

---

# Actor-Critic Fictitious Play in Simultaneous Move Multistage Games

---

Julien Perolat<sup>1</sup>  
Univ. Lille

Bilal Piot<sup>1</sup>  
Univ. Lille

Olivier Pietquin<sup>1</sup>  
Univ. Lille

## Abstract

Fictitious play is a game theoretic iterative procedure meant to learn an equilibrium in normal form games. However, this algorithm requires that each player has full knowledge of other players' strategies. Using an architecture inspired by actor-critic algorithms, we build a stochastic approximation of the fictitious play process. This procedure is on-line, decentralized (an agent has no information of others' strategies and rewards) and applies to multistage games (a generalization of normal form games). In addition, we prove convergence of our method towards a Nash equilibrium in both the cases of zero-sum two-player multistage games and cooperative multistage games. We also provide empirical evidence of the soundness of our approach on the game of Alesia with and without function approximation.

## 1 Introduction

Go, Chess, Checkers, Oshi-Zumo [10] are just a few example of Multistage games [7]. In these games, the interaction proceeds from stage to stage without looping back to a previously encountered situation. This model groups a broad class of multi-agent sequential decision processes where the interaction never goes back in the same state. This work focuses on Multi-Agent Reinforcement Learning [11] (MARL) in Multistage games. In this multi-agent environment, players evolve from state to state as a result of their mutual actions. During this interaction, all players re-

ceive a reward informing on how good was their action when performed in the state they were in. The goal of MARL is to learn a strategy that maximally accumulates rewards over time. Whilst the problem is fairly well understood when studying single agent RL, learning while independently interacting with other agents remains superficially explored. The range of open questions is so wide in that area [31] that it is worth giving a precise definition of our goal. In this paper, we follow a prescriptive agenda. We intend to find a learning algorithm that provably converges to Nash equilibrium in cooperative and in non-cooperative games. The goal is to find a strategy that can be executed independently by each player that corresponds to a Nash equilibrium. Many, if not most, approaches to address this problem consider a centralized learning procedure that produces an independent strategy for each player [21]. Centralized learning procedures are quite common and often perform better than decentralized learning procedures [13]. But these centralized learning procedure require synchronization between agents during learning (which is the main limitation of these methods). The agenda we follow in this paper is to propose a decentralized on-line learning method that provably converges to a Nash equilibrium in self-play. Decentralized algorithms, because they allow building identical independent agents that don't rely on anything but the observation of their state and reward, no central controller being required. On-line algorithms, on another hand, allow learning while playing and do not require prior computation of possible strategies.

This agenda is a fertile ground of interaction between traditional RL and game theory. Indeed, RL aims at building autonomous agents learning on-line in games against nature (where the environment is not interested in winning). For that reason, a wide variety of single agent RL algorithms have been adapted to multi-agent problems. But several major issues prevent direct use of standard RL with multi-agent systems. First, blindly applying single agent RL in a decentralized fashion implies that, from each agent's point of view, the other agents are part of the envi-

---

<sup>1</sup>now with DeepMind, London (UK)

ronment. Such an hypothesis breaks the crucial RL assumption that the environment is (at least almost) stationary [22]. Second, it introduces partial observability as each agent’s knowledge is restricted to its own actions and rewards while its behavior should depend on others’ strategies.

Decentralized procedures (unlike counterfactual regret minimization algorithms [34]) have been the topic of many studies in game theory and many approaches were proposed from policy hill climbing methods [8, 2] to evolutionary dynamics [33, 1] (related work will be detailed in Sec. 2). But those dynamics do not converge in all general-sum normal-form games, and, there exists a three-player normal form game [15] for which no first order uncoupled dynamics (*i.e.* most decentralized dynamics) can converge to a Nash equilibrium. Despite this counterexample, decentralize dynamics remain an important case to study because building a central controller for a multi-agent system is not always possible nor is observing the actions and rewards of every agent. Even if decentralized learning processes (as described in [15]) will never be guaranteed to converge in general, they should be at least guaranteed to converge in some interesting classes of games such as cooperative and zero-sum two-player games.

Fictitious play is a model-based process that learns Nash equilibria in normal form games. It has been widely studied and required assumptions were weakened over time [23, 17] since the original article of Robinson [28]. It has been extended to extensive form games (game trees) and, to a lesser extent, to function approximation [16]. However it is neither on-line nor decentralized except from the work of [23] which focuses on normal form games and [16] that has weak guarantees of convergence and focus on turn taking imperfect information games. Fictitious play enjoys several convergence guarantees [17] which makes it a good candidate for learning in simultaneous multistage stage games.

This paper contributes to fill a gap in the MARL literature by providing two online decentralized algorithms converging to a Nash equilibrium in multistage games both in the cooperative case and the zero-sum two-player case. Those two cases used to be treated as different agendas since the seminal paper of Shoham & *al.* [31] and we expect our work to serve as a milestone to reconcile them going further than normal form games [23, 17]. Our first contribution is to propose two novel on-line and decentralized algorithms inspired by actor-critic architectures for Markov Games, each of them working on two timescales. Those two algorithms perform the same actors’ update but use different methods for the critic. The first one performs

an off-policy control step whilst the second relies on a policy evaluation step. Although the actor-critic architecture is popular for its success in solving (continuous action) RL domains, we choose this architecture for a different reason. Our framework requires handling non-stationarity (because of adaptation of the other players) which is another nice property of actor-critic architectures. Our algorithms are stochastic approximations of two dynamical systems that generalize the work of [23] and [17] on the fictitious play process from normal form games to multistage games [7].

In the following, we first outline related work (in Sec. 2) and then describe the necessary background in both game theory and RL (Sec. 3) to introduce our first contribution, the two-timescale algorithms (Sec. 4). These algorithms are stochastic approximations of two continuous-time processes defined in Sec. 5. Then, we study (in Sec. 5) the asymptotic behavior of these continuous-time processes and show, as a second contribution, that they converge in self-play in cooperative games and in zero-sum two-player games. In Sec. 6, our third contribution proves that the algorithms are stochastic approximations of the two continuous-time processes. Finally, we perform an empirical evaluation (in Sec. 7).

## 2 Related Work

Decentralized reinforcement learning in games has been studied widely in the case of normal form games and includes regret minimization approaches [9, 12] or stochastic approximation algorithms [23]. However, to our knowledge, none of the previous methods have been extended to independent reinforcement learning in Markov Games or any intermediate models such as MSGs with guarantees of convergence both for cooperative and zero-sum case. Finding a single independent RL algorithm addressing both cases is still treated as separate agendas since the seminal paper [31].

***Q-Learning Like Algorithms:*** The adaptation of RL algorithms to the multi-agent setting was the first approach to address online learning in games. On-line algorithms like *Q*-learning [32] are often used in cooperative multi-agent learning environments but fail to learn a stationary strategy in simultaneous zero-sum two-player games. They fail in this setting because, in simultaneous zero-sum two-player games, it is not sufficient to use a greedy strategy to learn a Nash equilibrium. In [25], the *Q*-learning method is adapted to guarantee convergence to zero-sum two-player MGs. This method isn’t an independent learning algorithm anymore as each player needs to observe the action of the opponent. Other adaptation to *N*-player games were developed [18, 14]. However, all these methods

require the observation of the opponents' action and the two last ones are guaranteed to converge only under very conservative hypotheses. Moreover, these are not decentralized methods as each player needs to observe the reward of the others.

**Independent Policy Gradient:** Independent policy gradient methods is also an other attempt to address the problem of learning in games. In MGs and in normal form games, the policy hill climbing method [8] is probably the first approach with theoretical guarantees. But the guarantees of convergence are often limited to the two player two actions setting and fail to converge in zero-sum games. In this approach, each agent follows a gradient ascent on its expected outcome. The behavior of this algorithm can be improved using heuristics on the learning rate as reported in [8]. Many attempt were made to build on this approach such as the GIGA-WoLF algorithm but again, convergence properties are limited to the two player two actions case. It can also scale up using function approximation as in [2] which results in an actor-critic like method and fits in our setting.

**MCTS Algorithms in MSGs and Counterfactual Regret Minimization Algorithms (CFR):** A trove variety of tree search algorithms exists to compute Nash equilibria in MSGs. A review of those algorithms can be found in [7] and mainly focuses in zero-sum two-player simultaneous move multistage games. The convergence guarantees of those methods were first analyzed in [24] and then detailed in [19]. Those algorithms require that one can query the model at each stage of the game and as such do not belong to the family of independent RL methods.

An other family of algorithms related to our work is the CFR type of algorithms and its wide number of variant. Even if these algorithms have sample-based variations [20], they usually require each agent to be aware of the opponent's strategy and some game specific information. Only one variation of MCCFR was suggested in [20] would match our setting but is still unexplored in the literature and its unclear how function approximation could be used in this setting.

**Two-timescale algorithms in MDPs and in MGs:** Two-timescale algorithms have been the subject of a wide literature in MDPs that starts from the seminal work of [5]. These works mainly analyze the use of linear function approximation. However, when applied independently in a multiagent setting, those algorithms no longer have convergence guarantees (even without function approximation). In [27], the authors attempt to provide an on-line model-free algorithm with guarantees. This approach requires

each player to observe the reward of the others. The authors claim that their algorithm is guaranteed to converge to a Nash equilibrium but the proof is broken (details in App.K).

