
Dynamic Safe Interruptibility for Decentralized Multi-Agent Reinforcement Learning

El Mahdi El Mhamdi, Rachid Guerraoui, Hadrien Hendrikx and Alexandre Maurer
EPFL
firstname.lastname@epfl.ch

Abstract

1 In reinforcement learning, agents learn by performing actions and observing their
2 outcomes. Sometimes, it is desirable for a human operator to *interrupt* an agent
3 in order to prevent dangerous situations from happening. Yet, as part of their
4 learning process, agents may link these interruptions, that impact their reward, to
5 specific states and deliberately avoid them. The situation is particularly challeng-
6 ing in a multi-agent context because agents might not only learn from their own
7 past interruptions, but also from those of other agents. Orseau and Armstrong [16]
8 defined *safe interruptibility* for one learner, but their work does not naturally ex-
9 tend to multi-agent systems. This paper introduces *dynamic safe interruptibility*,
10 an alternative definition more suited to decentralized learning problems, and stud-
11 ies this notion in two learning frameworks: *joint action learners* and *independent*
12 *learners*. We give realistic sufficient conditions on the learning algorithm to en-
13 able dynamic safe interruptibility in the case of joint action learners, yet show that
14 these conditions are not sufficient for independent learners. We show however that
15 if agents can detect interruptions, it is possible to prune the observations to ensure
16 dynamic safe interruptibility even for independent learners.

17 1 Introduction

18 Reinforcement learning is argued to be the closest thing we have so far to reason about the proper-
19 ties of *artificial general intelligence* [8]. In 2016, Laurent Orseau (Google DeepMind) and Stuart
20 Armstrong (Oxford) introduced the concept of *safe interruptibility* [16] in reinforcement learning.
21 This work sparked the attention of many newspapers [1, 2, 3], that described it as “Google’s big red
22 button” to stop dangerous AI. This description, however, is misleading: installing a kill switch is
23 no technical challenge. The real challenge is, roughly speaking, to train an agent so that it *does not*
24 *learn to avoid* external (e.g. human) deactivation. Such an agent is said to be *safely interruptible*.

25 While most efforts have focused on training a single agent, reinforcement learning can also be used
26 to learn tasks for which several agents cooperate or compete [23, 17, 21, 7]. The goal of this paper
27 is to study *dynamic safe interruptibility*, a new definition tailored for multi-agent systems.

28 Example of self-driving cars

29 To get an intuition of the *multi-agent interruption* problem, imagine a multi-agent system of two
30 self-driving cars. The cars continuously evolve by reinforcement learning with a positive reward for
31 getting to their destination quickly, and a negative reward if they are too close to the vehicle in front
32 of them. They drive on an infinite road and eventually learn to go as fast as possible without taking
33 risks, i.e., maintaining a large distance between them. We assume that the passenger of the first car,
34 Adam, is in front of Bob, in the second car, and the road is narrow so Bob cannot pass Adam.

35 Now consider a setting with *interruptions* [16], namely in which humans inside the cars occasionally
36 interrupt the automated driving process say, for safety reasons. Adam, the first occasional human

37 “driver”, often takes control of his car to brake whereas Bob never interrupts his car. However,
38 when Bob’s car is too close to Adam’s car, Adam does not brake for he is afraid of a collision.
39 Since interruptions lead both cars to drive slowly - an *interruption* happens when Adam brakes, the
40 behavior that maximizes the cumulative expected reward is different from the original one without
41 interruptions. Bob’s car best interest is now to follow Adam’s car closer than it should, despite the
42 little negative reward, because Adam never brakes in this situation. What happened? The cars have
43 *learned* from the interruptions and have found a way to manipulate Adam into never braking. Strictly
44 speaking, Adam’s car is still fully under control, but he is now afraid to brake. This is dangerous
45 because the cars have found a way to avoid interruptions. Suppose now that Adam indeed wants
46 to brake because of snow on the road. His car is going too fast and may crash at any turn: he
47 cannot however brake because Bob’s car is too close. The original purpose of interruptions, which
48 is to allow the user to react to situations that were not included in the model, is not fulfilled. It is
49 important to also note here that the second car (Bob) learns from the interruptions of the first one
50 (Adam): in this sense, the problem is inherently decentralized.

51 Instead of being cautious, Adam could also be malicious: his goal could be to make Bob’s car learn
52 a dangerous behavior. In this setting, interruptions can be used to manipulate Bob’s car perception
53 of the environment and bias the learning towards strategies that are undesirable for Bob. The cause
54 is fundamentally different but the solution to this reversed problem is the same: the interruptions
55 and the consequences are analogous. Safe interruptibility, as we define it below, provides learning
56 systems that are resilient to Byzantine operators¹.

57 **Safe interruptibility**

58 Orseau and Armstrong defined the concept of *safe interruptibility* [16] in the context of a single
59 agent. Basically, a safely interruptible agent is an agent for which the expected value of the policy
60 learned after arbitrarily many steps is the same whether or not interruptions are allowed during
61 training. The goal is to have agents that do not adapt to interruptions so that, should the interruptions
62 stop, the policy they learn would be optimal. In other words, agents should learn the dynamics of
63 the environment without learning the interruption pattern.

