



# On the Combination of Game-Theoretic Learning and Multi Model Adaptive Filters

Michalis Smyrnakis<sup>1</sup>(✉) , Hongyang Qu<sup>2</sup> , Dario Bauso<sup>3,4</sup> ,  
and Sandor Veres<sup>2</sup> 

<sup>1</sup> Science and Technology Facilities Council, Daresbury, UK  
`michail.smyrnakis@stfc.ac.uk`

<sup>2</sup> Department of Automatic Control and Systems Engineering,  
University of Sheffield, Sheffield, UK  
`{h.qu,s.veres}@sheffield.ac.uk`

<sup>3</sup> Jan C. Willems Center for Systems and Control ENTEG,  
Faculty of Science and Engineering, University of Groningen Nijenborgh, Groningen,  
The Netherlands

<sup>4</sup> Dipartimento di Ingegneria, Università di Palermo,  
Viale delle Scienze, Palermo, Italy  
`d.bauso@rug.nl`

**Abstract.** This paper casts coordination of a team of robots within the framework of game theoretic learning algorithms. In particular a novel variant of fictitious play is proposed, by considering multi-model adaptive filters as a method to estimate other players' strategies. The proposed algorithm can be used as a coordination mechanism between players when they should take decisions under uncertainty. Each player chooses an action after taking into account the actions of the other players and also the uncertainty. Uncertainty can occur either in terms of noisy observations or various types of other players. In addition, in contrast to other game-theoretic and heuristic algorithms for distributed optimisation, it is not necessary to find the optimal parameters a priori. Various parameter values can be used initially as inputs to different models. Therefore, the resulting decisions will be aggregate results of all the parameter values. Simulations are used to test the performance of the proposed methodology against other game-theoretic learning algorithms.

**Keywords:** Game-theoretic learning · Distributed optimisation · Multi-model adaptive filters · Robot teams coordination · Fictitious play · Bayesian games · Potential games · State based games · Stochastic games

## 1 Introduction

Teams of robots can be used in many domains such as mine detection [69], medication delivery in medical facilities [16], formation control [45, 58] and exploration of unknown environments [31, 55]. A common feature shared by these

applications is that robots should either minimise a cost function or maximise a utility function in a distributed fashion. Thus, the resulting problem can be formulated as an distributed optimization one. Distributed optimisation arises also in several applications such as in smart grids [2,66], disaster management [25,57], robot team coordination [51,52], sensor networks [23,24,32,68], water distribution system optimisation [2,66] and scheduling problems [61].

In each of the aforementioned applications, the agents need to coordinate to achieve a common goal. If the desired task requires distributed optimisation of a utility or cost function, then the resulting problem turns into a game where each agent optimizes a portion of the common objective function based on local information. In such a scenario, game theory provides formal tools to assess the quality of the solution obtained.

As in [59], in this paper we address two kinds of uncertainty that can be arise in a game theoretic learning process that can be applied in robotic scenarios. The first one is related to uncertainty of measurements. Consider the case where, the robots do not have access to the other robots' states, or if the communication channel is noisy or involves faulty sensors, we say that the game has imperfect information, and the robots have to make a decision under uncertainty. Uncertainty can lead to wrong decision. For example, wrong decisions could be made when the positions of other robots are inferred by noisy odometry or noisy camera input. Another example in which noisy observation can have impact on the coordination process of a robot team is the case of a fleet of Unmanned Aerial Vehicles (UAVs) that need to take images of an area in order to identify the growth of the crops. The images should be taken from various angles. However, it is not always possible for a UAV to know the exact position and bearing of the other UAVs, and therefore, to make correct decisions about when to change photo-shooting angles.

The second form of uncertainty is related with various states or types that the environment or other robots can be. This uncertainty can be in the form of incomplete information. A team of UAVs can know the weather conditions for their tasks, "good" or "bad" weather for example, up to a certain probability. Consider also the case where, each member of a robotic team can be in various possible states regarding its battery life, affecting the possible choices they have. On another setup the choices robots are making in a current snapshot of their mission could affect their future choices. Stochastic, state based and Bayesian games are two categories of games which can be used in order to model distributed optimisation task under those types of uncertainty.

In this paper, we propose a novel game-theoretic learning algorithm, which can be used as a coordination mechanism among robots playing either complete information games with noisy observations or Bayesian games or Markov games. In detail, it is a synchronous algorithm, where Extended Kalman Filters Fictitious Play (EKFFP) [60] is combined with multi-model adaptive filters (MMAFs) [9]. The novelty of our algorithm is that the joint distribution of the uncertainty and the observed actions of other players' action are used to make decisions. Robots use multiple models to solve their optimisation task.

Each model is either a probabilistic representation of the noisy observations model of the other players' states or of the state of the world. Each model of MMAF represents a part of uncertainty and the final decision making is based on a weighted average over all the models.

Note here that a sequence of Bayesian games can also be used to describe partially observable games and decentralised partial observable Markov decision process (dec-POMDPs) [15]. Thus the proposed methodology can be used as a coordination mechanism for each individual Bayesian game (sub-tasks) till a solution to the dec-POMDP is found.

Another advantage of our algorithm, is that in contrast to EKFFP and various heuristic distributed optimisation algorithms, there is no need to tune any parameters. Therefore, there is no need to decide in advance the value of the EKFFP parameters. Instead, random valuations of these parameters can be used simultaneously. Each valuation predicts other robots' strategy from a different angle, which is represented by different models. Therefore, it is easier to adapt to evolution of other robots' strategy.

This article is an extended version of [59]. This version considers more types of uncertainty which are related with two new categories of games, namely stochastic and state based games. In addition, the basic game theoretic definitions and learning algorithms were explained in more detail. Also, additional information is provided regarding the implementation details of the proposed algorithm.

The rest of the paper is organised as follows. Next section contains a brief description of related work. In Sect. 3, a brief description of the game-theoretic notions that will be used in the rest of the paper are presented. In Sect. 4, the learning algorithm EKFFP is presented. Section 5 describes our fictitious play based algorithm, which integrates multi-model adaptive filters and extended Kalman filter. Section 6 discusses some implementation details of our algorithm and game-theoretic learning algorithms in general. In Sect. 7, we evaluate our algorithm in several case studies. In Sect. 8, we summarise our findings and present our future work.

## 2 Related Work

Distributed learning under noisy observation was considered in [27]. Particle swarm algorithms subjected to intrinsic noise was applied in [42] and [13]. In [21], noisy observations in a different context from this work were investigated, as no directly any knowledge about the noisy observation was used, such as their probability distribution. In [11], the uncertainty was dealt within a game theoretic framework under a simplified assumption that players use the same strategy through the iterations of the game.

Various approaches have been adopted to solve Bayesian games including: Bayesian action graph games [22], Multi-agent influence diagrams [26] and Newton method [19]. The difference between these approaches and the proposed algorithm is that their goal is to find the optimal strategy. However, the search for an optimal solution in cooperative Bayesian games is an NP-hard problem

[65], and thus, is not tractable. On the other hand, in our algorithm players take into account the other players' actions and their possible types when updating their desired action until they reach to a commonly accepted solution, which is usually an equilibrium point to the problem. In [15], an approximate solution was proposed based on an alternating maximisation algorithm, but this is not applicable when robots choose their actions simultaneously. Smooth fictitious play [43] and a variant of fictitious play for Bayesian games with continuous states [44] have been used to solve auctions in a competitive environment, which is not applicable in cooperative games on which this work is focused.

Various reinforcement learning techniques have been proposed as solutions to stochastic games. Examples include minimax Q and Q-learning in [30], decentralized Q-learning, distributed Q-learning, hysteretic Q-learning, and WOLF PHC in [34]. In [54] Classic fictitious play has also been used in order to solve a Markov games that used in order to model robotic arm manipulators. In [33] a log-linear process was proposed to solve state based games. When in [29] a two-memory better reply algorithm was proposed. In contrast, the proposed algorithm since a different model is used for each possible state of the game. Moreover, the model of each state is updated only when the state is active which allows to have a more accurate estimate of opponents' strategies in each state.

### 3 Game-Theoretic Definitions

This section contains a brief description of some game theoretic definitions that will be used in the rest of the paper.

#### 3.1 Normal Form Games

A game  $\Gamma$  in normal form is defined as a tuple

$$\Gamma = \langle \mathcal{I}, \{A^i\}_{i \in \mathcal{I}}, \{r^i\}_{i \in \mathcal{I}} \rangle,$$

where

- $\mathcal{I}$  is the set of indices of all players;
- $A^i$  is the set of all possible actions of player  $i$  and the set product  $A = \times_{i \in \mathcal{I}} A^i$  is the set of all joint actions;
- $r^i : A \rightarrow \mathbb{R}$  is the utility (reward) function of player  $i$ , which computes the reward that the player gains after a joint action is selected.

Joint action  $a$ , can be written as  $a = (a^i, a^{-i})$ , where  $a^{-i}$  is the joint action of all players but  $i$ . A strategy of player  $i$  is a probability distribution over its action space, and let  $\Delta^i$  denote the set of all the probability distributions over  $A^i$ . Each player uses a strategy  $\sigma^i \in \Delta^i$  to choose its action. Similarly to actions, a joint strategy  $\sigma \in \Delta$  is defined as an element of the set product  $\Delta = \times_{i \in \mathcal{I}} \Delta^i$ , and  $\sigma^{-i}$  is a joint strategy of all players but  $i$ .

In this paper we consider iterative games. In these games, a game is repeatedly played along a discrete sequence of time instances called rounds or iterations

when each player chooses their actions based on their strategies, the history of the observed joint actions in the played iterations of the game and the reward allocated to them.

A player uses a *pure* strategy when it deterministically chooses actions, therefore it puts all its mass function in a single action  $a^i \in A^i$  such that  $\sigma^i(a^i) = 1$ . When this is not the case, we call such a strategy, *mixed* strategy. The expected reward of player  $i$ , given its opponents' strategies  $\sigma^{-i}$ , is denoted by  $r^i(\sigma^i, \sigma^{-i})$ . If  $\sigma^i$  is a pure strategy with  $\sigma^i(a^i) = 1$ , the expected reward can be written as  $r^i(a^i, \sigma^{-i})$ .

The most common deterministic decision rule in game theory is the so-called *best response (BR)* by which players choose the actions which maximise their expected rewards. Formally, the action that a player  $i$  will choose, given its opponents strategies  $\sigma^{-i}$ , is

$$BR^i(\sigma^{-i}) = \operatorname{argmax}_{a^i \in A^i} r^i(a^i, \sigma^{-i}). \quad (1)$$

A normal form game  $\Gamma$  can be either a competitive or a coordination game, depending on its utility function. In competitive games, players have conflicted interests, while in coordination games, they maximise their reward when a common goal is achieved. Table 1 depicts the players' rewards of a zero-sum game, which is an canonical example of competitive games. There are two players in this game: the row player playing actions  $b^1$  and  $b^2$  and the column player playing  $a^1$  and  $a^2$ . Each entry in the table represents the reward they receive when the corresponding joint action is played. For example, entry 1, -1 means that the row player receives 1 and the column player receives -1 when they play joint action  $(b^1, a^1)$ . Therefore, the reward a player gains is what the other player loses. Similarly, Table 2 presents the rewards of a coordination game where both players receive the same reward.

