

# Real Time Symbolic Dynamic Programming for Hybrid MDPs

Luis G. R. Vianna and Leliane N. de Barros  
 IME - USP  
 São Paulo, Brazil

Scott Sanner  
 NICTA & ANU  
 Canberra, Australia

## Abstract

Symbolic Dynamic Programming (SDP) has been recently proposed as an exact and complete solution for finite-horizon Hybrid Markov Decision Processes (HMDPs), i.e. an MDP with mixed continuous and discrete state-space and actions with continuous parameters. In SDP, piecewise continuous functions are represented as eXtended Algebraic Decision Diagrams (XADDs). The XADD representation of piecewise functions naturally factors the continuous state-space into regions that share the same expression and allows an efficient symbolic update of the value for each region. If the initial state is known, using SDP to update the value of the whole state-space is unnecessary and costly. In this case, the optimal policy can be computed only for regions containing relevant states, as it is done in Real Time Dynamic Programming (RTDP) for MDPs. In this work, we propose combining SDP with RTDP into a new HMDP exact solution, named Real Time Symbolic Dynamic Programming (RTSDP). Unlike conventional RTDP, our proposed algorithm updates are not restricted to a single state but are generalised to the region where the current state belongs, which allows easier reuse of previous updates. Our algorithm is empirically tested on a nonlinear domain in which the action set is finite and a continuous actions domain in which actions have continuous parameters. We empirically show that, given an initial state, RTSDP can solve these HMDP problems faster and using less memory than SDP.

## Introduction

Real-world planning applications frequently involve continuous resources (e.g. fuel, energy and money) and can be modelled as Hybrid Markov Decision Processes (HMDPs). Symbolic dynamic programming (SDP) ((Sanner, Delgado, and de Barros 2011; Zamani, Sanner, and Fang 2012)) provides an exact solution for HMDPs with discrete noise, using the eXtended Algebraic Decision Diagram (XADD) representation. The XADD structure is a generalisation of the algebraic decision diagram (ADD) (Bahar et al. 1993) to represent functions of both discrete and continuous variables. Decision nodes may contain boolean variables or inequalities between expressions containing continuous variables. Terminal nodes contain an expression on the continu-

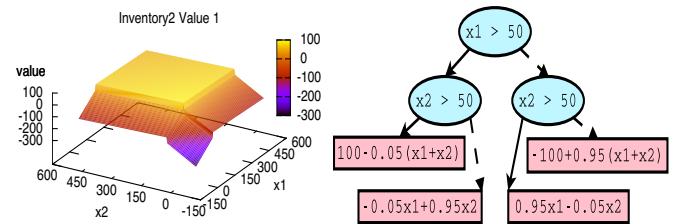


Figure 1: Illustration of a local reward function from the INVENTORY CONTROL domain: (left) graphical plot; (right) XADD representation. Blue ovals represent internal nodes, pink rectangles represent terminal nodes. The *true* branches are solid lines, *false* are dashed.

ous variables which could be evaluated to a real number. An example of XADD is shown in Figure 1.

**Example 1 - INVENTORY CONTROL (Scarf 2002)** An inventory control problem, consists of determining what items from the inventory should be ordered and how much to order of each. There is a continuous variables  $x_i$  for each item  $i$  and a single action order with continuous parameters  $dx_i$  for each item. There is a fixed limit  $L$  for the maximum number of items in the warehouse and a limit  $l_i$  for how much of an item can be ordered in a single action. The items are sold according to a stochastic demand, modelled by boolean variables  $d_i$  that represent whether demand is high ( $Q$  units) or low ( $q$  units) for item  $i$ . A reward is given for each sold item, but there are linear costs for items ordered and held in the warehouse. For item  $i$ , the reward is as follows:

$$R_i(x_i, d_i, dx_i) = \begin{cases} \text{if } d_i \wedge x_i + dx_i > Q : Q - 0.1x_i - 0.3dx_i \\ \text{if } d_i \wedge x_i + dx_i < Q : 0.7x_i - 0.3dx_i \\ \text{if } \neg d_i \wedge x_i + dx_i > q : q - 0.1x_i - 0.3dx_i \\ \text{if } \neg d_i \wedge x_i + dx_i < q : 0.7x_i - 0.3dx_i \end{cases}$$

A drawback of the SDP solution for HMDPs is the overhead of computing a complete optimal policy even for problems where the initial state is known and a partial solution would be sufficient. Though there exist previous solutions to HMDPs that use initial state information, they require monotonic resource consumption and aren't as general as SDP (Meuleau et al. 2009).

The Real Time Dynamic Programming (RTDP) algorithm is considered a state-of-the-art solver for MDPs ((Barto, Bradtke, and Singh 1995; Kolobov, Mausam, and Weld 2012)). It combines initial state information and a value function heuristic with asynchronous state value updates to generate an optimal partial policy closed for the relevant states, i.e., states reachable from the initial state following the optimal policy.

In this work, we propose a new HMDP solver, named Real Time Symbolic Dynamic Programming (RTSDP), which uses the symbolic representation of XADDs combined with RTDP to perform a heuristic search with asynchronous SDP, i.e. dynamic programming updating only some parts of the state space. As the XADD representation naturally factors the state-space into regions with the same continuous expressions, we perform *region based updates* that improve the value function not only for the current state, but *generalize this update to all states in the same region*. We claim that region updates promote better reusability of previous updates, improving the search efficiency.

Our algorithm is empirically tested on two challenging types of domains: nonlinear domains in which the action set is finite but the transition dynamic is nonlinear on state variables; and continuous actions domains in which actions have continuous parameters and the dynamics are linear on both state and parameter variables. We empirically show that, given an initial state, RTSDP can solve these HMDP problems faster and using less memory than SDP.

