

Multi-Robot Active Information Gathering with Periodic Communication

Mikko Lauri*, Eero Heinänen[†], Simone Frintrop*

Abstract—A team of robots sharing a common goal can benefit from coordination of the activities of team members, helping the team to reach the goal more reliably or quickly. We address the problem of coordinating the actions of a team of robots with periodic communication capability executing an information gathering task. We cast the problem as a multi-agent optimal decision-making problem with an information theoretic objective function. We show that appropriate techniques for solving decentralized partially observable Markov decision processes (Dec-POMDPs) are applicable in such information gathering problems. We quantify the usefulness of coordinated information gathering through simulation studies, and demonstrate the feasibility of the method in a real-world target tracking domain.

I. INTRODUCTION

Teams of robots are projected to be applied in a wide range of information gathering tasks, ranging from locating victims in search and rescue scenarios [1] to simultaneous localization and mapping [2]. To benefit most from the deployment of multiple robots, the robot team should coordinate its activities. Coordination becomes more challenging when communication between team members is limited. In this paper, we study how the team members can coordinate their activities in an information gathering task under periodic communication.

As a motivating example, consider the scenario in Figure 1 in which two robotic agents, here micro aerial vehicles (MAVs), are jointly estimating the state of a moving target. Each MAV can activate either a vision sensor that can detect the target at close range, or a radar sensor that can detect the target when it is further away. The MAVs can periodically communicate and share sensor data, forming a joint estimate of the target state. Between these periods, each MAV must act without knowledge of what the other is doing. Coordination benefits the MAVs: if the target is close to the first MAV but far from the second one, the first MAV should try to detect the target with its camera while the second MAV applies its radar sensor. Further, the MAVs can avoid simultaneously operating their radars, avoiding interference that may corrupt the data.

Related work. Controlling robot teams in information gathering tasks is often implemented by applying various relaxations to the control problem to avoid the high computational demands. A fully distributed algorithm applying gradient-based control with a mutual information reward is presented in [3], and in [4] the next best sensing locations for a robot

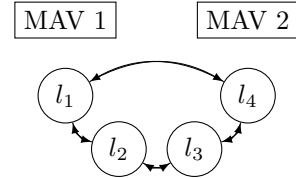


Fig. 1. Two micro aerial vehicles tracking a target between the locations l_i .

team are found via a Gauss-Seidel relaxation. In [5], the decentralized information gathering problem is linearized while also providing suboptimality guarantees. Periodic communication was considered through a problem constraint in [6], and by designing distributed data fusion techniques that can handle communication breaks in [7]. Distributed constraint optimization problems (DCOPs) have been applied to planning for coverage of multiple targets [8], [9]. Partially observable Markov decision processes (POMDPs) were applied to change detection in [10], with each robot reasoning about the beliefs of the others. Decentralized control for environmental monitoring and target tracking based on auctioning POMDP control policies is proposed in [11]. In [12], DCOPs and decentralized POMDPs were exploited for control in sensor networks, where the sensor nodes only have limited local interactions.

Common for many of the approaches above is that the underlying problem they attempt to tackle is a variant of a decentralized partially observable Markov decision process (Dec-POMDP) [13]. A Dec-POMDP is a generic model for cooperative multi-agent sequential decision making under uncertainty, where agents with a shared goal execute actions and perceive observations that provide partial information about the underlying hidden state of the system. Dec-POMDPs explicitly model uncertainty in sensor data and system dynamics, making them an appealing model for robotics problems. A solution to a Dec-POMDP is a control policy, computed centrally and distributed to the agents for execution. Although finding an optimal policy for a Dec-POMDP is NEXP-hard [14], recent work has been able to improve tractability at multi-robot control problems such as package delivery [15], [16].

Single-robot information gathering has also been addressed as a POMDP or a stochastic control problem, see [17], [18], [19], [20]. In these approaches, information gathering is explicitly addressed by defining an information-theoretic reward function, e.g., mutual information or entropy. Our aim here is to translate this into a multi-robot setting through Dec-POMDPs with information-theoretic rewards.

Contribution. Our contribution is summarized as follows. First, we introduce the ρ Dec-POMDP model, which extends

This work was supported by the Academy of Finland decision 268152, Optimal operation of observation systems in autonomous mobile machines.

* Department of Informatics, University of Hamburg, Hamburg, Germany, {lauri, frintrop}@informatik.uni-hamburg.de.

[†] Laboratory of Automation and Hydraulics, Tampere University of Technology, Tampere, Finland, eero.heinanen@student.tut.fi.

the Dec-POMDP to allow information-theoretic rewards. We show how existing Dec-POMDP solution methods can be extended to ρ Dec-POMDPs. Secondly, we quantify through experiments in target tracking domains the usefulness and feasibility of our approach to decentralized information gathering.

