# Joint Intrinsic Motivation for Coordinated Exploration in Multi-Agent Deep Reinforcement Learning

### Maxime Toquebiau
ECE Paris & Sorbonne Université, CNRS, Institut des
Systèmes Intelligents et de Robotique, ISIR, F-75005
Paris, France
maxime.toquebiau@gmail.com

### Nicolas Bredeche
Sorbonne Université, CNRS, Institut des Systèmes
Intelligents et de Robotique, ISIR, F-75005
Paris, France
nicolas.bredeche@sorbonne-universite.fr

### Faïz Benamar
Sorbonne Université, CNRS, Institut des Systèmes
Intelligents et de Robotique, ISIR, F-75005
Paris, France
faiz.ben_amar@sorbonne-universite.fr

### Jae-Yun Jun
ECE Paris
Paris, France
jaeyunjk@gmail.com

## ABSTRACT

Multi-agent deep reinforcement learning (MADRL) problems often encounter the challenge of sparse rewards. This challenge becomes even more pronounced when coordination among agents is necessary. In this paper, we propose a new approach for rewarding strategies where agents collectively exhibit novel behaviors. To achieve this, we introduce a measure of novelty that specifically considers diverse coordination patterns exhibited by a team of agents. We present JIM (Joint Intrinsic Motivation), a multi-agent intrinsic motivation method that follows the centralized learning with decentralized execution paradigm. By testing JIM with the state-of-the-art MADRL method QMIX, we demonstrate how joint exploration is crucial for solving tasks where the optimal strategy requires a high level of coordination.

## KEYWORDS

Multi-agent Systems, Deep Reinforcement Learning, Intrinsic Motivation

## 1 INTRODUCTION

One crucial aspect of human intelligence is its ability to act coincidentally with other human beings, to either cooperate or compete in a given task. This has led researchers to study reinforcement learning (RL) in the context of multi-agent systems (MAS), where multiple artificial agents interact with their environment and each other while concurrently learning to perform a task [10, 25]. However, having multiple agents in the environment makes the RL process significantly more difficult for several reasons [30]. In particular, the global reward depends on the actions of several independent agents, which makes the search for the optimal joint policy more complicated.

Recently, multi-agent deep reinforcement learning (MADRL) approaches have combined advancements in RL and deep learning to tackle long-standing problems in MAS such as credit assignment or partial observability [5, 10, 20]. These techniques are able to solve very complex multi-agent tasks such as autonomous driving [22] or real-time strategy video games [15]. However, major issues still remain with these approaches, such as the problem of relative

overgeneralization [27, 29] where agents struggle to find the optimal joint policy because local policies are attracted towards suboptimal areas of the search space. This makes most algorithms inefficient in tasks where the optimal strategy requires strong coordination among agents. Relative overgeneralization can be described as a problem of exploration of the joint-state space: as the success of the MAS depends on the coordination of multiple agents, exploring the joint-observation space is required to discover optimal joint behaviors. In this paper, we address the question of how to explore the joint-state space to efficiently discover superior coordinated strategies for solving the task at hand.

In single-agent RL, the problem of exploration has been studied to solve hard exploration tasks where positive reward signals are very sparse. One solution is to use intrinsic motivation [9, 16, 21] to incite agents to explore unknown parts of the environment. In addition to the environment reward, agents are given an auxiliary reward related to the novelty of encountered states. Maximizing this intrinsic reward leads agents to visit previously unexplored regions of the environment, ultimately discovering new solutions to the task. These methods have shown great success in helping RL agents solve hard exploration tasks [1, 17].

In the multi-agent setting, intrinsic objectives have also been studied to induce different kinds of behaviors in agents such as coordinated exploration [7], social influence [8, 26] or alignment with other agents' expectations [12]. However, previous works have only used local observations to generate intrinsic rewards. In the context of exploration, an intrinsic reward based only on local observations will lead to each agent exploring their own observation space without taking care of the current state of other agents. This can result in inefficient exploration in cooperative tasks where the success of the MAS depends on the coordination of all agents.

