
Cooperative Multi-Agent Reinforcement Learning with Sequential Credit Assignment

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Centralized training with decentralized execution is a standard paradigm for coop-
2 erative multi-agent reinforcement learning (MARL), with credit assignment being
3 a major challenge. In this paper, we propose a cooperative MARL method with
4 sequential credit assignment (SeCA) that deduces each agent’s contribution to the
5 team’s success one by one to learn better cooperation. We first present a sequential
6 MARL framework, under which we introduce a new counterfactual advantage to
7 evaluate each agent based on its preceding agents’ actions in a specific sequence.
8 As this credit assignment sequence tremendously impacts the performance, we
9 further present a sequence adjustment algorithm utilizing integrated gradients. It
10 dynamically modifies the sequence among agents according to their contribution
11 to the team. SeCA employs a network which either estimates the Q value for
12 training the centralized critic or deduces the proposed advantage of each agent for
13 decentralized policy learning. Our method is evaluated on a challenging set of
14 StarCraft II micromanagement tasks and achieves state-of-the-art performance.

15 1 Introduction

16 Cooperative multi-agent reinforcement learning (MARL) is a helpful tool in numerous applications
17 such as robot swarm control [9], autonomous vehicle coordination [3], network routing [36], and
18 productivity optimization [37]. This kind of problem where agents learn coordinated policies to
19 optimize the global reward has been extensively studied in recent years [7, 19, 18, 38, 8].

20 One natural way of addressing the cooperative MARL problem is the *centralized* approach, which
21 treats the team as a single actor with a joint action space. Although we can trivially apply single-agent
22 reinforcement learning algorithms to such settings, it usually does not scale well because the size of
23 the joint action space grows exponentially with the number of agents. Besides, it is not applicable
24 in real-world settings due to the inherent constraints on agent observability and communication.
25 An alternative approach is to learn *decentralized* policies by independently training agents based
26 on their local observations, but simultaneous exploration often brings non-stationarity that causes
27 unstable learning and difficulties in convergence. As a result, the majority of work on MARL
28 follows the *centralized training with decentralized execution* (CTDE) paradigm [17, 10, 22, 6], where
29 decentralized policies can access extra state information during training.

30 A crucial challenge of the CTDE paradigm in cooperative settings is to correctly deduce each agent’s
31 contribution to the team’s success, also known as the *multi-agent credit assignment problem* [4].
32 Existing methods can be classified as *implicit* and *explicit* credit assignment [39]. Previous implicit
33 methods often deduce all agents’ contributions by representing the global state-action value as an
34 aggregation of each agent’s state-action value [26, 22, 12, 24, 21, 29] and assigning the shared rewards
35 to each agent according to the joint action at one time. In this way, these methods avoid the complex
36 interaction analysis and instead fit these cooperation relationships by neural networks. However,

37 implicit methods often face limitations in expressiveness, and their extensions to continuous action
38 spaces may require additional strategies [39].

39 On the other hand, recognized explicit approaches calculate difference rewards [34] against a certain
40 reward baseline [28, 20, 6]. However, in cooperative MARL, evaluating any agent’s action requires
41 considering the actions of all agents, so it is often difficult to determine the impact a particular agent’s
42 behavior has on the team when we have not assessed other agents’ actions. In other words, we can
43 not say that a single agent’s action is bad if the team receives a small reward because the shared
44 reward is not decided only by this agent’s behavior. Maybe its action is actually good in that state.

45 This paper presents a sequential credit assignment SeCA to evaluate individual agent actions explicitly
46 and sequentially. Our motivation is to address the drawbacks of implicit methods that neglect the
47 cooperation between agents or simply leave it to neural networks and further improve explicit credit
48 assignment. In summary, we face two main challenges to learn a better explicit credit assignment: (1)
49 how to alleviate the problem that it is hard to accurately deduce the contribution of one agent without
50 previously assessing all the others’ action, and (2) how to evaluate agents better in an explicit way.

51 To deal with (1), we introduce a sequential MARL framework. As mentioned above, without assessing
52 the behaviors of other agents, we would never be able to evaluate a given agent’s action accurately.
53 However, we point out in this paper that some agents are less affected by such influences than others,
54 and we can first assign credit to them. For instance, evaluating a staff’s action needs to take the
55 CEO’s command or action into consideration, while the former has little importance in assessing the
56 CEO. Thus, we could evaluate the CEO first without considering the staff’s behavior and then analyze
57 the staff based on the CEO’s action. We fully consider the action coordination between agents and
58 explicitly deduce contribution to them one by one according to a particular order, so as to make up
59 for the disadvantage of implicit methods that the cooperation is only inexplicably fitted by neural
60 networks. Intuitively, the order significantly impacts the overall performance, so we further propose
61 an algorithm to adjust the sequence dynamically through integrated gradients [25].

62 As for (2), we compute an advantage function for each agent to attribute agent contributions explicitly.
63 COMA [6] is a representative method that computes a baseline for each agent to reason about
64 counterfactuals in which only one agent a ’s action changes, so its evaluation of a ’s action is based on
65 the joint action \mathbf{u}^{-a} of other agents. In other words, the policy gradient of COMA only encourages
66 agent a to learn in the direction that benefits the team while other agents are acting \mathbf{u}^{-a} , but the
67 others’ actions are not necessarily \mathbf{u}^{-a} when executing. Unlike COMA, we focus more on the action
68 coordination among agents and propose a new advantage under the proposed sequential framework.

69 We summarize the contributions of this paper as follows: (1) We propose a sequential MARL
70 framework in Section 3.2; (2) Under this framework, we introduce a sequential advantage function
71 for each agent to guide their learning explicitly in Section 3.3. We further prove that the sequential
72 credit assignment we proposed achieves additive advantage-decomposition. (3) We present a sequence
73 adjustment algorithm based on integrated gradients to modify the credit assignment order dynamically
74 in Section 3.4. This algorithm alleviates the impact caused by the sequence’s randomness and helps
75 achieve competitive performance on a challenging set of StarCraft II micromanagement tasks [23].

76 2 Related Work

77 Explicit credit assignment gives valuable insights into agent actions’ contributions to the shared
78 team reward and substantially promotes policy optimization. The representative method COMA [6]
79 utilizes a counterfactual baseline that marginalizes out a single agent’s action while keeping the other
80 agents’ actions fixed to calculate the advantage function. However, the advantage evaluates a single
81 agent’s action based on the other agents’ current behaviors and ignores different action combinations.
82 SQDDPG [30] distributes the global reward reflecting each agent’s contribution through Shapley
83 Value. Although SQDDPG provides a theoretically justified framework, its assumption on the
84 observability and convex game makes it impractical and performs poorly in complex environments.