**Table 1.** Zero sum game.

	$a^1$	$a^2$
$b^1$	1,-1	-1,1
$b^2$	-1,1	1,-1

**Table 2.** Coordination game.

	$a^1$	$a^2$
$b^1$	1,1	0,0
$b^2$	0,0	2,2

A joint strategy  $\tilde{\sigma} = (\tilde{\sigma}^i, \tilde{\sigma}^{-i})$  that satisfies

$$r^i(\tilde{\sigma}^i, \tilde{\sigma}^{-i}) \geq r^i(\sigma^i, \sigma^{-i}) \quad \forall i \in \mathcal{I}, \forall \sigma^i \in \Delta^i$$

is a Nash equilibrium [38]. Nash in [38] showed that every game has at least one equilibrium, i.e., there is at least one strategy  $\tilde{\sigma}$  where players do not benefit from deviating from it unilaterally. A Nash equilibrium can be either mixed or pure if the strategy  $\tilde{\sigma}$  is a mixed or a pure strategy respectively.

A class of games of particular interest is potential games. In [36], it was shown that distributed optimisation tasks can be cast as potential games. Hence, the search of an optimal solution in a distributed optimisation problem can be seen as searching for a Nash equilibrium in a potential game. A game  $\Gamma$  is a potential game if the rewards of all players can be replaced by a potential function  $\phi$  such that  $\forall a = (a^i, a^{-i}) \in A$  and  $\forall \bar{a} = (\bar{a}^i, \bar{a}^{-i}) \in A$ :

$$r^i(a^i, a^{-i}) - r^i(\bar{a}^i, \bar{a}^{-i}) = \phi(a^i, a^{-i}) - \phi(\bar{a}^i, \bar{a}^{-i}).$$

[36] showed that every potential game has at least one pure Nash equilibrium.

### 3.2 Bayesian Games

A Bayesian game, or game of incomplete information, is defined as a tuple

$$\mathcal{G} = \langle \mathcal{I}, \{\Theta^i\}_{i \in \mathcal{I}}, \{A^i\}_{i \in \mathcal{I}}, \{p(\theta^i)\}_{\theta^i \in \Theta, i \in \mathcal{I}}, \{r^i\}_{i \in \mathcal{I}} \rangle,$$

where

- $\mathcal{I}$  is the set of player indices;
- $\Theta^i$  is the set of types belonging to player  $i$  and  $\Theta = \times_{i \in \mathcal{I}} \Theta^i$ ;
- $A^i$  is set of possible actions of player  $i$ ;
- $r^i : A \rightarrow \mathbb{R}$  is the utility function of player  $i$ .

Each type of a player represents a possible internal state of the player. At any time, a player can only be in one of its types. The type of a player constitutes its private information, in the sense that each player  $i$  knows the type that it is in. In contrast, the other players only know the probability that player  $i$  can be in a certain type at that moment. If  $\theta^i \in \Theta^i$  is considered as the state of the environment, i.e., the type of a player is the state of the environment (world), then all players have the same types and thus  $\Theta^i = \Theta^j$ ,  $\forall i, j \in \mathcal{I}$ . The expected reward of a player in a Bayesian game is then estimated as:

$$r^i(\sigma^i, \theta^i) = \sum_{\theta^{-i} \in \Theta^{-i}} p(\theta^{-i}) r^i(\sigma^i, \theta^i, \sigma^{-i}, \theta^{-i}). \quad (2)$$

A Bayesian Nash equilibrium is defined as

$$\sigma^i \in \operatorname{argmax}_{a^i \in A} p(\theta^i | \theta^{-i}) r^i(\sigma^i, \theta^i, \sigma^{-i}, \theta^{-i})$$

Hence a Bayesian Nash equilibrium is a Nash equilibrium of the expanded game in which each player action space of pure strategies is the set of maps from  $\Theta^i$  to  $A^i$ .

As an example, consider a task allocation scenario: two robots 1 and 2 need to collaborate to finish two tasks: *easy* and *difficult*. The *difficult* task can be performed efficiently only if both robots work together on it. Robot 1 can do both tasks at the same efficiency, and hence it has only one type. Robot 2 has two types *A* and *B*. In type *A*, robot 2 can do the *easy* task with greater efficiency than the difficult task, while in type *B*, it performs both tasks with the same efficiency. Furthermore, robot 2 always knows its type, while robot 1 only knows that with probability  $p$  robot 2 is in type *A*, and with probability  $1 - p$  is in type *B*. In this game, the world has a unique state and thus does not affect robots' types. Table 3 illustrates an example of the utility function in this Bayesian game.

**Table 3.** Reward matrices of a two players Bayesian game.

Type A			Type B		
task	difficult	easy	task	difficult	easy
difficult	10,6	5,10	difficult	9,10	5,4
easy	8,5	6,8	easy	8,6	7,5

Each matrix of Table 3 represents the rewards that robots receive. The single type robot is the row player of the game, and the other robot is the column player. If robot 2 is in type *A*, then both robots receive the rewards in the left matrix, while when it is type *B*, they receive the reward in the right matrix. For example, the entry 5, 10 of the left matrix means that when robot 1 performs the *difficult* task and robot 2 does the *easy* task and is in type *A*. In this case, robot 1 receives 5 unit of reward and robot 2 receives 10 units. We reinforce here that at any time of game playing, robot 2 knows exactly which reward matrix is used, while robot 1 makes decisions based on the probability distribution of types of robot 2:  $p$  and  $1 - p$  for the left and right matrix respectively.

### 3.3 Stochastic Games

Stochastic games [53] is considering as a multi-agent extension of Markov Decision Processes. In these games the joint actions of the players can change the game that is played. Consider the following case which describes a stochastic game. Two players play repeatedly the zero sum game that is depicted in Table 1, but when the joint actions  $b^1, a^1$  or  $b^2, a^2$  is played then the game that the two players will play will change with probability  $p$  to the coordination game that is described in Table 2 and will play this game for the remaining repetitions of the game that should be played.

More formally a stochastic game  $G$  is defined as a tuple  $G = (Q, \mathcal{I}, A, P, \{r^i\}_{i \in \mathcal{I}})$ , where

- $Q$  is a finite set of states
- $\mathcal{I}$  is a finite set of players
- $A = \times_{i \in \mathcal{I}} A^i$  is a finite set of joint actions available for each state in  $Q$ .
- $P$  is a  $|Q| \times |A| \times |Q|$  transition probabilities matrix with  $|Q|$  and  $|A|$  denoting the cardinality of  $Q$  and  $A$  respectively. Thus,  $P(q, a, \tilde{q})$  will denote the probability that a state  $\tilde{q}$  will be visited from state  $q$  when the joint action  $a$  will be played.
- $r^i : Q \times A \rightarrow \mathbb{R}$  is the utility function of player  $i$ , with  $r^i(q, \{a_i\}_{i \in \mathcal{I}})$  denoting the reward that a player will gain if a joint action is played while the game is in state  $q$ .

The expected reward that players try to maximise in a stochastic game given a discount factor  $\beta \in [0, 1)$  and a mixed strategy  $\sigma$  is defined as:

$$\tilde{r}^i(q, \sigma) = \sum_{t=0}^{\infty} \beta^t \mathbb{E}(r^i(q, a) | q_0 = q).$$

The Nash equilibrium in the case of the stochastic games is defined as a joint strategy  $\tilde{\sigma}^i$  such as:

$$\tilde{r}^i(q, \tilde{\sigma}) \geq \tilde{r}^i(q, (\sigma^i, \tilde{\sigma}^{-i})) \forall \sigma^i \in \Delta^i \forall i \in \mathcal{I}.$$

As an example consider a simple scenario where two UAVs that will have to choose between two areas to monitor, area  $A$  and area  $B$ . Area  $A$  is a narrow and if both robots choose to monitor it there is a chance to crash. Therefore, it is possible that the robots wont be able to monitor any area for the rest of their mission. This problem can be formulated as a stochastic game with the following elements:  $Q = x, y$ , where state  $x$  represents the state that UAVs are functioning and  $y$  the state that the UAVs have crashed. The actions that the UAVs have in state  $x$  are  $A$  and  $B$  and they have no action in state  $y$  since they have crashed. The rewards of the robots are presented in Tables 4 and 5 for the states  $x$  and  $y$  respectively.

The transition probabilities when the game is in state  $x$  for each possible joint action is given by

$$P_x = \begin{matrix} & \begin{matrix} x & y \end{matrix} \\ \begin{matrix} (A,A) \\ (A,B) \\ (B,A) \\ (B,B) \end{matrix} & \begin{bmatrix} \frac{2}{3} & \frac{1}{3} \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \end{bmatrix} \end{matrix}.$$



**Table 4.** Reward of both players in state  $x$ . The left number of each tuple represents the rewards of the first robot (row player), and the second the rewards of the second robot (column player).

	A	B
A	(10,9)	(6,5)
B	(1,7)	(3,8)

**Table 5.** Reward of both players in state  $y$ .

	No action
No action	(0,0)

Each row of  $P_x$  represents a joint action and each column represents the possible new state of the game either  $x$  or  $y$ . When the game is on state  $y$  will remain in that state for ever. Under this set up this stochastic game has a pure strategy Markov equilibrium the joint action  $(A, B)$  which wouldn't be the equilibrium of the game if the game had only a single state state  $x$ .

**State-Based Games.** State-based games [33] are a simplified version of stochastic games. Similarly to stochastic games are defined as  $G = (Q, \mathcal{I}, A, P, \{r^i\}_{i \in \mathcal{I}})$ , but the players are myopic and thus they aim to maximise their current expected reward instead of discounted future rewards. Another important difference between stochastic and state-based games is that the in state based games the players are expected to have the same action in all the different states of the game. In these games the notion of recurrent state equilibrium have been proposed [33], which is defined as the tuple of joint action and state  $(a^*, q^*)$  such as:

- $q^* \in Q(a^*|q) \forall q \in Q(a^*|q^*)$ , where  $Q(a^*|q) \subseteq Q$  is the set of reachable states from initial state  $q$  when the joint action  $a$  is always selected.
- $r^i(q, a^*) \geq r^i(q, a^i a^{*, -i}) \forall a^i \in A^i, \forall q \in Q(a^*|q^*)$ ,  $\forall i \in \mathcal{I}$ .

