



# Practical solution techniques for first-order MDPs<sup>☆</sup>

Scott Sanner<sup>a,\*</sup>, Craig Boutilier<sup>b</sup>

<sup>a</sup> Statistical Machine Learning Group, National ICT Australia, Canberra, ACT, 0200, Australia

<sup>b</sup> Department of Computer Science, University of Toronto, Toronto, ON M5S 3H5, Canada

## ARTICLE INFO

### Article history:

Received 21 October 2007

Received in revised form 4 November 2008

Accepted 9 November 2008

Available online 24 November 2008

### Keywords:

MDPs

First-order logic

Planning

## ABSTRACT

Many traditional solution approaches to relationally specified decision-theoretic planning problems (e.g., those stated in the probabilistic planning domain description language, or PPDDL) ground the specification with respect to a specific instantiation of domain objects and apply a solution approach directly to the resulting ground Markov decision process (MDP). Unfortunately, the space and time complexity of these grounded solution approaches are polynomial in the number of domain objects and exponential in the predicate arity and the number of nested quantifiers in the relational problem specification. An alternative to grounding a relational planning problem is to tackle the problem directly at the relational level. In this article, we propose one such approach that translates an expressive subset of the PPDDL representation to a first-order MDP (FOMDP) specification and then derives a domain-independent policy without grounding at any intermediate step. However, such generality does not come without its own set of challenges—the purpose of this article is to explore practical solution techniques for solving FOMDPs. To demonstrate the applicability of our techniques, we present proof-of-concept results of our first-order approximate linear programming (FOALP) planner on problems from the probabilistic track of the ICAPS 2004 and 2006 International Planning Competitions.

Crown Copyright © 2008 Published by Elsevier B.V. All rights reserved.

## 1. Introduction

There has been an extensive line of research over the years aimed at exploiting structure in order to compactly represent and efficiently solve decision-theoretic planning problems modeled as Markov decision processes (MDPs) [12]. While traditional approaches from operations research typically use enumerated state and action models [62], these have proved impractical for large-scale AI planning tasks where the number of distinct states in a model can easily exceed the limits of primary and secondary storage on modern computers.

Fortunately, many MDPs can be compactly described by using a factored state and action representation and exploiting various independencies in the reward and transition functions [12]. The independencies and regularities laid bare by such representations can often be exploited in exact and approximate solution methods as well. Such techniques have permitted the practical solution of MDPs that would not have been possible using enumerated state and action models [22,36,38,75].

However, factored representations are only one type of structure that can be exploited in the representation of MDPs. Many MDPs can be described abstractly in terms of classes of domain objects and relations between those domain objects that may change over time. For example, a logistics problem specified in the probabilistic planning domain description

<sup>☆</sup> Parts of this article appeared in preliminary form in [S. Sanner, C. Boutilier, Approximate linear programming for first-order MDPs, in: Uncertainty in Artificial Intelligence (UAI-05), Edinburgh, Scotland, 2005, pp. 509–517; S. Sanner, C. Boutilier, Practical linear evaluation techniques for first-order MDPs, in: Uncertainty in Artificial Intelligence (UAI-06), Boston, MA, 2006].

\* Corresponding author.

E-mail address: ssanner@nicta.com.au (S. Sanner).

language (PPDDL) [89] may refer to domain objects such as boxes, trucks, and cities. If the objective is to deliver all boxes to their assigned destination cities then the locations of these boxes and trucks may change as a result of actions taken in pursuit of this objective. Since action templates such as loading or unloading a box are likely to apply generically to domain objects and can be specified independently of any ground domain instantiation (e.g., 4 trucks, 5 boxes, and 9 cities), this permits compact MDP descriptions by exploiting the existence of domain objects, relations over these objects, and the ability to express objectives and action effects using quantification.

Unfortunately, while relational specifications such as PPDDL permit very compact, domain-independent descriptions of a variety of MDPs, this compactness does not translate directly to effective solutions of the underlying planning problems. For example, one approach to solving a relational decision-theoretic planning problem might first construct sets of state variables and actions for all possible ground instantiations of each relation and action with respect to a specific domain (e.g., 4 trucks, 5 boxes, and 9 cities). Then this approach might apply known solution techniques to this ground factored representation of an MDP. Unfortunately, such an approach is domain-specific; and the size of the ground MDP grows polynomially in the number of domain objects, and exponentially in the predicate arity and the number of nested quantifiers in the problem specification. For sufficiently large domains and complex relational MDP specifications, grounding may not be a viable option.

An alternative approach to grounding is to apply a solution approach directly at the relational level. In this article, we discuss one such technique that translates an expressive subset of the relational PPDDL representation to a *first-order MDP (FOMDP)* [14] specification. A symbolic policy may then be derived with respect to this FOMDP, resulting in a domain-independent solution that exploits a purely lifted version of the Bellman equations and avoids grounding at any intermediate step. This stands in contrast to alternate first-order approaches discussed in Section 6.2 that induce symbolic representations of the solution from samples of the Bellman equation in ground problem instances.

Unfortunately, the use of first-order logical languages to describe our FOMDP specification and solution introduces the need for computationally expensive logical simplification and theorem proving. While this means that exact solutions are not tractable for many FOMDPs, there is often a high degree of regularity and structure present in many FOMDPs that can be exploited by the approximate (heuristic) solution techniques proposed in this article. To this end, this article continues the tradition of exploiting structure to find effective solutions for large MDPs.

After providing a review of MDPs and relevant solution techniques in Section 2 and the FOMDP formalism and its solution via symbolic dynamic programming [14] in Section 3, we make the following contributions to the practical solution of FOMDPs:

- (1) Section 3.2.2: We show how to translate a subset of PPDDL problems including universal and conditional effects to FOMDPs.
- (2) Section 4.1: We show how to exploit the logical structure of reward, value, and transition functions using first-order extensions of algebraic decision diagrams (ADDs) [4] for use in both exact and approximate FOMDP solutions.
- (3) Section 4.2: We apply additive decomposition techniques to universal reward specifications in a manner that leads to efficient solutions for our FOMDP representation and reasonable empirical performance on example problems.
- (4) Section 5.3: We show how to generalize the approximate linear programming technique for MDPs [19,36,72] to the case of FOMDPs by casting the optimization problem in terms of a first-order linear program.
- (5) Section 5.4: We define a linear program (LP) with first-order constraints and provide a constraint generation algorithm that utilizes a relational generalization of variable elimination [91] to exploit constraint structure in the efficient solution of this *first-order LP (FOLP)*.

To demonstrate the efficacy of our techniques, we present proof-of-concept results of our *first-order approximate linear programming (FOALP)* planner on problems from the probabilistic track of the ICAPS 2004 and 2006 International Planning Competitions in Section 5.6. Following this, we discuss a number of related first-order decision-theoretic planning approaches and discuss the relative advantages and disadvantages of each in Section 6. We conclude with a discussion of possible extensions to our techniques in Section 7.

## 2. Markov decision processes

Markov decision processes (MDPs) were first introduced and developed in the fields of operations research and economics [6,41,73]. The MDP has since been adopted as a model for decision-theoretic planning with fully observable state in the field of artificial intelligence [7,8,12] and as such provides the formal underpinning for the framework that we describe in this article. In this section, we describe various algorithmic approaches for making optimal sequential decisions in MDPs that we later generalize to the case of first-order MDPs. The following presentation derives from Puterman [62].

### 2.1. The MDP model and optimality criteria

Formally, a finite state and action MDP is specified by a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, h, \gamma \rangle$ .  $\mathcal{S}$  is a set of distinct states. An agent in an MDP can effect changes to its state by executing actions from the set  $\mathcal{A}$ . We base our initial presentation in this section on finite state and action MDPs; but in much of what follows, we will assume an infinite, discrete state and action

space. The standard techniques for MDPs discussed here can be generalized to countable or continuous state and action spaces [62].

The transition function  $\mathcal{T}$  is a family of probability distributions  $\mathcal{T}(s, a, s') = P(s'|a, s)$ , which denotes the probability that the world transitions from  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$  when action  $a \in \mathcal{A}$  was executed. This representation enforces the Markov property: the distribution over states  $s_{t+1}$  at time  $t + 1$  is independent of any previous state  $s_{t-i}$  and action  $a_{t-i}$ ,  $i \geq 1$ , given  $s_t$  and  $a_t$ .

The preferences of the agent are encoded in a reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . In addition to specifying single-step preferences, the agent must also specify how it trades off reward over the horizon  $h$  of remaining decision stages. In this article, we focus on the expected sum of discounted accumulated reward over an infinite horizon ( $h = \infty$ ) since this is most compatible with the (approximate) linear programming approach that we adopt later. In the calculation of discounted accumulated reward, we discount rewards  $t$  time steps into the future by a discount factor  $\gamma^t$  where  $\gamma \in [0, 1]$ . Throughout this article, we assume  $\gamma < 1$ . The use of  $\gamma < 1$  allows one to model the notion that an immediate reward  $r$  is worth more than the equivalent reward delayed one or more time steps in the future. Practically,  $\gamma < 1$  is required to ensure that the total expected reward is bounded in the case of infinite horizon MDPs.

A stationary policy takes the form  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , with  $\pi(s)$  denoting the action to be executed in state  $s$ . The *value* of policy  $\pi$  is the expected sum of discounted future rewards over horizon  $h$  given that  $\pi$  is executed. Its value function is given by:

$$V_\pi(s) = E_\pi \left[ \sum_{t=0}^h \gamma^t \cdot r^t \mid s_0 = s \right] \quad (1)$$

where  $r^t$  is a reward obtained at time  $t$ ,  $\gamma$  is a discount factor as defined above, and  $s_0$  is the initial starting state.

A *greedy policy*  $\pi_V$  with respect to a value function  $V$  is simply any policy that takes an action in each state that maximizes expected value with respect to  $V$ , defined as follows:

$$\pi_V(s) = \arg \max_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) V(s') \right\} \quad (2)$$

Thus, from any value function, we can derive a corresponding greedy policy that represents the best action choice with respect to that value estimation.

An *optimal policy*  $\pi^*$  in an infinite horizon MDP maximizes the value function for all states. An optimal policy  $\pi^*$  is any greedy policy with respect to the optimal value function  $V^*$  and likewise the optimal value function is the value of an optimal policy,  $V_{\pi^*}(s) = V^*(s)$ . We note that  $V^*$  satisfies the following fixed-point equality:

$$V^*(s) = \max_a \left\{ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} T(s, a, s') \cdot V^*(s') \right\} \quad (3)$$

Finding  $V^*$  constitutes finding an *exact solution* to an MDP. Throughout the article, we use the term *solution* more loosely to denote some attempt at approximating  $V^*$ , whether the approximation guarantees error bounds or is simply heuristic.

## 2.2. MDP solution algorithms

In this section we describe several exact and approximate solution techniques for MDPs that we later extend to the first-order case.

### 2.2.1. Value iteration

We begin our discussion of MDP solutions by providing two equations that form the basis of the stochastic dynamic programming algorithms used to solve MDPs.

We define  $V_\pi^0(s) = R(s, \pi(s))$  and then inductively define the *t-stage-to-go value function* for a policy  $\pi$  as follows:

$$V_\pi^t(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} T(s, \pi(s), s') \cdot V_\pi^{t-1}(s') \quad (4)$$

Based on this definition, Bellman's *principle of optimality* [6] establishes the following relationship between the optimal value function at stage  $t$  and the optimal value function at the previous stage  $t - 1$ :

$$V^t(s) = \max_{a \in \mathcal{A}} \left\{ R(s, a) + \gamma \sum_{s' \in \mathcal{S}} T(s, a, s') \cdot V^{t-1}(s') \right\} \quad (5)$$

The computation of  $V^t$  from  $V^{t-1}$  via this relationship is referred to as a *Bellman backup*. The *value iteration* algorithm consists of repeatedly performing Bellman backups to compute these *t-stage-to-go* value functions.

We note that the Bellman backup is often rewritten in the following two steps to separate out the backup of a value function through a single action and the maximization of this value over all actions:

$$Q^t(s, a) = R(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} T(s, a, s') \cdot V^{t-1}(s') \quad (6)$$

$$V^t(s) = \max_{a \in \mathcal{A}} \{Q^t(s, a)\} \quad (7)$$

Puterman [62] shows that terminating once the following Bellman error condition is met

$$\max_s |V^t(s) - V^{t-1}(s)| < \frac{\epsilon(1 - \gamma)}{2\gamma} \quad (8)$$

guarantees that the estimated value function  $V^t$  is  $\epsilon$ -optimal over an infinite horizon, that is, its value is within  $\epsilon$  of the optimal value:  $\max_s |V^t(s) - V^*(s)| < \epsilon$ .

We note that the value iteration approach requires time polynomial in the backup depth  $d$  and the number of states and actions, i.e.,  $O(|\mathcal{S}|^2 \cdot |\mathcal{A}| \cdot d)$ . Puterman [62] provides a proof that value iteration converges linearly.

### 2.2.2. Linear programming

An MDP can also be solved using the following linear program (LP):

$$\begin{aligned} \text{Variables:} \quad & V(s), \quad \forall s \in \mathcal{S} \\ \text{Minimize:} \quad & \sum_{s \in \mathcal{S}} V(s) \\ \text{Subject to:} \quad & 0 \geq R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)V(s') - V(s); \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \end{aligned} \quad (9)$$

Puterman [62] provides a proof that the solution to this LP is the optimal value function for an MDP.

### 2.2.3. Approximate linear programming

One general and popular approximate solution technique for MDPs is that of linear-value function approximation [36, 46, 47, 71, 72, 78]. Representing value functions as a linear combination of basis functions has many convenient computational properties, many of which will become evident as we incorporate relational structure in our MDP model. However, perhaps one of the most useful properties is that linear value function representations lead to MDP solutions requiring optimization with respect to linear objectives and linear constraints—that can be formulated as LPs.

In an  $n$ -state MDP, the exact value function can be specified as a vector in  $\mathbb{R}^n$ . This vector can be approximated by a value function  $\tilde{V}_{\tilde{w}}$  that is a linear combination of  $k$  fixed basis functions (or  $n$ -vectors), denoted  $b_i(s)$ :

$$\tilde{V}_{\tilde{w}}(s) = \sum_{i=1}^k w_i \cdot b_i(s) \quad (10)$$

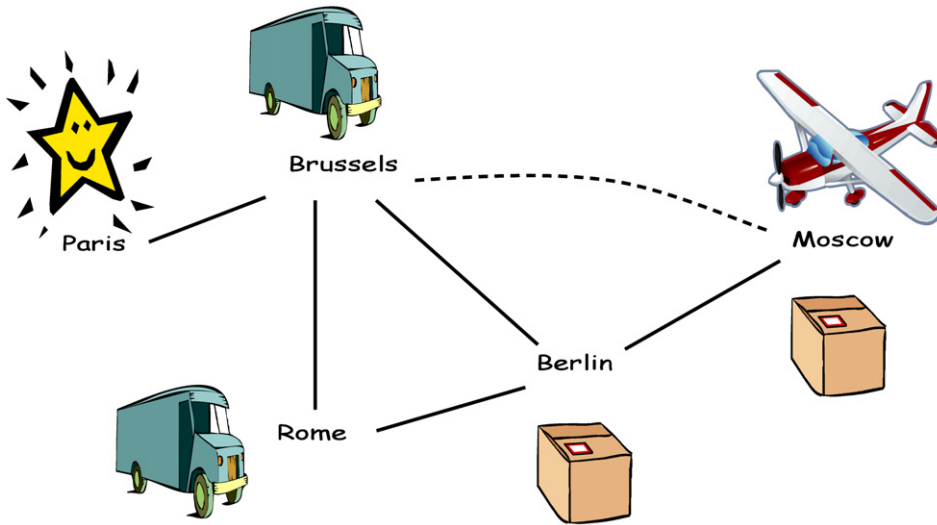
The linear subspace spanned by the basis set will generally not include the true value function, but one can use projection methods to minimize some error measure between the true value function and the linear combination of basis functions. The basis functions themselves can be specified by domain experts, constructed or learned in an automated fashion (e.g., [61]; [51]). We will consider first-order methods for automated basis function construction in Section 5 and related work in Section 6.

*Approximate linear programming (ALP)* is simply an extension of the linear programming solution of MDPs to the case where the value function is approximated. In a linear value function representation, the objective and constraints will be linear in the weights being optimized, leading to a direct LP formulation. Consequently, we arrive at the following variant of the previous exact LP solution:

$$\begin{aligned} \text{Variables:} \quad & \tilde{w} \\ \text{Minimize:} \quad & \sum_{s \in \mathcal{S}} \tilde{V}_{\tilde{w}}(s) \\ \text{Subject to:} \quad & 0 \geq R(s) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a)\tilde{V}_{\tilde{w}}(s') - \tilde{V}_{\tilde{w}}(s); \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \end{aligned} \quad (11)$$

## 2.3. Selecting an MDP solution approach

The choice of whether to use a linear programming or dynamic programming solution to MDPs is not always clear. Linear programming offers a simple one-shot solution, but it relies on efficient LP solvers. Dynamic programming is straightforward to implement, but may require a large number of iterations to converge. However, the choice of exact vs. approximate is almost invariably determined by the size of the state space. For sufficiently large state spaces, approximate solution techniques are the only viable option. But this last statement depends critically on how one measures the size of the state space.



**Fig. 1.** An example BoxWorld problem. Trucks may drive along solid lines and planes may fly along dashed lines. The goal in this instance is to get all boxes in Paris (indicated by the star).

Despite their promise, the exact and approximate solution techniques discussed above must represent the value function (and policy, if required) as vectors or functions over an explicitly enumerated state (and action) space. This is simply not feasible for large-scale AI planning problems. Fortunately, there are many representations (e.g., factored or relational) well suited to decision-theoretic planning that do not require explicit state or action enumeration in either the problem representation or the solution. To this end, we will be concerned with the exploitation of relational planning structure for the remainder of this article.

### 3. First-order MDPs

Given that relational representations seem natural for planning problems, it makes sense to attempt to exploit this relational structure at a first-order level without resorting to grounding. This is precisely the idea behind the first-order MDP model (FOMDP) and its symbolic dynamic programming solution [14], which we review in this section. For the remainder of this article, when we refer to a FOMDP without further qualification, we refer to the specific formalization presented in [14], although there are other possible first-order MDP formalizations and associated solution approaches (we discuss these alternatives in Section 6). The reader already familiar with the motivations for FOMDPs and the presentation and notation in [14] may wish to skip this section and proceed directly to the main contributions of this article in Sections 4 and 5.

#### 3.1. Motivation

Before we introduce FOMDPs and their solution, we begin with the basics of relational planning problem specifications and motivate the need for exploiting this structure at a lifted first-order level rather than at a ground propositional level.

##### 3.1.1. Relational planning specifications

We assume basic familiarity with unsorted first-order logic with equality. While we use a sorted notation for specifying object types of variables and predicate slots, we assume this sort information is compiled into an unsorted logical form where  $\forall \text{Sort} : c \phi(c)$  is rewritten as  $\forall c. \text{Sort}(c) \supset \phi(c)$  and likewise  $\exists \text{Sort} : c \phi(c)$  is rewritten as  $\exists c. \text{Sort}(c) \wedge \phi(c)$ . Assuming these transformations, we draw on the logical notation and semantics for unsorted first-order logic given in [16]. Specifically:

- **Predicate Symbols:** We assume a set of predicates  $P_i$  of each arity  $0 \leq i \leq m$  for some finite maximum  $m$ . We assume “=”  $\in P_2$  with its usual interpretation.
- **Function Symbols:** We assume a set of function symbols  $f_j$  of each arity  $0 \leq j \leq n$  for some finite maximum  $n$ .

In addition, we use a few notational conventions. All predicates (including unary predicates denoting domain object classes) are capitalized and all variables and constants are lowercased. We denote the types of predicate arguments using the notation  $\phi(\text{Sort}_1, \dots, \text{Sort}_k)$  for some predicate of arity  $k$ .<sup>1</sup>

<sup>1</sup> Logically, this requires a background theory axiom  $\forall x_1, \dots, x_k \phi(x_1, \dots, x_k) \supset \bigwedge_{i=1}^k \text{Sort}_i(x_i)$  for each predicate  $\phi(\text{Sort}_1, \dots, \text{Sort}_k)$ .

- 
- **Domain Object Types:** *Box, Truck, City*
  - **Relational ( $\mathcal{S}$ )tate Descriptors (with parameter sorts):**  
 $\text{BoxIn}(\text{Box}, \text{City}), \text{TruckIn}(\text{Truck}, \text{City}), \text{BoxOn}(\text{Box}, \text{Truck})$
  - **( $\mathcal{R}$ )eward:** if  $[\exists \text{Box} : b. \text{BoxIn}(b, \text{paris})]$  then 10 else 0
  - **( $\mathcal{A}$ )ctions (with parameter sorts) and ( $\mathcal{T}$ )ransition Function:**
    - $\text{load}(\text{Box} : b, \text{Truck} : t)$ :
      - Effects (probability 0.9):
        - \* when  $[\exists \text{City} : c. \text{BoxIn}(b, c) \wedge \text{TruckIn}(t, c)]$  then  $[\text{BoxOn}(b, t)]$
        - \*  $\forall \text{City} : c. \text{when} [\text{BoxIn}(b, c) \wedge \text{TruckIn}(t, c)]$  then  $[\neg \text{BoxIn}(b, c)]$
    - $\text{unload}(\text{Box} : b, \text{Truck} : t)$ :
      - Effects (probability 0.9):
        - \*  $\forall \text{City} : c. \text{when} [\text{BoxOn}(b, t) \wedge \text{TruckIn}(t, c)]$  then  $[\text{BoxIn}(b, c)]$
        - \* when  $[\exists \text{City} : c. \text{BoxOn}(b, t) \wedge \text{TruckIn}(t, c)]$  then  $[\neg \text{BoxOn}(b, t)]$
    - $\text{drive}(\text{Truck} : t, \text{City} : c)$ :
      - Effects (probability 1.0)
        - \* 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  $[\neg \text{TruckIn}(t, c_1)]$
    - $\text{noop}$ 
      - No effects.
- 

**Fig. 2.** A PPDDL-style representation of a simple variant of the BoxWORLD problem. The deterministic PDDL subset would exclude the probabilistic annotations of effects assuming that all effects occur with probability 1.0.

We can view many decision-theoretic planning problems as consisting of classes of domain objects and the changing relations that hold between those objects at different points in time. For example, in the BoxWORLD logistics problem [79] illustrated in Fig. 1, we have four classes of domain objects: *Box*, *City*, *Truck*, and *Plane*. For the relations that hold between them, we have  $\text{BoxIn}(\text{Box}, \text{City})$ ,  $\text{BoxOnTruck}(\text{Box}, \text{Truck})$ ,  $\text{TruckIn}(\text{Truck}, \text{City})$ ,  $\text{PlaneIn}(\text{Plane}, \text{City})$ ,  $\text{BoxOnPlane}(\text{Box}, \text{Plane})$ . In this framework, generic action templates such as loading or unloading a box from a truck or plane or driving trucks and flying planes between cities are likely to apply generically to domain objects and thus the planning problem can be specified independently of any ground domain instantiation.

One recent language for representing relational probabilistic planning problems is PPDDL [89]. At its core, PPDDL is a probabilistic extension of a subset of PDDL conforming to the deterministic ADL planning language [58]; ADL, in turn, introduced universal and conditional effects into the STRIPS representation [29]. To see the compactness of a relational representation, we provide a (P)PPDDL representation of the BoxWORLD problem in Fig. 2 where for simplicity, we omit the *Plane* class of objects and associated actions and relations and abbreviate  $\text{BoxOnTruck}(\text{Box} : b, \text{Truck} : t)$  as  $\text{BoxOn}(\text{Box} : b, \text{Truck} : t)$ .

General PPDDL specifications can be more compact for some problems than the PPDDL subset we refer to in this article. For example, in general PPDDL, universal and conditional effects and probabilities can be arbitrarily nested, thus allowing for exponentially more compact representations of probabilistic action effects than can be achieved with probabilities only at the top-level of effects [66]. In addition, there are some general PPDDL specifications that *cannot* be translated to the PPDDL subset described here. If the general PPDDL specification uses probabilistic effects *nested under* universal effects (e.g., each box falls off a truck with some independent probability), it is generally impossible to translate such a problem to the restricted PPDDL subset used here because it requires an indefinitely factored transition probability model that cannot be expressed with finite probability specifications restricted to the top level of effects. While we do not discuss such model-expressivity here, we refer the reader to Sanner and Boutilier [70] and Chapter 6 of Sanner [67] for a treatment of such issues in first-order MDPs.

While the meaning of the PPDDL representation in Fig. 2 is intended to be relatively straightforward, there are a few important points that should be explained. First, we assume that actions can be executed in all states so we do not encode explicit preconditions. While this assumption is not necessary, it does not have any effect on the value of an optimal policy in a domain that already has a *noop* action and it helps simplify our later notation. When an action executes, each probabilistic effect is realized independently according to the specified probability. For example, the *unload* action realizes its effects only 90% of the time, whereas the *drive* action deterministically realizes its effects on each execution.

