

Distributed Online Planning for Min-Max Problems in Networked Markov Games

Alexandros E. Tzikas¹, Jinkyoo Park², Mykel J. Kochenderfer¹, and Ross E. Allen³

Abstract—Min-max problems are important in multi-agent sequential decision-making because they improve the performance of the worst-performing agent in the network. However, solving the multi-agent min-max problem is challenging. We propose a modular, distributed, online planning-based algorithm that is able to approximate the solution of the min-max objective in networked Markov games, assuming that the agents communicate within a network topology and the transition and reward functions are neighborhood-dependent. This set-up is encountered in the multi-robot setting. Our method consists of two phases at every planning step. In the first phase, each agent obtains sample returns based on its local reward function, by performing online planning. Using the samples from online planning, each agent constructs a concave approximation of its underlying local return as a function of only the action of its neighborhood at the next planning step. In the second phase, the agents deploy a distributed optimization framework that converges to the optimal immediate next action for each agent, based on the function approximations of the first phase. We demonstrate our algorithm’s performance through formation control simulations.

Index Terms—distributed robot systems, networked robots, planning under uncertainty, Markov games, min-max optimization

I. INTRODUCTION

MULTI-agent reinforcement learning (MARL) involves multiple independent agents that operate in a common environment, each aiming to optimize a long-term reward by interacting with the environment and the other agents [1]. It has recently witnessed widespread success in many problems, such as the game of Go [2], autonomous driving [3], and simulated soccer [4].

MARL techniques can be applied to a large class of multi-agent decision problems. In this work, we will focus on the class of Markov games. The goal of each agent in the Markov game is to determine an optimal policy that optimizes its expected cumulative reward.

The Markov game can be extended to the distributed paradigm. In this setting, we assume that the agents communicate within a network topology. The agents can exchange information only with their neighbors, as defined by a communication graph. This paradigm is critical in applications where there is no central controller communicating with all agents. The absence of a central controller offers various benefits: i) the system does not have a single point of failure, ii) the distributed computation allows for deployment when the agents have limited communication capabilities or the number of agents is large. Usually, the goal of the agents in this setting is to maximize the expected cumulative team-average reward, by determining their optimal policies under the imposed communication restrictions [1].

Although other learning objectives have been explored in the literature for Markov games in the distributed setting [1], to the best of the authors’ knowledge, limited prior work [5] has investigated the problem where the network of agents aims to jointly maximize the worst-performing agent’s expected cumulative reward (termed the max-min or min-max problem from here onward) [1, pg. 19 & 24]. This is an important problem when fairness is required in the system. Most prior works consider a shared objective among the agents [6], a team-average reward [7], and a zero-sum set-up of agent objectives [8].

Various notions of fairness have been proposed in the literature, but in this work we focus on the notion of improving the performance of the worst-performing agent, which is known as the egalitarian objective [9], and has applicability in multi-robot systems [10]. Our contributions are as follows.

- We propose a distributed algorithm for the sequential decision-making problem of a multi-agent system with the egalitarian objective. Our algorithm consists of an online planning and a distributed optimization module, and solves the max-min Markov game problem under assumptions. We are able to combine online planning and distributed optimization by creating convex function approximations of the expected cumulative costs of the agents at each timestep. We assume the reward of each agent is dependent on the actions of its neighborhood and the portion of the state relevant to its neighborhood. We also assume that the transition probability for the portion of the state relevant to its neighborhood is only dependent on its neighborhood quantities.
- We demonstrate the performance of our algorithm on the formation control application. We pose the problem in min-max form and solve it through disciplined

¹A. E. Tzikas (corresponding author) and M. J. Kochenderfer are with the Department of Aeronautics and Astronautics, Stanford University, Stanford, CA 94305, U.S.A. {alexztzik, mykel}@stanford.edu

²J. Park is with the Department of Industrial and Systems Engineering, KAIST, Yuseong-gu, Daejeon 305-701, Republic of Korea. jinkyoo.park@kaist.ac.kr

³R. E. Allen is with the MIT Lincoln Laboratory, Lexington, MA 02421-6426, U.S.A. ross.allen@ll.mit.edu

Accepted to appear in the IEEE Robotics and Automation Letters.

© 2024 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

convex programming [11]. Our algorithm outperforms other baselines in this problem and its performance is comparable to the optimal action sequence.

This paper is organized as follows. In section II, we describe the prior related work. In section III, we include the problem statement, while in section IV we provide a qualitative description of our method. In section V we describe our algorithm and in section VI we test its performance in simulation. We conclude the paper in section VII.

II. RELATED WORK

A. Solving the Multi-Agent Min-Max Problem

Prior work has attempted to address the min-max multi-agent sequential decision-making problem, but there are limitations to the existing methods. A decentralized deep learning-based model to extract policies in this setting has been proposed [9]. The learned policy of every agent is only a function of the agent’s current observation, which means that inter-agent communication is not utilized in the decision-making step to allow for more informed decisions. The algorithm presented in this prior work also uses a minibatch of data at every timestep, which requires the simulation of the whole system. This is not always possible in a networked multi-agent system, because the agents can only communicate with their neighbors. In addition, the authors studied the problem with the assumption that the reward of each agent is not dependent on the actions of other agents: at each agent an advantage function, which is only a function of the agent’s quantities, is used. A variant of the problem has also been studied in the centralized setting, where the objective is to minimize a weighted sum of the worst agent cumulative cost and the cumulative cost of all agents [5]. This approach is limited: it requires a centralized controller and only solves a problem related to the min-max problem.

In the case of static distributed optimization, there exist algorithms that solve the static multi-agent max-min problem [12]. It has also been shown that the one-step min-max problem is equivalent, under specific conditions, to a min-sum optimization problem [12]. Such problems are solvable using many distributed optimization algorithms [13]–[15].

B. Online Planning Methods

In single-agent Markov decision problems, online methods are often used [16] when the optimal policy of the underlying problem is computationally expensive (or even intractable) to determine offline. Online methods choose the agent’s next action based on reasoning about states reachable from the current state and the future cumulative rewards that can potentially be received. Many variants of online planning exist in the literature. An algorithm called partially observable Monte Carlo planning with observation widening (POMCPOW) can handle continuous action and observation spaces by using weighted particle filtering in the Monte Carlo tree search (MCTS) [17].