85 Implicit methods are a more common way when addressing the credit assignment challenge. Among
86 them, LICA [39] is a policy-based method, which learns an end-to-end differentiable optimization
87 where it trains a hypernetwork that maps the state into a set of weights which, in turn, maps the
88 action policies into the Q estimate. On the other hand, value-based methods often represent the
89 global state-action value as an aggregation of the individual values. The value decomposition is linear

90 in the earlier work VDN [26], and it ignores the state information. QMIX [22] learns a non-linear
 91 mixing network with the global state and maps the individual state-action values into the joint Q value
 92 estimate. Although QMIX performs well in various environments, it still faces the mixing network’s
 93 monotonicity constraint limitation. QTRAN [24] further avoids the representation limitations by
 94 using linear constraints between individual utilities and the global state-action value. It guarantees
 95 optimal decentralization, but its constraints are computationally intractable, and the relaxations often
 96 lead to unsatisfied performance. QPLEX [29] decomposes Q values following the dueling structure,
 97 transferring the monotonicity condition from Q values to advantage values. QPD [35] leverages the
 98 integrated gradient attribution technique to decompose global Q values along trajectory paths based
 99 on the assumption that an agent’s local reward is linearly correlated with its contribution to the team.

100 3 Methods

101 3.1 Preliminaries

102 **Notations.** This work considers a fully cooperative multi-agent task with n agents $\mathcal{A} = \{1, \dots, n\}$
 103 as a Dec-POMDP [16] defined by a tuple $G = (S, U, P, r, Z, O, n, \gamma)$. The environment has a
 104 true state $s \in S$. Each agent a chooses an action u_t^a from its action space U at each timestep
 105 t and forms a joint action \mathbf{u}_t that induces a transition in the environment according to the state
 106 transition function $P(s_{t+1}|s_t, \mathbf{u}_t) : S \times U^n \times S \rightarrow [0, 1]$. The agents share the same reward function
 107 $r(s, \mathbf{u}) : S \times U^n \rightarrow \mathbb{R}$, and $\gamma \in [0, 1)$ is the discount factor. We consider partially observable
 108 scenarios in which agent a acquires its local observation $z^a \in Z$ drawn from $O(s_t, a) : S \times \mathcal{A} \rightarrow Z$.
 109 Each agent has an action-observation history $\tau^a \in T \equiv (Z \times U)^*$, on which it conditions a policy
 110 $\pi^a(u^a|\tau^a) : T \times U \rightarrow [0, 1]$. We denote joint quantities over agents in bold and joint quantities over
 111 agents other than a given agent a with the superscript $-a$.

112 **Integrated Gradients.** Many works aim to attribute the predictions of deep networks to their input
 113 features [1, 15, 2]. As one of them, integrated gradients [25] aggregates the gradients along the inputs
 114 that fall on the lines between the baseline \vec{b} and the input $\vec{x} = (x_1, \dots, x_j, \dots, x_d)$. It explains how
 115 much one feature affects the deep network output F while changing from $F(\vec{b})$ to $F(\vec{x})$ along a
 116 path between \vec{b} and \vec{x} . Given a path function $\tau(\alpha)$ with $\alpha \in [0, 1]$ specifying a path from baseline
 117 $\tau(0) = \vec{b}$ to the input $\tau(1) = \vec{x}$, then integrated gradients along the j^{th} dimension is acquired by:

$$c_j = \text{PathIG}_j^\tau(\vec{x}) ::= \int_0^1 \frac{\partial F(\tau(\alpha))}{\partial \tau_j(\alpha)} \frac{\partial \tau_j(\alpha)}{\partial \alpha} d\alpha, \quad (1)$$

118 where c_j represents x_j ’s contribution to the difference between baseline prediction $F(\vec{b})$ and $F(\vec{x})$.
 119 In this work, we leverage the integrated gradients technique to dynamically adjust the order of our
 120 proposed sequential credit assignment according to each agent’s contribution to the team.

121 3.2 Sequential MARL Framework

122 The relationship in a multi-agent system is complicated, as every agent makes decisions based on
 123 the environment interfered with by the other agents. If we model each agent as a node and model
 124 the cooperations between them as edges, the cooperative relationship will be built as a complicated
 125 web-like graph shown in Figure 1(a). Evaluating the actions of any agent should take into account
 126 the behaviors of other agents in this situation. It is hard to judge whether an agent’s current action is
 127 beneficial to the team when we have not evaluated other agents’ actions. If we cannot determine an
 128 analysis order, we can only analyze all the agents implicitly as most existing methods did, and the
 129 cooperation is often fitted only by deep neural networks, leading to unsatisfactory results.

130 This section presents a sequential framework for cooperative MARL, which aims to analyze agents’
 131 actions one by one. Our key assumption is that evaluations of some agents in a team are less affected
 132 than others. Thus we can study these less-affected agents first and then analyze the others based on
 133 the actions of these already-studied agents. For instance, when evaluating a staff’s action, the CEO’s
 134 decision plays a vital role because we have to judge whether the staff obeys the command or not. On
 135 the contrary, the staff intuitively has little impact on evaluating the CEO’s decision. In assessing the
 136 CEO, we often consider external factors such as market situation, modeled as state s in MARL.

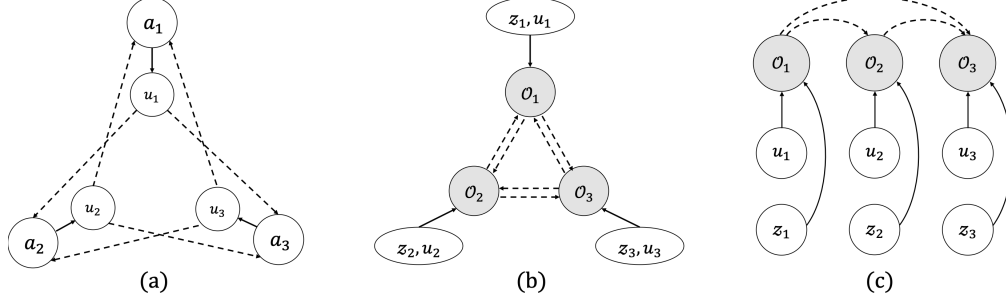


Figure 1: A toy example with three agents. (a) Agents affect each other as they choose actions based on the state interfered with by the others' actions. (b) The study on one agent will influence all the other agents' assessments in the original MARL framework. Agent's cooperation analyses are interrelated. (c) Each agent's cooperation study in the proposed sequential MARL framework. Dotted arrows representing correlations decrease from 6 in (b) to 3 in (c), reducing the complexity by half. This merit also holds for systems with other numbers of agents.