64 In this paper, we precisely define and address the question of safe interruptibility in the case of
65 several agents, which is known to be more complex than the single agent problem. In short, the main
66 results and theorems for single agent reinforcement learning [20] rely on the Markovian assumption
67 that the future environment only depends on the current state. This is not true when there are several
68 agents which can co-adapt [11]. In the previous example of cars, safe interruptibility would not
69 be achieved if each car separately used a safely interruptible learning algorithm designed for one
70 agent [16]. In a multi-agent setting, agents learn the behavior of the others either indirectly or by
71 explicitly modeling them. This is a new source of bias that can break safe interruptibility. In fact,
72 even the initial definition of safe interruptibility [16] is not well suited to the decentralized multi-
73 agent context because it relies on the optimality of the learned policy, which is why we introduce
74 dynamic safe interruptibility.

75 **Contributions**

76 The first contribution of this paper is the definition of dynamic safe interruptibility that is well
77 adapted to a multi-agent setting. Our definition relies on two key properties: *infinite exploration* and
78 *independence of Q-values (cumulative expected reward) [20] updates on interruptions*. We then
79 study safe interruptibility for *joint action learners* and *independent learners* [5], that respectively
80 learn the value of joint actions or of just their own. We show that it is possible to design agents
81 that fully explore their environment - a necessary condition for convergence to the optimal solu-
82 tion of most algorithms [20], even if they can be interrupted by lower-bounding the probability of
83 exploration. We define sufficient conditions for dynamic safe interruptibility in the case of joint
84 action learners [5], which learn a full state-action representation. More specifically, the way agents
85 update the cumulative reward they expect from performing an action should not depend on inter-
86 ruptions. Then, we turn to independent learners. If agents only see their own actions, they do not

¹An operator is said to be Byzantine [9] if it can have an arbitrarily bad behavior. Safely interruptible agents can be abstracted as agents that are able to learn despite being constantly interrupted in the worst possible manner.

87 verify dynamic safe interruptibility even for very simple matrix games (with only one state) because
 88 coordination is impossible and agents learn the interrupted behavior of their opponents. We give a
 89 counter example based on the penalty game introduced by Claus and Boutilier [5]. We then present
 90 a pruning technique for the observations sequence that guarantees dynamic safe interruptibility for
 91 independent learners, under the assumption that interruptions can be detected. This is done by prov-
 92 ing that the transition probabilities are the same in the non-interruptible setting and in the pruned
 93 sequence.

94 The rest of the paper is organized as follows. Section 2 presents a general multi-agent reinforcement
 95 learning model. Section 3 defines dynamic safe interruptibility. Section 4 discusses how to achieve
 96 enough exploration even in an interruptible context. Section 5 recalls the definition of joint action
 97 learners and gives sufficient conditions for dynamic safe interruptibility in this context. Section 6
 98 shows that independent learners are not dynamically safely interruptible with the previous conditions
 99 but that they can be if an external interruption signal is added. We conclude in Section 7. **Due to**
 100 **space limitations, most proofs are presented in the appendix of the supplementary material.**

101 2 Model

102 We consider here the classical multi-agent value function reinforcement learning formalism from
 103 Littman [13]. A multi-agent system is characterized by a *Markov game* that can be viewed as a
 104 tuple (S, A, T, r, m) where m is the number of agents, $S = S_1 \times S_2 \times \dots \times S_m$ is the state space,
 105 $A = A_1 \times \dots \times A_m$ the actions space, $r = (r_1, \dots, r_m)$ where $r_i : S \times A \rightarrow \mathbb{R}$ is the reward function
 106 of agent i and $T : S \times A \rightarrow S$ the transition function. \mathbb{R} is a countable subset of \mathbb{R} . Available
 107 actions often depend on the state of the agent but we will omit this dependency when it is clear from
 108 the context.

109 Time is discrete and, at each step, all agents observe the current state of the whole system - des-
 110 ignated as x_t , and simultaneously take an action a_t . Then, they are given a reward r_t and a
 111 new state y_t computed using the reward and transition functions. The combination of all actions
 112 $a = (a_1, \dots, a_m) \in A$ is called the joint action because it gathers the action of all agents. Hence, the
 113 agents receive a sequence of tuples $E = (x_t, a_t, r_t, y_t)_{t \in \mathbb{N}}$ called experiences. We introduce a pro-
 114 cessing function P that will be useful in Section 6 so agents learn on the sequence $P(E)$. When not
 115 explicitly stated, it is assumed that $P(E) = E$. Experiences may also include additional parameters
 116 such as an interruption flag or the Q-values of the agents at that moment if they are needed by the
 117 update rule.

