

# MRVF: MULTI-ROUND VALUE FACTORIZATION WITH GUARANTEED ITERATIVE IMPROVEMENT FOR MULTI-AGENT REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Value factorization restricts the joint action value in a monotonic form to enable efficient search for its optimum. However, the representational limitation of monotonic forms often leads to suboptimal results in cases with highly non-monotonic payoff. Although recent approaches introduce additional conditions on factorization to address the representational limitation, we propose a novel theory for convergence analysis to reveal that single-round factorizations with elaborated conditions are still insufficient for global optimality. To address this issue, we propose a novel Multi-Round Value Factorization (MRVF) framework that refines solutions round by round and finally obtains the global optimum. To achieve this, we measure the non-negative incremental payoff of a solution relative to the preceding solution. This measurement enhances the monotonicity of the payoff and highlights solutions with higher payoff, enabling monotonic factorizations to identify them. We evaluate our method in three challenging environments: non-monotonic one-step games, predator-prey tasks, and StarCraft II Multi-Agent Challenge (SMAC). Experiment results demonstrate that our MRVF outperforms existing value factorization methods, particularly in scenarios highly non-monotonic payoff.

## 1 INTRODUCTION

Multi-agent reinforcement learning (MARL) effectively addresses many real-world problems, such as robotics (Hüttenrauch et al., 2019), automated warehouses (Tao et al., 2024), and games (Bernier et al., 2019). MARL’s core issue is finding the optimal action in an action space that grows exponentially with the number of agents. To address this issue, value factorization methods represent the joint action value in a specific form. Early value factorization methods, such as VDN (Sunehag et al., 2017) and QMIX (Rashid et al., 2020b), represent the joint action value through a monotonic mix of individual action values. This monotonic relationship ensures that the optimal joint action can be obtained by searching through individual action spaces. However, because of the monotonic particularity, these methods struggle to obtain the global optimum in non-monotonic cases.

To mitigate the representational limitation in VDN and QMIX, the following methods propose new fitting functions (e.g., QPLEX (Wang et al., 2021) and ResQ (Shen et al., 2022)) or loss functions (e.g., QTRAN (Son et al., 2019) and WQMIX (Rashid et al., 2020a)). These designs can be treated as conditions on the current greedy action to enhance its convergence to the optimal action. However, the effectiveness of these conditions relies on a strong assumption: the current greedy action is the optimal action, which means that the optimal action has already been found. We find that when the current greedy action is among certain suboptimal actions, the greedy action may converge to such suboptimal actions, resulting in poor performance of these methods. **One of our contributions is that** we introduce a theoretical tool with a novel concept, **stable point**, to describe the convergence of greedy action under value factorization. Using this tool, we explain how existing methods converge to suboptimal solutions and provide specific cases where they fail to obtain the optimum.

Theoretically, we conclude that no matter how elaborate the conditions on greedy action are, monotonic factorization is almost impossible to obtain the optimal solution for non-monotonic cases. However, if a non-monotonic payoff is transformed into a monotonic one, we can easily obtain the optimal solution with monotonic factorization. To achieve this, **another contribution of ours is**

that we propose a multi-round value factorization (MRVF) framework. As illustrated in Figure 1, in each round, we replace the payoff with the payoff increment (non-negative), the payoff clipped by that of the preceding solution. This process gradually transforms the original payoff into a monotonic one round by round, thereby enabling monotonic factorization to achieve the optimal solution. Furthermore, using the theory of stable points, we prove that the solution obtained based on the increment is strictly improved from the previous round. This guarantees that the optimal solution can be obtained with a sufficient number of iterations.

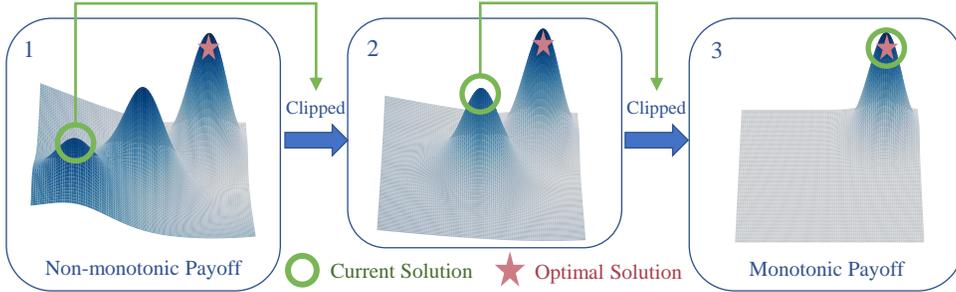


Figure 1: The process of transforming a non-monotonic payoff into a monotonic one round by round. At the beginning, MRVF obtains a suboptimal solution for the non-monotonic payoff. However, by clipping off these suboptimal values, we transform the payoff into a monotonic one that only leaves optimal values, thereby enabling monotonic factorization to find the optimal solution. The definition of monotonic payoff is presented in Appendix B.1.

To illustrate the superiority of MRVF, we conduct experiments on randomly generated one-step games that require significant coordination to get bonus, or otherwise a penalty. In addition, we evaluate MRVF on more challenging tasks such as the predator-prey task (Böhmer et al., 2020) and the StarCraft II MARL tasks (Whiteson et al., 2019). The experimental results show that MRVF outperforms existing value factorization methods, particularly in challenging scenarios that require high levels of agent coordination. The ablation study demonstrates the importance of strictly improving current solutions by using the payoff increment as the target value.

## 2 BACKGROUND

### 2.1 DEC-POMDP

In cooperative multi-agent systems, agents interact with the environment to achieve common objectives. This process can be modeled as a decentralized partially observable Markov decision process (Dec-POMDP) (Oroojlooy & Hajinezhad, 2023), defined by a tuple  $\langle \mathcal{S}, \mathcal{U}, P, O, R, \gamma, n \rangle$ , in which  $n$  is the number of agents,  $\mathcal{S}$  is the state space,  $\mathcal{U} = \mathcal{U}_1 \times \mathcal{U}_2 \times \dots \times \mathcal{U}_n$  is the action space,  $P(\mathbf{s}' | \mathbf{s}, \mathbf{u}) : \mathcal{S} \times \mathcal{U} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability between the states,  $O : \mathcal{S} \rightarrow \mathcal{O}$  is the joint observation function where  $\mathcal{O}$  is the joint observation space,  $r \sim R : \mathcal{S} \times \mathcal{U} \rightarrow \mathbb{R}$  is the reward,  $\gamma \in (0, 1]$  is the discount, agent have an action-observation history  $\tau \in \mathcal{T} \equiv (\mathcal{O} \times \mathcal{U})^*$ . Similarly to RL with a single agent, the objective of MARL is to find the policy  $\pi = (\pi_1, \pi_2, \dots, \pi_n) : \mathcal{T} \times \mathcal{U} \rightarrow [0, 1]$  that maximizes the joint action value  $Q_{jt}(\mathbf{s}, \mathbf{u}) = \mathbb{E}_\pi[\sum_t \gamma^t r_t | \mathbf{s}, \mathbf{u}]$ , the expectation of return. However, obtaining the global optimal action  $\mathbf{u}^* = \arg \max_{\mathbf{u}} Q_{jt}(\mathbf{s}, \mathbf{u})$  by searching the large action space is intractable for value-based learning.

### 2.2 MONOTONIC VALUE FACTORIZATION

Value factorization methods approximate the joint action value  $Q_{jt}$  with a specifically formed  $Q_{tot}$  by minimizing  $L_{tot}$ , the Mean Squared Error (MSE) between  $Q_{jt}$  and  $Q_{tot}$ . For example, VDN (Sunehag et al., 2017) factorizes  $Q_{tot}$  into the sum of individual  $Q_i$ :  $Q_{tot} = \sum_i Q_i$ . And QMIX (Rashid et al., 2020b) uses a monotonic function  $f_{mon}$  to represent the relationship:  $Q_{tot} = f_{mon}(Q_1, Q_2, \dots, Q_n)$  where  $\frac{\partial f_{mon}}{\partial Q_i} \geq 0, \forall i \in \{1, 2, \dots, N\}$ . The monotonicity of  $Q_{tot}$  with respect to  $Q_i$  ensures the alignment of their optimal actions, which is the Individual-Global-

108 Max (IGM) principle, defined as follows <sup>1</sup>:

$$109 \prod_{i=0}^n \arg \max_{u_i} Q_i(\tau_i, u_i) \subseteq \arg \max_{\mathbf{u}} Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) \quad (1)$$

110 where  $\prod_{i=0}^n \arg \max_{u_i} Q_i(\tau_i, u_i) = \arg \max_{u_1} Q_1(\tau_1, u_1) \times \cdots \times \arg \max_{u_n} Q_n(\tau_n, u_n)$ , and  $Q_i : \mathcal{T}_i \times$   
 111  $\mathcal{U}_i \rightarrow \mathbb{R}$  is the individual action value of agent  $i$ . Here we define  $\bar{\mathbf{u}} \in \prod_{i=0}^n \arg \max_{u_i} Q_i(\tau_i, u_i)$  as the  
 112 greedy action. If the IGM principle is satisfied, and  $Q_{\text{tot}}$  approximate  $Q_{\text{jt}}$  precisely, it is easy to get  
 113 the maximum of  $Q_{\text{jt}}$  by searching that of  $Q_i$ .

### 120 2.3 GRADIENT-FREE COMPONENTS IN VALUE FACTORIZATION

121 Value factorization methods use gradient descent to update gradient-based components, such as  
 122 individual Q values and the mixing network  $f_{\text{mon}}$ . However, not all components in the factorization  
 123 are gradient-based. Specifically, the components taking greedy action as a parameter are gradient-  
 124 free. To distinguish the meanings of  $\bar{\mathbf{u}}$  as the parameter and greedy action, we refer to  $\tilde{\mathbf{u}}$  as the  
 125 parameter of these components, where  $\bar{\mathbf{u}}$  is assigned to  $\tilde{\mathbf{u}}$ . Gradient-free components are common  
 126 in value factorization, for example, the decentralized  $\epsilon$ -greedy policy  $\pi$

$$127 \pi(\mathbf{u}|\boldsymbol{\tau}) = \left(\frac{\epsilon}{|\mathcal{U}|}\right)^{n-m} \left(1 - \epsilon + \frac{\epsilon}{|\mathcal{U}|}\right)^m \quad (2)$$

128 where  $m = |\{i|u_i = \tilde{u}_i\}|$ . Recent methods introduce additional terms involving  $\tilde{\mathbf{u}}$  into monotonic  
 129 factorization. For example, the weight  $w$  of the weighted  $L_{\text{tot}}$  in WQMIX (Rashid et al., 2020a) <sup>2</sup>

$$130 w(\boldsymbol{\tau}, \mathbf{u}) = \begin{cases} 1 & \hat{Q}_{\text{jt}}(\boldsymbol{\tau}, \mathbf{u}) > \hat{Q}_{\text{jt}}(\boldsymbol{\tau}, \tilde{\mathbf{u}}) \text{ or } \mathbf{u} = \tilde{\mathbf{u}} \\ \alpha < 1 & \text{otherwise} \end{cases} \quad (3)$$

131 where  $\hat{Q}_{\text{jt}}$  is an approximation of  $Q_{\text{jt}}$ , and the residual mask  $w_r$  in ResQ (Shen et al., 2022)

$$132 w_r(\boldsymbol{\tau}, \mathbf{u}) = \begin{cases} 0 & \mathbf{u} = \tilde{\mathbf{u}} \\ 1 & \text{otherwise} \end{cases} \quad (4)$$

## 140 3 SUBOPTIMALITY OF EXISTING SINGLE-ROUND FACTORIZATION

141 In this section, we discuss why the policy (or the greedy action) is trapped in local optima dur-  
 142 ing training, which results in existing value factorization methods obtain a suboptimal result. To  
 143 explain this, we need to analyze the convergence of greedy actions in value factorization. For the  
 144 convergence analysis, existing works only analyze the cases when  $\tilde{\mathbf{u}}$  (a parameter of gradient-free  
 145 components) is the optimal action, which is insufficient. In contrast, we provide an analytical ap-  
 146 proach to determine what outcomes that  $\bar{\mathbf{u}}$  (the greedy action) converges to for general cases of  $\tilde{\mathbf{u}}$ ,  
 147 and reveal situations where  $\bar{\mathbf{u}}$  converges to suboptimal actions under certain  $\tilde{\mathbf{u}}$ .

### 150 3.1 TRANSITION OF THE GREEDY ACTION

151 To introduce the analysis of the greedy action's convergence, we first illustrate how the greedy action  
 152 changes during training. We provide an example in Figure 2, where we observe that the transition  
 153 of greedy actions is similar to that of a Markov process:  $\tilde{\mathbf{u}}$  acts as the current "state" (controlling  
 154 the transitions) and  $\bar{\mathbf{u}}$  as the next "state". This process emerges because updates to gradient-free  
 155 parameters (indicated by  $\tilde{\mathbf{u}}$ ) and gradient-based parameters (indicated by  $\bar{\mathbf{u}}$ ) are asynchronous, with  
 156 their relationship shown in Figure 3.

157 Combining Figure 3, we illustrate the process in Figure 2 as follows:

- 158 (1) Initially,  $\tilde{\mathbf{u}} = (1, 3)$  (top right) generates an all-ones weight matrix.

159 <sup>1</sup>We use " $\subseteq$ " instead of " $=$ " in Equation (1), as  $Q_{\text{tot}}$  may have more maxima than individual  $Q$ s (Table 1).

160 <sup>2</sup>It stands for Centrally-Weighted QMIX.

- 162 (2) Minimizing  $L_{\text{tot}} = w * (Q_{\text{tot}} - Q_{\text{jt}})^2$  under the all-ones weight  $w$  yields the  $Q_{\text{tot}}$  (left table),  
 163 corresponding to  $\bar{\mathbf{u}} = (3, 3)$  (bottom right).  
 164 (3)  $\tilde{\mathbf{u}}$  is assigned the new value  $\bar{\mathbf{u}} = (3, 3)$ , thereby altering the weight matrix (shaded regions).  
 165 (4) Minimizing  $L_{\text{tot}}$  with the weight under  $\tilde{\mathbf{u}} = (3, 3)$  produces a new  $Q_{\text{tot}}$  (middle table).  
 166 (5)  $\bar{\mathbf{u}}$  keeps on changing (steps 2-4) until  $\bar{\mathbf{u}} = \tilde{\mathbf{u}}$  is reached (right table).

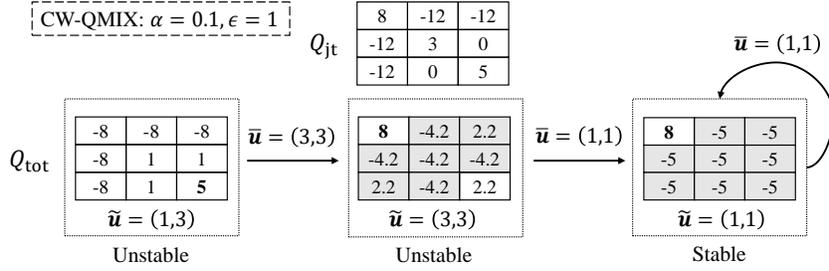


Figure 2: The transition of greedy action in WQMIX with  $\alpha = 0.1$  under uniform visitation ( $\epsilon = 1$ ). This process begins with  $\tilde{\mathbf{u}} = (1, 3)$  (top right) and stabilizes at  $\bar{\mathbf{u}} = (1, 1)$  (top left). The cells in  $Q_{\text{tot}}$  where  $w(s, \mathbf{u}) = \alpha$  are shaded, and the value corresponding to  $\bar{\mathbf{u}} = (\text{row}, \text{column})$  is in bold.

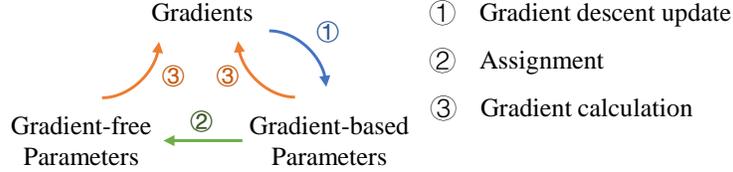


Figure 3: The relationship between gradient-based and gradient-free parameters. The gradient-free parameters are updated according to the values of gradient-based parameters rather than the gradients, leading to asynchronous updates between the two parameters.

