

GROUP-ORIENTED COOPERATION IN MULTI-AGENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Grouping is ubiquitous in natural systems and is essential for promoting efficiency
2 in team coordination. This paper introduces the concept of grouping into multi-
3 agent reinforcement learning (MARL) and provides a novel formulation of Group-
4 oriented MARL (GoMARL). In contrast to existing approaches that attempt to
5 directly learn the complex relationship between the joint action-values and indi-
6 vidual values, we empower groups as a bridge to model the connection between a
7 small set of agents and encourage cooperation among them, thereby improving the
8 efficiency of the whole team. In particular, we factorize the joint action-values as
9 a combination of group-wise values, which guide agents to improve their policies
10 in a fine-grained fashion. We propose a flexible grouping mechanism inspired by
11 variable selection and sparse regularization to generate dynamic groups and group
12 action-values. We further propose a hierarchical control for policy learning that
13 drives the agents in the same group to specialize in similar policies and possess
14 diverse strategies for various groups. Extensive experiments on a challenging set
15 of StarCraft II micromanagement tasks and Google Research Football scenarios
16 verify our method’s effectiveness and learning efficiency. Detailed component
17 studies show how grouping works and enhances performance.

18 1 INTRODUCTION

19 Cooperative multi-agent reinforcement learning (MARL) aims to coordinate multiple agents’ ac-
20 tions through shared team rewards and has become a helpful tool for solving multi-agent decision-
21 making problems, such as network routing (Ye et al., 2015), robot swarm control (Hüttenrauch et al.,
22 2017), crewless aerial vehicles (Xu et al., 2018), etc. Learning centralized policies is a natural way
23 to address the cooperative MARL problem. It treats the team as a single actor with a joint action
24 space. Although single-agent RL algorithms can be trivially transplanted to this setting, the global
25 information to train the model is usually unavailable during execution due to partial observability
26 or communication constraints, and the joint action space grows exponentially with the agent num-
27 ber (Gupta et al., 2017; OroojlooyJadid & Hajinezhad, 2019; Gronauer & Diepold, 2021). An alter-
28 native paradigm is to learn decentralized policies (Tan, 1993; Witt et al., 2020; Yu et al., 2021) by
29 independently training agents based on their local observations. However, simultaneous exploration
30 brings non-stationarity that causes unstable learning and convergence difficulties (Hernandez-Leal
31 et al., 2019; Zhang et al., 2021). The centralized training with decentralized execution (CTDE)
32 paradigm inherits the advantages of the above two paradigms and allows learning decentralized
33 policies in a centralized fashion (Oliehoek et al., 2008; Kraemer & Banerjee, 2016).

34 Although the CTDE paradigm solves many multi-agent problems and opens up the possibility for
35 agents to share information during training, it still faces two challenges. On the one hand, learning
36 efficient team cooperation by directly estimating the joint action-values from individual utilities is
37 exceptionally difficult, especially when the agent number is enormous. Most value function factor-
38 ization approaches have been only evaluated in domains with a handful of agents. The flat factoriza-
39 tion scheme is proven to lead to performance bottleneck (Phan et al., 2021), where it gets difficult to
40 provide sufficiently informative training signals for each agent. Although the state information pro-
41 vides the complete knowledge needed for learning, it is burdensome for agents to extract effective
42 guidance that facilitates cooperation. On the other hand, the guideline of sharing parameters is still
43 an open question. The full parameter-sharing mechanism limits the diversity of agents’ behavior
44 strategies (Li et al., 2021), leading agents’ policies to be similar or the same.

Both nature (Jeanson et al., 2005; Wittemyer & Getz, 2007) and MARL (Phan et al., 2021) have validated *grouping* as a means to promote efficient cooperation and break performance bottlenecks. MARL method VAST (Phan et al., 2021) approximates a factorization for agent sub-teams to overcome the performance bottleneck caused by the flat value factorization which directly assigns a centralized value function to each agent. Although grouping provides new possibilities for MARL, how to formulate a general grouping criterion for complex and diverse environments without any domain knowledge is still an open question and a matter of great interest. Most previous grouping works are proposed for well-structured tasks, *e.g.*, software engineering (Pavón & Gómez-Sanz, 2003; Bresciani et al., 2004), and typically predefine specific responsibilities or forms of task decomposition. VAST, on the other hand, artificially sets the team amount to half or a quarter of the number of agents. These methods all require apriori knowledge or settings that are potentially unavailable or unreasonable in practice and may discourage methods’ transferring to diverse environments.

To address the above challenges faced by the CTDE paradigm and group learning, we propose a novel formulation of Group-oriented MARL (GoMARL), a method for learning dynamic grouping without any domain knowledge or a priori setting. Instead of formulating a concrete grouping criterion, GoMARL implicitly learns dynamic groups with a novel sparsity-driven scheme. Concretely, it learns a dual hierarchy of value function factorization, where the learning weights of the decomposition from the group value to local utilities determine whether an agent is suitable for the current group. A sparsity regularization drives this end-to-end automatic group learning with the Least Absolute Shrinkage and Selection Operator (LASSO) (Tibshirani, 1996). Furthermore, GoMARL relies on well-designed architecture to transform the individual, group-related, and global information into the weights of corresponding networks in a manner reminiscent of hypernetworks (Ha et al., 2017), which flexibly adapts to the changes in group number and group size.

To alleviate the first problem faced by the CTDE paradigm, GoMARL adopts the dual-hierarchy value factorization to learn on a more focused and compact input representation. Different from VAST and other methods that learn value factorization with no information (learn blindly) or only with the state information (complete but hard to extract efficient guidance), we propose a fine-grained learning scheme to integrate information from the group perspective into the policy gradient by hypernetworks. It promotes *intra-group coordination with group information* and facilitates *inter-group cooperation with the global state* to further break the performance bottleneck. As for the second problem, parameter sharing prevents efficient cooperation due to policy similarity. However, our dynamic grouping naturally serves as an intermediary to encourage diversity while sharing parameters. We introduce a latent space to describe group information for hierarchical control to establish the connection between group and policy. Agents condition their behaviors on their group information embedded by a shared encoder, which is learned following specialization guidance, *i.e.*, imposing similarity within a group and diversity between groups. In this way, GoMARL synergizes groups with specialized policies, providing a *fully-sharing mechanism for learning diverse policies*.

We summarize our main contributions in this paper as follows:

1. We introduce an end-to-end adaptive group learning MARL solution with no domain knowledge or a priori setting, partitioning agents into dynamic groups to promote efficient cooperation.
2. Structurally, we present an architecture to estimate global and group-wise joint action-values with varying group numbers and agent numbers per group, realizing a dynamic end-to-end training.
3. Functionally, we propose agents that condition their behavior on latent group information to achieve group specialization and diversity while sharing all their parameters.

We test our method on a challenging set of StarCraft II micromanagement tasks (Samvelyan et al., 2019) and Google Research Football (Kurach et al., 2020) scenarios. GoMARL achieves superior performance with higher efficiency compared with notable baseline methods. We also provide detailed ablation studies to give insights into how grouping works and enhances learning performance.

2 BACKGROUND AND PRELIMINARIES

Dec-POMDP. This paper focuses on cooperative tasks with n agents $\mathcal{A} = \{a_1, \dots, a_n\}$ as a Dec-POMDP (Oliehoek & Amato, 2016) defined by a tuple $G = \langle S, U, P, r, Z, O, n, \gamma \rangle$. The environment has a true state $s \in S$. Each agent a chooses an action u_t^a from its action space U_a at timestep t and forms a joint action $\mathbf{u}_t \in (U_1 \times \dots \times U_n) \equiv U^n$ that induces a transition in the environment according to the state transition function $P(s_{t+1}|s_t, \mathbf{u}_t) : S \times U^n \times S \rightarrow [0, 1]$. $r(s, \mathbf{u}) : S \times U^n \rightarrow \mathbb{R}$

99 is the *reward* function yielding a global reward, and $\gamma \in [0, 1)$ is the discount factor. We consider
 100 partially observable scenarios in which agent a acquires its local *observation* $z^a \in Z$ drawn from
 101 $O(s_t, a) : S \times \mathcal{A} \rightarrow Z$. Each agent has an *action-observation history* $\tau^a \in T \equiv (U \times Z)^*$, on
 102 which it conditions a *policy* $\pi^a(u^a | \tau^a) : T \times U \rightarrow [0, 1]$. We denote joint quantities over agents in
 103 bold and joint quantities over agents other than a given agent a with the superscript $-a$.

104 **Value Function Factorization.** We consider cooperative MARL with centralized training and de-
 105 centralized execution paradigm, which has been a major focus in recent efforts (Foerster et al., 2018;
 106 Sunehag et al., 2018; Rashid et al., 2018; Iqbal & Sha, 2019; Mahajan et al., 2019). Some methods
 107 achieve CTDE through value function factorization, *i.e.*, factoring action-value functions into com-
 108 binations of per-agent utilities. The individual utility only depends on the local history of actions
 109 and observations, allowing agents to maximize their local utility functions independently. Among
 110 these attempts, the representative deep MARL approach QMIX (Rashid et al., 2018) improves the
 111 simple summation of individual utilities (Sunehag et al., 2018) by introducing a more expressive fac-
 112 torization: $Q^{tot} = f(Q^1(\tau^1, u^1; \theta_Q), \dots, Q^n(\tau^n, u^n; \theta_Q); \theta_f)$, where θ_f denotes the parameters of
 113 the monotonic mixing function generated by a hypernetwork (Ha et al., 2017).