### Hybrid Markov Decision Processes

An MDP (Puterman 1994) is defined by a set of states  $\mathcal{S}$ , actions  $\mathcal{A}$ , a probabilistic transition function  $\mathcal{P}$  and a reward function  $\mathcal{R}$ . In an HMDP, we consider a model factored in terms of state variables (Boutilier, Hanks, and Dean 1999), i.e. a state  $s \in \mathcal{S}$  is an assignment vector of variables  $(\vec{b}, \vec{x}) = (b_1, \dots, b_n, x_1, \dots, x_m)$ , where each  $b_i \in \{0, 1\}$  ( $1 \leq i \leq n$ ) is boolean and each  $x_j \in \mathbb{R}$  ( $1 \leq j \leq m$ ) is continuous. We assume the action set  $A$  to be composed of a finite set of parametric actions, i.e.  $A = \{a_1(\vec{y}), \dots, a_K(\vec{y})\}$ , where each  $\vec{y}$  is the vector of parameter variables. The functions of state and action variables can be compactly represented if structural independencies among variables are exploited in dynamic Bayesian network (DBN) (Dean and Kanazawa 1990).

We assume that the factored transition model satisfies the following properties: (i) the next state boolean variables are defined by a probabilistic distribution that only depends on previous state variables and action; (ii) the continuous variables are deterministically dependent on previous state variables, the chosen action and current state boolean variables; (iii) transition functions are piecewise polynomial on the continuous variables. These assumptions allow us to write the transition function in terms of state variables, i.e.:

$$\mathcal{P}(\vec{b}', \vec{x}' | a(\vec{y}), \vec{b}, \vec{x}) = \prod_{i=1}^n \mathcal{P}(b'_i | a(\vec{y}), \vec{b}, \vec{x}) \prod_{j=1}^m \mathcal{P}(x'_j | \vec{b}', a(\vec{y}), \vec{b}, \vec{x}). \quad (1)$$

Assumption (ii) implies the conditional probabilities for continuous variables  $\mathcal{P}(x'_j | \vec{b}', a(\vec{y}), \vec{b}, \vec{x})$  are Dirac Delta

functions, which correspond to deterministic transitions,  $\vec{x}' \leftarrow \vec{T}_{a(\vec{y})}(\vec{b}, \vec{x}, \vec{b}')$ . While this restricts stochasticity only to boolean variables, continuous transitions depend on their sampled values, thus allowing the representation general finite distributions. This is a common restriction in exact HMDP solutions ((Feng et al. 2004; Meuleau et al. 2009; Sanner, Delgado, and de Barros 2011)).

The reward function  $\mathcal{R}(\vec{b}, \vec{x}, a(\vec{y}))$  specifies the immediate reward obtained by taking action  $a(\vec{y})$  in state  $(\vec{b}, \vec{x})$  and also is a piecewise polynomial function on the continuous variables.

In this paper, we consider HMDPs with an initial state  $s_0 \in \mathcal{S}$  and a finite horizon  $H \in \mathbb{N}$ . A non-stationary policy  $\pi$  is a function which specifies the action  $a(\vec{y}) = \pi(s, h)$  to take in state  $s \in \mathcal{S}$  at step  $h \leq H$ . The solution of an HMDP planning problem is an optimal policy  $\pi^*$  that maximizes the expected accumulated reward,  $V^{\pi^*}(s_0)$ :

$$V^{\pi^*}(s_0) = \mathbb{E} \left[ \sum_{t=1}^H \mathcal{R}(s_t, a_t(\vec{y}_t)) \middle| s_{t+1} \sim \mathcal{P}(s' | s_t, a_t(\vec{y}_t)) \right], \quad (2)$$

where  $a_t = \pi^*(s_t, t)$  is action chosen by  $\pi^*$  at step  $t$ .

### Symbolic Dynamic Programming

Before explaining RTSDP, we will first describe the SDP algorithm from (Zamani, Sanner, and Fang 2012) and then explain how the XADD representation allows it to be efficiently performed.

#### SDP Algorithm

The symbolic dynamic programming (SDP) algorithm is a generalisation of the classical *value iteration* dynamic programming algorithm (Bellman 1957) for constructing optimal policies for MDPs. It proceeds by constructing a series of  $h$  stages-to-go optimal value functions  $V_h(\vec{b}, \vec{x})$ . The pseudocode of SDP is shown in Algorithm 1. Beginning with  $V_0(\vec{b}, \vec{x}) = 0$  (line 2), it obtains the *quality*  $Q_{a(\vec{y})}(\vec{b}, \vec{x})$  for each action  $a(\vec{y}) \in \mathcal{A}$  (lines 5-8) in state  $(\vec{b}, \vec{x})$  by regressing the expected reward for  $h - 1$  future stages  $V_{h-1}(\vec{b}, \vec{x})$  along with the immediate reward  $R(\vec{b}, \vec{x}, a(\vec{y}))$  as the following:

$$Q_{a(\vec{y})}(\vec{b}, \vec{x}) = R(\vec{b}, \vec{x}, a(\vec{y})) + \sum_{\vec{b}'} \int_{\vec{x}'} \left[ \prod_i \mathcal{P}(b'_i | a(\vec{y}), \vec{b}, \vec{x}) \cdot \prod_j \mathcal{P}(x'_j | \vec{b}', a(\vec{y}), \vec{b}, \vec{x}) \cdot V_{h-1}(\vec{b}', \vec{x}') \right]. \quad (3)$$

