

Multi-Agent Coordination in Adversarial Environments through Signal Mediated Strategies

Federico Cacciamani
Politecnico di Milano
federico.cacciamani@polimi.it

Marco Ciccone
Politecnico di Milano
marco.ciccone@polimi.it

Andrea Celli
Politecnico di Milano
andrea.celli@polimi.it

Nicola Gatti
Politecnico di Milano
nicola.gatti@polimi.it

ABSTRACT

Many real-world scenarios involve teams of agents that have to coordinate their actions to reach a shared goal. We focus on the setting in which a team of agents faces an opponent in a zero-sum, imperfect-information game. Team members can coordinate their strategies before the beginning of the game, but are unable to communicate during the playing phase of the game. This is the case, for example, in Bridge, collusion in poker, and collusion in bidding. In this setting, model-free RL methods are oftentimes unable to capture coordination because agents' policies are executed in a decentralized fashion. Our first contribution is a game-theoretic centralized training regimen to effectively perform trajectory sampling so as to foster team coordination. When team members can observe each other actions, we show that this approach provably yields equilibrium strategies. Then, we introduce a signaling-based framework to represent team coordinated strategies given a buffer of past experiences. Each team member's policy is parametrized as a neural network whose output is conditioned on a suitable exogenous signal, drawn from a learned probability distribution. By combining these two elements, we empirically show convergence to coordinated equilibria in cases where previous state-of-the-art multi-agent RL algorithms did not.

KEYWORDS

Team Games; Multi-Agent Reinforcement Learning; Coordination

ACM Reference Format:

Federico Cacciamani, Andrea Celli, Marco Ciccone, and Nicola Gatti. 2021. Multi-Agent Coordination in Adversarial Environments through Signal Mediated Strategies. In *Proc. of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*, Online, May 3–7, 2021, IFAAMAS, 10 pages.

1 INTRODUCTION

In many strategic interactions agents have similar goals and have incentives to team up, and share their final reward. In these settings, coordination between team members plays a crucial role. We focus on *ex ante coordination*, where team members have an opportunity to discuss and agree on tactics before the game starts, but will be unable to communicate during the game, except through their publicly-observed actions. Consider, as an illustration, a poker game

where some players are colluding against some identified target players and will share the final winnings after the game. Another instance of this problem is the card-playing phase of Bridge, in which two *defenders* have to coordinate their actions against the *declarer*, but they are prohibited from communicating by the rules of the game.

Finding an optimal equilibrium with *ex ante* coordination is NP-hard and inapproximable [10]. Celli and Gatti [10] introduced the first algorithm to compute optimal coordinated strategies for a team playing against an adversary. At its core, it is a column-generation algorithm exploiting a hybrid representation of the game, where team members play joint normal-form actions while the adversary employs sequence-form strategies [33]. It is crucial to observe that the number of joint normal-form actions of the team grows exponentially in the size of the game tree, which makes them impractical when dealing with games of medium/large size. More recently, Farina et al. [18] proposed a variation of the Fictitious Play algorithm, namely Fictitious Team-Play (FTP), to compute an approximate solution to the problem. Both approaches require to iteratively solve Mixed-Integer Linear Programs (MILP), which significantly limits the scalability of these techniques to large problems, with the biggest instances solved via FTP being in the order of 800 infosets per player. The biggest crux of these tabular approaches is the need for an explicit representation of the sequential game, which may not be exactly known to players, or could be too big to be stored in memory. For this reason, extremely large games are usually abstracted by bucketing similar states together. The problem with this approach is that abstractions procedures are domain-specific, require extensive domain knowledge, and fail to generalize to new scenarios (see, e.g., [8, 21, 22, 57]).

On the other side of the spectrum with respect to tabular equilibrium computation techniques there are Multi-Agent Reinforcement Learning (MARL) algorithms (see [9, 27] for a comprehensive tractation). These techniques do not require a complete knowledge of the environment and are sample-based by nature, but their application to imperfect-information adversarial team games presents a number of difficulties, such as dealing with private information and representing coordinated strategy spaces. The latter is of crucial importance since we will show that, even in simple settings, it is impossible to reach optimal coordination by exploiting completely decentralized policies, as it is customary in the RL literature.

Original Contributions. The contribution of the paper is two-fold. First, we study the problem of collecting meaningful histories

of play (i.e., trajectory sampling) of the team that can be used, in a subsequent phase, to compute strong coordinated strategies. Second, we propose a signal-mediated framework to represent and compute coordinated strategy from a buffer of past experience. We address the former problem showing that historical data can be efficiently collected via self-play on an approximated version of the team game, leveraging the notion of its *perfect-recall refinement*, which provides a natural way to perform centralized trajectory sampling when the objective is uncovering effective coordinated behaviors. Moreover, when the team members have symmetric observations, this approach allow us to prove that an optimal team coordinated strategy can be computed in polynomial time, on an equivalent two-player zero-sum game. Finally, we propose *signal mediated strategies* (SIMS) as a scalable way to capture coordination without the need for an explicit description of the underlying game. Specifically, SIMS represents a coordinated team strategy as the combination of a signaling policy (i.e., a distribution over signals) and one decentralized policy for each team member. First, a signal is sampled from the signaling policy and communicated to the team members. Then, each team member uses the signal to condition the output of a neural network encoding his/her decentralized policy. Therefore, in order to approximate an optimal coordinated strategy, team members have to learn from past experiences both the signaling policy and the *meaning* associated to each signal (i.e., a suitable parametrisation of the decentralised policies). We show that this is possible by testing our framework on a set of coordination games in which previous state-of-the-art multi-agent RL techniques could not reach an optimal coordinated strategy, and on an instance of a simple patrolling game defined over a grid-world.

2 RELATED WORKS

Learning how to coordinate multiple independent agents [4, 16] via Reinforcement Learning requires tackling multiple concurrent challenges, e.g., non-stationarity, alter-exploration and shadowed-equilibria [45]. There is a rich literature of algorithms proposed for learning cooperative behaviours among independent learners. Most of them are based on heuristics encouraging agents’ policies coordination [3, 5, 36, 37, 43, 44, 52].

Thanks to the recent successes of deep RL in single-agent environments [46, 59, 60], MARL is recently experiencing a new wave of interest and some old ideas have been adapted to leverage the power of function approximators [50, 51]. Several successful variants of the Actor-Critic framework based on the *centralized training/decentralized execution* paradigm have been proposed [19, 20, 42, 61]. These works encourage the emergence of coordination and cooperation, by learning a centralized Q -function that exploits additional information available only during training. Other approaches factorize the shared value function into an additive decomposition of the individual values of the agents [62], or combine them in a non-linear way [54], enforcing monotonic action-value functions. More recent works, showed the emergence of complex coordinated behaviours across team members in real-time games [28, 41], even with a fully independent asynchronous learning, by employing population-based training [29].

Player’s coordination is usually modeled from a game-theoretic perspective via the notion of *correlated equilibrium* (CE) [1], where

agents make decisions following a recommendation function, i.e., a *correlation device*. Learning a CE of *extensive-form games* (EFG) is a challenging problem as actions spaces grow exponentially in the size of the game tree. A number of works in the MARL literature address this problem (see, e.g., [14, 15, 23, 65]). Differently from these works, we are interested in the computation of TMECor [10].

In our work, we model the correlation device explicitly. By sampling a signal at the beginning of each episode, we show that the team members are capable of learning how to associate a precise meaning to a potentially uninformative signal. Our approach is closely related to the work by Chen et al. [13], which proposes a similar approach based on exogenous signals. Chen et al. [13] suggest that coordination can be encouraged by maximizing the mutual information between the signals and the joint policy.