114 **Variable Selection in Regression.** Finding critical explanatory variables (factors) in predicting the
 115 response variable is vital when solving regression problems. Each explanatory factor may be rep-
 116 resented by a group of derived input variables. Consider the general regression problem with M
 117 factors: $Y = \sum_{m=1}^M X_m \beta_m + \epsilon$, where Y is an $n \times 1$ response vector, X_m is an $n \times p_m$ matrix
 118 corresponding to the m -th factor, β_m is a coefficient vector of size p_m , and ϵ is a perturbation that
 119 follows a Gaussian distribution. The most considered model selection problem is a particular case
 120 when $p_1 = \dots = p_m = 1$. Model selection methods have been widely introduced (George &
 121 McCulloch, 1993; Shen & Ye, 2002; Efron et al., 2004). Among them, the Least Absolute Shrink-
 122 age and Selection Operator (LASSO) (Tibshirani, 1996) is the renowned and classic one, which is
 123 defined as $\hat{\beta}^{\text{LASSO}}(\lambda) = \arg \min_{\beta} (\|Y - X\beta\|^2 + \lambda \|\beta\|_{l_1})$. $\|\cdot\|_{l_1}$ stands for the vector l_1 -norm
 124 penalty *inducing sparsity in the solution*, and λ is a tuning parameter.

125 3 GROUP-ORIENTED LEARNING FRAMEWORK

126 This section introduces the Group-oriented Multi-agent Rein-
 127 forcement Learning (GoMARL) framework that integrates the
 128 group notion with MARL in a principled manner. Following
 129 the value function factorization approaches that represent the
 130 global action-value as an aggregation of individual utilities,
 131 GoMARL decomposes the joint action-value into group-wise
 132 values and trains agents by groups in a fine-grained scheme.
 133 Figure 1 illustrates the overview of the group-oriented learn-
 134 ing framework. It consists of an automatic grouping module,
 135 specialized agent networks generating local utilities $Q^i(\tau^i, \cdot)$,
 136 and a mixing network among groups. In Section 3.1, we intro-
 137 duce the automatic grouping module. It progressively divides
 138 the team into dynamic groups as training proceeds. Based on
 139 the grouping, we propose specialized agent networks that achieve similarity within each group and
 140 diversity among groups to generate the local Q^i in Section 3.2. Section 3.3 further presents the
 141 mixing of Q^i to estimate the group-wise and global action-values and the overall training procedure.

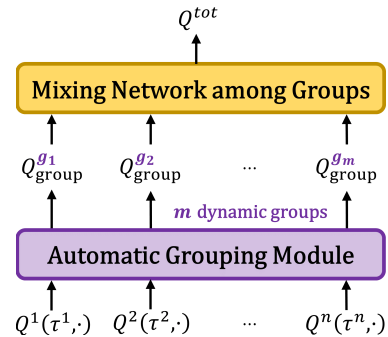


Figure 1: Overview of GoMARL.

142 3.1 AUTOMATIC GROUPING

143 **Definition 1** (Individual and Group). Given a cooperative task with n agents $\mathcal{A} = \{a_1, \dots, a_n\}$,
 144 we have a set of groups $\mathcal{G} = \{g_1, \dots, g_m\}$, $1 \leq m \leq n$. Each group g_i contains n_i ($1 \leq n_i \leq n$)
 145 different agents $g_i = \{a_1^i, \dots, a_{n_i}^i\} \subseteq \mathcal{A}$, where $g_i \cap g_j = \emptyset$, $i \neq j$, $\cup_i g_i = \mathcal{A}$, and $i, j \in [1, m]$. In
 146 this paper, we denote all variables with a *superscript* to describe the variable owner, *e.g.*, e^i and s^{g_j}
 147 are variables of agent a_i and group g_j , respectively.

148 The automatic grouping module aims to learn a mapping relationship $f_g : \mathcal{A} \mapsto \mathcal{G}$. Predefined
 149 explicit formulations and apriori knowledge are not necessities, as they are unavailable in complex
 150 environments. **Our key idea is to divide the system into dynamic groups in an end-to-end fashion**

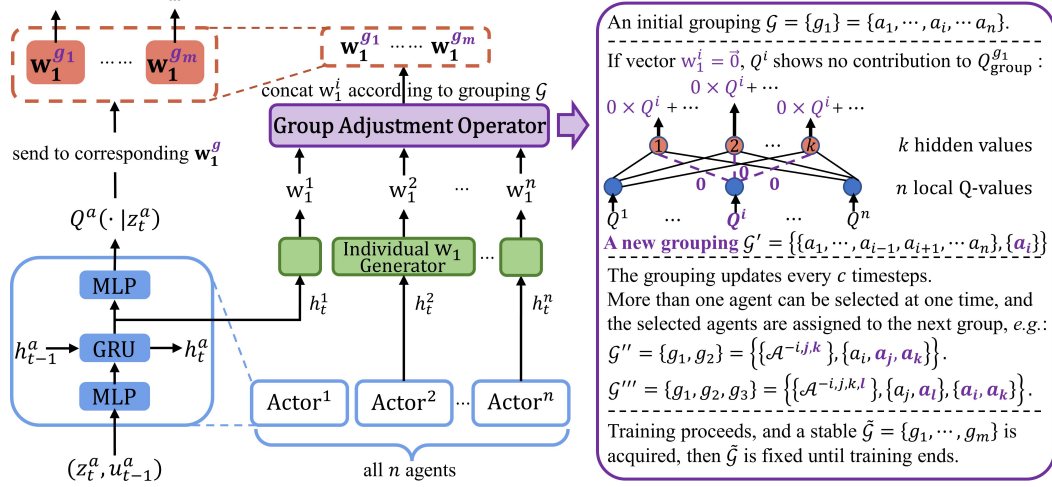


Figure 2: Schematic diagram of automatic grouping. The right side shows how grouping \mathcal{G} changes during training. In particular, we select the agents whose Q^i have little contribution to $Q_{\text{group}}^{g_j}$ to move out of the current group g_j . Based on the grouping \mathcal{G} , the group selection operator concatenates the weights w_1^i of agents in the same group to form the group-wise weights \mathbf{w}_1^g for mixing local utilities.

151 by maximizing the expected global return $Q_G^{\text{tot}}(s_t, \mathbf{u}_t) = \mathbb{E}_{s_{t+1:\infty}, \mathbf{u}_{t+1:\infty}} [\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t, \mathbf{u}_t; \mathcal{G}]$.
 152 Value function factorization approaches represent the joint action-value as an aggregation of the
 153 individual values, *i.e.*, a weighted sum of Q^i and biases. If the learned weights are restricted to be
 154 non-negative, as in the monotonic mixing (Rashid et al., 2018), each weight reflects the contribution
 155 of Q^i to Q^{tot} . We follow this setting and represent the group-wise Q_{group}^g as an aggregation of the
 156 individual values. Intuitively, if the learned weight of Q^i is small enough, agent a_i contributes
 157 a little to its current group. In other words, when an agent a_i takes action $u^i = \arg \max_u Q^i(\tau^i, u)$
 158 but does not contribute to its group value, it indicates that agent a_i does not belong to its current
 159 group. The right side of Figure 2 illustrates the schematic diagram of the grouping mechanism. In
 160 the beginning, all agents belong to the same group, and the grouping \mathcal{G} is gradually adjusted as the
 161 training proceeds. It is worth noting that the proposed dynamic grouping guideline does not require
 162 a predefined group number like VAST or any other domain knowledge.

163 GoMARL flexibly utilizes hypernetworks to accommodate the changeable group size (agent number
 164 per group). We construct individual w_1 generators $f_w^i(\cdot; \theta_{w_1}^i) : \tau^i \rightarrow w_1^i$ that map each agent a_i 's
 165 hidden state h^i with history information τ^i to a k -dimensional weight vector. All agents' w_1 are sent
 166 to a group adjustment operator \mathcal{O}_g and decide on a new grouping based on the proposed grouping
 167 guideline. According to the learned grouping \mathcal{G} , \mathcal{O}_g concatenates the w_1^i of agents in the same group
 168 to form a set of group-wise \mathbf{w}_1^g , *i.e.*, $\mathcal{O}_g : \{w_1^1, \dots, w_1^n\} \xrightarrow{\mathcal{G}} \{\mathbf{w}_1^{g_1}, \dots, \mathbf{w}_1^{g_m}\}$. This fixed network
 169 architecture can adapt to the grouping dynamics since each w_1^i is tied to agent a_i with hypernetwork
 170 f_w^i . No matter which group g_j agent a_i belongs to, Q^i engages in the mixing of $Q_{\text{group}}^{g_j}$ through w_1^i .

171 We achieve automatic group learning by proposing a sparsity-driven scheme. Specifically, GoMARL
 172 selects agents whose utilities contribute a little to their current group-values and adjusts their group-
 173 ing; thus, a regularization for sparsity on the w_1 generators $f_w^i(\cdot; \theta_{w_1}^i)$ is necessary. Applicable
 174 regularizers are various, and we provide a straightforward but effective attempt. We choose LASSO
 175 from among many feasible sparse regularizers due to its simplicity and not introducing additional
 176 learnable parameters. The w_1 generators are then trained by minimizing the following loss function:
 177

$$\mathcal{L}_g(\theta_{w_1}) = \mathbb{E}_{(z, \mathbf{u}, r, z') \sim \mathcal{B}} \sum_i \left(\|f_w^i(\tau^i(z^i, u^i); \theta_{w_1}^i)\|_{l_1} \right), \quad (1)$$

178 where \mathcal{B} is a replay buffer, and $\|\cdot\|_{l_1}$ stands for LASSO's l_1 -norm penalty. In practice, the weights'
 179 shrinkage to zero is infeasible. It is also challenging to determine a fixed threshold for groups
 180 of various sizes in diverse environments. We empirically utilize seventy percent of each group's
 181 average contribution (weights) to assess whether an agent is suitable for its current group. The
 182 experiments in Section 4 show that this setting is generic to different scenarios and testbeds. The
 183 grouping shifts every c timesteps, and each selected agent is assigned to the next group until it
 184 properly contributes to where it belongs. Appendix A shows insight into this grouping shift.