Probabilistic effects at the top-level of the effect specification consist of conjunctions of effects. Each individual effect can be *universal* and *conditional*. Universal effects denoted by universally quantified variables in the *then* clause permit the effect to apply to an arbitrary number of objects not explicitly named in the action's parameter list. Conditional effects denoted by *when* can be arbitrary first-order formulae specifying that the effects listed in the *then* clause hold in the post-action state if the *when* conditions hold in the pre-action state. When universally quantified variables are shared between the *when/then* clause pair, we refer to such effects as *universal conditional*. We note that each individual effect is only allowed to mention one positive or negative relation in the *then* portion of the clause. A conjunction of *then* effects can be easily specified as multiple effects with the same *when* condition. Disjunctive (i.e., non-deterministic) effects are prohibited in PPDDL. For example, when the  $\text{load}(b, t)$  action is executed, its effects are realized with probability 0.9. When these effects are realized, then for any city  $c$  that satisfies  $\text{BoxIn}(b, c) \wedge \text{TruckIn}(t, c)$  in the pre-action state,  $\text{BoxOn}(b, t) \wedge \neg \text{BoxIn}(b, c)$  will hold in

- 
- **Domain Object Instantiation:**
    - $\text{Box} = \{\text{box}_1, \text{box}_2, \text{box}_3\}$ ,  $\text{Truck} = \{\text{truck}_1, \text{truck}_2\}$ ,  $\text{City} = \{\text{paris}, \text{berlin}, \text{rome}\}$
  - **(S)ate-variable Atoms (i.e., binary state variables):**
    - $\text{BoxIn}$ :  
 $\{\text{BoxIn}(\text{box}_1, \text{paris}), \text{BoxIn}(\text{box}_2, \text{paris}), \text{BoxIn}(\text{box}_3, \text{paris}),$   
 $\text{BoxIn}(\text{box}_1, \text{berlin}), \text{BoxIn}(\text{box}_2, \text{berlin}), \text{BoxIn}(\text{box}_3, \text{berlin}),$   
 $\text{BoxIn}(\text{box}_1, \text{rome}), \text{BoxIn}(\text{box}_2, \text{rome}), \text{BoxIn}(\text{box}_3, \text{rome})\}$
    - $\text{TruckIn}$ :  
 $\{\text{TruckIn}(\text{truck}_1, \text{paris}), \text{TruckIn}(\text{truck}_1, \text{berlin}), \text{TruckIn}(\text{truck}_1, \text{rome}),$   
 $\text{TruckIn}(\text{truck}_2, \text{paris}), \text{TruckIn}(\text{truck}_2, \text{berlin}), \text{TruckIn}(\text{truck}_2, \text{rome})\}$
    - $\text{BoxOn}$ :  
 $\{\text{BoxOn}(\text{box}_1, \text{truck}_1), \text{BoxOn}(\text{box}_2, \text{truck}_1), \text{BoxOn}(\text{box}_3, \text{truck}_1),$   
 $\text{BoxOn}(\text{box}_1, \text{truck}_2), \text{BoxOn}(\text{box}_2, \text{truck}_2), \text{BoxOn}(\text{box}_3, \text{truck}_2)\}$
  - **(A)ctions:**
    - $\text{load}$ :  
 $\{\text{load}(\text{box}_1, \text{truck}_1), \text{load}(\text{box}_2, \text{truck}_1), \text{load}(\text{box}_3, \text{truck}_1)$   
 $\text{load}(\text{box}_1, \text{truck}_2), \text{load}(\text{box}_2, \text{truck}_2), \text{load}(\text{box}_3, \text{truck}_2)\}$
    - $\text{unload}$ :  
 $\{\text{unload}(\text{box}_1, \text{truck}_1), \text{unload}(\text{box}_2, \text{truck}_1), \text{unload}(\text{box}_3, \text{truck}_1),$   
 $\text{unload}(\text{box}_1, \text{truck}_2), \text{unload}(\text{box}_2, \text{truck}_2), \text{unload}(\text{box}_3, \text{truck}_2)\}$
    - $\text{drive}$ :  
 $\{\text{drive}(\text{truck}_1, \text{paris}), \text{drive}(\text{truck}_1, \text{berlin}), \text{drive}(\text{truck}_1, \text{rome})$   
 $\text{drive}(\text{truck}_2, \text{paris}), \text{drive}(\text{truck}_2, \text{berlin}), \text{drive}(\text{truck}_2, \text{rome})\}$
  - **(T)ransition Function:**  
 Follows directly from ground instantiation of PPDDL actions in Fig. 2.
  - **(R)eward:**  
 $\text{if } [\text{BoxIn}(\text{box}_1, \text{paris}) \vee \text{BoxIn}(\text{box}_2, \text{paris}) \vee \text{BoxIn}(\text{box}_3, \text{paris})] \text{ then } 10 \text{ else } 0$
- 

Fig. 3. One possible ground MDP instantiation of the BoxWORLD FOMDP.

the post-action state since both effects have equivalent *when* conditions. When these effects are not realized on 10% of the  $\text{load}(b, t)$  executions, no state changes occur and it is equivalent to a *noop* action.

One can easily see that this relationally specified domain-independent specification allows very compact MDP specifications when compared to a corresponding ground factored MDP representation. For example, consider instantiating the PPDDL problem in Fig. 2 to the ground factored MDP representation in Fig. 3 where we assume a problem instance with a domain instantiation of three boxes, three cities, and two trucks. While this is a trivially small domain instantiation, we note that its factored MDP representation requires 21 propositional atoms corresponding to over two million distinct states and 18 distinct actions that can be executed in each state. And the reward, which uses existential quantification in the relational PPDDL specification must be grounded to obtain the corresponding factored MDP representation. Clearly, for  $n$  objects, the grounded factor for the formula  $\exists \text{Box} : b. \text{BoxIn}(b, \text{paris})$  will contain  $|\text{Box}|$  state variables, but if the reward were changed to  $\forall \text{City} : c \exists \text{Box} : b. \text{BoxIn}(b, c)$ , the ground reward representation would contain  $|\text{Box}| \cdot |\text{City}|$  state variables—thus implying a combinatorial growth in the number of nested quantifiers.

In general, the number of ground atoms for a factored MDP representation will scale linearly in the number of relations, exponentially in the arity of each relation (assuming more than one domain object), and polynomially in the number of domain objects that fill each relation argument. To see this, let us assume for simplicity that all object class instantiations have  $k$  instances. Then a single unary relation would be represented by  $k$  ground atoms, a binary relation by  $k^2$  atoms, and an  $n$ -ary relation by  $k^n$  atoms. Similarly, the size of the grounding of any quantified formula is exponential in the number of nested quantifiers, linear in the number of relations, and exponential in the size of the domain object classes being quantified. Assuming  $k$  instances for all object classes and  $q$  nested (non-vacuous) quantifiers over formulae containing  $r$  relations, the resulting unsimplified ground representation of the formula would require  $rk^q$  ground atoms.

For sufficiently small predicate arities and levels of quantifier nesting (assuming these remain constant for a problem as the domain size varies), the space requirements for representing a ground MDP may be acceptable. Thus, if we have adequate space to permit the grounding of a relational MDP to obtain a factored MDP and we have the time to find an optimal solution to this factored MDP, then grounding gives us one approach to representing and solving relational MDPs for specific domain instances. However we note that while solving MDPs exactly is known to be polynomial in the number of states (see Section 2.2.2), the number of states is exponential in the number of ground atoms in a factored representation. This is Bellman's [6] well-known curse of dimensionality and since the number of ground atoms is at least linear in domain size, it implies that the exact solution methods discussed previously require time at least exponential in the domain size. This precludes the general possibility of exact solutions to grounded relational MDPs for all but the smallest domain sizes. While this suggests the use of approximation methods for solving grounded MDPs, there are useful lifted alternatives to representing and solving relational MDPs that we discuss next.

- 
- if  $(\exists b. \text{BoxIn}(b, \text{paris}))$   
then do *noop* (value = 100.00)
  - else if  $(\exists b^*, t^*. \text{TruckIn}(t^*, \text{paris}) \wedge \text{BoxOn}(b^*, t^*))$   
then do *unload*( $b^*, t^*$ ) (value = 89.0)
  - else if  $(\exists b, c, t^*. \text{BoxOn}(b, t^*) \wedge \text{TruckIn}(t, c))$   
then do *drive*( $t^*, \text{paris}$ ) (value = 80.0)
  - else if  $(\exists b^*, c, t^*. \text{BoxIn}(b^*, c) \wedge \text{TruckIn}(t^*, c))$   
then do *load*( $b^*, t^*$ ) (value = 72.0)
  - else if  $(\exists b, c_1^*, t^*, c_2. \text{BoxIn}(b, c_1^*) \wedge \text{TruckIn}(t^*, c_2))$   
then do *drive*( $t^*, c_1^*$ ) (value = 64.7)
  - else do *noop* (value = 0.0)
- 

**Fig. 4.** A decision-list representation of the expected discounted reward value for an exhaustive partitioning of the state space in the BoxWORLD problem. The optimal action is also shown for each partition where the optimal bindings of the action variables (denoted by a  $*$ ) correspond to any binding satisfying those variable names in the state formula.

### 3.1.2. Grounded vs. lifted solutions

In contrast to the grounded approach to representing relational MDPs as factored MDPs, it is important to point out that no matter how many domain objects there may be in an actual problem instance, the size of the PPDDL relational planning problem specification in Fig. 2 remains constant. Consequently, this invites the following question: if we can avoid a domain-dependent blowup in the representation of a relational MDP as in PPDDL, can we avoid a domain-dependent blowup in its solution too? Although we have yet to discuss the specifics of how we might find a domain independent solution to this PPDDL representation, in Fig. 4 we provide an optimal domain-independent value function and its corresponding policy for the relational PPDDL specification of the BoxWORLD problem in Fig. 2 (using discount factor  $\gamma = 0.9$ ).

The key features to note here are the state and action abstraction in the value and policy representation that are afforded by the first-order specification and solution of the problem. That is, this solution does not refer to any specific set of domain objects, say just  $\text{City} = \{\text{paris}, \text{berlin}, \text{rome}\}$ , but rather it provides a solution for *all possible domain object instantiations*. And while the BoxWORLD problem could not be represented as a grounded factored MDP for sufficiently large domain instantiations, much less solved, a domain-independent solution to this particular problem is quite simple and applies to domain instances of any size due to the power of state and action abstraction afforded by the first-order logical representation.

Thus, an alternative idea to grounding a relational MDP specification and solving it for a particular domain instance is to translate the PPDDL relational specification to a first-order MDP representation that is directly amenable to solutions via lifted symbolic dynamic programming. This approach obtains a solution that applies universally to all possible domain instantiations and has a time complexity that is independent of domain size. As we will see, the power of this lifted style of solution is that it exploits the existence of domain objects, relations over these objects, and the ability to express objectives and action effects using quantification.

## 3.2. Situation calculus background

Before we present the first-order MDP (FOMDP) formalism, we discuss the basics of the situation calculus, which in turn provides the logical foundations for our FOMDP representation. We begin by describing the necessary background material from the situation calculus and Reiter's default solution to the frame problem [64] required to understand FOMDPs. This includes a discussion of the basic ingredients of the situation calculus formulation: actions, situations, and fluents along with relevant axioms (e.g., unique names for actions and domain-specific axioms). Next we introduce effect axioms and explain how these can be derived from a PDDL specification. Then we show how effect axioms can be compiled into the successor-state axioms that underly the default solution to the frame problem of the situation calculus. We conclude by introducing the regression operator *Regr* that will prove crucial to our symbolic dynamic programming solution to first-order MDPs.

### 3.2.1. Basic ingredients

The situation calculus is a first-order language for axiomatizing dynamic worlds [52]. Its basic language elements consist of actions, situations and fluents:

- **Actions:** Actions are first-order terms consisting of an action function symbol and arguments. For example, an action for loading box  $b$  on truck  $t$  in the running BoxWORLD example is represented by  $\text{load}(b, t)$ .
- **Situations:** A situation is a first-order term denoting a specific state. The initial situation is usually denoted by  $s_0$  and subsequent situations resulting from action executions are obtained by applying the *do* function,  $\text{do}(a, s)$  representing the situation resulting from executing action  $a$  in situation  $s$ . For example, the situation resulting from loading box  $b$  on truck  $t$  in the initial situation  $s_0$  and then driving truck  $t$  to city  $c$  is given by the term  $\text{do}(\text{drive}(t, c), \text{do}(\text{load}(b, t), s_0))$ .
- **Fluents:** A fluent is a relation whose truth value varies from situation to situation. A fluent is simply a relation whose last argument is a situation term. For example, imagine an initial state  $s_0$  in which fluent  $\text{BoxOn}(b, t, s_0)$  is false, but fluents  $\text{TruckIn}(t, c, s_0)$  and  $\text{BoxIn}(b, c, s_0)$  are true. Then under the semantics of a deterministic version of the  $\text{load}(b, t)$



action (which we formally define in a moment),  $\text{BoxOn}(b, t, \text{do}(\text{load}(b, t), s_0))$  holds. We do not consider functional fluents in this exposition, but they are easily added to the language without adverse computational side effects [64].

### 3.2.2. From PDDL to a first-order logic domain theory

To axiomatize a PDDL domain theory in first-order logic, we must first consider how to describe the effects and non-effects of actions. We can begin by describing positive and negative effect axioms that characterize how fluents change as a result of actions. Note that in the following presentation, all relations that can change between states in PPDDL have been rewritten as fluents with an extra situation term. In addition, we assume all axioms are implicitly universally quantified.

- **Positive Effect Axioms:** positive effect axioms state which actions can explicitly make each fluent true; for example:

$$\begin{aligned} [\exists c. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] &\supset \text{BoxOn}(b, t, \text{do}(a, s)) \\ [\exists t. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] &\supset \text{BoxIn}(b, c, \text{do}(a, s)) \\ [\exists c_1. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] &\supset \text{TruckIn}(t, c, \text{do}(a, s)) \end{aligned}$$

- **Negative Effect Axioms:** negative effect axioms state which actions can explicitly make each fluent false; for example:

$$\begin{aligned} [\exists c. a = \text{unload}(b, t) \wedge \text{BoxOn}(b, t, s) \wedge \text{TruckIn}(t, c, s)] &\supset \neg \text{BoxOn}(b, t, \text{do}(a, s)) \\ [\exists t. a = \text{load}(b, t) \wedge \text{BoxIn}(b, c, s) \wedge \text{TruckIn}(t, c, s)] &\supset \neg \text{BoxIn}(b, c, \text{do}(a, s)) \\ [\exists c. a = \text{drive}(t, c) \wedge \text{TruckIn}(t, c_1, s)] &\supset \neg \text{TruckIn}(t, c_1, \text{do}(a, s)) \end{aligned}$$

In general, positive and negative effect axioms can be specified by considering all of the ways in which each action can affect each fluent. Fortunately, these axioms are easy to derive directly from the PDDL representation given in Fig. 2. In fact, one can verify that these effect axioms are simply syntactic rewrites of the PDDL effects where we have made the following transformations:

- (1) The action name from the PDDL effect is placed in an equality on the LHS of the  $\supset$ .
- (2) All universal quantifiers for universal effects are dropped as all unquantified variables are assumed to be universally quantified in the effect axioms.
- (3) The *when* conditions of the PDDL effect are conjoined on the LHS of the  $\supset$  with all fluents specified in terms of the situation  $s$ .
- (4) The *then* portion of the effect (which should be a single literal) is placed on the RHS of the  $\supset$  and is parameterized by the post-action situation  $\text{do}(a, s)$ . Whether the literal is negated or non-negated respectively determines whether the resulting axiom should be negative or positive.
- (5) Any free variables appearing only on the LHS of the  $\supset$  and not appearing free in the action term are explicitly existentially quantified in the LHS.

This takes care of specifying *what changes*, however we have not provided any axioms for specifying *what does not change*, i.e., the so-called *frame axioms*. Obviously, if we want to prove anything useful in our theory, we have to specify frame axioms. Otherwise, we would never be able to infer the properties of a successor or predecessor state for an action as simple as a *noop*. However, specifying exactly what does not change in a *compact* manner has been an extremely difficult problem to solve for the situation calculus—this is, of course, the infamous *frame problem*.

An especially elegant solution to the *frame problem* is that proposed by [63]. In this solution, we specify all positive and negative effects for a fluent, which conveniently, we have just done in our translation from PDDL above. We use the following normal form for positive effect axioms where  $F$  is a fluent and  $\gamma_F^+(\vec{x}, a, s)$  represents a first-order formula that, if true in  $s$ , results in  $F(\vec{x}, \text{do}(a, s))$  being true after action  $a(\vec{x})$  is executed in situation  $s$ :

$$\gamma_F^+(\vec{x}, a, s) \supset F(\vec{x}, \text{do}(a, s)) \quad (12)$$

Likewise, we use the following normal form for negative effect axioms where  $\gamma_F^-(\vec{x}, a, s)$  represents a first-order formula that if true in  $s$ , results in  $F(\vec{x}, \text{do}(a, s))$  being false after action  $a(\vec{x})$  is executed in situation  $s$ :

$$\gamma_F^-(\vec{x}, a, s) \supset \neg F(\vec{x}, \text{do}(a, s)) \quad (13)$$

We note that the potential difference between our previous presentation of positive and negative effect axioms and this normal form is that there is exactly *one* positive effect axiom for each positive fluent and *one* negative effect axiom for each negative fluent. This just happens to be the case in our example, but if it were otherwise, we could use the simple logical equivalence

$$[(C_1 \supset F) \wedge (C_2 \supset F)] \equiv [(C_1 \vee C_2) \supset F] \quad (14)$$

to rewrite any set of effect axioms derived from the PDDL subset of PPDDL into this normal form.

Next, we need to add in *unique name axioms* for all pairs of distinct action names  $A$  and  $B$  stating that

$$A(\vec{x}) \neq B(\vec{y}) \quad (15)$$

and also that identical actions have identical arguments:

$$A(x_1, \dots, x_k) = A(y_1, \dots, y_k) \supset x_1 = y_1 \wedge \dots \wedge x_k = y_k \quad (16)$$

From this normal form, unique names axioms, and *explanation closure axioms* that state these are the only effects that hold in our world model, Reiter showed that we can build *successor state axioms* (SSAs) that compactly encode both the effect and frame axioms for a fluent. The format of the successor state axiom for a fluent  $F$  is as follows:

$$\begin{aligned} F(\vec{x}, do(a, s)) &\equiv \Phi_F(\vec{x}, a, s) \\ &\equiv \gamma_F^+(\vec{x}, a, s) \vee F(\vec{x}, s) \wedge \neg \gamma_F^-(\vec{x}, a, s) \end{aligned} \quad (17)$$

For our running BoxWORLD example, we obtain the following SSAs:

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

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

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

While the notation might seem a bit cumbersome, the meaning of the axioms is quite intuitive. For example, the successor state axiom for  $BoxOn(b, t, \cdot)$  states that a box  $b$  is on a truck  $t$  after an action *iff* the action loaded box  $b$  on truck  $t$  or box  $b$  was already on truck  $t$  to begin with and the action did not unload it.

### 3.2.3. Regression

An important tool in the development of first-order MDPs is the ability to take a first-order state description  $\psi$  and “backproject” it through a deterministic action to see what conditions must have held prior to executing the action if  $\psi$  holds after executing the action. This is precisely the definition of *regression*. Fortunately, the SSAs lend themselves to a very natural specification definition of regression: if we want to regress a fluent  $F(\vec{x}, do(a, s))$  through an action  $a$ , we need only replace the fluent with its equivalent pre-action formula  $\Phi_F(\vec{x}, a, s)$ . In general, we can inductively define a regression operator  $Regr(\cdot)$  for all first-order formulae as follows [64]:

- $Regr(F(\vec{x}, do(a, s))) = \Phi_F(\vec{x}, a, s)$
- $Regr(\neg\psi) = \neg Regr(\psi)$
- $Regr(\psi_1 \wedge \psi_2) = Regr(\psi_1) \wedge Regr(\psi_2)$
- $Regr((\exists x)\psi) = (\exists x)Regr(\psi)$

Using the unique names assumption for actions and these regression rules, we can perform regression on any first-order logic formula. For example, if

$$\exists b. BoxIn(b, paris, do(unload(b^*, t^*), s))$$

holds then we can use the regression operator to determine what must have held in the pre-action situation  $s$ . Following is a derivation using the above rules:

$$\begin{aligned} &Regr(\exists b. BoxIn(b, paris, do(unload(b^*, t^*), s))) \\ &= \exists b. Repr(BoxIn(b, paris, do(unload(b^*, t^*), s))) \\ &= \exists b. \Phi_{BoxIn}(b, paris, unload(b^*, t^*), s) \\ &= \exists b. [[\exists t. unload(b^*, t^*) = unload(b, t) \wedge BoxOn(b, t, s) \wedge TruckIn(t, paris, s)] \\ &\quad \vee BoxIn(b, paris, s) \\ &\quad \wedge \neg [\exists t. unload(b^*, t^*) = load(b, t) \wedge BoxIn(b, paris, s) \wedge TruckIn(t, paris, s)]] \end{aligned}$$

At this point, we can use the unique names axioms for actions to simplify, and exploit rules for distributing quantifiers and renaming variables with respect to equality to obtain the following equivalent representation:

$$\begin{aligned} &= [\exists b, t. b = b^* \wedge t = t^* \wedge BoxOn(b, t, s) \wedge TruckIn(t, paris, s)] \\ &\quad \vee \exists b. BoxIn(b, paris, s) \\ &= [(\exists b. b = b^*) \wedge (\exists t. t = t^*) \wedge BoxOn(b^*, t^*, s) \wedge TruckIn(t^*, paris, s)] \\ &\quad \vee \exists b. BoxIn(b, paris, s) \end{aligned}$$

We will assume throughout the rest of this article that all object domains are non-empty.<sup>2</sup> This leads to the following fully simplified form of the regression:

$$\begin{aligned} \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, \text{do}(\text{unload}(b^*, t^*), s))) \\ = [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \vee \exists b. \text{BoxIn}(b, \text{paris}, s) \end{aligned} \quad (18)$$

This final result is very intuitive: it states that if there exists a box  $b$  in *paris* after unloading some box  $b^*$  from some truck  $t^*$ , then either the truck  $t^*$  was in *paris*, or a box was in *paris* to begin with.

### 3.3. FOMDP representation

Having defined the deterministic situation calculus translation of a simple PDDL model, we use this as a building block to obtain a first-order MDP (FOMDP) [14] from the restricted PPDDL syntax for relational MDPs that we introduced earlier. A FOMDP can be thought of as a universal MDP that abstractly defines the state, action, transition, and reward tuple  $\langle S, A, T, R \rangle$  for all possible domain instantiations (i.e., an infinite number of ground MDPs). In this subsection we formalize the building blocks of FOMDPs. We begin by introducing the *case* notation and operations and discuss the representation of the reward and value function as case statements. Then we describe how stochastic actions are represented by building on our previous situation calculus formalization. Once all of these components are defined, we will have everything needed to generalize the dynamic programming solution of MDPs from the ground case to the lifted case of symbolic dynamic programming for FOMDPs.

#### 3.3.1. Case representation of rewards, values, and probabilities

We introduce two useful variants of a *case notation* along with its logical definition to allow first-order specifications of the rewards, probabilities, and values required for FOMDPs:

$$\begin{aligned} (t = \text{case}[\phi_1, t_1; \dots; \phi_n, t_n]) &\equiv \left( t = \begin{array}{|c|} \hline \phi_1 : t_1 \\ \hline \vdots \\ \hline \phi_n : t_n \\ \hline \end{array} \right) \\ &\equiv \left( \bigvee_{i \leq n} \{\phi_i \wedge t = t_i\} \right) \end{aligned} \quad (19)$$

Here the  $\phi_i$  are *state formulae* where fluents in these formulae do not contain the term *do* and the  $t_i$  are terms. We note that in contrast to states, situations reflect the entire history of action occurrences. However, the specification of our FOMDP dynamics is Markovian and allows recovery of state properties from situation terms. For this reason, we can always represent the situation term using the free variable  $s$  without loss of generality. Often the  $t_i$  will be numerical constants and the  $\phi_i$  will partition state space.

We emphasize that the case notation for a logical formula (whether in the syntactic form  $t = \text{case}[\phi_1, t_1; \dots; \phi_n, t_n]$  or in the tabular form above) is simply a meta-logical notation used as a compact representation of the logical formula itself. In the meta-logical notation of cases, all formulae  $\phi_i$ , terms  $t_i$  and parameters of the case statement such as the situation term  $s$  refer to symbols of the underlying logical language. At a meta-logical level, a case statement may be viewed as a relation since the case “partition” formulae may overlap and may not be exhaustive. Case statements may be compared with (in)equalities and manipulated with arithmetic operations to produce other case statements (all at a meta-logical level).

To illustrate this notation concretely, we represent our BoxWORLD FOMDP reward function  $R(s)$  from our PPDDL representation in Fig. 2 as the following  $r\text{Case}(s)$  statement that reflects the immediate reward obtained in situation  $s$ :

$$r\text{Case}(s) = \begin{array}{|c|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) : 10 \\ \hline \neg \exists b. \text{BoxIn}(b, \text{paris}, s) : 0 \\ \hline \end{array} \quad (20)$$

For simplicity of presentation, we will assume the reward is not action dependent, but such dependencies can be introduced without difficulty. Throughout the text,  $R(s)$  will be used to represent a generic FOMDP reward case statement and  $r\text{Case}(s)$  will refer to the specific reward function. Thus, for BoxWORLD, we write  $R(s) = r\text{Case}(s)$  and wherever  $R(s)$  occurs, we can substitute the logical formula for  $r\text{Case}(s)$ .

Here we see that the first-order formulae in the case statement divide all possible ground states into two regions of constant-value: when there exists a box in *paris*, a reward of 10 is achieved, otherwise a reward of 0 is achieved. Likewise the value function  $V(s)$  that we derive through symbolic dynamic programming can be represented in exactly the same case format. Indeed,  $V^0(s) = R(s)$  in the first-order version of value iteration.

The case representation can also be used to specify transition probabilities (as we will see below). We first discuss the operations that can be performed on case statements.

<sup>2</sup> Logically, this requires a background theorem axiom for every object type *Sort* that states  $\exists o. \text{Sort}(o)$ . With this, we can use the simplification  $(\exists \text{Sort} : o. o = o^*) \supset \top$ .

### 3.3.2. Case operations

In this subsection we introduce various operations that can be applied to case statements providing both a formal logical definition and a graphical example that intuitively demonstrates the application of the case operation.

