

1 **Rev. 1** 1) Yes, the formulation of “robust Markov game” originated from [12]. But our contribution is to model the
2 “robust MARL” problem as such a game and develop algorithms to solve it, both of which are original. Introducing
3 the “nature” agent is exactly one contribution that makes the computation tractable and easier to understand in MARL
4 context. It is not simply a “trick”, but our effort to transfer [12] to robust MARL. We will add more clarification on this.
5 2) With general function approximation, PG is solving a non-convex optimization problem, the theoretical guarantee of
6 its optimality is very challenging. In fact, the optimality guarantee of PG methods, even in the single-agent setting, can
7 only be established in limited cases until recently. Note that the equilibrium guarantee of the original MADDPG is
8 non-existent yet. As our work is the first one that dealt with robust MARL under the framework of robust Markov games,
9 many other aspects, e.g., the algorithm design, solution existence, Q-learning convergence, and empirical validation,
10 need to be done first. 3) We had put experiment details in the supplement due to page limit. We will re-organize the
11 experiment session to add those details and include results for the cooperative navigation in the main content. We have
12 updated all the tables with 95% confidence intervals across 5 trials, 1000 evaluation episodes for each. Below is an
13 example for keep-away. The figures for the cooperative navigation in the supplement have also shown 95% CI across
14 five runs. We hope our response can help address your concerns on the novelty, and re-evaluate our submission.

15 **Rev. 2** Thanks for the very positive feedback, we will address the grammatical typos and the notations in Table
16 2, and add more clarification on the nature agent. 1) We will re-iterate the “distribution-free” point and add
17 clarification to avoid the possible confusion with “distributional uncertainties”. 2) For Keep-away, we have run
18 student’s t tests. All calculated p -values are below 0.01, suggesting that the improvements are statistically significant.

19 **Rev. 3** 1) We are not sure what was meant by “as I men-
20 tioned above”, as there was only the reviewer’s summary
21 above it. 2) The “nature” policy is approximated by
22 neural networks, as for other agents. We are confused
23 about the comment “eq 3.5 only minimizes the reward”,
24 since Eq. 3.5 has no minimization. There is no “original
25 reward” in our robust MARL setting, as the reward func-
26 tion is “uncertain”. Instead, the nature agent is always
27 “playing against” all agents, while the agents want to find
28 the policy that account for this adversarial uncertainty.
29 The MSE in line 13 is some heuristic to regularize the
30 nature policy’s output, with η being the weight coefficient.
31 We have mentioned the choice of η in Sec. F (appendix),
32 and will add more clarification in the main paper. 3) Clarity: We will address the typos. The “robust Markov
33 games” is fundamentally different from “Markov games with stochastic reward/transition”. The latter is nothing but the
34 standard MG model, since it allows the reward/transition to be stochastic. However, our robust MG model accounts for
35 the “model uncertainty” explicitly, like the single-agent setting in [1]. The mathematical description of uncertainty
36 is explicitly provided in the robust MG formulation. 4) Related work: We note that they are not that relevant, as
37 non-stationary MDP usually deals with an online/adaptive setting, with the focus of reducing the accumulated regret,
38 while we account for the “model uncertainty” from the beginning, and our algorithms are model-free. Most importantly,
39 our main focus is on the “multi-agent” side. 5) The MADDPG paper [19] has already shown that MADDPG
40 outperforms DDPG in all these particle environments. There is no added value to compare to DDPG again. 6) About
41 the variance, please check our response for reviewer 1, answer 3. Since we run evaluation over 1000 episodes for each
42 trial, the variance, i.e., CI, is fairly acceptable. 7) We have the same intuition that the confidence interval for high
43 noise level (7.0) should be larger than 5.0/6.0, we think we get the narrower CI only by chance. 8) We have started
44 running the suggested ablation experiments that the nature always selects the minimum reward in an uncertainty set.
45 But the training takes much longer for convergence and we could not obtain the final results before deadline. We will
46 add it to the final version. We notice that most of your comments are regarding the clarification, and the “weakness”
47 part seems not that fatal. We sincerely hope our response can help address your concerns, and re-evaluate our paper.

48 **Rev. 4** 1) Clarity and framing: our focus was “the agents may not have perfect information of the model” from each
49 agent’s perspective, including both *its own* and *others’* model. Sorry for the confusion in the traffic network example, we
50 will clarify on this point. 2) When there is no reward noise, non-R-MADDPG methods can obtain the exact reward,
51 while R-MADDPG still uses the nature agent to approximate the certain reward, which leads to extra approximation
52 errors. However, R-MADDPG does not always “perform notably worse” than the baselines. E.g., in the keep-away
53 experiments ($\lambda = 0$ or 1), R-MADDPG adversary performs better than MADDPG most of the time, and slightly worse
54 than M3DDPG. In practice, when entering a new environment, we suggest running both algorithms and compare. 3)
55 Please check our response to reviewer 1, answer 3 and reviewer 3, answer 6 on the variance concern. 4) We will
56 explain the “Without loss of generality...” sentence with more details in revision. 5) We will provide the proof of
57 Proposition 2.2 for completeness. As the main weakness was on clarify and framing, which can be addressed as per
58 your suggestion easily, we really appreciate the reviewer to re-evaluate our contribution, especially considering that our
59 work provides the first framework that accounts for “model-uncertainty” in MARL, with both theory and simulations.

Models		Uncertainty level (λ)			
AG	ADV	0.0	1.0	2.0	3.0
M3	M3	13.23 (0.12)	6.34 (0.18)	3.01 (0.14)	1.34 (0.09)
M3	MA	8.07 (0.18)	7.51 (0.17)	3.47 (0.14)	2.86 (0.13)
M3	RM	13.15 (0.12)	7.23 (0.18)	5.63 (0.17)	4.31 (0.16)
MA	M3	13.57 (0.11)	8.12 (0.18)	2.58 (0.13)	1.66 (0.09)
MA	MA	10.91 (0.16)	6.04 (0.17)	3.89 (0.15)	1.44 (0.09)
MA	RM	12.52 (0.13)	6.96 (0.17)	4.05 (0.15)	3.47 (0.14)
RM	M3	11.87 (0.15)	6.98 (0.18)	5.65 (0.16)	2.09 (0.11)
RM	MA	8.83 (0.17)	7.38 (0.18)	4.37 (0.16)	3.25 (0.14)
RM	RM	7.89 (0.18)	8.21 (0.18)	4.32 (0.16)	3.61 (0.14)

Table 1 in the paper updated with 95% confidence interval.