### 194 3.2 CONVERGENCE OF THE GREEDY ACTION

195 From the discussion above, the transition of  $\bar{\mathbf{u}}$  typically begins in a transient phase (steps 2-4) and  
 196 ends in a steady phase (step 5). The steady phase of this transition is what we are interested in,  
 197 which reflects the convergence result of greedy actions. Based on the above discussion, **we use**  
 198 **stable points to describe the convergence of  $\bar{\mathbf{u}}$** , defined as follows:

199 **Definition 1** (Stable Point). *In a specific value factorization framework, for a joint action value  $Q_{\text{jt}}$   
 200 and state  $s \in \mathcal{S}$ , if  $\tilde{\mathbf{u}}$  is a stable point of  $Q_{\text{jt}}(s, \cdot)$ , then for  $\tilde{\mathbf{u}} = \tilde{\mathbf{u}}$  in the optimization of minimizing  
 201  $L_{\text{tot}}$ , the converged  $\bar{\mathbf{u}}$  satisfies  $\bar{\mathbf{u}} = \tilde{\mathbf{u}}$ .*

202 Note that multiple stable points may exist in value factorization, any of which could be the final  
 203 result of greedy actions. To analyze the possible outcomes of greedy actions, we need to determine  
 204 whether an action is a stable point, which is described as follows:

- 205 (1) Set  $\tilde{\mathbf{u}} = \tilde{\mathbf{u}}$  to generate the gradient-free parameters.  
 206 (2) Minimize  $L_{\text{tot}}$  under  $\tilde{\mathbf{u}} = \tilde{\mathbf{u}}$  to obtain  $Q_{\text{tot}}$  (multiple solutions may exist) and corresponding  $\bar{\mathbf{u}}$ .  
 207 (3) If any  $\bar{\mathbf{u}}$  satisfies  $\bar{\mathbf{u}} = \tilde{\mathbf{u}}$ , then  $\tilde{\mathbf{u}}$  is a stable point.  
 208 (4) Furthermore, if  $\bar{\mathbf{u}} = \tilde{\mathbf{u}}$  holds for all solutions,  $\tilde{\mathbf{u}}$  is a strongly stable point. Otherwise,  $\tilde{\mathbf{u}}$  is a  
 209 weakly stable point.

210 The difference between mixer types lies in Step 2: For VDN-style mixers, we obtain  $Q_{\text{tot}}$  using  
 211 gradient descent. And for QMIX-style mixers, we obtain  $Q_{\text{tot}}$  based on the Ideal QMIX assumption:  
 212 The Ideal QMIX yields the optimal  $Q_{\text{tot}}$  that minimizes  $L_{\text{tot}}$  under any given  $\tilde{\mathbf{u}}$  (Appendix B.2).



where  $\max\{\hat{Q}_{jt}(\tau_t, \mathbf{u}_t^k) - \hat{Q}_{jt}(\tau_t, \bar{\mathbf{u}}_t^{k-1}), 0\}$  is the action-value increment.

Intuitively, the action-value increment enhances the monotonicity of  $\hat{Q}_{jt}$  round by round, thereby enabling monotonic factorization to find the optimal solution easily. **Theoretically, it ensures the strict improvement of greedy action when  $\bar{\mathbf{u}}_t^{k-1} \neq \mathbf{u}^*$ .** This is grounded in a property of monotonic factorization: the greedy action under QMIX will not converge to actions with the lowest target value of  $Q_{\text{tot}}$ , since these lowest-value actions are not stable points (Theorem 2). Therefore, with the the action-value increment,  $\bar{\mathbf{u}}_t^k$  only converge to actions that are superior to  $\bar{\mathbf{u}}_t^{k-1}$ , since all  $\mathbf{u}_t^k$  with  $\hat{Q}_{jt}(\tau_t, \mathbf{u}_t^k) \leq \hat{Q}_{jt}(\tau_t, \bar{\mathbf{u}}_t^{k-1})$  are assigned the minimum value, zero.

**Theorem 2.** Consider a non-constant joint action value  $Q_{jt}(\mathbf{s}, \cdot)$  ( $\forall \mathbf{s} \in \mathcal{S}, \forall C \in \mathbb{R}, Q_{jt}(\mathbf{s}, \cdot) \not\equiv C$ ) with finite state space  $\mathcal{S}$  and action space  $\mathcal{U}$ . For the ideal QMIX defined in Equation (12) and the centralized  $\epsilon$ -greedy policy  $\pi$  defined in Equation (8), we have  $\forall \mathbf{s} \in \mathcal{S}, \exists E \in (0, 1), \forall \epsilon \in (0, E), \forall \mathbf{u} \in \arg \min_{\mathbf{u}} Q_{jt}(\mathbf{s}, \mathbf{u}), \mathbf{u}$  is not a stable point. (Proof in Appendix C.1)

We illustrate in Table 2 how the target value of  $Q_{\text{tot}}$  shapes during the multi-round process. For the first round ( $k = 1$ ), we use the original  $\hat{Q}_{jt}$  as the target for  $Q_{\text{tot}}$ , as shown in Equation (6).

$$L_{\text{tot}} = E_{\tau_t, \mathbf{u}_t^k \sim \pi_k} [(Q_{\text{tot}}(\tau_t, \mathbf{u}_t^k) - \hat{Q}_{jt}(\tau_t, \mathbf{u}_t^k))^2], \quad k = 1 \quad (6)$$

8	-12	-12
-12	3	0
-12	0	<b>5</b>

(a)  $k = 1$ 

<b>3</b>	0	0
0	0	0
0	0	0

(b)  $k = 2$ 

0	0	0
0	0	0
0	0	0

(c)  $k = 3$ 

Table 2: This table shows the target of  $Q_{\text{tot}}$  in three round, in which the value corresponding to  $\bar{\mathbf{u}}$  is in bold when it is unique. (a): The target of  $Q_{\text{tot}}$  in the first round which is the original  $\hat{Q}_{jt}$ . (b): The target of  $Q_{\text{tot}}$  in the second round which is clipped with  $\hat{Q}_{jt}(\tau_t, \bar{\mathbf{u}}_t^1) = 5$ . (c): The target of  $Q_{\text{tot}}$  in the final round where any  $\mathbf{u}$  can be  $\bar{\mathbf{u}}^3$ .

The forward computation in step  $t$  generates the final action  $\mathbf{u}_t$  (or the greedy final action  $\bar{\mathbf{u}}_t$  during evaluation) that actually interacts with the environment. Note that the final action is not always the action in the final round. As shown in Table 2, the greedy action in the round  $k = 2$  is the optimal action. As a result, the greedy action can possibly converge to any action in round  $k = 3$ , which violates the Strict Improvement Condition. Therefore, early termination is necessary when the strict improvement is not achieved. Specifically, we check the strict improvement by comparing the  $\hat{Q}_{jt}$  values of the two greedy actions: When  $\hat{Q}_{jt}(\tau_t, \bar{\mathbf{u}}_t^k) \leq \hat{Q}_{jt}(\tau_t, \bar{\mathbf{u}}_t^{k-1})$ , the forward process in step  $t$  is terminated with an output  $\mathbf{u}_t^{k-1}$ , otherwise the process proceed to the next round.

## 4.2 OVERVIEW OF THE ARCHITECTURE

The architecture of our methods is shown in Figure 4. The architecture consists of three crucial parts: individual Q networks, a mixing network with positive weights, and a joint action value network. Individual Q networks receive the action-observation history  $\tau_t$  and the greedy action of the preceding round  $\bar{\mathbf{u}}_t^{k-1}$  (default action for the first round). Then, they output the actions of the current round  $\mathbf{u}_t^k$  through the  $\epsilon$ -greedy action selector. The mixing network receives individual Q values and outputs  $Q_{\text{tot}}$ . The joint action value network outputs  $\hat{Q}_{jt}$  to approximate the real  $Q_{jt}$ . We update  $\hat{Q}_{jt}$  by minimizing the temporal difference (TD) error in Equation (7).

$$L_{jt} = E_{\tau_t, \mathbf{u}_t \sim \pi} [(\hat{Q}_{jt}(\tau_t, \mathbf{u}_t) - (r_t + \gamma \hat{Q}_{jt}(\tau_{t+1}, \bar{\mathbf{u}}_{t+1})))^2] \quad (7)$$

where  $\bar{\mathbf{u}}_t$  is the greedy final action, and  $\hat{Q}_{jt}(\tau_{t+1}, \bar{\mathbf{u}}_{t+1})$  approximates  $\max_{\mathbf{u}_{t+1}} \hat{Q}_{jt}(\tau_{t+1}, \mathbf{u}_{t+1})$ .

Within step  $t$ , both the forward and backward processes take  $\tau_t$  as input, and in each round, the individual Q networks are fed with  $\tau_t$  and  $\bar{\mathbf{u}}_t^{k-1}$  to produce  $Q_i$  and  $\bar{\mathbf{u}}_t^k$ . The difference is that the forward process computes and compares  $\hat{Q}_{jt}$  to determine whether to terminate early, while the backward process does not involve early termination and computes  $\hat{Q}_{jt}$  only once using the recorded  $\mathbf{u}_t$  from the batch (not shown in Figure 4). In addition,  $Q_{\text{tot}}$  is computed only in the backward.

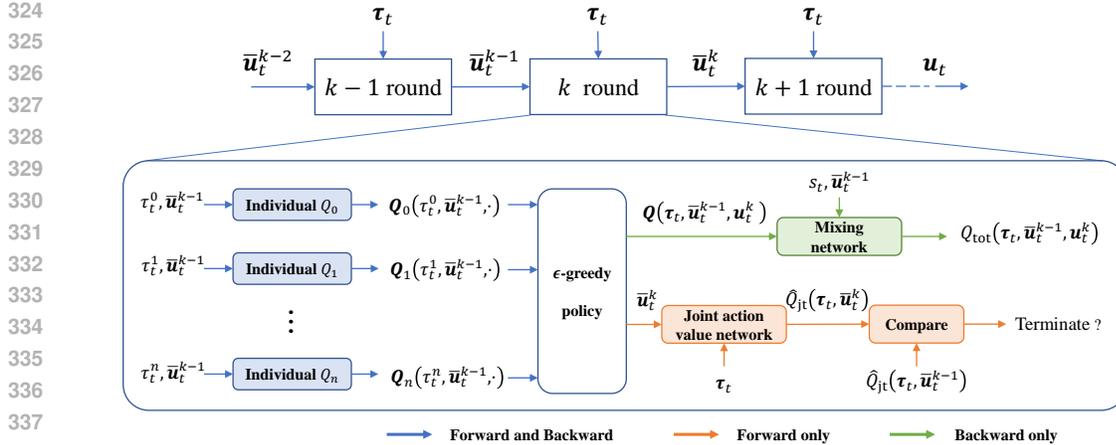


Figure 4: The architecture of multi-round value factorization framework.

### 4.3 SAMPLING

We approximate the expectation in  $L_{jt}$  and  $L_{tot}$  by sampling the action spaces. We design sampling strategies for multi-round frameworks to ensure stable training. For the final action  $u_t$  in  $L_{jt}$ , its sampling should not only cover the greedy final action and random actions, but also cover the greedy action in each round, which obtains an accurate target value in Equation (5). Therefore, during training, we select  $\bar{u}_t^k$  in a certain round  $k$  as the final action with probability  $p$ . For  $u_t^k$  in  $L_{tot}$ , its distribution  $\pi_k$  is the centralized  $\epsilon$ -greedy policy defined in Equation (8) with  $\tilde{u} = \bar{u}^k$ .

$$\pi(u|\tau) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{U}|} & u = \tilde{u} \\ \frac{\epsilon}{|\mathcal{U}|} & u \neq \tilde{u} \end{cases} \quad (8)$$

The reason for not using the decentralized one (defined in Equation (2)) is that, according to Theorem 2, the centralized  $\epsilon$ -greedy policy is required to guarantee the Strict Improvement Condition. More importantly, since  $Q_{tot}$  learns the actual payoff from  $\hat{Q}_{jt}$ , the centralized  $\epsilon$ -greedy policy allows the same actions to be sampled in  $L_{jt}$  and  $L_{tot}$  when sampling random actions, which prevents  $Q_{tot}$  from learning regions where  $\hat{Q}_{jt}$  exhibits underfitting.

## 5 EXPERIMENT

In this section, we evaluate the performance in one-step games, the predator-prey task (Böhmer et al., 2020) with increasing punishment, and a variety of SMAC (Whiteson et al., 2019) scenarios, among which one-step games and predator-prey tasks with large punishment are environments with highly non-monotonic payoff (defined in Appendix B.1). We report the median performance, with the shaded area denoting the 25%-75% percentile (0%-100% percentile for one-step games). Our experimental setup is provided in Appendix E, and ablation results are presented in Appendix F.

### 5.1 ONE-STEP GAME

From the discussion in Section 3.3, we find that existing methods fail to achieve optimality when, at the global optimum, any deviation by an agent leads to significant negative rewards. To demonstrate this limitation of existing methods, we evaluate their performance in the *risk-reward* game, where agents must reach a consensus to obtain rewards, or otherwise they incur punishments, and higher rewards for consensus correspond to harsher punishments for dissension. These games are randomly generated by averaging individual rewards (positive if consensus is reached, negative otherwise), and more details are provided in the Appendix E.1. The *risk-reward* game is an extension of the matrix game (Son et al., 2019) (Shen et al., 2022) that allows the participation of more than two agents.

As shown in Figure 5, in cases with 3 agents and 5 actions, monotonic factorizations (QMIX (Rashid et al., 2020b) and NA<sup>2</sup>Q (Liu et al., 2023)) obtain the smallest positive return, nearly zero, while

QTRAN (Wang et al., 2021) and WQMIX (Rashid et al., 2020a) only achieve a moderate positive return. In contrast, our method consistently attains the optimal return. When the number of agents increases to 5, the risk of obtaining positive returns increases, since consensus requires coordination among more agents. Existing methods rarely obtain positive returns, while our method still obtains the optimum with high probability. However, in cases with 5 agents and 8 actions, the positive rewards become too sparse to obtain. Despite this challenge, our method still obtains better results than others.

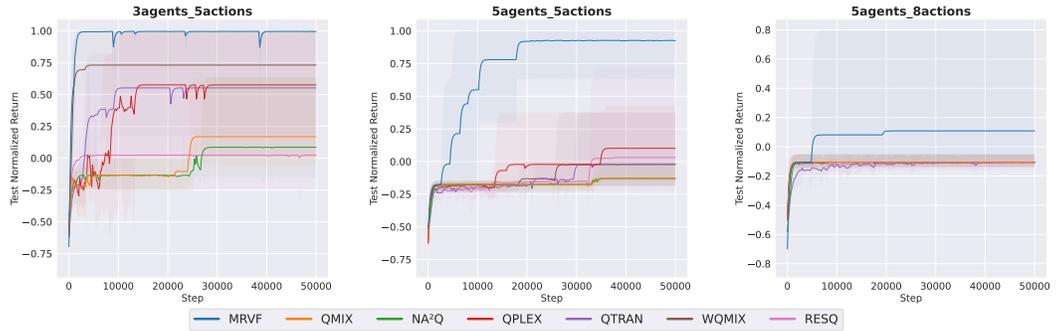


Figure 5: Test normalized return in the risk-reward games. The positive returns are normalized to  $[0, 1]$  (0 corresponds to the smallest positive return, and 1 corresponds to the largest). The negative returns are normalized to  $[-1, 0)$ . Five random cases for each setting.

### 5.2 PREDATOR PREY

We conduct experiments on the predator-prey task involving 8 agents, where capturing a prey requires cooperation between at least two agents. Agents are punished for capturing the prey alone. We vary the punishment from 0 to -5 to evaluate its impact on performance. The non-monotonic property of the payoff increases with the punishment, since agents must act more cautiously when choosing the "capture" action to ensure consensus and avoid punishments.

The results in Figure 6 show that all methods perform well when the punishment is zero. However, as the punishment intensifies, existing methods avoid the "capture" action to prevent punishments, resulting in  $return \approx 0$ . Among baseline approaches, WQMIX performs adequately only under moderate punishments, while the performance of ResQ shows significant fluctuations due to its weak stability. In contrast, our method consistently achieves the best performance, even under the strongest punishment of -5. Combined with the results in *risk-reward* games, we conclude that our method outperforms existing methods in environments with highly non-monotonic payoff.