We assume that the robots can periodically communicate to share information. In this sense, our approach is less general compared to e.g. [6] who require periodic communication via constraints on the problem. In contrast to e.g. [12], [11], [10], we do not assume any state transition or observation independence properties. Compared to existing approaches for multi-robot control via Dec-POMDPs, ours differs in that we explicitly minimize a measure of state estimate uncertainty.

Organization. Section II reviews the Dec-POMDP. In Section III, we introduce our extension, the ρ Dec-POMDP, and discuss its properties as applied to multi-robot information gathering. Section IV explains how solution algorithms for standard Dec-POMDPs may be applied to ρ Dec-POMDPs. Sections V and VI report results of simulation and real-world target tracking experiments. Section VII concludes the paper.

II. DECENTRALIZED COOPERATIVE DECISION MAKING

We study sequential decision-making by a team of robotic agents, formalized as a decentralized partially observable Markov decision process (Dec-POMDP). The definitions and concepts referred to in this section are based on [21] and [13].

Definition 1 (Dec-POMDP). *A Dec-POMDP is a tuple $\langle I, \mathcal{S}, \{\mathcal{A}_i\}_{i \in I}, \{\mathcal{Z}_i\}_{i \in I}, \mathbb{T}, \mathbb{O}, R, b_0 \rangle$, where*

- $I = \{1, 2, \dots, n\}$ is the set of n agents,
- \mathcal{S} is the finite state space,
- \mathcal{A}_i and \mathcal{Z}_i are the finite action and observation spaces of agent i such that $\mathcal{A} = \times_{i \in I} \mathcal{A}_i$ and $\times_{i \in I} \mathcal{Z}_i$ are the joint action and observation spaces, respectively,
- $\mathbb{T}: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a stochastic state transition model, such that $\mathbb{T}(s', a, s)$ gives the conditional probability of transitioning to state $s' \in \mathcal{S}$ when joint action $a \in \mathcal{A}$ is executed at state $s \in \mathcal{S}$,
- $\mathbb{O}: \mathcal{Z} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a probabilistic observation model, such that $\mathbb{O}(z', s', a)$ gives the conditional probability of perceiving joint observation $z' \in \mathcal{Z}$ in state $s' \in \mathcal{S}$ when the previous joint action was $a \in \mathcal{A}$,
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and
- b_0 is the initial joint belief state, encoding information about the state at time $t = 0$.

At each time step t , each agent i selects an action $a_t^i \in \mathcal{A}_i$, forming a joint action $a_t = (a_t^1, \dots, a_t^n)$. The agents obtain a shared reward according to R . The state then transitions according to \mathbb{T} , and each agent perceives an observation $z_{t+1}^i \in \mathcal{Z}_i$ such that the joint observation $z_{t+1} = (z_{t+1}^1, \dots, z_{t+1}^n)$ is distributed as specified by \mathbb{O} . The objective of the agents is to select actions such that the expected sum of rewards over a given horizon of time $h > 0$ is maximized. To characterize the solution, we give some supporting definitions.

Definition 2 (Local and joint history). *The local history of agent i at time t is $\theta_t^i = (a_0^i, z_1^i, \dots, a_{t-1}^i, z_t^i) \in \Theta_t^i =$*

$\times_t (\mathcal{A}_i \times \mathcal{Z}_i)$. The joint history $\theta_t \in \Theta_t = \times_i \Theta_t^i$ is the collection of each agent's local histories: $\theta_t = (\theta_t^1, \dots, \theta_t^n)$.

Definition 3 (Local and joint decision rule). *A local decision rule δ_t^i of agent i at time t is a function from its local histories to its individual actions: $\delta_t^i: \Theta_t^i \rightarrow \mathcal{A}_i$. A joint decision rule δ_t is a collection of local decision rules: $\delta_t = (\delta_t^1, \dots, \delta_t^n)$.*

Definition 4 (Joint policy). *A joint policy $\pi \in \Pi_h$ for horizon h is a collection of joint decision rules: $\pi = (\delta_0, \delta_1, \dots, \delta_{h-1})$.*

Thus, π is a mapping $\Theta_t \rightarrow \mathcal{A}$ for any t s.t. $0 \leq t < h$. The value $V(\pi)$ of a joint policy is given by

$$V(\pi) = \sum_{t=0}^{h-1} \sum_{\theta_t \in \Theta_t} P(\theta_t | \pi, b_0) R(\theta_t, \pi(\theta_t)). \quad (1)$$

Here, $P(\theta_t | \pi, b_0)$ is the probability of experiencing joint history θ_t when executing policy π starting from b_0 . Furthermore, $R(\theta_t, a_t) = \sum_{s_t \in \mathcal{S}} R(s_t, a_t) P(s_t | \theta_t, b_0)$ is the expected immediate reward of executing joint action a_t after experiencing joint history θ_t . An optimal solution is a joint policy π^* with a maximal value: $\pi^* = \operatorname{argmax}_{\pi \in \Pi_h} V(\pi)$.