137 We introduce a variable \mathcal{O}_i to help model this sequential MARL framework. This additional variable
 138 represents a random event that our cooperation study (e.g., credit assignment) on agent a_i is optimal or
 139 precise. Then the probability $p(\mathcal{O}_i)$ denotes the accuracy of our research on agent a_i . For illustration
 140 and understanding convenience, we discuss a simple multi-agent system with three agents as an
 141 example, in which agents are identified by a_i ($i \in \{1, 2, 3\}$). In original MARL, the evaluation of
 142 agent a_i will influence all the other agents' assessments. Thus events $\mathcal{O}_1, \mathcal{O}_2$ and \mathcal{O}_3 are mutually
 143 dependent, as shown in Figure 1(b). We calculate the probability of studying the system accurately
 144 by computing conditional probabilities:

$$p(\mathcal{O}_1, \mathcal{O}_2, \mathcal{O}_3) = p(\mathcal{O}_1) \cdot p(\mathcal{O}_2|\mathcal{O}_1) \cdot p(\mathcal{O}_3|\mathcal{O}_1, \mathcal{O}_2) \quad (2a)$$

$$\begin{aligned} & \vdots \\ & = p(\mathcal{O}_3) \cdot p(\mathcal{O}_2|\mathcal{O}_3) \cdot p(\mathcal{O}_1|\mathcal{O}_2, \mathcal{O}_3) \end{aligned} \quad (2b)$$

145 where $p(\mathcal{O}_j|\mathcal{O}_i)$ denotes the probability of agent a_j 's accurate analysis under the condition of
 146 conducting a precise study on agent a_i . It also indicates the accuracy of a_j 's analysis conditions on
 147 precisely assess a_i . We then conclude that:

$$p(\mathcal{O}_1, \mathcal{O}_2, \mathcal{O}_3) = p(\mathcal{O}_i) \cdot p(\mathcal{O}_j|\mathcal{O}_i) \cdot p(\mathcal{O}_k|\mathcal{O}_i, \mathcal{O}_j) \quad (3)$$

148 where $i, j, k \in \{1, 2, 3\}, i \neq j, k \neq i, j$.

149 We take Equ.(2a) as an example. To study the cooperation of this multi-agent system precisely (i.e.,
 150 big $p(\mathcal{O}_1, \mathcal{O}_2, \mathcal{O}_3)$), we can first analyze a_1 as accurately as possible (i.e., big $p(\mathcal{O}_1)$) and then go
 151 on to investigate a_2 and a_3 respectively with the best possible accuracy (i.e., big $p(\mathcal{O}_2|\mathcal{O}_1)$ and
 152 $p(\mathcal{O}_3|\mathcal{O}_1, \mathcal{O}_2)$) under the condition of preceding agents' precise analysis.

153 The sequential MARL framework reduces the complexity of the model with six dotted arrows that
 154 indicate correlations between agents' evaluations in Figure 1(b) by half, as those three dotted lines in
 155 Figure 1(c) show. Equ.(3) suggests that we can analyze the cooperation of a multi-agent system in
 156 any order, but from the CEO-Staff example, we can see that the difficulty of analyzing in various
 157 orders is not the same. Further discussion on the sequence will show in Section 3.4.

158 In general, we specify an order to analyze the cooperation in the sequential MARL framework. We
 159 fix an agent's actions after assessing it and study a particular agent based on the fixed actions of its
 160 preceding agents, reflecting the intuition that a CEO's decision has a strong influence on evaluating
 161 the staff in the example mentioned earlier. This sequential MARL framework significantly alleviates
 162 the correlations in studying agents and helps us assess their cooperation more directly.

163 3.3 Sequential Credit Assignment

164 Following the CTDE paradigm, we utilize a centralized critic for each actor to follow a gradient
 165 based on an advantage function A estimated from this critic:

$$g = \nabla_{\theta\pi} \log \pi(u|\tau_t^a) A. \quad (4)$$

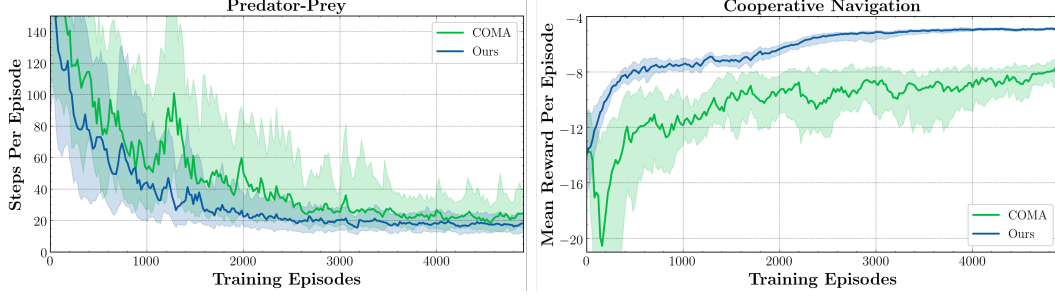


Figure 2: Performances between COMA’s counterfactual advantage and ours in two environments. (Left) *Predator-Prey*. Three predators cooperate to chase a faster prey that acts randomly in an area containing two obstacles. The game terminates when a predator captures the prey, and then a shared reward is given. The predators trained by our advantage capture the prey faster. (Right) *Cooperative Navigation* initializes three agents and three landmarks with random locations. Agents cooperate to cover all the landmarks, and the shared reward is the negative sum of displacements between each landmark and its nearest agent. Our method helps the team gain bigger rewards than COMA.

166 The advantage function A for each actor explicitly deduces how that particular agent contributes to
 167 the team. COMA [6] introduced a counterfactual baseline inspired by difference rewards [34]. For
 168 each agent a , COMA computes an advantage function that compares the Q-value for the action u^a to
 169 a counterfactual baseline that marginalizes out u^a while keeping the others’ actions \mathbf{u}^{-a} fixed:

$$A_{COMA}^a(s, \mathbf{u}) = Q(s, (u^a, \mathbf{u}^{-a})) - \sum_{u'^a} \pi^a(u'^a | \tau^a) \cdot Q(s, (\mathbf{u}^{-a}, u'^a)). \quad (5)$$

170 COMA avoids expensive calculations through careful network design. However, each agent’s
 171 contribution deduced by COMA is still imperfect. The evaluation of u^a is based on the fixed \mathbf{u}^{-a}
 172 in Equ.(5), so agent a will learn a policy that works better with \mathbf{u}^{-a} in this way. It ignores the joint
 173 actions $(u^a, \mathbf{u}^{-a'})$ with $\mathbf{u}^{-a'} \neq \mathbf{u}^{-a}$ that may lead to unexpected results when assessing u^a .