Consider the case where they should choose to cooperate or not in order to perform a task. Each time that they choose an action the state of the game can change. An example is the a state-based game presented in [29]. There are three states  $x, y$  and  $z$  and in each of them corresponds a different game, namely coordination game, prisoners' dielemma and maching pennies. The reward functions of these games are depicted in Tables.

**Table 6.** Reward of both players for the coordination game. The left number of each tuple represents the rewards of the first robot (row player), and the second the rewards of the second robot (column player).

	Cooperate	Not cooperate
Cooperate	(4,4)	(1,3)
Not cooperate	(3,1)	(2,2)

**Table 7.** Reward of both players for prisoners' dilemma. The left number of each tuple represents the rewards of the first robot (row player), and the second the rewards of the second robot (column player).

	Cooperate	Not cooperate
Cooperate	(2,2)	(0,3)
Not cooperate	(3,0)	(1,1)

**Table 8.** Reward of both players for matching pennies game. The left number of each tuple represents the rewards of the first robot (row player), and the second the rewards of the second robot (column player).

	Cooperate	Not cooperate
Cooperate	(-1,1)	(1,-1)
Not cooperate	(1,-1)	(-1,1)

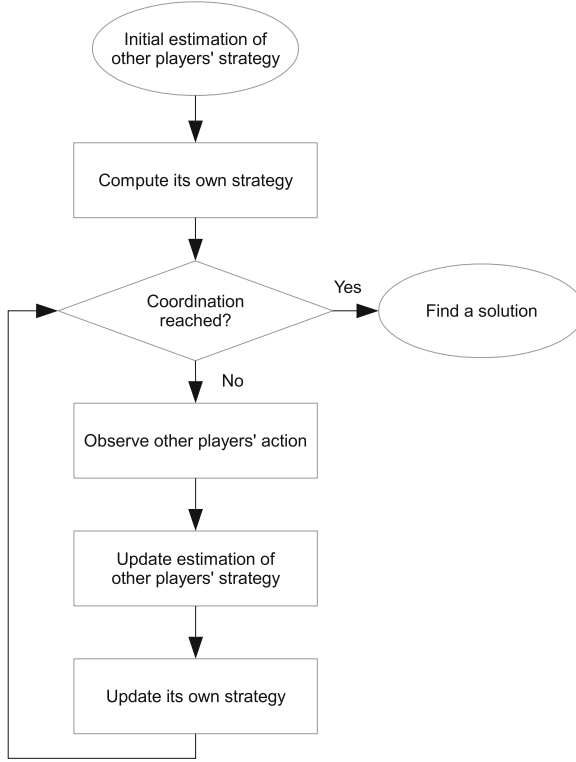
The transition probabilities from each state are given by the following transition matrices:

$$P_x = \begin{matrix} & \begin{matrix} x & y & z \end{matrix} \\ \begin{matrix} (A,A) \\ (A,B) \\ (B,A) \\ (B,B) \end{matrix} & \begin{bmatrix} 0 & \frac{1}{3} & \frac{2}{3} \\ 0 & 0 & 1 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix} \end{matrix}, P_y = \begin{matrix} & \begin{matrix} x & y & z \end{matrix} \\ \begin{matrix} (A,A) \\ (A,B) \\ (B,A) \\ (B,B) \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & \frac{1}{4} & \frac{3}{4} \\ 0 & 1 & 0 \\ \frac{3}{5} & \frac{2}{5} & 0 \end{bmatrix} \end{matrix}, P_z = \begin{matrix} & \begin{matrix} x & y & z \end{matrix} \\ \begin{matrix} (A,A) \\ (A,B) \\ (B,A) \\ (B,B) \end{matrix} & \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix} \quad (3)$$

where action  $A$  and  $B$  correspond to actions cooperate and not cooperate respectively. The recurrent equilibria for this state-based game are the action state pairs  $(BB, x)$  and  $(BB, y)$

## 4 Learning Algorithms

A distributed optimisation task can be cast as a game [67]. However, the formulation of the optimisation task as a game does not directly provide a solution to the game. A coordination mechanism between the robots is needed especially in cases where autonomy is a desirable property of the robot team. Game-theoretic learning algorithms can be used by robots to choose a joint action to solve the



**Fig. 1.** General procedure of game-theoretic learning algorithms.

game. The canonical example of game-theoretic learning algorithms is fictitious play (FP). Fictitious play is an iterative learning algorithm. Figure 1 illustrates the general procedure of game-theoretic learning algorithms.

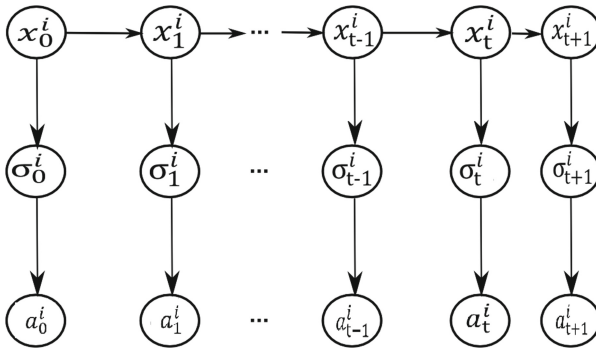
In each iteration  $t$ , each player estimates other players' strategies, and based on these estimates, chooses an action using the best response decision rule. At the initial iteration, i.e.,  $t = 0$ , every player  $i$  maintains some arbitrary, non-negative weights  $\kappa_t^{i \rightarrow j}$  for each other player  $j$  as the estimation of their strategy. In particular,  $\kappa_t^{i \rightarrow j}(a^j)$  is the weight for action  $a^j \in A^j$  of player  $j$ . At successive iterations, players update their weight functions based on other players' chosen actions. The update of player  $i$ 's weight function for player  $j$  is computed as follows [18]:

$$\kappa_t^{i \rightarrow j}(a^j) = \kappa_{t-1}^{i \rightarrow j}(a^j) + \begin{cases} 1 & \text{if } a^j = a_{t-1}^j \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

where  $a_{t-1}^j$  is the action that player  $j$  chooses at iteration  $t - 1$ . Based on these weights, player  $i$  then estimates player  $j$ 's strategy using the following equation:

$$\sigma_t^j(a^j) = \frac{\kappa_t^{i \rightarrow j}(a^j)}{\sum_{a^j \in A^j} \kappa_t^{i \rightarrow j}(a^j)}. \quad (5)$$

Fictitious play converges to the Nash equilibrium in many classes of games, such as  $2 \times 2$  games with generic payoffs [35], zero sum games [46], games that can be solved using iterative dominance [37],  $2 \times n$  games [4] and potential games [36]. However, this convergence can be very slow [18] because of the implicit assumption that all players use the same strategy throughout the game. In [56] and [60], variants of fictitious play were proposed, which were based on particle filters and extended Kalman filters respectively. These algorithms assume that players adapt their strategies through the iterations of the game. In both variants, the fictitious play process is described as a hidden Markov model (HMM). Each player maintains some unconstrained propensities<sup>1</sup>, which are responsible for their strategies. In each iteration of game playing, each player aims to predict other players' propensities, i.e., hidden layer of the HMM, by using the history of other players' actions, i.e., observations layer of the HMM. Figure 2 illustrates the evolution of propensities of player  $i$  through the iterations of the game and how they are related to strategies and actions.



**Fig. 2.** Propensities propagation throughout the game. The propensity  $x_t^i$  at time  $t$  depends only on the propensity at time  $t - 1$ . Moreover, the strategy  $\sigma_t^i$  of player  $i$  at time  $t$  depends only on the propensity of the same iteration, and the action  $a_t^i$  that the player chooses depends only on its strategy at the same iteration as well.

More formally, let  $x^i(a^i)$  denote the propensity of player  $i$  to play action  $a^i$ , and  $a_t^i$  the action of player  $i$  at the  $t$ -th iteration of the fictitious play process. The probability at which each player estimates about the propensity of other players is

$$p(x^j(a^j)|a_0^j, a_1^j, \dots, a_t^j) \quad \forall j \in \mathcal{I} \setminus \{i\}, \forall a^j \in A^j. \quad (6)$$

This probability was evaluated using particle filters in [56] and extended Kalman filters in [60]. In this paper, we only consider the variant of fictitious

<sup>1</sup> The strategies are probability distributions. Thus, when a dynamical model is used to propagate them, new estimates are not necessary to lay in the probability distributions space. For that reason, the intentions of players to choose an action, namely propensities, which are not bounded to probability distribution spaces, are used [56].

play based on extended Kalman filters. But the same methodology can be easily applied to the variant with particle filters. As each player estimates the propensity of every other individual player separately, only inference over a single “opponent” player, say player  $j$ , will be presented in the rest of the paper.

#### 4.1 Extended Kalman Filter Fictitious Play (EKFFP)

This variant of fictitious play is based on the assumption that players have no prior knowledge about other players’ strategies, and thus, an autoregressive model can be used to propagate the propensities [56]. In addition, inspired from the sigmoid functions that are used in neural networks to connect the weights and the observations, a Boltzman formula is used to relate the propensities with other players’ strategies [5]. The following state space model is used to describe EKFFP:

$$\begin{aligned} x_t^j(a^j) &= x_{t-1}^j(a^j) + \xi_{t-1}^j \\ I_{a_t^j=a^j}(a^j) &= h(x_t^j(a^j)) + \zeta_t^j \end{aligned} \tag{7}$$

where  $I_{a_t^j=a^j}(a^j)$  is the measurement equation, which relates the propensities to the actions of the players. The noise of the propensity process,  $\xi_{t-1}^j \sim N(0, \Xi)$ , which comprises the internal states, has zero *mean* and *covariance matrix*  $\Xi$ . The error,  $\zeta_t \sim N(0, Z)$ , of the observations has zero *mean* and *covariance matrix*  $Z$ . This error occurs because a discrete 0–1 process, such as the best response in Eq. (1) is represented through the continuous Boltzmann formula  $h(\cdot)$  in which  $\tau$  is a “temperature parameter”. The components of the vector  $h$ , are evaluated as:

$$h(x^j(a^j)) = \frac{\exp(x^j(a^j)/\tau)}{\sum_{a^k \in A^j} \exp(x^j(a^k)/\tau)}. \tag{8}$$

The behaviour of the EKFFP algorithm at the  $t$ -th iteration of a game can be described as follows. At first, player  $i$  uses the EKF process, which is based on the state space in Eq. (7), to predict other players’ propensities. Player  $i$  then using these estimates, evaluates player  $j$ ’s strategy of choosing an action  $a^j \in A^j$ ,  $\sigma_t^j(a^j)$ , as follows:

$$\sigma_t^j(a^j) = \frac{\exp(\bar{x}_t^j(a^k)/\tau)}{\sum_{a^k \in A^j} \exp(\bar{x}_t^j(a^k)/\tau)}, \tag{9}$$

where  $\bar{x}_t^j(a^k)$  is player  $i$ ’s prediction of the propensities of player  $j$  in order to choose action  $a^k \in A^j$  based on the state equations in Eq. (7) and using observations up to time  $t - 1$ . Player  $i$  then uses the estimates in Eq. (9) to choose an action using best response in Eq. (1). After all players have chosen an action, they use the EKF update process to correct their estimates about other players’ strategies in the light of the recently observed actions. Then, the next iteration of EKFFP starts with  $t = t + 1$ . The EKF estimations can be computed by any standard textbook procedure, such as in [49]. Algorithm 1 summarises the fictitious play algorithm when EKF is to predict other robots’ strategies.

