

# Structured Diversification Emergence via Reinforced Organization Control and Hierarchical Consensus Learning

Wenhao Li

East China Normal University  
Shanghai, China  
52194501026@stu.ecnu.edu.cn

Xiangfeng Wang\*

East China Normal University and  
SRIAS, Shanghai, China  
xfwang@cs.ecnu.edu.cn

Bo Jin\*

East China Normal University and  
SRIAS, Shanghai, China  
bjin@cs.ecnu.edu.cn

Junjie Sheng

East China Normal University  
Shanghai, China  
52194501003@stu.ecnu.edu.cn

Yun Hua

East China Normal University  
Shanghai, China  
52194501002@stu.ecnu.edu.cn

Hongyuan Zha

School of Data Science and AIRS, The  
Chinese University of Hong Kong,  
Shenzhen, China. zhahy@cuhk.edu.cn

## ABSTRACT

When solving a complex task, humans will spontaneously form teams and to complete different parts of the whole task, respectively. Meanwhile, the cooperation between teammates will improve efficiency. However, for current cooperative MARL methods, the cooperation team is constructed through either heuristics or end-to-end blackbox optimization. In order to improve the efficiency of cooperation and exploration, we propose a structured diversification emergence MARL framework named ROCHICO based on reinforced organization control and hierarchical consensus learning. ROCHICO first learns an adaptive grouping policy through the organization control module, which is established by independent multi-agent reinforcement learning. Further, the hierarchical consensus module based on the hierarchical intentions with consensus constraint is introduced after team formation. Simultaneously, utilizing the hierarchical consensus module and a self-supervised intrinsic reward enhanced decision module, the proposed cooperative MARL algorithm ROCHICO can output the final diversified multi-agent cooperative policy. All three modules are organically combined to promote the structured diversification emergence. Comparative experiments on four large-scale cooperation tasks show that ROCHICO is significantly better than the current SOTA algorithms in terms of exploration efficiency and cooperation strength.

## KEYWORDS

Cooperative MARL; Diversification; Organization Control

### ACM Reference Format:

Wenhao Li, Xiangfeng Wang, Bo Jin, Junjie Sheng, Yun Hua, and Hongyuan Zha. 2021. Structured Diversification Emergence via Reinforced Organization Control and Hierarchical Consensus Learning. In *Proc. of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*, Online, May 3–7, 2021, IFAAMAS, 15 pages.

## 1 INTRODUCTION

Multi-agent reinforcement learning (MARL) has been widely used and achieve fantastic performance in many application fields, like

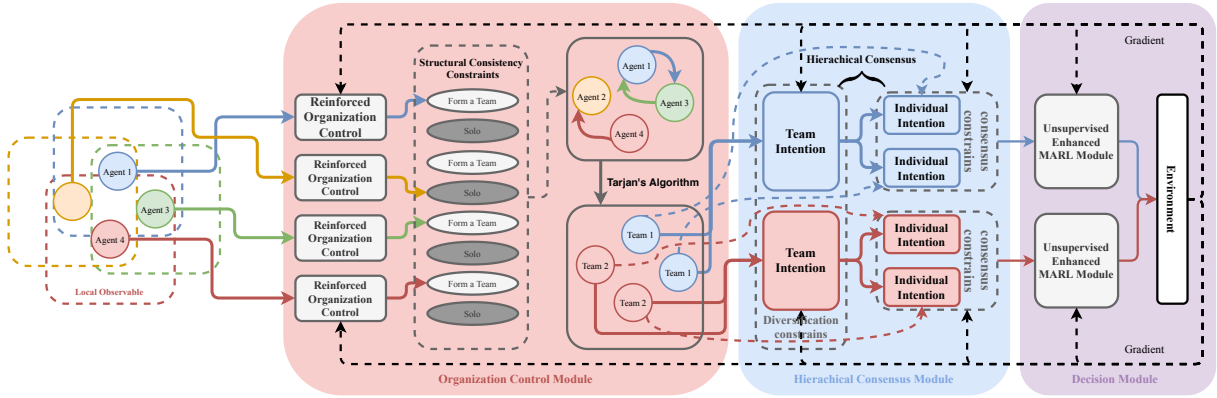
\*Corresponding authors

multiplayer games [21, 31], swarm robot control [27] etc. Most of the current MARL algorithms follow the centralized training and decentralized execution (CTDE, [29]) framework. In the centralized training phase, a decentralized policy needs to be learned for each agent through sharing local observations, parameters, or gradients among agents. However, these CTDE-based MARL algorithms have to consider each agent as an independent individual during the training procedure. Although these individuals can transmit information explicitly or implicitly to achieve collaboration, most of them are usually learned through an end-to-end blackbox scheme, which raises the difficulty to obtain meaningful communication protocols [45]. This makes it difficult for multi-agents to explore and collaborate effectively.

When humans perform tasks in an unknown environment, diverse teams instead of individuals, are usually used as the basic unit to make up for limited individual abilities. As for the nonstationarity of the external environment and the difficulty of the related task, the team can be restructured accordingly. In terms of MARL, this means that the agents must have the ability to dynamically team-up. To improve the efficiency of exploration in an unknown environment, the behavior of different teams need to be sufficiently diverse; Further considering the capacity limitations of a single agent, the agents from the same team need to cooperate closely to improve the efficiency of task completion. In this paper, we call this inter-team diversification and intra-team cooperation ability the *structured diversification emergence*.

The core assumption is that the behavior of an agent is determined by both the team goal and its perception of the environment, while diverse team goals and consistent environmental perception can lead to structured diversification emergence. Designing an efficient algorithm for the agent to learn structured diversification emergence ability, we need to answer the following important questions: 1) **Adaptive organization control**: how do the agents form teams spontaneously, and the team composition can dynamically change to adapt to the external environment? 2) **Structured diversification**: How to maintain diversity in the behavior of different teams and to form tight cooperation between the agents within the team? 3) **Behavior emergence**: How to combine the above two processes organically? Let us conduct an in-depth analysis.

*Adaptive organization control*. In a cooperative task, before taking actions, an individual will first assess whether her ability is sufficient to complete the task based on her observations of the



**Figure 1: The ROCHICO architecture. The algorithm is divided into three modules: organization control, hierarchical consensus, and decision. Each agent first makes a decision whether to team up with neighbors through the organization control module based on local observations. Next, the team intention generator and individual intention generator, trained based on the hierarchical consensus constraints, generate team intentions and individual intentions according to the teaming results. Finally, the decision-making module generates structured diversify policies based on team intentions and individual intentions.**

environment. If not, she will seek help from other individuals. Besides, since the environment is locally observable, a team would be more powerful by fetching the information from multiple individuals. Therefore, some existing works introduce teams in the MARL algorithms. However, these teams are constructed either through heuristics [51, 52] or through end-to-end blackbox optimization [4, 16, 39]. Compared with these methods to passively put individuals into the team, a more reasonable way is to let the individuals actively decide who to team up with.

*Structured diversification.* After the teams are formed, the most direct way to measure the diversity of team behavior (or the goal) is to compare the differences in policies between all agents in different teams. However, the number of agents in different teams are various, and the agent’s policy is usually a stochastic conditional probability distribution. As a result, for large state-action space, directly comparing the difference among these conditional distributions (or agent’s policies) is intractable. From another perspective, the behavior of agents in the same team may be generated by a potential team intention. If the intentions of different teams could be mapped into the same latent space, these intentions can be easily compared. Within each team, the team intention guides the behaviors of agents, and each agent’s perception of the surrounding environment can also influence agents’ behaviors. We can make a reasonable assumption: if two agents both have the same perception of the environment and the same goal, their behaviors should be closely coordinated. At the same time, an effective perception should be able to reconstruct the surrounding environment. This motivates us to impose consensus constraints on the agents’ perception of the environment, thus we can achieve tight collaboration within the team. A Similar idea can be found in Mao et al. [25], but it only imposes perceptual consensus constraints on the agent and other agents within its fixed neighborhood to encourage collaboration, which can’t leads to agents’ diverse behaviors.

In this paper, we propose a structured diversification emergence MARL algorithm, so-called ROCHICO, based on reinforced organization control and hierarchical consensus learning. As shown in

Figure 1, ROCHICO consists of three modules: organization control, hierarchical consensus and decision. First, the organization control module models the multi-agent organization control problem as a partially observable stochastic game (POSG), and introduces independent MARL to obtain a dynamic and autonomous teaming strategy. Second, the hierarchical consensus module obtains team intentions and individual intentions through contrastive learning and unsupervised learning. Structured diversification emerges by imposing hierarchical consensus constraints on hierarchical intentions. Finally, the decision module outputs the diversify cooperative policies based on results of reinforced organization control and hierarchical consensus learning to the environment, and feeds back external rewards to the other two modules, and decision module itself, so that all modules can be combined organically.

Our contributions mainly consist of the following folds: 1) We model the organization control of the multi-agent as an independent learning task, which enables the agent to autonomously and adaptively team up based on environmental feedback. 2) We impose hierarchical consensus constraints on hierarchical intentions obtained by the novel introduced contrastive learning and unsupervised learning auxiliary tasks to encourage the structured diversification emergence. 3) Performance experiments on various large-scale cooperative tasks show that ROCHICO is significantly better than the current SOTA algorithms in terms of exploration and cooperation efficiency.

