

---

# Regularized Softmax Deep Multi-Agent $Q$ -Learning

---

Ling Pan<sup>1</sup>, Tabish Rashid<sup>2</sup>, Bei Peng<sup>2</sup>, Longbo Huang<sup>1</sup>, Shimon Whiteson<sup>2</sup>

<sup>1</sup>Institute for Interdisciplinary Information Sciences, Tsinghua University

p117@mails.tsinghua.edu.cn, longbohuang@tsinghua.edu.cn

<sup>2</sup>Department of Computer Science, University of Oxford

{tabish.rashid, bei.peng, shimon.whiteson}@cs.ox.ac.uk

## Abstract

Tackling overestimation in  $Q$ -learning is an important problem that has been extensively studied in single-agent reinforcement learning, but has received comparatively little attention in the multi-agent setting. In this work, we empirically demonstrate that QMIX, a popular  $Q$ -learning algorithm for cooperative multi-agent reinforcement learning (MARL), suffers from a more severe overestimation in practice than previously acknowledged, and is not mitigated by existing approaches. We rectify this with a novel regularization-based update scheme that penalizes large joint action-values that deviate from a baseline and demonstrate its effectiveness in stabilizing learning. Furthermore, we propose to employ a softmax operator, which we efficiently approximate in a novel way in the multi-agent setting, to further reduce the potential overestimation bias. Our approach, Regularized Softmax (RES) Deep Multi-Agent  $Q$ -Learning, is general and can be applied to any  $Q$ -learning based MARL algorithm. We demonstrate that, when applied to QMIX, RES avoids severe overestimation and significantly improves performance, yielding state-of-the-art results in a variety of cooperative multi-agent tasks, including the challenging StarCraft II micromanagement benchmarks.

## 1 Introduction

In recent years, multi-agent reinforcement learning (MARL) has achieved significant progress [5, 24] under the popular centralized training with decentralized execution (CTDE) paradigm [27, 17, 8]. In CTDE, the agents must learn decentralized policies so that at execution time they can act based on only local observations, but the training itself is centralized, with access to global information. A critical challenge in this setting is how to represent and learn the joint action-value function [32].

In learning the value function, overestimation is an important challenge that stems from the max operator [38] typically used in the bootstrapping target. Specifically, the max operator in  $Q$ -learning [44] approximates the maximum *expected* value with the maximum *estimated* value. This can lead to overestimation as  $\mathbb{E} [\max_i X_i] \geq \max_i \mathbb{E} [X_i]$  due to noise [38, 13], where  $X_i$  is a random variable representing the  $Q$ -value of action  $i$  given a state. This overestimation error can accumulate during learning and hurt the performance of both value-based [41, 2, 36, 18] and actor-critic algorithms [9], and has been widely studied in the single-agent domain. However, overestimation can be even more severe in the multi-agent setting. For example, suppose there are  $n$  agents, each agent has  $K$  actions, and the  $Q$ -value for each action given a state is independently drawn from a uniform distribution  $U(0, 1)$ . Then,  $\max_i \mathbb{E} [X_i]$  is  $\frac{1}{2}$  while  $\mathbb{E} [\max_i X_i]$  is  $\frac{K^n}{K^n + 1}$ , which quickly approaches 1 (the maximum value of  $X_i$ ) as the size of the joint action space increases exponentially with the number of agents. Nonetheless, this problem has received much less attention in MARL.

QMIX [33] is a popular CTDE deep multi-agent  $Q$ -learning algorithm for cooperative MARL. It combines the agent-wise utility functions  $Q_a$  into the joint action-value function  $Q_{tot}$ , via a monotonic

mixing network to ensure consistent value factorization. Due to its superior performance, there have been many recent efforts to improve QMIX’s representation capability [35, 45, 31, 43]. However, the overestimation problem of the joint-action  $Q$ -function  $Q_{tot}$  can be exacerbated in QMIX due not only to overestimation in agents’  $Q_a$  but also the non-linear monotonic mixing network. Despite this, the role of overestimation in limiting its performance has been largely overlooked.

In this paper, we show empirically that overestimation in deep multi-agent  $Q$ -learning is more severe than previously acknowledged and can lead to divergent learning behavior in practice. In particular, consider double estimators, which can successfully reduce overestimation bias in the single-agent domain [41, 9]. Although QMIX applies Double DQN [41] to estimate the value function,<sup>1</sup> we find that it is ineffective in deep multi-agent  $Q$ -learning. As shown in Figure 1, value estimates of the joint-action  $Q$ -function can increase without bound in tasks from the multi-agent particle framework [21], yielding catastrophic performance degradation. Our experiments show that surprisingly, even applying Clipped Double  $Q$ -learning (a key technique from a state-of-the-art TD3 [9] algorithm) to the multi-agent setting does not resolve the severe overestimation bias in the joint-action  $Q$ -function. Therefore, alleviating overestimation in MARL is a particularly important and challenging problem.

To tackle this issue, we propose a novel update scheme that penalizes large joint-action  $Q$ -values. Our key idea is to introduce a regularizer in the Bellman loss. A direct penalty on the magnitude of joint-action  $Q$ -values can result in a large estimation bias and hurt performance. Instead, to better trade off learning efficiency and stability, we introduce a baseline into the penalty, thereby constraining the joint-action  $Q$ -values to not deviate too much from this baseline. Specifically, we use the discounted return as the baseline. By regularizing towards a baseline, we stabilize learning and effectively avoid the unbounded growth in our value estimates.

However, regularization is not enough to fully avoid overestimation bias in the joint-action  $Q$ -function due to the max operator in the target’s value estimate [38]. To this end, we propose to employ a softmax operator, which has been shown to efficiently improve value estimates in the single-agent setting [36, 30, 29]. Unfortunately, a direct application of the softmax operator is often too computationally expensive in the multi-agent case, due to the exponentially-sized joint action space. We therefore propose a novel method that provides an efficient and reliable approximation to the softmax operator based on a joint action subspace, where the gap between our approximation and its direct computation converges to 0 at an exponential rate with respect to its inverse temperature parameter. The computational complexity of our approximation scales only linearly in the number of agents, as opposed to exponentially for the original softmax operator. We show that our softmax operator can further improve the value estimates in our experiments. We refer to our method as *RES (Regularized Softmax) deep multi-agent  $Q$ -learning*, which utilizes the discounted return-based regularization and our approximate softmax operator.

To validate RES, we first prove that it can reduce the overestimation bias of QMIX. Next, we conduct extensive experiments in the multi-agent particle tasks [21], and show that RES simultaneously enables stable learning, avoids severe overestimation when applied to QMIX, and achieves state-of-the-art performance. RES is not tied to QMIX and can significantly improve the performance and stability of other deep multi-agent  $Q$ -learning algorithms, e.g., Weighted-QMIX [31] and QPLEX [43], demonstrating its versatility. Finally, to demonstrate its ability to scale to more complex scenarios, we evaluate it on a set of challenging StarCraft II micromanagement tasks [34]. Results show that RES-QMIX provides a consistent improvement over QMIX in all scenarios tested.

## 2 Background

**Decentralized partially observable Markov decision process (Dec-POMDP).** A fully cooperative multi-agent task can be formulated as a Dec-POMDP [26] represented by a tuple

<sup>1</sup>As mentioned in Appendix D.3 in [32] and open-source PyMARL [34] implementations.

$\langle A, S, U, P, r, Z, O, \gamma \rangle$ , where  $A \equiv \{1, \dots, n\}$  denotes the finite set of agents,  $s \in S$  is the global state, and  $\gamma \in [0, 1]$  is the discount factor. The action space for an agent is  $U$  with size  $K$ . At each timestep, each agent  $a \in A$  receives an observation  $z \in Z$  from the observation function  $O(s, a)$  due to partial observability, and chooses an action  $u_a \in U$ , which forms a joint action  $\mathbf{u} \in \mathbf{U} \equiv U^n$ . This leads to a transition to the next state  $s' \sim P(s'|s, \mathbf{u})$  and a joint reward  $r(s, \mathbf{u})$ . We assume that the reward function is bounded [39], i.e.,  $|r(s, \mathbf{u})| \leq R_{\max}$ . Each agent has an action-observation history  $\tau_a \in T \equiv (Z \times U)^*$ , based on which it constructs a policy  $\pi_a(u_a | \tau_a)$ . The goal is to find an optimal joint policy  $\pi = \langle \pi_1, \dots, \pi_n \rangle$ , whose joint action-value function is  $Q^\pi(s_t, \mathbf{u}_t) = \mathbb{E}[\sum_{i=0}^{\infty} \gamma^i r_{t+i}]$ .

**Deep multi-agent  $Q$ -learning.** Deep multi-agent  $Q$ -learning [37, 33, 35, 43] extends deep  $Q$ -learning [44], a popular value-based method for learning optimal action values, to the multi-agent setting. Given transitions  $(s, \mathbf{u}, r, s')$  sampled from the experience replay buffer  $\mathcal{B}$ , its objective is to minimize the mean squared error loss  $\mathcal{L}(\theta)$  on the temporal-difference (TD) error  $\delta = y - Q_{tot}(s, \mathbf{u})$ , where  $y = r + \gamma \max_{\mathbf{u}'} \bar{Q}_{tot}(s', \mathbf{u}')$  is the target value, and  $\bar{Q}_{tot}$  is the target network for the joint-action  $Q$ -function that is periodically copied from  $Q_{tot}$ . Parameters of  $Q_{tot}$  are denoted by  $\theta$  that are updated by  $\theta' = \theta - \alpha \nabla_\theta \mathcal{L}(\theta)$ , where  $\alpha$  is the learning rate. To mitigate the overestimation bias from the max operator in the target's value estimate [38], Double DQN [41, 13] estimates the target value as  $y = r + \gamma \bar{Q}_{tot}(s', \arg \max_{\mathbf{u}'} Q_{tot}(s', \mathbf{u}'))$  which decouples action selection and evaluation.

**Centralized training with decentralized execution (CTDE).** In CTDE [27, 17], agents are trained in a centralized way with access to the overall action-observation history and global state during training, but during execution have access only to their own local action-observation histories. The individual-global-max (IGM) property [35] is a popular concept to realize efficient CTDE as in Eq. (1), where  $Q_{tot}$  and  $Q_a$  denote the joint-action  $Q$ -function and agent-wise utilities respectively.

$$\arg \max_{\mathbf{u}} Q_{tot}(s, \mathbf{u}) = (\arg \max_{u_1} Q_1(s, u_1), \dots, \arg \max_{u_n} Q_n(s, u_n)). \quad (1)$$

The IGM property enables efficient decentralized execution, which ensures the consistency between greedy action selection in the local and joint  $Q$ -values. QMIX [33] is a popular CTDE method that combines  $Q_a$  into  $Q_{tot}$  via a non-linear monotonic function  $f_s$  that is state-dependent as in Eq. (2). To satisfy the IGM property,  $f_s$  is constrained to be monotonic in  $Q_a$ , i.e.,  $\frac{\partial f_s}{\partial Q_a} \geq 0, \forall a \in A$ , which is achieved by enforcing non-negative weights in  $f_s$ .

$$Q_{tot}(s, \mathbf{u}) = f_s(Q_1(s, u_1), \dots, Q_n(s, u_n)). \quad (2)$$

### 3 Overestimation in QMIX

In this section, we show empirically that the overestimation problem for deep multi-agent  $Q$ -learning can be more severe in practice than previously acknowledged. In particular, we demonstrate that state-of-the-art methods for tackling this issue in the single-agent domain can fail in our multi-agent setting.

We investigate the behavior of QMIX in the predator-prey task from the multi-agent particle environments [21], where 3 slower predators need to coordinate to capture a faster prey to solve the task. A detailed description of the task is included in Appendix E.1.1. We consider a fully cooperative setting where the prey is pre-trained by MADDPG [21] to avoid capture by the predators, and the predators need to learn to cooperate with each other in order to surround and capture the prey.