We begin by formally introducing the following binary  $\otimes$ ,  $\oplus$ , and  $\ominus$  operators on case statements [14]:

$$\text{case}[\phi_i, t_i : i \leq n] \otimes \text{case}[\psi_j, v_j : j \leq m] = \text{case}[\phi_i \wedge \psi_j, t_i \cdot v_j : i \leq n, j \leq m] \quad (21)$$

$$\text{case}[\phi_i, t_i : i \leq n] \oplus \text{case}[\psi_j, v_j : j \leq m] = \text{case}[\phi_i \wedge \psi_j, t_i + v_j : i \leq n, j \leq m] \quad (22)$$

$$\text{case}[\phi_i, t_i : i \leq n] \ominus \text{case}[\psi_j, v_j : j \leq m] = \text{case}[\phi_i \wedge \psi_j, t_i - v_j : i \leq n, j \leq m] \quad (23)$$

Intuitively, to perform an operation on case statements, we simply perform the corresponding operation on the cross-product of all case partitions of the operands. Letting each  $\phi_i$  and  $\psi_j$  denote generic first-order formulae, we can perform the “cross-sum”  $\oplus$  of case statements in the following manner:

$$\begin{array}{|c|c|} \hline \phi_1 : 10 \\ \hline \phi_2 : 20 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \psi_1 : 1 \\ \hline \psi_2 : 2 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \phi_1 \wedge \psi_1 : 11 \\ \hline \phi_1 \wedge \psi_2 : 12 \\ \hline \phi_2 \wedge \psi_1 : 21 \\ \hline \phi_2 \wedge \psi_2 : 22 \\ \hline \end{array}$$

Likewise, we can perform  $\ominus$ ,  $\otimes$ , and  $\max$  operations by, respectively, subtracting, multiplying, or taking the max of partition values. Note that for a binary operation involving a scalar and a case statement, a scalar value  $C$  may be viewed as  $\text{case}[\top, C]$  where  $\top$  is a tautology. We use the  $\oplus$  and  $\otimes$  operators to, respectively, denote summations and products of multiple case operands.

It is important to note that some partitions resulting from the application of the  $\oplus$ ,  $\ominus$ , and  $\otimes$  operators may be inconsistent; if we can identify such inconsistency, we simply discard such partitions. When the case partitions contain general first-order logic formulae, inconsistency detection is undecidable. However, for the symbolic dynamic programming algorithm discussed in this section, it is not required that all inconsistent partitions be discarded; failing to do so simply results in a non-minimal case representation that contains partitions not corresponding to any world state. In practice, we rely on time-limited incomplete theorem proving for inconsistency pruning.

We define a few additional operations on case statements, the first being the binary  $\cup$  operation:

$$\text{case}[\phi_i, t_i : i \leq n] \cup \text{case}[\psi_j, v_j : j \leq m] = \text{case}[\phi_1, t_1; \dots; \phi_n, t_n; \psi_1, v_1; \dots; \psi_m, v_m] \quad (24)$$

In this operation we simply construct the union of the partitions from each of the case statements; for example:

$$\begin{array}{|c|c|} \hline \phi_1 : 10 \\ \hline \phi_2 : 20 \\ \hline \end{array} \cup \begin{array}{|c|c|} \hline \psi_1 : 1 \\ \hline \psi_2 : 2 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \phi_1 : 10 \\ \hline \phi_2 : 20 \\ \hline \psi_1 : 1 \\ \hline \psi_2 : 2 \\ \hline \end{array}$$

Next we define two unary operations. The  $\exists \vec{x}. \text{case}(\vec{x})$  operation simply existentially quantifies the  $\text{case}(\vec{x})$  statement. Since  $\text{case}(\vec{x})$  is defined logically with a disjunction, we can distribute the  $\exists \vec{x}$  inside the disjunction:

$$\begin{aligned} \exists \vec{x}. \left( t = \begin{array}{|c|c|} \hline \phi_1(\vec{x}) : t_1 \\ \hline : : : \\ \hline \phi_n(\vec{x}) : t_n \\ \hline \end{array} \right) &\equiv \exists \vec{x}. \bigvee_{i \leq n} \{ \phi_i(\vec{x}) \wedge t = t_i \} \\ &\equiv \bigvee_{i \leq n} \{ \exists \vec{x}. \phi_i(\vec{x}) \wedge t = t_i \} \\ &\equiv \left( t = \begin{array}{|c|c|} \hline \exists \vec{x}. \phi_1(\vec{x}) : t_1 \\ \hline : : : \\ \hline \exists \vec{x}. \phi_n(\vec{x}) : t_n \\ \hline \end{array} \right) \end{aligned} \quad (25)$$

Normally we assume an implicit “ $t =$ ” for a case statement but show it above for logical clarity.

The second unary operation is denoted “casemax” (and not “max”) since it produces a case statement as opposed to a single numerical value. The result of casemax is a case statement where the maximal possible value of its case argument is assigned to each region of state space in the resulting case statement. Assuming that the case partitions are pre-sorted such that  $t_i > t_{i+1}$  and all partitions of equal value have been disjunctively merged we can formally define this operation as follows:

$$\text{casemax } \text{case}[\phi_1, t_1; \dots; \phi_n, t_n] = \text{case} \left[ \phi_i \wedge \bigwedge_{j < i} \neg \phi_j, t_i : i \leq n \right] \quad (26)$$

Following is a more intuitive graphical exposition of the same casemax operation:

$$\text{casemax} \begin{array}{|c|} \hline \phi_1 : t_1 \\ \hline \phi_2 : t_2 \\ \hline \vdots \\ \hline \phi_n : t_n \\ \hline \end{array} = \begin{array}{|c|} \hline \phi_1 : t_1 \\ \hline \phi_2 \wedge \neg \phi_1 : t_2 \\ \hline \vdots \\ \hline \phi_n \wedge \neg \phi_1 \wedge \neg \phi_2 \wedge \dots \wedge \neg \phi_{n-1} : t_n \\ \hline \end{array}$$

One can easily verify that if the partitions are sorted from the highest value  $t_1$  to the lowest  $t_n$ , then the highest value consistent with any state formula in the input case statement is assigned to the unique partition consistent with that state formulae in the resulting case statement. (If the  $\phi_i$  in the input are mutually exclusive, then the casemax results in a case statement logically equivalent to the original.) The application of casemax requires constructing new partition formulae, up to  $n$  times the length of the original formulae for a case statement with  $n$  partitions. Fortunately, the use of inconsistency detection discussed previously and first-order ADDs (FOADD) that we introduce in the next section will mitigate the impact of this blowup by respectively pruning inconsistent case partitions and simplifying the representation of case formulae.

It is important to point out that all of the case operators are purely symbolic in that the  $t_i$  case partition values are not necessarily restricted to constant numerical values, but can be arbitrary symbolic (possibly state-dependent) terms [14]. However, the casemax operator (as defined here) implicitly requires an ordering on the  $t_i$ . We assume for the rest of this section that the case values are numeric rather than symbolic, and apply the natural  $<$  operator for our ordering.

### 3.3.3. Stochastic actions and transition probabilities

To state the FOMDP transition function for an action, stochastic “agent” actions are decomposed into a *collection* of deterministic actions, each corresponding to a possible outcome of the stochastic action. We then use a case statement to specify a distribution according to which “Nature” may choose a deterministic action from this set whenever the stochastic action is executed. As a consequence we need only formulate SSAs using the deterministic *Nature’s choices* [2,15,59,64], thus obviating the need for a special treatment of stochastic actions in SSAs.

Letting  $A(\vec{x})$  be a stochastic action with Nature’s choices (i.e., deterministic actions)  $n_1(\vec{x}), \dots, n_k(\vec{x})$ , we represent the probability of  $n_i(\vec{x})$  given  $A(\vec{x})$  is executed in  $s$  by  $P(n_j(\vec{x}), A(\vec{x}), s)$ . Continuing with the translation of our simple PPDDL example, we note that the  $load(b, t)$  action has one set of effects that occurs with probability 0.9. We use the deterministic action  $loadS(b, t)$  to denote the successful occurrence of these effects, and we let the deterministic action  $loadF(b, t)$  denote the failure of these effects to execute. To do this, we must redefine our SSAs from the previous PDDL case: now  $load(b, t)$  is a stochastic action executed by the agent with  $loadS(b, t)$  and  $loadF(b, t)$  being possible outcomes (i.e., deterministic actions chosen by Nature). Similarly, we interpret the other two actions using  $unloadS(b, t)/unloadF(b, t)$  as the two deterministic outcomes for  $unload(b, t)$ , and  $driveS(t, c)/driveF(t, c)$  as the two deterministic outcomes for  $drive(t, c)$ . For completeness and correctness, we redefine our SSAs for BoxWORLD in terms of these new deterministic actions for the BoxWORLD FOMDP:

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

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

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

Here, we have simply replaced our previous deterministic action names from the PDDL version with the deterministic *success* versions of Nature’s choice actions that we will use in our FOMDP. Note that since the “failure” versions of the actions correspond to the “no effects” case, they obviously do not play any role in the SSAs. The frame assumption present in the SSAs ensures that what was not explicitly changed remains the same.

We can now specify a distribution  $P(n_j(\vec{x}), A(\vec{x}), s)$  over Nature’s choice deterministic outcome using case statements to specify families of distributions, where the partitions in the case statements correspond to different classes of states and stochastic action parameters on which the distributions are conditioned. We denote specific instances of  $P(n_j(\vec{x}), A(\vec{x}), s)$  with the case statement  $pCase(n_j(\vec{x}), A(\vec{x}), s)$  where  $\top$  is a tautology, for example:

$$\begin{aligned} pCase(loadS(b, t), load(b, t), s) &= \boxed{\top : 0.9} \\ pCase(loadF(b, t), load(b, t), s) &= \boxed{\top : 0.1} \end{aligned}$$

$$pCase(unloadS(b, t), unload(b, t), s) = \boxed{\top : 0.9} \quad (27)$$

$$pCase(unloadF(b, t), unload(b, t), s) = \boxed{\top : 0.1} \quad (28)$$

$$pCase(driveS(b, t), drive(b, t), s) = \boxed{\top : 1.0}$$

$$pCase(driveF(b, t), drive(b, t), s) = \boxed{\top : 0.0}$$

The above axiomatization does not fully illustrate the power of the FOMDP representation in that the probabilities are not state or action dependent, so we briefly digress to demonstrate a slightly more interesting variant. Suppose that the success of driving a truck to a city depends on whether the truck contains a box  $b$  with volatile material denoted by the predicate  $Volatile(b)$ . Then we can specify the family of distributions over Nature's choices for this stochastic action as follows:

$$\begin{aligned} pCase(driveS(t, c), drive(t, c), s) &= \begin{array}{|l|} \hline \exists b.BoxOn(b, t, s) \wedge Volatile(b) : 0.9 \\ \hline \neg(\exists b.BoxOn(b, t, s) \wedge Volatile(b)) : 1.0 \\ \hline \end{array} \\ pCase(driveF(t, c), drive(t, c), s) &= \begin{array}{|l|} \hline \exists b.BoxOn(b, t, s) \wedge Volatile(b) : 0.1 \\ \hline \neg(\exists b.BoxOn(b, t, s) \wedge Volatile(b)) : 0.0 \\ \hline \end{array} \end{aligned}$$

Here we see the transition probability of  $drive(t, c)$  can be easily conditioned on state properties of  $s$  and action parameters  $t$  and  $c$ .

It is important to note that the probabilities over all deterministic Nature's choices for a stochastic action sum to one:

$$\bigoplus_{j=1}^k P(n_j(\vec{x}), A(\vec{x}), s) = \boxed{\top : 1} ; \quad \forall \vec{x}, s$$

In addition, each  $P(n_j(\vec{x}), A(\vec{x}), s)$  should be a disjoint partitioning of state space such that no two case partitions ambiguously assign multiple probabilities to the same state. These two properties are crucial to having a well-defined probability distribution over all possible deterministic action outcomes for every possible state.

For this last example, the second property can be easily verified:

$$pCase(driveS(t, c), drive(t, c), s) \oplus pCase(driveF(t, c), drive(t, c), s) = \boxed{\top : 1}$$

### 3.4. Symbolic dynamic programming (SDP)

Symbolic dynamic programming (SDP) [14] is a dynamic programming solution to FOMDPs that produces a logical case description of the optimal value function. This is achieved through the symbolic operations of first-order decision-theoretic regression and maximization that perform the traditional dynamic programming Bellman backup at an abstract level without explicit enumeration of either the state or action spaces of the FOMDP. Among many possible applications, the use of SDP leads directly to a domain-independent value iteration solution to FOMDPs.

We will assume a constant numerical representation of values in order to explicitly perform the casemax during SDP in this article. However, we note that an appropriate generalization of casemax (cf., Chapter 6 of [67]) along with *Regr* of functional fluents [64] allows the definitions covered here to apply to general symbolic value representations using general terms rather than constants, hence the original use of “symbolic” in the name of the SDP algorithm.

#### 3.4.1. First-order decision-theoretic regression

Suppose we are given a value function  $V(s)$ . The first-order decision-theoretic regression (FODTR) [14] of this value function through an action  $A(\vec{x})$  yields a case statement containing the logical description of states and values that would give rise to  $V(s)$  after doing action  $A(\vec{x})$ . This is analogous to classical goal regression, the key difference being that action  $A(\vec{x})$  is stochastic. In MDP terms, the result of FODTR roughly corresponds to a Q-function (albeit one with free variables for the action parameters), which corresponds to the first half of a Bellman backup operation given in Eq. (6).<sup>3</sup>

We define the *first-order decision theoretic regression (FODTR)* as the situation calculus analog of Eq. (6) where we note that different successor states only arise through different Nature's choice deterministic actions:

$$FODTR[V(s), A(\vec{x})] = R(s) \oplus \gamma \cdot \left[ \bigoplus_{j=1}^k \{P(n_j(\vec{x}), A(\vec{x}), s) \otimes V(do(n_j(\vec{x}), s))\} \right] \quad (29)$$

<sup>3</sup> We do not use an action dependent reward  $R(s, A(\vec{x}))$ , but could substitute it if needed.

*FODTR* uses a meta-logical notation that takes as arguments  $V(s)$  representing the logical case statement for a value function with situation variable  $s$  and a parameterized stochastic action term  $A(\vec{x})$  with free variables  $\vec{x}$ . All subsequently defined operations on case statements in this article will be defined analogously.

The only problem with the  $FODTR[V(s), A(\vec{x})]$  operation as currently defined is that the formula  $V(do(n_j(\vec{x}), s))$  refers not to the current situation  $s$ , but to the future situation  $do(n_j(\vec{x}), s)$ , but this is easily remedied with regression:

$$FODTR[V(s), A(\vec{x})] = R(s) \oplus \gamma \cdot \left[ \bigoplus_{j=1}^k \{P(n_j(\vec{x}), A(\vec{x}), s) \otimes \text{Regr}(V(do(n_j(\vec{x}), s)))\} \right] \quad (30)$$

This is equivalent to the *FODTR* operation in Eq. (29) since the *Regr* operation preserves equivalence (by definition). Also on account of the equivalence preserving properties of *Regr*, we note that if  $V(s)$  partitions the state space then so must the resulting case statement for  $FODTR[V(s), A(\vec{x})]$ . Thus, from a logical description of  $V(s)$  we can derive one for its decision-theoretic regression  $FODTR[V(s), A(\vec{x})]$ . This is key to avoiding state and action enumeration in dynamic programming.

We denote an instance of the value function  $V(s)$  by the case statement  $vCase(s)$ . As defined previously, we also assume that the reward function  $R(s)$  and instances of Nature's choice probabilities  $P(n_j(\vec{x}), A(\vec{x}), s)$  are denoted by  $rCase(s)$  and  $pCase(n_j(\vec{x}), A(\vec{x}), s)$ , respectively.

As an example, let us compute the *FODTR* for  $vCase(s) = rCase(s)$  through the stochastic action  $A(\vec{x}) = \text{unload}(b^*, t^*)$  where  $rCase(s)$  is the *BoxWorld* reward as previously defined in Eq. (20). Since  $vCase(s)$  is logically defined, we can push the *Regr* operator into the individual  $vCase(s)$  partitions as follows:

$$\begin{aligned} FODTR[vCase(s), \text{unload}(b^*, t^*)] \\ = rCase(s) \oplus \gamma \left[ \bigoplus_{j=1}^k \left\{ pCase(n_j(\vec{x}), \text{unload}(b^*, t^*), s) \otimes \begin{array}{|l|} \hline \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(n_j(\vec{x}), s))) : 10 \\ \hline \text{Regr}(\neg \exists b. \text{BoxIn}(b, \text{paris}, do(n_j(\vec{x}), s))) : 0 \\ \hline \end{array} \right\} \right] \end{aligned}$$

Now, since the stochastic action  $A(\vec{x}) = \text{unload}(b^*, t^*)$ , we know that Nature's deterministic action choices  $n_j(\vec{x})$  range over  $\text{unloadS}(b^*, t^*)$  and  $\text{unloadF}(b^*, t^*)$ . We now substitute the  $pCase$  definitions for the deterministic actions  $\text{unloadS}(b^*, t^*)$  and  $\text{unloadF}(b^*, t^*)$  from Eqs. (27) and (28), respectively, obtaining:

$$\begin{aligned} FODTR[vCase(s), \text{unload}(b^*, t^*)] = rCase(s) \oplus \\ \gamma \left[ \left\{ \begin{array}{|l|} \hline \top : 0.9 \\ \hline \end{array} \otimes \begin{array}{|l|} \hline \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unloadS}(b^*, t^*), s))) : 10 \\ \hline \text{Regr}(\neg \exists b. \text{BoxIn}(b, \text{paris}, do(\text{unloadS}(b^*, t^*), s))) : 0 \\ \hline \end{array} \right\} \right. \\ \left. \oplus \left\{ \begin{array}{|l|} \hline \top : 0.1 \\ \hline \end{array} \otimes \begin{array}{|l|} \hline \text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unloadF}(b^*, t^*), s))) : 10 \\ \hline \text{Regr}(\neg \exists b. \text{BoxIn}(b, \text{paris}, do(\text{unloadF}(b^*, t^*), s))) : 0 \\ \hline \end{array} \right\} \right] \end{aligned}$$

We have already computed  $\text{Regr}(\exists b. \text{BoxIn}(b, \text{paris}, do(\text{unloadS}(b^*, t^*), s)))$  from Eq. (18) where the deterministic  $\text{unload}(b^*, t^*)$  from the PDDL case has been renamed to  $\text{unloadS}(b^*, t^*)$ . And by the properties of *Regr*, we know that  $\text{Regr}(\neg \phi) = \neg \text{Regr}(\phi)$  so we can easily obtain the regression of the negated partition in  $rCase(s)$ . It is important to note that if  $rCase(s)$  partitioned the post-action state space, the *Regr* operator preserves this partitioning in the pre-action state space. We note that

$$\text{Regr}(\phi(\vec{x}, do(\text{unloadF}(b^*, t^*), s))) = \phi(\vec{x}, s)$$

can be easily derived since  $\text{unloadF}(b^*, t^*)$  has no effects and is thus equivalent to a *noop* action. Making these substitutions, explicitly multiplying in the action probabilities and discount factor  $\gamma = 0.9$ , and explicitly writing out  $rCase(s)$ , we obtain the following (where, for readability, we use  $\neg$  to denote the conjunction of the negation of *all* partitions above the given partition in the case statement):

$$\begin{aligned} FODTR[vCase(s), \text{unload}(b^*, t^*)] = & \begin{array}{|l|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) : 10 \\ \hline \neg : 0 \\ \hline \end{array} \\ \oplus & \begin{array}{|l|} \hline [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] \\ \hline \vee \exists b. \text{BoxIn}(b, \text{paris}, s) : 8.1 \\ \hline \neg : 0 \\ \hline \end{array} \\ \oplus & \begin{array}{|l|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) : 0.9 \\ \hline \neg : 0 \\ \hline \end{array} \end{aligned}$$

Finally, explicitly carrying out the  $\oplus$ 's and simplifying yields the final result:

$$FODTR[vCase(s), \text{unload}(b^*, t^*)] = \begin{array}{|l|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) : 19.0 \\ \hline \neg \wedge [\text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] : 8.1 \\ \hline \neg : 0 \\ \hline \end{array} \quad (31)$$

The case statement resulting from *FODTR* contains free variables for the action parameters; in this case  $A(\vec{x}) = \text{unload}(b^*, t^*)$  so the free parameters are  $b^*$  and  $t^*$ . This result is intuitive: it states that if a box was already in *paris* then we get reward 19 (10 for the current reward and 9 for the discounted 1-step reward). Otherwise, if a box is not in *paris* in the current

state, but box  $b^*$  was on truck  $t^*$  in *paris* and the action was specifically  $unload(b^*, t^*)$ , then we get an expected future reward of 8.1 taking into account the success probability of unloading the box and the discount factor. Finally, if no box is in *paris* in the current state and we do not unload a box then we get 0 total reward.

This case statement represents the value of taking stochastic action  $unload(b^*, t^*)$  and acting so as to obtain the value given by  $rCase(s)$  thereafter. However, what we really need for symbolic dynamic programming is a logical description of a Q-function (recall Eq. (6)) that tells us the possible values that can be achieved for *any* action instantiation of  $b^*$  and  $t^*$ . This leads us to the following definition  $Q(A, s)$  of a first-order Q-function that makes use of the previously defined  $\exists \vec{x}$  unary case operator:

$$Q^t(A, s) = \exists \vec{x}. FODTR[V^{t-1}(s), A(\vec{x})] \quad (32)$$

We denote a specific instance of  $Q^t(A, s)$  by the case statement  $qCase^t(s, A)$ . We can think of  $qCase^t(s, A)$  as a logical description of the Q-function for action  $A(\vec{x})$  indicating the values that could be achieved by *any* instantiation of  $A(\vec{x})$ . By using the first-order case representation of states as well as action quantification via the  $\exists \vec{x}$  operation, FODTR effectively achieves *both action and state abstraction*.

Letting  $vCase^0(s) = rCase(s)$ , we can continue our running example to obtain a Q-function description for action  $unload$  where we have removed vacuous quantifiers. Technically,  $qCase^1(unload, s)$  would not be an exhaustive partitioning of the state space in that the 0 value partition from Eq. (32) is not the same one implied here from the  $\neg$  because the partition formulae above it have been quantified. However, throughout this article, we can exploit our assumption that all FOMDPs have a *noop* action to assume that the minimum value for any state is 0 (as opposed to being undefined). Thus we can always show the final 0 partition as  $\neg$  to indicate that any partitions not explicitly assigned a value by the above partitions are assigned a default value 0. Thus, we arrive at the following intuitive result:

$$qCase^1(unload, s) = \exists b^*, t^*. FODTR[vCase^0(s), unload(b^*, t^*)]$$

$\exists b. BoxIn(b, paris, s)$	: 19.0
$\exists b^*, t^*. [\neg \wedge BoxOn(b^*, t^*, s) \wedge TruckIn(t^*, paris, s)]$	: 8.1
$\neg$	: 0

In words, this states if the box was already in *paris* then we get a discounted reward of 19. Otherwise, if a box is not in *paris* in the current state, but there *exists* some box on a truck in *paris*, then we could unload it to get an expected discounted reward of 8.1. Finally, if there is no box on a truck to unload in *paris* and there is no box already in *paris* then we get 0 expected discounted reward. It is instructive to compare this description to the prior description of FODTR without existential action quantification—the difference is subtle, but important for action abstraction.

### 3.4.2. Symbolic maximization

At this point, we can decision-theoretically regress the value function through a *single* stochastic action to obtain a representation of its Q-function, but to complete the dynamic programming (Bellman backup) step in the spirit of Eq. (7) from Section 2, we need to know the maximum value that can be achieved by *any* action. For example, in the BoxWORLD FOMDP, our possible action choices are  $unload(b, t)$ ,  $load(b, t)$ , and  $drive(t, c)$  and our Q-function computations using Eq. (32) give us  $qCase^1(unload, s)$ ,  $qCase^1(load, s)$ , and  $qCase^1(drive, s)$ . In general, we will assume that we have  $m$  stochastic actions  $\{A_1(\vec{x}_1), \dots, A_m(\vec{x}_m)\}$  and a corresponding set of Q-functions  $\{qCase^t(A_1, s), \dots, qCase^t(A_m, s)\}$  derived from a common value function  $vCase^{t-1}(s)$ .

We might try to obtain a case description of the value function  $vCase^t(s)$  by simply applying the case  $\cup$  operator to merge all partitions of the Q-functions, i.e.,  $qCase^t(s, A_1) \cup \dots \cup qCase^t(s, A_m)$ . While this provides us with a description of possible values, it is not a value *function* because the state spaces of each Q-function may overlap, thus potentially assigning multiple values to the same underlying state. What we really want instead is to assign the *highest* possible value to each portion of state space. Fortunately, this is quite easy with the casemax operator. Thus we get the following equation for the symbolic maximization of Q-functions:

$$V^t(s) = \text{casemax}[Q^t(A_1, s) \cup \dots \cup Q^t(A_m, s)] \quad (33)$$

Recalling the way in which the casemax operation is computed from Eq. (26), every resulting instance  $vCase^t(s)$  of the value function  $V^t(s)$  will have the following case statement format where value case partition  $\psi_j$  corresponds to value  $v_j$  and  $v_i > v_{i+1}$ :

$$vCase^t(s) =$$

$\psi_1$	: $v_1$
$\psi_2 \wedge \neg \psi_1$	: $v_2$
$\vdots$	$\vdots$
$\psi_n \wedge \neg \psi_1 \wedge \neg \psi_2 \wedge \dots \wedge \neg \psi_{n-1}$	: $v_n$

This approach effectively gives us a decision-list representation of our value function (recall the optimal value function representation from Fig. 4). Thus, to determine the value for a state, we can simply traverse the list from highest to lowest value and take the value for the first case partition that is satisfied. The casemax operation guarantees that this value function will be a disjoint partitioning of the state space and our previous assumption that all actions are executable in all states ensures that this value function exhaustively assigns a value to all possible states (assuming  $vCase^{t-1}$  was exhaustive).