## 2 RELATED WORKS

### 2.1 Organization Control Mechanism

Organization control is defined as a mechanism or a process that enables a system to change its organization without explicit command during execution time [5]. The most relevant existing self-organization mechanisms to our work can be summarized as *task allocation* [6, 24, 33], *relation adaption* [9, 10, 19] and *coalition formation* [2, 26, 48]. For the reason that these methods are not RL-based and we will not expand here, more details can be found in Ye et al.

[49]. The task allocation refers to the agent actively allocates the task(s) to other agents because it cannot finish it by itself, which is different from classical task allocation in RL. Further, the key difference between our work with relation adaption or coalition formation methods, is that our organization control mechanism is obtained through independent learning, rather than based on heuristic techniques or on large communication driven negotiation.

In recent years, there are a few works introduce organization control into reinforcement learning. Abdallah and Lesser [1] uses RL to design the task allocation mechanism and transfer the learned knowledge across the different steps of organization control with heuristics mechanism. Zhang et al. [51, 52] integrate organization control into MARL to improve the convergence speed, which suffers large communication overhead because of the negotiation.

## 2.2 Behavior Diversification

Many cooperative tasks require agents to take different behaviors to achieve higher degree of task completion. Behavior diversification can be handcrafted or emerged through multi-agent interaction. Handcrafted diversification is widely studied as task allocation or role assignment. Heuristics mechanisms [3, 24, 36, 41] assign a specific task or a pre-defined role to each agent based on its goal (capability, visibility) or by searching. Shu and Tian [40] establishes a manager to assign suitable sub-tasks to rule-based workers with different preferences and skills. All these methods require that the sub-tasks and roles are pre-defined, while the worker agents are rule-based at the same time.

Recently, the emergent diversification was introduced to single-agent RL [7, 11–13] with the purpose to learn reusable diverse skills in complex and transferable tasks. In MARL, McKee et al. [28] introduces diversity into heterogeneous agents to learn more generalized policies for solving social dilemmas. Wang et al. [46] learns a role embedding encoder and a role decoder simultaneously. However, no mechanism guarantees the role decoder can generate different parameters and generate diversity policies accordingly, taking as input different role embeddings. Jiang and Lu [17] establishes an intrinsic reward for each agent through a well-trained probabilistic classifier. The intrinsic reward makes the agents more identifiable and promote diversification emergence. Based on Eysenbach et al. [7], learning low-level skills for each agent in hierarchical MARL is considered in Lee et al. [20], Yang et al. [47]. The high-level policy can utilize coordinated low-level diverse skills, but the high-level policy does not consider diversity.

## 3 PRELIMINARIES

**Cooperative POSGs.** POSG [14] is denoted as a seven-tuple based on the stochastic game (or Markov game) as

$$\langle \mathcal{X}, \mathcal{S}, \{\mathcal{A}^i\}_{i=1}^n, \{\mathcal{O}^i\}_{i=1}^n, \mathcal{P}, \mathcal{E}, \{\mathcal{R}^i\}_{i=1}^n \rangle,$$

where  $n$  denotes agents total number;  $\mathcal{X}$  represents the agent space;  $\mathcal{S}$  contains a finite set of states;  $\mathcal{A}^i$ ,  $\mathcal{O}^i$  and denote a finite action set and a finite observation set of agent  $i$  respectively;  $\mathcal{A} = \mathcal{A}^1 \times \mathcal{A}^2 \times \dots \times \mathcal{A}^n$  is the finite set of joint actions;  $\mathcal{P}(s'|s, \mathbf{a})$  denotes the Markovian state transition probability function;  $\mathcal{O} = \mathcal{O}^1 \times \mathcal{O}^2 \times \dots \times \mathcal{O}^n$  is the finite set of joint observations;  $\mathcal{E}(o|s)$  is the Markovian observation emission probability function;  $\mathcal{R}^i : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$

denotes the reward function of agent  $i$ . The game in POSG unfolds over a finite or infinite sequence of stages (or timesteps), where the number of stages is called *horizon*. In this paper, we consider the finite horizon case. The objective for each agent is to maximize the expected cumulative reward received during the game. For a cooperative POSG, we quote the definition in Song et al. [42],

$$\forall x \in \mathcal{X}, \forall x' \in \mathcal{X} \setminus \{x\}, \forall \pi_x \in \Pi_x, \forall \pi_{x'} \in \Pi_{x'}, \frac{\partial \mathcal{R}^{x'}}{\partial \mathcal{R}^x} \geq 0,$$

where  $x$  and  $x'$  are a pair of agents in agent space  $\mathcal{X}$ ;  $\pi_x$  and  $\pi_{x'}$  are the corresponding policies in the policy space  $\Pi_x$  and  $\Pi_{x'}$  respectively. Intuitively, this definition means that there is no conflict of interest for any pair of agents.

**QMIX.** The QMIX [35] algorithm, which is the follow-up work of the VDN [43], is the current SOTA of cooperative MARL. QMIX claims that the the  $Q$ -value functions before ( $Q_{tot}$ ) and after ( $Q_a$ ) decomposition should satisfy the following constraints:

$$\frac{\partial Q_{tot}}{\partial Q_a} \geq 0, \forall a \in A.$$

QMIX employs a hypernetwork-based mixing network to promote the two  $Q$ -function satisfying the above condition. Because of the non-linear mixing network, QMIX can outperform VDN.

## 4 ALGORITHMS

The overall framework of the proposed algorithm is shown in Figure 1, which can be divided into three modules. The organization control module receives local observations of all agents and makes adaptive teaming decisions use the traditional graph theory algorithm [44]. The hierarchical consensus module will then generate the team intention and the individual intention based on the obtained teaming results. The hierarchical consensus constraints are established on the above hierarchical intentions to promote the structured diversification emergence. Finally, the decision module outputs the structured diversification policies through a cooperative MARL algorithm. However, it should be noted that the organization control process is not differentiable, so we cannot train the overall model by the end-to-end scheme. We draw on the communication MARL algorithms [8, 18, 39], while passing the external rewards to the organization control module and decision module separately, so as to realize the joint training.

### 4.1 Organization Control Module

If we consider all agents as a graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ , each agent is a node<sup>1</sup>  $v \in \mathcal{V}$  in the graph  $\mathcal{G}$ . The edge  $e \in \mathcal{E}$  indicates whether the two agents connected by  $e$  belong to the same team. The main purpose of the organization control module is to determine the connections (edges) between agents. Then, we can naturally view the connected components as teams, while searching for connected components can be done efficiently by traditional graph theory algorithms [44].

In the beginning, if the problem in organization control module is modeled as a single-agent RL problem,  $\mathcal{O}(n^2)$  edges need to be determined if the graph contains  $n$  nodes, and the size of the action space is  $\mathcal{O}(2^{n^2})$ . This makes it impossible to scale to a larger multi-agent scenario (e.g., tens of agents). Therefore, we model

<sup>1</sup>Except for explicit emphasis, we will no longer distinguish between the two terms node and agent below.

the organization control problem as a MARL problem, while each node is considered as an agent. However, if each agent needs to determine its connections with all other nodes, the action space is still very large and be  $O(2^n)$ . Therefore, inspired by other team-based MARL algorithms, like [15, 16], we only consider the closest  $m$  other agents. This idea is intuitive but effective, because it is reasonable to team up the agents according to the adjacency of the spatial position for most tasks.

Formally, the organization control problem can be modeled as a cooperative POSG, which is denoted as:

$$\mathcal{M}_u := \left\langle \mathcal{X}_u, \mathcal{S}_u, \{\mathcal{A}_u^i\}_{i=1}^n, \{\mathcal{O}_u^i\}_{i=1}^n, \mathcal{P}_u, \mathcal{E}_u, \{\mathcal{R}_u^i\}_{i=1}^n \right\rangle.$$

We set the action of agent  $i$ ,  $a_u^i \in \mathcal{A}_u^i$ , as a  $m$ -dimension binary vector, denoting the connection action to  $m$ -nearest agents  $\{x^j \mid j \in \mathcal{N}_m(i)\}$  ( $m \ll n$ ) according to its local observation  $o_u^i$ , where the subscript  $u$  stands for 'unorganized'. This kind of MARL problem will suffer that agent  $i$  decides to connect with agent  $j$  but agent  $j$  does not want to connect with agent  $i$  or vice versa, which means:

$$a_u^i[j] \neq a_u^j[i], \quad \text{with } i \in \mathcal{N}_m(j), j \in \mathcal{N}_m(i).$$

Considering the positive effect of prosociality on promoting the cooperation of agents [32], we use the weakly connected graph  $\mathcal{G}_u = (\mathcal{V}_u, \mathcal{E}_u)$  to form teams, which is established by converting directed edges into undirected edges, i.e.,

$$e(i, j) = a_u^i[j] \vee a_u^j[i], \quad \text{with } i \in \mathcal{N}_m(j), j \in \mathcal{N}_m(i),$$

where  $\vee$  denotes "or" operation. Finally, the Tarjan's algorithm is employed for searching weakly connected components with worst-case time complexity as  $O(|\mathcal{V}_u| + |\mathcal{E}_u|)$  [44].

