

# Toward Error-bounded Algorithms for Infinite-Horizon DEC-POMDPs

Jilles S. Dibangoye  
Ecole des Mines de Douai  
Douai, France  
jilles.dibangoye@mines-  
douai.fr

Abdel-Iliah Mouaddib  
University of Caen  
Caen, France  
mouaddib@info.unicaen.fr

Brahim Chaib-draa  
Laval University  
Québec, Qc, Canada  
chaib@aid.ift.ulaval.ca

## ABSTRACT

Over the past few years, attempts to scale up infinite-horizon DEC-POMDPs are mainly due to approximate algorithms, but without the theoretical guarantees of their exact counterparts. In contrast,  $\varepsilon$ -optimal methods have only theoretical significance but are not efficient in practice. In this paper, we introduce an algorithmic framework ( $\beta$ -PI) that exploits the scalability of the former while preserving the theoretical properties of the latter. We build upon  $\beta$ -PI a family of approximate algorithms that can find (provably) error-bounded solutions in reasonable time. Among this family, **H-PI** uses a branch-and-bound search method that computes a near-optimal solution over distributions over histories experienced by the agents. These distributions often lie near structured, low-dimensional subspace embedded in the high-dimensional sufficient statistic. By planning only on this subspace, **H-PI** successfully solves all tested benchmarks, outperforming standard algorithms, both in solution time and policy quality.

## Categories and Subject Descriptors

I.2.11 [Distributed Artificial Intelligence]: Multiagent systems

## General Terms

Algorithms, Experimentation, Performance

## Keywords

artificial intelligence, decentralized pomdps, point-based solvers.

## 1. INTRODUCTION

In recent years, there has been increasing interest in finding scalable algorithms for solving multiple agent systems where agents cooperate to optimize a joint reward function while having different local information. To formalize and solve such problems, [5, 14] suggest similar models that enable a set of  $n$  agents to cooperate in order to control a partially observable Markov decision process (POMDP), namely decentralized partially observable Markov decision process (DEC-POMDP).

Unfortunately, finding either optimal or even near-optimal solutions of such a problem has been shown to be particularly hard

**Cite as:** Toward Error-bounded Algorithms for Infinite-Horizon DEC-POMDPs, J.S. Dibangoye, A.-I. Mouaddib and B. Chaib-draa, *Proc. of 10th Int. Conf. on Autonomous Agents and Multiagent Systems (AA-MAS 2011)*, Tumer, Yolum, Sonenberg and Stone (eds.), May, 2–6, 2011, Taipei, Taiwan, pp. 947–954.

Copyright © 2011, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

[5, 15]. Significant efforts have been devoted to developing near-optimal algorithms for DEC-POMDPs. These algorithms consist in searching in the entire space of all policies [17, 4]. State-of-the-art optimal or near-optimal methods such as **DEC-PI** [4], **DP** [11], **MAA\*** [17], **PBDP** [18], or mathematical programs [3, 6] suggest first performing the exhaustive enumeration of all possible policies before they prune dominated ones. This is both an advantage and a liability. On the one hand, it preserves the ability to eventually find an  $\varepsilon$ -optimal policy, which is a key property. Yet it makes these methods impractical even for small toy problems. This is mainly because they quickly run out of memory.

To tackle the memory bottleneck, a number of memory-bounded algorithms have been suggested, and proven to be remarkably scalable, but without the theoretical guarantees of their exact counterparts [4, 16, 7, 10, 2]. Memory-bounded algorithms use a fixed amount of memory, *i.e.*, the size of the solution is fixed prior to the execution of the algorithm. Infinite-horizon memory-bounded techniques such as **NLP** and **BPI** rely on mathematical programs that search the best possible policy for a fixed size [1]. On the other hand, finite-horizon memory-bounded methods including **MBDP** [16], **MBDP-OC** [7], **PBIP** [10], **PBIP-IPG** [2] and **CBPB** [12] are mainly point-based algorithms. They compute approximate policies over a bounded number of beliefs<sup>1</sup> by selecting only a few policies for each point. While applying finite-horizon algorithms to infinite-horizon cases is non-trivial, they provide good insights on approximation methods. However, both infinite and finite approaches lack theoretical guarantees on the approximation. So it would seem we are constrained to either solving small toy problems near-optimally, or solving large problems but possibly doing so badly.

In this paper, we introduce an algorithmic framework ( $\beta$ -PI) that builds upon both approximate and near-optimal techniques. This provides the ability to preserve the theoretical properties of the former, while exploiting the scalability of the latter. To do so,  $\beta$ -PI incorporates the error  $\beta$  the decision-maker can sacrifice at each time step of the execution stage for computational tractability. A theoretical analysis of  $\beta$ -PI provides error-bounds on solutions produced by many approximate techniques. We then exploit the **H-PI** algorithm that aims at reducing the error produced in approximate algorithms while improving the empirical performances [9]. In order to reduce the error-bound, **H-PI** relies on the concept of distributions over histories, *i.e.*, the sufficient statistic for the selection of a decision rule<sup>2</sup> in general DEC-POMDPs [8] – page 199. By planning over distributions over histories experienced by the agents, **H-PI** considerably tightens the error-bound produced. These distributions often lie near structured, low-dimensional sub-

<sup>1</sup>A belief is a probability distribution over the underlying states of the system.

<sup>2</sup>A decision rule is a mapping  $d: H \rightarrow A$  from history set to action set.

space embedded in the high-dimensional sufficient statistic. By maintaining only a single policy-node for each individual history, it circumvents the memory bottleneck. This is achieved by means of a branch-and-bound search method that tracks the best policy for each distribution over histories experienced by the agents. This paper also provides empirical results demonstrating the successful performance of **H-PI** algorithm on all tested benchmarks: outperforming standard algorithms, both in solution time and policy quality.

## 2. BACKGROUND AND RELATED WORK

We review the DEC-POMDP model and the associated notation, and provide a short overview of the state-of-the-art algorithms.

### 2.1 The DEC-POMDP Model

A  $n$ -agent DEC-POMDP model can be represented using a tuple  $(S, \{A^i\}, \{\Omega^i\}, h_0, P, O, R, \gamma)$ , where:  $S$  is a finite set of states;  $A^i$  denotes a finite set of actions available for agent  $i$ , and  $A = \otimes_{i=1}^n A^i$  is the set of joint actions, where  $a = (a^1, \dots, a^n)$  denotes a joint action;  $P(s'|s, a)$  is a Markovian transition function, that denotes the probability of transiting from state  $s$  to state  $s'$  when taking action  $a$ ;  $\Omega^i$  defines a finite set of observations available for agent  $i$ , and  $\Omega = \otimes_{i=1}^n \Omega^i$  is the set of joint-observations, where  $\omega = (\omega^1, \dots, \omega^n)$  is a joint observation;  $O(\omega|a, s')$  is an observation function, that denotes the probability of observing joint observation  $\omega$  given that joint action  $a$  was taken and led to state  $s'$ ;  $R(s, a)$  is a reward function, that denotes the reward signal received when executing action  $a$  in state  $s$ . The DEC-POMDP model is parameterized by:  $h_0$ , the initial joint history, *i.e.*, the team joint action/joint-observation trace. When the agents operate over an unbounded number of time-steps, the DEC-POMDP has a discount factor,  $\gamma \in [0, 1)$ . This model is referred as infinite-horizon DEC-POMDPs with discounted rewards.

#### 2.1.1 Sufficient Statistic

A key assumption of DEC-POMDPs is that during the online execution stage the true state of the world could not be sensed exactly and reliably: agents are imperfectly informed about the state of the world to differing degrees.

Given that the state is not directly observable, the agents can instead maintain a complete trace of all joint-observations and all joint-actions they ever executed during the offline planning stage, and use this to select their joint-actions. These joint-action/joint-observation traces are referred to as joint-history experienced by the agents. We formally define  $h_\tau^i := (a_0^i, \omega_1^i, a_1^i, \omega_2^i, \dots, a_{\tau-1}^i, \omega_\tau^i)$ ,  $h_\tau := (h_\tau^1, \dots, h_\tau^n)$ , and  $H_\tau := \{h_\tau\}$  to be the history of agent  $i$ , the joint history of the team, and the set of histories at time step  $\tau$ , respectively.