Solving the Markov game analytically is difficult, even as a centralized problem. A major challenge is the computational complexity. Online planning has been explored in the context

of centralized MARL with a cooperative objective. The idea of factored-value MCTS has been combined with the max-plus algorithm to obtain an anytime online planning algorithm that is computationally efficient because it takes advantage of the factored structure of the value function [18].

When we want to solve the Markov game in a distributed manner, the problem difficulty increases. An iterative, distributed solution method is then required that learns the optimal policies as the agents continue to interact with the environment. *By integrating MCTS into our multi-agent algorithm for the max-min problem, we obtain an anytime method that can provide action suggestions at the current timestep, without the need to solve the entire problem.*

Decentralized online planning algorithms have been proposed in the context of the Markov game [19], [20]. Czechowski and Oliehoek [20] introduce a decentralized online planning approach for the Markov game in the cooperative case, where each agent shares a common reward. Their approach is based on MCTS and allows only one agent to improve its policy at a time, using learned models for the other agents’ policies. The problem we are considering differs from this application in three significant ways: our problem objective is not equivalent to the collaborative objective explored in this prior work, we do not require cyclic communication among the agents for the action (policy) selection, and we do not require simulations with the current joint policy in order to model the agents’ behavior. In addition, a decentralized planning method based on MCTS has been proposed for scenarios where each agent makes sequential decisions in order to optimize a global objective that depends on the action sequences of all agents [19]. This algorithm cannot solve our problem of interest because it assumes that the global objective is known to all agents. In the max-min problem we are considering, each agent only knows its local reward function.

III. PROBLEM STATEMENT

In our formulation, we consider the framework of the Markov game, which is defined below.

Definition 1: A Markov game is defined by the tuple

$$(\mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, P, \{R^i\}_{i \in \mathcal{N}}, \gamma), \quad (1)$$

where $\mathcal{N} = \{1, \dots, N\}$ is the set agents, \mathcal{S} is the state space observed by all agents, \mathcal{A}^i is the action space of agent i . Let $\mathcal{A} = \times_{i \in \mathcal{N}} \mathcal{A}^i$ be the set of joint agent actions, where \times denotes the Cartesian product. Then $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ defines the transition probabilities by mapping each state-joint action pair to a distribution on the state space, which belongs to the set $\Delta(\mathcal{S})$. The reward function $R^i : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ determines the immediate reward received by agent i for a transition from (s, \mathbf{a}) to s' . The discount factor is $\gamma \in [0, 1)$.

Our objective is to maximize the worst-performing agent’s expected cumulative reward, under constrained communication at every timestep. The constrained communication is imposed by an undirected graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where \mathcal{E} is the set of communication links, i.e., $(i, j) \in \mathcal{E}$ if and only

if node i can communicate directly with node j . We assume that at every timestep agent i can send information to agent j if and only if agent j can send information to agent i , i.e. $(i, j) \in \mathcal{E} \iff (j, i) \in \mathcal{E}$. The neighborhood of agent i is denoted $\mathcal{N}_i = \{j \mid (i, j) \in \mathcal{E}\} \cup \{i\}$. The environment and the interaction of the agents with the environment are expressed through the tuple $(\mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, P, \{R^i\}_{i \in \mathcal{N}}, \gamma)$. Out of the entities in the tuple, we assume that \mathcal{A} , P , and γ are known to all agents, but R^i is only known to agent i . We further impose the restriction:

$$R^i(\mathbf{s}, \mathbf{a}, \mathbf{s}') = R^i(\mathbf{s}^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i}, \mathbf{s}'^{\mathcal{N}_i}), \forall i \in \mathcal{N}. \quad (2)$$

Eq. (2) states that the reward of agent i is only a function of the actions of the agents in the neighborhood of agent i , $\mathbf{a}^{\mathcal{N}_i} \in \mathcal{A}^{\mathcal{N}_i} = \times_{i \in \mathcal{N}_i} \mathcal{A}^i$, and of the portion of the state which is relevant to the neighborhood of agent i , $\mathbf{s}^{\mathcal{N}_i}$ and $\mathbf{s}'^{\mathcal{N}_i}$. $\mathbf{a}^{\mathcal{N}_i}$ and $\mathbf{s}^{\mathcal{N}_i}$ are subvectors of the vectors \mathbf{a} and \mathbf{s} , respectively. Eq. (2) is an extension to the common assumption that the reward of agent i is dependent only on the actions of its neighbors [21]. We note that nevertheless, every agent has access to the complete current state \mathbf{s} , because we do not deal with partial observability in this work.

We also assume that

$$P(\mathbf{s}'^{\mathcal{N}_i} \mid \mathbf{s}, \mathbf{a}) = P(\mathbf{s}'^{\mathcal{N}_i} \mid \mathbf{s}^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i}). \quad (3)$$

The portion of the state relevant to the neighborhood of agent i is conditionally dependent only on the current timestep's quantities of the neighborhood. This can be encountered in multi-robot systems where the state is the concatenation of the state of each robot and the robots have independent dynamics.

We finally introduce a new entity for every agent i , $I_t^{\mathcal{N}_i} \in \mathcal{I}_t^{\mathcal{N}_i}$, where $\mathcal{I}_t^{\mathcal{N}_i}$ is the set of possible pieces of information available to agent i from its neighborhood at timestep t . It could for example include the action sequence of agent i 's neighbors up to the previous timestep or incorporate knowledge about the neighbors' neighborhoods from previous timesteps, which are propagated to agent i through the network. Note that $I_t^{\mathcal{N}_i} \in \mathcal{I}_t^{\mathcal{N}_i}$ can only depend on information available at or before timestep t . Assuming an initial state \mathbf{s}_0 , the expected cumulative reward of agent i is:

$$\tilde{R}^i(\mathbf{s}_0) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t R^i(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \mid a_t^i \sim \pi^i(\cdot \mid \mathbf{s}_t, I_t^{\mathcal{N}_i}), \mathbf{s}_0 \right], \quad (4)$$

where the subscript t refers to the timestep and $\pi^i(\cdot \mid \cdot)$ denotes the policy of agent i . The optimization problem that we aim to solve in a distributed manner is formulated as follows:

$$\begin{aligned} \min_{\pi^i(\cdot \mid \cdot), \forall i \in \mathcal{N}} \quad & \max_{i \in \mathcal{N}} -\tilde{R}^i(\mathbf{s}_0) \\ \text{s.t.} \quad & \pi^i : \mathcal{S} \times \mathcal{I}_t^{\mathcal{N}_i} \rightarrow \Delta(\mathcal{A}^i), \forall i \in \mathcal{N}. \end{aligned} \quad (5)$$

We will call problem (5) the main problem (min-max or max-min interchangeably) from now on. We can impose

additional restrictions on the policies, such as requiring that they be deterministic.

