

18<sup>th</sup> International Conference on Automated Planning and Scheduling  
September 14-18. 2008 Sydney, Australia

ICAPS-08 Tutorial on

## **First-Order Planning Techniques**

Organizers

Scott Sanner, NICTA (Australia)  
Kristian Kersting, Fraunhofer IAIS (Germany)



## ICAPS 2008 Tutorial

# Techniques for First-order Planning



**Scott Sanner**  
NICTA

## Motivation



**Kristian Kersting**  
Fraunhofer IAIS

with an invited presentation by [Saket Joshi](#), Tufts University



## Tutorial Outline

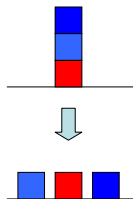
### • Motivation

- Deductive First-order Planning
  - Situation Calculus
  - Symbolic Dynamic Programming
  - Relational Bellman Algorithm (ReBeL)
  - First-order Decision Diagrams (FODDs)
- Inductive First-order Planning
- Conclusion

## 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

## Benefits of Relational Languages

- STRIPS, PDDL, GGP are relational languages...
  - Refer to relational *fluents*:
    - e.g., *Bin(?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.,  $(\text{forall } (?b ?c) ((\text{Destination } ?b ?c) \Rightarrow (\text{Bin } ?b ?c)))$
  - Are *domain-independent* and often compact!

## How to Solve?

- Relational planning *problem*:



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

- Solve *ground problem* for *each domain instance*?
  - 3 trucks: 2 planes: 3 boxes:
- Or solve lifted specification for *all domains* at once?

## Full Specification: BoxWorld

- **Relational Fluents:** *BozIn(Boz, City)*, *TruckIn(Truck, City)*, *BozOn(Boz, Truck)*
- **Goal:**  $[\exists \text{Boz} : b. \text{BozIn}(b, \text{paris})]$
- **Actions:**
  - *load(Boz : b, Truck : t):*
    - \* **Effects:**
      - when  $[\exists \text{City} : c. \text{BozIn}(b, c) \wedge \text{TruckIn}(t, c)]$  then  $[\text{BozOn}(b, t)]$
      - $\forall \text{City} : c. \text{when } [\text{BozIn}(b, c) \wedge \text{TruckIn}(t, c)]$  then  $[\neg \text{BozIn}(b, c)]$
  - *unload(Boz : b, Truck : t):*
    - \* **Effects:**
      - $\forall \text{City} : c. \text{when } [\text{BozOn}(b, t) \wedge \text{TruckIn}(t, c)]$  then  $[\text{BozIn}(b, c)]$
      - when  $[\exists \text{City} : c. \text{BozOn}(b, t) \wedge \text{TruckIn}(t, c)]$  then  $[\neg \text{BozOn}(b, t)]$
  - *drive(Truck : t, City : c):*
    - \* **Effects:**
      - when  $[\exists \text{City} : c_1. \text{TruckIn}(t, c_1)]$  then  $[\text{TruckIn}(t, c)]$
      - $\forall \text{City} : c_1. \text{when } [\text{TruckIn}(t, c_1)]$  then  $[\neg \text{TruckIn}(t, c_1)]$

## Solving Ground BoxWorld

- Apply planner to BoxWorld grounded w.r.t. domain, e.g.,
  - **Domain Object Instantiation:**
    - $Box = \{box_1, box_2, box_3\}$ ,  $Truck = \{truck_1, truck_2\}$ ,  $City = \{paris, berlin, rome\}$
  - **Ground Fluents (i.e., binary state variables):**
    - $BozIn: \{BozIn(box_1, paris), BozIn(box_2, paris), BozIn(box_3, paris), BozIn(box_1, berlin), BozIn(box_2, berlin), BozIn(box_3, berlin), BozIn(box_1, rome), BozIn(box_2, rome), BozIn(box_3, rome)\}$
    - $TruckIn: \{TruckIn(truck_1, paris), TruckIn(truck_1, berlin), TruckIn(truck_1, rome), TruckIn(truck_2, paris), TruckIn(truck_2, berlin), TruckIn(truck_2, rome)\}$
    - $BoxOn: \{BoxOn(box_1, truck_1), BoxOn(box_2, truck_1), BoxOn(box_3, truck_1), BoxOn(box_1, truck_2), BoxOn(box_2, truck_2), BoxOn(box_3, truck_2)\}$
  - **Ground Actions:**
    - $load: \{load(box_1, truck_1), load(box_2, truck_1), load(box_3, truck_1), load(box_1, truck_2), load(box_2, truck_2), load(box_3, truck_2)\}$
    - $unload: \{unload(box_1, truck_1), unload(box_2, truck_1), unload(box_3, truck_1), unload(box_1, truck_2), unload(box_2, truck_2), unload(box_3, truck_2)\}$
    - $drive: \{drive(truck_1, paris), drive(truck_1, berlin), drive(truck_1, rome), drive(truck_2, paris), drive(truck_2, berlin), drive(truck_2, rome)\}$
  - **Goal:**  $[BozIn(box_1, paris) \vee BozIn(box_2, paris) \vee BozIn(box_3, paris)]$

#states exponential in #state-vars!

Exponential #state-vars in arity

Exponential #actions in arity

Exponential in #nested quantifiers

## A First-order Solution to BoxWorld

- Derive solution deductively at lifted PDDL level:

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

Optimal for any domain instantiation!

- Great, but how do I obtain this solution?

## Tutorial Overview

- Foundational theory for exploiting first-order structure in planning
  - deterministic and probabilistic
  - representations and implementation
- The first part covers a *deductive* approach
  - plan solely based on model
  - no simulations or sampled data
    - requires grounding
- The second part reviews *inductive* approaches

## ICAPS 2008 Tutorial

### Techniques for First-order Planning

### Deterministic Planning in the Situation Calculus

Scott Sanner

NICTA

## Tutorial Outline

- Motivation
- **Deductive First-order Planning**
  - **Situation Calculus**
  - Symbolic Dynamic Programming
  - Relational Bellman Algorithm (ReBeL)
  - First-order Decision Diagrams (FODDs)
- Inductive First-order Planning
- Conclusion

## 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 b/w situations
  - e.g.,  $BoxOn(b, t, s)$ ,  $TruckIn(t, c, s)$ ,  $BoxIn(t, c, s)$

## Situation Calculus: PDDL to Effects

- Recall **BoxWorld** PDDL specification...

- $\text{load}(\text{Box} : b, \text{Truck} : t)$ :
  - Effects:
    - \*  $\text{when } [\exists \text{City} : c. \text{BoxIn}(b, c) \wedge \text{TruckIn}(t, c)] \text{ then } [\text{BoxOn}(b, t)]$
    - \*  $\forall \text{City} : c. \text{when } [\text{BoxIn}(b, c) \wedge \text{TruckIn}(t, c)] \text{ then } [\neg \text{BoxIn}(b, c)]$
- $\text{unload}(\text{Box} : b, \text{Truck} : t)$ :
  - Effects:
    - \*  $\forall \text{City} : c. \text{when } [\text{BoxOn}(b, t) \wedge \text{TruckIn}(t, c)] \text{ then } [\text{BoxIn}(b, c)]$
    - \*  $\text{when } [\exists \text{City} : c. \text{BoxOn}(b, t) \wedge \text{TruckIn}(t, c)] \text{ then } [\neg \text{BoxOn}(b, t)]$
- $\text{drive}(\text{Truck} : t, \text{City} : c)$ :
  - Effects:
    - \*  $\text{when } [\exists \text{City} : c_1. \text{TruckIn}(t, c_1)] \text{ then } [\text{TruckIn}(t, c)]$
    - \*  $\forall \text{City} : c_1. \text{when } [\text{TruckIn}(t, c_1)] \text{ then } [\neg \text{TruckIn}(t, c_1)]$

## Situation Calculus: PDDL to Effects

- Translate to **positive** and **negative** effect axioms

- $\text{load}(\text{Box} : b, \text{Truck} : t)$ :
  - Effects:
    - \*  $[\exists c. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] \supset \text{BoxOn}(b, t, \text{do}(a, s))$
    - \*  $[\exists t. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] \supset \neg \text{BoxIn}(b, c, \text{do}(a, s))$
- $\text{unload}(\text{Box} : b, \text{Truck} : t)$ :
  - Effects:
    - \*  $[\exists t. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] \supset \text{BoxIn}(b, c, \text{do}(a, s))$
    - \*  $[\exists c. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] \supset \neg \text{BoxOn}(b, t, \text{do}(a, s))$
- $\text{drive}(\text{Truck} : t, \text{City} : c)$ :
  - Effects:
    - \*  $[\exists c_1. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] \supset \text{TruckIn}(t, c, \text{do}(a, s))$
    - \*  $[\exists c. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] \supset \neg \text{TruckIn}(t, c_1, \text{do}(a, s))$

## Situation Calculus: PDDL to Effects

- Now, merge into **positive** effect axioms

$$\gamma_F^+(\vec{x}, a, s) \supset F(\vec{x}, \text{do}(a, s))$$

and **negative** effect axioms

$$\gamma_F^-(\vec{x}, a, s) \supset \neg F(\vec{x}, \text{do}(a, s))$$

- Use **rule** to combine multiple effects

$$[(C_1 \supset F) \wedge (C_2 \supset F)] \equiv [(C_1 \vee C_2) \supset F]$$

## Frame Problem

- Now we have **positive** and **negative** effects

$$\begin{aligned} \gamma_{\text{BoxIn}}^+(\vec{x}, a, s) &\supset \text{BoxIn}(\vec{x}, \text{do}(a, s)) & \gamma_{\text{BoxIn}}^-(\vec{x}, a, s) &\supset \neg \text{BoxIn}(\vec{x}, \text{do}(a, s)) \\ \gamma_{\text{TruckIn}}^+(\vec{x}, a, s) &\supset \text{TruckIn}(\vec{x}, \text{do}(a, s)) & \gamma_{\text{TruckIn}}^-(\vec{x}, a, s) &\supset \neg \text{TruckIn}(\vec{x}, \text{do}(a, s)) \\ \gamma_{\text{BoxOn}}^+(\vec{x}, a, s) &\supset \text{BoxOn}(\vec{x}, \text{do}(a, s)) & \gamma_{\text{BoxOn}}^-(\vec{x}, a, s) &\supset \neg \text{BoxOn}(\vec{x}, \text{do}(a, s)) \end{aligned}$$

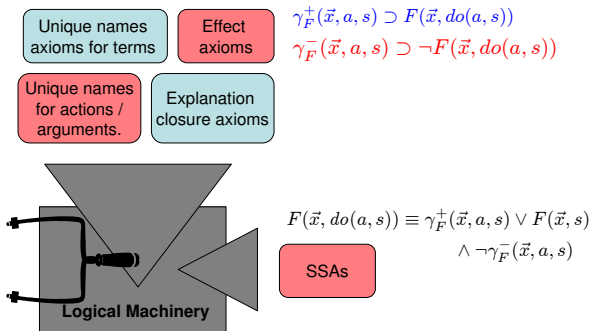
so we have compactly specified what changes.

- How to compactly specify what does not change?

- Infamous **Frame Problem**
- Intuition:
  - “what does not change, remains same”
  - this is Reiter’s **Default Solution**
  - but we have to logically formalize it...

## Successor State Axioms (SSAs)

- Default solution to frame problem given as SSAs:



## SSAs

- Shorthand:

$$\begin{aligned} F(\vec{x}, \text{do}(a, s)) &\equiv \Phi_F(\vec{x}, a, s) \\ &\equiv \gamma_F^+(\vec{x}, a, s) \vee F(\vec{x}, s) \wedge \neg \gamma_F^-(\vec{x}, a, s) \end{aligned}$$

- Reality check:

What changes and does not change!

$$\begin{aligned} \text{BoxOn}(b, t, \text{do}(a, s)) &\equiv \Phi_{\text{BoxOn}}(b, t, a, s) \\ &\equiv [\exists c. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] \\ &\quad \vee \text{BoxOn}(b, t, s) \wedge \neg [\exists c. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] \end{aligned}$$

$$\begin{aligned} \text{BoxIn}(b, c, \text{do}(a, s)) &\equiv \Phi_{\text{BoxIn}}(b, c, a, s) \\ &\equiv [\exists t. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] \\ &\quad \vee \text{BoxIn}(b, c, s) \wedge \neg [\exists t. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] \end{aligned}$$

$$\begin{aligned} \text{TruckIn}(t, c, \text{do}(a, s)) &\equiv \Phi_{\text{TruckIn}}(t, c, a, s) \\ &\equiv [\exists c_1. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] \\ &\quad \vee \text{TruckIn}(t, c, s) \wedge \neg [\exists c_1. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] \end{aligned}$$

## Regression

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

## Regression Example

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

$$\begin{aligned} &\text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unload}(b^*, t^*), s))) \\ &= \exists b. \Phi_{\text{BoxIn}}(b, \text{paris}, \text{unload}(b^*, t^*), s) \\ &= \exists b. [\exists t. \text{unload}(b^*, t^*) = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, \text{paris}, s)] \\ &\quad \vee \text{BoxIn}(b, \text{paris}, s) \\ &\quad \wedge \neg [\exists t. \text{unload}(b^*, t^*) = \text{load}(b, t) \wedge \text{BoxIn}(b, \text{paris}, s) \wedge \text{TruckIn}(t, \text{paris}, s)] \\ &= [\exists b, t. b = b^* \wedge t = t^* \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, \text{paris}, s)] \\ &\quad \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \quad \text{make non-empty domain assumption} \\ &= [(\exists b. b = b^*) \wedge (\exists t. t = t^*) \wedge \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \\ &\quad \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \quad \text{note free vars } b^*, t^*; \text{ why?} \\ &= [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \end{aligned}$$

## Regression Example

- But what action instantiation of  $\text{unload}(b^*, t^*)$  leads to:
 
$$\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unload}(b^*, t^*), s))$$
  - Just have to existentially quantify  $b^*, t^*$ 
    - Can obtain instances via query extraction w.r.t. state KB
- $$\begin{aligned} &\exists b^*, t^*. \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unload}(b^*, t^*), s))) \\ &= \exists b^*, t^*. [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \\ &\quad \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \\ &= [\exists b^*, t^*. \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \\ &\quad \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \end{aligned}$$

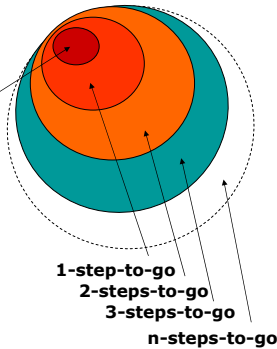
First-order state & action abstraction!  
Don't have to enumerate all states,  $b^*, t^*$ !