### 3 Background

**Markov Games (MGs):** Formally, an MG is a tuple  $\langle N, S, \mathbf{A}, p(s'|s, \mathbf{a}), \mathbf{r}(s, \mathbf{a}) \rangle$ . In each state  $s \in S$  of the MG, each of the  $N$ -players simultaneously takes an action  $a^i \in A^i$ . The joint action  $\mathbf{a}$  is  $(a^1, \dots, a^N)$  or  $(a^i, \mathbf{a}^{-i})$  (where  $\mathbf{a}^{-i}$  is the tuple  $(a^1, \dots, a^{i-1}, a^{i+1}, \dots, a^N)$ ) and the joint action set  $\mathbf{A}$  is  $A^1 \times \dots \times A^N = A^i \times \mathbf{A}^{-i}$ . As a result of this joint action  $\mathbf{a}$ , player  $i$  collects a reward  $r^i(s, \mathbf{a})$  ( $\mathbf{r}(s, \mathbf{a}) = (r^1(s, \mathbf{a}), \dots, r^N(s, \mathbf{a}))$ ) and move to a next state  $s'$  following the probability  $p(s'|s, \mathbf{a})$ . In such a setting, the usual goal is to find a strategy for each player that maximizes some long term reward. A strategy for player  $i$  is a function from state  $s$  to a distribution over actions  $A^i$  and is written  $\pi^i(\cdot|s)$ . The joint strategy will be written  $\boldsymbol{\pi} = (\pi^1, \dots, \pi^N) = (\pi^i, \boldsymbol{\pi}^{-i})$ . When the strategy of each player is fixed, the MG behaves like a Markov chain of kernel  $\mathcal{P}_{\boldsymbol{\pi}}(s'|s) = E_{\mathbf{a} \sim \boldsymbol{\pi}}[p(s'|s, \mathbf{a})]$  (note that  $\mathcal{P}_{\boldsymbol{\pi}}$  can be seen as a squared matrix of size  $|S| \times |S|$ ). The reward associated with that Markov chain is averaged over the joint strategy  $r_{\boldsymbol{\pi}}^i(s) = E_{\mathbf{a} \sim \boldsymbol{\pi}}[r^i(s, \mathbf{a})]$ . When the strategy of the  $i^{\text{th}}$  player's opponents is fixed, the process is reduced to an MDP of kernel  $p_{\boldsymbol{\pi}^{-i}}(s'|s, a^i) = E_{\mathbf{a}^{-i} \sim \boldsymbol{\pi}^{-i}}[p(s'|s, a^i, \mathbf{a}^{-i})]$  and reward function  $r_{\boldsymbol{\pi}^{-i}}^i(s, a^i) = E_{\mathbf{a}^{-i} \sim \boldsymbol{\pi}^{-i}}[r^i(s, a^i, \mathbf{a}^{-i})]$ .

**Multistage Game:** We consider games that can be modeled as trees (see Figure 1). In these games, the interaction between players start in an initial state  $\tilde{s}$  and proceeds in stages as illustrated in Figure 1 and terminates at state  $\Omega$  (a formal definition is given in Section A of the appendix).

**Value Function:** In a  $\gamma$ -discounted multistage game, (with  $\gamma \in (0, 1)$ ), each player's goal is to maximize the rewards it accumulates starting from any state  $s$ . This value is defined as follows:

$$v_{\boldsymbol{\pi}}^i(s) = E\left[\sum_{t=0}^{+\infty} \gamma^t r_{\boldsymbol{\pi}}^i(s_t) \mid s_0 = s, s_{t+1} \sim \mathcal{P}_{\boldsymbol{\pi}}(\cdot|s_t)\right], \quad (1)$$

and,

$$v_{\boldsymbol{\pi}}^i = \left(\mathcal{I} + \gamma \mathcal{P}_{\boldsymbol{\pi}} + \dots + \gamma^{|\mathcal{S}|} \mathcal{P}_{\boldsymbol{\pi}}^{|\mathcal{S}|}\right) r_{\boldsymbol{\pi}}^i. \quad (2)$$

Note that for a multistage game, the value function is well defined even for  $\gamma = 1$  since the process has an

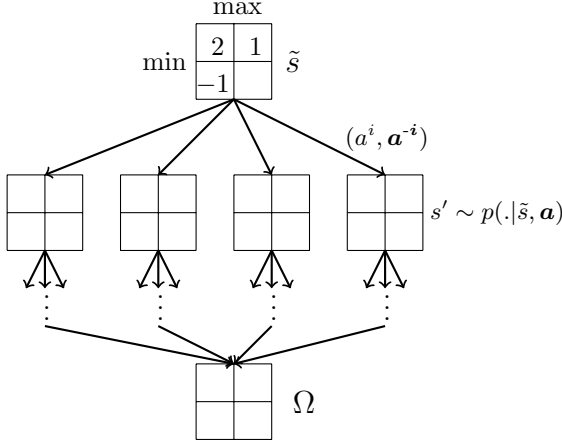


Figure 1: Example of a two-player zero-sum multistage game with deterministic dynamics.

absorbing state  $\Omega$  where the reward function is null. For  $t > |S|$ , we have that  $\mathcal{P}_\pi^t r_\pi^i = 0$  since after  $|S|$  steps (at most) the process ends up in the absorbing state  $\Omega$ . The matrix  $\mathcal{I}$  is the identity application.

**State-Action Value Function:** we can also define a state-action value function per player, that is the value of performing action  $a^i$  in state  $s$  and then following its strategy:

$$Q_\pi^i(s, a^i) = r_{\pi^{-i}}^i(s, a^i) + \gamma \sum_{s'} p_{\pi^{-i}}(s'|s, a^i) v_\pi^i(s'). \quad (3)$$

**Logit Choice Function:** when given a choice between actions  $a \in A$  based on an expected value  $Q(a)$  (in the following  $Q(a)$  will be the state-action value function  $Q_\pi^i(s, a^i)$  or  $Q_{\sigma, \pi^{-i}}^{*i}(s, a^i)$  defined below), it is common not to take the action that maximizes  $Q(a)$  but to choose suboptimal actions so as to favor exploration. One common choice of suboptimal actions is to pick them according to the logit choice function  $B_\eta(Q)(a) = \frac{\exp(\eta^{-1}Q(a))}{\sum_{\tilde{a}} \exp(\eta^{-1}Q(\tilde{a}))}$ . We will write  $C_\eta(Q) = \sum_{\tilde{a} \in A} Q(\tilde{a}) B_\eta(Q)(\tilde{a})$  the expected outcome if actions are taken according to the logit choice function. This definition can be generalized to the notion of choice probability function (see Section B in the appendix) and will be written  $B_\sigma(Q)$  and  $C_\sigma(Q)$  where  $\sigma$  is a function in  $\mathbb{R}^{\Delta A}$ . Here  $\Delta A$  is the set of distributions over actions. In the case of the logit choice function,  $\eta$  is linked to  $\sigma$  through the following formula  $\sigma(\pi) = \sum_a \eta \pi(a) \ln(\pi(a))$ .

**Best Responses:** If the strategy  $\pi^{-i}$  of all opponents is fixed, the value of the best response is  $v_{\sigma, \pi^{-i}}^{*i}$  and is recursively defined as follows (starting from  $\Omega$  and going in increasing order with respect to  $\phi$ ):

$v_{\sigma, \pi^{-i}}^{*i}(s) = C_\sigma(r_{\pi^{-i}}^i(s, a^i) + \sum_{s' \in S} p_{\pi^{-i}}(s'|s, a^i) v_{\sigma, \pi^{-i}}^{*i}(s'))$ . From the definition of that value function, we define the corresponding  $Q$ -function:

$$Q_{\sigma, \pi^{-i}}^{*i}(s, a^i) = r_{\pi^{-i}}^i(s, a^i) + \gamma \sum_{s'} p_{\pi^{-i}}(s'|s, a^i) v_{\sigma, \pi^{-i}}^{*i}(s'). \quad (4)$$

**Bellman Operators:** We define the two following Bellman operators on the  $Q$ -function:

$$\begin{aligned} [T_\pi^i Q^i](s, a^i) &= r_{\pi^{-i}}^i(s, a^i) \\ &+ \gamma \sum_{s'} p_{\pi^{-i}}(s'|s, a^i) E_{b^i \sim \pi^i(\cdot|s')} [Q^i(s', b^i)], \end{aligned} \quad (5)$$

and

$$\begin{aligned} [T_{\sigma, \pi^{-i}}^{*i} Q^i](s, a^i) &= r_{\pi^{-i}}^i(s, a^i) \\ &+ \gamma \sum_{s'} p_{\pi^{-i}}(s'|s, a^i) C_\sigma(Q^i(s', \cdot)). \end{aligned} \quad (6)$$

The value  $Q_\pi^i$  is the fixed point of the operator  $T_\pi^i$  and  $Q_{\sigma, \pi^{-i}}^{*i}$  is the fixed point of the operator  $T_{\sigma, \pi^{-i}}^{*i}$ . Furthermore, a strategy  $\tilde{\pi}^i$  is greedy with respect to a  $Q$ -function  $Q^i$  if  $T_{\sigma, \pi^{-i}}^{*i} Q^i = T_{\tilde{\pi}^i, \pi^{-i}}^i Q^i$ .