The probability mass function (pmf) $P(\cdot | \theta_t, b_0)$ is a sufficient statistic for the state given the joint history $\theta_t = (a_0, z_1, \dots, a_{t-1}, z_t)$ and the joint belief state b_0 . We denote $P(\cdot | \theta_t, b_0) \equiv b_t$, and by \mathcal{B} the space of all such pmfs over the state space \mathcal{S} . Joint belief states are computed recursively applying the Bayesian filtering operator $\tau: \mathcal{B} \times \mathcal{A} \times \mathcal{Z} \rightarrow \mathcal{B}$:

$$b_t = \tau(b_{t-1}, a_{t-1}, z_t) \equiv \frac{1}{\eta(z_t | b_{t-1}, a_{t-1})} \cdot \sum_{s_{t-1} \in \mathcal{S}} \mathbb{T}(s_t, a_{t-1}, s_{t-1}) b_{t-1}(s_{t-1}), \quad (2)$$

where $\eta(z_t | b_{t-1}, a_{t-1})$ is the normalizing factor equal to the prior probability of perceiving z_t when joint action a_{t-1} is taken in joint belief state b_{t-1} . For later use, define

$$b_t = \tau(\theta_t, b_0) \equiv \tau(\tau(\dots, a_{t-2}, z_{t-1}), a_{t-1}, z_t) \quad (3)$$

as an equivalent shorthand notation for expressing the joint belief state b_t as function of the joint history θ_t and b_0 .

III. DECENTRALIZED INFORMATION GATHERING

We next present our extension of the Dec-POMDP model for information-theoretic rewards. Additionally, we motivate and discuss our assumption of periodic communication.

A. Information-theoretic rewards for Dec-POMDPs

Consider an information gathering task for a robot team. It seems clear that an optimal policy should lead the robots to act such as to reach a state estimate with low uncertainty. Reward functions that only depend on the state and action do not appear to be fully compatible with such objectives [22]. For example, consider a robot team collecting information about a physical process which they are unable to affect through their actions, such as monitoring underwater ocean currents. In this

case, the underlying state is not meaningful for the robots' objective, nor are the possible actions by themselves.

In single-agent POMDPs, information gathering has been addressed, e.g., by augmenting the action space to include new actions that reward information-gathering [23], [24], or by applying information-theoretic rewards [22]. The first approach results in a multiplication of the size of the agents' action spaces \mathcal{A}_i . As the number of possible policies for an agent i in a Dec-POMDP is $|\mathcal{A}_i|^{((|\mathcal{A}_i||\mathcal{Z}_i|^h - 1)/(|\mathcal{A}_i||\mathcal{Z}_i| - 1))}$ [21], this approach does not seem promising to translate to Dec-POMDPs. We consider the second approach, which corresponds to setting a reward function that depends on the joint belief state instead of the true underlying state of the system.

In a Dec-POMDP, an individual agent usually cannot compute a joint belief state based only on its local history. However, during the centralized planning phase the possible joint histories θ_t are available, and each joint history corresponds to some joint belief state as indicated by Eq. (3). Thus, information theoretic reward functions may be evaluated and applied in Dec-POMDPs while computing a joint policy. Inspired by [22], we call the resulting model a ρ Dec-POMDP.

Definition 5 (ρ Dec-POMDP). *A ρ Dec-POMDP is a tuple $\langle I, \mathcal{S}, \{\mathcal{A}_i\}_{i \in I}, \{\mathcal{Z}_i\}_{i \in I}, \mathbb{T}, \mathbb{O}, \rho, b_0 \rangle$, where $I, \mathcal{S}, \{\mathcal{A}_i\}_{i \in I}, \{\mathcal{Z}_i\}_{i \in I}, \mathbb{T}, \mathbb{O}$, and b_0 are as in the Dec-POMDP, and the reward function is $\rho : \mathcal{B} \times \mathcal{A} \rightarrow \mathbb{R}$.*

Choosing $\rho(b, a) = \sum_{s \in \mathcal{S}} R(s, a)b(s)$, the definition above subsumes the standard Dec-POMDP (Definition 1). As in Eq. (1), the value of a joint policy π in a ρ Dec-POMDP is

$$V(\pi) = \sum_{t=0}^{h-1} \sum_{\theta_t \in \Theta_t} P(\theta_t | \pi, b_0) \rho(\tau(\theta_t, b_0), \pi(\theta_t)), \quad (4)$$

where $\tau(\theta_t, b_0)$ is the joint belief state computed by Eq. (3).

An appropriate reward function encourages the agents to reach joint histories corresponding to joint belief states with low uncertainty. This is quantified via uncertainty functions.