3 PRELIMINARIES

In this section we provide a brief overview of extensive-form games (see also the textbook by Shoham and Leyton-Brown [58]).

An extensive-form games \mathcal{G} is a tree-form model of sequential interactions involving a set of players \mathcal{P} . A node v of the tree is defined by all the information on the current state of the game. For instance, in a poker game, a node is determined by the history of actions up to that point, and by the hand of each player. The set of actions available to the relevant player at a node v is denoted by $\mathcal{A}(v)$. Leaf nodes are called *terminal nodes*. We denote the set of terminal nodes by \mathcal{Z} . Each player $i \in \mathcal{P}$ as a payoff function $u_i : \mathcal{Z} \rightarrow \mathbb{R}$ which specifies her final reward for reaching a certain leaf. Exogenous stochasticity is represented via a *chance* player (denoted by C) which selects actions with a fixed known probability distribution. Given $z \in \mathcal{Z}$, we denote by $\rho_C(z)$ the probability with which the chance player plays so as to reach z .

Private information is modeled through the notion of **information states** (a.k.a. information sets). An information state s_i of player i comprises all nodes of the tree which are indistinguishable to i . Taken together, all information states of player i form a partition of the nodes where i has to act. We denote the set of all information states of player i as \mathcal{S}_i . Given $s \in \mathcal{S}_i$, for any pair of nodes $v, w \in s$, nodes v and w must have the same set of available actions. As is customary in the related literature, we assume *perfect recall*, i.e., no player forgets what he/she knew earlier in the game.

In this setting, we distinguish two fundamental paradigms for strategy representation. A **behavioral strategy profile** for player i is a collection specifying a point in the strategy simplex for each information state in \mathcal{S}_i . Formally, for any $s \in \mathcal{S}_i$, $\pi_i[s] \in \Delta(\mathcal{A}(s))$ specifies the probability distribution according to which player i selects an action at s . The second strategy representation is based on the notion of *normal-form plan*, which is a vector specifying an action $a_s \in \mathcal{A}(s)$ for each information state $s \in \mathcal{S}_i$. A *reduced normal-form plan* p_i is a normal-form plan where irrelevant information is removed: it specifies an action only for information states that can be reached following the actions specified by p_i higher up in the game tree. We denote the set of reduced-normal-form plans of i as P_i . Given a leaf $z \in \mathcal{Z}$, we denote by $P_i(z) \subseteq P_i$ the set of reduced-normal-form plans in which player i plays so as to reach z . A **normal-form strategy** for player i is a probability distribution $\mu_i \in \Delta(P_i)$, where $\mu_i[p_i]$ is the probability with which player i

will play according to the actions specified by p_i . We say that two strategies of player i are *realization equivalent* if they force the same distribution over the leaves of the game, *i.e.*, the probability of player i playing so as to reach any $z \in \mathcal{Z}$ is the same under both strategies.

4 CHALLENGES OF TEAM COORDINATION

A **team** is a set of players sharing the same objectives [2, 64]. We study games where a team faces an opponent in a zero-sum interaction. In order to simplify the presentation, we describe our results for the case of a team composed by two agents, denoted by T1 and T2, playing against an opponent O. The extension of our framework to the case with multiple team members and opponents is straightforward. We are interested in settings where team members can communicate and agree on a coordinated strategy before the beginning of the game, but are unable to communicate during the playing phase. Two examples of this setting are collusion in poker and collusion during bidding, where T1 and T2 will share their earning at the end of the game and do not want to be detected, and Bridge, where T1 and T2 are members of the same team, and the rules of the game do not allow them to communicate during the game.

4.1 TMECor and Coordinated Strategies

The most appropriate notion of equilibrium for this setting is the *team-maxmin equilibrium with coordination device (TMECor)* introduced by Celli and Gatti [10]. A powerful, game-theoretic way to think about coordination is through the notion of **coordination device**. Intuitively, before the game starts, team members can identify a set of joint normal-form plans within $P_{T1} \times P_{T2}$. Then, just before the play, the coordination device draws one of such plans according to a suitable probability distribution $\mu_T \in \Delta(P_{T1} \times P_{T2})$, and team members will act as specified in the selected joint plan. A probability distribution over $P_{T1} \times P_{T2}$ is called a **coordinated strategy**. A TMECor is a Nash equilibrium (NE) of the game where team members play their best coordinated strategy. Let $u_T : \mathcal{Z} \rightarrow \mathbb{R}$ be the shared payoff function of the team. Then, computing a TMECor amounts to solving the following optimization problem:

$$\begin{aligned} \max_{\mu_T} \min_{\mu_O} & \sum_{z \in \mathcal{Z}} \sum_{\substack{p_{T1} \in P_{T1}(z) \\ p_{T2} \in P_{T2}(z) \\ p_O \in P_O(z)}} \mu_T[p_{T1}, p_{T2}] \mu_O[p_O] u_T(z) \\ \text{s.t.} & \mu_T \in \Delta(P_{T1} \times P_{T2}) \\ & \mu_O \in \Delta(P_O) \end{aligned} \quad (1)$$

By taking the dual of the inner minimization problem, Problem (1) can be reformulated a linear programming problem. The main difficulty is then managing the coordinated strategy μ_T , as its dimension grows exponentially in the size of the game tree. In an ϵ -TMECor, neither the team nor the opponent can gain more than ϵ by deviating from their strategy, assuming that the other does not deviate.

4.2 Common Pitfalls of Coordination

By sampling a recommendation from the joint probability distribution μ_T the coordination device introduces a correlation between

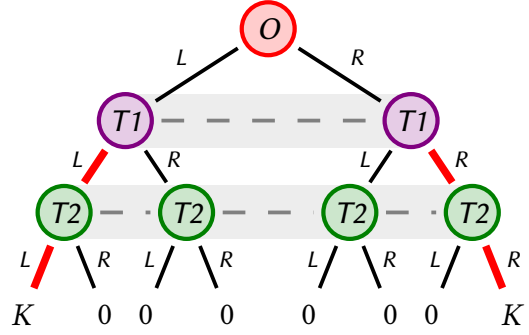


Figure 1: Coordination game. Nodes within the same information state are connected through the grey dotted lines. Leaf nodes display the payoff of the team. The payoff of the opponent is $u_O(\cdot) = -u_T(\cdot)$.

team members' actions that would otherwise be impossible to capture through behavioral strategies. This is illustrated by the following simple example.

EXAMPLE 1 (COORDINATION GAME). *The game is played by a team of two players (as usual denoted by T1 and T2) and an opponent O. Each player of the game has two available actions: left (L) and right (R). Players have to select one of the two actions without having observed the choice of the other players. Team members receive a payoff of K only if they both guess correctly the action taken by the opponent and mimic that action. For example, when the opponent plays L, the team is rewarded K if and only if both T1 and T2 play L. Otherwise they have payoff equal to 0. Team's rewards are depicted in the leaves of the tree in Figure 1.*

If team members did not have the opportunity of coordinating their strategies, then the best they could do is selecting an action randomly. This corresponds to the Nash equilibrium of the game without coordination, where T1, T2, and O play a uniform behavioral strategy. This leads to an expected return for the team of K/4. When coordination is possible, team members can skew their joint strategy to play only the reduced-normal-form plans (L, L) and (R, R) (displayed in red in Figure 1). This allows them to avoid playing pairs of actions that would surely result in a 0 payoff, independently of O's action. The TMECor of the game is reached when O plays with equal probability L and R, and the team plays according to a joint coordinated strategy such that $\mu_T[(L, L)] = \mu_T[(R, R)] = .5$. At the TMECor the team has an expected utility of K/2: the team can double its returns by adopting coordinated strategies.