**Operators on the Value Function:** we define the counterparts of those operators on value functions as:

$$[T_\pi^i v^i](s) = r_\pi^i(s) + \gamma \sum_{s'} p_\pi(s'|s) v^i(s') = r_\pi^i + \gamma \mathcal{P}_\pi v^i \quad (7)$$

and

$$[T_{\sigma, \pi^{-i}}^{*i} v^i] = C_\sigma(r_{\pi^{-i}}^i(s, \cdot) + \gamma \sum_{s'} p_{\pi^{-i}}(s'|s, \cdot) v^i(s')). \quad (8)$$

From these definitions, we have the two following value functions:

$$v_{\sigma, \pi^{-i}}^{*i}(s) = C_\sigma(Q_{\sigma, \pi^{-i}}^{*i}(s, \cdot)) \quad (9)$$

and

$$v_\pi^i(s) = E_{b^i \sim \pi^i(\cdot|s)} [Q_\pi^i(s, b^i)]. \quad (10)$$

**Smooth Nash Equilibrium:** The goal in this setting is to find a strategy  $\pi^i$  for each player that recursively (in increasing order with respect to function  $\phi(\cdot)$ ) fulfills the following condition:  $\forall i, E_{a^i \sim \pi^i(\cdot|s)} [Q_\pi^i(s, \cdot)] = C_\sigma(Q_{\sigma, \pi}^i(s, \cdot))$ . As a consequence, we have that  $Q_{\sigma, \pi^{-i}}^{*i} = Q_\pi^i$  and  $v_\pi^i = v_{\sigma, \pi^{-i}}^{*i}$ . Furthermore, in the case of a zero-sum two-player game, we have  $v_\pi^i = -v_\pi^{-i}$  (where  $v_\pi^{-i}$  is the value function of the opponent)

## 4 Actor-Critic Fictitious Play

We present here our first contribution: the actor-critic fictitious play algorithm for MGs. This is an on-line and decentralized process. At each time-step  $n$ , all players are in state  $s_n$ , players choose independently an action  $a_n^i$  according to their current strategy  $\pi_n^i(\cdot|s_n)$  and observe independently one from another a reward signal  $r_n^i = r^i(s_n, \mathbf{a}_n)$ . The process is decentralized, meaning players do not observe others' actions nor their rewards. Then, the game moves to the following state  $s_{n+1}$ . If the process reaches the absorbing state  $\Omega$ , we simply restart from the beginning ( $\tilde{s}$ ).

---

### Algorithm 1 On-line Actor Critic Fictitious Play

---

**Input:** An initial strategy  $\pi_0^i$  and an initial value  $v_0^i = 0$ . Two learning rates  $\{\alpha_n\}_{n \geq 0}$ ,  $\{\beta_n\}_{n \geq 0}$  satisfying assumption **A 3** and an initial state  $s_0 = \tilde{s}$ .

**for**  $n=1,2,\dots$  **do**

Agent  $i$  draws action  $a_n^i \sim \pi^i(\cdot|s_n)$ .

Agent  $i$  observes reward  $r_n^i = r^i(s_n, \mathbf{a}_n)$ .

Every player observes the next state  $s_{n+1} \sim p(\cdot|s_n, \mathbf{a}_n)$ .

**actor step**

$\pi_{n+1}^i(s_n, \cdot) = (1 - \beta_n)\pi_n^i(s_n, \cdot) + \beta_n B_\sigma(Q_n^i(s_n, \cdot))$

**critic step**

Either an off-policy control step:

$Q_{n+1}^i(s_n, a_n^i) = (1 - \alpha_n)Q_n^i(s_n, a_n^i) + \alpha_n (r_n^i + C_\sigma(Q_n^i(s_{n+1}, \cdot)))$

Or a policy evaluation step:

$Q_{n+1}^i(s_n, a_n^i) = (1 - \alpha_n)Q_n^i(s_n, a_n^i) + \alpha_n (r_n^i + E_{b \sim \pi_n^i(\cdot|s_{n+1})}(Q_n^i(s_{n+1}, b)))$

**if**  $s_{n+1} = \Omega$  **then**

$s_{n+1} = \tilde{s}$ .

**end if**

**end for**

**Return** The joint strategy  $\pi$  and values  $v^i$  for all  $i$ .

---

The learning algorithm performs two updates. First, it updates the players' strategy (actors' update). The strategy  $\pi_{n+1}^i$  is a mixture between the current strategy  $\pi_n^i$  and either a local best response  $B_\sigma(Q_{\pi_n^i}^i(s_n, \cdot))$  or a global best response  $B_\sigma(Q_{\sigma, \pi_n^i}^{*i}(s_n, \cdot))$ . The actors' update is performed according to a slow timescale  $\beta_n$ . Second, it performs the critics' update which evaluates the current strategy. It happens on a fast timescale  $\alpha_n$  on which we can consider that the strategy of every player is stationary. If at the actors' step we want to act according to a local best response dynamics, the critic step will perform a policy evaluation step. If we want to perform a global best response dynamics, the critic will perform an off-policy evaluation step. Thus, at the slow timescale, the  $Q$ -function  $Q_n^i$  has almost

converged to  $Q_{\sigma, \pi_n^i}^{*i}$  or to  $Q_{\pi_n^i}^i$ . Therefore, we obtain canonically two algorithms.

In the next section, as a second contribution, we introduce the two dynamical processes corresponding to these algorithms and we prove that they possess desirable properties (*i.e.* rationality and convergence in self play for zero-sum two-player and cooperative games). The proofs rely non trivial techniques to propagate the Lyapunov stability property of the Fictitious play process of these ODE on the tree structure of the multistage games. Then in Sec. 6, as a third contribution, we show formally that these two algorithms (Algo.1) are stochastic approximations of those dynamical processes. This is again non-trivial as we need to prove the convergence of a two-time scale discrete scheme depending on a Markov chain.

## 5 Fictitious play in Markov Games

In this section, we propose novel definitions for two perturbed best response dynamics in the case of MGs. These dynamics are defined as a set of Ordinary Differential Equations (ODE) that generalizes the continuous time Fictitious play process to MSGs [17]. Then, we prove the convergence of these processes in multistage games. To do so, we build on the work of [17] on the stability of the Fictitious play process in normal form games. By induction on the tree structure of multistage games, we prove that our processes have stable attractors in zero-sum two-player and in cooperative multistage games. Later in Sec. 6, we will prove that our actor-critic algorithms track the solutions of that ODE and thus are guaranteed to converge in both settings. The first one considers a local best response dynamics:

$$\dot{\pi}_t^i(\cdot|s) = d_{\pi_t}(s)[B_\sigma(Q_{\pi_t}^i(s, \cdot)) - \pi_t^i(\cdot|s)] \quad (11)$$

The second process considers a global best response:

$$\dot{\pi}_t^i(\cdot|s) = d_{\pi_t}(s)[B_\sigma(Q_{\sigma, \pi_t^i}^{*i}(s, \cdot)) - \pi_t^i(\cdot|s)] \quad (12)$$

Two properties are usually desirable for such a stochastic process which are *rationality* and *convergence* in self-play [8]. Rationality implies that if other players converge to a stationary strategy, the learning algorithm will converge to a best response strategy. The convergence property received many definitions in the MARL literature and usually ensures convergence of the algorithm against a class of other algorithms or, as studied here, the convergence in self-play.

**Rationality:** For the rest of this paragraph, let us study a fixed player  $i$ . First, if we consider the case

where other players are stationary (i.e.  $\pi_t^{-i} = \pi^{-i}$ ), the strategy of player  $i$  will converge to  $B_\sigma(Q_{\sigma, \pi^{-i}}^{*i}(s, \cdot))$  for the second process Eq. (12) (since  $B_\sigma(Q_{\sigma, \pi^{-i}}^{*i}(s, \cdot))$  does not depend on  $\pi_t^i$  the solution of the second process converges exponentially toward the best response strategy). For the first process, let us show that if the strategy  $\pi^i$  follows the dynamics described by (11), it converges to a best response strategy. The proof of this property requires the following two technical lemmas.

**( $T, \delta$ )-perturbation:** Let us consider the following ODE where  $h(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a Lipschitz continuous function:  $\dot{z}(t) = h(z(t))$  (13)

We consider the case where ODE (13) has an asymptotically stable attractor set  $J$ . A bounded measurable function  $y(\cdot) : \mathbb{R}^+ \rightarrow \mathbb{R}^n$  is a  $(T, \delta)$ -perturbation of (13) ( $T, \delta > 0$ ) if there exist  $0 = T_0 < T_1 < T_2 < \dots < T_n$  with  $T_{i+1} - T_i \geq T$  and solutions  $z^j(t)$ ,  $t \in [T_j, T_{j+1}]$  of (13) for all  $j \geq 0$  such that  $\sup_{t \in [T_j, T_{j+1}]} \|z^j(t) - y(t)\| < \delta$ . The following technical lemma can be proved [6].