Figures 2 and 3(a) show the performance of QMIX and its estimated values of the joint-action  $Q$ -function during training respectively. QMIX suffers from catastrophic performance degradation and a severe overestimation bias, where its value estimates grow without bound. Figures 3(b) and 3(c) illustrate the mean agent-wise utilities  $Q_a$  and their gradients  $\partial f_s / \partial Q_a$  over agents respectively. We see that for QMIX, the gradients  $\partial f_s / \partial Q_a$  increase rapidly and continuously during learning. Figure 3(d) shows the learned weight (matrix) of the monotonic mixing network  $f_s$ . The weights similarly grow larger during training and further amplify the overestimation in each  $Q_a$  when computing  $Q_{tot}$ , leading to a severely overestimated joint-action  $Q$ -function that continues to grow without bound.

To reduce overestimation, Ackermann et al. [1] extend the state-of-the-art TD3 algorithm [9], which addresses overestimation in the single-agent case, to the multi-agent setting. We obtain QMIX (CDQ)

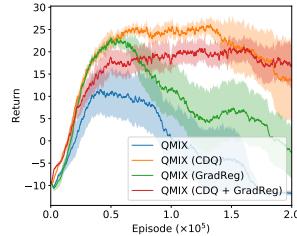


Figure 2: Return of QMIX and its variants in the predator-prey environment.

by applying the key Clipped Double  $Q$ -learning technique from TD3 [9] to per-agent utilities in QMIX for value estimation.<sup>2</sup> As shown in Figure 3(b), QMIX (CDQ) slows the increase of agent-wise utilities. However, Figure 3(a) shows that it still does not eliminate the severe overestimation bias in the joint-action  $Q$ -function, and subsequently suffers a performance degradation as shown in Figure 2. This is due to the gradients of the monotonic mixing network,  $\partial f_s / \partial Q_a$ , which still continuously increase as shown in Figure 3(c), leading to the large and increasing value estimates in Figure 3(a).

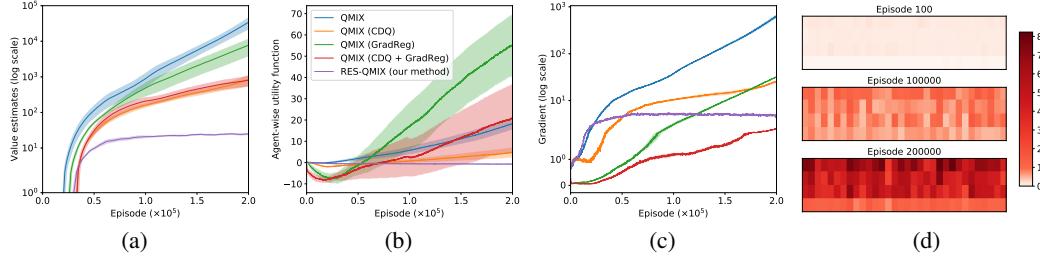


Figure 3: Learning statistics during training in predator-prey. (a) Value estimates in log scale. (b) Mean agent-wise utility functions  $Q_a$  over agents. (c) Mean gradients  $\frac{\partial f_s}{\partial Q_a}$  over agents in log scale. (d) The learned weight (matrix) in the monotonic mixing network  $f_s$  of QMIX, where a darker color represents a larger value (a detailed description is in Appendix A.1).

Another way to avoid overly large value estimates is to limit the gradients themselves. We propose Gradient Regularization (GradReg) to prevent the gradients  $\partial f_s / \partial Q_a$  from growing too large by regularizing the gradient using a quadratic penalty. Specifically, the loss function of QMIX (GradReg) is defined as  $\mathcal{L}_{\text{GradReg}}(\theta) = \mathbb{E}_{(s, u, r, s') \sim \mathcal{B}} [\delta^2 + \lambda (\partial f_s / \partial Q_a)^2]$ , where  $\delta$  is the TD error defined in Section 2 and  $\lambda$  is the regularization coefficient. However, although QMIX (GradReg) with a tuned coefficient for the regularizer prevents the gradients from growing overly large during the early stage of learning (Figure 3(c)), the agent-wise utility functions grow larger instead as shown in Figure 3(b). This then still leads to a particularly large value estimate as shown in Figure 3(a). Finally, a direct combination of both CDQ and GradReg also fails to avoid the problem.

From these experiments, we see that even applying the state-of-the-art TD3 [9] method from the single-agent literature can fail in our multi-agent setting. In addition, it is also insufficient to tackle the problem by regularizing the magnitude of the gradients  $\partial f_s / \partial Q_a$ . These results show that mitigating overestimation in MARL is particularly important and challenging, and requires novel solutions.

## 4 Regularized Softmax (RES) Deep Multi-Agent $Q$ -Learning

Our analysis shows that preventing gradients from being large as in QMIX (GradReg) and/or utilizing clipped double estimators as in QMIX (CDQ) are not sufficient to prevent overly large value estimates. It is important to note that these methods take an *indirect* approach to reduce the overestimation of the joint-action  $Q$ -value, and as shown do not solve the severe overestimation problem. Therefore, we propose a novel way to directly penalize large joint-action  $Q$ -values to avoid value estimate explosion.

Directly penalizing large joint-action  $Q$ -values can push their values towards 0, yielding a large estimation bias. Instead, we introduce a baseline into the penalty to better trade off learning efficiency and stability. Specifically, the new learning objective penalizes joint-action  $Q$ -values deviating from a baseline  $b(s, u)$ , by adding a regularizer to the loss:  $\lambda (Q_{\text{tot}}(s, u) - b(s, u))^2$ , where we use the mean squared error loss and  $\lambda$  is the regularization coefficient. A potential choice for the baseline is the  $N$ -step return, i.e.,  $b(s, u) = \sum_{t=0}^{N-1} \gamma^t r_t + \gamma^N \max_{u'} Q_{\text{tot}}(s_N, u')$ . Figure 4(a) shows the value estimates of Regularized (RE)-QMIX ( $b = N$ -step return) using the  $N$ -step return as the baseline (with best  $N, \lambda$ ). However, as it still involves  $Q_{\text{tot}}$ , the value estimation can still grow extremely large despite being discounted, yielding estimates that almost coincide with those of QMIX.

Therefore, we propose to use the discounted return starting from state  $s$  as the baseline to regularize learning, i.e.,  $b(s, u) = R_t(s, u) = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$  for the  $t$ -th timestep. The regularized learning objective directly penalizes large joint-action  $Q$ -values and allows the value estimates to remain

<sup>2</sup>Applying CDQ on  $Q_{\text{tot}}$  leads to larger value estimates than applying it on  $Q_a$  (details are in Appendix A.2).

grounded by real data. As shown in Figure 4(a), RE-QMIX effectively avoids particularly large value estimations and stabilizes learning.

However, while the regularized learning objective is effective in resolving value explosion, it can still overestimate due to the use of the max operator when computing target value estimates [38]. Figure 4(b) shows the estimated and true values of RE-QMIX (with best  $\lambda$ ). The estimated values are computed by averaging over 100 states sampled from the replay buffer at each timestep, and we estimate true values by averaging the discounted returns which are obtained by following the greedy policy with respect to the current  $Q_{tot}$  starting from the sampled states. As shown in Figure 4(b), although RE-QMIX avoids unbounded growth of its value estimates, it does not fully avoid the overestimation bias, which still estimates the target value according to the max operator [38].

To further mitigate overestimation bias in the joint-action  $Q$ -function, we adopt the softmax operator, which has shown great potential in reducing overestimation bias in place of the max operator in single-agent domains [36, 30, 29]. Specifically, the softmax operator is defined in Eq. (3), where  $\beta \geq 0$  is the inverse temperature parameter.

$$\text{sm}_{\beta, \mathbf{U}}(Q_{tot}(s, \cdot)) = \sum_{\mathbf{u} \in \mathbf{U}} \frac{e^{\beta Q_{tot}(s, \mathbf{u})}}{\sum_{\mathbf{u}' \in \mathbf{U}} e^{\beta Q_{tot}(s, \mathbf{u}')}} Q_{tot}(s, \mathbf{u}). \quad (3)$$

When  $\beta$  approaches  $\infty$  or 0, softmax reduces to the max or mean operator respectively. Unfortunately, computing Eq. (3) in the multi-agent case can be computationally intractable as the size of the joint action space grows exponentially with the number of agents. In addition, as the action space in the multi-agent case is much larger than that in the single-agent case, some joint-action  $Q$ -value estimates  $Q_{tot}(s, \mathbf{u})$  can be unreliable due to a lack of sufficient training. Thus, directly taking them all into consideration for computing the softmax operator in Eq. (3) can result in inaccurate value estimates.

We propose a novel method to approximate the computation of the softmax operator efficiently and reliably. Specifically, we first obtain the maximal joint action  $\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} Q_{tot}(s, \mathbf{u})$  with respect to  $Q_{tot}$  according to the individual-global-max (IGM) property discussed in Section 2. Then, for each agent  $a$ , we consider  $K$  joint actions, by changing only agent  $a$ 's action while keeping the other agents' actions  $\hat{\mathbf{u}}_{-a}$  fixed. We denote the resulting action set of agent  $a$  by  $U_a = \{(u_a, \hat{\mathbf{u}}_{-a}) | u_a \in U\}$ . Finally, we form a joint action subspace  $\hat{\mathbf{U}} = U_1 \cup \dots \cup U_n$ , where each  $U_a$  contributes  $K$  actions, and use  $\hat{\mathbf{U}}$  in Eq. (3) for computing the approximate softmax operator. In Theorem 1, we show that the gap between our approximation in  $\hat{\mathbf{U}}$  and its direct computation in the joint action space  $\mathbf{U}$  converges to 0 at an exponential rate with respect to  $\beta$ . The proof can be found in Appendix B.1.

**Theorem 1.** Let  $\mathbf{u}^*$  and  $\mathbf{u}'$  denote the optimal joint actions in  $\mathbf{U}$  and  $\mathbf{U} - \hat{\mathbf{U}}$  with respect to  $Q_{tot}$ , respectively. The difference between our approximate softmax operator and its direct computation in the whole action space satisfies:  $\forall s \in \mathcal{S}, |\text{sm}_{\beta, \hat{\mathbf{U}}}(Q_{tot}(s, \cdot)) - \text{sm}_{\beta, \mathbf{U}}(Q_{tot}(s, \cdot))| \leq \frac{2R_{\max}}{1-\gamma} \frac{|\mathbf{U} - \hat{\mathbf{U}}|}{|\mathbf{U} - \hat{\mathbf{U}}| + \exp(\beta(Q_{tot}(s, \mathbf{u}^*) - Q_{tot}(s, \mathbf{u}')))}$ , where  $|\mathbf{U} - \hat{\mathbf{U}}|$  denotes the size of the set  $\mathbf{U} - \hat{\mathbf{U}}$ .

**Discussion.** Since its computational complexity is linear instead of exponential in the number of agents, our approximation is feasible to compute even when the number of agents grows large, as opposed to the original softmax operator. According to the IGM property,  $\hat{\mathbf{U}}$  consists of joint actions that are close to the maximal joint action  $\hat{\mathbf{u}}$ . Thus, it is more likely to contain joint actions with more accurate and reliable value estimates, especially when using  $\epsilon$ -greedy exploration, since they are more likely to be sampled and therefore trained. Due to this reduced reliance on *unreliable* joint-action  $Q$ -values, our approximate softmax operator actually outperforms its counterpart using a direct computation in the joint action space  $\mathbf{U}$ , as well as a random sampling scheme, as shown in Section 5.1.3. The full algorithm for our approximate softmax is in Appendix B.2. As shown in

Figure 4(b), RES (Regularized Softmax)-QMIX fully addresses the overestimation bias and achieves better value estimates than RE-QMIX.<sup>3</sup>

The loss of Regularized Softmax (RES) Deep Multi-Agent  $Q$ -Learning  $\mathcal{L}_{\text{RES}}(\theta)$  is defined as:

$$\mathbb{E}_{(s, \mathbf{u}, r, s') \sim \mathcal{B}} \left[ \delta_{\text{sm}_\beta}^2 + \lambda (Q_{\text{tot}}(s, \mathbf{u}) - R_t(s, \mathbf{u}))^2 \right], \quad (4)$$

where  $\delta_{\text{sm}_\beta} = r + \gamma \text{sm}_{\beta, \hat{\mathbf{U}}}(\bar{Q}_{\text{tot}}(s', \cdot)) - Q_{\text{tot}}(s, \mathbf{u})$ , and  $\lambda$  denotes the coefficient of the regularizer.