Despite their theoretical superiority with respect to behavioral strategies, coordinated strategies have a major downside in practice: they require an exponential number of bits to be represented. This is because the set of joint reduced-normal-form plans grows exponentially in the size of the game tree. Hence, previous work on the topic largely focuses on providing more manageable representations of the coordinated strategy space (see, *e.g.*, Celli and Gatti [10], Farina et al. [18]). Here, we take a radically different approach by proposing a model-free framework to approximate coordinated strategies in a RL fashion. Classical multi-agent RL algorithms employ decentralized policies, which can be described

as behavioral strategies. Therefore, they are unable to attain the optimal coordinated outcome even in simple settings such as the one depicted in the previous example. We show how to reach a middle-ground between decentralized policies (*i.e.*, behavioral strategies) and coordinated strategies.

5 CENTRALIZED TRAINING FOR IMPERFECT-INFORMATION TEAM GAMES

Our framework subdivides the training procedure for approximating a TMECor in two separate phases: in the first phase the game trajectories (*i.e.*, sequences of state-action pairs) are collected in self-play and stored in a buffer \mathcal{M} . In the second phase, the coordinated team strategy, that will be used at test time, is learned from \mathcal{M} via supervised learning. This approach was already exploited by known algorithms for the two-player, zero-sum setting such as Neural Fictitious Self Play (NFSP) [26] and Deep-CFR [7]. However, the team coordination setting presents a number of additional challenges: first, unlike in the two-player, zero-sum setting, it is not clear how to collect trajectory samples in such a way to guarantee convergence when two or more team members are coordinating against an opponent. Second, the coordinated strategy to be used at test time must be able to capture coordination without having to represent the exponentially large μ_T . In this section we propose a solution for the former problem, the latter is discussed in the remainder of the paper (Section 6).

The former problem amounts to populating the buffer \mathcal{M} with meaningful trajectories to learn a coordinated strategy. We do that through a centralized training phase in self-play during which we let each team member share with the other its private information. We can provide a useful interpretation of this procedure by considering team members T1 and T2 as a single meta-player T. In the original game \mathcal{G} , T may have imperfect recall. For example, in Figure 1, T would not remember his/her first move when choosing its second action. Equivalently, in card-playing game with private cards, T would have to periodically forget about T1’s cards and regain memory of T2’s hand, and vice versa. By letting T1 and T2 sharing their private information, we are making sure that T has perfect recall. This is because information sharing produces finer grained information states. If we denote by \mathcal{G}^* the game resulting from this process, we can say that \mathcal{G}^* is a **perfect-recall refinement** of \mathcal{G} for the meta-player T. Following Lanctot et al. [35], we define a perfect-recall refinement for T as follows:

Definition 1 (Perfect-recall refinement). *Given a game \mathcal{G} with an imperfect-recall player T, \mathcal{G}^* is a perfect-recall refinement of \mathcal{G} if T has perfect recall in \mathcal{G}^* and \mathcal{G} is an abstraction of \mathcal{G}^* , that is if, for any pair of nodes v, w of T it holds $\mathcal{A}(v) \subseteq \mathcal{A}^*(v)$ and $v, w \in s_T$, then there exists an information state $s_T^* \in \mathcal{S}_T^*$ such that $v, w \in s_T^*$.*

As an illustrative example, the perfect-recall refinement of the game in Figure 1 is obtained by splitting the information state of T2 into two distinct information states: one following action L of T1, and the other following from action R of T1. In a perfect-recall refinement team members share the same observations on the state of the game. Specifically, in \mathcal{G}^* team members have imperfect information that can be due only to either partial observability of the actions of the opponent (as it happens in the perfect recall refinement of Figure 1), or to private information of the opponent

due to a chance moves higher up in the game tree. The key observation is that in \mathcal{G}^* either T1 and T2 both observe an action (of any player, chance included), or they both do not. We say that \mathcal{G}^* has **symmetric observations** for the team. In this setting, the underlying reason for which the meta-player T has imperfect recall is the limited observability within the team: that is, T1 not being able to observe every T2’s action and vice versa.

5.1 Coordination in Games with Team Symmetric Observability

In the class of games in which \mathcal{G} already has symmetric observations for the team, we show that our approach provably yields a TMECor. To show this we need to reintroduce the notion of A-loss recall [30, 32]. A player has A-loss recall if he/she has perfect recall, or if his/her losses of memory can be traced back to forgetting his/her own actions.

Definition 2 (Symmetric observability). *A game \mathcal{G} has symmetric observability for the team if and only if the meta-player T has A-loss recall.*

We observe that a perfect-recall refinement \mathcal{G}^* always has symmetric observability for the team, but the converse is not true in general. A practical example where this condition holds is the game of Goofspiel [56]. In this setting, both team members cannot observe the opponent’s move up until the end of each turn, and do not have any private information but the action they just played.

In the following, \mathcal{G}^* is always treated as a two-player, zero-sum game between the meta-payer T and O. Let $\pi = (\pi_T, \pi_O)$ be an arbitrary behavioral strategy profile of \mathcal{G}^* . In order to prove our theoretical results we need the following auxiliary definitions.

Definition 3. *Two games \mathcal{G} and \mathcal{G}' differing only for their information partitions ($\{\mathcal{S}_i\}_{i \in \mathcal{P}}$ and $\{\mathcal{S}'_i\}_{i \in \mathcal{P}}$, respectively) are μ -equivalent if for any player i and for any normal-form strategy μ of i in \mathcal{G} , there exists a realization equivalent normal-form strategy μ' in \mathcal{G}' , and vice versa.*

Given a node v , let $X_i(v)$ be the set of information state-action pairs of player i on the path from the root of the tree to v . We will make use of the *inflation* operation [17, 30, 49], which we define as follows:

Definition 4 (Immediate inflation). *Let \mathcal{S}_i and \mathcal{S}'_i be two information partitions of player i . We say that \mathcal{S}'_i is an immediate inflation of \mathcal{S}_i iff there exists $s \in \mathcal{S}_i$ and $s', s'' \in \mathcal{S}'_i$ such that: (i) the set of nodes comprised by s is equal to the set of nodes comprised by s' and s'' (*i.e.*, $s = s' \cup s''$), and (ii) for each $v \in s'$ and $w \in s''$ there exists $\bar{s} \in \mathcal{S}_i \cap \mathcal{S}'_i$ such that $(\bar{s}, a) \in X_i(s')$, $(\bar{s}, b) \in X_i(s'')$ for some actions $a \neq b$.*

Definition 5 (Inflation). *Given a player i , an information partition \mathcal{S}'_i is an inflation of \mathcal{S}_i iff it is obtained by successive applications of immediate inflation operations to \mathcal{S}_i .*

When an inflation of \mathcal{S}_i has no further immediate inflations, it is called *complete inflation*.

By leveraging the notion of inflation we can prove the following result, which constitutes a strong motivation for our approach to the computation of a TMECor.

