counting DecPOMDPs: Use histograms for policies → blows up state space
- Policy-counted DecPOMDP: Count how many agents follow each representative policy
- Dynamic programming possible per partition in all cases
  - Dynamic programming for isomorphic and policy-counted DecPOMDPs polynomial in the number of agents

# Outline: Decision Making – Structure

## *Structure by Groups in the Agent Set*

- Agent types
- Partitioned decPOMDPs

## ***Structure by Relations in the State Space***

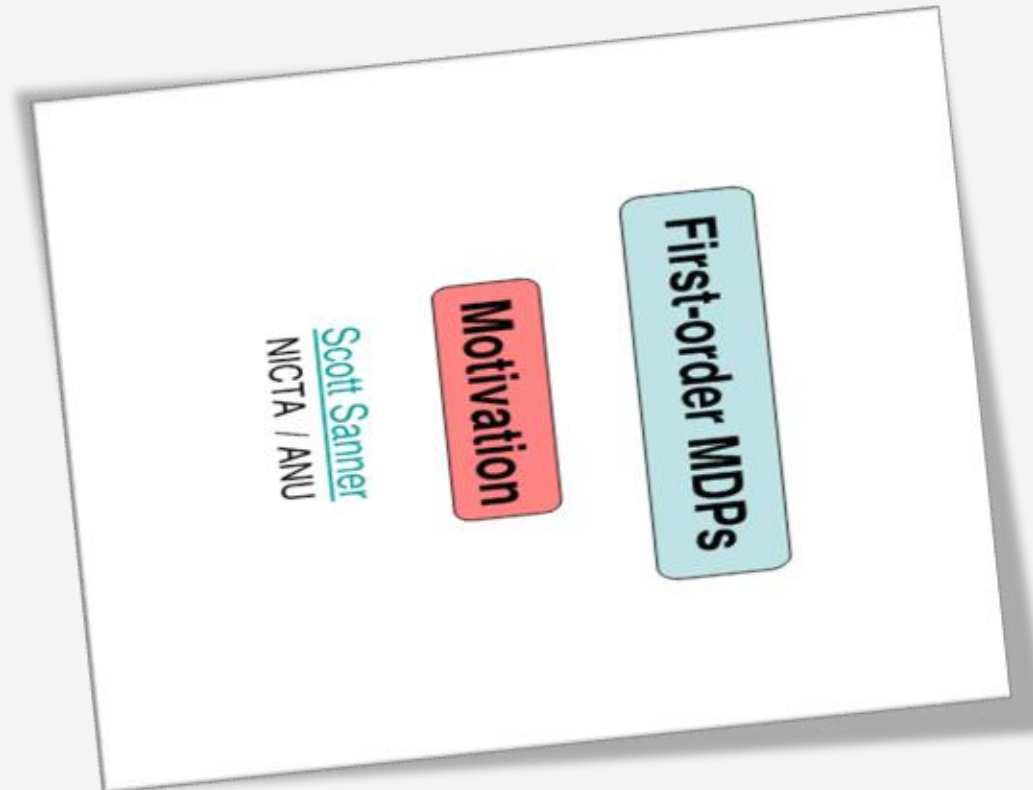
- Situation calculus
- First-order MDPs

## *Structure by Features in the State Space*

- Dynamic Bayesian networks
- Factored MDPs

# Acknowledgement

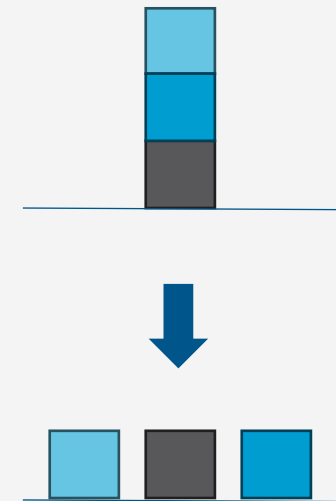
- Thanks to Scott Sanner!



## Motivation: Planning Languages

- Common languages:
  - STRIPS
  - PDDL
    - More expressive than STRIPS
    - For example, universal and conditional effects:

```
(:action put-all-blue-blocks-on-table
  :parameters ( )
  :precondition ( )
  :effect (forall (?b)
            (when (Blue ?b)
              (not (OnTable ?b))))))
```
- General Game Playing (GGP)
  - One or more agents

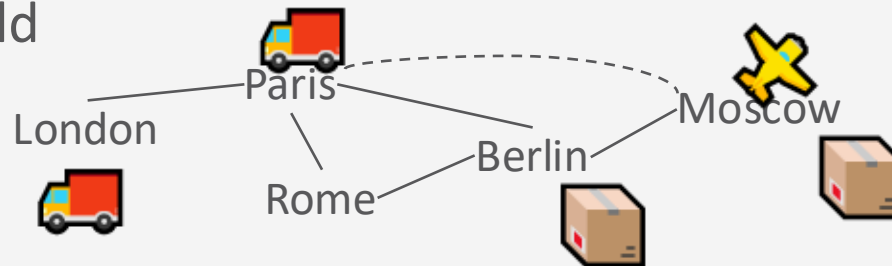


## Motivation: Benefits of Relational Languages

- STRIPS, PDDL, GGP are relational languages...
  - Refer to relational fluents:
    - E.g., *BoxIn(? b, ? c)*, *OnTable(? b)*
  - Specify relations between objects
  - Change over time
- Use first-order logic to specify...
  - Action preconditions
  - Action effects
  - Goals / rewards
    - E.g., `(forall (?b ?c) ((Destination ?b ?c) => (BoxIn ?b ?c)))`
- Are domain-independent and often compact!

## Motivation: How to Solve?

- Relational planning problem
  - E.g., box world



```

(:action load-box-on-truck-in-city
  :parameters (?b - box ?t - truck ?c - city)
  :precondition (and (BoxIn ?b ?c) (TruckIn ?t ?c))
  :effect (and (On ?b ?t) (not (BoxIn ?b ?c))))
  
```

- Solve *ground* problem for each domain instance?

- E.g., instance with 3 trucks , 2 planes , 3 boxes 

- Or solve lifted specification for *all* domains at once?

## Box World: Full (Relational) Specification

- Relational fluents:  $BoxIn(Box, City)$ ,  $TruckIn(Truck, City)$ ,  $BoxOn(Box, Truck)$
- Goal:  $[\exists Box : b. BoxIn(b, paris)]$
- Actions:
  - $load(Box : b, Truck : t)$ :
    - Effects:
      - when  $[\exists City : c. BoxIn(b, c) \wedge TruckIn(t, c)]$  then  $[BoxOn(b, t)]$
      - $\forall City : c.$  when  $[BoxIn(b, c) \wedge TruckIn(t, c)]$  then  $[\neg BoxIn(b, c)]$
  - $unload(Box : b, Truck : t)$ :
    - Effects:
      - $\forall City : c.$  when  $[BoxOn(b, t) \wedge TruckIn(t, c)]$  then  $[BoxIn(b, c)]$
      - when  $[\exists City : c. BoxOn(b, t) \wedge TruckIn(t, c)]$  then  $[\neg BoxOn(b, t)]$
  - $drive(Truck : t, City : c)$ :
    - Effects:
      - when  $[\exists City : c_1. TruckIn(t, c_1)]$  then  $[TruckIn(t, c)]$
      - $\forall City : c_1.$  when  $[TruckIn(t, c_1)]$  then  $[\neg TruckIn(t, c_1)]$

## Solving Ground Box World

- Apply planner to Box World grounded with respect to domain, e.g.,
  - Domain object instantiations:
    - $Box = \{box_1, box_2, box_3\}$ ,  $Truck = \{truck_1, truck_2\}$ ,  $City = \{paris, berlin, rome\}$
  - Ground fluents:
    - $BoxIn: \{BoxIn(box_1, paris), BoxIn(box_2, paris), BoxIn(box_3, paris), BoxIn(box_1, berlin), BoxIn(box_2, berlin), BoxIn(box_3, berlin), BoxIn(box_1, rome), BoxIn(box_2, rome), BoxIn(box_3, rome)\}$
    - $TruckIn: \{TruckIn(truck_1, paris), TruckIn(truck_2, paris), TruckIn(truck_1, berlin), TruckIn(truck_2, berlin), TruckIn(truck_1, rome), TruckIn(truck_2, rome)\}$
    - $BoxOn: \{BoxOn(box_1, truck_1), BoxOn(box_2, truck_1), BoxOn(box_3, truck_1), BoxOn(box_1, truck_2), BoxOn(box_2, truck_2), BoxOn(box_3, truck_2)\}$
  - Ground actions:
    - $load: \{load(box_1, truck_1), load(box_2, truck_1), load(box_3, truck_1), load(box_1, truck_2), load(box_2, truck_2), load(box_3, truck_2)\}$
    - $unload: \{unload(box_1, truck_1), unload(box_2, truck_1), unload(box_3, truck_1), unload(box_1, truck_2), unload(box_2, truck_2), unload(box_3, truck_2)\}$
    - $drive: \{drive(truck_1, paris), drive(truck_2, paris), drive(truck_1, berlin), drive(truck_2, berlin), drive(truck_1, rome), drive(truck_2, rome)\}$
  - Goal:  $[BoxIn(box_1, paris) \vee BoxIn(box_2, paris) \vee BoxIn(box_3, paris)]$

Number of fluents  
exponential in arity

Number of actions  
exponential in arity

Goal description exponential in  
number of nested quantifiers

## A First-order Solution to Box World

- Derive solution deductively at lifted PDDL level  $\rightarrow$  Optimal for any domain instantiation!

```
if ( $\exists b. \text{BoxIn}(b, \text{paris})$ ) then  
  do noop  
else if ( $\exists b^*, t^*. \text{TruckIn}(t^*, \text{paris}) \wedge \text{BoxOn}(b^*, t^*)$ ) then  
  do unload( $b^*, t^*$ )  
else if ( $\exists b, c, t^*. \text{BoxOn}(b, t^*) \wedge \text{TruckIn}(t, c)$ ) then  
  do drive( $t^*, \text{paris}$ )  
else if ( $\exists b^*, c, t^*. \text{BoxIn}(b^*, c) \wedge \text{TruckIn}(t^*, c)$ ) then  
  do load( $b^*, t^*$ )  
else if ( $\exists b, c_1^*, t^*, c_2. \text{BoxIn}(b, c_1^*) \wedge \text{TruckIn}(t^*, c_2)$ ) then  
  do drive( $t^*, c_1^*$ )  
else do noop
```

- Great, but how do I obtain this solution?