Given  $Q_{a(\vec{y})}(\vec{b}, \vec{x})$  for each  $a(\vec{y}) \in \mathcal{A}$ , we can define the  $h$  stages-to-go value function as the quality of the best action in each state  $(\vec{b}, \vec{x})$  as follows (line 7):

$$V_h(\vec{b}, \vec{x}) = \max_{a(\vec{y}) \in \mathcal{A}} \left\{ Q_{a(\vec{y})}(\vec{b}, \vec{x}) \right\}. \quad (4)$$

The optimal value function  $V_h(\vec{b}, \vec{x})$  and optimal policy  $\pi_h^*$  for each stage  $h$  are calculated after  $h$  steps. Note that the policy is obtained from the maximisation in Equation 4,  $\pi_h^*(\vec{b}, \vec{x}) = \arg \max_{a(\vec{y})} Q_{a(\vec{y})}(\vec{b}, \vec{x})$  (line 8). The regression step (Procedure 1.1) computes Equation 3 using the XADD data structure.

**Algorithm 1:** SDP(HMDP M, H)  $\rightarrow (V^h, \pi_h^*)$  (Zamani, Sanner, and Fang 2012)

---

```

1  $V_0 := 0, h := 0$ 
2 while  $h < H$  do
3    $h := h + 1$ 
4   foreach  $a \in A$  do
5      $Q_{a(\vec{y})} \leftarrow \text{Regress}(V_{h-1}, a(\vec{y}))$ 
6    $V_h \leftarrow \max_{a(\vec{y})} Q_{a(\vec{y})}$  // Maximize all  $Q_{a(\vec{y})}$ 
7    $\pi_h^* \leftarrow \arg \max_{a(\vec{y})} Q_{a(\vec{y})}$ 
8 return  $(V_h, \pi_h^*)$ 

Procedure 1.1: Regress( $V_{h-1}, a(\vec{y})$ )
1  $Q = \text{Prime}(V_{h-1})$  // Rename all symbolic variables -
  //  $b_i \rightarrow b'_i$  and all  $x_i \rightarrow x'_i$ 
2 foreach  $b'_i$  in  $Q$  do
3   // Discrete marginal summation, also denoted  $\sum_{b'}$ 
4    $Q \leftarrow [Q \otimes P(b'_i | \vec{b}, \vec{x}, a(\vec{y}))] |_{b'_i=1}$ 
       $\oplus [Q \otimes P(b'_i | \vec{b}, \vec{x}, a(\vec{y}))] |_{b'_i=0}$ 
5 foreach  $x'_j$  in  $Q$  do
6   // Continuous marginal integration
7    $Q \leftarrow \text{subst}_{x'_j=T_a(\vec{y})(\vec{b}, \vec{x}, \vec{b}')} Q$ 
8 return  $Q \oplus R(\vec{b}, \vec{x}, a(\vec{y}))$ 

```

---

## XADD representation

The XADD data structure (Sanner, Delgado, and de Barros 2011) is used to efficiently represent and manipulate piecewise polynomial functions. An XADD is a directed acyclic graph with two kinds of nodes, internal and terminal. Internal nodes are decision nodes that contain two edges, namely *true* and *false* branches, and a decision that can be either a boolean variable or a polynomial inequality on the continuous variables. Terminal nodes contain a polynomial expression on the continuous variables. XADDs represent piecewise functions by using decision nodes to separate regions and terminal nodes for the local functions (i.e. a expression valid inside a region).

For example, Figure 1(left) shows the plot of a piecewise function and Figure 1(right) its XADD representation. Every path in an XADD defines a partial assignment of boolean variables and a continuous region that satisfies all inequalities in its internal nodes.

The SDP algorithm uses four kinds of XADD operations: (i) Algebraic operations (sum  $\oplus$ , product  $\otimes$ ) are naturally defined within each region, so the resulting partition is a cross-product of the original partitions and within each region the operation is straightforward. (See Figure 2a), (ii) Substitution operations (replacing or instantiating variables) affect both internal and terminal expressions and any decision that is no longer symbolic is replaced

by the corresponding branch.(See Figure 2b), (iii) Comparison operations (maximisation) cannot be performed directly between terminal expressions so a new decision node comparing terminal expressions is created.(See Figure 2c), (iv) Parameter maximization operations modify expressions substituting the maximal values and modify the structure adding comparisons between different branches (See Figure 2d). These XADD operations have been defined previously and we refer to (Sanner, Delgado, and de Barros 2011; Zamani, Sanner, and Fang 2012) for details.

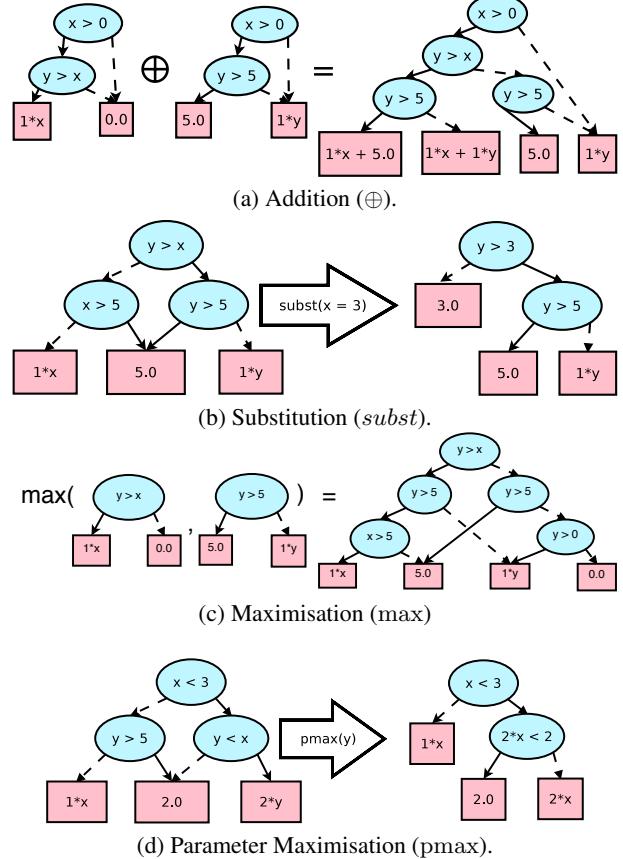


Figure 2: XADD Operations.