Theorem 1. *Given a game \mathcal{G} with symmetric observability, for any $\pi = (\pi_T, \pi_O)$ which is an NE of \mathcal{G}^* there exists a pair of realization equivalent strategies (μ_T, μ_O) which is a TMECor of \mathcal{G} , and vice versa.*

PROOF. Given \mathcal{G} , let T be the team meta-player. Formally, T’s information states are such that

$$S_T = S_{T1} \cup S_{T2}.$$

Then, the set of coordinated strategies of the team $\Delta(P_{T1} \times P_{T2})$ is equal to the set of normal-form strategies of T, i.e., $\Delta(P_T)$. We are left with a two-player, zero-sum game between O and T. O has perfect recall and, by Def. 1, T has A-loss recall. Then, \mathcal{G} is A-loss.

Case ①: from μ_T to π_T . By Theorem 5.A of Kaneko and Kline [30] we have that since the information partition of \mathcal{G} satisfies the A-loss condition, then the complete inflation of \mathcal{G} coincides with the perfect-recall refinement \mathcal{G}^* . Since the inflation procedure preserves the same μ -equivalence class, we have that \mathcal{G} and \mathcal{G}^* are μ -equivalent. Hence, if we denote by μ_T a TMECor of \mathcal{G} , by Definition 3 there exists a normal-form strategy μ^* of the team meta-player in \mathcal{G}^* which is realization equivalent to μ_T . By Kuhn’s theorem [34], every normal-form strategy of \mathcal{G}^* has an equivalent behavioral strategy: there exists a behavioral strategy π_T^* of the team meta-player in \mathcal{G}^* which is realization equivalent to μ_T^* , which implies realization equivalence to μ_T . By definition of T, $\Delta(P_{T1} \times P_{T2}) = \Delta(P_T)$. Then, since O’s information partition is left unchanged going from \mathcal{G} to \mathcal{G}^* , if μ_T is an NE with strategy space $\Delta(P_{T1} \times P_{T2})$, then μ_T^* and π_T^* are NE of \mathcal{G}^* .

Case ②: from π_T to μ_T . The proof follows the same points of the previous case. \square

Hence, Theorem 1 justifies the introduction of our centralized training regiment over the perfect-recall refinement of each game, since in many cases this implies sampling trajectories from true equilibrium strategies. Then, we have the following key result, which is in striking contrast from with impossibility results by Celli and Gatti [10].

Theorem 2. *For any game \mathcal{G} with symmetric observability, a TMECor can be computed in polynomial time.*

PROOF. The result immediately follows from Theorem 1 with the following remarks:

- An NE of a two-player, zero-sum game (i.e., \mathcal{G}^*) can be computed in polynomial time via linear programming by exploiting, for example, the *sequence-form representation* of the game [33].
- The complete inflation of an arbitrary game \mathcal{G} can be computed in polynomial time [12, Theorem 3.3].
- Given π_T in \mathcal{G}^* it is possible to compute the reach probability associated to each leaf node $z \in \mathcal{Z}$ in polynomial-time. From there, a realization equivalent normal-form strategy for T can be computed in polynomial time [11, Theorem 4]. This is a TMECor strategy for the team in \mathcal{G} . \square

Even in games where T does not have A-loss recall, collecting trajectories on a perfect-recall refinement allows team members to populate the buffer of past experience with meaningful trajectories, which are the result of coordinated play in the ideal setting in which they are able to share information. As we will show in

the experimental evaluation, this is essential to compute strong coordinated strategies at test time. The crucial problem of performing strategy mapping between \mathcal{G}^* and \mathcal{G} is entirely handled by the SIMS framework (Section 6).

5.2 Centralized Training: Trajectory Sampling on Perfect-recall Refinement

A perfect-recall refinement of a game, as defined in Def. 1, provides the team players with the ability to observe the actions played by their team members at each step of the game. This simulates the best-case scenario in which agents can communicate during game-play: hence, we can collect meaningful trajectories on the relaxed game \mathcal{G}^* using any trajectory sampling algorithm and still exploit the extra information available at training time. We test our framework using NFSP [26] and QMIX [54] (see Appendix for further details). Going back to the coordination game with horizon 2 in Figure 1, the team players’ observations can be split in two separate information issues: the state of the game and the teammate action. Consider, for example, an episode of the game in which T1 plays action L and T2 plays action R. Then, the observations of the players will be: $\sigma_{T1} : \{o : R, a_{T2} : \text{None}\}$ and $\sigma_{T2} : \{o : R, a_{T1} : L\}$, where the played actions are $a = \{a_{T1} : L, a_{T2} : R\}$ and σ_{T1}, σ_{T2} are both R because in the original game there is only one information set for each player.

Before storing the collected trajectory into the buffer, players’ observations are purged from the extra knowledge of the teammate’s action since this information is not available at execution time. In the example, only T1’s action is purged from T2’s observation as for the multistage nature of the original game the only player that does not observe the action of the other is T2. Note that in multi-stage games like the coordination game of the example, in order to purge the observation of T2 from the information obtained by observing T1’s action is enough to set the observation of T2 equal to the one of T1. In the next section, we explain how to use the collected trajectories to learn the team strategy via signal coordination.

6 SIMS: SIGNAL MEDIATED STRATEGIES

In this section, we focus on the problem of representing team coordinated strategies and, in doing so, we implicitly solve the problem of mapping strategies of \mathcal{G}^* back to the original game \mathcal{G} .

As noted by Farina et al. [18], any coordinated strategy μ_T can always be represented as the convex combination of a finite set of behavioral strategy profiles (π_{T1}, π_{T2}) of the team. Therefore, we exploit a set of exogenous signals to condition team members’ decentralized policies. Specifically, the decentralized policy followed by each team member at test time is conditioned on a signal which is sampled just before the beginning of the episode. The questions here are: (i) how to properly learn the probability distribution over signals in order to optimally balance different decentralized policies? (ii) how to make sure team members don’t just ignore signals?

Let \mathcal{M} be a memory of trajectories collected from game-play interactions, following the centralized sampling procedure described in Section 5.2. Each sample in \mathcal{M} is a pair (o, t) containing a game observation and its target action. Once the experience buffer is full of meaningful coordinated trajectories, we can try to imitate the

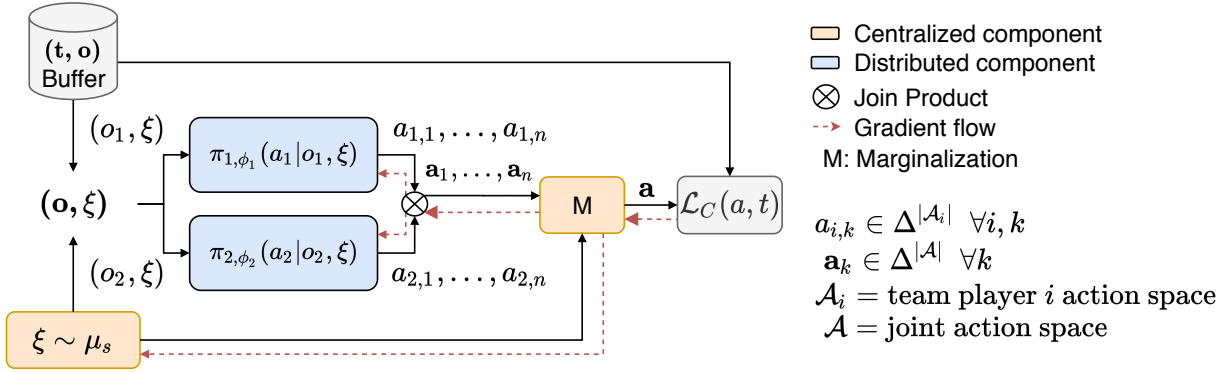


Figure 2: Diagram block of the training of SIMS to learn a coordinated strategy in a game with 2 team players.

team players’ behaviour via supervised learning. In practice, coordination is achieved by conditioning the team players’ policies on an exogenous signal drawn from a learnable distribution, explicitly implementing a *coordination device*.