# Situation Calculus

- Logic formalism designed for representing and reasoning about dynamic domains
  - First introduced by John McCarthy in 1963
- Basic elements
  - Actions that can be performed in the world
  - Situations
  - Fluents that describe the state of the world
- Domain
  - Action precondition axioms, one for each action
  - Successor state axioms, one for each fluent
  - Axioms describing the world in various situations
  - Foundational axioms of the situation calculus: situations are histories, induction on situations

## Situation Calculus: Ingredients

- Actions
  - First-order terms with action parameters
  - E.g.,  $load(b, t)$ ,  $unload(b, t)$ ,  $drive(t, c)$
- Situations
  - Term that encodes action history
  - E.g.,  $s$ ,  $s_0$ ,  $do(load(b, t), s)$ ,  $do(load(b, t), drive(t, c), s)$
- Fluents
  - Relation whose truth value varies between situations
  - E.g.,  $BoxOn(b, t, s)$ ,  $TruckIn(t, c, s)$ ,  $Box(t, c, s)$
- Effects?

## Situation Calculus: PDDL to Effects

- Translate action effects into positive and negative effect axioms
  - $load(Box : b, Truck : t)$ :
    - when  $[\exists City : c. BoxIn(b, c) \wedge TruckIn(t, c)]$  then  $[BoxOn(b, t)]$
    - $\forall City : c. when [BoxIn(b, c) \wedge TruckIn(t, c)]$  then  $[\neg BoxIn(b, c)]$
  - $unload(Box : b, Truck : t)$ :
    - $\forall City : c. when [BoxOn(b, t) \wedge TruckIn(t, c)]$  then  $[BoxIn(b, c)]$
    - when  $[\exists City : c. BoxOn(b, t) \wedge TruckIn(t, c)]$  then  $[\neg BoxOn(b, t)]$
  - $drive(Truck : t, City : c)$ :
    - when  $[\exists City : c_1. TruckIn(t, c_1)]$  then  $[TruckIn(t, c)]$
    - $\forall City : c_1. when [TruckIn(t, c_1)]$  then  $[\neg TruckIn(t, c_1)]$

## Situation Calculus: PDDL to Effects

- Use rule to combine multiple effects  $C_1 \Rightarrow F, C_2 \Rightarrow F$  over the same fluent  $F$  into effect axioms:  $\gamma_F^+(\vec{x}, a, s) \Rightarrow F(\vec{x}, do(a, s)), \gamma_F^-(\vec{x}, a, s) \Rightarrow F(\vec{x}, do(a, s))$ 
  - Rule:  $[(C_1 \Rightarrow F) \wedge (C_2 \Rightarrow F)] \equiv [(C_1 \vee C_2) \Rightarrow F]$
  - As a sort of shorthand notation
    - E.g.,
      - $[\exists c. a = load(b, t) \wedge BIn(b, c, s) \wedge TIn(t, c, s)] \Rightarrow BOn(b, t, do(a, s)) \rightarrow \gamma_{BOn}^+(\vec{x}, a, s) \Rightarrow BOn(\vec{x}, do(a, s))$
      - $[\exists c. a = unload(b, t) \wedge BOn(b, t, s) \wedge TIn(t, c, s)] \Rightarrow \neg BOn(b, t, do(a, s)) \rightarrow \gamma_{BOn}^-(\vec{x}, a, s) \Rightarrow \neg BOn(\vec{x}, do(a, s))$
      - $[\exists t. a = unload(b, t) \wedge BOn(b, t, s) \wedge TIn(t, c, s)] \Rightarrow BIn(b, c, do(a, s)) \rightarrow \gamma_{BIn}^+(\vec{x}, a, s) \Rightarrow BIn(\vec{x}, do(a, s))$
      - $[\exists t. a = load(b, t) \wedge BIn(b, c, s) \wedge TIn(t, c, s)] \Rightarrow \neg BIn(b, c, do(a, s)) \rightarrow \gamma_{BIn}^-(\vec{x}, a, s) \Rightarrow \neg BIn(\vec{x}, do(a, s))$
      - $[\exists c_1. a = drive(t, c) \wedge TIn(t, c_1, s)] \Rightarrow TIn(t, c, do(a, s)) \rightarrow \gamma_{TIn}^+(\vec{x}, a, s) \Rightarrow TIn(\vec{x}, do(a, s))$
      - $[\exists c. a = drive(t, c) \wedge TIn(t, c_1, s)] \Rightarrow \neg TIn(t, c_1, do(a, s)) \rightarrow \gamma_{TIn}^-(\vec{x}, a, s) \Rightarrow \neg TIn(\vec{x}, do(a, s))$

## Frame Problem

- Positive and negative effect axioms specify what changes
  - $\gamma_{BOn}^+(\vec{x}, a, s) \Rightarrow BOn(\vec{x}, do(a, s))$        $\gamma_{BOn}^-(\vec{x}, a, s) \Rightarrow \neg BOn(\vec{x}, do(a, s))$
  - $\gamma_{BIn}^+(\vec{x}, a, s) \Rightarrow BIn(\vec{x}, do(a, s))$        $\gamma_{BIn}^-(\vec{x}, a, s) \Rightarrow \neg BIn(\vec{x}, do(a, s))$
  - $\gamma_{TIn}^+(\vec{x}, a, s) \Rightarrow TIn(\vec{x}, do(a, s))$        $\gamma_{TIn}^-(\vec{x}, a, s) \Rightarrow \neg TIn(\vec{x}, do(a, s))$
- Assume completeness regarding these effect axioms:
  - That is, assume that  $\gamma_F^+(\vec{x}, a, s) \Rightarrow F(\vec{x}, do(a, s))$ ,  $\gamma_F^-(\vec{x}, a, s) \Rightarrow \neg F(\vec{x}, do(a, s))$  characterise all the conditions under which an action  $a$  changes the value of fluent  $F$
  - Formalise as explanation closure axioms
    - $\neg F(\vec{x}, s) \wedge F(\vec{x}, do(a, s)) \Rightarrow \gamma_F^+(\vec{x}, a, s) \quad \equiv \quad \neg F(\vec{x}, s) \wedge \neg \gamma_F^+(\vec{x}, a, s) \Rightarrow \neg F(\vec{x}, do(a, s))$ 
      - If  $F$  was false and was made true by doing action  $a$ , then condition  $\gamma_F^+$  must have been true
    - $F(\vec{x}, s) \wedge \neg F(\vec{x}, do(a, s)) \Rightarrow \gamma_F^-(\vec{x}, a, s) \quad \equiv \quad F(\vec{x}, s) \wedge \neg \gamma_F^-(\vec{x}, a, s) \Rightarrow F(\vec{x}, do(a, s))$ 
      - If  $F$  was true and was made false by doing action  $a$  then condition  $\gamma_F^-$  must have been true

## Successor State Axioms (SSAs)

- Inputs / Requirements
  - Unique names for actions / arguments
  - Positive and negative effect axioms
    - $\gamma_F^+(\vec{x}, a, s) \Rightarrow F(\vec{x}, do(a, s)), \gamma_F^-(\vec{x}, a, s) \Rightarrow \neg F(\vec{x}, do(a, s))$
  - Explanation closure axioms
    - $\neg F(\vec{x}, s) \wedge F(\vec{x}, do(a, s)) \Rightarrow \gamma_F^+(\vec{x}, a, s), F(\vec{x}, s) \wedge \neg F(\vec{x}, do(a, s)) \Rightarrow \gamma_F^-(\vec{x}, a, s)$
  - Integrity:  $\neg \exists \vec{x}, a, s. \gamma_F^+(\vec{x}, a, s) \wedge \gamma_F^-(\vec{x}, a, s)$
- SSA for each  $F$ :
  - $F(\vec{x}, do(a, s)) \equiv \gamma_F^+(\vec{x}, a, s) \vee (F(\vec{x}, s) \wedge \neg \gamma_F^-(\vec{x}, a, s))$
  - Shorthand:
    - $F(\vec{x}, do(a, s)) \equiv \Phi_F(\vec{x}, a, s)$

## Successor State Axioms (SSAs): Example

- SSA for each  $F$ :  $F(\vec{x}, do(a, s)) \equiv \gamma_F^+(\vec{x}, a, s) \vee (F(\vec{x}, s) \wedge \neg\gamma_F^-(\vec{x}, a, s))$ 
  - Shorthand:  $F(\vec{x}, do(a, s)) \equiv \Phi_F(\vec{x}, a, s)$
- $BoxOn(b, t, do(a, s)) \equiv \Phi_{BoxOn}(b, t, a, s)$   
 $\equiv [\exists c. a = load(b, t) \wedge BoxIn(b, c, s) \wedge TruckIn(t, c, s)]$   
 $\vee (BoxOn(b, t, s) \wedge \neg[\exists c. a = unload(b, t) \wedge BoxOn(b, t, s) \wedge TruckIn(t, c, s)])$
- $BoxIn(b, c, do(a, s)) \equiv \Phi_{BoxIn}(b, c, a, s)$   
 $\equiv [\exists t. a = unload(b, t) \wedge BoxOn(b, t, s) \wedge TruckIn(t, c, s)]$   
 $\vee (BoxIn(b, c, s) \wedge \neg[\exists t. a = load(b, t) \wedge BoxIn(b, c, s) \wedge TruckIn(t, c, s)])$
- $TruckIn(t, c, do(a, s)) \equiv \Phi_{TruckIn}(t, c, a, s)$   
 $\equiv [\exists c_1. a = drive(t, c) \wedge TruckIn(t, c_1, s)]$   
 $\vee (TruckIn(t, c, s) \wedge \neg[\exists c_1. a = drive(t, c) \wedge TruckIn(t, c_1, s)])$

# Regression

- Idea: Use SSAs to regress from goal towards a (possibly only partially defined) initial state
  - A bit like lifted backward search
- Regression
  - If  $\phi$  held after action  $a$ , then *regression* is the  $\phi'$  that held before action  $a$
  - Exploit following properties
    - $Regr(\neg\psi) = \neg Regr(\psi)$
    - $Regr(\psi_1 \wedge \psi_2) = Regr(\psi_1) \wedge Regr(\psi_2)$
    - $Regr((\exists x)\psi) = (\exists x)Regr(\psi)$
    - $Regr\left(F(\vec{x}, do(a, s))\right) = \Phi_F(\vec{x}, a, s)$