Definition 6 (Uncertainty function [25]). *Any non-negative, concave function $g : \mathcal{B} \rightarrow \mathbb{R}^+$ is applicable as an uncertainty function.*

An example of an uncertainty function is the Shannon entropy $H(b) = -\sum_s b(s) \log_2 b(s)$ [26]. An uncertainty function has a small value for joint belief states near the degenerate case where all probability mass is concentrated on one state, and a greater value near the uniform distribution.

Throughout the remainder of the paper, we will consider ρ Dec-POMDP models with reward functions defined as

$$\rho(b, a) = \sum_{s \in \mathcal{S}} R(s, a)b(s) - \alpha g(b), \quad (5)$$

where $R(s, a)$ models the rewards dependent on states and actions only, and g is an uncertainty function encoding the information gathering objective¹. The term $\alpha > 0$ sets the balance of state-dependent and information gathering rewards.

¹We apply a minus sign here to penalize for uncertainty in b .

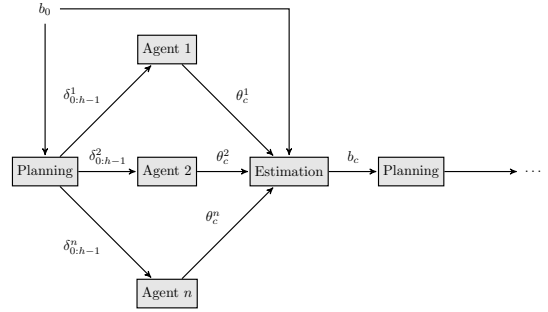


Fig. 2. Decentralized information gathering with communication period c . The agents act according to local decision rules $\delta_{0:h-1}^i$ for c time steps ($c \leq h$), and then share their local histories θ_c^i . A joint belief state b_c may then be estimated, and the information applied to plan subsequent actions.

B. Periodic communication

Deploying a team of robots to collect information is useful if the collected information will eventually be combined to form a joint state estimate. We assume the robots share data via periodic communication every c time steps. This leads to a scheme outlined in Figure 2. Given an initial joint belief state, a joint policy is computed for a horizon h , and distributed to the agents for execution. After c time steps, the agents communicate to share their local histories, enabling estimation of a new joint belief state. The new information may be applied to revise the joint policy. The communication period c need not be equal to the horizon h .

For $c = 1$, information is immediately shared by all agents and the problem reduces to a multi-agent POMDP [27]. The case $c > 1$ is similar to the k -steps delayed communication case [13]: at time t all agents know θ_{t-k} . However, for periodic communication this information is obtained only during the communication intervals, instead of at every time step.

One might argue that it is sufficient to set a non-zero information reward g only for the last time instant of the planning horizon. While appropriate in some applications, a counterargument can be presented for applications such as the robot team monitoring ocean currents. Here, we wish to estimate the state of the process accurately at all time steps to learn about the process dynamics. In such cases, minimizing the “average uncertainty” via Eq. (5) is reasonable.

IV. PLANNING IN ρ DEC-POMDPs

A planning algorithm for standard Dec-POMDPs can be applied to ρ Dec-POMDPs if it does not rest on the assumption that the reward is state and action dependent only. For instance, approaches that cast the problem as a search through the joint policy space Π_h can be applied once reward function evaluation is modified appropriately. For example, variants of generalized multi-agent A* algorithms [21] are applicable.

We applied the multi-agent A* (MAA*) algorithm [28]. MAA* is a heuristic search algorithm similar to the classic A* algorithm. Starting with $t = 1$, MAA* constructs a search tree over partial joint policies $\phi_t = (\delta_0, \delta_1, \dots, \delta_{t-1})$ with $t < h$ that specify how the agents act until time step t . Thus,

each node in the search tree represents a partial joint policy. The search tree is expanded to cover the partial joint policies ϕ_{t+1} by appending to ϕ_t any possible joint decision rule: $\phi_{t+1} = (\phi_t, \delta_t)$. The tree is expanded in a best-first order, as determined by a value estimate of the nodes currently in the tree. The value estimate $\hat{V}(\phi_t)$ of a node with partial joint policy ϕ_t is computed in two parts: by exact evaluation of ϕ_t plus an estimate of the value of the remaining $h-t$ time steps via a heuristic function H_{h-t} : $\hat{V}(\phi_t) = V(\phi_t) + H_{h-t}(\phi_t)$. Here, $V(\phi_t)$ is the value of ϕ_t computed via Eq. (4), and $H_{h-t}(\phi_t)$ is the heuristic value. If H_{h-t} overestimates the true expected reward over the remaining $h-t$ time steps, MAA* returns an optimal policy [28].