We define the sufficient statistic at step  $\tau$ ,  $\mu_\tau \in \Delta H_\tau$  to be a probability distribution of the team over joint-histories  $H_\tau$  [8]. Furthermore,  $\mu_\tau$  at time step  $\tau$  is calculated recursively, using only the distribution over histories one time step earlier,  $\mu_{\tau-1}$ , along with the most recent joint decision rule  $d_{\tau-1}: H_{\tau-1} \rightarrow A$ :

$$\mu_\tau(h_\tau) = \mu_{\tau-1}(h_{\tau-1}) \cdot p(h_{\tau-1}, d_{\tau-1}(h_{\tau-1}), \omega_\tau | h_{\tau-1})$$

for all  $h_{\tau-1} \in H_{\tau-1}$  and  $\omega_\tau \in \Omega$ , where joint history  $h_\tau$  is given by joint history one step earlier, along with its corresponding joint-action and a given joint-observation, *i.e.*,  $(h_{\tau-1}, d_{\tau-1}(h_{\tau-1}), \omega_\tau)$ . Notice that  $p(h, a, \omega | h') = \sum_{s, s'} O(\omega | a, s') P(s' | a, s) \mu_{\tau-1}(h')$  and  $\mu_0(h) = 1$  if and only if  $h = h_0$ . Finally, the distribution over individual history  $h^i$  is defined by  $\mu^i(h^i) = \sum_h p(h^i | h) \mu(h)$ , where  $p(h^i | h) = 1$  if and only if there exists  $h^j$  such that  $h^i h^j =$

$h$ , otherwise  $p(h^i | h) = 0$ . If not all joint histories are reachable,  $\mu_\tau$  yields a positive probability only for reachable histories denoted  $\bar{H}_\tau$ . Unfortunately,  $\bar{H}_\tau$  can get very large as time goes on. More precisely, set of histories increases exponentially with increasing time step  $\bar{H}_\tau = \mathcal{O}(|A^i| |\Omega^i|^{n\tau})$ . For this reason, we want to plan only over a small set of distributions over histories experienced by the agents. These distributions often lie near structured, low-dimensional subspace embedded in the high-dimensional sufficient statistic.

#### 2.1.2 Optimization Criterion

The goal of DEC-POMDP planning is to find a sequence of joint-actions  $\{a_0, \dots, a_\tau\}$  maximizing the expected sum of rewards  $\mathbb{E}[\sum_{\tau=0}^{\infty} \gamma^\tau R(s_\tau, a_\tau)]$ . Given the initial belief, the goal is to find a joint-policy  $\delta$  that yields the highest expected reward. Throughout the paper, a policy  $\delta$  of the team is represented as a deterministic joint-policy graph. That is, a vector  $(\delta^1, \delta^2, \dots, \delta^n)$  of individual policy graphs as illustrated in Figure 1. We note  $X :=$

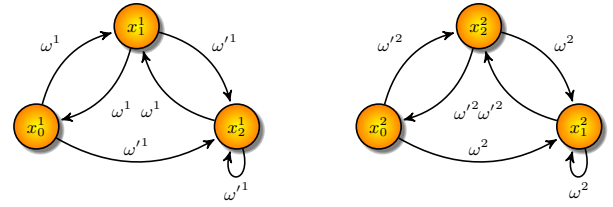


Figure 1: Deterministic policy graphs for two agents.

$X^1 \times \dots \times X^n$  the set of joint-policy nodes  $x = (x^1, \dots, x^n)$ . A policy for a single agent  $i$  is therefore represented as a deterministic policy graph  $\delta^i$ , where  $X^i = \{x^i\}$  denotes a set of policy nodes. Solving such a problem usually relies on successive approximations of the joint-policy graph. After  $\tau$  consecutive iterations, the solution consists of a set of hyperplanes  $\Lambda_\tau = \{v^x\}$ , together with the corresponding joint-policy graph  $\delta_\tau$ . The value function at iteration  $\tau$  can be formulated as:

$$v_\tau(\mu_\tau) = \sup_{\substack{x_{1\tau}^1, x_{2\tau}^1, \dots, x_{|H_\tau|}^1 \\ x_{1\tau}^n, x_{2\tau}^n, \dots, x_{|H_\tau|}^n}} \sum_{k=1}^{|H_\tau|} \mu_\tau(h_{k\tau}) \cdot v^{x_{k\tau}}(h_{k\tau}) \quad (1)$$

where  $x_{k\tau} = (x_{k\tau}^1, x_{k\tau}^2, \dots, x_{k\tau}^n)$  denotes the joint-policy node associated to joint-history  $h_{k\tau}$ . The resulting set of joint-policy nodes  $X_\tau := \{x_{k\tau}\}_{k=1, \dots, m}$  represents the next step joint-policy graph  $\delta_\tau$ . When joint-policy node  $x$  is associated to joint-action  $a$ , for all  $s \in S$ , hyperplane  $v^x$  follows:

$$\begin{aligned} v^x(h_\tau) &= \sum_s p(s | h_\tau) (R(s, a) + \gamma \sum_\omega v^{a, \omega, x'}(s)) \\ p(s | h_\tau) &= \frac{\sum_{s'} O(\omega | a, s') P(s' | a, s) \mu_{\tau-1}(h_{\tau-1})}{p(h_\tau, a, \omega | h_{\tau-1})} \end{aligned}$$

where  $p(s | h_\tau)$  denotes the probability of being in state  $s$  given history  $h_\tau$ , this probability distribution is also referred to as belief state. The estimate value  $v^{a, \omega, x'}$  denotes the value of taking joint-action  $a$  followed by a joint-observation  $\omega$  conditional transition to joint-policy node  $x$ , and is given by:

$$v^{a, \omega, x'}(s) = \sum_{s'} P(s' | s, a) O(\omega | a, s') v^{x'}(s')$$

It is worth noticing that Equation (1) denotes the supremum over all next step joint-policy graphs that selects both: on the one hand,

the best<sup>3</sup> hyperplane  $v^{x_{k\tau}}$  for each joint-history  $h_{k\tau}$ ; and a single policy-node  $x_{k\tau}^i$  for each individual history  $h_{k\tau}^i$ , on the other hand, thus preserving the ability to control the system in a distributed manner.

Throughout the paper, we will use superscripts either to name agent, e.g.,  $x^i$ ,  $\Omega^i$  or to distinguish estimate values between joint-policy graphs and joint-policy nodes, e.g.,  $v^\delta$ ,  $v^x$ , or  $v^{a,\omega,x}$ . In addition, we use subscripts to indicate time step or iteration, e.g.,  $v_\tau$ ,  $\Lambda_\tau$ ,  $\delta_\tau$ . Finally, except specific indications,  $\|\cdot\|$  denotes the Chebyshev norm.

## 2.2 Related Work

In this section, we discuss near-optimal as well as approximate approaches to solving infinite horizon DEC-POMDPs with discounted rewards.

**DEC-PI** was the first attempt to compute a near-optimal policy for infinite-horizon DEC-POMDPs with discounted rewards. It builds over a series of exhaustive backups a vector of stochastic policy graphs, one for each agent [4]. However, the number of policy nodes generated by the exhaustive backups is double exponential in the number of iterations. In order to reduce the number of policy nodes generated, **DEC-PI** does pruning by using iterated elimination of dominated policies after each backup, without losing the ability to eventually converge to a near-optimal policy. Performing this pruning, however, can be expensive, in addition the number of policy nodes still grows exponentially.

To alleviate these problems, we can use a heuristic search technique (referred to as **I-MAA\***), which uses a best-first search in the space of deterministic joint-policy graphs with a fixed size. **I-MAA\*** prunes dominated joint-policy graphs at earlier construction stages. This is done by calculating a heuristic for the joint-policy graph given known parameters and filling in the remaining parameters one at a time in a best-first fashion. Both **DEC-PI** and **I-MAA\*** provide good guarantees on the solution quality, but they do not scale beyond small toy problems. This is mainly due to the explosion in memory. As such, researchers have turned their attention to a family of approximate memory-bounded algorithms.

Thus, a version of **DEC-PI** namely **DEC-BPI** was introduced to keep the policy size bounded over the iterations of the **DEC-PI** algorithm [4]. **DEC-BPI** iterates through the policy nodes of each agent's stochastic policy graph and attempts to find an improvement. A linear program searches for a probability distribution over actions and transitions into the agent's current policy graph that increases the value of the policy graph for any initial belief state and any initial policy node of the other agents' policy graph. If an improvement is discovered, the policy node is updated based on the probability distributions found. Each policy node of each agent is examined in turn and the algorithm terminates when no policy graph can be further improved. While this algorithm allows memory to remain fixed, it provides only a locally optimal solution. Unfortunately, this locally optimal solution can be arbitrarily far from the actual optimal solution.

In an attempt to address some of these problems, a set of optimal policy graphs given a fixed size with nonlinear program was defined in the **NLP** algorithm. Because it is often intractable to solve this **NLP** optimally, a locally optimal solver is used. Unlike **DEC-BPI**, this approach allows initial belief state to be used so smaller policy graphs may be generated and the improvement takes place in one step. While concise policy graphs with high value can be produced, large policy graphs, which may be required for some prob-

lems, cannot be produced by the current locally optimal solvers. Even more importantly, these locally optimal solvers do not provide any error-bound on their solutions. So it would seem we are constrained to either solving small toy problems near-optimally, or solving large problems but possibly doing so badly.