## Regression: Example

- Given:  $\exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s))$
- Regress through  $\text{unload}(b^*, t^*)$

$$\begin{aligned} & \text{Regr} \left( \exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s)) \right) \\ &= \exists b. \text{Regr} \left( \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s)) \right) \end{aligned}$$

$$= \exists b. \Phi_{\text{BoxIn}}(b, \text{paris}, \text{unload}(b^*, t^*), s)$$

$$= \exists b. [\exists t. \text{unload}(b^*, t^*) = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, \text{paris}, s)]$$

$$\vee (\text{BoxIn}(b, \text{paris}, s))$$

$$\wedge \neg [\exists t. \text{unload}(b^*, t^*) = \text{load}(b, t) \wedge \text{BoxIn}(b, \text{paris}, s) \wedge \text{TruckIn}(t, \text{paris}, s)]$$

$$= [\exists b, t. b = b^* \wedge t = t^* \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s)$$

$$= [(\exists b. b = b^*) \wedge (\exists t. t = t^*) \wedge \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)]$$

$$\vee \exists b. \text{BoxIn}(b, \text{paris}, s)$$

$$= [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s)$$

- If  $\phi$  held after action  $a$ , then *regression* is the  $\phi'$  that held before action  $a$
- Exploit following properties
  - $\text{Regr}(\neg\psi) = \neg\text{Regr}(\psi)$
  - $\text{Regr}(\psi_1 \wedge \psi_2) = \text{Regr}(\psi_1) \wedge \text{Regr}(\psi_2)$
  - $\text{Regr}((\exists x)\psi) = (\exists x)\text{Regr}(\psi)$
  - $\text{Regr}(F(\vec{x}, \text{do}(a, s))) = \Phi_F(\vec{x}, a, s)$

Cannot be made true

$$\rightarrow \phi \wedge \neg[\perp] \equiv \phi \wedge \top \equiv \phi$$

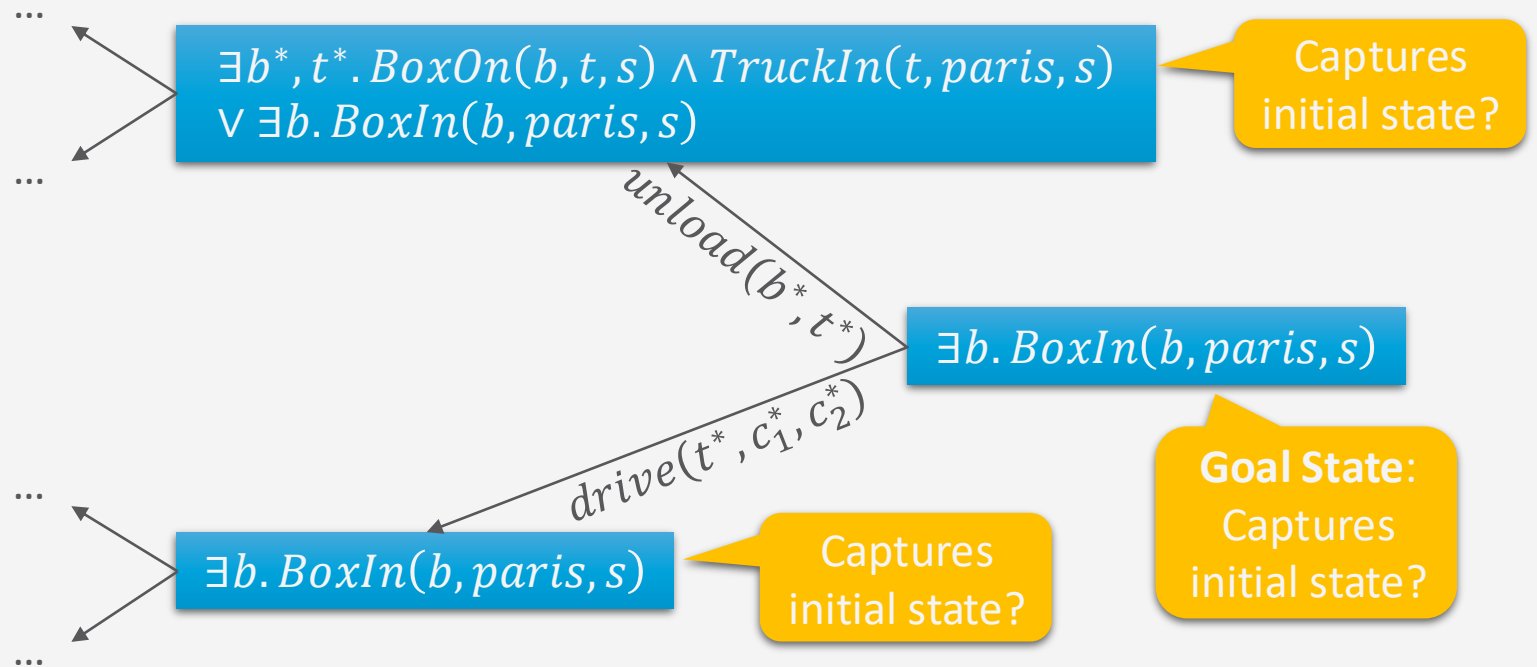
Make non-empty domain assumption for  $b, t$

## Regression: Example

- Given:  $\exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s))$
- Regress through  $\text{unload}(b^*, t^*)$ 
  - $\text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s)))$   
 $= [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s)$
- To get action instantiations of  $\text{unload}(b^*, t^*)$ , query knowledge base (KB, i.e., planning domain)
  - Existentially quantify  $b^*, t^*$  and obtain instances via query extraction w.r.t. KB
    - KB consists of first-order state and action abstraction  $\rightarrow$  do not have to enumerate all states,  $b^*, t^*$
    - $\exists b^*, t^*. \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s)))$   
 $= \exists b^*, t^*. [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s)$   
 $= [\exists b^*, t^*. \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s)$

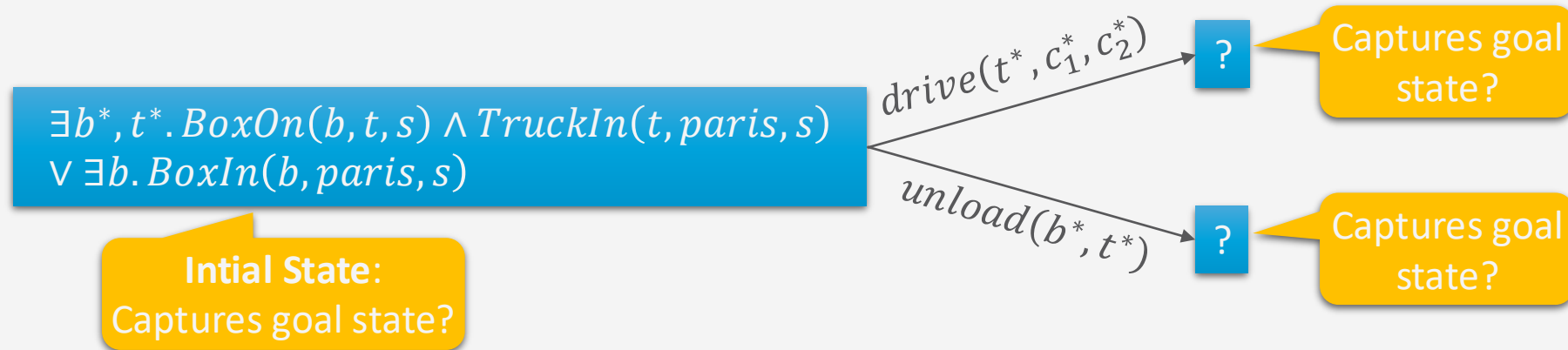
# Regression Planning

- Define abstract goal state
  - E.g.,  $\exists b. \text{BoxIn}(b, \text{paris}, s)$
  - Check if regression through action sequence holds in initial state
- Goal regression planning
  - Provide initial state, actions
    - Initial state description can be partial
  - Use regression to tell whether goal will hold



## Progression and Forward Search?

- Can we do lifted forward-search planning?



- Progression not first-order definable! (Reiter, 2001)
- Could progress ground state
  - But this does not exploit first-order structure

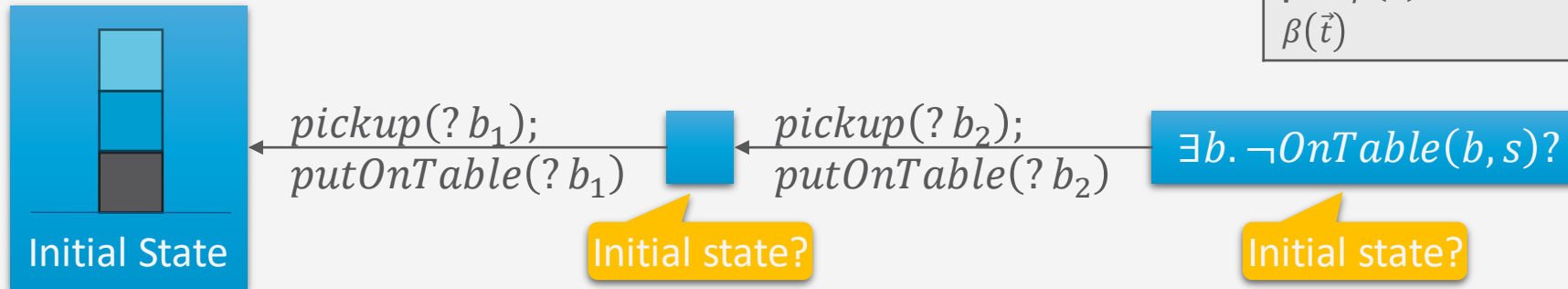
## Golog: Restricted Plan Search

- **ALGOL** in **LOGic**
  - Search the space of sequential action plans
  - Regress actions to initial state to test reachability
  - Restrict action space with program:

$\alpha$	primitive action
$\phi?$	condition test
$(\delta_1, \delta_2)$	sequence
<b>if <math>\phi</math> then <math>\delta</math> endif</b>	conditional
<b>while <math>\phi</math> then <math>\delta</math> endwhile</b>	loop
$(\delta_1   \delta_2)$	nondeterministic choice of actions
$\pi \vec{x} [\delta]$	nondeterministic choice of arguments
$\delta^*$	nondeterministic iteration
<b>proc <math>\beta(\vec{x}) \delta</math> endProc</b>	procedure call definition
$\beta(\vec{t})$	procedure call

## Golog: Example

- Golog program
  - $(\pi b [\neg \text{OnTable}(b, s)?, \text{pickup}(b), \text{putOnTable}(b)])^*, \forall b. \text{OnTable}(b, s)?$
- Diagram of Golog planning