In addition to the external rewards  $\{r_e^i\}_{i=1}^n$  from the environment, we also introduce additional intrinsic rewards to train the agents. Specifically, in order to strengthen training stability without causing excessive fluctuation on the graph structure, the novel *structural consistency intrinsic reward*  $r_u^i$  for each agent  $i$  is defined by:

$$r_u^i = \frac{1}{|\mathcal{N}_m(i)|} \cdot \text{GED} \left( \mathcal{G}_u(\mathcal{N}_m(i) \vee i), \mathcal{G}'_u(\mathcal{N}_m(i) \vee i) \right),$$

where  $\text{GED}(\cdot, \cdot)$  represents the *graph edit distance* [37].  $\mathcal{G}_u(\mathcal{N}_m(i) \vee i)$  and  $\mathcal{G}'_u(\mathcal{N}_m(i) \vee i)$  represent the sub-graph only contains node  $i$  and its  $m$ -nearest neighbors before and after take action  $a_u^i$  respectively. The total reward for each agent  $i$  is:

$$r_{u+}^i = r_e^i + \alpha_u r_u^i,$$

where  $\alpha_u$  indicates the strength of constraint for structural consistency. For the organization control problem, the goal is to maximize the summation of all agents' expected accumulated rewards, which can be denoted as follows:

$$\max \mathcal{J}_u = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tau_u} [r_{u+}^i(\tau_u)].$$

Further to raise training efficiency, we use independent learning combined with DQN and parameter sharing mechanism. [30] has shown the significant performance of independent learning on multi-agent tasks. Specifically, we can minimize following TD(0) error for each agent  $i$ , i.e.,

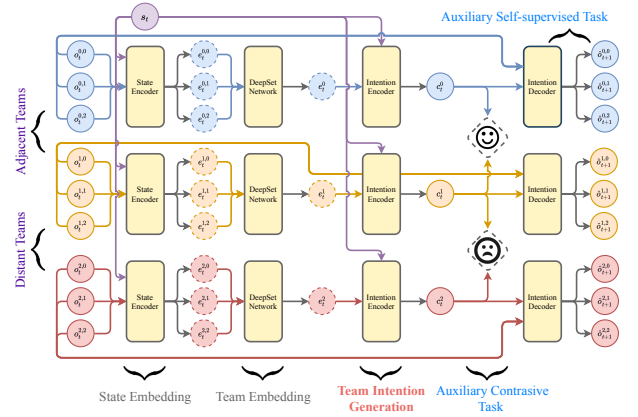
$$\mathcal{L}_u^i(\theta_u) = \mathbb{E}_{(o_u^i, a_u^i, r_{u+}^i, o_u^{i'}) \sim D} \left[ \left( Q_{\theta_u}(o_u^i, a_u^i) - y \right)^2 \right].$$

where  $y = r_{u+}^i + \gamma \max_{a_u^{i'}} Q_{\bar{\theta}_u}(o_u^{i'}, a_u^{i'})$  and  $\theta_u, \bar{\theta}_u$  parameterize  $Q$  function and target  $Q$  function separately. Finally, the overall objective function is:

$$\min \mathcal{L}_u^Q(\theta_u) = \sum_i L_u^i(\theta_u).$$

## 4.2 Hierarchical Consensus Module

The goal of the hierarchical consensus module is to achieve efficient multi-agent exploration and cooperation, that is, structured diversification emergence. To this end, we put forward the concepts of hierarchical intentions, i.e., team intentions and individual intentions. Then the structured diversification emergence can be achieved through contrastive learning and hierarchical consensus learning. The hierarchical consensus module is composed of two sub-modules: team intention generator (Figure 2) and individual intention generator (Figure 3).



**Figure 2: The network structure of the team intention generator in the training phase. Two teams are close to each other and one team is far apart. The team intention is based on the representations of all agents in the team, which are aggregated through a DeepSet network, and further is calculated through the intention encoder. The loss function consists of both contrastive loss and self-supervised loss.**

**4.2.1 Team Intention Generator.** The team intention should have the following characteristics: 1) the team intention must reflect the team behavior (or goal), so that it should be generated based on the joint observation of all agents within the team; 2) to improve the exploration efficiency of the team in an unknown environment, the team intentions among different teams must be diverse; 3) the team intention can reflect the team behaviors in a short period time, so we can predict the future observation of the agents within the team based on the team intention. In addition to the second one, the agents in each team make decisions independently, so they must access the local information of the other agents in the same team to achieve diverse behaviors. Considering the difficulty of communication learning [22], we use the global state of environment additionally to assist the generation of team intention, similar as Rashid et al. [34]. The global state is only used in the training phase since we use a CTDE framework.

Formally, in order to generate the team intention for each team, the state encoder  $f_\mu(\cdot)$  parameterized by  $\mu$  receives the joint observation  $\mathbf{o}_t^k = (o_t^{k,1}, o_t^{k,2}, \dots, o_t^{k,n_k})$  of all  $n_k$  agents in team  $k$ , and together with the global state  $s_t$  (such as the minimap of the environment) at timestep  $t$ . The generated agent state embeddings of each team  $(e_t^{k,1}, e_t^{k,2}, \dots, e_t^{k,n_k})$  are feed into a DeepSet network [50], i.e.,  $f_\nu(\cdot)$  parameterized by  $\nu$  to generate the team embedding  $e_t^k$ :

$$e_t^k = f_\nu \left( f_\mu \left( o_t^{k,1} \right), f_\mu \left( o_t^{k,2} \right), \dots, f_\mu \left( o_t^{k,n_k} \right) \right).$$

Finally, the team intention encoder  $f_\omega(\cdot)$  parameterized by  $\omega$  receives team embedding  $e_t^k$  and global state  $s_t$  again to generate the team intention  $c_t^k$ :

$$c_t^k = f_\omega \left( e_t^k, s_t \right).$$

In order to generate diverse team intentions, we model the team intention generation as a *contrastive learning* problem. First, the average spatial position  $[\bar{x}_t^k, \bar{y}_t^k]$  of the team  $k$  at each timestep  $t$  can be calculated based on the spatial position of the all agents within the team  $\{(x_t^{k,1}, y_t^{k,1}), \dots, (x_t^{k,n_k}, y_t^{k,n_k})\}$ . Combined with the current timestamp,  $t_k$ , a spatiotemporal feature representation  $e_{ts}^k = [\bar{x}_t^k, \bar{y}_t^k, t_k]$  of the team  $k$  can be obtained. The Euclidean distance in the spatiotemporal space is used to measure the distance between team  $k$  and  $l$ , i.e.,

$$d(k, l) = \|e_{ts}^k - e_{ts}^l\|_2^2.$$

A simple version directed graph  $\mathcal{G}_\ell(\mathcal{V}_\ell, \mathcal{E}_\ell)$  can be established in team level, where  $\mathcal{V}_\ell$  denotes the node set of  $\mathcal{G}_\ell$  and each node  $v^k \in \mathcal{V}_\ell$  represents a team  $k$ . The  $\mathcal{E}_\ell$  denotes the edge set and the two teams connected by edge  $e^k \in \mathcal{E}_\ell$  have similar intentions. For each edge in edge set  $\mathcal{E}_\ell$  of  $\mathcal{G}_\ell$ , we have

$$e(k, l) := \left\{ \mathbf{1}[d(k, l) = \min_v d(k, v)] \right\} \vee \left\{ \mathbf{1}[d(l, k) = \min_v d(l, v)] \right\}.$$

There is an edge connection between  $k, l$  two nodes iff anyone is the nearest neighbor of the other. Similar to the organization control module, we use Tarjan's algorithm to find all the weakly connected components in the directed graph  $\mathcal{G}_\ell$ . Teams that belong to the same weakly connected component will be assigned the same label. The team intention  $c_t^k$  is set to be the feature and the Tarhan's algorithm label  $y_t^k$  is set to be the label, which constructs a supervised training set  $\{(c_t^k, y_t^k)\}$ . After combining with the triplet loss [38], we could construct a contrastive learning problem by minimize following objective for each team  $k$ :

$$\mathcal{L}_\ell^k(\mu, \nu, \omega) = \mathbb{E} \left[ \max \left( 0, \|c_t^k - c_t^u\|_2^2 - \|c_t^k - c_t^v\|_2^2 + m \right) \right],$$

where  $m$  is a margin parameter,  $k, u$  share the same label ( $y_t^k = y_t^u$ ) and  $k, v$  have different labels ( $y_t^k \neq y_t^v$ ).  $u$  and  $v$  are sampled from the weak connected component same as  $k$  and different from  $k$  respectively.