### 3.4.3. First-order value iteration

One should note that the SDP equations given here are exactly the lifted versions of the classical dynamic programming solution to MDPs given previously in Eqs. (6) and (7) from Section 2. Since those equations were used in part to define a value iteration algorithm, we can use the lifted versions to define a *first-order value iteration* algorithm where  $\epsilon$  is our error tolerance:

- (1) Initialize  $V^0(s) = R(s)$ ,  $t = 1$ .
- (2) Compute  $V^t(s)$  given  $V^{t-1}(s)$  using Eqs. (32) and (33).
- (3) If the following Bellman error inequality holds

$$\|V^t(s) \ominus V^{t-1}(s)\|_\infty \leq \frac{\epsilon(1-\gamma)}{2\gamma} \quad (34)$$

then terminate and return  $V^t(s)$ , otherwise go to step 2.

Here, we define  $\|V^t(s) \ominus V^{t-1}(s)\|_\infty$  as the maximal absolute value of any consistent partition in the case statement resulting from  $V^t(s) \ominus V^{t-1}(s)$ .

For example, applying first-order value iteration to the 0-stage-to-go value function (i.e.,  $vCase^0(s) = rCase(s)$ , given previously in Eq. (20)) yields the following simplified 1- and 2-stage-to-go value functions in the BoxWorld problem domain:

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

After sufficient iterations of first-order value iteration, the  $t$ -stage-to-go value function converges, giving the optimal value function (and as we derive in a moment, an optimal policy) from Fig. 4.

Boutilier et al. [14] provide a proof that SDP and thus every step of value iteration produces a correct logical description of the value function. From this, we can lift the error bounds from the ground MDP case in Eq. (8) to show domain-independent error bounds on the first-order abstracted value estimate:

**Corollary 3.4.1.** *Let  $V^*(s)$  be the optimal value function for a FOMDP. Terminating according to the criteria given in Step 3 of first-order value iteration guarantees  $\|V^t(s) - V^*(s)\|_\infty < \epsilon$  for any domain instantiation (even infinite) of the FOMDP.*

More generally, as a direct result of this corollary, we can derive domain-independent error bounds for the first-order representation of the value function produced by any first-order MDP algorithm (see Section 6 for other first-order algorithms).

**Corollary 3.4.2.** *Let  $\hat{V}(s)$  be an arbitrary first-order case representation of a value function. Let  $\hat{V}'(s)$  be the result of applying Eqs. (32) and (33) to  $\hat{V}(s)$  for a FOMDP. Let  $\epsilon = \frac{2}{1-\gamma} \|\hat{V}'(s) \ominus \hat{V}(s)\|_\infty$ . Then  $\|\hat{V}(s) - V^*(s)\|_\infty < \epsilon$  for any domain instantiation of the FOMDP.*

The difference of  $\gamma$  between the bounds of Corollaries 3.4.1 and 3.4.2 occurs because the former refers to a bound on  $V^t(s)$ , while the latter refers to a bound on  $\hat{V}(s) = V^{t-1}(s)$  and value iteration is known to contract the error by  $\gamma$  on each iteration.

### 3.4.4. Policy representation

Given a value function, it is important to be able to derive a first-order greedy policy representation from it, just as we did in the ground case in Section 2. This policy can then be used to directly determine actions to apply when acting in a ground instantiation of the FOMDP, or it can be used to define first-order versions of (approximate) policy iteration [69].

Fortunately, given a value function  $V(s)$ , it is easy to derive a greedy policy from it. Assuming we have  $m$  parameterized actions  $\{A_1(\vec{x}), \dots, A_m(\vec{x})\}$ , we can formally derive the policy  $\pi(s)[\cdot]$  using the  $\cdot$  to denote the value representation from which the policy is derived as follows (taking into account a few modifications to the case operators that we discuss in a moment):

$$\pi(s)[V(s)] = \text{casemax} \left( \bigcup_{i=1..m} \exists \vec{x}. FODTR[V(s), A_i(\vec{x})] \right) \quad (35)$$

We often refer to a specific instance of  $\pi(s)$  with the case statement  $\pi Case(s)$ . For bookkeeping, we require that each partition  $\langle \phi, t \rangle$  in  $\exists \vec{x}. FODTR[V(s), A_i(\vec{x})]$  maintain a mapping to the action  $A_i$  that generated it, which we denote as  $\langle \phi, t \rangle \rightarrow$

$A_i$ . Then, given a particular world state  $s$ , we can evaluate  $\pi\text{Case}(s)$  to determine which maximal policy partition  $\langle \phi, t \rangle \rightarrow A_i$  is satisfied by  $s$  and thus, which action  $A_i$  should be applied. If we retrieve the bindings of the existentially quantified action variables  $\exists \vec{x}$  in that satisfying policy partition, we can easily determine the parameterization of action  $A_i$  that should apply according to the policy.

To make this concrete, we derive a simple greedy policy for the BoxWorld FOMDP assuming the value function  $V(s) = r\text{Case}(s)$  and that we only have two actions  $\text{unload}(b^*, t^*)$  and  $\text{noop}$ . Noting that we have already computed  $\text{FODTR}[r\text{Case}(s), \text{unload}(b^*, t^*)]$  in Eq. (31) and that  $\text{FODTR}[r\text{Case}(s), \text{noop}]$  will just be  $r\text{Case}(s)$  with 10 replaced by 19, we obtain the following policy:

$$\begin{aligned} \pi\text{Case}[r\text{Case}(s)] &= \text{casemax}(\{\exists b^*, t^*. \text{FODTR}[r\text{Case}(s), \text{unload}(b^*, t^*)]\} \cup \{\text{FODTR}[r\text{Case}(s), \text{noop}]\}) \\ &= \begin{array}{|l|l|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) & : 19.0 \longrightarrow \text{noop} \\ \hline \neg \text{“} \wedge [\exists b^*, t^*. \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] : & 8.1 \longrightarrow \text{unload}(b^*, t^*) \\ \hline \neg \text{“} & : 0 \longrightarrow \text{noop} \\ \hline \end{array} \end{aligned}$$

For a more interesting policy, we refer the reader back to the optimal value function and policy for BoxWorld given in Fig. 4.

Technically, we note that there may be an infinite number of actions that can be applied since there are an infinite number of ground instantiations of  $\text{unload}(b^*, t^*)$  depending on the domain instantiation. Thus, this policy representation manages to compactly represent the selection of an optimal action amongst an infinite set.

#### 4. Practical FOMDP solution techniques

The last section reviewed a symbolic dynamic programming (SDP) algorithm theoretically capable of producing an  $\epsilon$ -optimal value function for a FOMDP that does not require theorem proving to detect inconsistent case partitions or logical simplification to maintain compact representations of case partition formulae. However, in practice, both theorem proving and simplification are needed to control the representational blowup of the value function occurring at each step of value iteration.

To this end, the first half of this section introduces a practical first-order extension of the algebraic decision diagram (ADD) [4] data structure, the *first-order ADD* (FOADD), for maintaining case statements in a simplified, non-redundant format that facilitates theorem proving for inconsistency detection. We show how FOADDs can be used to exploit structure in SDP for FOMDPs in much the same manner that ADDs have been used to exploit structure in dynamic programming for MDPs [38]. We conclude with an illustrative empirical results demonstrating that FOADDs enable an automated solution to basic FOMDPs. We will discuss related work on first-order decision diagrams (FODDs) [81], also applied to FOMDPs, in Section 6.

In the second half, we introduce an additive decomposition approach for approximately solving FOMDPs with universal reward specifications. This approach is motivated in part by previous decomposition methods and enables the application of FOMDP solution techniques to a reward specification that otherwise renders SDP solution approaches intractable.

##### 4.1. Representation and solution with first-order ADDs

An *algebraic decision diagram* (ADD) [4] is a data structure for compactly representing a function from  $\mathbb{B}^n \rightarrow \mathbb{R}$  using a directed acyclic graph. ADDs have been used to compactly model transition functions, rewards, and value functions in factored MDPs [12]. Moreover, value iteration defined in terms of ADD operations has yielded substantial improvements in time and space complexity over enumerated state representations [38].

To extend these ideas to the first-order framework, we define methods for breaking down first-order case partition formulae into their boolean propositional components and create a compact *first-order ADD* (FOADD) representation of case statements. Then we can apply known ADD algorithms to perform the  $\otimes$ ,  $\oplus$ , and  $\ominus$  case operations. Once we have shown how to do this, we end with a discussion of the practical use of FOADDs and a small example of a FOADD application to SDP.

###### 4.1.1. FOADD construction and operations

The first aspect of FOADDs concerns how to construct them automatically from a case representation. Since ADDs are propositional, we need some method of finding propositional structure in first-order formulae. We can do this by permuting quantifiers at the same level of nesting (e.g.,  $[\exists x, y. \phi] \equiv [\exists y, x. \phi]$ ) and by distributing quantifiers as deeply into case formulae as possible using the following rewrite rule templates ( $\diamond$  indicates variables other than those explicitly quantified):

$$[\exists x. A(x, \diamond) \vee B(x, \diamond)] \longrightarrow [(\exists x. A(x, \diamond)) \vee (\exists x. B(x, \diamond))] \quad (36)$$

$$[\forall x. A(x, \diamond) \wedge B(x, \diamond)] \longrightarrow [(\forall x. A(x, \diamond)) \wedge (\forall x. B(x, \diamond))] \quad (37)$$

$$[\exists x. A(x, \diamond) \wedge B(\diamond)] \longrightarrow [(\exists x. A(x, \diamond)) \wedge (B(\diamond))] \quad (38)$$

$$[\forall x. A(x, \diamond) \vee B(\diamond)] \longrightarrow [(\forall x. A(x, \diamond)) \vee (B(\diamond))] \quad (39)$$

**(a) Given case statement:**

$$\text{case} = \begin{array}{|c|c|} \hline \exists x. [A(x) \vee \forall y. A(x) \wedge B(x) \wedge \neg A(y)] & 1 \\ \hline \neg \text{''} & 0 \\ \hline \end{array}$$
**(b) Push down quantifiers, expose propositional structure:**

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

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

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

$$\text{case} = \begin{array}{|c|c|} \hline a \vee (b \wedge \neg a) & 1 \\ \hline \neg \text{''} & 0 \\ \hline \end{array}$$
**(c) Convert to first-order (A)ADD:**

**Fig. 5.** An example conversion from a case statement to a compact FOADD representation demonstrating first-order CSI.

We also perform equality simplification using the non-empty domain assumption with the following two rules:

$$[\exists x. x = y \wedge A(x, \diamond)] \longrightarrow A(y, \diamond) \quad (40)$$

$$[\forall x. x \neq y \vee A(x, \diamond)] \longrightarrow A(y, \diamond) \quad (41)$$

The first rule is fairly straightforward while the second rule follows simply from the negation of the first rule with renaming. In practice, we iteratively apply simplification rules (40), (41) followed by rewrite rules (36)–(39) working from the innermost to the outermost quantifiers until no more rewrites can be applied. While other orders may give different (potentially smaller) results, we find that this deterministic approach is generally sufficient to expose most propositional structure in first-order formulae.

We provide the following example application of these rewrite and simplification rules to demonstrate their power:

$$\begin{aligned} & \exists x, z. [x = y \wedge A(x, \diamond) \wedge B(y, z)] \\ & \equiv \exists x. [x = y \wedge A(x, \diamond) \wedge (\exists z. B(y, z))] \quad [\text{Apply rewrite rule (38) for } z] \\ & \equiv (\exists x. x = y \wedge A(x, \diamond)) \wedge (\exists z. B(y, z)) \quad [\text{Apply rewrite rule (38) for } x] \\ & \equiv A(y, \diamond) \wedge (\exists z. B(y, z)) \quad [\text{Apply simplification rule (40) for } x] \end{aligned}$$

To build a FOADD, we first apply these rules to expose the propositional structure of a first-order formula. Consider the example in Fig. 5(a,b); we start with

$$\exists x. [A(x) \vee \forall y. A(x) \wedge B(x) \wedge \neg A(y)] \quad (42)$$

and apply rewrite rule (39) for  $y$  followed by (36) for  $x$  to obtain

$$[\exists x. A(x)] \vee ([\exists x. A(x) \wedge B(x)] \wedge [\forall y. \neg A(y)]). \quad (43)$$

Once we have pushed quantifiers as far down as possible, we extract the propositional structure of the formula by considering propositional connectives over quantified formulae as follows:

$$\boxed{\exists x. A(x)} \vee \left( \boxed{[\exists x. A(x) \wedge B(x)]} \wedge \boxed{\forall y. \neg A(y)} \right) \quad (44)$$

Each of these boxes represents a formula that we cannot further decompose into propositional components. Consequently, we treat each of these boxed formulae as propositions. To do this, we maintain a table of mappings from propositional variables  $p$ , naming each first-order formula, to first-order formulae  $\psi$ :  $\{p \rightarrow \psi\}$ . To convert a new formula  $\phi$  in a case statement to a propositional variable, we examine each formula-to-proposition mapping in our table. If  $\phi \equiv \psi$  for some  $\psi$  in the table, we return its corresponding proposition  $p$ ; if  $\phi \equiv \neg\psi$ , we return  $\neg p$ ; otherwise, we add a new proposition label  $q$  and add the mapping  $q \rightarrow \phi$  to our table and return  $q$ . In our example, having built the table shown in Fig. 5(b), we can convert the formula to its propositional counterpart:

$$a \vee (b \wedge \neg a) \quad (45)$$

At this point, we can build an ADD from a case statement whose formulae are purely propositional. What makes this ADD first-order is the additional proposition to first-order formula mapping that gives each proposition a first-order definition. Standard ADDs can exploit *context-specific independence (CSI)* [13] (i.e., where the value of a function is independent of an input variable given the assignment to other variables). There is, however, an additional form of CSI that we can exploit in FOADDs—*first-order CSI*. This first-order CSI follows from the structured and potentially overlapping nature of the propositional variables. For instance, in our example,  $\neg a \supset \neg b$ , so as we traverse its FOADD representation, we can force the decision node for  $b$  in the context of  $a$ . This is shown in Fig. 5(c).

The options for detecting first-order CSI include:

- (a) Do not perform any first-order CSI detection at all.
- (b) Maintain information about all pairwise implications in the propositional mapping table and detect just this pairwise first-order CSI during the application of FOADD operations.
- (c) Perform full simplification for all decision nodes in the context of the conjunction of all decisions made for parent nodes during all operations on the FOADD.

Obviously (a) requires no additional computation, but can give rise to FOADDs with potentially dead paths. In contrast, (c) requires substantial computation in return for extensive simplification. In practice, we find (b) to offer the most reasonable tradeoff between computation and simplification; time-limited theorem proving, although incomplete, suffices to identify many pairwise node implications that lead to substantial first-order CSI pruning. It is trivial to extend the ADD algorithm to do this additional consistency check in the presence of parent decisions when performing the standard ADD *Apply* and *Reduce* operations. However, if (b) or (c) are used, it is not sound to reorder the ADD nodes since the first-order context of these prunings may change and thus may no longer be valid after node reordering.

Once we convert a case statement to an FOADD, we can apply the  $\otimes$ ,  $\oplus$ , and  $\ominus$  case operations to FOADDs by making direct use of the ADD *Apply* operations of multiplication, addition, and subtraction [4]. We can reuse standard ADD operations for FOADDs since they are just ADDs with augmented variable definitions in the propositional mapping table. Thus, the only practical difference between ADD and FOADD operations is that these augmented variable definitions may lead to additional pruning of structure due to first-order CSI.

In general, FOADDs may be treated as ADDs, except for the requirement to consult the propositional mapping table in the following circumstances:

1. when constructing a FOADD;
2. when converting a FOADD back to a case representation or evaluating a ground state; or
3. when exploiting first-order CSI using method (b) or (c) above, we may consult this table during the ADD *Reduce* and *Apply* procedures.

#### 4.1.2. Practical considerations

Replacing case statements with FOADDs in the representation and solution of FOMDPs has the potential to exploit a great deal of structure that naturally occurs in these representations. First, the disjunctive nature of positive effects in the regression of FOMDP formulae introduces a number of disjunctions during the application of algorithms such as SDP. Second, the existential quantification of the action variables in these formulae introduce existential quantifiers that can be distributed through the disjunctions introduced by *Regr*. Consequently, every SDP step introduces structure that can be directly exploited by the previously described methods for exposing propositional structure of first-order formulae. As such, our approach to representing FOADDs is well-suited to FOMDPs as we demonstrate below with a small example.

However, if we were to define a complete SDP algorithm for FOMDPs that only uses FOADDs, we would need to define special unary FOADD operations such as *Regr*, *casemax*, and  $\exists \vec{x}$  used in the SDP algorithm. While *Regr* can be easily defined (note that a FOADD is just a compact representation of a case statement and thus *Regr* can still be applied), it changes the logical meaning of the FOADD nodes since they have a first-order definition. In general, maintaining a canonical representation after performing *Regr* on a FOADD requires expensive node reordering operations. The application of  $\exists \vec{x}$  and *casemax* also generally require expensive node reordering operations. For these reasons, we do not apply *Regr*, *casemax*, or  $\exists \vec{x}$  to FOADDs in practice, instead opting for a pragmatic use of FOADDs that exploits their strengths.<sup>4</sup>

The primary advantage of FOADDs is the provision of efficient binary operations and formula simplification through the breakdown of propositional structure and the elimination of redundancy that occurs during their construction. In doing this simplification, FOADDs remove a lot of burden from the theorem prover, which must otherwise detect inconsistency with highly redundant representations. Thus, in our SDP algorithms, we use FOADDs where they are most useful and efficient—binary operations and logical simplification—and revert to the case representation to perform the unary operations of *Regr*, *casemax*, and  $\exists \vec{x}$  that can be expensive due to the need for internal node rotations. This approach leads to a viable SDP algorithm, to which we now turn.

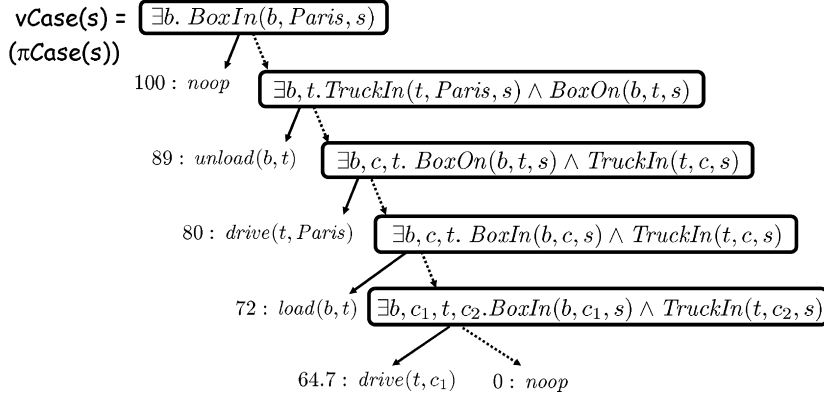
<sup>4</sup> While we do not discuss *Regr*, *casemax*, and  $\exists \vec{x}$  for FOADDs further here, the reader is referred to [67] for additional information on how one might perform these operations efficiently.

• **BoxWorld FOADD reward representation:**

$$r\text{Case}(s) = \boxed{\exists b, c. \text{BoxIn}(b, \text{Paris}, s)}$$

$\swarrow$  10       $\searrow$  0

• **BoxWorld FOADD optimal value function and policy:**



**Fig. 6.** An example FOADD representation of the reward in BoxWorld and the FOADD representation of the optimal value function and policy for this domain.

#### 4.1.3. Symbolic dynamic programming with FOADDs

The use of FOADDs in the somewhat hybrid manner discussed above allows the development of a practical SDP algorithm.

We have implemented a fully automated first-order value iteration algorithm and tested it on several examples to develop a sense of its effectiveness. One problem tested is the running BoxWorld FOMDP example. The FOADDs for the reward, optimal value function and policy are given in Fig. 6. For the variable ordering, we simply maintain the order of formulae as they were added to the variable mapping table in the FOADD during the SDP algorithm. We use the Vampire theorem prover [65] for detecting equivalence and inconsistency. The total running time for this solution until convergence within tolerance  $1e-4$  was 15.7 s on a 2 Ghz Pentium with 2 Gb of RAM. Unsurprisingly, the final FOADD for this problem gives exactly the decision list structure that we would expect for the BoxWorld problem as shown in Fig. 4.

We have also used our FOADD value iteration algorithm to solve other variants of the BoxWorld problem, including the version given in [14] with an extra fluent for *Rain*(*s*) and action probabilities conditioned on this fluent. We also used a BoxWorld reward with the following structure:

$$R(s) = \begin{array}{|l|l|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) \wedge \text{TypeA}(b) & : 10 \\ \hline \neg \text{---} \wedge \exists b. \text{BoxIn}(b, \text{paris}, s) \wedge \neg \text{TypeA}(b) & : 5 \\ \hline \neg \text{---} & : 0 \\ \hline \end{array} \quad (46)$$

Here in addition to the *Rain*(*s*) fluent, we have also added a non-fluent predicate *TypeA*(*b*) to distinguish types of boxes and varying rewards for each type of box. The FOADDs for these solutions are too large to display, but we note that after a small number of steps of value iteration, the value function FOADD stopped growing indicating that all relevant state partitions had been identified. Value iteration continued with this quiesced FOADD until all values at the leaves converged. The respective solution times to convergence within tolerance  $1e-4$  for these more complex problems were 70.4 s and 489 s on a 2 Ghz Pentium with 2 Gb of RAM. For comparison, the ReBel algorithm [43] produced the same solution for the first FOMDP variant with the *Rain*(*s*) fluent in <6 s on a 3.1 Ghz machine. ReBel's specialization for a less expressive subset of FOMDPs (still capturing BoxWorld, however) results in a substantial performance edge. We discuss differences between ReBel and the work in this article in Section 6.

There appear to be at least two general criteria for problem domains to demonstrate finitely-sized optimal value functions with the current case representation as occurred in these examples: (1) the non-zero reward case partitions must be existentially quantified and (2) the FOMDP dynamics must not introduce transitive structure that cannot be finitely bounded by domain axioms. As this last requirement is vague, we provide an example. In the BoxWorld problem covered in this section, we implicitly assume that all cities are accessible from each other via the *drive* action. If instead we had some underlying road topology indicated by *Conn*(*City* : *c*<sub>1</sub>, *City* : *c*<sub>2</sub>) that restricted the *drive* action and we did not know this topology in terms of prior knowledge specified as domain axioms, then the SDP algorithm would likely need to generate representations for all possible topologies, thus likely leading to a value function of infinite “size.” Infinite-sized value functions can also occur when condition (1) is violated as we discuss in the next subsection. We discuss potential research directions to mitigate these observed deficiencies of the case representation in Section 7.1.

Unfortunately, the FOADD solution approach has failed to scale to more complex problems used in the planning community (particularly problems from the ICAPS 2004 and 2006 International Planning Competitions) since they typically use more complex rewards, including those with universal quantifiers. Whereas problems with existentially quantified rewards may exhibit a finite-size optimal value function, this is rarely the case with universal rewards. Thus additional techniques are required to handle this problem, as we discuss next.

#### 4.2. Decomposing universal rewards

In first-order domains, we are often faced with *universal reward expressions* that assign some positive value to the world states satisfying a formula of the general form  $\forall y \phi(y, s)$ , and 0 otherwise. For instance, in our BoxWORLD problem, we may define a reward as having *all* boxes  $b$  at their assigned destination city  $c$  given by  $\text{Dst}(b, c)$ :

$$R(s) = \frac{\forall b, c. \text{Dst}(b, c) \supset \text{BoxIn}(b, c, s) : 1}{\neg : 0} \quad (47)$$

One difficulty with such rewards is that our case statements provide a piecewise-constant representation of the value function. However, with universal rewards, the value function typically depends on the *number* of domain objects of interest. In our example, value at a state depends on the number of boxes not at their proper destination (since this can impact the minimum number of steps it will take to obtain the reward). So a  $t$ -stage-to-go value function in this case would have the following characteristic structure (where we use English in place of first-order logic for readability):

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

Obviously, since there are  $t$  distinct values in an optimal  $t$ -stage-to-go value function, the piecewise-constant case representation requires a minimum of  $t$  case partitions to represent this value function. And when we combine these counting dynamics with other interacting processes in the FOMDP, we often see an uncontrollable combinatorial blowup in the number of case partitions of value functions for FOMDPs with universally defined rewards. As noted by [33], effectively handling universally quantified rewards is one of the most pressing issues in the practical solution of FOMDPs.

To address this problem we adopt a decompositional approach, motivated in part by techniques for additive rewards in MDPs [11,53,61,74]. We divide our solution into off-line and on-line components where the on-line component requires a finite-domain assumption in order to execute the policy.

##### 4.2.1. Offline generic goal solution

Intuitively, given a goal-oriented reward that assigns positive reward if  $\forall \vec{y} G(\vec{y}, s)$  is satisfied, and zero otherwise, we can decompose it into a set of ground goals  $\{G(\vec{y}_1), \dots, G(\vec{y}_n)\}$  for all possible  $\vec{y}_j$  in a ground domain of interest. If we reach a state where all ground goals are true, then we have satisfied  $\forall \vec{y} G(\vec{y}, s)$ .

Of course, our methods solve FOMDPs without knowledge of the specific domain, so the set of ground goals that will be faced at run-time is unknown. Thus, in the offline FOMDP solution, we assume a *generic* ground goal  $G(\vec{y}^*)$  for a “generic” object vector  $\vec{y}^*$ . Assuming that our universal reward takes an implicative form as it does in our reworked BoxWORLD example, the conditions in the antecedent ( $\text{Dst}(b, c)$ ) indicate the goal objects of interest (all pairs  $\langle b, c \rangle$  satisfying  $\text{Dst}(b, c)$ ) and the consequent of the implication indicates the specific goal  $G(\vec{y}, s)$  to be achieved for these objects ( $\text{BoxIn}(b, c, s)$ ).