## 3. A NEW NEAR-OPTIMAL ALGORITHM

In this section, we present a new near-optimal algorithm ( **$\beta$ -PI**) for solving infinite-horizon DEC-POMDPs with discounted rewards.  **$\beta$ -PI** has only theoretical significance and is not efficient in practice. However, its framework serves as a foundation to derive methods that are both error-bounded and very efficient in practice, as discussed in the next section (Section 4).

**$\beta$ -PI** algorithm (Figure 2) consists of a two-step method: the policy-evaluation (step 2); and the policy-improvement (step 3). At each iteration  $\tau$ , the policy-evaluation estimates the set of hyperplanes  $\Lambda_\tau$  of the current joint-policy graph  $\delta_\tau$ , while the policy improvement updates set  $\Lambda_\tau$  into set  $\Lambda_{\tau+1}$ . Thereafter, it transforms the current joint-policy graph  $\delta_\tau$  into an improved one  $\delta_{\tau+1}$  through comparison of  $\Lambda_\tau$  and  $\Lambda_{\tau+1}$ . Finally, duplicated, dominated, and unreachable old policy nodes are pruned.

### Policy Iteration Algorithm $\beta$ -PI

1. Set parameters  $(\beta, \varepsilon)$  and joint-policy graph  $\delta_0$ .
2. (*Policy Evaluation*) Obtain  $\Lambda_\tau$  by evaluating  $\delta_\tau$ .
3. (*Policy Improvement*) Transform  $\delta_\tau$  to  $\delta_{\tau+1}$  through

$$\Lambda_{\tau+1} = (\beta\text{-}\mathbb{H}) \cdot \Lambda_\tau$$

4. If  $\|v_{\tau+1} - v_\tau\| < \varepsilon(1 - \gamma)/2\gamma$ , stop and return  $\delta_{\tau+1}$ . Otherwise set  $\tau = \tau + 1$  and return to step 2.

Figure 2:  $\beta$ -PI Algorithm.

The above procedure is similar to classical  $\varepsilon$ -optimal policy-iteration algorithms [4], when parameter  $\beta = 0$ . This parameter denotes the decision-maker's preference on the quality of the returned solution. Indeed,  **$\beta$ -PI** is designed to return a solution with error bounded by  $\beta$  when compared to the  $\varepsilon$ -optimal solution, so as to satisfy the decision-maker's preferences. To this end,  **$\beta$ -PI** replaces the exhaustive backup operator performed in classical policy-iteration algorithms by a new backup operator ( $\beta\text{-}\mathbb{H}$ ) that builds up the improved value function with error bounded by  $\beta$ .

### 3.1 Backup Operator

A backup operator aims at computing an improved set  $\Lambda_{\tau+1} = \{v^{x'}\}$  given set  $\Lambda_\tau$ . Each joint-policy node  $x'$  corresponds to a joint-action choice  $a$  (resp.  $a^i$ ) followed by a joint-observation choice  $\omega$  (resp.  $\omega^i$ ) conditional transition to joint-policy node  $x$  (resp.  $x^i$ ). As a result, one can represent a joint-policy node  $x'$  as a set of action-observation-node trios  $\{(a^i - \omega^i - x^i)\}_{i=1, \dots, n}$ .

$$\Lambda_{\tau+1} = (\beta\text{-}\mathbb{H}) \cdot \Lambda_\tau$$

1.  $\forall (a^i, \omega^i)$  compute set  $X_{\tau+1}^{a^i, \omega^i} = \mathbf{IEDT}(a^i, \omega^i, X_\tau^i)$ .
2.  $\forall a$ , compute set  $X_{\tau+1}^a = \bigotimes_{i=1}^n (\bigotimes_{\omega^i \in \Omega^i} X_{\tau+1}^{a^i, \omega^i})$ .
3. Compute set  $\Lambda_{\tau+1} = \{v^{x'} \mid x' \in \cup_{a \in A} X_{\tau+1}^a\}$ .

Figure 3: Backup Operator  $\beta\text{-}\mathbb{H}$ .

Following this idea, backup operator  $\beta\text{-}\mathbb{H}$ , described in Figure 3, prunes those dominated trios, but without actually constructing joint-policy nodes  $x'$  exhaustively (step 1). Next, it creates set

<sup>3</sup>The best hyperplane is not necessarily the maximal hyperplane for a given joint-history, this is why the sup operator is outside the  $\sum$ , and states the major difference with POMDPs both in terms of complexity and value function expression.

$X_{\tau+1}^a$  ( $\forall a \in A$ ), a cross-product over agents and individual observations, which includes one trio  $(a^i - \omega^i - x^i)$  from each  $X_{\tau+1}^{a^i, \omega^i}$  (step 2). Finally, it takes the union of  $X^a$  sets and creates the improved value function represented by set  $\Lambda_{\tau+1}$  (step 3).

The intuition behind  $\beta$ - $\mathbb{H}$  is that rather than adding all possible trios  $A^i \times \Omega^i \times X_{\tau}^i$  as suggested in classical algorithms [4], we only add trios that would be part of non  $\beta$ -dominated joint-policy node  $x'$ . That is, trio  $(a^i - \omega^i - x^i)$  is pruned if its value goes up over  $\beta$  by changing policy node  $x^i$  by another one for some distribution  $\zeta(\cdot)$  of states  $s$ ; policy node  $x^j$  of the other agents; action  $a^j$ ; and observation  $\omega^j$ . For the sake of simplicity we use the abbreviation  $\rho^j = (s, a^j, \omega^j, x^j)$ . Our  $\beta$ -dominance criterion is therefore formulated as follows: maximize  $\xi$ , s.t.:  $\forall y^i \in X_{\tau+1}^{a^i, \omega^i} \setminus \{x^i\}$ , in

$$\sum_{\rho^j} \zeta(\rho^j) v^{a, \omega, x}(s) + \xi + \beta \leq \sum_{\rho^j} \zeta(\rho^j) v^{a, \omega, y}(s) \quad (2)$$

where  $\sum_{\rho^j} \zeta(\rho^j) = 1$  and  $\zeta(\rho^j) \geq 0$  and  $a = a^i a^j$ ,  $\omega = \omega^i \omega^j$ ,  $x = x^i x^j$  and  $y = y^i x^j$ . We provide in Figure 4 an example on how the concept of  $\beta$ -dominance can be used in practice.

---

#### Algorithm 1 Iterative Elimination of $\beta$ -Dominated Trios

---

```

1: procedure PRUNE( $a^i, \omega^i, X_{\tau}^i$ )
2:   Initialize:  $X_{\tau+1}^{a^i, \omega^i} \leftarrow X_{\tau}^i$ ;
3:   repeat
4:     for  $i = 1, 2, \dots, n$  do
5:       for  $x^i \in X_{\tau+1}^{a^i, \omega^i}$  do
6:         Maximize  $\xi$ , s.t.:  $\forall y^i \in X_{\tau+1}^{a^i, \omega^i} \setminus \{x^i\}$ , in Eq. (2).
7:         if ( $\xi \leq 0$ ) then Remove  $y^i$  from  $X_{\tau+1}^{a^i, \omega^i}$ ;
8:   until no changes occur

```

---

We are now ready to present the iterative elimination of  $\beta$ -dominated trios algorithm – called **IEDT** Algorithm 1. This algorithm loops over each trio  $(a^i, \omega^i, x^i)$  and tests whether there exists a probability distribution  $\zeta$  such that  $(a^i, \omega^i, x^i)$   $\beta$ -dominates any other trio, e.g.,  $(a^i, \omega^i, y^i)$ . Those trios are kept in sets  $X_{\tau+1}^{a^i, \omega^i}$ . Then it repeats this procedure for all agents, until no more changes occur. **DEC-PI** and **DP** introduce a similar pruning mechanism namely iterative elimination of dominated policies. However, ours remains fundamentally different. The key difference lies in when this pruning takes place and what we actually prune. **DEC-PI** and **DP** iterative elimination procedures take place after each exhaustive backup, i.e., they first generate all possible policies before they prune dominated ones. **IEDT**, however, takes place before the exhaustive backup, providing the ability to prune all  $\beta$ -dominated trios before we actually build the next step policies. As a result, all policies generated by **IEDT** are non  $\beta$ -dominated policies, as discussed below.

## 3.2 Theoretical Analysis

This section states and proves the theoretical properties of our  $\beta$ -**PI** algorithm, including: error-bound and convergence. We first show that our  $\beta$ -**PI** algorithm does not prune trios that would be part of a non  $\beta$ -dominated joint-policy node.