---

**Algorithm 1.** Extended Kalman filter fictitious play [60].
 

---

- 1: **while**  $t < \text{max iterations}$  **do**
  - 2:   **for all**  $j \in \mathcal{I} \setminus \{i\}$  **do**
  - 3:     Predict other players' propensities for the next iteration  $t + 1$  using the state equations in Equation (7).
  - 4:   Use the beliefs about other players' strategies in Equation (9) and choose an action using BR in Equation (1).
  - 5:   Observe other players' actions
  - 6:   **for all**  $j \in \mathcal{I} \setminus \{i\}$  **do**
  - 7:     Update estimates of player  $j$ 's propensities using extended Kalman Filtering to obtain  $\bar{x}_t^j(a^k)$ .
  - 8:    $t = t + 1$
- 

## 5 Multi-model Adaptive Filter EKFFP (Source: [59])

### 5.1 Multi-model Adaptive Filters

The EKFFP process requires the definition of the covariance matrices  $\Xi$  and  $Z$  for the random variables  $\xi$  and  $\zeta$  respectively. The performance of the learning algorithm is affected by the values of these covariance matrices. In [60] specific values were proposed for those covariance matrices, although these values are not optimal for all games. In this work, we propose a new approach that uses many models, each of which represents a pair of covariance matrices  $\Xi$  and  $Z$ . This approach then uses a weighted sum of these models in order to obtain an estimate of other players' propensities, instead of estimating the propensities from a single pair of covariance matrices. For Bayesian games, players can have a propensity estimate for each state of the nature or each type of other players.

The framework that allows many models to be considered under the EKFFP process is multiple model adaptive filters [6, 9, 10]. Let  $L$  be the set of all models that are used. Instead of estimating the propensity in Eq. (6), each player should estimate

$$p(x^j(a^j), l | a_0^j, a_1^j, \dots, a_t^j), \quad (10)$$

where  $l \in L$  is one of the possible models, each of which either refers to a pair of covariance matrices for potential games, or the state of the nature or another player's type  $\Theta^i$  in Bayesian games.

To simplify notations, we use  $\tilde{a}_t^j$  to denote  $(a_0^j, a_1^j, \dots, a_t^j)$ . The estimate of other players' propensities in Eq. (10) can be written as:

$$p(x^j(a^j), l | \tilde{a}_t^j) = p(l | \tilde{a}_t^j) p(x^j(a^j) | l, \tilde{a}_t^j), \quad (11)$$

where  $p(x^j(a^j) | l, \tilde{a}_t^j)$  for a given  $l$  is the standard EKF estimate of other player's propensity and  $p(l | \tilde{a}_t^j)$  can be seen as the weight factor of each model.

Using Bayes rule, the probability  $p(l|\tilde{a}_t^j)$  can be written as:

$$p(l|\tilde{a}_t^j) = \frac{p(\tilde{a}_t^j|l)p(l)}{\sum_{l \in L} p(\tilde{a}_t^j|l)p(l)}, \quad (12)$$

where  $p(l)$  is the prior distribution of the model  $l$ . The probability  $p(\tilde{a}_t^j|l)$  can be written as

$$\begin{aligned} rlp(\tilde{a}_t^j|l) &= p(a_t^j, a_{t-1}^j, \dots, a_0^j|l) \\ &= p(a_t^j, a_{t-1}^j, \dots, a_1^j|a_0^j, l)p(a_0^j|l) \\ &\quad \vdots \\ &= p(a_t^j|\tilde{a}_{t-1}^j, l)p(a_{t-1}^j|\tilde{a}_{t-2}^j, l) \cdots p(a_0^j|l). \end{aligned} \quad (13)$$

The propensities are described by the hidden Markov model, so they are conditionally independent, and thus, Eq. (13) can be written as:

$$p(\tilde{a}_t^j|l) = \prod_{q=0}^{q=t} p(a_q^j|l). \quad (14)$$

## 5.2 Multi-model Adaptive Filters EKFFP (MMAF-EKFFP)

Let  $|L|$  denote the cardinality of set  $L$ , and  $x_{t,l}^j(a^j)$  the propensity of player  $j$ , playing action  $a^j$  at the  $t^{\text{th}}$  iteration under model  $l$ . In the multi-model adaptive filters EKFFP process, each player uses  $|L|$  models for the propensity of each other player. In particular, each model is a state model:

$$\begin{aligned} x_{t,l}^j(a^j) &= x_{t-1,l}^j(a^j) + \xi_{t-1,l}^j \\ I_{a_t^j=a^j}^l(a^j) &= h(x_{t,l}^j(a^j)) + \zeta_{t,l}^j. \end{aligned} \quad (15)$$

For each of these models, player  $i$  uses the EKF process to predict player  $j$ 's propensity, i.e.,  $\tilde{x}_{t,l}^j(a^k)$ , in order to choose an action  $a^k \in A^j$  under model  $l$ . This prediction is weighted using Eq. (12). The estimate of player  $j$ 's propensity to choose action  $a^k \in A^j$  is then the sum of the weighted estimates of each model:

$$\bar{x}_{t,a^k}^j = \sum_{l \in L} p(l|a^k) \tilde{x}_{t,l}^j(a^k), \quad (16)$$

where  $p(l|a^k)$  is evaluated using Eq. (12) and (14). Then Eq. (9) can be applied to evaluate player  $j$ 's strategy. Each model of the estimates of player  $j$ 's propensity are updated using the standard EKF process under the light of the new action player  $j$  has chosen. Algorithm 2 and Fig. 3 summarise the MMAF-EKFFP algorithm.

---

**Algorithm 2.** MMAF-EKFFP, [59].
 

---

```

1: while  $t < \text{max iterations}$  do
2:   for all  $j \in \mathcal{I} \setminus \{i\}$  do
3:     for all  $l \in L$  do
4:       For each model  $l$  predict other players' propensities for the next iteration
          $t + 1$ , using Equation (15).
5:       Evaluate  $\bar{x}_t^j(a^k)$  using Equation (16)
6:       Compute the beliefs about other players' strategies using Equation (9), and
         choose an action using BR in Equation (1)
7:       Observe other players' actions
8:     for all  $j \in \mathcal{I} \setminus \{i\}$  do
9:       for all  $l \in L$  do
10:        Update each model's estimate of other players' propensities using extended
          Kalman filter to obtain  $\bar{x}_t^j(a^k)$ 
11:    $t = t + 1$ 

```

---

## 6 Implementation Details

In this section, we discuss some implementation details in robotics of the MMAF-EKFFP algorithm and of game-theoretic learning algorithms in general.

### 6.1 Utility Functions

Utility functions have been used in robotics as a metric of robots' performance in various applications such as [8, 40, 62–64, 70]. Utility functions can incorporate various aspects of a coordination task such as the cost a robot pays to perform an action, the reward that will be produced if a task is completed successfully, and the aptness of a robot to perform a specific task. In Sect. 7.1, an example of a utility function for a task allocation problem [1] is presented.

Methodologies for designing a utility function is out of the scope of this paper. Nonetheless, we mention here that wonderful life utility (WLU) [67] is a methodology for constructing a utility function, which allows coordination tasks to be cast as potential games. An important property of potential games, especially useful in robotic applications, is that they have at least one pure Nash equilibrium. Therefore, there is at least one optimum joint action where robots will not deviate from it unilaterally.

### 6.2 Robots' Decisions

As any iterative learning algorithm, our algorithm assumes that a specific game can be iteratively be played and a final decision is reached either when the algorithm converges to an equilibrium or when the maximum number of iterations is reached. In robotics, this can be seen as the following coordination mechanism. The robots in each iteration choose an action which they intent to play and makes aware the other robots for its intentions, i.e., by communicating this



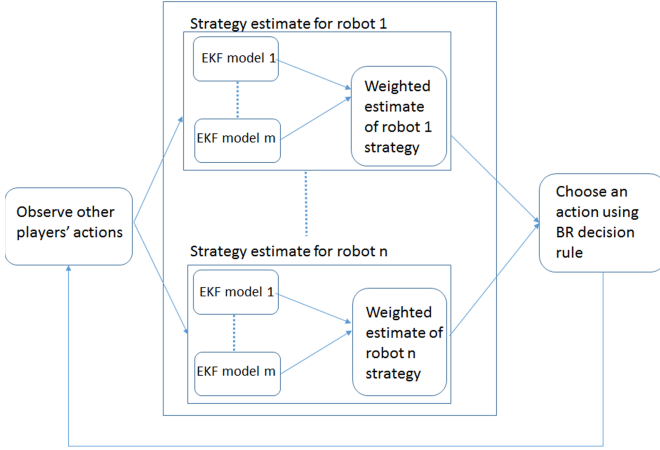


Fig. 3. The MMAF EKFP process [59].

intention to the other robots. Based on this information, they update their estimates about other players strategies and update the action they are intended to choose. The action that the team of robots will execute will be the joint action of the final iteration of the learning algorithm.

### 6.3 Complexity

The computational complexity of the EKF algorithm is upper bounded by the complexity of inverting a matrix  $\mathcal{O}(n^3)$  [7], where  $n$  is the rank of the inverted matrix. Therefore the additional computational complexity of EKFFP when it is compared with classic FP is upper bounded by  $\mathcal{O}((|\mathcal{I}| - 1)|A^k|^3)$ , where  $A^k$  ( $k \in \mathcal{I}$ ) is the largest set of actions among all players. Similarly, for MMAF-EKFFP it is  $\mathcal{O}(|\mathcal{M}|(|\mathcal{I}| - 1)|A^k|^3)$ , where  $|\mathcal{M}|$  denotes the number of models which are used. The difference of the two algorithms is of a multiplicative magnitude of  $\mathcal{M}$ . This computational difference can be vanished if the computations of each model  $m \in \mathcal{M}$  are executed in parallel.

### 6.4 Stopping Criteria

A general and crucial issue with learning algorithms is about stopping criteria, namely, the criteria for which the robots stop the iterative game playing stage (coordination stage) and make a decision. An obvious choice is when a maximum number of iterations is reached. This number should depend on the size of the game, number of robots and their available actions, the problem of interest and constraints which arise from it. For example, in cases where the communication is expensive or robots should make decisions in real time, the maximum number of iterations cannot be arbitrarily large.