## Recap

- We translated PDDL to SitCalc theory
  - converted PDDL effects to SitCalc effect axioms
  - derived SSAs from effect axioms
    - using default solution to Frame Problem
- Introduced regression operator
  - extracted action instantiation to achieve goal
- Let the planning begin...

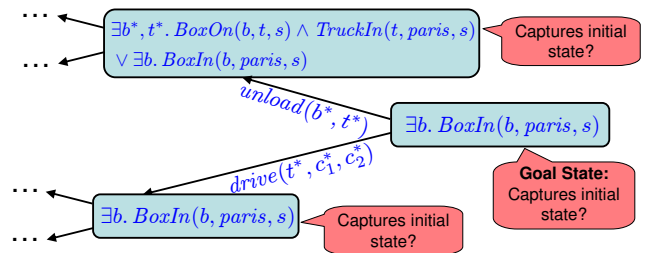
## Regression Planning

- Define abstract goal state, e.g.,
 
$$\exists b. \text{BoxIn}(b, \text{paris}, s)$$
- Check if regression through action sequence holds in initial state



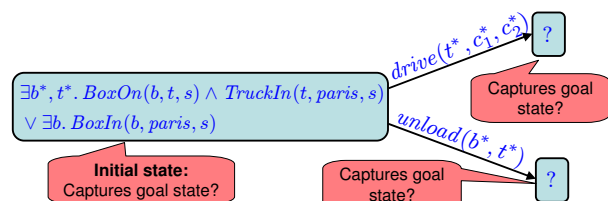
## First-order Goal-regression

- We can now do goal regression planning!
  - Provide initial state and sequence of actions
  - Use regression,  $\exists$  to tell whether goal will hold



## Progression and Forward-search?

- Can we do lifted *forward-search planning*?



- *Progression* not first-order definable! (Reiter, 01)
- Could progress *ground* state
  - But this does not exploit first-order structure

## Golog: Restricted Plan Search

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

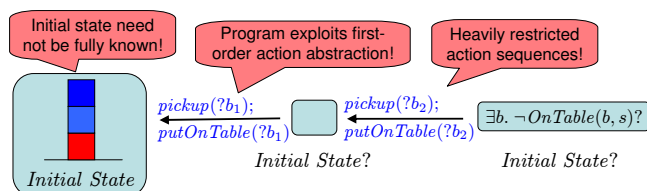
$\alpha$	primitive action
$\phi?$	condition test
$(\delta_1, \delta_2)$	sequence
if $\phi$ then $\delta_2$ endIf	conditional
while $\phi$ then $\delta$ endWhile	loop
$(\delta_1   \delta_2)$	nondeterministic choice of action
$\pi \bar{x} [\delta]$	nondeterministic choice of arguments
$\delta^*$	nondeterministic iteration
proc $\beta(\bar{x}) \delta$ endProc	procedure call definition
$\beta(\bar{t})$	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:



## For Further Reading

- Knowledge in Action:**  
In-depth coverage of SitCalc default solution, applications (Reiter, 2001)
- Golog**  
(Levesque, Reiter, Lesperance, Lin, Journal Logic Programming, 1997)
- Extensions**
  - ConGolog: concurrent Golog (de Giacomo, Lesperance, Levesque, AIJ-00)
  - DT-Golog: decision-theoretic, covered next (Soutchanski, Boutilier, Reiter, Thrun, AAAI-20)

For MDPs, covered next.

## Conclusion

- Situation Calculus**
  - First-order specification of action theory
  - Default solution addresses Frame Problem
    - Effective approach to PDDL-expressive planning
- Supports Regression Planning**
  - Initial state need not be fully specified
  - Can restrict action space with Golog program
  - Exploits state & action abstraction
    - Avoids enumerating all state & action instances!

## ICAPS 2008 Tutorial

### Techniques for First-order Planning

### FOMDPs and Symbolic Dynamic Programming

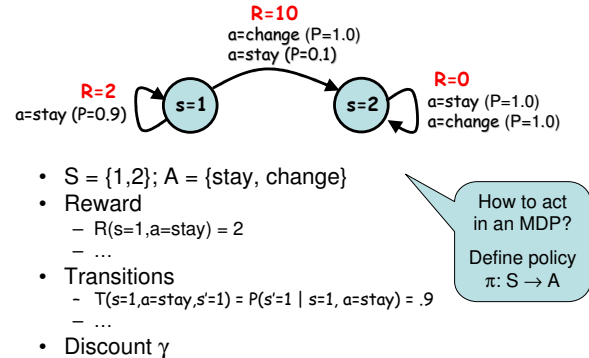
[Scott Sanner](#)

NICTA

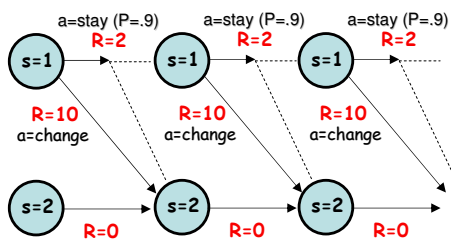
## Tutorial Outline

- **Motivation**
- **Deductive First-order Planning**
  - Situation Calculus
  - **Symbolic Dynamic Programming**
  - Relational Bellman Algorithm (ReBeL)
  - First-order Decision Diagrams (FODDs)
- Inductive First-order Planning
- Conclusion

## MDPs $\langle S, A, T, R, \gamma \rangle$

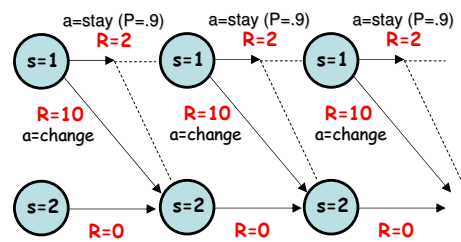


## What's the best Policy?



- Immediate vs. long-term gain?

## What's the best Policy?



- Must define reward criterion to optimize!
  - Discount factor  $\gamma$  important ( $\gamma=.9$  vs.  $\gamma=.1$ )

## MDP Policy, Value, & Solution

- Define value of a policy  $\pi$ :
 
$$V_{\pi}(s) = E_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t \cdot r_t \mid s = s_0 \right]$$
- Tells how much value you expect to get by following  $\pi$  starting from state  $s$
- MDP Optimal Solution:
  - Find optimal policy  $\pi^*$  that maximizes value
  - Fortunately:  $\exists \pi^*. \forall s, \pi. V_{\pi^*}(s) \geq V_{\pi}(s)$

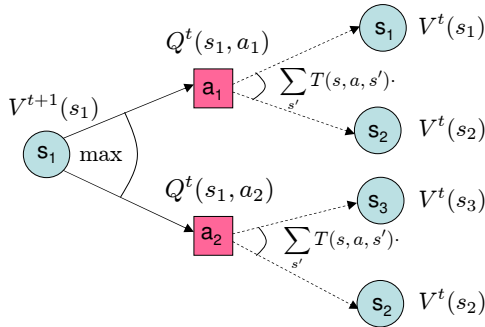
## Value Iteration: from finite to $\infty$ decisions

- Given optimal  $t-1$ -stage-to-go value function
- How to act optimally with  $t$  decisions?
  - Take action  $a$  then act so as to achieve  $V^{t-1}$  thereafter
 
$$Q^t(s, a) := R(s, a) + \gamma \cdot \sum_{s' \in S} T(s, a, s') \cdot V^{t-1}(s')$$
  - What is expected value of best action  $a$  at decision stage  $t$ ?
 
$$V^t(s) := \max_{a \in A} \{Q^t(s, a)\}$$
  - At  $\infty$  horizon, get same value ( $=V^*$ )
 
$$\lim_{t \rightarrow \infty} \max_s |V^t(s) - V^{t-1}(s)| = 0$$
    - $\pi^*$  says act same at each decision stage for  $\infty$  horizon!

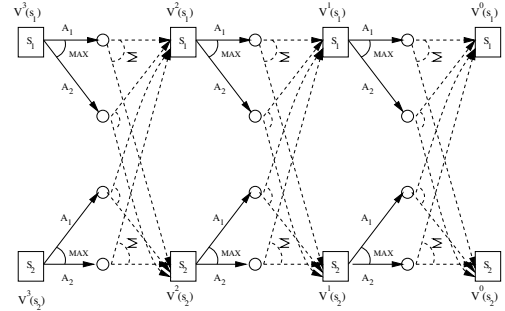


## Single Dynamic Programming Step

- Graphical view:



## Synchronous DP Updates (Value Iteration)



## Value Function $\rightarrow$ Policy

- Can derive policy from value function  $V$
- Given arbitrary value  $V$  (optimal or not)...
  - A *greedy policy*  $\pi_V$  takes action in each state that maximizes expected value w.r.t.  $V$ :

$$\pi_V(s) = \arg \max_a \left\{ R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s') \right\}$$

- If can act so as to obtain  $V$  after doing action  $a$  in state  $s$ ,  $\pi_V$  guarantees  $V(s)$  in expectation

## How to Specify & Solve “First-order MDPs”?

Following:  
[Boutilier, Reiter, Price, IJCAI-01]

## First-order (FO)MDPs: Case Statement

- $\langle S, A, T, R \rangle$  for FOMDPs defined in terms of *cases*
  - E.g., express *reward* in *BoxWorld* FOMDP as...

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

- Operators:** Define unary, binary case operations
  - E.g., can take “cross-sum”  $\oplus$  (or  $\otimes$ ,  $\ominus$ ) of cases...

$$\begin{array}{|c|c|} \hline \phi & 10 \\ \hline \neg \phi & 20 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \phi & 3 \\ \hline \neg \phi & 4 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \phi \wedge \phi & 13 \\ \hline \phi \wedge \neg \phi & 14 \\ \hline \neg \phi \wedge \phi & 23 \\ \hline \neg \phi \wedge \neg \phi & 24 \\ \hline \end{array}$$

## Stochastic Actions & FODTR

- Stochastic actions using deterministic SitCalc:
  - User’s stochastic action:  $A(x) = load(b, t)$
  - Nature’s choice:  $n(x) \in \{loadS(b, t), loadF(b, t)\}$
  - Probability distribution over Nature’s choice:

$$P(loadS(b, t) \mid load(b, t)) = \begin{array}{|l|l|} \hline snow(s) & .1 \\ \hline \neg snow(s) & .6 \\ \hline \end{array}$$

$$P(loadF(b, t) \mid load(b, t)) = \begin{array}{|l|l|} \hline snow(s) & .9 \\ \hline \neg snow(s) & .4 \\ \hline \end{array}$$

- First-order decision-theoretic regression
  - FODTR = *expectation* of regression:
 
$$FODTR[vCase(s), A(x)] = E_{P(n(x)|A(x))} [Regr[vCase(s), n(x)]]$$

## Q-functions and Backups

- FODTR almost gives us a Q-function

$$FODTR[vCase(unload(b,t))] = \begin{array}{c|c} On(b,t,s) & 5 \\ \hline \neg On(b,t,s) & 0 \end{array}$$

- FODTR specific to action variables
- Also need to add reward, discount

- Specify a backup operator for this

$$B^{unload}[vCase(s)] = rCase(s) \oplus \gamma \begin{array}{c|c} \exists b,t. On(b,t,s) & 5 \\ \hline \exists b,t. \neg On(b,t,s) & 0 \end{array}$$

- Yields a first-order Q-function

## Symbolic Dynamic Programming

- What value if 0-stages-to-go?

– Obviously  $V^0(s) = rCase(s)$

- What value if 1-stage-to-go?

