



Institute of Science and Technology

Optimal cost indefinite-horizon reachability in goal DEC-POMDPs

1 Anonymous, 2 Anonymous

Deposited at: 12 Dec 2018 11:53

ISSN: 2664-1690

IST Austria (Institute of Science and Technology Austria)
Am Campus 1
A-3400 Klosterneuburg, Austria

Copyright © 2015, by the author(s).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.



Institute of Science and Technology

Optimal Cost Indefinite-Horizon Reachability in Goal DEC-POMDPs

1 Anonymous and 2 Anonymous

Technical Report No. IST-2015-326-v1+1
Deposited at 19 Feb 2015 10:25
<http://repository.ist.ac.at/326/1/main.pdf>

IST Austria (Institute of Science and Technology Austria)
Am Campus 1
A-3400 Klosterneuburg, Austria

Copyright © 2012, by the author(s).

All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

Optimal Cost Indefinite-Horizon Reachability in Goal DEC-POMDPs

Anonymous1 and Anonymous2

February 19, 2015

Abstract

DEC-POMDPs extend POMDPs to a multi-agent setting, where several agents operate in an uncertain environment independently to achieve a joint objective. DEC-POMDPs have been studied with finite-horizon and infinite-horizon discounted-sum objectives, and there exist solvers both for exact and approximate solutions. In this work we consider Goal-DEC-POMDPs, where given a set of target states, the objective is to ensure that the target set is reached with minimal cost. We consider the indefinite-horizon (infinite-horizon with either discounted-sum, or undiscounted-sum, where absorbing goal states have zero-cost) problem. We present a new method to solve the problem that extends methods for finite-horizon DEC-POMDPs and the RTDP-Bel approach for POMDPs. We present experimental results on several examples, and show our approach presents promising results.

1 Introduction

POMDPs and DEC-POMDPs. Partially-observable Markov decision processes (POMDPs) are standard models for problems related to probabilistic planning, where an agent tries to optimize an objective in an uncertain environment [16, 23]. There are a wide range of applications of POMDPs ranging from reinforcement learning [18], to software verification [11] to robot motion planning [17]. However, in many scenarios there is not a single agent, but a set of agents whose joint goal is to optimize an objective in an uncertain environment. Each agent has a different view of the state space of the system, and must choose a *local* policy based on her own view, such that the *joint* policy optimizes the objective function. For example, a classic scenario is when two robots are necessary to achieve a task (like moving a large box) and each of them has a specific view of the environment [26, 7]. *Decentralized POMDPs* (DEC-POMDPs) provide the appropriate model for such scenarios, which are extension of POMDPs to multi-agent setting, where each agent has her personal view of the entire system [26, 7, 2, 22].

Indefinite-horizon objectives. In this work we consider DEC-POMDPs with indefinite-horizon objectives. We consider DEC-POMDPs with a set of target (or goal) states, and a cost function that assigns a cost to each transition. Given a discount factor $0 < \gamma \leq 1$, the cost of a path is the discounted-sum of the costs of the transitions in the path. If $\gamma = 1$, then we refer to the objective as the undiscounted-sum objective, and if $\gamma < 1$, then we refer as discounted-sum objective. The costs are accumulated until the target is reached (i.e., once a target state is reached, from then on the costs are zero). The objective function we consider is to ensure that the target set is reached with probability 1 and the accumulated cost according to the discounted-sum (or undiscounted-sum)

is minimized. The objectives we consider are the classical optimization criteria in the setting of POMDPs and DEC-POMDPs. Note that while a finite-horizon objective requires to optimize the cost for a given finite number of steps, and an infinite-horizon objective requires to optimize over the whole length of the paths, an indefinite-horizon objective requires to optimize until the target set is reached.

Restriction. Most problems related to optimizations in POMDPs are undecidable [24, 21]. To develop a practical approach to solve POMDPs, restricted class of POMDPs are considered, where from every state there is a path to a target state, and such POMDPs are referred to as Goal-POMDPs [9]. Several examples of POMDPs that arise in practice can be modeled as Goal-POMDPs [9, 19]. Hence in this work we consider DEC-POMDPs with the restriction that from every state there is a path to a target state, and refer to them as Goal-DEC-POMDPs.

Previous results. One of the most efficient and practical approach to solve Goal-POMDPs is the RTDP-Bel approach that extends the classical RTDP approach [5]. There are several methods, exact and approximate to solve finite-horizon and infinite-horizon DEC-POMDPs [26, 1] (see Related work for detailed discussion).

Our contributions. In this work we present an approach to solve Goal-DEC-POMDPs with indefinite-horizon objectives. Our approach combines and extends two classical algorithmic approaches, namely, the RTDP-Bel approach for POMDPs, and the principle of policy-iteration algorithm for MDPs (which has been considered for DEC-POMDPs as well [4]). Our approach is as follows: we consider an agent, and fix the policies of all other agents, and treat the problem as a POMDP and obtain an improved policy for the agent using a RTDP-Bel approach to solve the POMDP. We repeat this process for all agents. In principle, we fix policies for the agents, and use the RTDP-Bel approach for the policy improvements. The two challenges for this approach are: (i) how to fix policies in DEC-POMDPs and interpret the problem as POMDP; and (ii) then translate the policy from the POMDP to the DEC-POMDP. For the first issue, we compute a probability distribution over states that mimics the belief-state for the agent, even though in the multi-agent setting there is no good notion of belief-state. For the second issue, we merge belief-states in the POMDP that corresponds to probability distribution over states in the DEC-POMDP, and thus interpret the policy from the POMDP in the DEC-POMDP. In addition to our approach, we present several heuristics to make the approach practical. Our heuristics mainly achieve state-space reduction for the POMDP construction, and are based on (i) merging belief-states that are close in terms of probability distribution; and (ii) even if the DEC-POMDP reaches a point where no policy is returned by the RTDP-Bel approach, instead of playing a naive policy, we play based on policies that correspond to similar beliefs in the POMDP. We have implemented our approach along with the heuristics, and we present experimental results on several DEC-POMDP examples from the literature that demonstrates the effectiveness of our approach.