If the optimal solutions for the task of interest, and therefore the Nash equilibria, are known, then the coordination state can be terminated before the maximum number of iterations is reached. It can be stopped when the joint action is one of the Nash equilibria or the reward of the selected joint action is not less than a constant  $\epsilon$  from the equilibrium reward.

Even in cases where no Nash equilibria of the game are known, it is still possible to have a stopping criterion before the maximum number of iterations is reached. Coordination can be established if the learning algorithm has converged to a specific joint action. Convergence to a specific joint action can be defined as the repeated choice for  $c$  iterations of his particular joint action. Note here that there is no fixed  $c$  which can be used for all games. The size of the game and its nature should be taken into account when  $c$  is defined.

Depending on the application and the need to converge to a Nash equilibrium, different combinations of the aforementioned criteria can be used.

## 6.5 Example of a Sequence of Bayesian Games

In [14], it was shown that dec-POMDPs can be cast as a sequence of Bayesian games. In this section, we present a process of implementing a sequence of Bayesian games to solve a coordination task. Consider the task allocation problem between two robots with rewards shown in Table 3 in Sect. 3. Depending to the type of robot 2, there are two pure Nash equilibria where robots can converge. The first one is (*easy, easy*) when robot 2 is of type *A* and (*difficult, difficult*) when robot 2 is of type *B*. In order to accomplish their mission robots should finish both the easy and the difficult task. Independently of the action the robots will choose for this game, they will have accomplish only half of the necessary tasks. Therefore the players will have to play a sequence of games in order to finish both tasks, easy and difficult. The first game is the one with rewards depicted in Table 3. Another game should be defined in order to complete the unselected task, after making a decision based on the first game. An example of such a game is defined in Table 9. There the robots should choose if they will both try to finish the remaining task or only one should try or none of them should try. Note here that the game depicted in Table 9 makes the two robots coordinate and choose to do the remaining task together. Nonetheless, depending on the nature of the problem another reward matrix can be used in order to allow robots having different behaviours.

**Table 9.** Reward matrices of the Bayesian game for the remaining action.

Type A			Type B		
	do	not do		do	not do
do	10,10	-5,0	do	10,5	-5,0
not do	-5,0	0,0	not do	0,-10	0,0

## 7 Simulation Results

### 7.1 Results in Potential Games (Source: [59])

In this section the performance of the proposed algorithm, MMAF-EKFFP, is compared against the one of EKFFP in a resource allocation task. This is the vehicle-target assignment game which was introduced in [1]. In this potential game,  $N$  robots and  $M$  targets are placed in an area. The goal of each robot is to engage a target in order to destroy it. The actions of each robot are simply the choice of a target to engage. Each robot can choose only one target to engage, but a target can be engaged by many robots. The probability robot  $i$  has to destroy a target  $m$  is assumed to be independent of the probability another robot  $j$  has to destroy the same target. The probability that a target  $m$  will be destroyed is computed as:

$$1 - \prod_{i:a^i=m} (1 - p_{im}),$$

where  $p_{im}$  is the probability at which the robot  $i$  can destroy target  $m$ .

The reward that robots will share is the sum of the rewards each target  $m$  will produce if a specific joint action  $a$  is selected:

$$r_{global}(a) = \sum_{m \in M} r_m(a), \quad (17)$$

where  $r_m(a)$ , is defined as the product of target's  $m$  value  $V_m$  and the probability it can be destroyed by the robots which engage it. More formally, we can express the utility that is produced by target  $m$  as follows.

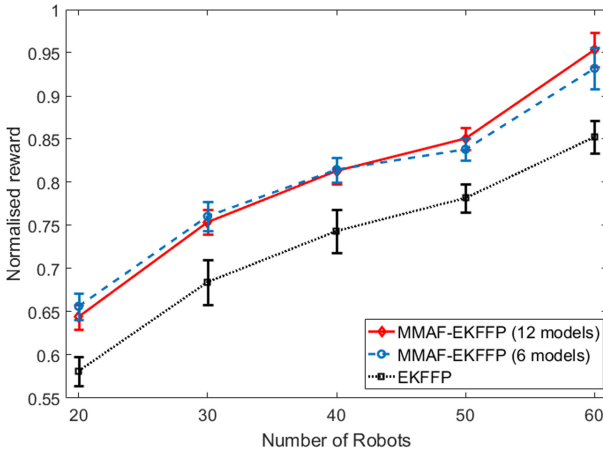
$$r_m(a) = V_m(1 - \prod_{i \in \mathcal{I}:a^i=m} (1 - p_{im})). \quad (18)$$

Multiple cases of the vehicle target assignment game were considered, with varying number of robots. In particular 20 targets and  $N = \{20, 30, 40, 50, 60\}$  robots were considered. The simulations were run in a computer with dual Intel Xeon E5-2643 v2 processors (3.50 GHz, 6 cores) and 384 GB memory. The probability that each robot  $i$  can destroy a target  $m$  was set to be proportional to their euclidean distance. For each case, we run MMAF-EKFFP with six and twelve models respectively. Each model represents a pair of covariance matrices  $\Xi$  and  $Z$ . In all cases, the values of  $\Xi$  and  $Z$  were uniformly sampled from the interval  $(0, 0.1]$ . Comparisons are also made with the classic EKFFP with  $\Xi$  and  $Z$  defined as in [60]. We reinforce here that there does not exist a universal combination of  $\Xi$  and  $Z$  which maximises the performance of EKFFP for all games [60]. The parameters which are reported in [60] are suggestive. Therefore, the question which arises is which pair of parameters  $\Xi$  and  $Z$  maximise the performance of EKFFP for the games of interest. MMAF-EKFFP provides a solution to this problem, since many combinations of  $\Xi$  and  $Z$  can be used as part of different models. Then the players will select actions based on the weighted sum of these models. The results presented in this paper are averaged

over 200 runs for each case. In each run, target values  $V_m$  were uniformly chosen from  $(0, 10]$ . The target and the robot positions were uniformly chosen from  $[0, 1]$ . This results in different scales of reward function. In order to make comparison feasible, the following normalised version of the global utility was used:

$$r_{total} = \frac{r_{global}}{\sum_m V_m}.$$

As it is depicted in Fig. 4, MMAF-EKFFP with multiple models performs better than the classic EKFFP, which uses only a single model.



**Fig. 4.** Average reward over 200 runs of the vehicle target assignment game, [59].

In addition, MMAF-EKFFP with 6 and 12 models perform similarly when cases with up to 40 robots are considered. When more players are considered, having more models improves the results of the algorithm.

It is expected that MMAF-EKFFP has heavier computational cost than EKFFP as the number of models is increased. Nonetheless as it is shown in Table 10, the running time of MMAF-EKFFP is in the same level as EKFFP because it is implemented using parallel computing, taking advantage of the multiple cores which many development platforms already have in robotics. In this experiment, each model is processed by an individual thread, which then runs in an individual physical computation core.

**Remark.** This game can be also solved by message passing algorithms like maxsum [17]. However, the size of search table and the size of the messages that need to be exchanged between agents render the application of these methods prohibitive complex in real time applications. In particular, if each robot were using maxsum algorithm, then it should update and transmit  $20^{50}$  messages. This is because a robot can engage any of the 20 available targets and the reward of a robot  $i$  depends on the joint action  $a^{-i}$  of all the other robots.

**Table 10.** Time in seconds that 20 iterations needed for various number of models EKF fictitious play, from [59].

	20 robots	30 robots	40 robots	50 robots
EKFFP	0.0649 ± 0.0042	0.1133 ± 0.0130	0.1548 ± 0.0109	0.2321 ± 0.0227
MMAF-EKFFP (6 models sequential)	0.3243 ± 0.0263	0.5909 ± 0.2247	1.1428 ± 0.1618	1.6509 ± 0.3116
MMAF-EKFFP (12 models sequential)	0.7773 ± 0.0327	1.3294 ± 0.1102	2.2283 ± 0.2177	3.0378 ± 0.2063
MMAF-EKFFP (6 models parallel)	0.0568 ± 0.0013	0.0906 ± 0.0034	0.1326 ± 0.0059	0.1918 ± 0.0105
MMAF-EKFFP (12 models parallel)	0.0616 ± 0.0030	0.1091 ± 0.0065	0.1619 ± 0.0071	0.2326 ± 0.0192

### 7.2 Games with Noisy Rewards (Source: [59])

This section contains the simulation results in a scenario where the robots receive noisy observations of the other robots’ actions. Noisy observations denote observations which for some reason is incorrect. Consider the case where two UAVs monitor a part of a field in order to decide if fertilisation is needed. The best results are obtained when one UAV flies at the top of the area and the other flies at the sides of the area. This scenario can be modelled as a two player game, which is depicted in Table 11. If both UAVs choose to go at the top of the area of interest, they will collide and they will receive some negative rewards. On the other hand, if they both choose to fly at the side of the area, then they will gain no reward because the quality of the images they gather will be poor. Finally, some positive reward is generated only if the two UAVs make different decisions. Note that the natural choice of the two UAVs is to choose to fly at the side of the area. They change their decision only when they believe with probability greater than 0.8 that the other UAV will choose to go at the side of the area of interest. In addition, each UAV can observe correctly the intention of the other UAV with probability  $p$  in every iteration of the coordination process.

**Remark.** In this paper, we assume that each robot knows the probability distribution of noisy observations, either by having some prior knowledge about its sensors’ specification or using some methods to estimate that distribution as the one proposed in [47] to estimate this distribution.

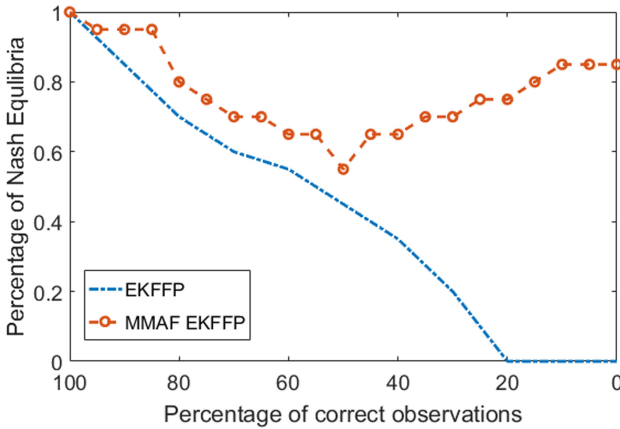
**Table 11.** UAVs’ rewards for the game with noisy observations.

