

1 **Reviewer 1.** We appreciate R1’s recognition of the novelty of our contribution to MARL and the potential impact on a  
 2 range of problems. We address R1’s two concerns below. **Regarding our chosen baselines**, we note that baselines  
 3 we include represent three major existing categories: 1) policy gradient and actor-critic with discrete or continuous  
 4 “give-reward” actions are direct applications of conventional RL (which have been applied to multi-agent incentivization  
 5 in recent work (Lupu et al. 2020)); 2) LOLA is an archetype of second-order approaches; 3) Inequity Aversion (IA)  
 6 draws domain knowledge from models in evolutionary biology and sociology to alter individual rewards. Hence we  
 7 believe the existing baselines allow for fair benchmarking of our new approach. While Social Influence (Jaques et  
 8 al. 2019) does not have open-source code, we can make an indirect comparison by noting that IA reaches a score  
 9 around 250 by  $1.6 \times 10^8$  steps (Figure 3a in IA (Hughes et al. 2018)), outperforming Social Influence that reaches  
 10 score of 200 by  $3 \times 10^8$  steps (Figure 1a in Jaques et al. 2019) in the original Cleanup map with 5 agents. Hence,  
 11 the fact that LIO outperforms IA in our experiments implies that LIO compares favorably with Social Influence.  
 12 We clarify that LIO technically does not involve a second-order gradient, as the gradients  
 13 are w.r.t. separate parameters  $\theta$  (policy) and  $\eta$  (incentive function). **Regarding scalability**,  
 14 we clarify that the bi-level optimization does not necessarily imply difficulty in scaling  
 15 up, because the learning of incentives is conducted in a *pairwise* manner: in equation  
 16 (7) for a fixed reward-giver agent  $i$ , each term of the summation corresponds to the pair  
 17  $(i, j)$  for recipient  $j \neq i$ . Figure 1 shows that LIO attains the global optimum collective  
 18 reward of 17 ( $= 2 \times 10 - 3$ ) in the Escape Room game with  $N = 5$  agents, out of which  
 19  $M = 3$  agents are incentivized to cooperate despite penalties of  $-1$  each. Scaling up to  
 20 large populations poses new questions regarding population-level phenomenon (such as  
 21 social norms) that modulate the impact of incentives; we leave this to future work.

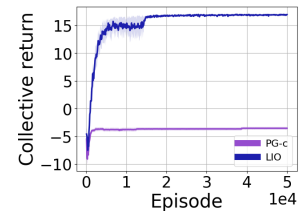


Figure 1: Escape Room ( $N = 5, M = 3$ )

22 **Reviewer 2.** We appreciate R2’s positive feedback on our quantitative results and we are glad that our behavioral  
 23 analysis of the learned incentive functions provided insight. We believe this work is a suitable contribution to the  
 24 NeurIPS community, in addition to the broad area of multi-agent learning, as we tackle the open question of emergent  
 25 cooperation from decentralized learning charted out in Hughes et al. (NeurIPS 2018) by building on general meta-  
 26 gradient methods (Xu et al., NeurIPS 2018). In our revision, we will elaborate on the wide range of new research  
 27 questions generated by our work, including theoretical analysis of dynamically-changing incentive functions, new  
 28 population-level effects in a scaled up context, and adaptive ways to account for the currency of incentives. Regarding  
 29 Figure 6b where the agent gives nonzero reward for “fire cleaning beam but miss” after 40k steps, one reason is that the  
 30 agent’s actual partner in training already converged to the behavior of consistently cleaning waste successfully (LIO in  
 31 Figure 6a), so it may have “forgotten” the difference between successful and unsuccessful usage of the cleaning beam.  
 32 As demonstrated more clearly in the Escape Room results (e.g. Figures 5b and 5d), this can be avoided by choosing a  
 33 sufficiently large lower bound on the exploration rate by all agents, so that all agents pose the risk of deviating from  
 34 cooperative behavior, which forces LIO to maintain correct incentivization.

35 **Reviewer 3.** We thank R3 for recognizing our contribution to the general class of opponent-shaping algorithms. We  
 36 address each concern as follows. 1) Our definition of “decentralized” focuses solely on the inability to optimize social  
 37 welfare directly, which is the crux of social dilemmas, and which holds regardless of access to the global state (e.g.,  
 38 Prisoner’s Dilemma is fully observable). Hence our definition does not mention the degree of observability. 2) We  
 39 explain in Appendix A that the coefficient  $\alpha$  in the cost for incentivization is indeed an important hyperparameter  
 40 and provide intuition for how to choose it in practice. We agree that a sweep over  $\alpha$  can provide more insight. More  
 41 broadly, we believe there is room to develop an adaptive scheme to trade off between small  $\alpha$ , which allows time to  
 42 learn the effect of incentives, and large  $\alpha$ , which penalizes redundant incentivization. 3) In Figure 5d, by episode 50k,  
 43 the cooperator still receives nonzero incentives between 0.5 and 1, but the winner’s received incentives has noticeably  
 44 converged to zero. 4) We became aware of Hostallero et al. (AAMAS 2020) after submission, and we believe there is a  
 45 crucial methodological difference from our work. They use the temporal difference error of agent  $k$ ’s  $Q$ -function to  
 46 modify the rewards of  $k$ ’s peers, so that those peer agents have incentive to take actions that lead to more favorable  
 47 reward than average for agent  $k$ . This is a *passive* approach from agent  $k$ ’s viewpoint, and is closer to the method in  
 48 Hughes et al. (2018), since each agent’s original reward is modified by a hand-designed function based on other agents’  
 49 performance. In contrast, a LIO agent *actively* differentiates through the recipient’s learning step to update the learned  
 50 incentive function, which it uses to change the recipient’s total reward. We will include this in the revision.

51 Hostallero, D. et al. (2020). Inducing Cooperation through Reward Reshaping based on Peer Evaluations in Deep Multi-Agent  
 52 Reinforcement Learning. *AAMAS*.

53 Hughes, E. et al. (2018). Inequity aversion improves cooperation in intertemporal social dilemmas. *NeurIPS*.

54 Jaques, N. et al. (2019). Social influence as intrinsic motivation for multi-agent deep reinforcement learning. *ICML*.

55 Lupu, A. and Precup, D. (2020). Gifting in multi-agent reinforcement learning. *AAMAS*.

56 Xu, Z. et al. (2018). Meta-gradient reinforcement learning. *NeurIPS*.