Signal Mediated Strategies (SIMS). Let be \mathcal{T} the set of team players $(T_1, T_2, \dots, T_{|\mathcal{T}|})$, with observation space \mathcal{S} , and action space \mathcal{A}_{T_i} . For simplicity, we assume that the game observation o is the same for all team players and t is a tuple $(t_1, \dots, t_{|\mathcal{T}|})$ that specifies the target action for each team member. In the rest of the section and in Figure 2, for the sake of clarity, we drop the prefix T from the player notation and denote each player only by its index.

Each team member’s policy is defined as $\pi_{i, \phi_i} : \mathcal{S} \times \mathcal{E} \rightarrow \Delta^{|\mathcal{A}_i|}$, where \mathcal{E} is the space of signals and ϕ_i are the parameters of a deep neural network that parameterizes the policy. Crucially, each policy is conditioned both on the observation o (the state of the game) and on an exogenous signal $\xi \in \mathcal{E}$, sampled at the beginning of the game. The signals are sampled according to a categorical distribution $\mu_s = \text{Cat}(n, \theta)$, where n is the number of available signals and their probabilities are learned from experience, parameterized by θ . We sample all the n signals from the distribution μ_s and compute the action distributions $a_{j,k}$ from π_{j, ϕ_j} for each team member and for each signal $1 \leq k \leq n$. The marginal action distributions of the team players are then multiplied via joint product to obtain the team action distribution \mathbf{a}_k for each signal ξ_k :

$$\mathbf{a}_k = \prod_{j \in \mathcal{T}} \pi_j(\cdot | o, z_k, \theta_j) = \prod_{j \in \mathcal{T}} a_{j,k} \quad (2)$$

We marginalize the joint action distribution over all the n signals $\mathbf{a} = \sum_{i \leq n} \mu_s[\xi_i] \mathbf{a}_i$, and compute a classification loss \mathcal{L}_C with respect to the target action t :

$$\mathcal{L}_C(\mathbf{a}, \theta) = \text{CrossEntropy}(\mathbf{a}, t | \theta). \quad (3)$$

We also introduce an additional entropy regularization term $\mathcal{L}_{E,s} = \sum_{t \in \mathcal{T}} H(a_{t,k})$, to enforce pure strategies over the actions probabilities for each distinct signal. The overall supervised learning loss to be minimized is:

$$\mathcal{L}(\mathbf{a}_1, \dots, \mathbf{a}_n, \theta) = \mathcal{L}_C + \beta \sum_{k \leq n} \mu_s[\xi_k] \mathcal{L}_{E,s}(\mathbf{a}_k) \quad (4)$$

where we dropped the arguments of \mathcal{L}_C and $\mathcal{L}_{E,d}$ for clarity, and β is the coefficient of the regularization term. A block diagram of the SIMS framework is shown in Figure 2.

7 EXPERIMENTAL EVALUATION

In this section, we empirically evaluate our framework against some of the state-of-the-art multi-agent RL algorithms available in the literature. First, we provide some details on our experimental setting, then we provide the main experimental results.

7.1 Experimental Setting

We use as benchmarks for our experimental evaluation different instances of the coordination game in Figure 1 and an instance of the *patrolling game* shown in Figure 4.

Coordination game. We consider two different instances of the coordination game in Figure 1. Specifically, *coord-2* and *coord-4* are coordination games with horizon 2 and 4 respectively, *i.e.*, each team member plays once or twice before reaching a terminal node. In each game, the team receives a positive payoff only if both players guess correctly the action of the opponent and mimic its action. Otherwise, all players have payoff equal to 0. The final team payoff is split equally among the team members. Note that these game are known to be the worst-case instances in terms of difficulty of coordination [10], which makes them ideal to understand whether team members are effectively coordinating or not. In particular, we consider the setting in which there is an imbalance in the team’s payoffs, *i.e.*, instead of receiving K for playing both *left* or *right*, team members receive K and $K/2$, respectively. Imbalanced payoffs make the process of learning an optimal coordinated strategy more challenging with respect to the balanced setting. Intuitively, this is because a uniform probability distribution over signals is no longer enough to reach an optimal strategy. In all the experiments we set $K = 100$. We also tested other combinations of imbalanced payoffs, but we omit them since they yield similar results.

Patrolling game. We focus on a simple patrolling game defined over a grid-world. The setting is described in Figure 4. Both team members play as the defending agents and, starting from the central position in the grid, they have to reach one of the four sites that must be defended. The game evolves synchronously and at every

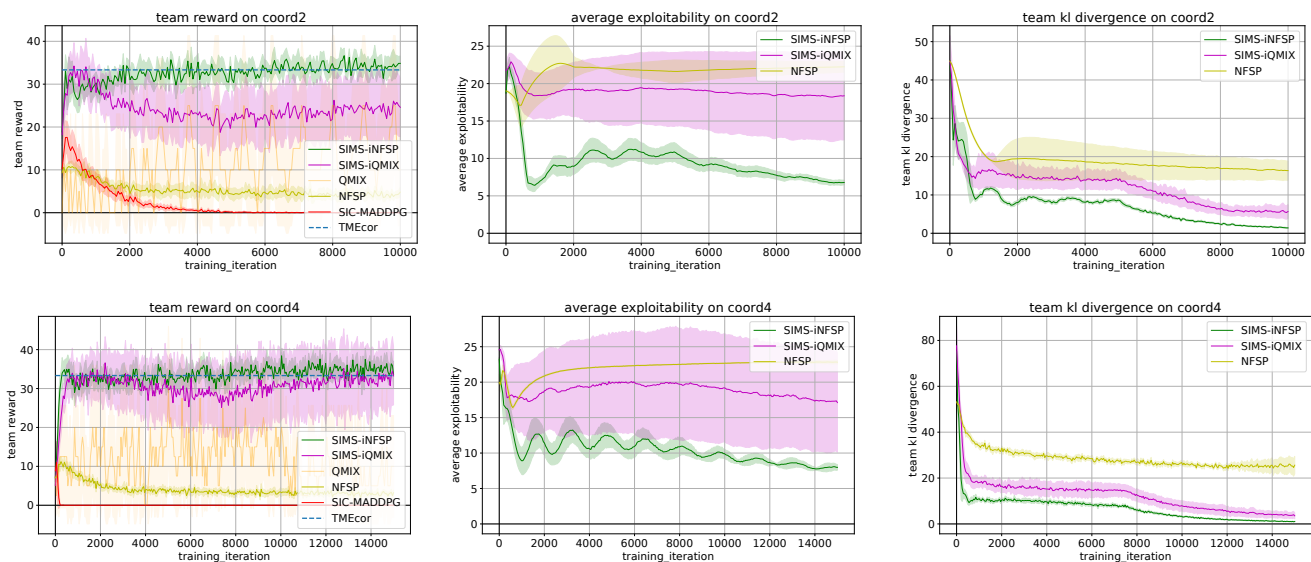


Figure 3: First row: performance of SIMS-iNFSP on a coordination game with horizon 2 (coord-2). Second row: performance on coordination game with horizon 4 (coord-4). Left column: Average Team reward. Center column: Average exploitability. Right column: Kullback-Leibler divergence between the joint average strategy and the optimal TMECor strategy.

round each agent can either move in one of the four directions or remain on the same cell. After three time instants, the opponent decides which one of the sites to attack (without observing the team members’ positions). The defense is considered successful if both agents are located in the exact cell attacked by the opponent, unsuccessful otherwise. In case of a successful defense the team gets as a reward +1 and the opponent -1, while in case of a successful attack the team gets -1 and the opponent +1. The game is denoted as *patrolling_4_3* to indicate 4 sites to defend and attacks taking place after 3 time steps.