185 3.2 SPECIALIZED AGENT NETWORKS SHARING ALL THE PARAMETERS

186 Policy decentralization with shared parameters is widely utilized to improve scalability and learning
 187 efficiency. However, agents tend to behave similarly when sharing parameters, preventing effective
 188 exploration and complex cooperative policies. In addition, it is also undesirable to entirely forgo
 189 shared parameters in pursuit of diverse strategies since proper sharing accelerates learning. To make
 190 full use of the grouping, we introduce a hierarchical control for policy learning that drives the agents
 191 in the same group to specialize in similar policies and possess diverse strategies across groups.

192 As shown in Figure 3, we construct an agent
 193 info encoder $f_e(\cdot; \theta_e)$ to embed agents' hid-
 194 den states. The acquired agent info summarizes
 195 *agent's history from the group perspective*. To
 196 achieve this group-related view, we train the en-
 197 coder network to be an extractor by a group spe-
 198 cialization regularizer, where the agent info of
 199 agents from the same group is similar. To avoid
 200 all agents' info collapsing in a similar manner,
 201 the regularizer also encourages diversity between
 202 agents from different groups. Formally, we
 have the following Similarity-Diversity objective to train the info encoder:

$$\begin{aligned} \text{minimize } \mathcal{L}_{SD}(\theta_e) &= \mathbb{E}_{\mathcal{B}} \left(\sum_{i \neq j} I(i, j) \cdot \text{cosine}(f_e(h^i; \theta_e), f_e(h^j; \theta_e)) \right), \\ \text{where } I(i, j) &= \begin{cases} -1, & a^i, a^j \in g^k. \\ 1, & a^i \in g^k, a^j \in g^l, k \neq l. \end{cases} \end{aligned} \quad (2)$$

203 The SD-loss trains the encoder to extract agent info e^i that is recognizable to agents' group. This
 204 identifiable agent info is then fed into a decoder network $f_d(\cdot; \theta_d)$ to generate the parameters of
 205 the agent network's upper MLP. This decoder hypernetwork is trained by the TD-loss introduced in
 206 Section 3.3. In this way, the value-based agents condition their behavior on their agent info with
 207 group-related information embedded, achieving specialized policies through hierarchical control.

208 The proposed agent network has two merits. First, different from promising approaches like CDS (Li
 209 et al., 2021) that promote diversity by sharing only a fraction of the network, GoMARM enables
 210 diversified policies while *sharing all the parameters*, hybridizing the efficient learning of parameter
 211 sharing and diversity for complex cooperation. Second, we utilize hypernetworks f_d to integrate
 212 the informative extracted agent info e^{a_i} into the gradients to provide group-related information.
 213 Unlike all previous methods that either learn with no information (learn blindly) or only with state
 214 s (complete but hard to extract effective guidance), GoMARM achieves fine-grained learning with
 215 *agent info from the group perspective* in the gradient for efficient policy learning. Concretely, the
 216 partial derivative $\frac{\partial Q^{tot}}{\partial \theta_\pi}$ for updating the policy parameters θ_π of the GRU and the bottom MLP is:

$$\frac{\partial Q^{tot}}{\partial \theta_\pi} = \frac{\partial Q^{tot}}{\partial Q^a} \frac{\partial Q^a}{\partial \theta_\pi} = \frac{\partial Q^{tot}}{\partial Q^a} \frac{\partial Q^a}{\partial v^a} \frac{\partial v^a}{\partial \theta_\pi} = f_d(e^i) \cdot \frac{\partial Q^{tot}}{\partial Q^a} \frac{\partial v^a}{\partial \theta_\pi}, \quad (3)$$

217 where v is the representation after the GRU. Eqn.(3) shows that e^i is deeply involved in agent a_i 's
 218 policy updating, facilitating group-related guidance on policy learning for better cooperation.

219 3.3 OVERALL LEARNING FRAMEWORK

220 We next introduce the estimation of the group-wise Q_{group}^g and global joint action-value function
 221 Q^{tot} . Section 3.1 presents the generation of each agent's w_1^i that determines the grouping \mathcal{G} . Al-
 222 though w_1^i enables a weighted mixing of local utilities of agents in group g_i to generate the group
 223 action-value $Q_{\text{group}}^{g_i}$, the naive mixture of $\sum_{a_j \in g_i} w_1^j Q^j + b$ lacks the guidance of group-related in-
 224 formation to reflect the action-value in a *specific group-wise state*. Therefore, we build a two-layer
 225 mixing structure that embeds the *group state* into the weights w_2^g of the second layer to generate
 226 Q_{group}^g . In particular, group g_i 's state s^{g_i} is a *fusion* of the agent info e^j presented in Section 3.2 for
 227 all $a_j \in g_i$. To cohesively summarize the group state based on the agent info of all agents in this
 228 group, we apply the max pooling operation (Ranzato et al., 2007) over each dimension of the agent
 229 info e^{g_i} to generate the group state s^{g_i} describing the current group status. The pooling operation
 230 also ensures adaptability to the changeable agent number per group. We further build a group-wise
 231 w_2 generator to map the group state into the weights of the second mixing layer. *Equipped with this*

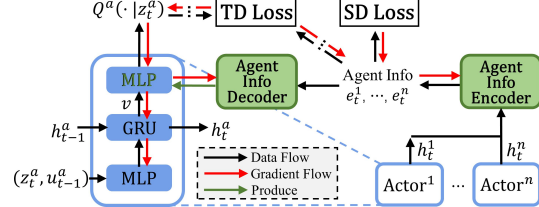


Figure 3: Architecture of the agent networks.

232 w_2 generator hypernetwork, s^g is integrated into the gradients to provide group-related information
 233 for mixing Q_{group}^g (similar to Eqn.(3)) and facilitates efficient *intra-group coordination*.

234 The overall learning framework of GoMARL is
 235 shown in Figure 4, and the bottom part illus-
 236 trates the generation of s^g and w_2 . The two-
 237 layer mixing estimates the group-wise value by
 238 both w_1 that decides the grouping and w_2 that
 239 carries group state information. Just as our
 240 fixed network can dynamically adapt to vari-
 241 ous group sizes, the hypernetwork w_2 gener-
 242 ator ties each group g_j to its $w_2^{g_j}$ and enables
 243 our architecture to flexibly adapt to the group
 244 number changes as well. We estimate the joint
 245 value Q^{tot} by mixing all the Q_{group}^g in a similar fashion. Concretely, the two layers of the total mix-
 246 ing network are generated by two hypernetworks, respectively taking group states s^{g_i} and the global
 247 state s as inputs. Similar to Eqn.(3), the group and global state are deeply involved in the value gra-
 248 dient and guide the *inter-group cooperation*. Each layer’s biases are produced in the same manner as
 249 the corresponding weights. The architecture of the total mixing network is akin to the group mixing
 250 network and is omitted in Figure 4. The formulation gives the TD-loss of the estimated Q^{tot} :

$$\mathcal{L}_{TD}(\theta) = \mathbb{E}_{\mathcal{B}} \left[\left(r + \gamma \max_{\mathbf{u}'} \bar{Q}^{tot}(s', \mathbf{u}') - Q^{tot}(s, \mathbf{u}) \right)^2 \right], \quad (4)$$

251 where \bar{Q}^{tot} is a target network with periodic updates. The overall learning objective is:

$$\mathcal{L}(\theta) = \mathcal{L}_{TD}(\theta) + \lambda_g \mathcal{L}_g(\theta_{w_1}) + \lambda_{SD} \mathcal{L}_{SD}(\theta_e), \quad (5)$$

252 where $\theta = (\theta_\pi, \theta_e, \theta_d, \theta_w)$, θ_w denotes the parameters of hypernetworks producing all mixing
 253 weights and biases, λ_g and λ_{SD} are two scaling factors.

254 Although containing two mixing networks to respectively estimate Q_{group}^g and Q^{tot} , GoMARL’s total
 255 mixing-net size is quite close to or even smaller than the size of QMIX’s single mixing network, as
 256 verified in experiments. This is mainly attributed to the input dimension reduction of the weights
 257 generator hypernetworks. The commonly used QMIX’s mixing network takes the global state s as
 258 the input of all the hypernetworks. In contrast, we take specific information of the grouping (*i.e.*,
 259 agents’ hidden states h , group state s^g , and the global state s is only used in the top hypernetwork)
 260 as inputs. This parameter reduction offsets the increase of an extra mixing network. Compared
 261 with QMIX’s mixing only with the state information, our dual-mixing structure allows fine-grained
 262 mixing of Q_i or Q_{group}^g to be guided by more detailed information (*i.e.*, local history, group state,
 263 and global state) embedded in the gradients, facilitating intra- and inter-group coordination. We
 264 further prove that GoMARL’s dual-hierarchy factorization maintains decentralizability by satisfying
 265 the IGM (Individual-Global-Max) for Q^{tot} and Q^i for arbitrary grouping \mathcal{G} in Appendix B.

266 4 EXPERIMENTS

267 **Baselines.** We compare GoMARL with prominent baselines to verify its effectiveness and effi-
 268 ciency. Hu et al. (2021) fairly compared existing MARL methods without code-level optimizations
 269 and reported that QMIX (Rashid et al., 2018) and QPLEX (Wang et al., 2021a) are the top two
 270 value function factorization methods. The authors also finetuned QMIX (denoted as Ft-QMIX in
 271 our paper), which attains higher win rates than the vanilla QMIX. Therefore, we compare GoMARL
 272 with Ft-QMIX, and QPLEX to show its performance as a value factorization method (QMIX is
 273 also shown). Besides, the baselines also include role-based methods (ROMA (Wang et al., 2020),
 274 RODE (Wang et al., 2021b)) and the representative credit assignment method RIIT (Hu et al., 2021).
 275 The latter combines effective modules of noticeable methods and has recently gotten much attention.