In addition, in order to enhance the impacts of team intention on the team behavior, besides using the team intention as the input of the following individual intention generator, we introduce another self-supervised task. The intention decoder  $f_\xi(\cdot)$  parameterized by  $\xi$  receives team intention  $c_t^k$  of team  $k$  at current timestep  $t$  as input, and output the prediction of joint observation

$\hat{o}_{t+1}^k = (o_{t+1}^{k,1}, o_{t+1}^{k,2}, \dots, o_{t+1}^{k,n_k})$  at next timestep. Formally, we minimize following regression objective function for each team  $k$  so as to formulate the self-supervised task, i.e.,

$$\mathcal{L}_\ell^k(\xi) = \mathbb{E} \left[ \frac{1}{n_k} \sum_{i=1}^{n_k} \left\| f_\xi(o_t^{k,i}, c_t^k) - o_{t+1}^{k,i} \right\|_2^2 \right].$$

The network structure is shown in Figure 2 and the overall problem of the team intention generator can be formulated as follows

$$\min \mathcal{L}_\ell^{tg}(\mu, \nu, \omega, \xi) = \mathbb{E} \left[ \sum_k \mathcal{L}_\ell^k(\mu, \nu, \omega) + \lambda_{tg} \cdot \mathcal{L}_\ell^k(\xi) \right].$$

**4.2.2 Individual Intention Generator.** With the team intention as the guidance for the diverse behavior between teams, the hierarchical consensus module aims to realize the structured diversification emergence through individual intentions consensus. The key idea to achieve better cooperation within the team is that all agents in the same team should have a consistent cognition of the surrounding environment and the task. The individual intention generation is divided into the following three steps:

- 1). The individual encoder  $g_\phi$  parameterized by  $\phi$  encodes the local observation  $o_t^{k,i}$  (of agent  $i$  in team  $k$  at timestep  $t$ ) into individual embedding  $h_t^{k,i}$  (yellow rounded square in Figure 3(b)), which only contains agent-specific information;
- 2). We regard agents in the same team as nodes in a new fully-connected graph. A GNN  $g_\psi$  parameterized by  $\psi$  then is introduced to further aggregate all individual embeddings  $\mathbf{h}_t^k = (h_t^{k,1}, \dots, h_t^{k,n_k})$  and extracts the individual cognition  $\chi_t^{k,i}$  (blue rounded square in Figure 3(b)) by

$$\chi_t^{k,i} = g_\psi \left( \sum_{j \in k} h_t^{k,j} \right);$$

- 3). The variational encoder  $g_\varphi$  (shown in Figure 3(c)) receives the concatenated feature  $(\chi_t^{k,i}, c_t^k)$  (consist of individual cognition  $\chi_t^{k,i}$  and team intention  $c_t^k$  of team  $k$ ) as input, and output the individual intention  $\zeta_t^{k,i}$  (green rounded square in Figure 3(b)) through<sup>2</sup>

$$\zeta_t^{k,i} = g_\varphi(\chi_t^{k,i}, c_t^k).$$

Then, we impose hierarchical consensus constraints on the generated individual intentions, which leads to optimizing the following function for each agent  $i$  in team  $k$

$$\mathcal{L}_\ell^i(\phi, \psi, \varphi) = \mathbb{E} \left[ \frac{1}{n_k - 1} \sum_{j \neq i} \text{KL} \left[ q_\varphi(\zeta_t^{k,i} | o_t^{k,i}) \| q_\varphi(\zeta_t^{k,j} | o_t^{k,j}) \right] \right].$$

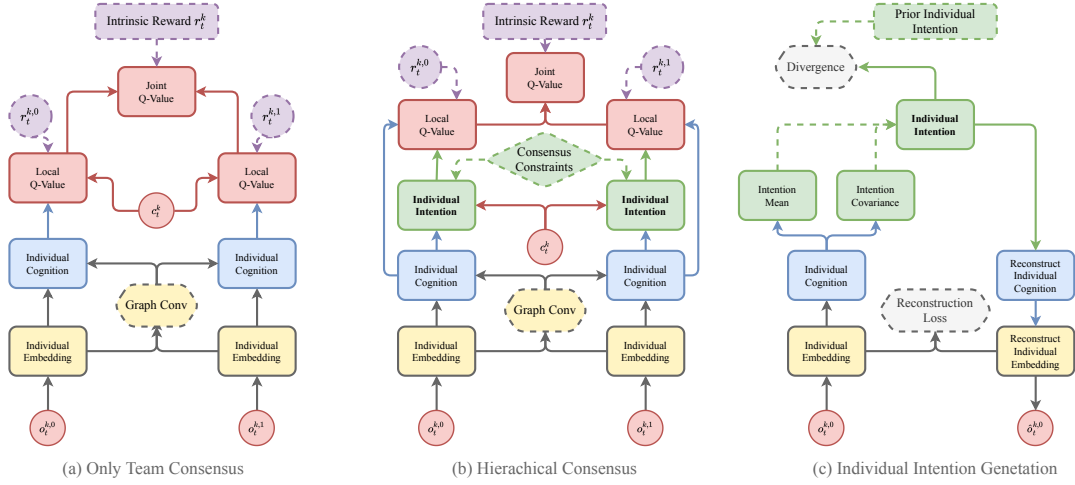
Recall the loss function of the variational autoencoder

$$\mathcal{L}_\ell^i(\varphi) = \mathbb{E} \left[ \|o_t^{k,i} - \hat{o}_t^{k,i}\|_2^2 + \text{KL} \left[ q_\varphi(\zeta_t^{k,i} | o_t^{k,i}) \| p(\zeta_t^{k,i}) \right] \right].$$

The prior distribution  $p(\zeta_t^{k,i})$  in second term can be replaced by our above consensus constraints, which leads to

$$\mathcal{L}_\ell^i(\phi, \psi, \varphi) = \mathbb{E} \left[ \|o_t^{k,i} - \hat{o}_t^{k,i}\|_2^2 + \frac{1}{n_k - 1} \sum_{j \neq i} \text{KL} \left[ q_\varphi(\zeta_t^{k,i} | o_t^{k,i}) \| q_\varphi(\zeta_t^{k,j} | o_t^{k,j}) \right] \right].$$

<sup>2</sup>Since the reparameterization trick is used in order to enable the backpropagation during implementation, it is reasonable to write the equal sign here.



**Figure 3: The network structure of the hierarchical consensus module and decision module. (a) The baseline structure that only uses team intention to generate final policies. (b) The individual intention is generated based on the team intention and consensus constraints. (c) The individual intention generator of ROCHICO algorithm is a standard variational autoencoder.**

The overall problem of individual intention generation can be formulated as

$$\min \mathcal{L}_\ell^{ig}(\phi, \psi, \varphi) := \sum_{i=1}^n \mathcal{L}_\ell^i(\phi, \psi, \varphi).$$

### 4.3 Decision Module

Before outputs the agents' policies, it should be noted that the team intention to affect the final policies indirectly. To affect the agents' behaviors more directly and efficiently, we introduce another intrinsic team reward and utilize the QMIX [34] technique to assign the rewards to all agents. Formally, the intrinsic team reward  $r_{t,\ell}^k$  is defined as:

$$r_{t,\ell}^k = \sum_{j \neq k} \|c_t^k - c_t^j\|_2^2.$$

We additionally use the external reward  $r_\ell^i$  to learn the local  $Q$  function, as a result the local  $Q$  function will be trained based on two different reward signals  $r_\ell^i$  and  $r_{t,\ell}^k$ . To make the behavior of the agents in each team have a certain diversity, we firstly combine the individual intention  $\zeta_t^{k,i}$  with the individual cognition  $\chi_t^{k,i}$  as the input of the QMIX algorithm, and generates the local  $Q$ -value. Further we optimize the local  $Q$ -value function by minimize following TD(0) error:

$$\mathcal{L}_\ell^i(\theta_\ell^i) = \mathbb{E}_{(\zeta_t^{k,i}, \chi_t^{k,i}, a_{t,\ell}^i) \sim D} \left[ \left( Q_{\theta_\ell^i}(\zeta_t^{k,i}, \chi_t^{k,i}, a_{t,\ell}^i) - y^i \right)^2 \right].$$

where  $y^i = r_\ell^i + \gamma \max_{a_{t+1,\ell}^i} Q_{\theta_\ell^i}(\zeta_{t+1}^{k,i}, \chi_{t+1}^{k,i}, a_{t+1,\ell}^i)$  and  $\theta_\ell^i, \bar{\theta}_\ell^i$  parameterize local  $Q$  function and local target  $Q$  function of agent  $i$  respectively. For each team  $k$ , the team joint  $Q$ -value can be denoted

as

$$Q_{\theta_\ell^k}(\zeta_t^k, \chi_t^k, a_{t,\ell}^k) = Q_{\theta_\ell^k} \left( Q_{\theta_\ell^k}(\zeta_t^{k,1}, \chi_t^{k,1}, a_{t,\ell}^{k,1}), \dots, Q_{\theta_\ell^k}(\zeta_t^{k,n_k}, \chi_t^{k,n_k}, a_{t,\ell}^{k,n_k}) \right),$$

and we also optimize the team joint  $Q$ -value function by minimize following TD(0) error:

$$\mathcal{L}_\ell^k(\theta_\ell^k) = \mathbb{E}_{(\zeta_t^k, \chi_t^k, a_{t,\ell}^k) \sim D} \left[ \left( Q_{\theta_\ell^k}(\zeta_t^k, \chi_t^k, a_{t,\ell}^k) - y^k \right)^2 \right].$$

where  $y^k = r_{t,\ell}^k + \gamma \max_{a_{t+1,\ell}^k} Q_{\theta_\ell^k}(\zeta_{t+1}^k, \chi_{t+1}^k, a_{t+1,\ell}^k)$  and  $\theta_\ell^k, \bar{\theta}_\ell^k$  parameterize joint  $Q$ -value function and joint target  $Q$ -value function of team  $k$  respectively. The overall learning problem of decision module is

$$\min \mathcal{L}_\ell^Q(\{\theta_\ell^i\}, \{\bar{\theta}_\ell^i\}) = \sum_i \mathcal{L}_\ell^i(\theta_\ell^i) + \lambda_{QMIX} \sum_k \mathcal{L}_\ell^k(\theta_\ell^k).$$

## 5 EXPERIMENTS

### 5.1 Environments

We evaluate the algorithm performances on four large-scale cooperative environments, including *Pacmen*, *Block*, *Pursuit* and *Battle*. More details can be found in appendix.