It is easy to construct a generic instance of a reward function  $R_{G(\vec{y}^*)}(s)$  given a single goal. In our BoxWORLD example we would introduce the distinguished constants  $b^*$  and  $c^*$  to denote our goal objects of interest  $G(b^*, c^*)$ :

$$r\text{Case}_{G(b^*, c^*)}(s) = \frac{\text{BoxIn}(b^*, c^*, s) : 1}{\neg \text{BoxIn}(b^*, c^*, s) : 0} \quad (48)$$

Given this simple reward, it is easy to derive a value function  $V_{G(\vec{y}^*)}(s)$  for this FOMDP using SDP or the approximate FOMDP solution algorithms that we introduce in subsequent sections.  $V_{G(\vec{y}^*)}(s)$  and its corresponding policy assume that  $\vec{y}^*$  is the only object vector of interest satisfying relevant type constraints and goal preconditions in the domain. In our running BoxWORLD example, the optimal  $v\text{Case}_{G(b^*, c^*)}(s)$  would look very similar to Fig. 4 (or 6) with some differences owing to the fact that our reward is defined in terms of constants  $b^*$  and  $c^*$  rather than existentially quantified variables  $b$  and  $c$ .

We next derive Q-functions for each action  $A_i(\vec{x})$  from the value function  $V_{G(\vec{y}^*)}(s)$  for the “generic” domain:

$$Q_{G(\vec{y}^*)}(A_i, s) = \exists \vec{x}. \text{FODTR}[V_{G(\vec{y}^*)}(s), A_i(\vec{x})] \quad (49)$$

For our running BoxWORLD example, we would derive  $q\text{Case}_{G(b^*, c^*)}(A_i, s)$  for  $A_i \in \{\text{unload}, \text{load}, \text{drive}\}$ .

---

**input** : (1) For each action template  $A_i(\vec{x})$  a set of (non-disjoint) Q-functions  $qCase_{G(\vec{y}^*)}(A_i, s)$  for a specific ground instantiation  $\vec{y}^*$  of a goal  $G$ .  
 (2) A set of  $n$  unsatisfied goals  $\{G(\vec{y}_1), \dots, G(\vec{y}_n)\}$  to achieve.  
 (3) A ground state  $s$  to find the best action for.

**output** : The optimal ground action  $A_i(\vec{c})$  to execute with respect to the given state and additive decomposition of unsatisfied goals:  $A_i(\vec{c}) = \arg \max_{i, \vec{c}} \sum_{j=1}^n qCase_{G(\vec{y}_j)}(A_i(\vec{c}), s)$

**begin**

// In hash table  $h$ , entries map ground actions to corresponding value:  $A(\vec{x}) \rightarrow v$ .  
 Initialize empty hash table  $h$ ;

// Now, compute additive values for all matching ground actions

**foreach** (action  $A_i$ ) **do**

**foreach** (goal  $G(\vec{y}_j)$ ) **do**

Replace all occurrences of  $\vec{y}^*$  in  $qCase_{G(\vec{y}^*)}(A_i, s)$  with  $\vec{y}_j$ ;

**foreach** (case partition  $\langle \exists \vec{x} \phi(\vec{x}), t \rangle \in qCase_{G(\vec{y}_j)}(A_i, s)$ ) **do**

**foreach** (ground binding  $\vec{x} = \vec{c}$  satisfying  $\exists \vec{x} \phi(\vec{x})$ ) **do**

**if** ( $A_i(\vec{c}) \rightarrow v$  is already in  $h$  for some  $v$ ) **then**

Update  $h$  to contain  $A_i(\vec{c}) \rightarrow (v + \frac{t}{n})$ ;

**else**

Update  $h$  to contain  $A_i(\vec{c}) \rightarrow \frac{t}{n}$ ;

// Assume  $h$  tracks its maximal entry:  $A_i(\vec{c}) \rightarrow v$ .  
 Return the maximal  $A_i(\vec{c})$  from  $h$ ;

**end**

---

**Algorithm 1.**  $EvalPolicy(\{qCase_{G(\vec{y}^*)}(\cdot, \cdot)\}, \{G(\vec{y}_1), \dots, G(\vec{y}_n)\}, s) \rightarrow A_i(\vec{c})$ .

#### 4.2.2. Online policy evaluation

With the offline solution (i.e., Q-function for each action) of a generic goal FOMDP in hand, we address the online problem of action selection for a specific domain instantiation given at run-time. We assume a set of ground goals  $\{G(\vec{y}_1), \dots, G(\vec{y}_n)\}$  corresponding to a specific finite domain given at run-time. If we assume that (typed) domain objects are treated uniformly in the uninstantiated FOMDP, as is the case in many logistics and planning problems, then we obtain the Q-function for any goal  $G(\vec{y}_j)$  by replacing all ground terms  $\vec{y}^*$  in  $qCase_{G(\vec{y}^*)}(A_i, s)$  with the respective terms  $\vec{y}_j$  to obtain  $qCase_{G(\vec{y}_j)}(A_i, s)$ .

Returning to our running example, from the value function  $vCase_{G(b^*, c^*)}(s)$  we derived a Q-function  $qCase_{G(\vec{y}^*)}(A_i, s)$  for each action  $A_i$ . If at run-time, we are given the three goals  $Dst(b_1, paris)$ ,  $Dst(b_2, berlin)$ , and  $Dst(b_3, rome)$ , then we would substitute these goals into our Q-functions to obtain three goal-specific Q-functions for each action  $A_i$ :

$$\{qCase_{G(b_1, paris)}(A_i, s), qCase_{G(b_2, berlin)}(A_i, s), qCase_{G(b_3, rome)}(A_i, s)\} \quad (50)$$

Action selection requires finding an action that maximizes value with respect to the original universal reward. Following [11,53], we do this by treating the *sum of the Q-values* of any action in the subgoal MDPs as a measure of its Q-value in the joint (original) MDP. Specifically, we assume that each goal contributes uniformly and additively to the reward, so the Q-function for an entire set of ground goals  $\{G(\vec{y}_1), \dots, G(\vec{y}_n)\}$  determined by our domain instantiation is just  $\sum_{j=1}^n \frac{1}{n} qCase_{G(\vec{y}_j)}(A_i, s)$ . Action selection (at run-time) in any ground state is realized by choosing the action with maximum additive Q-value. Naturally, we do not want to explicitly create the joint Q-function, but instead use an efficient “scoring” technique that evaluates potentially useful actions by iterating through the individual Q-functions as described in Algorithm 1.

While this additive and uniform decomposition may not be appropriate for all domains with goal-oriented universal rewards (and certainly offers no performance guarantees on account of its heuristic nature), we have found it to provide reasonable results for domains such as BoxWORLD as we empirically demonstrate in the next section. While our approach only currently handles rewards with universal quantifiers, this reflects the form of many planning problems. Nonetheless, this technique could be extended for more complex universal rewards, the general open question being how to assign credit among the constituents of such a reward.

## 5. Linear-value approximation for FOMDPs

Perhaps the greatest difficulty with the symbolic dynamic programming (SDP) approach and practical extensions discussed in the last section is that the size of the value function case representation grows polynomially on each iteration and

thus exponentially in terms of the number of iterations.<sup>5</sup> Similar growth can occur for the first-order formulae representing the state partitions themselves. Once these formulae become too large to practically detect equivalence or inconsistency, all hope of obtaining a compact representation of the value function is lost as the number of partitions in the case representation grow unboundedly with no practical means for simplification or pruning. Indeed, the SDP approaches above, using both FOADDs and universal reward decomposition, are incapable of producing value functions and policies competitive with other planners from the ICAPS 2004 and 2006 International Probabilistic Planning Competitions [32,50].

Given that approximate solution techniques such as linear value approximation [19,36,71] have allowed MDP solutions to scale far beyond the limits of exact algorithms, at the same time offering reasonable error guarantees, this suggests generalizing linear value approximation techniques to FOMDPs. In this section, we generalize the LP methods for ground MDPs, discussed in Section 2, to the first-order case. This reduces the task of solving an FOMDP to that of obtaining good weights for a set of basis functions that approximates the optimal value function. This requires the generalization of linear programs to handle first-order constraints and further requires efficient extensions of solution methods such as constraint generation and variable elimination in cost networks to exploit the first-order structure of these constraints.

To develop a completely automated linear-value approximation approach to FOMDPs we must address the issue of automatic basis function construction; to do this, we adapt techniques proposed by [33]. With appropriate domain axioms defining legal states, our techniques provide fully first-order, non-grounded solutions to FOMDPs derived from PPDDL and can compete with planners from the ICAPS 2004 and ICAPS 2006 International Probabilistic Planning Competitions.

### 5.1. Benefits of linear-value approximation

Linear-value approximation for FOMDPs is attractive for several reasons:

- Given that much of the computation in linear value approximation reduces to solving LPs, this reduces the algorithm design space to the setup and solution of linear programs.
- Since the size of linear-value approximations is fixed, it can be used to moderate the complexity of the resulting solution algorithm. This leads to a flexible solution approach that trades off approximation accuracy and computation.
- Linear value approximation does not require extensive logical simplification in practice, just weight projections that make use of a theorem prover. This is a tremendous advantage over exact techniques that require substantial simplification in order to maintain a compact representation.
- Linear value approximation have yielded reasonable empirical performance for ground and factored MDPs, suggesting promise for its application to FOMDPs.
- If we do not use additive reward decomposition techniques of Section 4.2 (which approximate the FOMDP model), then we can derive domain-independent error bounds on our resulting value function using Corollary 3.4.2.

### 5.2. First-order linear-value representation

We represent a value function as a weighted sum of  $k$  first-order basis functions, denoted  $b_i(s)$ , each ideally containing a small number of formulae that provide a first-order abstraction of state space:

$$V(s) = \bigoplus_{i=1}^k w_i \cdot b_i(s) \quad (51)$$

Throughout this section, we assume that each individual basis function  $b_i(s)$  is represented by a case statement that is an exhaustive and disjoint partitioning of state space. This property will be useful when we define the backup operators next. However, two basis functions may assign non-zero values to overlapping regions of state space; in fact this can be quite useful for representing additively decomposable values.

Such a linear value function representation can often provide a reasonable approximation of the exact value function, especially given the additive structure inherent in many real-world problems. For example, as argued in previous sections, many planning problems have additive reward functions or multiple goals, both of which lend themselves to approximation via linearly additive basis functions. Unlike exact solution methods where value functions can grow exponentially in size during the solution process and must be logically simplified, here we maintain the value function in a compact form that requires no simplification, just discovery of good weights.

As an example, consider approximation of the value function for our BoxWorld FOMDP from the last section, using the following basis functions (we refer to specific instances of  $b_i(s)$  as  $bCase_i(s)$ ):

$$bCase_1(s) = \begin{array}{|l} \exists b. \text{BoxIn}(b, \text{paris}, s) : 1 \\ \neg : 0 \end{array}$$

<sup>5</sup> In the worst case, a single case operation can yield a quadratic blowup in the number of case partitions in terms of the maximum number of case partitions in its operands.



$$\begin{aligned}
bCase_2(s) &= \begin{array}{|l|} \hline \exists b, t. BoxOn(b, t, s) : 1 \\ \hline \neg, " : 0 \\ \hline \end{array} \\
bCase_3(s) &= \begin{array}{|l|} \hline \exists b, t. TruckIn(t, paris, s) \wedge BoxOn(b, t, s) : 1 \\ \hline \neg, " : 0 \\ \hline \end{array}
\end{aligned} \tag{52}$$

Then each instance of  $V(s)$  (denoted by  $vCase(s)$ ) has the form:

$$vCase(s) = [w_1 \cdot bCase_1(s)] \oplus [w_2 \cdot bCase_2(s)] \oplus [w_3 \cdot bCase_3(s)] \tag{53}$$

Each basis function is relatively small and represents a portion of state space to which we would expect to assign some positive value in order to approximate the BoxWORLD value function.

### 5.2.1. Backup operators

Suppose we are given a value function  $V(s)$ . Backing up this value function through an action  $A(\vec{x})$  yields a case statement containing the logical description of states that would give rise to  $V(s)$  after doing action  $A(\vec{x})$ , as well as the values thus obtained.

However, due to the free variables in action  $A(\vec{x})$ , there are in fact two types of backups that we can perform. The first,  $B^{A(\vec{x})}[\cdot]$ , regresses a value function through an action and produces a case statement with *free variables* for the action parameters. The second,  $B^A[\cdot]$ , existentially quantifies over the free variables  $\vec{x}$  in  $B^{A(\vec{x})}[\cdot]$ . Thus, the application of  $B^A[\cdot]$  results in a case description of the regressed value function indicating the values that could be achieved by *any* instantiation of  $A(\vec{x})$  in the pre-action state.

The definition of  $B^{A(\vec{x})}[\cdot]$  is almost the same as the *first-order decision theoretic regression* (FODTR) operator from Eq. (30), except that we do not explicitly add in the reward. Slightly modifying our definitions from Section 3.3.3, we let  $n_1(\vec{x}), \dots, n_q(\vec{x})$  be the set of Nature's deterministic actions for stochastic action  $A(\vec{x})$ . Then we define  $B^{A(\vec{x})}[\cdot]$  as follows:

$$B^{A(\vec{x})}[V(s)] = \gamma \left[ \bigoplus_{j=1}^q \{P(n_j(\vec{x}), A(\vec{x}), s) \otimes Repr(V(do(n_j(\vec{x}), s)))\} \right] \tag{54}$$

Defining  $B^{A(\vec{x})}[\cdot]$  in this way without the reward makes it a linear operator. Thus, if we apply this operator to our linear-value function representation, it distributes to each first-order basis function:

$$B^{A(\vec{x})}[V(s)] = B^{A(\vec{x})} \left[ \bigoplus_{i=1}^k w_i \cdot b_i(s) \right] = \bigoplus_{i=1}^k w_i \cdot B^{A(\vec{x})}[b_i(s)] \tag{55}$$

Having defined  $B^{A(\vec{x})}[\cdot]$ , we now use it to define  $B^A[\cdot]$ :<sup>6</sup>

$$B^A[V(s)] = \exists \vec{x}. \{B^{A(\vec{x})}[V(s)]\} \tag{56}$$

Unfortunately, if we apply  $B^A[\cdot]$  to our linear-value function representation, we see that  $B^A[\cdot]$  is not necessarily linear:

$$B^A[V(s)] = B^A \left[ \bigoplus_{i=1}^k w_i \cdot b_i(s) \right] = \exists \vec{x}. \left\{ \bigoplus_{i=1}^k w_i \cdot B^{A(\vec{x})}[b_i(s)] \right\} \tag{57}$$

The difficulty is that the existential quantification of  $B^A[\cdot]$  jointly constrains the backup of all basis functions that contain the existentially quantified variable as a free variable.

These problems can be mitigated, however. We begin with a few definitions.

**Definition 5.2.1.** We say that a deterministic action  $n_j(\vec{x})$  *affects* a fluent  $F$  if there is a positive or negative effect axiom that contains  $a = n_j(\vec{x})$  in the body of the axiom and  $F$  in the head (cf., Section 3.2.2). We say that a stochastic action  $A(\vec{x})$  *affects* a fluent  $F$  if at least one of Nature's choices  $n_j(\vec{x})$  for  $A(\vec{x})$  affects  $F$ . Finally, a formula  $\phi$  is *affected by* a stochastic action  $A(\vec{x})$  iff  $\phi$  contains a fluent affected by  $A(\vec{x})$ . Since a case statement is defined as a logical formula, this definition extends to case statements in the obvious way.

**Property 5.2.2.** When a basis function case statement  $b_i(s)$  is affected by a stochastic action  $A(\vec{x})$ ,  $B^{A(\vec{x})}[b_i(s)]$  will contain the action arguments  $\vec{x}$  as free variables. The inverse of this property is also true: if a stochastic action  $A(\vec{x})$  does not affect a basis function  $b_i(s)$ ,  $B^{A(\vec{x})}[b_i(s)]$  will not contain the action arguments as free variables.

<sup>6</sup> For simplicity, we assume that the reward is independent of the action arguments  $\vec{x}$ , allowing us to exclude the reward from the  $\exists \vec{x}$  operation of  $B^A$ . If required, such dependencies could be added with appropriate adjustments to our definitions.

To exploit this property, we let  $I_A^+$  denote the set of indices  $i$  for basis functions  $b_i(s)$  that are affected by an action  $A(\vec{x})$  (so that for all  $i \in I_A^+$ ,  $B^{A(\vec{x})}[b_i(s)]$  contains at least one of the free variables  $\vec{x}$ ). Likewise, we let  $I_A^-$  denote the set of indices of basis functions  $b_i(s)$  not affected by an action (so that for all  $i \in I_A^-$ ,  $B^{A(\vec{x})}[b_i(s)]$  contains none of the free variables  $\vec{x}$ ). We can exploit the fact that the  $\exists \vec{x}$  is vacuous for case statements not containing free variables  $\vec{x}$  and remove these terms from the scope of the  $\exists \vec{x}$  quantification. This yields the following form for  $B^A$ :

$$B^A \left[ \bigoplus_i w_i b_i(s) \right] = \left( \bigoplus_{i \in I_A^-} w_i B^{A(\vec{x})}[b_i(s)] \right) \oplus \exists \vec{x}. \left( \bigoplus_{i \in I_A^+} w_i B^{A(\vec{x})}[b_i(s)] \right) \quad (58)$$

Consequently, if no fluent occurs in more than a few basis functions and no action affects more than a few fluents then we can reasonably expect the result of applying  $B^A$  to retain some additive structure. The first property can be controlled by the appropriate design of basis functions. The second is true of typical planning domains.

As a concrete example to demonstrate the backup operators and the exploitation of additive structure, let us compute  $B^{drive}[\cdot]$  for our previously specified linear-value function from Eq. (53):

$$\begin{aligned} B^{drive}[vCase(s)] &= \exists t^*, c^*. B^{drive(t^*, c^*)}[vCase(s)] \\ &= \exists t^*, c^*. B^{drive(t^*, c^*)}[w_1 \cdot bCase_1(s) \oplus w_2 \cdot bCase_2(s) \oplus w_3 \cdot bCase_3(s)] \\ &= \exists t^*, c^*. \{ w_1 \cdot B^{drive(t^*, c^*)}[bCase_1(s)] \oplus w_2 \cdot B^{drive(t^*, c^*)}[bCase_2(s)] \oplus w_3 \cdot B^{drive(t^*, c^*)}[bCase_3(s)] \} \\ &= \exists t^*, c^*. \left\{ w_1 \cdot \begin{array}{c|c} \exists b. BoxIn(b, paris, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \oplus w_2 \cdot \begin{array}{c|c} \exists b, t. BoxOn(b, t, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \right. \\ &\quad \left. \oplus w_3 \cdot \begin{array}{c|c} \exists b, t. [t = t^* \wedge c^* = paris \wedge \exists c_1 TruckIn(t, c_1, s)] & \\ \hline \vee TruckIn(t, paris, s) \wedge BoxOn(b, t, s) & : 0.9 \\ \hline \neg & : 0 \end{array} \right\} \quad (59) \end{aligned}$$

Here, we note that the first and second basis functions are not affected by the  $drive(t^*, c^*)$  action and thus their backup through this action is equivalent to a backup through a *noop*. Since the third basis function is affected by the action  $drive(t^*, c^*)$  and this introduces the action parameters  $t^*$  and  $c^*$  into the result of its backup, we can push the quantifiers in to just this third case statement:

$$\begin{aligned} B^{drive}[vCase(s)] &= \exists t^*, c^*. B^{drive(t^*, c^*)}[vCase(s)] \\ &= w_1 \cdot \begin{array}{c|c} \exists b. BoxIn(b, paris, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \oplus w_2 \cdot \begin{array}{c|c} \exists b, t. BoxOn(b, t, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \\ &\quad \oplus w_3 \cdot \exists t^*, c^*. \left\{ \begin{array}{c|c} \exists b, t. [t = t^* \wedge c^* = paris \wedge \exists c_1 TruckIn(t, c_1, s)] & \\ \hline \vee TruckIn(t, paris, s) \wedge BoxOn(b, t, s) & : 0.9 \\ \hline \neg & : 0 \end{array} \right\} \end{aligned}$$

Finally, we carry out the explicit  $\exists t^*, c^*$  operation on the third case statement where we distribute the quantifiers inside the case partitions and simplify. This allows us to remove the  $\exists t^*, c^*$  by rewriting equalities and exploiting the non-empty domain assumption:

$$\begin{aligned} B^{drive}[vCase(s)] &= \exists t^*, c^*. B^{drive(t^*, c^*)}[vCase(s)] \\ &= w_1 \cdot \begin{array}{c|c} \exists b. BoxIn(b, paris, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \oplus w_2 \cdot \begin{array}{c|c} \exists b, t. BoxOn(b, t, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \\ &\quad \oplus w_3 \cdot \begin{array}{c|c} \exists b, t. [(\exists c_1. TruckIn(t, c_1, s)) \vee TruckIn(t, paris, s)] \wedge BoxOn(b, t, s) : 0.9 & \\ \hline \neg & : 0 \end{array} \quad (60) \end{aligned}$$

This example demonstrates best case performance for  $B^A[\cdot]$ , where an action only affects one basis function thus allowing the other basis functions to be removed from the scope of the  $\exists \vec{x}$  operator. Then the  $\exists \vec{x}$  operator can be easily applied to a single case statement without incurring a representational blowup that would otherwise occur if the  $\exists \vec{x}$  ranged over a sum of case statements and the explicit “cross-sum”  $\oplus$  was required.

Of course, in many cases, more than one basis function will be affected by an action. For example, if we had computed  $B^{unload}[vCase(s)]$ , all three basis functions would have been affected by the action and we would have had to explicitly compute the “cross-sum”  $\oplus$  of the backups of all three basis functions. While this effectively counteracts many of the benefits of linear-value approximation since additive structure can no longer be exploited, we will see that by generating our basis functions in a restricted manner, we can often manage to avoid computing the explicit  $\oplus$ , even when *all* basis functions are affected by an action. We will discuss this further when we discuss basis function generation.

### 5.3. First-order approximate linear programming

We now generalize the approximate linear programming (ALP) approach for MDPs (see Eq. (11)) to first-order MDPs. If we simply substitute appropriate notation, we arrive at the following formulation of first-order ALP (FOALP) [68]:

$$\begin{aligned}
 \text{Variables: } & w_i; \quad \forall i \leq k \\
 \text{Minimize: } & \sum_s \bigoplus_{i=1}^k w_i \cdot b_i(s) \\
 \text{Subject to: } & 0 \geq R(s) \oplus B^A \left[ \bigoplus_{i=1}^k w_i \cdot b_i(s) \right] \ominus \bigoplus_{i=1}^k w_i \cdot b_i(s); \quad \forall A, s
 \end{aligned} \tag{61}$$

As with ALP, our variables are the weights of our basis functions and our objective is to minimize the sum of values over all states  $s$ . We have one constraint for each stochastic action  $A$  (e.g., in BoxWORLD,  $A \in \{\text{unload}, \text{load}, \text{drive}\}$ ) and each state  $s$ . One advantage of FOALP over SDP is that it does not require a casemax, thus avoiding the representational blowup incurred by this step in SDP.<sup>7</sup> Unfortunately, while the objective and constraints in ALP for a ground MDP range over a finite number of states, this direct generalization to the FOALP approach for FOMDPs requires dealing with infinitely (or indefinitely) many states  $s$ .

Since we are summing over infinitely many states in the FOALP objective, it is ill-defined. Thus, we redefine the FOALP objective in a manner that preserves the intention of the original approximate linear programming solution for MDPs. In ALP (see Eq. (11)), the objective equally weights each state and minimizes the sum of the value function over all states. However, if we look at the case partitions  $\langle \phi_i(s), t_i \rangle$  of each basis function  $b_i(s)$  case statement, each case partition serves as an aggregate representation of ground states assigned equal value. Consequently, rather than count ground states in our FOALP objective—of which there will generally be an infinite number per partition—we suppose that each basis function partition is chosen because it represented a potentially useful partitioning of state space, and thus weight each case partition equally. Consequently, we rewrite the FOALP objective as follows:

$$\sum_s \bigoplus_{i=1}^k w_i \cdot b_i(s) = \bigoplus_{i=1}^k w_i \sum_s b_i(s) \sim \bigoplus_{i=1}^k w_i \sum_{\langle \phi_j, t_j \rangle \in b_i} \frac{t_j}{|b_i|}$$

We use  $|b_i|$  to indicate the number of partitions in the  $i$ th basis function. This approach can be seen as aggregating states within a basis function partition into one abstract state and then weighting each abstract state uniformly in importance. For the case of 0–1 indicator basis functions as in Eq. (52), this yields a simple objective of  $\sum_{i=1}^k w_i$ . Of course, this solution requires approximating the original objective and thus FOALP does not represent an exact generalization of the ground ALP approach to the first-order case. Nonetheless, we show that this approximation still leads to reasonable results in our empirical evaluation.

With the issue of the infinite objective resolved, this leaves us with one final problem—the infinite number of constraints (i.e., one for every state  $s$ ). Fortunately, we can work around this since case statements are finite. Since the value  $t_i$  for each case partition  $\langle \phi_i(s), t_i \rangle$  is constant over all situations satisfying the  $\phi_i(s)$ , we can explicitly sum over the  $\text{case}_i(s)$  statements in each constraint to yield a single case statement representation of the constraints. The key observation here is that the finite number of constraints represented in the single “flattened” case statement hold iff the original infinite set of constraints in Eq. (61) hold.

