

Agent-Temporal Credit Assignment for Optimal Policy Preservation in Sparse Multi-Agent Reinforcement Learning

Aditya Kapoor

Tata Consultancy Services, Mumbai

Sushant Swamy

Birla Institute of Technology and Science, Goa

Kale-ab Tessera

University of Edinburgh

Mayank Baranwal

Tata Consultancy Services, Mumbai

Mingfei Sun

University of Manchester

Harshad Khadilkar

Indian Institute of Technology, Bombay

Stefano V. Albrecht

University of Edinburgh

Abstract

In multi-agent environments, agents often struggle to learn optimal policies due to sparse or delayed global rewards, particularly in long-horizon tasks where it is challenging to evaluate actions at intermediate time steps. We introduce Temporal-Agent Reward Redistribution (TAR²), a novel approach designed to address the agent-temporal credit assignment problem by redistributing sparse rewards both temporally and across agents. TAR² decomposes sparse global rewards into time-step-specific rewards and calculates agent-specific contributions to these rewards. We theoretically prove that TAR² is equivalent to potential-based reward shaping, ensuring that the optimal policy remains unchanged. Empirical results demonstrate that TAR² stabilizes and accelerates the learning process. Additionally, we show that when TAR² is integrated with single-agent reinforcement learning algorithms, it performs as well as or better than traditional multi-agent reinforcement learning methods.

1 Introduction

In cooperative multi-agent reinforcement learning (MARL), multiple autonomous agents learn to interact and collaborate to execute tasks in a shared environment by maximizing a global return (Busoniu et al., 2008). MARL has shown considerable potential in solving decentralized partially observable Markov decision processes (Dec-POMDPs) (Oliehoek & Amato, 2016; Amato, 2024; Zhang et al., 2021), where each agent has access to only local information (partial observation) and need to select actions based on their local action-observation (or sometimes only observation) histories to maximize the global (team) return. Applications of MARL include complex video games such as StarCraft-II (Vinyals et al., 2019), Defense of the Ancients (DOTA) (Berner et al., 2019), Google Football (Kurach et al., 2020), and Capture the Flag (CTF) (Jaderberg et al., 2019), and real world applications of warehouse logistics (Krnjaic et al., 2022; Agrawal et al., 2023), e-commerce (Shelke et al., 2023; Baer et al., 2019), robotics (Sartoretti et al., 2019; Damani et al., 2021), and routing problems (Zhang et al., 2018; Vinitzky et al., 2020; Zhang et al., 2023b). These applications illustrate the potential of MARL to develop sophisticated strategies and behaviors through coordinated teamwork and collaboration.

Despite these successes, cooperative multi-agent systems face a significant challenge in credit assignment, which is crucial for learning effective policies Foerster et al. (2018). Credit assignment

in multi-agent systems encompasses two main aspects: *temporal credit assignment* and *agent credit assignment*. Temporal credit assignment involves decomposing sparse, delayed rewards into intermediate time steps within a multi-agent trajectory. The assignment of agent credit focuses on discerning the contribution of each agent to these decomposed temporal rewards (Albrecht et al., 2024). Addressing both aspects is essential for effective learning in cooperative multi-agent systems.

Significant progress has been made in addressing the credit assignment problem with methods such as VDN(Sunehag et al., 2017), QMIX(Rashid et al., 2020), QTRAN(Son et al., 2019), COMA(Foerster et al., 2018), and PRD (Freed et al., 2021). However, these methods are primarily designed to deal with agent credit assignment and may not be suitable for environments with sparse or delayed rewards (Papoudakis et al., 2020; De Witt et al., 2020). Additionally, the representations required for effective credit assignment might not align with those needed for learning Q-values or critics. Recent advances in temporal credit assignment have introduced dense reward functions in single-agent settings (Arjona-Medina et al., 2019; Ren et al., 2021; Liu et al., 2019; Gangwani et al., 2020) and multi-agent settings (Xiao et al., 2022; She et al., 2022). These methods aim to learn a Markovian proxy reward function that replaces the environment’s sparse rewards with dense rewards. Motivated by this progress, we aim to address the combined challenge of agent and temporal credit assignment in multi-agent tasks with sparse or delayed rewards.

In this paper, we propose Temporal-Agent Reward Redistribution (TAR²), a novel approach to address the problem of agent-temporal credit assignment by learning a reward redistribution function that decomposes sparse environment rewards to each time step of the multi-agent trajectory and further redistributes the temporally decomposed rewards to each agent based on their contributions. We theoretically prove that there exists a class of reward redistribution functions that can be formulated as potential-based reward shaping (Ng, 1999; Devlin & Kudenko, 2011), under which the optimal policies are preserved in the original reward function of the environment. TAR² extends AREL’s (Xiao et al., 2022) reward model that uses a temporal attention module to analyze the influence of state-action tuples along trajectories, followed by an agent attention module to identify the relevance of other agents for each agent. This alternation between the two attention modules allows the reward function to identify agent-specific state-action tuples that are key to the sparse environment rewards received by the multi-agent system. Thus, TAR² learns agent-specific temporal rewards and enables the use of single agent reinforcement learning (RL) algorithms such as IQL (Tan, 1997), IAC (Foerster et al., 2018), and IPPO (Schulman et al., 2017; De Witt et al., 2020) to solve multi-agent tasks. This approach separates the problem of credit assignment from learning Q-functions and critics, leveraging the simplicity and scalability of single-agent RL algorithms in complex environments. We empirically demonstrate the sample efficiency of our approach against competitive state-of-the-art baselines on SMACLite (Michalski et al., 2023).

In summary, our contribution is three folds:-

- We introduce TAR², a reward redistribution model that can be used to assign temporal and agent credit assignment to densify the reward signal.
- We theoretically show that TAR² is equivalent to potential-based reward shaping which ensures that the optimal policy learned using TAR² is also optimal under the environment’s original reward function.
- We theoretically prove that any intermediate policy gradient updates under TAR² and the environment’s original reward function share the same direction, ensuring that the trajectory of policy updates remains consistent across both reward functions.
- We empirically validate our method on Alice & Bob, various battle environments in SMACLite (Michalski et al., 2023) and various environment configurations of Google football (Kurach et al., 2020).