In this paper, we introduce a novel multi-agent exploration approach called Joint Intrinsic Motivation (JIM) which can be combined with any MADRL algorithm that follows the centralized training with decentralized execution paradigm (CTDE). JIM exploits centralized information to motivate agents to explore new coordinated behaviors. In order to compute joint novelty, JIM combines two previous state-of-the-art approaches: NovelD [32] for exploring unknown parts of the environment, and E3B [6] for having more diverse trajectories. Experimental studies show that combining JIM

Maxime Toquebiau, Nicolas Bredeche, Faïz Benamar, and Jae-Yun Jun

with the state-of-the-art algorithm QMIX [20] helps to overcome the problem of relative overgeneralization.

## 2 RELATED WORKS

In recent years, deep reinforcement learning techniques have been used in the context of MAS to tackle long-standing issues in multi-agent learning. Successful single-agent RL approaches have been adapted to the CTDE framework [10, 31], using a centralized value function to guide the training of decentralized policies. Recent studies have investigated the problem of credit assignment [5] in MADRL, i.e, distributing the global reward among agents based on their participation. Value factorization methods also do this implicitly [24], combining the output of local value functions into a centralized one that predicts the current value of the system. In particular, QMIX [20] uses a separate network to predict the Q-value of the joint action, given the output of local Q-values and the global state of the environment. QMIX has established itself as a long-standing state-of-the-art approach, despite its inherent limitations that several works have tried to surpass [19, 23]. However, MADRL algorithms have been shown to suffer from the problem of relative overgeneralization [28, 29]. So far, few works have addressed this problem: Wei et al. [28] propose maximum entropy RL to explore the joint-action space, and MAVEN [13] augments QMIX using a hierarchical policy to guide the exploration of joint behaviors.

A promising approach to overcome relative overgeneralization is to intrinsically motivate agents to explore their environment, ultimately discovering the optimal reward signals. In single-agent RL, curiosity has been defined to help agents solve hard exploration tasks [9, 16, 21] by rewarding the visitation of states considered as novel. For measuring novelty, several methods have used the error of trainable prediction models. The Intrinsic Curiosity Module (ICM) [17] trains a model of environment dynamics and uses the prediction error as a measure of novelty. Random Network Distillation (RND) [2] uses a target network that produces a random encoding of the state and trains a predictor network to generate the same encoding, the prediction error being the measure of novelty. The idea behind these two approaches is that the prediction models will yield low novelty for states similar to what they have trained on while producing high novelty for unknown parts of the environment. RIDE [18] and NovelD [32] use respectively ICM and RND to compute a reward from the difference of novelty between the next state and the current state, pushing the agents to always seek novel states. Similarly, NGU [1] and E3B [6] use clustering techniques to reward states that are distant from previous states. Finally, a similar approach is proposed by AGAC [4] which trains an adversarial policy to predict the main policy's output, the latter being rewarded with the former's prediction error.

In MADRL, recent works have demonstrated the effectiveness of intrinsic rewards in promoting desirable behaviors in groups of agents. One example is social influence [8, 26] that rewards agents for performing actions that have a significant impact on other agents. Ma et al. [12] propose an intrinsic reward based on the average alignment with other agents' expectations, promoting more predictable behaviors in agents. Lupu et al. [11] propose to reward policies that perform diverse trajectories in comparison to a population of agents, which is shown to help train agents to be more versatile. Du et al. [3]

use intrinsic objectives as a credit assignment technique. Finally, Iqbal and Sha [7] propose an approach for coordinated exploration using several metrics for estimating the novelty of observations that depend on all agents' past experiences. However, their model is computationally expensive and does not address the exploration of the joint-observation space, which can be problematic for hard exploration tasks where relative overgeneralization can occur.

In this paper, we address the challenge of relative overgeneralization by rewarding agents for exploring the joint-observation space. In the following sections, we will present the necessary formal background and an overview of the proposed algorithm that implements joint intrinsic motivation.