To better understand the effectiveness of our method, we provide a theoretical justification that connects the loss function of RES and a new Bellman operator in Theorem 2.

**Theorem 2.** *Given the same sample distribution, the update of the RES method is equivalent to the update using  $\mathcal{L}(\theta) = \mathbb{E}_{(s, \mathbf{u}, r, s') \sim \mathcal{B}} [(y - Q_{\text{tot}}(s, \mathbf{u}))^2]$  with learning rate  $(\lambda + 1)\alpha$ , which estimates the target value according to  $y = \frac{r + \gamma \text{sm}_{\beta, \hat{\mathbf{U}}}(\bar{Q}_{\text{tot}}(s', \cdot))}{\lambda + 1} + \frac{\lambda R_t(s, \mathbf{u})}{\lambda + 1}$ .*

Its proof is in Appendix C. In the special case where  $\lambda = 0$ , it is equivalent to changing the max operator to our approximate softmax in the target value estimation. When  $\lambda > 0$ , it can be thought of as learning with a weighted combination of the 1-step TD target with our approximate softmax and the discounted return, which allows the bootstrapping target to remain grounded by real data.<sup>4</sup>

Let  $\mathcal{T}$  be the value estimation operator which estimates the value of the next state  $s'$ . From Theorem 2, we have  $\mathcal{T}_{\text{QMIX}} = \max_{\mathbf{u}'} \bar{Q}_{\text{tot}}(s', \mathbf{u}')$ ,  $\mathcal{T}_{\text{RE-QMIX}} = \frac{\max_{\mathbf{u}'} \bar{Q}_{\text{tot}}(s', \mathbf{u}') + \lambda R_{t+1}(s')}{\lambda + 1}$ , and  $\mathcal{T}_{\text{RES-QMIX}} = \frac{\text{sm}_{\beta, \hat{\mathbf{U}}}(\bar{Q}_{\text{tot}}(s', \cdot)) + \lambda R_{t+1}(s')}{\lambda + 1}$ . We now analyze the relationship between the value estimation bias of these operators in Theorem 3 (the proof is in Appendix D), and show that RES-QMIX can reduce the overestimation bias when built upon QMIX. It is also worth noting that since RES is general, it can be readily applied to other deep multi-agent  $Q$ -learning algorithms, as investigated in Section 5.1.4.

**Theorem 3.** *Let  $B(\mathcal{T}) = \mathbb{E}[\mathcal{T}(s')] - \max_{\mathbf{u}'} Q_{\text{tot}}^*(s', \mathbf{u}')$  be the bias of value estimates of  $\mathcal{T}$  and the true optimal joint-action  $Q$ -function  $Q_{\text{tot}}^*$ . Given the same assumptions as in [41] for the joint-action  $Q$ -function, where there exists some  $V_{\text{tot}}^*(s') = Q_{\text{tot}}^*(s', \mathbf{u}')$  for different joint actions,  $\sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s')) = 0$ , and  $\frac{1}{|\mathcal{U}|} \sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s'))^2 = C$  ( $C > 0$ ) with  $\bar{Q}_{\text{tot}}$  an arbitrary joint-action  $Q$ -function, then  $B(\mathcal{T}_{\text{RES-QMIX}}) \leq B(\mathcal{T}_{\text{RE-QMIX}}) \leq B(\mathcal{T}_{\text{QMIX}})$ .*

From the proof of Theorem 3 and Figure 4(a), we find that RE-QMIX can reduce the overestimation bias of QMIX following the same standard assumptions in [41, 36]. In addition, Theorem 3 also shows that the bias of value estimates of RES-QMIX is no larger than that of RE-QMIX, which validates our RES method.

## 5 Experiments

We conduct a series of experiments in the multi-agent particle tasks [21] to answer: (i) How much can our RES method improve over QMIX? (ii) How does RES-QMIX compare against state-of-the-art methods in performance and value estimates? (iii) How sensitive is RES to important hyperparameters and what is the effect of each component? (iv) Can RES be applied to other algorithms? We also evaluate our method on the challenging SMAC benchmark [34] to demonstrate its scalability.

### 5.1 Multi-Agent Particle Environments

The predator-prey (PP) task was introduced in Section 3. Physical deception (PD) involves 2 cooperating agents and 1 adversary, whose goal is to reach a single target landmark from a total of two landmarks, while the adversary is unaware of the target. World (W) involves 4 slower agents who must coordinate to catch 2 faster adversaries that desire to eat food. In covert communication (CC), one agent must send a private message to the other over a public channel with a private key, while an adversary tries to reconstruct the message without the key. In the above tasks, we consider a fully cooperative setting where the adversaries are pre-trained by MADDPG [21] for  $10^4$  episodes.

<sup>3</sup>Note that applying the softmax operator on agent-wise utilities results in a larger underestimation bias and significantly underperforms RES, a detailed discussion is included in Appendix B.3.

<sup>4</sup>The result of QMIX with learning rate  $(\lambda + 1)\alpha$  is in Appendix C.1 (also suffers from performance drop).

We compare RES against state-of-the-art value factorization algorithms including VDN [37], QMIX [33], QTRAN [35], Weighted QMIX (including CW-QMIX and OW-QMIX) [31], and QPLEX [43] using the PyMARL [34] implementations and setup, and the popular actor-critic algorithm MADDPG [21]. For RES-QMIX, we fix the inverse temperature  $\beta$  to be 0.05 while the regularization coefficient  $\lambda$  is selected based on a grid search over  $\{1e-2, 5e-2, 1e-1, 5e-1\}$  as investigated in Section 5.1.3. Each algorithm is run with five random seeds, and is reported in mean  $\pm$  standard deviation. A detailed description of the tasks and implementation details is in Appendix E.1.

### 5.1.1 Performance Comparison

We first investigate how much of a performance improvement RES-QMIX achieves over QMIX. Performance comparison in different environments is shown in Figure 5, and the mean normalized return averaged over different environments can be found in Appendix E.2. The results show that RES-QMIX significantly outperforms QMIX in performance and efficiency, and achieves stable learning behavior, without any catastrophic performance degradation. We then investigate the performance of RES-QMIX in comparison with state-of-the-art baselines. As shown in Figure 5, RES-QMIX significantly outperforms all baselines and achieves state-of-the-art performance in PP, PD, and W tasks. In the CC environment, RES-QMIX matches the best performers including VDN, QTRAN and QPLEX, while QMIX and Weighted QMIX suffer from oscillation. We also demonstrate the robust performance of RES-QMIX in stochastic environments in Appendix E.3.

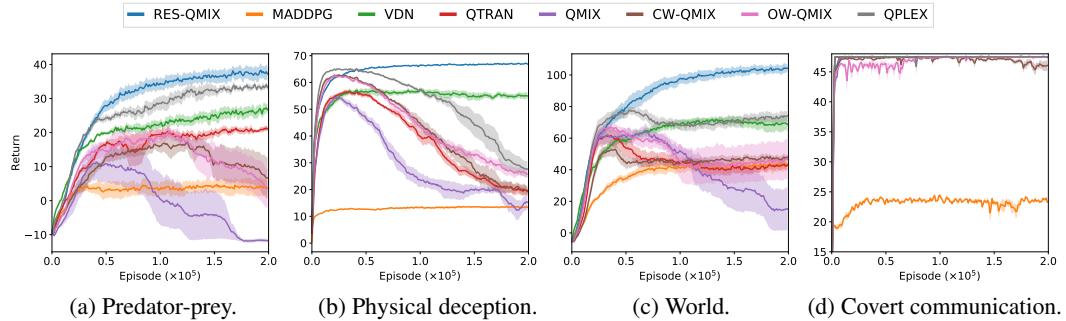


Figure 5: Performance comparison.

### 5.1.2 Value Estimation Analysis

To understand the reasons for RES-QMIX’s better performance, we analyze the bias of value estimation (the difference between estimated values and corresponding true values, which are computed in the same way as in Section 4) in predator-prey. To facilitate a fair comparison among different categories of algorithms including value-based and actor-critic methods, we normalize the bias of each method by  $100 * \frac{\text{estimated value} - \text{true value}}{|\text{true value}|}$ .

Figure 6 shows that both QMIX and Weighted QMIX suffer from large and rapidly increasing overestimation bias, yielding the severe performance degradation shown in Figure 5(a). Value estimates of VDN, which is based on a linear decomposition of  $Q_{tot}$ , increase more slowly at the end of training. However, it still incurs large overestimation as in QTRAN and QPLEX. Unlike all other value factorization methods, MADDPG learns an unfactored critic that directly conditions on the full state and joint action. It is less sample efficient, which indicates that value factorization is important in these tasks [16, 7, 15]. Thus, MADDPG results in a lower return and value estimates (shown in Appendix E.4) compared to all other value factorization methods, but still overestimates. RES-QMIX achieves the smallest bias and fully mitigates the overestimation bias of QMIX,<sup>5</sup> resulting in stable performance and outperforming all other methods.

<sup>5</sup>Overestimation has been shown to be problematic and can accumulate during learning, and a research direction starting from Double DQN [13, 41, 9] aims to reduce it. As shown in [40], there is no unbiased estimator in general settings.

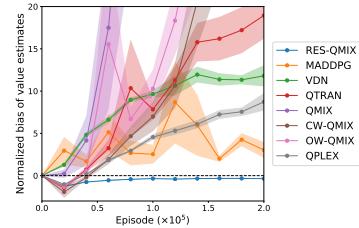


Figure 6: Normalized bias of value estimates.

### 5.1.3 Ablation Study

We now analyze how sensitive RES-QMIX is to some important hyperparameters, and the effect of each component in our method including the regularizer, the softmax operator, and the approximation scheme in the predator-prey task. Full results for all environments can be found in Appendix E.5.

**The effect of the regularization coefficient  $\lambda$ .** Figure 7(a) shows the performance of RES-QMIX with varying  $\lambda$ . RES-QMIX is sensitive to this hyperparameter, which serves a critical role in trading off stability and efficiency. A small value of  $\lambda$  fails to avoid the performance degradation, while a large value of  $\lambda$  focuses more on learning based on the regularization term instead of the bootstrapping target, and affects learning efficiency. There exists an intermediate value that provides the best trade-off. In fact, we find that this is the only parameter that needs to be tuned for RES-QMIX.

**The effect of the inverse temperature  $\beta$ .** As shown in Figure 7(b), RES-QMIX is insensitive to the hyperparameter  $\beta$ , where the performance remains competitive for a wide range of  $\beta$ . Thus, we fix  $\beta$  in RES-QMIX to be 0.05, which performs the best, in all multi-agent particle environments.

**The effect of each component.** By comparing RES-QMIX against RE-QMIX and S-QMIX in Figure 7(c), we see that the regularization component is critical for stability, while combining with our softmax operator further improves learning efficiency.

**The effect of the approximation scheme.** We compare our proposed approximation scheme for the softmax operator with a random sampling scheme and a direct computation in the joint action space. Specifically, RES-QMIX (RS) randomly samples the same number of joint actions as in RES-QMIX, while RES-QMIX (DC) directly computes softmax in the joint action space. The runtime for QMIX, RES-QMIX, and RES-QMIX (DC) is 4.8, 5.2, and 10.7 hours respectively, which shows that RES-QMIX only requires a small amount of extra computation compared to QMIX, and is much more computationally efficient than RES-QMIX (DC). Additionally, from Figure 7(c), it outperforms RES-QMIX (DC), showing that our approximation scheme is not only a necessary compromise for computational efficiency, but also improves performance as discussed in Section 4. RES-QMIX also significantly outperforms its counterpart with random sampling, which also validates its effectiveness.