IV. MOTIVATING THE PROPOSED APPROACH

We highlight the challenges in solving the multi-agent min-max problem using online planning by showing how it differs from the single-agent case. With this analysis, our approach to solving the problem is motivated.

A. Online Planning Module

This section highlights the differences between performing online planning in the single-agent min-max problem and the multi-agent min-max case. In the single-agent setting, the min-max criterion simply becomes the min criterion. The subsection on the single-agent case describes already known facts from the literature. The subsection on the multi-agent case motivates our novel algorithm by following a parallel structure to the previous subsection. It is an informal description of the reasoning behind our proposed algorithm.

1) *The Single-Agent Case:* Much of the success of online planning methods in single-agent reinforcement learning can be attributed to the Bellman equation [16], which pertains to the underlying framework, i.e., the Markov decision process (MDP). An MDP is the degenerate case of problem (5) for $\mathcal{N} = \{1\}$. According to the Bellman equation, the optimal solution to an MDP satisfies the principle of dynamic programming. Also, if an optimal policy exists for an MDP, a deterministic optimal policy exists. Namely, an optimal action $a_{\mathbf{s}}$, when in state \mathbf{s} , in order to maximize the expected cumulative reward, is

$$a_{\mathbf{s}} = \operatorname{argmax}_{a \in \mathcal{A}^1} Q^*(\mathbf{s}, a) = \operatorname{argmax}_{a \in \mathcal{A}^1} R(\mathbf{s}, a) + \gamma \sum_{\mathbf{s}' \in \mathcal{S}} P(\mathbf{s}' \mid \mathbf{s}, a) V^*(\mathbf{s}'), \quad (6)$$

where $V^*(\mathbf{s}')$ denotes the optimal cumulative reward when starting at state \mathbf{s}' and following an optimal policy for the MDP, $Q^*(\mathbf{s}, a)$ denotes the return when taking action a at \mathbf{s} and then following an optimal policy, and $R(\mathbf{s}, a)$ is the expected reward when transitioning from \mathbf{s} with action a . The solution to an instance of the MDP is a sequence of actions (obtained via an optimal policy) that maximizes the expected cumulative reward.

The online planning methods attempt to obtain an action sequence with a return realization close to the optimal one. They work as follows. At the first timestep, we assume that the state is \mathbf{s}_0 . Then, as a surrogate to eq. (6), the online planning method creates a search tree with root node \mathbf{s}_0 and first-level branches pertaining to L sampled actions $\{\bar{a}^{(1)}, \dots, \bar{a}^{(L)}\}$ that are to be taken at \mathbf{s}_0 . We assume that for these, the online planning method is able to obtain a good approximation of the value $Q^*(\mathbf{s}_0, \bar{a}^{(l)})$, which we denote $\hat{Q}_{\mathbf{s}_0, \bar{a}^{(l)}}$, by simulating future action trajectories. Then, the algorithm chooses as $a_{\mathbf{s}_0}$ (the action that will actually be taken at \mathbf{s}_0) the action $\bar{a}^{(\operatorname{argmax}_{1 \leq l \leq L} \hat{Q}_{\mathbf{s}_0, \bar{a}^{(l)}})}$. The agent performs $a_{\mathbf{s}_0}$, transitioning to state \mathbf{s}_1 at the next timestep, and receiving reward r_1 . At this time, the goal remains

to optimize the cumulative reward since the first timestep. Hence, the agent should find a_{s_1} according to

$$a_{s_1} = \operatorname{argmax}_{a \in \mathcal{A}^1} r_1 + \gamma Q^*(s_1, a) = \operatorname{argmax}_{a \in \mathcal{A}^1} Q^*(s_1, a), \quad (7)$$

and we see that we arrive at the same optimization problem we would get if we directly applied Bellman's equation at s_1 , i.e., the past sequence of actions and rewards is irrelevant to the current action selection. Hence, the same online planning procedure as above is followed for every timestep.

The main question is how can the online planning method obtain good estimates of $Q^*(s, a)$, which are denoted $\tilde{Q}_{s,a}$, when $Q^*(s, a)$ internally assumes the agent follows the optimal policy. Online methods, such as MCTS, update the value or candidacy for optimality of an action a at a given state s in the search tree using sampled trajectory paths and their cumulative rewards starting at that state-action pair. The value of a at s is determined by the empirical average return of paths starting at this state-action pair, when at later states in the path, s' , the action to be followed is chosen using the upper confidence bound (UCB1) exploration heuristic [16]. UCB1 uses the approximate Q values, $\tilde{Q}_{s',\cdot}$, computed so far in the search tree for node s' , and selects as next action the one that maximizes partly the $\tilde{Q}_{s',\cdot}$ value and partly an exploration bonus. Based on the received returns, MCTS updates the estimate $\tilde{Q}_{s,a}$. This indicates that the value of an action at state node s is determined with respect to a proxy to the best policy (since maximizing actions are taken in the path). Therefore, the online planning method can be seen as an approximate policy iteration algorithm [16], where approximate Q values are computed based on simulated trajectories (proxy of policy evaluation), and then actions are chosen in internal levels based on the UCB1 metric (which is a proxy for policy improvement), to finally find a good estimate of Q^* . Hence, single-agent online planning methods produce trajectory samples using action sequences that follow the basic rule of optimality for MDPs, eq. (6), which in turn follows the principle of dynamic programming. This explains why $\tilde{Q}_{s,\bar{a}^{(l)}}$ is a good approximation of $Q^*(s, a^{(l)})$ with the UCB1 metric.

