

1 We appreciate the careful reviews!

2 **Reviewer 1:** (1) “Can this framework get rid of the knowledge of T ?” Some information about T is necessary, since
3 the time horizon is an important ingredient in optimally balancing exploration and exploitation. That said, precise
4 knowledge of T is not necessary. The framework (penalties, policies, and upper bounds) can naturally incorporate
5 the unknown T within the *Bayesian setting*: i.e., the horizon T is also a random variable whose prior distribution is
6 known. As a simple case, if T is independent of the DM’s actions, we can reformulate the objective function of the inner
7 problem as $\sum_{t=1}^{\infty} \gamma_t (r_t(\mathbf{a}_{1:t}; \omega) - z_t(\mathbf{a}_{1:t}; \omega))$ where the discount factor $\gamma_t \triangleq \mathbb{P}[T \geq t]$ is the survivor probability, and
8 $r_t(\cdot)$ and $z_t(\cdot)$ are the reward and penalty terms used in the paper. Alternatively, we can treat the random variable T like
9 the random reward realizations – sample T from its prior distribution while a penalty function (additionally) penalizes
10 for the gain from knowing T (you can imagine that the outcome ω now includes the realization of T). Structural results
11 such as weak duality and strong duality will continue to hold.

12 (2) “In most cases, IRS has computational complexity that is linear (or even polynomial) in T for each round.” To
13 be fair, the computational complexity of IRS.FH is independent of T (like TS). Moreover, we can construct a dual
14 feasible penalty function that mixes IRS.FH and IRS.V-ZERO, which induces an algorithm whose complexity is
15 $O(K \min\{T, T_0\}^2)$ for some predefined constant T_0 (in the inner problem, IRS.V-ZERO-like penalties are applied
16 for the initial $\lfloor T_0/K \rfloor$ pulls and then IRS.FH-like penalties are applied for the later pulls). Such a practical variant
17 has performance that does not scale with T beyond T_0 . In our experiments, this variant works well even for moderate
18 values of T_0 . Other heuristic variations of IRS policies with bounded T dependence are also possible.

19 (3) “Any theoretical results for IRS.V-EMAX?” We have the additional result that $W^{\text{IRS.V-EMAX}}(T, \mathbf{y}) \leq$
20 $W^{\text{TS}}(T, \mathbf{y})$, this could be added to Theorem 2. We don’t have a suboptimality analysis for IRS.V-EMAX yet.

21 **Reviewer 2: Comparison with information-directed sampling (IDS).** It is also remarkable to us that IDS performs
22 well without the knowledge of T . The IRS policies outperform when the finite horizon creates a considerable tension in
23 the exploitation-exploration trade-off. As illustrated in the numerical experiments, this would be the case when the
24 arms are dissimilar and the time horizon is short relative to the number of arms.

25 We believe that our algorithms have other advantages over IDS. First, IDS requires a significant amount of work
26 that is specific to the application; e.g., it requires to compute the expected change in entropy, which is typically
27 obtained by a numerical integration after doing some distribution-specific reformulation. For example, if some arms
28 yield normally-distributed rewards and the others yield Bernoulli-distributed rewards, implementing IDS will be very
29 challenging. Moreover, IDS’s computational complexity is $O(K^2)$ per decision, whereas IRS.FH and IRS.V-ZERO are
30 linear in K . Finally, this framework can naturally deal with other constraints apart from the time-horizon one; see the
31 answer for Reviewer 3 below.

32 **Reviewer 3: Additional discussion on our contributions.** Respectfully, we believe that our generalization of TS
33 to finite-horizon problems is novel and has not appeared in the literature. A common heuristic for the finite-horizon
34 setting would be *posterior reshaping*, mentioned in Chappelle and Li (2011), which reduces the variance of the posterior
35 distribution with an ad hoc parameter. Another relevant work is Russo, Tse and Van Roy (2017), in which the authors
36 propose *satisficing Thompson sampling* for the discounted infinite-horizon setting, which also introduces an auxiliary
37 parameter to control the degree of exploitation explicitly. In contrast to these heuristic proposals, this paper provides a
38 principled method that does not require any tuning or additional parameters, and suggests a unified framework that
39 includes TS and the Bayesian optimal policy as special cases. Also note that the decision making procedure of every
40 IRS policy is recursive like TS: i.e., the decision at a certain moment depends only on the posterior distribution and the
41 remaining horizon at that moment.

42 We absolutely agree with the fact that the stochastic MAB with independent arms has already been studied extensively.
43 That said, to the extent that this problem is practically interesting, we provide methods that are competitive with
44 commonly employed solution methods such as TS. Moreover, even though this paper focuses on this simplest setting,
45 our framework applies for more complicated settings. Consider the following examples: (a) Correlated arms in a
46 finite-horizon setting (e.g., $R_{a,n} \sim \mathcal{N}(\mathbf{x}_a^\top \boldsymbol{\theta}, \sigma^2)$ where $\boldsymbol{\theta}$ is shared across the arms): IRS.V-ZERO can be immediately
47 implemented by adopting the DP algorithm discussed in §B.2. (b) MAB with the delayed reward realization: IRS.FH
48 can be immediately implemented by simulating the DM’s learning process in the presence of delay. (c) MAB with a
49 budget constraint (in which each arm consumes a certain amount of budget and the DM wants to maximize the total
50 reward within a limited budget): all IRS algorithms can be implemented by solving a budget-constrained optimization
51 problem instead of a horizon-constrained optimization problem. In these extensions, we obtain not only the online
52 decision making policies but also their performance bounds as in this paper. Generally speaking, our framework
53 provides a systemic way of improving TS by taking into account the exploitation-exploration trade-off more carefully,
54 particularly in the presence of some constraint that incurs incomplete learning. We believe that this feature is novel in
55 the literature and also very crucial in practice.