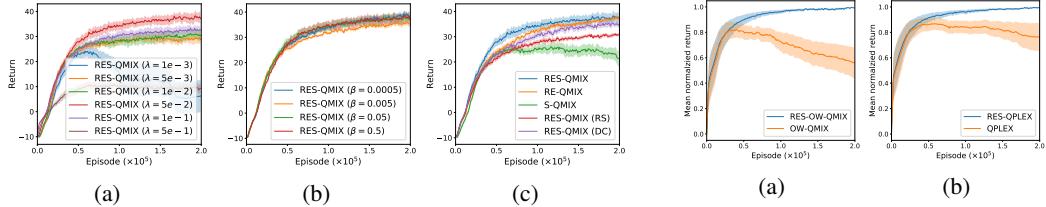


Figure 7: Ablation study. (a) Varying  $\lambda$ . (b) Varying  $\beta$ . (c) The effect of each component (the regularizer, the softmax operator, and the approximation scheme).

### 5.1.4 Applicability to Other Algorithms

The proposed RES method is general and can be readily applied to other  $Q$ -learning based MARL algorithms. To demonstrate its versatility, we apply it to two recent algorithms, Weighted QMIX [31] and QPLEX [43]. For Weighted QMIX, we use OW-QMIX as the base algorithm because it outperforms CW-QMIX as shown in Figure 5. The improvement in mean normalized return of RES-based methods over their vanilla counterparts in different environments is summarized in Figure 8, while the full comparison of learning curves in each environment is in Appendix E.6. As demonstrated, RES-Weighted QMIX and RES-QPLEX provide consistent improvement in performance and stability over Weighted QMIX and QPLEX, respectively, demonstrating the versatility of our RES method.

## 5.2 StarCraft II Micromanagement Benchmark

To demonstrate the ability of our method to scale to more complex scenarios, we also evaluate it on the challenging StarCraft II micromanagement tasks [34] by comparing RES-QMIX to QMIX using SC2 version 4.10 for five random seeds. Both algorithms are implemented using the PyMARL framework [34] with default hyperparameters as detailed in Appendix E.1. We evaluate the method

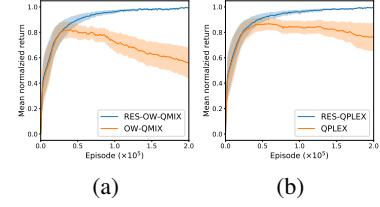


Figure 8: Mean normalized return of RES when applied to (a) Weighted QMIX and (b) QPLEX.

on maps with different difficulties ranging from easy (`2s3z`, `3s5z`), hard (`2c_vs_64zg`), to super hard (`MMM2`) as classified in [34], with a detailed description in Appendix E.1.2.

Figures 9(a)-(d) show the test win rate. RES-QMIX provides a consistent performance improvement over QMIX. Additional comparison results with the most competitive baseline QPLEX, as in the multi-agent particle tasks, are provided in Appendix E.8, where RES-QMIX also outperforms QPLEX. Comparison results of value estimates are in Appendix E.9 due to space limitation, where RES-QMIX achieves smaller value estimates than QMIX and validates its effectiveness. Additionally, we show that RES is still effective even without double estimators. The results are presented in Appendix E.10, showing that RES-QMIX (single) outperforms both QMIX and QMIX (single) on all maps we tested.

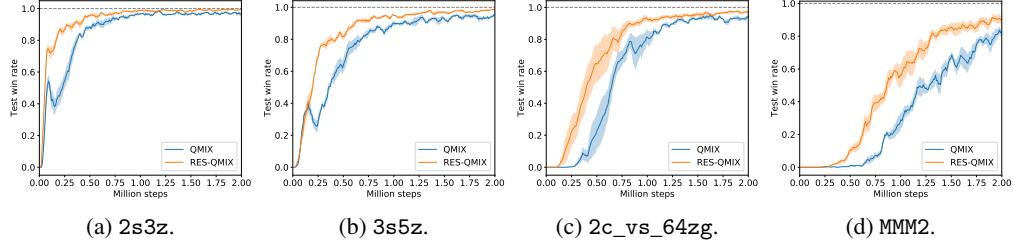


Figure 9: Comparison of RES-QMIX and QMIX in StarCraft II micromanagement tasks.

## 6 Related Work

Centralized training with decentralized execution (CTDE) [27, 17] is a popular paradigm in MARL, where a large number of recent work focuses on training agents under this paradigm. MADDPG [21] is an off-policy actor-critic algorithm that extends DDPG [19] to the multi-agent setting. Another line of research focuses on value-based  $Q$ -learning algorithm, which learns a centralized but factored joint action-value function  $Q_{tot}$ . For instance, VDN factorizes the joint-action  $Q$ -function into a linear combination of agent-wise utilities, and QMIX [33] improves the representation capability using a non-linear monotonic mixing network. There has also been a number of approaches improving QMIX’s representation capability including QTRAN [35], Weighted QMIX [31], QPLEX [43] and its exploration ability [23, 11].

Methods to alleviate overestimation have been extensively studied in reinforcement learning, with a focus on the single-agent case [13, 41, 36, 18, 6, 9, 29]. Double  $Q$ -learning [13, 41] uses a pair of independent estimators to reduce overestimation in  $Q$ -learning. In [36, 30, 29], it is shown that the softmax operator is effective for reducing estimation bias in the single-agent setting. However, it is challenging to compute the softmax operator in the multi-agent case due to the exponentially-sized joint action space, for which we propose an efficient approximation scheme. Ackermann et al. [1] extend TD3 [9] to the multi-agent setting while Gan et al. [10] propose a soft Mellowmax operator [10] to tackle overestimation in MARL. However, the analysis in [10] is based on the assumption that the gradients  $\partial f_s / \partial Q_a$  are bounded, which might not hold in practice as we have shown in Figure 3. Both methods fail to tackle the severe overestimation problem as shown in Figure 2 and Appendix E.7.

To improve learning efficiency, He et al. [14] propose a learning objective based on lower and upper bounds of the optimal  $Q$ -function. As the upper bound is based on  $N$ -step returns, it cannot avoid the severe overestimation as investigated in Section 4. Self-imitation learning (SIL) [25] aims to improve exploration via leveraging the agent’s past good experiences, which only performs SIL updates when the estimated value function is smaller than the discounted return. However, this is not effective for reducing overestimation bias in our case, as the SIL part pays little attention to experiences whose estimated values are greater than the discounted return.<sup>6</sup> The update rule of RES shown in Theorem 2 is related to the mixed Monte Carlo update [3, 28]. However, our primary motivation is to reduce severe overestimation as opposed to speed up learning.

<sup>6</sup>In Appendix A.3, we also compare RE-QMIX with a variant that only activates the penalty when the estimated  $Q$ -value is larger than the discounted return; its value estimates are close to those of RE-QMIX.

## 7 Conclusion

Overestimation is a critical problem in RL, and has been extensively studied in the single-agent setting. In this paper, we showed that it presents a more severe practical challenge in MARL than previously acknowledged, and solutions in the single-agent domain fail to successfully tackle this problem. We proposed the RES method based on a novel regularization-based update and the softmax operator with an efficient and reliable approximation in a novel way. Extensive experiments showed that RES-QMIX significantly reduces overestimation and outperforms state-of-the-art baselines. RES is general and can be applied to other deep multi-agent  $Q$ -learning methods, and can also scale to a set of challenging StarCraft II micromanagement tasks. The analysis and development of RES shed light on how to design better value estimation in MARL. Interesting directions for future work include theoretical study for the phenomenon discovered in Section 3, an adaptive scheduling for the regularization coefficient, and the application of RES to other MARL methods.

## References

- [1] J. Ackermann, V. Gabler, T. Osa, and M. Sugiyama. Reducing overestimation bias in multi-agent domains using double centralized critics. *arXiv preprint arXiv:1910.01465*, 2019.
- [2] O. Anschel, N. Baram, and N. Shimkin. Averaged-dqn: Variance reduction and stabilization for deep reinforcement learning. In *International Conference on Machine Learning*, pages 176–185. PMLR, 2017.
- [3] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, volume 29, pages 1471–1479, 2016.
- [4] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- [5] Y. Cao, W. Yu, W. Ren, and G. Chen. An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial informatics*, 9(1):427–438, 2012.
- [6] J. D. Co-Reyes, Y. Miao, D. Peng, Q. V. Le, S. Levine, H. Lee, and A. Faust. Evolving reinforcement learning algorithms. In *International Conference on Learning Representations*, 2021.
- [7] C. S. de Witt, B. Peng, P.-A. Kamienny, P. Torr, W. Böhmer, and S. Whiteson. Deep multi-agent reinforcement learning for decentralized continuous cooperative control. *arXiv preprint arXiv:2003.06709*, 2020.
- [8] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [9] S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596, 2018.
- [10] Y. Gan, Z. Zhang, and X. Tan. Stabilizing q learning via soft mellowmax operator. *arXiv preprint arXiv:2012.09456*, 2020.
- [11] T. Gupta, A. Mahajan, B. Peng, W. Böhmer, and S. Whiteson. Uneven: Universal value exploration for multi-agent reinforcement learning. *arXiv preprint arXiv:2010.02974*, 2020.
- [12] D. Ha, A. Dai, and Q. V. Le. Hypernetworks. In *International Conference on Learning Representations*, 2017.
- [13] H. Hasselt. Double q-learning. In *Advances in Neural Information Processing Systems*, volume 23, pages 2613–2621, 2010.
- [14] F. S. He, Y. Liu, A. G. Schwing, and J. Peng. Learning to play in a day: Faster deep reinforcement learning by optimality tightening. In *International Conference on Learning Representations*, 2016.
- [15] S. Iqbal, C. A. S. de Witt, B. Peng, W. Böhmer, S. Whiteson, and F. Sha. Ai-qmix: Attention and imagination for dynamic multi-agent reinforcement learning. *arXiv preprint arXiv:2006.04222*, 2020.

- [16] S. Iqbal and F. Sha. Actor-attention-critic for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 2961–2970. PMLR, 2019.
- [17] L. Kraemer and B. Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190:82–94, 2016.
- [18] Q. Lan, Y. Pan, A. Fyshe, and M. White. Maxmin q-learning: Controlling the estimation bias of q-learning. In *International Conference on Learning Representations*, 2020.
- [19] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [20] Z. Liu, X. Li, B. Kang, and T. Darrell. Regularization matters in policy optimization—an empirical study on continuous control. *arXiv preprint arXiv:1910.09191*, 2019.
- [21] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in neural information processing systems*, pages 6379–6390, 2017.
- [22] M. C. Machado, M. G. Bellemare, E. Talvitie, J. Veness, M. Hausknecht, and M. Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.
- [23] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [24] L. Matignon, L. Jeanpierre, and A.-I. Mouaddib. Coordinated multi-robot exploration under communication constraints using decentralized markov decision processes. In *AAAI 2012*, pages p2017–2023, 2012.
- [25] J. Oh, Y. Guo, S. Singh, and H. Lee. Self-imitation learning. In *International Conference on Machine Learning*, pages 3878–3887. PMLR, 2018.
- [26] F. A. Oliehoek, C. Amato, et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- [27] F. A. Oliehoek, M. T. Spaan, and N. Vlassis. Optimal and approximate q-value functions for decentralized pomdps. *Journal of Artificial Intelligence Research*, 32:289–353, 2008.
- [28] G. Ostrovski, M. G. Bellemare, A. Oord, and R. Munos. Count-based exploration with neural density models. In *International conference on machine learning*, pages 2721–2730. PMLR, 2017.
- [29] L. Pan, Q. Cai, and L. Huang. Softmax deep double deterministic policy gradients. In *Advances in Neural Information Processing Systems*, volume 33, pages 11767–11777, 2020.
- [30] L. Pan, Q. Cai, Q. Meng, W. Chen, and L. Huang. Reinforcement learning with dynamic boltzmann softmax updates. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, pages 1992–1998, 7 2020.
- [31] T. Rashid, G. Farquhar, B. Peng, and S. Whiteson. Weighted qmix: Expanding monotonic value function factorisation for deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 10199–10210, 2020.
- [32] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research*, 21(178):1–51, 2020.
- [33] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 4295–4304, 2018.
- [34] M. Samvelyan, T. Rashid, C. S. de Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr, J. Foerster, and S. Whiteson. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- [35] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 5887–5896, 2019.
- [36] Z. Song, R. Parr, and L. Carin. Revisiting the softmax bellman operator: New benefits and new perspective. In *International Conference on Machine Learning*, pages 5916–5925. PMLR, 2019.