**Lemma 1.** *Given  $\epsilon, T > 0$ , there exists a  $\bar{\delta} > 0$  such that for  $\delta \in (0, \bar{\delta})$ , every  $(T, \delta)$ -perturbation of Eq. (13) converges to  $J^\epsilon$  (def: the  $\epsilon$ -neighborhood of  $J$ ).*

The next lemma is used all over this section. It proves that if, instead of following a given process  $\dot{x}(t) = f(x(t))$ , one follows a perturbed version of that process  $\dot{y}(t) = f_t(y(t))$  one converges to the set of stable equilibria of the unperturbed process if  $f_t(\cdot)$  converges toward  $f(\cdot)$ .

**Lemma 2.** *Let's study the two following ODE  $\dot{x}(t) = f(x(t))$  and  $\dot{y}(t) = f_t(y(t))$  and let's assume that  $\limsup \|f(\cdot) - f_t(\cdot)\| \rightarrow 0$  uniformly. Furthermore, let's assume that  $f$  is Lipschitz continuous. Then, any bounded solution of  $\dot{y}(t) = f_t(y(t))$  converges to the set of attractors of  $\dot{x}(t) = f(x(t))$ .*

*Proof.* The full proof of this lemma is left in appendix C.  $\square$

Now, we show that if player  $i$  follows the direction given by the local best response dynamics Eq. (11), it converges to a best response to  $\pi^{-i}$ .

**Proposition 1.** *In any game, if the strategy  $\pi^{-i}$  of the opponents is fixed, the process (11) converges to a best response.*

*Proof.* The proof of this proposition is left in appendix G  $\square$

Now, let's consider the case where the strategy of other players  $\pi_t^{-i}$  converges to a stationary strategy  $\pi^{-i}$ . In

that case, one can show that, for a fixed  $i$ , the solutions of ODE (11) and (12) are the same as if  $\pi_t^{-i}$  was fixed.

**Proposition 2.** *If  $\pi_t^{-i} \rightarrow \pi^{-i}$ , then the solutions of ODE (11) and (12) are the same as the one of:*

$$\dot{\pi}_t^i(\cdot|s) = d_{\pi_t^i, \pi^{-i}}(s)[B_\sigma(Q_{\pi_t^i, \pi^{-i}}^i(s, \cdot)) - \pi_t^i(\cdot|s)]$$

and

$$\dot{\pi}_t^i(\cdot|s) = d_{\pi_t^i, \pi^{-i}}(s)[B_\sigma(Q_{\sigma, \pi^{-i}}^{*i}(s, \cdot)) - \pi_t^i(\cdot|s)]$$

*Proof.* The proof of this proposition is left in appendix H  $\square$

Thus, if all other players converge to a stationary strategy, the strategy of player  $i$  converges to a best response to  $\pi^{-i}$ .

**Convergence in Multistage Games:** Concerning the convergence in self-play of those two fictitious play processes (ODE (11) and (12)), one must first be aware that they will not converge in all multistage games. Even in normal form games, the fictitious play process is known to be unstable in some cases. For example, the fictitious play process is known to oscillate in some two-player games [30]. More surprisingly, [15] proved in the case of normal form games that there exists no uncoupled dynamics which guarantees Nash convergence. In their paper, they consider dynamics for each player which only depends on its own reward and possibly on the strategy of all players. They present an example in which this class of dynamical systems has no stable equilibrium. We were able to prove the convergence of our process in the case of zero-sum two-player and cooperative multistage games.

**Proposition 3.** *In a zero-sum two player multistage game, the process (11) and the process (12) converge to a smooth Nash equilibrium  $\pi$ .*

*Proof.* The proof of this proposition is left in appendix I  $\square$

**Remark 1.** *Proposition 3 can be adapted to study the cooperative case (i.e. when players receive the same reward  $\forall i, j$   $r^i(s, \mathbf{a}) = r^j(s, \mathbf{a})$ ). The proof in the cooperative case also works by induction. The corresponding proposition and proof can be found in the appendix E.*

## 6 Stochastic Approximation with Two-Timescale

In this section, we provide a novel stochastic approximation theorem taking advantage of two independent techniques. First, we use two timescales because the

process we are studying (defined on the strategy of each player called the *actor*) has a complex dependency on the strategy of all players through the  $Q$ -function (the *critic*). In an on-line setting, one can't have access to the  $Q$ -function of a given strategy (either  $Q_\pi^i$  or  $Q_{\sigma, \pi^i}^{*i}$ ) as it is an asymptotic solution of the critic process. Thus formally, we look for an asymptotically stable solution of  $\dot{y}(t) = g(\mu(y(t)), y(t))$  (actors' update) where  $\mu(y)$  is the stable solution of another process  $\dot{x}(t) = f(x(t), y)$  (critics' update). In our case  $\mu(y)$  is either  $Q_\pi^i$  or  $Q_{\sigma, \pi^i}^{*i}$ ,  $y$  is the strategy and  $x$  the action-value function. Instead of waiting until the subroutine converges (the critic part) to iterate over the main routine (the actor part), [6] gives an elegant solution to that class of problems by using two timescales to update simultaneously  $x_n$  and  $y_n$ . The process  $x_n$  will be updated according to a "fast" timescale  $\alpha_n$ . On that timescale,  $y_n$  behaves as if it was stationary and thus,  $x_n$  is an estimate of  $\mu(y_n)$  (i.e.  $\|\mu(y_n) - x_n\| \rightarrow 0$ ). The process  $y_n$  will move on a "slower" timescale  $\beta_n$ . On that slower timescale, one can treat the process  $x_n$  as  $\mu(y_n)$  and thus,  $y_n$  will converge to a stable solution of  $\dot{y}(t) = g(\mu(y(t)), y(t))$ .

Second we use an averaging technique. Our process is on-line and follows a Markovian dynamics of stationary distribution  $d_{\pi_t}$  controlled by the policy  $\pi_n$ . All policies and  $Q$ -functions are only updated on state  $s_n$ . Thus, the two processes  $x_n$  and  $y_n$  also depend on a Markov process  $Z_n$  controlled by  $y_n$ . Again, [3] shows that, in the case of a simple timescale, the process  $y_{n+1} = y_n + \alpha_n f(y_n, Z_n)$  tracks the solution of  $\dot{y}(t) = f(y(t), d_{y(t)})$  where  $\bar{f}(y(t), d_{y(t)})$  is the average of  $f$  over the stationary distribution of  $Z_n$  (the distribution  $d_y$ ). Formally,  $\bar{f}(y, d_y) = \sum_{z \in \mathcal{Z}} f(y, z) d_y(z)$ . In a nutshell, we generalize that result on stochastic approximation with a controlled Markov noise to two timescale stochastic approximation. In our case, we need to study the following recursion:

$$x_{n+1} = x_n + \alpha_n f(x_n, y_n, Z_n), \quad (14)$$

$$y_{n+1} = y_n + \beta_n g(x_n, y_n, Z_n), \quad (15)$$

Where:

**A 1.** *Functions  $f$  and  $g$  are jointly continuous in their arguments and Lipschitz in their two first arguments uniformly with respect to the third.*

**A 2.** *The controlled Markov process  $Z_n$  takes its value in a discrete space  $\mathcal{Z}$  controlled by variable  $y_n$ . The variable  $Z_{n+1}$  follows the kernel  $p(\cdot | Z_n, y_n)$  which is uniformly continuous in  $y_n$ . Furthermore, let us suppose that if  $y_n = y$ , the Markov chain  $Z_n$  has a unique invariant distribution  $d_y(\cdot)$ .*

We note  $\bar{f}(x, y, d_y) = \sum_{z \in \mathcal{Z}} f(x, y, z) d_y(z)$  the function  $f$  averaged over the stationary distribution of

the Markov chain defined in the previous assumption. Similarly,  $\bar{g}(x, y, d_y) = \sum_{z \in \mathcal{Z}} g(x, y, z) d_y(z)$ .