Thus, the SDP Bellman update (Equations 3 & 4) is performed synchronously as a sequence of XADD operations<sup>1</sup>, in two steps:

- Regression of all states:

$$Q_a^{DD}(\vec{b}, \vec{x}, \vec{y}) = R_a^{DD}(\vec{b}, \vec{x}, \vec{y}) \oplus \sum_{\vec{b}'} P_a^{DD}(\vec{b}, \vec{x}, \vec{y}, \vec{b}') \otimes \left[ \text{subst}_{(x' = T_a^{DD}(\vec{b}, \vec{x}, \vec{y}, \vec{b}'))} V_{h-1}'^{DD}(b', x') \right]. \quad (5)$$

- Maximization for all states w.r.t. parametrised actions:

$$V_h^{DD}(\vec{b}, \vec{x}) = \max_a \left( \text{pmax}_{\vec{y}} Q_a^{DD}(\vec{b}, \vec{x}, \vec{y}) \right). \quad (6)$$

<sup>1</sup>  $F^{DD(vars)}$  is a notation to indicate that  $F$  is an XADD symbolic function of  $vars$ .

As these operations are performed the complexity of the value function increases and so does the number of regions and nodes in the XADD. A goal of this work is to propose an efficient solution that avoids exploring unreachable or unnecessary regions.

## Real Time Symbolic Dynamic Programming

In this section we describe the RTSDP algorithm (Section 4.2), an extension of the RTDP algorithm for MDPs (Section 4.1) and the first asynchronous exact solution for HMDPs.

### The RTDP algorithm

The RTDP algorithm (Barto, Bradtke, and Singh 1995) (Algorithm 2) uses an initial admissible heuristic value function  $V$  and performs trials to improve it iteratively. A trial, described in Procedure 2.1, is a simulated execution of a policy interleaved with local Bellman updates on the visited states. Starting from a known initial state, an update is performed, an action is chosen and a next state is sampled and the trial continues from this new state. Trials are stopped when they reach the horizon  $H$ .

The update operation (Procedure 2.2) uses the current value function  $V$  to estimate the future reward and evaluates every action  $a$  on the current state in terms of its quality  $Q_a$ . The action with greatest quality, called greedy action  $a_g$ , is chosen and the state's value is set to  $Q_{a_g}$ .

RTDP performs updates on one state at a time and only on relevant states, i.e. states that are reachable under the greedy policy. Nevertheless, if the initial heuristic is admissible, there is an optimal policy whose relevant states are visited infinitely often in RTDP trials. Thus the value function converges to  $V^*$  on these states and an optimal partial policy closed for the initial state is obtained.

---

**Algorithm 2:** RTDP( $MDP\ M, s_0, H, V \rightarrow V_H^*(s_0)$

---

- 1 **while**  $\neg Solved(M, s_0, V) \wedge \neg TimeOut()$  **do**
- 2    $V \leftarrow MakeTrial(M, s_0, V, H)$
- 3 **return**  $V$

**1 Procedure 2.1:** MakeTrial( $MDP\ M, s, V, h$ )

- 2 **if**  $s \in GoalStates$  **or**  $h == 0$  **then**
- 3   **return**  $V$
- 4  $V_h(s), a_g \leftarrow Update(M, s, V, h)$
- 5 // Sample next state from  $P(s, a, s')$
- 6  $s^{new} \leftarrow Sample(s, a_g, M)$
- 7 **return**  $MakeTrial(M, s^{new}, V, h - 1)$

**1 Procedure 2.2:** Update( $MDP\ M, s, V, h$ )

- 2 **foreach**  $a \in A$  **do**
- 3    $Q_a \leftarrow R(s, a) + \sum_{s' \in S} P(s, a, s') \cdot [V_{h-1}(s')]$
- 4  $a_g \leftarrow \arg \max_a Q_a$
- 5  $V_h(s) \leftarrow Q_{a_g}$
- 6 **return**  $V_h(s), a_g$

**1 Procedure 2.3:** Sample( $MDP\ M, s, a_g$ )

- 2  $s' \sim P(s, a_g)$
- 3 **return**  $s'$

---

### The RTSDP Algorithm

We will now generalise the RTDP algorithm for HMDPs. The main RTDP (Algorithm 2) and `makeTrial` (Procedure 2.1) functions are trivially extended as they are independent of the state or action representation. The procedure that requires a non-trivial extension is the `Update` function (Procedure 2.2). Considering the continuous state space, it is inefficient to represent values for individual states, therefore we define a new type of update, named *region update*, which is one of the contributions of this work.