2) *The Multi-Agent Case:* In order to solve the main problem (5) online, we follow the reasoning in the previous subsection. The goal is to maximize the worst agent's expected cumulative reward since the beginning of time. At the first timestep, the state is s_0 and we use online planning to create the functions $\hat{Q}_0^i(s_0, \mathbf{a}_0^{\mathcal{N}^i})$ at every agent i . These functions are to be good approximations of the local return functions at the optimal solution of problem (5). They should express the local cumulative return, assuming that the agents follow the optimal policies of problem (5) after timestep 0 and at timestep 0 take \mathbf{a}_0 . Similarly, functions $\hat{Q}_t^i(s_t, \mathbf{a}_t^{\mathcal{N}^i})$ should express the local cumulative return starting at t , assuming that the agents follow the optimal policies of problem (5) before and after timestep t , and at timestep t take \mathbf{a}_t . Then, the problem

$$\max_{\mathbf{a}_0 \in \times_{i \in \mathcal{N}} \mathcal{A}^i} \min_{i \in \mathcal{N}} \hat{Q}_0^i(s_0, \mathbf{a}_0^{\mathcal{N}^i}) \quad (8)$$

serves as an approximation of problem (5) at the initial timestep. The agents solve problem (8) in a distributed manner and take the action computed as the maximizer. Then each agent i receives a local reward r_1^i and the state becomes s_1 at timestep 1. At this timestep, the agents will again apply online planning to continue solving the original problem (5). They now create new function approximations $\gamma \hat{Q}_1^i(s_1, \mathbf{a}_1^{\mathcal{N}^i}) + r_1^i$, which should express the local cumulative rewards assuming that the agents follow the optimal policies of problem (5) before and after timestep 1 and at timestep 1 take \mathbf{a}_1 . The problem

$$\max_{\mathbf{a}_1 \in \times_{i \in \mathcal{N}} \mathcal{A}^i} \min_{i \in \mathcal{N}} r_1^i + \gamma \hat{Q}_1^i(s_1, \mathbf{a}_1^{\mathcal{N}^i}) \quad (9)$$

serves as the new proxy of problem (5) in timestep 1 and is again solved using a distributed algorithm. The same process continues at every timestep.

For this method to work, the online planning approach needs to obtain function approximations $\hat{Q}_t^i(\cdot, \cdot)$, which are good approximations of the returns of the agents when following the optimal policies of problem (5). In our algorithm, we assume that, at every timestep t , each agent i creates a search tree for its own local reward, with root node s_t , which is used to compute $\hat{Q}_t^i(s_t, \mathbf{a}_t^{\mathcal{N}^i})$. The function $\hat{Q}_t^i(s_t, \mathbf{a}_t^{\mathcal{N}^i})$ is a concave function approximation of the values $\hat{Q}_{t,s_t,\bar{\mathbf{a}}^{\mathcal{N}^i}^{(l)}}^i$ for the sampled neighborhood action tuples $\bar{\mathbf{a}}^{\mathcal{N}^i}^{(l)}$ at the first level of the tree search of agent i at timestep t , which has root node s_t . In this case, using the UCB1 criterion for the expansion of the search tree of each agent need not be the approach that provides the best approximation of performance, $\hat{Q}_t^i(s_t, \mathbf{a}_t^{\mathcal{N}^i})$, for the optimal policy of problem (5). The reasons for this are twofold: the principle of dynamic programming does not necessarily hold for the solution of problem (5), as shown in the following theorem, and the UCB1 in the search tree is performed only with respect to the local reward. Nevertheless, UCB1 is still used and future work will investigate other criteria.

Theorem 1: For a problem of the form (5), the principle of dynamic programming does not necessarily hold.

Proof: The proof is based on a simple counter-example, where a sequence of actions that are non-optimal for their respective min-max subproblems, constitute an optimal action sequence for the min-max problem involving all timesteps. The details of the proof are included in the appendix. ■

Overall, there are two significant challenges in applying online planning to solve the multi-agent min-max problem that do not appear in the single-agent case. In the multi-agent min-max scenario it is important to keep in memory the cumulative reward received up to the current timestep for every agent. This is because the dynamic programming principle does not hold. In addition, while decisions can be made with $\tilde{Q}_{\cdot,\cdot}$ values for sampled actions in the single agent setting, we require a function approximation $\hat{Q}_t^i(\cdot, \cdot)$ in the multi-agent scenario, because of the different trajectories sampled by each agent performing online planning.

B. Distributed Optimization Module

As mentioned in the previous section, when solving the general version of problem (5), each agent performs online planning to compute an approximation of its local cumulative reward and then the agents determine their current actions by solving the approximate max-min problem at the current timestep. This is done by deploying a distributed optimization framework. In our algorithm, we choose to use the method presented by Srivastava et al. [12]. This is a distributed stochastic subgradient method. We decide to use a subgradient method, because it is generally applicable: any convex function has a nonempty and bounded subdifferential in the interior of its domain. The distributed optimization framework by Srivastava et al. [12] solves the problem

$$\min_{\mathbf{a} \in \times_{i \in \mathcal{N}} \mathcal{A}^i} \max_{i \in \mathcal{N}} f^i(\mathbf{a}), \quad (10)$$

where each f^i is convex and known only to agent i . The closed and convex set $\times_{i \in \mathcal{N}} \mathcal{A}^i$ ($\mathbf{a} \in \times_{i \in \mathcal{N}} \mathcal{A}^i$) is known to all agents. We now look at the method from the perspective of agent i . The estimate of the optimal point and optimal value of problem (10) of agent i , at iteration k , are denoted $\tilde{\alpha}_k^i$ and $\tilde{\eta}_k^i$, respectively. The method proceeds with two steps at every iteration $k \geq 0$. At first, it forms an intermittent adjustment, and then agent i takes a step towards minimizing its own function f^i . This is written:

$$\begin{bmatrix} \tilde{\alpha}_k^i \\ \tilde{\eta}_k^i \end{bmatrix} = \sum_{j \in \mathcal{N}_i} w_k^{ij} \begin{bmatrix} \alpha_k^j \\ \eta_k^j \end{bmatrix}, \quad \mathbf{v}_k^i = \begin{bmatrix} \tilde{\alpha}_k^i \\ \tilde{\eta}_k^i \end{bmatrix} - \frac{\beta_k}{N} \begin{bmatrix} \mathbf{0}_n \\ 1 \end{bmatrix}, \quad (11)$$

$$\begin{bmatrix} \alpha_{k+1}^i \\ \eta_{k+1}^i \end{bmatrix} = \Pi_{\mathcal{A} \times \mathbb{R}} \left[\mathbf{v}_k^i - \beta_k r^i \left(g^i + \begin{bmatrix} \epsilon_k^i \\ 0 \end{bmatrix} \right) \right]. \quad (12)$$