**A 3.** *The sequences  $\{\alpha_n\}_{n \geq 0}$  and  $\{\beta_n\}_{n \geq 0}$  are two positives decreasing step-size sequences satisfying:  $\sum_{n \geq 0} \alpha_n = \sum_{n \geq 0} \beta_n = \infty$ ,  $\sum_{n \geq 0} \alpha_n^2 < \infty$ ,  $\sum_{n \geq 0} \beta_n^2 < \infty$  and  $\beta_n = o(\alpha_n)$*

**A 4.** *We need  $\sup_n \|x_n\| < \infty$  and  $\sup_n \|y_n\| < \infty$*

**A 5.** *For any constant  $y$ , the ODE:  $\dot{x}(t) = \bar{f}(x(t), y, d_y)$ , has a globally asymptotically stable equilibrium  $\mu(y)$ . That stable equilibrium  $\mu(\cdot)$  is a Lipschitz continuous function.*

**A 6.** *The ODE  $\dot{y}(t) = \bar{g}(\mu(y(t)), y(t), d_{y(t)})$  has a globally asymptotically stable equilibrium  $y^*$*

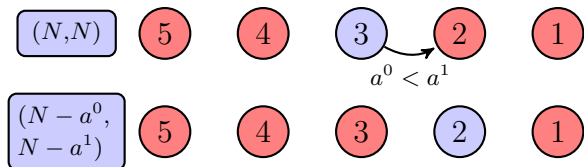
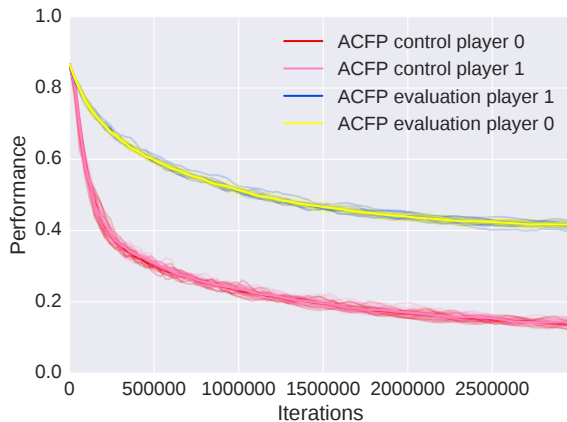
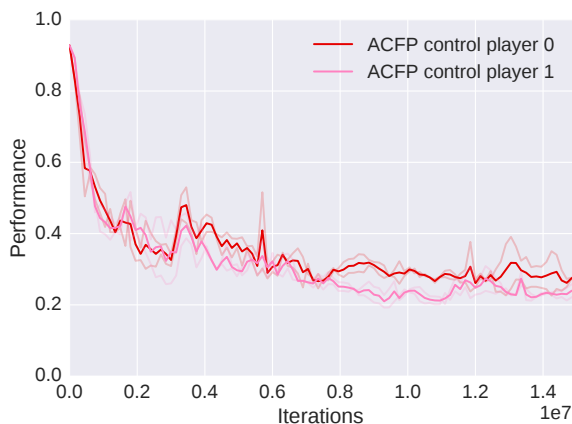
**Theorem 1.**  *$(x_n, y_n) \rightarrow (\mu(y^*), y^*)$  almost surely.*

**Analysis of Actor-Critic Fictitious Play:** The analysis of our two algorithms is a direct application of the previous theorem and is left in appendix (App. J). In a nutshell, the critic  $Q_n$  is updated on a fast timescale (i.e.  $x_n$ ) and we will prove that it's recursion satisfy assumption A 1, A 2, A 3, A 4, and A 5. The actor  $\pi_n$  is updated on a slow timescale (i.e.  $y_n$ ) and satisfy assumption A 1, A 2, A 3, A 4, and A 6. Thanks to the use of choice probability function, assumption A 1 is satisfied on the slow timescale since a choice probability function is lipschitz with respect to the  $Q$ -function.

## 7 Experiment on Alesia

Alesia [26], which resembles the Oshi-Zumo game in [10, 7] (meaning "the pushing sumo"), is a two-player board game where each player starts with  $N$  coins. They are positioned in the middle of a board with  $2K + 1$  different positions (see Figure 2). At each round, each player chooses secretly a certain amount of coins to bid (let's say  $(a^0, a^1)$ ). If the player's budget is not null, s/he has to bet at least one coin. If player 0's bid is larger (respectively smaller) than the one of player 1, player 0 (respectively player 1) moves toward his side of the board. If the bids are equal, then players will stay on their current position. This process continues until players have no coins left or until one player reaches one side of the board. The final position determines the winner. The game ends with a draw if no one succeeded to reach his side of the board. If player 1 (respectively player zero) reaches his side of the board, the reward is  $(-1, 1)$  (respectively  $(1, -1)$ ).

This game is challenging in many aspects. First, rewards are sparse and are received at the end of the game. Furthermore, strategies need to be stochastic [10, 7]. The performance criterion used to compare algorithms is the difference of the value of the


 Figure 2: Alesia rules for  $K = 2$ .

 Figure 3: Performance of the strategy  $\pi^i$  ( $y$ -axis) along iterations ( $x$ -axis).

 Figure 4: Performance of the strategy  $\pi^i$  ( $y$ -axis) along iterations ( $x$ -axis).

joint strategy  $v_{\pi}^i(\tilde{s})$  and the value of the best response  $v_{0, \pi^{-i}}^i(\tilde{s})$  without any perturbation  $\sigma$  (*i.e.*  $\sigma = 0$ )

We ran experiments for a game with learning rates of the form  $\alpha_n = \frac{\alpha_0 a_\alpha}{a_\alpha + n b_\alpha}$  and  $\beta_n = \frac{\beta_0 a_\beta}{a_\beta + n b_\beta}$ . We used a logit choice function with  $\eta$  decaying with the number of iterations  $\eta_n = \frac{\eta_0 a_\eta}{a_\eta + n}$ . The two step-sizes  $\alpha_n$  and  $\beta_n$  were chosen to satisfy the conditions of Theorem 1. The experiments were ran with an initial budget of  $N = 10$  and a board of size  $K = 3$ . The step-size parameters were  $(\alpha_0, a_\alpha, b_\alpha) = (0.1, 10^4, 0.8)$ ,

$(\beta_0, a_\beta, b_\beta) = (0.01, 10^4, 1.0)$  and  $(\eta_0, a_\eta) = (0.1, 10^6)$ . The results are displayed in figure 3 and shows that both algorithms learn a strategy that tend to be closer to its own best response over learning (meaning that both algorithms learn a Nash equilibrium). An empirical finding is that ACFP with a control critic converges faster than the one with a policy evaluation critic on that game.

**Function approximation:** We also ran these experiments with function approximations both on the strategy and on the  $Q$ -function. The main issue in Alesia is that the number of actions available to each player depends on the remaining budget of the player. For each player, we aggregated states for which the budget of the player is fixed. The aggregated states are consecutive states with respect to the budget of the other players'. The results reported in Figure 4 show that our method is robust to this form of function approximation. We ran experiment with an initial budget of  $N = 15$  on which the algorithm was able to learn a minimax strategy. The step-size parameters were  $(\alpha_0, a_\alpha, b_\alpha) = (0.1, 10^5, 0.9)$ ,  $(\beta_0, a_\beta, b_\beta) = (0.01, 10^5, 1.0)$  and  $(\eta_0, a_\eta) = (0.05, 10^8)$ .

## 8 Conclusion

To summarize, we defined two novel algorithms. These algorithms are on-line and decentralized procedures that provably converge in two classes of multistage games: zero-sum two-player and cooperative games. Compared to previous family of methods [8, 25], these algorithms can be applied to a larger number of settings without heuristics. This paper opens several interesting research avenues. The first is the study of the stability of the two dynamical processes in MGs. The proof given for multistage games propagates the convergence property from the end of the game to the initial state. Since MGs have a graph structure, this technique can't be applied. Thus, other stability arguments must be found to ensure convergence. The main contributions of this paper were algorithmic and theoretical but we also proposed an empirical evaluation that provides good evidence that our methods are applicable to real problems and could scale up thanks to function approximation. Whilst there is a wide literature dealing with the off-policy evaluation part or the off-policy control [32], the update of the strategy using general function approximation remains challenging. The main issue is that the strategy space should support convex combination with a best response. Boosting methods such as in [29] could offer a solution. The convex combination of the strategy and the best response could also be done thanks to a classification neural network.



## Acknowledgments

We wish to thank Marc Lanctot for the helpful conversations, his comments definitely helped us improving the final version of the paper. This work has received partial funding from CPER Nord-Pas de Calais/FEDER DATA Advanced data science and technologies 2015-2020 and the European Commission H2020 framework program under the research Grant number 687831 (BabyRobot).

## References

- [1] N. Akchurina. Multiagent Reinforcement Learning: Algorithm Converging to Nash Equilibrium in General-Sum Discounted Stochastic Games. In *Proc. of AAMAS*, 2009.
- [2] B. Banerjee and J. Peng. Adaptive policy gradient in multiagent learning. In *Proc. AAMAS*. ACM, 2003.
- [3] Vivek S Borkar. Stochastic approximation with controlled markovnoise. *Systems & control letters*, 2006.
- [4] Vivek S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge University Press, 2009.
- [5] V.S. Borkar. Stochastic Approximation with Two Time Scales. *Systems & Control Letters*, 29(5):291–294, 1997.
- [6] V.S. Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 1997.
- [7] B. Božanský, V. Lisý, M. Lanctot, J. Čermák, and M.H.M. Winands. Algorithms for computing strategies in two-player simultaneous move games. *Artificial Intelligence*, 237:1 – 40, 2016.
- [8] M. Bowling and M. Veloso. Rational and Convergent Learning in Stochastic Games. In *Proc. of IJCAI*, 2001.
- [9] Sbastien Bubeck and Nicol Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- [10] M. Buro. *Solving the Oshi-Zumo Game*, pages 361–366. Springer US, 2004.
- [11] L. Busoni, R. Babuska, and B. De Schutter. A Comprehensive Survey of Multiagent Reinforcement Learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, 2008.
- [12] Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006.
- [13] Jakob N. Foerster, Yannis M. Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Proc. of NIPS*, 2016.
- [14] A. Greenwald, K. Hall, and R. Serrano. Correlated Q-learning. In *Proc. of ICML*, 2003.
- [15] S. Hart and A. Mas-Colell. Uncoupled dynamics do not lead to nash equilibrium. *The American Economic Review*, 2003.
- [16] J. Heinrich, M. Lanctot, and D. Silver. Fictitious self-play in extensive-form games. In *Proc. of ICML*, 2015.
- [17] J. Hofbauer and W.H. Sandholm. On the global convergence of stochastic fictitious play. *Econometrica*, 70(6):2265–2294, 2002.
- [18] J. Hu and M. P. Wellman. Nash Q-Learning for General-Sum Stochastic Games. *Journal of Machine Learning Research*, 4:1039–1069, 2003.
- [19] Vojtěch Kovařík and Viliam Lisý. Analysis of hannan consistent selection for monte carlo tree search in simultaneous move games. *arXiv preprint arXiv:1509.00149*, 2015.
- [20] Marc Lanctot, Kevin Waugh, Martin Zinkevich, and Michael Bowling. Monte carlo sampling for regret minimization in extensive games. In *Proc. of NIPS*, 2009.
- [21] Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, karl Tuyls, Julien Perolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent reinforcement learning. In *Proc. of NIPS*, 2017.
- [22] Guillaume J. Laurent, Laëticia Matignon, and N. Le Fort-Piat. The world of independent learners is not Markovian. *Int. J. Know.-Based Intell. Eng. Syst.*, pages 55–64, 2011.
- [23] D.S. Leslie and E.J. Collins. Generalised weakened fictitious play. *Games and Economic Behavior*, 56(2):285–298, 2006.
- [24] Viliam Lisý, Vojta Kovarik, Marc Lanctot, and Branislav Bosansky. Convergence of monte carlo tree search in simultaneous move games. In *Proc. of NIPS*, 2013.
- [25] M. L. Littman. Markov Games as a Framework for Multi-Agent Reinforcement Learning. In *Proc. of ICML*, 1994.