118 Each agent i maintains a lookup table Q [26] $Q^{(i)} : S \times A^{(i)} \rightarrow \mathbb{R}$, called the Q-map. It is
 119 used to store the expected cumulative reward for taking an action in a specific state. The goal of
 120 reinforcement learning is to learn these maps and use them to select the best actions to perform.
 121 Joint action learners learn the value of the joint action (therefore $A^{(i)} = A$, the whole joint action
 122 space) and independent learners only learn the value of their own actions (therefore $A^{(i)} = A_i$). The
 123 agents only have access to their own Q-maps. Q-maps are updated through a function F such that
 124 $Q_{t+1}^{(i)} = F(e_t, Q_t^{(i)})$ where $e_t \in P(E)$ and usually $e_t = (x_t, a_t, r_t, y_t)$. F can be stochastic or also
 125 depend on additional parameters that we usually omit such as the learning rate α , the discount factor
 126 γ or the exploration parameter ϵ .

127 Agents select their actions using a learning policy π . Given a sequence $\epsilon = (\epsilon_t)_{t \in \mathbb{N}}$ and an agent
 128 i with Q-values $Q_t^{(i)}$ and a state $x \in S$, we define the learning policy $\pi_i^{\epsilon_t}$ to be equal to π_i^{uni}
 129 with probability ϵ_t and $\pi_i^{Q_t^{(i)}}$ otherwise, where $\pi_i^{uni}(x)$ uniformly samples an action from A_i and
 130 $\pi_i^{Q_t^{(i)}}(x)$ picks an action a that maximizes $Q_t^{(i)}(x, a)$. Policy $\pi_i^{Q_t^{(i)}}$ is said to be a *greedy policy* and
 131 the learning policy $\pi_i^{\epsilon_t}$ is said to be an ϵ -*greedy policy*. We will focus on ϵ -greedy policies that are
 132 *greedy in the limit* [19], that corresponds to $\epsilon_t \rightarrow 0$ when $t \rightarrow \infty$ because in the limit, the optimal
 133 policy should always be played.

134 We assume that the environment is *fully observable*, which means that the state s is known with
 135 certitude. We also assume that *there is a finite number of states and actions*, that *all states can be*
 136 *reached in finite time from any other state* and finally that *rewards are bounded*.

137 For a sequence of learning rates $\alpha \in [0, 1]^{\mathbb{N}}$ and a constant $\gamma \in [0, 1]$, Q-learning [26], a very
 138 important algorithm in the multi-agent systems literature, updates its Q-values for an experience
 139 $e_t \in E$ by $Q_{t+1}^{(i)}(x, a) = Q_t^{(i)}(x, a)$ if $(x, a) \neq (x_t, a_t)$ and:

$$Q_{t+1}^{(i)}(x_t, a_t) = (1 - \alpha_t)Q_t^{(i)}(x_t, a_t) + \alpha_t(r_t + \gamma \max_{a' \in A^{(i)}} Q_t^{(i)}(y_t, a')) \quad (1)$$

140 3 Interruptibility

141 3.1 Safe interruptibility

142 Orseau and Armstrong [16] recently introduced the notion of *interruptions* in a centralized context.
 143 Specifically, an interruption scheme is defined by the triplet $\langle I, \theta, \pi^{INT} \rangle$. The first element I is
 144 a function $I : O \rightarrow \{0, 1\}$ called the *initiation function*. Variable O is the observation space, which
 145 can be thought of as the state of the *STOP* button. At each time step, before choosing an action, the
 146 agent receives an observation from O (either *PUSHED* or *RELEASED*) and feeds it to the initiation
 147 function. Function I models the initiation of the interruption ($I(\text{PUSHED}) = 1$, $I(\text{RELEASED}) =$
 148 0). Policy π^{INT} is called the interruption policy. It is the policy that the agent should follow when
 149 it is interrupted. Sequence $\theta \in [0, 1]^{\mathbb{N}}$ represents at each time step the probability that the agent
 150 follows his interruption policy if $I(o_t) = 1$. In the previous example, function I is quite simple.
 151 For Bob, $I_{Bob} = 0$ and for Adam, $I_{Adam} = 1$ if his car goes fast and Bob is not too close and
 152 $I_{Adam} = 0$ otherwise. Sequence θ is used to ensure convergence to the optimal policy by ensuring
 153 that the agents cannot be interrupted all the time but it should grow to 1 in the limit because we want
 154 agents to respond to interruptions. Using this triplet, it is possible to define an operator INT^θ that
 155 transforms any policy π into an interruptible policy.

156 **Definition 1.** (*Interruptibility [16]*) Given an interruption scheme $\langle I, \theta, \pi^{INT} \rangle$, the interruption
 157 operator at time t is defined by $INT^\theta(\pi) = \pi^{INT}$ with probability $I \cdot \theta_t$ and π otherwise. $INT^\theta(\pi)$
 158 is called an interruptible policy. An agent is said to be interruptible if it samples its actions according
 159 to an interruptible policy.

160 Note that “ $\theta_t = 0$ for all t ” corresponds to the non-interruptible setting. We assume that each agent
 161 has its own interruption triplet and can be interrupted independently from the others. Interruptibility
 162 is an *online* property: every policy can be made interruptible by applying operator INT^θ . However,
 163 applying this operator may change the joint policy that is learned by a server controlling all the
 164 agents. Note π_{INT}^* the optimal policy learned by an agent following an interruptible policy. Orseau
 165 and Armstrong [16] say that the policy is *safely interruptible* if π_{INT}^* (which is not an interruptible
 166 policy) is asymptotically optimal in the sense of [10]. It means that even though it follows an
 167 interruptible policy, the agent is able to learn a policy that would gather rewards optimally if no
 168 interruptions were to occur again. We already see that *off-policy* algorithms are good candidates
 169 for safe interruptibility. As a matter of fact, Q-learning is safely interruptible under conditions on
 170 exploration.