**LEMMA 1.** *Any joint-policy node  $x'$  that includes  $\beta$ -dominated trio  $(a^i, \omega^i, x^i)$  is  $\beta$ -dominated by some probability distribution over some joint-policy nodes.*

**PROOF.** We will prove this result for 2 agents (although it holds for more), and from the agent 1's perspective. Let trio  $(a^i, \omega^i, x^i)$  be a  $\beta$ -dominated trio. We now show that any joint-policy node  $x'$

that includes  $(a^i, \omega^i, x^i)$  is  $\beta$ -dominated by some probability distribution over a set of joint-policy nodes  $\{y'\}$  that are identical to  $x'$  except instead of including  $(a^i, \omega^i, x^i)$ , some probability distribution over trios  $\{(a^i, \omega^i, y^i)\}$  is chosen. That is,  $\forall s, \forall \omega^2, \forall x^2$ ,

$$v^{a, \omega^1 \omega^2, x^1 x^2}(s) \leq \sum_{y^1 \neq x^1} p(y^1) \cdot v^{a, \omega^1 \omega^2, y^1 x^2}(s) + \beta \quad (3)$$

$$\sum_{\omega \in \Omega} v^{a, \omega, x^1 x^2}(s) \leq \sum_{\omega \in \Omega} \sum_{y^1 \neq x^1} p(y^1) \cdot v^{a, \omega, y^1 x^2}(s) + \beta \quad (4)$$

$$v^{x'}(s) - \beta \leq \sum_{\omega^1, y^1} p(\omega^1, y^1) \left( R(s, a) + \sum_{\omega \in \Omega} v^{a, \omega, y^1 x^2}(s) \right) \quad (5)$$

$$v^{x'}(s) - \beta \leq \sum_{y'} p(y') \cdot v^{y'}(s) \quad (6)$$

where  $\omega = (\omega^1, \omega^2)$  and  $a = (a^1, a^2)$ . The inequality (3) results from inequality (2) where trio  $(a^i, \omega^1, x^1)$  is supposed stochastically  $\beta$ -dominated; in the inequality (4), we consider the sum over all joint observations  $\omega \in \Omega$ ; in inequality 5, we add the immediate reward, and build joint-policy nodes  $x'$  and  $\{y'\}$ . The cross-product between the probability distribution  $p(y^1)$  and the uniform probability distribution over  $\Omega^1$  enables us to build a probability distribution  $p(\omega^1, y^1)$  and we build in a similar way a probability distribution  $p(y')$  so that  $x'$  is stochastically dominated by  $\{y'\}$ . This ends the proof.  $\square$

We now show that our  $\beta$ -**PI** algorithm returns a near-optimal solution. To do so, we apply the *Banach Fixed-Point Theorem* [13] to prove that  $\beta$ - $\mathbb{H}$  is a contraction mapping on the space  $\mathcal{V}$  of bounded functions on  $S$  with supremum norm. The proof that  $\beta$ -**PI** is near-optimal follows from properties of norms and contraction mappings.

**THEOREM 1.** *Our  $\beta$ -**PI** algorithm returns a near-optimal solution for any initial history, with error bounded by  $(\varepsilon/2 + \beta/(1 - \gamma))$ .*

**PROOF.** First, we prove that  $\beta$ - $\mathbb{H}$  is a contraction mapping on the space of value functions  $\mathcal{V}$  for any positive scalar  $\beta$ . Since  $S$  is discrete,  $\beta$ - $\mathbb{H}$  maps  $\mathcal{V}$  into  $\mathcal{V}$ .

Let  $v$  and  $u$  be estimate values of value functions  $V$  and  $U$  in  $\mathcal{V}$  respectively. Fix  $h_0 \in H_0$ , assume that  $(\beta$ - $\mathbb{H})v(\mu_0) \geq (\beta$ - $\mathbb{H})u(\mu_0)$ , and let  $x_{h_0} = \arg \max_{x: v^x \in (\beta$ - $\mathbb{H})\mathcal{V}} \sum_s v^x(h_0)$ . Denote  $a_{h_0}$  the joint-action associated to joint-policy node  $x_{h_0}$ , and  $\xi = (\beta$ - $\mathbb{H})v(\mu_0) - (\beta$ - $\mathbb{H})u(\mu_0)$ . Then,  $0 \leq \xi \leq \sum_s (R(s, a_{h_0}) + \gamma \sum_{s', \omega} P(s'|s, a_{h_0}) O(\omega|a_{h_0}, s') v^x(s')) p(s|h_0) - \sum_s (R(s, a_{h_0}) + \gamma \sum_{s', \omega} P(s'|s, a_{h_0}) O(\omega|a_{h_0}, s') u^y(s')) p(s|h_0)$ . Finally,

$$\begin{aligned} \xi &\leq \gamma \sum_{s, s', \omega} P(s'|s, a_{h_0}) O(\omega|a_{h_0}, s') p(s|h_0) [v^x(s') - u^y(s')] \\ &\leq \gamma \sum_{s, s', \omega} P(s'|s, a_{h_0}) O(\omega|a_{h_0}, s') p(s|h_0) \|v - u\| \\ &= \gamma \|v - u\| \end{aligned}$$

Repeating this argument in the case  $(\beta$ - $\mathbb{H})u(\mu_0) \geq (\beta$ - $\mathbb{H})v(\mu_0)$  implies that  $|(\beta$ - $\mathbb{H})v(\mu_0) - (\beta$ - $\mathbb{H})u(\mu_0)| \leq \gamma \|v - u\|$  for all initial distributions  $\mu_0 \in \Delta \bar{H}_0$ . Taking the supremum over  $\mu_0$  in the above expression gives the result.

We now able to show that  $\beta$ -**PI** returns a joint-policy graph  $\delta$  that is near-optimal. Suppose that  $\|v_{\tau+1} - v_{\tau}\| \leq \varepsilon(1 - \gamma)/2\gamma$  holds for some iteration. Then, the overall error in  $\beta$ -**PI** is bounded by  $\|v^{\delta} - v_{\tau+1}\| + \|v_{\tau+1} - v^*\|$ . Since  $\delta$  is a fixed point of  $(\beta$ - $\mathbb{H})$ ,

the first expression is bounded as follows:  $\|v^\delta - v_{\tau+1}\|$

$$\begin{aligned} &= \|(\beta\text{-}\mathbb{H}) \cdot v^\delta - v_{\tau+1}\| \\ &\leq \|(\beta\text{-}\mathbb{H}) \cdot v^\delta - (\beta\text{-}\mathbb{H}) \cdot v_{\tau+1}\| + \|(\beta\text{-}\mathbb{H}) \cdot v_{\tau+1} - v_{\tau+1}\| \\ &\leq \gamma \|v^\delta - v_{\tau+1}\| + \|(\beta\text{-}\mathbb{H}) \cdot v_{\tau+1} - (\beta\text{-}\mathbb{H}) \cdot v_\tau\| \\ &\leq \gamma \|v^\delta - v_{\tau+1}\| + \gamma \|v_{\tau+1} - v_\tau\| \end{aligned}$$

where inequalities follow because  $(\beta\text{-}\mathbb{H})$  is a contraction mapping on  $\mathcal{V}$ . Rearranging terms yields:

$$\|v^\delta - v_{\tau+1}\| \leq \frac{\gamma}{1-\gamma} \|v_{\tau+1} - v_\tau\|.$$

Then, the second expression follows because  $(0\text{-}\mathbb{H})$  and  $(\beta\text{-}\mathbb{H})$  are contraction mappings on  $\mathcal{V}$ :  $\|v_{\tau+1} - v^*\|$

$$\begin{aligned} &\leq \|(\beta\text{-}\mathbb{H})v_{\tau+1} - (0\text{-}\mathbb{H})v_{\tau+1}\| + \|(0\text{-}\mathbb{H})v_{\tau+1} - v^*\| \\ &\leq \beta + \|(0\text{-}\mathbb{H})v_{\tau+1} - (0\text{-}\mathbb{H})v^*\| \\ &\leq \beta + \gamma \|v_{\tau+1} - v^*\| \end{aligned}$$

Rearranging terms yields:  $\|v_{\tau+1} - v^*\| \leq \beta/(1-\gamma)$ . Thus when  $\|v_{\tau+1} - v_\tau\| \leq \varepsilon(1-\gamma)/2\gamma$  holds, the first expression is bounded by  $\varepsilon$  and the second expression is bounded by  $\beta/(1-\gamma)$ , so that the error produced by  $\beta\text{-}\mathbf{PI}$  is bounded by  $\varepsilon/2 + \beta/(1-\gamma)$ .  $\square$

To better understand the significance of the error-bound in  $\beta\text{-}\mathbf{PI}$ , let's consider its terms. The first term  $\varepsilon/2$  denotes the error produced by the stopping criterion in  $\beta\text{-}\mathbf{PI}$  algorithm in Figure 2, step 4. This criterion stops the algorithm before a fixed point of  $\beta\text{-}\mathbb{H}$  has been found. It is required when optimal joint-policies do not exist, so we seek  $\varepsilon$ -optimal joint-policies for  $\beta = 0$ . This criterion also guarantees that  $\beta\text{-}\mathbf{PI}$  terminates after a finite number of iterations. The second term  $\beta/(1-\gamma)$  defines the error the decision-maker's preference produced by pruning all  $\beta$ -dominated hyperplanes. Decreasing the decision-maker's parameter  $\beta$  reduces the error-bound and increases the solution size, but it is usually worthwhile to avoid an explosion of the solution size.