276 **Experimental setup.** All methods are trained with 8 parallel runners for 10M steps. We evaluate
 277 them every 10K steps with 32 episodes and report the 1st, median, and 3rd quartile win rates across
 278 5 random seeds. The detailed settings, including fair comparison, are introduced in Appendix C.1.

279 4.1 PERFORMANCE ON STARCRAFT II MICROMANAGEMENT TASKS

280 All methods are evaluated on a set of challenging SMAC maps that vary in difficulty by *Hard*
 281 (3s_vs_5z, 5m_vs_6m, 8m_vs_9m) and *Super Hard* (corridor, MMM2, 3s5z_vs_3s6z). *Easy*

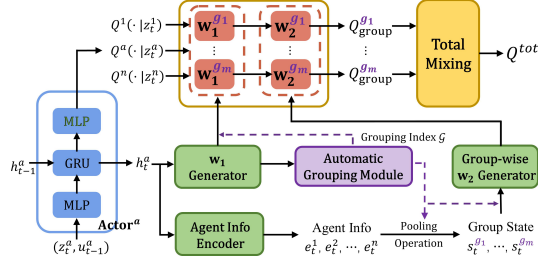


Figure 4: Overall learning framework.

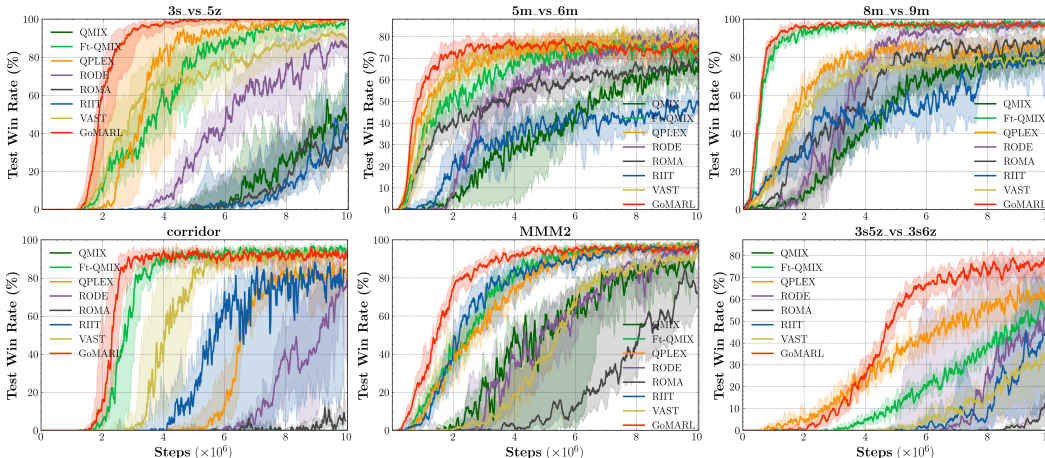


Figure 5: Performance comparison of GoMARL with baseline methods on SMAC. (Update VAST)

282 scenarios are not chosen as they can be easily solved, and the performances of all methods are close.
 283 The chosen tasks involve homogeneous and heterogeneous teams with asymmetric battles, allowing
 284 a holistic study of all methods. Appendix C.3 introduces the traits of these scenarios.

285 **Parameter size for value mixing.** GoMARL maintains two mixing networks to respectively es-
 286 timate Q_{group}^g and Q^{tot} . If GoMARL has more parameters for value mixing is a natural question.
 287 Appendix C.2 gives the parameter sizes for value mixing of all compared methods. When the agent
 288 amount exceeds 5, our dual-hierarchy architecture has the least parameters among all the baselines.

289 **Overall Performance.** The comparison of GoMARL against baseline algorithms on the SMAC
 290 tasks is shown in Figure 5. As the results show, each baseline method only achieves satisfactory
 291 performance on some of the challenging benchmarks with specific properties they specialize in,
 292 e.g., RIIT performs well on MMM2 but converges much slower in other tasks, QPLEX’s leaning is not
 293 efficient in 8m_vs_9m and corridor. RODE attains better overall performance than ROMA, but
 294 its hyperparameters are sensitive to scenarios; thus, its learning is not stable as QPLEX. Finetuned
 295 QMIX significantly over-performs the vanilla QMIX and has more efficient learning than other
 296 baselines. Our method has a similar win rate as Ft-QMIX in 8m_vs_9m and corridor. However,
 297 its superiority in both efficiency and effectiveness can be clearly validated in all the other challenging
 298 tasks. Next, we conduct detailed ablation studies and component analyses on GoMARL’s modules
 299 to illustrate how our method improves learning efficiency and enhances performance.

300 **Ablation study and component analysis.** GoMARL contains two key components: (1) an auto-
 301 matic grouping mechanism that progressively divides the team into proper groups as training pro-
 302 ceeds; and (2) specialized agent networks that generate diversified policies according to the group-
 303 ing while sharing all the parameters. We conduct ablation studies on the three *Super Hard* tasks
 304 (corridor and MMM2, and 3s5z_vs_3s6z) to show how each module influences performance.

305 (1) *Ablation of the proposed grouping mechanism.* We validate our grouping mechanism by com-
 306 paring GoMARL’s dynamic grouping with other intuitive groupings in the top row of Figure 6. All
 307 compared methods utilized the same architecture as GoMARL to reflect how grouping itself influ-
 308 ences performance. Setting all agents as a group in corridor converges faster but has significant
 309 variances, as the shadow illustrates. This may be due to the hard exploration of efficient cooper-
 310 ation without grouping guidance. Another two intuitive groupings, each agent a group and an equal
 311 division into two groups $\{\{a_1, a_2, a_3\}, \{a_4, a_5, a_6\}\}$, have minor variance. However, their learning
 312 efficiency is affected since inappropriate groupings fail to promote cooperative behaviors. MMM2
 313 contains heterogeneous agents; thus, a natural grouping is to keep the agents of the same type in a
 314 group. As the results show, this natural grouping over-performs setting each agent in a group and
 315 all agents in a group. On the other hand, our mechanism dynamically adjusts the grouping and con-
 316 verges faster. Setting all agents as a group in 3s5z_vs_3s6z also has a great variance. It is difficult
 317 for agents to learn complex cooperation if each is set to be a group, as the bad performance in both
 318 win rate and efficiency shows. 3s5z_vs_3s6z is another scenario with heterogeneous agents, and
 319 the natural grouping of homogeneous agents as a group is studied. It performs nearly the same as
 320 our approach because our dynamic grouping mechanism learns the grouping exactly according to
 321 the agents’ type in this scenario. We also analyze this map with a grouping containing three groups
 322 $\{\{a_1, a_4, a_5\}, \{a_2, a_6, a_7\}, \{a_3, a_8\}\}$ (1s2z, 1s2z, 1s1z); however, this balanced grouping fails to

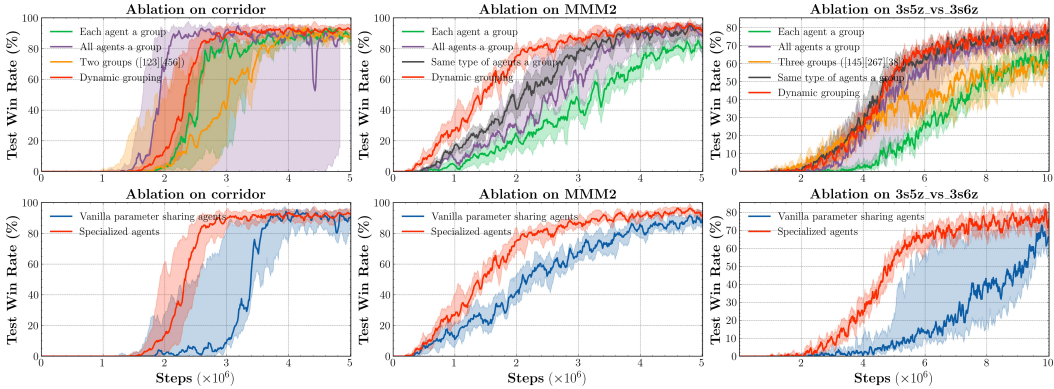


Figure 6: Ablations of the grouping mechanism (top) and specialized agent networks (bottom).



Figure 7: The final learned policies that fit the grouping $\{\{a_1, a_2\}, \{a_3, a_4, a_5, a_6\}\}$ on corridor.

323 form effective cooperation. We can see from these results that multi-agent systems with the same
 324 learning framework perform diversely with various grouping. Appropriate grouping facilitates ef-
 325 ficient cooperation and accelerates learning. The proposed grouping mechanism can automatically
 326 learn adaptive grouping in different tasks and assist GoMARL with superior and stable performance.

327 (2) *Learned grouping analysis.* To better demonstrate whether the learned grouping makes sense,
 328 we further visualize the final trained strategy in one corridor battle, as illustrated in Figure 7.
 329 Six allied Zealots fight twenty-four Zerglings on this super-hard map. The massive disparity in
 330 unit numbers between the two sides implies that the whole team cannot launch an attack together.
 331 The only winning strategy is to sacrifice a small number of agents who leave the team and attract the
 332 attention of most enemies. Taking this opportunity, our large force eliminates the rest of the enemies.
 333 The surviving agents then use the same tactic to attract several enemies to the team every time and
 334 kill them together. In our visualization, Agent 1 and Agent 2 sacrifice themselves to attract most
 335 enemies and bring enough time for the team to eliminate the remaining enemies. Other agents fight
 336 as a small group and successfully kill all the surviving enemies. Our dynamic grouping mechanism
 337 learns a two-group setting in this battle, where Agent 1 and Agent 2 are in the same group while the
 338 others are set in another group. This grouping is explicable in light of the combat situation, and this
 339 reasonable grouping guidance contributes to the superior performance of our approach.