- Heavily restricted action sequences
- Program exploits first-order action abstraction
- Initial state need not be fully known

$\alpha$	primitive action
$\phi?$	condition test
$(\delta_1, \delta_2)$	sequence
<b>if <math>\phi</math> then <math>\delta</math> endif</b>	conditional
<b>while <math>\phi</math> then <math>\delta</math> endwhile</b>	loop
$(\delta_1   \delta_2)$	nondeterministic choice of actions
$\pi \vec{x} [\delta]$	nondeterministic choice of arguments
$\delta^*$	nondeterministic iteration
<b>proc <math>\beta(\vec{x}) \delta</math> endProc</b>	procedure call definition
$\beta(\vec{t})$	procedure call

## Frame Problem

- Frame problem: How to (*compactly*) specify what does not change?
  - Intuition: “What does not change, remains the same.”
    - Reiter’s so-called Default Solution
  - Not so easy to specify
    - Moving one thing might move another thing, even though the other thing is never directly touched
    - How to distinguish between relevant and irrelevant side effects? And use that efficiently?
- Default solution to frame problem given as successor state axioms (SSAs), which we construct next

## Interim (Interim) Summary

- Situation calculus to describe a relational world
  - Can convert PDDL (and state-variable domains) into effect axioms
  - Derive SSAs from effect axioms
    - Using default solution to frame problem
- Regression operator
  - Extract action instantiation to achieve goal
- Regression planning
  - Initial state need not be fully specified
  - Exploit state and action abstraction
    - Avoid enumerating all state and action instances
- Frame problem

Next step: Extend this idea for decision-theoretic planning with uncertain action outcomes

$$rCase(s) = \begin{array}{|l|l|} \hline \forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s) & 1 \\ \hline \neg(\forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s)) & 0 \\ \hline \end{array}$$

# First-order MDPs

Structure by Relations in the State Space

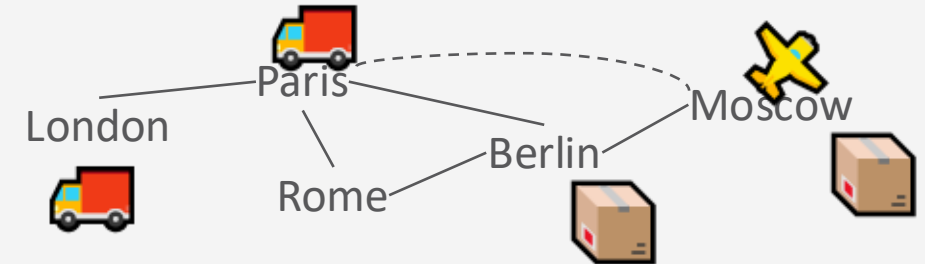
## First-order MDP (FOMDP)

- Components of MDP in an FOMDP specified as a collection of *case statements*
  - E.g., express reward in Box World FOMDP as

$$rCase(s) = \begin{array}{|l|l|} \hline \forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s) & 1 \\ \hline \neg(\forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s)) & 0 \\ \hline \end{array}$$

- Operators: define unary and binary case operations
  - E.g., cross-sum  $\oplus$  (or  $\ominus$ ,  $\otimes$ ) of cases

$$\begin{array}{|c|c|} \hline \phi & 10 \\ \hline \neg\phi & 20 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \varphi & 3 \\ \hline \neg\varphi & 4 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \phi \wedge \varphi & 10 + 3 \\ \hline \phi \wedge \neg\varphi & 10 + 4 \\ \hline \neg\phi \wedge \varphi & 20 + 3 \\ \hline \neg\phi \wedge \neg\varphi & 20 + 4 \\ \hline \end{array}$$



# Stochastic Actions and First-order Decision-theoretic Regression (FODTR)

- Stochastic actions using deterministic situation calculus
  - User's stochastic action, e.g.,  $A(x) = load(b, t)$
  - Nature's choice, e.g.,  $n(x) \in \{loadS(b, t), loadF(b, t)\}$
  - Probability distribution over nature's choice, e.g.,

Probability distribution → Adds up to 1 over success and failure choice

$$0.1 + 0.9 = 1$$

$$0.6 + 0.4 = 1$$

$$P(loadS(b, t) | load(b, t)) =$$

$snow(s)$	0.1
$\neg snow(s)$	0.6

$$P(loadF(b, t) | load(b, t)) =$$

$snow(s)$	0.9
$\neg snow(s)$	0.4

- First-order decision-theoretic regression (FODTR)
  - FODTR = *expectation* of regression:

$$FODTR[vCase(s), A(\vec{x})] = \mathbf{E}_{P(n(\vec{x}) | A(\vec{x}))} [Regr(vCase(s), n(\vec{x}))]$$

## FODTR & Q-Functions

- Result of FODTR is a case statement encoding a first-order Q-function

$$FODTR[vCase(s), A(\vec{x})] = R(s) \oplus \gamma \bigoplus_{j=1}^k P(n_j(\vec{x}), A(\vec{x}), s) \otimes Regr \left( V \left( do \left( n_j(\vec{x}) \right), s \right) \right)$$

- E.g.,

$$FODTR[vCase(s), unload(b^*, t^*)]$$

$$= rCase(s) \oplus \gamma \bigoplus_{j=1}^k pCase(n_j(\vec{x}), unload(b^*, t^*), s)$$

$$\otimes \begin{array}{|l|l|} \hline Regr \left( \exists b. BoxIn(b, paris, do(n_j(\vec{x}), s)) \right) & 10 \\ \hline Regr \left( \neg \exists b. BoxIn(b, paris, do(n_j(\vec{x}), s)) \right) & 0 \\ \hline \end{array}$$

$$rCase(s) = \begin{array}{|l|l|} \hline \exists b. BoxIn(b, paris, s) & 10 \\ \hline \neg(\exists b. BoxIn(b, paris, s)) & 0 \\ \hline \end{array}$$

$$pCase(loadS(b, t), load(b, t), s) = \begin{array}{|l|l|} \hline \top & 0.9 \\ \hline \end{array}$$

$$pCase(unloadS(b, t), unload(b, t), s) = \begin{array}{|l|l|} \hline \top & 0.9 \\ \hline \end{array}$$

$$pCase(driveS(c, t), drive(c, t), s) = \begin{array}{|l|l|} \hline \top & 1 \\ \hline \end{array}$$

## FODTR & Q-Functions

$$FODTR[vCase(s), unload(b^*, t^*)]$$

$$= rCase(s) \oplus \gamma \bigoplus_{j=1}^k pCase(n_j(\vec{x}), unload(b^*, t^*), s) \otimes$$

$Regr(\exists b. BoxIn(b, paris, do(n_j(\vec{x}), s)))$	10
$Regr(\neg \exists b. BoxIn(b, paris, do(n_j(\vec{x}), s)))$	0.0

$$= rCase(s) \oplus \gamma \left[ \begin{array}{l} \text{T} \quad 0.9 \quad \otimes \\ \text{T} \quad 0.1 \quad \otimes \end{array} \right]$$

$Regr(\exists b. BoxIn(b, paris, do(unloadS(b^*, t^*), s)))$	10
$Regr(\neg \exists b. BoxIn(b, paris, do(unloadS(b^*, t^*), s)))$	0.0

$Regr(\exists b. BoxIn(b, paris, do(unloadF(b^*, t^*))))$	10
$Regr(\neg \exists b. BoxIn(b, paris, do(unloadF(b^*, t^*))))$	0.0

## FODTR & Q-Functions

$$FODTR[vCase(s), unload(b^*, t^*)]$$

$$= rCase(s) \oplus \gamma \bigoplus_{j=1}^k pCase(n_j(\vec{x}), unload(b^*, t^*), s) \otimes$$

$Regr\left(\exists b. BoxIn(b, paris, do(n_j(\vec{x}), s))\right)$	10
$Regr\left(\neg\exists b. BoxIn(b, paris, do(n_j(\vec{x}), s))\right)$	0.0

$$= \begin{array}{|c|c|} \hline \exists b. BoxIn(b, paris, s) & 10 \\ \hline \neg & 0.0 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \exists b. BoxIn(b, paris, s) & 9.0 \\ \hline \neg & 0.0 \\ \hline \end{array}$$

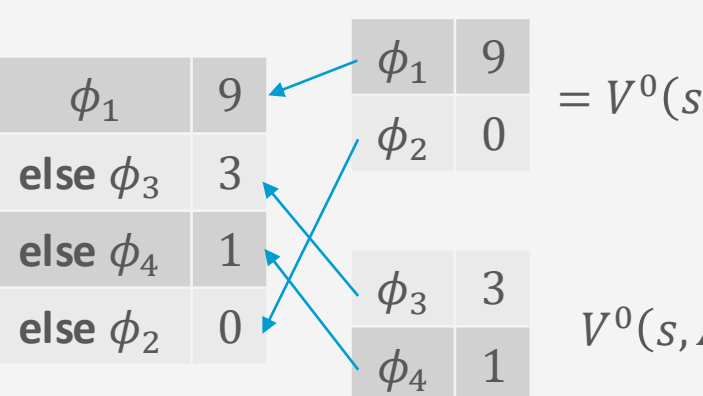
$$\oplus \begin{array}{|c|c|} \hline \exists b^*, t^*. BoxOn(b, t, s) \wedge TruckIn(t, paris, s) \vee \exists b. BoxIn(b, paris, s) & 8.1 \\ \hline \neg & 0.0 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|} \hline \exists b. BoxIn(b, paris, s) & 19.0 & \rightarrow noop \\ \hline \neg \wedge [\exists b, t. BoxOn(b, t, s) \wedge TruckIn(t, paris, s)] & 8.1 & \rightarrow unload(b^*, t^*) \\ \hline \neg & 0.0 & \rightarrow noop \\ \hline \end{array}$$

# Symbolic Dynamic Programming (SDP)

- What value if 0-stages-to-go?
  - Immediate reward:  $V^0(s) = rCase(s)$
- What value if 1-state-to-go?
  - We know value for each action  $\rightarrow$  Take maximum for each state

$$V^1(s) = \max_s \begin{cases} \begin{matrix} \phi_1 & 9 \\ \phi_2 & 0 \end{matrix} = V^0(s, A_1) \\ \begin{matrix} \phi_3 & 3 \\ \phi_4 & 1 \end{matrix} = V^0(s, A_2) \end{cases}$$

$$V^1(s) = \begin{matrix} \phi_1 & 9 \\ \text{else } \phi_3 & 3 \\ \text{else } \phi_4 & 1 \\ \text{else } \phi_2 & 0 \end{matrix}$$