## 3 BACKGROUND

### 3.1 Dec-POMDP

To describe cooperative multi-agent tasks, we use the setting of decentralized POMDP (Dec-POMDP) [14], defined as a tuple $\langle S, A, T, O, O, R, n, \gamma \rangle$ with $n$ the number of agents. $S$ describes the set of global states $s$ of the environment. $O$ is the set of joint observations, with one joint observation $\mathbf{o} = \{o_1, ..., o_n\} \in O$, and $A$ the set of joint actions, with one joint action $\mathbf{a} = \{a_1, ..., a_n\} \in A$. $T$ is the transition function defining the probability $P(s'|s, \mathbf{a})$ to transition from state $s$ to next state $s'$ with the joint action $\mathbf{a}$. $O$ is the observation function defining the probability $P(\mathbf{o}|\mathbf{a}, s')$ to observe the joint observation $\mathbf{o}$ after taking joint action $\mathbf{a}$ and ending up in $s'$. $R : O \times A \to \mathbb{R}$ is the reward function producing at each time step the reward shared by all agents. Finally, $\gamma \in [0, 1)$ is the discount factor controlling the importance of immediate rewards against future gains.

### 3.2 Intrinsic rewards

In Section 2, we introduced intrinsic motivation as a way to incite agents to actively explore their environment. To this end, at each time step $t$, agents receive an augmented reward $r_t = r_t^e + \beta r_t^{int}$, where $r_t^e$ is the extrinsic reward given by the environment, $r_t^{int}$ is the intrinsic reward and $\beta$ is a hyperparameter controlling the weight of the intrinsic reward in the agents' objective.

In this section, we describe three methods of intrinsic rewards from the literature that we will use later in Section 4.2.

*Random Network Distillation.* In Random Network Distillation (RND), Burda et al. [2] compute novelty using two neural networks with the same architecture: a target network $\phi$ and a predictor network $\phi'$. The target's parameters are initialized randomly and fixed. It takes as input the state $s_t$ and produces a random embedding $\phi(s_t)$. The predictor is trained to output the same embedding, minimizing the Euclidean distance:

$$RND_t(s_t) = \|\phi(s_t) - \phi'(s_t)\|_2 \qquad (1)$$

This distance is used as a measure of the novelty of state $s_t$ and is given as an intrinsic reward to agents.

*NovelD.* Zhang et al. [32] build upon RND to devise a novelty criterion termed NovelD. It is defined as follows:

$$N(s_t, s_{t+1}) = \max[RND(s_{t+1}) - \alpha RND(s_t), 0] \times$$
$$\times \mathbb{I}\{N_e(s_{t+1}) = 1\} \quad (2)$$

with $\alpha$ a scaling factor and $N_e$ an episodic count of visited states. The first part is the core of the novelty criterion. It uses RND to reward agents for positive gains in novelty between the current and the next states. The second part is an episodic restriction that ensures the reward is given only when state $s_{t+1}$ is observed for the first time in this episode. This restriction limits the use of NovelD to discrete state spaces as it relies on an explicit count of visited states.

*E3B.* With E3B, Henaff et al. [6] propose an episodic bonus based on the position of the observed state with respect to an ellipse that fits all states previously encountered in the current episode. Formally, it is computed as follows:

$$b(s_t) = \psi(s_t)^\top C_{t-1}^{-1} \psi(s_t), \tag{3}$$

with

$$C_{t-1} = \sum_{i=1}^{t-1} \psi(s_i)\psi(s_i)^\top + \lambda I, \tag{4}$$

where $I$ is the identity matrix and $\lambda$ a scalar coefficient. $\psi$ is an embedding network trained using an inverse dynamics model [17]: embeddings of following states $\psi(s_t)$ and $\psi(s_{t+1})$ are used by a separate neural network trained to predict the action $a_t$ taken between these states. As a result of this training process, $\psi$ encodes parts of the observation that are controllable by the agents (please refer to [6] for details). Intuitively, $b$ can be understood as a generalization of a count-based episodic bonus for a continuous state space. States that are close to previously encountered states in the current episode will yield low bonuses, whereas states that are very different will produce high bonuses.

