

1 We thank all the reviewers for their valuable feedback. We first address some common concerns.

2 **Deterministic systems.** Although deterministic systems seem restrictive in theory, in practice, lots of RL problems are  
3 indeed deterministic. Also, this assumption makes the problem more tractable. We will emphasize that we focus on  
4 deterministic systems in the next version, and will also leave extending our result to stochastic environments as an open  
5 problem.

6 **Practicality.** We stress that our goal is to design provably efficient algorithms for RL with general reward functions.  
7 In this paper, we focus on giving sufficient and necessary conditions that admit efficient algorithms, and we believe  
8 our algorithmic insights (discretization, augmenting state space) can be applied in practice, which we are currently  
9 exploring. However, this is not the focus of the current work.

10 ——— **To Reviewer #1** ———

11 **Section 5.** The main goal of the algorithm in Section 5 is to motivate the discretization procedure used in the more  
12 complicated algorithm in Section 6, and thus we do not focus on optimizing the approximation guarantee. In the next  
13 version we will provide pseudocode in Section 5 to make the description of the algorithm formal.

14 **Correctness of the complexity result / The algorithm in Section 6 is also rather brute force.** The number of  
15 multisets of cardinality  $k$ , with elements taken from a finite set of cardinality  $n$ , is  $\binom{n+k-1}{k} \leq (k+1)^n$ . This bound  
16 can be found on the wikipedia page of multiset. For our case,  $k = H$  and  $n = \Theta(\log(1/\delta)/\delta)$ , and thus the number of  
17 possible multisets is at most  $H^{\Theta(\log(1/\delta)/\delta)}$ . Notice that "multisets" are different from "sequences", i.e., for multisets  
18 we do not care about orders of elements. Indeed, the number of sequences of length  $k$ , with elements taken from a  
19 finite set of cardinality  $n$ , could be as large as  $n^k$ . This also explains why our algorithm is not brute force: we carefully  
20 discretize reward values to make the number of possible elements small, and we exploit the symmetric of the objective  
21 function so that we only need to deal with multisets instead of sequences to avoid an exponential dependency on  $H$ .

22 ——— **To Reviewer #2** ———

23 We are grateful to the reviewer for providing detailed comments on the writing of our paper, and we will revise the  
24 paper according to the reviewer's comment in the next version.

25 **What prevents these techniques from being applied in stochastic environments?** Consider the following case in  
26 the symmetric norm setting (Section 5 in our paper), which suggests that the stochastic case is fundamentally more  
27 difficulty: there are two actions at the initial state, and all further actions do not affect the rewards. If the first action is  
28 chosen, then half of the reward values will be 1 and half of the reward values will be 0. If the second action is chosen,  
29 then all reward values will be a fair coin (0 or 1 with equal probability). To find a near-optimal policy, the agent must  
30 carefully compare the expected objective values of both choices and thus cannot be handled by our algorithm. However,  
31 if rewards are deterministic, one can simply return the action with more reward values of 1 which is optimal. Thus,  
32 even for the setting that rewards are stochastic and transitions are deterministic, the problem becomes much harder.

33 ——— **To Reviewer #3** ———

34 **Finite-horizon problems / exponential in the horizon length.** We would like to remind the reviewer that the running  
35 time of our algorithm is polynomial in the planning horizon  $H$  instead of being exponential in  $H$ . See Theorem 4.1 for  
36 the precise statement. One can reduce discounted MDPs to finite-horizon MDPs by considering the first  $\tilde{O}(1/(1-\gamma))$   
37 levels, and the  $H$  dependency in the complexity of our algorithm will be replaced by  $\tilde{O}(1/(1-\gamma))$ .

38 **Improve the presentation.** We will revise the paper, make the proofs more readable and improve the presentation in  
39 general according to the reviewer's comments in the next version. In particular, we will explain the high-level ideas at  
40 the beginning of Section 4 and 5 instead of at the end.

41 ——— **To Reviewer #4** ———

42 **Improvements to the presentation / care around the use of language / Broader Impact.** We are grateful to the  
43 reviewer for providing detailed comments on the writing of our paper, and we will revise the paper according to the  
44 reviewer's comments in the next version.

45 **"query" / "exponential number of values".** We will make it clear in the next version that we are proving lower bounds  
46 on the number of times that the agent evaluates the objective function  $f$ . Here by a "query", we mean that the agent  
47 evaluates the objective function  $f$  on some specific input.

48 **Literature on non-Markovian reward.** We are grateful to the reviewer for providing a comprehensive list of papers  
49 on non-Markovian reward, and we are planning to add them into the related work section in the next version.