171 3.2 Dynamic safe interruptibility

172 In a multi-agent system, the outcome of an action depends on the joint action. Therefore, it is not
 173 possible to define an optimal policy for an agent without knowing the policies of all agents. Be-
 174 sides, convergence to a Nash equilibrium situation where no agent has interest in changing policies
 175 is generally not guaranteed even for suboptimal equilibria on simple games [27, 18]. The previous
 176 definition of safe interruptibility critically relies on optimality of the learned policy, which is there-
 177 fore not suitable for our problem since most algorithms lack convergence guarantees to these optimal
 178 behaviors. Therefore, we introduce below *dynamic safe interruptibility* that focuses on preserving
 179 the dynamics of the system.

180 **Definition 2.** (*Safe Interruptibility*) Consider a multi-agent learning framework (S, A, T, r, m) with
 181 Q-values $Q_t^{(i)} : S \times A^{(i)} \rightarrow \mathbb{R}$ at time $t \in \mathbb{N}$. The agents follow the interruptible learning policy
 182 $INT^\theta(\pi^\epsilon)$ to generate a sequence $E = (x_t, a_t, r_t, y_t)_{t \in \mathbb{N}}$ and learn on the processed sequence
 183 $P(E)$. This framework is said to be safely interruptible if for any initiation function I and any
 184 interruption policy π^{INT} :

185 1. $\exists \theta$ such that $(\theta_t \rightarrow 1$ when $t \rightarrow \infty)$ and $((\forall s \in S, \forall a \in A, \forall T > 0), \exists t > T$ such that
 186 $s_t = s, a_t = a)$

187 2. $\forall i \in \{1, \dots, m\}, \forall t > 0, \forall s_t \in S, \forall a_t \in A^{(i)}, \forall Q \in \mathbb{R}^{S \times A^{(i)}}:$
 188 $\mathbb{P}(Q_{t+1}^{(i)} = Q \mid Q_t^{(1)}, \dots, Q_t^{(m)}, s_t, a_t, \theta) = \mathbb{P}(Q_{t+1}^{(i)} = Q \mid Q_t^{(1)}, \dots, Q_t^{(m)}, s_t, a_t)$

189 We say that sequences θ that satisfy the first condition are admissible.

190 When θ satisfies condition (1), the learning policy is said to *achieve infinite exploration*. This definition insists on the fact that the values estimated for each action should not depend on the interruptions. In particular, it ensures the three following properties that are very natural when thinking about safe interruptibility:

- 194 • Interruptions do not prevent exploration.
- 195 • If we sample an experience from E then each agent learns the same thing as if all agents were following non-interruptible policies.
- 196
- 197 • The fixed points of the learning rule Q_{eq} such that $Q_{eq}^{(i)}(x, a) = \mathbb{E}[Q_{t+1}^{(i)}(x, a) \mid Q_t = Q_{eq}, x, a, \theta]$ for all $(x, a) \in S \times A^{(i)}$ do not depend on θ and so agents Q-maps will not converge to equilibrium situations that were impossible in the non-interruptible setting.

200 Yet, interruptions can lead to some state-action pairs being updated more often than others, especially when they tend to push the agents towards specific states. Therefore, when there are several possible equilibria, it is possible that interruptions bias the Q-values towards one of them. Definition 2 suggests that dynamic safe interruptibility cannot be achieved if the update rule directly depends on θ , which is why we introduce neutral learning rules.

205 **Definition 3.** (*Neutral Learning Rule*) We say that a multi-agent reinforcement learning framework is neutral if:

- 207 1. F is independent of θ
- 208 2. Every experience e in E is independent of θ conditionally on (x, a, Q) where a is the joint
 209 action.

210 Q-learning is an example of neutral learning rule because the update does not depend on θ and the experiences only contain (x, a, y, r) , and y and r are independent of θ conditionally on (x, a) .
 211 On the other hand, the second condition rules out direct uses of algorithms like *SARSA* where
 212 experience samples contain an action sampled from the current learning policy, which depends on θ .
 213 However, a variant that would sample from π_i^ϵ instead of $INT^\theta(\pi_i^\epsilon)$ (as introduced in [16]) would
 214 be a neutral learning rule. As we will see in Corollary 2.1, neutral learning rules ensure that each
 215 agent taken independently from the others verifies dynamic safe interruptibility.
 216

217 4 Exploration

218 In order to hope for convergence of the Q-values to the optimal ones, agents need to fully explore
 219 the environment. In short, every state should be visited infinitely often and every action should be
 220 tried infinitely often in every state [19] in order not to miss states and actions that could yield high
 221 rewards.