**Pacmen.** 64 agents initialized at the maze center and 64 dots scatter randomly at four corners of the squared map. Agents get the reward by eat dots. The dots are distributed in different corners, the agent needs to team up and travel to different corners to eat more dots.

**Block.** There are 32 blockers and 32 blockees who have superior speed than the blockers. There also are 64 foods initialized at one side of the squared map. Blockers and blockees are only rewarded by eat foods. Since blockee runs faster than blocker, blocker needs to learn diverse policies to block blockees and eat food Simultaneously.

**Pursuit.** There are 64 predators and 64 preys who have superior speed than the predators. Since the prey runs faster than the predators, the predators need to learn to round up through structured diversification policies.

**Battle.** 64 agents learn to fight against 64 enemies who have superior abilities than agents. As the hit point of enemy is 10 (more than single agent’s damage), agents need to continuously cooperate to kill the enemy. All environments are implemented by MAgent [53].

## 5.2 Baselines

IDQN is chosen as the baseline. Due to the connection between ROCHICO with QMIX and NCC-Q, we also compare these two methods as baselines. However, both QMIX and NCC-Q are not designed for large-scale scenarios. Therefore, we first randomly split all agents to multiple teams, and use QMIX and NCC-Q algorithms in each one. See Appendix for detailed hyperparameter settings.

## 5.3 Performance Comparison

*Pacmen* and *Block* need to pay more attention to the division of labor of agents than *Pursuit* and *Battle* environments, so the diversity of policies and the degree of collaboration between agents should have a greater impact on the final performance. As can be seen from Figures 4(a) and Figures 4(b), since IDQN does not take into account the cooperation between agents, the weak individual ability limits the overall task completion. Compared with IDQN, QMIX, and NCC-Q, which encourage the cooperation between agents explicitly, achieve a certain degree of diversification of policies with better performance. In addition, NCC-Q explicitly imposes consensus constraints on the policies within the team, thereby it can achieve better collaboration. However, QMIX and NCC-Q both use predefined teaming strategies, they cannot dynamically adapt to the non-stationary environment. ROCHICO outperforms all the other three methods: it can perform adaptive teaming because of reinforced organization control; and it can achieve tight cooperation even when the teams change dynamically because of the hierarchical consensus learning.

For *Pursuit* and *Battle* environments, more attentions are paid on the flexibility of agents’ collaboration policies. ROCHICO can achieve the best performance in these two more complex environments, as seen from Figure 4(c) and Figure 4(d). The ability of the individual agent is more important in *Pursuit*, as a result, the performance of IDQN is better than QMIX and NCC-Q; the flexible switching ability of the individual agent is more important in *Battle*, the performance of IDQN becomes poor.

## 5.4 Ablation Study

In this part, we first conduct an ablation analysis on the three important components of the ROCHICO algorithm: hierarchical consensus constraints, structural consistency intrinsic reward in the organization control module, and the intrinsic reward in the decision module. Here we choose the *Pursuit* environment because all comparisons have the most significant difference in this environment. The detailed comparisons can be summarized as follows:

1). We can see from Figure 4(e) that removing the structural consistency intrinsic reward (i.e., ROCHICO-C) has the least impact on the algorithm performance. This is because the edge between two

nodes in the organization control module is determined by the OR operation, which makes it difficult for the edge existing at the previous timestep to disappear due to randomness at the next timestep. This indirectly realizes a certain degree of regularity for the stability of the graph structure.

2). It can be seen from the Figure 4(e) that after removing the hierarchical consistency constraint (ROCHICO-G, shown in Figure 3(a)), the performance of the ROCHICO algorithm has dropped significantly. Although IDQN performs better than QMIX and NCC-Q in Figure 4(c), this does not mean that cooperation is unnecessary in *Pursuit* environment. QMIX and NCC-Q use a predefined teaming strategy instead of the adaptive strategy in ROCHICO, which limits the capabilities of the agents.

3). After removing the intrinsic reward in decision module (ROCHICO-I, replaced by sum of local rewards of all agents in the team), the performance of ROCHICO-I has a great decline from ROCHICO. This shows that use pf team intention to promote diversification indirectly is difficult. Through utilizing the intrinsic rewards based on intention difference between teams, the diversified policy emergence will become easier.

	ROCHICO-C	ROCHICO-G	ROCHICO-I	ROCHICO-1	ROCHICO-3	ROCHICO(-2)
<i>Pacmen</i>	953(±56)	889(±78)	803(±45)	893(±52)	767(±74)	988(±40)
<i>Block</i>	687(±45)	635(±38)	580(±52)	685(±56)	620(±44)	736(±46)
<i>Pursuit</i>	2252(±420)	2012(±385)	1588(±394)	2486(±403)	2403(±345)	2548(±380)
<i>Battle</i>	221(±24)	199(±33)	169(±23)	194(±26)	160(±32)	232(±27)

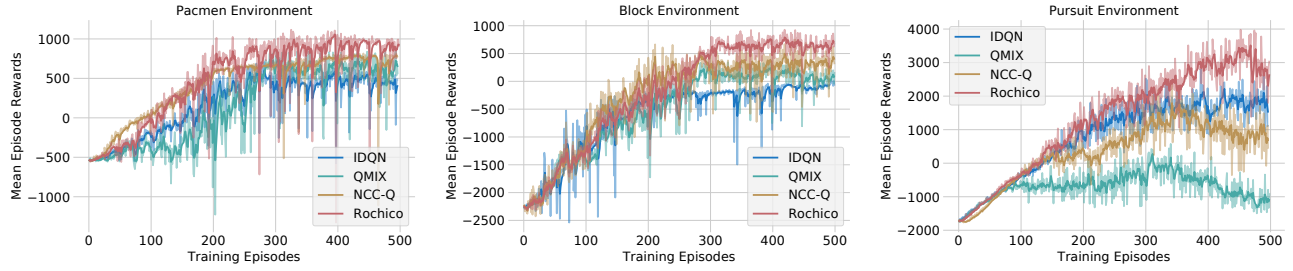
**Table 1: The average episode rewards of ablation algorithms in test environment. The mean and standard variance are calculated under 5 random seeds.**

The ablation analysis on the decision of the organization control module is shown in Figure 4(f). The number suffix  $k = \{1, 2, 3\}$  behind the ROCHICO indicates that each agent needs to decide at the same time whether to form a team with the nearest  $k$  agents.  $k = 2$  is the default setting of ROCHICO algorithm. It can be seen from Figure 4(f) that ROCHICO-2(ROCHICO) outperforms both ROCHICO-1 and ROCHICO-3. ROCHICO-1 has a smaller range which will weaken the connection between agents, and less effective collaboration between agents can be established; ROCHICO-3 has a larger range which will make the teams too fixed, and unable to flexibly adapt to the non-stationary environment. The average episode rewards of the above ablations in the test environment are shown in Table 1.

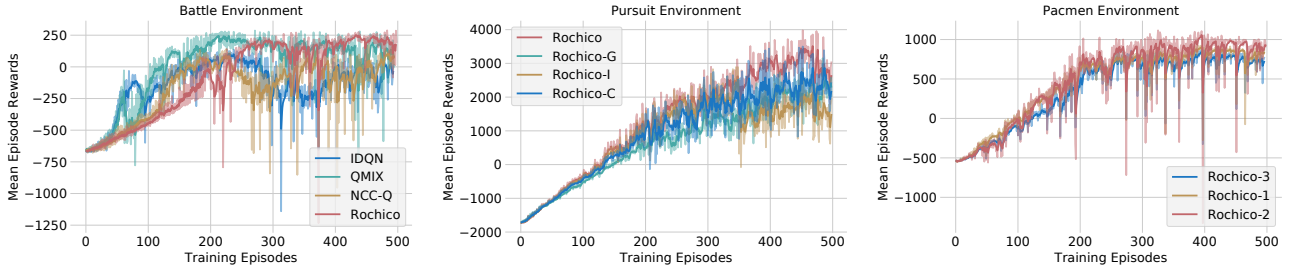
## 5.5 Emergence Behavior Analysis

Figure 5 shows the changing of the averaged team number in the training process. The curves corresponding to all environments show a downward trend in the training process, which means that the organization control module does learn teaming strategies that can promote cooperation. It can be seen from the Figure 5 that the *Pursuit* environment has the largest averaged team number. The *Pursuit* environment is more focused on the ability of the agent itself, which is also consistent with the previous analysis.

Further, we select the *Pacmen* environment to deeply analyze the teaming strategy of the organization control module and the individual intentions at different stages for task completion. For the convenience of the presentation, we conduct the experiment in the



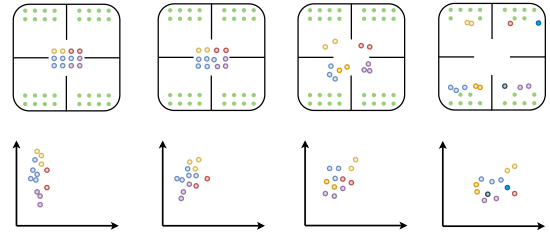
(a) The performance comparison in Pacmen environment. (b) The performance comparison in Block environment. (c) The performance comparison in Pursuit environment.