174 To analyze each agent a ’s contribution more objectively, we consider the influence of all joint actions
 175 with u^a . Considering all potential action combinations, we calculate a counterfactual advantage for
 176 each agent’s action, derived by computing the expectation on all the actions of other agents:

$$A^a(s, \mathbf{u}) = \mathbb{E}_{\mathbf{u}^{-a}} [Q(s, (u^a, \mathbf{u}^{-a}))] - \mathbb{E}_{\mathbf{u}^{-a}} \left[\sum_{u'^a} \pi^a(u'^a | \tau^a) \cdot Q(s, (\mathbf{u}^{-a}, u'^a)) \right]. \quad (6)$$

177 Under our proposed sequential MARL framework, we carry out credit assignment according to a
 178 specific order, and there is no need to consider all the possible joint actions. After assessing agent a ,
 179 we fix its action and evaluate agents after it based on a ’s fixed action, so the following agents’ credit
 180 assignments do not have to compute the expectation on u^a anymore.

181 We now give the detailed sequential credit assignment for a team with n agents identified by
 182 $a_i (i \in \{1, \dots, n\})$ under one specific sequence $\{a_1, a_2, \dots, a_n\}$, and it can also be concluded from the
 183 rest $(n! - 1)$ orders in the same way. Here we denote $\mathbf{u}_{a_1}^{a_i-1} = [u^{a_1}, u^{a_2}, \dots, u^{a_{i-1}}] (i = 2, 3, \dots, n)$.

184 As for agent $a_i (i \neq 1)$ in the sequence, the contribution of its leading agents a_1, a_2, \dots, a_{i-1} has
 185 been deduced. We fix the leading agents’ actions and assess agent a_i ’s action based on $\mathbf{u}_{a_1}^{a_i-1}$, so
 186 there is no need to calculate the expectations on $[u^{a_1}, u^{a_2}, \dots, u^{a_{i-1}}]$, simplifying Equ.(6) to:

$$\begin{aligned} A^{a_i}(s, \mathbf{u}) &= \sum_{u'^{a_{i+1}}} \dots \sum_{u'^{a_n}} \pi^{a_{i+1}}(u'^{a_{i+1}} | \tau^{a_{i+1}}) \dots \pi^{a_n}(u'^{a_n} | \tau^{a_n}) \cdot Q(s, (\mathbf{u}_{a_1}^{a_i}, u'^{a_{i+1}}, \dots, u'^{a_n})) \\ &\quad - \sum_{u'^{a_i}} \dots \sum_{u'^{a_n}} \pi^{a_i}(u'^{a_i} | \tau^{a_i}) \dots \pi^{a_n}(u'^{a_n} | \tau^{a_n}) \cdot Q(s, (\mathbf{u}_{a_1}^{a_i-1}, u'^{a_i}, \dots, u'^{a_n})). \end{aligned} \quad (7)$$

187 Then the first agent a_1 ’s advantage is:

$$\begin{aligned} A^{a_1}(s, \mathbf{u}) &= \sum_{u'^{a_2}} \dots \sum_{u'^{a_n}} \pi^{a_2}(u'^{a_2} | \tau^{a_2}) \dots \pi^{a_n}(u'^{a_n} | \tau^{a_n}) \cdot Q(s, (u^{a_1}, u'^{a_2}, \dots, u'^{a_n})) \\ &\quad - \sum_{u'^{a_1}} \dots \sum_{u'^{a_n}} \pi^{a_1}(u'^{a_1} | \tau^{a_1}) \dots \pi^{a_n}(u'^{a_n} | \tau^{a_n}) \cdot Q(s, (u'^{a_1}, u'^{a_2}, \dots, u'^{a_n})) \end{aligned} \quad (8)$$

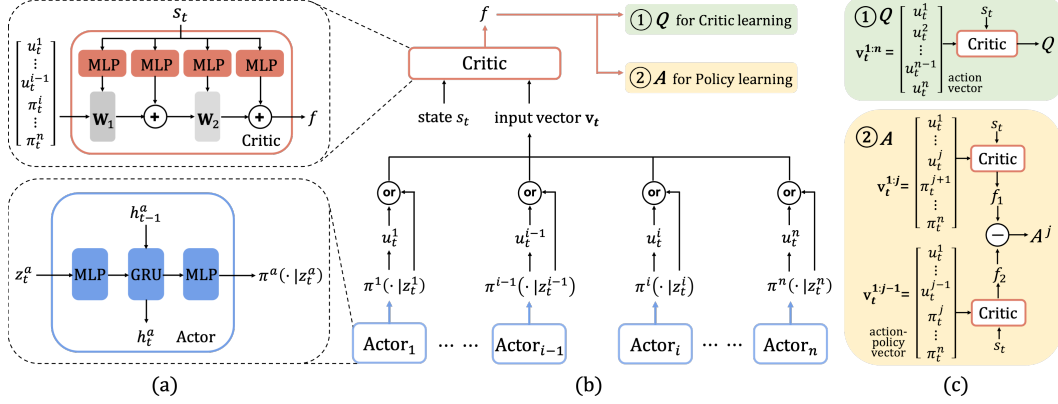


Figure 3: (a) A centralized mixing critic network that maps the state into a set of weights (top) and the decentralized agent network structure (bottom). (b) The overall SeCA architecture. (c) Critic learning (top) and policy learning (bottom) flow. View in color if possible for better understanding.

188 To illustrate the effectiveness of our sequential counterfactual advantage, we conduct a simple
 189 but illuminating test in two common multi-agent particle environments [11], *Predator-Prey* and
 190 *Cooperative Navigation*. We train both methods with 5 random seeds, and agents are trained for 5000
 191 episodes. We provide detailed information on the environments and experiments in the Appendix. As
 192 shown in Figure 2, our sequential advantage functions help agents handle the task faster and better.

193 Our sequential advantage for each agent achieves an additive decomposition of the total advantage
 194 function, which to some extent explains the soundness and superiority of our advantage over COMA’s.