	Top	Sideways
Top	-4, -4	0, 1
Sideways	1, 0	0, 0

Figure 5 shows the results of MMAF-EKFFP and EKFFP for the game with rewards depicted in Table 11. As it is shown here the percentage of times that

MMAF-EKFFP converged to a Nash equilibrium is always greater than 50%. This is significantly better than the results reported in [11], where the results of Generalised Weakened fictitious play (GWFP) [28] and Filtered Fictitious Play (FFP) [11]. The probability of converging to a Nash equilibrium reached zero when the 50% or more of the observations were faulty for FFP. For the case of GWFP the probability of converging to a Nash equilibrium reached zero when the 30% or more of the observations were faulty. Note here that even EKFFP performs better than FFP and GWFP since the probability of converging to a Nash equilibrium reached zero when the 80% or more of the observations were faulty.

The minimum probability of convergence to a Nash equilibrium for the MMAF-EKFFP algorithm is observed when the 50% of the observations were faulty. This is because the observations had exactly the same chance to be correct or faulty. When players know that the probability of a faulty observation is greater than a correct one, they use this information and increase their chances to converge to an equilibrium point.



**Fig. 5.** Probability to converge to Nash equilibrium as a function of the percentage of the correctly observed opponents' actions, [59].

### 7.3 Results in Bayesian Games (Source: [59])

The majority of the methods that are used to solve Bayesian games, such as Agent Security via Approximate Policies (ASAP) [41], Mixed Integer Programming Nash (MIP Nash) [50], brute force search methods [39] or Multiple Linear Programs [12] tries to find the Bayes Nash equilibrium with the highest reward. In order to solve the Bayesian game in Table 3, we can transform it into a strategic form game using Harsanyi's transformation [20], and search for the Nash equilibrium of this new game. An example of Harsanyi's transformation is depicted in Table 12, where the Bayesian game of Table 3 has been cast as a

strategic form game, with probabilities  $p$  and  $1 - p$  of the second player being of type  $A$  or  $B$ . In this game, robot 1 is the row player and robot 2 the column player.

Note here that as the second robot can be of any of the two types in the strategic form game, its actions consist of all the possible combinations of its actions in the two games. Therefore, in the strategic form representation of Table 3, the second robot has four possible actions  $E_A D_B, E_A E_B, D_A D_B, D_A E_B$ , where  $E_A D_B$  denotes selecting the easy task if it is of type A and the difficult task if it is of type B etc.

The game with rewards in Table 12 have three Bayes-Nash equilibria. Two pure Bayes-Nash equilibria and one mixed. The joint actions (*difficult*,  $D_A D_B$ ) and (*easy*,  $D_A E_B$ ) are pure Bayes-Nash equilibria. The mixed Bayes-Nash equilibrium exists when  $p > \frac{1}{5}$ : robot 1 chooses the difficult task with probability  $\frac{5p-1}{5-p}$  and robot 2 chooses  $D_A D_B$  with probability  $\frac{2}{3}$ . Nonetheless, even finding the Bayes-Nash equilibria does not answer to the question which action the robots should choose if they are playing any of the two games. For example if the type of robot 2 is A, then the equilibrium of the actual game which will be played is easy;easy, as it can be seen from Table 3. On the other hand, if the type of the robot 2 is of type B then the equilibrium of the actual game which will be played is difficult, difficult.

**Table 12.** Strategic form game’s rewards, of the Bayesian game of Table 3 as a function of  $p$ .

		Robot 2			
		$E_A D_B$	$E_A E_B$	$D_A D_B$	$D_A E_B$
Robot 1	Difficult	$9 + p, 10 - 4p$	$5 + 5p, 4 + 2p$	$9 - 4p, 10$	$5, 4 + 6p$
	Easy	$8, 6 - p$	$7 + p, 5$	$8 + 2p, 6 - 2p$	$7 - p, 5 + 3p$

On the other hand, the proposed algorithm allows robots to learn what the other robots are doing, and evaluate the probability of choosing a particular sequence of actions given that they are of a specific type. Then based on this knowledge, they choose the action that maximises their expected reward conditional to the possible types of the other players.

Now we study the performance of MMAF-EKFFP in two games with incomplete information. The first game had at least one pure strategy Nash equilibrium, and the second has only a mixed strategy Nash equilibrium.

The first game is the one depicted in Tables 3 and 9. MMAF-EKFFP always converged to the pure Nash equilibrium of the game, and therefore, the two robots always chose to work on the same task (easy or difficult) jointly.

The second example which is considered comes from a security problem [41]. In security games, a group of security robots try to secure some areas of interest and an attacker robot tries to invade these areas. In [3, 48], the security robots cannot be physically present in all the areas of interest at the same time. Instead,

they can choose among various patrol routes. Security robots choose the route and the areas they will patrol based on the importance of the areas or the likelihood at which an attacker robot will be appear etc. The security problem can be cast as a two player game [41]. If there are  $\mathcal{M}$  areas of interest, then the action of the security robots will be the  $d$ -tuple ( $d \leq m$ ) of the areas which the robots will patrol. The order by which the robots visit the areas is also taking into account. For instance, in the case of three areas  $\{1, 2, 3\}$ , each patrol order, e.g.,  $1 \leftarrow 2 \leftarrow 3$  or  $1 \leftarrow 3 \leftarrow 2$ , is a different action. The attacker robot can choose any of the  $\mathcal{M}$  available areas to invade. Assume that the security robots can choose from  $\mathcal{K}$  patrolling routes. The reward of the security robots  $r_i(k)$  ( $k \in \mathcal{K}$ ) and that of attacking robot  $r_j(m)$  ( $m \in \mathcal{M}$ ) are estimated respectively as follows.

$$r_i(k) = \begin{cases} -u_i(a^i) & \text{if } a^j \notin a^i \\ p_{a^i}c_i + (1 - p_{a^i})(-u_i) & \text{if } a^j \in a^i, \end{cases} \quad (19)$$

where  $u_i(a^i)$  is the value of area  $a^i$  to the security robots,  $p_{a^i}$  is the probability that the security robots can catch the attacker in the  $a^i$  area,  $c_i$  is the reward to the security robots if the attacker is caught:

$$r_j = \begin{cases} u_j & \text{if } a^j \notin a^i \\ -p_{a^j}c_j + (1 - p_{a^j})(u_j) & \text{if } a^j \in a^i, \end{cases} \quad (20)$$

where  $u_i(a^j)$  is the value of area  $a^j$  to the attacker robot,  $p_{a^j}$  is the probability that the security robots can catch the attacker when patrolling area  $a^j$ , and  $c_j$  is the cost to the attacker robot if it is caught.

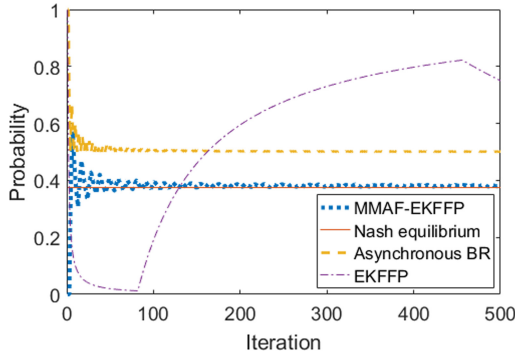
Consider the case where two areas are available and the actions available to security robot are  $1 \leftarrow 2$  and  $2 \leftarrow 1$  respectively. In addition, assume that the attacking robot can be of two types  $A$  and  $B$ . The above reward function, when similar parameters to [41] are used, can be modelled as a Bayesian game as it is depicted in Table 13.

**Table 13.** Reward matrices of the attacking and security robot for two different types of attacking robot. The top table is the game played when the attacking robot is of type  $A$ .

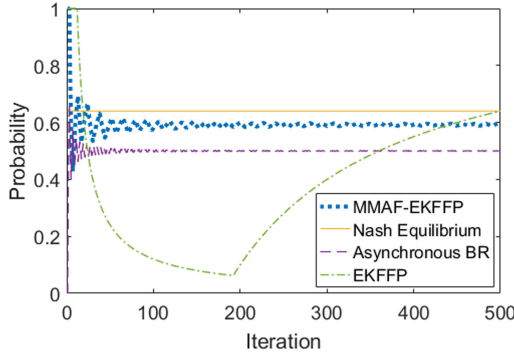
	Areas 1,2	Areas 2,1
Area 1	-1,0.5	-0.125,-0.125
Area 2	-0.375,0.125	-1,0.5
	Areas 1,2	Areas 2,1
Area 1	-1.2,0.4	-0.025,-0.025
Area 2	-0.275,0.125	-1.2,0.4



When the attacker is type  $A$ , the optimal policy for the security robots, which leads to a Nash equilibrium, is to choose a mixed strategy and play action  $1 \leftarrow 2$  with probability 0.58 and  $2 \leftarrow 1$  with probability 0.42. The mixed strategy for the attacking robot is to choose area 1 with probability 0.375 and area 2 with probability 0.625. If attacker is of type  $B$ , the security robots' optimal policy is to choose a mixed strategy by playing action  $1 \leftarrow 2$  with probability 0.35 and  $2 \leftarrow 1$  with probability 0.65. The mixed strategy for the attacking robot is to choose area 1 with probability 0.39 and area 2 with probability 0.61. The security robots assume that the attacking robot is of type  $A$  with probability 0.9 and of type  $B$  with probability 0.1. Figure 6 illustrates the probability with which the attacking robot, chose Area 1 when it was of type  $A$ . The performance of MMAF-EKFFP was compared with EKFFP and the asynchronous best response algorithm which proposed in [15] in order to solve Bayesian games. As it can be seen from Fig. 6, MMAF-EKFFP converged to the Nash equilibrium of the game while the two other algorithms failed to converged to a value close to the Nash equilibrium. Similarly, Fig. 7 depicts the probability with which the security robot chose action  $2 \leftarrow 1$ . As it can be observed from Fig. 6, MMAF-EKFFP converged to the Nash equilibrium, while the other algorithms failed to converged to a value close to the Nash equilibrium. Similar results were obtained for various combinations of the probabilities with which the attacking robot could be of type  $A$  or  $B$ . In particular, the differences in the results between the case where the attacking robot is of type  $A$  with probability 0.9 and the case when it is of type  $B$  with probability 0.1, are less than 0.01 from the reported results in Figs. 6 and 7.



**Fig. 6.** Probability of the attacking robot to choose area 1 when it is of type  $A$  as a function of the number of iterations. The solid line represents the Nash equilibrium, the squared-line the MMAF-EKFFP algorithm, the dashed line the Asynchronous Best response and the dotted line the EKFFP algorithm, [59].



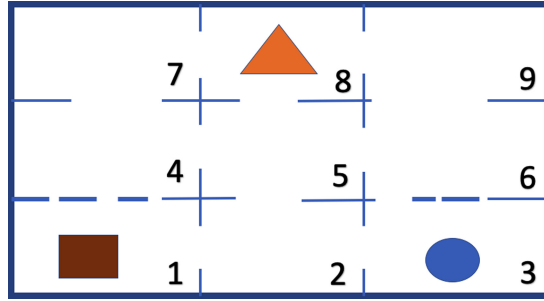
**Fig. 7.** Probability of the security robot to choose action  $2 \leftarrow 1$  when the attacking robot is of type  $B$  as a function of the number of iterations. The solid line represents the Nash equilibrium, the squared-line the MMAF-EKFFP algorithm, the dashed line the Asynchronous Best response and the dotted line the EKFFP algorithm, [59].