## 4 ALGORITHM

In this section, we introduce the Joint Intrinsic Motivation (JIM) exploration criterion for coordinated multi-agent exploration. Firstly, we describe the motivation behind our approach by providing a detailed description of the problem of relative overgeneralization. Then, we define the intrinsic reward used for motivating agents to explore a continuous state-space environment in a coordinated fashion. Finally, we explain how this reward is used in a multi-agent setting with JIM.

### 4.1 The challenge of coordinated actions

Addressing hard exploration environments is challenging because of the few positive reward signals that exist to guide the agent's learning process. This becomes even worse with MAS as the completion of a task depends on the actions of multiple independent agents. When strong coordination is needed, agents will struggle to find the optimal strategy and settle for an easier suboptimal joint strategy, which is a problem known as relative overgeneralization [27, 29]. Figure 1a provides an example of a social dilemma game where relative overgeneralization occurs. The optimal strategy requires both agents to choose action A. But if only one agent chooses action A, the payoff is very bad. Therefore, agents will independently prefer to take actions B or C, as action A most often leads to sub-optimal outcomes.

In MAS, this can be seen as a problem of ill-coordinated exploration. As success depends on coordinated behaviors, exploration of joint policies is required in order to discover which ones lead to

optimal returns. In the example of Figure 1a, exploring independent strategies will lead to ultimately choosing suboptimal actions as they individually may yield better expected returns. On the other hand, we argue that uniformly exploring joint actions would enable agents to choose optimal joint strategies more often and consequently learn more efficient individual behaviors. The approach described in the following two sections implements an algorithm that efficiently rewards agents for exploring the joint-observation space, in order to consistently find optimal strategies.

### 4.2 Double-timescale Intrinsic Reward

We define a novelty metric that combines two exploration criteria working at different timescales:

- A **life-long exploration criterion** that captures how novel is the current observation with respect to all observations since the beginning of training.
- An **episodic exploration criterion** that captures the difference between the current observation and all previous observations in the current episode.

Intuitively, the life-long reward motivates agents to search for never-experienced parts of the environment. Meanwhile, the episodic bonus induces more diverse trajectories. These two elements will feed each other and reinforce agents to efficiently explore their environment.

Concretely, for each transition from state $s_t$ to the next state $s_{t+1}$, we define the double-timescale intrinsic reward as follows:

$$r_t(s_t, a_t, s_{t+1}) = N_l(s_t, s_{t+1}) \times \sqrt{2b(s_{t+1})} \tag{5}$$

The first term $N_l(s_t, s_{t+1})$ corresponds to the life-long novelty and the second term $\sqrt{2b(s_{t+1})}$ corresponds to the episodic novelty.

The life-long novelty $N_l$ is inspired from NovelD [32] (see Section 3.2):

$$N_l(s_t, s_{t+1}) = \max[RND(s_{t+1}) - \alpha RND(s_t), 0], \tag{6}$$

with $\alpha$ a scaling factor and *RND* the novelty measure. We remove the episodic restriction of the original approach as it relies on an episodic count of visited states. This makes it impractical in a continuous state space, as one state is very unlikely to be visited twice. Instead, we scale the life-long novelty $N_l$ with the elliptical episodic bonus $b$ from E3B [6] (see also Section 3.2). This bonus replaces the episodic restriction by scaling $N_l$ up or down, depending on the novelty of the current state compared to what has been observed in the current episode. As $b$ provides very large bonuses and decreases very fast, we use $\sqrt{2b(s_{t+1})}$ to both smooth out large values and increase small ones.