2 Related Works

In this section, we review various techniques proposed to address temporal and agent credit assignment in both single-agent and multi-agent systems. We begin by discussing potential-based reward shaping, a method that provides theoretical guarantees of sample-efficient learning of optimal policies in single-agent (Ng, 1999) and multi-agent (Xiaosong Lu, 2011; Devlin & Kudenko, 2011) settings. The use of potential based shaping methods have shown to accelerate the learning process.

2.1 Temporal Credit Assignment

Temporal credit assignment focuses on decomposing sparse or episodic environment rewards into dense reward functions by attributing credit to each time step in an episode.

RUDDER (Arjona-Medina et al., 2019) and its variants (Patil et al., 2020) use contribution analysis to break down episodic rewards into per-time-step rewards by computing the difference between predicted returns at successive time steps. Similarly, Zhang et al. (2023a) perform return-equivalent contribution analysis. Liu et al. (2019) leverage auto-regressive architectures from natural language processing, such as Transformers (Vaswani et al., 2017), to attribute credit to every state-action tuple in the trajectory. Methods by Efroni et al. (2021) and Ren et al. (2021) learn proxy reward functions via trajectory smoothing based on reinforcement learning algorithms that utilize least-squares error. Harutyunyan et al. (2019) introduce a new family of algorithms that use new information to assign credit in hindsight. Han et al. (2022) redesign the value function to predict returns for both historical and current steps by approximating these decompositions. Zhu et al. (2023) propose a bi-level optimization framework to learn a reward redistribution for effective policy learning. These methods have been primarily developed for single-agent settings and may not scale well to multi-agent reinforcement learning (MARL) due to the exponential growth in the joint observation-action space.

In multi-agent settings, recent works have also addressed temporal credit assignment. IRCR (Gangwani et al., 2020) developed a count-based method to learn a proxy reward function for both single and multi-agent settings. AREL (Xiao et al., 2022) uses attention networks to perform reward redistribution, while She et al. (2022) employ an attention encoder network followed by a decoder to address both agent and temporal credit assignment in delayed reward settings.

2.2 Agent Credit Assignment

Most prior works focus on agent credit assignment in multi-agent systems. Devlin et al. (2014) and Foerster et al. (2018) employ difference rewards to assess each agent’s contribution to the global reward. Value Decomposition Networks (VDN) (Sunehag et al., 2017) decompose the joint value function into agent-specific value functions, assuming additivity. Rashid et al. (2020) introduce monotonicity constraints on the joint Q function to learn individual Q values for each agent. Son et al. (2019) generalize this approach to decompose joint Q functions into agent-specific Q functions. Wang et al. (2020) leverage Shapley values to model the joint Q function for agent credit assignment. Zhou et al. (2020) propose an entropy-regularized actor-critic method to efficiently explore multi-agent credit assignment. Freed et al. (2021) use Transformer attention mechanisms in the critic of an actor-critic method to identify relevant agent subgroups for effective multi-agent credit assignment. However, these techniques do not address temporal credit assignment and are therefore inadequate for learning optimal policies in episodic or highly delayed reward settings.

In summary, while significant progress has been made in addressing either agent or temporal credit assignment, the combined challenge of both remains underexplored. Our proposed Temporal-Agent Reward Redistribution (TAR²) aims to fill this gap by effectively handling both agent and temporal credit assignment, enabling efficient learning in multi-agent environments with sparse or delayed rewards.

3 Background

In this section, we introduce our problem setup within the framework of a decentralized partially observable Markov decision process (Dec-POMDP) (Oliehoek & Amato, 2016; Amato, 2024). We describe how agents operate under partial observability and must make decisions based on local observations and histories. We then discuss the episodic multi-agent reinforcement learning (MARL) setting, where agents receive rewards only at the end of each episode, posing a significant challenge for credit assignment. To address this, we explore potential-based reward shaping for multi-agent systems (Ng, 1999; Devlin & Kudenko, 2011), a technique that reshapes the reward function to facilitate learning while preserving the optimal policy. Finally, we analyze how imperfect credit assignment impacts the variance of the policy gradient in multi-agent systems. We show that improper credit distribution among agents leads to high variance in advantage estimates, which in turn exacerbates the learning process and hinders the convergence to optimal policies.

3.1 Decentralized Partially Observable Markov Decision Processes (Dec-POMDP)

A Dec-POMDP is represented by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{T}, \mathcal{O}, \mathcal{N}, \mathcal{R}_\zeta, \rho_0, \gamma)$ where $s \in \mathcal{S}$ is the environment state space, $a \in \mathcal{A}$ is the joint action space denoted by $\mathcal{A} := \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_N$ and $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition function. $r_{\text{global}, t} \sim \mathcal{R}_\zeta(s_t, a_t) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the global reward shared among agents at every timestep of the trajectory. ρ_0 is the initial state distribution and $\gamma \in [0, 1]$ is the discount factor. $\pi = \prod_{i=1}^N \pi_i$ is the joint policy of the multi-agent system which comprises of independent agent policies π_i . Each agent $i \in \{1 \dots N\}$ receives an observation $o_i \in \mathcal{O}_i$ from the observation function $\mathcal{T}(s, i) : \mathcal{S} \times \mathcal{N} \rightarrow \mathcal{O}$. Because the state is not directly observable, it is typically beneficial for each agent to remember a history of its observations or observations-actions. $h_t \in \mathcal{H} := \mathcal{H}_1 \dots \mathcal{H}_N$ is the set of agent observation (-action) histories up to the current time step t where $h_{i,t} \in \mathcal{H}_i$ and defined as $h_{i,t} = \{o_{i,1}, a_{i,1}, \dots, o_{i,t}\}$ denotes agent i 's history and $h_{-i,t}$ is the history of all other agents except agent i . At each time step every agent selects an action $a_i \in \mathcal{A}_i$ according to its policy $\pi_i : H_i \times \mathcal{A}_i \rightarrow [0, 1]$. $\tau = \{o_{0,1}, a_{0,1}, \dots, o_{0,N}, a_{0,N}, \dots, o_{|\tau|,1}, a_{|\tau|,1}, \dots, o_{|\tau|,N}, a_{|\tau|,N}\}$ is the multi-agent trajectory where $|\tau|$ is the horizon length of the trajectory. The goal of the agents is to determine their individual optimal policies that achieve maximum global return $E_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} \left[\sum_{t=1}^{|\tau|} \gamma^t r_t \right]$.