Baselines. We evaluate SIMS in combination with two different algorithms for collecting trajectories, namely **SIMS-iNFSP** and **SIMS-iQMIX**. The prefix “i” indicates that trajectories have been collected using the relaxed game (*i.e.*, the perfect-recall refinement) and the extra information about the action of the teammate has been discarded when saved into the buffer as described in Section 5.2 – “i” stands for *inflation operator* from Def. 4 and Def. 5. In order to improve the performances of the trajectory sampling performed by iNFSP, we adopt parameter sharing between team members. We show the necessity of centralized training by sampling the trajectories on the relaxed game and comparing the performance against versions of the NFSP and QMIX algorithms where trajectories have been collected on the original game. Finally, we also perform experiments using SIC-MADDPG [13], another framework that models the coordination device explicitly. SIC-MADDPG extends the actor-critic framework proposed in [42] for competitive-cooperative environments by adding signals coordination and ensuring that the signal is taken into account with a mutual information regularizer. MADDPG and hence SIC-MADDPG take advantage of the centralized training by observing extra information such as the state of the game and the actions of other players. In a coordination game, this

is equivalent to observing only the actions chosen by other players, as the game can be considered inherently stateless. Furthermore, we empirically show that oftentimes SIC-MADDPG and QMIX cannot reach convergence to a TMECor, even in settings in which iNFSP is guaranteed to converge.

Implementation details. All the policies and value networks are parameterized by Multi-Layer Perceptron (MLPs) with two fully-connected layers of 128 units each and *ReLU* activation. The batch size has been set to 128 and the optimizer used is *Adam* [31] for all experiments, with a learning rate $lr = 10^{-3}$ and default β_1 and β_2 parameters. The replay buffers used by all algorithms have a size of $2 \cdot 10^4$, while the reservoir buffers in NFSP have a size of 10^5 . The mixing layer used by QMIX has also 128 units with *ReLU* activation. In all the considered games, there exists an optimal coordinated strategy employing only two signals. In SIMS, we used a signal space composed of 5 signals for the coordination games and 4 signals for the patrolling game. In order to improve stability of the strategy computation, we use a linear scheduling for the entropy regularizer coefficient β in Eq. 4. In particular, we keep fixed the parameter at β_{init} for the first half of the training, and then linearly increase it up to β_{end} for the second half of the training phase. For all our experiments we set $\beta_{init} = 0$ and $\beta_{end} = 1$. Moreover, we accumulate the gradients of the signals distribution’s parameters and perform a back-propagation step every $N_{sig} = 20$ iterations. Both the introduction of the scheduling for β and the accumulation of gradients for the signal distribution significantly stabilize the training. All the algorithms are evaluated for 100 episodes every 50 training iterations by averaging across 10 different runs. We used PyTorch [53] and the multi-agent environment abstraction from RLLib [40].

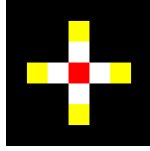


Figure 4: Patrolling game instance. Yellow cells are the facilities that must be defended, the red cell is the starting cell for both defending agents.

7.2 Results

Now, we discuss the results on the experimental settings described above, showing the success of the SIMS framework in achieving team coordination.

Coordination game. In coordination games, decentralized strategies are not expressive enough to describe the optimal coordinated behaviours of the team members. This can be observed by the unsatisfactory performance of NFSP, QMIX and SIC-MADDPG, when compared to the optimal TMECor. On the other hand, SIMS-*i*NFSP is able to compute and represent an optimal coordinated strategy for the team. As shown in Figure 3 (Right), the team reward of the joint policy computed through SIMS-*i*NFSP is close to the optimal one reached at the equilibrium. This is further confirmed by Figure 3 (Center), as the exploitability of the strategy obtained through SIMS-*i*NFSP decreases towards zero. Informally, the exploitability of a strategy gives a measure of how much worse that strategy performs versus a best response of the opponent, compared to an equilibrium strategy. Finally, the leftmost column of Figure 3 reports the Kullback-Leibler divergence between the joint policy and the TMECor strategy profile during training. The analysis of SIMS-*i*QMIX and its comparison with the other algorithms show two important aspects. On one side it shows how the proposed decentralized training paradigm outperforms the centralized training paradigms of QMIX and SIC-MADDPG, that together with NFSP fail in being able to capture coordination between team members. This happens because, by relaxing the environment, we are giving the possibility to the team members to always avoid playing uncoordinated actions, resulting in an increase of the average reward obtained, and in a decrease in the Kullback-Leibler divergence with the optimal strategy. Comparing SIMS-*i*NFSP and SIMS-*i*QMIX further stresses the importance of collecting meaningful trajectories. While the trajectory sampling performed by *i*NFSP has theoretical guarantees of converging to an equilibrium in some specific settings (see Section 5.1), this is not the case for *i*QMIX. This results in high instability of SIMS-*i*QMIX’s performances and higher average exploitability, that is a higher distance from the equilibrium. We also observed that policy gradient approaches struggle to find a pair of strategies with satisfactory performances. We conjecture that this is due to the sparsity of the rewards, a key characteristic of coordination games that inevitably weakens the gradients and prevents the algorithm from learning.

Patrolling game. The analysis of the patrolling game offers the opportunity to visualize the level of coordination achieved via

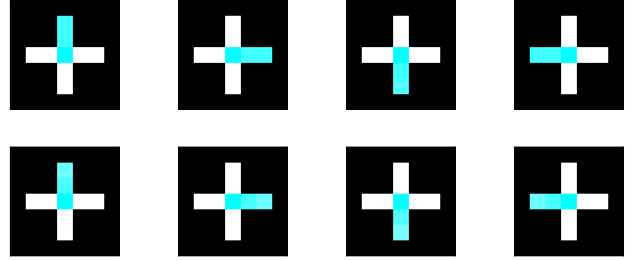


Figure 5: First row: heatmaps of player T1 on *patrolling_4_3*. Second row: heatmaps of player T2 on *patrolling_4_3*. Every column represents a single signal.

SIMS-*i*NFSP through the heatmaps that we report in Figure 5. The heatmaps clearly show that team members learn to associate the same shared meaning to each signal. Specifically, each signal (*i.e.*, each column of Figure 5) is associated with the same site by both team members, resulting in an optimal coordination between them. In Figure 6, we plot the average reward achieved by SIMS-*i*NFSP against the average reward of NFSP and the reward at the TMECor. Also in this case, SIMS-*i*NFSP is able to guarantee to team members the optimal reward, resulting in better performances than NFSP.

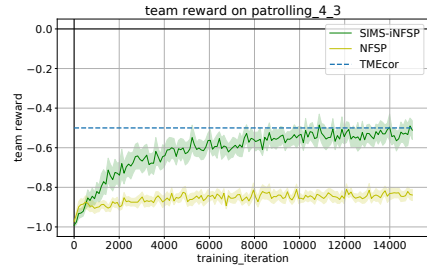


Figure 6: Team reward on a patrolling game with 4 sites to defend and attack after 3 time steps.