Combining these two rewards makes it possible to take the benefits of both. NovelD pushes agents to explore regions of the state space that are not well-known to agents. Meanwhile, the elliptical episodic bonus favors diverse trajectories, inducing agents to always seek new observations during a single episode. As the agents explore their environment, the prediction error of RND (see equation (1)) slowly decreases. Thus, the life-long novelty decreases as well, tending toward zero, allowing agents to progressively focus on the extrinsic reward. Finally, the proposed intrinsic reward does not rely on any explicit count of visited states. As a consequence, it can be used in continuous state spaces.

|   | $A$ | $B$ | $C$ |
|---|---|---|---|
| $A$ | 10 | −5 | −5 |
| $B$ | −5 | 7 | 7 |
| $C$ | −5 | 7 | 7 |

**(a)**                    **(b)**

**Figure 1: Two examples of relative overgeneralization: (a) payoff matrix of a social dilemma game, (b) heat-map of the reward function of the `rel_overgen` environment for $D = 40$ and $\delta = 30$.**

## 4.3 The Joint Intrinsic Motivation algorithm

Building from the intrinsic reward introduced previously, we propose the Joint Intrinsic Motivation (JIM) algorithm to incite MADRL agents to explore the joint-observation space. At each time step, all agents receive the same global reward $r_t = r_t^e + \beta r_t^{JIM}$, where $r_t^e$ is the extrinsic reward given by the environment, $r_t^{JIM}$ is our joint exploration criterion and $\beta$ is a hyper-parameter controlling the weight of the intrinsic reward. The exploration criterion in JIM uses the double-timescale intrinsic reward defined earlier to compute the novelty of the joint observation:

$$r_t^{JIM}(\mathbf{o}_t, \mathbf{a}_t, \mathbf{o}_{t+1}) = N_l(\mathbf{o}_t, \mathbf{o}_{t+1}) \times \sqrt{2b(\mathbf{o}_{t+1})}, \tag{7}$$

where $\mathbf{o}_t = \{o_t^i\}_{0 \le i \le N}$, i.e., the concatenation of all local observations. Rather than only exploring their local-observation space, agents will be rewarded for finding new combinations of observations with other agents of the system.

As JIM uses joint observations for computing the intrinsic reward, it can be associated with any MADRL algorithm that fits in the CTDE paradigm. These algorithms usually employ a centralized value function [10, 20, 31] that looks at the joint observation to predict the value of the agents' actions. Such centralized value functions will be able to associate rewards provided by JIM to new configurations in the joint observation space, thus inducing the agents to actively search for these configurations.

One could note that the joint observation has two notable drawbacks: the number of dimensions grows exponentially with the number of agents and there is a risk of capturing redundant information. These issues are both alleviated with JIM as we use dimensionality reduction techniques. $N_l$ and $b$ use respectively $\phi$ and $\psi$ (see Section 3.2) as embedding networks to encode the joint observation into a more condensed representation that contains only the required information.

## 5 IMPLEMENTATION DETAILS

As previously said, JIM can be used to augment any MADRL approach that fits in the CTDE paradigm. In the experiments presented in the next section, we use JIM with QMIX [20]. We use the default QMIX architecture and hyperparameters, as presented in the original paper. Both embedding networks $\psi$ and $\phi$ (see Section 3.2) are feed-forward neural networks with respectively one and three hidden layers of dimension 128. They output encodings of dimension 64. The hyperparameter $\alpha$ (see equation (6)) is set to 0.2 and $\lambda$ (see equation (4)) to 0.1. Finally, the weight $\beta$ of the intrinsic reward in the total reward is always set to 1. All our code is available online[1].

## 6 EXPERIMENTS

In this Section, we evaluate the ability of JIM to address the problem of relative overgeneralization. In order to do so, we implement the JIM exploration criterion within the state-of-the-art QMIX algorithm for MADRL and benchmark the resulting algorithm on different versions of a toy problem where the problem of relative overgeneralization can be artificially tuned.