- [37] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *AAMAS*, pages 2085–2087, 2018.
- [38] S. Thrun and A. Schwartz. Issues in using function approximation for reinforcement learning. In *Proceedings of the Fourth Connectionist Models Summer School*, pages 255–263. Hillsdale, NJ, 1993.
- [39] J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE transactions on automatic control*, 42(5):674–690, 1997.
- [40] H. Van Hasselt. Estimating the maximum expected value: an analysis of (nested) cross validation and the maximum sample average. *arXiv preprint arXiv:1302.7175*, 2013.
- [41] H. Van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [42] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- [43] J. Wang, Z. Ren, T. Liu, Y. Yu, and C. Zhang. Qplex: Duplex dueling multi-agent q-learning. In *International Conference on Learning Representations*, 2021.
- [44] C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- [45] Y. Yang, J. Hao, B. Liao, K. Shao, G. Chen, W. Liu, and H. Tang. Qatten: A general framework for cooperative multiagent reinforcement learning. *arXiv preprint arXiv:2002.03939*, 2020.

## A Details of Results in Section 3 and Additional Results

### A.1 Description of Figure 3(d) in the Main Text

Figure 3(d) in the main text shows the learned weights in the monotonic mixing network  $f_s$  of QMIX during learning, where the details of the network structure can be found in Section E.1. The weights of the monotonic mixing network are obtained from hypernetworks [12], whose outputs are reshaped into a matrix with appropriate size. Specifically, the outputs of the first and second layers are reshaped to  $(\text{num\_agents}, \text{embed\_dim})$  and  $(\text{embed\_dim}, 1)$  respectively, which corresponds to the first three rows (as there are three agents in the environment) and the final row in Figure 3(d) in the main text with  $\text{embed\_dim} = 32$  corresponding to the columns. A darker color represents a larger value.

### A.2 Additional Results of QMIX (CDQ)

In QMIX (CDQ), we apply CDQ on agent-wise utility functions  $Q_a$  in QMIX by  $\min(\bar{Q}_a(s', u'_a), Q_a(s', u'_a))$ , where  $u'_a = \arg \max_{u'_a} Q_a(s', u'_a)$ . Comparison results between applying CDQ on  $Q_a$  and the joint-action  $Q$ -function  $Q_{tot}$  (QMIX (CDQ-joint)) are shown in Figures 10(a)-(b). As demonstrated, QMIX (CDQ-joint) underperforms QMIX (CDQ), which also results in larger value estimates.

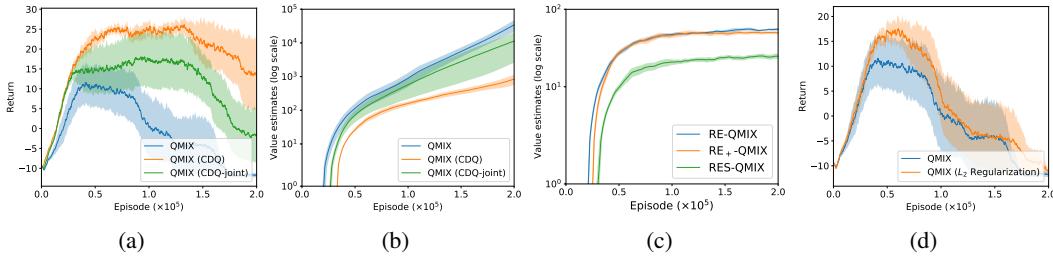


Figure 10: (a) and (b) Comparison of performance and value estimates of QMIX (CDQ) and QMIX (CDQ-joint). (c) Value estimates (in log scale) of RE-QMIX, RE<sub>+</sub>-QMIX, and RES-QMIX. (d) Performance comparison of QMIX and QMIX ( $L_2$  Regularization).

### A.3 Additional Results of a Clipped Version of RE-QMIX

Figure 10(c) shows comparison results for value estimates of RE-QMIX, RE<sub>+</sub>-QMIX, and RES-QMIX. Specifically, RE<sub>+</sub> uses a clipped version of the regularizer in RE-QMIX by  $\lambda((Q_{tot}(s, \mathbf{u}) - R_t(s, \mathbf{u}))_+)^2$ , where  $\lambda$  denotes the coefficient and  $(\cdot)_+ = \max(\cdot, 0)$ . As a result, RE<sub>+</sub>-QMIX only penalizes joint-action  $Q$ -values when  $Q_{tot}(s, \mathbf{u}) > R_t(s, \mathbf{u})$ . As shown in Figure 10(c), the value estimates for RE<sub>+</sub>-QMIX are very close to those of RE-QMIX. This is because the regularizer is active in most cases (94.4% on average during training).

### A.4 Additional Results of QMIX ( $L_2$ Regularization)

Figure 10(d) shows the performance comparison of QMIX and QMIX ( $L_2$  Regularization) with best regularization coefficient. The loss function of QMIX ( $L_2$  Regularization) is obtained by adding  $\lambda \|\theta\|_2^2$  to the original loss function [20], where  $\lambda$  is the coefficient. As shown, applying  $L_2$  regularization fails to avoid performance degradation.

## B Our Approximate Softmax Operator

### B.1 Proof of Theorem 1

**Theorem 1.** Let  $\mathbf{u}^*$  and  $\mathbf{u}'$  denote the optimal joint actions in  $\mathbf{U}$  and  $\mathbf{U} - \hat{\mathbf{U}}$  with respect to  $Q_{tot}$ , respectively. The difference between our approximate softmax operator and its direct computation in the whole action space satisfies:  $\forall s \in \mathcal{S}, |\text{sm}_{\beta, \hat{\mathbf{U}}}(Q_{tot}(s, \cdot)) - \text{sm}_{\beta, \mathbf{U}}(Q_{tot}(s, \cdot))| \leq \frac{2R_{\max}}{1-\gamma} \frac{|\mathbf{U} - \hat{\mathbf{U}}|}{|\mathbf{U} - \hat{\mathbf{U}}| + \exp(\beta(Q_{tot}(s, \mathbf{u}^*) - Q_{tot}(s, \mathbf{u}')))}$ , where  $|\mathbf{U} - \hat{\mathbf{U}}|$  denotes the size of the set  $\mathbf{U} - \hat{\mathbf{U}}$ .

*Proof.* By definition, we have that

$$|\text{sm}_{\beta, \hat{\mathbf{U}}}(Q_{tot}(s, \cdot)) - \text{sm}_{\beta, \mathbf{U}}(Q_{tot}(s, \cdot))| \quad (5)$$

$$= \left| \frac{\sum_{\mathbf{u} \in \mathbf{U}} \exp(\beta Q_{tot}(s, \mathbf{u})) Q_{tot}(s, \mathbf{u})}{\sum_{\bar{\mathbf{u}} \in \mathbf{U}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))} - \frac{\sum_{\mathbf{u} \in \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \mathbf{u})) Q_{tot}(s, \mathbf{u})}{\sum_{\bar{\mathbf{u}} \in \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))} \right| \quad (6)$$

For simplicity, we denote  $a = \sum_{\bar{\mathbf{u}} \in \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))$ ,  $b = \sum_{\bar{\mathbf{u}} \in \mathbf{U} - \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))$ ,  $c = \sum_{\mathbf{u} \in \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \mathbf{u})) Q_{tot}(s, \mathbf{u})$ , and  $d = \sum_{\mathbf{u} \in \mathbf{U} - \hat{\mathbf{U}}} \exp(\beta Q_{tot}(s, \mathbf{u})) Q_{tot}(s, \mathbf{u})$ . Then, Eq. (6) =  $\frac{b}{a+b} \left| \frac{c}{a} - \frac{d}{b} \right|$ .

We have that

$$\left| \frac{c}{a} - \frac{d}{b} \right| = |\text{sm}_{\beta, \hat{\mathbf{U}}}(Q_{tot}(s, \cdot)) - \text{sm}_{\beta, \mathbf{U} - \hat{\mathbf{U}}}(Q_{tot}(s, \cdot))| \quad (7)$$

$$\leq \left| \max_{\mathbf{u} \in \mathbf{U}} Q_{tot}(s, \mathbf{u}) - \min_{\mathbf{u} \in \mathbf{U}} Q_{tot}(s, \mathbf{u}) \right| \quad (8)$$

$$\leq \left| \max_{\mathbf{u} \in \mathbf{U}} Q_{tot}(s, \mathbf{u}) \right| + \left| \min_{\mathbf{u} \in \mathbf{U}} Q_{tot}(s, \mathbf{u}) \right| \quad (9)$$

$$\leq 2 \max_{\mathbf{u} \in \mathbf{U}} |Q_{tot}(s, \mathbf{u})| \quad (10)$$

$$\leq \frac{2R_{\max}}{1-\gamma}. \quad (11)$$

We also have that

$$\frac{b}{a+b} = \frac{\sum_{\bar{\mathbf{u}} \in U - \hat{U}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))}{\sum_{\bar{\mathbf{u}} \in U} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))} \quad (12)$$

$$\leq \frac{|\mathbf{U} - \hat{\mathbf{U}}| \exp(\beta Q_{tot}(s, \mathbf{u}'))}{|\mathbf{U} - \hat{\mathbf{U}}| \exp(\beta Q_{tot}(s, \mathbf{u}')) + \sum_{\bar{\mathbf{u}} \in \hat{U}} \exp(\beta Q_{tot}(s, \bar{\mathbf{u}}))} \quad (13)$$

$$\leq \frac{|\mathbf{U} - \hat{\mathbf{U}}| \exp(\beta Q_{tot}(s, \mathbf{u}'))}{|\mathbf{U} - \hat{\mathbf{U}}| \exp(\beta Q_{tot}(s, \mathbf{u}')) + \exp(\beta Q_{tot}(s, \mathbf{u}^*))} \quad (14)$$

$$= \frac{|\mathbf{U} - \hat{\mathbf{U}}|}{|\mathbf{U} - \hat{\mathbf{U}}| + \exp(\beta(Q_{tot}(s, \mathbf{u}^*) - Q_{tot}(s, \mathbf{u}'))} \quad (15)$$

Therefore, we obtain that

$$|\text{sm}_{\beta, \hat{U}}(Q_{tot}(s, \cdot)) - \text{sm}_{\beta, \mathbf{U}}(Q_{tot}(s, \cdot))| \leq \frac{2R_{\max}}{1 - \gamma} \frac{|\mathbf{U} - \hat{\mathbf{U}}|}{|\mathbf{U} - \hat{\mathbf{U}}| + \exp(\beta(Q_{tot}(s, \mathbf{u}^*) - Q_{tot}(s, \mathbf{u}')))}. \quad (16)$$

As a result, the gap converges to 0 in an exponential rate with respect to the inverse temperature parameter  $\beta$ .  $\square$

## B.2 Algorithm and More Details for Computing the Approximate Softmax Operator

The full algorithm for computing the approximate softmax operator is in Algorithm 1.

---

### Algorithm 1 Approximate softmax operator

---

- 1: Obtain the maximal joint action w.r.t.  $Q_{tot}$ :  $\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} Q_{tot}(s, \mathbf{u})$
- 2: Construct  $U_a$ :  $U_a = \{(u_a, \hat{\mathbf{u}}_{-a}) | u_a \in U\}$
- 3: Construct the joint action subspace:  $\hat{U} = U_1 \cup \dots \cup U_n$
- 4: Compute the approximate softmax operator:

$$\text{sm}_{\beta, \hat{U}}(Q_{tot}(s, \cdot)) = \sum_{\mathbf{u} \in \hat{U}} \frac{e^{\beta Q_{tot}(s, \mathbf{u})}}{\sum_{\mathbf{u}' \in \hat{U}} e^{\beta Q_{tot}(s, \mathbf{u}')}} Q_{tot}(s, \mathbf{u})$$


---

Figure 11 also illustrates the computation of our approximate softmax operator in the joint action subspace  $\hat{U}$ . The left part in Figure 11 corresponds to the maximal joint action w.r.t.  $Q_{tot}$ . The middle part demonstrates the joint action subspace  $\hat{U}$  (whose size is  $nK$ ) for computing our approximate softmax operator, where the yellow block means that the corresponding action can be one of the actions in  $U$ . The right part shows the joint action space  $\mathbf{U}$  (whose size is  $K^n$ ), with the orange block showing that the action can be one of the actions in  $U$  except for the local greedy action. As discussed in the main text, the action space in the multi-agent setting is much larger than that in the single-agent case, and some joint-action  $Q$ -value estimates  $Q_{tot}(s, \mathbf{u})$  can be unreliable due to a lack of sufficient training. As a result, directly taking them all into consideration for computing the softmax operator as in the single-agent case [36, 29] can result in inaccurate value estimates. As shown in Figure 11, according to the individual-global-max (IGM) property discussed in Section 2 in the main text,  $\hat{U}$  consists of joint actions that are close to the maximal joint action  $\hat{\mathbf{u}}$ , which is more likely to contain joint actions with more accurate and reliable value estimates. As a result, our softmax operator provides an efficient and reliable approximation. Theorem 1 in the main text provides a theoretical guarantee for  $\hat{U}$ , and Section 5.1.4 in the main text validates its effectiveness by comparing it with other choices.