195 **Claim 1.** *The proposed sequential credit assignment achieves additive advantage-decomposition.*

196 *Proof.* See Appendix A. □

197 Facing the same problem as COMA that those evaluations are expensive, we model the first term
 198 in Equ.(7) as a function f_ϕ of $(u^{a_1}, u^{a_2}, \dots, u^{a_i}, \pi^{a_{i+1}}, \dots, \pi^{a_n})$ to address this issue, and the second
 199 term is a similar function of $(u^{a_1}, u^{a_2}, \dots, u^{a_{i-1}}, \pi^{a_i}, \dots, \pi^{a_n})$. Thus, we rewrite Equ.(7) as:

$$A^{a_i} = f_\phi(s; u^{a_1}, u^{a_2}, \dots, u^{a_i}, \pi^{a_{i+1}}, \dots, \pi^{a_n}) - f_\phi(s; u^{a_1}, u^{a_2}, \dots, u^{a_{i-1}}, \pi^{a_i}, \dots, \pi^{a_n}). \quad (9)$$

200 Here f_ϕ is a function evaluating agents’ action-policy vectors, where $f_\phi(u^{a_1}, u^{a_2}, \dots, u^{a_n}) = Q$ and
 201 $f_\phi(\pi^{a_1}, \pi^{a_2}, \dots, \pi^{a_n}) = V$. We design the complete setup for SeCA, which is illustrated in Figure 3.

202 **Critic Learning.** We train critic f_ϕ on-policy to estimate Q , utilizing a practical variant of TD(λ) [27]
 203 adapted for use with deep neural networks. In particular, the critic parameter ϕ is updated by minibatch
 204 gradient descent to minimize the following loss:

$$\mathcal{L}_t(\phi) = \left(y_t^{(\lambda)} - f_\phi(s_t, \mathbf{u}_t) \right)^2, \text{ where } y_t^{(\lambda)} = r_t + \gamma \left(\lambda y_{t+1}^{(\lambda)} + (1 - \lambda) f_\phi(s_{t+1}, \mathbf{u}_{t+1}) \right). \quad (10)$$

205 We utilize a target critic f_{ϕ^-} [14] to improve learning stability and update $\phi^- \leftarrow \phi$ periodically. The
 206 critic learning flow is shown at the top of Figure 3(c). The input for critic training is the state s and
 207 the action vector $\mathbf{u} = [u^1, u^2, \dots, u^n]$ denoted as $\mathbf{v}^{1:n}$.

208 **Policy Learning.** We optimize each agent a ’s policy parameter θ_a by maximizing the following
 209 objective, which contains our proposed advantage function and an entropy regularization term \mathcal{H} :

$$g^a = \mathbb{E}_{\tau \sim \pi} [\nabla_{\theta_a} \log \pi^a(u^a | \tau^a) A^a(s, \mathbf{u}) + \mathcal{H}(\pi^a(\cdot | \tau^a))], \quad (11)$$

210 where the derivative of the adaptive entropy regularization term $\mathcal{H}(\pi^a(\cdot | \tau^a))$ [39] with respect to the
 211 i -th action probability p_i^a is given by:

$$d\mathcal{H}_i := -\xi \cdot (\log p_i^a + 1) / H(\pi^a(\cdot | \tau^a)), \text{ where } H(\pi^a(\cdot | \tau^a)) = \mathbb{E}_{u^a \sim \pi^a} [-\log \pi^a(u^a | \tau^a)]. \quad (12)$$

212 We share parameters among agents, and the gradient we use to train the actor shared by all agents is:

$$g = \mathbb{E}_{\tau \sim \pi} \left[\sum_a (\nabla_{\theta_a} \log \pi^a(u^a | \tau^a) A^a(s, \mathbf{u}) + \mathcal{H}(\pi^a(\cdot | \tau^a))) \right]. \quad (13)$$

213 The inputs of the centralized critic f_ϕ to compute the advantage function are the state s and two
 214 action-policy vectors $\mathbf{v}^{1:i} = [u^1, \dots, u^i, \pi^{i+1}, \dots, \pi^n]$ and $\mathbf{v}^{1:i-1} = [u^1, \dots, u^{i-1}, \pi^i, \dots, \pi^n]$. The
 215 bottom of Figure 3(c) demonstrates the policy learning flow.

216 3.4 Sequence Adjustment Through Integrated Gradients

217 We apply integrated gradients to adjust the credit assignment sequence dynamically. Reviewing the
 218 enlightening and straightforward CEO-Staff example discussed in Section 3.2, we can evaluate the
 219 staff’s behavior based on the CEO’s decision, but assessing the CEO does not require much attention
 220 to the staff’s action. Therefore, we would analyze the CEO first and then evaluate the staff based on
 221 the CEO’s current action. However, this example is not generalized for two reasons: (1) There are
 222 often multiple agents taking the same role in a system with superior-subordinate relationships, and the
 223 sequence of these agents is hard to determine; (2) Not all scenarios have such superior-subordinate
 224 relationships. The agents often do not need to follow others’ commands in many applications.

225 We generalize the CEO-Staff example to propose a universal model. Instead of focusing on the roles
 226 among the agents as in [31, 32], we are more interested in agents’ contributions. Although the CEO
 227 and the staff have a superior-subordinate relationship, they are essentially employees of an enterprise.
 228 The staff plays an auxiliary role and acts based on the CEO’s decision. The staff’s work is meaningful
 229 only if the CEO’s decision is correct. Therefore, we often intuitively assume that an enterprise’s
 230 leader is paid more and contributes more. Based on this, we transform the roles of the CEO and staff
 231 into employees with different contributions to the enterprise. In the sequential MARL framework, we
 232 first assign credit to the agent with a higher contribution to the team.

233 The attribution method is a powerful way to determine the influence of input features’ each component
 234 on the network output value [2]. Among them, integrated gradients [25] leverages path integral to
 235 aggregate gradients along the inputs that fall on the lines between the baseline and the input, which
 236 is a natural tool for measuring each agent’s contribution. QPD [35] utilizes the integrated gradient
 237 attribution technique to decompose shared rewards along trajectory paths, revealing how much each
 238 agent’s observation and action contributes to the global Q value. However, it remains unclear whether
 239 individual Q value should be linearly correlated to or approximated by the agent’s contribution, as in
 240 the case of QPD. The proper connection between agents’ contributions and their individual Q values
 241 in a cooperative team is worth well studied for the community.