While the synchronous SDP update (Equations 5 and 6) modifies the expressions for all paths of  $V^{DD(\vec{b}, \vec{x})}$ , the region update only modifies a single path, the one that corresponds to the region containing the current state  $(\vec{b}_c, \vec{x}_c)$ , denoted by  $\Omega$ . This can be achieved by using a “mask”, which restricts the update operation to region  $\Omega$ . The “mask” is an XADD symbolic indicator of whether any  $(\vec{b}, \vec{x})$  belongs to  $\Omega$ , i.e.:

$$I[(\vec{b}, \vec{x}) \in \Omega] = \begin{cases} 0, & \text{if } (\vec{b}, \vec{x}) \in \Omega \\ +\infty, & \text{otherwise} \end{cases}, \quad (7)$$

where 0 indicates valid regions and  $+\infty$  invalid ones. The mask is applied in an XADD with a sum ( $\oplus$ ) operation, the valid regions stay unchanged whereas the invalid regions change to  $+\infty$ . The mask is applied to probability, transition and reward functions restricting these operations to  $\Omega$ , i.e.:

$$T_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}', \vec{x}')} = T_a^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}', \vec{x}')} \oplus I[(b, x) \in \Omega] \quad (8)$$

$$P_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}')} = P_a^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}')} \oplus I[(b, x) \in \Omega] \quad (9)$$

$$R_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y})} = R_a^{DD(\vec{b}, \vec{x}, \vec{y})} \oplus I[(b, x) \in \Omega] \quad (10)$$

Finally we can define the region update using the symbolic XADD operations with these restricted functions. In the following steps:

- Regression of region  $\Omega$ :

$$Q_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y})} = R_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y})} \oplus \sum_{\vec{b}'} P_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}')} \otimes \left[ subst_{(x' = T_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y}, \vec{b}')} )} V_{h-1}'^{DD(b', x')} \right]. \quad (11)$$

- Maximization on region  $\Omega$  w.r.t. parametrised actions:

$$V_{h,\Omega}^{DD(\vec{b}, \vec{x})} = \max_a \left( \underset{\vec{y}}{\operatorname{pmax}} Q_{a,\Omega}^{DD(\vec{b}, \vec{x}, \vec{y})} \right). \quad (12)$$

- Update of  $V$  restricted to region  $\Omega$ :

$$V_h^{DD(b, x)} = \begin{cases} V_{h,\Omega}^{DD(b, x)}, & \text{if } (\vec{b}, \vec{x}) \in \Omega \\ V_h^{DD(b, x)}, & \text{otherwise} \end{cases} \quad (13)$$

This case update is performed by a XADD minimisation:

$$V_h^{DD(b, x)} \leftarrow \min(V_{h,\Omega}^{DD(b, x)}, V_h^{DD(b, x)}), \quad (14)$$

where  $V_{h,\Omega}$  is  $+\infty$  for the invalid regions, and is lower on the valid regions, so that the minimum is indeed the correctly updated  $V_h$ .

This new update procedure is described in Algorithm 3. Note that the current state is denoted by  $(\vec{b}_c, \vec{x}_c)$  to distinguish it from symbolic variables  $(\vec{b}, \vec{x})$ .

---

**Algorithm 3:**

```

Region-update(HMDP,  $(\vec{b}_c, \vec{x}_c), V, h \rightarrow V_h, a_g, \vec{y}_g$ )


---


1 //GetRegion extracts XADD path for  $(\vec{b}_c, \vec{x}_c)$ )
2  $I[(\vec{b}, \vec{x}) \in \Omega] \leftarrow \text{GetRegion}(V_h, (\vec{b}_c, \vec{x}_c))$ 
3  $V'^{DD(b',x')} \leftarrow \text{Prime}(V_{h-1}^{DD(b,x)})$  // Prime variable names
4 foreach  $a(\vec{y}) \in A$  do
5    $P_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y},\vec{b}')} \leftarrow P_a^{DD(\vec{b},\vec{x},\vec{y},\vec{b}')} \oplus I[(b,x) \in \Omega]$ 
6    $T_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y},\vec{b}',\vec{x}')} \leftarrow T_a^{DD(\vec{b},\vec{x},\vec{y},\vec{b}',\vec{x}')} \oplus I[(b,x) \in \Omega]$ 
7    $R_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y})} \leftarrow R_a^{DD(\vec{b},\vec{x},\vec{y})} \oplus I[(b,x) \in \Omega]$ 
8    $Q_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y})} \leftarrow R_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y})} \oplus \sum_{\vec{b}'} P_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y},\vec{b}')} \otimes$ 
     $\left[ \text{subst}_{(x' = T_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y},\vec{b}')} )} V'^{DD(b',x')} \right]$ 
9    $\vec{y}_g^a \leftarrow \arg \max_{\vec{y}} Q_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y})}$  //Maximizing parameter
10   $V_{h,\Omega}^{DD(\vec{b},\vec{x})} \leftarrow \max_a (\max_{\vec{y}} Q_{a,\Omega}^{DD(\vec{b},\vec{x},\vec{y})})$ 
11   $a_g \leftarrow \arg \max_a (Q_{a,\Omega}^{DD(\vec{b},\vec{x})}(\vec{y}_g^a))$ 
12 //The value is updated through a minimisation.
13   $V_h^{DD(b,x)} \leftarrow \min(V_h^{DD(b,x)}, V_{h,\Omega}^{DD(b,x)})$ 
14 return  $V_h(b_c, x_c), a_g, \vec{y}_g^{a_g}$ 


---



```

## Empirical Evaluation

In this section, we evaluate RTSDP for two domains: INVENTORY CONTROL (Example 1) and MARS ROVER (Example 2). These domains show different challenges for planning as INVENTORY CONTROL contains continuously parametrised actions whereas MARS ROVER has a Nonlinear reward function. We show the efficiency of our algorithm by comparing it with synchronous SDP. The space efficiency is measured by the number of nodes in the XADD solution. The time efficiency is measured by the convergency of the value of the initial state  $s_0$ .