1.1 Related work

POMDPs. There are several works for discounted POMDPs [20, 28, 25], as well as for Goal-MDPs and Goal-POMDPs [9, 19, 12]. In this work we extend the Goal-POMDP solution to Goal-DEC-POMDPs.

DEC-POMDPs. The problem of DEC-POMDPs with various objectives has been an active research area, see [27] for a survey. Some important works related to DEC-POMDPs are as follows: (i) exact algorithms for finite-horizon objectives [15, 31];

(ii) approximate methods for finite-horizon objectives [22, 26]; (iii) approximate methods for infinite-horizon objectives [1, 7, 30]. The undecidability for infinite-horizon objectives was established in [6]. A translation of DEC-POMDPs to continuous-state MDPs was established in [13]. The error-bounded approximations for infinite-horizon discounted-sum objectives was considered in [14]. An approach based on fixing policies in turns for finite-horizon objectives was considered in [4]. While [4] does not consider indefinite-horizon objectives, we show how to extend the policy iteration along with RTDP-Bel approach to DEC-POMDPs with indefinite-horizon objectives. One of the most closely related work is [2] that also consider reaching goals in DEC-POMDPs, and our approach is very different: [2] considers an approach where policies for all the agents are constructed simultaneously, and in contrast we consider extension of policy iteration (for single agent with other policies fixed) and RTDP-Bel approach to solve the problem.

2 Goal POMDPs and DEC-POMDPs

In this section we present the definitions of POMDPs, DEC-POMDPs, strategies, objectives and other notions that are necessary to present our results.

We use \mathbb{N} and \mathbb{R} to denote the sets of positive integers and real numbers, respectively. Given a finite set X , we denote by $\mathcal{P}(X)$ the set of subsets of X , i.e., $\mathcal{P}(X)$ is the power set of X . A probability distribution f on X is a function $f : X \rightarrow [0, 1]$ such that $\sum_{x \in X} f(x) = 1$, and we denote by $\mathcal{D}(X)$ the set of all probability distributions on X . For $f \in \mathcal{D}(X)$ we denote by $\text{Supp}(f) = \{x \in X \mid f(x) > 0\}$ its support.

DEC-POMDPs. A *Decentralized Partially Observable Markov Decision Process (DEC-POMDP)* is defined as a tuple $D = (I, S, \{\mathcal{A}_i\}, \delta, \{\mathcal{Z}_i\}, \mathcal{O}, c, \lambda)$, where (i) I is a finite set of agents; (ii) S is a set of states; (iii) \mathcal{A}_i is a finite set of actions for every agent $i \in I$; (iv) $\delta : S \times \prod_{i \in I} \mathcal{A}_i \rightarrow \mathcal{D}(S)$ is a probabilistic transition function that given the current state s and the set of actions $\vec{a} \in \prod_{i \in I} \mathcal{A}_i$ for every agent gives the probability distribution over the successor states; (v) \mathcal{Z}_i is a finite set of observations for every agent $i \in I$; (vi) $\mathcal{O} : S \times \prod_{i \in I} \mathcal{A}_i \rightarrow \mathcal{D}(\prod_{i \in I} \mathcal{Z}_i)$ gives the probability of observing $\vec{z} \in \prod_{i \in I} \mathcal{Z}_i$ given the current state and the played actions; (vii) $c : S \times \prod_{i \in I} \mathcal{A}_i \rightarrow \mathbb{R}$ is a function that given a state and the played actions assigns a cost; and (viii) $\lambda \in \mathcal{D}(S)$ is the initial state distribution. We will for simplicity use $\vec{\mathcal{A}}$ as a shortcut for $\prod_{i \in I} \mathcal{A}_i$ and similarly $\vec{\mathcal{Z}}$ for $\prod_{i \in I} \mathcal{Z}_i$.

POMDPs. A *Partially Observable Markov Decision Process (POMDP)* is a special type of DEC-POMDP where the set of agents I is restricted to be a singleton.

Plays and cones. A *play* (or a path) in a DEC-POMDP is an infinite sequence $\rho = (s_0, \vec{a}_0, s_1, \vec{a}_1, s_2, \vec{a}_2, \dots)$ of states and action sets $\vec{a}_i \in \vec{\mathcal{A}}$ such that $s_0 \in \text{Supp}(\lambda)$ and for all $i \geq 0$ we have $\delta(s_i, \vec{a}_i)(s_{i+1}) > 0$. We write Ω for the set of all plays. For a finite prefix $w \in (S \cdot \vec{\mathcal{A}})^* \cdot S$ of a play, we denote by $\text{Cone}(w)$ the set of plays with w as the prefix (i.e., the cone or cylinder of the prefix w), and denote by $\text{Last}(w)$ the last state of w .