### 7.4 Results on Stochastic and State-Based Games

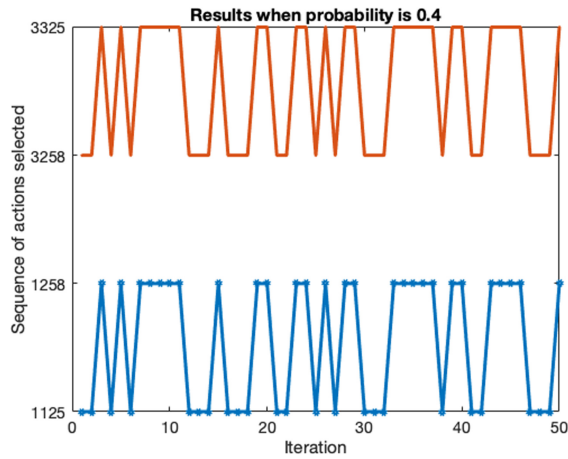
In this section the results obtained for two scenarios of state based games are presented. The first is the game that was described in Tables 6, 7 and 8, with the corresponding transition probabilities of Eq. 3. This state space game was played for 100 replications and in each replication 100 iterations of the state-based game where played. In this game, the players that used EKFMMAF fictitious play, where playing the recurrent equilibrium of the game. In particular, 57% of the times they were selecting the (Not cooperate, Not cooperate) action of the coordination game, and 43% of times they were selecting the (Not cooperate, Not cooperate) of the prisoners’ dilemma game.

Consider the game that it is depicted in Fig. 8. In this game two robots located in squares 1 and 3 respectively should reach the square with number 8 where their target is. When a robot is in a particular square this represents a specific state. Even though the available actions of the players are the same in all states, different behaviours can emerge when a robots are in different squares.

The robots can move either one square forward, or one square backwards, or one square at the left or one square at the right. If they both choose to move to the same square, case of possible collision, then only one is can go to the new square with probability  $\frac{1}{2}$ , while the other one stays at its current cell. In addition, the robots can transition from squares 1 and 3 to squares 4 and 6 respectively with probability  $p$ . One iteration of the game finishes when at least one robot reaches the square number 8.



**Fig. 8.** State-based game of two robots. The robots (brown square and blue cycle) should reach their target on square number eight (orange triangle). (Color figure online)



**Fig. 9.** Sequence of actions when probability  $p = 0.4$ . (Color figure online)

**Table 14.** Percentage of times that each robot reached its target.

$p$	% of times only robot 1 reached cell #8	% of times only robot 2 reached cell #8	% of times both robots reached cell #8
0.4	0.54	0.46	0
0.5	0.44	0.56	0
0.6	0.48	0.04	0.48
0.7	0.22	0.02	0.76
0.8	0.02	0.04	0.94

Table 14 shows the percentage of times that either only a single robot or both robots managed to reach their target. As expected better coordination achieved when  $p \geq 0.6$ . In particular, as it is depicted to Figs. 9 and 10, both players were choosing to towards the square number two, since it would give them better chances to reach their target. When  $p \geq 0.6$  the robots manage to coordinate and one of was selecting as its first action square 2 while the other was choosing squares 4 and 6 respectively, Figs. 11, 12 and 13.

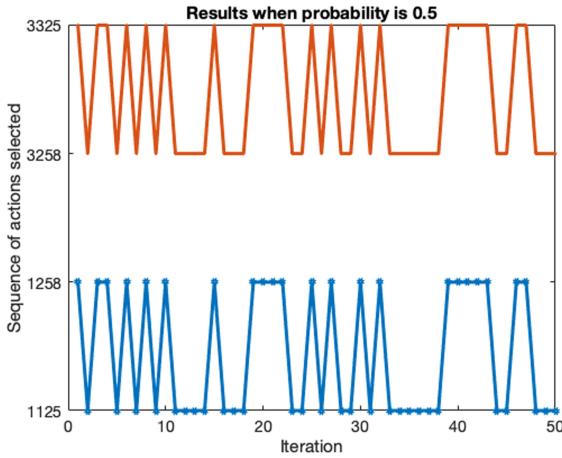


Fig. 10. Sequence of actions when probability  $p = 0.5$ . (Color figure online)

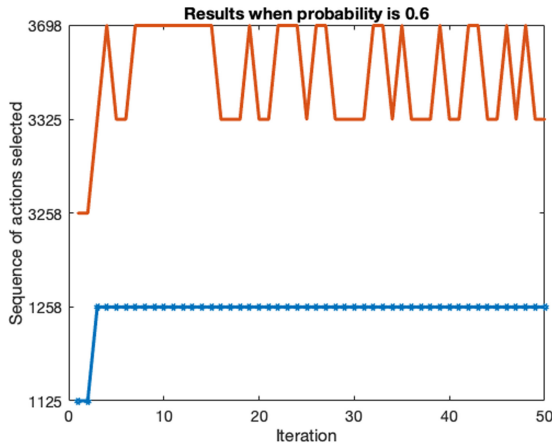
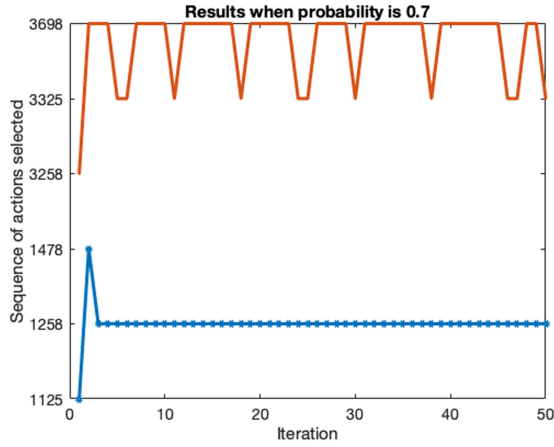
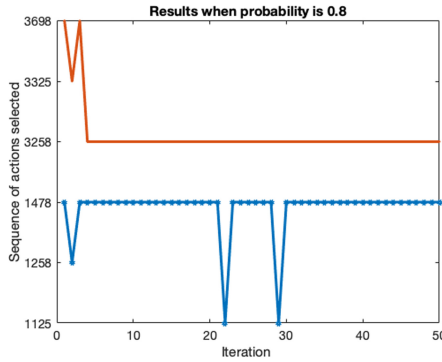


Fig. 11. Sequence of actions when probability  $p = 0.6$ . (Color figure online)



**Fig. 12.** Sequence of actions when probability  $p = 0.7$ . (Color figure online)



**Fig. 13.** Sequence of actions when probability  $p = 0.8$ . (Color figure online)

## 8 Conclusions and Future Works

A new game-theoretic learning algorithm based on EKFFP and multi-model adaptive filters has been proposed. This new algorithm can take into account various forms of uncertainty.

The performance of the proposed algorithm was tested in various games, that are related to robotic applications. The experimental results showed that it can provide better solution than the classic EKFFP algorithm and can be run as fast as the latter if implemented using parallel computing.

As future work the proposed algorithm will be applied in robotic platforms with various computing capabilities, in scenarios that can be formulated as games. This will test the applicability of the proposed methodology in real life scenarios.

## References

1. Arslan, G., Marden, J.R., Shamma, J.S.: Autonomous vehicle-target assignment: a game-theoretical formulation. *J. Dyn. Syst. Meas. Contr.* **129**(5), 584–596 (2007)
2. Ayken, T., Imura, J.i.: Asynchronous distributed optimization of smart grid. In: 2012 Proceedings of SICE Annual Conference (SICE), pp. 2098–2102. IEEE (2012)
3. Beard, R.W., McLain, T.W.: Multiple UAV cooperative search under collision avoidance and limited range communication constraints. In: Proceedings of 42nd IEEE Conference on Decision and Control, vol. 1, pp. 25–30. IEEE (2003)
4. Berger, U.: Fictitious play in 2xn games. *J. Econ. Theory* **120**(2), 139–154 (2005)
5. Bishop, C.M.: *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford (1995)
6. Blair, W., Bar-Shalom, T.: Tracking maneuvering targets with multiple sensors: does more data always mean better estimates? *IEEE Trans. Aerosp. Electron. Syst.* **32**(1), 450–456 (1996)
7. Bonato, V., Marques, E., Constantinides, G.A.: A floating-point extended Kalman filter implementation for autonomous mobile robots. *J. Signal Process. Syst.* **56**(1), 41–50 (2009)
8. Botelho, S., Alami, R.: M+: a scheme for multi-robot cooperation through negotiated task allocation and achievement. In: 1999 IEEE International Conference on Robotics and Automation, 1999, Proceedings, vol. 2, pp. 1234–1239 (1999)
9. Brown, R.G., Hwang, P.Y.: *Introduction to Random Signals and Applied Kalman Filtering: With Matlab Exercises and Solutions*. Wiley, New York, 1 (1997)
10. Caputi, M.J.: A necessary condition for effective performance of the multiple model adaptive estimator. *IEEE Trans. Aerosp. Electron. Syst.* **31**(3), 1132–1139 (1995)
11. Chapman, A.C., Williamson, S.A., Jennings, N.R.: Filtered fictitious play for perturbed observation potential games and decentralised POMDPs. CoRR abs/1202.3705 (2012). <http://arxiv.org/abs/1202.3705>
12. Conitzer, V., Sandholm, T.: Choosing the best strategy to commit to. In: ACM Conference on Electronic Commerce (2006)
13. Di Mario, E., Navarro, I., Martinoli, A.: Distributed learning of cooperative robotic behaviors using particle swarm optimization. In: Hsieh, M.A., Khatib, O., Kumar, V. (eds.) *Experimental Robotics*. STAR, vol. 109, pp. 591–604. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-23778-7\\_39](https://doi.org/10.1007/978-3-319-23778-7_39)
14. Emery-Montemerlo, R.: Game-theoretic control for robot teams. Ph.D. thesis, The Robotics Institute. Carnegie Mellon University (2005)
15. Emery-Montemerlo, R., Gordon, G., Schneider, J., Thrun, S.: Game theoretic control for robot teams. In: Proceedings of the 2005 IEEE International Conference on Robotics and Automation, ICRA 2005, pp. 1163–1169. IEEE (2005)
16. Evans, J.M., Krishnamurthy, B.: HelpMate®, the trackless robotic courier: a perspective on the development of a commercial autonomous mobile robot. In: de Almeida, A.T., Khatib, O. (eds.) *Autonomous Robotic Systems*. LNCIS, vol. 236, pp. 182–210. Springer, London (1998). <https://doi.org/10.1007/BFb0030806>
17. Farinelli, A., Rogers, A., Jennings, N.: Agent-based decentralised coordination for sensor networks using the max-sum algorithm. *Auton. Agents Multi-Agent Syst.* **28**, 337–380 (2014)
18. Fudenberg, D., Levine, D.: *The Theory of Learning in Games*. The MIT Press, Cambridge (1998)
19. Govindan, S., Wilson, R.: A global newton method to compute Nash equilibria. *J. Econ. Theory* **110**(1), 65–86 (2003)