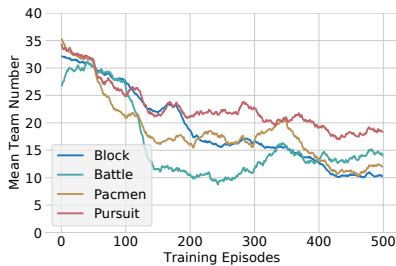
(d) The performance comparison in Battle environment. (e) The ablation study in Pursuit environment. (f) The ablation study in self-organization phase.

**Figure 4: The performance comparison and ablation study of the Rochico algorithm.**

*Pacmen* with only 12 agents. The different colors in Figure 6 represent different teams (green represents food). The following two-dimensional scatter plot is the visualized result of the t-SNE [23] algorithm by reducing the dimensionality of the individual intentions at the corresponding timestep. It can be seen from Figure 6 that the points belonging to the same team show obvious aggregation, while the points belonging to different teams are far apart. As the task progresses, if the difference between the sub-tasks completed by different teams becomes larger, the corresponding points will be farther away.



**Figure 6: The relationship between team pattern and team intention in *Pacmen* environment.**



**Figure 5: The averaged team number changing.**

## 6 CONCLUSION

In this paper, in order to improve the efficiency of multi-agent exploration and collaboration in complex tasks, we propose a MARL framework ROCHICO based on reinforced organization control and hierarchical consensus learning. In the organization control module, we model the multi-agent organization control problem as a

cooperative POSG and use an independent MARL algorithm to output an adaptive teaming strategy. In the hierarchical consensus module, based on the auxiliary tasks of contrastive learning and self-supervised learning, the exploration efficiency and the collaboration efficiency of multi-agents are improved through hierarchical consensus learning. The comparison between ROCHICO and current SOTA cooperative MARL algorithms in four large-scale cooperative multi-agent environments shows that our algorithm can complete complex tasks more efficiently through richer policies diversity and tighter agents collaboration.

## 7 ACKNOWLEDGE

This work was supported in part by National Key Research and Development Program of China (No. 2020AAA0107400), NSFC (No. 12071145), STCSM (No. 18DZ2270700 and 20511101100), the Open Research Projects of Zhejiang Lab (NO.2021KE0AB03) and a grant from Shenzhen Institute of Artificial Intelligence and Robotics for Society.



## REFERENCES

- [1] Sherief Abdallah and Victor R. Lesser. 2007. Multiagent reinforcement learning and self-organization in a network of agents. In *AAMAS*.
- [2] Georgios Chalkiadakis, Edith Elkind, Evangelos Markakis, Maria Polukarov, and Nicholas R. Jennings. 2010. Cooperative Games with Overlapping Coalitions. *J. Artif. Intell. Res.* 39 (2010), 179–216.
- [3] Mehdi Dastani, Virginia Dignum, and Frank Dignum. 2003. Role-assignment in open agent societies. In *AAMAS*.
- [4] Christian Schroeder de Witt, Jakob Foerster, Gregory Farquhar, Philip Torr, Wendelin Boehmer, and Shimon Whiteson. 2019. Multi-Agent Common Knowledge Reinforcement Learning. In *NeurIPS*.
- [5] Giovanna Di Marzo Serugendo, Marie-Pierre Gleizes, and Anthony Karageorgos. 2005. Self-organization in multi-agent systems. *Knowledge Engineering Review* 20, 2 (2005), 165–189.
- [6] Daniela Scherer Dos Santos and Ana LC Bazzan. 2012. Distributed clustering for group formation and task allocation in multiagent systems: A swarm intelligence approach. *Applied Soft Computing* 12, 8 (2012), 2123–2131.
- [7] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. 2019. Diversity is All You Need: Learning Skills without a Reward Function. In *ICLR*.
- [8] Jakob Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. 2016. Learning to communicate with deep multi-agent reinforcement learning. In *NeurIPS*.
- [9] Matthew E Gaston and Marie Desjardins. 2005. Agent-organized networks for dynamic team formation. In *AAMAS*.
- [10] Robin Glinton, Katia P. Sycara, and Paul Scerri. 2008. Agent Organized Networks Redux. In *AAAI*.
- [11] Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine. 2018. Composable deep reinforcement learning for robotic manipulation. In *ICRA*.
- [12] Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. 2017. Reinforcement learning with deep energy-based policies. In *ICML*.
- [13] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *ICML*.
- [14] Eric A Hansen, Daniel S Bernstein, and Shlomo Zilberstein. 2004. Dynamic programming for partially observable stochastic games. In *AAAI*.
- [15] Jiechuan Jiang, Chen Dun, Tiejun Huang, and Zongqing Lu. 2019. Graph Convolutional Reinforcement Learning. In *ICLR*.
- [16] Jiechuan Jiang and Zongqing Lu. 2018. Learning attentional communication for multi-agent cooperation. In *NeurIPS*.
- [17] Jiechuan Jiang and Zongqing Lu. 2020. The Emergence of Individuality in Multi-Agent Reinforcement Learning. *arXiv preprint arXiv:2006.05842* (2020).
- [18] Daewoo Kim, Sangwoo Moon, David Hostallero, Wan Ju Kang, Taeyoung Lee, Kyunghwan Son, and Yung Yi. 2018. Learning to Schedule Communication in Multi-agent Reinforcement Learning. In *ICLR*.
- [19] Ramachandra Kota, Nicholas Gibbins, and Nicholas R. Jennings. 2012. Decentralized approaches for self-adaptation in agent organizations. *ACM Trans. Auton. Adapt. Syst.* 7 (2012), 1:1–1:28.
- [20] Youngwoon Lee, Jingyun Yang, and Joseph J Lim. 2020. Learning to Coordinate Manipulation Skills via Skill Behavior Diversification. In *ICLR*.
- [21] Wenhao Li, Bo Jin, Xiangfeng Wang, Junchi Yan, and Hongyuan Zha. 2020. F2A2: Flexible Fully-decentralized Approximate Actor-critic for Cooperative Multi-agent Reinforcement Learning. *arXiv preprint arXiv:2004.11145* (2020).
- [22] Ryan Lowe, Jakob Foerster, Y-Lan Boureau, Joelle Pineau, and Yann Dauphin. 2019. On the Pitfalls of Measuring Emergent Communication. In *AAMAS*.
- [23] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, Nov (2008), 2579–2605.
- [24] Kathryn Sarah Macarthur, Ruben Stranders, Sarvapali Ramchurn, and Nicholas Jennings. 2011. A distributed anytime algorithm for dynamic task allocation in multi-agent systems. In *AAAI*.
- [25] Hangyu Mao, Wulong Liu, Jianye Hao, Jun Luo, Dong Li, Zhengchao Zhang, Jun Wang, and Zhen Xiao. 2020. Neighborhood Cognition Consistent Multi-Agent Reinforcement Learning. In *AAAI*.
- [26] Milan Mares. 2000. Fuzzy coalition structures. *Fuzzy Sets Syst.* 114 (2000), 23–33.
- [27] Laëtitia Matignon, Laurent Jeanpierre, and Abdel-Ilhah Mouaddib. 2012. Coordinated multi-robot exploration under communication constraints using decentralized Markov decision processes. In *AAAI*.
- [28] Kevin R McKee, Ian Gemp, Brian McWilliams, Edgar A Duñez-Guzmán, Edward Hughes, and Joel Z Leibo. 2020. Social Diversity and Social Preferences in Mixed-Motive Reinforcement Learning. In *AAMAS*.
- [29] Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. 2008. Optimal and approximate Q-value functions for decentralized POMDPs. *Journal of Artificial Intelligence Research* 32 (2008), 289–353.
- [30] Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V Albrecht. 2020. Comparative Evaluation of Multi-Agent Deep Reinforcement Learning Algorithms. *arXiv preprint arXiv:2006.07869* (2020).
- [31] Peng Peng, Quan Yuan, Ying Wen, Yaodong Yang, Zhenkun Tang, Haitao Long, and Jun Wang. 2017. Multiagent bidirectionally-coordinated nets for learning to play starcraft combat games. *arXiv preprint arXiv:1703.10069* (2017).
- [32] Alexander Peysakhovich and Adam Lerer. 2018. Prosocial Learning Agents Solve Generalized Stag Hunts Better than Selfish Ones. In *AAMAS*.
- [33] Sarvapali D Ramchurn, Alessandro Farinelli, Kathryn S Macarthur, and Nicholas R Jennings. 2010. Decentralized coordination in robocup rescue. *Comput. J.* 53, 9 (2010), 1447–1461.
- [34] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *ICML*.
- [35] Tabish Rashid, Mikayel Samvelyan, C. S. Witt, Gregory Farquhar, Jakob N. Foerster, and S. Whiteson. 2018. QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. In *ICML*.
- [36] Pedro V Sander, Denis Peleshchuk, and Barbara J Grosz. 2002. A scalable, distributed algorithm for efficient task allocation. In *AAMAS*.
- [37] Alberto Sanfeliu and King-Sun Fu. 1983. A distance measure between attributed relational graphs for pattern recognition. *IEEE transactions on systems, man, and cybernetics* 3 (1983), 353–362.
- [38] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *CVPR*.
- [39] Junjie Sheng, Xiangfeng Wang, Bo Jin, Junchi Yan, Wenhao Li, Tsung-Hui Chang, Jun Wang, and Hongyuan Zha. 2020. Learning Structured Communication for Multi-agent Reinforcement Learning. *arXiv preprint arXiv:2002.04235* (2020).
- [40] Tianmin Shu and Yuandong Tian. 2019. M<sup>3</sup> RL: Mind-aware Multi-agent Management Reinforcement Learning. In *ICLR*.
- [41] Mark Sims, Daniel Corkill, and Victor Lesser. 2008. Automated organization design for multi-agent systems. *Autonomous agents and multi-agent systems* 16, 2 (2008), 151–185.
- [42] Yuhang Song, Jianyi Wang, Thomas Lukasiewicz, Zhenghua Xu, Mai Xu, Zihan Ding, and Lianlong Wu. 2020. Arena: A General Evaluation Platform and Building Toolkit for Multi-Agent Intelligence. In *AAAI*.
- [43] Peter Sunehag, G. Lever, A. Gruslys, W. Czarniecki, V. Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, K. Tuyls, and T. Graepel. 2018. Value-Decomposition Networks For Cooperative Multi-Agent Learning. In *AAMAS*.
- [44] Robert Tarjan. 1972. Depth-first search and linear graph algorithms. *SIAM journal on computing* 1, 2 (1972), 146–160.
- [45] Zheng Tian, Shihao Zou, Tim Warr, Lisheng Wu, and Jun Wang. 2018. Learning to communicate implicitly by actions. *arXiv preprint arXiv:1810.04444* (2018).
- [46] Tonghan Wang, Heng Dong, Victor Lesser, and Chongjie Zhang. 2020. Multi-Agent Reinforcement Learning with Emergent Roles. In *ICML*.
- [47] Jiachen Yang, Igor Borovikov, and Hongyuan Zha. 2020. Hierarchical Cooperative Multi-Agent Reinforcement Learning with Skill Discovery. In *AAMAS*.
- [48] Dayong Ye, Minjie Zhang, and Danny Sutanto. 2013. Self-Adaptation-Based Dynamic Coalition Formation in a Distributed Agent Network: A Mechanism and a Brief Survey. *IEEE Transactions on Parallel and Distributed Systems* 24 (2013), 1042–1051.
- [49] Dayong Ye, Minjie Zhang, and Athanasios V Vasilakos. 2016. A survey of self-organization mechanisms in multiagent systems. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 47, 3 (2016), 441–461.
- [50] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. 2017. Deep sets. In *NeurIPS*.
- [51] Chongjie Zhang, Sherief Abdallah, and Victor Lesser. 2009. Integrating organizational control into multi-agent learning. In *AAMAS*.
- [52] Chongjie Zhang, Victor R Lesser, and Sherief Abdallah. 2010. Self-organization for coordinating decentralized reinforcement learning. In *AAMAS*.
- [53] L Zheng, J Yang, H Cai, W Zhang, J Wang, and Y Yu. 2018. MAgent: A many-agent reinforcement learning platform for artificial collective intelligence. In *AAAI*.

