

1 We thank the reviewers for the positive feedback and their interest in our work! Below we address some questions.

2 **Response to Reviewer 2: Empirical evaluation:** Interestingly, we actually did an empirical evaluation in the earlier  
3 phase of the project. We found that UCBZero can actually learn the navigation policies to all states in a 20 x 20 random  
4 2D maze with the same number of samples required for UCB-H to learn the navigation policy to a particular target  
5 state. Both algorithms are well-tuned for hyperparameters. We didn't include it in the submission because after all the  
6 algorithm is only for tabular RL and is not of much empirical value to real complex problems. But we indeed hope that  
7 our work provides some theoretical justification to the multi-task RL problem showing that efficient multi-task learning  
8 is at least possible in the tabular setting. **Variables used before defined:** Thank you very much for pointing that out.  
9 We will make sure to define them earlier in the paper in the revision.

10 **Response to Reviewer 3:** Thanks for the many technical questions! We are happy to clarify them. **Additional related**  
11 **works:** [Hazan et. al. 2019] is an interesting work that we are not aware of, which explicitly aims at visiting each  
12 state as uniformly as possible. In that sense, it is more similar to the reward-free paper of [Jin et. al. 2020], which  
13 explicitly aims at visiting each state enough time as a subgoal of the exploration phase. Our result, however, shows  
14 that the approach taken by Jin et. al. is too pessimistic (suffering from an additional factor of  $S$ ) if the eventual goal  
15 is still to perform tasks. The other two papers [Even Dar et. al. 2009, Shani et. al. 2020] studies the problem of  
16 regret minimization under adversarial rewards, which is very different from our setting, where the rewards are still  
17 stochastic but only unseen during exploration, and the goal is still to perform best policy identification, rather than  
18 regret minimization because the regret is not even defined in our setting without a reward function. Overall, we think  
19 it's not very precise to define UCBZero as an algorithm that merely tries to visit all the state-action pairs in an MDP.  
20 Arguably, all provably efficient RL algorithms must visit all essential states enough times in order to guarantee either  
21 small regret or small sample complexity. The goal of UCBZero is to perform multi-task learning even when the reward  
22 is absent during exploration. We did show, however, that as a bi-product, UCBZero visits each state enough times  
23 to allow other forms of downstream learning tasks. **UCRL2 with reward zero and nonoptimal dependency on H:**  
24 These are great questions! And they are connected. Since UCBZero is essentially a zero-reward version of UCB-H,  
25 our dependence on  $H$  essentially inherits from UCB-H, which is known to be suboptimal because of its model-free  
26 nature. At a higher-level, our result is an example that the zero-reward version of a standard provably efficient RL  
27 algorithm can be adapted to the task-agnostic RL setting and achieve good sample complexity. It is therefore very  
28 reasonable to conjecture that a more efficient RL algorithm, e.g. UCRL2, if adapted in a similar way, can achieve a  
29 better dependency on  $H$ . This is a great direction for future work. **Necessity of  $\log N$ :** Again, great question! The key  
30 observation is that we assume the rewards are stochastic, so even if you have lots of data on each state-action pair, given  
31 infinitely many stochastic reward functions, there will be one reward function whose instantiation deviates from its  
32 mean, with a large probability. As a very simple sample, consider that you want to estimate the mean of  $N$  Bernoulli  
33 random variables simultaneously. The number of times you need to sample each Bernoulli random variable will scale  
34 as  $\log N$  due to union bound. In fact, this result is profound. It shows that the reward-free RL task defined in [Jin et. al.  
35 2020] is impossible when the reward functions are stochastic. The bound must scale with  $\log N$ . All the questions are  
36 very good, and we could see why readers may have these confusions. We will make sure to discuss all of the above  
37 points in the revision. Thanks again for the very valuable feedback!

38 **Response to Reviewer 4: Empirical works in task-agnostic RL:** Yes, this is actually one of the big motivations  
39 behind this paper. We've seen a lot of empirical works on this topic (e.g. an ICLR workshop last year), but very little  
40 theoretical discussion, and our work aims to close the gap. We have one paragraph in the related work sections dedicated  
41 to the empirical work and mentioned algorithms such as Go-Explore and Hindsight Experience Replay (HER), but we  
42 are happy to add the missing ones you mentioned to the paragraph. **Dependency on  $H$  and  $N$ :** The dependency on  $H$   
43 can potentially be improved upto  $H^2$ , matching the lower bound. See also response 2 to reviewer 3. The dependency  
44 on  $N$ , however, is shown to be unavoidable by our lower-bound. This is due to our more realistic assumption that  
45 only instantiations of reward are available rather than full reward function assumed in [Jin et. al. 2020]. In most  
46 empirical problems, however,  $N$  is small. For example, if one wants to learn to navigate to all states,  $N = S$ , and  
47  $\log S$  is negligible compared to the other  $S$  term. **Notation in Alg 1:** Sorry for the confusion. So both  $b_t$  and  $\alpha_t$  means  
48 to be functions of  $t$ .  $t = ++ N(s, a)$  is used immediately in the next lines as inputs to  $b_t$  and  $\alpha_t$ . We will make it  
49 clearer and add the suggested comment lines in the revision!  **$N$  dependent v.s.  $N$  independent:** Sorry again for the  
50 confusion. All tasks are always independent of each other. Here "dependent" means whether the sample complexity  
51 depends on the number of tasks  $N$ , and " $N$  independent" means that the bound doesn't scale with  $N$  and therefore still  
52 holds finite even when the number of tasks goes to infinity. Thank you for pointing out the many typos and for the  
53 polishing suggestions! We will make sure to follow them in the revision.

54 **Response to Reviewer 5:** Thank you for the very positive feedback on our work! We agree that the algorithm is mainly  
55 of theoretical interest and the current gap between the empirical success of deep RL and the theoretical understanding  
56 is still large. We hope that our paper can provide some intuition to the (fast-growing) empirical task-agnostic RL  
57 community!