Unfortunately, it is worth noticing that the number of preserved hyperplanes in  $\beta\text{-}\mathbf{PI}$  would be very large in many practical cases. This is mainly because it is likely that there exists a probability distribution for which many hyperplanes  $\beta$ -dominate any other. Indeed, there are infinitely many possible probability distributions. In the next section, we provide two enhancements that overtake the limitations of  $\beta\text{-}\mathbf{PI}$  to scale up while preserving its ability to bound the error produced.

## 4. ERROR-BOUNDED ALGORITHMS

First, we want to plan only over distributions of histories experienced by the agents  $B := \{\mu\}$  as a means of reducing the infinitely many possible probability distributions considered into  $\beta\text{-}\mathbf{PI}$ . Then, we want to arbitrarily increase parameter  $\beta$  such that at some point only one policy node will be preserved for each individual history as a means of reducing the solution size.

Thereafter, those distributions can be used to build non  $\beta$ -dominated hyperplanes at each iteration of  $\beta\text{-}\mathbf{PI}$  algorithm. In particular, the linear program (inequality 2) can be replaced by a series of comparisons over joint-histories  $h$  where distributions  $\mu \in B$  is positive, without affecting the ability to find a near-optimal solution with respect to  $B$ . More formally, if inequality

$$v^{a,\omega^i,\omega^j,x^i,x^j}(h) \leq v^{a,\omega^i,\omega^j,y^i,y^j}(h) + \beta$$

holds for any  $(\omega^j, x^j)$ , then trio  $(a^i, \omega^i, x^i)$  is  $\beta$ -dominated by trio  $(a^i, \omega^i, y^i)$  for  $h$ .

To better understand the pruning procedure using joint-histories, consider the example illustrated in Figure 4. This figure shows the

steps of pruning  $\beta$ -dominated trios for a problem with 2 agents; 2 individual observations  $\{\omega^i, \omega'^i\}$  and policy-nodes  $\{x^i, x'^i\}$ ; a joint-history  $h = h^1 h^2$ ; joint-action  $a = a^1 a^2$ , and  $\beta = 0.9$ . The set of trios are represented in a form of a Bayesian game Figure (4.A). Figure (4.B) illustrates the first pruning of  $\beta$ -dominated trios for each agent (red lines). As an example, trio  $(a^1, \omega^1, x^1)$  is pruned since it is  $\beta$ -dominated by trio  $(a^1, \omega^1, x'^1)$ . The pruning process continues until no more pruning occurs. Figure (4.D) shows that for joint-history  $h$ , joint-action  $a$  and  $\beta = 0.9$ , there is only one possible non  $\beta$ -dominated hyperplane (each agent has only one possible policy node for each observation – a single policy node for each individual history). This remark is crucial to bound the error produced by algorithms that keep only one hyperplane for each joint-history  $h$  or its corresponding belief state  $p(\cdot|h)$ , such as point-based solvers **MBDP** and **PBIP**. That is, there exists a possible large scalar  $\beta_B$  such that there is only one non  $\beta_B$ -dominated hyperplane for each joint-history  $h$ .

The next section presents **H-PI**, which provides an efficient and scalable derivation of  $\beta\text{-}\mathbf{PI}$ , while preserving the ability to bound the error produced.

### 4.1 Heuristic PI Algorithm

The heuristic policy-iteration (**H-PI**) algorithm replaces backup operator  $\beta\text{-}\mathbb{H}$  in  $\beta\text{-}\mathbf{PI}$  by a more scalable backup operator denoted  $\beta\text{-}\mathbb{H}_B$ . This operator performs the backup only over a set of distributions over histories  $\mu \in B$ , by means of a branch-and-bound search in the space of non  $\beta$ -dominated policy nodes.

As **H-PI** proceeds in the same way for every iteration, and each distribution  $\mu \in B$ , we therefore restrict our attention to the following problem that occurs at each iteration and for each  $\mu$ : the problem of assigning trio  $(a^i, \omega^i, x^i)$  for each individual history  $h^i$  where  $\mu^i(h^i) > 0$ , and this for every individual observation  $\omega^i \in \Omega^i$  and all agents  $i = 1, 2, \dots, n$ , and such that the resulting joint-policy nodes  $\{x_h\}_{h: \mu(h)>0}$  are the best possible. That is, the corresponding set of hyperplanes  $\{v^{x_h}\}_{h: \mu(h)>0}$   $\beta$ -dominates any other at  $\mu$ . More formally,  $\beta\text{-}\mathbb{H}_B$  computes  $\{v^{x_h}\}_{h: \mu(h)>0}$ , such that for any other set of hyperplanes  $\{v^{x'_h}\}_h$  the following holds:  $\sum_h \mu(h)v^{x_h}(h) + \beta \geq \sum_h \mu(h)v^{x'_h}(h)$ .

The idea behind **H-PI** is to build a search tree in which nodes  $\theta$  are sets of partially specified mappings  $\{(d^i, \sigma^i)\}_i$ , where  $d^i: H^i \rightarrow A^i$  is a mapping from individual history set  $H^i = \{h^i | \mu^i(h^i) > 0\}$  to individual action set  $A^i$ ; and  $\sigma^i: H^i \times \Omega^i \rightarrow \cup_{a^i} X^{a^i, \omega^i}$  is a mapping from pairs of individual history and observation to non  $\beta$ -dominated policy nodes.

Notice that by assigning a value to a variable we often constrain the possible assignments of the other variables. To better understand this, let's consider the assignment of value  $a^i$  to variable  $d^i(h^i)$ , as a result variables  $\sigma^i(h^i, \omega^i)$  are constrained to choose their values in  $X^{a^i, \omega^i}$ . Thereafter, it is likely that trios that were non  $\beta$ -dominated before the assignment become  $\beta$ -dominated after. For this reason, **H-PI** interleaves each search node expansion step with an iterative elimination of  $\beta$ -dominated trios for each expanded nodes in the search tree. This provides **H-PI**'s first pruning mechanism. The second one prunes nodes based on upper and lower bounds.

We define the upper-bound based on the decomposition of the exact estimate into two estimates. The first estimate,  $G(\theta, \mu)$ , is the exact estimate coming from variables where  $\theta$  is constrained<sup>4</sup>. The second estimate,  $H(\theta, \mu)$ , is the upper-bound value coming from variables where  $\theta$  is not constrained. That is,  $\hat{v}(\theta, \mu) = G(\theta, \mu) +$

<sup>4</sup>A partially specified mapping  $\phi: \mathcal{X} \rightarrow \mathcal{Y}$  is said to be constrained at  $x \in \mathcal{X}$  if  $\phi(x)$  has been assigned a value  $y \in \mathcal{Y}$ , otherwise it is said to be non constrained.

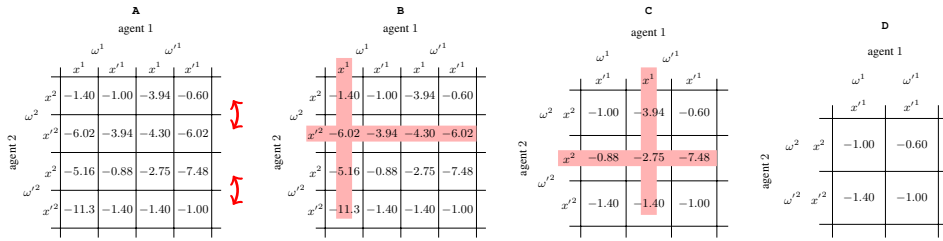


Figure 4: Illustration of the  $\beta$ -dominance criterion, where  $\beta = .9$ .