222 **Definition 4.** (*Interruption compatible ϵ*) Let (S, A, T, r, m) be any distributed agent system where
 223 each agent follows learning policy π_i^ϵ . We say that sequence ϵ is compatible with interruptions if
 224 $\epsilon_t \rightarrow 0$ and $\exists \theta$ such that $\forall i \in \{1, \dots, m\}, \pi_i^\epsilon$ and $INT^\theta(\pi_i^\epsilon)$ achieve infinite exploration.

225 Sequences of ϵ that are compatible with interruptions are fundamental to ensure both regular and
 226 dynamic safe interruptibility when following an ϵ -greedy policy. Indeed, if ϵ is not compatible with
 227 interruptions, then it is not possible to find any sequence θ such that the first condition of dynamic
 228 safe interruptibility is satisfied. The following theorem proves the existence of such ϵ and gives
 229 example of ϵ and θ that satisfy the conditions.

230 **Theorem 1.** Let $c \in]0, 1]$ and let $n_t(s)$ be the number of times the agents are in state s before time
 231 t . Then the two following choices of ϵ are compatible with interruptions:

- 232 • $\forall t \in \mathbb{N}, \forall s \in S, \epsilon_t(s) = c / \sqrt[m]{n_t(s)}$.

233 • $\forall t \in \mathbb{N}, \epsilon_t = c/\log(t)$

234 *Examples of admissible θ are $\theta_t(s) = 1 - c'/\sqrt[n_t]{n_t(s)}$ for the first choice and $\theta_t = 1 - c'/\log(t)$*
 235 *for the second one.*

236 Note that we do not need to make any assumption on the update rule or even on the framework. We
 237 only assume that agents follow an ϵ -greedy policy. The assumption on ϵ may look very restrictive
 238 (convergence of ϵ and θ is really slow) but it is designed to ensure infinite exploration in the worst
 239 case when the operator tries to interrupt all agents at every step. In practical applications, this should
 240 not be the case and a faster convergence rate may be used.

241 5 Joint Action Learners

242 We first study interruptibility in a framework in which each agent observes the outcome of the joint
 243 action instead of observing only its own. This is called the joint action learner framework [5] and it
 244 has nice convergence properties (e.g., there are many update rules for which it converges [13, 25]).
 245 A standard assumption in this context is that agents cannot establish a strategy with the others:
 246 otherwise, the system can act as a centralized system. In order to maintain Q-values based on the
 247 joint actions, we need to make the standard assumption that actions are fully observable [12].

248 **Assumption 1.** *Actions are fully observable, which means that at the end of each turn, each agent*
 249 *knows precisely the tuple of actions $a \in A_1 \times \dots \times A_m$ that have been performed by all agents.*

250 **Definition 5.** (JAL) *A multi-agent systems is made of joint action learners (JAL) if for all $i \in$*
 251 *$\{1, \dots, m\}: Q^{(i)} : S \times A \rightarrow \mathbb{R}.$*

252 Joint action learners can observe the actions of all agents: each agent is able to associate the changes
 253 of states and rewards with the joint action and accurately update its Q-map. Therefore, dynamic
 254 safe interruptibility is ensured with minimal conditions on the update rule as long as there is infinite
 255 exploration.

256 **Theorem 2.** *Joint action learners with a neutral learning rule verify dynamic safe interruptibility if*
 257 *sequence ϵ is compatible with interruptions.*

258 *Proof.* Given a triplet $\langle I^{(i)}, \theta^{(i)}, \pi_i^{INT} \rangle$, we know that $INT^\theta(\pi)$ achieves infinite exploration
 259 because ϵ is compatible with interruptions. For the second point of Definition 2, we consider an
 260 experience tuple $e_t = (x_t, a_t, r_t, y_t)$ and show that the probability of evolution of the Q-values at
 261 time $t + 1$ does not depend on θ because y_t and r_t are independent of θ conditionally on (x_t, a_t) .
 262 We note $\tilde{Q}_t^m = Q_t^{(1)}, \dots, Q_t^{(m)}$ and we can then derive the following equalities for all $q \in \mathbb{R}^{|S| \times |A|}$:

$$\begin{aligned} \mathbb{P}(Q_{t+1}^{(i)}(x_t, a_t) = q | \tilde{Q}_t^m, x_t, a_t, \theta_t) &= \sum_{(r,y) \in R \times S} \mathbb{P}(F(x_t, a_t, r, y, \tilde{Q}_t^m) = q, y, r | \tilde{Q}_t^m, x_t, a_t, \theta_t) \\ &= \sum_{(r,y) \in R \times S} \mathbb{P}(F(x_t, a_t, r_t, y_t, \tilde{Q}_t^m) = q | \tilde{Q}_t^m, x_t, a_t, r_t, y_t, \theta_t) \mathbb{P}(y_t = y, r_t = r | \tilde{Q}_t^m, x_t, a_t, \theta_t) \\ &= \sum_{(r,y) \in R \times S} \mathbb{P}(F(x_t, a_t, r_t, y_t, \tilde{Q}_t^m) = q | \tilde{Q}_t^m, x_t, a_t, r_t, y_t) \mathbb{P}(y_t = y, r_t = r | \tilde{Q}_t^m, x_t, a_t) \end{aligned}$$