- [26] J. Perolat, B. Scherrer, B. Piot, and O. Pietquin. Approximate Dynamic Programming for Two-Player Zero-Sum Markov Games. In *Proc. of ICML*, 2015.
- [27] HL Prasad, Prashanth LA, and Shalabh Bhatnagar. Two-Timescale Algorithms for Learning Nash Equilibria in General-Sum Stochastic Games. In *Proc. of AAMAS*, 2015.
- [28] J. Robinson. An iterative method of solving a game. *Annals of mathematics*, pages 296–301, 1951.
- [29] B. Scherrer and M. Geist. Local Policy Search in a Convex Space and Conservative Policy Iteration as Boosted Policy Search. In *Proc. of ECML*, 2014.
- [30] L.S. Shapley. Some topics in two-person games. *Advances in game theory*, 1964.
- [31] Y. Shoham, R. Powers, and T. Grenager. If multi-agent learning is the answer, what is the question? *Artificial Intelligence*, 171(7):365 – 377, 2007.
- [32] C. J. C. H. Watkins and P. Dayan. Q-learning. *Machine Learning*, 8(3):279–292, 1992.
- [33] J. W. Weibull. *Evolutionary game theory*. MIT press, 1997.
- [34] Martin Zinkevich, Michael Johanson, Michael H Bowling, and Carmelo Piccione. Regret minimization in games with incomplete information. In *Proc. of NIPS*, 2007.

## A Multistage Game:

We consider games that can be modeled as trees (see Fig. 1). The state  $\tilde{s}$  is the root of the game tree and each state can be reached with some non-zero probability with at least a deterministic joint strategy. Furthermore, we enforce the fact that once a state is visited, it can never be visited again except state  $\Omega$ , the end of the game. Formally, a multistage game is an MG with an initial state  $\tilde{s} \in S$  and an absorbing state  $\Omega \in S$ . To enforce that a state can be only visited once, we define an ordering bijection  $\phi : S \rightarrow \{1, \dots, |S|\}$  where  $|S|$  is the cardinal of  $S$ ,  $\phi(\tilde{s}) = |S|$ ,  $\phi(\Omega) = 1$  and such that  $\forall s, s' \in S \times S \setminus \{(\Omega, \Omega)\}$ ,  $\phi(s) \leq \phi(s') \Rightarrow \forall \pi, \mathcal{P}_\pi(s'|s) = 0$ . Moreover the reward in state  $\Omega$  is null (meaning  $\mathbf{r}(\Omega, \cdot) = 0$ ) and no player has more than one action available in that absorbing state. Furthermore, for any state  $s \in S \setminus \{\tilde{s}, \Omega\}$  there exists a deterministic strategy  $\pi$  and a time  $t \leq |S|$  such that  $\mathcal{P}_\pi^t(s|\tilde{s}) > 0$ . This condition means that every states can be reached from state  $\tilde{s}$  with at least a deterministic strategy.

## B Generalized best response

In the previous section, we defined the logit choice function. One can generalize this notion with the concept of choice probability function (see [17]). Imagine a player is given a choice of an action  $a \in A$  based on the outcome  $Q(a)$ . Instead of choosing based on the outcome  $Q(a)$ , he also obtains an additive random payoff  $\epsilon_a$  and chooses with respect to  $Q(a) + \epsilon_a$ . The vector  $(\epsilon_a)_{a \in A}$  is a positive random variable taking it's value in  $\mathbb{R}^{|A|}$  and does not depend on  $Q(a)$ . These choice probability functions allow us to get smooth Bellman operators and ease the analysis of our algorithms. A choice probability function is defined as follows:  $B_\sigma(Q)(a) = P(\operatorname{argmax}_{\tilde{a} \in A} [Q(\tilde{a}) + \epsilon_{\tilde{a}}] = a)$ . This definition is equivalent to the following one where  $\sigma(\cdot)$  is a deterministic perturbation:  $B_\sigma(Q) = \operatorname{argmax}_{\tilde{\pi} \in \Delta(A)} [E_{\tilde{a} \sim \tilde{\pi}} [Q(\tilde{a})] - \sigma(\tilde{\pi})]$ . The function  $C_\sigma(Q)$  is the average value considering the choice probability function  $B_\sigma(Q)$ :  $C_\sigma(Q) = \sum_{\tilde{a} \in A} Q(\tilde{a}) B_\sigma(Q)(\tilde{a})$ . This perturbation is said admissible [17] if for all  $y$ ,  $\nabla^2 \sigma(y)$  is positive definite on  $\mathbb{R}_0 = \{z \in \mathbb{R}^n : \sum_j z_j = 0\}$ , and if  $\|\nabla \sigma(y)\|$  goes to infinity as  $y$  approaches the boundary of  $\Delta A$ . For example,  $B_\sigma(Q)(a) = \frac{\exp(\eta^{-1} Q(a))}{\sum_{\tilde{a}} \exp(\eta^{-1} Q(\tilde{a}))}$  if the  $\sigma(\cdot)$  perturbation is defined as  $\sigma(\pi) = \sum_a \eta \pi(a) \ln(\pi(a))$ . In that case, the noise  $\epsilon_{\tilde{a}}$  follows a Gumbel distribution. From now on, we consider that  $\sigma(\cdot)$  is fixed.

## C Proof of Lemma 2

**Lemma 2:** *Let's study the two following ODE  $\dot{x}(t) = f(x(t))$  and  $\dot{y}(t) = f_t(y(t))$  and let's assume that  $\limsup \|f(\cdot) - f_t(\cdot)\| \rightarrow 0$  uniformly. Furthermore, let's assume that  $f$  is Lipschitz continuous. Then, any solution of  $\dot{y}(t) = f_t(y(t))$  converges to the set of attractors of  $\dot{x}(t) = f(x(t))$ .*

*Proof.* let us fix  $t$  such as for all  $T > 0$ ,  $\|f(\cdot) - f_{t+T}(\cdot)\| \leq \epsilon$ . if  $x(t) = y(t)$ , then:

$$x(t+T) - y(t+T) = \int_t^{t+T} [f(x(\tau)) - f(y(\tau))] d\tau + \int_t^{t+T} [f(y(\tau)) - f_\tau(y(\tau))] d\tau \quad (16)$$

$$\|(x - y)(t+T)\| \leq \int_t^{t+T} \|f(x(\tau)) - f(y(\tau))\| d\tau + \int_t^{t+T} \|f(y(\tau)) - f_\tau(y(\tau))\| d\tau \quad (17)$$

$$\leq K \int_t^{t+T} \|(x - y)(\tau)\| d\tau + T\epsilon \quad (18)$$

Let's write  $g(T) = \|(x - y)(t+T)\|$ . We have  $g(T) \leq K \int_0^T g(\tau) d\tau + T\epsilon$

Let's write  $h(T) = T\epsilon + K \int_0^T g(\tau) d\tau + \frac{\epsilon}{K}$  and  $h'(T) = \epsilon + Kg(T) \leq \epsilon + K(h(T) - \frac{\epsilon}{K}) \leq Kh(T)$

With the differential form of the Grönwall lemma:

$$h(\delta) \leq h(0) \exp(KT) = \frac{\epsilon}{K} \exp(KT)$$

And we have  $\|(x - y)(t + T)\| = g(T) \leq \frac{\epsilon}{K} (\exp(KT) - 1)$

This inequality means that, given  $T$ , for all  $\delta$  we can choose  $t_0$  large enough such that any trajectory of  $\dot{y}(t) = f_{t_0+t}(y(t))$  is a  $(T, \delta)$ -perturbation of  $\dot{x}(t) = f(x(t))$ . Then, lemma 1 concludes the proof.  $\square$

## D Proof of Theorem 1

$$x_{n+1} = x_n + \alpha_n f(x_n, y_n, Z_n), \quad (19)$$

$$y_{n+1} = y_n + \beta_n g(x_n, y_n, Z_n), \quad (20)$$

Where:

**A 1.** Functions  $f$  and  $g$  are jointly continuous in their arguments and Lipschitz in their two first arguments uniformly with respect to the third,

**A 2.** The controlled Markov process  $Z_n$  takes its value in a discrete space  $\mathcal{Z}$  controlled by variable  $y_n$ . The variable  $Z_{n+1}$  follows the transition kernel  $p(\cdot | Z_n, y_n)$  which is uniformly continuous in  $y_n$ . Furthermore, let us suppose that if  $y_n = y$ , the Markov chain  $Z_n$  has a unique invariant distribution  $d_y(\cdot)$ .