## 6.1 Environment definition

We design a simple test environment that expands the example of relative overgeneralization shown in Figure 1a. In this environment called `rel_overgen`, two agents can move on a discrete one-dimensional axis with $D$ possible positions. Each agent is denoted by its position, namely x and y. At each time step, agents observe their position as a one-hot vector $o_t^x = \{o_t^{x,j} = 1$ if x $= j$, 0 otherwise$\}_{0 \le j < D}$ and can choose between three actions: move in one direction or the other or stay in position. They receive a reward corresponding to their combined position:

$$r_t^e(x, y; \delta) = \max\Big(R^+ - \frac{\delta}{D}\big[(x - r_x^+)^2 + (y - r_y^+)^2\big],$$
$$R^- - \frac{1}{8D}\big[(x - r_x^-)^2 + (y - r_y^-)^2\big]\Big) \tag{8}$$

The result of this formula is displayed in Figure 1b. The reward combines two hyperboles in opposite corners: one narrow that culminates at $R^+$ at position $(r_x^+, r_y^+)$, and another much wider that plateaus at $R^-$ at position $(r_x^-, r_y^-)$. We set the optimal reward $R^+$ to 12 and the suboptimal $R^-$ to 0. The width of the optimal reward spike is controlled by the parameter $\delta$: a higher $\delta$ value yields a narrower spike.

The goal of the agents is to find where to go to maximize their aggregated rewards. The wide suboptimal hyperbole will probably attract agents to minimize their loss. The optimal reward spike is difficult to find because it covers a small portion of the state space, but it guarantees much greater returns. We can vary the difficulty of the task by changing the width of this optimal reward spike: the narrower the spike, the harder it is to find.

In this environment, we expect MADRL to struggle to find the optimal reward spike. Exploring local states could help but would not be sufficient to consistently solve the task. As the dimension $D$ of the local-state space is fairly small, novelty rewards will quickly vanish and will not help agents to find the optimal reward spike. Exploring the joint-observation space adequately is required in order to consistently find optimal rewards. As JIM will reward exploration until all combined positions $(x, y)$ are visited several times, agents will visit the optimal reward spike more often, thus helping them to learn the optimal coordinated strategy.

---

[1]https://github.com/MToquebiau/Joint-Intrinsic-Motivation

**Figure 2: Performance of variants of QMIX in the `rel_overgen` environment, with three levels of difficulty. On top, we show the heat-maps representing the reward function in each version of the environment, where the difficulty is dictated by the width coefficient of the optimal reward spike $\delta$ (as defined in equation (8)). Increasing $\delta$ leads to a slightly narrower optimal reward spike. Below is shown the performance during training of QMIX with no intrinsic reward (QMIX), local intrinsic motivation (QMIX+LIM) and joint intrinsic motivation (QMIX+JIM) (mean and standard deviation shown for 15 runs each). We see that a slight decrease in the size of the optimal reward spike results in a considerable increase in the difficulty of the task.**

## 6.2 Results

The results shown in Figure 2 confirm this hypothesis. We show the performance of QMIX [20] in rel_overgen with no intrinsic reward (QMIX) and with two intrinsic rewards: our approach JIM (QMIX+JIM) and a local version of our intrinsic motivation where each agent generates its own intrinsic reward based on its local observations (QMIX+LIM[2]). Further, we show the performance in three levels of difficulty dictated by the width of the optimal reward spike. We plot the mean and standard deviation for 15 independent runs each.

The results clearly demonstrate the importance of exploring the joint-state space. QMIX alone manages to get some good performance on the easy scenario, but the large standard deviation shows its inconsistency. In the harder scenarios, QMIX's performance degrades strongly, never finding any positive reward in the hardest case. JIM clearly improves the performance. In the easy scenario, QMIX+JIM consistently goes for the optimal reward spike. In the harder settings, it still performs well, even in the "very hard" scenario where the optimal reward spike covers only 0.013% of all combined positions. The results of QMIX+LIM show that exploring the local-observation space helps agents find the optimal reward spike more often. However, it performs worse than JIM as it does not insure that all combined positions are sufficiently explored. This shows that exploring the joint-observation space is crucial to allow agents to discover optimal coordinated behaviors.