### Algorithm 2 Heuristic Backup Operator

---

```

1: procedure  $\beta$ - $\mathbb{H}_B(\mu)$ 
2:   Initialize: Incumbent :=  $v(\mu)$ ; Live :=  $\{\theta_0\}$ 
3:   repeat
4:     Select  $\theta_k \in$  Live with the highest  $\hat{v}(\theta_k, \mu)$ 
5:     Live := Live  $\setminus \{\theta_k\}$ 
6:     Branch on  $\theta_k$  generating  $\theta_{k_1}, \dots, \theta_{k_m}$ 
7:     for  $1 \leq p \leq m$  do
8:       if  $\hat{v}(\theta_{k_p}, \mu) >$  Incumbent  $+$   $\beta$  then
9:         if  $\theta_{k_p}$  is completely defined then
10:           Incumbent :=  $\hat{v}(\theta_{k_p}, \mu)$ 
11:           Solution :=  $\theta_{k_p}$ 
12:         else Live := Live  $\cup \{\theta_{k_p}\}$ 
13:   until Live =  $\emptyset$ 
14:   return Solution

```

---

$H(\theta, \mu)$ , where:

$$G(\theta, \mu) = \sum_h R(h, d(h)) + \gamma \sum_{(h, \omega)} \mu(h) v^{d(h), \omega, \sigma(h, \omega)}(h)$$

where  $d(h)$  and  $\sigma(h, \omega)$  are constrained,

$$H(\theta, \mu) = \sum_h \max_a R(h, a) + \gamma \sum_{(h, \omega)} \mu(h) \max_{v^a, \omega, x} v^{a, \omega, x}(h)$$

where  $d(h)$  and  $\sigma(h, \omega)$  are not constrained. Notice that  $R(h, d(h))$  is given by  $\sum_s \mu(h) p(s|h) R(s, d(h))$ .

**H-PI** search starts with a pool of live nodes with a partially specified mapping  $\theta$ , where none of the variables are specified, see Algorithm 2. Moreover, the value hereof is used as the value (called incumbent) of the current best solution, (line 2). At each iteration of the search, a node  $\theta$  that yields the highest upper-bound is selected for exploration from the pool of live nodes, (lines 4-5). Then, a branching is performed: two or more children of the node are constructed through the specification of a single variable, (line 6). Furthermore, for each of the generated child nodes  $\theta_k$ , the upper-bound is computed. In this case, the current node corresponds to a completely specified node, its upper-bound is its exact value at  $\mu$ , the value hereof is compared to the incumbent, and the best solution and its value are kept, (lines 8-11). If its upper-bound is not better than the incumbent, the node is discarded, since no completely specified descendant nodes of that node can be better than the incumbent. Otherwise, the possibility of a better solution in the descendant nodes cannot be ruled out, and the node is then joined to the pool of live nodes, (line 12). When the search tree has been completely explored, the algorithm starts a new search tree with a new distribution over histories  $\mu$ , until all have been processed, and this at each iteration.

**H-PI** may be considered as an extension and generalization of either near-optimal search methods such as **I-MAA\***, or point-based search techniques for solving finite-horizon DEC-POMDPs including **MBDP**, **PBIP**. Indeed, **H-PI** is designed to provide error-bounds on the solution produced as does near-optimal methods. **H-PI** meets

this requirement either by planning only over a small set of distributions  $\mu$  or by using parameter  $\beta$ , or doing both. In addition, **H-PI** is able to scale up through the selection of a small set of distributions  $\mu$ , and by planning only over a small number of histories among those where  $\mu(h) > 0$ . In particular, when we plan separately over histories  $h$  where  $\mu(h) > 0$ , we actually perform a point-based search method as does **MBDP**, **PBIP**. Even within the latter case, **H-PI** remains fundamentally different with respect to other point-based search methods. The key difference lies in how the heuristic function is computed. While finite-horizon DEC-POMDP heuristic functions are all based only on the current state of assignments of values to variables, **H-PI** performs an additional step of iterative elimination of  $\beta$ -dominated trios after each node expansion step, thus tightening its heuristic function. The following provide theoretical guarantees on the solution produced by **H-PI**.

## 4.2 Convergence and Error-bound

For any set of distributions over histories  $B$  and iteration  $\tau$ , **H-PI** produces an estimate  $v_\tau$  with the corresponding set of hyperplanes  $\Lambda_\tau$ . The error between  $v_\tau$  and the true value function  $v^*$  is bounded. The bound depends on four parameters: the density  $\varepsilon_B$  of the set of distributions over histories  $B$ , where  $\varepsilon_B$  is the maximum distance from any legal distribution  $\mu$  to  $B$ , that is:  $\varepsilon_B = \max_{\mu' \in \Delta_{\bar{H}}} \min_{\mu \in B} \|\mu' - \mu\|_1$ ; the distance  $\beta_B$  between hyperplanes that compose  $v_\tau$ , where  $\beta_B$  is the maximum Chebyshev distance of any pair of hyperplanes into  $\Lambda_\tau$ , that is:  $\beta_B = \max_{v^x, v^y \in \Lambda_\tau} \|v^x - v^y\|$ ; the probability  $\mu_B = 1 - \min_{\tau} \sum_{h \in H_\tau} \mu_\tau(h)$  that a history is visited during the online execution stage, but not taken into account during the offline planning stage; and the Chebyshev distance  $\|r\| = \max_{s, a} |R(s, a)|$  over the rewards  $R(s, a)$ , which defines maximum possible rewards that occur after a one step decision.

That is, by keeping all non-dominated hyperplanes over a denser sampling of distribution set  $\Delta_{\bar{H}}$ ,  $v_\tau$  converges to  $v^*$ , the true value function. Cutting off **H-PI** iterations at any sufficiently large time step, we know that the divergence between  $v_\tau$  and the optimal value function  $v^*$  is bounded. The following lemma states and proves a bound on the error  $\|(\beta_B\text{-}\mathbb{H}_B)v_\tau - (0\text{-}\mathbb{H})v_\tau\|$  produced by one application of the backup operator  $\beta_B\text{-}\mathbb{H}_B$ .

**LEMMA 2.** *The error  $\eta_{prune}$  produced by  $(\beta_B\text{-}\mathbb{H}_B)$  when performing the value function backup over  $B$  instead of  $\Delta_{\bar{H}}$ , is bounded by:  $\eta_{prune} \leq \mu_B \cdot (\beta_B + \varepsilon_B \|r\| / (1 - \gamma))$ .*

**PROOF.** First, we note that applying a similar argument to that used to derive that  $\beta\text{-}\mathbb{H}$  is a contraction mapping, we prove that  $\beta_B\text{-}\mathbb{H}_B$  is also a contraction mapping on  $\mathcal{V}$ . Let  $v$  be a value function in  $\mathcal{V}$ , and  $(0\text{-}\mathbb{H}_B)$  be the backup operator that plans only over distributions  $\mu \in B$  but keeps all non dominated hyperplanes for each distribution  $\mu \in B$ . Using the triangle inequality, we know that the error  $\|(\beta_B\text{-}\mathbb{H}_B)v - (0\text{-}\mathbb{H})v\|$  produced by  $\beta_B\text{-}\mathbb{H}_B$  is bounded by  $\|(\beta_B\text{-}\mathbb{H}_B)v - (0\text{-}\mathbb{H}_B)v\| + \|(0\text{-}\mathbb{H}_B)v - (0\text{-}\mathbb{H})v\|$ . We thus propose a two-fold step method that bounds the two expressions above.

On the one hand, we establish the error  $\phi_1 = \|(\beta_B\text{-}\mathbb{H}_B)v - (0\text{-}\mathbb{H}_B)v\|$  made by preserving only one non  $\beta$ -dominated policy node for each individual history. Let  $h$  be a joint history where  $\beta_B\text{-}\mathbb{H}_B$  makes it worst error. This is achieved by pruning away policy-node  $x^i$  and hyperplane  $v^{x^i x^j}$ . Let  $v^{x_h^i x_h^j}$  be the hyperplane that is maximal for  $h$ . By pruning  $v^{x^i x^j}$ ,  $\beta_B\text{-}\mathbb{H}_B$  makes an error of at most  $\mu(h)[v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h)]$ . Furthermore, we know that  $v^{x^i x^j}(h) \leq v^{x_h^i x_h^j}(h)$ . Therefore,  $\phi_1$

$$\leq \mu(h)[v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h)] \quad (7)$$

$$= \mu(h)[v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h) + (v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h))] \quad (8)$$

$$\leq \mu(h)[v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h) + v^{x_h^i x_h^j}(h) - v^{x^i x^j}(h)] \quad (9)$$

$$= \mu(h)[v^{x_h^i x_h^j} - v^{x^i x^j}] \cdot p(h) \quad (10)$$

$$\leq \|v^{x_h^i x_h^j} - v^{x^i x^j}\| \cdot \mu(h) \quad (11)$$

$$\leq \beta_B \cdot \mu_B \quad (12)$$

The equation (8) results from adding zero ( $v^{x^i x^j}(h) - v^{x_h^i x_h^j}(h)$ ) to equation (7). In inequality (9), we replace the third expression  $v^{x^i x^j}$  on the right hand side by  $v^{x_h^i x_h^j}$ , since  $v^{x_h^i x_h^j}$  is maximal for  $h$ . Rearranging terms in equation (9) yields equation (10) where  $p(h)$  is the matrix form of  $p(s|h)$ . Applying the Chebyshev norm and the definition of  $\beta_B$  result in inequalities (11) and (12), respectively.