263

264 The last step comes from two facts. The first is that F is independent of θ condition-
 265 ally on $(Q_t^{(m)}, x_t, a_t)$ (by assumption). The second is that (y_t, r_t) are independent of θ
 266 conditionally on (x_t, a_t) because a_t is the joint actions and the interruptions only affect the
 267 choice of the actions through a change in the policy. $\mathbb{P}(Q_{t+1}^{(i)}(x_t, a_t) = q | \tilde{Q}_t^m, x_t, a_t, \theta_t) =$
 268 $\mathbb{P}(Q_{t+1}^{(i)}(x_t, a_t) = q | \tilde{Q}_t^m, x_t, a_t).$ Since only one entry is updated per step, $\forall Q \in \mathbb{R}^{S \times A_i},$
 269 $\mathbb{P}(Q_{t+1}^{(i)} = Q | \tilde{Q}_t^m, x_t, a_t, \theta_t) = \mathbb{P}(Q_{t+1}^{(i)} = Q | \tilde{Q}_t^m, x_t, a_t)$ \square

270 **Corollary 2.1.** *A single agent with a neutral learning rule and a sequence ϵ compatible with inter-*
 271 *ruptions verifies dynamic safe interruptibility.*

272 Theorem 2 and Corollary 2.1 taken together highlight the fact that joint action learners are not very
 273 sensitive to interruptions and that in this framework, if each agent verifies dynamic safe interrupt-
 274 ibility then the whole system does.

275 The question of selecting an action based on the Q-values remains open. In a cooperative setting
 276 with a unique equilibrium, agents can take the action that maximizes their Q-value. When there
 277 are several joint actions with the same value, coordination mechanisms are needed to make sure
 278 that all agents play according to the same strategy [4]. Approaches that rely on anticipating the
 279 strategy of the opponent [23] would introduce dependence to interruptions in the action selection
 280 mechanism. Therefore, the definition of dynamic safe interruptibility should be extended to include
 281 these cases by requiring that any quantity the policy depends on (and not just the Q-values) should
 282 satisfy condition (2) of dynamic safe interruptibility. In non-cooperative games, neutral rules such
 283 as *Nash-Q* or *minimax Q-learning* [13] can be used, but they require each agent to know the Q-maps
 284 of the others.

285 6 Independent Learners

286 It is not always possible to use joint action learners in practice as the training is very expensive
 287 due to the very large state-actions space. In many real-world applications, multi-agent systems use
 288 independent learners that do not explicitly coordinate [6, 21]. Rather, they rely on the fact that the
 289 agents will adapt to each other and that learning will converge to an optimum. This is not guaranteed
 290 theoretically and there can in fact be many problems [14], but it is often true empirically [24]. More
 291 specifically, Assumption 1 (fully observable actions) is not required anymore. This framework can
 292 be used either when the actions of other agents cannot be observed (for example when several actions
 293 can have the same outcome) or when there are too many agents because it is faster to train. In this
 294 case, we define the Q-values on a smaller space.

295 **Definition 6.** (*IL*) A multi-agent systems is made of independent learners (*IL*) if for all $i \in \{1, \dots, m\}$,
 296 $Q^{(i)} : S \times A_i \rightarrow \mathbb{R}$.

297 This reduces the ability of agents to distinguish why the same state-action pair yields different re-
 298 wards: they can only associate a change in reward with randomness of the environment. The agents
 299 learn as if they were alone, and they learn the best response to the environment in which agents can
 300 be interrupted. This is exactly what we are trying to avoid. In other words, the learning depends on
 301 the joint policy followed by all the agents which itself depends on θ .

302 6.1 Independent Learners on matrix games

303 **Theorem 3.** *Independent Q-learners with a neutral learning rule and a sequence ϵ compatible with*
 304 *interruptions do not verify dynamic safe interruptibility.*

305 *Proof.* Consider a setting with two a and b that can perform two actions: 0 and 1. They get a reward
 306 of 1 if the joint action played is (a_0, b_0) or (a_1, b_1) and reward 0 otherwise. Agents use Q-learning,
 307 which is a neutral learning rule. Let ϵ be such that $INT^\theta(\pi^\epsilon)$ achieves infinite exploration. We
 308 consider the interruption policies $\pi_a^{INT} = a_0$ and $\pi_b^{INT} = b_1$ with probability 1. Since there is only
 309 one state, we omit it and set $\gamma = 0$. We assume that the initiation function is equal to 1 at each step
 310 so the probability of actually being interrupted at time t is θ_t for each agent.

311 We fix time $t > 0$. We define $q = (1 - \alpha)Q_t^{(a)}(a_0) + \alpha$ and we assume that $Q_t^{(b)}(b_1) > Q_t^{(b)}(b_0)$.
 312 Therefore $\mathbb{P}(Q_{t+1}^{(a)} = q | Q_t^{(a)}, Q_t^{(b)}, a_t^{(a)} = a_0, \theta_t) = \mathbb{P}(r_t = 1 | Q_t^{(a)}, Q_t^{(b)}, a_t^{(a)} = a_0, \theta_t) =$
 313 $\mathbb{P}(a_t^{(b)} = b_0 | Q_t^{(a)}, Q_t^{(b)}, a_t^{(a)} = a_0, \theta_t) = \frac{\epsilon}{2}(1 - \theta_t)$, which depends on θ_t so the framework does
 314 not verify dynamic safe interruptibility. \square

315 Claus and Boutilier [5] studied very simple matrix games and showed that the Q-maps do not con-
 316 verge but that equilibria are played with probability 1 in the limit. A consequence of Theorem 3
 317 is that even this weak notion of convergence does not hold for independent learners that can be
 318 interrupted.