8 CONCLUSIONS AND FUTURE WORK

In this paper we propose to leverage the notion of perfect-recall refinement of a game to perform centralized trajectory sampling for a team of agents whose goal is to learn a coordinated strategy. Moreover, we introduce a supervised learning framework (SIMS) for computing and representing joint coordinated strategies of a team from a buffer of past experiences.

We provide a preliminary experimental evaluation which shows promising performance of our framework. Future works will be devoted to further testing the ideas we presented on more challenging benchmarks where the abstraction power of deep RL techniques can be fully appreciated. For example, we plan to test our techniques on predator-prey environments such as Wolpack [38]. We are also planning more extensive tests on different trajectory-sampling algorithms that may be used to populate the buffer of past experiences used by the team.

A ALGORITHMS FOR TRAJECTORY SAMPLING

Trajectory sampling plays an important role for the convergence of the players’ strategies to the equilibrium of the game. We give a brief description of two offline RL algorithms that can be used in our framework to collect trajectories on the relaxed game.

Neural Fictitious Self-Play (NFSP). *Fictitious Play* [6] is a game-theoretic self-play algorithm in which players iteratively play their best responses against their opponents’ average strategies. In certain classes of games, e.g. two-player zero-sum [55] and many-player potential games [48], the average strategies are guaranteed to converge to Nash Equilibria (NE), or to approximate ϵ -NE [39], under approximate best responses and perturbed average strategy updates.

The original formulation of Fictitious Play is defined for normal-form games, resulting in exponential complexity for extensive-form games. Heinrich et al. [25] introduced exact (XFP) and machine learning-based (FPS) versions of the model that are implemented in behavioural strategies. The algorithms are realization equivalent to their normal-form counterpart, inheriting its convergence guarantees to an exact and ϵ -NE, respectively while reducing the complexity from exponential to linear in time and space.

A third variant, used in this work, *Neural Fictitious Self-Play* (NFSP) [26], combines FSP with neural function approximators. In NFSP, agents interact with each others generating datasets of experience in self-play. Each agent collects and stores its experienced transition tuples, $(s_t, a_t, r_{t+1}, s_{t+1})$, in a memory, \mathcal{M}_{RL} , while its own best response behaviour, (s_t, a_t) , is stored in a separate memory, \mathcal{M}_{SL} . Each agent uses off-policy reinforcement learning by training a DQN [47], $Q(s, a|\theta^Q)$, to predict action values, from the memory \mathcal{M}_{RL} of its opponents’ behaviour. The agent’s approximate best response strategy is defined as $\beta = \epsilon$ -greedy(Q), which selects a random action with probability ϵ and chooses the action that maximizes the predicted action values otherwise. A separate neural network, $\Pi(s, a|\theta^\Pi)$, is trained to imitate the agent’s own past best response behaviour using supervised classification on the data in \mathcal{M}_{SL} . This network maps states to action probabilities and defines the agent’s average strategy, the one that is guaranteed to converge to the equilibrium strategy.

QMIX. QMIX [54] is an offline value-based method that can train decentralised policies in a centralised end-to-end fashion inducing agents’ coordination. Learning how to coordinate multiple agents toward cooperative behaviours is hard due to numerous challenges. One can train fully decentralized policies disregarding agents’ interactions as in *Independent Q-Learning (IQL)* [63], but this may not converge because of non-stationary caused by others agents learning and exploration. On the other hand, centralised learning of joint actions can naturally handle coordination problems and avoids non-stationarity, but is hard to scale, as the joint action space grows exponentially in the number of agents. Similarly to *Value Decomposition Networks (VDNs)* [62], QMIX lies between IQL and centralised Q -Learning. By estimating a joint action-values Q_{tot} as a non-linear combination of per-agent values Q_a , conditioned only on local observations, QMIX can represent complex centralised action-value functions with a factored representation

that scales well in the number of agents. QMIX enforces that a global arg max performed on Q_{tot} yields the same result as a set of individual *argmax* operations performed on each Q_a . This allows each agent to compute decentralized policies by choosing greedy actions with respect to its Q_a . Agents can also be conditioned on their action-observation history to deal with partial observability in the environment by using recurrent neural networks to model their value functions [24]. For more detailed descriptions of the presented algorithms please refer to the original papers.

REFERENCES

- [1] Robert J Aumann. 1974. Subjectivity and correlation in randomized strategies. *Journal of mathematical Economics* 1, 1 (1974), 67–96.
- [2] Nicola Basilico, Andrea Celli, Giuseppe De Nittis, and Nicola Gatti. 2017. Team-maxmin equilibrium: efficiency bounds and algorithms. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- [3] Daan Bloembergen, Michael Kaisers, and Karl Tuyls. 2010. Lenient frequency adjusted Q-learning. In *Proc. of 22nd Belgium-Netherlands Conf. on Artif. Intel.*
- [4] Craig Boutilier. 1999. Sequential optimality and coordination in multiagent systems. In *IJCAI*, Vol. 99. 478–485.
- [5] Michael Bowling and Manuela Veloso. 2001. Rational and convergent learning in stochastic games. In *International joint conference on artificial intelligence*, Vol. 17. Lawrence Erlbaum Associates Ltd, 1021–1026.
- [6] George W Brown. 1951. Iterative solution of games by fictitious play. *Activity analysis of production and allocation* 13, 1 (1951), 374–376.
- [7] Noam Brown, Adam Lerer, Sam Gross, and Tuomas Sandholm. 2018. Deep counterfactual regret minimization. *arXiv preprint arXiv:1811.00164* (2018).
- [8] Noam Brown and Tuomas Sandholm. 2018. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science* 359, 6374 (2018), 418–424.
- [9] Lucian Busoni, Robert Babuska, and Bart De Schutter. 2008. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 38, 2 (2008), 156–172.
- [10] Andrea Celli and Nicola Gatti. 2018. Computational results for extensive-form adversarial team games. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [11] Andrea Celli, Alberto Marchesi, Tommaso Bianchi, and Nicola Gatti. 2019. Learning to correlate in multi-player general-sum sequential games. In *Advances in Neural Information Processing Systems*. 13076–13086.
- [12] Andrea Celli, Giulia Romano, and Nicola Gatti. 2019. Personality-Based Representations of Imperfect-Recall Games. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*. 1868–1870.
- [13] Liheng Chen, Hongyi Guo, Haifeng Zhang, Fei Fang, Yaoming Zhu, Ming Zhou, Weinan Zhang, Qing Wang, and Yong Yu. 2019. Signal Instructed Coordination in Team Competition. *arXiv preprint arXiv:1909.04224* (2019).
- [14] Ludek Cigler and Boi Faltings. 2011. Reaching correlated equilibria through multi-agent learning. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*. International Foundation for Autonomous Agents and Multiagent Systems, 509–516.
- [15] Ludek Cigler and Boi Faltings. 2013. Decentralized anti-coordination through multi-agent learning. *Journal of Artificial Intelligence Research* 47 (2013), 441–473.
- [16] Caroline Claus and Craig Boutilier. 1998. The dynamics of reinforcement learning in cooperative multiagent systems. *AAAI/IAAI 1998*, 746-752 (1998), 2.
- [17] Norman Dalkey. 1953. Equivalence of information patterns and essentially determinate games. *Contributions to the Theory of Games* 2 (1953), 217–244.
- [18] Gabriele Farina, Andrea Celli, Nicola Gatti, and Tuomas Sandholm. 2018. Ex ante coordination and collusion in zero-sum multi-player extensive-form games. In *Advances in Neural Information Processing Systems*. 9638–9648.
- [19] Jakob Foerster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. 2016. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*. 2137–2145.
- [20] Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [21] Andrew Gilpin and Tuomas Sandholm. 2006. A competitive Texas Hold’em poker player via automated abstraction and real-time equilibrium computation. In *Proceedings of the National Conference on Artificial Intelligence*, Vol. 21. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 1007.
- [22] Andrew Gilpin and Tuomas Sandholm. 2007. Better automated abstraction techniques for imperfect information games, with application to Texas Hold’em poker. In *Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems*. 1–8.
- [23] Amy Greenwald and Keith Hall. 2003. Correlated Q-Learning. In *ICML*. 242–249. <http://www.aaai.org/Library/ICML/2003/icml03-034.php>