## SUPPLEMENTARY MATERIAL

### A. Environments

**Pacmen.** This scenario is a fully cooperative task, where  $N$  agents initialized at the maze center and  $M$  dots scatter randomly at four corner rooms. Agents get the reward by eating dots. Each agent has a local observation that contains a circle view with a radius 7 centered at the agent itself. The moving or attacking range of the agent is the only 1 neighbor grid. The reward is  $-0.01$  for moving,  $+0.5$  for attacking the dot,  $-0.1$  for attacking a blank grid, and  $+5$  for eat a dot. Since the dots are distributed in different corners, the agent needs to be grouped automatically and travel to different corners to eat more dots.

**Block.** This scenario is a fully cooperative task, where  $N$  blockers and  $L$  blockees who have superior speed than the blocker. There also are  $M$  foods initialized on one side of the map. Blockers and blockees are only rewarded by eating foods. The moving or attacking range of the blockers is the 5 neighbor grids and the blockees have larger local observed range. The blockers have very large hit points and therefore are considered indestructible. Blockees could be killed by blockers. The reward is 0 for moving,  $-0.2$  for attacking,  $-1$  for being killed, and  $+5$  for eat one food. Since the blockee runs faster than the blocker, the blocker needs to learn to use diverse strategies to block blockees and eat food at the same time.

**Pursuit.** This scenario is a fully cooperative task, where  $N$  predators and  $L$  preys who have superior speed than the predators. The moving or attacking range of the agent is the 4 neighbor grids and the predator has a larger local observed range. The reward is 0 for moving,  $-0.2$  for attacking the prey, 1 for killing,  $-1$  for being killed, and  $-0.2$  for attacking a blank grid. Since the prey runs faster than the predator, the predator needs to learn to round up through the division of labor and cooperation.

**Battle.** This scenario is a fully cooperative task, where  $N$  agents learn to fight against  $L$  enemies who have superior abilities than the agents. The moving or attacking range of the agent is the 4 neighbor grids, however, the enemy can move to one of 12 nearest grids or attack one of 8 neighbor grids. Each agent/enemy has 10 hit points. The reward is  $-0.005$  for moving,  $+5$  for attacking the enemy,  $-0.1$  for being killed, and  $-0.1$  for attacking a blank grid. As the hit point of the enemy is 10, agents have to continuously cooperate to kill the enemy. Therefore, the task is much more challenging than Pursuit in terms of learning to cooperate. The schematic diagrams of all environments are shown in Figure 7.

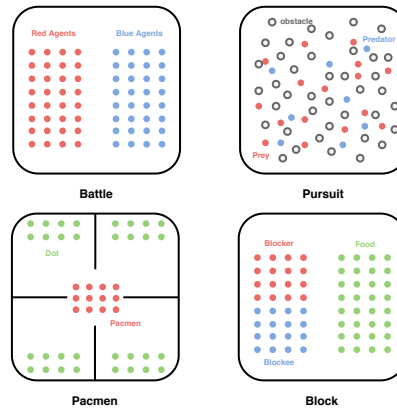


Figure 7: The simulated environments.

### B. Hyperparameters

The detailed hyperparameter settings of all algorithms are shown in Table 2, Table 3, Table 4 and Table 5.

### C. Notations

The notations of three module in ROCHICO are shown in Table 6, Table 7 and Table 8.

Hyperparameter	Value
Batch Size	512
Experience Replay Buffer Size	1,024,000
Target Network Update Frequency	1 times per 1000 samples
Train Frequency	4 times per batch
Train Episodes	500
Max Episode Length	250
<i>epsilon</i> Piecewise Decay	[0, 200, 400] episodes, [1, 0.2, 0.05]
Layer Number of CNN	2
Kernel Size of Each Layer	$3 \times 3$
Kernel Number of Each Layer	32
Hidden Size of MLP	[256, 512]
Activation Function	ReLU
Double Q-Learning	True
Dueling Q-Learning	True
Learning Rate	0.0001
$\gamma$	0.99
Hidden Sizes of MLP in Team Intention Generator	[32, 32, 32]
Kernel Size of CNN in Team Intention Generator	$3 \times 3$
Kernel Number of CNN in Team Intention Generator	32
Hypernet Kernel Size of CNN	$3 \times 3$
Hypernet Layer Number of MLP	2
Hypernet Hidden Size of MLP	64
Layer Number of GCN	1
Hidden Size of GCN	32
Hidden Sizes of VAE Encoder	[32, 32]
Hidden Sizes of VAE Decoder	[32, 32]
Dimension of Hierarchical Intentions	32

**Table 2: The hyperparameters setting of ROCHICO algorithm.**

Hyperparameter	Value
Batch Size	512
Experience Replay Buffer Size	4,096,000
Target Network Update Frequency	1 times per 1000 samples
Train Frequency	4 times per batch
Train Episodes	500
Max Episode Length	250
<i>epsilon</i> Piecewise Decay	[0, 200, 400] episodes, [1, 0.2, 0.05]
Layer Number of CNN	2
Kernel Size of Each Layer	$3 \times 3$
Kernel Number of Each Layer	32
Hidden Size of MLP	[256, 512]
Activation Function	ReLU
Double Q-Learning	True
Dueling Q-Learning	True
Learning Rate	0.0001
$\gamma$	0.99

**Table 3: The hyperparameters setting of IDQN algorithm.**

Hyperparameter	Value
Batch Size	512
Experience Replay Buffer Size	4,096,000
Target Network Update Frequency	1 times per 1000 samples
Train Frequency	4 times per batch
Train Episodes	500
Max Episode Length	250
<i>epsilon</i> Piecewise Decay	[0, 200, 400] episodes, [1, 0.2, 0.05]
Layer Number of CNN	2
Kernel Size of Each Layer	$3 \times 3$
Kernel Number of Each Layer	32
Hidden Size of MLP	[256, 512]
Activation Function	ReLU
Double Q-Learning	True
Dueling Q-Learning	True
Learning Rate	0.0001
$\gamma$	0.99
Hypernet Kernel Size of CNN	$3 \times 3$
Hypernet Layer Number of MLP	2
Hypernet Hidden Size of MLP	64