340 (3) *Ablation of the specialized agent networks.* As shown in the bottom row of Figure 6, our spe-
 341 cialized agent networks greatly improve the learning efficiency. Although GoMARL’s agents also
 342 share all the parameters like vanilla parameter-sharing agents, our proposed info encoder embeds
 343 group-related information into the agent info with similarity and diversity regularizer. The agent
 344 info decoder further produces each agent’s upper MLP to generate diversified policies. Compared to
 345 the vanilla paramete-sharing mechanism that limits diversity, our agents have various styles of be-
 346 haviors related to their group, which encourages extensive exploration and accelerates learning. In
 347 addition, our agents share all the parameters (GRU, bottom MLP, agent info encoder, and decoder)
 348 to obtain their policies, preserving the advantages in scalability of the vanilla parameter-sharing
 349 mechanism. Detailed studies on the inner component of λ_{SD} and the performance improvement of
 350 baseline methods with our agents are shown in Appendix D to further demonstrate the effectiveness.

351 4.2 PERFORMANCE ON GOOGLE RESEARCH FOOTBALL

352 We also test GoMARL on two challenging Google Research Football Academy offensive scenarios,
 353 3_vs_1_with_keeper and counterattack_easy. Agents in this environment need to coordi-
 354 nate timing and positions for organizing offense to seize fleeting opportunities, and only scoring
 355 leads to rewards. The tasks in GRF are far more difficult than those in SMAC, making many meth-
 356 ods that work well in SMAC invalid. Therefore, a secondary test on GRF is better proof of our
 357 approach’s effectiveness. Appendix C.4 details the basic information about the GRF environment.

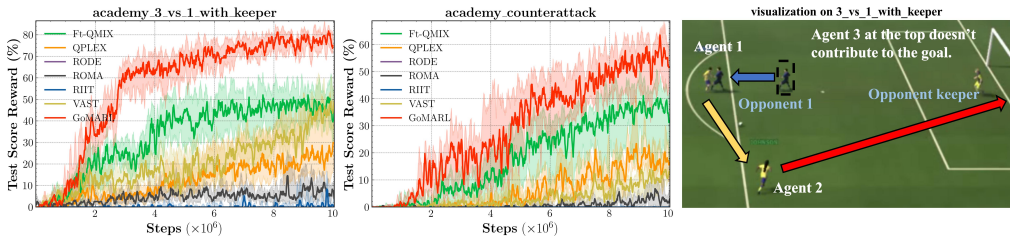


Figure 8: Comparison on GRF tasks (Update VAST) and a visualization of 3_vs_1_with_keeper.

358 **Performance.** The left side of Figure 8 shows the performance comparison of all methods. Only
 359 Ft-QMIX and GoMARL achieve more than 30% of the score reward in both scenarios. QPLEX
 360 and ROMA only score with a small probability, while other methods fail to give a single goal.
 361 GoMARL maintains its superior overall performance with outstanding learning efficiency. The
 362 excellent performance in the second testbed further demonstrates the transferability of our method.

363 **Visualizations.** The trained strategies are visualized to check if the learned grouping makes sense.
 364 Figure 8 (right) visualizes the 3_vs_1_with_keeper, and the visualization of counterattack
 365 is illustrated in Appendix E. The cooperation of two players is enough to score a goal in
 366 3_vs_1_with_keeper. As illustrated, Agent 1 passes the ball to Agent 2 (shown in yellow arrow)
 367 when the opponent player rushing (blue arrow) to tackle. GoMARL agents master the cooperative
 368 timing and position, and Agent 2 smoothly receives the ball and shoots at the best timing (red arrow).
 369 Our method learns the grouping of $\{\{a_1, a_2\}, \{a_3\}\}$ that is explicable to the game.

370 5 RELATED WORK

371 This paper focuses on cooperative MARL problems with the CTDE paradigm (Oliehoek & Amato,
 372 2016). GoMARL utilizes value function factorization, an approach aiming to address the multi-
 373 agent credit assignment problem. Early attempts at value function factorization (Schneider et al.,
 374 1999; Russell & Zimdars, 2003) require apriori knowledge to predefine specific responsibilities
 375 or design suitable team reward decompositions. Deep MARL methods learn value factorization
 376 with no domain knowledge by treating agents as independent factors. VDN (Sunehag et al., 2018)
 377 learns a linear decomposition into a sum of local utility functions used for greedy action selection.
 378 QMIX (Rashid et al., 2018) enlarges the functions that can be represented by the mixing network
 379 but still faces the monotonicity limitation. QTRAN (Son et al., 2019) further improves the ex-
 380 pressivity by using constraints between individual utilities and the global action-value; however,
 381 the constraints are computationally intractable, and the relaxations often lead to unsatisfied perfor-
 382 mances. MAVEN (Mahajan et al., 2019) learns latent embeddings to integrate a diverse ensemble of
 383 monotonic approximations. LICA (Zhou et al., 2020) learns end-to-end differentiable policy opti-
 384 mization to remove the expressivity constraint. VMIX (Su et al., 2021) combines A2C with QMIX
 385 to extend the monotonicity constraint to value networks. QPLEX (Wang et al., 2021a) decomposes
 386 Q values with a dueling structure, transferring the monotonicity condition to advantage values.

387 Focus on group development enables agents to maintain diversified policies and promote efficient
 388 collaborations. Early works are proposed only for tasks with a clear structure and train agents
 389 with similar traits to specialize in specific sub-tasks (DeLoach & Garcia-Ojeda, 2010; Bonjean
 390 et al., 2014). Recent works attempt to learn individual connections implicitly. VAST (Phan et al.,
 391 2021) learns value factorization for sub-teams based on a priori setting on group number. Wang
 392 et al. (2019) considers pairwise mutual influence to encourage interdependence between agents.
 393 ROMA (Wang et al., 2020) learns dynamic roles that depend on the context each agent observes.
 394 Iqbal et al. (2021) randomly groups agents into related and unrelated groups, allowing agents to
 395 explore only specific entities. RODE (Wang et al., 2021b) decomposes the joint action spaces and
 396 integrates the action effects into the role policies to boost learning. CDS (Li et al., 2021) incorporates
 397 agent-specific modules to promote sharing among agents while keeping necessary diversity.

398 6 CONCLUSION

399 Grouping like natural systems is essential to promote efficient cooperation of multi-agent systems.
 400 Instead of predefining grouping utilizing apriori knowledge, this paper proposes an automatic group-
 401 ing mechanism that gradually learns reasonable grouping as training proceeds. Based on the dy-
 402 namic grouping, we further encourage specialization in policies to promote individual similarity
 403 and group diversity, achieving efficient intra- and inter-group cooperation. With these novelties, our
 404 method, GoMARL, achieves impressive performance on both SMAC and GRF benchmarks.

REFERENCES

- Noelie Bonjean, Wafa Mefteh, Marie-Pierre Gleizes, Christine Maurel, and Frédéric Migeon. Adelfe 2.0, handbook on agent-oriented design processes, 2014.
- Paolo Bresciani, Anna Perini, Paolo Giorgini, Fausto Giunchiglia, and John Mylopoulos. Tropos: An agent-oriented software development methodology. In *International Conference on Autonomous Agents and Multi-Agent Systems*, volume 8, pp. 203–236. Springer, 2004.
- Scott A DeLoach and Juan Carlos Garcia-Ojeda. O-mase: A customisable approach to designing and building complex, adaptive multi-agent systems. *International Journal of Agent-Oriented Software Engineering*, 4(3):244–280, 2010.
- Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *AAAI Conference on Artificial Intelligence*, pp. 2974–2982, 2018.
- Edward I George and Robert E McCulloch. Variable selection via gibbs sampling. *Journal of the American Statistical Association*, 88(423):881–889, 1993.
- Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: A survey. *Artificial Intelligence Review*, pp. 1–49, 2021.
- Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multi-Agent Systems*, pp. 66–83, 2017.
- David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. In *International Conference on Learning Representations*, 2017.
- Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. A survey and critique of multiagent deep reinforcement learning. In *International Conference on Autonomous Agents and Multi-Agent Systems*, pp. 750–797, 2019.
- Jian Hu, Siyang Jiang, Seth Austin Harding, Haibin Wu, and Shih-wei Liao. Riit: Rethinking the importance of implementation tricks in multi-agent reinforcement learning. *arXiv preprint arXiv:2102.03479*, 2021.
- Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. Guided deep reinforcement learning for swarm systems. *arXiv preprint arXiv:1709.06011*, 2017.
- Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 2961–2970. PMLR, 2019.
- Shariq Iqbal, Christian A Schroeder De Witt, Bei Peng, Wendelin Böhmer, Shimon Whiteson, and Fei Sha. Randomized entity-wise factorization for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 4596–4606, 2021.
- Raphaël Jeanson, Penelope F Kukuk, and Jennifer H Fewell. Emergence of division of labour in halictine bees: Contributions of social interactions and behavioural variance. *Animal Behaviour*, 70(5):1183–1193, 2005.
- Jeewon Jeon, Woojun Kim, Whiyoung Jung, and Youngchul Sung. Maser: Multi-agent reinforcement learning with subgoals generated from experience replay buffer. In *International Conference on Machine Learning*, pp. 10041–10052, 2022.
- Jiechuan Jiang and Zongqing Lu. The emergence of individuality. In *International Conference on Machine Learning*, pp. 4992–5001, 2021.
- Landon Kraemer and Bikramjit Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.

- Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zajac, Olivier Bachem, Lasse Espeholt, Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research football: A novel reinforcement learning environment. In *AAAI Conference on Artificial Intelligence*, pp. 4501–4510, 2020.
- Chenghao Li, Tonghan Wang, Chengjie Wu, Qianchuan Zhao, Jun Yang, and Chongjie Zhang. Celebrating diversity in shared multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 34, pp. 3991–4002, 2021.
- Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, pp. 7611–7622, 2019.
- Frans A Oliehoek and Christopher Amato. *A concise introduction to decentralized pomdps*. Springer, 2016.
- Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.
- Afshin OroojlooyJadid and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *arXiv preprint arXiv:1908.03963*, 2019.
- Juan Pavón and Jorge Gómez-Sanz. Agent oriented software engineering with ingenias. In *International Central and Eastern European Conference on Multi-Agent Systems*, pp. 394–403. Springer, 2003.
- Thomy Phan, Fabian Ritz, Lenz Belzner, Philipp Altmann, Thomas Gabor, and Claudia Linnhoff-Popien. Vast: Value function factorization with variable agent sub-teams. volume 34, 2021.
- Marc’Aurelio Ranzato, Y-Lan Boureau, Yann Cun, et al. Sparse feature learning for deep belief networks. In *Advances in Neural Information Processing Systems*, volume 20, 2007.
- Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 4295–4304, 2018.
- Stuart J Russell and Andrew Zimdars. Q-decomposition for reinforcement learning agents. In *International Conference on Machine Learning*, pp. 656–663, 2003.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. *CoRR*, abs/1902.04043, 2019.
- Jeff G Schneider, Weng-Keen Wong, Andrew W Moore, and Martin A Riedmiller. Distributed value functions. In *International Conference on Machine Learning*, pp. 371–378, 1999.
- Xiaotong Shen and Jianming Ye. Adaptive model selection. *Journal of the American Statistical Association*, 97(457):210–221, 2002.
- Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 5887–5896, 2019.
- Jianyu Su, Stephen Adams, and Peter Beling. Value-decomposition multi-agent actor-critics. In *AAAI Conference on Artificial Intelligence*, pp. 11352–11360, 2021.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *International Conference on Autonomous Agents and Multi-Agent Systems*, pp. 2085–2087, 2018.
- Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *International Conference on Machine Learning*, pp. 330–337, 1993.

- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling multi-agent q-learning. In *International Conference on Learning Representations*, pp. 1–9, 2021a.
- Tonghan Wang, Jianhao Wang, Yi Wu, and Chongjie Zhang. Influence-based multi-agent exploration. In *International Conference on Learning Representations*, 2019.
- Tonghan Wang, Heng Dong, Victor Lesser, and Chongjie Zhang. Roma: Multi-agent reinforcement learning with emergent roles. In *International Conference on Machine Learning*, pp. 9876–9886, 2020.
- Tonghan Wang, Tarun Gupta, Anuj Mahajan, Bei Peng, Shimon Whiteson, and Chongjie Zhang. Rode: Learning roles to decompose multi-agent tasks. In *International Conference on Learning Representations*, pp. 1–9, 2021b.
- Christian Schroeder Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- G Wittemyer and Wayne M Getz. Hierarchical dominance structure and social organization in african elephants, *loxodonta africana*. *Animal Behaviour*, 73(4):671–681, 2007.
- Siyang Wu, Tonghan Wang, Chenghao Li, and Chongjie Zhang. Containerized distributed value-based multi-agent reinforcement learning. *arXiv preprint arXiv:2110.08169*, 2021.
- Zhao Xu, Yang Lyu, Quan Pan, Jinwen Hu, Chunhui Zhao, and Shuai Liu. Multi-vehicle flocking control with deep deterministic policy gradient method. In *International Conference on Control and Automation*, pp. 306–311. IEEE, 2018.
- Dayong Ye, Minjie Zhang, and Yun Yang. A multi-agent framework for packet routing in wireless sensor networks. *Sensors*, 15(5):10026–10047, 2015.
- Chao Yu, Akash Velu, Eugene Vinitsky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of Reinforcement Learning and Control*, pp. 321–384, 2021.
- Lulu Zheng, Jiarui Chen, Jianhao Wang, Jiamin He, Yujing Hu, Yingfeng Chen, Changjie Fan, Yang Gao, and Chongjie Zhang. Episodic multi-agent reinforcement learning with curiosity-driven exploration. *Advances in Neural Information Processing Systems*, 34:3757–3769, 2021.
- Meng Zhou, Ziyu Liu, Pengwei Sui, Yixuan Li, and Yuk Ying Chung. Learning implicit credit assignment for cooperative multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 11853–11864, 2020.

405 A THE INSIGHT OF GROUP SHIFTING

406 The proposed automatic grouping mechanism dynamically adjusts the team division as the training
 407 proceeds. In the beginning, all agents belong to the same group. As introduced in Section 3.1, we
 408 examine each w_1^i ($i \in [1, n]$) every c timesteps to check if the grouping needs adjustment. If there
 409 are agents who take actions but contribute little to their group, it indicates that these agents do not
 410 belong to their current group. All these selected agents are assigned to the next group (a new one
 411 for agents in the last group) until they appropriately contribute to where they belong.

412 We utilize the example in Figure 2 (also Figure A1 for your reading convenience) to illustrate the
 413 insight of this grouping shift, *i.e.*, this shifting scheme guarantees that agents belonging to the same
 414 group will not be misclassified into different groups during the dynamic adjustments.

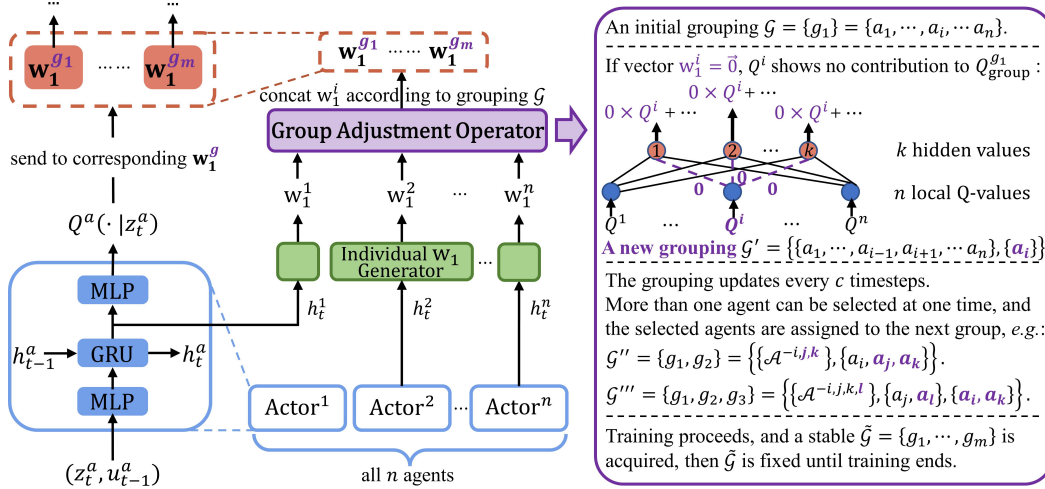


Figure A1: Schematic diagram of automatic grouping. Right box shows how grouping \mathcal{G} changes during training. In particular, we select the agents whose Q^i have little contribution to Q_{group}^g to move out of the current group g_j . Based on the grouping \mathcal{G} , the group selection operator concatenates the weights w_1^i of agents in the same group to form the group-wise weights w_1^g for mixing local utilities.

415 In this example, agent a_i is first selected after $\alpha_1 \cdot c$ timesteps' training, and the initial grouping
 416 $\mathcal{G} = \{a_1, a_2, \dots, a_n\}$ shifts to $\mathcal{G}' = \{\{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n\}, \{a_i\}\}$. Subsequently, after
 417 another training period, at timestep $\alpha_2 \cdot c$, two agents a_j and a_k in the first group were selected
 418 simultaneously to be moved out of the first group because of their small contribution. At this point,
 419 they were automatically placed in the second group, *i.e.*, the group where agent a_i is located. $\mathcal{G}'' =$
 420 $\{\{\mathcal{A}^{-i,j,k}\}, \{a_i, a_j, a_k\}\}$. Instead of placing a_j and a_k in a brand new group, our shifting ensures
 421 that a_j and a_k have the opportunity to train with a_i together as a group, determining if a_i , a_j , and
 422 a_k (or two of them) are supposed to be in a group. Later on, agent a_l in the first group and agent
 423 a_i, a_k in the second group are chosen at timestep $\alpha_3 \cdot c$. The selection of a_i and a_k indicates that
 424 agent a_j cannot cooperate well with them. Therefore, a_i and a_k are set in a new group (they should
 425 not go back to the first group since they have been proved inappropriate for the first group), while
 426 a_j stays alone in the second group. Agent a_l from the first group is assigned automatically to the
 427 second group (*i.e.*, $\{a_j\}$) in this example. A natural question is why this agent is not sent to the
 428 third group $\{a_i, a_k\}$. There is no need to worry about this under our shifting. If a_l is supposed
 429 to be with a_i and a_k , it will be selected afterward since a_j fails to form efficient cooperation with
 430 $\{a_i, a_k\}$. Therefore, a_j also cannot cooperate well with a_l , and a_l will be selected later on and set
 431 in group $\{a_i, a_k\}$. Equipped with the shifting that assigns selected agents to its following group,
 432 our proposed automatic grouping mechanism ensures that each grouping is tested by a fixed period
 433 of cooperative attempts. Therefore, each agent can be grouped with appropriate agents that can
 434 efficiently cooperate.

435 **B DUAL-HIERARCHY FACTORIZATION SATISFYING IGM**

436 The IGM (Individual-Global-Max) principle proposed by QTRAN (Son et al., 2019) is defined on
 437 the correspondence between individual greedy actions and joint greedy actions. Formally, if there
 438 exist individual action-value functions $[Q^i]_{i=1}^n$ that satisfy:

$$\arg \max_{\mathbf{u}} Q^{tot}(\boldsymbol{\tau}, \mathbf{u}) = \begin{pmatrix} \arg \max_{u^1} Q^1(\tau^1, u^1) \\ \vdots \\ \arg \max_{u^n} Q^n(\tau^n, u^n) \end{pmatrix}, \quad (\text{A1})$$

439 where $Q^{tot}(\boldsymbol{\tau}, \mathbf{u})$ is the joint action-value function with joint action observation histories $\boldsymbol{\tau}$ and joint
 440 action \mathbf{u} , we say that $[Q^i]_{i=1}^n$ satisfies IGM for Q^{tot} .