– We know value for each action

$$V^1(s) = \begin{array}{c|c} \phi_1 & 9 \\ \hline \text{else } \phi_2 & 3 \\ \hline \text{else } \phi_3 & 0 \\ \hline \text{else } \phi_4 & 0 \end{array} \begin{array}{c} \swarrow \\ \searrow \\ \swarrow \\ \searrow \end{array} \begin{array}{c|c} \phi_1 & 9 \\ \hline \phi_2 & 0 \\ \hline \phi_3 & 3 \\ \hline \phi_4 & 1 \end{array} = B^{A1}[rCase(s)]$$

$$= B^{A2}[rCase(s)]$$

– Now just need max for every state

- Value iteration: (BoutReiPr, IJCAI-01)

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

## First-order ADDs

[Sanner, Thesis]

- Want to compactly represent:

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

- Push down quantifiers, expose prop. structure:

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

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

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

- Convert to first-order ADD

$$case = \begin{array}{c} a \\ \swarrow \quad \searrow \\ 1 \quad \quad 0 \end{array} \begin{array}{c} \text{First-order CSI!} \\ \swarrow \quad \searrow \\ 1 \quad \quad 0 \end{array}$$

## Results for SDP with FOADDs

- Replace case with FO(A)ADDs, e.g. *BoxWorld*

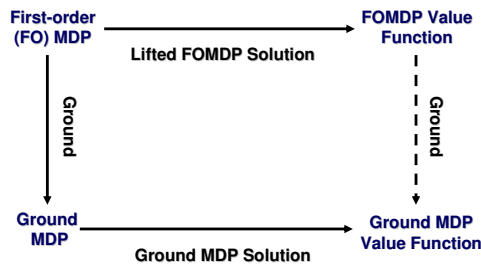
$$rCase(s) = \begin{array}{c|c} \exists b. BIn(b, Paris, s) & \\ \hline & 10 \quad 0 \end{array}$$

- Use FO(A)ADD ops for structured SDP (using  $\gamma=0.9$ )...

$$vCase(s) = \begin{array}{c|c} \exists b. BIn(b, Paris, s) & \\ \hline & 100 : noop \quad \exists b, t. TIn(t, Paris, s) \wedge On(b, t, s) \\ \hline & 89 : unload(b, t) \quad \exists b, t. On(b, t, s) \\ \hline & 80 : drive(t, Paris) \quad \exists b, c. BoxIn(b, c, s) \wedge \exists t. TIn(t, c, s) \\ \hline & 72 : load(b, t) \quad \dots \end{array}$$

## Correctness of SDP

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



## Related Purely Deductive Approaches

- Value Iteration:

– **ReBel algorithm**

(Kersting, van Otterlo, de Raedt, IJML-04)

– **FOVIA algorithm for fluent calculus**

(Karabaev & Skvortsova, UAI-05)

– **First-order decision diagrams (FODDs)**

(Wang, Joshi, Khaldon, IJCAI-07; JK, ICAPS-08; WJK, JAIR-08)

- Approximate Linear Programming (ALP)

– **First-order ALP (FOALP)**

(Sanner & Boutilier, UAI-05)

- Policy Iteration

– **Approximate policy iteration (FOAPI)**

(Sanner & Boutilier, UAI-06)

– **Modified policy iteration with FODDs**

(Wang & Joshi, UAI-07)

- Factored FOMDPs – FOMDP extension

– **Factored SDP and Factored FOALP**

(Sanner & Boutilier, ICAPS-07)

Kristian covers this.

Saket covers this.

3rd place in ICAPS IPPC5 (after FPG, FF-Replan)

## Conclusions

- MDP: model of decision-theoretic planning
  - Common solution is dynamic programming
- “FOMDPs” are one model for lifted decision-theoretic planning
  - Use SitCalc specified action theory
  - Use case to represent reward, probabilities
  - Symbolic dynamic programming = lifted DP
  - State & action abstraction for MDPs & DP

ICAPS 2008 Tutorial

**Techniques for  
First-order Planning**

**Relational Bellman  
Algorithm (ReBeL)**

[Kristian Kersting](#)

Fraunhofer IAIS

• Thanks to Prasad Tadepalli, Alan Fern, Kurt Driessens, Martijn Van Otterlo,

## Tutorial Outline

- **Motivation**
- **Deductive First-order Planning**
  - Situation Calculus
  - Symbolic Dynamic Programming
  - **Relational Bellman Algorithm (ReBeL)**
  - First-order Decision Diagrams (FODDs)
- Inductive First-order Planning
- Conclusion

## ReBeL

[Kersting, Van Otterlo, De Raedt ICML04]

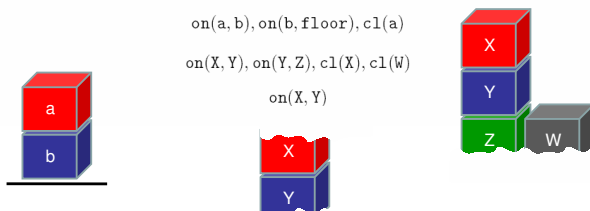
A **relational Bellman** algorithm

- Instance of SDP
- Sacrifices expressivity/compactness for „simplicity“
- Abstract state = existentially quantified conjunction of atoms (logical query) with equality constraints
- „Direct treatment“ of probabilistic actions
- Basic data structure = decision lists

## ReBeL's Abstract States

**Definition 2** An *abstract state* is a conjunction  $Z$  of logical atoms, i.e., a logical query.

Abstract states represent sets of states. More formally, a state is an interpretation, i.e. a set of grounds facts.



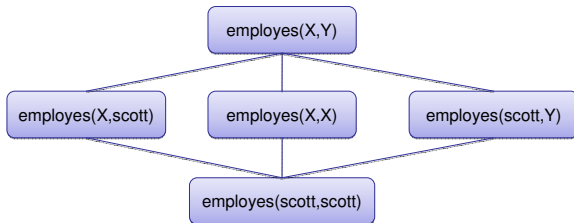
## Aside: Subsumption

- Recall that a clause like  
**C: grandparent(X,Z):- parent(X,Y), parent(Y,Z)**  
 is **C: {grandparent(X,Z), ~parent(X,Y), ~parent(Y,Z)}**

Clause **C1** **subsume** **C2** iff there exists a substitution  $\theta$  s.t.  $\theta C1$  is a **subset of** **C2**

- Thus for the following pair of clauses, **C1** subsumes **C2**:  
**C1: mem(A,[B|C]):- mem(A,C).**  
**C2: mem(0,[1,0]):- nat(0), nat(1), mem(0,[0]).**
- **Note that C1 "looks" more general.**

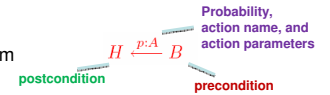
## Aside: Subsumption



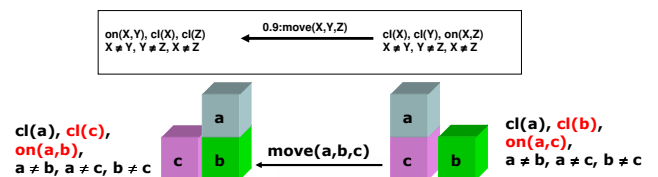
- Subsumption induces a „generality“ lattice

## ReBeL's Abstract Actions

An action is an expression of the form



[Semantics]: If current state  $b$  subsumed by  $B$ , i.e.,  $b \preceq_{\theta} B$ , then taking action  $A$  will result in  $[b \setminus B\theta] \cup H\theta$  with probability  $p$ , and  $b$  remains unchanged otherwise.

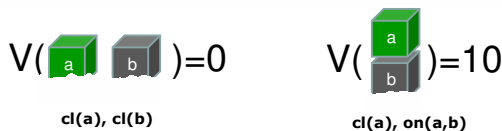


## ReBeL's Reward Model

**Definition 3** An *abstract state value function*  $V_t$  is a finite list of value rules of the form  $v(c) \leftarrow B$  where  $B$  is an abstract state and  $c \in \mathbb{R}$  is the value of all states subsumed by  $B$ .

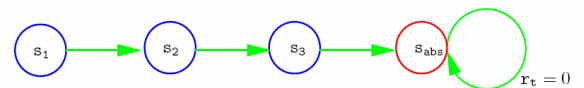
$v(10) \leftarrow \text{on}(a, b).$   
 $v(0) \leftarrow \text{true}.$

For abstract state  $Z$ ,  $V_t(Z) = \max\{c \mid v(c) \leftarrow B \in V_t \text{ and } Z \preceq_{\theta} B\}$ .



State-action values are modeled by Q-rules:  $q(c) \leftarrow B, A$

## Absorbing States



- In RL, episodic tasks can be encoded using absorbing states
  - transition only to themselves
  - generate zero reward
- In ReBeL, we encode absorbing states using **artificial** deterministic functions

$\text{on}(a, b) \xleftarrow{1.0:\text{absorbing}} \text{on}(a, b).$

## Integrity Constraints ...

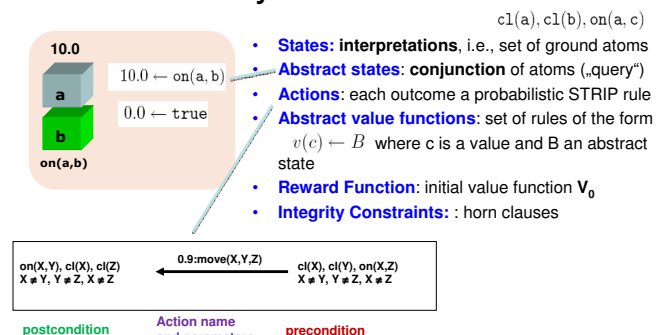
... are „simply“ Horn clauses

$\text{false} \leftarrow \text{on}(X, Y), \text{cl}(Y)$   
 $X \neq Y \leftarrow \text{on}(X, Y)$

Note that actions are constraint, too. They cannot yield illegal states.

- We used Frühwirth's *Constraint Handling Rules* and Buntine's *Generalized Subsumption* (i.e., subsumption w.r.t. to a background theory)

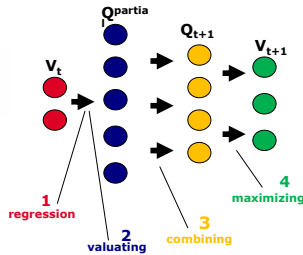
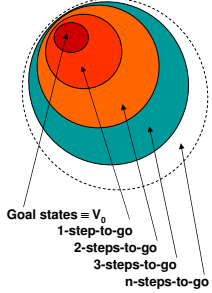
## Summary of ReBeL's MDPs



**Theorem:**  
Every ReBeL MDP induces a (possibly infinite) MDP

## Recap: SDP

Given an initial, abstract reward model  $R(\equiv V_0)$ ,  
Compute the  $t$ -steps-to-go abstract state value functions  $V_t$ ,  $t = 1, 2, \dots$



### Value Iteration

repeat for each  $t$   
for each state do

$$V_{t+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_t(s')] \quad \text{until 'convergence'}$$



## Step 1: Regression

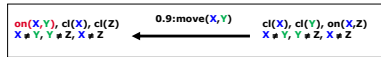


- Abstract actions define how states change
- In order to do the update from state Z, we need to consider all states that can reach Z by applying the action
- Intuition: 'Revers' the action effects
- Compute the weakest precondition of Z given the action

on(a,b)



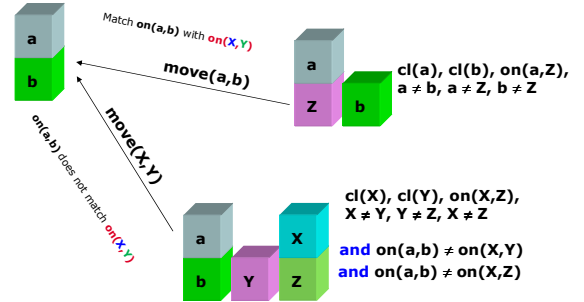
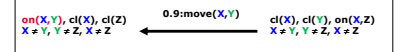
## Step 1: Regression



Two cases:

- move caused on(a,b):** We have been in abstract state  $S_1 \equiv (cl(a), cl(b), on(a, Z), a \neq b, b \neq Z, a \neq Z)$
- move did not cause on(a,b):** We moved X on Y but not a on b. So, we have been in abstract state  $T \equiv (cl(X), cl(Y), on(X, Z), on(a, b), X \neq Y, Y \neq Z, X \neq Z)$  satisfying  $c \equiv (on(X, Y) \neq on(a, b) \wedge on(X, Z) \neq on(a, b)) \equiv (X \neq a \vee Y \neq b \vee Z \neq b)$  which guarantees that applying move(X,Y) will not affect on(a,b). Simplifies to  $S_2 \equiv (T \wedge X \neq a)$ ,  $S_3 \equiv (T \wedge Y \neq b)$  and  $S_4 \equiv (T \wedge Z \neq b)$ .

## Step 1: Regression



## Step 1: Regression

- on(a,b) is simple effect. In general, we have multiple effects, i.e., conjunctions of atoms.

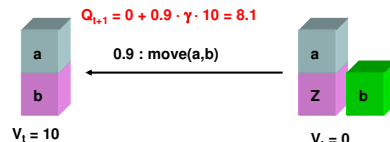
If  $Post \stackrel{A}{\Leftarrow} Pre$  and the 'next' state is Z  
we take all possible overlapping subsets of Post and Z  
apply the resulting substitution  $\theta$  on Post and Pre  
add 'unexplained' effects in Z to Pre $\theta$   
calculate and add constraints

on(a, b), on(c, d) : on(a, b) was the effect and on(c, d) was already true  
on(a, b), on(c, d) : on(c, d) was the effect and on(a, b) was already true  
on(a, b), on(c, d) : on(X, Y) was the effect and on(a, b), on(c, d) were already true  
and constraints make sure that on(X, Y) is not on(a, b) nor on(c, d).

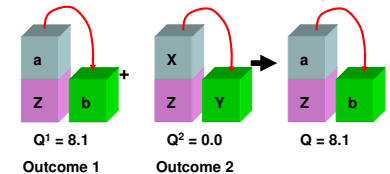
## Steps 2&3: Valuation & Combination ...

... that means computing the Q rules

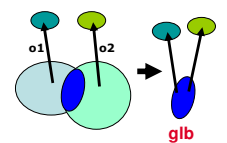
From regression, for each state  $S'$  in  $V_t$  we obtain  $(S, A)$ -pairs such that doing A in S results in  $S'$ .



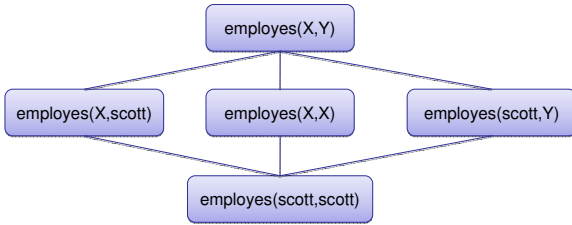
Step 2:  
Computing values for a single outcome



Step 3:  
Combining outcomes



### Aside: Greatest Lower Bound (GLB)



- Recall that subsumption induces a lattice (over the reduced clauses). Hence, for each node, we can compute the greatest lower bound.
- GLB: **set of ground states, in which both abstract states hold.**

### Q-Rules after first „Iteration“

$q(10) \leftarrow \text{on}(a, b), a \neq b, \text{absorb.} \quad \langle 1 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, X \neq a, a \neq b, \text{move}(X, Y). \quad \langle 2 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, Y \neq b, a \neq b, \text{move}(X, Y). \quad \langle 3 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, a \neq b, Z \neq b, \text{move}(X, Y). \quad \langle 4 \rangle$   
 $q(9) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), a \neq b, b \neq Z, a \neq Z, \text{move}(a, b) \quad \langle 5 \rangle$   
 $q(0) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), X \neq Y, Y \neq Z, X \neq Z, \text{move}(X, Y). \quad \langle 6 \rangle$

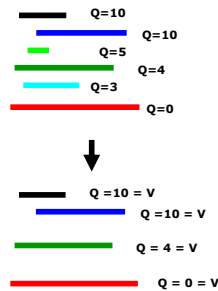
### Step 4: Maximizing ...