**Example 2 (MARS ROVER NonLinear Domain) from (Sanner, Delgado, and de Barros 2011)** This domain describes a planetary rover moving in a two dimensional plane  $(x, y)$ , and searching for  $k$  marked points from where it must take pictures. There are also  $k$  boolean  $tp_i$  variables, one for each point, indicating whether the picture of point  $i$  was already taken. There is a single move action in this domain that simply reduces the distance from the rover to a specific point by  $\frac{1}{3}$  of the current distance. For all experiments, this target point was set to  $(0, 0)$ . The intent of this action is to represent the fact that a rover may move progressively more slowly as it approaches a target position in order to reach the position with high accuracy. There are  $k$  take-pic $_i$  actions that provide reward according to nonlinear expressions over the continuous  $x$  and  $y$  variables, e.g. it is quadratically proportional to the distance from the picture point. Hence

for various points, the rover has to trade-off whether to take each picture at its current position or to get a larger reward by first moving and potentially getting closer before taking the picture. The reward for the take-pic actions for a point in  $(0,0)$  is:

$$R_{\text{take-pic}_i}(x, y, tp_i) = \begin{cases} \text{if } \neg tp_i \wedge x^2 + y^2 < 4 : 4 - x^2 + y^2 \\ \text{otherwise} : 0, \end{cases}$$

which shows that reward  $r_i > 0$  is received at the first time the rover chooses take-pic and its distance to  $(0,0)$  is less than 2.

In Figure 3 the SDP plot shows the optimal initial state value for different stages-to-go, i.e.  $V_h(s_0)$  is the value of the  $h$ -th point in the plot. The RTSDP plot is the estimated initial state value at the end of each trial, until the optimal value  $V_H(s_0)$  is found. Note that RTSDP reaches this optimal value much earlier than SDP and its value function has far fewer nodes. Thus, our solver can achieve the optimal value for the initial state much faster and using a significantly smaller amount of memory space than synchronous SDP.

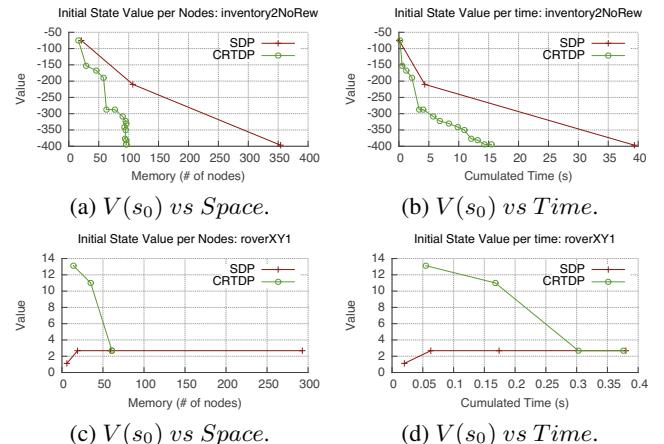


Figure 3: CRTDP and SDP performance comparison on two domains INVENTORY CONTROL (*top*) and MARS ROVER (*bottom*). Initial state value per memory (*left*) and time (*right*). SDP points are complete iterations (horizontal) and RTSDP points are complete trials.

Figures 4 and 5 show the value functions for different stages-to-go generated by RTSDP and SDP when solving the INVENTORY CONTROL instance 2, with  $H = 3$  and MARS ROVER instance 1, with  $H = 4$ . The SDP solutions (*right*) are the optimal value functions ( $V_h^{DD(\vec{b},\vec{x})}$ ) shown as a 3D plot (continuous state variables X state value) and as an XADD. The RTSDP solutions (*left*) are the value functions obtained after a number of trials, such that the initial state value has converged to the optimal. The RTSDP solver runs all trials from the initial state, so the value function for  $H$  stages-to-go,  $V_H$  (*bottom*) is always updated on the region containing the initial state. As the updates refine the partition, the region containing  $s_0$  shrinks. Note that once a region no longer contains the initial state its heuristic value

remains the same (See Figures 4i, 4j, 5i and 5j). In both domains, the use of reachability to prune updates is clear, there are plateaus on the unreachable regions that correspond to the initial heuristic, while the regions containing the initial state have a value similar to the optimal.

Another important comparison is how the solution size (number of XADD nodes) changes with the number of stages-to-go. For the SDP (Figures 4d, 4h, 4l, 5d, 5h and 5l), the complexity of the value function and solution size always increases quickly with horizon. However, for RTSDP (Figures 4b, 4f, 4j, 5b, 5f and 5j) there is different pattern: on one hand, as horizon increases the complexity of the solution in general does increase; On the other hand, if the number of performed actions is small, the number of reachable regions is also small and therefore the number of different terminal nodes and overall XADD size is not so large for the greater horizons.

A peculiarity of this algorithm is that even though its solution is focused to obtain the initial state correct value, it provides useful information on many other states and stages-to-go, what makes it possible to be reused.

## Related Work

Previous exact solutions for HMDP were often limited to hyper-rectangular partitions ((Feng et al. 2004; Li and Littman 2005)). SDP (Sanner, Delgado, and de Barros 2011; Zamani, Sanner, and Fang 2012) allows general piecewise functions, but generates an expensive to compute complete policy.

The most relevant existing heuristic search solution for HMDPs is given in (Meuleau et al. 2009); however, its requirement of monotonic use of resources and hyper-rectangular regions is more restrictive than our model. As far as we know, this is the only work to use a continuous region based Bellman update operation, and that uses heuristic search to compute a partial policy that generalizes beyond visited states.