- [24] Matthew Hausknecht and Peter Stone. 2015. Deep Recurrent Q-Learning for Partially Observable MDPs. In *AAAI Fall Symposium on Sequential Decision Making for Intelligent Agents (AAAI-SDMIA15)* (Arlington, Virginia, USA).
- [25] Johannes Heinrich, Marc Lanctot, and David Silver. 2015. Fictitious self-play in extensive-form games. In *International Conference on Machine Learning*, 805–813.
- [26] Johannes Heinrich and David Silver. 2016. Deep reinforcement learning from self-play in imperfect-information games. *arXiv preprint arXiv:1603.01121* (2016).
- [27] Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. 2019. A survey and critique of multiagent deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems* 33, 6 (2019), 750–797.
- [28] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. 2019. Human-level performance in 3D multiplayer games with population-based reinforcement learning. *Science* 364, 6443 (2019), 859–865.
- [29] Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. 2017. Population based training of neural networks. *arXiv preprint arXiv:1711.09846* (2017).
- [30] Mamoru Kaneko and J Jude Kline. 1995. Behavior strategies, mixed strategies and perfect recall. *International Journal of Game Theory* 24, 2 (1995), 127–145.
- [31] Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*.
- [32] J Jude Kline. 2002. Minimum memory for equivalence between ex ante optimality and time-consistency. *Games and Economic Behavior* 38, 2 (2002), 278–305.
- [33] Daphne Koller, Nimrod Megiddo, and Bernhard Von Stengel. 1996. Efficient computation of equilibria for extensive two-person games. *Games and economic behavior* 14, 2 (1996), 247–259.
- [34] HW Kuhn. 1953. Extensive Games and the Problem of Information, Contributions to the Theory of Games II. *Annals of Mathematics Studies* 28 (1953), 193–216.
- [35] Marc Lanctot, Richard Gibson, Neil Burch, Martin Zinkevich, and Michael Bowling. 2012. No-regret learning in extensive-form games with imperfect recall. In *Proceedings of the 29th International Conference on Machine Learning*, 1035–1042.
- [36] Martin Lauer and Martin Riedmiller. 2000. An algorithm for distributed reinforcement learning in cooperative multi-agent systems. In *In Proceedings of the Seventeenth International Conference on Machine Learning*. Citeseer.
- [37] Martin Lauer and Martin Riedmiller. 2004. Reinforcement learning for stochastic cooperative multi-agent systems. In *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 3*. Citeseer, 1516–1517.
- [38] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. 2017. Multi-agent reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 464–473.
- [39] David S Leslie and Edmund J Collins. 2006. Generalised weakened fictitious play. *Games and Economic Behavior* 56, 2 (2006), 285–298.
- [40] Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gonzalez, Michael Jordan, and Ion Stoica. 2018. RLlib: Abstractions for Distributed Reinforcement Learning (*Proceedings of Machine Learning Research, Vol. 80*), Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 3053–3062. <http://proceedings.mlr.press/v80/liang18b.html>
- [41] Siqi Liu, Guy Lever, Nicholas Heess, Josh Merel, Saran Tunyasuvunakool, and Thore Graepel. 2019. Emergent Coordination Through Competition. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=BkG8sjR5Km>
- [42] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*. 6379–6390.
- [43] Laëtitia Matignon, Guillaume Laurent, and Nadine Le Fort-Piat. 2008. A study of FMQ heuristic in cooperative multi-agent games.
- [44] Laëtitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. 2007. Hysteretic q-learning: an algorithm for decentralized reinforcement learning in cooperative multi-agent teams. In *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 64–69.
- [45] Laetitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. 2012. Independent reinforcement learners in cooperative markov games: a survey regarding coordination problems. *The Knowledge Engineering Review* 27, 1 (2012), 1–31.
- [46] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. 2013. Playing Atari with Deep Reinforcement Learning. *CoRR abs/1312.5602* (2013). [arXiv:1312.5602](http://arxiv.org/abs/1312.5602)
- [47] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529.
- [48] Dov Monderer and Lloyd S Shapley. 1996. Fictitious play property for games with identical interests. *Journal of economic theory* 68, 1 (1996), 258–265.
- [49] Akira Okada. 1987. Complete inflation and perfect recall in extensive games. *International Journal of Game Theory* 16, 2 (1987), 85–91.
- [50] Shayegan Omidshafiei, Jason Pazis, Christopher Amato, Jonathan P How, and John Vian. 2017. Deep decentralized multi-task multi-agent reinforcement learning under partial observability. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2681–2690.
- [51] Gregory Palmer, Karl Tuyls, Daan Bloembergen, and Rahul Savani. 2018. Lenient Multi-Agent Deep Reinforcement Learning. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems* (Stockholm, Sweden) (AAMAS ’18). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 443–451. <http://dl.acm.org/citation.cfm?id=3237383.3237451>
- [52] Liviu Panait and Sean Luke. 2005. Cooperative multi-agent learning: The state of the art. *Autonomous agents and multi-agent systems* 11, 3 (2005), 387–434.
- [53] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*. 8026–8037.
- [54] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholm, Sweden, 4295–4304. <http://proceedings.mlr.press/v80/rashid18a.html>
- [55] Julia Robinson. 1951. An iterative method of solving a game. *Annals of mathematics* (1951), 296–301.
- [56] Sheldon M Ross. 1970. *Goofspiel—the game of pure strategy*. Technical Report. CALIFORNIA UNIV BERKELEY OPERATIONS RESEARCH CENTER.
- [57] Tuomas Sandholm. 2010. The state of solving large incomplete-information games, and application to poker. *Ai Magazine* 31, 4 (2010), 13–32.
- [58] Yoav Shoham and Kevin Leyton-Brown. 2008. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press.
- [59] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature* 529, 7587 (2016), 484.
- [60] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharrshan Kumaran, Thore Graepel, et al. 2018. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science* 362, 6419 (2018), 1140–1144.
- [61] Sainbayar Sukhbaatar, Rob Fergus, et al. 2016. Learning multiagent communication with backpropagation. In *Advances in Neural Information Processing Systems*. 2244–2252.
- [62] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Flores Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. 2017. Value-Decomposition Networks For Cooperative Multi-Agent Learning. *CoRR abs/1706.05296* (2017). [arXiv:1706.05296](http://arxiv.org/abs/1706.05296)
- [63] Ming Tan. 1993. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*. 330–337.
- [64] Bernhard von Stengel and Daphne Koller. 1997. Team-Maxmin Equilibria. *Games and Economic Behavior* 21, 1 (1997), 309 – 321.
- [65] Chongjie Zhang and Victor Lesser. 2013. Coordinating the multi-agent reinforcement learning with limited communication. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*. International Foundation for Autonomous Agents and Multiagent Systems, 1101–1108.