1 **We appreciate all reviewers for their feedback!** We're glad that they find our methods well presented (**R2,3,4**), moti-
2 vated (**R3,4**), and contextualized (**R2,4**), novel (**R1,2,4**), simple and practical (**R3**), and experiments well-designed (**R2**).
3 **Reviewer #1** Q1. Definition of $Q$. The critic aims to estimate the joint action-value based on the **action probabilities** (AP).
4 As discussed in L121-128, our intuition is to train policies directly towards optimal cooperation with **full differentiability**,
5 and we use sampled actions (**special cases of AP with probability 1**) for critic training because the target values (defined
6 by action-specific rewards $r(s_t, u_t)$) over arbitrary AP are hard to estimate. In fact, similar ideas were explored for single RL
7 settings [Wierstra,Schmidhuber,ECML'07; Weber et al, AISTATS'19] with proper justification. We'll add more discussion
8 and citations accordingly. Q2. $k$ iterations critic update. Yes, $k$ is intended to give better critic estimation and tuning LR is
9 equivalent; it was included as a **practical generalization**: training till convergence could take long and risk overfitting, and
10 tuning $k$ instead of LR may avoid overshooting. Q3. $k$ for other PG baselines? Yes, e.g. $k$=2 for LICA/MADDPG and $k$=1 for
11 others work best empirically in SC II. We'll revise to avoid confusion. Q4. Different $\lambda$ for MLP vs mixing critic. We observed
12 that the **architecture change alone resulted in more stochastic joint actions**, and as clarified in L296, the choices of $\lambda$
13 for MLP critic ensure a fair comparison of **policy stochasticity** (Fig.2(b)) against the best LICA run ($\lambda$=0.09). We found
14 that setting $\lambda$=0.09 for MLP critic clearly results in over-regularization and gives even worse performance. Q5. Need more
15 runs/inconsistency with SMAC paper. We want to point out that our results on all maps except `2c_vs_64zg` are consistent
16 with previous work (e.g. [3,20,21]); for `2c_vs_64zg` specifically, our investigation suggests that the inconsistency is due
17 to a **mismatch in SC II gameplay version**: we base our experiments on the latest SMAC repo which uses **v4.10**, while
18 SMAC paper seems to use **v4.6** (commit history); critically, **v4.7** added changes that made Colossi units more powerful,
19 changing the dynamics of `2c_vs_64zg`. Nevertheless, we'll add more runs for SCII as suggested. Q6. Compare with
20 MAVEN. As suggested, we added comparisons on 2 **Super Hard** maps in Fig.A/B. With same #iterations, **LICA performs**
21 **considerably better**. Q7. Why $t$ in $s_t$ for Eq.2? Optimizing expected returns over different $t$ is rather standard and often
22 implied under various notational choices; e.g. see [4,3,28] and their implementation. Q8. Eqn for per-agent policy gradients.
23 Due to full differentiability (L145), the PG for agent $a \propto \sum_t \nabla_{\theta_a} p_t^a \nabla_{p_t^a} Q^\pi \left( s_t, p_t^1, ..., p_t^a, ..., p_t^n \right)$ with $p_t^a = \pi_{\theta_a}^a (\cdot | z_t^a)$;
24 we'll update accordingly. Q9. Details of MPE. For Fig.3(b,c), we use 200 steps (L214), -1 reward for every pairwise
25 collision, and we report the **mean reward over all timesteps and agents** in each episode. We'll clarify the metrics in the
26 paper; see also our base repos [13,28]. Q10. Add discussions for QMIX/MADDPG. Thanks! We'll update accordingly.
27 **Reviewer #2** *Thanks for recognizing our work!* Q1. LICA in continuous domains. While this is a future extension, we
28 emphasize that LICA doesn't pose extra constraints on top of previous work [4,9,13] that readily handles continuous actions.
29 **Reviewer #3** Q1. Benefits/novelty of mixing critic. Let us consider the generalization where both MLP critic ($C_{MLP}$) and
30 mixing critic ($C_{Mix}$) operate on representations of states and actions $f_s(s)$, $f_a(a)$. Then, in both cases, we have $\frac{\partial Q}{\partial a} = \frac{\partial Q}{\partial h} \frac{\partial h}{\partial a}$,
31 where $h = f_s(s) + f_a(a)$ for $C_{MLP}$ and $h = f_s(s) f_a(a)$ for $C_{Mix}$ is the **first mixed representation** of $s,a$ before activation (i.e.