## B.3 Discussion of the Softmax Operator on Agent-Wise Utilities

As discussed in the main text, we propose to employ the softmax operator to further mitigate the overestimation bias in the joint-action  $Q$ -function by  $\text{sm}_{\beta, \hat{U}}(Q_{tot}(s, \cdot))$ . One may be interested

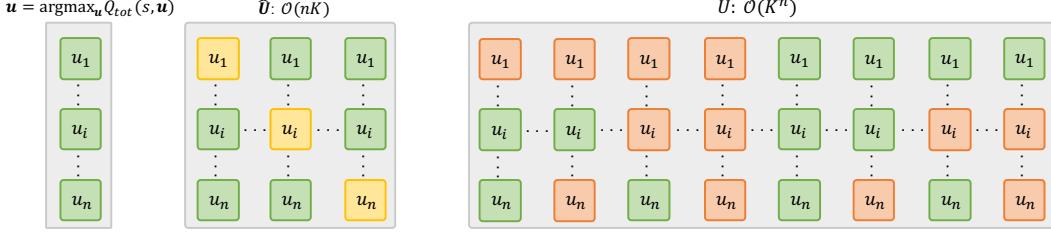


Figure 11: Left: The maximal joint action w.r.t.  $Q_{tot}$ . Middle: Illustration of the joint action subspace  $\hat{U}$  for computing our approximate softmax operator. Right: Illustration of the joint action space  $U$ .

in its application on the agent-wise utility functions by  $f_s(\text{sm}_{\beta,U}(Q_1(s,\cdot)), \dots, \text{sm}_{\beta,U}(Q_n(s,\cdot)))$ . The results are shown in Figure 12, where we refer to the method as RE-QMIX (softmax on  $Q_a$ ). From Figure 12(b), we can see that this results in overly pessimistic value estimates and a larger underestimation bias shown in Figure 12(c). In addition, as shown in Figure 12(a), it also significantly underperforms RES-QMIX, demonstrating the necessity of a careful design of the softmax operator in deep multi-agent  $Q$ -learning methods in MARL.

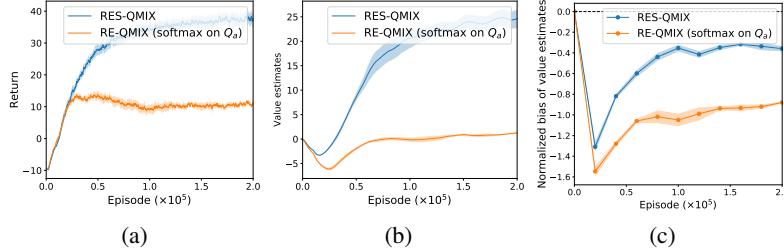


Figure 12: Comparison between RES-QMIX and RE-QMIX with the softmax operator on agent-wise utilities. (a) Performance. (b) Value estimates. (c) Normalized bias of value estimates.

## C Proof of Theorem 2

**Theorem 2.** *Given the same sample distribution, the update of the RES method is equivalent to the update using  $\mathcal{L}(\theta) = \mathbb{E}_{(s,u,r,s') \sim \mathcal{B}} [(y - Q_{tot}(s, u))^2]$  with learning rate  $(\lambda + 1)\alpha$ , which estimates the target value according to  $y = \frac{r + \gamma \text{sm}_{\beta,\hat{U}}(\bar{Q}_{tot}(s', \cdot))}{\lambda + 1} + \frac{\lambda R_t(s, u)}{\lambda + 1}$ .*

*Proof.* Let  $\theta$  and  $\bar{\theta}$  denote parameters of  $Q_{tot}$  and the target network  $\bar{Q}_{tot}$  respectively. The gradient for the learning objective of the RES method is

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{RES}}(\theta) &= -2(r + \gamma \text{sm}_{\beta,\hat{U}}(\bar{Q}_{tot}(s', \cdot | \bar{\theta})) - Q_{tot}(s, u | \theta)) \nabla_{\theta} Q_{tot}(s, u | \theta) \\ &\quad + 2\lambda(Q_{tot}(s, u | \theta) - R_t(s, u)) \nabla_{\theta} Q_{tot}(s, u | \theta). \end{aligned} \quad (17)$$

Then, we have

$$\theta' = \theta + 2\alpha(\lambda + 1) \left( \frac{r + \gamma \text{sm}_{\beta,\hat{U}}(\bar{Q}_{tot}(s', \cdot | \bar{\theta}))}{\lambda + 1} + \frac{\lambda R_t(s, u)}{\lambda + 1} - Q_{tot}(s, u | \theta) \right) \nabla_{\theta} Q_{tot}(s, u | \theta), \quad (18)$$

where  $\alpha$  is the learning rate. Therefore, it is equivalent to using a learning rate  $\alpha' = (\lambda + 1)\alpha$ , and estimating the target value  $y$  by  $\frac{r + \gamma \text{sm}_{\beta,\hat{U}}(\bar{Q}_{tot}(s', \cdot))}{\lambda + 1} + \frac{\lambda R_t(s, u)}{\lambda + 1}$ .  $\square$

## C.1 Additional Results

Figure 13 shows the comparison results for QMIX, RES-QMIX, and QMIX with learning rate (lr)  $(\lambda + 1)\alpha$ . We see that solely using a larger learning rate still fails to tackle the problem and cannot avoid performance degradation.

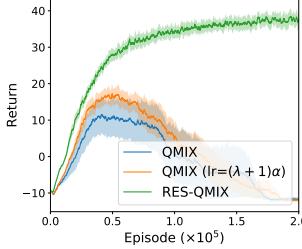


Figure 13: Comparison results of QMIX, QMIX ( $\text{lr}=(\lambda + 1)\alpha$ ), and RES-QMIX.

## D Proof of Theorem 3

**Theorem 3** Let  $B(\mathcal{T}) = \mathbb{E}[\mathcal{T}(s')] - \max_{\mathbf{u}'} Q_{\text{tot}}^*(s', \mathbf{u}')$  be the bias of value estimates of  $\mathcal{T}$  and the true optimal joint-action  $Q$ -function  $Q_{\text{tot}}^*$ . Given the same assumptions as in [41] for the joint-action  $Q$ -function, where there exists some  $V_{\text{tot}}^*(s')$  such that  $V_{\text{tot}}^*(s') = Q_{\text{tot}}^*(s', \mathbf{u}')$  for different joint actions,  $\sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s')) = 0$ , and  $\frac{1}{|\mathcal{U}|} \sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s'))^2 = C$  ( $C > 0$ ) with  $\bar{Q}_{\text{tot}}$  an arbitrary joint-action  $Q$ -function, then  $B(\mathcal{T}_{\text{RES-QMIX}}) \leq B(\mathcal{T}_{\text{RE-QMIX}}) \leq B(\mathcal{T}_{\text{QMIX}})$ .

*Proof.* For the left-hand-side, by definition, we have that

$$\text{sm}_{\beta, \hat{\mathcal{U}}}(\bar{Q}_{\text{tot}}(s', \cdot)) \leq \max_{\mathbf{u}' \in \hat{\mathcal{U}}} \bar{Q}_{\text{tot}}(s', \mathbf{u}') = \max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}'). \quad (19)$$

Therefore,  $\mathcal{T}_{\text{RES-QMIX}} \leq \mathcal{T}_{\text{RE-QMIX}}$ , and  $B(\mathcal{T}_{\text{RES-QMIX}}) \leq B(\mathcal{T}_{\text{RE-QMIX}})$ .

For the right-hand-side, we have that

$$B(\mathcal{T}_{\text{RE-QMIX}}) - B(\mathcal{T}_{\text{QMIX}}) = \mathbb{E} \left[ \frac{\lambda}{\lambda + 1} \left( R_{t+1}(s') - \max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}') \right) \right] \quad (20)$$

$$= \frac{\lambda}{\lambda + 1} \left( V_{\text{tot}}^\pi(s') - \mathbb{E} \left[ \max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}') \right] \right) \quad (21)$$

$$\leq \frac{\lambda}{\lambda + 1} \left( \max_{\mathbf{u}' \in \mathcal{U}} Q_{\text{tot}}^*(s', \mathbf{u}') - \mathbb{E} \left[ \max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}') \right] \right), \quad (22)$$

where  $V_{\text{tot}}^\pi(s')$  is the expected discounted return starting from state  $s'$  under the current behavior policy  $\pi$ , and Eq. (22) follows from the fact that its value is no larger than that of the optimal policy.

Then, it suffices to prove that  $\mathbb{E} [\max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}')] \geq \max_{\mathbf{u}' \in \mathcal{U}} Q_{\text{tot}}^*(s', \mathbf{u}')$ .

Assume that the joint-action  $Q$ -function follows the same assumptions of Theorem 1 in [41], i.e., there exists some  $V_{\text{tot}}^*(s')$  such that  $V_{\text{tot}}^*(s') = Q_{\text{tot}}^*(s', \mathbf{u}')$  for different joint actions,  $\sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s')) = 0$ , and  $\frac{1}{|\mathcal{U}|} \sum_{\mathbf{u}'} (\bar{Q}_{\text{tot}}(s', \mathbf{u}') - V_{\text{tot}}^*(s'))^2 = C$  for some  $C > 0$ . The assumption means that value estimates are correct on average, but there exists estimation error in the joint-action  $Q$ -value estimates for some joint actions. Then, following the analysis in [41], we have that  $\max_{\mathbf{u}' \in \mathcal{U}} \bar{Q}_{\text{tot}}(s', \mathbf{u}') - \max_{\mathbf{u}' \in \mathcal{U}} Q_{\text{tot}}^*(s', \mathbf{u}') \geq 0$ .

Therefore, RE-QMIX can reduce the overestimation bias of QMIX, and we have that  $B(\mathcal{T}_{\text{RE-QMIX}}) \leq B(\mathcal{T}_{\text{QMIX}})$ .  $\square$

## E Experimental Details

### E.1 Experimental Setup

**Tasks.** The multi-agent particle environments [21] are based on an open-source implementation,<sup>7</sup> where the global state is the concatenation of observations of all agents. For StarCraft Multi-Agent Challenge (SMAC),<sup>8</sup> we use the latest version 4.10.<sup>9</sup>

**Baselines.** Value factorization methods including VDN, QMIX, and QTRAN are implemented using the PyMARL [34] framework,<sup>10</sup> while we use authors' open-source implementations for Weighted QMIX<sup>11</sup> and QPLEX.<sup>12</sup> The actor-critic method MADDPG is based on an open-source implementation,<sup>13</sup> which uses the Gumbel-Softmax trick to tackle discrete action spaces [21]. We use default hyperparameters and setup as in [34], where we list as in Table 1. The only exception is that we tune the target-update-interval for value factorization methods in multi-agent particle environments, as these algorithms fail to learn in these environments with default hyperparameters that are originally fine-tuned for SMAC environments. Specifically, the target-update-interval is 800 for QMIX, QTRAN, Weighted QMIX and QPLEX (which updates the target network every 800 episodes), while it remains 200 for VDN as the default value works well. Each agent network is a deep recurrent Q-network (DRQN) consisting of 3 layers: a fully-connected layer, a GRU layer with a 64-dimensional hidden state, and a fully connected layer with ReLU activation. The mixing network consists of a 32-dimensional hidden layer using ELU activation, and the hypernetwork [12] consists of two layers with 64-dimensional hidden state using ReLU activation. Note that QPLEX uses a different architecture for the mixing network with more parameters, which consists of two modules (with default hyperparameters as in [43]): a transformation network and a dueling mixing network with a multi-head attention module [42]. By default, all  $Q$ -learning based MARL algorithms (including VDN, QMIX, QTRAN, Weighted QMIX, QPLEX and our RES method) use double estimators as in Double DQN [41] to estimate the target value as pointed out in Appendix D.3 in [32]. All experiments are run on P100 GPU. The code will be released upon publication of the paper.