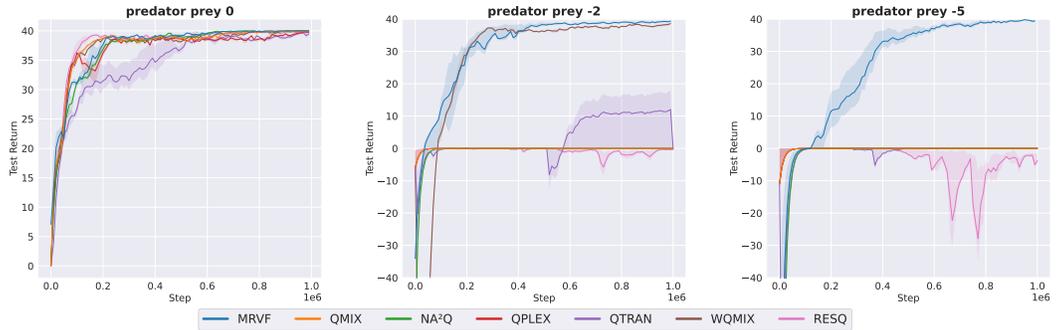


Figure 6: Test return in the predator prey tasks with punishments 0 (left), -2 (middle), and -5 (right). The non-monotonicity of the payoff increases with the punishment.

### 5.3 STARCRAFT II MULTI-AGENT CHALLENGE

We conduct experiments on the SMAC benchmark in scenarios of varying difficulty. Monotonic value factorization methods, including QMIX (Rashid et al., 2020b) and NA²Q (Liu et al., 2023),

are particularly well-suited for the SMAC because changes in an individual agent’s action within a single step have little impact on overall performance (see Appendix E.3 for details). Nevertheless, as shown in Figure 7, MRVF even achieves an advantage over monotonic value factorization methods, particularly in *3s\_vs\_5z*, *2c\_vs\_64zg* and *bane\_vs\_bane* scenarios. Meanwhile, the final performance of MRVF is the best in most scenarios, as shown in Table 3. Combined with the results in the *risk-reward* games and the predator prey task, we conclude that MRVF performs the best in both monotonic and non-monotonic scenarios.

In contrast, other methods (QPLEX (Wang et al., 2021), ResQ (Shen et al., 2022)), QTRAN (Son et al., 2019) and WQMIX (Rashid et al., 2020a)) often converge to suboptimal stable points. This leads to a high variance in performance within a scenario and across scenarios. Within the same scenario, although they occasionally achieve strong performance (as seen in the upper bounds of the shaded regions in Figure 7), their average performance lags behind. In addition, across different scenarios, their performance may vary drastically. For example, WQMIX performs well in *bane\_vs\_bane* but poorly in both *5m\_vs\_6m* and *3s\_vs\_5z*. Therefore, MRVF achieves significant improvements in robustness and performance over them in almost all scenarios.

Methods	MRVF	QMIX	NA <sup>2</sup> Q	QPLEX	QTRAN	WQMIX	RESQ
Best	<b>6</b>	3	4	1	1	1	3
Worst	<b>0</b>	1	0	3	5	2	1

Table 3: The number of SMAC scenarios where each method performs the best and the worst. We evaluate the final performance (averaged over steps after 1.5M) with a tolerance of  $\pm 0.05$ .

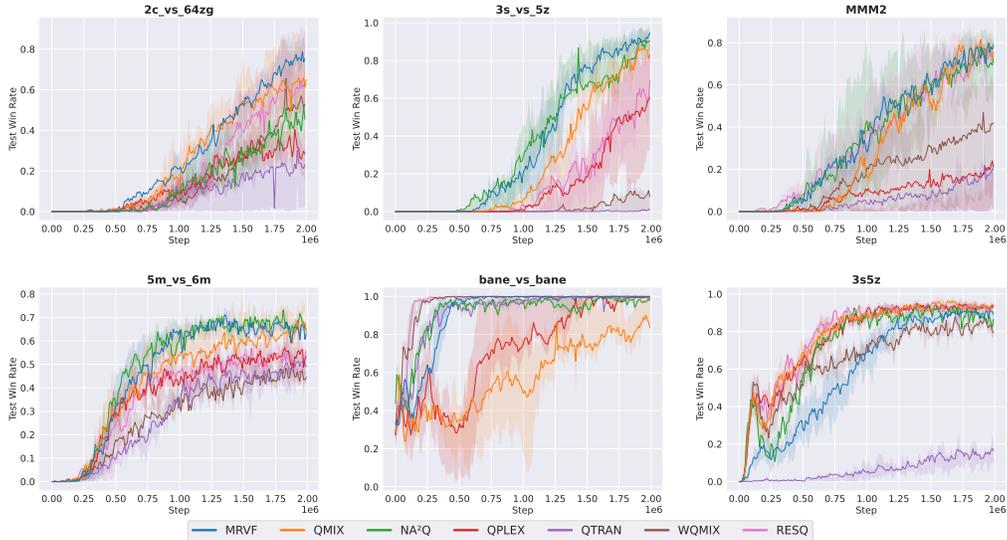


Figure 7: Test win rate in the SMAC benchmarks.

## 6 CONCLUSION

In this paper, we introduce a novel theoretical tool for studying the convergence of greedy action under value factorization, and propose MRVF, a novel framework for cooperative multi-agent reinforcement learning. In Section 3, we use this theoretical tool to derive the condition for global optimality in single-round factorization, and provide examples to demonstrate why existing methods struggle to satisfy this condition. In Section 4, we propose the condition for global optimality in multi-round factorization, which is strictly improving the greedy action round by round. To satisfy this condition, we design the forward and backward computation of MRVF. In addition, we design new sampling strategies suitable for multi-round factorization to ensure training stability. Experiments on non-monotonic one-step games, predator-prey tasks, and the StarCraft II Multi-Agent Challenge show the superior performance and robustness of MRVF compared to state-of-the-art value factorization methods.

## REFERENCES

- 486  
487  
488 Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic  
489 language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.
- 490  
491 Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy  
492 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large  
493 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 494  
495 Wendelin Böhmer, Vitaly Kurin, and Shimon Whiteson. Deep coordination graphs. In *International  
496 Conference on Machine Learning*, pp. 980–991. PMLR, 2020.
- 497  
498 Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of  
499 gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- 500  
501 Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and  
502 Joelle Pineau. Tarmac: Targeted multi-agent communication. In *International Conference on  
503 machine learning*, pp. 1538–1546. PMLR, 2019.
- 504  
505 Gang Ding, Zeyuan Liu, Zhirui Fang, Kefan Su, Liwen Zhu, and Zongqing Lu. Multi-agent coordi-  
506 nation via multi-level communication. *Advances in Neural Information Processing Systems*, 37:  
507 118513–118539, 2024.
- 508  
509 Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Mahajan,  
510 Jakob Foerster, and Shimon Whiteson. Smacv2: An improved benchmark for cooperative multi-  
511 agent reinforcement learning. *Advances in Neural Information Processing Systems*, 36:37567–  
512 37593, 2023.
- 513  
514 Jakob Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. Learning  
515 to communicate with deep multi-agent reinforcement learning. *Advances in neural information  
516 processing systems*, 29, 2016.
- 517  
518 Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.  
519 Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial  
520 intelligence*, volume 32, 2018.
- 521  
522 Wei Fu, Chao Yu, Zelai Xu, Jiaqi Yang, and Yi Wu. Revisiting some common practices in cooper-  
523 ative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp.  
524 6863–6877. PMLR, 2022.
- 525  
526 Shengchao Hu, Li Shen, Ya Zhang, and Dacheng Tao. Learning multi-agent communication from  
527 graph modeling perspective. In *International Conference on Learning Representations*, 2024.
- 528  
529 Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. Deep reinforcement learning for  
530 swarm systems. *Journal of Machine Learning Research*, 20(54):1–31, 2019.
- 531  
532 Chuming Li, Jie Liu, Yinmin Zhang, Yuhong Wei, Yazhe Niu, Yaodong Yang, Yu Liu, and Wanli  
533 Ouyang. Ace: Cooperative multi-agent q-learning with bidirectional action-dependency. In *Pro-  
534 ceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 8536–8544, 2023.
- 535  
536 Zichuan Liu, Yuanyang Zhu, and Chunlin Chen. NA<sup>2</sup>Q: Neural attention additive model for inter-  
537 pretable multi-agent q-learning. In *International Conference on Machine Learning*, pp. 22539–  
538 22558. PMLR, 2023.
- 539  
Anuj Mahajan, Tabish Rashid, Mikayel Samvelyan, and Shimon Whiteson. Maven: Multi-agent  
variational exploration. *Advances in neural information processing systems*, 32, 2019.
- D Mguni, T Jafferjee, J Wang, O Slumbers, N Perez-Nieves, F Tong, L Yang, J Zhu, and Y Yang.  
Ligs: Learnable intrinsic-reward generation selection for multi-agent learning. In *ICLR 2022-10th  
International Conference on Learning Representations*, volume 2022. ICLR, 2022.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word represen-  
tations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- 540 Yaru Niu, Rohan R Paleja, and Matthew C Gombolay. Multi-agent graph-attention communication  
541 and teaming. In *AAMAS*, volume 21, pp. 20th, 2021.
- 542 Afshin Oroojlooy and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement  
543 learning. *Applied Intelligence*, 53(11):13677–13722, 2023.
- 544 Ling Pan, Tabish Rashid, Bei Peng, Longbo Huang, and Shimon Whiteson. Regularized softmax  
545 deep multi-agent q-learning. *Advances in Neural Information Processing Systems*, 34:1365–1377,  
546 2021.
- 547 Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: expanding  
548 monotonic value function factorisation for deep multi-agent reinforcement learning. In *Proceed-*  
549 *ings of the 34th International Conference on Neural Information Processing Systems*, pp. 10199–  
550 10210, 2020a.
- 551 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob Foerster,  
552 and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement  
553 learning. *The Journal of Machine Learning Research*, 21(1):7234–7284, 2020b.
- 554 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
555 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 556 Siqi Shen, Mengwei Qiu, Jun Liu, Wei-quan Liu, Yongquan Fu, Xinwang Liu, and Cheng Wang.  
557 Resq: A residual q function-based approach for multi-agent reinforcement learning value factor-  
558 ization. *Advances in Neural Information Processing Systems*, 35:5471–5483, 2022.
- 559 Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran: Learning  
560 to factorize with transformation for cooperative multi-agent reinforcement learning. In *Internat-*  
561 *ional conference on machine learning*, pp. 5887–5896. PMLR, 2019.
- 562 Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropaga-  
563 tion. *Advances in neural information processing systems*, 29, 2016.
- 564 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max  
565 Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition  
566 networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- 567 Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings*  
568 *of the tenth international conference on machine learning*, pp. 330–337, 1993.
- 569 Lesong Tao, Miao Kang, Jinpeng Dong, Songyi Zhang, Ke Ye, Shitao Chen, and Nanning Zheng.  
570 Poaql: A partially observable altruistic q-learning method for cooperative multi-agent reinforce-  
571 ment learning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp.  
572 15076–15082. IEEE, 2024.
- 573 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
574 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
575 *tion processing systems*, 30, 2017.
- 576 Lipeng Wan, Zeyang Liu, Xingyu Chen, Xuguang Lan, and Nanning Zheng. Greedy based value  
577 representation for optimal coordination in multi-agent reinforcement learning. In *International*  
578 *Conference on Machine Learning*, pp. 22512–22535. PMLR, 2022.
- 579 Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling  
580 multi-agent q-learning. In *International Conference on Learning Representations*, 2021.
- 581 Jianhong Wang, Yuan Zhang, Yunjie Gu, and Tae-Kyun Kim. Shaq: Incorporating shapley value  
582 theory into multi-agent q-learning. *Advances in Neural Information Processing Systems*, 35:  
583 5941–5954, 2022.
- 584 Tonghan Wang, Jianhao Wang, Chongyi Zheng, and Chongjie Zhang. Learning nearly decompos-  
585 able value functions via communication minimization. In *International Conference on Learning*  
586 *Representations*, 2020.

594 S Whiteson, M Samvelyan, T Rashid, CS De Witt, G Farquhar, N Nardelli, TGJ Rudner, CM Hung,  
595 PHS Torr, and J Foerster. The starcraft multi-agent challenge. In *Proceedings of the International*  
596 *Joint Conference on Autonomous Agents and Multiagent Systems, AAMAS*, pp. 2186–2188, 2019.  
597

598 Yaodong Yang, Jianye Hao, Ben Liao, Kun Shao, Guangyong Chen, Wulong Liu, and Hongyao  
599 Tang. Qatten: A general framework for cooperative multiagent reinforcement learning. *arXiv*  
600 *preprint arXiv:2002.03939*, 2020.

601 Jianing Ye, Chenghao Li, Jianhao Wang, and Chongjie Zhang. Towards global optimality in cooper-  
602 ative marl with the transformation and distillation framework. *arXiv preprint arXiv:2207.11143*,  
603 2022.

604 Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The  
605 surprising effectiveness of ppo in cooperative multi-agent games. *Advances in neural information*  
606 *processing systems*, 35:24611–24624, 2022.

607

608 Yifan Zang, Jinmin He, Kai Li, Haobo Fu, Qiang Fu, Junliang Xing, and Jian Cheng. Automatic  
609 grouping for efficient cooperative multi-agent reinforcement learning. *Advances in neural infor-*  
610 *mation processing systems*, 36:46105–46121, 2023.

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

## A RELATED WORK

Reinforcement learning has been extensively studied in single-agent scenarios. However, when extended to multi-agent scenarios, a key issue is that the action space grows exponentially with the number of agents. To address this issue, existing approaches can be broadly categorized into two paradigms. The first paradigm treats other agents as part of the environment and focuses on learning individual value functions or policies. The second paradigm considers all agents as a unified entity and learns joint value functions with specific forms.

### A.1 OVERVIEW OF EXISTING WORK

**Independent Learning Methods** Among the first category of methods, IQL (Tan, 1993) is one of the earliest methods where agents independently learn individual action values. However, this method does not account for the dynamic policies of other agents. Subsequent work addresses this limitation by learning individual policies with joint action values. For example, COMA (Foerster et al., 2018) introduces a counterfactual baseline in actor-critic framework, while MAPPO (Yu et al., 2022) adapts PPO (Schulman et al., 2017) to multi-agent settings.

**Value Factorization Methods** The second category of methods factorizes the joint value function into individual components. Monotonic factorization is an intuitive way to satisfy the IGM principle. For example, VDN (Sunehag et al., 2017) represents the joint action value as the sum of individual action values. QMIX (Rashid et al., 2020b) represents the joint action value with a monotonic network that uses a positive weighted linear to deal with individual action values. Qat-tan (Yang et al., 2020) uses multi-head attention (Vaswani et al., 2017) to generate mixing weights. NA<sup>2</sup>Q (Liu et al., 2023) improves interpretability through Taylor expansion-based monotonic factorization. However, monotonic factorizations may suffer from representational limitations. To address this issue, QTRAN (Son et al., 2019) and WQMIX (Rashid et al., 2020a) pay more attention to the estimation of greedy action values. Specifically, QTRAN introduces a compensatory base to relax the non-optimal errors, while WQMIX reduces the weight on approximation of non-optimal parts. Other methods like QPLEX (Wang et al., 2021) and ResQ (Shen et al., 2022) design complete factorizations of the joint action value. QPLEX uses the maximum of a mixer value and a non-positive advantage function. ResQ combines a monotonic function with a non-positive function. GVR (Wan et al., 2022) takes a different approach by making  $u^*$  the only stable point, although this requires numerous approximations.

**Communication-based Methods** Other researchers study communication mechanisms in multi-agent systems. CommNet (Sukhbaatar et al., 2016) enables multi-round communication for sharing observations. DIAL (Foerster et al., 2016) designs gradient-based messages exchanged between agents. To decide with whom to communicate, TarMAC (Das et al., 2019) uses a signature-based soft attention mechanism, and MAGIC (Niu et al., 2021) uses graph neural networks. NDQ (Wang et al., 2020) combines value factorization and communication by mixing individual  $Q$ s through communication. ACE (Li et al., 2023) models communication as intermediate processes in the Markov Decision Process. In addition to parallel decision making, recent methods study sequential decision making through communication. PG-AR (Fu et al., 2022) randomizes the order of decision making. SeqComm (Ding et al., 2024) use intentions to prioritize agents in sequential decision-making.

**Other MARL Methods** Recent works introduce additional innovations in MARL. RES (Pan et al., 2021) addresses value overestimation using regularized softmax losses. MAVEN (Mahajan et al., 2019) enhances exploration efficiency through randomized latent vectors. LIGS (Mguni et al., 2022) proposes learnable intrinsic rewards to improve multi-agent cooperation. SHAQ (Wang et al., 2022) integrates the Shapley value by modeling  $Q$  functions as marginal contributions of agents.