Strategies (or policies). A *local strategy* (or a policy) for agent $i \in I$ is a recipe to extend prefixes of plays and is a function $\sigma_i : (\mathcal{Z}_i \cdot \mathcal{A}_i)^* \cdot \mathcal{Z}_i \rightarrow \mathcal{D}(\mathcal{A}_i)$ that given a finite history of observations observed by agent i and actions played by agent i selects a probability distribution over the actions \mathcal{A}_i . For agent $i \in I$ we denote by σ_i^U the strategy that plays all the actions \mathcal{A}_i uniformly at random¹, i.e., for all histories $\rho =$

¹For simplicity and wlog, we consider that all actions in \mathcal{A}_i are available in all states.

$(z_0, a_0, z_1, a_1, \dots, a_{n-1}, z_n)$ and actions $a \in \mathcal{A}_i$ we have $\sigma_i^U(\rho)(a) = 1/|\mathcal{A}_i|$. Given a local strategy σ_i for every agent $i \in I$, we denote by σ_I the joint strategy $\prod_{i \in I} \sigma_i$ and σ_I^U for $\prod_{i \in I} \sigma_i^U$.

Probability and Expectation Measures. Given a joint strategy σ_I and an initial state distribution λ , the unique probability measure obtained given σ_I is denoted as $\mathbb{P}_\lambda^{\sigma_I}(\cdot)$. Let $\mathcal{K} : (S \cdot \bar{\mathcal{A}})^* \cdot S \rightarrow \mathcal{D}((\bar{\mathcal{Z}} \cdot \bar{\mathcal{A}})^* \cdot \bar{\mathcal{Z}})$ be the probability distribution that given a finite prefix w of a play gives a probability distribution over the joint observation and action sequence ρ observed by the joint strategy. We first define the measure $\mu_\lambda^{\sigma_I}(\cdot)$ on cones. For $w = s$, where $s \in S$, we have $\mu_\lambda^{\sigma_I}(\text{Cone}(w)) = \lambda(s)$ and for $w' = w \cdot \bar{a} \cdot s$ we have $\mu_\lambda^{\sigma_I}(\text{Cone}(w')) = \mu_\lambda^{\sigma_I}(\text{Cone}(w)) \cdot \sum_{\rho \in \text{Supp}(\mathcal{K}(w))} \mathcal{K}(w)(\rho) \cdot \sigma_I(\rho)(\bar{a}) \cdot \delta(\text{Last}(w), \bar{a})(s)$. By Carathéodory's extension theorem, the function $\mu_\lambda^{\sigma_I}(\cdot)$ can be uniquely extended to a probability measure $\mathbb{P}_\lambda^{\sigma_I}(\cdot)$ over Borel sets of infinite plays [8]. We denote by $\mathbb{E}_\lambda^{\sigma_I}[\cdot]$ the expectation measure associated with the strategy σ_I .

Objectives. We consider the following objectives:

(Un)Discounted Sum Objectives. The traditional optimization objective in (DEC-)POMDPs is defined as follows: given a discount factor $0 < \gamma \leq 1$ and a play $\rho = (s_0, \bar{a}_0, s_1, \bar{a}_1, s_2, \bar{a}_2, \dots)$ the discounted *infinite-horizon* sum of the the play ρ is the value $\text{Sum}_\gamma^\infty(\rho) = \sum_{i=0}^\infty \gamma^i \cdot c(s_i, \bar{a}_i)$. A discounted *finite-horizon* sum, where a horizon length k is given, only considers the prefix of length k of a play, i.e., $\text{Sum}_\gamma^k(\rho) = \sum_{i=0}^k \gamma^i \cdot c(s_i, \bar{a}_i)$. The *undiscounted sum* is a special case where the discount factor γ is required to be 1.

Reachability Objectives. Goal (DEC-)POMDPs are additionally equipped with a *reachability objective*. A reachability objective is a measurable set of plays $\varphi \subseteq \Omega$ defined by a set of absorbing *goal states* $G \subseteq S$ as follows: $\text{Reach}(G) = \{(s_0, \bar{a}_0, s_1, \bar{a}_1, s_2, \bar{a}_2, \dots) \in \Omega \mid \exists i \geq 0 : s_i \in G\}$. In the framework of goal (DEC-)POMDPs it is additionally required that for all states $s \in G$ and actions $\bar{a} \in \bar{\mathcal{A}}$, (i) we have $\delta(s, \bar{a})(s) = 1$ (i.e., s is absorbing); and (ii) the cost function assigns $c(s, \bar{a}) = 0$; i.e., once a play reaches a goal state it stays there and stops to accumulate costs. Observe that the cost restriction on the absorbing goal states of goal (DEC-)POMDPs allows the indefinite-horizon objectives to be expressed as infinite-horizon objectives.

Almost-sure winning. Given a (DEC-)POMDP D with a reachability objective $\text{Reach}(G)$ a joint strategy σ_I is *almost-sure winning* iff $\mathbb{P}_\lambda^{\sigma_I}(\text{Reach}(G)) = 1$. We denote by $\text{Almost}_D(G, \lambda)$ the set of all joint almost-sure winning strategies from the initial state distribution λ .

Optimal cost under almost-sure winning. Given a goal (DEC-)POMDP D with a reachability objective $\text{Reach}(G)$ and a cost function c we are interested in minimizing the expected discounted (or undiscounted) sum of costs before reaching the goal set G , while ensuring that the goal set is reached almost-surely. Formally, the value of an almost-sure winning joint strategy $\sigma_I \in \text{Almost}_D(G, \lambda)$ is the expectation $\text{Val}(\sigma_I) = \mathbb{E}_\lambda^{\sigma_I}[\text{Sum}_\gamma^\infty]$, and the objective is to minimize $\text{Val}(\sigma_I)$ over almost-sure winning strategies.