Table 1: Hyperparameters.

Hyperparameter	Value
Discount factor	0.99
Replay buffer size	5000 episodes
Batch size	32 episodes
Warmup steps	50000
Optimizer	RMSprop
Learning rate	$5 \times 10^{-4}$
Initial $\epsilon$	1.0
Final $\epsilon$	0.05
Linearly annealing steps for $\epsilon$	50k
Double DQN update	True

**RES.** RES is also implemented based on the PyMARL [34] framework with the same network structures and hyperparameters as discussed above. For our RES method, to estimate the target joint-action  $Q$ -value using the softmax operator based on double estimators, the softmax weighting is computed in the joint action subspace based on current joint-action  $Q$ -network  $Q_{tot}$  in Eq. (23):

$$sm_{\beta, \hat{U}}(\bar{Q}_{tot}(s, \cdot)) = \sum_{u \in \hat{U}} \frac{e^{\beta Q_{tot}(s, u)}}{\sum_{u' \in \hat{U}} e^{\beta Q_{tot}(s, u')}} \bar{Q}_{tot}(s, u). \quad (23)$$

<sup>7</sup><https://github.com/shariqibal2810/multiagent-particle-envs>

<sup>8</sup><https://github.com/oxwhirl/smac>

<sup>9</sup>Note that the results reported in [34] use SC2.4.6.2.69232, and performance is not always comparable across versions.

<sup>10</sup><https://github.com/oxwhirl/pymarl>

<sup>11</sup><https://github.com/oxwhirl/wqmix>

<sup>12</sup><https://github.com/wjh720/QPLEX>

<sup>13</sup><https://github.com/shariqibal2810/maddpg-pytorch>

As discussed in Section 5.1.3 in the main text, we only need to tune  $\lambda$  while keeping  $\beta$  fixed (whose performance is competitive within a wide range of values).

In multi-agent particle environments, for RES-QMIX, the inverse temperature  $\beta$  of the softmax operator is fixed to be 0.05, where the regularization coefficient  $\lambda$  is selected based on a grid search over  $\{1e-2, 5e-2, 1e-1, 5e-1\}$ . Specifically,  $\lambda$  is  $1e-2$  for covert communication (CC),  $5e-2$  for predator-prey (PP) and world (W), and  $5e-1$  for physical deception (PD). As for RES-Weighted-QMIX, it is based on OW-QMIX [31] with an optimistic weighting, as it outperforms its counterpart CW-QMIX as shown in Figure 5 in the main text. The parameter  $\beta$  is 0.1 while  $\lambda$  is  $5e-2$  for PP and W, and  $5e-1$  for PD and CC. As for RES-QPLEX, we set  $\beta = 0.1$  and  $\lambda = 5e-2$  for all environments. For RES-QMIX in SMAC, the regularization coefficient  $\lambda$  is  $1e-2$  for 3m and the super hard map MMM2 while being  $5e-2$  for the remaining maps. The hyperparameter  $\beta$  for the softmax operator is fixed to be 5.0 for all maps.

### E.1.1 Multi-agent Particle Environments

Table 2 summarizes detailed information for our tested environments based on the multi-agent particle framework with discrete action space, where observation dim and action dim correspond to dimension of observation space and action space for the agents respectively. In world,<sup>14</sup> unmovable entities also include forests (which are accessible) and foods besides the inaccessible landmarks. Figure 14 shows the illustration of the multi-agent particle tasks.

Table 2: Information of environments in the multi-agent particle framework.

Name	#Agents	#Adver-saries	#Land-marks	State dim	Observation dim	Action dim
Predator-prey (PP)	3	1	2	16	62	5
Physical deception (PD)	2	1	2	10	28	5
World (W)	4	2	1	34	200	5
Covert communication (CC)	2	1	2	8	20	4

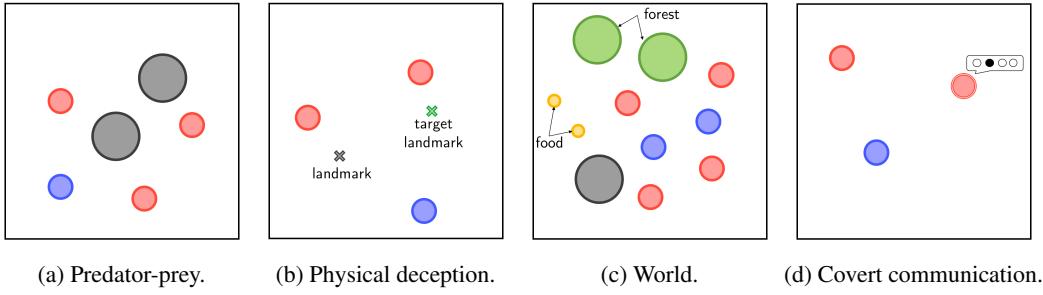


Figure 14: Illustration of the multi-agent particle tasks.

### E.1.2 StarCraft Multi-Agent Challenge (SMAC)

Table 3 describes the tested maps in SMAC, and Figure 15 shows the illustration for the maps.

## E.2 Mean Normalized Return in Multi-Agent Particle Environments

Figure 16 shows the mean normalized return of different algorithms averaged over all multi-agent particle environments. As shown, RES-QMIX significantly outperforms state-of-the-art methods in performance and efficiency.

<sup>14</sup>This is a variant of the environment simple\_world\_comm, where we modify the environment to allow for Discrete action space instead of MultiDiscrete action space.

Table 3: Information about agents, enemies and difficulty of tested maps in SMAC.

Name	Agents	Enemies	Difficulty
3m	3 Marines	3 Marines	Easy
2s3z	2 Stalkers and 3 Zealots	2 Stalkers and 3 Zealots	Easy
3s5z	3 Stalkers and 5 Zealots	3 Stalkers and 5 Zealots	Easy
2c_vs_64zg	2 Colossi	64 Zerglings	Hard
MMM2	1 Medivac, 2 Marauders, and 7 Marines	1 Medivac, 3 Marauders, and 8 Marines	Super hard



Figure 15: Illustration of the StarCraft II micromanagement tasks.

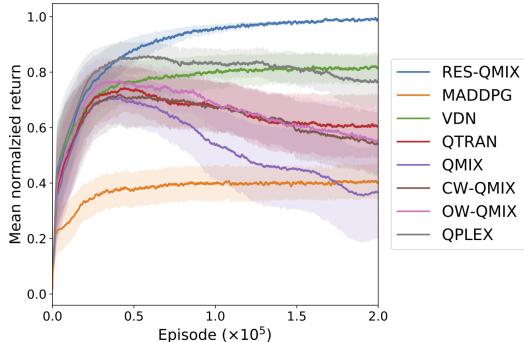


Figure 16: Mean normalized return averaged over all multi-agent particle environments.

### E.3 Performance of RES-QMIX in Stochastic Environments

In this section, we investigate the robustness of RES-QMIX in the stochastic environment. To support stochasticity in the environment, we conduct evaluation in a variant of the predator-prey task based on sticky actions [22], which has been widely used in Atari games [4] to inject stochasticity to the environments. Specifically, at each timestep, the action selected by the agent will be executed with probability  $1 - p$ , and with probability  $p$  the last action taken by the agent will be executed (where the stickiness  $p$  is typically set to be 0.25 in [22]).

Figure 17 shows the performance of RES-QMIX in predator-prey with sticky actions in different levels. As shown, RES-QMIX is robust to different levels of stochasticity and significantly outperforms QMIX in all cases.

### E.4 Comparison of Value Estimates in Predator-Prey

Figure 18 shows comparison of true values and estimated values (which are obtained in the same way as in Section 4 in the main text) of different algorithms in predator-prey. As discussed in the main text,

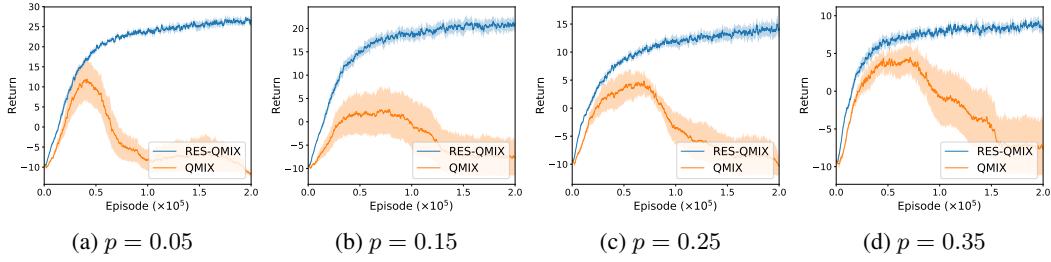


Figure 17: Performance comparison of RES-QMIX and QMIX in stochastic variants of the predator-prey task.

value estimates for QMIX and Weighted QMIX increase rapidly, which leads to large overestimation bias (Figure 6 in the main text) and severe performance degradation (Figure 5(a) in the main text). Value estimates of VDN (which is based on a linear decomposition of the joint-action  $Q$ -function) increase more slowly at the end of training, but still incurs large overestimation bias as in QTRAN and QPLEX. Unlike all other value factorization methods, MADDPG learns an unfactored critic that directly conditions on the full state and joint action. It is less sample efficient, which indicates that value factorization is important in these tasks. Thus, MADDPG results in a lower return (Figure 5(a) in the main text) and value estimates compared to all other value factorization methods, but still overestimates. RES-QMIX achieves the smallest bias and fully mitigates the overestimation bias of QMIX, resulting in stable performance and outperforming all other methods (as shown in Figure 5(a) in the main text).

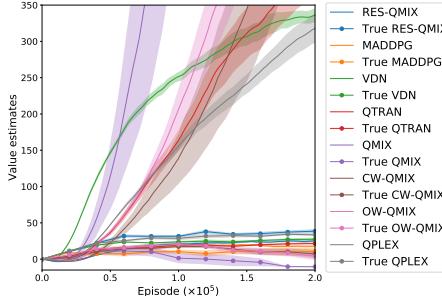


Figure 18: Comparison of true values and estimated values of different algorithms in predator-prey.

## E.5 Full Ablation Study

Figures 19 shows the ablation study of the effect of each component in our approach in different environments by comparing RES-QMIX, RE-QMIX, S-QMIX, and QMIX. As shown, the regularization component is critical for stability, while combining with our softmax operator further improves learning efficiency (or avoids performance oscillation).

## E.6 Full Learning Curves of RES-QMIX/Weighted QMIX/QPLEX

Our RES method is general and can be readily applied to different  $Q$ -learning based MARL algorithms. To demonstrate its versatility, we apply it to three different deep multi-agent  $Q$ -learning algorithms: QMIX, Weighted QMIX, and QPLEX. Full comparison of learning curves of RES over QMIX, Weighted QMIX, and QPLEX in the multi-agent particle environments are shown in Figure 20, 21, and 22, respectively. As shown, our RES-based methods significantly outperform corresponding baseline methods.

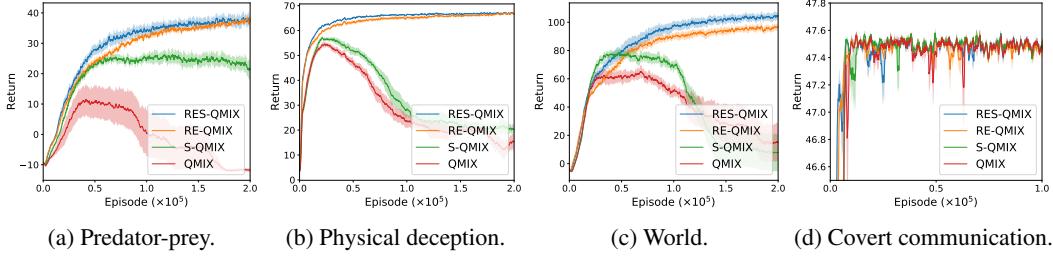


Figure 19: Ablation study of the effect of each component.