- Value iteration
  - Obtain  $V^{n+1}$  from  $V^n$  until  $(V^{n-1} \ominus V^n) < \epsilon$

## Value Iteration for $t = 1, 2$ of the Box World Example

- With increasing iterations, the sequence of actions considered gets longer

$vCase^1(s) =$	$\exists b. BoxIn(b, paris, s)$	19.0
	$\neg \text{“} \wedge [\exists c. BoxOn(b, t, s) \wedge TruckIn(t^*, paris, s)]$	8.1
	$\neg \text{“}$	0.0
$vCase^2(s) =$	$\exists b. BoxIn(b, paris, s)$	26.1
	$\neg \text{“} \wedge [\exists b, t. BoxOn(b, t, s) \wedge TruckIn(t, paris, s)]$	15.4
	$\neg \text{“} \wedge [\exists b, c, t. BoxOn(b, t, s) \wedge TruckIn(t, c, s)]$	7.3
	$\neg \text{“}$	0.0

```

if ( $\exists b. BoxIn(b, paris)$ ) then
  do noop
else if ( $\exists b^*, t^*. TruckIn(t^*, paris) \wedge BoxOn(b^*, t^*)$ )
  do unload( $b^*, t^*$ )
else if ( $\exists b, c, t^*. BoxOn(b, t^*) \wedge TruckIn(t, c)$ ) then
  do drive( $t^*, paris$ )
else if ( $\exists b^*, c, t^*. BoxIn(b^*, c) \wedge TruckIn(t^*, c)$ ) then
  do load( $b^*, t^*$ )
else if ( $\exists b, c_1^*, t^*, c_2. BoxIn(b, c_1) \wedge TruckIn(t^*, c_2)$ )
  do drive( $t^*, c_1^*$ )
else do noop
  
```

# First-order Algebraic Decision Diagrams (FOADDs)

- We want to compactly represent arbitrary case statements
  - E.g.,

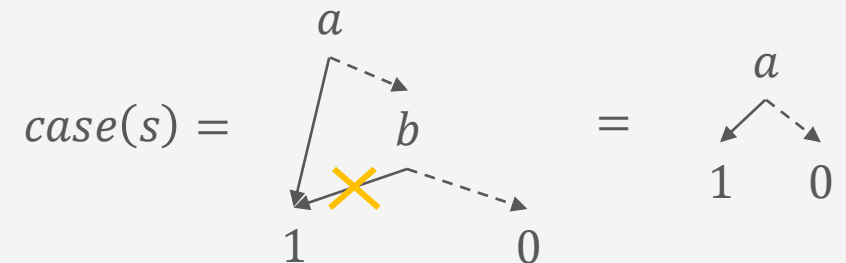
$$case(s) = \begin{array}{|l|l|} \hline \exists x. [A(x) \vee \forall y. A(x) \wedge B(x) \wedge \neg A(y)] & 1 \\ \hline \neg(\exists x. [A(x) \vee \forall y. A(x) \wedge B(x) \wedge \neg A(y)]) & 0 \\ \hline \end{array}$$

- Push down quantifiers, expose propositional structure  $\rightarrow$  convert into FOADD

$$[\exists x. A(x)] \vee ([\exists x. A(x) \wedge B(x)] \wedge [\forall y. \neg A(y)])$$

Variable	Variable $\Leftrightarrow$ FOL KB
$a$	$\equiv [\exists x. A(x)]$
$b$	$\equiv [\exists x. A(x) \wedge B(x)]$

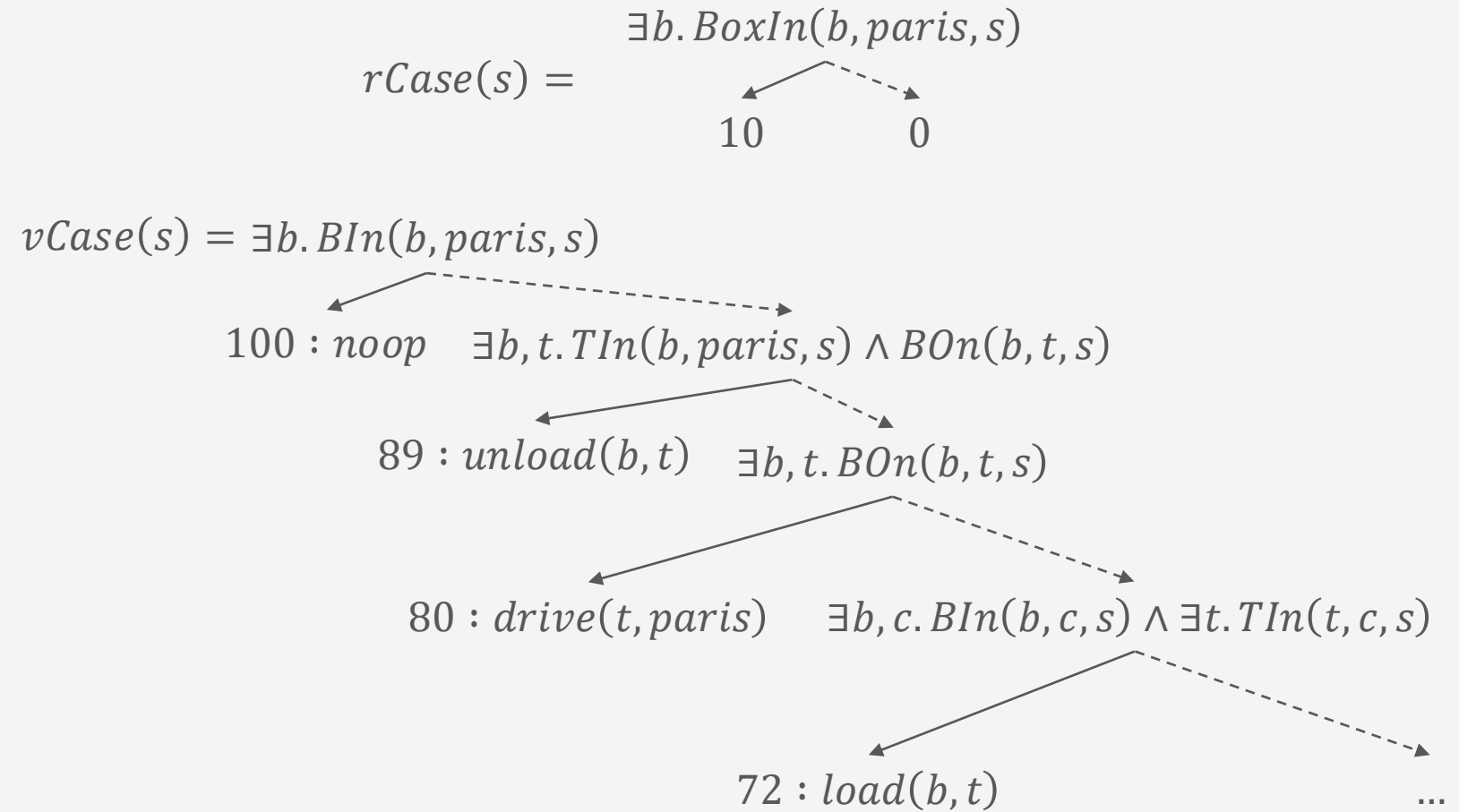
$$case(s) = \begin{array}{|l|l|} \hline a \vee (b \wedge \neg a) & 1 \\ \hline \neg(a \vee (b \wedge \neg a)) & 0 \\ \hline \end{array}$$



First-order context-specific independence

## Results for SDP with FOADDs

- Encode case statements with FOADDs
  - Solid line: true case
  - Dotted line: false case
- Use FOADD operations for structured SDP
  - E.g., Box World
    - Using  $\gamma = 0.9$



Factored SDP for factored FOMDPs  
[Sanner and Boutilier, 2007]

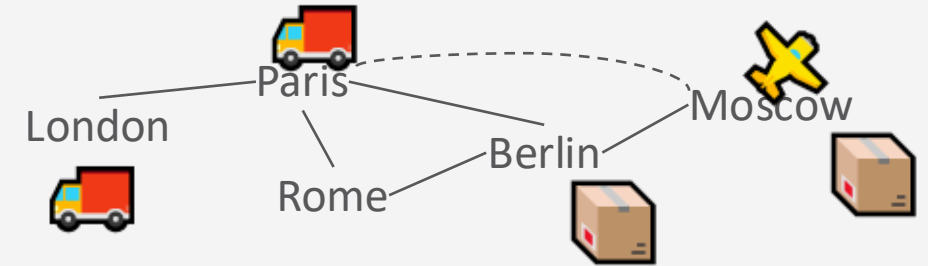
## Correctness of SDP

- Show SDP for FOMDPs is correct w.r.t. ground MDP



# Caveats of First-order Planning

- Many problems have topologies
  - E.g., reachability constraints in logistics
- If topology not fixed a priori
  - First-order solution must consider  $\infty$  topologies
  - In general case, leads to  $\infty$  values / policies
    - Universal rewards
    - Value function must distinguish  $\infty$  cases
    - Policy will also likely be  $\infty$




$$rCase(s) = \begin{array}{|l|l|} \hline \forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s) & 1 \\ \hline \neg(\forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s)) & 0 \\ \hline \end{array}$$
  

$$V^t(s) = \begin{array}{|l|l|} \hline \forall b, c. Dest(b, c) \Rightarrow BoxIn(b, c, s) & 1 \\ \hline One\ box\ not\ at\ destination & \gamma \\ \hline Two\ boxes\ not\ at\ destination & \gamma^2 \\ \hline \vdots & \vdots \\ \hline t - 1\ boxes\ not\ at\ destination & \gamma^{t-1} \\ \hline \end{array}$$

## Caveats of First-order Planning

- Unreachable states
  - PDDL domains often under-constrained
    - E.g., logistics: one box cannot be in two cities at once
  - Constraints implicitly obeyed in initial state
    - Action effects cannot violate constraints
      - Reachable legal states are small subset of all states
    - But (P)PDDL does not constrain legal states
- If no background theory to restrict legal states
  - First-order planning must solve for all states
    - When initial state unknown
  - Where majority of states are actually illegal
- First-order planning w/o initial state solves more difficult problem than search-based solutions
  - Initial state contains implicit constraint information
  - Reachable state space is small subset of all states