441 VDN (Sunehag et al., 2018) and QMIX (Rashid et al., 2018) are renowned methods that attempt to
 442 factorize Q^{tot} assuming additivity and monotonicity respectively that satisfies the IGM. VAST (Phan
 443 et al., 2021) combines the value factorization operators of VDN and QTRAN to heritage the IGM
 444 property. Here, we formulate the following Proposition 1 to show that GoMARL also maintains
 445 decentralizability by satisfying the IGM principle for Q^{tot} and Q^i .

446 **Proposition 1.** *Given a multi-agent system $M = \langle \mathcal{A}, S, U, P, r, Z, O; \mathcal{G} \rangle$, where each agent $a_i \in \mathcal{A}$
 447 with local utility Q^i belongs to only one group $g_j \in \mathcal{G}$ for GoMARL’s dual-hierarchy factorization.
 448 IGM is satisfied for Q^{tot} and $[Q^i]$ for each agent $a_i \in g_j, g_j \in \mathcal{G}$.*

449 *Proof.* The dynamic grouping mechanism in Section 3.1 regards the mixing weights of the individ-
 450 ual utilities as the contribution of Q^i to the group value. Therefore, it restricts the learned weights to
 451 be non-negative, as in the monotonic mixing that satisfies IGM. Thus, the maximization of $[Q^i]_{i=1}^n$
 452 maximizes $Q_{group}^{g_j}$ such that $\mathbf{u}^{g_j} = [u^i]_{i \in g_j}$, where $\mathbf{u}^{g_j} = \arg \max_{\mathbf{u}^{g_j} \in [U]_{a_i \in g_j}} Q_{group}^{g_j}(\boldsymbol{\tau}^{g_j}, \mathbf{u}^{g_j})$
 453 and $u^i = \arg \max_{u^i} Q^i(\tau^1, u^1)$. As introduced in Section 3.3, GoMARL’s factorization of Q^{tot}
 454 into $Q_{group}^{g_j}$ adopts a similar fashion with the factorization of $Q_{group}^{g_j}$ into Q^{tot} . Therefore, the
 455 factorization of Q^{tot} into $Q_{group}^{g_j}$ also satisfies the IGM such that:

$$\mathbf{u} = [\mathbf{u}^{g_j}]_{g_j \in \mathcal{G}} = [[u^i]_{i \in g_j}]_{g_j \in \mathcal{G}} = [u^i]_{i \in \mathcal{A}}. \quad (\text{A2})$$

456 Therefore, GoMARL ensures that the greedy local action set of all agents $\mathcal{A} = \{a_1, \dots, a_n\}$ maxi-
 457 mizes Q^{tot} for time-varying grouping \mathcal{G} according to the IGM principle in Eqn.(A1). \square

458 **C EXPERIMENT DETAILS**

459 **C.1 DETAILED EXPERIMENTAL SETUP FOR FAIR COMPARISON**

460 We compare all methods in six StarCraft II micromanagement tasks (SMAC) (Samvelyan et al.,
 461 2019) and two challenging Google Research Football (GRF) (Kurach et al., 2020) scenarios. Meth-
 462 ods are trained with 8 parallel runners for 10M steps in both testbeds. We evaluate each method every
 463 10K steps with 32 episodes and report the 1st, median, and 3rd quartile win rates across 5 random
 464 seeds. The detailed setting of GoMARL’s hyperparameters is introduced in our source code.

465 Many algorithms introduce implementation tricks when they are implemented. These code-level
 466 optimizations were studied in depth in Witt et al. (2020); Yu et al. (2021); Hu et al. (2021) and were
 467 shown to have a significant impact on algorithm performance. Considering that different baseline
 468 algorithms may use some of these code-level optimizations and thus cause unfair experimental com-
 469 parisons, this paper conducts experiments under strict control on tricky code-level implementations
 470 (e.g., reward clipping/scaling/normalization, gradient clipping, observation normalization, learning
 471 rate annealing, death agent masking, etc.), ensuring the comparisons are as fair as possible. Besides,
 472 specific parameter tuning in diverse scenarios is unfair to compare methods, so our experiments use
 473 fixed parameters for all methods in all environments.

474 To further ensure fair comparisons, our experiments are based on the PyMARL2 framework pro-
 475 posed for the purpose of fairly comparing algorithms, which is the fairest open-sourced framework
 476 we could find. Although RODE and ROMA are not included in this framework, we utilized the
 477 same PyMARL2 settings for them to ensure fairness. The performance of RODE and ROMA in
 478 this paper is one of the best among public papers, better than their performances in papers from the
 479 same authors or research team (Wu et al., 2021; Zheng et al., 2021) and other research (Jiang & Lu,
 480 2021; Jeon et al., 2022), and even better than their original papers in some environments, with abso-
 481 lutely no unfair comparisons. Please refer to PyMARL2’s open-source implementation¹ for further
 482 training details and fair comparison settings.

483 C.2 PARAMETER SIZE FOR VALUE MIXING

484 GoMARL maintains two mixing networks to estimate Q_{group}^g and Q^{tot} respectively. Whether Go-
 485 MARL has more parameters for value mixing is a natural question. We analyze that our input
 486 dimension reduction of the hypernetworks offsets the increase of an extra mixing network in Sec-
 487 tion 3.3. Here, we give the detailed mixing network size (GoMARL’s two mixing networks are
 488 counted) of all the methods in Table A1. The results of GoMARL is the average size of the 5 runs
 489 since each run may learn a slightly different grouping with diverse group amounts. Our two-mixing
 490 architecture has fewer parameters than the baseline QMIX when there are a large number of agents.

Table A1: The size of parameters for value mixing

Maps	(Ft-)QMIX	QPLEX	RODE	ROMA	RIIT	GoMARL
3s_vs_5z	21.601K	72.482K	43.202K	13.281K	37.986K	26.530K
5m_vs_6m	31.521K	107.574K	63.042K	25.377K	51.362K	31.554K
8m_vs_9m	53.313K	197.460K	106.626K	63.393K	93.986K	51.427K
corridor	68.929K	303.808K	137.858K	81.537K	122.882K	53.859K
MMM2	84.929K	342.248K	169.858K	134.401K	177.282K	74.244K
3s5z_vs_3s6z	63.105K	243.156K	126.21K	81.345K	118.466K	61.028K

491 C.3 DETAILED INFORMATION ABOUT SMAC AND ITS SCENARIOS

492 A group of units controlled by decentralized agents cooperates to defeat the enemy agent system con-
 493 trolled by handcrafted heuristics in each SMAC micromanagement problem. Each agent’s partial
 494 observation comprises the attributes (such as `health`, `location`, `unit_type`) of all units shown
 495 up in its view range. The global state information includes all agents’ positions and `health`, and al-
 496 lied units’ last actions and `cooldown`, which is only available to agents during centralized training.
 497 The agents’ discrete action space consists of `attack[enemy_id]`, `move[direction]`, `stop`,
 498 and `no-op` for the dead agents only. Particular unit `Medivac` has no action `attack[enemy_id]`
 499 but has `heal[enemy_id]`. Agents can only attack enemies within their shooting range. Proper
 500 micromanagement requires agents to maximize the damage to the enemies and take as little damage
 501 as possible in combat, so they need to cooperate with each other or even sacrifice themselves. We
 502 follow the default setup of SMAC in our experiments, and more settings, including rewards and de-
 503 tailed observation/state information, can be acquired from the original paper Samvelyan et al. (2019)
 504 or implementation².

505 Based on baseline algorithms’ performances, the scenarios in SMAC are broadly grouped into three
 506 categories: *Easy*, *Hard*, and *Super Hard*. The key to winning some *Hard* or *Super Hard* battles is
 507 mastering specific micro techniques, such as *focus fire*, *kiting*, *avoid overkill*, et cetera. The battles
 508 can be symmetric or asymmetric, and the group of agents can be homogeneous or heterogeneous.
 509 Here we provide some characteristics of each *Super Hard* scenario to help gain insights into the
 510 good or poor performance of the methods:

- 511 • `corridor` is a *Super Hard* map that need extensive exploration. Six allied Zealots fight
 512 twenty-four Zerglings on this super hard map. The massive disparity in unit numbers be-

¹<https://github.com/hijkzzz/pymarl2>

²<https://github.com/oxwhirl/smac>

513 tween the two sides implies that the whole team cannot launch an attack together. The only
514 winning strategy is to sacrifice a small number of agents who leave the team and attract the
515 attention of most enemies. Taking this opportunity, our large force eliminates the rest of
516 the enemies. The surviving agents then use the same tactic to attract several enemies to the
517 team every time and kill them together.

- 518 • `MMM2` is a representative *Super Hard* asymmetric battle between two heterogeneous teams
519 with three kinds of units. One Medivac, two Marauders, and seven Marines have to battle
520 against a team with one more Marine. Marauder has greater attack damage and health than
521 Marine but with a longer attack `cooldown`. Medivac has no damage but can heal any
522 other agent in the team. This map with three kinds of units and many agents requires more
523 cooperation between agents, so we picked this map for our ablation studies.
- 524 • `3s5z_vs_3s6z` is another *Super Hard* map that requires breaking the bottleneck of explo-
525 ration, where three Stalkers and five Zealots battle against three Stalkers and six Zealots.