### A.2 RELATIONSHIP TO EXISTING WORK

**Relationship to Weighted QMIX** Both Weighted QMIX (Rashid et al., 2020a) and the single-round value factorization in our method aim to enhance convergence toward actions superior to a given action. Weighted QMIX achieves this by reducing the fitting weight for inferior actions,

702 whereas our method assigns them the minimum target value. However, Weighted QMIX cannot  
 703 guarantee strict improvement over the given action, as suboptimal actions may be stable points. In  
 704 contrast, our method ensures strict improvement by preventing actions with minimum target value  
 705 from becoming stable points. Moreover, even when Weighted QMIX obtains the optimal action, it  
 706 may fail to stabilize at this optimum because the optimal action might not be a stable point. Although  
 707 our method similarly risks deviating from the optimal action, we mitigate this by early termination  
 708 when no further improvement is detected.

709  
 710 **Relationship to GVR** The main idea of GVR (Wan et al., 2022) is to establish the optimal action  
 711 as a stable point while rendering other actions non-stable. However, this approach requires a prior  
 712 knowledge of the optimal action, which is fundamentally infeasible in MARL since finding the opti-  
 713 mal action is the primary objective. Consequently, GVR inevitably relies on approximations of ideal  
 714 conditions, which not only increase the complexity but may also convert certain suboptimal actions  
 715 into stable points. In contrast to GVR, our method converts inferior actions into unstable points  
 716 in each iteration. Determining inferior actions is computationally straightforward, and by assign-  
 717 ing them the minimum target value, existing monotonic value factorizations can effectively exclude  
 718 these inferior actions. Therefore, our method strictly improves the current action through iterations,  
 719 providing a more reliable approximation to the optimal action. Furthermore, GVR only defines the  
 720 concept of stable points within the policy distribution, and its analysis is limited to a particular form  
 721 of QMIX. In contrast, we not only explicitly identify the origin of stable points (characterized by the  
 722 lag between changes in non-differential variables and gradient-based parameters) but also provide a  
 723 formal definition. In addition, we provide a comprehensive analysis on existing value factorization  
 724 methods, unifying their suboptimal performance under a common explanation: the presence of at  
 725 least one non-optimal stable point.

726  
 727 **Relationship to Communication** Our method can also be interpreted as a multi-round commu-  
 728 nication algorithm within a parallel execution framework. In each iteration, agents exchange pre-  
 729 decision information and generate new decisions through communication. The key challenge in such  
 730 frameworks is ensuring that post-communication actions strictly improve upon pre-communication  
 731 actions - failure to achieve this prevents agents from reaching consensus, resulting in suboptimal  
 732 solutions. In our method, by rendering actions inferior to pre-communication ones as non-stable  
 733 points, our approach not only guarantees the improvement of post-communication actions but also  
 734 approaches the optimal action within sufficient communication rounds. In contrast, sequential exe-  
 735 cution frameworks essentially extend single-step action selection across agents, where each agent’s  
 736 action depends on higher-priority agents’ choices. Such frameworks not only demonstrate lower  
 737 efficiency but also introduce additional concerns regarding policy convergence and suboptimality  
 738 caused by execution ordering (Ding et al., 2024).

## 738 B DEFINITION

739  
 740 In this section, we will define some important concepts that have been used in previous work but  
 741 have not yet been rigorously defined.

### 742 B.1 MONOTONIC PAYOFF

743  
 744 Note that the definition of monotonicity for payoffs differs from that in functions of real variables,  
 745 as there is no ordering relationship defined in the action space, the domain of the payoff. We define  
 746 an order relation on the actions based on the function  $F$  as follows:

747 **Definition 2** (Order Relation  $\succeq$ ). *Let  $u_i^1, u_i^2 \in \mathcal{U}_i$  be two actions of agent  $i$ . For the function  $F$ , if*  
 748 *the order relation between  $u_i^1$  and  $u_i^2$  is  $u_i^1 \succeq u_i^2$ , then*

$$749 \quad F(u_i^1, \mathbf{u}_{-i}) \geq F(u_i^2, \mathbf{u}_{-i}), \forall \mathbf{u}_{-i} \in \mathcal{U}_{-i} \quad (9)$$

750 where  $\mathbf{u}_{-i} = (u_1, \dots, u_{i-1}, u_{i+1}, \dots, u_n)$  is the joint action without  $u_i$ .

751  
 752 The function  $F$  can be replaced with  $Q_{jt}$ , a measurement of payoff. However, the order relation  
 753 based on  $Q_{jt}$  exists over the entire action space only if  $Q_{jt}$  is monotonic, which is defined as follows:

754 **Definition 3** (Monotonic  $Q_{jt}$ ).  *$\forall s \in \mathcal{S}$ , if for every agent  $i$  and any two actions  $u_i^1, u_i^2 \in \mathcal{U}_i$ , there*  
 755 *exists an order relation that is either  $u_i^1 \succeq u_i^2$  or  $u_i^2 \succeq u_i^1$ , then  $Q_{jt}$  is monotonic.*

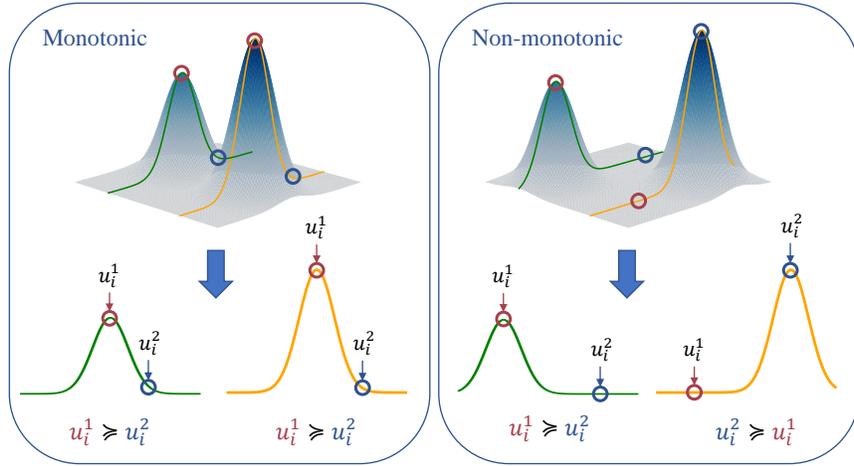


Figure 8: An illustration of Definition 3. We intercept curves corresponding to different  $\mathbf{u}_{-i}$  and check whether their monotonicity with respect to  $u_i$  is consistent.

Table 4 shows the payoff matrices with different monotonicity. Monotonic factorization methods can easily obtain the optimal solution in monotonic cases ((a) and (b) in Table 4) and sometimes in non-monotonic cases ((c) in Table 4). However, monotonic factorizations struggle to obtain the optimal solution in highly non-monotonic cases ((d) in Table 4).

9	8	7
6	5	4
3	2	1

(a) Monotonic

8	9	7
3	2	1
5	6	4

(b) Monotonic

9	0	0
0	5	0
0	0	0

(c) Non-monotonic

9	-9	-9
-9	5	0
-9	0	0

(d) Highly non-monotonic

Table 4: Payoff matrices with different monotonicity, where the value in row  $r$  and column  $c$  denotes the payoff of joint action  $\mathbf{u} = (r, c)$ . (a): A monotonic payoff matrix. (b) A monotonic payoff matrix generated by rearranging the rows and columns of (a). (c) A non-monotonic payoff matrix, where the relation between any two rows or columns is partially ordered. (d) A highly non-monotonic payoff matrix, where the relation between any two rows or columns is largely unordered.

## B.2 IDEAL QMIX

The objective of QMIX is to find individual  $Q$ s and  $Q_{\text{mon}}$  that minimize the MSE between  $Q_{\text{tot}}$  and  $Q_{\text{jt}}$ , which is defined in Equation (10) for certain  $\mathbf{s} \in \mathcal{S}^3$ .

$$L_{\text{tot}}(\mathbf{s}) = \sum_{\mathbf{u}} \pi(\mathbf{u}|\mathbf{s})(Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) - Q_{\text{jt}}(\mathbf{s}, \mathbf{u}))^2 \quad (10)$$

Since the relationship between individual  $Q$ s and  $Q_{\text{mon}}$  is monotonic, the magnitude relationship of  $Q_{\text{mon}}$  is constrained by that of individual  $Q$ s, which is  $\forall i \in \{1, 2, \dots, n\}, \forall u_i, v_i \in \mathcal{U}_i, \mathbf{u}_{-i} \in \mathcal{U}_{-i} = \mathcal{U}_1 \times \dots \times \mathcal{U}_{i-1} \times \mathcal{U}_{i+1} \times \dots \times \mathcal{U}_n$

$$(Q_{\text{mon}}(\mathbf{s}, \mathbf{u}_{-i}, u_i) - Q_{\text{mon}}(\mathbf{s}, \mathbf{u}_{-i}, v_i))(Q_i(\mathbf{s}, u_i) - Q_i(\mathbf{s}, v_i)) \geq 0 \quad (11)$$

Combining the objective defined in Equation (10) and the constrain defined in Equation (11), we get the ideal optimization problem that QMIX solves. In practice, QMIX uses the network as a monotonic function  $f_{\text{mon}}$  to represent the constraint between individual  $Q$ s and  $Q_{\text{mon}}$ , and applies stochastic gradient descent (SGD) to minimize the objective. In this way,  $Q_i$  is updated based on its gradient which involves  $\frac{\partial f_{\text{mon}}}{\partial Q_i}$ . However, since how  $\frac{\partial f_{\text{mon}}}{\partial Q_i}$  varies is unknown, to avoid this confounding factor, we assume that QMIX can obtain the optimal combination of individual  $Q$ s and

<sup>3</sup>To simplify the discussion, we focus on the fully observable environment by replacing  $\tau$  with  $\mathbf{s}$  in  $Q_{\text{tot}}$

$Q_{\text{mon}}$  that satisfies Equation (11), and consequently minimize  $L_{\text{tot}}$ . We name QMIX based on this assumption as the ideal QMIX which is defined as follows <sup>4</sup>:

**Definition 4** (Ideal QMIX). *Let  $\mathcal{Q}$  be the set of pairs  $(\mathbf{Q}, Q_{\text{mon}})$ , where the individual action values  $\mathbf{Q} = (Q_1, Q_2, \dots, Q_n)$  and  $Q_{\text{mon}}$  satisfy Equation (11). For a joint action value  $Q_{\text{jt}}$  and policy  $\pi$ ,  $\mathbf{Q}$  and  $Q_{\text{mon}}$  of the ideal QMIX converge to  $\mathbf{Q}^*$  and  $Q_{\text{mon}}^*$  which satisfy  $\forall \mathbf{s} \in \mathcal{S}$*

$$\mathbf{Q}^*, Q_{\text{mon}}^* \in \arg \min_{(\mathbf{Q}, Q_{\text{mon}}) \in \mathcal{Q}} \sum_{\mathbf{u}} \pi(\mathbf{u}|\mathbf{s}) (Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) - Q_{\text{jt}}(\mathbf{s}, \mathbf{u}))^2 \quad (12)$$

where  $Q_{\text{tot}} = Q_{\text{mon}}$  in QMIX.

As for other value factorization methods such as WQMIX (Rashid et al., 2020a) and ResQ (Shen et al., 2022) where  $Q_{\text{mon}}$  appears in their design, we also assume they share the same property as the ideal QMIX, which aligns with the assumption used in their papers.

From Definition 4, the ideal QMIX obtains a monotonic representation of  $Q_{\text{jt}}$  with the minimum  $L_{\text{tot}}$ . We illustrate this process in matrix cases: Given the order of rows and columns, we can find the monotonic matrix with local minimum  $L_{\text{tot}}$  by solving the constraint optimization problem. To find the monotonic matrix with minimum  $L_{\text{tot}}$ , we compare  $L_{\text{tot}}$  for all possible orders. We give an example of the ideal QMIX shown in Table 5.

	$L_{\text{tot}}=36.22$	$L_{\text{tot}}=45.56$																											
<table border="1" style="border-collapse: collapse;"> <tr><td>8</td><td>-12</td><td>-12</td></tr> <tr><td>-12</td><td>3</td><td>0</td></tr> <tr><td>-12</td><td>0</td><td>5</td></tr> </table>	8	-12	-12	-12	3	0	-12	0	5	<table border="1" style="border-collapse: collapse;"> <tr><td>-8</td><td>-8</td><td>-8</td></tr> <tr><td>-8</td><td>1</td><td>1</td></tr> <tr><td>-8</td><td>1</td><td><b>5</b></td></tr> </table>	-8	-8	-8	-8	1	1	-8	1	<b>5</b>	<table border="1" style="border-collapse: collapse;"> <tr><td><b>8</b></td><td>-5</td><td>-5</td></tr> <tr><td>-5</td><td>-5</td><td>-5</td></tr> <tr><td>-5</td><td>-5</td><td>-5</td></tr> </table>	<b>8</b>	-5	-5	-5	-5	-5	-5	-5	-5
8	-12	-12																											
-12	3	0																											
-12	0	5																											
-8	-8	-8																											
-8	1	1																											
-8	1	<b>5</b>																											
<b>8</b>	-5	-5																											
-5	-5	-5																											
-5	-5	-5																											
(a) $Q_{\text{jt}}$	(b) $Q_{\text{tot}}^*$	(c) $Q_{\text{tot}}^*$ for $\bar{\mathbf{u}} = \mathbf{u}^*$																											

Table 5: (a): A non-monotonic payoff matrix, where the value presented in row  $r$  and column  $c$  denotes the  $Q_{\text{jt}}$  value of joint action  $(r, c)$ . (b): The monotonic representation with the global minimum  $L_{\text{tot}}$ . (c) The optimal monotonic representation constrained by  $\bar{\mathbf{u}} = \mathbf{u}^*$ . Since the  $L_{\text{tot}}$  of (b) is less than the  $L_{\text{tot}}$  of (c),  $Q_{\text{tot}}$  of the ideal QMIX under uniform visitation ( $\epsilon = 1$ ) will converge to (b), which fails to get the global optimal action.

## C PROOF

### C.1 THEOREM ON MRVF

**Theorem 1.** *For a joint action value  $Q_{\text{jt}}$  with a finite action space  $\mathcal{U}$  and  $\mathbf{s} \in \mathcal{S}$ , consider a sequence  $\{\bar{\mathbf{u}}^k | \bar{\mathbf{u}}^k \in \mathcal{U}\}$ , if  $\forall k > 1, Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^k) > Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^{k-1})$  when  $\bar{\mathbf{u}}^{k-1} \neq \mathbf{u}^*$ , then  $\exists K > 0, \bar{\mathbf{u}}^K = \mathbf{u}^*$ .*

*Proof.* Let  $\mathcal{U}_+^k = \{\mathbf{u} | Q_{\text{jt}}(\mathbf{s}, \mathbf{u}) > Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^k), \mathbf{u} \in \mathcal{U}\}$ . Assume that  $\forall k > 0, \bar{\mathbf{u}}^k \neq \mathbf{u}^*$ . Since  $\forall k > 1, Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^k) > Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^{k-1})$  when  $\bar{\mathbf{u}}^{k-1} \neq \mathbf{u}^*$ ,  $\mathcal{U}_+^k \subset \mathcal{U}_+^{k-1}$ . Therefore, we have  $|\mathcal{U}_+^{k-1}| - |\mathcal{U}_+^k| \geq 1$ . Then, we have  $\forall k > 0, |\mathcal{U}_+^k| \leq |\mathcal{U}_+^1| - (k-1)$ . If we choose  $K > |\mathcal{U}| + 1$ , since  $\mathcal{U}_+^1 \subseteq \mathcal{U}$ , we get  $|\mathcal{U}_+^K| < |\mathcal{U}_+^1| - |\mathcal{U}| \leq 0$ . Since  $\forall k > 0, |\mathcal{U}_+^k| \geq 0$ , the assumption is false. Therefore, we prove that  $\exists K > 0, \bar{\mathbf{u}}^K = \mathbf{u}^*$ .  $\square$

**Lemma 1.** *For the converged  $Q_{\text{tot}}$  of the idea QMIX defined in Equation (12) and  $\epsilon$ -greedy policy  $\pi$  that is continuous for  $\epsilon \in [0, 1]$  and satisfies Equation (13), we have  $\lim_{\epsilon \rightarrow 0^+} Q_{\text{tot}}^*(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{jt}}(\mathbf{s}, \tilde{\mathbf{u}})$  for  $\mathbf{s} \in \mathcal{S}$ .*

$$\lim_{\epsilon \rightarrow 0^+} \pi(\mathbf{u}|\mathbf{s}) = \begin{cases} 1 & \mathbf{u} = \tilde{\mathbf{u}} \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

<sup>4</sup>The policy  $\pi$  is omitted in the definition by Rashid et al. (2020a). However, it plays a crucial role in shaping  $Q_{\text{tot}}$  and will be extensively analyzed in our proof.

