

1 We thank all reviewers for their comments. Minor comments will be addressed in the final version.

## 2 Reviewer 1

3 **Clear description of the setting** We want to emphasize first that our problem setting is a standard restless bandit  
4 setting with a few specific choices.  $P_k^{\text{active}}$  and  $P_k^{\text{passive}}$  are the transition matrices of the arm  $k$  when it is pulled or not,  
5 respectively.  $X_t$  is a  $K$  dimensional binary vector such that the  $k^{\text{th}}$  component  $X_{t,k}$  represents the reward of the arm  
6  $k$ . Since the learner only observes the rewards of pulled arms, only the  $N$  components  $X_{t,A_t}$  will be available to the  
7 learner. These notions are defined in lines 49 - 58. We will make our description clearer in the final version.

8 **Messages of the experiments** The first experiment empirically checks the Bayesian regret of our algorithm is indeed  
9  $\tilde{O}(\sqrt{T})$ . The second experiment shows the algorithm still works in the frequentist setting. Figure 3 (left) illustrates  
10 how the value function of the policy  $\pi_l$  chosen in an episode converges to the baseline value for a variety of competitor  
11 mappings (the best fixed action, the myopic policy, and the Whittle index policy). Figure 3 (right) shows the posterior  
12 weights on the true parameters monotonically increase. We will describe the details of our experiments more carefully  
13 and make figures more readable.

## 14 Reviewer 2

15 **1. Motivating application of the episodic setting** Yes, the assumption of periodic restart of the system is somewhat  
16 limiting, and the regret analysis in the infinite time horizon is an interesting open question. Analyzing the finite  
17 horizon case should be an intermediate step towards solving this question. Moreover, the episodic case itself has a few  
18 motivating applications. For example, in the dynamic channel access problem that we consider in our experiment, the  
19 channel provider might reset their system every night when network traffic is low for maintenance related reasons. After  
20 the reset, every channel should be available for use, which can be thought as the beginning of a new episode.

21 **2. Super-time-instant** It is indeed possible to tackle the problem by considering each deterministic policy as an arm.  
22 However, this would result in very large (possibly infinite)  $K$ , the number of arms, and the existing bounds become  
23 vacuous as they depend (polynomially) on  $K$ . The bound in Dai et al. [2011] is meaningful since there are only two  
24 competing policies. This perspective still conveys interesting points, and we will add the comparison in the final version.

25 **3. More complete picture in intro** We totally agree that existing results, including ours, are just limited in different  
26 aspects. We will clarify this point more clearly. Nevertheless, we want to point out that this is the *first* paper that  
27 analyzes Thompson sampling in *multi-armed* restless bandit problems.

28 **4. The optimal policy depends on  $L$**  Yes, the optimal policy will depend on the episode length. It will change the  
29 baseline value in our regret definition in (2), but the same regret bound will still apply. It is one of our main contributions  
30 that the regret bound applies regardless of the choice of the benchmark.

## 31 Reviewer 3

32 **Finer analysis within each episode** First of all, the point raised by the reviewer is completely true. The episode length  
33  $L$  should remain small to make the regret bound meaningful. We mainly considered the case where  $L$  is fixed as a  
34 constant and the number of resets,  $m$ , increases arbitrarily so that the posterior distribution concentrates sufficiently  
35 around the truth. A fundamental reason why we did not do the finer analysis within the episode is because our algorithm  
36 fixes a policy  $\pi_l$  and runs it throughout the episode  $l$ . If we get to do the finer analysis, then that means our algorithm  
37 changes the policy more often, which comes with an extra cost. For example, in the regret analysis by Ouyang et  
38 al. [“Learning Unknown Markov Decision Processes: A Thompson Sampling Approach,” NIPS 2017], who analyze  
39 Thompson sampling in non-episodic fully observable MDPs, the bound includes  $K_T$ , the number of different policies  
40 that Thompson sampling runs up to time  $T$ .

41 **Tightness of the regret bound** As pointed out in the remark right after Theorem 5, our result reproduces the regret  
42 bound of  $O(\sqrt{KT \log T})$  in the stationary MAB problem, whose lower bound is shown to be  $\Omega(\sqrt{KT})$ . This suggests  
43 that our bound is optimal in  $K$  and  $T$  up to a logarithmic factor. When  $L = 1$ , the problem becomes a combinatorial  
44 bandit problem (of choosing a set of  $N$  active arms out of a total of  $K$ ) in which case the best known regret bound is  
45  $O(\sqrt{KN^3 T \log K})$  (e.g., see “Combinatorial bandits” by Cesa-Bianchi and Lugosi [2012]. Their bound is actually  
46  $O(\sqrt{KNT \log K})$ , but they normalize the loss to be in  $[0, 1]$ , whereas our reward is in  $[0, N]$ ). Our bound agrees with  
47 their bound up to logarithmic terms. Finally, the optimal dependence on  $L$  remains open. We will add the discussion of  
48 tight dependence in the final version.