To define a heuristic function, the Dec-POMDP is relaxed, e.g., into a centralized POMDP, or into a fully observable centralized Markov decision process (MDP) [21]. The heuristic function is obtained, e.g., in the case of POMDP relaxation, by finding the optimal value of the POMDP. We apply heuristic functions obtained via a POMDP relaxation. This preserves the uncertainty aspect in the problem, and the optimal value can be computed even for information-theoretic reward functions by existing techniques [22].

V. SIMULATION EXPERIMENTS

In this section, we study the usefulness of coordinating information gathering activities in a target tracking domain. To the best of our knowledge, there are currently no methods directly comparable to ours that address the same problem. The related approaches [12], [10] require transition and observation independence, and the auctioning method proposed in [11] and the multi-robot Dec-POMDP studies [15], [16] use state and action based rewards. Instead, we compare optimal ρ Dec-POMDP policies to hand-tuned heuristic policies.

A. Cooperative target tracking domain

In the scenario of Figure 1, the state consists of the target location $l \in \{l_1, l_2, l_3, l_4\} = L$ and a binary variable describing the target status; neutral (0) or hostile (1). The target does not change its status, but moves in a different pattern depending on the status. A neutral target stays in place with probability p_0 , and moves to either neighbouring location with probability $(1-p_0)/2$. A hostile target stays in place with probability p_1 . We set $p_0 = 0.85$ and $p_1 = 0.6$.

Both MAVs have two actions, a_c and a_r , referring to applying a camera or a radar, respectively. The observations $Z_1 = Z_2 = L$ correspond to perceiving the target at any of the locations. The observations do not provide information on the target status, which has to be inferred based on a series of observations. Sensor accuracy depends on the distance to the target. At l_1 , the target is at distance 0 from MAV 1, and at distance 3 from MAV 2, at l_2 , the distances are 1 and 2, and so on. We model the sensors by Gaussian distributions with mean at the target location, and standard deviation dependent on the target status and increasing as a function of the distance. For a sensor mode j , given a distance d , its standard deviation is $\sigma_j(d) = \sigma_{j,0} \cdot 2^{(d/d_{j,0})}$, where $\sigma_{j,0}$ is the nominal standard

TABLE I
SENSOR PARAMETERS IN THE SIMULATION.

Sensor mode j	Neutral target		Hostile target	
	$d_{j,0}$	$\sigma_{j,0}$	$d_{j,0}$	$\sigma_{j,0}$
camera	0.6	0.3	0.7	0.75
radar	1.0	0.2	1.0	0.45
radar (interference)	2.0	1.0	1.5	1.2

deviation, and $d_{j,0}$ is the half efficiency distance of the sensor. The parameters we applied are summarized in Table I.

The reward is as in Eq. (5), with $\alpha = 1$ and g as Shannon entropy. In $R(s, a)$, for each agent applying the radar action a_r , there is a reward of -0.1, and additionally, if the target is hostile, an additional reward of -1 or -0.1 if the target is at distance 0 or 1, respectively. This term models the higher costs of applying the radar, and the risk of revealing the MAVs' own location to a hostile target if radar is engaged at short range.

B. Experimental results

We varied the communication interval c and the initial joint belief state b_0 . We compare the value of an optimal policy of the ρ Dec-POMDP² to five heuristic policies. Policy 1 is a risk-averse strategy where both MAVs only apply cameras. In policy 2, the first MAV always applies its camera, and the second MAV always applies its radar. Policy 3 is the same as policy 2, reversing the roles. Policy 4 is a turn-taking policy where the first MAV starts by applying its camera while the second MAV applies its radar, and on subsequent time steps they switch sensors; policy 5 is the same with reversed roles.

We set b_0 uniform with respect to the target's location, and varied the initial probability that the target is neutral. Figure 3 shows the values of the policies for $h = 3$ as function of the probability that the target is neutral. If it is very likely that the target is hostile, the heuristic policy of only applying the cameras is close to optimal as it avoids the risk of additional costs for radar use on hostile targets. The cameras only policy has a much lower value when it is more likely that the target is neutral. In this case, the other heuristic policies all work equally well, and are close to optimal. None of the heuristic policies can consistently reach near-optimal performance. An important advantage of an optimal policy is that it adapts to changing situations, unlike the fixed heuristic policies.

We then set b_0 uniform, and varied the communication interval c between 1 and 3. At every c th time step, a new policy for horizon $h = 3$ was computed and then applied for the next c decisions. We ran 50 simulations on the task, each for 51 decisions, recording the rewards for each policy. Table II shows the average total rewards with 95% confidence intervals for an optimal policy for $h = 3$ while varying c , each heuristic policy, and a policy of choosing random actions. Policies 2 and 3 are grouped together as "fixed roles", and 4 and 5 as "turn-taking" as there was no significant difference

²We apply MAA* from the MADP toolbox www.fransoliehoek.net/madp/, extended by us to handle information theoretic rewards.