In the equations above, w_k^{ij} denotes the weight that agent i assigns to its neighbor j at timestep k , β_k is the step size, ϵ_k^i the random subgradient error, $r^i > 1$, g^i a subgradient of $\tilde{g}^i(\boldsymbol{\alpha}, \eta) = \tilde{g}^i(\mathbf{z}) = \max\{0, f^i(\boldsymbol{\alpha}) - \eta\}$ at \mathbf{v}_k^i , and n denotes the dimensionality of $\mathbf{a} \in \times_{i \in \mathcal{N}} \mathcal{A}^i$. Here, $\Pi_{\mathcal{A} \times \mathbb{R}}$ denotes the Euclidean projection on the set $\mathcal{A} \times \mathbb{R}$.

Under noise, subgradient, communication topology, and assigned weights assumptions, it is guaranteed that the iterates $\mathbf{z}_k^i = (\alpha_k^i, \eta_k^i)$ converge to a common $\mathbf{z}^* = (\alpha^*, \eta^*)$ such that α^*, η^* are an optimal point and the optimal value of problem (10), respectively [12].

V. PROPOSED ALGORITHM

Our proposed algorithm is Algorithm 1. At every planning step, it passes through two phases. In the first phase, each agent i creates, through online planning, a convex function approximation of its local expected cumulative cost, which is a function of only the immediate next action of the agents in its neighborhood. In the current instance of the algorithm, the online planning module uses the POMCPOW algorithm [17] (with the observation equaling the state). The convex function approximator used is the adaptive max-affine partitioning algorithm (AMAP) [22], which offers a balance between the regression speed and model quality. The steps mentioned above correspond to lines 5–8 in Algorithm

1. In the second phase, the subgradient-based distributed optimization algorithm [12] is deployed to solve the min-max problem involving the convex function approximations from the first phase and determine the agents' next action (lines 9–12). In Algorithm 1, we give our proposed algorithm from the perspective of agent i . Because of properties (2–3), agent i only needs to create a search tree with actions of its neighbors and the portion of the state relative to its neighborhood.

Algorithm 1 Proposed algorithm from the perspective of agent i at timestep t

- 1: **Inputs:**
 - 2: planning timestep t
 - 3: current environment state \mathbf{s}_t
 - 4: cumulative reward up to timestep $t - 1$: \bar{R}_{t-1}^i
 - 5: obtain L samples of future expected cumulative reward from state $\mathbf{s}_t^{\mathcal{N}_i}$: using online planning (e.g., MCTS), sample an action tuple $\mathbf{a}^{\mathcal{N}_i, (l)}$ and compute $\hat{Q}_{t, \mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i, (l)}}^i$, $l \in 1, \dots, L$
 - 6: negate the sampled cumulative rewards to get future cumulative cost samples $-\hat{Q}_{t, \mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i, (l)}}^i$
 - 7: create a convex function approximation $-\hat{Q}_t^i(\mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i})$ based on the samples $-\hat{Q}_{t, \mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i, (l)}}^i$
 - 8: add $-\bar{R}_{t-1}^i$ to $-\gamma^t \hat{Q}_t^i(\mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i})$ and get $f_t^i(\mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i})$
 - 9: initialize the estimates $\alpha_0^i \in \mathcal{A}$, $\eta_0^i \in \mathbb{R}$
 - 10: **for** $k = 0, \dots, K$ **do**
 - 11: $\alpha_{k+1}^i, \eta_{k+1}^i \leftarrow$ one step of algorithm (12)
 - 12: perform the action pertaining to agent i in α_{K+1}^i
 - 13: collect immediate reward r_t^i and $\bar{R}_t^i \leftarrow \bar{R}_{t-1}^i + \gamma^t r_t^i$
-

The number of samples L controls the accuracy of the cost function approximation. The number of iterations K determines the consensus among the agent iterates and their distance to the minimum of the optimization problem. The agent iterates determine the next action for the agents.

A. Computational Complexity

The per-agent computational complexity of the proposed algorithm is completely determined by the three submodules included in Algorithm 1: the online planning module (POMCPOW) in line 5, the convex function approximator [22] in line 7, and the distributed optimization framework in lines 10–11. In Algorithm 1, POMCPOW uses L planning calls. With a max tree depth Υ , it has complexity $\mathcal{O}(LY \log L)$ [17]. The AMAP model improvement step has complexity $\mathcal{O}(\max\{\phi_{\text{LSPA}}, D_i\} \max\{H, D_i\} D_i L)$, when H hyperplanes are used, D_i is the argument dimensionality of $-\hat{Q}_t^i(\mathbf{s}_t^{\mathcal{N}_i}, \cdot)$, and ϕ_{LSPA} denotes the number of iterations of a sub-algorithm [22]. Assuming an ensemble of M models and multiple model improvements, the AMAP total complexity is $\mathcal{O}(L^{D_i/D_i+4} M \max\{\phi_{\text{LSPA}}, D_i\} \max\{H, D_i\} D_i L + M L D_i^2)$. The second term comes from the initial least-squares problems. Finally K iterations of the distributed optimization algorithm are performed, with the complexity

of each dominated by the complexity of eq. (12). The complexity of eq. (12) can greatly vary depending on the structure of \mathcal{A} , but certain cases allow for analytic solutions [23]. Therefore we suppose a generic complexity for eq. (12), $\mathcal{O}(\Pi)$. In total, our algorithm's complexity is the sum of the above complexities. We however note that the method used for each of the submodules of Algorithm 1 can be independently replaced, e.g., online planning can be performed with an algorithm other than POMCPOW.