In our work we assume that the (DEC-)POMDP do not contain states from which the goal states are no longer reachable, i.e., for every initial state distribution λ the joint strategy σ_I^U belongs to the set of almost-sure winning strategies $\text{Almost}_D(G, \lambda)$.

3 Approximate Algorithm for Goal POMDPs

In this section we summarize the key concepts of the RTDP-Bel algorithm [9] that is an approximate solver for goal POMDPs. In the following section we show how to extend

RTDP-Bel for solving goal DEC-POMDPs, inspired by an approach from the work [4] for finite-horizon objectives.

The state of the art algorithm for solving goal POMDPs is RTDP-Bel [9], that is an adaptation of the older RTDP algorithm [5] to goal POMDPs. It uses the well-known approach of turning a POMDP into a completely observable continuous MDP over belief states [3, 29].

Formally, a *belief state* is a probability distribution $b \in \mathcal{D}(S)$, that represents the current information the agent has about the current state.

Belief state updates. Given the current belief state b , the action a played by the agent, and the observation z observed, it is straightforward to compute the newly reached information state b' [10] using the transition and observation probabilities. We write $b' = \text{Update}(b, a, z)$ for the belief state update.

Belief state discretization. In order to bound the size of the MDP, the algorithm uses a discretization method Disc that maps belief states to *qbeliefs*. A *qbelief* is a function $q : S \rightarrow \mathbb{N}$. Given a belief state b the corresponding discretized qbelief $q = \text{Disc}(b)$ is defined as follows: for every state $s \in S$ the value of $q(s)$ is $\text{ceil}(D \cdot b(s))$, where D is a positive integer parameter of the method. Note, that the discretization method preserves the support of the belief states.

Algorithm. The algorithm runs simulations of the POMDP as follows. Every simulation starts in the initial belief state λ , and terminates when a goal state or an upper bound on the length of the simulations is reached. The information learned from the simulations is stored in a hashtable that maps qbeliefs to an estimation of the expected cost until a goal state is reached, i.e., the value of the qbelief.

During the simulation the algorithm maintains the current belief state b . In every step of the simulation for every action $a \in \mathcal{A}$ and observation $z \in \mathcal{Z}$ the new belief state $b' = \text{Update}(b, a, z)$ is computed that is discretized into a qbelief $q' = \text{Disc}(b')$. The next action to be played by the strategy is the one that minimizes the estimations of the expected costs stored in the hashtable for the qbeliefs q' . The data from the simulation are then used to update the hashtable, and to improve the strategy for the following simulations.

There are theoretical downsides of using the discretization method: the convergence is not guaranteed, and the values of the learned strategy may oscillate. Additionally, the learned expected cost is not guaranteed to be an upper bound on the optimal cost. However, the learned strategy can be evaluated (by fixing the strategy and obtaining the value in the Markov chain) and the obtained value is guaranteed to be an upper bound on the optimal cost. Despite the theoretical disadvantages, in practice the RTDP-Bel method works very well. In the experimental sections of [9, 12] it is shown that RTDP-Bel scales to examples that could not be solved exactly before. Finally, on smaller examples, the computed values match the values obtained from an exact solver. Hence in practice RTDP-Bel algorithm is the state-of-the-art solver for goal POMDPs.

4 Approximate Algorithm for Goal DEC-POMDPs

In this part we present our new approach that extends the RTDP-Bel algorithm to provide joint strategies for DEC-POMDPs. We build on the approach presented in [4] for finite-horizon objectives, which is based on the principle of policy iteration. Informally, we fix policies of all but one agent, obtain a POMDP, and use RTDP-Bel for the policy improvement in a POMDP.

Algorithm overview. Intuitively, we initialize every agent $i \in I$ with an initial strategy σ_i^0 and fix all the initial strategies of the agents in the DEC-POMDP except for one. As the strategies for all but one agent, say i , are fixed, the resulting model is a POMDP. By running RTDP-Bel on the POMDP we obtain a new strategy σ_i^1 for agent i . In the next round we choose another agent j and fix all the current strategies of the remaining agents, construct a POMDP, and improve the current strategy of agent j . We iteratively cycle through agents and improve their strategies until the desired value of the joint strategy is reached. For simplicity of presentation and without loss of generality we will consider only two agent DEC-POMDPs in the following text.²

Main challenges. There are two key issues relevant to our algorithmic approach. First, is fixing strategies for agents in a DEC-POMDP and obtain a POMDP as a result. In the case of two agents i and j , the belief state update of the agent i whose strategy is fixed depends on the action played by agent j . Thus unless all actions are fixed, belief update interpretation is difficult and we take an optimistic (or cooperative) approach since the agents have a joint objective. Second, we obtain a strategy in the reduced POMDP, and the challenge is to interpret the strategy of the POMDP in the DEC-POMDP. Below we present the solution for both the issues.