... that means computing the next value function by maximizing the Q-rules

$$V_{t+1}(s) = \max_a Q_{t+1}(s, a)$$

**10 : move(X,Y) ← cl(X), cl(Y)**  
**5 : move(a,b) ← cl(a), cl(b)**

- Sort the Q-rules in descending order
- Remove top  $q(d) \leftarrow B, A$  from the rules
- if no other more general rule
- then add  $v(d) \leftarrow B$  to  $V_{t+1}$  and remove all Qrules that are subsumed
- continue until no more rules



### Step 4: Maximizing ...

$q(10) \leftarrow \text{on}(a, b), a \neq b, \text{absorb.} \quad \langle 1 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, X \neq a, a \neq b, \text{move}(X, Y). \quad \langle 2 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, Y \neq b, a \neq b, \text{move}(X, Y). \quad \langle 3 \rangle$   
 $q(10) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), \text{on}(a, b), X \neq Y, Y \neq Z, X \neq Z, a \neq b, Z \neq b, \text{move}(X, Y). \quad \langle 4 \rangle$   
 $q(9) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), a \neq b, b \neq Z, a \neq Z, \text{move}(a, b) \quad \langle 5 \rangle$   
 $q(0) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), X \neq Y, Y \neq Z, X \neq Z, \text{move}(X, Y). \quad \langle 6 \rangle$

### Step 4: Maximizing ...

$q(9) \leftarrow \text{cl}(a), \text{cl}(b), \text{on}(a, Z), a \neq b, b \neq Z, a \neq Z, \text{move}(a, b) \quad \langle 5 \rangle$   
 $q(0) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), X \neq Y, Y \neq Z, X \neq Z, \text{move}(X, Y). \quad \langle 6 \rangle$

### Step 4: Maximizing ...

$q(0) \leftarrow \text{cl}(X), \text{cl}(Y), \text{on}(X, Z), X \neq Y, Y \neq Z, X \neq Z, \text{move}(X, Y). \quad \langle 6 \rangle$

Finally, we get ...

$$V_{t+1}$$

$$\begin{aligned} v(10) &\leftarrow \text{on}(a, b), a \neq b. \\ v(9) &\leftarrow \text{on}(a, Y), \text{cl}(b), \text{cl}(a), b \neq Y, a \neq Y, a \neq b. \\ v(0) &\leftarrow \text{on}(X, Y), \text{cl}(X), \text{cl}(Z), Z \neq Y, X \neq Y, X \neq Z. \end{aligned}$$

## Blocks World

Predicates:  $\text{on}/2$  and  $\text{cl}/1$

$$\begin{aligned} \text{on}(X, Y), \text{cl}(X), \text{cl}(Z), &\xleftarrow{1.0:\text{move}(X,Y)} \text{cl}(Y), \text{cl}(X), \text{on}(X, Z), \\ X \neq Y, Y \neq Z, X \neq Z & X \neq Y, Y \neq Z, X \neq Z \end{aligned}$$

Note that the number of blocks is unspecified

Value function,  $V_0$  e.g.:

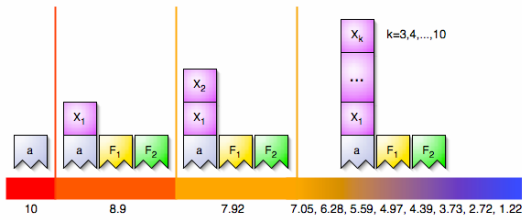
$$\begin{aligned} v(10) &\leftarrow \text{cl}(a). \\ v(0) &\leftarrow \text{true}. \end{aligned}$$

IC & sample constraint

$$\begin{aligned} \text{cl}(a) &\xleftarrow{1.0:\text{absorbing}} \text{cl}(a). \\ X \neq Y &\leftarrow \text{on}(X, Y). \end{aligned}$$

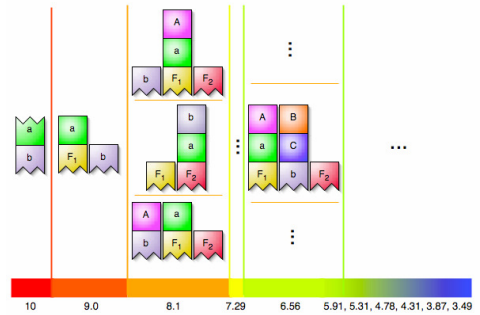
In the experiments, the discount  $\gamma = 0.9$ .

## Blocks World: $\text{cl}(a)$



**Observation 1** Abstraction does not guarantee convergence in infinite domains because an infinite number of abstract states can be required.

## Blocks World: $\text{on}(a,b)$



## ReBel: Logistics Domain



$$\begin{aligned} \text{bin}(B, C), \text{tin}(T, C), R &\xleftarrow{pr:\text{unload}(B,T)} \text{on}(B, T), \text{tin}(T, C), R \\ \text{on}(B, T), \text{tin}(T, C), R &\xleftarrow{pr:\text{load}(B,T)} \text{bin}(B, C), \text{tin}(T, C), R \\ \text{tin}(T, C'), C \neq C' &\xleftarrow{1.0:\text{drive}(T,C')} \text{tin}(T, C) \end{aligned}$$

( $pr$  is 0.9 if  $R$  is rain and 0.7 if  $R$  is not\_rain)

Reward model

$$\begin{aligned} v(10) &\leftarrow \text{bin}(b, p) \\ v(0) &\leftarrow \text{true}. \end{aligned}$$

IC and absorbing

$$\begin{aligned} \text{false} &\leftarrow \text{rain}, \text{not\_rain} \\ \text{bin}(b, p) &\xleftarrow{1:\text{absorbing}} \text{bin}(b, p) \end{aligned}$$

$\gamma = 0.9$

## ReBel: Logistics Domain



abstract states	1	2	3	4	5	6	7	8	9	10
bin(b, p).	10.000	10.000	10.000	10.000	10.000	10.000	10.000	10.000	10.000	10.000
tin(a, p), on(b, A), not_rain.	8.100	8.820	8.895	8.901	8.901	8.901	8.901	8.901	8.901	8.901
tin(a, p), on(b, A), rain.	8.800	8.001	8.460	8.584	8.618	8.627	8.629	8.630	8.630	8.630
tin(a, B), on(b, A), not_rain.	7.290	7.946	8.005	8.010	8.011	8.011	8.011	8.011	8.011	8.011
tin(a, B), on(b, A), rain.	5.670	7.201	7.614	7.726	7.756	7.764	7.766	7.767	7.767	7.767
tin(a, B), bin(b, B), not_rain.		5.905	6.968	7.111	7.128	7.130	7.131	7.131	7.131	7.131
tin(a, B), bin(b, B), rain.		3.572	5.501	6.282	6.563	6.658	6.689	6.699	6.702	6.702
tin(a, B), bin(b, C), not_rain.			5.314	6.271	6.400	6.416	6.417	6.418	6.418	6.418
tin(a, B), bin(b, C), rain.			3.215	4.951	5.854	5.907	5.903	6.020	6.029	6.029
tin(a, B).	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Convergence on both structural and value level.

(In about 2 min.)

## Conclusions

- ReBeL is an instance of SDP
  - Avoids the full state and action enumeration of classical approaches
  - Lifted solution applies to any (ground) instance
  - Basic tool: Constraint-Logic Programming
- Sacrifices expressivity/compactness for „simplicity“
- Employs constraint logic programming
- Background Knowledge is not a feature, but a necessity
- Convergence: Structural and Value level

ICAPS 2008 Tutorial

## Techniques for First-order Planning

## First-order Decision Diagrams (FODDs)

[Saket Joshi](#)

Tufts University

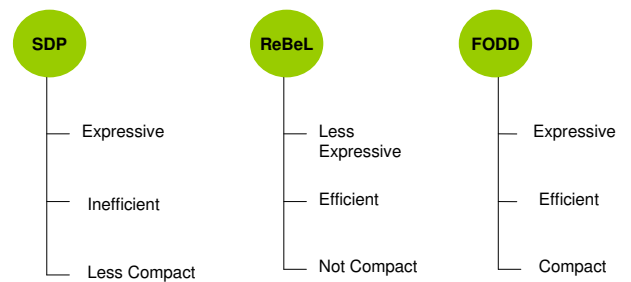
• Thanks to Chenggang Wang and Roni Khardon

## Tutorial Outline

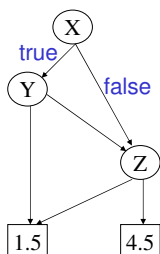
- **Motivation**
- **Deductive First-order Planning**
  - Situation Calculus
  - Symbolic Dynamic Programming
  - Relational Bellman Algorithm (ReBeL)
  - **First-order Decision Diagrams (FODDs)**
- Inductive First-order Planning
- Conclusion

[Wang, Joshi, Khardon, IJCAI 07, JAIR08]

## Motivation for FODDs

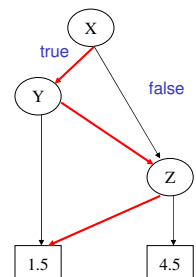


## Propositional ADDs - Syntax



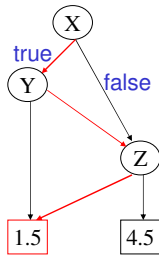
## Propositional ADDs - Semantics

- Variable valuation
  - {X = true, Y = false, Z = true}
- Function output: 1.5





## Propositional ADDs – Normal Form



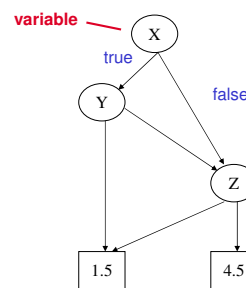
Propositional ADD

For a given variable ordering, every function has a unique representation.

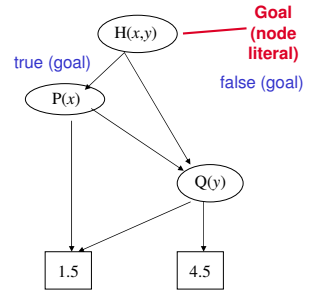
But, how do we deal with logical atoms ?

## First-Order Decision Diagrams

**Idea:** nodes are existentially quantified atoms (goals)



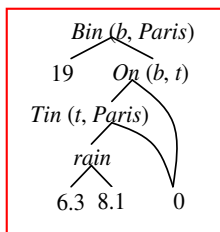
•Propositional ADD



•First-order decision diagram

## First-Order Decision Diagrams (FODDs) – Semantics based on Single Paths

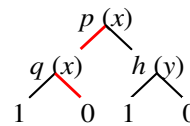
Same as for decision trees and their relational variants such as TILDE [Blockeel, De Raedt 98].



1.  $\exists b, \text{Bin}(b, \text{Paris}) : 19$
2.  $\neg \exists b, \text{Bin}(b, \text{Paris}) \wedge \exists b, t, \text{On}(b, t) \wedge \text{Tin}(t, \text{Paris}) \wedge \text{rain} : 6.3$
3.  $\neg \exists b, \text{Bin}(b, \text{Paris}) \wedge \exists b, t, \text{On}(b, t) \wedge \text{Tin}(t, \text{Paris}) \wedge \neg \text{rain} : 8.1$
4. Otherwise : 0

## First Order Decision Diagrams (FODDs) – Semantics based on Multiple Paths

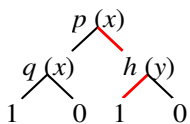
Semantics defined in terms of variable valuations  $\text{MAP}_B(I, \zeta)$



- Domain: {1, 2, 3}
- Interpretation I: {p(1), q(2), h(3)}
- $\zeta_1 = \{x/1, y/1\}$
- $\text{MAP}_B(I, \zeta_1) = 0$

## First Order Decision Diagrams (FODDs) – Semantics based on Multiple Paths

Semantics defined in terms of variable valuations  $\text{MAP}_B(I, \zeta)$



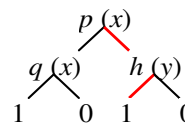
- Domain: {1, 2, 3}
- Interpretation I: {p(1), q(2), h(3)}
- $\zeta_1 = \{x/2, y/3\}$
- $\text{MAP}_B(I, \zeta_1) = 1$

## First Order Decision Diagrams (FODDs) – Semantics based on Multiple Paths

Semantics defined in terms of variable valuations  $\text{MAP}_B(I, \zeta)$

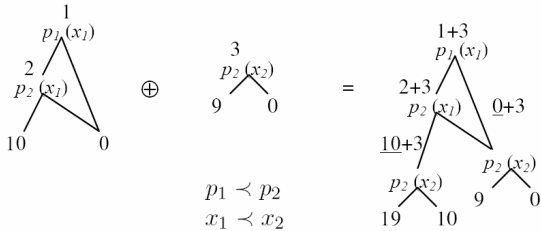
$$\text{MAP}_B(I) = \max_{\xi} \{\text{MAP}_B(I, \xi)\}$$

**- take the maximum over all paths -**



- Domain: {1, 2, 3}
- Interpretation I: {p(1), q(2), h(3)}
- $\text{MAP}_B(I) = 1$

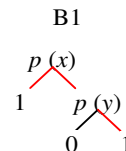
## First Order Decision Diagrams (FODDs) – Combination



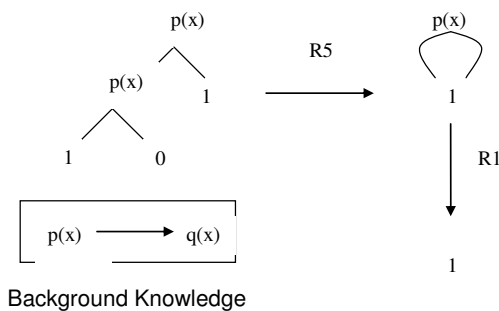
- Assume a fixed order among predicates and parameters
- Choose lower label as new root and combine sub-diagrams recursively.
- Stop when combining two leaves; perform numerical operation

## First Order Decision Diagrams (FODDs) – Reductions

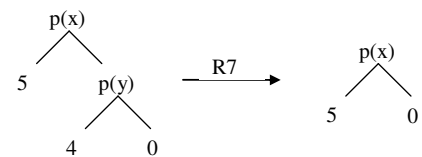
Is this the most compact representation?