VI. EXPERIMENTS

In order to evaluate the performance of our algorithm, we pose the formation control problem as a min-max problem. In this setting, the agents aim to converge to states such that desired relative states are satisfied between them [24].¹

A. Experiment Setup

We consider a Markov game where the state consists of the x, y positions of the agents, i.e., $\mathbf{s} = (s^1, \dots, s^N) \in \mathcal{S} = \times_{i \in \mathcal{N}} \mathbb{R}^2$, where $s^i = (x^i, y^i)$ the position of agent i . The dynamics of each agent i , from timestep t to $t + 1$, with action a_t^i , evolve according to

$$s_{t+1}^i = s_t^i + \begin{bmatrix} 1 & 0 \\ -1 & 2 \end{bmatrix} a_t^i, \quad a_t^i = (a, b), \quad a, b \in [0, 1]. \quad (13)$$

We assume deterministic dynamics, in order to allow comparisons with various techniques. The instantaneous reward for agent i is

$$R^i(\mathbf{s}^{\mathcal{N}_i}, \mathbf{a}^{\mathcal{N}_i}, \mathbf{s}'^{\mathcal{N}_i}) = \sum_{j \neq l \in \mathcal{N}_i} \|(s^j - s^l) - (d^j - d^l)\|, \quad (14)$$

where d^j denotes the desired position of agent j in a specific formation that satisfies the relative states. Our goal is to solve problem (5), which also leads the agents to satisfy their desired relative states. A reward close to 0 signifies satisfaction of the desired relative states. We assume $\gamma = 1$ and a finite horizon of $T = 150$.

We consider three graph topologies $\mathcal{G}_1, \mathcal{G}_2$, and \mathcal{G}_3 . In \mathcal{G}_1 , five agent communicate over an almost fully connected graph (only agent 1 is not a neighbor of agent 5). In \mathcal{G}_2 , five agent have the topology 1-2-3-4-5-1. In \mathcal{G}_3 , eight agents communicate over the topology 1-2-3-4-5-6-7-8.

B. Implementation Details for the Proposed Algorithm

Our proposed approach has the following parameters in Algorithm 1: $L = 100$, and $K = 1000$. For the online framework we deploy POMCPOW [17], since both our action and state spaces are continuous. Every agent i has access to the positions of its neighbors, $\mathbf{s}^{\mathcal{N}_i}$. Internally in POMCPOW of agent i , the observation equals the next simulated state for the agents in \mathcal{N}_i . For the POMCPOW implementation at each agent, we use the action and observation progressive widening values of $k_a = 2.0$, $\alpha_a = 0.5$, $k_o = 0.0$, and $\alpha_o = 0.5$ in Listing (1) of [17]. It also holds that the n of

¹The source code for the experiments can be found at: https://github.com/alextzik/distr_online_maxmin_markov_game.

Listing (1) in [17] equals L in our case. The search depth is 5. The rollout planner chosen for our algorithm is based on the affine formation control system

$$\frac{ds^i(t)}{dt} = - \sum_{j \in \mathcal{N}_i} (s^i(t) - s^j(t)) - (d^i - d^j), \quad \forall i, \quad (15)$$

for which it is guaranteed that $s^i(t) \rightarrow d^i + \xi$ as $t \rightarrow \infty$, where ξ is a constant displacement vector [24]. Therefore, if the rollout policy is called at a given depth in the search tree of an agent, it first solves the system of differential equations (15) with the state of the current node in the search tree as the initial condition and then chooses for the agents the action that takes them closer to the converged states of eq. (15), in the sense of a one-step look-ahead, which is an easily solvable convex problem. Note that the rollout policy is used within the POMCPOW of each agent r . Hence, in POMCPOW of agent r , eq. (15) only involves the agents i such that $i \in \mathcal{N}_r$, and agent r can assume a fully connected topology in the rollout planner. After simulating the actions for the agents, the search moves to a new node, if the maximum depth has not been reached. Note that different initializations of eq. (15) lead to different converged positions, which means that simply using the rollout planner should not solve our problem. In addition, the converged state for agent i can be different in the individual systems solved by the agents $j \in \mathcal{N}_i$.

C. Baseline Algorithm Description

Using the rollout planner at every timestep with depth 1 constitutes a simple baseline. The agents solve eq. (15), for the true graph topology and involving every $i \in \mathcal{N}$, in a distributed manner and then each agent takes the action that in the one step look-ahead notion brings it closer to its converged state for eq. (15). This is the ‘rollout baseline’ below. Assuming Euler discretization with step Δt and total time T_f , each timestep of the rollout baseline at agent i is $\mathcal{O}(\frac{T_f}{\Delta t} |\mathcal{N}_i|)$. The second baseline is the ‘POMCPOW baseline’, where each agent uses Algorithm 1 up to line 5, with complexity $\mathcal{O}(L \Upsilon \log L)$. It then follows the action pertaining to itself from the sample, $\hat{Q}_{t, \mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}_t^{\mathcal{N}_i}, (t)}^i$ with the maximum value. Obviously, the proposed algorithm is the most computationally intensive among the compared methods, but with significantly better min-max behavior, as shown below.

In the scenario that we are considering, we can obtain an upper bound on performance, by finding the optimal point and value of problem (5). Since the proposed problem is deterministic, the following centralized convex problem is solved by the optimal open-loop sequence of actions for the agents

$$\min_{\mathbf{a}_0^i, \dots, \mathbf{a}_{T-1}^i, \mathbf{s}_0^i, \dots, \mathbf{s}_T^i, \forall i \in \mathcal{N}} \max_{i \in \mathcal{N}} - \sum_{t=0}^{T-1} R^i(\mathbf{s}_t^{\mathcal{N}_i}, \mathbf{a}_t^{\mathcal{N}_i}, \mathbf{s}_{t+1}^{\mathcal{N}_i}). \quad (16)$$

Problem (16) serves as the main comparison to our work. It is referred to ‘optimal’ below. We have chosen this simplified

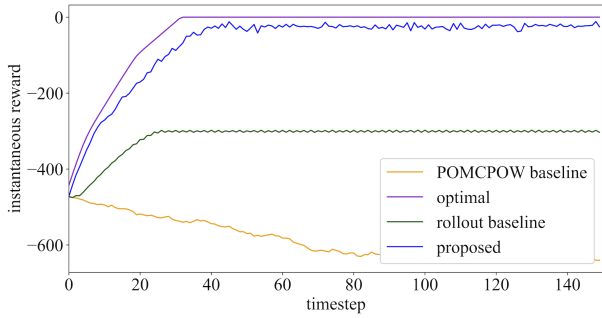


Fig. 1: Instantaneous reward of the worst-performing agent for \mathcal{G}_1 . Five agents communicate over an almost fully connected communication network (only agent 1 is not a neighbor of agent 5). We observe that our proposed method performs better than the baselines on the max-min criterion.

application in order to make the upper bound of performance tractable.