*Proof.* Assume that  $\exists \Delta > 0, \forall E > 0, \exists \epsilon \in (0, E), |Q_{\text{tot}}^*(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \tilde{\mathbf{u}})| \geq \Delta$ . The loss of  $Q_{\text{tot}}^*$  defined in Equation (12) is

$$\begin{aligned} L_{\text{tot}}^* &= \sum_{\mathbf{u}} \pi(\mathbf{u}|\mathbf{s})(Q_{\text{tot}}^*(s, \mathbf{u}) - Q_{\text{jt}}(s, \mathbf{u}))^2 \\ &= \sum_{\mathbf{u} \neq \tilde{\mathbf{u}}} \pi(\mathbf{u}|\mathbf{s})(Q_{\text{tot}}^*(s, \mathbf{u}) - Q_{\text{jt}}(s, \mathbf{u}))^2 + \pi(\tilde{\mathbf{u}}|\mathbf{s})(Q_{\text{tot}}^*(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \tilde{\mathbf{u}}))^2 \\ &\geq \pi(\tilde{\mathbf{u}}|\mathbf{s})\Delta^2 \end{aligned} \quad (14)$$

Let  $Q'_{\text{tot}}(s, \cdot) \equiv \frac{\Delta}{2} + Q_{\text{jt}}(s, \tilde{\mathbf{u}})$ . The loss of  $Q'_{\text{tot}}$  is

$$L'_{\text{tot}} = \sum_{\mathbf{u} \neq \tilde{\mathbf{u}}} \pi(\mathbf{u}|\mathbf{s})\left(\frac{\Delta}{2} + Q_{\text{jt}}(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \mathbf{u})\right)^2 + \pi(\tilde{\mathbf{u}}|\mathbf{s})\left(\frac{\Delta}{2}\right)^2 \quad (15)$$

Consider  $L_{\text{tot}}^* - L'_{\text{tot}}$ , we have

$$L_{\text{tot}}^* - L'_{\text{tot}} \geq \pi(\tilde{\mathbf{u}}|\mathbf{s})\left(\frac{\Delta}{2}\right)^2 - \sum_{\mathbf{u} \neq \tilde{\mathbf{u}}} \pi(\mathbf{u}|\mathbf{s})\left(\frac{\Delta}{2} + Q_{\text{jt}}(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \mathbf{u})\right)^2 \quad (16)$$

Notice that

$$\lim_{\epsilon \rightarrow 0^+} \left(\pi(\tilde{\mathbf{u}}|\mathbf{s})\left(\frac{\Delta}{2}\right)^2 - \sum_{\mathbf{u} \neq \tilde{\mathbf{u}}} \pi(\mathbf{u}|\mathbf{s})\left(\frac{\Delta}{2} + Q_{\text{jt}}(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \mathbf{u})\right)^2\right) = \left(\frac{\Delta}{2}\right)^2 > 0 \quad (17)$$

Therefore,  $\exists E > 0, \forall \epsilon \in (0, E), L_{\text{tot}}^* > L'_{\text{tot}}$  which is against with condition that  $Q_{\text{tot}}^*$  satisfies Equation (12). Thus, the assumption is false. We have  $\forall \Delta > 0, \exists E > 0, \forall \epsilon \in (0, E), |Q_{\text{tot}}^*(s, \tilde{\mathbf{u}}) - Q_{\text{jt}}(s, \tilde{\mathbf{u}})| < \Delta$  which means  $\lim_{\epsilon \rightarrow 0^+} Q_{\text{tot}}^*(s, \tilde{\mathbf{u}}) = Q_{\text{jt}}(s, \tilde{\mathbf{u}})$ .  $\square$

**Theorem 2.** Consider a non-constant joint action value  $Q_{\text{jt}}(s, \cdot)$  ( $\forall s \in \mathcal{S}, \forall C \in \mathbb{R}, Q_{\text{jt}}(s, \cdot) \neq C$ ) with finite state space  $\mathcal{S}$  and action space  $\mathcal{U}$ . For the ideal QMIX defined in Equation (12) and the centralized  $\epsilon$ -greedy policy  $\pi$  defined in Equation (8), we have  $\forall s \in \mathcal{S}, \exists E \in (0, 1), \forall \epsilon \in (0, E), \forall \mathbf{u} \in \arg \min_{\mathbf{u}} Q_{\text{jt}}(s, \mathbf{u}), \mathbf{u}$  is not a stable point.

*Proof.* For  $s \in \mathcal{S}$ , Assume that  $\forall E_s \in (0, 1], \exists \epsilon \in (0, E_s), \exists Q_{\text{tot}}^*$  that  $Q_{\text{jt}}(s, \bar{\mathbf{u}}^*) = \min_{\mathbf{u}} Q_{\text{jt}}(s, \mathbf{u})$  and  $\tilde{\mathbf{u}} = \bar{\mathbf{u}}^*$  ( $\bar{\mathbf{u}}^*$  is a stable point). We define  $\mathcal{U}^+ = \{\mathbf{u} | Q_{\text{jt}}(s, \mathbf{u}) > \min_{\mathbf{u}} Q_{\text{jt}}(s, \mathbf{u})\}$  and  $\mathcal{U}^- = \mathcal{U} - \mathcal{U}^+$ . Since  $\lim_{\epsilon \rightarrow 0^+} Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) = Q_{\text{jt}}(s, \bar{\mathbf{u}}^*)$  (Lemma 1 and  $\tilde{\mathbf{u}} = \bar{\mathbf{u}}^*$ ), we have  $\exists E_s^0 \in (0, 1], \forall \epsilon \in (0, E_s^0), Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) < \min_{\mathbf{u} \in \mathcal{U}^+} Q_{\text{jt}}(s, \mathbf{u})$ . Thus,  $\forall \mathbf{u} \in \mathcal{U}^+$ , since  $Q_{\text{tot}}^*(s, \mathbf{u}) \leq Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) < Q_{\text{jt}}(s, \mathbf{u})$ , consider a term of  $L_{\text{tot}}^*$ , we have

$$\begin{aligned} L_{\text{tot}}^*(\mathbf{u}) &= \pi(\mathbf{u}|\mathbf{s})(Q_{\text{tot}}^*(s, \mathbf{u}) - Q_{\text{jt}}(s, \mathbf{u}))^2 \\ &> \frac{\epsilon}{|\mathcal{U}|}(Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \mathbf{u}))^2 \end{aligned} \quad (18)$$

Therefore, for the loss of  $Q_{\text{tot}}^*$  defined in Equation (12), we have

$$\begin{aligned} L_{\text{tot}}^* &= \sum_{\mathbf{u}} \pi(\mathbf{u}|\mathbf{s})(Q_{\text{tot}}^*(s, \mathbf{u}) - Q_{\text{jt}}(s, \mathbf{u}))^2 \\ &> \frac{\epsilon}{|\mathcal{U}|} \sum_{\mathbf{u} \in \mathcal{U}^+} (Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \mathbf{u}))^2 + (1 - \epsilon + \frac{\epsilon}{|\mathcal{U}|})(Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \bar{\mathbf{u}}^*))^2 \end{aligned} \quad (19)$$

Consider certain  $\mathbf{u}^+ \in \mathcal{U}^+$ . We define  $Q'_{\text{tot}}$  as

$$Q'_{\text{tot}}(s, \mathbf{u}) = \begin{cases} Q_{\text{jt}}(s, \mathbf{u}^+) & \mathbf{u} = \mathbf{u}^+ \\ Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) & \text{otherwise} \end{cases} \quad (20)$$

The loss of  $Q'_{\text{tot}}$  defined in Equation (12) is ( $\tilde{\mathbf{u}} = \bar{\mathbf{u}}^*$ )

$$\begin{aligned} L'_{\text{tot}} &= \frac{\epsilon}{|\mathcal{U}|} \left( \sum_{\mathbf{u} \in \mathcal{U}^+ - \{\mathbf{u}^+\}} (Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \mathbf{u}))^2 + \sum_{\mathbf{u} \in \mathcal{U}^-} (Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \bar{\mathbf{u}}^*))^2 \right) \\ &\quad + (1 - \epsilon + \frac{\epsilon}{|\mathcal{U}|})(Q_{\text{tot}}^*(s, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(s, \bar{\mathbf{u}}^*))^2 \end{aligned} \quad (21)$$

Comparing  $L_{\text{tot}}^*$  and  $L'_{\text{tot}}$ , we have

$$L_{\text{tot}}^* - L'_{\text{tot}} > \frac{\epsilon}{|\mathcal{U}|} ((Q_{\text{tot}}^*(\mathbf{s}, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(\mathbf{s}, \mathbf{u}^+))^2 - |\mathcal{U}^-| (Q_{\text{tot}}^*(\mathbf{s}, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^*))^2) \quad (22)$$

Notice that

$$f(x) = (x - Q_{\text{jt}}(\mathbf{s}, \mathbf{u}^+))^2 - |\mathcal{U}^-| (x - Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^*))^2 \quad (23)$$

is a continuous function and  $\lim_{x \rightarrow Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^*)} f(x) = (Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^*) - Q_{\text{jt}}(\mathbf{s}, \mathbf{u}^+))^2 > 0$ . Since

$\lim_{\epsilon \rightarrow 0^+} Q_{\text{tot}}^*(\mathbf{s}, \bar{\mathbf{u}}^*) = Q_{\text{jt}}(\mathbf{s}, \bar{\mathbf{u}}^*)$ , we have  $\exists E_s^1 \in (0, 1], \forall \epsilon \in (0, E_s^1), f(Q_{\text{tot}}^*(\mathbf{s}, \bar{\mathbf{u}}^*)) > 0$ . Therefore, we find  $E_s = \min\{E_s^0, E_s^1\} \in (0, 1]$  that  $\forall \epsilon \in (0, E_s) L_{\text{tot}}^* > L'_{\text{tot}}$ , which is against with the condition that  $Q_{\text{tot}}^*$  satisfies Equation (12). Thus,  $\exists E_s \in (0, 1], \forall \epsilon \in (0, E_s), \forall \mathbf{u} \in \arg \min_{\mathbf{u}} Q_{\text{jt}}(\mathbf{s}, \mathbf{u})$ ,  $\mathbf{u}$  is not a stable point. Finally, we complete the proof by taking  $E = \min_{s \in \mathcal{S}} E_s$ .  $\square$

## C.2 THEOREM ON RESQ

We will prove that  $\mathbf{u}^*$  is the unique yet weakly stable point of ResQ (Shen et al., 2022). Here, we present the expression of  $Q_{\text{tot}}$  of ResQ.

$$Q_{\text{tot}}(\mathbf{s}, \mathbf{u}) = Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) + w_{\text{r}}(\mathbf{s}, \mathbf{u}) * Q_{\text{r}}(\mathbf{s}, \mathbf{u}) \quad (24)$$

where  $Q_{\text{mon}} = f_{\text{mon}}(Q_1, Q_2, \dots, Q_n)$ ,  $Q_{\text{r}} \leq 0$  and

$$w_{\text{r}}(\mathbf{s}, \mathbf{u}) = \begin{cases} 0 & \mathbf{u} = \tilde{\mathbf{u}} \\ 1 & \text{otherwise} \end{cases} \quad (25)$$

We provide a specific instance for the weak stability of ResQ in Table 6. Then, we prove that the weak stability occurs in general cases in Theorem 3.

8	-12	-12
-12	3	0
-12	0	5

(a)  $Q_{\text{jt}}$

8	8	8
8	8	9
8	8	8

(b)  $Q_{\text{mon}}$

0	-20	-20
-20	-5	-9
-20	-8	-4

(c)  $w_{\text{r}} * Q_{\text{r}}$

Table 6: This table shows the weak stability of ResQ for a specific  $Q_{\text{jt}}$  in (a), where we let  $\tilde{\mathbf{u}} = \mathbf{u}^*$  initially to check whether ResQ will stabilize at  $\mathbf{u}^*$ . We highlight the values related to  $\tilde{\mathbf{u}}$  in red and  $\bar{\mathbf{u}}$  in blue.  $Q_{\text{tot}}$  is the sum of  $Q_{\text{mon}}$  in (b) and  $w_{\text{r}} * Q_{\text{r}}$  in (c), which is identical to  $Q_{\text{jt}}$  ( $L_{\text{tot}} = 0$ ). However, since  $\bar{\mathbf{u}} \neq \mathbf{u}^*$  may happens,  $\mathbf{u}^*$  is not a strongly stable point for ResQ.

**Lemma 2.** For a joint action value  $Q_{\text{jt}}$  and  $\mathbf{s} \in \mathcal{S}$ ,  $\forall \tilde{\mathbf{u}} \in \mathcal{U}$ , we can find  $Q_{\text{tot}}$  with  $L_{\text{tot}} = 0$  and  $\bar{\mathbf{u}} \neq \tilde{\mathbf{u}}$ .

*Proof.* We can find a  $Q_{\text{mon}}$  that satisfies

- (1)  $Q_{\text{mon}} \geq Q_{\text{jt}}$
- (2)  $Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{jt}}(\mathbf{s}, \tilde{\mathbf{u}})$
- (3)  $\bar{\mathbf{u}} \neq \tilde{\mathbf{u}}$

Here we give a specific  $Q_{\text{mon}}$  that satisfies these conditions:

$$Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) = \begin{cases} Q_{\text{jt}}(\mathbf{s}, \tilde{\mathbf{u}}) & \mathbf{u} = \tilde{\mathbf{u}} \\ \max_{\mathbf{u}} Q_{\text{jt}}(\mathbf{s}, \mathbf{u}) + \Delta & \text{otherwise} \end{cases} \quad (26)$$

where  $\Delta > 0$ . We let  $Q_{\text{r}} = Q_{\text{jt}} - Q_{\text{mon}}$ . Since  $Q_{\text{mon}} \geq Q_{\text{jt}}$ , we have  $Q_{\text{r}} \leq 0$ .

$\forall \mathbf{u} \neq \tilde{\mathbf{u}}$ , we have  $Q_{\text{jt}}(\mathbf{s}, \mathbf{u}) = Q_{\text{tot}}(\mathbf{s}, \mathbf{u}) = Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) + 1 * Q_{\text{r}}(\mathbf{s}, \mathbf{u})$ . For  $\tilde{\mathbf{u}}$ , since  $Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{jt}}(\mathbf{s}, \tilde{\mathbf{u}})$ , we have  $Q_{\text{jt}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{tot}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) + 0 * Q_{\text{r}}(\mathbf{s}, \tilde{\mathbf{u}})$ .

Therefore,  $\forall \mathbf{u} \in \mathcal{U}$ , we have  $Q_{\text{jt}}(\mathbf{s}, \mathbf{u}) = Q_{\text{tot}}(\mathbf{s}, \mathbf{u})$ , which means  $L_{\text{tot}} = 0$ . Thus, we find  $Q_{\text{tot}}$  with  $L_{\text{tot}} = 0$  and  $\bar{\mathbf{u}} \neq \tilde{\mathbf{u}}$ .  $\square$

**Theorem 3.** For a joint action value  $Q_{jt}$  and  $\mathbf{s} \in \mathcal{S}$ ,  $\forall \mathbf{u}^* \in \arg \max_{\mathbf{u}} Q_{jt}(\mathbf{s}, \mathbf{u})$ ,  $\mathbf{u}^*$  is a weakly stable point of  $\text{Res}Q$  (Shen et al., 2022), while  $\forall \mathbf{u}^- \notin \arg \max_{\mathbf{u}} Q_{jt}(\mathbf{s}, \mathbf{u})$  is not a stable point of  $\text{Res}Q$ .

*Proof.* First, we illustrate that  $\mathbf{u}^*$  is a stable point. **Assume** that  $\tilde{\mathbf{u}} = \mathbf{u}^*$ . We can find a  $Q_{\text{mon}}$  that satisfies

- (1)  $Q_{\text{mon}} \geq Q_{jt}$
- (2)  $Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{jt}(\mathbf{s}, \tilde{\mathbf{u}})$
- (3)  $\bar{\mathbf{u}} = \mathbf{u}^*$

where  $Q_{\text{mon}} \equiv Q_{jt}(\mathbf{s}, \mathbf{u}^*)$  is a specific case.

