

1 We thank all the reviewers for their valuable feedback and appreciating our contributions. We first address some  
2 common concerns.

3 **The proof applies only to deterministic systems / Deterministic systems seem highly restrictive.** Despite deter-  
4 ministic systems seem restrictive in theory, in practice, lots of RL problems are indeed deterministic. Moreover,  
5 all known algorithms that work under the assumption that the optimal  $Q$ -function is linear require deterministic or  
6 near-deterministic systems [Wen and Van Roy, 2013, Du et al., 2019].

7 **The proof depends on the assumption on the gap optimality / The model has to be very correct.** In this paper,  
8 we show that unless the gap  $\rho = \Omega(\sqrt{\dim_E} \delta)$  where  $\delta$  is the approximation error, any algorithm requires exponential  
9 number of samples even just to find a near-optimal policy (see Proposition 1.2). Therefore, such an assumption is  
10 necessary for any algorithm with polynomial sample complexity.

11 Please find our response to each individual reviewer below.

12 — **To Reviewer #1** —

13 **The algorithm for the general case requires an oracle.** When the number of actions is finite (as in Atari games),  
14 the agent can possibly enumerate all actions and find  $f_1, f_2 \in \mathcal{F}$  separately for each action by running continuous  
15 optimization algorithms that can handle constraints (e.g. projected gradient ascent). When the action space is continuous  
16 (as in control tasks), the agent could directly optimize  $a, f_1, f_2$  by running continuous optimization algorithms (as done  
17 in practice). Moreover, we would like to note that our paper is concerned with the statistical efficiency, and the oracle  
18 does not require any new sample (it solves an optimization problem based on existing samples).

19 Compared to the “Know-What-It-Knows” oracle, our uncertainty oracle just requires solving an optimization problem,  
20 while it is even unclear whether the “Know-What-It-Knows” oracle can be implemented statistically efficiently for  
21 general function classes. We will make the comparison clearer in the next version.

22 **The range of the return is assumed to lie in  $[0,1]$ .** This is a standard regularity assumption in RL theory, and is  
23 required to make sure that the empirical mean of the reward values concentrates around their expectation by taking  
24 enough samples. Such assumption is required for the algorithm in the supplementary material (Section D). In general,  
25 if the summation of the reward values is in  $[0, C]$ , then the sample complexity of the algorithm in Section D will be  
26 increased by a factor of  $C^2$ . Note that the required assumption that  $\rho = \Omega(\sqrt{\dim_E} \delta)$  keeps unchanged even if one  
27 changes the range of the reward values. E.g., if one scales all reward values by a factor of  $C$ , then the ratio between  $\rho$   
28 and  $\delta$  remains unchanged.

29 — **To Reviewer #2** —

30 We would like to thank the reviewer for the positive feedbacks.

31 — **To Reviewer #3** —

32 **Cannot operate in the scenario where there exists more than one optimal policy.** We disagree that our algorithm  
33 does not work in the scenario where there exists more than one optimal policy. Consider the case that for some state  $s$ ,  
34 there are three actions  $a_1, a_2$  and  $a_3$ . If  $Q^*(s, a_1) = Q^*(s, a_2) = 1$  and  $Q^*(s, a_3) = 0$ , then by Definition 3.1, the gap  
35 would be 1 and our algorithm still works. However in this case, it is clear that there could be more than one optimal  
36 policy.