Suggests need for hybrid  
first-order / search-based approaches



## A Note on First-order Modelling in Reinforcement Learning

- Novel propositional situations worth exploring may be instances of a well-known context in the relational setting → *exploitation* promising
  - E.g., household robot learning water-taps
    - Having opened one or two water-taps in a kitchen, one can expect other water-taps in kitchens to work similarly
      - ⇒ Priority for exploring water-taps in kitchens in general reduced
      - ⇒ Information gathered likely to carry over to water-taps in other places
      - ❖ **Hard to model in propositional setting: each water-tap is novel**

## Interim Summary

- FOMDPs are one model for lifted decision-theoretic planning
  - Exploit state and action abstraction for MDPs
- Use situation calculus specified action theory
- Use case statements to represent reward, probabilities
- Symbolic dynamic programming = lifted DP
  - Use FOADDs to compactly represent case statements
  - First-order context-specific independence to compactify FOADDs

# Outline: Decision Making – Structure

## *Structure by Groups in the Agent Set*

- Agent types
- Partitioned decPOMDPs

## *Structure by Relations in the State Space*

- Situation calculus
- First-order MDPs

## ***Structure by Features in the State Space***

- Dynamic Bayesian networks
- Factored MDPs

## State Space

- Often has structure
  - Not just a black box with  $n$  states, leading to a probability distribution  $T(s', s, a) = P(s' | s, a)$
  - Can be structured by relations (FO MDPs), regularities (robot navigation; harbour management)
- Generally,  $n$  different features that describe a state space
  - Encode in  $n$  individual random variables  $S_i$  with respective domains  $\text{dom}(S_i) = \{v_1, \dots, v_{d_i}\}$ 
    - State space size then describable as  $|S| = \prod_i d_i \leq d^n, d = \max_i d_i$ 
      - I.e., exponential in the number of random variables
- Given (conditional) independences between different  $S_i$ , factorisation of probability distributions in model possible
  - Applicable to MDPs, POMDPs, DecPOMDPs, partitioned DecPOMDPs
  - Most work exists for factored MDPs (also the simplest case to consider)

## Factorisation in General

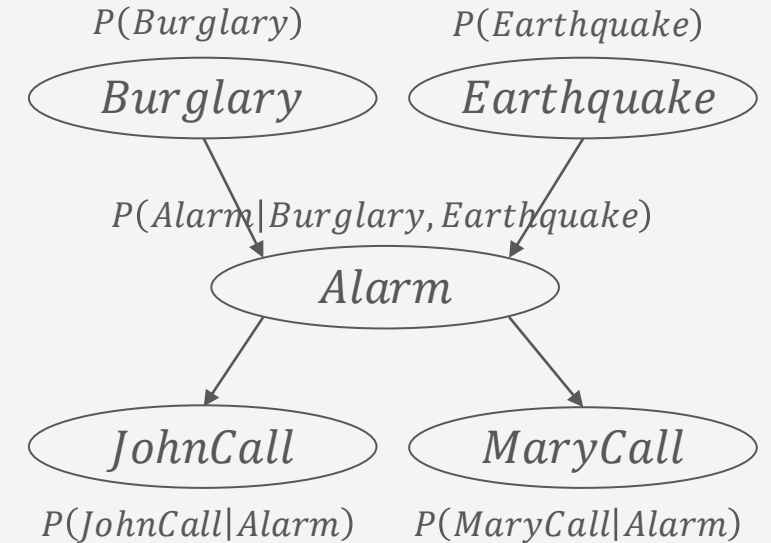
- (Conditional) independences:
  - $A \perp B$  ( $A, B$  independent)  $\Leftrightarrow P(A, B) = P(A) \cdot P(B)$
  - $A \perp B \mid C$  ( $A, B$  conditionally independent given  $C$ )  $\Leftrightarrow P(A, B \mid C) = P(A \mid C) \cdot P(B \mid C)$ 
    - Alternate version:  $A \perp B \mid C \Leftrightarrow P(A \mid B, C) = P(A \mid C)$
- (Conditional) independences allow for factorising a distribution into smaller factors
  - In general: Factorisation of a full joint probability distribution  $P(S_1, \dots, S_n)$  into  $m$  factors over subsets  $\mathcal{C}$  of random variables that form  $P(S_1, \dots, S_n)$  after multiplication (and normalisation):

$$P(S_1, \dots, S_n) = \frac{1}{Z} \prod_{j=1}^m \phi(\mathcal{C}_j)$$

- Where  $\mathcal{C}_j$  refers to sets of random variables that are mutually dependent on each other
- Memory complexity:  $\mathcal{O}(d^n)$  vs.  $\mathcal{O}(m \cdot d^{|\mathcal{C}_{max}|})$

# Probabilistic Graphical Models (PGMs)

- PGMs use a graph structure to represent dependences
- Common formalism: **Bayesian network** (BN)  $B$ 
  - Directed acyclic graph
    - Nodes: random variables  $S_i$
    - Edges: if  $S_i$  depends on  $S_j$ , edge  $S_j \rightarrow S_i$
  - Factors: conditional probability distributions (CPDs)  $\forall i P(S_i | \text{pa}(S_i))$ 
    - Roots:  $\text{pa}(S_i) = \emptyset \rightarrow$  Prior distributions  $P(S_i)$
    - Usually not depicted in graph; have to be denoted somewhere
    - Semantics:  $P(S_1, \dots, S_n) = \prod_{i=1}^n P(S_i | \text{pa}(S_i))$
- Not further considered here:
  - Undirected version with potential functions  $\phi$  as factors:
    - Factor graphs, Markov networks
    - Same semantics, different graphical representation



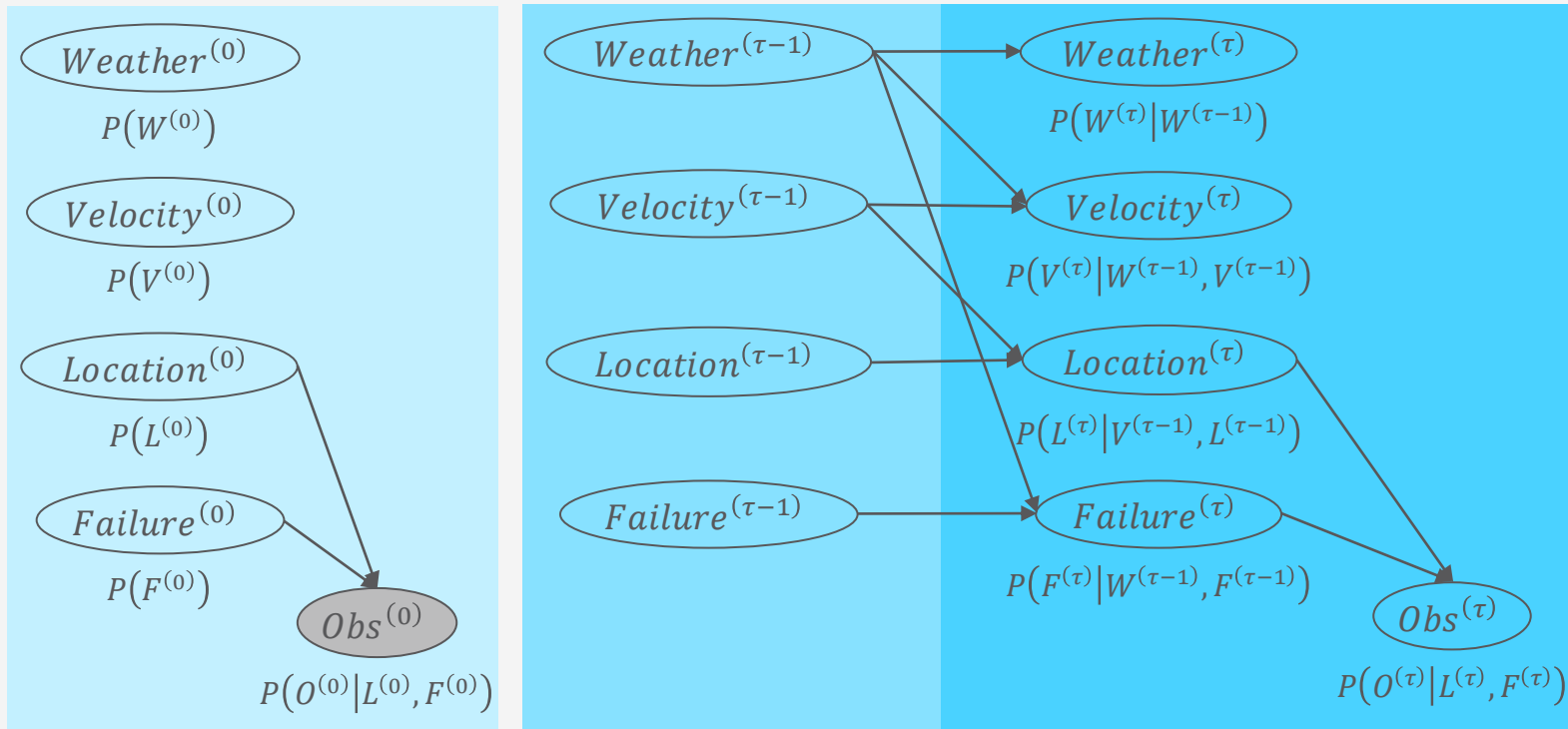
Full joint probability distribution size:  $d^5$   
 Sizes of CPDs:  $d + d + d^3 + d^2 + d^2$   
 Given  $d = 2$ :  $2^5 = 32$  vs. 20  
 (As probabilities add to 1:  
 size - 1 for each probability distribution in each CPD,  
 i.e.,  $1 + 1 + 4 + 2 + 2 = 10$ )