Fixing a strategy in a DEC-POMDP. Given a goal DEC-POMDP $D = (\{i, j\}, S, (\mathcal{A}_i, \mathcal{A}_j), \delta, (\mathcal{Z}_i, \mathcal{Z}_j), \mathcal{O}, c, \lambda)$, let σ_i be a strategy for agent i that is for simplicity of the type $\sigma_i : \mathcal{D}(S) \rightarrow \mathcal{D}(\mathcal{A}_i)$, i.e., given a belief state gives a probability distribution over the actions to be played next. For a DEC-POMDP D with a set of goal states G , we denote by $P = D \upharpoonright \sigma_i$ the goal POMDP P that is obtained from DEC-POMDP D by fixing strategy σ_i of agent i . The POMDP $P = (\{j\}, S', \{\mathcal{A}_j\}, \delta', \{\mathcal{Z}_j\}, \mathcal{O}', c', \lambda')$ is defined as follows:

- The states are $S' = \{(s, b) \mid s \in S, b \in \mathcal{D}(S)\}$, i.e., a state of D and a probability distribution that mimics the belief state of agent i .
- The transition probability $\delta'((s, b), a_j)((s', b'))$ for $a_j \in \mathcal{A}_j$ is defined as the product of the follows. Let $P(b, \vec{a}) : \mathcal{D}(S) \times \overline{\mathcal{A}} \rightarrow \mathcal{D}(\mathcal{D}(S))$ denote the probability distribution of updating to belief state b' , given the current belief state is b and actions \vec{a} are played. The transition probability is given as follows:

$$\begin{aligned} \delta'((s, b), a_j)((s', b')) = \\ \sum_{a_i \in \mathcal{A}_i} \sigma_i(b)(a_i) \cdot \delta(s, (a_i, a_j))(s') \cdot P(b, (a_i, a_j))(b') \end{aligned}$$

i.e., for every action a_i of agent i the components are the probability that agent i plays action a_i , the transition probability of moving to state s' , and the probability of updating to belief state b' . Note that due to the presence of multiple agents, the notion of a belief state is not well defined, as it depend on the action of agent j . The intuitive explanation of the distribution stored in the state, is the belief state of agent i if agent j cooperates and plays the expected action.

- The observation function \mathcal{O}' is a projection of the function \mathcal{O} , conditioned on strategy σ_i . Formally:

$$\begin{aligned} \mathcal{O}'((s, b), a_j)(z_j) = \\ \sum_{a_i \in \mathcal{A}_i} \sigma_i(b)(a_i) \cdot \sum_{z_i \in \mathcal{Z}_i} \mathcal{O}(s, a_i, a_j)((z_i, z_j)) \end{aligned}$$

²In case of more agents, it is intuitively possible to merge the remaining agents into a single agent in every iteration and instead of a set of strategies consider the joint strategy of the merged agents. The presented methods can be easily extended to handle joint strategies.

- The cost c' is weighted over the actions played by strategy σ_i , i.e., the function is:

$$c((s, b), a_j) = \sum_{a_i \in \mathcal{A}_i} \sigma_i(b)(a_i) \cdot c(s, (a_i, a_j))$$

- The initial belief state for agent i is λ , therefore the distribution is defined $\lambda'((s, b)) = \lambda(s)$ for $b = \lambda$ and 0 otherwise.

The goal states of the POMDP $P = D \upharpoonright \sigma_i$ are all states (s, b) , where s is an goal state of DEC-POMDP D . The output POMDP P is then used as an input for the RTDP-Bel algorithm. The tool outputs a deterministic strategy σ_j for the POMDP, i.e., given a belief state the strategy returns a single action to be played next. In the next part we show how to interpret a strategy σ_j of POMDP P in the DEC-POMDP D .

Interpreting the POMDP strategy. Let the current two strategies for agents i and j be σ_i^k and σ_j^ℓ , respectively. We fix a strategy σ_i^k in the DEC-POMDP D and obtain a strategy σ_j for the POMDP $P = D \upharpoonright \sigma_i^k$. In this part we present how to interpret the strategy σ_j in the context of the DEC-POMDP and obtain an improved strategy $\sigma_j^{\ell+1}$ for agent j . The strategy σ_j is of type $\mathcal{D}(S \times \mathcal{D}(S)) \rightarrow \mathcal{D}(\mathcal{A}_j)$, whereas the strategy $\sigma_j^{\ell+1}$ needs to be of type $\mathcal{D}(S) \rightarrow \mathcal{D}(\mathcal{A}_j)$.

We define a transformation function `Reduce` that reduces belief states b' from the DEC-POMDP to belief states b of the POMDP as follows: Given a belief $b' \in \mathcal{D}(S \times \mathcal{D}(S))$ in P , the reduced belief $b = \text{Reduce}(b')$, where $b \in \mathcal{D}(S)$ in D , for a state $s \in S$ is defined as follows:

$$b(s) = \sum_{\bar{b} \in \mathcal{D}(S)} b'(s, \bar{b})$$

In other words, given a state s we sum the probabilities of all states, where the first component is s , i.e., the states differ only in the belief state of agent i .

Algorithm 1 GOAL DEC-POMDPS

Input: goal DEC-POMDP D , Integer k

Output: A joint strategy σ_I

▷ Initialize the agents with initial strategies

for $i \in I$ **do**
 $\sigma_i^0 \leftarrow \sigma_i^U$

▷ Iteratively improve the strategies

for $0 \leq u < k$ **do**
for $j \in I$ **do**
POMDP $P \leftarrow D \upharpoonright (\prod_{i \in I \setminus j} \sigma_i^u)$
 $\sigma_j \leftarrow \text{RTDP-Bel}(P)$
 $\sigma_j^{u+1} \leftarrow \text{InterpretOnDEC-POMDP}(\sigma_j)$

$\sigma_I \leftarrow \prod_{i \in I} \sigma_i^k$
return σ_I

Given a belief state b in the DEC-POMDP D , the new strategy $\sigma_j^{\ell+1}(b)$ is defined as follows:

1. Compute all belief states b' such that $b = \text{Reduce}(b')$, where the corresponding qbeliefs $q = \text{Disc}(b')$ are stored in the hashtable. If there are no such beliefs play all the available actions uniformly at random.
2. For every belief state b' compute the discretized qbelief $q = \text{Disc}(b')$ and utilize the hashtable stored by RTDP-Bel that for any given qbelief stores its value, i.e., the expected sum of costs before the goal state is reached.
3. The strategy $\sigma_j^{\ell+1}(b)$ is defined to play the action according to the optimal value over the considered qbeliefs stored in the hashtable with probability 99%. With the remaining probability 1% play all the other actions from \mathcal{A}_j uniformly at random. Note that this ensures that all actions are played with positive probability, and as a consequence we can show (in item 1 of Theorem 1) that the POMDP we obtain is a goal POMDP.