3.2 Return Decomposition in Episodic Multi-Agent Reinforcement Learning

In most MARL systems, each agent receives a global reward $r_{\text{global}, t}$ after executing the joint action a_t in state s_t . However, in episodic MARL setups, agents only receive a global reward signal from the environment at the end of the trajectory, known as the episodic reward or trajectory return $r_{\text{global}, \text{episodic}}$. The objective in such environments is to maximize the trajectory return, $E_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} (r_{\text{global}, \text{episodic}}(\tau))$. Delayed reward settings introduce significant bias and variance (Ng, 1999) during the learning process, exacerbating sample inefficiency.

3.3 Potential-based reward shaping

Ng (1999) presented a single-agent reward shaping method to address the credit assignment problem by introducing a potential-based shaping reward to the environment. The combination of the shaping reward with the original reward can enhance the learning performance of a reinforcement learning algorithm and accelerate the convergence to the optimal policy. Devlin & Kudenko (2011) and Xiaosong Lu (2011) extended potential-based reward shaping to multi-agent systems.

Theorem 1. *Given an n -player discounted stochastic game $M = (S, A_1, \dots, A_n, T, \gamma, R_1, \dots, R_n)$, we define a transformed n -player discounted stochastic game $M' = (S, A_1, \dots, A_n, T, \gamma, R_1 + F_1, \dots, R_n + F_n)$, where $F_i \in S \times S$ is a shaping reward function for player i . We call F_i a potential-based shaping function if F_i has the form:*

$$F_i(s, s') = \gamma \Phi_i(s') - \Phi_i(s),$$

where $\Phi_i : S \rightarrow \mathbb{R}$ is a potential function. Then, the potential-based shaping function F_i is a necessary and sufficient condition to guarantee the Nash equilibrium policy invariance such that:

- **(Sufficiency)** If F_i ($i = 1, \dots, n$) is a potential-based shaping function, then every Nash equilibrium policy in M' will also be a Nash equilibrium policy in M (and vice versa).
- **(Necessity)** If F_i ($i = 1, \dots, n$) is not a potential-based shaping function, then there may exist a transition function T and reward function R such that the Nash equilibrium policy in M' will not be the Nash equilibrium policy in M .

In summary, potential-based reward shaping ensures that Nash equilibrium policies are preserved, enhancing learning without altering the strategic dynamics. This principle underpins our proposed reward redistribution method, which we will validate in the following sections, demonstrating its effectiveness in multi-agent reinforcement learning.

3.4 Impact of Faulty Credit Assignment on Policy Gradient Variance

To understand the impact of imperfect credit assignment, we analyze the effect of other agents on the policy gradient update of agent i . Consider an actor-critic gradient estimate for a multi-agent system in a Dec-POMDP setting, computed using a state-action sample from an arbitrary timestep t . We make no assumptions about the policy parameters of the agents in the multi-agent system. Ideally, the policy gradient update for agent i should be computed using

$$\hat{\nabla}_{\theta_i} J(\theta, h) = \nabla_{\theta_i} \log \pi_i(a_i | h_i) \mathbb{E}_{\neg h_i, \neg a_i} [A(h, a)]$$

Computing $\mathbb{E}_{\neg h_i, \neg a_i} [A(h, a)]$ is challenging to compute due to the high dimensionality, dependency on other agents, and the computational complexity involved in accurately modeling and estimating the interdependent histories and actions of multiple agents. However, in practice multi-agent policy-gradient algorithms like MAPPO (Yu et al., 2022), MADDPG (Lowe et al., 2017) etc employ $A(h, a)$ to compute the policy gradient update for each agent. As a result, the credit assignment problem manifests as high variance in advantage estimates, leading to slower learning because of noisier policy gradient estimates.

Multi-agent policy gradient methods approximate the true *advantage* by computing \hat{A} , which is actually a stochastic advantage estimation of taking a joint action a while observing the joint agent history h , and following the joint policy π . The advantage function is typically defined as $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$, where $Q^\pi(s, a)$ and $V^\pi(s)$ are the state-action value function and state-value function, respectively (Sutton & Barto, 1998). However, in practice the state-action value function and state-value function are approximated using $\hat{Q}^\pi(h, a)$ and $\hat{V}^\pi(h)$ in Dec-POMDPs and there are many ways to compute \hat{A} , generally all involving some error, as the true value functions are unknown (Sutton & Barto, 1998; Schulman et al., 2015).. Intuitively, this is the centralized advantage function which measures how much better it is to select a joint action a than a random action from the joint policy π , while in state s . Besides, to update the policy of agent i , we need to compute the advantage of selecting action a_i , taking into account the specific contribution and context of agent i within the multi-agent system. Perfect credit assignment would be possible if the advantage function could be computed perfectly for each agent, as it directly measures how a particular action of an agent impacted the total reward obtained by the group.

The conditional variance of $\text{Var}(\hat{\nabla}_{\theta_i} J | h, a)$, given h and a , is proportional to the variance of \hat{A} . While this statement typically implies multiple samples are considered to estimate the variance accurately, we can gain insight into the variance introduced by the contributions of other agents by initially

focusing on a single sample scenario.

$$\text{Var}(\hat{\nabla}_\theta J|h, a) = (\nabla_\theta \log \pi(a_i|h_i)) (\nabla_\theta \log \pi(a_i|h_i))^T \text{Var}(\hat{A}|h, a).$$

It is therefore evident that the variance of the policy gradient update is directly proportional to the variance of the advantage estimate: $\text{Var}(\hat{\nabla}_\theta J|h, a) \propto \text{Var}(\hat{A}|h, a)$.