On the other hand, we establish the error  $\phi_2 = \|(0\text{-}\mathbb{H}_B)v - (0\text{-}\mathbb{H})v\|$  produced by  $(0\text{-}\mathbb{H}_B)$  by planning only over  $B$  instead of  $\Delta\bar{H}$ . Let  $\mu' \in \Delta\bar{H} \setminus B$  be the distribution where  $\beta_B\text{-}\mathbb{H}_B$  makes its worst error, and  $\mu \in B$  be the closest sampled distribution to  $\mu'$ . Let  $u$  be the value function that would be maximal at  $\mu'$ . Let  $v$  be the value function that is maximal at  $h$ . By failing to include hyperplanes that compose  $u$  in its solution set,  $(0\text{-}\mathbb{H}_B)$  makes an error of at most  $u(\mu') - v(\mu')$ . In addition, we know that  $v(\mu) \geq u(\mu)$ . So,  $\phi_2$

$$\leq u(\mu') - v(\mu') \quad (13)$$

$$= u(\mu') - v(\mu') + (u(\mu) - u(\mu')) \quad (14)$$

$$\leq u(\mu') - v(\mu') + v(\mu) - u(\mu) \quad (15)$$

$$= (u - v) \cdot (\mu' - \mu) \quad (16)$$

$$\leq \|u - v\| \cdot \|\mu' - \mu\|_1 \cdot \mu_B \quad (17)$$

$$\leq \varepsilon_B \cdot \mu_B \cdot \|r\|/(1 - \gamma) \quad (18)$$

The equation (14) results from adding zero ( $u(\mu) - u(\mu')$ ). In inequality (15), we replace the third expression  $u(\mu)$  on the right hand side by  $v(\mu)$ , since  $v$  is maximal at  $\mu$ . Rearranging terms in equation (15) yields equation (16). Inequality (17) follows from Hölder inequality and inequality (18) results from the definition of  $\varepsilon_B$ . This ends the proof.  $\square$

**THEOREM 2.** *For any distribution set  $B$  and any iteration  $\tau$ , the error of  $\mathbf{H}\text{-}\mathbf{PI}$  algorithm,  $\|v_\tau - v^*\|$ , is bounded by:  $\eta_\tau \leq \varepsilon/2 + (\beta_B/(1 - \gamma) + \|r\|\varepsilon_B/(1 - \gamma)^2) \mu_B$ .*

**PROOF.** The overall error  $\eta_\tau$  in  $\mathbf{H}\text{-}\mathbf{PI}$  at iteration  $\tau$  is bounded by  $\|v_\tau - v_\tau^*\| + \|v_\tau^* - v^*\|$ . Because  $\beta_B\text{-}\mathbb{H}_B$  is a contraction mapping, when the stopping criterion  $\|v_\tau - v_{\tau-1}\| \leq \varepsilon(1 - \gamma)/\gamma$  holds, the second term  $\|v_\tau^* - v^*\|$  is bounded by  $\varepsilon/2$ . The remainder of this proof states and demonstrates a bound on the first term  $\eta_\tau = \|v_\tau - v_\tau^*\|$  as follows:  $\eta_\tau = \|(\beta_B\text{-}\mathbb{H}_B)v_{\tau-1} - (0\text{-}\mathbb{H})v_{\tau-1}^*\|$

$$\leq \|(\beta_B\text{-}\mathbb{H}_B)v_{\tau-1} - (0\text{-}\mathbb{H})v_{\tau-1}\| + \|(0\text{-}\mathbb{H})v_{\tau-1} - (0\text{-}\mathbb{H})v_{\tau-1}^*\|$$

This follows from the definition of backup operators  $\beta_B\text{-}\mathbb{H}_B$  and  $0\text{-}\mathbb{H}$ , as well as the norm properties. We note that the first term on the right hand side of the last inequality is in fact error estimate  $\eta_{\text{prune}}$ . Moreover, as  $0\text{-}\mathbb{H}$  is a contraction mapping, the second

term on the right hand side of the last inequality is bounded by  $\gamma\|v_{\tau-1} - v_{\tau-1}^*\|$ . Replacing these terms yields:

$$\eta_\tau \leq \eta_{\text{prune}} + \gamma\|v_{\tau-1} - v_{\tau-1}^*\| \quad (19)$$

Then, the error-bound follows as a consequence of Lemma 2, the definition of  $\eta_{\tau-1} = \|v_{\tau-1} - v_{\tau-1}^*\|$  and series sum properties:

$$\begin{aligned} \eta_\tau &\leq \eta_{\text{prune}} + \gamma\eta_{\tau-1} \\ &\leq \left(\beta_B + \frac{\|r\|\varepsilon_B}{1-\gamma}\right) \cdot \mu_B + \gamma\eta_{\tau-1} \\ &\leq \left(\frac{\beta_B}{1-\gamma} + \frac{\|r\|\varepsilon_B}{(1-\gamma)^2}\right) \mu_B \end{aligned}$$

This ends the proof.  $\square$

This result is rather intuitive. Indeed, the error produced by the  $\mathbf{H}\text{-}\mathbf{PI}$  relies on three terms. The first  $\varepsilon/2$  denotes the error produced by cutting off  $\mathbf{H}\text{-}\mathbf{PI}$  iterations when the stopping criterion is reached. The second term  $\beta_B/(1 - \gamma)$  represents the error produced by adding only non  $\beta$ -dominated hyperplanes – and in some case only a single hyperplane for each joint history. In other words, by pruning all  $\beta_B$ -dominated hyperplanes for each joint history. The last term  $\varepsilon_B\|r\|/(1 - \gamma)^2$  illustrates the error produced by planning only over a small set  $B$ . The overall error states the relationship between exact  $\mathbf{PI}$ ,  $\beta\text{-}\mathbf{PI}$  and  $\mathbf{H}\text{-}\mathbf{PI}$  algorithms.

This result synthesizes error-bounds for policy iteration algorithms with respect to three criteria: the backup operator used; the distribution set, and the pruning criterion. This general error-bound can be used to bound the error produced by any algorithm designed within  $\beta\text{-}\mathbf{PI}$ 's algorithmic framework. In particular, when we plan only over a single joint history  $h$  at a time using  $\mathbf{H}\text{-}\mathbf{PI}$  – as does point-based algorithms including  $\mathbf{MBDP}$  [16],  $\mathbf{MBDP}\text{-OC}$  [7],  $\mathbf{PBIP}$  [10], and  $\mathbf{PBIP}\text{-IPG}$  [2], the error is bounded by  $\varepsilon/2 + \frac{\beta_B}{(1-\gamma)} + \frac{\|r\|\varepsilon_B}{(1-\gamma)^2} \mu_B$ . However, if we plan over the entire distribution  $\mu$ ,  $\mathbf{H}\text{-}\mathbf{PI}$  yields a tighter error-bound, *i.e.*,  $\varepsilon/2 + \frac{\|r\|}{(1-\gamma)^2} \cdot \varepsilon_B \cdot \mu_B$ .

This error-bound also suggests that  $\mathbf{H}\text{-}\mathbf{PI}$  can tighten the error even more when its distributions set  $B$  is uniformly dense in the set of reachable distributions  $\Delta\bar{H}$ . Selecting the best distribution set in this sense would require the generation of all possible distributions given all possible decision rule and the distributions at hand. As it current stands, we do not address this problem. The selection of our distribution set, is based on trajectories of distributions. We create trajectories based on the current value function. Each such trajectory starts with the initial distribution  $\mu_0$ , we then executes the greedy decision rule specified by the current value function, and finally select the successor distribution.

## 5. EMPIRICAL EVALUATIONS

We now evaluate the performance of  $\mathbf{H}\text{-}\mathbf{PI}$  in comparison with other recent approximate solvers, such as  $\mathbf{NLP}$  and  $\mathbf{BPI}$ . Experiments have been run on Intel Core Duo 1.83GHz CPU processor with 1Gb main memory.

### 5.1 Results

We begin by demonstrating the advantage of  $\mathbf{H}\text{-}\mathbf{PI}$  with respect to  $\mathbf{NLP}$  and  $\mathbf{BPI}$ . As we can see in graphs in Figure 5,  $\mathbf{H}\text{-}\mathbf{PI}$  outperforms  $\mathbf{NLP}$  and  $\mathbf{BPI}$  in all tested DEC-POMDP domains and in both computation time and solution quality.