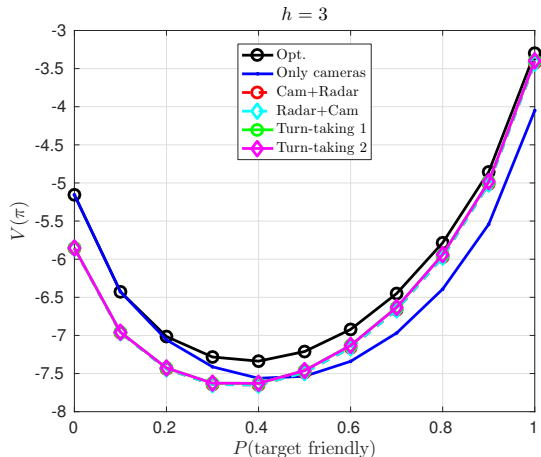


Fig. 3. Values of optimal and heuristic policies as a function of the initial probability of the target being neutral. The policies corresponding to the red, cyan, green and magenta lines have almost equal values and overlap each other.

TABLE II
AVERAGE REWARDS WITH 95% CONFIDENCE INTERVALS.

Policy	Comm. interval	Reward
Optimal	1	-89.9 ± 1.2
	2	-90.1 ± 1.6
	3	-89.8 ± 1.2
Cameras only	-	-96.6 ± 2.0
Fixed roles	-	-95.0 ± 1.5
Turn-taking	-	-90.7 ± 1.4
Random	-	-104.2 ± 2.1

between them. Here the optimal policy performs equally well regardless of the communication interval, while the turn-taking policy also performs well. The lack of improvement for lower c is explained by the fact that the actual horizon of the task is equal to 51, the number of decisions to be taken, which is much longer than the planning horizon applied.

VI. COOPERATIVE TRACKING

We set up a target tracking experiment as shown in Figure 4. The target in the center of the figure was programmed to move randomly in the area. The markers on the target can be detected by the robot on the left applying its laser range finder, and by the robot on the right applying its camera.

To study cooperative target tracking, we applied a Kalman filter (KF) to estimate the target’s position and velocity while limiting the amount of input data. A schematic of the experimental setup is shown in Figure 5. At each 1-second time interval, both observers select a detection sector to focus their attention on, as indicated by the labels a_1 through a_5 . We only input an observation of the target to the KF if it was within the selected detection sector. Each detection sector was 15 degrees wide, with a maximum range of 2.5 and 3 meters, respectively. Selecting overlapping detection sectors could result in interference corrupting the data.

As the KF state estimate is continuous, at each time instant

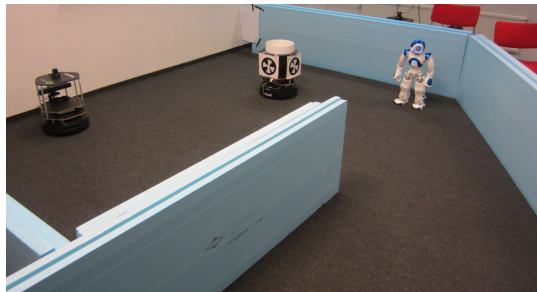


Fig. 4. The robots on the left and right track the target robot at the center.

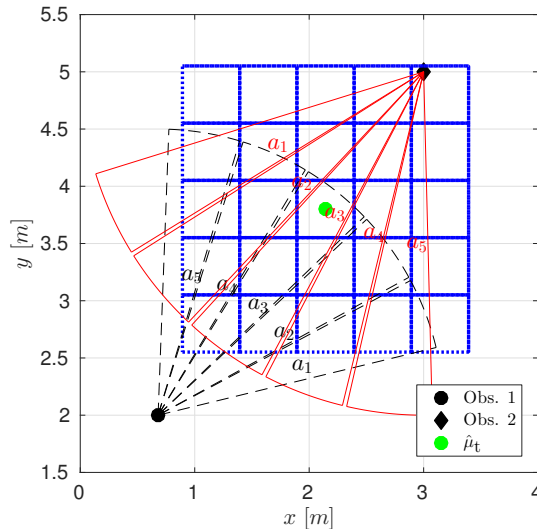


Fig. 5. One of five detection sectors a_1 through a_5 is chosen by each of the two observers. The Kalman filter position estimate is mapped to a finite grid, shown in blue, centred at the estimate mean $\hat{\mu}_t$.

we mapped the tracking task into a finite ρ Dec-POMDP. The KF state estimate was discretized to a 5-by-5 grid centred at the mean $\hat{\mu}_t$ of the current target position estimate. The grid cell size was adaptively tuned so the grid covered the 3-sigma range of the estimate. The detection sectors were defined always setting the middle sector’s center to point towards $\hat{\mu}_t$. For target motion we assumed a Gaussian velocity distribution, with mean equal to the velocity estimate from the KF, and covariance matching the expected velocity of the target. The observations indicated if the target was detected or not within the selected detection sector. There was a nominal false negative probability of 0.15, and a false positive probability of 0.05. If overlapping detection sectors were selected, these probabilities were increased in proportion to the area of overlap, up to a maximum value of 0.5. The reward was as in Eq. (5), R all zero, $\alpha = 1$, and g as Shannon entropy.