Let us analyze the variance of the advantage function in cooperative multi-agent setting,

$$\mathcal{A}(h, a) = \mathcal{Q}(h, a) - \mathcal{V}(h)$$

Let us assume that agent i 's reward contribution at an arbitrary time-step t is denoted by $r_{i,t}$ and the episodic reward described in subsection 3.2 can be derived using $r_{\text{global,episodic}}(\tau) = \sum_{t=1}^{|\tau|} \sum_{i=1}^N r_{i,t}|h, a$. From this we can rewrite $\mathcal{Q}(h, a)$ and $\mathcal{V}(h)$ as $\mathcal{Q}(h, a) = \mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{i=1}^N r_{i,t}|h, a]$ and $\mathcal{V}(h) = \mathbb{E}_\pi [\mathcal{Q}(h, a)]$

$$\mathcal{A}(h, a) = \mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{i=1}^N r_{i,t}|h, a] - \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{i=1}^N r_{i,t}|h]]$$

Based on the linearity of expectations on $\sum_{t=1}^{|\tau|} \sum_{i=1}^N r_{j,t} = \sum_{t=1}^{|\tau|} r_{i,t} + \sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}$

$$\begin{aligned} \mathcal{A}(h, a) &= (\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h, a] + \mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h, a]) \\ &\quad - (\mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h]] + \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h]]) \\ \mathcal{A}(h, a) &= (\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h, a] - \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h]]) \\ &\quad - (\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h, a] - \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h]]) \end{aligned}$$

The advantage estimate considering only the contribution of agent i is $\mathcal{A}_i = \mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h, a] - \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} r_{i,t}|h]]$ the only advantage term that should be considered while calculating the policy gradient update for agent i whereas the advantage estimate due to other agent contributions $\mathcal{A}_{\neg i} = \mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h, a] - \mathbb{E}_\pi [\mathbb{E}_{s_0 \sim \rho_0, s \sim \mathcal{P}, a_i \sim \pi_i} [\sum_{t=1}^{|\tau|} \sum_{j \neq i}^N r_{j,t}|h]]$ act as noise. Thus,

$$\mathcal{A}(h, a) = \mathcal{A}_i(h, a) + \mathcal{A}_{\neg i}(h, a)$$

Using variance of the sum of random variables,

$$\text{Var}(\mathcal{A}(h, a)) = \text{Var}(\mathcal{A}_i(h, a)) + \text{Var}(\mathcal{A}_{\neg i}(h, a)) + 2\text{Cov}(\mathcal{A}_i(h, a), \mathcal{A}_{\neg i}(h, a))$$

To express the equation in terms of variance, we use the Cauchy-Schwarz inequality, which states that for any two random variables \mathcal{A}_i and $\mathcal{A}_{\neg i}$:

$$\text{Cov}(\mathcal{A}_i(h, a), \mathcal{A}_{\neg i}(h, a)) \leq \sqrt{\text{Var}(\mathcal{A}_i(h, a))\text{Var}(\mathcal{A}_{\neg i}(h, a))}$$

By substituting this inequality, we get an upper bound on our equation,

$$\begin{aligned}\text{Var}(\mathcal{A}(h, a)) &\leq \text{Var}(\mathcal{A}_i(h, a)) + \text{Var}(\mathcal{A}_{\neg i}(h, a)) + 2\sqrt{\text{Var}(\mathcal{A}_i(h, a))\text{Var}(\mathcal{A}_{\neg i}(h, a))} \\ \text{Var}(\mathcal{A}(h, a)) &\leq (\sqrt{\text{Var}(\mathcal{A}_i(h, a))} + \sqrt{\text{Var}(\mathcal{A}_{\neg i}(h, a))})^2\end{aligned}$$

The above equation shows that the variance of the policy gradient update grows approximately linearly with the number of agents in the multi-agent system. This increase in variance reduces the signal-to-noise ratio of the policy gradient, necessitating more updates for effective learning. Proper credit assignment can mitigate this issue by enhancing the signal-to-noise ratio, thereby facilitating more sample-efficient learning.

4 Method

4.1 Definition of reward redistribution function

In this paper, we address the challenge of temporal and agent credit assignment in fully cooperative multi-agent systems with episodic global rewards. Our goal is to learn a reward redistribution function that preserves the optimal policy of the original reward function of the environment. We aim to achieve this by defining a reward redistribution function that decomposes the episodic trajectory reward, also known as trajectory return, to each agent based on their contribution to the team’s outcome at every time step.

We premise that the joint history and action at the final timestep of the multi-agent trajectory serve as a good proxy for predicting the episodic global reward. Consequently, we predict the per timestep reward for each agent by assessing the contribution of its state-action tuple towards generating this joint history of the multi-agent system.

Assumption 1. *The reward redistribution function $r_{i,t}$, which assigns the reward received by agent i at time step t , is determined by analyzing the importance of each state-action tuple against the joint history and action at the final timestep that predicts the episodic global reward.*

$$r_{\text{global}, \text{episodic}}(\tau) = \sum_{i=1}^N \sum_{t=1}^{|\tau|} r_{i,t}(h_{i,t}, a_{i,t}, h_{\neg i,t}, a_{\neg i,t}, h_{|\tau|}, a_{|\tau|}), \quad (1)$$

This assumption aligns with the framework adopted by prior works (Xiao et al., 2022; Ren et al., 2021; Efroni et al., 2021) which also assumes that the episodic return has some structure in nature, e.g., a sum-decomposable form.

4.2 Assembling the reward function

We propose a method to redistribute the trajectory returns temporally, assigning credit to each time step in the multi-agent trajectory. Subsequently, the temporally redistributed rewards are further decomposed across agents based on their individual contributions, ensuring that the relationship expressed in (1) is maintained.

Given the dual nature of credit assignment—attributing relative credit to: (1) each time step in the multi-agent trajectory, and (2) each agent at every time step—we can derive the relationship between the trajectory return, $r_{\text{global}, \text{episodic}}$, and the redistributed reward received by agent i at time step t , $r_{i,t}$.

To formalize, we assume the following process:

Assumption 2. *The trajectory return $r_{\text{global}, \text{episodic}}$ is redistributed temporally to obtain rewards for each time step, $r_{\text{global}, t}$, such that:*

$$r_{\text{global}, \text{episodic}}(\tau) = \sum_{t=1}^{|\tau|} r_{\text{global}, t}(h_t, a_t, h_{|\tau|}, a_{|\tau|}).$$

Subsequently, these temporally redistributed rewards are decomposed across agents:

$$r_{\text{global},t}(h_t, a_t) = \sum_{i=1}^N r_{i,t}(h_{i,t}, a_{i,t}, h_{\neg i,t}, a_{\neg i,t}, h_{|\tau|}, a_{|\tau|}),$$

ensuring that:

$$r_{\text{global},\text{episodic}}(\tau) = \sum_{i=1}^N \sum_{t=1}^{|\tau|} r_{i,t}(h_{i,t}, a_{i,t}, h_{\neg i,t}, a_{\neg i,t}, h_{|\tau|}, a_{|\tau|}).$$

This two-step decomposition process—first temporally and then across agents—ensures that each agent’s contribution at each time step is appropriately credited, thereby providing a robust framework for reward redistribution in multi-agent systems.