As we explain above,  $\mathbf{H}\text{-}\mathbf{PI}$  plans only over a small set of selected distributions experienced by the agents during the offline planning stage. Such distributions often lie near a structured, low-dimensional subspace. For example, in the boxpushing domain, we only consider 6 distributions since the domain is very structured

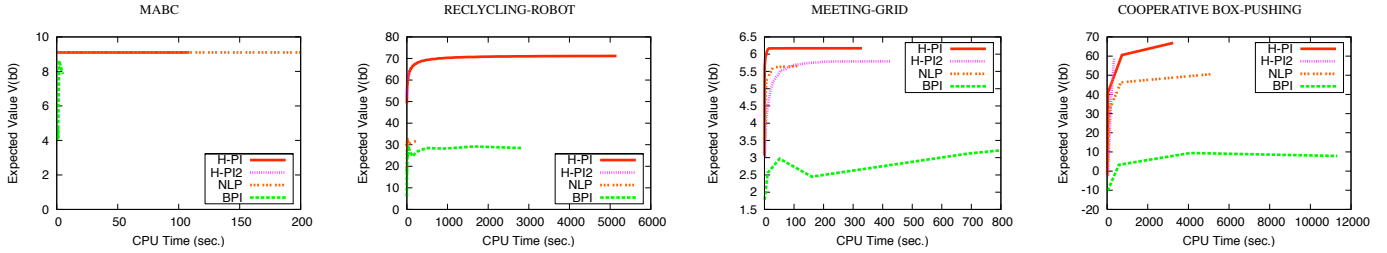


Figure 5: Performance results for DEC-POMDP benchmark problems from the literature

– more precisely often there is only a single possible next observation for a given history, this considerably limits the number of possible next distributions. While **NLP** and **BPI** compute the policy based on a continuum, by planning over this low-dimensional subspace, **H-PI** saves considerable computational efforts – thus increasing its ability to find good solutions very quickly. Furthermore, it builds the best possible joint-policy graph – assigning a single policy node for each individual history. This tightens the size of the solution produced. Notice, however, that the size of the solution is not bounded by  $B$ . Indeed, in the infinite-horizon case the policy nodes are interconnected – that is by keeping policy node  $x$  we also keep policy nodes that are reachable starting from  $x$ . Finally, because of all its enhancements, **H-PI** does not need to bound the size of the solution produced – it is able to provide larger policy graphs for problems that require such policies so as to achieve reasonable performances. For example, in the recycling robot domain, **H-PI** produces twice the expected value produced using either **NLP** or **BPI** but it requires policies that are 250 times larger than those in either **NLP** or **BPI**. It is worthwhile to notice that policy graphs produced by **NLP** and **BPI** are stochastic – thus even though the number of policy nodes is reduced, the equivalent deterministic policy graph would be much larger.

We continue to study the performance of **H-PI** with respect to the distribution set  $B$ . When we plan other belief states corresponding to histories where  $\mu(h) > 0$ , **H-PI** is referred to as **H-PI<sub>2</sub>**, otherwise we use **H-PI**. When evaluating the performance of **H-PI** in comparison with **H-PI<sub>2</sub>** – see graphs Figure 5, we note that **H-PI<sub>2</sub>** is faster but keeps too much joint policy nodes – this limits its performances in comparison to **H-PI**. This is mainly because, **H-PI<sub>2</sub>** often keeps many policy nodes of each individual history, while **H-PI** only keeps a single one. Moreover, as we already discussed, by planning over belief states rather than distributions over histories the error-bound is larger. This explains why the expected value produced by **H-PI** is always superior to the one produced using **H-PI<sub>2</sub>**. However, by increasing the number of belief states considered we may increase the expected value. As illustrated in Figure 6, as the number of belief states grows the solution quality improves and the computation time also grows (and vice versa). Even more importantly, at some point increasing the number of belief states do not provide significant improvement in the solution quality. These observations support the theoretical results on the error produced by **H-PI**, that is a denser sampling of the set  $\Delta\bar{H}$  produces more distributions and results in a tighter error bound. It also highlights the impetus of using a sampling method that selects good distributions or belief states.

## 6. CONCLUSION

We have introduced a new algorithmic framework ( $\beta$ -**PI**) that exploits the scalability of the approximate methods while preserving the theoretical properties of the near-optimal techniques. In particular, it provides the ability to bound the error produced when we approximate the solution using the sufficient statistic in gen-

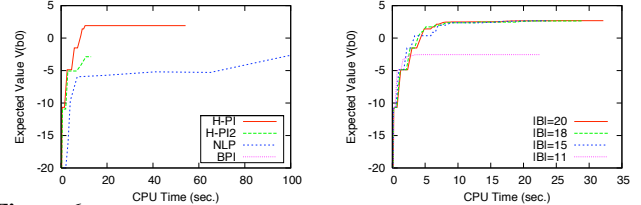


Figure 6: Performance results of the **H-PI** algorithm for the multi-agent tiger domain, and different belief space sizes.

eral DEC-POMDPs. We introduce a heuristic derivation of  $\beta$ -**PI**, namely **H-PI**. We have demonstrated how **H-PI** outperforms state-of-the-art infinite-horizon DEC-POMDP solvers in all tested domains. In this paper we identify the general requirements from a  $\beta$ -**PI** solver, and suggested a possible implementation for DEC-POMDPs. In the future, we will investigate the integration of methods for the selection of good distributions. We also intend to apply  $\beta$ -**PI** to factored domains such as *fire-fighting* or *network distributed sensors* [14], reducing the dimensionality of the sufficient statistic – thus enabling us to scale to even larger domains.

## 7. REFERENCES

- [1] C. Amato, D. S. Bernstein, and S. Zilberstein. Optimizing memory-bounded controllers for decentralized pomdps. In *UAI*, 2007.
- [2] C. Amato, J. S. Dibangoye, and S. Zilberstein. Incremental policy generation for finite-horizon dec-POMDPs. in *ICAPS*, 2009.
- [3] R. Aras and A. Dutech. An investigation into Mathematical Programming for Finite Horizon Decentralized POMDPs. in *JAIR*, 2010. to appear.
- [4] D. S. Bernstein, C. Amato, E. A. Hansen, and S. Zilberstein. Policy iteration for decentralized control of Markov Decision Processes. *JAIR*, 34:89–132, 2009.
- [5] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control of Markov Decision Processes. *Math. Oper. Res.*, 27(4), 2002.
- [6] A. Boularias and B. Chaib-draa. Exact dynamic programming for decentralized POMDPs with lossless policy compression. In *ICAPS*, pages 20–27, 2008.
- [7] A. Carlin and S. Zilberstein. Value-based observation compression for DEC-POMDPs. In *AAMAS*, 2008.
- [8] J. S. Dibangoye. *Contribution à la résolution des problèmes décisionnels de Markov centralisés et décentralisés: algorithmes et théorie*. PhD thesis.
- [9] J. S. Dibangoye, B. Chaib-draa, and A.-I. Mouaddib. Policy iteration algorithms for DEC-POMDPs with discounted rewards. in *MSDM*, 2009.
- [10] J. S. Dibangoye, A.-I. Mouaddib, and B. Chaib-draa. Point-based incremental pruning heuristic for solving finite-horizon DEC-POMDPs. in *AAMAS*, 2009.
- [11] E. A. Hansen, D. S. Bernstein, and S. Zilberstein. Dynamic programming for partially observable stochastic games. In *AAAI*, pages 709–715, 2004.
- [12] A. Kumar and S. Zilberstein. Point-based backup for decentralized pomdps: complexity and new algorithms. In *AAMAS*, pages 1315–1322, 2010.
- [13] M. L. Putterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, New York, NY, 1994.
- [14] D. V. Pynadath and M. Tambe. Multiagent teamwork: analyzing the optimality and complexity of key theories and models. In *AAMAS*, pages 873–880, 2002.
- [15] Z. Rabinovich, C. V. Goldman, and J. S. Rosenschein. The complexity of multiagent systems: the price of silence. In *AAMAS*, pages 1102–1103, 2003.
- [16] S. Seuken and S. Zilberstein. Formal models and algorithms for decentralized decision making under uncertainty. *JAAIAS*, 17(2):190–250, 2008.
- [17] D. Szer and F. Charpillet. An optimal best-first search algorithm for solving infinite horizon DEC-POMDPs. In *ECML*, pages 389–399, 2005.
- [18] D. Szer and F. Charpillet. Point-based dynamic programming for DEC-POMDPs. In *AAAI*, pages 16–20, July 2006.