

1 We would like to thank the reviewers for their thorough evaluations and for bringing to our attention some missing
2 citations and typos. We answer specific questions raised by the reviewers, below.

3 **Performance w.r.t. number of agents (R1).** We discuss how our methods scale with the number of agents in Appendix
4 H. Specifically, Figure 16 shows that our method’s benefits hold but that the underlying algorithm (MADDPG in this
5 case) fails to handle many agents; this has also been shown in [12].

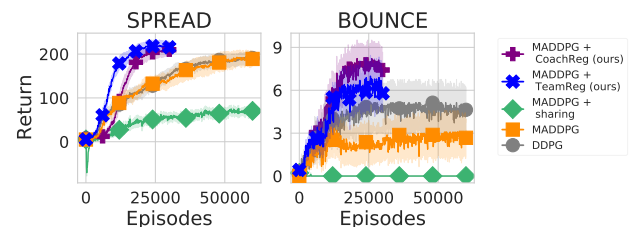
6 **Related work and novelty (R2, R3).** (To R2 and R3) We are grateful for bringing to our attention some relevant work
7 in hierarchical RL. Importantly, however, the novelty of CoachReg does not lie in training sub-policies (which are
8 obtained here through a simple and novel masking procedure) but rather in co-evolving synchronized sub-policies
9 across multiple agents. This is indeed closer to the joint exploration work of Mahajan et al. (2019)¹ (as pointed out
10 by R2). Yet, a major difference is that MAVEN’s situational-prediction occurs only on the first timestep and requires
11 synchronized random seeds across the agents *at test time*, whereas with CoachReg agents explicitly learn a set of
12 subpolicies, of which they choose one to execute at every timestep. Each agent chooses without using a common
13 sampling procedure and execution is therefore fully decentralized. We will update our manuscript with these more
14 explicit clarifications. (To R2) Thank you for referring us to Hong et al.² whose method is very close to “MADDPG +
15 agent-modelling”, the TeamReg ablation that we compare against in Section 7, Figure 5. As we discuss in Related Work
16 (L211-214), agent-modelling (through cross-entropy prediction) is now a widely used MARL component and Hong et
17 al., like [11], uses it as an auxiliary task to learn richer representations. Similarly, TeamReg relies on agent-modelling
18 (L141-146) but our contribution with TeamReg is to instead use it to explicitly influence other agents’ behavior toward
19 being predictable rather than just learning a representation (L145-148 and L219-222). To our knowledge, this is a novel
20 contribution and has not been considered in prior work.

21 **Positioning of the paper and missing keywords (R2).** While the high level positioning of this work in the Centralized
22 Training Decentralized Execution framework (CTDE) is already made clear throughout the paper (L20, 38, 42, 237,
23 266, 323), we will highlight it in the abstract as suggested. However, we do not believe that “the planning setting”
24 (usually referring to making use of a transition model rather than the agents model) or “self-play” (where an agent,
25 short of having an opponent to train with, plays against itself) are relevant keywords for our work.

26 **Importance to the broader community, reflection, motivation and transfer (R2).** As highlighted by other reviewers,
27 our work makes significant contributions to the research community: at a high level we question the widespread
28 assumption that centralized training always outperforms decentralized training, proposing a definition for coordinated
29 behavior (based on behavior predictability), in order to improve upon it. We propose two (2) novel practical coordination
30 promoting methods that are applicable to any CTDE algorithm and evaluate them on three (3) different baselines based
31 on the prevalent MADDPG algorithm, as well as two (2) ablated versions of our methods.

32 **MADDPG baseline (R3).** We disagree with the premise that MADDPG is a weak baseline and argue that the evaluation
33 setting plays a major role in allowing valid and insightful experimental results. Several recent works have pointed out
34 that hyperparameter tuning often plays a fundamental role in determining which algorithms best perform at a given task
35 (Henderson et al. (2018)³, Colas et al. (2019)⁴). In our work, we make a substantial effort to offer fair and significant
36 comparisons by allowing our three (3) baselines (DDPG, MADDPG, MADDPG + sharing) and two (2) ablations
37 (MADDPG + agent modelling, MADDPG + policy masks) a full hyperparameter tuning, yielding a competitive suite
38 of baselines. To substantiate the importance of such re-tuning, we provide here additional experiments reporting the
39 improvements of our tuned MADDPG over MADDPG with the original hyperparameters configuration from [22].
40 The improvements are 900 % (SPREAD), 1300 % (BOUNCE), 700 % (COMPROMISE) and 400 % (CHASE) and
41 highlight the important performance gains allowed by our evaluation procedure to the baselines.

42 **Conclusiveness of the results (R4).** As requested, we ex-
43 tended the training of the baselines and the inlined figure
44 shows that our methods still outperform them. Addition-
45 ally, we believe that our evaluation is sound, conclusive
46 and substantiates our claims (as concluded by R1). A key
47 question here is “do the proposed coordination-inducing
48 methods improve performance of the CTDE framework?”.
49 We answer this by examining the impact of our proposed
50 ideas on the widely used MADDPG CTDE algorithm and we perform an ablation study to probe each element of our
51 contribution in more detail. We apply a careful experimental methodology (Tables 2 through 11) to both continuous and
52 discrete action environments of varying complexity, requiring *significant* computing resources: e.g., the retraining in
53 Table 1 of our submission alone require 120 CPU-days and prevented us from extending the number of training steps
54 illustrated in this rebuttal.



¹ Mahajan et al., MAVEN : Multi-Agent Variational Exploration (2019) ² Hong et al. A Deep Policy Inference... (2018)

³ Henderson et. al., Deep RL that matters. (2018) ⁴ Colas et al., A Hitchhiker’s Guide to Statistical Comparisons of RL ... (2019)