242 Here we avoid detailed analysis on the relationship between agents’ contributions and their individual
 243 rewards. Instead, we use integrated gradients to measure agents’ contributions to the state transition
 244 and adjust the credit assignment sequence based on their contributions. In particular, we estimate
 245 agent a ’s contribution c^a in the trajectory path $\tau_{t_1}^{t_2}$ from time t_1 to t_2 based on its policy vector π^a :

$$c^a = \sum_{x_j \in \pi^a} \text{PathIG}_j^{\tau_{t_1}^{t_2}}(\pi^a), \quad (14)$$

246 where x_j is j -th dimension of the policy vector π^a . The computation for PathIG is shown in Equ.(1).
 247 We compute each agent’s contribution c to the state transition from s_{t_1} to s_{t_2} and analyze the agent
 248 with higher c first. We further study the adjustment frequency and its effectiveness in Section 4.2

249 4 Experiments and Analysis

250 4.1 Experimental Setup

251 We consider a challenging set of cooperative StarCraft II maps from the SMAC benchmark [23]
 252 classified as *Easy*, *Hard*, and *Super Hard* scenarios according to the baseline algorithms’ performance.
 253 The inherent differences among various methods and their training procedure (e.g., on/off-policy
 254 learning for value-based/policy-based methods) bring difficulties when comparing methods in a
 255 reasonably fair manner without introducing additional components (e.g., importance sampling [13, 33]
 256 for off-policy methods). To attribute any poor performance of policy-based methods to potential
 257 algorithmic limitations or poor training conditions (in particular, high variance due to small batch
 258 sizes or insufficient gradient steps), we follow [5, 39], training all methods with 32 parallel runners
 259 to generate trajectories and using batches of 32 episodes. We evaluate each method every 320K
 260 steps with 32 episodes and report the 1st, median, and 3rd quartile win rates across 5 random seeds.
 261 Detailed information about the scenarios and the experimental setup is shown in the Appendix.

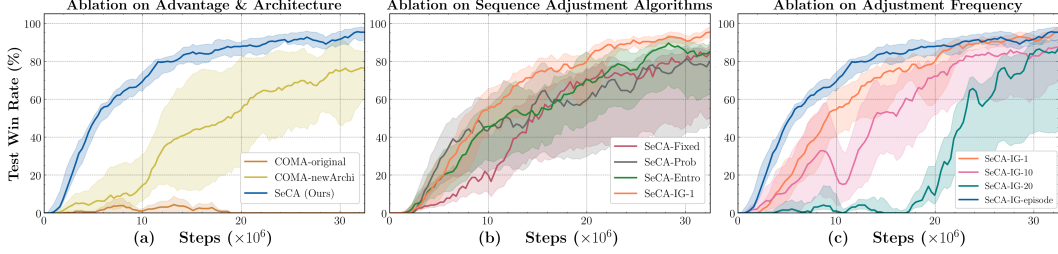


Figure 4: Ablations for SeCA’s key elements on scenario MMM2 (*Super Hard*). (a) investigates the effects of our sequential advantage and network architecture. (b) validates our sequence adjustment through integrated gradients. (c) shows the test win percentage with various adjustment frequencies.

262 4.2 Ablation Studies

263 We first carry out ablation experiments on a *Super Hard* map MMM2 to validate key elements of SeCA.

264 **Proposed Advantage and Architecture.** In Section 3.3, we compare our sequential advantage
 265 with COMA’s in two simple multi-agent particle environments and show our superiority in Figure 2.
 266 Afterward, we introduce a f_ϕ approximation and a corresponding network architecture. Here we apply
 267 the same approximation and architecture for COMA’s counterfactual advantage (COMA-newArchi)
 268 and compare it with the original COMA and our method SeCA to show the effects of our advantage
 269 function, approximation, and network architecture. The result is illustrated in Figure 4(a). COMA
 270 performs poorly on this Super Hard map but acquires significant improvement with our approximation
 271 and architecture. Our sequential advantage further accelerates and stabilizes the training.

272 **Sequence Adjustment Algorithm.** SeCA’s credit assignment sequence is dynamic. We compare
 273 our method with some intuitive adjustments to validate its effects. One could first evaluate agents
 274 with higher current-action probability (SeCA-Prob) or lower policy entropy (SeCA-Entro), as these
 275 agents are more confident in their acts, and we can assess other agents based on their behaviors. Since
 276 SeCA-Prob and Entro get a new order at each step, to be fair, we set the path length in Equ.(14) to one,
 277 i.e., consider agents’ contributions based on the transition from s_t to s_{t+1} (SeCA-IG-1). Figure 4(b)
 278 illustrates that SeCA-Prob and Entro learn better than the fixed method (SeCA-Fixed), but Prob has a
 279 larger variance than Entro. Fixed is better than expected, which we believe is because that the fixed
 280 sequence acquires adequate training. Our integrated-gradients-adjustment performs the best in win
 281 rates and stability, and the others have inferior performance and incredibly high variance.

282 **Sequence Adjustment Frequency.** We next consider how the sequence adjustment frequency in
 283 SeCA-IG affects the performance. Except per step adjustment (i.e., SeCA-IG-1), one could also
 284 update the sequence after a stage or an episode. If we change the credit assignment order for every
 285 episode during training (SeCA-IG-episode), then $\tau_{t_1}^{t_2}$ in Equ.(14) represents a whole episode. As for
 286 stage adjustment, it is hard to define a stage in these tasks, and the stage length varies in diverse maps.
 287 Here we set stage length to 10 and 20, respectively denoted as SeCA-IG-10 and SeCA-IG-20. As the
 288 results in Figure 4(c) show, IG-1 and IG-episode have similar final win rates. However, IG-episode
 289 converges more quickly with smaller variance. The reason for IG-10(20)’s mediocre performance
 290 and high variance may be because the stage length needs to be dynamically adjusted. Inappropriate
 291 adjustment frequency fails to adapt to the stage changes in the task and causes insufficient training
 292 for each sequence. We utilized SeCA-IG-episode in other experiments and will investigate dynamic
 293 stage learning in the future to improve stage adjustment.