Our approach is summarized as Algorithm 1. In the first step all the agents are initialized with a strategy that always plays all available actions uniformly at random. The algorithm is parametrized with an integer k , that specifies how many times should a strategy of an individual agent be improved.

Theorem 1 *The following assertions hold:*

1. *During the computation of Algorithm 1, the constructed POMDPs are goal POMDPs.*
2. *The value of the computed joint strategy σ_I is an upper bound on value of the optimal-cost joint strategy.*

For the first item, note that for every agent $i \in I$ and every $0 \leq u \leq k$ the strategy σ_i^u satisfies that for every belief state b we have $\text{Supp}(\sigma_i^u(b)) = \mathcal{A}_i$. Therefore, any path in the DEC-POMDP can be mapped to a path in the POMDP that is obtained by fixing the strategies of the agents. Since we consider goal DEC-POMDP, it follows that in the POMDP there is a path from every state to a target state, and hence we obtain a goal POMDP. For the second item, note that the value learned from RTDP-Bel is not guaranteed to be an upper bound on the optimal value, however, we can evaluate the joint strategy, and the value of the computed joint strategy is an upper bound on the optimal value. \square

Remark 1 *Note that in our approach we store state and belief state pairs, instead of states and qbelief pairs. We remark that the simpler approach to replace belief states with qbeliefs in the states of the POMDPs, that are obtained by fixing the strategies of the agents, might not yield satisfactory results. Intuitively, the stored belief state may evolve for a sequence of steps, even when the corresponding qbelief remains unchanged. The qbelief is changed only after a certain number of steps. This behavior is not captured if qbeliefs instead of the belief states are considered.*

5 Heuristics

In this section we present a number of possible heuristics to make to presented approach practical. Intuitively, any joint strategy σ_I , that satisfies that the individual strategies of the agents always play all the available actions with positive probability, is almost-sure winning for the reachability objective. Therefore, the we focus on heuristics that reduce the size of the state space of the POMDPs constructed in Algorithm 1, to speed up the construction of the POMDP and run RTDP-Bel on smaller instances.

Merge states with close belief states. Given a DEC-POMDP D and a strategy σ_i for agent $i \in I$ the constructed POMDP $P = D \upharpoonright \sigma_i$ has a continuous set of states due to the

fact that the states s are elements in $S \times \mathcal{D}(S)$. We exploit the fact, that if two probability distributions are similar, then it suffices to keep one of them as representative and reduce the size of the POMDP. Given a threshold $\delta > 0$ we say that two belief states $b, b' \in \mathcal{D}(S)$ are δ -close if the following two conditions hold: (1) $\text{Supp}(b) = \text{Supp}(b')$; and (2) for all $s \in S : |b(s) - b'(s)| \leq \delta$. In the constructed POMDP, whenever a new state (s, b) is to be added, it is checked that whether there exists a state (s, b') , where b and b' are δ -close. If it is the case, then transitions to (s, b) are redirected to (s, b') , otherwise a new state (s, b) is added. This heuristic ensures that the constructed POMDP has a finite state set, and can be used to adjust the size of the POMDP.

Relax the interpretation of the strategy. As defined in Section 4 (item 1 of Interpreting the POMDP strategy), whenever no matching qbelief is found in the hashtable from the RTDP-Bel computation, a default naive option of playing all the available actions uniformly at random is chosen. In this heuristic we consider instead other qbeliefs stored in the hashtable, that are close to the qbelief that is missing (in a similar sense as in the previous heuristic). The resulting distribution over the actions will still play all the available actions with positive probability, but the resulting distribution is a weighted sum based on the closeness of the other qbeliefs and their probability distributions over played actions.

Initial strategies of the agents. It is possible to initialize individual agents with different initial strategies, that can affect the performance of the computation. Without prior knowledge of the model, the strategy that plays all the available actions uniformly at random performs well, as it explores the largest state space of the DEC-POMDP with high probability.

More aggressive reductions of the state space. For larger examples one can significantly reduce the state space of the POMDP by deploying more aggressive reductions. One can view the belief state component of the state as an information that the strategy uses to make decisions. There are various possibilities to reduce the size, and we give a few examples:

- Remove states s from the belief state b (assign $b(s) = 0$), when their probability $b(s)$ is below a certain threshold and distribute the probability mass to the other states.
- For a given integer n keep only the n highest probable states in the belief state and remove all the others.

6 Implementation and Experimental Results

We have implemented in Java our approach presented in Section 4 together with the heuristics from Section 5. We have modified the existing implementation of RTDP-Bel in C++ to export the strategy together with the content of the hashtable storing the values and actions for individual qbeliefs to a file.