## 7 CONCLUSION

In this paper, we present the JIM algorithm, which employs an exploration criterion to reward teams of cooperating agents for exploring the space of joint observations. JIM can be integrated to enhance any Multi-Agent Deep Reinforcement Learning algorithm and can be applied to problems with continuous state-action spaces. By combining JIM with the state-of-the-art QMIX algorithm, we demonstrate its efficiency. Our results show that QMIX with JIM outperforms both the original QMIX algorithm and QMIX with single-agent intrinsic rewards. Notably, JIM enables the discovery of optimal coordinated behaviors that would be hard to find otherwise as they necessitate a high level of coordination between agents. Our short-term objective for this work is to validate the results in a more realistic scenario involving multiple agents addressing the same cooperative task. Preliminary results indicate that JIM allows agents to efficiently explore more complex continuous environments to discover optimal coordinated strategies.

## ACKNOWLEDGMENT

## REFERENCES

[1] Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven Kapturowski, Olivier Tieleman, Martin Arjovsky, Alexander Pritzel, Andrew Bolt, and Charles Blundell. 2020. Never Give Up: Learning Directed Exploration Strategies. In *International Conference on Learning Representations*. https://openreview.net/forum?id=Sye57xStvB

---

[2]LIM: Local Intrinsic Motivation

[2] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. 2019. Exploration by Random Network Distillation. In *International Conference on Learning Representations*.

[3] Yali Du, Lei Han, Meng Fang, Tianhong Dai, Ji Liu, and Dacheng Tao. 2019. LIIR: Learning Individual Intrinsic Reward in Multi-Agent Reinforcement Learning. In *Advances in Neural Information Processing Systems*. https://proceedings.neurips.cc/paper/2019/hash/07a9d3fed4c5ea6b17e80258dee231fa-Abstract.html

[4] Yannis Flet-Berliac, Johan Ferret, Olivier Pietquin, Philippe Preux, and Matthieu Geist. 2021. Adversarially Guided Actor-Critic. In *International Conference on Learning Representations*. https://hal.inria.fr/hal-03167169

[5] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. 2018. Counterfactual Multi-Agent Policy Gradients. *Proceedings of the AAAI Conference on Artificial Intelligence*. https://ojs.aaai.org/index.php/AAAI/article/view/11794

[6] Mikael Henaff, Roberta Raileanu, Minqi Jiang, and Tim Rocktäschel. 2022. Exploration via Elliptical Episodic Bonuses. In *Advances in Neural Information Processing Systems*. https://openreview.net/forum?id=Xg-yZos9qJQ

[7] Shariq Iqbal and Fei Sha. 2019. Coordinated Exploration via Intrinsic Rewards for Multi-Agent Reinforcement Learning. https://openreview.net/forum?id=rkltE0VKwH

[8] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro A. Ortega, DJ Strouse, Joel Z. Leibo, and Nando de Freitas. 2019. Social Influence as Intrinsic Motivation for Multi-Agent Deep Reinforcement Learning. In *Proceedings of the 36th International Conference on Machine Learning*.

[9] Joel Lehman and Kenneth O Stanley. 2011. Abandoning objectives: Evolution through the search for novelty alone. *Evolutionary computation* 19, 2, 189–223.