Let’s define a function $w_t \sim W_\omega(h_t, a_t, h_{|\tau|}, a_{|\tau|}) : (\mathcal{H} \times \mathcal{A}) \times (\mathcal{H}_{|\tau|} \times \mathcal{A}_{|\tau|}) \rightarrow \mathbb{R}$ that redistributes the rewards across the temporal axis of the multi-agent trajectory. Thus, we can express the multi-agent temporal reward at an arbitrary time step t as

$$r_{\text{global},t} = w_t r_{\text{global},\text{episodic}}(\tau)$$

Similarly, let’s define a function $w'_{i,t} \sim W_\kappa(h_{i,t}, a_{i,t}, h_{\neg i,t}, a_{\neg i,t}) : \mathcal{H}_i \times \mathcal{A}_i \times \mathcal{H}_{\neg i} \times \mathcal{A}_{\neg i} \rightarrow \mathbb{R}$ that redistributes the temporal rewards at an arbitrary time step t across agents. Hence, now we can express the reward that agent i receives as

$$r_{i,t} = w'_{i,t} r_{\text{global},t}$$

Finally, deriving the relationship between $r_{i,t}$ and $r_{\text{global},\text{episodic}}(\tau)$ for an arbitrary time-step t

$$r_{i,t} = w'_{t,i} w_t r_{\text{global},\text{episodic}}(\tau)$$

Based on the definition of reward redistribution function

$$\begin{aligned} \sum_{i=1}^N \sum_{t=1}^{|\tau|} r_{i,t}(h_{i,t}, a_{i,t}, h_{\neg i,t}, a_{\neg i,t}, h_{|\tau|}, a_{|\tau|}) &= r_{\text{global},\text{episodic}}(\tau) \\ \left(\sum_{i=1}^N \sum_{t=1}^{|\tau|} w'_{t,i} w_t r_{\text{global},\text{episodic}}(\tau) \right) &= r_{\text{global},\text{episodic}}(\tau) \\ \sum_{i=1}^N \sum_{t=1}^{|\tau|} w'_{t,i} w_t &= 1 \\ \sum_{t=1}^{|\tau|} \left(\sum_{i=1}^N w'_{t,i} \right) \times w_t &= 1 \end{aligned}$$

The solution for the above equation is,

$$\sum_{i=1}^N w'_{t,i} = 1 \tag{2}$$

$$\sum_{t=1}^{|\tau|} w_t = 1 \tag{3}$$

To construct the new reward function $\mathcal{R}_{\omega,\kappa}$, we incorporate the original reward \mathcal{R}_ζ and the redistributed credit $r_{i,t}$ that each agent i receives at time step t . This redistributed credit reflects the agent’s relevance to the final outcome of the multi-agent system. The relevance of each state-action tuple is determined by the reward redistribution functions W_ω and W_κ , which effectively assign the trajectory rewards temporally and across agents.

$$\begin{aligned} R_{\omega,\kappa}(s_t, a_t, s_{t+1}) &= R_\zeta(s_t, a_t, s_{t+1}) + r_{i,t} \\ R_{\omega,\kappa}(s_t, a_t, s_{t+1}) &= R_\zeta(s_t, a_t, s_{t+1}) + w_{t,i} w_t r_{\text{global,episodic}}(\tau) \end{aligned} \quad (4)$$

This formalizes the reward redistribution process, ensuring that each agent receives a reward proportionate to its contribution to the overall team performance at each timestep.

4.3 Optimal Policy Preservation

To ensure that the optimal policy learned using the densified reward function is also optimal in the environment’s original reward function, we establish the following theorem:

Theorem 2. *Let’s consider two Dec-POMDPs as defined in subsection 3.1, $\mathcal{M}_{\text{env}} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{T}, \mathcal{O}, \mathcal{N}, \mathcal{R}_\zeta, \rho_0, \gamma)$ and $\mathcal{M}_{\text{rrf}} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{T}, \mathcal{O}, \mathcal{N}, \mathcal{R}_{\omega,\kappa}, \rho_0, \gamma)$. The only distinction between \mathcal{M}_{env} and \mathcal{M}_{rrf} are the reward functions. If π_θ^* is the optimal policy in \mathcal{M}_{rrf} then π_θ^* is also optimal in \mathcal{M}_{env} .*

Proof. We know that π_{theta}^* is optimal in \mathcal{M}_{rrf} . For π_{theta}^* to be optimal in \mathcal{M}_{env} , we need to show that $\mathcal{R}_{\omega,\kappa} = \mathcal{R}_\zeta + \mathcal{F}(s_t, a_t, s_{t+1})$ where $\mathcal{F}(s_t, a_t, s_{t+1})$ is a potential based shaping function which is a necessary and sufficient condition for optimal policy preservation 3.3.

It is therefore sufficient to show that the equation (4) takes the form $\mathcal{R}_{\omega,\kappa}(s_t, a_t, s_{t+1}) = \mathcal{R}_\zeta(s_t, a_t, s_{t+1}) + \gamma\phi(s_{t+1}) - \phi(s_t)$. Comparing this format to equation (4), assuming $\gamma = 1$ we arrive at $\phi(s_{t+1}) - \phi(s_t) = w'_{t,i} w_t r_{\text{global,episodic}}(\tau)$. This relation holds for $\phi(s_t) = r_{\text{global,episodic}}(\tau) (\sum_{t'=0}^t w'_{t',i} w_{t'})$ \square

This result ensures that if an arbitrary policy π_θ when trained using the reward function $\mathcal{R}_{\omega,\kappa}$ in \mathcal{M}_{rrf} converges to an optimal policy π_θ^* then π_θ^* will also be optimal for the original reward function \mathcal{R}_ζ in \mathcal{M}_{env} .