Key difference to existing tools. Several benchmarks presented in the literature are naturally expressed as goal DEC-POMDPs. However, in the literature the infinite-horizon discounted-sum objective is considered, rather than indefinite-horizon (discounted-sum or undiscounted-sum) objectives. Thus our setup is different from most existing tools. Since in the literature the objective is infinite-horizon discounted-sum, for the benchmarks the objectives are expressed as follows: whenever a target state of the DEC-POMDP is reached a negative cost (reward) is acquired and the DEC-POMDP is restarted back to the initial distribution. Note that in this approach undiscounted-sum objectives are not allowed, and the discounted-sum costs are accumulated even after

Meeting in the Grid problem $ S = 16, \mathcal{A} = 3, \mathcal{Z} = 2, I = 2$		
Num of iter.	Val	Tot. time
0	10.28	-
1	6.5	4.1s
2	6.5	8.2s
3	4.87	9.64s
4	4.87	13.7s
5	4.87	15.1s
6	4.87	17.6s

Table 1: Results for the *Meeting in the Grid problem*

the goal has been reached. Our algorithm allows to specify goal states in the input file and can handle undiscounted-sum (as well as discounted-sum) indefinite horizon problems.

Experimental setup. For every example we run a number of iterations of strategy improvements. The constructed strategies are then evaluated by constructing a Markov chain and the value of the strategies is computed by running simulations within the RTDP-Bel algorithm. The time needed to evaluate the strategies after the last iteration finishes is very small and is not included in the running times. We have experimented with a number of examples from the literature on an Intel Core i7 (2.5 GHz) CPU equipped with 16GB of RAM. We have extended the `.dpomdp` file format with the possibility to specify goal states. For every POMDP constructed we run 20k simulations with RTDP-Bel, to construct a strategy. We use the built-in standard cut-off length 250 that bounds the length of simulations. In none of the examples our algorithm was using more than 1.5 GB of RAM.

Meeting in a grid problem. The problem is based on the Meeting in the Grid problem of [7] and used later in [2, 14]. We require two agents that move in a 2×2 grid to meet at the same grid position. The state where the two agents share a single grid position is a goal state. Initially the agents are placed uniformly at random among all the positions in the grid, where they do not share the same position. We assign a cost of 1 to every move of the agents until they meet in the grid. The objective is to minimize the undiscounted-sum of costs, i.e., the expected number of steps the agents need to meet at a grid position (in other words, this is the stochastic shortest path problem). In the Table 1 we report the obtained values for individual iterations, where iteration 0 corresponds to the naive joint strategy, where all the agents play all available action uniformly at random. The values of the strategies after the 6th iteration do not change much.

Decentralized Tiger problem. The problem is an extension of the well-known tiger problem POMDP [17] to the multi-agent setting and is used extensively in the literature [13, 2]. The problem models two agents that are standing in a hallway with two doors. Behind one of the doors is a tiger, behind the other a treasure. Therefore there are two states: the tiger is behind the left door or behind the right door. Both agents have three actions at their disposal: open the left door, open the right door, and listen. They cannot observe the action of the other agent. Whenever both agents open jointly a door with the treasure the model reaches a goal state. Otherwise the model is restarted (tiger is placed again randomly). Every move before reaching the goal state has cost 1, we consider the undiscounted-sum objective. The results are presented in Table 4.

Unlike in the previous example, one iteration was sufficient to obtain strategies that do not improve with more iterations (in this example the optimal value is 2).

Decentralized Tiger problem $ S = 3, \mathcal{A} = 3, \mathcal{Z} = 2, I = 2$		
Num of iter.	Val	Tot. time
0	17.37	-
1	2.05	0.37s
2	2.04	0.58s

Table 2: Results for the *Decentralized Tiger problem*

Box-pushing problem. The cooperative box-pushing problem is a well-known robotics problem introduced by [KZ97]. Two agents have to cooperate to move a big object (the box) that they could not move on their own. Even though the robots cannot communicate with each other, they have to achieve a certain degree of coordination to move the box at all. The problem was considered as a DEC-POMDP in [26]. We define the goal states as the situation when some of the boxes is pushed in its goal position.

Comparison with other DEC-POMDP solvers. The other existing related tools for DEC-POMDPs are the *goal Directed* algorithm introduced in [2] and the state-of-the-art infinite-horizon discounted-sum DEC-POMDP solver introduced in [13]. These implementations are not (publicly) available for comparison. We compare our results with the results reported from [13, 2]. Note that since we have indefinite-horizon undiscounted-sum objectives, and the cost assignment are different, the computed values are different. We have also modified DEC-POMDPs by adding new goal states. Nevertheless, the remaining states, transitions, and observations are similar. Also the results were computed on different but similar platforms. Hence the time comparison maybe approximate at best.

Rover example. In our experimental result we compare with the examples studied in [2] which also considers. Along with the examples we presented, [2] also considers another example namely, Mars Rover. For Mars Rover our algorithm performs poorly, mainly, due to the construction of the POMDP which is quite large. Both the construction of the POMDP and RTDP-Bel running on the POMDP was impractical.

Box Pushing problem $ S = 100, \mathcal{A} = 4, \mathcal{Z} = 5, I = 2$		
Num of iter.	Val	Tot. time
0	556.68	-
1	414.82	73.4s
2	252.52	136.3s
3	252.52	199.2s
4	252.52	260.5s

Table 3: Results for the *Box Pushing problem*

Algorithm	Dec. Tiger		Meet in the Grid		Box pushing	
	Val.	Time (s)	Val.	Time (s)	Val.	Time (s)
BFS	-14.1	12007	4.2	17	-2	1696
DEC-PBI	-52.6	102	3.6	2227	9.4	4094
NLP	-1.1	6174	5.7	117	54.23	1824
Goal Directed	5.0	75	5.6	4	149.9	199
FB-HSVI	13.5	6	-	-	199.4	15.2
Our algorithm		0.4		9.6		136

Table 4: Time comparison of the existing algorithms

7 Conclusions

In this work we presented a new approach, based on RTDP-Bel for goal POMDPs and the principle of policy iteration algorithm, for goal DEC-POMDPs. Along with our approach we presented a number of heuristics. One bottleneck of our approach is that the POMDP constructed can be large, and an important direction of future work would be to explore methods for state-space reduction of the POMDP. In DEC-POMDP an important concern is to obtain small-size policies, which was considered in [2]. Whether our approach can be extended to obtain small policies is another direction of future work.