## Conclusion

The proposed Continuous RTDP algorithm can use heuristic search to speed up the solution of DC-MDPs in exchange for a solution focused on the relevant states. This is the first algorithm to use asynchronous symbolic dynamic programming and heuristic search. For future work, there are many possibilities of improvement: better convergency and anytime solution quality are likely to be obtained with the labelling (Bonet and Geffner 2003) and online (Kolobov, Mausam, and Weld 2012) modifications. Also, a very important factor to improve performance that should be explored is to use more informative heuristics that can greatly reduce the number of relevant regions.

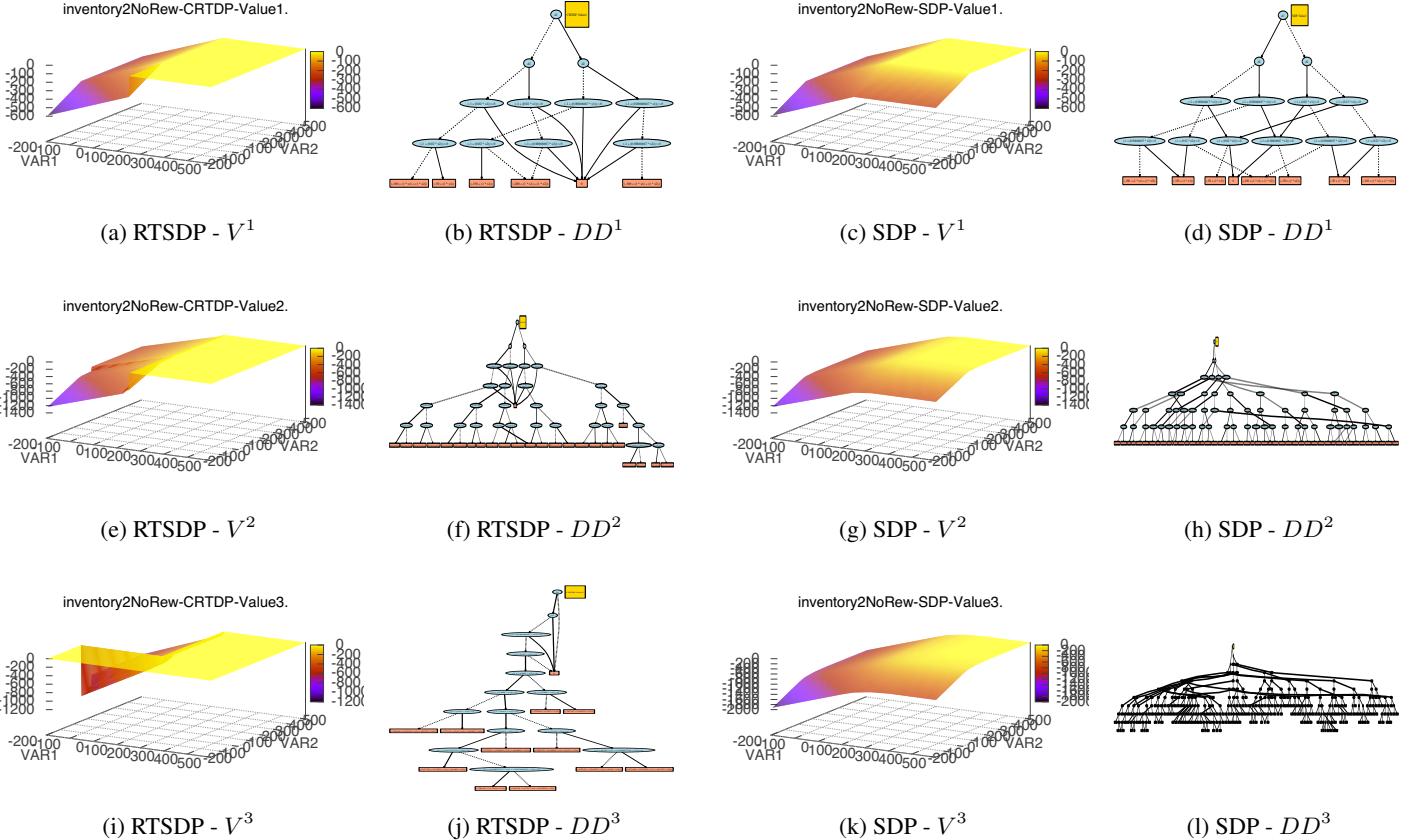


Figure 4: Partial value functions generated by RTSDP and SDP on INVENTORY CONTROL, with  $H = 3$  for different number of steps-to-go.

## References

- Bahar, R. I.; Frohm, E. A.; Gaona, C. M.; Hachtel, G. D.; Macii, E.; Pardo, A.; and Somenzi, F. 1993. Algebraic decision diagrams and their applications. In Lightner, M. R., and Jess, J. A. G., eds., *ICCAD*, 188–191. IEEE Computer Society.
- Barto, A.; Bradtke, S.; and Singh, S. 1995. Learning to act using real-time dynamic programming. *Artificial Intelligence* 72(1-2):81–138.
- Bellman, R. E. 1957. *Dynamic Programming*. USA: Princeton University Press.
- Bonet, B., and Geffner, H. 2003. Labeled RTDP: Improving the convergence of real-time dynamic programming. In *ICAPS-03*, 12–21. AAAI Press.
- Boutilier, C.; Hanks, S.; and Dean, T. 1999. Decision-theoretic planning: Structural assumptions and computational leverage. *Journal of Artificial Intelligence Research* 11:1–94.
- Dean, T., and Kanazawa, K. 1990. A model for reasoning about persistence and causation. *Computational Intelligence* 5(3):142–150.
- Feng, Z.; Dearden, R.; Meuleau, N.; and Washington, R. 2004. Dynamic programming for structured continuous