294 4.3 Comparisons with State-of-the-arts

295 We compare SeCA with some competitive algorithms, including the representative explicit credit
 296 assignment method COMA, the policy-based implicit method LICA, the common-used baseline
 297 QMIX and QTRAN. Methods are evaluated on 6 scenarios, including 2 *Easy* ones (2s3z, 1c3s5z),
 298 2 *Hard* ones (2c_vs_64zg, 3s_vs_5z), and 2 *Super Hard* ones (MMM2, 3s5z_vs_3s6z). We train
 299 all methods for 32 million steps in *Easy* maps and 64 million steps in *Hard* and *Super Hard* maps.
 300 These scenarios involve homogeneous and heterogeneous teams, symmetric and asymmetric battles,
 301 allowing a holistic study on all methods. Our experiments are based on the latest PyMARL [23]

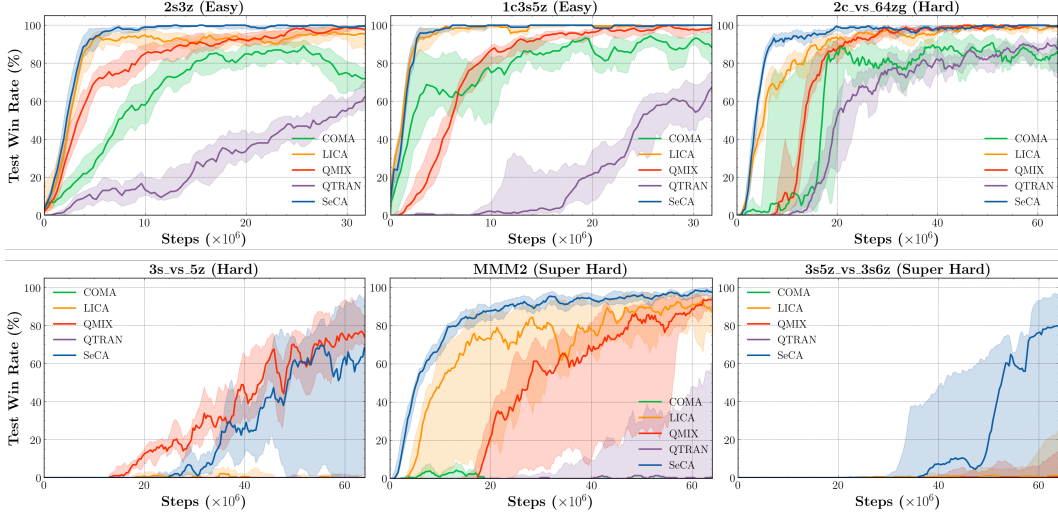


Figure 5: The comparison of SeCA against various baseline algorithms on six SMAC maps.

302 utilizing SC2.4.10. Performance is not always comparable between versions, so the results may be
 303 subtly different from the original papers.

304 As we can see in Figure 5, SeCA demonstrates its robustness by achieving good performances in
 305 scenarios with various characteristics. All methods except COMA and QTRAN solve two *Easy*
 306 scenarios, and SeCA performs better in convergence speed and stability. SeCA’s advantage is further
 307 extended in the *Hard* map *2c_vs_64zg*, and it converges significantly faster than other methods.
 308 Although classified only as *Hard*, *3s_vs_5z* invalidates most algorithms except QMIX and SeCA, as
 309 Stalkers have to learn dispersing and making enemies give chase while maintaining enough distance
 310 ("kiting" technique) in this map. SeCA has a higher variance than QMIX. This is possibly because
 311 the Stalkers’ scattering prioritizes individual performance over cooperation which is more in line
 312 with QMIX’s monotonicity constraint. Nevertheless, SeCA’s performance improvements on the
 313 *Super Hard* scenarios *MMM2* and *3s5z_vs_3s6z* demonstrate the effectiveness of our method. LICA’s
 314 performance in *3s5z_vs_3s6z* here is different from the original paper, as the original results for
 315 this map are obtained by using a different entropy coefficient, which is explained in its open-source
 316 implementation.¹ This parameter tuning is unfair when comparing methods, so all experiments in this
 317 paper use the fixed entropy coefficient. We also visualize the learned sequences in different battles of
 318 *3s_vs_5z* to provide insights into our sequence adjustment in the Appendix.

319 We are supposed to compare our method with QPD that also utilizes integrated gradients to show
 320 our improvement. However, QPD modifies the original SMAC environment to acquire additional
 321 information for policy training, which is mentioned in its open-source implementation.² Therefore, it
 322 is unfair to compare QPD’s learning curves in the modified environment with other methods, and
 323 QPD’s authors did not provide methods’ learning curves comparison in the original paper. We follow
 324 them, providing a win rate table in the Appendix to show our superiority over QPD.

325 5 Conclusions and Future Work

326 This paper presents SeCA, a cooperative MARL framework with sequential credit assignment.
 327 SeCA computes counterfactual advantage functions to evaluate each agent based on the actions of
 328 the preceding agents under a specific sequence. The sequence is adjusted dynamically according
 329 to agents’ contributions to the team deduced by integrated gradients. SeCA accelerates policy
 330 convergence and improves the final performance over existing recognized methods in practice. In the
 331 future, we will further investigate stage learning in an episode and adjust the sequence per stage to
 332 improve SeCA and achieve adaptive cooperation in various task situations.

¹<https://github.com/mzho7212/LICA>

²<https://github.com/QPD-NeurIPS2019/QPD>

333 References

- 334 [1] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding
335 of gradient-based attribution methods for deep neural networks. In *International Conference on*
336 *Learning Representations*, 2018.
- 337 [2] Guillem Brasó Andilla. Attribution methods for deep convolutional networks.
- 338 [3] Yongcan Cao, Wenwu Yu, Wei Ren, and Guanrong Chen. An overview of recent progress in
339 the study of distributed multi-agent coordination. *IEEE Transactions on Industrial Informatics*,
340 9(1):427–438, 2012.
- 341 [4] Yu-han Chang, Tracey Ho, and Leslie Kaelbling. All learning is local: Multi-agent learning in
342 global reward games. In *Advances in Neural Information Processing Systems*, pages 808–814,
343 2004.
- 344 [5] Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. Liir: Learning individual
345 intrinsic reward in multi-agent reinforcement learning. In *Advances in Neural Information*
346 *Processing Systems*, pages 4403–4414, 2019.
- 347 [6] Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.
348 Counterfactual multi-agent policy gradients. In *AAAI Conference on Artificial Intelligence*,
349 pages 2974–2982, 2018.
- 350 [7] Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control
351 using deep reinforcement learning. In *International Conference on Autonomous Agents and*
352 *Multiagent Systems*, pages 66–83, 2017.
- 353 [8] Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. A survey and critique of multiagent
354 deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33(6):750–797,
355 2019.
- 356 [9] Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. Guided deep reinforcement
357 learning for swarm systems. In *AAMAS Autonomous Robots and Multirobot Systems (ARMS)*
358 *Workshop*, 2017.
- 359 [10] Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for
360 decentralized planning. *Neurocomputing*, 190:82–94, 2016.
- 361 [11] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent
362 actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information*
363 *Processing Systems*, pages 6382–6393, 2017.
- 364 [12] Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-
365 agent variational exploration. In *Advances in Neural Information Processing Systems*, pages
366 7611–7622, 2019.
- 367 [13] A Rupam Mahmood, Hado van Hasselt, and Richard S Sutton. Weighted importance sampling
368 for off-policy learning with linear function approximation. In *Advances in Neural Information*
369 *Processing Systems*, pages 3014–3022, 2014.
- 370 [14] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
371 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
372 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 373 [15] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and
374 understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.
- 375 [16] Frans A Oliehoek and Christopher Amato. *A concise introduction to decentralized POMDPs*.
376 Springer, 2016.
- 377 [17] Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value
378 functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353,
379 2008.

