

1 We sincerely thank all the reviewers, and feel really honored to receive such positive and constructive comments. We  
2 are in the process of incorporating many of the changes into the final version.

3 **Reviewer 1** 1) We will improve the broader impact section by emphasizing the implications of our theoretical  
4 results on applications. 2) We will mention total variation distance in the appendix, and correct the typo on “Corollary  
5 3.3”. Thanks for the careful reading.

6 **Reviewer 2** 1) Note that the smooth planning oracle is not needed throughout the paper, and is thus not the “primary  
7 concept” in our paper. It is only used in Sec. 3.2 for finding  $\epsilon$ -NE policies with near-optimal sample complexity. We  
8 have justified this in lines 209-231, and will add some demonstration on simple games in the final version. 2) Yes, by  
9 saying “rarely”, R-MAX is definitely one of the very few ones. We have discussed R-MAX in lines 82-83. We will  
10 change the wording in the final version. 3) By saying “especially model-free ones...” this sentence, we simply meant  
11 that the convergence techniques for single-agent Q-learning/model-free methods cannot be applied directly in MARL.  
12 The works on Q-learning in games you mentioned exactly conquered this issue, with non-trivial efforts. We will add  
13 more clarifications on this. 4) We will address all the grammatical comments/typos in the final version. Thanks a lot  
14 for the careful examining.

15 **Reviewer 3** 1) We will add more comparisons with [4,17] on the lower bound proof, in the final version. 2) We  
16 thought the “generative model” setting is a standard one. Thanks for pointing this, and we will add the clarification. We  
17 believe an adaptive policy guided sampling may indeed decrease the sample complexity of “model-based” methods, if  
18 more structure about the “model” is known beforehand, e.g., which state transition is more significant. Otherwise, for  
19 general models without special structure, estimating the whole transition model (approximately) seems inevitable, and  
20 our matching upper-bound seems to be hardly improvable.

21 **Reviewer 4** On the smooth planning oracle: First, technically, our “leave-one-out” proof technique requires some  
22 “smoothness” of the change of either the “value-function” or the “policy” (for proving  $\epsilon$ -approximate NE value and  
23  $\epsilon$ -NE policy, respectively), when one state is made “absorbing”. When showing the near-optimal sample complexity  
24 to find  $\epsilon$ -approximate NE value, the change in value function is indeed smooth, so no “smooth planning oracle” is  
25 needed; when showing the near-optimal sample complexity to find  $\epsilon$ -NE policy, the small change in the Markov  
26 game cannot guarantee the change of the “best-response” policy of the opponent to be small, too. In fact, a small  
27 change in the value may correspond to a drastic change in the best-response policy. Thus, we introduce the smooth  
28 planning oracle to address this. Second, computationally, as we have discussed in lines 209-231, such an oracle can be  
29 readily satisfied in practice, especially in the regularized setting. Regularized setting makes solving the matrix game a  
30 strongly-convex-strongly-concave problem, and thus improves the computational efficiency, over the un-regularized  
31 ones. Whether it is possible to achieve “near-optimal sample complexity” for achieving “ $\epsilon$ -NE policy” without the  
32 smooth oracle may require totally different proof techniques, and is left as our future work.