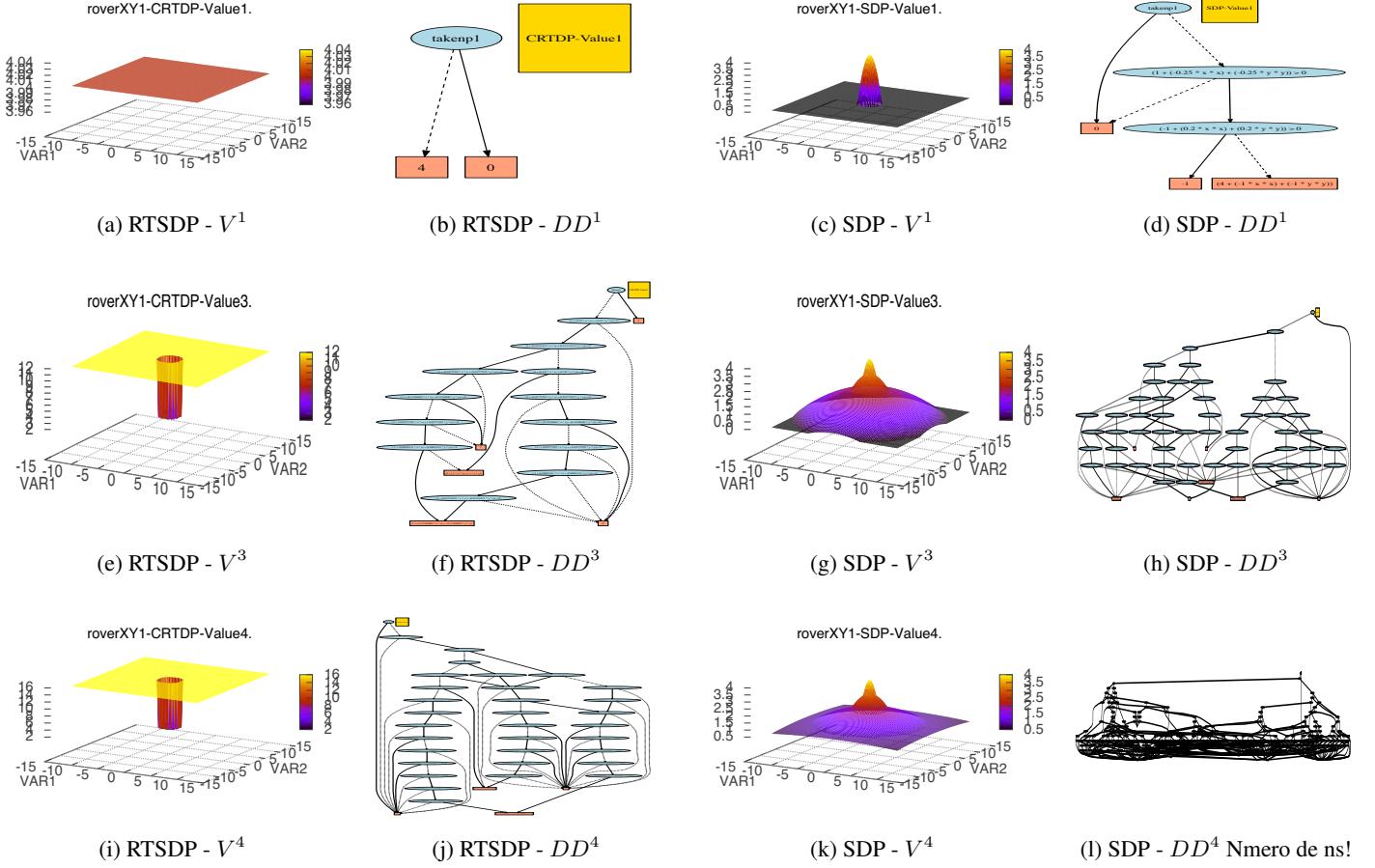


Figure 5: Partial value functions generated by RTSDP and RTSDP on MARS ROVERNonlinear, with  $H = 4$  for different actions-to-go.

markov decision problems. In *Proceedings of the UAI-04*, 154–161.

Kolobov, A.; Mausam; and Weld, D. 2012. LRTDP Versus UCT for online probabilistic planning. In *AAAI Conference on Artificial Intelligence*.

Li, L., and Littman, M. L. 2005. Lazy approximation for solving continuous finite-horizon MDPs. In *AAAI*, 1175–1180.

Meuleau, N.; Benazera, E.; Brafman, R. I.; Hansen, E. A.; and Mausam. 2009. A heuristic search approach to planning with continuous resources in stochastic domains. *Journal Artifical Intelligence Research (JAIR)* 34:27–59.

Puterman, M. L. 1994. *Markov decision processes*. New York: John Wiley and Sons.

Sanner, S.; Delgado, K. V.; and de Barros, L. N. 2011. Symbolic dynamic programming for discrete and continuous state MDPs. In *Proceedings of UAI-2011*.

Scarf, H. E. 2002. Inventory theory. *Operations Research* 50(1):186–191.

Zamani, Z.; Sanner, S.; and Fang, C. 2012. Symbolic dynamic programming for continuous state and action MDPs. In Hoffmann, J., and Selman, B., eds., *AAAI*. AAAI Press.