To understand this, consider the constraints for the *drive* action in FOALP, substituting our previously defined basis functions  $b\text{Case}_i(s)$  from Eq. (52) for  $b_i(s)$ , the results of the  $B^{\text{drive}}$  operator for these basis functions from Eq. (60), and the reward definition for BoxWORLD given by  $r\text{Case}(s)$  in Eq. (20) for  $R(s)$ . We substitute all of these directly into the constraint of the form in Eq. (61) above to obtain:

$$\begin{aligned}
 0 \geq & \frac{\exists b. \text{BoxIn}(b, \text{paris}, s) : 10}{\neg : 0} \oplus w_1 \cdot \frac{\exists b. \text{BoxIn}(b, \text{paris}, s) : 0.9}{\neg : 0} \\
 & \oplus w_2 \cdot \frac{\exists b, t. \text{BoxOn}(b, t, s) : 0.9}{\neg : 0} \\
 & \oplus w_3 \cdot \frac{\exists b, t. [(\exists c_1. \text{TruckIn}(t, c_1, s)) \vee \text{TruckIn}(t, \text{paris}, s)] \wedge \text{BoxOn}(b, t, s) : 0.9}{\neg : 0} \\
 & \ominus w_1 \cdot \frac{\exists b. \text{BoxIn}(b, \text{paris}, s) : 1}{\neg : 0} \ominus w_2 \cdot \frac{\exists b, t. \text{BoxOn}(b, t, s) : 1}{\neg : 0} \\
 & \ominus w_3 \cdot \frac{\exists b, t. \text{TruckIn}(t, \text{paris}, s) \wedge \text{BoxOn}(b, t, s) : 1}{\neg : 0}; \quad \forall s
 \end{aligned} \tag{62}$$

<sup>7</sup> The reasons for this are the same as for the lack of a max in the ground case as discussed in Section 2.2.3.

Next we perform an explicit  $\oplus$  and  $\ominus$  for some of the case statements, simplify the resulting partitions, and distribute the weights into the partition values:

$$\begin{aligned}
 0 \geq & \begin{array}{|c|c|} \hline \exists b. \text{BoxIn}(b, \text{paris}, s) : 10 - 0.1 \cdot w_1 & \\ \hline \neg \text{“} & : 0 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline \exists b, t. \text{BoxOn}(b, t, s) : -0.1 \cdot w_2 & \\ \hline \neg \text{“} & : 0 \\ \hline \end{array} \\
 \oplus & \begin{array}{|c|c|} \hline \exists b, t. \text{TruckIn}(t, \text{paris}, s) \wedge \text{BoxOn}(b, t, s) & : -0.1 \cdot w_3 \\ \hline \neg \wedge \exists b, t, c_1. \text{TruckIn}(t, c_1, s) \wedge \text{BoxOn}(b, t, s) : 0.9 \cdot w_3 & \\ \hline \neg \text{“} & : 0 \\ \hline \end{array} ; \quad \forall s
 \end{aligned} \tag{63}$$

To maintain our representation in a compact and perspicuous form, we define the following propositional renamings for the first-order formulae in these case statements:<sup>8</sup>

$$\begin{aligned}
 \phi_1(s) &\equiv \exists b. \text{BoxIn}(b, \text{paris}, s) \\
 \phi_2(s) &\equiv \exists b, t. \text{BoxOn}(b, t, s) \\
 \phi_3(s) &\equiv \exists b, t. \text{TruckIn}(t, \text{paris}, s) \wedge \text{BoxOn}(b, t, s) \\
 \phi_4(s) &\equiv \exists b, t, c_1. \text{TruckIn}(t, c_1, s) \wedge \text{BoxOn}(b, t, s)
 \end{aligned}$$

Finally, we fully expand the  $\oplus$  to obtain an explicit representation of *all* FOALP constraints for the *drive* action in our BoxWORLD example:

$$\begin{array}{|c|c|} \hline \phi_1(s) \wedge \phi_2(s) \wedge \phi_3(s) & : 0 \geq 10 - 0.1 \cdot w_1 + -0.1 \cdot w_2 + -0.1 \cdot w_3 \\ \hline \phi_1(s) \wedge \phi_2(s) \wedge \neg \phi_3(s) \wedge \phi_4(s) & : 0 \geq 10 - 0.1 \cdot w_1 + -0.1 \cdot w_2 + 0.9 \cdot w_3 \\ \hline \phi_1(s) \wedge \phi_2(s) \wedge \neg \phi_3(s) \wedge \neg \phi_4(s) & : 0 \geq 10 - 0.1 \cdot w_1 + -0.1 \cdot w_2 \\ \hline \phi_1(s) \wedge \neg \phi_2(s) \wedge \phi_3(s) & : 0 \geq 10 - 0.1 \cdot w_1 + -0.1 \cdot w_3 \\ \hline \phi_1(s) \wedge \neg \phi_2(s) \wedge \neg \phi_3(s) \wedge \phi_4(s) & : 0 \geq 10 - 0.1 \cdot w_1 + 0.9 \cdot w_3 \\ \hline \phi_1(s) \wedge \neg \phi_2(s) \wedge \neg \phi_3(s) \wedge \neg \phi_4(s) & : 0 \geq 10 - 0.1 \cdot w_1 + -0.1 \cdot w_2 \\ \hline \neg \phi_1(s) \wedge \phi_2(s) \wedge \phi_3(s) & : 0 \geq -0.1 \cdot w_2 + -0.1 \cdot w_3 \\ \hline \neg \phi_1(s) \wedge \phi_2(s) \wedge \neg \phi_3(s) \wedge \phi_4(s) & : 0 \geq -0.1 \cdot w_2 + 0.9 \cdot w_3 \\ \hline \neg \phi_1(s) \wedge \phi_2(s) \wedge \neg \phi_3(s) \wedge \neg \phi_4(s) & : 0 \geq -0.1 \cdot w_2 \\ \hline \neg \phi_1(s) \wedge \neg \phi_2(s) \wedge \phi_3(s) & : 0 \geq -0.1 \cdot w_3 \\ \hline \neg \phi_1(s) \wedge \neg \phi_2(s) \wedge \neg \phi_3(s) \wedge \phi_4(s) & : 0 \geq 0.9 \cdot w_3 \\ \hline \neg \phi_1(s) \wedge \neg \phi_2(s) \wedge \neg \phi_3(s) \wedge \neg \phi_4(s) & : 0 \geq 0 \\ \hline \end{array} ; \quad \forall s \tag{64}$$

Here, if we had detected that any partition formula had been inconsistent, we would have removed it and the corresponding constraint.

While we note that technically there are an infinite number of constraints (one for every possible state  $s$ ), there are only a finite number of *distinct* constraints. In fact, the case representation conveniently partitions the state space into regions with the same constraint. Thus, to solve the FOALP problem, we could enumerate all consistent constraints for every action and then directly solve the resulting LP. In addition to the above constraints for the *drive* action in BoxWORLD, this approach would require us to carry out a similar procedure for the *unload*, *load*, and *noop* actions; however, once we did this, we would have all of the constraints necessary for solving the FOALP first-order linear program specification.

However, as the number of basis functions increases, the number of constraints can grow exponentially in the number of case statements in the constraint. To tackle this problem, we examine the underlying optimization problem in the next section.

#### 5.4. First-order linear programs

We can restate the FOALP problem as the optimal solution to a general *first-order linear program* (FOLP) for which we provide a generic solution. A FOLP is nothing more than a standard linear program where the constraints are written in terms of a sum of case statements whose case partition values may be specified as linear combinations of the weights. Efficiently solving FOLPs poses a number of difficulties—and we tackle these difficulties next.

##### 5.4.1. General formulation

A FOLP is specified as follows:

<sup>8</sup> One will note that the renaming of first-order formulae with “propositional” variables is in the same spirit as FOADDs. Consequently, we note that FOADDs prove to be an efficient method for representing and performing operations on the constraints that occur in FOALP.

---

```

input      : (1) A set  $C = \{case_1, \dots, case_n\}$ .
              (2) An ordering  $\langle R_1 \dots R_n \rangle$  of all relations in  $C$ .
output     : (1) The maximum value  $v$  achievable.
              (2) A set  $S = \{\langle \phi_i, t_i \rangle \in case_i\}$  for  $i = 1 \dots n$  s.t.  $v = t_1 + \dots + t_n$ .
begin
  // Convert  $C$  into CNF
  for ( $i = 1 \dots n$ ) do
    foreach ( $\langle \phi_i, t_i \rangle \in case_i(s)$ ) do
       $\perp$  Convert  $\phi_i$  to a set of CNF formulae.
  foreach (relation  $R \in \langle R_1 \dots R_n \rangle$  (in order)) do
    // Divide  $C$  into two sets of cases based on whether they contain  $R$ 
     $C_R^+ := \{case_i | case_i \in C \wedge \exists j. (\langle \phi_j, t_j \rangle \in case_i) \wedge \phi_j \text{ contains relation } R\}$ 
     $C_R^- := C \setminus C_R^+$ 
    // Build explicit “cross-sum”  $\oplus$  of its cases & convert to CNF
     $case_R^+ := \bigoplus_{case_i \in C_R^+} case_i$ 
    foreach ( $\langle \phi_j, t_j \rangle \in case_R^+$ ) do
       $\perp$  Convert  $\phi_j$  to CNF.

    foreach ( $\langle \phi_j, t_j \rangle \in case_R^+$  in order from highest to lowest value) do
      Resolve all clauses in  $\phi_j$  on relation  $R$  until quiescence or inference limit.
      // All resolvents on  $R$  derived so further resolution on these clauses
      // cannot lead to the empty clause — thus clauses can be removed
      Remove all clauses in  $\phi_j$  containing  $R$ 
      // Remove inconsistent partitions (i.e., those containing empty clause)
      if ( $\emptyset \in \phi_j$ ) then
         $\perp$  Remove  $\langle \phi_j, t_j \rangle$  from  $case_R^+$  and continue with next  $\langle \phi_j, t_j \rangle$ .
      // Remove  $\theta$ -subsumed partitions that are dominated
      foreach ( $\langle \phi_i, t_i \rangle \in case_R^+$  where  $t_i > t_j$ ) do
        if ( $\phi_j \preceq_\theta \phi_i$ ) then
           $\perp$  Remove  $\langle \phi_j, t_j \rangle$  from  $case_R^+$  and continue with next  $\langle \phi_j, t_j \rangle$ .
       $C := \{case_R^+\} \cup C_R^-$ 
   $v := 0$ ;  $S := \emptyset$ 
  foreach (maximal value partition  $\langle \phi_j, t_j \rangle$  of each case  $\in C$ ) do
     $\perp v := v + t_j$ ;  $S := S \cup$  all partitions from input  $C$  contributing to  $\langle \phi_j, t_j \rangle$ 
  Return  $\langle v, S \rangle$ .
end

```

---

**Algorithm 2.**  $FOMax(C, \langle R_1 \dots R_n \rangle) \longrightarrow \langle S, v \rangle$ .

Variables:  $w_1, \dots, w_k$ ;

Minimize:  $\sum_{i=1}^k c_i w_i$

Subject to:  $0 \geq case_{1,1}(\vec{w}, s) \oplus \dots \oplus case_{1,l(1)}(\vec{w}, s); \quad \forall s$

:

$0 \geq case_{m,1}(\vec{w}, s) \oplus \dots \oplus case_{m,l(n)}(\vec{w}, s); \quad \forall s$

(65)

The  $k$  variables  $\vec{w} = \langle w_1, \dots, w_k \rangle$  and objective weights  $\vec{c} = \langle c_1, \dots, c_k \rangle$  are defined as in a typical LP, the main difference being the form of the constraints. Here we have  $m$  different constraints of varying length  $l(j)$  (i.e., the number of case statements in constraint  $j$ ,  $1 \leq j \leq n$ ). We allow the  $t_i$  in each partition  $\langle \phi_i, t_i \rangle$  of  $case(\vec{w}, s)$  to be linearly dependent on the weights  $\vec{w}$  (e.g.,  $t_i = 3w_1 + 2w_2$ ). We note that the first-order LP for FOALP can be cast in this general form. As previously discussed in our FOALP example, we could simply compute the explicit “cross-sum”  $\oplus$  to flatten out each constraint  $j$  into a single case statement as in Eq. (64). However, this could be inefficient as it scales exponentially in the number of summed case statements. Fortunately, we can extend constraint generation methods used in factored MDPs [71] to the first-order case as we show next.

#### 5.4.2. First-order cost network maximization

In the constraint generation approach to solving a FOLP, the most important operation is to find a most-violated constraint given a current solution (i.e., setting of weights  $\vec{w}$ ). In this section, we formulate this problem as maximization over a first-order generalization of a cost network [23] represented as follows:

$$0 \geq \max_s [\text{case}_1(\vec{w}, s) \oplus \dots \oplus \text{case}_n(\vec{w}, s)] \quad (66)$$

The use of  $\max_s$  indicates that we are only interested in the single value (and corresponding case partitions contributing to this value) that maximizes the RHS. casemax would be less efficient here since it would exhaustively enumerate *all* values and constraints when we only require the single maximal value and constraint.

To determine the  $\max_s$  with this form of the constraints, we define the *FOMax* algorithm (see Algorithm 2) to carry out this computation. It is similar to *variable elimination* [90] or *bucket elimination* [23] (which makes a stronger connection to resolution), except that we use a simple ordered version of first-order resolution in place of propositional ordered resolution. Thus, we term this generalized variable elimination technique used by *FOMax* to be *relation elimination*.

Ostensibly, relation elimination and the technique of *first-order variable elimination* (FOVE) [20,21,60] appear similar since they both deal with lifted versions of variable elimination. However, they fundamentally apply to different problems: FOVE does not permit quantified formulae in its representation, while relational elimination permits full first-order logic in its representation; furthermore, FOVE permits the representation of indefinite products and sums whereas relation elimination only permits finite products and sums. Here we require full first-order logic, but not indefinite products or sums. While it is beyond the scope of this article to delve into a detailed discussion, we note that both relation elimination and FOVE can be combined when required; this occurs, for example, in FOALP approaches to factored FOMDP solutions (cf., [70] and Chapter 6 of [67]).

We provide a concrete example of *FOMax* and relation elimination in Fig. 7. Relation elimination proceeds analogously to variable elimination, except that we choose a relation  $R$  to eliminate at every step rather than a propositional variable. Elimination order can affect the time and space requirements of *FOMax* since eliminating  $R$  requires the “cross-sum”  $\oplus$  of all case statements containing  $R$ , incurring a polynomial blowup in the number of case statements being summed. In practice, we greedily eliminate the relation  $R$  at each step that minimizes this representational blowup, although this is not guaranteed to provide an optimal order.

On any elimination step of *FOMax*, once all of the case statements containing  $R$  have been explicitly “cross-summed,” the next step is to determine whether any case partitions are inconsistent (via resolution) or  $\theta$ -subsumed and dominated in value (using the generalized  $\theta$  subsumption operator  $\leq_\theta$  [18] with respect to our background theory, similar to the approach used by ReBel [43]); in both cases, these partitions may be removed since they will never contribute to the maximally consistent partition. Once all relations have been eliminated, maximal case partitions and their values extracted

Suppose we are given the following hypothetical constraint specification for a first-order linear program:

$$0 \geq \max_s \left( \begin{array}{c|c|c} \forall b, c. Dst(b, c) \supset BoxIn(b, c, s) : 10 & \exists b, c. Dst(b, c) \wedge \neg BoxIn(b, c, s) : w_1 & \exists t, c. TruckIn(t, c, s) : w_2 \\ \hline \neg \text{“} & \neg \text{“} & \neg \text{“} \\ \hline & : -w_1 & : 0 \end{array} \right) \oplus$$

Assume our last LP solution gave  $w_1 = 2$  and  $w_2 = 1$ . We can compute the most violated constraint (if one exists) by evaluating the weights in the constraint and applying *FOMax*. We begin by converting all first-order formulae to CNF where  $c_1, \dots, c_6$  are Skolemized constants. Formulae are negated prior to Skolemization and once in CNF are only resolved with each other (i.e., never negated).

$$0 \geq \max_s \left( \begin{array}{c|c|c} \{\neg Dst(b, c) \vee BoxIn(b, c, s)\} : 10 & \{Dst(c_3, c_4), \neg BoxIn(c_3, c_4, s)\} : 2 & \{TruckIn(c_5, c_6, s)\} : 1 \\ \hline \{Dst(c_1, c_2), \neg BoxIn(c_1, c_2, s)\} : 0 & \{\neg Dst(b, c) \vee BoxIn(b, c, s)\} : -2 & \{\neg TruckIn(t, c, s)\} : 0 \end{array} \right) \oplus$$

Assume elimination order  $BoxIn, Dst, TruckIn$ . First we eliminate  $BoxIn$ : we take the cross-sum  $\oplus$  of case statements containing  $BoxIn$ , repeatedly resolve clauses in each partition on  $BoxIn$  until quiescence, and remove all clauses containing  $BoxIn$  (indicated by struck-out text):

$$0 \geq \max_s \left( \begin{array}{c|c|c} \{\neg Dst(b, c) \vee \cancel{BoxIn(b, c, s)}, Dst(c_3, c_4), \neg \cancel{BoxIn(c_3, c_4, s)}, \neg Dst(c_3, c_4), \emptyset\} : 12 & & \\ \hline \{\neg Dst(b, c) \vee \cancel{BoxIn(b, c, s)}\} : 8 & & \{TruckIn(c_5, c_6, s)\} : 1 \\ \hline \{Dst(c_1, c_2), Dst(c_3, c_4), \neg \cancel{BoxIn(c_1, c_2, s)}, \neg \cancel{BoxIn(c_3, c_4, s)}\} : 2 & & \{\neg TruckIn(t, c, s)\} : 0 \\ \hline \{\neg Dst(b, c) \vee \cancel{BoxIn(b, c, s)}, Dst(c_1, c_2), \neg \cancel{BoxIn(c_1, c_2, s)}, \emptyset\} : -2 & & \end{array} \right) \oplus$$

Because the partitions valued 12 and  $-2$  contain the empty clause  $\emptyset$  (i.e., they are inconsistent), we can remove them. And because the partition of value 8 dominates the partition of value 2 (i.e.,  $2 < 8$  and the empty clause set of the value 8 partition trivially  $\theta$ -subsumes the clauses of the value 2 partition), we can remove it as well. This yields the following simplified result:

$$0 \geq \max_s \left( \begin{array}{c|c} \{\} : 8 & \{TruckIn(c_5, c_6, s)\} : 1 \\ \hline & \{\neg TruckIn(t, c, s)\} : 0 \end{array} \right) \oplus$$

From here it is obvious that the *Dst* elimination step will have no effect and the *TruckIn* elimination step will yield a maximal consistent partition with value 9. Since this is a positive value and thus a violation of the original constraint, we can generate the new linear constraint  $0 \geq 10 + -w_1 + w_2$  based on the original constituent partitions that led to this maximal constraint violation.

Fig. 7. An example use of *FOMax* to find the maximally violated constraint during first-order constraint generation.

from the remaining sum of case statements are used to generate the maximal value (and case partitions contributing to this value).

We note that the ordered resolution strategy we use in *FOMax* is not refutation-complete: it may loop indefinitely at an intermediate relation elimination step before finding a latter relation with which to resolve a contradiction. This is an unavoidable consequence of the fact that refutation resolution for general first-order theories is semi-decidable. From a practical standpoint, it is necessary to bound the number of resolutions performed at each relation elimination step (100 clauses per elimination step in our experiments) to prevent non-termination of *FOMax* due to an infinite number of resolutions. This incomplete theorem proving approach may generate unnecessary constraints corresponding to unsatisfiable regions of state space; while these constraints serve to overconstrain the set of feasible solutions, this has not led to infeasibility problems in practice. Furthermore, we often omit the generalized  $\theta$ -subsumption test  $\leq_\theta$  since the savings from this simplification does not outweigh its computational cost. This does not affect completeness since simplification is not required for inconsistency detection.

Finally, we remark that if the resolution procedure does finitely terminate before the inference limit is reached on every step of *FOMax*, then the conjunction of case partition formulae returned by *FOMax* is guaranteed to be satisfiable as a consequence of the completeness of refutation resolution. Research on decidable resolution procedures for expressive subsets of first-order logic [54] may pave the way for stronger completeness guarantees for *FOMax* in future work.

#### 5.4.3. First-order constraint generation

We can use the *FOMax* algorithm to find the maximal constraint violation when we have constraints of the form in Eq. (66). This allows us to define the following first-order constraint generation algorithm where we have specified some solution tolerance  $\epsilon$ :

- (1) Initialize LP with  $i = 0$ ,  $\vec{w}^i = \vec{0}$ , and empty constraint set.
- (2) For each constraint in the cost-network form of Eq. (66), find the maximally violated constraint  $C$  (if one exists) using the *FOMax* algorithm applied to the constraint instantiated with  $\vec{w}^i$ .
- (3) If  $C$ 's constraint violation is larger than  $\epsilon$ , add  $C$  to the LP constraint set, otherwise return  $\vec{w}^i$  as solution.
- (5) Solve LP with new constraints to obtain  $\vec{w}^{i+1}$ , goto step 2.

In first-order constraint generation, we initialize our LP with an initial setting of weights, but no constraints. Note that the initial weights  $\vec{w}^0 = \vec{0}$  will violate at least one constraint in a FOMDP with non-zero reward. Then we alternate between generating constraints based on maximal constraint violations at the current solution and re-solving the LP with these additional constraints. This process repeats until no constraints are violated and we have found the optimal solution. In practice, this approach typically generates *far* fewer constraints than the full exhaustive enumeration approach given by Eq. (64). To provide intuitions for this, we refer back to the example of finding the most violated constraint in Fig. 7.

Using first-order constraint generation, we now have a solution to the first-order LP from Eq. (65), thus providing a general solution for FOALP. At this point, the only step for FOALP that we have not automated is the generation of basis functions, which we discuss next.

#### 5.5. Automatic generation of basis functions

The effective use of linear approximations requires a “good” set of basis functions, one that spans a space containing a good approximation to the true value function. Previous work has addressed the issue of basis function generation in ground MDPs [51,57], while other work has addressed the inductive generation of first-order features or basis functions from sampled experience [83,86]. Here we consider a deductive first-order basis function generation method that draws on the work of [33]. Specifically, they use regressions of the reward as candidate basis functions for learning a value function. This technique has allowed them to generate fully or  $t$ -stage-to-go optimal policies for a range of BLOCKSWORLD problems.

We leverage a similar approach for generating candidate basis functions using regression, except that rather than use these candidate basis functions to learn a value function, we fit their weights without sampling or grounding by using FOALP. Algorithm 3 provides an overview of our basis function generation algorithm. The motivation for this approach is as follows: if some portion of state space  $\phi$  has value  $v > \tau$  in an existing approximate value function for some non-trivial threshold  $\tau$ , then this suggests that states that can reach this region (i.e., found by  $\text{Regr}(\phi)$  through some deterministic action) should also have reasonable value. However, since we have already assigned value to  $\phi$ , we want the new basis function to focus on the area of state space not covered by  $\phi$ ; thus we negate  $\phi$  and conjoin it with  $\text{Regr}(\phi)$ .

As a small example, given the initial *weighted* basis function  $bCase_1(s) = w_1 \cdot rCase(s)$  from BoxWORLD,

$$bCase_1(s) = w_1 \cdot \frac{\begin{array}{c} \exists b. \text{BoxIn}(b, \text{paris}, s) : 10 \\ \neg \text{“} : 0 \end{array}}{\quad}, \quad (67)$$

we derive the following *weighted* basis function from  $bCase_1(s)$  when considering deterministic action  $A_i = \text{unloadS}(b^*, t^*)$  during basis function generation:

$$bCase_2(s) = w_2 \cdot \frac{\begin{array}{c} \neg[\exists b. \text{BoxIn}(b, \text{paris}, s)] \wedge [\exists c. \text{BoxOn}(b^*, t^*, s) \wedge \text{TruckIn}(t^*, \text{paris}, s)] : 1 \\ \neg \text{“} : 0 \end{array}}{\quad} \quad (68)$$

---

**input** : (1) A FOMDP specification.  
           (2) a value threshold  $\tau$   
           (3) an iteration limit  $n$

**output** : A set  $B$  of basis functions  $bCase_i(s)$  and corresponding weights  $w_i$ .

**begin**

*// Note:  $rCase(s)$  may be a sum of cases, so we can start with many basis functions.*  
 $B := \{rCase(s)\}$

**for** ( $i = 1 \dots n$ ) **do**

**foreach** ( $bCase_i(s) \in B$ ) **do**

**foreach** ( $\langle \phi_i(s), t_i \rangle \in bCase_i(s)$ ) **do**

**foreach** (deterministic action  $n_j(\vec{x})$ ) **do**

$$B := B \cup \begin{array}{ll} \neg \phi_i \wedge \exists \vec{x} \text{ Repr}(\phi_i(do(n_j(\vec{x}), s))) & : 1 \\ \neg \text{“} & : 0 \end{array}$$

Solve for the weights  $\vec{w}$  using FOALP.

**foreach** ( $bCase_i(s) \in B$ ) **do**

**if** ( $w_i < \tau$ ) **then**

Discard  $bCase_i(s)$  from  $B$  and ensure it is not regenerated.

**if** (no new basis functions generated on this iteration) **then**

Return  $B, \vec{w}$ .

Return  $B, \vec{w}$ .

**end**

---

**Algorithm 3.**  $BasisGen(FOMDP, \tau, n) \rightarrow B$ .

If one examines the form of these two basis functions, the inherent “orthogonality” between the new basis functions and the ones from which they were derived allows for significant computational optimizations. For example, since the top partition of  $bCase_1(s)$  takes the form  $\phi_1$  and the top partition of  $bCase_2(s)$  takes the form  $\neg \phi_1 \wedge \phi_2$ , these two partitions are mutually exclusive and could never jointly contribute to the value of a state. Thus, when two basis functions are orthogonal in this manner, we can efficiently perform an explicit “cross-sum”  $\oplus$  on them to obtain a single *compact* case statement representing both weighted basis functions:

$$bCase_{1,2}(s) = bCase_1(s) \oplus bCase_2(s)$$

$$= \begin{array}{ll} \exists b. BoxIn(b, paris, s) & : w_1 \cdot 10 \\ \neg [\exists b. BoxIn(b, paris, s)] \wedge [\exists c. BoxOn(b^*, t^*, s) \wedge TruckIn(t^*, paris, s)] & : w_2 \\ \neg \text{“} & : 0 \end{array} \quad (69)$$

This style of basis function generation also has many computational advantages for FOALP. To see this, we return to our original discussion concerning the fact that the  $B^A[\cdot]$  operator as defined in Eq. (58) will not be able to preserve additive structure when all basis functions in the linear-value function representation are affected by the stochastic action  $A(\vec{x})$ . Recalling Property 5.2.2, if all basis functions are affected by  $A(\vec{x})$ , then the backup  $B^A[\cdot]$  of a sum of basis functions will require their explicit “cross-sum” since they will all have free variables  $\vec{x}$  causing them to be summed with  $\exists \vec{x}$  is applied. However, in the best case, if the explicit “cross-sum” was already pre-computed for orthogonal basis functions by merging them, then this blowup will not occur.