[10] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2017. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*. https://dl.acm.org/doi/abs/10.5555/3295222.3295385

[11] Andrei Lupu, Brandon Cui, Hengyuan Hu, and Jakob Foerster. 2021. Trajectory Diversity for Zero-Shot Coordination. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 7204–7213. https://proceedings.mlr.press/v139/lupu21a.html

[12] Zixian Ma, Rose E Wang, Li Fei-Fei, Michael S. Bernstein, and Ranjay Krishna. 2022. ELIGN: Expectation Alignment as a Multi-Agent Intrinsic Reward. In *Advances in Neural Information Processing Systems*. https://openreview.net/forum?id=uPyNR2yPoe

[13] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. 2019. MAVEN: Multi-Agent Variational Exploration. In *Advances in Neural Information Processing Systems*, Vol. 32. https://proceedings.neurips.cc/paper/2019/file/f816dc0acface7498e10496222e9db10-Paper.pdf

[14] Frans A Oliehoek and Christopher Amato. 2016. *A concise introduction to decentralized POMDPs*. Springer. https://www.ccis.northeastern.edu/home/camato/publications/OliehoekAmato16book.pdf

[15] OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d. O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. 2019. Dota 2 with Large Scale Deep Reinforcement Learning. arXiv:1912.06680

[16] Pierre-Yves Oudeyer and Frederic Kaplan. 2007. What is Intrinsic Motivation? A Typology of Computational Approaches. *Frontiers in neurorobotics*. https://doi.org/10.3389/neuro.12.006.2007

[17] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. 2017. Curiosity-driven Exploration by Self-supervised Prediction. *Proceedings of the 34th International Conference on Machine Learning*.

[18] Roberta Raileanu and Tim Rocktäschel. 2020. RIDE: Rewarding Impact-Driven Exploration for Procedurally-Generated Environments. In *International Conference on Learning Representations*. https://openreview.net/forum?id=rkg-TJBFPB

[19] Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. 2020. Weighted QMIX: Expanding monotonic value function factorisation for deep multiagent reinforcement learning. In *Advances in Neural Information Processing Systems*. https://proceedings.neurips.cc/paper/2020/hash/73a427badebe0e32caa2e1fc7530b7f3-Abstract.html

[20] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *Proceedings of the 35th International Conference on Machine Learning*. http://proceedings.mlr.press/v80/rashid18a.html

[21] Jürgen Schmidhuber. 1991. A possibility for implementing curiosity and boredom in model-building neural controllers. In *Proceedings of the international conference on simulation of adaptive behavior: From animals to animats*.

[22] Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. 2016. Safe, Multi-Agent, Reinforcement Learning for Autonomous Driving. arXiv:1610.03295

[23] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. 2019. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*.

[24] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. 2018. Value-Decomposition Networks For Cooperative Multi-Agent Learning Based On Team Reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*.

[25] Ming Tan. 1993. Multi-Agent Reinforcement Learning: Independent versus Cooperative Agents. In *Proceedings of the Tenth International Conference on Machine Learning*.

[26] Tonghan Wang, Jianhao Wang, Yi Wu, and Chongjie Zhang. 2020. Influence-Based Multi-Agent Exploration. In *International Conference on Learning Representations*. https://openreview.net/forum?id=BJgy96EYvr

[27] Ermo Wei and Sean Luke. 2016. Lenient learning in independent-learner stochastic cooperative games. *The Journal of Machine Learning Research* 17, 2914–2955.

[28] Ermo Wei, Drew Wicke, David Freelan, and Sean Luke. 2018. Multiagent Soft Q-Learning. In *2018 AAAI Spring Symposium Series*.

[29] Rudolf Paul Wiegand. 2003. *An Analysis of Cooperative Coevolutionary Algorithms*. Ph.D. Dissertation. George Mason University. http://l.academicdirect.org/Horticulture/GAs/Refs/PhD_Wiegand&Jong_2003.pdf

[30] Michael Wooldridge. 2009. *An introduction to multiagent systems*. John wiley & sons.

[31] Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. 2021. The Surprising Effectiveness of PPO in Cooperative, Multi-Agent Games. arXiv:2103.01955

[32] Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E Gonzalez, and Yuandong Tian. 2021. NovelD: A Simple yet Effective Exploration Criterion. In *Advances in Neural Information Processing Systems*. https://proceedings.neurips.cc/paper/2021/file/d428d070622e0f4363fceae11f4a3576-Paper.pdf