We let  $Q_r = Q_{jt} - Q_{\text{mon}}$ . Since  $Q_{\text{mon}} \geq Q_{jt}$ , we have  $Q_r \leq 0$ .

$\forall \mathbf{u} \neq \tilde{\mathbf{u}}$ , we have  $Q_{jt}(\mathbf{s}, \mathbf{u}) = Q_{\text{tot}}(\mathbf{s}, \mathbf{u}) = Q_{\text{mon}}(\mathbf{s}, \mathbf{u}) + 1 * Q_r(\mathbf{s}, \mathbf{u})$ . For  $\tilde{\mathbf{u}}$ , since  $Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{jt}(\mathbf{s}, \tilde{\mathbf{u}})$ , we have  $Q_{jt}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{tot}}(\mathbf{s}, \tilde{\mathbf{u}}) = Q_{\text{mon}}(\mathbf{s}, \tilde{\mathbf{u}}) + 0 * Q_r(\mathbf{s}, \tilde{\mathbf{u}})$ .

Therefore, we have  $L_{\text{tot}} = 0$ , which is the lower bound of loss. However, according to Lemma 2, we can find  $Q'_{\text{tot}}$  with  $L'_{\text{tot}} = L_{\text{tot}}$  while  $\bar{\mathbf{u}} \neq \mathbf{u}^*$ . Thus,  $\mathbf{u}^*$  is a weakly stable point of  $\text{Res}Q$ .

Second, we illustrate that  $\mathbf{u}^-$  is not a stable point. Let  $\tilde{\mathbf{u}} = \mathbf{u}^-$ , **Assume** that there exists  $Q_{\text{tot}}$  satisfying  $\bar{\mathbf{u}} = \mathbf{u}^-$  and achieving the minimum loss. According to Lemma 2, we have  $L_{\text{tot}} = 0$ , otherwise  $L_{\text{tot}}$  is not the minimum.

Since  $L_{\text{tot}} = 0$ , we have  $\forall \mathbf{u} \in \mathcal{U}$ ,  $Q_{\text{tot}}(\mathbf{s}, \mathbf{u}) = Q_{jt}(\mathbf{s}, \mathbf{u})$ . For  $\mathbf{u}^-$ , we have  $Q_{\text{tot}}(\mathbf{s}, \mathbf{u}^-) = Q_{jt}(\mathbf{s}, \mathbf{u}^-)$ . For  $\mathbf{u}^*$ , we have  $Q_{\text{tot}}(\mathbf{s}, \mathbf{u}^*) = Q_{jt}(\mathbf{s}, \mathbf{u}^*)$ .

However, since  $Q_{\text{tot}}(\mathbf{s}, \mathbf{u}^-) \geq Q_{\text{tot}}(\mathbf{s}, \mathbf{u}^*)$  due to  $\bar{\mathbf{u}} = \mathbf{u}^-$ , we get a contradiction where  $Q_{jt}(\mathbf{s}, \mathbf{u}^-) \geq Q_{jt}(\mathbf{s}, \mathbf{u}^*)$ . Therefore,  $\mathbf{u}^-$  is not a stable point. □

## D STABILITY ANALYSIS OF QPLEX

In this section, we provide a suboptimal case of QPLEX in matrix games. Here, we present the expression of  $Q_{\text{tot}}$  of QPLEX.

$$Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) = \sum_i \max_{u_i} Q_i(\tau_i, u_i) + \sum_i w_i(\boldsymbol{\tau}, \mathbf{u}) * (Q_i(\tau_i, u_i) - \max_{u_i} Q_i(\tau_i, u_i)) \quad (27)$$

where  $w_i(\boldsymbol{\tau}, \mathbf{u}) \geq 0$ .

For QPLEX, stable point occurs when the gradients of  $Q$  and  $w$  satisfy that  $\forall i \in \{1, 2, \dots, N\}$

$$\frac{\partial L_{\text{tot}}}{\partial Q_i} = \mathbf{0}, \quad \frac{\partial L_{\text{tot}}}{\partial w_i} \geq 0 \quad (28)$$

where  $\frac{\partial L_{\text{tot}}}{\partial w_i(\boldsymbol{\tau}, \mathbf{u})} > 0$  only when  $w_i(\boldsymbol{\tau}, \mathbf{u}) = 0$ .

The suboptimal cases of QPLEX are presented in Table 7 and Table 8, where there are three stable points but only one is optimal. For readers to verify those stable points, we expand Equation (28) into Equation (29) and substitute  $\tilde{\mathbf{u}}$  with  $\bar{\mathbf{u}}$ .

$$\frac{\partial L_{\text{tot}}}{\partial Q_i(\tau_i, u_i)} = \begin{cases} \sum_{\mathbf{u} \in \mathcal{U}(u_i = u_i)} (Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) - Q_{jt}(\boldsymbol{\tau}, \mathbf{u})) w_i(\boldsymbol{\tau}, \mathbf{u}) & u_i \neq \bar{u}_i \\ \sum_{\mathbf{u} \in \mathcal{U}(u_i = u_i)} (Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) - Q_{jt}(\boldsymbol{\tau}, \mathbf{u})) w_i(\boldsymbol{\tau}, \mathbf{u}) + \\ \sum_{\mathbf{u} \in \mathcal{U}} (Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) - Q_{jt}(\boldsymbol{\tau}, \mathbf{u})) (1 - w_i(\boldsymbol{\tau}, \mathbf{u})) & u_i = \bar{u}_i \end{cases} \quad (29)$$

$$\frac{\partial L_{\text{tot}}}{\partial w_i(\boldsymbol{\tau}, \mathbf{u})} = (Q_{\text{tot}}(\boldsymbol{\tau}, \mathbf{u}) - Q_{jt}(\boldsymbol{\tau}, \mathbf{u})) (Q_i(\tau_i, u_i) - Q_i(\tau_i, \bar{u}_i))$$

where  $\mathcal{U}(u_i = u_i) = \{u | u \in \mathcal{U}, u_i = u_i\}$ .

8	-12	-12
-12	3	0
-12	0	5

(a)  $Q_{jt}$

$Q_2$	<b>2.77</b>	0.31	0.37	
$Q_1$	<b>5.23</b>	<b>8</b>	-12	-12
	0.49	-12	3	0
	0.33	-12	0	5

$Q_2$	1.82	<b>3.09</b>	0.70	
$Q_1$	1.42	5.5	-12	-12
	<b>2.41</b>	-12	<b>5.5</b>	0
	1.09	-12	0	5

$Q_2$	1.29	0.66	<b>3.41</b>	
$Q_1$	1.38	6.5	-12	-12
	0.68	-12	3	0
	<b>3.09</b>	-12	0	<b>6.5</b>

(b)  $Q_{tot}^*$                       (c)  $Q_{tot}^*$                       (d)  $Q_{tot}^*$

<b>0.25</b>	1.51	1.71
4.21	0.48	0.10
4.08	0.59	0.09

0.00	17.64	5.01
8.53	<b>7.65</b>	1.41
5.30	4.16	0.15

0.00	6.70	10.86
3.35	0.26	2.71
5.10	1.03	<b>1.92</b>

(e)  $w_1^*$                       (f)  $w_1^*$                       (g)  $w_1^*$

<b>2.01</b>	8.13	8.34
0.85	1.12	3.12
4.62	2.08	1.05

0.00	9.07	5.22
13.71	<b>7.94</b>	2.29
8.23	3.76	0.13

0.00	2.58	2.50
4.93	1.05	0.23
8.72	2.36	<b>1.47</b>

(h)  $w_2^*$                       (i)  $w_2^*$                       (j)  $w_2^*$

Table 7: The example of QPLEX in which  $u^*$  is not the only one stable points. The quantities of  $Q_{tot}^*$  and  $w_i^*$  listed in the same column correspond to a particular stable point, where we highlight  $\bar{u}^*$  in bold.

8	0	0
0	3	0
0	0	5

(a)  $Q_{jt}$

$Q_2$	<b>4.14</b>	1.85	0.73	
$Q_1$	<b>3.86</b>	<b>8</b>	0	0
	0.44	0	3	0
	0.54	0	0	5

$Q_2$	2.57	<b>2.75</b>	0.99	
$Q_1$	2.59	5.5	0	0
	<b>2.75</b>	0	<b>5.5</b>	0
	1.11	0	0	5

$Q_2$	2.86	0.57	<b>3.15</b>	
$Q_1$	3.06	6.5	0	0
	0.82	0	3	0
	<b>3.35</b>	0	0	<b>6.5</b>

(b)  $Q_{tot}^*$                       (c)  $Q_{tot}^*$                       (d)  $Q_{tot}^*$

<b>0.94</b>	1.31	1.10
2.34	0.83	0.91
2.41	1.49	0.83

0.00	33.25	1.48
8.82	<b>13.44</b>	0.44
3.00	3.35	0.16

0.00	2.91	22.31
2.13	0.54	2.57
4.35	1.53	<b>14.12</b>

(e)  $w_1^*$                       (f)  $w_1^*$                       (g)  $w_1^*$

<b>1.44</b>	3.50	2.34
0.76	0.95	1.43
0.66	1.34	0.07

0.00	19.52	2.99
31.74	<b>20.50</b>	3.13
3.35	1.73	0.14

0.00	2.20	0.10
3.89	0.83	1.94
22.74	2.53	<b>6.15</b>

(h)  $w_2^*$                       (i)  $w_2^*$                       (j)  $w_2^*$

Table 8: The example in which QPLEX has two suboptimal stable points, while QMIX can easily obtain the optimum. The quantities of  $Q_{tot}^*$  and  $w_i^*$  listed in the same column correspond to a particular stable point, where we highlight  $\bar{u}^*$  in bold.

## E EXPERIMENTAL SETUP

We adopt a setup similar to PyMARL (Whiteson et al., 2019). We use the implementation of QMIX, QTRAN, QPLEX, WQMIX, ResQ, and NA<sup>2</sup>Q (Liu et al., 2023) from their open-source repositories<sup>5 6 7 8</sup>. These codes are released under the Apache License V2.0.

All algorithms apply TD(0) to update  $Q_{\text{tot}}$  including ResQ, as TD( $\lambda$ ) result in catastrophic performance in *5m\_vs\_6m* while no measurable improvement in other scenarios. The hyperparameters for all algorithms and environments are presented in Table 9. For WQMIX, we set  $\alpha = 0.1$  for all environments. For MRVF, we use three rounds of value factorization, with a probability  $p = 0.2$  for sampling preselected actions.

For hardware, we run experiments on an NVIDIA 3090 GPU for risk-reward games, predator-prey tasks, and SMAC scenarios (excluding *bane\_vs\_bane*). In addition, we use an NVIDIA A800 GPU for MRVF in *2c\_vs\_64zg* and *MMM2* scenarios.

Hyperparameter	Value	Description
Batch Size	32	Number of episodes per update
Replay buffer size	5000	Maximum number of stored episodes
Target update interval	200	Frequency of updating the target network
Initial $\epsilon$	1.0	The initial $\epsilon$ in the $\epsilon$ -greedy policy
Final $\epsilon$	0.05	The final $\epsilon$ in the $\epsilon$ -greedy policy
Anneal steps for $\epsilon$	50,000	Number of steps for linearly decay of $\epsilon$
Discount $\gamma$	0.99	Discount of future return
Test interval	10,000	Frequency of test evaluation
Test episodes	32	Number pf episodes to test

Table 9: The hyperparameters for experiments

**Network structure** We provide details of the network structures presented in Figure 4. The individual Q network and the mixing network ( $Q_{\text{tot}}$ ) are shown in Figure 9, and the joint action value network ( $\hat{Q}_{\text{jt}}$ ) is shown in Figure 10.

The individual Q network processes the observation using a linear layer followed by a GRU (Chung et al., 2014) with 64 hidden dimensions. The mixing network processes the state with 2-layer linear network (64 hidden dimension) and processes the individual Qs with 2-layer linear network, where the weights and bias (generated by the state’s linear processor) have 32 hidden dimension. Inspired by the embedding module (Bengio et al., 2003) (Mikolov et al., 2013), we generate the features of  $\bar{u}^k$  by selecting them from the continuous features (128 dimensions) of the observation or state. For the concatenated feature, we double the hidden dimensions to 128.

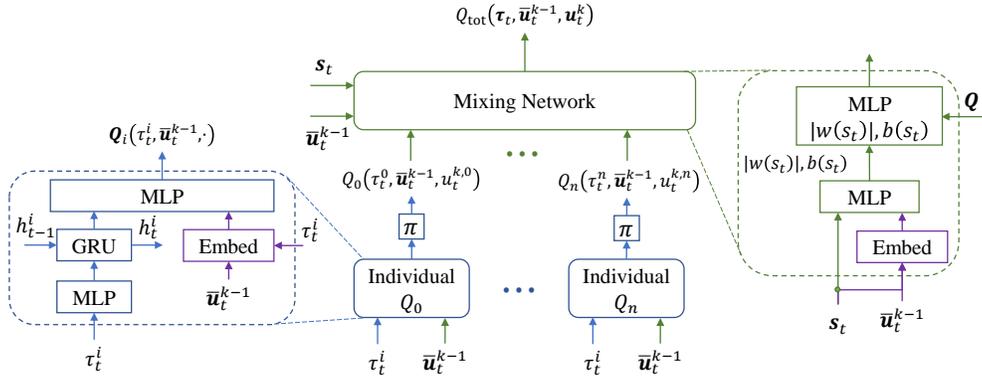
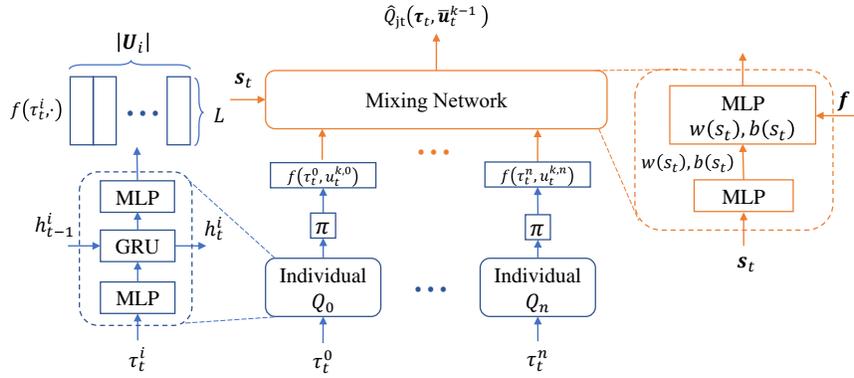
The network structure of  $\hat{Q}_{\text{jt}}$  is similar to that of  $Q_{\text{tot}}$ , with two key differences: First, we use a 3-layer linear network to process the state. Second, the output of the individual network has 64 dimensions. We input the state to  $Q_{\text{jt}}$  in all SMAC scenarios except *3s5z* and *3s\_vs\_5z*, because masking observations of dead units loses critical information for return prediction in complex scenarios. In the *bane\_vs\_bane* scenario, we replace the embedding feature with the action index to avoid GPU out-of-memory errors.

<sup>5</sup><https://github.com/oxwhirl/pymarl>

<sup>6</sup><https://github.com/oxwhirl/wqmix>

<sup>7</sup><https://github.com/xmu-rl-3dv/ResQ>

<sup>8</sup><https://github.com/zichuan-liu/NA2Q>

Figure 9: The network structure of  $Q_{\text{tot}}$ .Figure 10: The network structure of  $Q_{\text{jt}}$ .

### E.1 ONE-STEP GAME

We randomly generate  $Q_{\text{jt}}$  of the risk-reward game with the following steps: First, randomly generate the individual reward vectors  $r_i \geq 0$ . Second, we randomly generate a bijection  $\mathcal{U} \rightarrow \mathcal{U}$  that maps  $u_i$  to another action  $v_i$  for all agents  $i$ , which is the exchange of rows and columns for a matrix. Finally, we let  $Q_{\text{jt}}(\mathbf{u}) = \text{sign}(\mathbf{v}) * \sum_{i=0}^n r_i(v_i)$  where  $\text{sign}(\mathbf{v}) = 1$  if all agents choose the same mapped action, and  $\text{sign}(\mathbf{v}) = -1$  otherwise. Here, we give an example of the risk-reward matrix in Table 10.