The ρ Dec-POMDP optimal policy was computed with $h = 3$ and applied to select detection sectors over the next $c = 3$ time steps. We compared this to a policy where the first robot repeated a_1, a_2, \dots, a_5 , while the second robot repeated this sequence in reversed order, and a random policy. We modelled interference due to overlapping detection sectors by corrupting

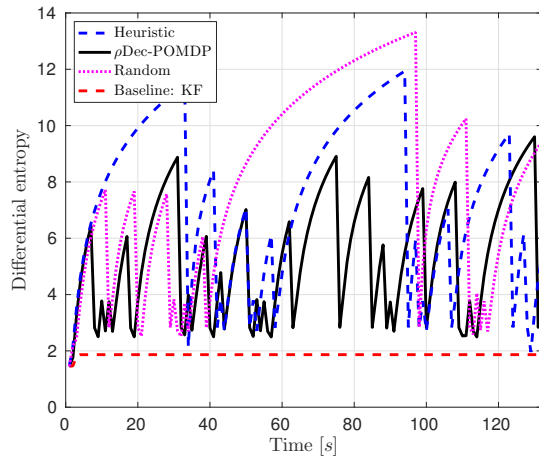


Fig. 6. Differential entropy of the position estimate (lower is less uncertain) as function of time.

the observations according to a probability proportional to the area of overlap. A baseline was computed applying all data in a KF, without regard to the detection sectors.

Figure 6 shows the differential entropy [26] of the position estimate as function of time. Due to its adaptivity, the ρ Dec-POMDP method performs better than the heuristic or random policy, maintaining a lower average entropy. Interference occurred during 2 time steps for ρ Dec-POMDP, and during 41 time steps for the heuristic policy, and 5 time steps for the random policy. Thus, the ρ Dec-POMDP policy avoids the risk of corrupting observations due to selecting overlapping detection sectors. Compared to the baseline KF estimate, the ρ Dec-POMDP policy had a sum of squared position error of 143.4, the heuristic policy 178.4, and the random policy 519.7.

VII. CONCLUSION

For modelling information gathering by a robot team, we presented ρ Dec-POMDP, extending the Dec-POMDP to information-theoretic rewards. A ρ Dec-POMDPs may be solved applying existing algorithms for Dec-POMDPs, with modified reward function evaluation. We verified the feasibility of our approach for cooperative target tracking. Due to the adaptivity of ρ Dec-POMDP policies, they can outperform heuristic approaches. Future work includes extended empirical evaluation, and possibly combining ρ Dec-POMDPs with distributed state estimation to relax communication assumptions.

REFERENCES

- [1] S. Balakirsky, S. Carpin, A. Kleiner, M. Lewis, A. Visser, J. Wang, and V. A. Ziparo, "Towards heterogeneous robot teams for disaster mitigation: Results and performance metrics from RoboCup rescue," *J. Field Robot.*, vol. 24, no. 11-12, pp. 943-967, 2007.
- [2] S. Saeedi, M. Trentini, M. Seto, and H. Li, "Multiple-robot simultaneous localization and mapping: A review," *J. Field Robot.*, vol. 33, no. 1, pp. 3-46, 2016.
- [3] B. J. Julian, M. Angermann, M. Schwager, and D. Rus, "Distributed robotic sensor networks: An information-theoretic approach," *Int. J. Robot. Res.*, vol. 31, no. 10, pp. 1134-1154, 2012.
- [4] K. Zhou and S. I. Roumeliotis, "Multirobot active target tracking with combinations of relative observations," *IEEE T. Robot.*, vol. 27, no. 4, pp. 678-695, Aug. 2011.