D. Performance Evaluation

We present results for the instantaneous reward of the worst-performing agent (in terms of return, which is the max-min objective) as a function of planning step in Figures 1–3. In Figure 1, five agents communicate over \mathcal{G}_1 for all timesteps. In Figure 2, five agents communicate over a switching network topology: every 10 timesteps, the topology switches between \mathcal{G}_1 and \mathcal{G}_2 . In Figure 3, eight agents communicate over the fixed topology \mathcal{G}_3 .

Our proposed algorithm performs better than the baselines by a wide margin. The POMCPOW baseline cannot achieve optimal performance, because it does not lead to coordination as explained above. The rollout baseline is also not able to converge to the desired configuration, because the converged states keep changing at the solution of eq. (15) at every timestep. The oscillatory behavior of reward for the two baselines in Figure 2 is attributed to the changing topology, which alters the number of terms in the reward sum. Nevertheless, the reward within the timesteps of constant topology is not increasing, which indicates a lack of coordination for the reasons explained above. Our algorithm is able to reach an ϵ -suboptimal configuration, but cannot exactly reach the optimal formation, because the agents’ function approximations, used in the min-max optimization problem to obtain the next actions, are not minimized by the zero action at later steps. This explains the observed oscillatory performance. We expect smaller oscillations as L and tree density in POMCPOW and H in AMAP increase. Our proposed algorithm is able to obtain convergence rates similar to the open-loop optimal action sequence, constrained however by the connectedness of the underlying communication topology.

VII. LIMITATIONS AND FUTURE WORK

Our approach is superior to the baselines presented, but has limitations. We observe that the online planning and

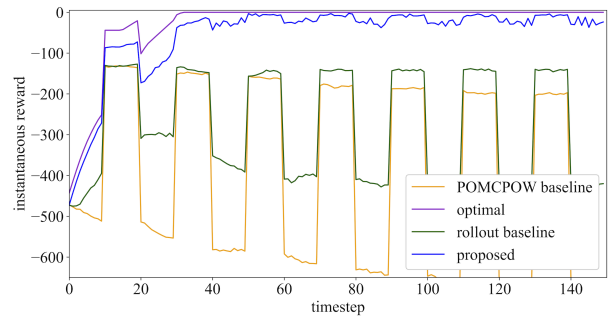


Fig. 2: Instantaneous reward of the worst-performing agent for a switching topology of five agents. Every 10 timesteps the topology changes between \mathcal{G}_1 , an almost fully connected topology, and \mathcal{G}_2 , which is a cyclic graph. We observe that our proposed method performs better than the baselines on the max-min criterion.

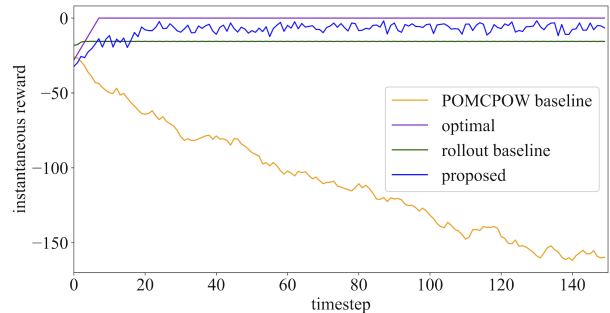


Fig. 3: Instantaneous reward of the worst-performing agent for \mathcal{G}_3 . Eight agents communicate over the topology 1-2-3-4-5-6-7-8. The worst agent cumulative reward in our method is -1188, while the worst agent return in the POMCPOW baseline is -2338. We observe that our proposed method performs better than the baselines with respect to the max-min criterion.

function approximator modules of our algorithm lead to oscillatory behavior. Tuning the parameters and exploring ways to mitigate this would be important before deployment in actual robots. Because computation is limited in robotic applications, an efficient implementation of the submodules of Algorithm 1 is also needed. In the current implementation, the online planning framework is the bottleneck in performance.

From a theoretical perspective, as discussed in Section IV-A.2, the UCB1 metric need not be the best exploration metric for the min-max problem. Theoretical analysis of the reasons behind its sub-optimality, and determination of a better metric are potential avenues of research. Our algorithm is also only approximate. Future work will aim towards a distributed algorithm that can provably converge to the optimal policies for the agents in the context of problem (5). Finally, we assume that every agent knows the complete feasible set \mathcal{A} for the actions of all agents beforehand.

Assume the case of 2 neighboring agents interacting over the course of 2 timesteps (final state at $t = 3$). Also assume $\mathcal{I}_t^{N_i}$ is the empty set and the dynamics are deterministic: $s_{t+1} = s_t + a_t^1 + a_t^2$. The initial state is $s_1 = 0$, $a_t^i \in \{1, 2\}$, and $\gamma = 1$. The allowed policies are deterministic.