## Strong Reduction

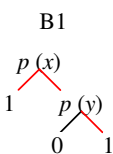


## Weak Reduction



Under any interpretation, for every valuation reaching value 4, there is another valuation reaching value 5

## First Order Decision Diagrams (FODDs) – Reductions



B2

1

No!!! Several reduction operators have been developed

Reductions are key to practical FODDs. The rest is SDP

$$MAP_{B1}(I) = MAP_{B2}(I) \text{ for any } I$$

## Value Iteration using FODDs

- Instance of SDP

$$FODTR[vCase, A(\vec{x})] = rCase \oplus \gamma [\oplus_j \{pCase(n_j(\vec{x})) \otimes Regr[vCase, A(\vec{x})]\}]$$

$$qCase^t(A(\vec{x})) = \max \exists \vec{x}. FODTR[vCase^{t-1}, A(\vec{x})]$$

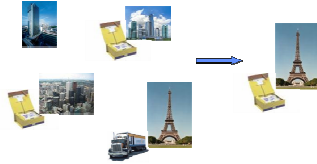
$$vCase^t = \max_{a \in \{A_1(\vec{x}_1), \dots, A_m(\vec{x}_m)\}} qCase^t(a)$$

- Let's go through an example

## FODDs for MDPs - Logistics Domain Revisited

### • Predicates

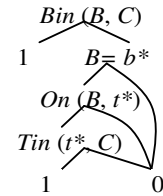
Bin (Box, City),  
Tin (Truck, City),  
On (Box, Truck)



### • Actions

Load(box,truck) {LoadS, LoadF}  
Unload(box,truck) {UnloadS, UnloadF}  
Drive(truck,city)

## Encoding the Domain Dynamics - Truth value diagrams (TVDs)



The TVD for  $\text{Bin}(B, C)$  under  $\text{unloadS}(b^*, t^*)$

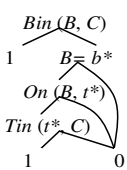
- For every action A schema and predicate schema P, a TVD is a FODD with  $\{0,1\}$  leaves
- Gives the truth value of the predicate P in the next state when A is executed in the current state.
- It captures the truth values for all instances of P
- Only the variables in A and P can appear in the corresponding TVD
- Can express universal effects

## Truth value diagrams (TVDs)

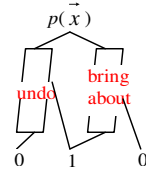
predicate is true

if it **was true before** and is not **undone** by the action

or **was false before** and is **brought about** by the action



The TVD for  $\text{Bin}(B, C)$  under  $\text{unloadS}(b^*, t^*)$

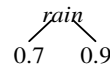


TVD template

Multi-path semantics is beneficial. With single path semantics, a TVD would have to specify all possible ways a predicate can become true

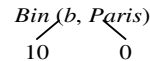
## Encoding Nature's Choice and Reward Function

### • Action choice probabilities



Probability of UnloadS being chosen given Unload is executed.

### • Reward and value functions



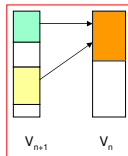
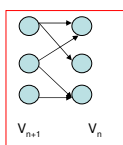
## Value Iteration with FODDs

### • The value iteration algorithm

$$V_{n+1}(s) = \max_{a \in A} [r(s) + \gamma \sum_{s' \in S} \Pr(s' | s, a) V_n(s')]$$

### • The first-order value iteration

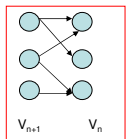
$$1. T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$



## Value Iteration with FODDs

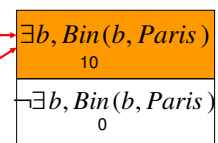
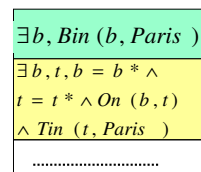
### • The value iteration algorithm

$$V_{n+1}(s) = \max_{a \in A} [r(s) + \gamma \sum_{s' \in S} \Pr(s' | s, a) V_n(s')]$$



### • The first-order value iteration

$$1. T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$



UnloadS

## Value Iteration with FODDs

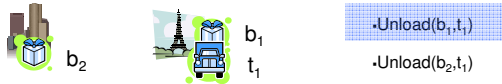
- The value iteration algorithm

$$V_{n+1}(s) = \max_{a \in A} [r(s) + \gamma \sum_{s' \in S} \Pr(s' | s, a) V_n(s')]$$

- The first-order value iteration

$$1. T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$

$$2. Q_{n+1}^A = R \oplus \gamma \otimes \text{obj} - \max(T_{n+1}^{A(\vec{x})}(V_n))$$



## Value Iteration with FODDs

- The value iteration algorithm

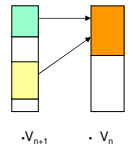
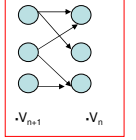
$$V_{n+1}(s) = \max_{a \in A} [r(s) + \gamma \sum_{s' \in S} \Pr(s' | s, a) V_n(s')]$$

- The first-order value iteration

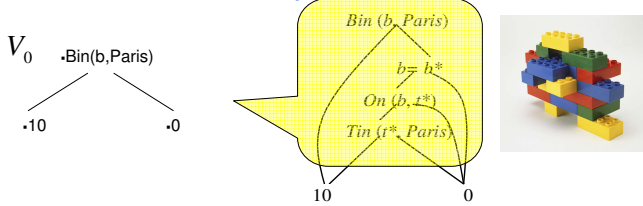
$$1. T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$

$$2. Q_{n+1}^A = R \oplus \gamma \otimes \text{obj} - \max(T_{n+1}^{A(\vec{x})}(V_n))$$

$$3. V_{n+1} = \max_A Q_{n+1}^A$$



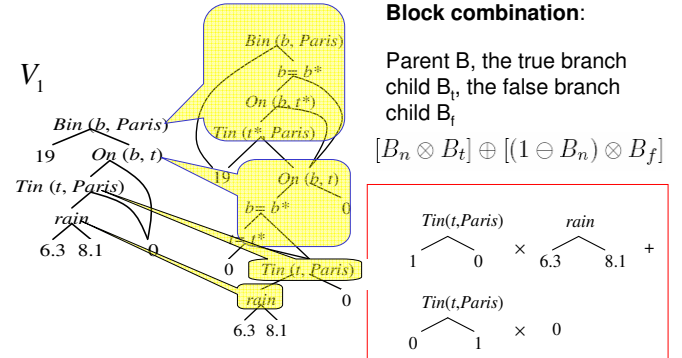
## Value Iteration with FODDs - Regression by block replacement



$$T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$

Replace each node with the corresponding TVD, with the outgoing edges connected to 0 and 1 leaves of the TVD

## Value Iteration with FODDs - Regression by block replacement



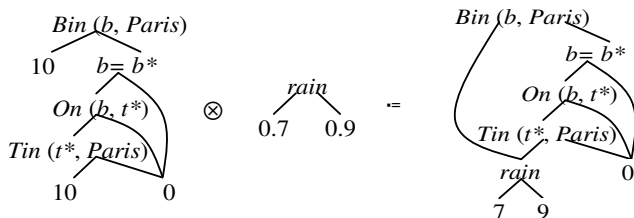
Block combination:

Parent B, the true branch child B<sub>t</sub>, the false branch child B<sub>f</sub>

$$[B_n \otimes B_t] \oplus [(1 \ominus B_n) \otimes B_f]$$

## Value Iteration with FODDs - Adding regression results

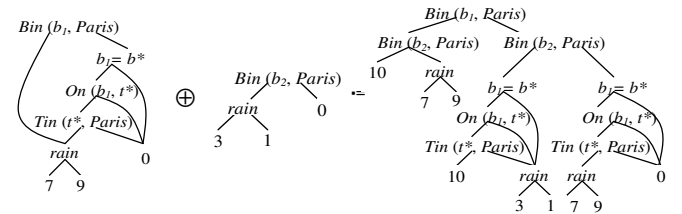
Each regression result multiplied with the corresponding choice probability



$$T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$

## Value Iteration with FODDs - Adding regression results

Standardize apart different regression results before adding

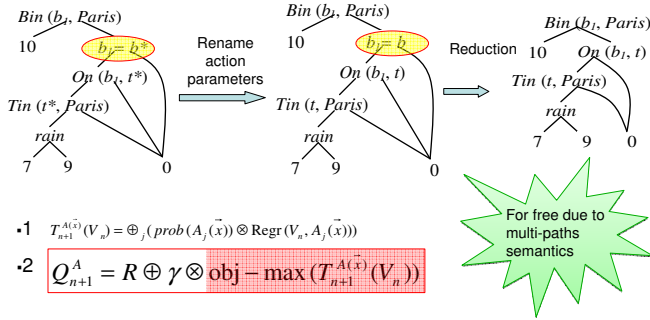


$$T_{n+1}^{A(\vec{x})}(V_n) = \oplus_j (\text{prob}(A_j(\vec{x})) \otimes \text{Regr}(V_n, A_j(\vec{x})))$$

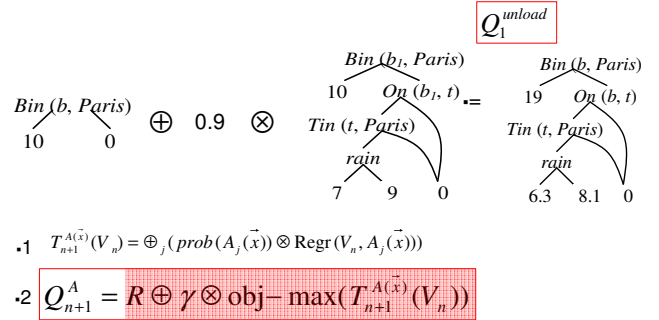
- Finally, reduce the resulting FODD

## Value Iteration with FODDs - Object maximization

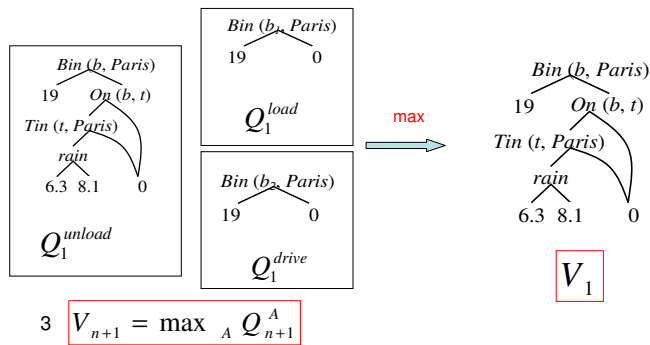
Maximizing over the action parameters to get the maximum value achievable by an instance of this action



## Value Iteration with FODDs



## Value Iteration with FODDs - Maximizing over actions



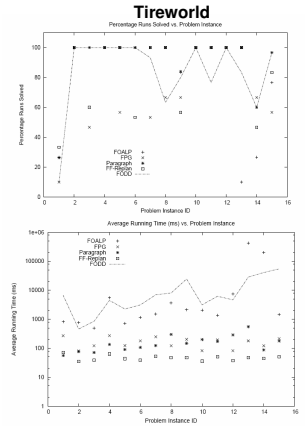
## ICAPS IPPC Results

### Logistics:

- Without approx.
  - slower than ReBel
- with approx.
  - comparable

	Coverage	Time (ms)	Reward
GPT	100%	2220	57.66
Policy Iteration with policy language bias	46.66%	60466	36
Re-Engg NMRDPP	10%	290830	-387.7
FODD-Planner	100%	65000	47.3

Fileworld main results



## FODDs for Relational MDPs - Summary

- FODDs compactly represent functions such as truth values, Q-values etc. over logical spaces
- Complete set of operators to reduce, multiply, add, etc. FODDs
  - direct implementation of SDP

## FODDs for Relational MDPs - Summary

- Approximation a la SPUDD possible: merge ...
  - substructures with similar values
  - Leaves, which are within a certain distance,
  - ....
- Policy iteration approach exists
  - Does not implement the same algorithm as original PI; instead it incorporates an element of policy improvement
  - Theorem: the sequence of value functions obtained from relational modified policy iteration converges monotonically to the optimal value function.
- Initial approach on partially observed, relational MDPs