Of course, since different actions generate different non-orthogonal basis functions from the same “parent” basis function, it will not generally hold that all basis functions are pairwise orthogonal to each other. Nonetheless, if we can exploit the mutual orthogonality of *subsets* of the basis functions to efficiently carry-out their explicit “cross-sum”, then we can still achieve an exponential time speedup relative to the worst-case of the  $B^A[\cdot]$  operator that requires the explicit computation of the “cross-sum”. To see how subsets of basis functions can be efficiently summed, we refer back to Eq. (69), which provides an example sum of two orthogonal basis functions. In general, any mutually orthogonal subset of basis functions can be merged in this way.

As a consequence, we can exploit properties of orthogonal basis function generation in FOALP to mitigate exponential space and time scaling in the number of basis functions, where worst-case exponential scaling arises at various points due to the need to explicitly compute the “cross-sum” of the linear-value representation. While we do not claim this method of basis function generation will be appropriate for all domains, we will demonstrate that it works reasonably well for the stochastic planning problems evaluated in the next section.



## 5.6. Empirical results

We evaluated FOALP on PPDDL planning problems from the ICAPS 2004 [50] and ICAPS 2006 [32] International Probabilistic Planning Competitions (IPPC). We divide the discussion of results according to each competition in order to reflect the differences in the competition setup, the data collected, and the specific planners that entered each competition.

We used the Vampire theorem prover and the CPLEX 9.0 LP solver<sup>9</sup> in our FOALP implementation and applied *BasisGen* (Algorithm 3) to our FOMDP translation of these PPDDL domains, generated as described in Section 3.2.2. We additively decomposed universal rewards using the technique described in Section 4.2; we note that doing so prevents us from obtaining any approximation guarantees on the solution generated by FOALP.

We provided FOALP with additional background theory axioms that were not encoded in the PPDDL source: if a fluent was intended to have functional arguments in PPDDL (PPDDL does not make provisions for specifying this property explicitly), we provide a background axiom stating this. So, for example, in our running BoxWorld example, we would provide the following functional constraint axioms:

$$\begin{aligned} \forall b, c_1, c_2, s. \text{BoxIn}(b, c_1, s) \wedge \text{BoxIn}(b, c_2, s) &\supset c_1 = c_2 \\ \forall t, c_1, c_2, s. \text{TruckIn}(t, c_1, s) \wedge \text{TruckIn}(t, c_2, s) &\supset c_1 = c_2 \\ \forall b, t_1, t_2, s. \text{BoxOn}(b, t_1, s) \wedge \text{BoxOn}(b, t_2, s) &\supset t_1 = t_2 \end{aligned}$$

In words, these axioms state that a box can only be in one city, a truck can only be in one city, and a box can only be on one truck. Any search-based or inductive planner that is given an initial state respecting these constraints (which was always the case in the competition instances) would never have to consider such erroneous states violating these constraints since they are unreachable from non-erroneous states satisfying these constraints. However, FOALP has no initial state knowledge in its offline solution phase and will produce extremely poorly approximated value functions if it cannot rule out such erroneous states as being inconsistent.

The need for these constraints may be viewed as a major drawback of the FOALP approach and was the reason that, although FOALP entered the ICAPS 2006 Probabilistic Planning Competition, it did not compete on 6 of the 10 problem domains (since these 6 problem domains were released at the start of the competition and rules prevented the planners from being modified beyond this point). On the other hand, we note that functional constraints on fluents represent a minimal type of problem knowledge often easily encoded by the person specifying a PPDDL problem; the constraints for BoxWorld are a good example. As an aid to future non-grounding planners, we recommend that the capability to specify functional constraints on fluents be incorporated in future versions of the PPDDL specification. If such constraints are known to hold on all initial states, automated techniques based on reachability analysis could also be used to prove such constraints hold as well.

In the following sections, we present proof-of-concept results comparing FOALP to other planners across a sampling of problems where FOALP has been able to generate policies for IPPC problems.

### 5.6.1. ICAPS 2004 probabilistic planning competition problems

We applied FOALP to the BoxWorld logistics and BlocksWorld probabilistic planning problems from the ICAPS 2004 IPPC [50]. In the BoxWorld logistics problem, the domain objects consist of trucks, planes, boxes, and cities. The number of boxes and cities varied in each problem instance, but there were always 5 trucks and 5 planes. Trucks and planes are restricted to particular routes between cities in a problem instance-specific manner. The goal in BoxWorld was to deliver all boxes to their destination cities and there were costs associated with each action. The transition functions allowed for trucks and planes to stochastically end up in destinations other than that intended by the execution of their respective drive and fly actions. BlocksWorld is just a stochastic version of the standard domain where blocks are moved between the table and other stacks of blocks to form a goal configuration. In this version, a block may be dropped with some probability while picking it up or placing it on a stack.

We stopped our offline basis function generation algorithm after iteration 7 in *BasisGen* (Algorithm 3) taking less than 2 hours for both problems on a 2 Ghz Pentium with 2 Gb of RAM; iteration 8 could not complete due to memory constraints. We note that if we were not using the “orthogonal” basis function generation described in Section 5.5, we would not get past iteration 2 of basis function generation (the system does not terminate within 10 hours at iteration 3); thus, these optimizations have substantially increased the number of basis functions for which FOALP is a viable solution option.

We compared FOALP to the three other top-performing planners on these problems: *NMRDPP* is a temporal logic planner with human-coded control knowledge [77]; *mGPT* is an RTDP-based planner [10]; (*Purdue-Humans* is a human-coded planner, *Classy* is an inductive first-order policy iteration planner, and *FF-Replan* [85] (2004 version) is a deterministic replanner based on FF [39]. Results for all of these planners are given in Fig. 8.

Since FOALP was only able to complete 7 iterations of basis function generation, this effectively limits the lookahead horizon of our basis functions to 7 steps. A lookahead of 8 would be required to properly plan in the final BoxWorld problem instance and thus FOALP failed on this instance. It is important to note that in comparing FOALP to the other planners,

<sup>9</sup> <http://www.ilog.com/products/cplex/>.

Problem	Competing Probabilistic Planners					FOALP
	NMRDPP	mGPT	Humans	Classy	FF-Replan	
<i>bx c10 b5</i>	438	184	419	376	425	433
<i>bx c10 b10</i>	376	0	317	0	346	366
<i>bx c10 b15</i>	0	–	129	0	279	0
<i>bw b5</i>	495	494	494	495	494	494
<i>bw b11</i>	479	466	480	480	481	480
<i>bw b15</i>	468	397	469	468	0	470
<i>bw b18</i>	352	–	462	0	0	464
<i>bw b21</i>	286	–	456	455	459	456

**Fig. 8.** Cumulative reward of 5 planning systems and FOALP (100 run avg.) on the BoxWORLD and BLOCKSWORLD probabilistic planning problems from the ICAPS 2004 IPPC (– indicates no data). BoxWORLD problems are indicated by a prefix of *bx* and followed by the number of cities *c* and boxes *b* used in the domain. BLOCKSWORLD problems are indicated by a prefix of *bw* and followed by the number of blocks *b* used in the domain.

NMRDPP and Humans used hand-coded control knowledge. FF-Replan was a very efficient search-based deterministic planner that had a significant advantage because near-optimal policies in these specific goal-oriented problems can be obtained by assuming that the highest probability action effects occur deterministically and making use of classical search-based planning techniques. The only autonomous *fully stochastic* planners were mGPT and Classy (itself an inductive first-order planning approach), and FOALP performs comparably to both of these planners and outperforms them by a considerable margin on some problem instances.

### 5.6.2. ICAPS 2006 probabilistic planning competition problems

We now present results for FOALP on three problem domains from the ICAPS 2006 IPPC [32]: BLOCKSWORLD, TIREWORLD, and ELEVATORS.<sup>10</sup> In BLOCKSWORLD, there are blocks and a table and the goal is to stack and unstack blocks from each other in an effort to achieve a goal configuration of the blocks with respect to the table. TIREWORLD is a relatively simple problem where the goal is to drive from a goal city to a destination city, while being able to pick up a spare tire in some cities. One stochastic outcome of driving between cities is that a tire may go flat and can only be fixed when a spare tire is present. Thus, routes with cities that contain spare tires are preferred to other routes that do not. Finally, ELEVATORS is a problem with a grid-like state space. The horizontal dimension of the grid corresponds to positions on a floor and the vertical dimension corresponds to different floors. There may be elevators at each position that can move vertically between floors. An agent can occupy one position on one floor and can move left or right between positions or can move into or out of an elevator if it is at the appropriate floor or position. Any elevator can be moved up or down independently of whether the agent resides in it. There can be gates at certain positions, which probabilistically teleport the agent back to the start position of floor 1, position 1. Finally, there are a number of coins at different known positions and the goal is for the agent to retrieve them all.

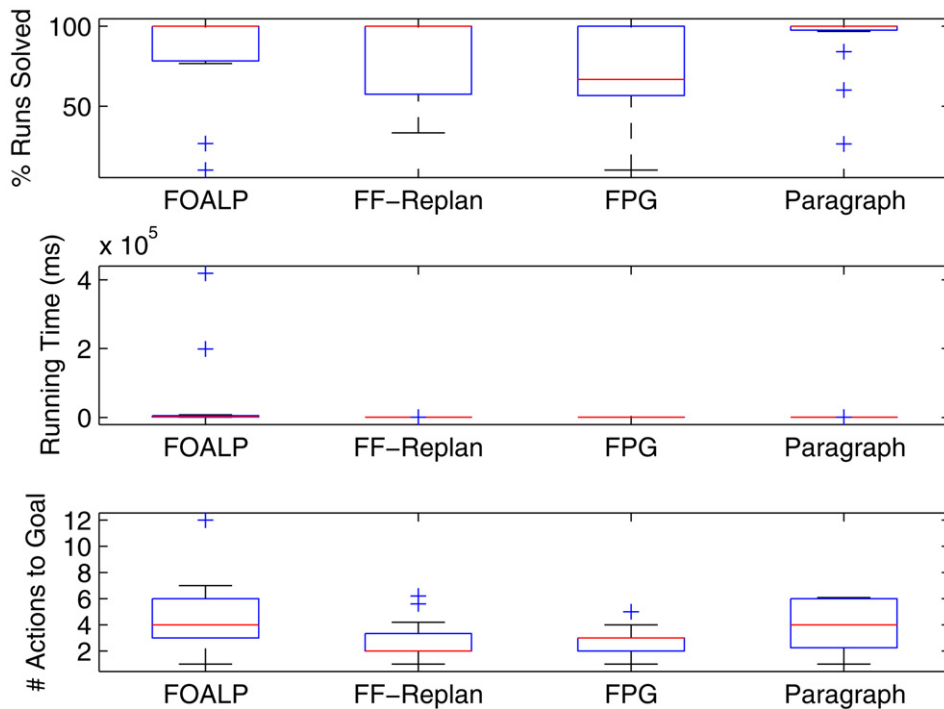
In all of the following results, *BasisGen* (Algorithm 3) was run for a four-hour fixed time limit on a 2 Ghz Pentium with 2 Gb of RAM to generate solutions for successively larger sets of basis functions. At the four-hour mark, we halted the solution process and used the largest (most recent) set of basis functions and weights for which FOALP had successfully terminated. Since the offline solution time of 4 hours can be amortized over an indefinite number of instances for a given problem, we do not report this in the online policy evaluation times in the following results.

In Figs. 9, 10, and 11, we provide data for FOALP and competing planners that specifies the number of problem instances solved, the online solution generation time, and the average number of actions required to reach the goal in each successful problem. We compare to the following planners that entered the competition<sup>11</sup>: (1) *FPG* [17], which uses policy gradient search in a factored representation of the Q-functions; (2) *sfdp* [76], which uses ADD-based dynamic programming [38] with reachability constraints based on initial state knowledge; (3) *Paragraph* [49], which uses a probabilistic extension of Graphplan [9] for probabilistic planning; (4) *FF-Replan* [88] (2006 version) is a deterministic replanner based on FF [39]. We note that all planners in this competition aside from FOALP are ground planners in that they use a propositional representation of a PPDDL problem for a specific domain instantiation.

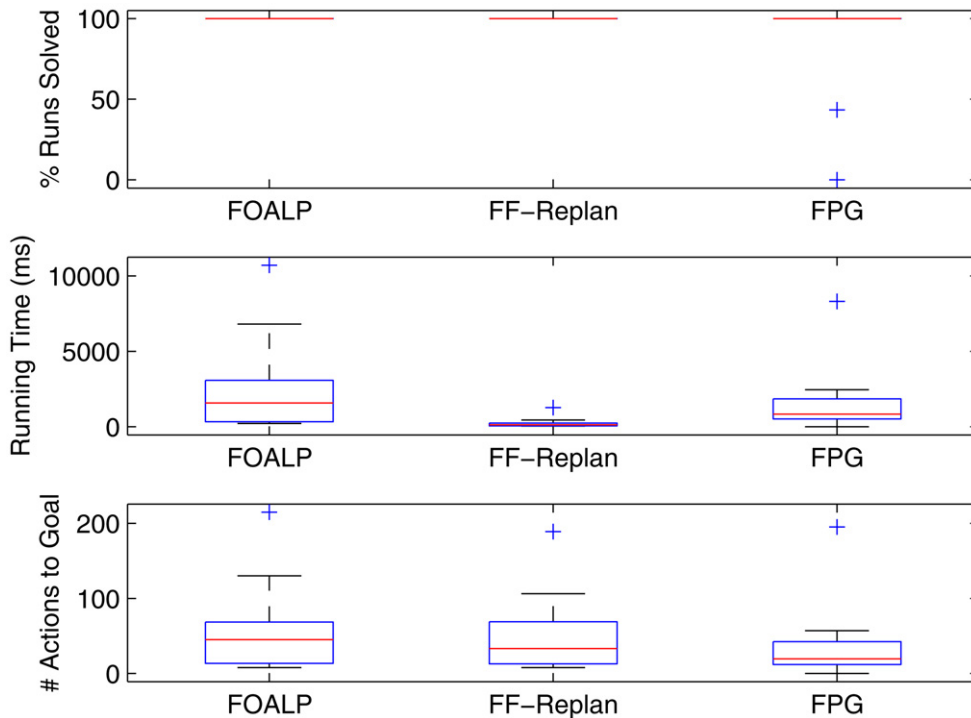
The results vary by problem, so we explain each in turn. In TIREWORLD, FOALP's policy allowed it to solve most problems although its policy was suboptimal in the number of actions and % problems solved in comparison to FF-Replan. In this case, it appears that the approximation inherent in the FOALP approach fared poorly in comparison to a deterministic replanner like FF-Replan that could perform nearly optimally on this problem. FOALP's slow policy evaluation on this problem is due to the transitive nature of the road connection topology and the lack of optimization in FOALP's logical policy evaluator. In ELEVATORS, the top three planners including FOALP all performed comparably with the deterministic replanner performing

<sup>10</sup> In the ICAPS 2006 IPPC, FOALP ran on the three problems reported here as well as EXPLODING-BLOCKSWORLD (not reported here). We do not report the EXPLODING-BLOCKSWORLD results since the competition version of the FOALP planner was restricted to use only the BLOCKSWORLD subset of the EXPLODING-BLOCKSWORLD problem description. In this section, we only show results for problems where FOALP was able to generate a policy for the full problem description.

<sup>11</sup> Not all planners ran on all of the problems in the competition. Furthermore, some planners did not provide results on all problem instances, this is noted for each result plot.

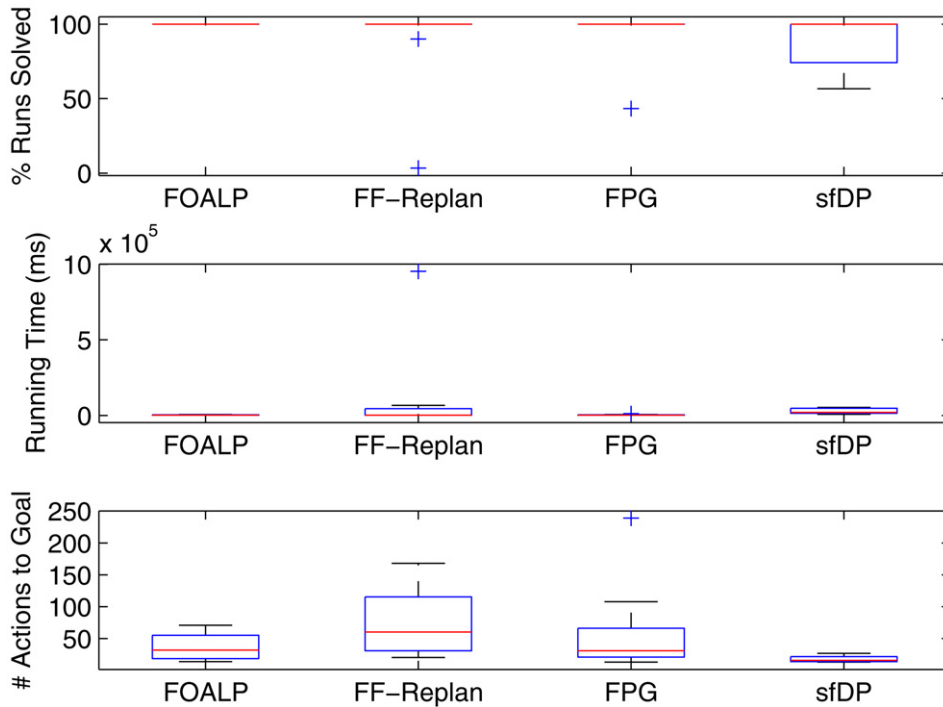


**Fig. 9.** A boxplot of performance of four planners on 15 instances of the TIREWORLD problem domain from the probabilistic track of the ICAPS 2006 IPPC. sfDP did not produce results for this problem; all other planners reported results for all instances.



**Fig. 10.** A boxplot of performance of three planners on 15 instances of the ELEVATORS problem domain from the probabilistic track of the ICAPS 2006 IPPC. sfDP and Paragraph did not produce results for these problems; FF-Replan and FPG did not report results for 2 and 3 problem instances, respectively.

consistently faster than the others, again due to the suitability of this domain for deterministic replanning and the relative speed of that approach. The goals in this domain are highly decomposable and FOALP thus benefited substantially from its additive goal decomposition approach. In BLOCKSWORLD, FOALP shows the best performance, solving more problems,



**Fig. 11.** A boxplot of performance of four planners on 15 instances of the BlocksWorld problem domain from the probabilistic track of the ICAPS 2006 IPPC. Paragraph did not produce results for these problems; FF-Replan, FPG, and sfDP did not report results for 1, 5, and 10 problem instances, respectively.

taking less time on the hard instances (FPG did not report results for the 5 hardest instances, thus skewing its results), and reaching the goal with the fewer actions (sfDP did not report results for the 10 hardest instances, thus skewing its results). In this case, FOALP's performance owes to two advantages: (1) first-order abstraction in BlocksWorld considerably helps the system avoid much of the combinatorial complexity that the ground planners face, and (2) the additive goal decomposition, although not optimal for all BlocksWorld problems, performed very well on these problem instances.

### 5.6.3. Summary of results

In summary, the first-order representation of FOALP seems to offer robust performance across a range of domain instance sizes and problems. However, as discussed at the end of Section 4.1.3, the case representation used by FOALP is a limiting factor in its performance due to its inability to exploit value structure in problems requiring reasoning about universal rewards (for which suboptimal additive reward decomposition techniques were used) or transitive reachability (for which the deficiency is quite clear from the TIREWORLD results). We discuss potential research directions to mitigate these observed deficiencies in Section 7.1.

## 6. Related work

In this section, we review work related to that presented in this article across two important dimensions: deductive first-order decision-theoretic planners based on symbolic dynamic programming (SDP), and inductive lifted decision-theoretic planners based on learning first-order representations of value functions, control knowledge, or policies from grounded domain instantiations.

### 6.1. Variants of symbolic dynamic programming

There have been a variety of alternative exact approaches to solving relationally specified MDPs without grounding in the spirit of SDP. Each of these approaches apply an SDP-like algorithm to their own first-order MDP representation. Like SDP, these algorithms all have guarantees on domain-independent error bounds for the value functions they produce and can produce exact domain-independent value functions when they exist. However, all of these approaches are restricted to solve less expressive variants of relational MDPs than SDP as we describe below.

First-order value iteration (FOVIA) [40,42] and the Relational Bellman algorithm (ReBel) [43] are value iteration algorithms that solve a restricted subclass of relational MDPs, most notably disallowing combined *universal conditional* effects (as defined in Section 3.1.1). Since universal conditional effects are a powerful planning formalism underlying the ADL extension to STRIPS, it can be argued that this is a significant limitation of these alternate SDP approaches. Both have provided

fully automated proof-of-concept results; we were able to directly compare SDP with FOADDs and ReBel on the BoxWORLD problem in Section 4.1.3. ReBel's specialization for a less expressive subset of FOMDPs (still capturing BoxWORLD, however) results in a substantial performance edge for this problem although both produce the same, exact solution. Results for ReBel and FOVIA are not available for the specific versions of the planning competition domains that we examined in Section 5.

First-order decision diagrams (FODDs) [81] have been introduced to compactly represent case statements and to permit efficient application of symbolic dynamic programming operations to another restricted class of relational MDPs via value iteration [80] and policy iteration [82].

Since FODDs are very similar in spirit to the FOADDs we defined in Section 4.1, we enumerate some of the major differences between these two formalisms:

- (1) FODDs disallow explicit universal quantification. This prevents FODDs from being applied to relational MDPs with universal preconditions or alternating quantifiers in their effects, although importantly, they *can* handle universal conditional effects.
- (2) Unlike FOADDs, which are maintained in a canonical form, FODDs are maintained in a sorted format, but are not guaranteed to be in a canonical form. As such, they rely on a range of simplification rules to maintain compact representations. This approach has the advantage that some diagrams without a strict order can be exponentially more compact than diagrams with a strict order [80]. However, rather than having a well-defined simplification algorithm leading to a canonical form, simplification in FODDs is somewhat open-ended and heuristic.
- (3) There is no need to reorder internal decision nodes after *Regr* in FODDs in order to maintain a canonical form. In this way, *Regr* is more efficient in FODDs than in FOADDs. This results in value and policy iteration algorithms that can be performed completely in terms of FODDs, unlike the current FOADD representation.
- (4) FODDs assume an implicit semantics where the maximal value is assumed for all instantiations of the free variables, thus precluding the need to perform explicit  $\exists x$  and  $\text{casemax}$ . In FOADDs, such operations would need to be performed explicitly. As such, the use of FODDs can lead to very compact representations for decision-theoretic planning, but this semantics may interfere with extensions of FODDs to handle universally quantified formulae.

Consequently, FODDs represent an interesting alternative in the design space of data structures for the compact representation of case statements. Nonetheless, the major limitation with respect to the work we present in this article is their limitations w.r.t. representing some forms of universal quantification. Ideally the best approach would be to combine the advantages of FOADDs with those of FODDs. This is a non-trivial problem, however, and an interesting future research direction.

## 6.2. Alternative lifted approaches to decision-theoretic planning

There are many alternative approaches to first-order decision-theoretic planning that reason inductively about sample domain instances and sample trajectories to produce lifted value functions or policies. This stands as an alternative to reasoning symbolically about actions and rewards directly at a first-order level without grounding as done in this article.

In one class of approaches, sampled experience from grounded domain instantiations is used to directly induce relational representations of value or Q-functions in a reinforcement learning approach. This can be done with pure reinforcement learning using relational decision or regression trees to learn a value or Q-function [25], combining this with supervised guidance [24], or using Gaussian processes and graph kernels over relational structures to learn a value or Q-function [31].

A second approach uses experience sampled from ground domain instantiations to induce first-order policy representations. In one version, policies can be learned directly from sampled experience trajectories generated using other planners [84]. In a different vein, policies can be learned in an approximate policy iteration framework [87] that combines trajectory sampling with policy updates derived from these trajectories. In this approach, sample experience trajectories can be generated using planning heuristics [26] and/or random walks on problem sizes that are adaptively scaled as planner performance improves [27].

A third inductive approach (that could also be used in conjunction with FOALP) allows first-order features to be learned from experience rather than symbolically deriving them directly from the relational MDP specification as described in Section 5.5. In one approach, heuristic control knowledge represented in a first-order taxonomic syntax can be learned from solution trajectories on a given problem [86]. In another recent approach, relational basis functions can be learned from sampled trajectories and then used in an approximate value iteration framework [83].

Since the approaches in this subsection also produce first-order value functions or policies, it is important to compare and contrast them with the symbolic deductive approach we adopt. In this approach, our *ideal* objectives are threefold:

- (1) Obtaining domain-independent exact or bounded approximate solutions where possible while exploiting natural relational and first-order planning structure.
- (2) Avoiding potential pitfalls of value functions and policies specific to biases from (small) sampled domain instantiations.
- (3) Avoiding an intractable representational blowup by grounding in the solution algorithm.