## Dynamic Bayesian Networks

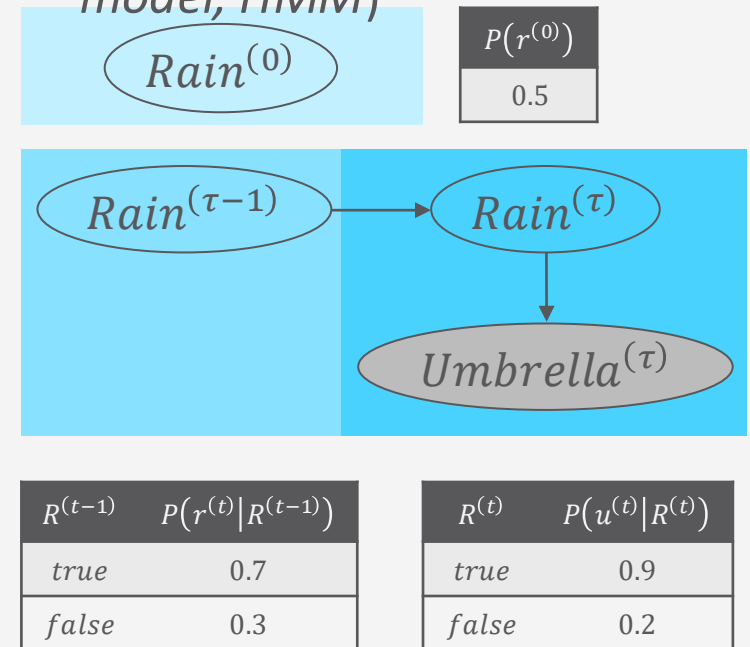
- MDP models a sequential, i.e., temporal, stationary, Markovian probabilistic setting
  - Factorisation also needs to encode a sequential, stationary, Markovian probabilistic setting
- Popular modeling formalism used:  
**Dynamic BN** (DBN) is a two-tuple  $(B^{(0)}, B^{(\rightarrow)})$ 
  - Template variables  $S_i$  indexed by time step  $\tau$  in BNs
    - Can be instantiated for particular time steps  $t$
  - BN  $B^{(0)}$  for time step 0 to encode
    - If set to uniform distributions or using DBN for fix point calculations, can be safely ignored
  - BN  $B^{(\rightarrow)}$  for time step  $\tau$  with connections from time step  $\tau - 1$  (copy pattern)
    - Markov-1 → Only connections from  $\tau - 1$  to  $\tau$
    - Stationary →  $B^{(\rightarrow)}$  identical for all  $t \in \{1, \dots\}$
  - Semantics: unroll for  $T$  time steps and multiply

# Dynamic Bayesian Networks: Example

- Left: vehicle localization task, where a moving car tries to track its current location using the data obtained from a, possibly faulty, sensor



- Right: Toy example of a special case of a DBN with one latent and one observable variable (*hidden Markov model, HMM*)



## Factored MDPs

- MDP with its state space  $S$  structured according to  $S_1, \dots, S_n$ , which in general means that
  - Transition probability distribution  $T(S', S, A) = P(S' | S, A)$  is given by
 
$$T(S'_1, \dots, S'_n, S_1, \dots, S_n, A) = P(S'_1, \dots, S'_n | S_1, \dots, S_n, A)$$
    - Or using the template notation:  $T(S^{(\tau)}, S^{(\tau-1)}, A^{(\tau-1)}) = P(S^{(\tau)} | S^{(\tau-1)}, A^{(\tau-1)})$  is given by
 
$$T(S_1^{(\tau)}, \dots, S_n^{(\tau)}, S_1^{(\tau-1)}, \dots, S_n^{(\tau-1)}, A^{(\tau-1)}) = P(S_1^{(\tau)}, \dots, S_n^{(\tau)} | S_1^{(\tau-1)}, \dots, S_n^{(\tau-1)}, A^{(\tau-1)})$$
    - Note that the overall size of  $T$  does not increase as the state space size is identical
  - Given that  $S_1, \dots, S_n$  represent features of (hopefully weakly) connected parts of a system,  $T$  can be factored according to (conditional) independences  $\rightarrow$  often represented using a DBN
  - Factorisation of  $T$ :

$$T(S', S, A) = P(S'_1, \dots, S'_n | S_1, \dots, S_n, A) = \prod_{i=1}^n P(S'_i | \text{pa}(S'_i)) =: T_B$$

## Factored MDPs: Actions and Rewards

- To be correct, the DBN just described is a standard DBN extended with random variable nodes for actions, whose assignment we want to determine, and reward nodes to denote that a reward function outputs a reward depending on the state (and action)
  - BN extended with so-called decision and utility nodes called **influence or decision diagram**

*Side note:* Since the state in MDPs is fully observable, every random variable in a DBN is observable, which is not the general case for DBNs, where usually there is a set of latent variables, which are never observed and as such often queried, and a set of evidence variables, which are usually observed (save for sensor failures).

## Factored MDPs: Actions and Rewards

- What about rewards?  
If the reward remains a function over the complete state space without any factorisation, we have not gained much
- But remember: Multi-attribute utility theory
  - Reward function with preference independence between subsets of random variables  
→ additive reward function
  - Factorisation of  $R$ :

$$R(S) = R(S_1, \dots, S_n) = \sum_{j=1}^m R_j(C_j)$$

- Best case  $R(S_1, \dots, S_n) = \sum_{i=1}^n R_i(S_i)$
- Compare factorisation of  $T$ :  $T(S', S, A) = P(S'_1, \dots, S'_n | S_1, \dots, S_n, A) = \prod_{i=1}^n P(S'_i | \text{pa}(S'_i))$

## Factored MDPs: Space Complexity

- With a structured state space, representation size down
  - Given
    - State space with  $n$  features and a maximum domain size of  $d$
    - DBN over  $n$  features and a maximum domain size of  $d$ , with  $c = \max_{i \in \{1, \dots, n\}} |\text{pa}(S_i)| + 1$
    - Given action space of size  $a$
  - Space complexity
    - Transition function  $T(S', S, A)$ :  $\mathcal{O}(d^n \cdot a)$  vs.  $\mathcal{O}(n \cdot d^c \cdot a)$
    - Reward function  $R(S)$ :  $\mathcal{O}(d^n)$  vs.  $\mathcal{O}(n \cdot d^c)$

## Solving Factored MDPs

- Bellman equation:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s' \in \text{dom}(S)} P(s'|a, s) U(s')$$

- Becomes

$$U(s_1, \dots, s_n) = \sum_{j=1}^m R_j(\mathbf{C}_j) + \gamma \max_{a \in A(s_1, \dots, s_n)} \sum_{s'_1 \in \text{dom}(S_1)} \dots \sum_{s'_n \in \text{dom}(S_n)} \prod_{i=1}^N P(s_i^{(\tau)} | \text{pa}(s_i^{(\tau)})) U(s'_1, \dots, s'_n)$$

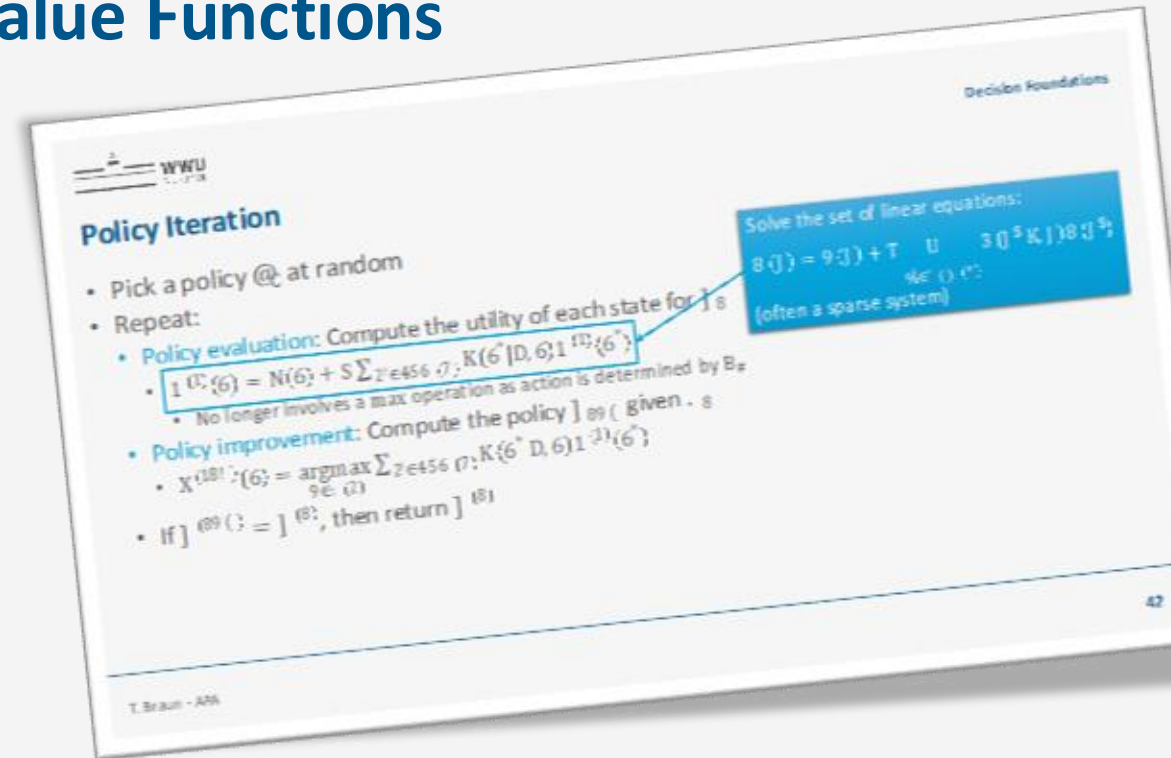
- Unfortunately, a factored MDP does not induce a factored value function  $U$ 
  - One way to go: concentrate on value functions that have a factored representation
    - Approximate the unfactored value function with a factored one

# Linear Value Functions

- **Linear value function**  $\mathcal{V}$  over a set of **basis functions**  $H = \{h_1, \dots, h_k\}$ 
  - Function  $\mathcal{V}$  that can be written as  $\mathcal{V}(s_1, \dots, s_n) = \sum_{j=1}^k w_j \cdot h_j(s_1, \dots, s_n)$  for some coefficients  $w = (w_1, \dots, w_k)'$ 
    - Let  $\mathcal{H}$  be the linear subspace of  $\mathbb{R}^n$  spanned by  $H$
    - Let  $H$  be an  $n \times k$  matrix whose columns are the  $k$  basis functions viewed as vectors
    - Then,  $\mathcal{V}$  can be written as  $Hw$
  - Equivalent expressive power to, e.g., single layer neural network
    - Features corresponding to the basis functions
    - Optimise the coefficients  $w$  to obtain a good approximation for true value function
  - Separates the problem of defining a reasonable space of features and the induced space  $\mathcal{H}$ , from the problem of searching within the space
    - Former problem is typically purview of domain experts, latter is focus of analysis + algorithmic design