From now on, we define  $\bar{f}(x, y, d_y) = \sum_{z \in \mathcal{Z}} f(x, y, z) d_y(z)$  which is the function  $f$  averaged over the stationary distribution of the Markov chain defined in the previous assumption. Similarly,  $\bar{g}(x, y, d_y) = \sum_{z \in \mathcal{Z}} g(x, y, z) d_y(z)$ .

**A 3.** The sequences  $\{\alpha_n\}_{n \geq 0}$  and  $\{\beta_n\}_{n \geq 0}$  are two positives decreasing step-size sequences satisfying:  $\sum_{n \geq 0} \alpha_n = \infty$ ,  $\sum_{n \geq 0} \beta_n = \infty$ ,  $\sum_{n \geq 0} \alpha_n^2 < \infty$ ,  $\sum_{n \geq 0} \beta_n^2 < \infty$  and  $\beta_n = o(\alpha_n)$

**A 4.** We need  $\sup_n \|x_n\| < \infty$  and  $\sup_n \|y_n\| < \infty$

**A 5.** For any constant  $y$ , the ODE:

$$\frac{dx(t)}{dt} = \bar{f}(x(t), y, d_y), \quad (21)$$

has a globally asymptotically stable equilibrium  $\mu(y)$ . That stable equilibrium  $\mu(\cdot)$  is a Lipschitz continuous function.

**A 6.** The ODE:

$$\frac{dy(t)}{dt} = \bar{g}(\mu(y(t)), y(t), d_{y(t)}) \quad (22)$$

has a globally asymptotically stable equilibrium  $y^*$

**Theorem 1.**  $(x_n, y_n) \rightarrow (\mu(y^*), y^*)$  almost surely (a.s.).

*Proof.* First, rewrite Eq. (15) as:

$$y_{n+1} = y_n + \alpha_n \left[ \frac{\beta_n}{\alpha_n} g(x_n, y_n, Z_n) \right].$$

Since the function  $g$  is Lipschitz in the two first arguments (A 1), since  $\mathcal{Z}$  is discrete (A 2) and since  $x_n$  and  $y_n$  are bounded (A 4), then we have that  $\frac{\beta_n}{\alpha_n} g(x_n, y_n, Z_n) \rightarrow 0$  a.s.. Then from corollary 8<sup>1</sup> of chapter 6.3 of [4] (first presented in [3]) it follows that  $(x_n, y_n)$  converges to an internally chain transitive invariant set of the ODE  $\dot{y}(t) = 0$  and  $\dot{x}(t) = \bar{f}(x(t), y(t), d_{y(t)})$ . In other words,  $\|x_n - \mu(y_n)\| \rightarrow 0$  a.s..

<sup>1</sup>As written in chapter 2 [4] in all stochastic approximation scheme studied in the book (except those in chapter 9) can be added a noise  $\epsilon_n$  that converges to 0 a.s.

Second, let's write (15) as:

$$y_{n+1} = y_n + \alpha_n [g(\mu(y_n), y_n, Z_n) + (g(x_n, y_n, Z_n) - g(\mu(y_n), y_n, Z_n))]. \quad (23)$$

Since  $g(\cdot, \cdot, \cdot)$  is Lipschitz in the two first variables uniformly with respect to the third one, we have that  $\|g(x_n, y_n, Z_n) - g(\mu(y_n), y_n, Z_n)\| \leq K \|x_n - \mu(y_n)\|$  (for some  $K$ ). Then Eq. (15) can be rewritten as:

$$y_{n+1} = y_n + \alpha_n [g(\mu(y_n), y_n, Z_n) + \epsilon_n], \quad (24)$$

where  $\epsilon_n \rightarrow 0$  a.s.. Again, from corollary 8 of chapter 6.3 of [4], we have that  $y_n \rightarrow y^*$  and  $x_n \rightarrow \mu(y^*)$   $\square$

## E Convergence in Cooperative Multistage Games

In a cooperative game (as defined in [11]), the reward signal is as follows:  $\forall s, i, j, r^i(s, \mathbf{a}) = r^j(s, \mathbf{a})$ . Thus we have the property that  $\forall s, i, j, v_{\pi}^i(s) = v_{\pi}^j(s)$

**Proposition 4.** *In a cooperative two player multistage game, the process (11) and the process (12) converge to a smooth Nash equilibrium  $\pi$ .*

*Proof.* The proof of this result works again by induction on the set  $S_n = \{s \in S | \phi(s) \leq n\}$ . Let's suppose that for all states  $s$  in  $S_n$ , the process converges to a smooth Nash equilibrium. This means that for all states in  $s \in S_n$  and for all players  $i$ , the strategies  $\pi_t^i(\cdot | s)$  converge to  $\pi^i(\cdot | s)$  such as  $v_{\pi}^i(s) = v_{\sigma, \pi^{-i}}^{*i}(s)$ . We also have that  $\forall i, j, v_{\pi}^i(s) = v_{\pi}^j(s)$ .

Let  $\hat{s}$  be the state such that  $\phi(\hat{s}) = n + 1$ . Then, we define:

$$M^i(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s' | \hat{s}, a^i, \mathbf{a}^{-i}) v_{\pi}^i(s')$$

$$M_t^i(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s' | \hat{s}, a^i, \mathbf{a}^{-i}) v_{\pi_t}^i(s')$$

And:

$$M_t^{*i}(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s' | \hat{s}, a^i, \mathbf{a}^{-i}) v_{\sigma, \pi_t^{-i}}^{*i}(s')$$

Since the strategy  $\pi_t^i(\cdot | s)$  converges to  $\pi^i(\cdot | s)$ , we have  $v_{\pi_t}^i(s)$  and  $v_{\sigma, \pi_t^{-i}}^{*i}(s)$  that converges to  $v_{\pi}^i(s)$  for all  $s \in S_n$  and finally  $M_t^i(\hat{s}, a^i, \mathbf{a}^{-i})$  and  $M_t^{*i}(\hat{s}, a^i, \mathbf{a}^{-i})$  converges to  $M^i(\hat{s}, a^i, \mathbf{a}^{-i})$ . Then, lemma 2 and results of convergence of stochastic fictitious play for  $N$ -player potential games from [17] guarantees that the process in state  $\hat{s}$  will converge to a smooth Nash equilibrium of the normal form game defined by  $M^i(\hat{s}, a^i, \mathbf{a}^{-i})$ . Finally, we have  $v_{\pi}^i(\hat{s}) = v_{\sigma, \pi^{-i}}^{*i}(\hat{s})$  and  $\forall i, j, v_{\pi}^i(\hat{s}) = v_{\pi}^j(\hat{s})$ .  $\square$

## F Remark on the Stationary Distribution

The distribution  $d_{\pi}$  is the stationary distribution of the Markov process defined by  $s' \sim p_{\pi}(\cdot | s)$  and if  $s' = \Omega$ , we restart the process from  $\tilde{s}$ . One can show that if we start the process with a strategy  $\pi_0$  that gives a non-zero probability for each action and each player, the distribution over states  $d_{\pi_t}(s)$  is never zero since we consider that there is at least a deterministic strategy that reaches any state. Furthermore, since we consider smooth best response dynamics, the smooth best response will assign to each action some probability which is bounded away from zero (since the  $Q$ -function are bounded). Thus, we can always consider that we will visit each state with some minimal probability (i.e.  $d_{\pi_t}(s) \geq \delta > 0$ ).