In *practice*, the approaches advocated in this article are unable to effectively achieve objective (1): the heuristics (such as universal reward decomposition from Section 4.2) required to apply our techniques to planning competition problems prevent the derivation of bounds. Objective (2) may be met in practice, although the approximations required for practical applications introduce their own representational biases. Finally, objective (3) may also be satisfied in practice, although the domain-independent approach introduces its own representational blowup by effectively planning for every possible domain instantiation.

In comparison, inductive first-order approaches outlined above share a goal similar to (1) in exploiting natural relational planning structure in a domain-independent manner, but cannot claim to support (2) since they must sample. Theoretical complexity results by [44,45] indicate that (3) is indeed possible to achieve for inductive approaches in some settings. We further note that in practice, the bias and computational complexity inherent in sampling a small set of possible ground domain instantiations of an MDP is not generally problematic since policies that work on one domain instantiation often generalize to similar or larger domains given an appropriate representation language [87].

So we may then ask: which first-order approach is better, inductive or deductive?<sup>12</sup> Empirically, recent results [83] show that inductive first-order approaches outperform FOALP. Is this the final answer? Hopefully not; but clearly there is still a great deal of work to be done in order to make first-order deductive approaches fully competitive with recent state-of-the-art first-order inductive approaches. Perhaps even more promising though is the potential to combine advances among both approaches; [33] do this in work that combines inductive logic programming with first-order decision-theoretic regression, showing that optimal policies can be induced from few training samples if using deductive methods to generate candidate policy structure. Such approaches offer the hope of combining the best of both worlds while sharing the goal of exploiting first-order structure in relational decision-theoretic planning problems.

## 7. Future directions and concluding remarks

In this article, we have motivated the need to exploit relational structure in decision-theoretic planning problems. To this end, we have provided a thorough review of the FOMDP representation of [14] and showed how to translate an expressive subset of PPDDL to this particular FOMDP representation. We reviewed the solution of FOMDPs via symbolic dynamic programming and contributed additional practical solution techniques based on the use of first-order ADDs (FOADDs), additive value decomposition of universal rewards, and first-order approximate linear programming (FOALP). Combining all of these ideas, we have provided proof-of-concept results from the probabilistic track of the ICAPS 2004 and 2006 International Planning Competitions.

We outline some interesting directions for future work, and offer some concluding remarks on decision-theoretic planning in the framework of FOMDPs.

### 7.1. Future directions

There are a number of open issues raised by our work that merit further exploration. We enumerate a few of them:

- (1) An interesting approach for the practical application of FOMDPs to decision-theoretic planning is to combine their approximate offline solution with online methods for enhancing their performance. We need only look at the range of successful planners used in planning competitions for ideas. Perhaps one of the most useful approaches would be to use offline methods for solving FOMDPs to generate a first-order approximated value function. Then we could use such a value function as a heuristic seed for online search methods such as RTDP [5,22]. Another approach would be to consider domain-specific control knowledge encoded as temporal logic constraints as in TLPlan [3], program constraints as in Golog [48] (both TLPlan and Golog are deterministic planners) or decision-theoretic extensions such as DT-Golog [15]. We discuss the use of program constraints further in a moment.
- (2) We did not explore approximate extensions of value iteration for FOMDPs. Given the success of the APRICODD planner [75] that performs approximate value iteration using ADDs, this approach is quite appealing for first-order approximate value iteration using FOADDs. When the FOADD representing the value function becomes too large, we can simply prune out nodes in the FOADD in an effort to reduce the size of the value function while minimizing the approximation error.
- (3) One promising use of FOMDPs is at the highest level of an abstraction hierarchy for agent-based decision-theoretic planning. [22] demonstrate that an MDP model can be approximated to a structure that is efficiently solvable and that error bounds can be obtained on the resulting optimal policy in the abstracted model with respect to the optimal policy in the non-abstracted version. If we lift such results to FOMDPs, then this offers a very appealing paradigm for their use: we can approximate a general FOMDP model to a level that we know we can solve efficiently while obtaining error bounds on the performance of the optimal policy in this approximated model. Or, further afield, we can use a solution

<sup>12</sup> To clarify, we use the term inductive to refer to any algorithm with an inductive component. However, it should be noted that all of the inductive approaches mentioned above incorporate some form of deduction by sampling from the Bellman equations then using induction to obtain a symbolic representation from these samples. In contrast the SDP and FOALP approaches advocated in this article can be viewed as pure symbolic deduction since they deduce their value representations from a lifted version of the Bellman equation.

to this approximated model as guidance for other more computationally expensive algorithms like ground heuristic search or as seed values [22] or shaped rewards [55] for value iteration in the non-abstracted MDP model.

In addition to these immediate open problems posed by our techniques, we have only touched on the surface of FOMDPs and the vast array of stochastic decision processes and symbolic solution methods that are possible. There remain a number of promising directions for the exploitation of structure in relationally-specified decision-theoretic planning problems that we briefly describe here:

- (1) One of the original goals in the FOMDP and symbolic dynamic programming frameworks [14] was to allow for very general symbolic representations. While most current FOMDP research has assumed a constant numerical representation of the values in case statement partitions, there are many situations where we might obtain non-constant values in our case statements, e.g., compactly representing value functions in FOMDPs with universal rewards that depend on the count of objects satisfying a property in a given situation, or in the context of modeling continuous state properties, perhaps combined with discrete state properties in a first-order generalization of *hybrid MDPs* [34,37]. However, as the case statement is generalized to handle non-constant numerical representations, case operators like the casemax must be appropriately generalized to efficiently handle such value representations (see Section 6.2.3 of [67] for one example of such a casemax generalization). Furthermore, theorem provers must also be capable of reasoning about counting properties or (constrained) continuous variables in such symbolic case statement enhancements in order to detect the inconsistency of state partitions.
- (2) In many FOMDPs there is an element of underlying topological graph structure. For example, in logistics planning, this graph structure may involve the accessibility of different cities via roads and flight routes. Currently, this graph structure is not exploited by our solution methods. Yet its regularity, if known *a priori*, could likely be exploitable by solution methods that could “compile” out this graph structure. This approach would be far more advantageous than relying on the first-order case representation to extract relevant graph properties using the cumbersome specification of transitively composed relations (i.e.,  $\exists c_1, c_2. Road(c_1, c_2) \wedge \exists c_3. Road(c_2, c_3) \wedge \exists c_4. Road(c_3, c_4) \wedge \dots$ ).
- (3) We often have a predefined set of constraints on the behavior of an agent and we need to optimize the agent's policy with respect to those constraints. If we can specify the program constraints in the form of a Golog program [48], then we can generalize the hierarchy of abstract machines (HAM) architecture [1,56] to the case of solving FOMDPs with respect to Golog program constraints. Such a solution would permit the (approximately) optimal execution of an incompletely specified program over all possible domain-instantiations. Various approaches in the decision-theoretic DT-Golog framework [15,28] have provided an initial investigation into these ideas.

The above suggestions are but a few of the many possible extensions to the work presented in this article and first-order decision-theoretic planning in general.

## 7.2. Concluding remarks

For a few years immediately succeeding the publication of the symbolic dynamic programming solution [14] to relationally specified MDPs, this domain-independent non-grounding approach was disparaged as being unrealistic for practical applications due to the complexity of value functions or due to the need for logical simplification and theorem proving [30, 35,84]. While these are all in fact significant obstacles to be overcome in the practical application of first-order MDPs to decision-theoretic planning, this article has aimed to show that these obstacles are not insurmountable. It has provided a substantial step in the direction of demonstrating that with careful attention paid to the first-order representation and algorithms specifically designed to exploit that representation, non-grounded lifted solutions are viable in practice as we demonstrated with our proof-of-concept results from the ICAPS 2004 and 2006 International Planning Competitions. Our hope is that this article lays the foundations for further exploration of these non-grounding approaches and permits the integration of these ideas with other lines of research in decision-theoretic planning.

## Acknowledgements

We are grateful to the three anonymous reviewers for their extensive comments and suggestions: these have vastly improved the presentation and discussion of our work. We are also grateful to Kee Siong Ng who provided many suggestions and corrections. This work was supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada. This research was conducted while the first author (now with NICTA) was at the Department of Computer Science, University of Toronto. NICTA is funded by the Australian Government's Backing Australia's Ability and the Centre of Excellence programs.

## References

- [1] D. Andre, S. Russell, Programmable reinforcement learning agents, in: *Advances in Neural Information Processing Systems (NIPS-01)*, vol. 13, 2001, pp. 78–85.
- [2] F. Bacchus, J.Y. Halpern, H.J. Levesque, Reasoning about noisy sensors in the situation calculus, in: *International Joint Conference on Artificial Intelligence (IJCAI-95)*, Montreal, 1995, pp. 1933–1940.

- [3] F. Bacchus, F. Kabanza, Using temporal logics to express search control knowledge for planning, *Artificial Intelligence* 116 (1–2) (2000) 123–191.
- [4] R.I. Bahar, E. Frohm, C. Gaona, G. Hachtel, E. Macii, A. Pardo, F. Somenzi, Algebraic decision diagrams and their applications, in: *IEEE/ACM International Conference on CAD*, 1993, pp. 428–432.
- [5] A.G. Barto, S.J. Bradtko, S.P. Singh, Learning to act using real-time dynamic programming, Tech. Rep. UM-CS-1993-002, U. Mass. Amherst, 1993.
- [6] R.E. Bellman, *Dynamic Programming*, Princeton University Press, Princeton, NJ, 1957.
- [7] D.P. Bertsekas, *Dynamic Programming*, Prentice Hall, Englewood Cliffs, NJ, 1987.
- [8] D.P. Bertsekas, J.N. Tsitsiklis, *Neuro-Dynamic Programming*, Athena Scientific, Belmont, MA, 1996.
- [9] A.L. Blum, M.L. Furst, Fast planning through graph analysis, in: *IJCAI 95*, Montreal, 1995, pp. 1636–1642.
- [10] B. Bonet, H. Geffner, mGPT: A probabilistic planner based on heuristic search, in: *Online Proceedings for The Probabilistic Planning Track of IPC-04*: <http://www.cs.rutgers.edu/~mlittman/topics/ipc04-pt/proceedings/>, 2004.
- [11] C. Boutilier, R.I. Brafman, C. Geib, Prioritized goal decomposition of Markov decision processes: Toward a synthesis of classical and decision theoretic planning, in: *International Joint Conference on Artificial Intelligence (IJCAI-97)*, Nagoya, 1997, pp. 1156–1162.
- [12] C. Boutilier, T. Dean, S. Hanks, Decision-theoretic planning: Structural assumptions and computational leverage, *Journal of Artificial Intelligence Research (JAIR)* 11 (1999) 1–94.
- [13] C. Boutilier, N. Friedman, M. Goldszmidt, D. Koller, Context-specific independence in Bayesian networks, in: *Uncertainty in Artificial Intelligence (UAI-96)*, Portland, OR, 1996, pp. 115–123.
- [14] C. Boutilier, R. Reiter, B. Price, Symbolic dynamic programming for first-order MDPs, in: *International Joint Conference on Artificial Intelligence (IJCAI-01)*, Seattle, 2001, pp. 690–697.
- [15] C. Boutilier, R. Reiter, M. Soutchanski, S. Thrun, Decision-theoretic, high-level agent programming in the situation calculus, in: *AAAI-00*, Austin, TX, 2000, pp. 355–362.
- [16] R. Brachman, H. Levesque, *Knowledge Representation and Reasoning*, Morgan Kaufmann Publishers Inc., San Francisco, CA, 2004.
- [17] O. Buffet, D. Aberdeen, The factored policy gradient planner (ipc-06 version), in: *Proceedings of the Fifth International Planning Competition*, 2006.
- [18] W. Buntine, Generalized subsumption and its application to induction and redundancy, *Artificial Intelligence* 36 (1988) 375–399.
- [19] D. de Farias, B.V. Roy, The linear programming approach to approximate dynamic programming, *Operations Research* 51 (6) (2003) 850–865.
- [20] R. de Salvo Braz, E. Amir, D. Roth, Lifted first-order probabilistic inference, in: *19th International Joint Conference on Artificial Intelligence (IJCAI-2005)*, Edinburgh, UK, 2005, pp. 1319–1325.
- [21] R. de Salvo Braz, E. Amir, D. Roth, MPE and partial inversion in lifted probabilistic variable elimination, in: *National Conference on Artificial Intelligence (AAAI-06)*, Boston, USA, 2006.
- [22] R. Dearden, C. Boutilier, Abstraction and approximate decision-theoretic planning, *Artificial Intelligence* 89 (12) (1997) 219–283.
- [23] R. Dechter, Bucket elimination: A unifying framework for reasoning, *Artificial Intelligence* 113 (1999) 41–85.
- [24] K. Driessens, S. Dzeroski, Integrating experimentation and guidance in relational reinforcement learning, in: *International Conference on Machine Learning (ICML)*, 2002, pp. 115–122.
- [25] S. Dzeroski, L. DeRaedt, K. Driessens, Relational reinforcement learning, *Machine Learning Journal (MLJ)* 43 (2001) 7–52.
- [26] A. Fern, S. Yoon, R. Givan, Approximate policy iteration with a policy language bias, in: *Advances in Neural Information Processing Systems 16 (NIPS-03)*, December 2003.
- [27] A. Fern, S. Yoon, R. Givan, Learning domain-specific control knowledge from random walks, in: *International Conference on Planning and Scheduling (ICAPS-04)*, June 2004, pp. 191–199.
- [28] A. Ferrein, C. Fritz, G. Lakemeyer, Extending DTGolog with options, in: *18th International Joint Conference on Artificial Intelligence (IJCAI-2003)*, Acapulco, Mexico, 2003, pp. 144–151.
- [29] R.E. Fikes, N.J. Nilsson, STRIPS: A new approach to the application of theorem proving to problem solving, *AI Journal* 2 (1971) 189–208.
- [30] N.H. Gardiol, L.P. Kaelbling, Envelope-based planning in relational MDPs, in: *Advances in Neural Information Processing Systems 16 (NIPS-03)*, Vancouver, CA, 2004, pp. 1040–1046.
- [31] T. Gartner, K. Driessens, J. Ramon, Graph kernels and Gaussian processes for relational reinforcement learning, *Machine Learning Journal (MLJ)* 64 (2006) 91–119.
- [32] A. Gerevini, B. Bonet, B. Givan (Eds.), *Online Proceedings for The Fifth International Planning Competition IPC-05*: <http://www ldc.usb ve/ bonet/ ipc5/ docs/ ipc-2006- booklet. pdf. gz>, Lake District, UK, 2006.
- [33] C. Gretton, S. Thiebaux, Exploiting first-order regression in inductive policy selection, in: *Uncertainty in Artificial Intelligence (UAI-04)*, Banff, Canada, 2004, pp. 217–225.
- [34] C. Guestrin, M. Hauskrecht, B. Kveton, Solving factored MDPs with continuous and discrete variables, in: *20th Conference on Uncertainty in Artificial Intelligence*, 2004, pp. 235–242.
- [35] C. Guestrin, D. Koller, C. Gearhart, N. Kanodia, Generalizing plans to new environments in relational MDPs, in: *18th International Joint Conference on Artificial Intelligence (IJCAI-2003)*, Acapulco, Mexico, 2003, pp. 1003–1010.
- [36] C. Guestrin, D. Koller, R. Parr, S. Venkatarman, Efficient solution methods for factored MDPs, *Journal of Artificial Intelligence Research (JAIR)* 19 (2002) 399–468.
- [37] M. Hauskrecht, B. Kveton, Linear program approximations for factored continuous-state Markov decision processes, in: *Advances in Neural Information Processing Systems 16*, 2004, pp. 895–902.
- [38] J. Hoey, R. St-Aubin, A. Hu, C. Boutilier, SPUDD: Stochastic planning using decision diagrams, in: *Uncertainty in Artificial Intelligence (UAI-99)*, Stockholm, 1999, pp. 279–288.
- [39] J. Hoffmann, B. Nebel, The FF planning system: Fast plan generation through heuristic search, *Journal of Artificial Intelligence Research (JAIR)* 14 (2001) 253–302.
- [40] S. Hölldobler, E. Karabaev, O. Skvortsova, FluCaP: A heuristic search planner for first-order mdps, *Journal of Artificial Intelligence Research (JAIR)* 27 (2006) 419–439.
- [41] R.A. Howard, *Dynamic Programming and Markov Processes*, MIT Press, 1960.
- [42] E. Karabaev, O. Skvortsova, A heuristic search algorithm for solving first-order MDPs, in: *Uncertainty in Artificial Intelligence (UAI-05)*, Edinburgh, Scotland, 2005, pp. 292–299.
- [43] K. Kersting, M. van Otterlo, L. de Raedt, Bellman goes relational, in: *International Conference on Machine Learning (ICML-04)*, ACM Press, 2004, pp. 465–472.
- [44] R. Khordon, Learning action strategies for planning domains, *Artificial Intelligence* 113 (1–2) (1999) 125–148.
- [45] R. Khordon, Learning to take actions, *Machine Learning* 35 (1) (1999) 57–90.
- [46] D. Koller, R. Parr, Computing factored value functions for policies in structured MDPs, in: *International Joint Conference on Artificial Intelligence (IJCAI-99)*, Stockholm, 1999, pp. 1332–1339.
- [47] D. Koller, R. Parr, Policy iteration for factored MDPs, in: *Uncertainty in Artificial Intelligence (UAI-00)*, Stockholm, 2000, pp. 326–334.
- [48] H.J. Levesque, R. Reiter, Y. Lespérance, F. Lin, R. Scherl, GOLOG: A logic programming language for dynamic domains, *Journal of Logic Programming* 31 (1–3) (1997) 59–83.



- [49] I. Little, Paragraph: A Graphplan-based probabilistic planner, in: Proceedings of the Fifth International Planning Competition, 2006.
- [50] M.L. Littman, H.L.S. Younes (Eds.), Online Proceedings for The Probabilistic Planning Track of IPC-04: <http://www.cs.rutgers.edu/mlittman/topics/ipc04-pt/proceedings/>, Vancouver, Canada, 2004.
- [51] S. Mahadevan, Samuel meets Amarel: Automating value function approximation using global state space analysis, in: National Conference on Artificial Intelligence (AAAI-05), Pittsburgh, 2005, pp. 1000–1005.
- [52] J. McCarthy, Situations, actions and causal laws, Tech. rep., Stanford University, 1963, reprinted, in: M. Minsky (Ed.), Semantic Information Processing, MIT Press, Cambridge, MA, 1968, pp. 410–417.
- [53] N. Meuleau, M. Hauskrecht, K.-E. Kim, L. Peshkin, L.P. Kaelbling, T. Dean, C. Boutilier, Solving very large weakly coupled Markov decision processes, in: National Conference on Artificial Intelligence (AAAI-98), Madison, WI, 1998, pp. 165–172.
- [54] B. Motik, Reasoning in description logics using resolution and deductive databases, Ph.D. thesis, Universität Karlsruhe (TH), Karlsruhe, Germany, January 2006.
- [55] A.Y. Ng, D. Harada, S. Russell, Policy invariance under reward transformations: theory and application to reward shaping, in: Proc. 16th International Conf. on Machine Learning, Morgan Kaufmann, San Francisco, CA, 1999, pp. 278–287.
- [56] R. Parr, S. Russell, Reinforcement learning with hierarchies of machines, in: M.M.K. Jordan, S. Solla (Eds.), Advances in Neural Information Processing Systems 10, MIT Press, Cambridge, MA, 1998, pp. 1043–1049.
- [57] R. Patrascu, P. Poupart, D. Schuurmans, C. Boutilier, C. Guestrin, Greedy linear value-approximation for factored Markov decision processes, in: National Conference on Artificial Intelligence (AAAI-02), Edmonton, 2002, pp. 285–291.
- [58] E.P.D. Pednault, ADL: Exploring the middle ground between STRIPS and the situation calculus, in: KR, 1989, pp. 324–332.
- [59] D. Poole, The independent choice logic for modelling multiple agents under uncertainty, Artificial Intelligence 94 (1–2) (1997) 7–56.
- [60] D. Poole, First-order probabilistic inference, in: IJCAI, 2003, pp. 985–991.
- [61] P. Poupart, C. Boutilier, R. Patrascu, D. Schuurmans, Piecewise linear value function approximation for factored MDPs, in: National Conference on Artificial Intelligence (AAAI-02), Edmonton, 2002, pp. 292–299.
- [62] M.L. Puterman, Markov Decision Processes: Discrete Stochastic Dynamic Programming, Wiley, New York, 1994.
- [63] R. Reiter, The frame problem in the situation calculus: A simple solution (sometimes) and a completeness result for goal regression, in: V. Lifschitz (Ed.), Artificial Intelligence and Mathematical Theory of Computation (Papers in Honor of John McCarthy), Academic Press, San Diego, CA, 1991, pp. 359–380.
- [64] R. Reiter, Knowledge in Action: Logical Foundations for Specifying and Implementing Dynamical Systems, MIT Press, 2001.
- [65] A. Riazanov, A. Voronkov, The design and implementation of vampire, AI Communications 15 (2) (2002) 91–110.
- [66] J. Rintanen, Expressive equivalence of formalisms for planning with sensing, in: 13th International Conference on Automated Planning and Scheduling, 2003, pp. 185–194.
- [67] S. Sanner, First-order decision-theoretic planning in structured relational environments, Ph.D. thesis, University of Toronto, Toronto, ON, Canada, March 2008.
- [68] S. Sanner, C. Boutilier, Approximate linear programming for first-order MDPs, in: Uncertainty in Artificial Intelligence (UAI-05), Edinburgh, Scotland, 2005, pp. 509–517.
- [69] S. Sanner, C. Boutilier, Practical linear evaluation techniques for first-order MDPs, in: Uncertainty in Artificial Intelligence (UAI-06), Boston, MA, 2006.
- [70] S. Sanner, C. Boutilier, Approximate solution techniques for factored first-order MDPs, in: 17th International Conference on Automated Planning and Scheduling (ICAPS-07), 2007, pp. 288–295.
- [71] D. Schuurmans, R. Patrascu, Direct value approximation for factored MDPs, in: Advances in Neural Information Processing 14 (NIPS-01), Vancouver, 2001, pp. 1579–1586.
- [72] P. Schweitzer, A. Seidmann, Generalized polynomial approximations in Markovian decision processes, Journal of Mathematical Analysis and Applications 110 (1985) 568–582.
- [73] L.S. Shapley, Stochastic games, Proceedings of the National Academy of Sciences 39 (1953) 327–332.
- [74] S.P. Singh, D. Cohn, How to dynamically merge Markov decision processes, in: Advances in Neural Information Processing Systems (NIPS-98), MIT Press, Cambridge, MA, 1998, pp. 1057–1063.
- [75] R. St-Aubin, J. Hoey, C. Boutilier, APRICODD: Approximate policy construction using decision diagrams, in: Advances in Neural Information Processing 13 (NIPS-00), Denver, 2000, pp. 1089–1095.
- [76] F. Teichteil, P. Fabiani, Symbolic stochastic focused dynamic programming with decision diagrams, in: Proceedings of the Fifth International Planning Competition, 2006.
- [77] S. Thiebaux, C. Grettton, J. Slaney, D. Price, F. Kabanza, Decision-theoretic planning with non-Markovian rewards, Journal of Artificial Intelligence Research 25 (January 2006) 17–74.
- [78] J.N. Tsitsiklis, B. Van Roy, Feature-based methods for large scale dynamic programming, Machine Learning 22 (1996) 59–94.
- [79] M. Veloso, Learning by analogical reasoning in general problem solving, Ph.D. thesis, Carnegie Mellon University, August 1992.
- [80] C. Wang, S. Joshi, R. Khardon, First order decision diagrams for relational MDPs, in: Twentieth International Joint Conference on Artificial Intelligence (IJCAI-07), Hyderabad, India, 2007, pp. 1095–1100.
- [81] C. Wang, S. Joshi, R. Khardon, First order decision diagrams for relational MDPs, Journal of Artificial Intelligence Research (JAIR) 31 (2008) 431–472.
- [82] C. Wang, R. Khardon, Policy iteration for relational MDPs, in: Uncertainty in Artificial Intelligence (UAI-07), Vancouver, Canada, 2007.
- [83] J. Wu, R. Givan, Discovering relational domain features for probabilistic planning, in: 17th International Conference on Automated Planning and Scheduling (ICAPS 2007), 2007, pp. 344–351.
- [84] S. Yoon, A. Fern, R. Givan, Inductive policy selection for first-order Markov decision processes, in: Uncertainty in Artificial Intelligence (UAI-02), Edmonton, 2002, pp. 569–576.
- [85] S. Yoon, A. Fern, R. Givan, Learning reactive policies for probabilistic planning domains, in: Online Proceedings for The Probabilistic Planning Track of IPC-04: <http://www.cs.rutgers.edu/mlittman/topics/ipc04-pt/proceedings/>, 2004.
- [86] S. Yoon, A. Fern, R. Givan, Learning measures of progress for planning domains, in: 20th National Conference on Artificial Intelligence, July 2005, pp. 1217–1222.
- [87] S. Yoon, A. Fern, R. Givan, Approximate policy iteration with a policy language bias: Learning to solve relational Markov decision processes, Journal of Artificial Intelligence Research (JAIR) 25 (2006) 85–118.
- [88] S. Yoon, A. Fern, R. Givan, FF-Replan: A baseline for probabilistic planning, in: 17th International Conference on Automated Planning and Scheduling (ICAPS-07), 2007, pp. 352–359.
- [89] H.L.S. Younes, M.L. Littman, D. Weissman, J. Asmuth, The first probabilistic track of the international planning competition, Journal of Artificial Intelligence Research (JAIR) 24 (2005) 851–887.
- [90] N.L. Zhang, D. Poole, A simple approach to Bayesian network computations, in: Proc. of the Tenth Canadian Conference on Artificial Intelligence, 1994, pp. 171–178.
- [91] N.L. Zhang, D. Poole, Exploiting causal independence in Bayesian network inference, Journal of Artificial Intelligence Research (JAIR) 5 (1996) 301–328.