526 C.4 DETAILED INFORMATION ABOUT GRF AND ITS SCENARIOS

527 Google Research Football (GRF) includes several scenarios which can be commonly found in foot-
528 ball games. Agents need to coordinate timing and positions for organizing offense to seize fleeting
529 opportunities, and only scoring leads to rewards. Each agent’s partial observation contains the ab-
530 solute positions and moving direction of the ego-agent, relative positions and moving directions of
531 other agents, and the ball. The global state information includes all agents’ and ball’s absolute po-
532 sitions/directions. Agents have a discrete action space of 19, including moving in eight directions,
533 sliding, shooting, and passing. Proper cooperation requires agents to pass and shoot effectively.
534 Other details can be acquired from the original paper (Kurach et al., 2020) or its implementation³.

535 Here we provide an introduction of the two scenarios we utilized to compare methods to help gain
536 insights into the good or poor performance of the methods:

- 537 • `3_vs_1_with_keeper` contains three of our agents and two opponents players (a de-
538 fender and a keeper). Our agents try to score from the edge of the penalty area. One of
539 them stands in the middle, while the others are located on both sides of the area. Initially,
540 the agent at the center keeps the ball and directly faces the defender.
- 541 • `academy_counterattack_easy` contains four of our agents and two opponent players
542 (a defender and a keeper). Agents are initialized far from the penalty area and stand evenly
543 in an arc centered on the goal. The second agent from the top initially keeps the ball and
544 has to pass it to a teammate at the appropriate time to avoid interception.

545 D EXTRA COMPONENT STUDIES ON THE SPECIALIZED AGENT NETWORK

546 In Section 4.1, we conducted ablation experiments on the proposed specialized agent networks. As
547 shown in the bottom row of Figure 6, our specialized agent networks significantly improve learning
548 efficiency. Compared to agent networks utilizing the vanilla parameter-sharing mechanism that
549 limits policy diversity, our agents have various styles of behaviors related to their group, which
550 encourages extensive exploration and accelerates learning.

551 Besides, we study the inner component of the specialized agent networks to provide deeper in-
552 sight. Policy specialization in GoMARL is driven by a specialization regularizer, *i.e.*, the similarity-
553 diversity objective to train the info encoder in Eqn.(2). The influence of the scaling factor λ_{SD} on
554 the performance is shown in Figure A2. According to these component studies, we set $\lambda_{SD} = 0.03$
555 in GoMARL for all the other experiments.

556 Furthermore, the proposed specialized agent network is highly transferable. To further validate
557 its effectiveness, we perform the specialized agents on other baseline methods (Ft-QMIX, QPLEX,
558 RODE, RIIT) to see if they can perform better. ROMA is not included since ROMA’s agent networks
559 are produced by its learned roles, and the replacement will destroy the main idea of the method. As
560 shown in Figure A3, The performance of all methods is further improved when equipped with our
561 specialized agent network. Specifically, the learning efficiency of Ft-QMIX is boosted. The variance

³<https://github.com/google-research/football>

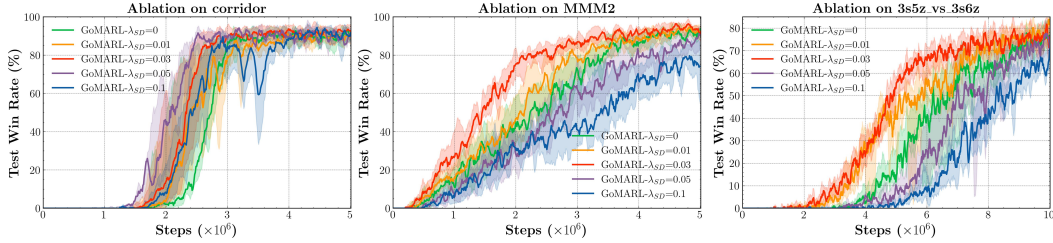


Figure A2: The influence of the scaling factors λ_{SD} on the performance.

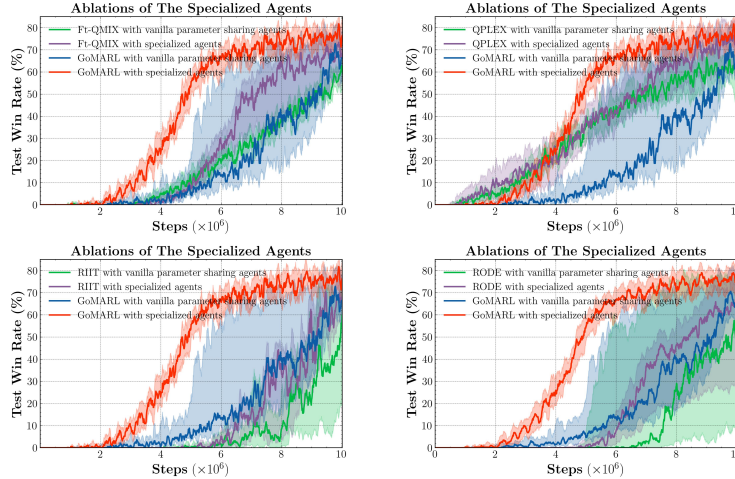


Figure A3: Improvement of baselines with our specialized agents to further prove its effectiveness.

562 of RIIT is markedly reduced, and the win rate is increased by about 10%. The improvement of
 563 RODE is the most obvious, both the learning efficiency and win rate are enhanced, and the variance
 564 is very clearly reduced. QPLEX’s improvement is not very obvious; however, it obtains a slightly
 565 higher learning speed and achieves similar performance to GoMARL at the end of training.

566 Most importantly, even equipped with our specialized network, all baseline algorithms fail to surpass
 567 GoMARL in terms of learning efficiency and final performance. GoMARL with vanilla parameter
 568 sharing looks much inferior to GoMARL with specialized agents, so it is worth questioning whether
 569 the dynamic grouping module is effective. This experiment fully illustrates that although dynamic
 570 group learning may reduce the learning efficiency to a certain extent *in the early stage of training*,
 571 however, *in the middle and late training stages*, the learned grouping will have a significant effect
 572 when equipped with the proposed specialized agent network. Therefore, the two main modules of
 573 GoMARL, the dynamic grouping module and specialized agents, are both crucial.

574 E EXTRA VISUALIZATION OF GOOGLE RESEARCH FOOTBALL MATCH

575 In Section 4.2, we visualize a match of 3_vs_1_with_keeper to prove the rationality of the group-
 576 ing GoMARL learned. Here we give another visualization of the counterattack scenario to
 577 show a more complex strategy and its corresponding learned grouping.

578 As shown in Figure A4(a), Agent 2 holds the ball at the beginning and faces an opponent rushing
 579 toward him. Agent 2 passes the ball to Agent 1 to prevent the ball from being stolen. Subsequently, in
 580 (b), Agent 1 carries the ball and tries to break through. However, the opponent goalkeeper blocks his
 581 attacking route, and Agent 1 chooses to continue passing after a short carry. Agent 3 makes a good
 582 run and catches the ball smoothly but shoots quickly in Figure A4(c) since the opponent is close.

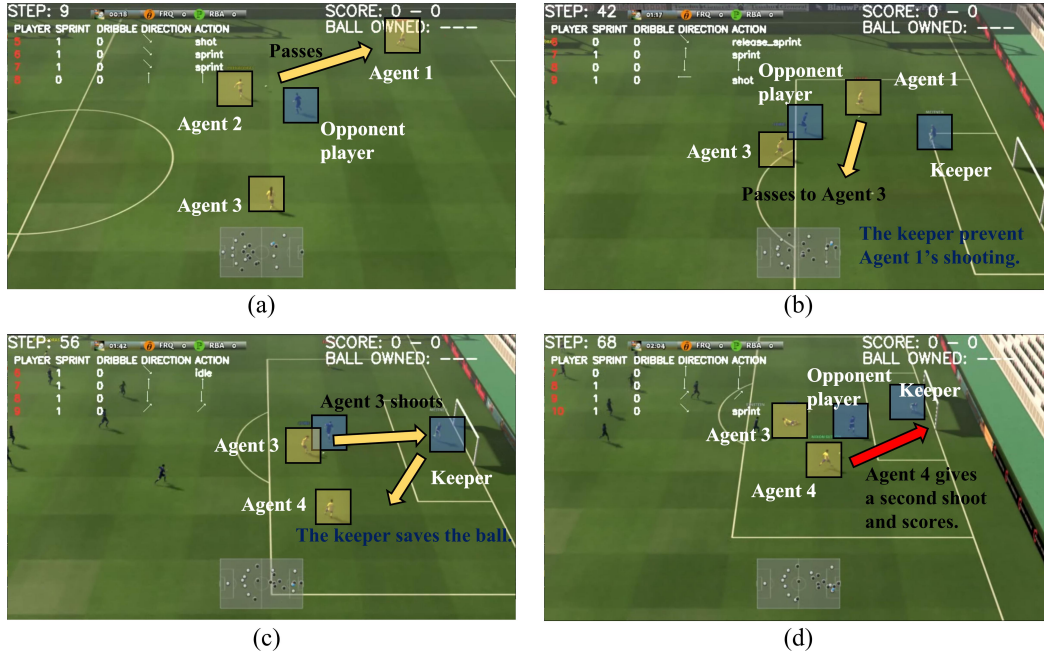


Figure A4: A visualization of learned policies on academy_counterattack_easy. Yellow arrows show the motion of the ball. The red arrow illustrates the scoring shoot.

583 The goalkeeper easily saves this hasty attack. Agent 4 in (d), who learned excellent coordination
 584 with Agent 3, stops the ball and immediately adds another shot to create the goal.

585 During this complex goal, GoMARL's automatic grouping module learned a reasonable group-
 586 ing $\{\{Agent1, Agent2\}, \{Agent3, Agent4\}\}$, in which the first group successfully brought the ball
 587 into the penalty area through smooth coordination, while the second group finally created the goal
 588 through skillful cooperation of shooting and a second shooting.