32 after concat+linear for $C_{MLP}$ and before $\sigma(\cdot)$ for $C_{Mix}$, Fig.1(b)). Since $g(h)=Q$ is non-linear/non-interpretable in both cases,
33 the crucial difference is thus that $\frac{\partial Q}{\partial a} = \frac{\partial Q}{\partial h} \frac{\partial h}{\partial f_a} \frac{\partial f_a}{\partial a} = \frac{\partial Q}{\partial h} \frac{\partial f_a}{\partial a}$ for $C_{MLP}$ and $\frac{\partial Q}{\partial a} = \frac{\partial Q}{\partial h} \frac{\partial h}{\partial f_a} \frac{\partial f_a}{\partial a} = \frac{\partial Q}{\partial h} f_s(s) \frac{\partial f_a}{\partial a}$ for $C_{Mix}$, i.e. $C_{Mix}$
34 **adds an extra, direct state representation**. ...do not necessarily lead to better credit assignment (CA): While better CA is
35 not *guaranteed*, we argue **better utilization** of state provides a basis for better CA. Rightness of $\frac{\partial Q}{\partial a}$ ...determined by accuracy
36 of $Q(s,a)$... $C_{Mix}$ just learns a better $Q(s,a)$? we argue that the **composition** of $\frac{\partial Q}{\partial a}$ in $C_{Mix}$ is the key factor, and a better $Q(s,a)$,
37 if any, would rather be a result of it. $C_{MLP}$ also contains state...: We intend to convey that $C_{Mix}$ has a better utilization of $s$ and
38 will revise all inaccuracies in Sec 3.2. Discussion (3,4)...aren't contributions: We'll revise accordingly; note that they remain
39 valid and were discussed as LICA's *properties* rather than novelties. Concat after MLP for $C_{MLP}$: As suggested, we ran a com-
40 parison in Fig.C where MLPs are added before concat; results confirm our earlier analysis which covers this case. Q2. Could
41 LICA converge to stable policies? While we cannot provide a full analysis here, we emphasize that our empirical evidence
42 across different $\lambda$'s, scenarios, complexity (Fig.4(a-f)), and environments with repeated runs (Fig.3/4) suggests that policies
43 eventually reach a stochasticity equilibrium (Fig.2(b,c)); this may in fact sustain smoother object landscapes and aid policy
44 convergence [1]. Q3. Compare with MAAC. By design, the simplicity of the quoted **1-step** game obviates most key aspects
45 that differentiate on/off-policy learning (future estimation, separate target/behavior nets, replay buffers) and focuses only on
46 the **mechanism for credit assignment**. However, we appreciate your suggestion and will add this discussion accordingly.
47 **Reviewer #4** Q1. Improvements in MPE. We stress that compared to the



48 previous SOTA [28], our method achieved similar gains despite approach-
49 ing the limits of the selected envs. Q2. Complex settings w/ uneven mix of
50 'individual performance' and 'cooperation'. In fact, MMM2 (**Super Hard**,
51 Fig.4(f), Supp L20-23, and demo) is *precisely* one such setting where our method has **sizable advantage over others**.
52 Winning heavily relies on the performance of the 1 healer unit and cooperation of the 9 attack units. Q3. SC II: further
53 training/more complex settings. We emphasize that many previous work mainly focuses on **Easy** maps (e.g. [3,4,20])
54 and lacks diversity in map choices (e.g. [3,4,20,14,ROMA ICML'20]); on our diverse maps (L252-254), **we achieved similar**
55 **or significantly more gains** compared to previous work **with similar #iterations**. At **R4**'s request, we also added results
56 on 2 extra **Super Hard** maps (6h_vs_8z,3s5z_vs_3s6z) in Fig. A/B, showing **sizable gains over previous methods**.
57 Q4. It reads more like a report. We respectfully disagree. On top of **R2**'s recognition and our above response, we'd also
58 highlight our comparison against SOTA in 2019 [3,25] and our extensive component studies (Sec 4.3, Supp A2, Fig.2) that
59 are equally or more comprehensive compared to previous work (e.g. [3,4,20,28,14]).