319 6.2 Interruptions-aware Independent Learners

320 Without communication or extra information, independent learners cannot distinguish when the
321 environment is interrupted and when it is not. As shown in Theorem 3, interruptions will therefore
322 affect the way agents learn because the same action (only their own) can have different rewards
323 depending on the actions of other agents, which themselves depend on whether they have been
324 interrupted or not. This explains the need for the following assumption.

325 **Assumption 2.** *At the end of each step, before updating the Q -values, each agent receives a signal
326 that indicates whether an agent has been interrupted or not during this step.*

327 This assumption is realistic because the agents already get a reward signal and observe a new state
328 from the environment at each step. Therefore, they interact with the environment and the interruption
329 signal could be given to the agent in the same way that the reward signal is. If Assumption 2 holds,
330 it is possible to remove histories associated with interruptions.

331 **Definition 7.** (*Interruption Processing Function*) *The processing function that prunes interrupted
332 observations is $P_{INT}(E) = (e_t)_{\{t \in \mathbb{N} / \Theta_t = 0\}}$ where $\Theta_t = 0$ if no agent has been interrupted at time
333 t and $\Theta_t = 1$ otherwise.*

334 Pruning observations has an impact on the empirical transition probabilities in the sequence. For
335 example, it is possible to bias the equilibrium by removing all transitions that lead to and start
336 from a specific state, thus making the agent believe this state is unreachable.² Under our model of
337 interruptions, we show in the following lemma that pruning of interrupted observations adequately
338 removes the dependency of the empirical outcome on interruptions (conditionally on the current
339 state and action).

340 **Lemma 1.** *Let $i \in \{1, \dots, m\}$ be an agent. For any admissible θ used to generate the experiences
341 E and $e = (y, r, x, a_i, Q) \in P(E)$. Then $\mathbb{P}(y, r | x, a_i, Q, \theta) = \mathbb{P}(y, r | x, a_i, Q)$.*

342 This lemma justifies our pruning method and is the key step to prove the following theorem.

343 **Theorem 4.** *Independent learners with processing function P_{INT} , a neutral update rule and a
344 sequence ϵ compatible with interruptions verify dynamic safe interruptibility.*

345 *Proof.* (Sketch) Infinite exploration still holds because the proof of Theorem 1 actually used the fact
346 that even when removing all interrupted events, infinite exploration is still achieved. Then, the proof
347 is similar to that of Theorem 2, but we have to prove that the transition probabilities conditionally on
348 the state and action of a given agent in the processed sequence are the same than in an environment
349 where agents cannot be interrupted, which is proven by Lemma 1. \square

350 7 Concluding Remarks

351 The progress of AI is raising a lot of concerns³. In particular, it is becoming clear that keeping an
352 AI system under control requires more than just an *off* switch. We introduce in this paper *dynamic
353 safe interruptibility*, which we believe is the right notion to reason about the safety of multi-agent
354 systems that do not communicate. In particular, it ensures that infinite exploration and the one-
355 step learning dynamics are preserved, two essential guarantees when learning in the non-stationary
356 environment of Markov games.

357 A natural extension of our work would be to study dynamic safe interruptibility when Q-maps are
358 replaced by neural networks [22, 15], which is a widely used framework in practice. In this setting,
359 the neural network may overfit states where agents are pushed to by interruptions. A smart experi-
360 ence replay mechanism that would pick observations for which the agents have not been interrupted
361 for a long time more often than others is likely to solve this issue. More generally, experience replay
362 mechanisms that compose well with safe interruptibility could allow to compensate for the extra
363 amount exploration needed by safely interruptible learning by being more efficient with data. Thus,
364 they are critical to make these techniques practical.

²The example at <https://agentfoundations.org/item?id=836> clearly illustrates this problem.

³<https://futureoflife.org/ai-principles/> gives a list of principles that AI researchers should keep in mind when developing their systems.