- [5] N. Atanasov, J. L. Ny, K. Daniilidis, and G. J. Pappas, "Decentralized active information acquisition: Theory and application to multi-robot SLAM," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, May 2015, pp. 4775-4782.
- [6] G. A. Hollinger and S. Singh, "Multirobot coordination with periodic connectivity: Theory and experiments," *IEEE T. Robot.*, vol. 28, no. 4, pp. 967-973, Aug. 2012.
- [7] G. A. Hollinger, S. Yerramalli, S. Singh, U. Mitra, and G. S. Sukhatme, "Distributed data fusion for multirobot search," *IEEE T. Robot.*, vol. 31, no. 1, pp. 55-66, Feb. 2015.
- [8] M. Jain, M. Taylor, M. Tambe, and M. Yokoo, "DCOPs Meet the Real World: Exploring Unknown Reward Matrices with Applications to Mobile Sensor Networks," in *Intl. Joint Conf. on AI (IJCAI)*, 2009.
- [9] R. Zivan, H. Yedidsion, S. Okamoto, R. Grinton, and K. Sycara, "Distributed constraint optimization for teams of mobile sensing agents," *Auton. Agent. Multi-Ag.*, vol. 29, no. 3, pp. 495-536, 2015.
- [10] J. Renoux, A. I. Mouaddib, and S. L. Gloanec, "A decision-theoretic planning approach for multi-robot exploration and event search," in *Proc. Intelligent Robots and Systems (IROS)*, Sept. 2015, pp. 5287-5293.
- [11] J. Capitan, M. T. Spaan, L. Merino, and A. Ollero, "Decentralized multi-robot cooperation with auctioned POMDPs," *Int. J. Robot. Res.*, vol. 32, no. 6, pp. 650-671, 2013.
- [12] R. Nair, P. Varakantham, M. Tambe, and M. Yokoo, "Networked distributed POMDPs: A synthesis of distributed constraint optimization and POMDPs," in *Proc. Natl. Conf. on AI (AAAI)*, 2005, pp. 133-139.
- [13] F. A. Oliehoek and C. Amato, *A Concise Introduction to Decentralized POMDPs*. Springer, 2016.
- [14] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein, "The Complexity of Decentralized Control of Markov Decision Processes," *Math. Oper. Res.*, vol. 27, no. 4, pp. 819-840, Nov. 2002.
- [15] S. Omidshafiei, A. a. Agha-mohammadi, C. Amato, and J. P. How, "Decentralized control of Partially Observable Markov Decision Processes using belief space macro-actions," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, May 2015, pp. 5962-5969.
- [16] S. Omidshafiei, A. a. Agha-mohammadi, C. Amato, S. Y. Liu, J. P. How, and J. Vian, "Graph-based Cross Entropy method for solving multi-robot decentralized POMDPs," in *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, May 2016, pp. 5395-5402.
- [17] N. Atanasov, B. Sankaran, J. Le Ny, G. Pappas, and K. Daniilidis, "Nonmyopic View Planning for Active Object Classification and Pose Estimation," *IEEE T. Robot.*, vol. 30, no. 5, pp. 1078-1090, 2014.
- [18] B. Charrow, V. Kumar, and N. Michael, "Approximate representations for multi-robot control policies that maximize mutual information," *Auton. Robot.*, vol. 37, no. 4, pp. 383-400, 2014.
- [19] V. Indelman, L. Carlone, and F. Dellaert, "Planning in the continuous domain: a generalized belief space approach for autonomous navigation in unknown environments," *Int. J. Robot. Res.*, vol. 34, no. 7, pp. 849-882, 2015.
- [20] M. Lauri and R. Ritala, "Planning for robotic exploration based on forward simulation," *Robot. Auton. Syst.*, vol. 83, pp. 15 - 31, 2016.
- [21] F. A. Oliehoek, M. T. J. Spaan, and N. Vlassis, "Optimal and approximate Q-value functions for decentralized POMDPs," *J. Artif. Intell. Res.*, vol. 32, no. 1, pp. 289-353, 2008.
- [22] M. Araya-López, O. Buffet, V. Thomas, and F. Charpillet, "A POMDP Extension with Belief-dependent Rewards," in *Advances in Neural Information Processing Systems 23*, 2010, pp. 64-72.
- [23] M. T. J. Spaan, T. S. Veiga, and P. U. Lima, "Active cooperative perception in network robot systems using POMDPs," in *Proc. IEEE/RSJ Conf. Intelligent Robots and Systems (IROS)*, Oct. 2010, pp. 4800-4805.
- [24] —, "Decision-theoretic planning under uncertainty with information rewards for active cooperative perception," *Auton. Agent. Multi-Ag.*, vol. 29, no. 6, pp. 1157-1185, 2015.
- [25] M. H. DeGroot, *Optimal Statistical Decisions*. Hoboken, New Jersey: John Wiley & Sons, Inc., 2004, Wiley Classics Library edition.
- [26] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Hoboken, NJ: John Wiley & Sons, Inc., 2006.
- [27] D. V. Pynadath and M. Tambe, "The communicative multiagent team decision problem: Analyzing teamwork theories and models," *J. Artif. Intell. Res.*, vol. 16, pp. 389-423, 2002.
- [28] D. Szer, F. Charpillet, and S. Zilberstein, "MAA*: A Heuristic Search Algorithm for Solving Decentralized POMDPs," in *Proc. 21st Conf. on Uncertainty in Artificial Intelligence (UAI)*, 2005, pp. 576-583.