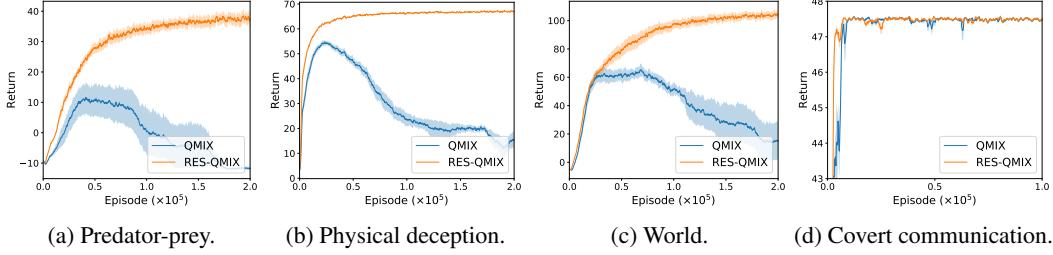


Figure 20: Performance comparison of QMIX and RES-QMIX.

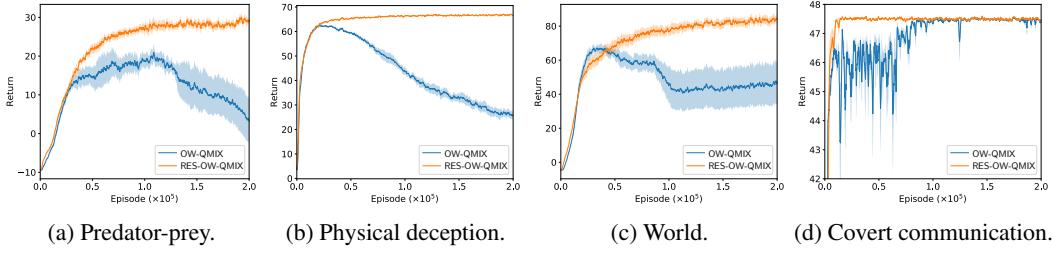


Figure 21: Performance comparison of Weighted QMIX and RES-Weighted QMIX.

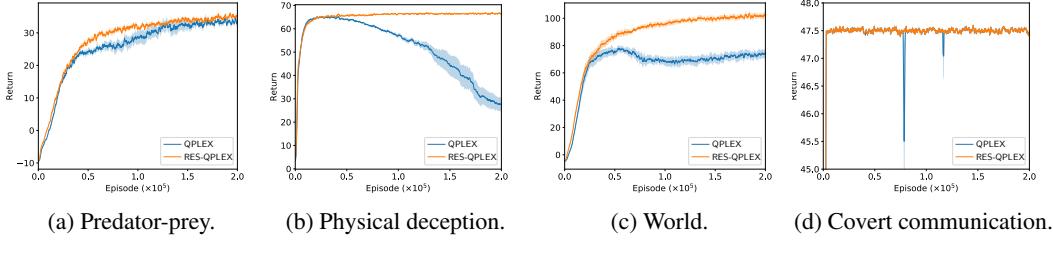


Figure 22: Performance comparison of QPLEX and RES-QPLEX.

## E.7 Performance Comparison with SM2

Gan et al. [10] propose the soft Mellowmax (SM2) operator to tackle overestimation in reinforcement learning. Figures 23 (a)-(d) show the comparison results of RES-QMIX, SM2-QMIX (with fine-tuned hyperparameters) and QMIX in each environment, and Figure 23(e) summarizes the mean normalized return. As shown, SM2-QMIX fails to tackle the severe overestimation problem, and also significantly underperforms RES-QMIX.

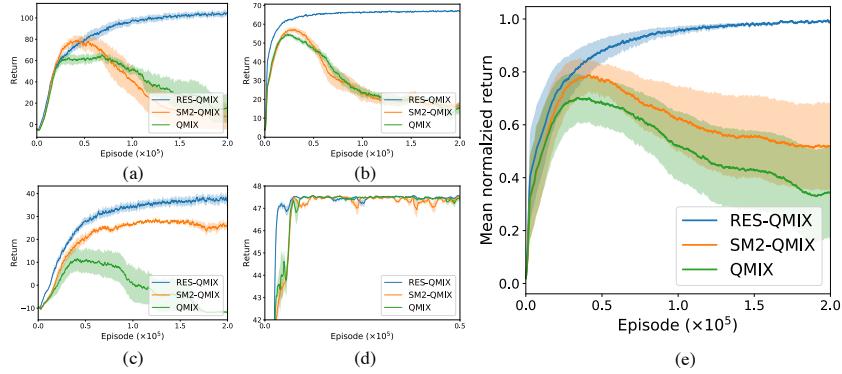


Figure 23: Performance comparison of RES-QMIX, SM2-QMIX and QMIX. (a) Predator-prey. (b) Physical deception. (c) World. (d) Covert communication. (e) Mean normalized return.

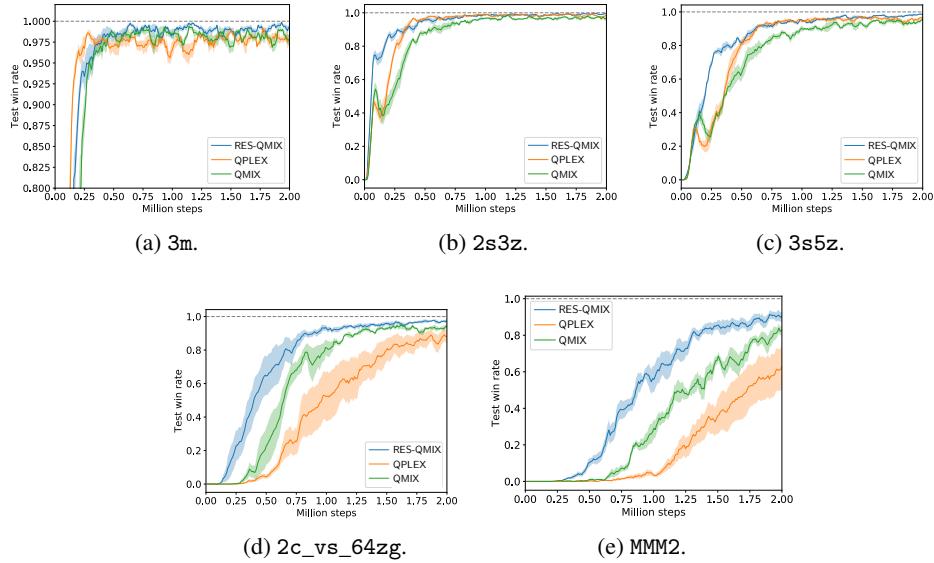


Figure 24: Comparison of test win rate of QMIX, RES-QMIX and QPLEX.

## E.8 Performance Comparison of RES-QMIX with Other Baseline Method in SMAC

In the main text, we analyze how much of a performance improvement RES-QMIX achieves over QMIX in StarCraft II micromanagement tasks in Figure 9. We also compare RES-QMIX in StarCraft II micromanagement tasks against QPLEX, which is the most competitive algorithm in multi-agent particle environments. Experimental settings are the same as in Section E.1.

Figure 24 shows the test win rate, where RES-QMIX significantly outperforms QPLEX in final test win rate and sample efficiency in all but one environments. The only exception is the easy map 3m, where QPLEX is more sample efficient but underperforms RES-QMIX at the end of training. Specifically, the final test win rate is 98.4% and 97.1% for RES-QMIX and QPLEX in 3m respectively. It is also worth noting that QPLEX uses more parameters due to the multi-head attention module for the hypernetwork compared with RES-QMIX and QMIX.

## E.9 Comparison of Value Estimates of RES-QMIX and QMIX in SMAC

Comparison results of value estimates of RES-QMIX and QMIX in StarCraft II micromanagement tasks are shown in Figure 25, where RES-QMIX achieves smaller value estimates than QMIX, which demonstrates its effectiveness.

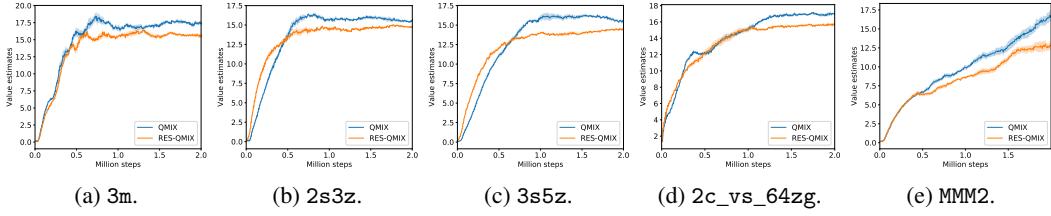


Figure 25: Value estimates of RES-QMIX and QMIX in StarCraft II micromanagement tasks.

### E.10 Performance of RES-QMIX with Single Estimator in SMAC

Comparison of test win rate of RES-QMIX (single), QMIX (single), and QMIX is shown in Figure 26. Experimental setup is the same as in Section E.1, with the only exception of removing double estimators, i.e., it does not estimate the target value using double estimators as in Double DQN [41]. Environment-specific hyperparameters for RES-QMIX (single) in StarCraft II micromanagement tasks include  $\lambda$  and  $\beta$ . For the super hard map MMM2,  $\lambda$  is  $1e-1$  while being  $5e-2$  for the remaining maps. The parameter  $\beta$  is 0.5 for 3m, 5 for 2s3z, 3s5z and MMM2, and 50 for 2c\_vs\_64zg. For RES-QMIX (single), the computation of the softmax operator is  $\text{sm}_{\beta, \hat{U}}(\bar{Q}_{\text{tot}}(s', \cdot))$  following Eq. (1) in the main text.

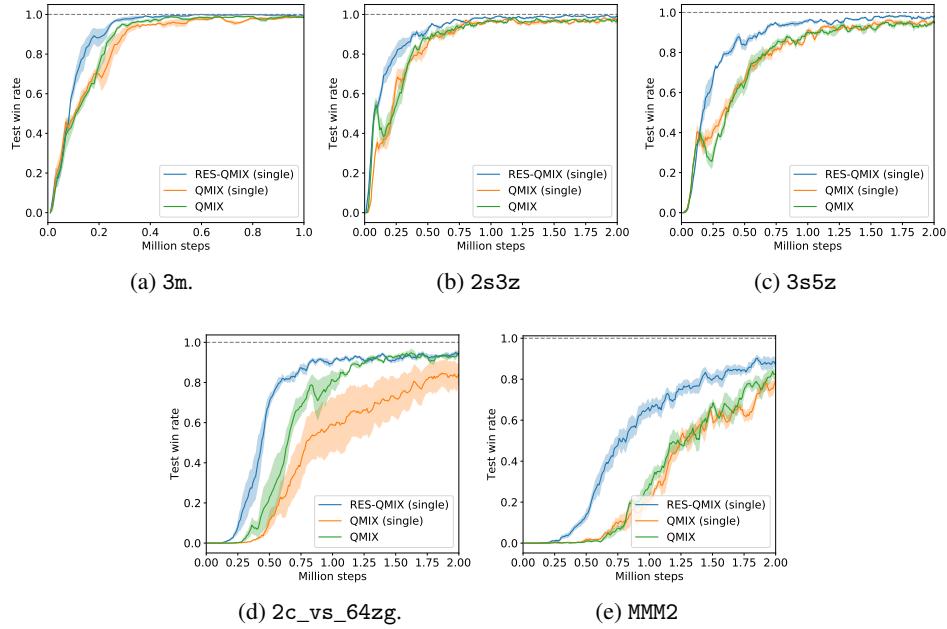


Figure 26: Performance comparison of RES-QMIX (single), QMIX (single), and QMIX.

From Figure 26, we can see that RES-QMIX (single) significantly outperforms QMIX and QMIX (single), demonstrating that our RES method is still effective even without double estimators.