

1 We thank the reviewers for their feedback on our sub-
 2 mission. We have fixed a config error when running OW-
 3 QMIX on Predator Prey (there is now very little vari-
 4 ance). We have also run additional experiments demon-
 5 strating significantly better performance of weighted
 6 QMIX on another *hard* SMAC map *bane_vs_bane*.

7 **Reviewer 1:** > "Why not include QTRAN, MADDPG,
 8 or MASAC on any of the SMAC experiments?"

9 We have already included QTRAN, MADDPG and
 10 MASAC on the SMAC experiments in Figure 2. Due to
 11 their relatively poor performance there we did not run them on the *Super-Hard* SMAC maps (Figures 3 and 4) due to
 12 the large computational cost of those experiments.

13 > "How does using a QMIX approximation to Q^* work?... And won't this cause the policy induced by Q^* to be different?"

14 In order to enable tractable maximisation, we use our QMIX approximation (Q_{tot}) to Q^* to suggest the best joint action.
 15 In general, the greedy action for Q^* and Q_{tot} can differ during training. However, Theorems 1 and 2 prove that given
 16 sufficient training (and an appropriate α) they will be the same.

17 > "...the modified architecture uses a hypernetwork layer—does this mean it is restricted to positive weights like QMIX?"

18 Yes, the first layer of \hat{Q}^* 's mixing network is restricted to non-negative weights, but \hat{Q}^* is not restricted to being
 19 monotonic due to subsequent layers.

20 **Reviewer 2:** > "The weighted QMIX only modifies QMIX by using a weighting function to get the Q_{tot} , and the two
 21 kinds of weighting functions seem too simple, so the contribution seems incremental."

22 We disagree that the simplicity of the weighting function makes our approach too incremental. The use of a weighting
 23 function in order to train a monotonic approximation to a learned unrestricted Q^* is a significant algorithmic change
 24 over QMIX. Additionally, we have proven that the two weighting functions we have considered are guaranteed to ensure
 25 the maximal joint action is correct (given sufficient training and an appropriate α) in contrast to QMIX which can fail to
 26 recover the optimal joint action for the simple matrix game in Table 2. Furthermore, the framework we have introduced
 27 for analysing Weighted QMIX can be used to analyse QTRAN and explain its empirical performance.

28 > "Extra computation cost also restricts its scalability."

29 Compared to QMIX, during training we must perform inference and train an additional model (with the same complexity
 30 as QMIX). This does **not** restrict the scalability of Weighted QMIX compared to QMIX, as demonstrated by our
 31 experiments on *bane_vs_bane* featuring 24 agents. We will include a discussion of the two papers you have provided.

32 **Reviewer 3:** > "The proof of your theory lacks discussion of POMDP settings." We deliberately restricted our
 33 theoretical analysis to the MMDP setting in order to avoid the additional complexity of partial observability. The
 34 MMDP setting allows for a cleaner presentation that focuses on our main goal of analysing the effect of the limited
 35 representation of QMIX on the learned Q_{tot} (and thus the learned policy).

36 > "...performance of QMIX+ \hat{Q}^* has a significance difference ... in Figure 8... The use of weighting is not that convinced."

37 On 3s5z the weighting does not affect performance, but on 5m_vs_6m it has a significant effect and on Predator Prey
 38 every method without the weighting is unable to solve the task, showing that it is crucial to our method.

39 > "In Section 6.2.3, the performance of the Weighted QMIX method is unacceptable." That's the point: Section 6.2.3
 40 aims to show the limitations of our method, which we believe is important for identifying areas for future research.

41 > "The authors argue that the complexity introduced by \hat{Q}^* is responsible for the regression in performance."

42 Figures 4 and 5 demonstrate a clear performance difference for Weighted QMIX when **only** changing the architecture
 43 used to represent \hat{Q}^* . This provides evidence that the poor performance is due to the architecture used to represent \hat{Q}^* .

44 > "... α of weighing Function, although the author gives a basis for selection in the appendix, the value of α seems not to
 45 be verified." Our theoretical results only show that there exists an α which works in **all** cases. They are not intended to
 46 provide a method for selecting an appropriate α . We will discuss the selection of α for experiments in the main paper.

47 > "... agent's ordering of actions is pointed to be important in representing value functions. But the proposed architecture
 48 seems to be incapable for dealing that case." \hat{Q}^* is capable of representing **any** joint action Q -value function.

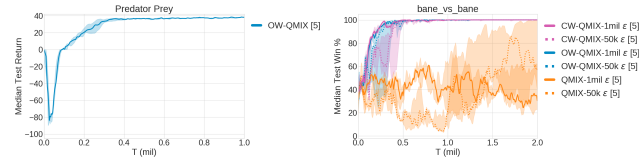
49 > "... sampling uniformly from a replay buffer does not strictly lead to a uniform weighting schema. However in the
 50 realization of Weighted QMIX you provided in Section 5, the loss in Equation 8 also suffers from the same problem."

51 A uniform weighting is an assumption we make to simplify our analysis (to make it clearer). It is not required for the
 52 Deep RL realisation of QMIX or Weighted QMIX.

53 > "How does the content of lines 163-167 relate to context?" Lines 163-167 explain why the failure modes discussed in
 54 Section 3.1 are problematic. In particular, they show fundamental limitations of QMIX that cannot be addressed without
 55 a significant algorithmic change, even in the idealised setting we consider. We will fix points 4, 5, and 6 on Clarity.

56 > "How will input s to the Mix network portion be handled during execution?"

57 During decentralised execution only the agent parts of Q_{tot} are required. We do not need access to the state, the mixing
 58 network, or \hat{Q}^* during execution.



(Left) Corrected experimental results for OW-QMIX on Predator Prey. (Right) Weighted QMIX vs QMIX on *bane_vs_bane* for ϵ annealed over 50k and 1mil.