At timestep 2 we can only have $s_2 \in \{2, 3, 4\}$. If:

$$\min_{i \in \{1, 2\}} R^i(2, (1, 1)) = \min\{90, 100\} > \quad (17)$$

$$\min_{i \in \{1, 2\}} R^i(2, \mathbf{a}), \quad \forall \mathbf{a} \neq (1, 1), \quad (18)$$

$$R^1(2, \mathbf{a}) = 200, R^2(2, \mathbf{a}) = 50, \forall \mathbf{a} \neq (1, 1), \quad (19)$$

then the optimal (with respect to the max-min criterion) policy for each agent in the subproblem starting at state 2 is action 1. Similarly, assume:

$$\min_{i \in \{1, 2\}} R^i(3, (1, 2)) = \min\{90, 100\} > \quad (20)$$

$$\min_{i \in \{1, 2\}} R^i(3, \mathbf{a}), \quad \forall \mathbf{a} \neq (1, 2), \quad (21)$$

$$R^1(3, \mathbf{a}) = 200, R^2(3, \mathbf{a}) = 50, \forall \mathbf{a} \neq (1, 2) \quad (22)$$

and

$$\min_{i \in \{1, 2\}} R^i(4, (2, 1)) = \min\{90, 100\} > \quad (23)$$

$$\min_{i \in \{1, 2\}} R^i(4, \mathbf{a}), \quad \forall \mathbf{a} \neq (2, 1), \quad (24)$$

$$R^1(4, \mathbf{a}) = 200, R^2(4, \mathbf{a}) = 50, \forall \mathbf{a} \neq (2, 1). \quad (25)$$

Now, assume $R^1(0, \mathbf{a}) = 5$, $R^2(0, \mathbf{a}) = 200$ for all \mathbf{a} . Then, to solve the min-max problem (5), the agents should choose any action when in state $s_1 = 0$ and for the state s_2 the system arrives at, they should choose any action other than the one optimal for the min-max subproblem starting at that state. In other words, suboptimal policies for the subproblems that could arise in timestep 2 have better (greater) overall max-min objective value in the complete problem than the optimal policies for these subproblems. This completes the proof.

ACKNOWLEDGMENTS

This material is based upon work supported by the Under Secretary of Defense for Research and Engineering under Air Force Contract No. FA8702-15-D-0001. Any opinions, findings, conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the Under Secretary of Defense for Research and Engineering.

Toyota Research Institute (TRI) also provided funds to assist the authors with this research, but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

For the first author, this work is also partially funded through the Alexander S. Onassis Foundation Scholarship program.

- [1] K. Zhang, Z. Yang, and T. Başar, “Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms,” *Handbook of Reinforcement Learning and Control*, 2021.
- [2] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, “Mastering the Game of Go with Deep Neural Networks and Tree Search,” *Nature*, 2016.
- [3] S. Shalev-Shwartz, S. Shammah, and A. Shashua, “Safe, Multi-Agent, Reinforcement Learning for Autonomous Driving,” *arXiv preprint arXiv:1610.03295*, 2016.
- [4] S. Liu, G. Lever, J. Merel, S. Tunyasuvunakool, N. Heess, and T. Graepel, “Emergent Coordination Through Competition,” in *International Conference on Learning Representations*, 2018.
- [5] C. Zhang and J. A. Shah, “Fairness in Multi-Agent Sequential Decision-Making,” in *Advances in Neural Information Processing Systems*, 2014.
- [6] M. Lauer and M. A. Riedmiller, “An Algorithm for Distributed Reinforcement Learning in Cooperative Multi-Agent Systems,” in *International Conference on Machine Learning*, 2000.
- [7] S. Kar, J. M. Moura, and H. V. Poor, “QD-Learning: A Collaborative Distributed Strategy for Multi-Agent Reinforcement Learning through Consensus+Innovations,” *IEEE Transactions on Signal Processing*, vol. 61, no. 7, pp. 1848–1862, 2013.
- [8] M. L. Littman, “Markov Games as a Framework for Multi-Agent Reinforcement Learning,” in *Machine Learning Proceedings*. Elsevier, 1994, pp. 157–163.
- [9] M. Zimmer, C. Glanois, U. Siddique, and P. Weng, “Learning Fair Policies in Decentralized Cooperative Multi-Agent Reinforcement Learning,” in *International Conference on Machine Learning*, 2021.
- [10] A. Kishimoto and N. Sturtevant, “Optimized Algorithms for Multi-Agent Routing,” in *International Joint Conference on Autonomous Agents and Multiagent Systems*, 2008.
- [11] M. Grant, S. Boyd, and Y. Ye, “Disciplined Convex Programming,” in *Global Optimization: From Theory to Implementation*. Springer, 2006, pp. 155–210.
- [12] K. Srivastava, A. Nedić, and D. Stipanović, “Distributed Min-Max Optimization in Networks,” in *International Conference on Digital Signal Processing*, 2011, pp. 1–8.
- [13] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers,” *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [14] O. Shorinwa, T. Halsted, J. Yu, and M. Schwager, “Distributed Optimization Methods for Multi-Robot Systems: Part 1—A Tutorial,” *IEEE Robotics & Automation Magazine*, 2024.
- [15] —, “Distributed Optimization Methods for Multi-Robot Systems: Part 2—A Survey,” *IEEE Robotics & Automation Magazine*, 2024.
- [16] M. J. Kochenderfer, T. A. Wheeler, and K. H. Wray, *Algorithms for Decision Making*. MIT Press, 2022.
- [17] Z. Sunberg and M. Kochenderfer, “Online Algorithms for POMDPs with Continuous State, Action, and Observation Spaces,” in *International Conference on Automated Planning and Scheduling*, vol. 28, 2018, pp. 259–263.
- [18] S. Choudhury, J. K. Gupta, P. Morales, and M. J. Kochenderfer, “Scalable Anytime Planning for Multi-Agent MDPs,” in *International Conference on Autonomous Agents and Multiagent Systems*, 2021, pp. 341–349.
- [19] G. Best, O. M. Cliff, T. Patten, R. R. Mettu, and R. Fitch, “DecMCTS: Decentralized Planning for Multi-Robot Active Perception,” *The International Journal of Robotics Research*, vol. 38, no. 2-3, pp. 316–337, 2019.
- [20] A. Czechowski and F. A. Oliehoek, “Decentralized MCTS via Learned Teammate Models,” in *International Joint Conference on Artificial Intelligence*, 2021, pp. 81–88.
- [21] Y. Yi, G. Li, Y. Wang, and Z. Lu, “Learning to Share in Multi-Agent Reinforcement Learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*.
- [22] G. Balázs, “Convex Regression: Theory, Practice, and Applications,” Ph.D. dissertation, University of Alberta, 2016.
- [23] S. P. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [24] M. Mesbahi and M. Egerstedt, “Graph Theoretic Methods in Multi-agent Networks,” in *Graph Theoretic Methods in Multiagent Networks*. Princeton University Press, 2010.