$r_2 \backslash r_1$	3	4	5
0	<b>3</b>	-4	-5
1	-4	<b>5</b>	-6
2	-5	-6	<b>7</b>

(a)

$r_2 \backslash r_1$	5	4	3
2	<b>7</b>	-6	-5
0	-5	-4	<b>3</b>
1	-6	<b>5</b>	-4

(b)

Table 10: An example of the risk-reward matrix. (a): the risk-reward matrix before action mapping, where positive rewards are in the diagonal. (b): the risk-reward matrix after action mapping, which scatters positive rewards.

We constrain the values of  $Q_{\text{jt}}$  to the interval  $[-10, 10]$ . To achieve this, we first generate the reward vectors uniformly from  $[0, 1]$ , then divide each component by  $n$ , and finally multiply the result by

1188 10. The normalized return in Figure 5 is defined as

$$1189 \text{ return}(\mathbf{u}) = \begin{cases} \frac{Q_{jt}(\mathbf{u}) - \min_{\mathbf{u}^+ \in \mathcal{U}^+} Q_{jt}(\mathbf{u}^+)}{\max_{\mathbf{u}^+ \in \mathcal{U}^+} Q_{jt}(\mathbf{u}^+) - \min_{\mathbf{u}^+ \in \mathcal{U}^+} Q_{jt}(\mathbf{u}^+)} & Q_{jt}(\mathbf{u}) \geq 0 \\ \frac{Q_{jt}(\mathbf{u}) - \min_{\mathbf{u}^- \in \mathcal{U}^-} Q_{jt}(\mathbf{u}^-)}{\min_{\mathbf{u}^- \in \mathcal{U}^-} Q_{jt}(\mathbf{u}^-)} & Q_{jt}(\mathbf{u}) < 0 \end{cases} \quad (30)$$

1194 where  $\mathcal{U}^+ = \{\mathbf{u}^+ | \mathbf{u}^+ \in \mathcal{U}, Q_{jt}(\mathbf{u}^+) \geq 0\}$  and  $\mathcal{U}^+ \cup \mathcal{U}^- = \mathcal{U}$ .

## 1198 E.2 PREDATOR PREY

1200 The predator-prey task (Böhmer et al., 2020) involves a multi-agent scenario where 8 predators  
1201 cooperate to capture 8 prey. Each predator can choose from six possible actions: four directional  
1202 movements (up, down, left, right), staying still, or attempting to "capture" prey.

### 1204 Capture Conditions

- 1205 • A predator can only select the "capture" action when occupying the same position with  
1206 prey.
- 1207 • Successful capture requires at least two predators simultaneously choosing "capture" on  
1208 the same prey.
- 1209 • Successful capture yields +10 reward.
- 1210 • If only one predator chooses "capture" for a prey, the team receives a punishment.
- 1211 • The captured prey and successful predators are immediately removed from the environ-  
1212 ment. Observations of removed agents are masked with zeros.

## 1215 E.3 STARCRAFT II MULTI-AGENT CHALLENGE

1217 The SMAC (SC2.4.10 version) involves two opposing teams competing to defeat each other. Players  
1218 control one team's units with the objective of eliminating all enemy units. The reward system  
1219 consists of three components:

### 1221 Reward

- 1222 • **Kill Reward:** +10 for eliminating an enemy unit by reducing its HP to zero.
- 1223 • **Win Reward:** +200 for defeating all enemy units and winning the game.
- 1224 • **Damage:** The amount of change in an enemy's HP.

1227 In the original PyMARL implementation (Whiteson et al., 2019), rewards are enforced as positive by  
1228 taking absolute values. However, in scenarios where enemies can regenerate health/shields, agents  
1229 might artificially inflate rewards by allowing recovery and reinflicting damage, rather than focusing  
1230 on winning. To address this, we adopt the reward without taking its absolute values.

1231 We illustrate why monotonic value factorizations excel in SMAC from two aspects. First, the in-  
1232 dividual action has slight impact on returns. Take  $3s\_vs\_5z$  scenario for example, the enemy unit,  
1233 Zealot, has 100 health and 50 recoverable shields. Thus, the maximum return is

$$1234 \underbrace{150 * 5}_{Damage} + \underbrace{10 * 5}_{Kill} + \underbrace{200}_{Win} = 1000 \quad (31)$$

1237 In particular, damages constitutes the major component of return (75% proportion). Within a step,  
1238 an individual agent's action only affects its damage, and its action will not cause a sharp change  
1239 of the state. Consequently, a slight deviation from the optimal joint action is insufficient to cause a  
1240 significant drop in return. Therefore, monotonic value factorization advantages in SMAC, as without  
1241 the sharp drop in return near the global optimum, monotonic value factorization can attain the global  
optimum.

Second, choosing the optimal individual action also performs the best in average. We consider a "focus fire" case that frequently appears in SMAC. In this case, two agents fight with three enemies denoted as "A", "B", and "C", and focusing fire on the same target is the trick to win. Those enemies are identical units, except that Enemy A currently has less HP. We present the  $Q_{jt}$  of this case in Table 11. Although  $Q_{jt}$  in Table 11 is non-monotonic, the average return of choosing action A is superior to that of choosing other actions. In this case, monotonic value factorization can obtain the optimal value.

	A	B	C	$\emptyset$
A	2.2	1.1	1.1	-0.9
B	1.1	2	1	-1
C	1.1	1	2	-1
$\emptyset$	-0.9	-1	-1	-2

Table 11:  $Q_{jt}$  in the "focus fire" case of the SMAC environment: The action is denoted by its attack target, and the action  $\emptyset$  means no attack. The data in the table only represents the relative advantages of each joint action rather than actual values. The agents receive a higher return (+2) when both focus fire on the same target, a moderate return (+1) when attacking different targets, and lower returns for not attacking. Additionally, attacking low-health units (Enemy A) grants a +0.1 bonus to the returns.

## F ABLATION STUDY

### F.1 EXPERIMENT AND ANALYSIS

We analyze the impact of multi-round iteration and strictness in improving greedy action (Equation (5)). We design MRVF-single, which uses a single round, and MRVF-non-strict, which also uses the original  $Q_{jt}$  as the target value of  $Q_{tot}$  (Equation (6)) for round  $k > 1$ . We present the results in risk-reward games (Figure 11), predator-prey tasks (Figure 12), and SMAC (Figure 13).

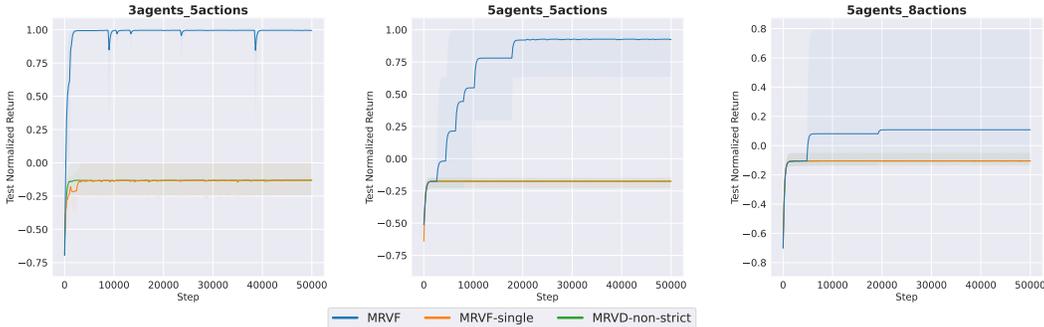
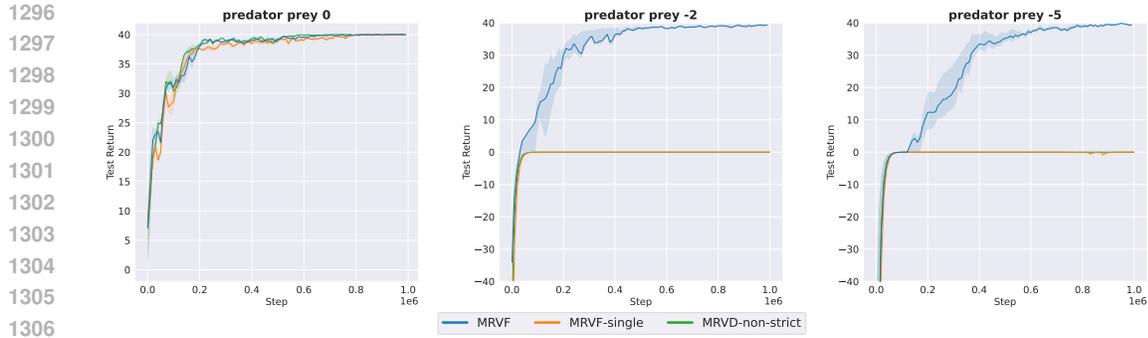


Figure 11: Test normalized return in the risk-reward games.

The results presented in Figure 11 and Figure 12 demonstrate that the strict improvement contributes to better performance in highly non-monotonic scenarios. The explanation for this is that MRVF-non-strict cannot guarantee improvement in the greedy action compared to the preceding round, which prevents it from obtaining the optimal action within a given number of rounds.

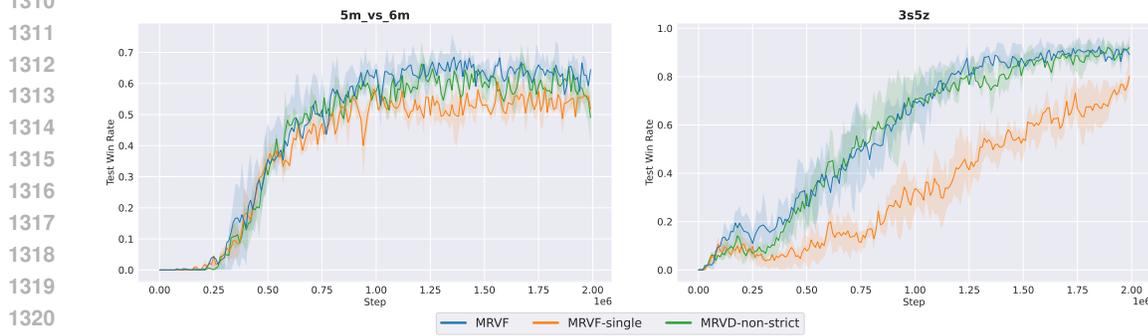
As for approximately monotonic scenarios like SMAC, Figure 13 shows a slight improvement compared to single-round factorization. In these scenarios, the improvement comes from additional information in decision-making. Specifically, by feeding  $\bar{u}_t^{k-1}$  into the individual Q networks, each agent knows the intentions of others and could infer their states accordingly. Without global information, this additional information contributes to better decision-making. In addition, due to the approximate monotonicity of the payoff, the first round of MRVF already produces a satisfactory solution, which explains the slightness of the improvement.



1307 Figure 12: Test return in the predator-prey tasks with punishments 0 (left), -2 (middle), and -5 (right).

1308

1309



1322 Figure 13: Test win rate in the SMAC benchmarks.

1323

1324

## 1325 F.2 DISCUSSION ON ROUNDS

1326 We will discuss how many rounds MRVF requires to obtain the optimal solution under scenarios with different levels of non-monotonicity. As discussed in Section F.1, a second round is required to obtain the optimal solution in highly non-monotonic environments. To support our claim, we calculate the proportion of  $u_t^k = u_t$  for each round  $k$  throughout test episodes, as shown in Figure 14.

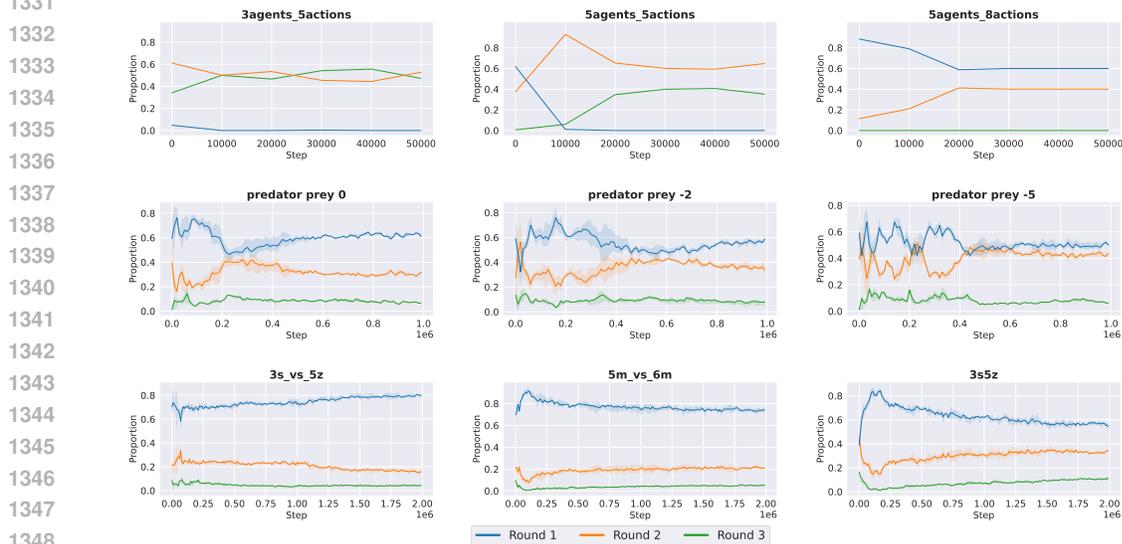


Figure 14: Proportion of final actions generated in each round throughout test episodes.

From Figure 14, we find that in the predator-prey environment, as the penalty increases, the proportion of second-round decisions rises from 30% to nearly 50%. In contrast, for approximately monotonic scenarios such as SMAC, the proportion of second-round decisions typically remains around 20%, and the performance gap between multi-round and single-round is small. In addition, a third round is generally unnecessary unless the reward is so sparse that the second round fails to find the optimal action. As shown in Figure 14, third-round decisions account for only about 5% across all predator-prey and SMAC scenarios. Nevertheless, to balance performance and efficiency, we recommend starting with a maximum of three rounds and adjusting it based on the proportion.

## G LIMITATIONS

Although MRVF achieves significant improvements over existing methods in highly non-monotonic scenarios, it faces the issue of high computational complexity: MRVF requires calculating  $Q_{\text{tot}}$  and  $Q_{\text{jt}}$  in each round. Although we reduce computational complexity by reusing parts of network outputs (e.g., individual features in  $Q_{\text{jt}}$ ) and early terminating the multi-round iterations for evaluation, the computational complexity of MRVF is still  $O(K(N + 1))$ , while that of QMIX is at  $O(N)$ . In addition, as shown in Equation (6) and Equation (5), the update of  $Q_{\text{tot}}$  depends on the results of  $Q_{\text{jt}}$ . As a result, the update of  $Q_{\text{tot}}$  is delayed compared to standard TD learning, which leads to a slow rise of the win rate in *3s\_vs\_5z* (Figure 7). However, learning directly from  $Q_{\text{jt}}$  is inevitable for environments that are not retraceable. Otherwise, if exploration of branching trajectories is allowed, TD targets could be used in the update of  $Q_{\text{tot}}$ .

We summarize the reasons why MRVF may fail to achieve the optimal solution in practice: First, the approximation errors of  $\hat{Q}_{\text{jt}}$  propagate to  $Q_{\text{tot}}$ , which may prevent the  $Q_{\text{tot}}$  from reaching the global optimum. Second, monotonic factorization is not ideal. Since we use a neural network as the mixing network, the approximation errors introduced by the network make it difficult to obtain the  $Q_{\text{tot}}$  that minimizes  $L_{\text{tot}}$ . Third, overly sparse rewards may further amplify these approximation errors. As shown in Figure 5, MRVF struggles to achieve optimal performance in *5agents\_8actions*.

## H SUPPLEMENTARY EXPERIMENT AND ANALYSIS

In this section, we provide analysis and experiments that are not included in the main part of this paper.

### H.1 STARCRAFT II MULTI-AGENT CHALLENGE

In Section 3.3, we analyze the stability of QPLEX (Wang et al., 2021), WQMIX (Rashid et al., 2020a) and ResQ (Shen et al., 2022). We reveal that suboptimal stable points are the main cause of the suboptimality in highly non-monotonic scenarios including risk-reward games and predator-prey tasks with large punishments. However, in Figure 7, these methods show inferior performance compared to monotonic value factorizations in SMAC scenarios. We provide the following explanation for this phenomenon.