4.4 Policy Gradient Update Equivalence with Reward Redistribution

In this subsection, we establish that the policy gradient update for an arbitrary agent k , derived from the reward redistribution function, shares the same direction but exhibits a smaller magnitude than the policy gradient update in the environment’s original reward function. Furthermore, this ensures that the policy update trajectory towards the optimal policy for an arbitrary initial policy is preserved for every agent.

Proposition 1. *Let π_θ be the parametric policy in a decentralized execution paradigm, where the joint policy is expressed as a product of individual agent policies $\pi_\theta = \prod_{k=1}^N \pi_{\theta_k}$ (Oliehoek & Amato, 2016; Amato, 2024). The policy gradient update for an arbitrary agent k under the reward redistribution function $R_{\omega,\kappa}$ follows the same direction as the policy gradient update under the environment’s original reward function R_ζ , preserving the policy update trajectory towards the individual optimal policies.*

Proof. Consider the policy gradient update for agent k under the reward redistribution function:

$$\nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^{|\tau|} r_{k,t} \right] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} [\delta(\tau) r_{\text{global,episodic}}(\tau)] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{i=1}^N \sum_{t=1}^T \delta_t r_{i,t} \right],$$

where τ is the multi-agent trajectory attained from the joint policy π_θ and $\delta : \mathcal{H} \times \mathcal{A} \times \mathcal{H}_{|\tau|} \times \mathcal{A}_{|\tau|} \rightarrow \mathbb{R} \geq 0$ is a non-negative scalar-valued function conditioned on the trajectory. For brevity, we drop the detailed notation of the variables. From the definition of the reward redistribution function in Assumption 1, we have:

$$\begin{aligned} \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} [r_{\text{global,episodic}}(\tau)] &= \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{i=1}^N \sum_{t=1}^T r_{i,t} \right] \\ &= \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^T r_{k,t} + \sum_{i \neq k} \sum_{t=1}^T r_{i,t} \right] \\ &= \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^T r_{k,t} \right] + \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{i \neq k} \sum_{t=1}^T r_{i,t} \right]. \end{aligned}$$

Given the definitions $\sum_{i=1}^N w'_{t,i} = 1$, equation (2), and $\sum_{t=1}^{|\tau|} w_t = 1$, equation (3)), we rewrite the above equation as:

$$\nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} [r_{\text{global,episodic}}(\tau)] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^{|\tau|} r_{k,t} \right] + \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\left(\sum_{t=1}^{|\tau|} w_t (1 - w'_{k,t}) \right) r_{\text{global,episodic}}(\tau) \right]$$

$$\text{Let } (w_t (1 - w'_{k,t})) = M_t,$$

$$\nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} [r_{\text{global,episodic}}(\tau)] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^{|\tau|} r_{k,t} \right] + \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\left(\sum_{t=1}^{|\tau|} M_t \right) r_{\text{global,episodic}}(\tau) \right]$$

$$\nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\left(1 - \sum_{t=1}^{|\tau|} M_t \right) r_{\text{global,episodic}}(\tau) \right] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^{|\tau|} r_{k,t} \right]$$

$$\text{Comparing with } \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} \left[\sum_{t=1}^{|\tau|} r_{k,t} \right] = \nabla_{\theta_k} \mathbb{E}_{\pi_{\theta_k}} [\delta(\tau) r_{\text{global,episodic}}(\tau)], \delta(\tau) = 1 - \sum_{t=1}^{|\tau|} M_t.$$

Since $1 \geq (1 - w'_{k,t}) \geq 0$, $1 \geq w_t \geq 0$, and $\sum_{t=1}^{|\tau|} w_t = 1$, the term $\sum_{t=1}^{|\tau|} w_t \times (1 - w'_{k,t})$ represents a weighted sum, ensuring that:

$$1 \geq \delta(\tau) = 1 - \sum_{t=1}^{|\tau|} M_t \geq 0.$$

Thus, training a policy with the reward redistribution function is equivalent to training with the environment's original reward function, as it preserves the direction of each agent's policy gradient update. This ensures that the policy evolves similarly in both settings and that the policy update trajectory for an arbitrary initial policy is preserved. \square

4.5 Temporal-Agent Reward Redistribution (TAR²) architectural details

We build upon the architecture proposed by Xiao et al. (2022) by extending its capabilities to decompose the episodic reward (trajectory return) not only temporally but also at the agent level. This advancement allows our model to predict $r_{i,t}$, effectively learning implicit temporal and agent weights. These weights satisfy the relationship $r_{i,t} = w'_{t,i} w_t r_{\text{global,episodic}}(\tau)$, thereby enhancing the accuracy and granularity of credit assignment in multi-agent systems.

5 Experimental Setup

We demonstrate the effectiveness of our approach TAR² with single-agent and multi-agent reinforcement learning algorithms against some competitive baselines in the 5m_vs_6m battle scenario of the SMACLite (Michalski et al., 2023) environment.

5.1 Baselines

In order to validate the effectiveness of our reward redistribution mechanism we compare its performance with many other forms of reward functions. We train all the baseline reward functions with IPPO (De Witt et al., 2020; Schulman et al., 2017) and MAPPO (Yu et al., 2022) and report them in Fig 1

Episodic rewards: This is the episodic reward setting where each agent receives a global reward signal at the end of the trajectory.

Dense temporal rewards: In this setting, each agent receives the original global dense reward signal described in subsection 5.2.

Dense AREL temporal rewards: This setting employs AREL reward redistribution that temporally assigns rewards to the multi-agent trajectory as described in (Xiao et al., 2022).

Dense IRCR temporal rewards: In this setting, each agent receives a global reward at every time step following this equation $r_{global,t} = R_{episodic}(\tau)/|\tau|$ (Gangwani et al., 2020). This baseline also exemplifies a unique form of reward redistribution in our case where the state and action tuple of each agent is of equal importance and hence each of them receive the same global reward.