# Approximate Policy Iteration with Linear Value Functions

- Restrict policy iteration algorithm to only use value functions  $\mathcal{V}$  within the provided  $\mathcal{H}$ 
  - Policy improvement as before
  - Policy evaluation changes
    - Whenever policy iteration takes a step that results in a  $\mathcal{V}$  outside of  $\mathcal{H}$ , project result back into  $\mathcal{H}$  by finding a value function within  $\mathcal{H}$  closest to  $\mathcal{V}$
- Projection operator  $\Pi$ 
  - Mapping  $\Pi : \mathbb{R}^n \rightarrow \mathcal{H}$
  - $\Pi$  is said to be a projection w.r.t. a norm  $\|\cdot\|$  if  $\Pi\mathcal{V} = Hw^*$  such that  $w^* \in \arg \min_w \|Hw - \mathcal{V}\|$
  - $\Pi$  is the linear combination of the basis functions that is closest to  $\mathcal{V}$  w.r.t. chosen norm



# Approximate Policy Iteration with Linear Value Functions

- Policy evaluation for a policy  $\pi^{(t)}$ 
  - Value function — the value of acting according to the current policy  $\pi^{(t)}$  — is approximated through a linear combination of basis functions
- Given  $\pi^{(t)}$ , i.e., actions are fixed,
  - $T(S', S, A) = T(S', S, \pi^{(t)}) = T(S', S)$
- Policy evaluation can be written in terms of matrices and vectors
  - $\mathcal{V}$  and  $R$  as  $n$ -dimensional vectors and  $T$  as an  $n \times n$ -dimensional matrix, denoted  $V, R, T$
  - Then,  $\mathcal{V} = R + \gamma T \mathcal{V}$ 
    - System of linear equations with one equation for each state  $\rightarrow$  approximate solution within  $\mathcal{H}$ :  
$$w^{(t)} = \arg \min_w \|Hw - (R + \gamma THw)\| = \arg \min_w \|(H - \gamma TH)w^{(t)} - R\|$$
    - Problem: How to choose  $\|\cdot\|$  wisely, i.e., providing error bounds?

## Approximate Policy Iteration with Linear Value Functions

- Convergence and error analysis for MDPs use max-norm ( $\mathcal{L}_\infty$ )  
→ Tie projection operator to  $\mathcal{L}_\infty$  norm
- Minimising the  $\mathcal{L}_\infty$  norm studied in optimisation literature as the problem of finding the Chebyshev solution to an overdetermined linear system of equations
  - I.e., finding  $w^*$  such that  $w^* \in \arg \min_w \|Cw - b\|_\infty$ 
    - $C = (H - \gamma TH)$ ,  $b = R$
  - Algorithm due to Stiefel (1960) solves problem by linear programming:
    - Variables:  $w_1, \dots, w_k, \phi$ ;
    - Minimise:  $\phi$ ;
    - Subject to:  $\phi \geq \sum_{j=1}^k c_{ij} \cdot w_j - b_i$  and  $\phi \geq b_i - \sum_{j=1}^k c_{ij} \cdot w_j$ ,  $i = 1, \dots, n$ .
- At solution  $(w^*, \phi^*)$ ,  $w^*$  is the Chebyshev solution and  $\phi^*$  is the  $\mathcal{L}_\infty$  projection error

Only  $k + 1$  variables but  $2n$  constraints:  
Impractical in general but in factored MDPs  
with linear value functions, constraints can  
be represented efficiently → tractable

## Factored Value Functions

- Efficient computation of value function using  $h_1, \dots, h_k$  ( $\mathbf{s} = s_1, \dots, s_n$ ) using Q value function

$$Q(\mathbf{s}, a) = R(\mathbf{s}, a) + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) \mathcal{V}(\mathbf{s}') = R(\mathbf{s}, a) + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) \sum_i w_i h_i(\mathbf{s}')$$

- Define  $G(\mathbf{s}, a)$  with  $g_i(\mathbf{s}, a) := \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) h_i(\mathbf{s}')$

$$G(\mathbf{s}, a) := \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) \sum_i w_i h_i(\mathbf{s}') = \sum_i w_i \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) h_i(\mathbf{s}') = \sum_i w_i g_i(\mathbf{s}, a)$$

- Can compute each basis function separately
- Factored (linear) value function
  - Linear function over the basis set  $h_1, \dots, h_k$  where scope of each basis function  $h_i$  restricted to some subset of variables  $\mathcal{C}_i \subset \mathcal{S}$
  - Goal: scopes of  $h_1, \dots, h_k$  correspond to cliques in graph of DBN representing transition model  $T$

## Factored Value Functions: Use in Q Value Function

- Consider  $g(\mathbf{s}, a) := \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' | \mathbf{s}, a) h(\mathbf{s}') = T_B h$ 
  - $P(\mathbf{s}' | \mathbf{s}, a)$  factored as a DBN  $T_B$
  - $h$  has restricted scope over  $\mathcal{C}$
- Sum over  $\mathcal{C}'$  conditioned on ancestors  $\mathbf{R} = \text{anc}(\mathcal{C}')$  of  $\mathcal{C}'$  in  $T_B$

$$\begin{aligned}
 g_i(\mathbf{s}, a) &= \sum_{\mathbf{s}' \in \mathcal{S}'} P(\mathbf{s}' | \mathbf{s}, a) h_i(\mathbf{s}') = \sum_{\mathbf{s}' \in \mathcal{S}'} P(\mathbf{s}' | \mathbf{s}, a) h_i(\mathbf{c}') \\
 &= \sum_{\mathbf{c}' \in \mathcal{C}'} P(\mathbf{c}' | \mathbf{s}, a) h_i(\mathbf{c}') \underbrace{\sum_{\mathbf{r}' \in \mathcal{S}' \setminus \mathcal{C}'} P(\mathbf{r}' | \mathbf{s}, a)}_{= 1} = \sum_{\mathbf{c}' \in \mathcal{C}'} P(\mathbf{c}' | \mathbf{r}, a) h_i(\mathbf{c}')
 \end{aligned}$$

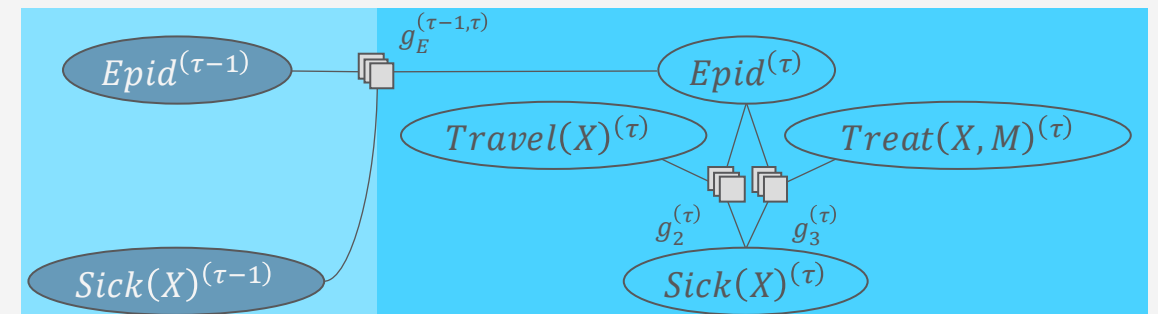
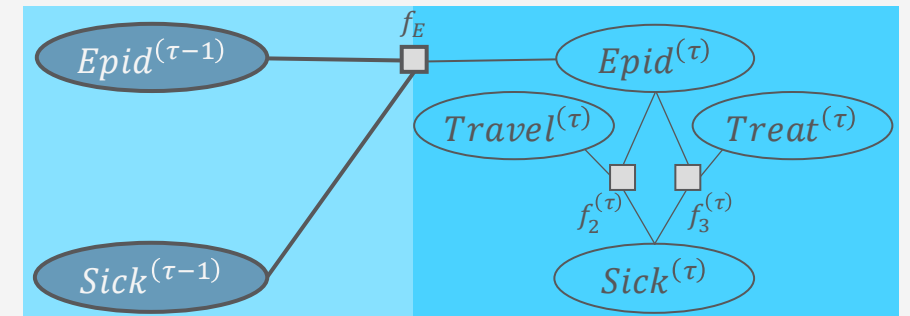
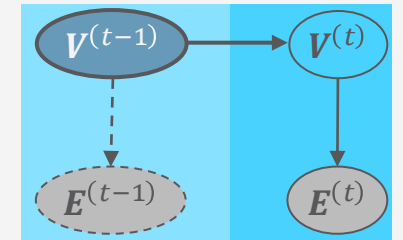
- Depends on the number of values  $\mathbf{R}$  can take, which depends on  $\mathcal{C}'$  and complexity of dynamics represented in  $T_B$ , i.e., connectivity of graph  $B$

## Factored Value Functions: Use in LP with Exponentially Many Constraints

- Constraints of form  $\phi \geq \sum_i w_i c_i(\mathbf{s}) - b(\mathbf{s}), \forall \mathbf{s} \in \mathcal{S}$ 
  - $\phi, w_1, \dots, w_k$  free variables
  - $\mathbf{s}$  ranges over all states
- Can be replaced by one equivalent non-linear constraint  $\phi \geq \max_{\mathbf{s}} \sum_i w_i c_i(\mathbf{s}) - b(\mathbf{s})$
- Tackle problem of representing non-linear constraint by
  - Computing maximum assignment for a fixed set of weights
    - Simpler problem: Given fixed weights  $w_i$ , compute  $\phi^* = \max_{\mathbf{s}} \sum_i w_i c_i(\mathbf{s}) - b(\mathbf{s})$
  - Representing non-linear constraint by small set of linear constraints using a construction called factored LP

## Relations and Symmetries in a Factored State Space

- Can add additional assumptions about relations and symmetries in the DBN
  - DBN may describe objects and relations among them
  - Objects may be considered indistinguishable
- Possible to group (and then count) objects in the transition and reward functions
  - Basis function for each group



## Interim Summary: *Structure by Features in the State Space*

- State space characterised by set of attributes
  - (Conditional) independences allow for factorisation of functions in MDP
  - Probabilistic graphical models represent such factorisations
- Factored MDP: MDP with a DBN as a representation of the transition model
  - Reduction in space complexity
  - Factored transition function does not lead to factored value function
- Factored (linear) value functions over a set of basis functions
  - Enable computing policy evaluation efficiently
- Approximate policy iteration
  - Project results outside of subspace spanned by basis functions back into subspace

## Outline: Decision Making – Structure

### *Structure by Groups in the Agent Set*

- Agent types
- Partitioned decPOMDPs

### *Structure by Relations in the State Space*

- Situation calculus
- First-order MDPs

### *Structure by Features in the State Space*

- Dynamic Bayesian networks
- Factored MDPs

⇒ Next: Human-awareness