- 380 [18] Afshin OroojlooyJadid and Davood Hajinezhad. A review of cooperative multi-agent deep
381 reinforcement learning. *arXiv preprint arXiv:1908.03963*, 2019.
- 382 [19] Gregory Palmer, Karl Tuyls, Daan Bloembergen, and Rahul Savani. Lenient multi-agent deep
383 reinforcement learning. In *International Conference on Autonomous Agents and MultiAgent*
384 *Systems*, pages 443–451, 2018.
- 385 [20] Scott Proper and Kagan Tumer. Modeling difference rewards for multiagent learning. In
386 *International Conference on Autonomous Agents and Multi-Agent Systems*, pages 1397–1398,
387 2012.
- 388 [21] Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding
389 monotonic value function factorisation. In *Advances in Neural Information Processing Systems*,
390 pages 10199–10210, 2020.
- 391 [22] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster,
392 and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent
393 reinforcement learning. In *International Conference on Machine Learning*, pages 4295–4304,
394 2018.
- 395 [23] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas
396 Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon
397 Whiteson. The starcraft multi-agent challenge. *CoRR*, abs/1902.04043, 2019.
- 398 [24] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran:
399 Learning to factorize with transformation for cooperative multi-agent reinforcement learning.
400 In *International Conference on Machine Learning*, pages 5887–5896, 2019.
- 401 [25] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
402 *International Conference on Machine Learning*, pages 3319–3328, 2017.
- 403 [26] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi,
404 Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-
405 decomposition networks for cooperative multi-agent learning based on team reward. In *In-*
406 *ternational Conference on Autonomous Agents and MultiAgent Systems*, pages 2085–2087,
407 2018.
- 408 [27] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*,
409 3(1):9–44, 1988.
- 410 [28] Kagan Tumer and Adrian Agogino. Distributed agent-based air traffic flow management. In
411 *International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 1–8,
412 2007.
- 413 [29] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling
414 multi-agent q-learning. *arXiv preprint arXiv:2008.01062*, 2020.
- 415 [30] Jianhong Wang, Yuan Zhang, Tae-Kyun Kim, and Yunjie Gu. Shapley q-value: A local reward
416 approach to solve global reward games. In *AAAI Conference on Artificial Intelligence*, pages
417 7285–7292, 2020.
- 418 [31] Tonghan Wang, Heng Dong, Victor Lesser, and Chongjie Zhang. Roma: Multi-agent reinforce-
419 ment learning with emergent roles. In *International Conference on Machine Learning*, pages
420 9876–9886, 2020.
- 421 [32] Tonghan Wang, Tarun Gupta, Anuj Mahajan, Bei Peng, Shimon Whiteson, and Chongjie Zhang.
422 Rode: Learning roles to decompose multi-agent tasks. *arXiv preprint arXiv:2010.01523*, 2020.
- 423 [33] Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Remi Munos, Koray Kavukcuoglu,
424 and Nando de Freitas. Sample efficient actor-critic with experience replay. *arXiv preprint*
425 *arXiv:1611.01224*, 2016.
- 426 [34] David H Wolpert and Kagan Tumer. Optimal payoff functions for members of collectives. In
427 *Modeling complexity in economic and social systems*, pages 355–369. World Scientific, 2002.

- 428 [35] Yaodong Yang, Jianye Hao, Guangyong Chen, Hongyao Tang, Yingfeng Chen, Yujing Hu,
 429 Changjie Fan, and Zhongyu Wei. Q-value path decomposition for deep multiagent reinforcement
 430 learning. In *International Conference on Machine Learning*, pages 10706–10715, 2020.
- 431 [36] Dayong Ye, Minjie Zhang, and Yun Yang. A multi-agent framework for packet routing in
 432 wireless sensor networks. *Sensors*, 15(5):10026–10047, 2015.
- 433 [37] Wang Ying and Sang Dayong. Multi-agent framework for third party logistics in e-commerce.
 434 *Expert Systems with Applications*, 29(2):431–436, 2005.
- 435 [38] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A
 436 selective overview of theories and algorithms. *arXiv preprint arXiv:1911.10635*, 2019.
- 437 [39] Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit
 438 credit assignment for cooperative multi-agent reinforcement learning. In *Advances in Neural*
 439 *Information Processing Systems*, 2020.

440 Checklist

- 441 1. For all authors...
- 442 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 443 contributions and scope? [Yes]
- 444 (b) Did you describe the limitations of your work? [Yes] We discussed it in the experiment
 445 analysis in Section 4.3 and future work in Section 5.
- 446 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 447 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 448 them? [Yes]
- 449 2. If you are including theoretical results...
- 450 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 451 (b) Did you include complete proofs of all theoretical results? [Yes] We provided the proof
 452 of our Claim in the supplemental material.
- 453 3. If you ran experiments...
- 454 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 455 mental results (either in the supplemental material or as a URL)? [Yes] We provided
 456 our code and instructions in the supplemental material.
- 457 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 458 were chosen)? [Yes] We described the training details in the supplemental material.
- 459 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 460 ments multiple times)? [Yes] See Figure 2, 4 and 5.
- 461 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 462 of GPUs, internal cluster, or cloud provider)? [Yes] We described it in the supplemental
 463 material.
- 464 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 465 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 466 (b) Did you mention the license of the assets? [Yes]
- 467 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 468 We provided our code in the supplemental material.
- 469 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 470 using/curating? [N/A]
- 471 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 472 information or offensive content? [N/A]
- 473 5. If you used crowdsourcing or conducted research with human subjects...
- 474 (a) Did you include the full text of instructions given to participants and screenshots, if
 475 applicable? [N/A]

476
477
478
479

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]