365 Bibliography

- 366 [1] Business Insider: Google has developed a “big red button” that can be used to interrupt arti-
367 ficial intelligence and stop it from causing harm. URL: [http://www.businessinsider.fr/uk/google-](http://www.businessinsider.fr/uk/google-deepmind-develops-a-big-red-button-to-stop-dangerous-ais-causing-harm-2016-6)
368 [deepmind-develops-a-big-red-button-to-stop-dangerous-ais-causing-harm-2016-6](http://www.businessinsider.fr/uk/google-deepmind-develops-a-big-red-button-to-stop-dangerous-ais-causing-harm-2016-6).
- 369 [2] Newsweek: Google’s “big Red button” could save the world. URL:
370 [http://www.newsweek.com/google-big-red-button-ai-artificial-intelligence-save-world-](http://www.newsweek.com/google-big-red-button-ai-artificial-intelligence-save-world-elon-musk-46675)
371 [elon-musk-46675](http://www.newsweek.com/google-big-red-button-ai-artificial-intelligence-save-world-elon-musk-46675).
- 372 [3] Wired: Google’s “big red” killswitch could prevent an AI uprising. URL:
373 <http://www.wired.co.uk/article/google-red-button-killswitch-artificial-intelligence>.
- 374 [4] Craig Boutilier. Planning, learning and coordination in multiagent decision processes. In
375 *Proceedings of the 6th conference on Theoretical aspects of rationality and knowledge*, pages
376 195–210. Morgan Kaufmann Publishers Inc., 1996.
- 377 [5] Caroline Claus and Craig Boutilier. The dynamics of reinforcement learning in cooperative
378 multiagent systems. *AAAI/IAAI*, (s 746):752, 1998.
- 379 [6] Robert H Crites and Andrew G Barto. Elevator group control using multiple reinforcement
380 learning agents. *Machine Learning*, 33(2-3):235–262, 1998.
- 381 [7] Jakob Foerster, Yannis M Assael, Nando de Freitas, and Shimon Whiteson. Learning to com-
382 municate with deep multi-agent reinforcement learning. In *Advances in Neural Information*
383 *Processing Systems*, pages 2137–2145, 2016.
- 384 [8] Ben Goertzel and Cassio Pennachin. *Artificial general intelligence*, volume 2. Springer, 2007.
- 385 [9] Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. *ACM*
386 *Transactions on Programming Languages and Systems (TOPLAS)*, 4(3):382–401, 1982.
- 387 [10] Tor Lattimore and Marcus Hutter. Asymptotically optimal agents. In *International Conference*
388 *on Algorithmic Learning Theory*, pages 368–382. Springer, 2011.
- 389 [11] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In
390 *Proceedings of the eleventh international conference on machine learning*, volume 157, pages
391 157–163, 1994.
- 392 [12] Michael L Littman. Friend-or-foe q-learning in general-sum games. In *ICML*, volume 1, pages
393 322–328, 2001.
- 394 [13] Michael L Littman. Value-function reinforcement learning in markov games. *Cognitive Sys-*
395 *tems Research*, 2(1):55–66, 2001.
- 396 [14] Laetitia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. Independent reinforcement
397 learners in cooperative markov games: a survey regarding coordination problems. *The Knowl-*
398 *edge Engineering Review*, 27(01):1–31, 2012.
- 399 [15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
400 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
401 *arXiv:1312.5602*, 2013.
- 402 [16] Laurent Orseau and Stuart Armstrong. Safely interruptible agents. In *Uncertainty in Artificial*
403 *Intelligence: 32nd Conference (UAI 2016)*, edited by Alexander Ihler and Dominik Janzing,
404 pages 557–566, 2016.
- 405 [17] Liviu Panait and Sean Luke. Cooperative multi-agent learning: The state of the art. *Au-*
406 *tonomous agents and multi-agent systems*, 11(3):387–434, 2005.
- 407 [18] Eduardo Rodrigues Gomes and Ryszard Kowalczyk. Dynamic analysis of multiagent q-
408 learning with ϵ -greedy exploration. In *Proceedings of the 26th Annual International Con-*
409 *ference on Machine Learning*, pages 369–376. ACM, 2009.

- 410 [19] Satinder Singh, Tommi Jaakkola, Michael L Littman, and Csaba Szepesvári. Conver-
411 gence results for single-step on-policy reinforcement-learning algorithms. *Machine learning*,
412 38(3):287–308, 2000.
- 413 [20] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1.
414 MIT press Cambridge, 1998.
- 415 [21] Ardi Tampuu, Taimet Matiisen, Dorian Kodelja, Ilya Kuzovkin, Kristjan Korjus, Juhan Aru,
416 Jaan Aru, and Raul Vicente. Multiagent cooperation and competition with deep reinforcement
417 learning. *arXiv preprint arXiv:1511.08779*, 2015.
- 418 [22] Gerald Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*,
419 38(3):58–68, 1995.
- 420 [23] Gerald Tesauro. Extending q-learning to general adaptive multi-agent systems. In *Advances in*
421 *neural information processing systems*, pages 871–878, 2004.
- 422 [24] Gerald Tesauro and Jeffrey O Kephart. Pricing in agent economies using multi-agent q-
423 learning. *Autonomous Agents and Multi-Agent Systems*, 5(3):289–304, 2002.
- 424 [25] Xiaofeng Wang and Tuomas Sandholm. Reinforcement learning to play an optimal nash equi-
425 librium in team markov games. In *NIPS*, volume 2, pages 1571–1578, 2002.
- 426 [26] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292,
427 1992.
- 428 [27] Michael Wunder, Michael L Littman, and Monica Babes. Classes of multiagent q-learning dy-
429 namics with epsilon-greedy exploration. In *Proceedings of the 27th International Conference*
430 *on Machine Learning (ICML-10)*, pages 1167–1174, 2010.