## ICAPS 2008 Tutorial

### Techniques for First-order Planning

### Inductive First-Order Planning

Kristian Kersting

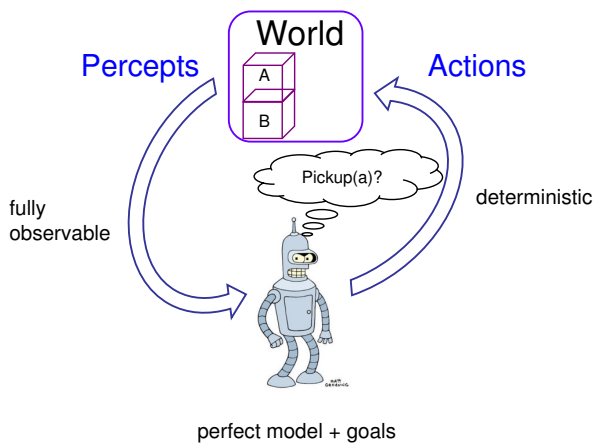
Fraunhofer IAIS

• Thanks to Prasad Tadepalli, Alan Fern, Kurt Driessens, Martijn Van Otterlo,

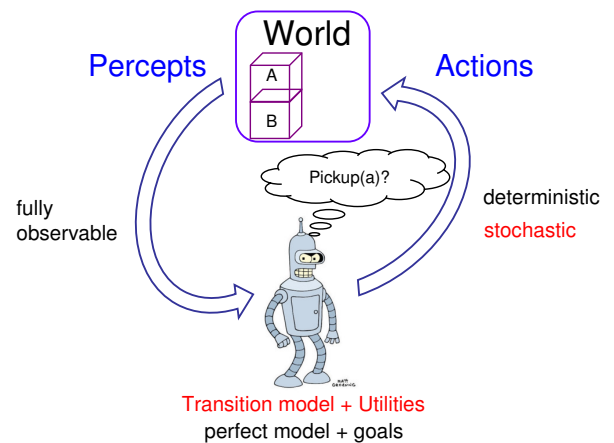
## Tutorial Outline

- **Motivation**
- Deductive First-order Planning
  - Situation Calculus
  - Symbolic Dynamic Programming
  - Relational Bellman Algorithm (ReBeL)
  - First-order Decision Diagrams (FODDs)
- **Inductive First-order Planning**
- Conclusion

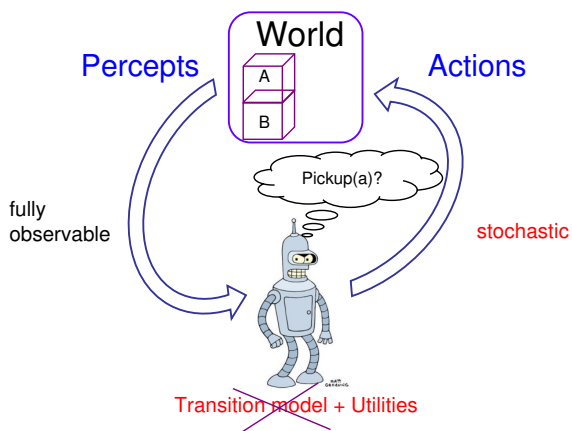
So far: Classical Planning ...



... and relational MDPs



Now: Reinforcement Learning



## Reinforcement Learning

- **No knowledge of environment**
  - Can only act in the world and observe states and reward
- Situated agent: learner **must decide** what actions to take at each step
- Must solve **credit/blame assignment** to actions
  - What actions are responsible for success?
- **Tradeoff** between **exploiting** what is known vs. **exploring** something new
- Must act with limited amount of reasoning

## Model-Based vs. Model-Free RL

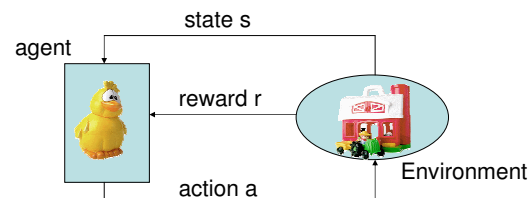
- **Model-based approach to RL:**

- learn the MDP model, or an approximation of it
- use it for policy evaluation or to find the optimal policy
- can use value iteration or policy iteration treating the learned model as if it is correct
- can do sample-based update of the value function

- **Model-free approach to RL:**

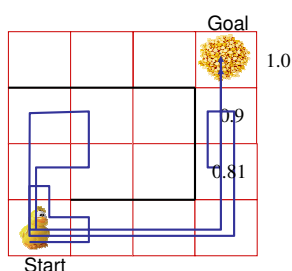
- derive the optimal policy without explicitly learning the model
- learn an action-based value function or Q-function
- we will focus on this in this tutorial !!!

## Reinforcement Learning



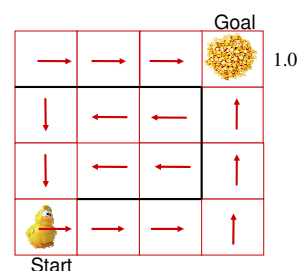
**Agent's goal:** Choose actions to maximize total reward

## Reinforcement Learning



... and so on ...

## Reinforcement Learning



Learn a **policy** mapping **states** to optimal **actions**

## Value Function

Compute

$$V_{i+1}(s) = \sum_a \pi(s, a) \sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V_i(s'))$$

- Prefer actions such that the maximal expected one-step ahead value.
- Could also be done using Q-Values, i.e, values assigned to state action pairs (e.g. Q-Learning)

			Goal
0.729	0.81	0.9	1.0
0.431	0.387	0.349	0.9
0.478	0.431	0.387	0.81
Start	0.531	0.591	0.656
			0.729

## Q-Learning: Model-Free RL

(1) Learn the optimal Q function.

(2) Act greedily with respect to Q(s,a).

- Q(s,a) is the expected value of taking action a in state s and then following the optimal policy thereafter.

- Optimal Q-function satisfies  $Q^*(s) = \max_{a'} Q(s, a')$   
s.t.  $Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') V^*(s')$

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

- After taking action a in state s and reaching s':  
 $Q(s, a) \leftarrow Q(s, a) + \alpha (R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a))$

(noisy) sample of Q-value  
based on next state

## Q-Learning

1. Start with initial Q-function (e.g. all zeros)
2. Take action according to an **explore/exploit policy** (should converge to greedy policy)
3. Perform update  

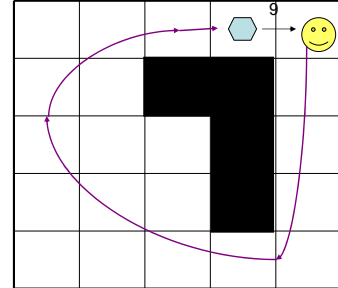
$$Q(s,a) \leftarrow Q(s,a) + \alpha(R(s,a) + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

$Q(s,a)$  is current estimate of optimal Q-function.
4. Goto 2

- Does not require model since we learn Q directly!
- Uses **explicit  $|S| \times |A|$  table** to represent Q
- Explore/exploit policy directly uses Q-values
  - ▲ E.g. use  $\epsilon$ -greedy or Boltzmann exploration.

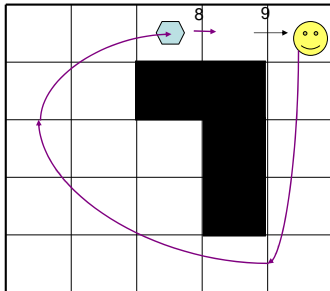
## A Grid Example

Rewards:  
 10 for reaching the goal state  
 -1 for every action.  
 $\alpha$  is set to 1 for simplicity.  
 Update:  $Q(s,a) = r + \max_b (Q(s',b))$



## A Grid Example

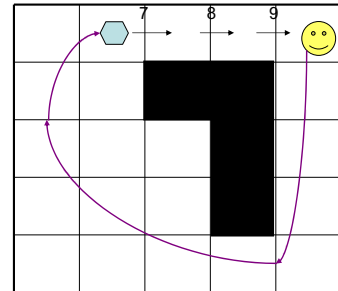
Rewards:  
 10 for reaching the goal state  
 -1 for every action.  
 $\alpha$  is set to 1 for simplicity.  
 Update:  $Q(s,a) = r + \max_b (Q(s',b))$



## A Grid Example

Choose an action  $a = \operatorname{argmax}_a Q(s,a)$  reaching  $s'$

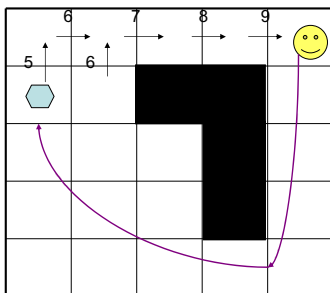
Update  $Q(s,a) = r + \max_b Q(s',b)$



## A Grid Example

Choose an action  $a = \operatorname{argmax}_a Q(s,a)$  reaching  $s'$

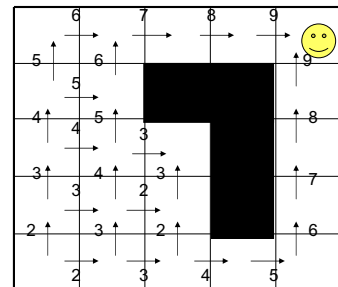
Update  $Q(s,a) = r + \max_b Q(s',b)$



## A Grid Example

Choose an action  $a = \operatorname{argmax}_a Q(s,a)$  reaching  $s'$

Update  $Q(s,a) = r + \max_b Q(s',b)$



The values converge to the optimal Q-values under GLIE policy



## Large-Scale Problems

- Typical state spaces in AI domains are exponentially large
  - Bellman's curse of dimensionality
- Learning a model and utility function
  - Can be difficult to learn good models for large complex environments
  - But if we can learn a model then learning utility function is simpler than learning  $Q(s,a)$
  - Also can reuse the model for “related problems”
- Learning Q-function
  - Simpler to implement since we don't need to worry about representing and learning a model
  - But Q-functions can be substantially more complex than utility functions (they must somehow make up for not having the model)

## Large Relational State Spaces

- When a problem has a large state space we can **not** longer represent the V or Q functions as explicit tables
  - **Generally the case for RMDPs** with a non-trivial numbers of objects
- Even if we had enough memory
  - Never enough training data!
  - Learning takes too long
- What to do??

## Relational RL

- RMDPs with fixed set of objects can be “propositionalized”
  - Describe states via traditional feature vectors that list values of all properties and relations.
    - $[on(a,b)=true, on(b,a)=false, ontable(a)=false, ontable(b)=true]$
- Can then directly apply feature-based RL
  - Loses the relational structure provided by objects
  - Policies can't be applied directly to new object domains
  - Can be difficult to learn from such large feature vectors
- **Relational RL attempts to learn value functions resp. policies that directly exploit relational structure:**
  - Faster learning w.r.t. propositionalization even with fixed # of objects
  - Learn policies that generalize across object domains

## Roadmap for RRL

- **Function Approximation**
- Relational Value Function Learning
  - Propositionalization
  - Relational Regression
- Relational Policy Learning
  - Approximate Policy Iteration
  - Nonparametric Policy Gradient

## Function Approximation

- Never enough training data!
  - Must **generalize** what is learned from one situation to other “similar” new situations

### Basic Idea:

1. Represent Q-function using a compact representation
  - Function encoding size much smaller than table
2. Learn function from experience instead of table

- Q-function updates arising from experience in one state can influence Q-estimate in other similar states
  - Facilitates generalization of experience
- We will first consider feature-based approximation

## Feature Based Function Approx.

- Define a set of n **state-action features**  $f_1(s,a), \dots, f_n(s,a)$ 
  - The features are used as our representation of state-action pairs
  - State-action pairs with similar features will be considered similar
  - In RRL s and a are relational states and actions
- **Example Representation:** linear approximator
 
$$\hat{Q}_\theta(s,a) = \theta_0 + \theta_1 f_1(s,a) + \theta_2 f_2(s,a) + \dots + \theta_n f_n(s,a)$$
- More generally one can use any form of function approximator in terms of these features
  - Regression trees, Kernel regression, Neural networks, etc.

## Q-Learning for Linear Approximators

1. Start with initial parameter values
2. Take action according to an **explore/exploit policy**
3. Perform Q-update for each parameter

$$\theta_i \leftarrow ?$$

4. Goto 2

### Aside: Gradient Descent for Squared Error

- **Given:** sequence of states-action pairs with target Q-values  $\langle s_1, a_1, q(s_1, a_1) \rangle, \langle s_2, a_2, q(s_2, a_2) \rangle, \dots$
- **Goal:** minimize the sum of squared errors between our estimated function and each target value:

$$E_j = \frac{1}{2} (\hat{Q}_\theta(s_j, a_j) - q(s_j, a_j))^2$$

squared error of example j      our estimated value for j'th state      target value for j'th state

- After seeing j'th state the **stochastic gradient descent rule** tells us to update all parameters by:

$$\theta_i \leftarrow \theta_i - \alpha \frac{\partial E_j}{\partial \theta_i}, \quad \frac{\partial E_j}{\partial \theta_i} = \frac{\partial E_j}{\partial \hat{Q}_\theta(s_j, a_j)} \frac{\partial \hat{Q}_\theta(s_j, a_j)}{\partial \theta_i}$$

learning rate

### Aside: continued

$$\theta_i \leftarrow \theta_i + \alpha \frac{\partial E_j}{\partial \theta_i} = \theta_i + \alpha \underbrace{(q(s_j, a_j) - \hat{Q}_\theta(s_j, a_j))}_{\frac{\partial E_j}{\partial \hat{Q}_\theta(s_j, a_j)}} \frac{\partial \hat{Q}_\theta(s_j, a_j)}{\partial \theta_i}$$

depends on form of approximator

- For a linear approximation function:

$$\hat{Q}_\theta(s, a) = \theta_1 f_1(s, a) + \theta_2 f_2(s, a) + \dots + \theta_n f_n(s, a)$$

$$\frac{\partial \hat{Q}_\theta(s_j)}{\partial \theta_i} = f_i(s_j, a_j)$$

- Thus the update becomes:

$$\theta_i \leftarrow \theta_i + \alpha (q(s_j, a_j) - \hat{Q}_\theta(s_j, a_j)) f_i(s_j, a_j)$$