References

- [1] C. Amato, D. S. Bernstein, and S. Zilberstein. Optimizing Memory-Bounded Controllers for Decentralized POMDPs. In *UAI*, pages 1–8, 2007.
- [2] C. Amato and S. Zilberstein. Achieving goals in decentralized POMDPs. In *AAMAS*, 2009.
- [3] K. J. Åström. Optimal control of Markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications*, 10(1):174–205, 1965.
- [4] B. Banerjee, J. Lyle, L. Kraemer, and R. Yellamraju. Sample Bounded Distributed Reinforcement Learning for Decentralized POMDPs. In *AAAI*, 2012.
- [5] A. G. Barto, S. J. Bradtke, and S. P. Singh. Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72(1):81–138, 1995.
- [6] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The Complexity of Decentralized Control of Markov Decision Processes. *Math. Oper. Res.*, 27(4):819–840, 2002.
- [7] D. S. Bernstein, E. A. Hansen, and S. Zilberstein. Bounded policy iteration for decentralized POMDPs. In *IJCAI*, pages 52–57, 2005.
- [8] P. Billingsley, editor. *Probability and Measure*. Wiley-Interscience, 1995.
- [9] B. Bonet and H. Geffner. Solving POMDPs: RTDP-Bel vs. Point-based Algorithms. In *IJCAI*, pages 1641–1646, 2009.

- [10] A. R. Cassandra, L. P. Kaelbling, and M. L. Littman. Acting optimally in partially observable stochastic domains. In *AAAI*, volume 94, pages 1023–1028, 1994.
- [11] P. Cerný, K. Chatterjee, T. A. Henzinger, A. Radhakrishna, and R. Singh. Quantitative synthesis for concurrent programs. In *Proc. of CAV*, LNCS 6806, pages 243–259. Springer, 2011.
- [12] K. Chatterjee, M. Chmelik, R. Gupta, and A. Kanodia. Optimal Cost Almost-sure Reachability in POMDPs. In *AAAI*, 2015.
- [13] J. S. Dibangoye, C. Amato, O. Buffet, and F. Charpillet. Optimally Solving Dec-POMDPs as Continuous-State MDPs. In *IJCAI*, 2013.
- [14] J. S. Dibangoye, O. Buffet, and F. Charpillet. Error-Bounded Approximations for Infinite-Horizon Discounted Decentralized POMDPs. In *ECML/PKDD*, pages 338–353, 2014.
- [15] E.A. Hansen, D. S. Bernstein, and S. Zilberstein. Dynamic programming for partially observable stochastic games. In *AAAI*, volume 4, pages 709–715, 2004.
- [16] H. Howard. *Dynamic Programming and Markov Processes*. MIT Press, 1960.
- [17] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1):99–134, 1998.
- [18] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: A survey. *J. of Artif. Intell. Research*, 4:237–285, 1996.
- [19] A. Kolobov, Mausam, D.S. Weld, and H. Geffner. Heuristic search for generalized stochastic shortest path MDPs. In *ICAPS*, 2011.
- [20] H. Kurniawati, D. Hsu, and W.S. Lee. SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. In *Robotics: Science and Systems*, pages 65–72, 2008.
- [21] O. Madani, S. Hanks, and A. Condon. On the undecidability of probabilistic planning and related stochastic optimization problems. *Artif. Intell.*, 147(1-2):5–34, 2003.
- [22] R. Nair, M. Tambe, M. Yokoo, D. Pynadath, and S. Marsella. Taming decentralized POMDPs: Towards efficient policy computation for multiagent settings. In *IJCAI*, pages 705–711, 2003.
- [23] C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of Markov decision processes. *Mathematics of Operations Research*, 12:441–450, 1987.
- [24] A. Paz. *Introduction to probabilistic automata (Computer science and applied mathematics)*. Academic Press, 1971.
- [25] J. Pineau, G. Gordon, S. Thrun, et al. Point-based value iteration: An anytime algorithm for POMDPs. In *IJCAI*, volume 3, pages 1025–1032, 2003.
- [26] S. Seuken and S. Zilberstein. Improved Memory-Bounded Dynamic Programming for Decentralized POMDPs. In *UAI*, pages 344–351, 2007.

- [27] S. Seuken and S. Zilberstein. Formal models and algorithms for decentralized decision making under uncertainty. *AAMAS*, 17(2):190–250, 2008.
- [28] T. Smith and R. Simmons. Heuristic search value iteration for POMDPs. In *UAI*, pages 520–527. AUAI Press, 2004.
- [29] E. J. Sondik. The optimal control of partially observable Markov processes over the infinite horizon: Discounted costs. *Operations Research*, 26(2):282–304, 1978.
- [30] D. Szer and F. Charpillet. An optimal best-first search algorithm for solving infinite horizon DEC-POMDPs. In *ECML*, pages 389–399. Springer, 2005.
- [31] D. Szer, F. Charpillet, and S. Zilberstein. MAA*: A Heuristic Search Algorithm for Solving Decentralized POMDPs. In *UAI*, pages 576–590, 2005.