## G Proof of Proposition 1

*Proof.* We prove the result by induction on the set  $S_n = \{s \in S | \phi(s) \leq n\}$ . First, the property is true in state  $\Omega$  since by definition, there is only one action available per player. Suppose that for all  $s \in S_n$ , the process converges to a best response (i.e.  $v_{\pi}^i(s) = v_{\sigma, \pi^{-i}}^{*i}$ ). Let  $\hat{s}$  be the state of order  $n + 1$  (i.e.  $\phi(\hat{s}) = n + 1$ ). From the definition of the  $Q$ -function (Eq. (3) and (4)), we have that in  $\hat{s}$  the  $Q_{\pi_t^i, \pi^{-i}}^i(\hat{s}, \cdot)$  converges to  $Q_{\pi}^i(\hat{s}, \cdot)$  uniformly over the actions. Moreover, since for all  $s \in S_n$ ,  $v_{\pi}^i(s) = v_{\sigma, \pi^{-i}}^{*i}$  we have that  $Q_{\pi}^i(\hat{s}, \cdot) = Q_{\sigma, \pi^{-i}}^{*i}(\hat{s}, \cdot)$ . From lemma 2, we get the convergence of  $\pi^i(\cdot | \hat{s})$  to a best response if the distribution  $d_{\pi}$  is non-null in all states (see Section F).  $\square$

## H Proof of Proposition 2

*Proof.* The proof comes from the fact that  $B_\sigma(Q_{\pi^i, \pi^{-i}}^i(s, \cdot))$  converges to  $B_\sigma(Q_{\pi^i, \pi^{-i}}^i(s, \cdot))$  uniformly with respect to  $\pi^i$  (since the  $Q$ -function is polynomial in  $\pi$  and  $B_\sigma(\cdot)$  is Lipschitz with respect to the  $Q$ -function. Moreover, since the state space is finite, simple convergence imply the uniform convergence) and the fact that  $B_\sigma(Q_{\pi^i, \pi^{-i}}^i(s, \cdot))$  is Lipschitz with respect to  $\pi^i$ .  $\square$

## I Proof of Proposition 3

*Proof.* The proof of this result works by induction on the set  $S_n = \{s \in S | \phi(s) \leq n\}$ . Again, the property is true in state  $\Omega$ . Let's suppose that for all states  $s$  in  $S_n$ , the process converges to a smooth Nash equilibrium. This means that for all states in  $s \in S_n$  and for all players  $i$ , the strategies  $\pi_t^i(\cdot|s)$  converge to  $\pi^i(\cdot|s)$  such as  $v_\pi^i(s) = v_{\sigma, \pi^{-i}}^{*i}(s)$ . We also have that  $v_\pi^i(s) = -v_\pi^{-i}(s)$ .

Let  $\hat{s}$  be the state such that  $\phi(\hat{s}) = n + 1$ . Then, we define:

$$M^i(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s'| \hat{s}, a^i, \mathbf{a}^{-i}) v_\pi^i(s')$$

$$M_t^i(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s'| \hat{s}, a^i, \mathbf{a}^{-i}) v_{\pi_t^i}^i(s')$$

And:

$$M_t^{*i}(\hat{s}, a^i, \mathbf{a}^{-i}) = r^i(\hat{s}, a^i, \mathbf{a}^{-i}) + \sum_{s' \in S} p(s'| \hat{s}, a^i, \mathbf{a}^{-i}) v_{\sigma, \pi_t^{-i}}^{*i}(s')$$

Since the strategy  $\pi_t^i(\cdot|s)$  converges to  $\pi^i(\cdot|s)$ , we have  $v_{\pi_t^i}^i(s)$  and  $v_{\sigma, \pi_t^{-i}}^{*i}(s)$  that converges to  $v_\pi^i(s)$  for all  $s \in S_n$  and finally  $M_t^i(\hat{s}, a^i, \mathbf{a}^{-i})$  and  $M_t^{*i}(\hat{s}, a^i, \mathbf{a}^{-i})$  converges to  $M^i(\hat{s}, a^i, \mathbf{a}^{-i})$ . Then, lemma 2 and results of convergence of stochastic fictitious play from [17] guarantees that the process in state  $\hat{s}$  will converge to a smooth Nash equilibrium of the normal form game defined by  $M^i(\hat{s}, a^i, \mathbf{a}^{-i})$ . Finally, we have  $v_\pi^i(\hat{s}) = v_{\sigma, \pi^{-i}}^{*i}(\hat{s})$  and  $v_\pi^i(\hat{s}) = -v_\pi^{-i}(\hat{s})$ .  $\square$

## J Analysis of Actor-Critic Fictitious Play

The two on-line algorithms we present here (Algo 1) are stochastic approximations of the two dynamical systems (11) and (12). In an on-line setting, the interaction between players proceeds as follows. At a step  $n$ , players are in state  $s_n$  and individually take an action  $a_n^i \sim \pi_n^i(\cdot|s_n)$ . Each of them receives a reward  $r^i(s_n, a_n^1, \dots, a_n^N)$  and the process moves from state  $s_n$  to state  $s_{n+1} \sim p(\cdot|s_n, a_n^1, \dots, a_n^N)$ . In addition, we consider the controlled Markov process  $Z_n = (s_n, a_n^1, \dots, a_n^N, s_{n+1})$  (controlled by the strategy  $\pi_n$ ). Finally, if  $s_{n+1} = \Omega$ , we restart from state  $\hat{s}$  and  $Z_{n+1} = (\hat{s}, \dots)$ . We define the following two-timescale processes (with  $\alpha_n$  and  $\beta_n$  defined as in A 3):

$$\pi_{n+1}^i = \pi_n^i + \beta_n \mathbf{1}_{s_n} \mathbf{1}_{s_n}^\top \left[ B_\sigma(Q_n^i(s_n, \cdot)) - \pi_n^i(\cdot|s_n) \right] \quad (25)$$

And:

$$Q_{n+1}^i = Q_n^i + \alpha_n \mathbf{1}_{(s_n, a_n^i)} \mathbf{1}_{(s_n, a_n^i)}^\top \left[ r^i(s_n, \mathbf{a}_n) + C_\sigma \left( Q_n^i(s_{n+1}, \cdot) \right) - Q_n^i(s_n, a_n^i) \right] \quad (26)$$

Or:

$$Q_{n+1}^i = Q_n^i + \alpha_n \mathbf{1}_{(s_n, a_n^i)} \mathbf{1}_{(s_n, a_n^i)}^\top \left[ r^i(s_n, \mathbf{a}_n) + E_{b^i \sim \pi^i(\cdot|s_{n+1})} \left( Q_n^i(s_{n+1}, b^i) \right) - Q_n^i(s_n, a_n^i) \right] \quad (27)$$

Assumption A 1 and assumption A 2 are verified. Regarding assumption A 4, the strategy  $\pi_n^i$  is obviously bounded since it remains in the simplex. If  $Q_0^i = 0$ , the state-action value function is bounded as follows  $\|Q_n^i(s, a)\| \leq \phi(s) R_{\max}$  (this can be shown by recursion). The two faster processes (Eq. (26) and (27)) track the two following ODE:

$$\dot{Q}_t^i(s, a^i) = d_\pi(s) \pi^i(\cdot|s) \left( \left[ T_{\sigma, \pi^{-i}}^{*i} Q_t^i \right] (s, a^i) - Q_t^i(s, a^i) \right)$$

And

$$\dot{Q}_t^i(s, a^i) = d_\pi(s) \pi^i(\cdot|s) \left( \left[ T_\pi^i Q_t^i \right] (s, a^i) - Q_t^i(s, a^i) \right)$$

Those two equations admit as an attractor  $Q_{\sigma^*, \pi^{-i}}^*$  and  $Q_\pi^i$ . Then, the strategy recursion follows either ODE (11) (if the subroutine is defined by Eq. (27)) or ODE (12) (if the subroutine is defined by (26)).

## K On the Guarantees of Convergence of OFF-SGSP and ON-SGSP

The main contribution to the field of MARL related to our work is the paper of Prasad & al [27]. These algorithms are not decentralized but one of them, ON-SGSP, is an on-line and model-free algorithm. This paper proposes two algorithms OFF-SGSP and ON-SGSP which are stochastic approximations of a dynamical system described in section 8. The authors claims that these algorithms converges to a Nash equilibrium of the game. The proof of the stability given in lemma 11 is wrong. In the following we point out the issue with this proof.

Using their notations:

- $G$  is the set of Nash equilibrium (the feasible set of the optimization problem),
- $K$  is the limit set of the dynamical system (and  $G \subset K$ )
- $K_1$  is the set of limits points of the dynamical system which are Nash equilibrium  $K \cap G$
- $K_2$  is the complementary of  $K_1$  in  $K$  (i.e.  $K_2 = K \setminus K_1$ )

Lemma 11 shows that  $K_2$  contains only unstable equilibrium and conclude that both processes converges to  $K_1$  and thus to a Nash equilibrium since  $K_1$  is not empty (this fact is proven early in the paper).

Unfortunately, that proof contains a mistake. The proof proceeds as follows: They show that if  $\pi^* \in K_2$  then there exists  $a^i, x$  and  $i$  such that  $g_{x, a^i}^i(\mathbf{v}_\pi^i, \pi^{-i}) > 0$  and they conclude that consequently  $\frac{\partial f(\mathbf{v}_\pi, \pi)}{\partial \pi^i} < 0$ . However, there is no direct link between the sine of  $g_{x, a^i}^i(\mathbf{v}_\pi^i, \pi^{-i})$  and  $\frac{\partial f(\mathbf{v}_\pi, \pi)}{\partial \pi^i}$  since  $f(\mathbf{v}_\pi, \pi) = \sum_{i=1}^N \sum_{x \in S} \sum_{z \in A^i(x)} \pi^i(x, z) g_{x, z}^i(\mathbf{v}_\pi^i, \pi^{-i})$ .

This imply that both processes might converge in  $K_2$  (i.e. not to a Nash equilibrium).