## Q-Learning for Linear Approximators

1. Start with initial parameter values
  2. Take action according to an **explore/exploit policy**
  3. Perform Q-update for each parameter
- $$\theta_i \leftarrow \theta_i + \alpha (R(s, a) + \underbrace{\beta \max_{a'} \hat{Q}_\theta(s', a')}_{\text{Predicted Target } q(s,a)} - \underbrace{\hat{Q}_\theta(s, a)}_{\text{Current estimate}}) f_i(s, a)$$
4. Goto 2

## Roadmap for RRL

- Function Approximation
- **Relational Value Function Learning**
  - Propositionalization
  - Relational Regression
- Relational Policy Learning
  - Approximate Policy Iteration
  - Nonparametric Policy Gradient

## Relational RL via Feature Engineering

- **What if our states and actions are relational?**
- One approach: **engineer a fixed set of “relational features”**
  - Each feature returns a value for any relational state-action pair
  - Features should be well defined regardless of the number of objects
    - Unlike naïve propositionalization approach
  - Ideally feature values should be similar for similar relational states
- With such a feature representation, can use **feature-based Q-learning** to learn **a relational Q-function**.
  - Success relies critically on ability to define appropriate features
  - Often requires significant effort and insight into problem

## Example: Tactical Battles in Wargus

- Wargus is real-time strategy (RTS) game
  - Tactical battles are a key aspect of the game



5 vs. 5



10 vs. 10

- RL Task:** learn a policy to control  $n$  friendly agents in a battle against  $m$  enemy agents
  - Policy should be applicable to tasks with different sets and numbers of agents
  - That is, policy should be relational

## Example: Tactical Battles in Wargus

- Relational States:** contain information about the locations, health, and current activity of all friendly and enemy agent
- Relational Actions:** Attack(F,E)
  - causes friendly agent F to attack enemy E
- Policy Structure:** each decision cycle loop through each friendly agent F and use a learned Q-function to select enemy to attack
  - I.e. select enemy E for F that maximizes  $Q(s, \text{Attack}(F, E))$
- $Q(s, \text{Attack}(F, E))$  is relational since any agents can be substituted for F and E
  - We used a linear function approximator with Q-learning

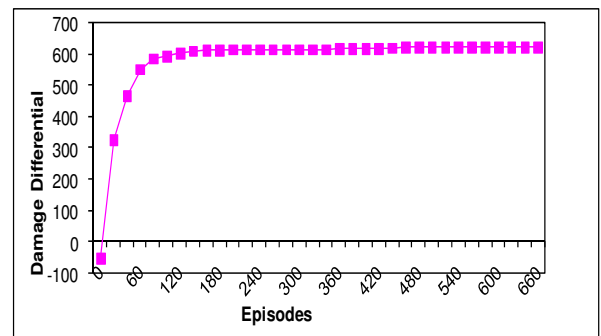
## Example: Tactical Battles in Wargus

$$\hat{Q}_{\theta}(s, a) = \theta_1 + \theta_1 f_1(s, a) + \theta_2 f_2(s, a) + \dots + \theta_n f_n(s, a)$$

- Engineered a set of relational features
  - $\{f_1(s, \text{Attack}(F, E)), \dots, f_n(s, \text{Attack}(F, E))\}$
- Example Features:**
  - # of other friendly agents that are currently attacking E
  - Health of friendly agent F
  - Health of enemy agent E
  - Difference in health values
  - Walking distance between F and E
  - Is E the enemy agent that F is currently attacking?
  - Is F the closest friendly agent to E?
  - Is E the closest enemy agent to E?
  - ...
- Features are well defined for any number of agents

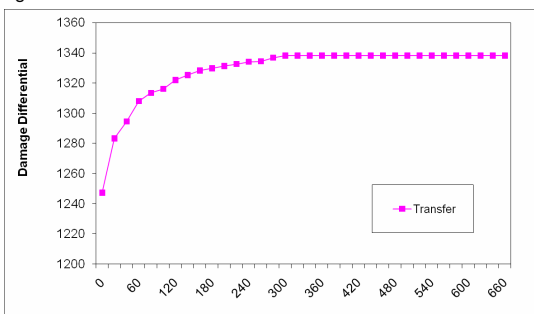
## Example: Tactical Battles in Wargus

- Linear Q-learning in 5 vs. 5 battle



## Example: Tactical Battles in Wargus

- Initialize Q-function for 10 vs. 10 to one learned for 5 vs. 5
  - Initial performance is very good which demonstrates relational generalization from 5 vs. 5 to 10 vs. 10



## Roadmap for RRL

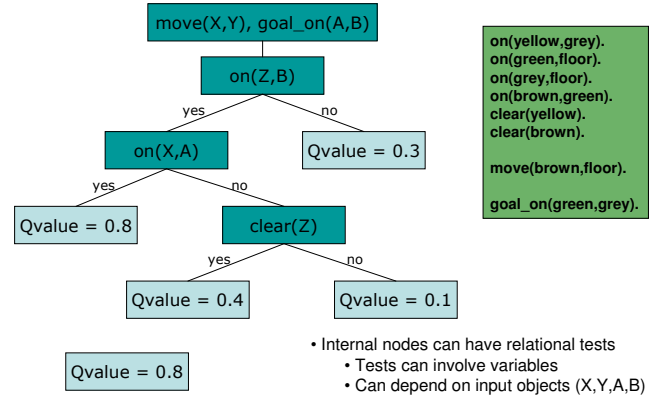
- Function Approximation
- Relational Value Function Learning**
  - Propositionalization
  - Relational Regression**
- Relational Policy Learning
  - Approximate Policy Iteration
  - Nonparametric Policy Gradient

## Relational Regression

- The previous approach relies on feature engineering to reduce a relational problem to a propositional one
  - Requires significant effort and trial-error
- Can we learn a relational value function automatically without propositionalization?**
- Several relational regression algorithms for batch supervised learning
  - Relational regression trees, Gaussian processes, Nearest neighbors

## Relational Regression Trees

[Blockeel, De Raedt AIJ 101]



## Original RRL algorithm

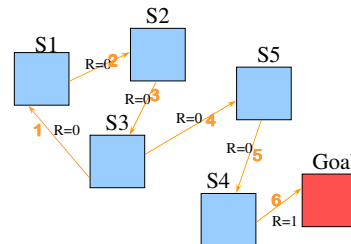
[Dzeroski et al. MLJ 43(1/2)]

- TILDE is a batch learning approach, but RL is an incremental process
- RRL will accumulate training data and call TILDE periodically

initialize an empty **example-set**  
**while** (true)  
 • generate episode through the use of a **standard Q-learning** algorithm using the current tree as Q-function  
 • **generate example** ( $s_i, a_i, q_i$ ) for each state-action pair encountered  
 • **add** the examples **to** the **example-set**  
 • run **tilde** on the knowledge-base

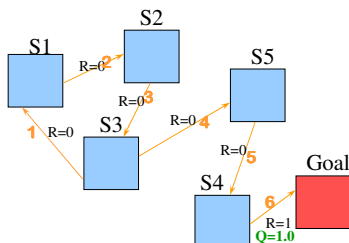
## RRL Example

- Use current policy until a goal-state is reached



## RRL Example

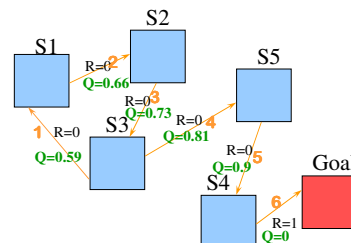
- Back propagate Q-value estimates



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

## RRL Example

- Store (state, action, qvalue) triples as batch training set
- Give to TILDE to learn a tree



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

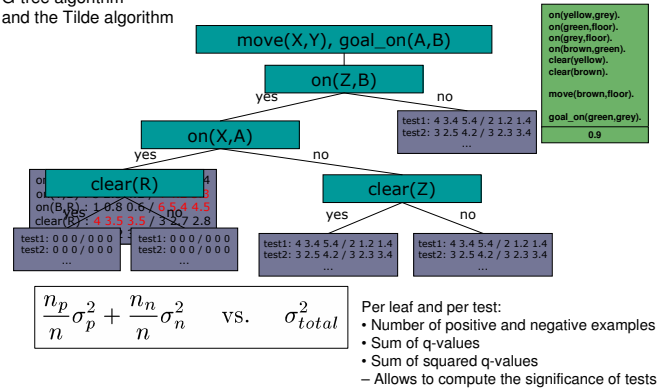
state(s3). action(a). qvalue(0.56).	state(s3). action(d). qvalue(0.81).
state(s1). action(b). qvalue(0.66).	state(s5). action(e). qvalue(0.90).
state(s2). action(c). qvalue(0.73).	state(s4). action(f). qvalue(1.00).

## Extension RRL-TD

- Problems with RRL approach
  - Example set increases with every episode
  - No value-update of old examples
  - Trees are rebuild from scratch each episode
- Build trees incrementally**
  - G-tree algorithm [Chapman & Kaelbling, IJCAI 91] is an online, incremental regression tree learner for propositional data
  - Straightforward to extend to relational setting

## The TG algorithm

Based on the G-tree algorithm and the Tilde algorithm



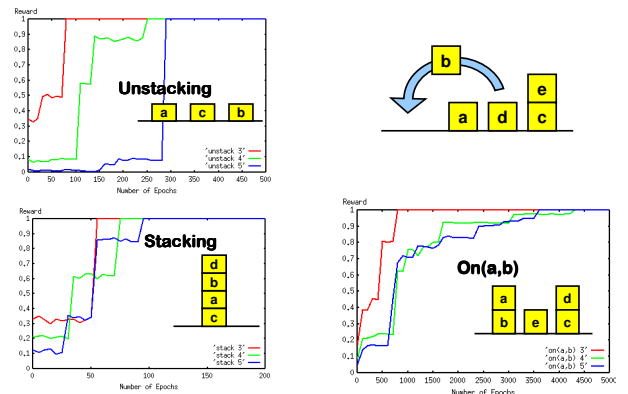
[Driessens et al. ECML01]

## RRL-TG algorithm

- No need to accumulate examples generated during RRL
- Simply update tree incrementally

initialize tree to a **single root node**  
**while** (true)  
 generate episode through the use  
 of a **standard Q-learning** algorithm  
 using the current tree as Q-function  
**generate example**  $(s_i, a_i, q_i)$  for each  
 state-action pair encountered  
**update tree** using the **TG-algorithm** and  
 the generated examples

## Experiments. Blocks World

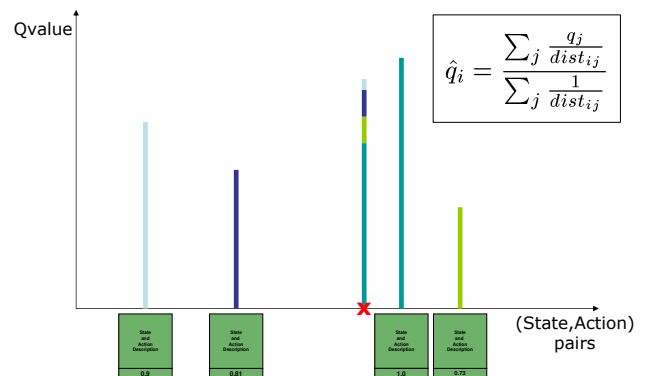


## Timings

		3 blocks	4 blocks	5 blocks
Batch RRL	Stack (30 epochs)	6.16 min	62.4 min	306 min
	Unstack (30 epochs)	8.75 min	Not Stated	Not Stated
	On(a,b) (30 epochs)	20 min	Not Stated	Not Stated
RRL-TG	Stack (200 epochs)	19.2 sec	26.5 sec	39.3 sec
	Unstack (500 epochs)	1.10 min	1.92 min	2.75 min
	On(a,b) (5000 epochs)	25.0 min	57 min	102 min

[Driessens, Ramon ICML03]

## Instance Based Regression



- Requires a distance metric or kernel between relational state-action pairs

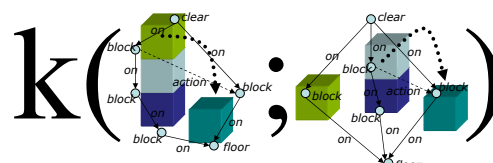
## Instance-Based Regression

- Instance based regression typically
  - Stores past examples (s,a,q)
  - Given a new example (s',a') interpolate wrt example set to estimate q'
- Require a “kernel”  $K(x,x')$  that measures distances between any examples x and x'
- A variety of algorithms exist that turn a kernel function into a regression method
  - Gaussian Processes, Support Vector Regression, K-NN methods

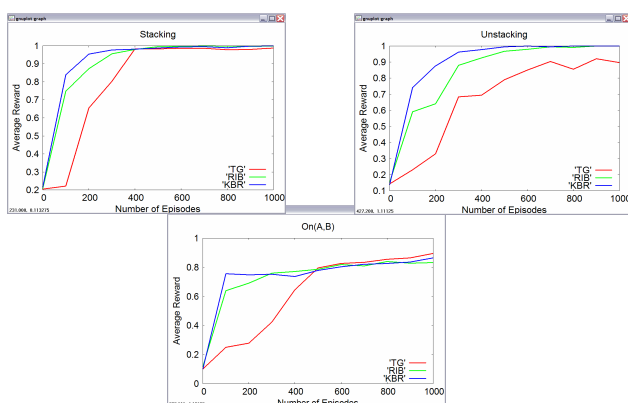
Kernel for relational data?