20. Harsanyi, J.C., Selten, R.: A generalized Nash solution for two-person bargaining games with incomplete information. *Manage. Sci.* **18**(5-part-2), 80–106 (1972)
21. Hennig, P.: Fast probabilistic optimization from noisy gradients. In: ICML, no. 1, pp. 62–70 (2013)
22. Jiang, A.X., Leyton-Brown, K.: Bayesian action-graph games. In: *Advances in Neural Information Processing Systems*, pp. 991–999 (2010)
23. Kho, J., Rogers, A., Jennings, N.R.: Decentralized control of adaptive sampling in wireless sensor networks. *ACM Trans. Sen. Netw.* **5**(3), 1–35 (2009)
24. Kho, J., Rogers, A., Jennings, N.R.: Decentralized control of adaptive sampling in wireless sensor networks. *ACM Trans. Sen. Netw. (TOSN)* **5**(3), 19 (2009)
25. Kitano, H., et al.: RoboCup rescue: search and rescue in large-scale disasters as a domain for autonomous agents research. In: 1999 IEEE International Conference on Systems, Man, and Cybernetics, 1999. IEEE SMC 1999 Conference Proceedings, vol. 6, pp. 739–743. IEEE (1999)
26. Koller, D., Milch, B.: Multi-agent influence diagrams for representing and solving games. *Games Econ. Behav.* **45**(1), 181–221 (2003)
27. Kostelnik, P., Hudec, M., Šamulka, M.: Distributed learning in behaviour based mobile robot control. In: Sinac, P. (Ed) *Intelligent Technologies-Theory and Applications IOP press* (2002)
28. Leslie, D.S., Collins, E.J.: Generalised weakened fictitious play. *Games Econ. Behav.* **56**(2), 285–298 (2006)
29. Li, C., Xing, Y., He, F., Cheng, D.: A strategic learning algorithm for state-based games. *Automatica* **113**, 108615 (2020)
30. Littman, M.L.: Markov games as a framework for multi-agent reinforcement learning. In: *Machine Learning Proceedings 1994*, pp. 157–163. Elsevier (1994)
31. Madhavan, R., Fregene, K., Parker, L.: Distributed cooperative outdoor multirobot localization and mapping. *Auton. Rob.* **17**(1), 23–39 (2004)
32. Makarenko, A., Durrant-Whyte, H.: Decentralized data fusion and control in active sensor networks. In: *Proceedings of 7th International Conference on Information Fusion*, vol. 1, pp. 479–486 (2004)
33. Marden, J.R.: State based potential games. *Automatica* **48**(12), 3075–3088 (2012)
34. Matignon, L., Laurent, G.J., Le Fort-Piat, N.: Independent reinforcement learners in cooperative Markov games: a survey regarding coordination problems. *Knowl. Eng. Rev.* **27**(1), 1–31 (2012)
35. Miyasawa, K.: On the convergence of learning process in a 2x2 non-zero-person game (1961)
36. Monderer, D., Shapley, L.: Potential games. *Games Econ. Behav.* **14**, 124–143 (1996)
37. Nachbar, J.: Evolutionary’ selection dynamics in games: convergence and limit properties. *Int. J. Game Theory* **19**, 59–89 (1990)
38. Nash, J.: Equilibrium points in n-person games. *Proc. Nat. Acad. Sci. U.S.A.* **36**, 48–49 (1950)
39. Oliehoek, F.A., Spaan, M.T., Dibangoye, J.S., Amato, C.: Heuristic search for identical payoff Bayesian games. In: *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems*, vol. 1, pp. 1115–1122. International Foundation for Autonomous Agents and Multiagent Systems (2010)
40. Parker, L.: ALLIANCE: an architecture for fault tolerant multirobot cooperation. *IEEE Trans. Robot. Autom.* **14**(2), 220–240 (1998)
41. Paruchuri, P., Pearce, J.P., Tambe, M., Ordonez, F., Kraus, S.: An efficient heuristic for security against multiple adversaries in Stackelberg games. In: *AAAI Spring Symposium: Game Theoretic and Decision Theoretic Agents*, pp. 38–46 (2007)

42. Pugh, J., Martinoli, A.: Distributed scalable multi-robot learning using particle swarm optimization. *Swarm Intell.* **3**(3), 203–222 (2009)
43. Rabinovich, Z., Gerding, E., Polukarov, M., Jennings, N.R.: Generalised fictitious play for a continuum of anonymous players. In: Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI), 01 July 2009, pp. 245–250 (2009)
44. Rabinovich, Z., Naroditskiy, V., Gerding, E.H., Jennings, N.R.: Computing pure Bayesian-Nash equilibria in games with finite actions and continuous types. *Artif. Intell.* **195**, 106–139 (2013)
45. Raffard, R.L., Tomlin, C.J., Boyd, S.P.: Distributed optimization for cooperative agents: application to formation flight. In: 43rd IEEE Conference on Decision and Control, 2004. CDC. vol. 3, pp. 2453–2459. IEEE (2004)
46. Robinson, J.: An iterative method of solving a game. *Ann. Math.* **54**, 296–301 (1951)
47. Rosen, D.M., Leonard, J.J.: Nonparametric density estimation for learning noise distributions in mobile robotics. In: 1st Workshop on Robust and Multimodal Inference in Factor Graphs, ICRA (2013)
48. Ruan, S., Meirina, C., Yu, F., Pattipati, K.R., Popp, R.L.: Patrolling in a stochastic environment. In: 10 International Symposium on Command and Control (2005)
49. Sakka, S.: Bayesian Filtering and Smoothing. Cambridge University Press, Cambridge (2013)
50. Sandholm, T., Gilpin, A., Conitzer, V.: Mixed-integer programming methods for finding nash equilibria. In: Proceedings of the National Conference on Artificial Intelligence, vol. 20 (2005)
51. Semsar-Kazerooni, E., Khorasani, K.: Optimal consensus algorithms for cooperative team of agents subject to partial information. *Automatica* **44**(11), 2766–2777 (2008)
52. Semsar-Kazerooni, E., Khorasani, K.: Multi-agent team cooperation: a game theory approach. *Automatica* **45**(10), 2205–2213 (2009)
53. Shapley, L.S.: Stochastic games. *Proc. Nat. Acad. Sci.* **39**(10), 1095–1100 (1953)
54. Sharma, R., Gopal, M.: Fictitious play based Markov game control for robotic arm manipulator. In: Proceedings of 3rd International Conference on Reliability, Infocom Technologies and Optimization, pp. 1–6 (2014)
55. Simmons, R., Apfelbaum, D., Burgard, W., Fox, D., Moors, M., Thrun, S., Younes, H.: Coordination for multi-robot exploration and mapping. In: AAAI/IAAI, pp. 852–858 (2000)
56. Smyrnakis, M., Leslie, D.S.: Dynamic opponent modelling in fictitious play. *Comput. J.* **53**, 1344–1359 (2010)
57. Smyrnakis, M., Galla, T.: Decentralized optimisation of resource allocation in disaster management. In: Preston, J., et al. (eds.) *City Evacuations: An Interdisciplinary Approach*, pp. 89–106. Springer, Heidelberg (2015). [https://doi.org/10.1007/978-3-662-43877-0\\_5](https://doi.org/10.1007/978-3-662-43877-0_5)
58. Smyrnakis, M., Kladis, G.P., Aitken, J.M., Veres, S.M.: Distributed selection of flight formation in UAV missions. *J. Appl. Math. Bioinform.* **6**(3), 93–124 (2016)
59. Smyrnakis, M., Qu, H., Bauso, D., Veres, S.M.: Multi-model adaptive learning for robots under uncertainty. In: Rocha, A.P., Steels, L., van den Herik, H.J. (eds.) Proceedings of the 12th International Conference on Agents and Artificial Intelligence, ICAART 2020, Valletta, Malta, 22–24 February 2020, vol. 1, pp. 50–61 (2020)
60. Smyrnakis, M., Veres, S.: Coordination of control in robot teams using game-theoretic learning. *Proc. IFAC* **14**, 1194–1202 (2014)



61. Stranjak, A., Dutta, P.S., Ebden, M., Rogers, A., Vytelingum, P.: A multi-agent simulation system for prediction and scheduling of aero engine overhaul. In: AAMAS 2008: Proceedings of the 7th International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 81–88, May 2008
62. Timofeev, A., Kolushev, F., Bogdanov, A.: Hybrid algorithms of multi-agent control of mobile robots. In: International Joint Conference on Neural Networks, 1999. IJCNN 1999, vol. 6, pp. 4115–4118 (1999)
63. Tsalatsanis, A., Yalcin, A., Valavanis, K.P.: Dynamic task allocation in cooperative robot teams. *Robotica* **30**, 721–730 (2012)
64. Tsalatsanis, A., Yalcin, A., Valavanis, K.: Optimized task allocation in cooperative robot teams. In: 17th Mediterranean Conference on Control and Automation, 2009. MED 2009, pp. 270–275 (2009)
65. Tsitsiklis, J.N., Athans, M.: On the complexity of decentralized decision making and detection problems. *IEEE Trans. Autom. Control* **30**(5), 440–446 (1985)
66. Voice, T., Vytelingum, P., Ramchurn, S.D., Rogers, A., Jennings, N.R.: Decentralised control of micro-storage in the smart grid. In: AAAI, pp. 1421–1427 (2011)
67. Wolpert, D., Tumer, K.: A survey of collectives. In: Tumer, K., Wolpert, D. (eds.) *Collectives and the Design of Complex Systems*, pp. 1–42. Springer, New York. [https://doi.org/10.1007/978-1-4419-8909-3\\_1](https://doi.org/10.1007/978-1-4419-8909-3_1)
68. Zhang, P., Sadler, C.M., Lyon, S.A., Martonosi, M.: Hardware design experiences in ZebraNet. In: Proceedings of SenSys 2004, pp. 227–238. ACM (2004)
69. Zhang, Y., Schervish, M., Acar, E., Choset, H.: Probabilistic methods for robotic landmine search. In: Proceedings of the 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2001), pp. 1525–1532 (2001)
70. Zlot, R., Stentz, A., Dias, M.B., Thayer, S.: Multi-robot exploration controlled by a market economy. In: IEEE International Conference on Robotics and Automation, 2002. Proceedings. ICRA 2002, vol. 3 (2002)