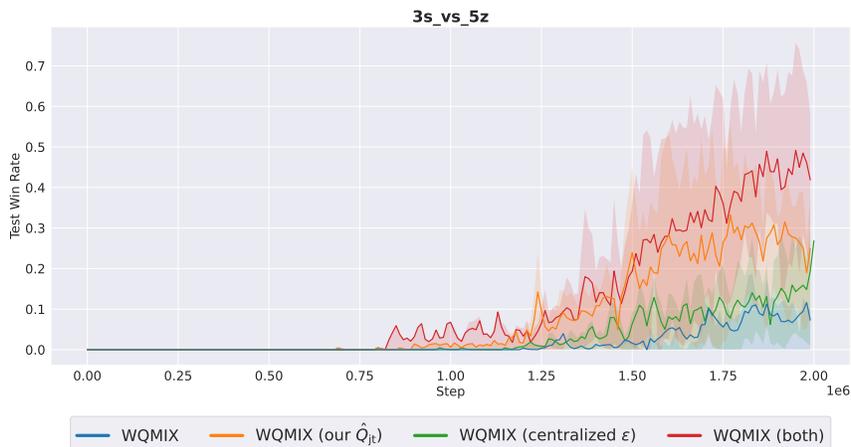
**QPLEX and ResQ** From Table 8 in Appendix D, in cases where QMIX can easily obtain the global optimum, QPLEX may still fail due to multiple suboptimal stable points. In addition, for ResQ, Theorem 3 shows that  $u^*$  is a weakly stable point in general cases. The weak stability of ResQ is particularly severe because any action can potentially become  $\bar{u}$  (Lemma 2). Consequently, QPLEX and ResQ do not perform well in SMAC.

**WQMIX** WQMIX is designed for highly non-monotonic scenarios, and from Figure 5 and Figure 6, it outperforms monotonic value factorization methods. Although WQMIX is supposed to perform similarly to QMIX in weakly non-monotonic or monotonic scenarios, it shows poor performance in SMAC (Figure 7). We provide three factors resulting in this phenomenon.

- (1) **The joint action-value network:** Individual features constitute only a small portion of the mixer’s input in the joint action-value network. Consequently, the joint action-value network hardly distinguishes different actions, since action information is primarily embedded in these individual features.

- 1404 (2) **Decentralized  $\epsilon$ -greedy policy:** The decentralized  $\epsilon$ -greedy policy (defined in Equa-  
 1405 tion (2)) combined with the weight in WQMIX enhances the convergence to the local  
 1406 optimum, since the decentralized  $\epsilon$ -greedy policy explores locally and the weight discards  
 1407 lower values around the local optimum.  
 1408 (3) **Weight:** When  $\alpha < 1$ , the weight reduces sample efficiency, as downweighting samples is  
 1409 similar to discarding them.  
 1410

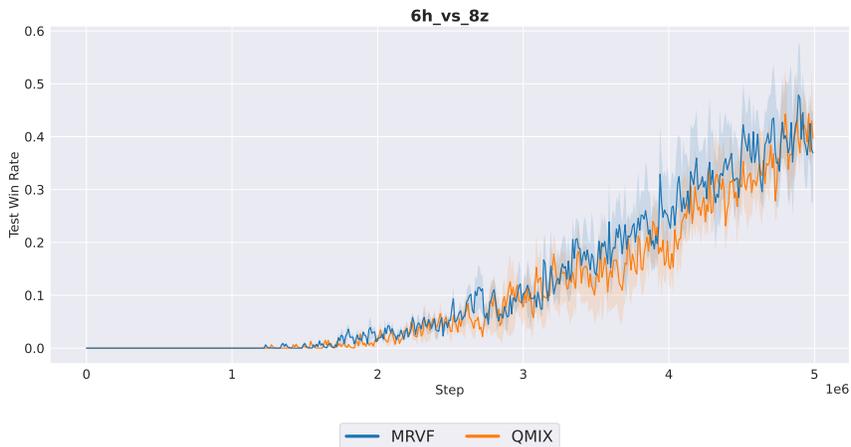
1411 We conduct additional experiments to support our analysis. We replace the network of  $\hat{Q}_{jt}$  with  
 1412 that of ours and replace the decentralized  $\epsilon$ -greedy policy with the centralized one. The result in  
 1413 Figure 15 demonstrates that these three factors do result in the performance gap between WQMIX  
 1414 and QMIX.  
 1415



1427  
 1428  
 1429  
 1430  
 1431 Figure 15: Test win rate in  $3s\_vs\_5z$ . WQMIX ( $\hat{Q}_{jt}$ ) represents WQMIX replacing the network of  
 1432  $\hat{Q}_{jt}$  with that of ours. WQMIX (centralized  $\epsilon$ ) represents WQMIX replacing the policy with the  
 1433 centralized one. WQMIX (both) represents WQMIX with both modifications.  
 1434

## 1435 H.2 SUPER-HARD SCENARIO OF SMAC

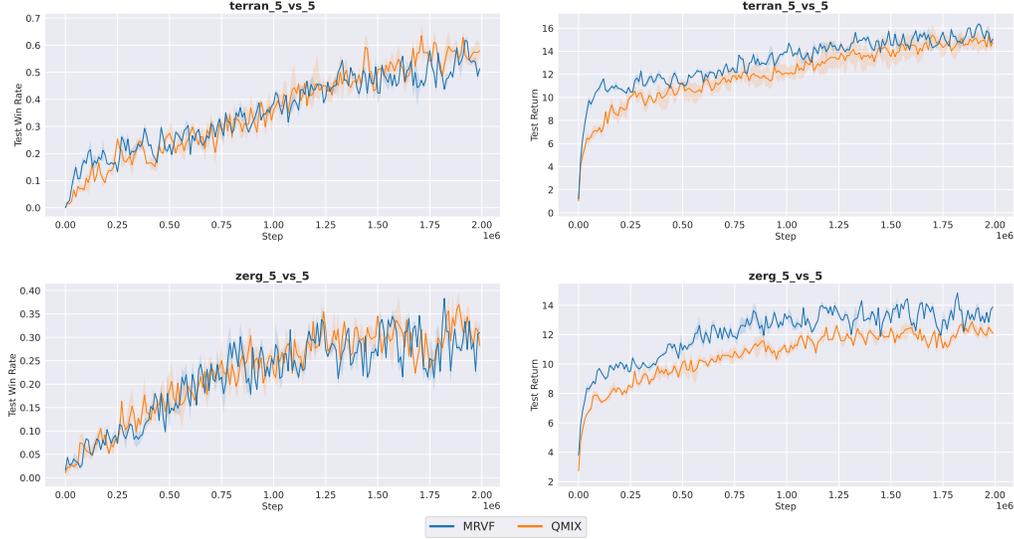
1436  
 1437 We conduct additional experiments on the  $6h\_vs\_8z$  scenario, which is a particularly challenging  
 1438 scenario where the "focus fire" is a key tactic for victory (Whiteson et al., 2019). We adopt a setup  
 1439 similar to the  $bane\_vs\_bane$  scenario, except that we set anneal steps for  $\epsilon$  to  $10^6$ . The results in Fig-  
 1440 ure 16 show a slight difference between MRVF and QMIX, supporting our analysis in Appendix E.3.  
 1441



1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457 Figure 16: Test win rate in  $6h\_vs\_8z$  in SMAC benchmarks.

### 1458 H.3 SMACv2

1459  
 1460 Results on SMACv2 (Ellis et al., 2023) are shown in Figure 17. Overall, the performance gap  
 1461 between MRVF and QMIX is small. This is because SMACv2 inherits the reward design of SMAC;  
 1462 as shown in Appendix E.3, its payoff is nearly monotonic. However, we observe a mismatch between  
 1463 the return and the win rate in SMACv2, suggesting that the reward design may need to account for  
 1464 environmental stochasticity. Due to this mismatch, we present the SMACv2 experiments only in the  
 1465 appendix.



1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485 Figure 17: Test win rate (left) and return (right) in SMACv2 benchmarks.

### 1486 H.4 PREDATOR PREY

1487  
 1488 From the results in Figure 6, QMIX fails to obtain the optimal action when the penalty is -2. How-  
 1489 ever, an analysis based on ideal QMIX will reveal that the optimal action is a stable point for any  
 1490  $\epsilon \in (0, 1]$ . In fact, this does not indicate a theoretical flaw. An explanation is that suboptimal actions  
 1491 are also stable points as  $\epsilon \rightarrow 0$ .

1492  
 1493 We will provide an analysis of the stable points when  $\epsilon \rightarrow 0$ . Taking two agents as an example, the  
 1494 reward function for the predator-prey environment with a penalty of -2 is given in Equation (32).

$$1495 \text{reward} = \begin{bmatrix} 10 & -2 & \dots & -2 \\ -2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -2 & 0 & \dots & 0 \end{bmatrix} \quad (32)$$

1496 where the reward is a matrix of shape  $6 \times 6$ .

1497 Consider  $\tilde{\mathbf{u}}$  that neither agent executes the "capture" action, which yields a reward of 0. As  $\epsilon \rightarrow 0$ ,  
 1498  $Q_{\text{tot}}^*$  that minimizes  $L_{\text{tot}}$  is given in Equation (33).

$$1499 \lim_{\epsilon \rightarrow 0} Q_{\text{tot}}^* |_{\tilde{\mathbf{u}} \neq \mathbf{u}^*} = \begin{bmatrix} -2 & -2 & \dots & -2 \\ -2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -2 & 0 & \dots & 0 \end{bmatrix}, L_{\text{tot}}^* = 0 + 0 * o(\epsilon) + 12^2 * o(\epsilon^2) \quad (33)$$

1500  
 1501 where  $\bar{\mathbf{u}}^* = \tilde{\mathbf{u}}$ , and both correspond to the agents that perform the "non-capture" action. Therefore,  
 1502 "non-capture" action pairs are stable points in predator prey with penalty -2.

Since suboptimal actions are in the majority, it is difficult for the greedy action to converge to the optimal one. However, because suboptimal actions are stable points only when  $\epsilon$  is small, we can enhance the convergence to the optimal action by slowing down the decay of  $\epsilon$ . To support our analysis, we conduct additional experiments on QMIX with  $\epsilon$  fixed at 1 and QMIX with  $\epsilon$  decaying in one million steps, as shown in Figure 18. From the result in predator prey -2, agents successfully learn to capture prey in both settings, yet they still fail when the penalty is -5. In addition, increasing the weight of the optimal action in the policy can enhance the convergence to the optimal action. For example, the weighting mechanism in WQMIX serves a similar purpose. However, these methods still fail when the penalty is set to -5.

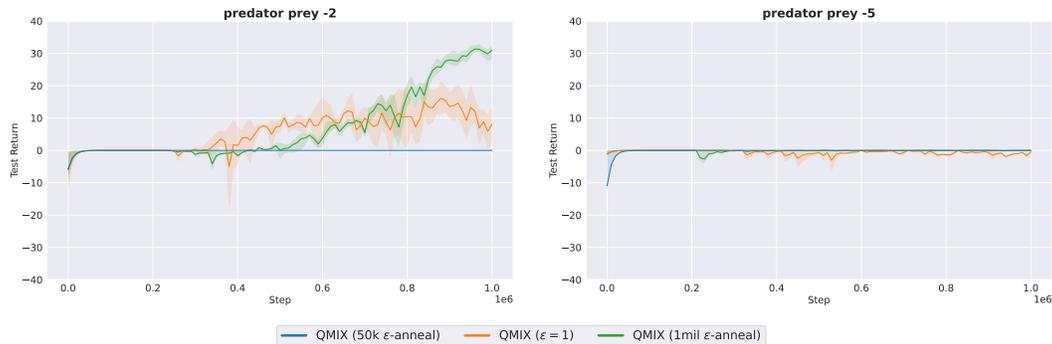


Figure 18: Test return in the predator prey tasks with punishments -2 (left), and -5 (right). QMIX (50k  $\epsilon$ -anneal) represents QMIX with  $5 * 10^4$  anneal steps for  $\epsilon$ . QMIX (1mil  $\epsilon$ -anneal) represents QMIX with  $10^6$  anneal steps for  $\epsilon$ . QMIX ( $\epsilon = 1$ ) represents QMIX with  $\epsilon = 1$  constantly.

## H.5 COMPARISON WITH OTHER MARL BASELINES

In Section 5, we compare MRVF with algorithms such as WQMIX (Rashid et al., 2020a) that claim to address the representation limitation in value factorization. However, the results show that none of these methods fully resolves this limitation. Other value-factorization approaches, including NA2Q (Liu et al., 2023) and GoMARL (Zang et al., 2023), are not designed to address the representation limitation. On-policy methods such as MAPPO (Yu et al., 2022) suffer from similar issues, converging to suboptimal Nash Equilibria (Ye et al., 2022), which are analogous to suboptimal stable points in value factorization. Consequently, all of these methods perform poorly in environments with highly non-monotonic payoffs, as illustrated in Figure 19.

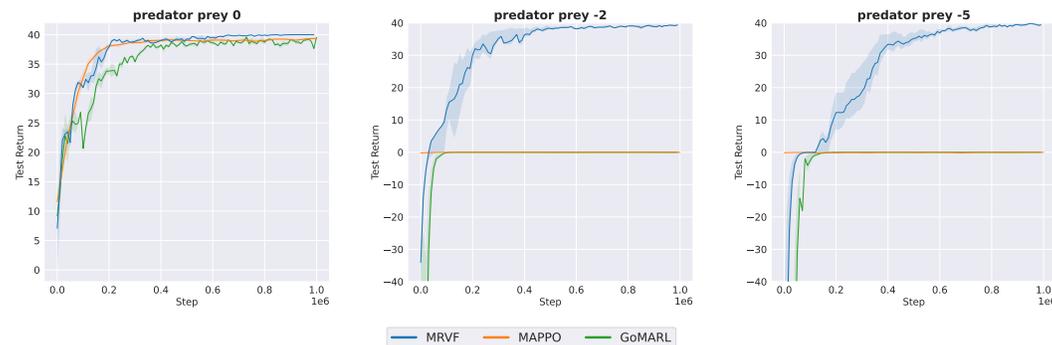


Figure 19: Test return in the predator prey tasks with punishments 0 (left), -2 (middle), and -5 (right).

The comparison results on the SMAC environment are shown in Figure 20. It can be observed that the win rates of the on-policy methods (including MAPPO (Yu et al., 2022) and CommFormer (Hu et al., 2024)) increase more slowly than that of MRVF, and their final performance is usually lower. Moreover, the performance of the on-policy methods depends on hyperparameters, as in Table 12, the network and its inputs for  $3s\_vs\_5z$  differ from those used in the other scenarios.

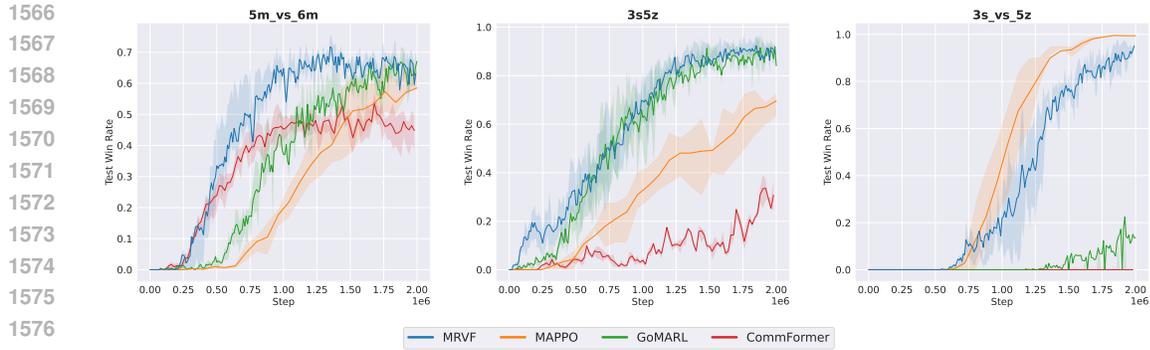


Figure 20: Test win rate in the SMAC benchmarks.

We use the implementation of GoMARL, MAPPO and CommFormer are from their open-source repositories <sup>9 10 11</sup>. The hyperparameter settings for the on-policy methods are listed in Table 12. Please refer to (Yu et al., 2022) and (Hu et al., 2024) for further details.

	epoch	clip	network	stacked frames	training threads	rollout threads
3s_vs_5z	15	0.05	mlp	4	1	8
5m_vs_6m	10	0.05	rnn	1	1	8
3s5z	5	0.2	rnn	1	1	8
predator prey	10	0.2	rnn	1	1	8

Table 12: The hyperparameters for on-policy MARL

<sup>9</sup><https://github.com/zyfsjycc/GoMARL>

<sup>10</sup><https://github.com/marlbenchmark/on-policy>

<sup>11</sup><https://github.com/charleshsc/CommFormer>