## Graph Kernels

- Relational can be viewed as graphs
- A number of kernels exist for graphs



## Some Experimental Results

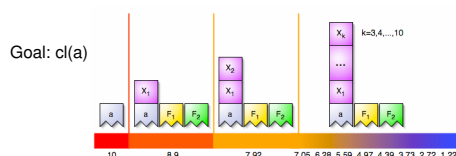


## Roadmap for RRL

- Function Approximation
- Relational Value Function Learning
  - Propositionalization
  - Relational Regression
- Relational Policy Learning
  - Approximate Policy Iteration
  - Nonparametric Policy Gradient

## Direct Policy Learning

- Value functions can often be much more complex to represent than the corresponding policy



- When policies have much simpler representations than the corresponding value functions, **direct search in policy space can be a good idea**

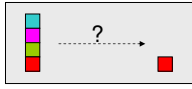
**Policy:** put each block on top of a on the floor

## Roadmap for RRL

- Function Approximation
- Relational Value Function Learning
  - Propositionalization
  - Relational Regression
- Relational Policy Learning
  - Approximate Policy Iteration
  - Nonparametric Policy Gradient

## Challenge Problem

Consider the following class of stochastic blocks world problems:



**Goal:** clear off blocks in the goal

- Optimal policy is:

- obvious to us
- simple and compact
- independent of number of blocks

## Policy for Simple Domains

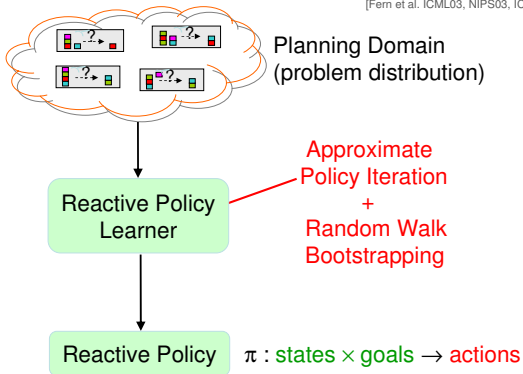
A compact policy for this domain:

- If holding a block, put it down on the table, else...
- Pick up a clear block above a block that is clear in the goal.

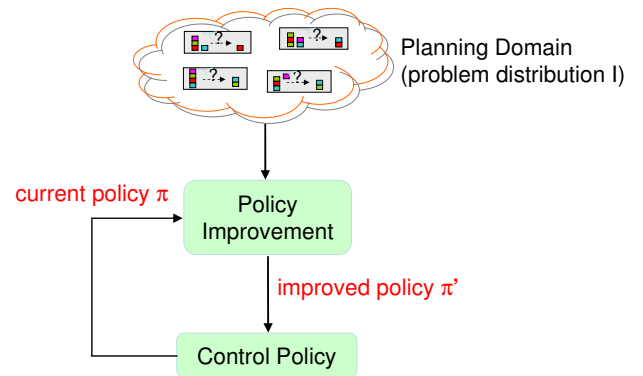


## Learning Domain-Specific Policies

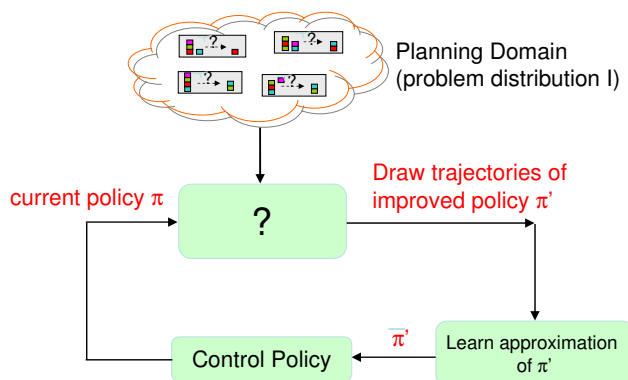
[Fern et al. ICML03, NIPS03, ICAPS04]



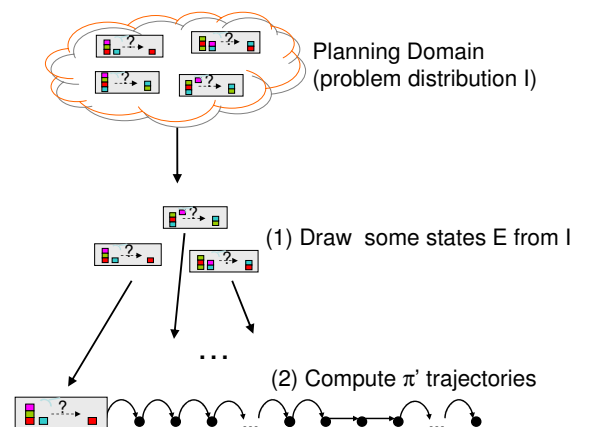
## Policy Iteration



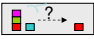
## Approximate Policy Iteration



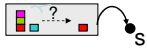
## Drawing Trajectories from improved policy $\pi'$



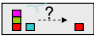
## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

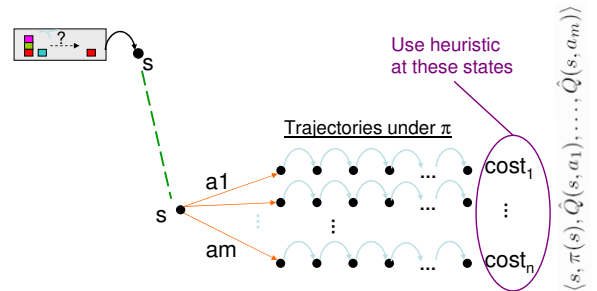
Output: a trajectory under improved policy  $\pi'$




## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

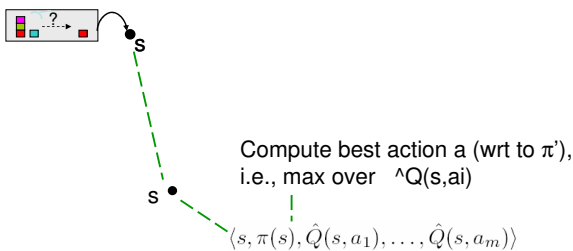
Output: a trajectory under improved policy  $\pi'$




## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

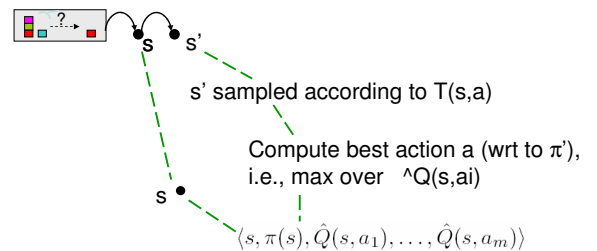
Output: a trajectory under improved policy  $\pi'$



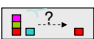
## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

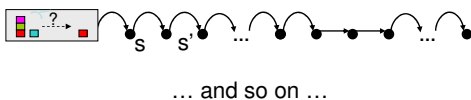
Output: a trajectory under improved policy  $\pi'$



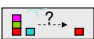
## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

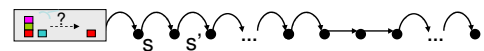
Output: a trajectory under improved policy  $\pi'$



## Computing $\pi'$ Trajectories from $\pi$

Given: current policy  $\pi$  and problem 

Output: a trajectory under improved policy  $\pi'$



This way, we generate examples of state-action pairs that are (approx.) sampled w.r.t.  $\pi'$  because they are (approx.) evaluated taking next "value iteration" into account



## Policy Learning Algorithm

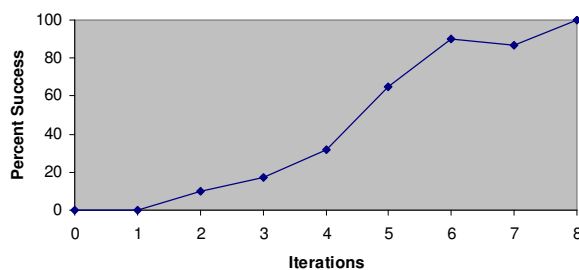
- Training data: generate  $\pi'$  trajectories and save state-action pairs  $\langle s_1, A_1 \rangle, \langle s_2, A_2 \rangle, \langle s_3, A_3 \rangle$ , etc.
- Use a **Rivest-style decision-list learning** approach (induce one rule at a time, in order).
  - While there are training instances remaining,
    - Find a good rule
    - Remove instances covered by the rule
  - **Similar to ReBeL's maximization step**
- Rules found by heuristically guided beam search
  - heuristic combines consistency and coverage
  - rules searched from small to large

## Experiments

- Evaluate on the seven domains from AIPS-2000 + TL-PLAN planning competition.
  - 5 domains : can represent good policies
  - 2 domains : can not represent good policies
- Compared against state-of-the-art planner FF
  - FF's heuristic is very good for most of these domains.
- API equal or better FF's performance when we can represent policies.

## Experimental Results

Blocks World (20 blocks)



## Domains with Good Policies

Success Percentage

	Blocks	Elevator	Schedule	Briefcase	Gripper
API	100	100	100	100	100
FF-Plan	28	100	100	0	100

Typically our solution lengths are comparable to FF's.

## Roadmap for RRL

- Function Approximation
- Relational Value Function Learning
  - Propositionalization
  - Relational Regression
- **Relational Policy Learning**
  - Approximate Policy Iteration
  - **Nonparametric Policy Gradient**

## Policy Gradients with Function Approximation (Sutton et al.)

- Parameterized policy:  $\pi(s, a, \theta)$
- Gradient search w.r.t. world-value

$$\begin{aligned}
 \frac{\partial p(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_{s,a} d^\pi(s) \pi(s, a, \theta) Q^\pi(s, a) \\
 &= \sum_{s,a} d^\pi(s) Q^\pi(s, a) \frac{\partial \pi(s, a, \theta)}{\partial \theta}
 \end{aligned}$$

## Non-Parametric Policy Gradient

- Express policy using an arbitrary potential function

$$\pi(s, a, \Psi) = \frac{e^{\Psi(s, a)}}{\sum_b e^{\Psi(s, b)}}$$

- Use a Functional Gradient (Friedmann)

$$\Delta \Psi \approx \alpha \frac{\partial \rho}{\partial \Psi}$$

## Functional Gradient Boosting

- (inspired by Friedman et al. and Dietterich et al.)
- Regular Parameterized Gradients

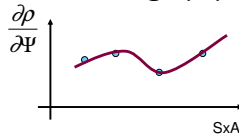
$$F(\theta) \rightarrow \theta_m = \theta_0 + \delta_1 + \delta_2 + \delta_3 + \dots + \delta_m$$

- Functional Gradients

$$\Psi \rightarrow \Psi_m = \Psi_0 + \Delta_1 + \Delta_2 + \Delta_3 + \dots + \Delta_m$$

## Functional Gradient Boosting (2)

- Evaluate gradient locally
  - Wherever there is data
- Use your favorite algorithm to generalize
  - Propositional, relational, continuous



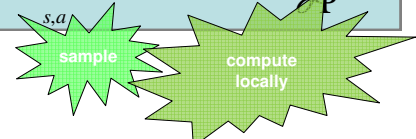
$$\Psi \rightarrow \Psi_m = \Psi_0 + \Delta_1 + \Delta_2 + \Delta_3 + \dots + \Delta_m$$

Infinitely many parameters

## In Practice

- Following Sutton et al.

$$\begin{aligned} \frac{\partial \rho}{\partial \Psi} &= \frac{\partial}{\partial \Psi} \sum_{s, a} d^\pi(s) \pi(s, a) Q^\pi(s, a) \\ &= \sum_{s, a} d^\pi(s) Q^\pi(s, a) \frac{\partial \pi(s, a)}{\partial \Psi} \end{aligned}$$



## Local Evaluation

$$Q^\pi(s, a)$$

Simple: Monte-Carlo estimate  
Future Work: actor-critic

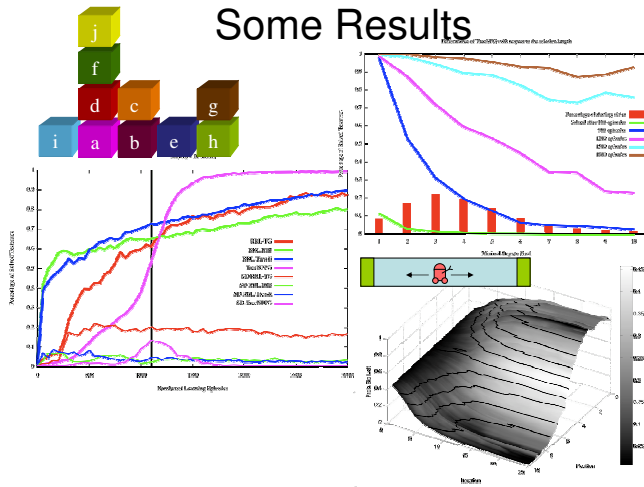
$$\pi(s, a) = \frac{e^{\Psi(s, a)}}{\sum_b e^{\Psi(s, b)}}$$

$$\begin{aligned} \frac{\partial \pi(s, a)}{\partial \Psi(s, a)} &= \pi(s, a)(1 - \pi(s, a)) \\ \frac{\partial \pi(s, a)}{\partial \Psi(s, b)} &= -\pi(s, a)\pi(s, b) \end{aligned}$$

## Gradient Tree Boosting

- Generate behavior traces following  $\pi(s, a, \Psi) = \frac{e^{\Psi(s, a)}}{\sum_b e^{\Psi(s, b)}}$
- For each encountered state  $s$   
For each available action in that state  
Generate a learning example:  
for chosen action  $a$   $\langle s, a, Q(s, a)\pi(s, a)(1 - \pi(s, a)) \rangle$   
for other actions  $b$   $\langle s, b, -Q(s, a)\pi(s, a)\pi(s, b) \rangle$
- Learn a tree:  $\Delta_{n+1}$
- $\Psi_{n+1} = \Psi_n + \Delta_{n+1}$

## Some Results



## Conclusions: Inductive FO Planning

- Model-free RL algorithms solve MDPs without knowing the transition model
- Relational RL algorithms attempt to exploit the relational structure
  - ▶ Faster learning, policies generalize across object domains
- Can exploit relational regression and classification methods for model-free RRL
  - ▶ Feature engineering
  - ▶ Relational regression trees (RRL-TG, NPPG)
  - ▶ Relational decision lists (API)
  - ▶ Nonparametric Policy Gradient
    - ▶ Unifies finite, continuous, and relational RL