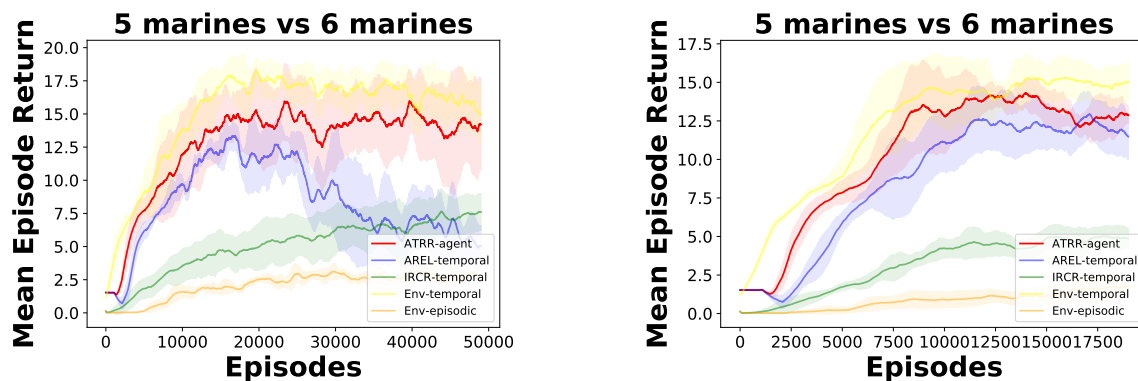
Dense TAR² agent rewards (ours): This is the reward setting proposed in the subsection 4.5.

5.2 Environment

StarCraft Multi-Agent Challenge Lite (SMACLite): The StarCraft Multi-Agent Challenge (SMAC) (Samvelyan et al., 2019) is an RL environment based on the StarCraft II real-time strategy game, in which a team of agents fights against an opposing team controlled by the game engine’s centralized hard-coded AI. We specifically consider the the lightweight and open-source SMACLite version (Michalski et al., 2023). We consider a battle scenario, 5m_vs_6m, where 5 agent-controlled marines battle 6 enemy marines. In this battle situation, the dense reward received by a particular agent while attacking an enemy unit is the difference in the health and shield points removed from that enemy unit in that particular timestep. If a particular agent kills an enemy unit, it receives a reward of 10. Upon defeating the entire enemy team, a reward of (200 / number of agents) is given to each surviving agent. The returns are then normalized such that the maximum possible group return is 20. However, we accumulate the dense reward for each multi-agent trajectory and provide it as a feedback only at the end of the episode.

6 Results and Discussion

As presented in Figure 1, our method TAR²-agent outperform other reward function baselines except for the environment’s original dense reward setting as described in subsection 5.2. This particular baseline is an oracle since it has been manually designed to achieve the objective of this specific environment. While training TAR², we used the same hyperparameters as proposed in (Xiao et al., 2022) with a slight modification to the training procedure. Since AREL (Xiao et al., 2022) was trained with off-policy reinforcement learning algorithms like QMIX (Rashid et al., 2020), they seemed to not require a warm-up period to train the reward function alone. Since in our experiments we train single and multi-agent on-policy policy gradient algorithms, we empirically discovered that a warm-up period (2000 episodes) performed better.



(a) Performance of IPPO in different reward settings.

(b) Performance of MAPPO in different reward settings.

Figure 1: Average agent episodic rewards with standard deviation for task 5m_vs_6m.

7 Conclusion and future work

This paper studied the multi-agent agent-temporal credit assignment problem in MARL tasks with episodic rewards. We proposed a agent-temporal reward redistribution (TAR²) function that theoretically guarantees the preservation of the optimal policy under the original reward function. Our experimental results demonstrate that TAR² outperforms all baselines, showing faster convergence speed.

In future work, we want to explore the agent-temporal reward redistribution by utilizing the attention weights generated by the temporal and agent attention blocks during a forward pass since they naturally fit well in the proposed framework. We want to also demonstrate the effectiveness of our approach against more competitive state-of-the-art baselines and across a variety of other MARL environments of varying difficulty. An interesting line of investigation would be to see the transfer-learning capabilities of such models 1) with more agents than it was trained with 2) across different environments with similar objectives.

References

- Aakriti Agrawal, Amrit Singh Bedi, and Dinesh Manocha. Rta: An attention inspired reinforcement learning method for multi-robot task allocation in warehouse environments. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1393–1399, 2023. doi: 10.1109/ICRA48891.2023.10161310. 1
- Stefano V. Albrecht, Filippos Christianos, and Lukas Schäfer. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024. URL <https://www.marl-book.com>. 2
- Christopher Amato. (a partial survey of) decentralized, cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2405.06161*, 2024. 1, 4, 9
- Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 3
- Schirin Baer, Jupiter Bakakeu, Richard Meyes, and Tobias Meisen. Multi-agent reinforcement learning for job shop scheduling in flexible manufacturing systems. In *2019 Second International Conference on Artificial Intelligence for Industries (AI4I)*, pp. 22–25, 2019. doi: 10.1109/AI4I46381.2019.00014. 1

- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019. 1
- Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008. doi: 10.1109/TSMCC.2007.913919. 1
- Mehul Damani, Zhiyao Luo, Emerson Wenzel, and Guillaume Sartoretti. Primal _2: Pathfinding via reinforcement and imitation multi-agent learning-lifelong. *IEEE Robotics and Automation Letters*, 6(2):2666–2673, 2021. 1
- Christian Schroeder De Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020. 2, 11
- Sam Devlin and Daniel Kudenko. Theoretical considerations of potential-based reward shaping for multi-agent systems. In *Adaptive Agents and Multi-Agent Systems*, 2011. URL <https://api.semanticscholar.org/CorpusID:1116773>. 2, 3, 4
- Sam Devlin, Logan Yliniemi, Daniel Kudenko, and Kagan Tumer. Potential-based difference rewards for multiagent reinforcement learning. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pp. 165–172, 2014. 3
- Yonathan Efroni, Nadav Merlis, and Shie Mannor. Reinforcement learning with trajectory feedback. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 7288–7295, 2021. 3, 7
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018. 1, 2, 3
- Benjamin Freed, Aditya Kapoor, Ian Abraham, Jeff Schneider, and Howie Choset. Learning cooperative multi-agent policies with partial reward decoupling. *IEEE Robotics and Automation Letters*, 7(2):890–897, 2021. 2, 3
- Tanmay Gangwani, Yuan Zhou, and Jian Peng. Learning guidance rewards with trajectory-space smoothing. *Advances in Neural Information Processing Systems*, 33:822–832, 2020. 2, 3, 11
- Beining Han, Zhizhou Ren, Zuofan Wu, Yuan Zhou, and Jian Peng. Off-policy reinforcement learning with delayed rewards. In *International Conference on Machine Learning*, pp. 8280–8303. PMLR, 2022. 3
- Anna Harutyunyan, Will Dabney, Thomas Mesnard, Mohammad Gheshlaghi Azar, Bilal Piot, Nicolas Heess, Hado P van Hasselt, Gregory Wayne, Satinder Singh, Doina Precup, et al. Hindsight credit assignment. *Advances in neural information processing systems*, 32, 2019. 3
- Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al. Human-level performance in 3d multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019. 1
- Aleksandar Krnjaic, Raul D Steleac, Jonathan D Thomas, Georgios Papoudakis, Lukas Schäfer, Andrew Wing Keung To, Kuan-Ho Lao, Murat Cubuktepe, Matthew Haley, Peter Börsting, et al. Scalable multi-agent reinforcement learning for warehouse logistics with robotic and human co-workers. *arXiv preprint arXiv:2212.11498*, 2022. 1

- Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zajac, Olivier Bachem, Lasse Espeholt, Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research football: A novel reinforcement learning environment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 4501–4510, 2020. 1, 2
- Yang Liu, Yunan Luo, Yuanyi Zhong, Xi Chen, Qiang Liu, and Jian Peng. Sequence modeling of temporal credit assignment for episodic reinforcement learning. *arXiv preprint arXiv:1905.13420*, 2019. 2, 3
- Ryan Lowe, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in neural information processing systems*, 30, 2017. 5
- Adam Michalski, Filippos Christianos, and Stefano V Albrecht. Smaclite: A lightweight environment for multi-agent reinforcement learning. *arXiv preprint arXiv:2305.05566*, 2023. 2, 11
- AY Ng. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the 16th International Conference on Machine Learning*, pp. 278, 1999. 2, 3, 4
- Frans A. Oliehoek and Chris Amato. A concise introduction to decentralized pomdps. In *Springer-Briefs in Intelligent Systems*, 2016. URL <https://api.semanticscholar.org/CorpusID:3263887>. 1, 4, 9
- Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. Benchmarking multi-agent deep reinforcement learning algorithms in cooperative tasks. In *NeurIPS Datasets and Benchmarks*, 2020. URL <https://api.semanticscholar.org/CorpusID:235417602>. 2
- Vihang P Patil, Markus Hofmarcher, Marius-Constantin Dinu, Matthias Dorfer, Patrick M Blies, Johannes Brandstetter, Jose A Arjona-Medina, and Sepp Hochreiter. Align-rudder: Learning from few demonstrations by reward redistribution. *arXiv preprint arXiv:2009.14108*, 2020. 3
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020. 2, 3, 11
- Zhizhou Ren, Ruihan Guo, Yuan Zhou, and Jian Peng. Learning long-term reward redistribution via randomized return decomposition. *arXiv preprint arXiv:2111.13485*, 2021. 2, 3, 7
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019. 11
- Guillaume Sartoretti, Justin Kerr, Yunfei Shi, Glenn Wagner, TK Satish Kumar, Sven Koenig, and Howie Choset. Primal: Pathfinding via reinforcement and imitation multi-agent learning. *IEEE Robotics and Automation Letters*, 4(3):2378–2385, 2019. 1
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015. 5
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. 2, 11
- Jennifer She, Jayesh K Gupta, and Mykel J Kochenderfer. Agent-time attention for sparse rewards multi-agent reinforcement learning. *arXiv preprint arXiv:2210.17540*, 2022. 2, 3
- Omkar Shelke, Pranavi Pathakota, Anandsingh Chauhan, Harshad Khadilkar, Hardik Meisheri, and Balaraman Ravindran. Multi-agent learning of efficient fulfilment and routing strategies in e-commerce. *arXiv preprint arXiv:2311.16171*, 2023. 1

- Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International conference on machine learning*, pp. 5887–5896. PMLR, 2019. 2, 3
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017. 2, 3
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, MA, 1998. 5
- Ming Tan. Multi-agent reinforcement learning: Independent versus cooperative agents. In *International Conference on Machine Learning*, 1997. URL <https://api.semanticscholar.org/CorpusID:268857333>. 2
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- Eugene Vinitsky, Nathan Lichtle, Kanaad Parvate, and Alexandre Bayen. Optimizing mixed autonomy traffic flow with decentralized autonomous vehicles and multi-agent rl. *arXiv preprint arXiv:2011.00120*, 2020. 1
- Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, L. Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom Le Paine, Caglar Gulcehre, Ziyun Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575:350 – 354, 2019. URL <https://api.semanticscholar.org/CorpusID:204972004>. 1
- Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. Shapley q-value: A local reward approach to solve global reward games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 7285–7292, 2020. 3
- Baicen Xiao, Bhaskar Ramasubramanian, and Radha Poovendran. Agent-temporal attention for reward redistribution in episodic multi-agent reinforcement learning. *arXiv preprint arXiv:2201.04612*, 2022. 2, 3, 7, 10, 11
- Sidney N. Givigi Jr Xiaosong Lu, Howard M. Schwartz. Policy invariance under reward transformations for general-sum stochastic games. *Journal of Artificial Intelligence Research*, 41:397–406, 2011. 3, 4
- Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information Processing Systems*, 35:24611–24624, 2022. 5, 11
- Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pp. 5872–5881. PMLR, 2018. 1
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of reinforcement learning and control*, pp. 321–384, 2021. 1

- Yudi Zhang, Yali Du, Biwei Huang, Ziyan Wang, Jun Wang, Meng Fang, and Mykola Pechenizkiy. Grd: A generative approach for interpretable reward redistribution in reinforcement learning. *arXiv preprint arXiv:2305.18427*, 2023a. 3
- Yulin Zhang, William Macke, Jiaxun Cui, Sharon Hornstein, Daniel Urieli, and Peter Stone. Learning a robust multiagent driving policy for traffic congestion reduction. *Neural Computing and Applications*, pp. 1–14, 2023b. 1
- Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit credit assignment for cooperative multi-agent reinforcement learning. *Advances in neural information processing systems*, 33:11853–11864, 2020. 3
- Tianchen Zhu, Yue Qiu, Haoyi Zhou, and Jianxin Li. Towards long-delayed sparsity: learning a better transformer through reward redistribution. IJCAI '23, 2023. ISBN 978-1-956792-03-4. doi: 10.24963/ijcai.2023/522. URL <https://doi.org/10.24963/ijcai.2023/522>. 3