**Table 4: The hyperparameters setting of QMIX algorithm.**

Hyperparameter	Value
Batch Size	512
Experience Replay Buffer Size	1,024,000
Target Network Update Frequency	1 times per 1000 samples
Train Frequency	4 times per batch
Train Episodes	500
Max Episode Length	250
<i>epsilon</i> Piecewise Decay	[0, 200, 400] episodes, [1, 0.2, 0.05]
Layer Number of CNN	2
Kernel Size of Each Layer	$3 \times 3$
Kernel Number of Each Layer	32
Hidden Size of MLP	[256, 512]
Activation Function	ReLU
Double Q-Learning	True
Dueling Q-Learning	True
Learning Rate	0.0001
$\gamma$	0.99
Hypernet Kernel Size of CNN	$3 \times 3$
Hypernet Layer Number of MLP	2
Hypernet Hidden Size of MLP	64
Layer Number of GCN	1
Hidden Size of GCN	32
Hidden Sizes of VAE Encoder	[32, 32]
Hidden Sizes of VAE Decoder	[32, 32]
Dimension of Latent Variable	32

**Table 5: The hyperparameters setting of NCC - Q algorithm.**

Symbols	Descriptions
$\mathcal{G}$	The directed graph constructed by all agents in the environment.
$\mathcal{V}$	The node set of graph $\mathcal{G}$ and each node represents a agent.
$\mathcal{E}$	The edge set of graph $\mathcal{G}$ and edges are determined by agents' policies.
$v$	A node in node set $\mathcal{V}$ .
$n$	The mode of node set $\mathcal{V}$ .
$M_{org}$	The POSG used to model organization control problem.
$\mathcal{X}_u$	The agent space of $M_{org}$ .
$\mathcal{S}_u$	The state space of $M_{org}$ .
$\mathcal{A}_u^i$	The action space of agent $i$ belongs to agent space $\mathcal{X}_u$ .
$\mathcal{O}_u^i$	The obsercation space of agent $i$ belongs to agent space $\mathcal{X}_u$ .
$\mathcal{P}_u$	The transition model of $M_{org}$ .
$\mathcal{E}_u$	The emission probability model of $M_{org}$ .
$\mathcal{R}_u^i$	The reward function of agent $i$ belongs to agent space $\mathcal{X}_u$ .
$a_u^i, a_u^{i'}$	Two consecutive actions belongs to action space $\mathcal{A}_u^i$ .
$o_u^i, o_u^{i'}$	The consecutive observations belongs to observation space $\mathcal{O}_u^i$ .
$m$	The upper bound of nearest neighbors of any agent in agent space $\mathcal{X}_u$ .
$N_m(i)$	The $m$ -nearest neighbors set of agent $i$ belongs to $\mathcal{X}_u$ .
$d(i)$	The index of agent $i$ of action belongs to action space $\mathcal{A}_u^i$ .
$\mathcal{G}_u$	The corresponding undirected graph of $\mathcal{G}$ .
$\mathcal{V}_u$	The node set of graph $\mathcal{G}_u$ and each node represents a agent.
$\mathcal{E}_u$	The edge set of graph $\mathcal{G}_u$ and edges are determined by agents' policies.
$e(i, j)$	The edge between node $i$ and node $j$ in $\mathcal{G}_u$ .
$r_e^i$	The external reward of agent $i$ .
$r_u^i$	The structural consistency intrinsic reward of agent $i$ .
$r_{u+}^i$	The total reward of agent $i$ .
$\alpha_u$	The strength of contrait for structural consistency.
$\text{GED}(\cdot)$	The graph edit distance
$\mathcal{G}_u(N_m(i) \vee i)$	The sub-graph only contains node $i$ and its $m$ -nearest neighbors before take action $a_u^i$ .
$\mathcal{G}_u^r(N_m(i) \vee i)$	The sub-graph only contains node $i$ and its $m$ -nearest neighbors after take action $a_u^i$ .
$\mathcal{J}_u$	The optimization goal of organization control module.
$\tau_u$	The sample trajectory.
$\mathcal{L}_u^i(\cdot)$	The TD(0) error of agent $i$ .
$\theta_u, \bar{\theta}_u$	The parameters of $Q$ function and target $Q$ function.
$D$	The experience replay buffer.
$y$	The TD backup.
$\gamma$	The discount factor.
$\mathcal{L}_u^Q(\theta_u)$	The overall objective function of organization control module.

**Table 6: Notations in the organization control module.**

Symbols	Descriptions
$f_{\mu}(\cdot)$	The state encoder parameterized by $\mu$ .
$k, u, v$	The team indicators.
$n_k$	The agent number of team $k$ .
$o_t^{k,i}$	The observation of agent $i$ in team $k$ at timestep $t$ .
$\mathbf{o}_t^k$	The joint observation of all agents in team $k$ at timestep $t$ .
$e_t^{k,i}$	The observation embedding of agent $i$ in team $k$ at timestep $t$ .
$f_v(\cdot)$	The DeepSet network parameterized by $v$ .
$e_t^k$	The team embedding of team $k$ at timestep $t$ .
$f_{\omega}(\cdot)$	The team intention encoder parameterized by $\omega$ .
$s_t$	The global state of timestep $t$ .
$c_t^k$	The team intention of team $k$ at timestep $t$ .
$(x_t^{k,i}, y_t^{k,i})$	The spatial position of agent $i$ in team $k$ at timestep $t$ .
$[\bar{x}_t^k, \bar{y}_t^k]$	The average spatial position of team $k$ at timestep $t$ .
$t_k$	The timestep of team $k$ .
$e_{ts}^k$	The spatiotemporal feature representation of team $k$ .
$d(k, l)$	The Euclidean distance in the spatiotemporal space between team $k$ and $l$ .
$\mathcal{G}_{\ell}(\mathcal{V}_{\ell}, \mathcal{E}_{\ell})$	The undirected graph constructed by team (similar as undirected graph in organization control module).
$e(k, l)$	The edge between team $k$ and team $l$ .
$y_t^k$	The label of team $k$ .
$\mathcal{L}_{\ell}^k(\mu, v, \omega)$	The contrastive learning objective function of hierarchical consensus learning module.
$m$	The margin parameter in contrastive learning objective function.
$f_{\xi}(\cdot)$	The team intention decoder parameterized by $\xi$ .
$\hat{o}_{t+1}^{k,i}$	The reconstructed observation by team intention decoder of agent $i$ in team $k$ at timestep $t$ .
$\hat{\mathbf{o}}_{t+1}^k$	The joint reconstructed observation by team intention decoder of all agents in team $k$ at timestep $t$ .
$\mathcal{L}_{\ell}^k(\xi)$	The prediction loss of team intention decoder.
$\mathcal{L}_{\ell}^{tg}(\mu, v, \omega, \xi)$	The loss function of team intention generator.
$\lambda_{tg}$	The temperature parameter of loss function of team intention generator.
$g_{\phi}(\cdot)$	The individual encoder parameterized by $\phi$ .
$h_t^{k,i}$	The individual embedding of agent $i$ in team $k$ at timestep $t$ .
$g_{\psi}(\cdot)$	The graph convolutional network parameterized by $\psi$ .
$\mathbf{h}_t^k$	The joint embedding of all agents in team $k$ at timestep $t$ .
$\chi_t^{k,i}$	The individual cognition of agent $i$ in team $k$ at timestep $t$ .
$g_{\varphi}(\cdot)$	The variational encoder parameterized by $\varphi$ .
$\zeta_t^{k,i}$	The individual intention of agent $i$ in team $k$ at timestep $t$ .
$\mathcal{L}_{\ell}^i(\phi, \psi, \varphi)$	The hierarchical consensus loss of hierarchical consensus module.
$\mathcal{L}_{\ell}^i(\varphi)$	The loss function of variation autoencoder.
$p(\zeta_t^{k,i})$	The prior distribution of individual intention of agent $i$ in team $k$ at timestep $t$ .
$\mathcal{L}_{\ell}^{ig}(\phi, \psi, \varphi)$	The loss function of individual intention generator.

**Table 7: Notations in the hierarchical consensus module.**

<b>Symbols</b>	<b>Descriptions</b>
$r_{t,\ell}^k$	The intrinsic reward of team $k$ at timestem $t$ .
$\mathcal{L}_\ell^i(\theta_\ell^i)$	The TD(0) error of local $Q$ function of agent $i$ in team $k$ at timestep $t$ .
$Q_{\theta_\ell^i}(\cdot), Q_{\bar{\theta}_\ell^i}$	The local $Q$ function and local target $Q$ function of agent $i$ in team $k$ , parameterized by $\theta_\ell^i$ and $\bar{\theta}_\ell^i$ respectively.
$y^i$	The TD backup of agent $i$ in team $k$ .
$\mathcal{L}_\ell^k(\theta_\ell^k)$	The TD(0) error of joint $Q$ function of team $k$ at timestep $t$ .
$Q_{\theta_\ell^k}(\cdot), Q_{\bar{\theta}_\ell^k}$	The joint $Q$ function and joint target $Q$ function of team $k$ , parameterized by $\theta_\ell^k$ and $\bar{\theta}_\ell^k$ respectively.
$y^k$	The TD backup of team $k$ .
$\mathcal{L}_\ell^Q(\{\theta_\ell^i\}, \{\theta_\ell^k\})$	The loss function of decision module.
$\lambda_{QMIX}$	The temperature paramter of loss function of decision module.

**Table 8: Notations in the decision module.**