

---

# Sample-Mean Anchored Thompson Sampling for Offline-to-Online Learning with Distribution Shift

---

Anonymous Authors<sup>1</sup>

## Abstract

Offline-to-online learning aims to improve online decision-making by leveraging offline logged data. A central challenge in this setting is the distribution shift between offline and online environments. While some existing works attempt to leverage shifted offline data, they largely rely on UCB-type algorithms. Thompson sampling (TS) represents another canonical class of bandit algorithms, well known for its strong empirical performance and naturally suited to offline-to-online learning through its Bayesian formulation. However, unlike UCB indices, posterior samples in TS are not guaranteed to be optimistic with respect to the true arm means. This makes indices constructed from purely online and hybrid data difficult to compare and complicates their use. To address this issue, we propose sample-mean anchored TS (Anchor-TS), which introduces a novel median-based anchoring rule that defines the arm index as the median of an online posterior sample, a hybrid posterior sample, and the online sample mean. The median anchoring systematically corrects bias induced by distribution shift by mitigating over-estimation for suboptimal arms and under-estimation for optimal arms, while exploiting offline information to obtain more accurate estimates when the shift is small. We establish theoretical guarantees showing that the proposed algorithm safely leverages offline data to accelerate online learning, and quantifying how the degree of distribution shift and the size of offline data affect the resulting regret reduction. Extensive experiments demonstrate consistent improvements of our algorithm over baselines.

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

## 1. Introduction

Stochastic multi-armed bandits (MAB) provide a fundamental framework for sequential decision-making under uncertainty (Lai & Robbins, 1985; Lattimore & Szepesvári, 2020), where a learning agent repeatedly selects actions to learn the unknown environment and maximize cumulative reward. The upper confidence bound (UCB) (Auer et al., 2002) and Thompson sampling (TS) (Agrawal & Goyal, 2017) are two canonical algorithm types to solving the problem: the former relies on optimism-driven confidence intervals, and the latter selects actions by sampling from posterior distributions. Despite their importance, classical bandit algorithms typically start from scratch. However, learning purely from online interactions can be costly or risky. Many applications provide access to offline logged data collected by historical policies, which motivates the offline-to-online learning setting that leverages such data to accelerate subsequent online learning (Nair et al., 2020; Lee et al., 2022; Shivaswamy & Joachims, 2012; Wagenmaker & Pacchiano, 2023).

Recently, there has been increasing interest in offline-to-online algorithms under distribution shift (Cheung & Lyu, 2024; Yin & Fang, 2025; He et al., 2025; Qu et al., 2025), where the offline reward distribution differs from the online environment due to system updates, non-stationarity, or sim-to-real gaps. A recurring insight in this line of work is that achieving performance no worse than purely online learning requires some form of prior knowledge about the distribution shift. Existing approaches mainly adopt UCB-type algorithms. A representative strategy is to construct two UCBs: a hybrid UCB that incorporates offline data together with a bias correction term, and a purely online UCB based solely on online observations (Cheung & Lyu, 2024). Actions are selected according to the minimum of the two, which guarantees regret no worse than that of pure online UCB.

Thompson sampling (TS) represents another canonical class of bandit algorithms and has been extensively studied across a wide range of online learning settings (Agrawal & Goyal, 2012; Daniel et al., 2018). It is well known for its strong empirical performance, often outperforming UCB-based methods in practice, and its Bayesian formulation makes it particularly appealing for offline-to-online learning as

055 unbiased offline data can be naturally incorporated into the  
 056 prior distribution (Oetomo et al., 2023; Agnihotri et al.,  
 057 2024).

058 Despite these advantages, the posterior-sampling nature of  
 059 TS makes it fundamentally more challenging under distribu-  
 060 tion shift. In UCB-based methods, the index is optimistic,  
 061 and regret can be attributed to inaccurate estimates of subop-  
 062 timal arms. By contrast, TS relies on posterior samples that  
 063 may fall on either side of the true arm means, and its regret  
 064 depends jointly on the accuracy of both the optimal and  
 065 suboptimal arms. This principal difference has important  
 066 implications for incorporating offline data. For UCB-based  
 067 methods, conservatively taking the minimum of a purely on-  
 068 line index and a hybrid index provides a safe guard against  
 069 distribution shift. For TS, however, no analogous compari-  
 070 son rule exists: taking the minimum risks underestimating  
 071 the optimal arm, while taking the maximum may overesti-  
 072 mate suboptimal arms. This lack of a principled mechanism  
 073 for comparing posterior samples makes it nontrivial to ex-  
 074 ploit biased offline data for TS.  
 075

076 To address this challenge, we propose sample-mean an-  
 077 chored Thompson sampling (Anchor-TS), a modified TS  
 078 algorithm based on a median-of-indices aggregation scheme.  
 079 At each round, Anchor-TS assigns to each arm an index de-  
 080 fined as the median of three statistics: (i) an online posterior  
 081 sample, (ii) a hybrid posterior sample that integrates offline  
 082 and online data, and (iii) the online sample mean, which  
 083 serves as an unbiased anchor. This median-based aggrega-  
 084 tion induces an inherent robustness–efficiency trade-off.  
 085 When the offline data are reliable, the hybrid posterior sam-  
 086 ple typically exhibits reduced variance and remains close  
 087 to the anchor, thereby accelerating learning. Conversely,  
 088 under severe bias in the offline component, the hybrid sam-  
 089 ple tends to behave as an outlier and is filtered out by the  
 090 median, so that the resulting index is primarily driven by  
 091 online observations. Furthermore, to prevent the underesti-  
 092 mation of the optimal arms’ online sample mean from being  
 093 amplified under median aggregation when its offline reward  
 094 is smaller, we additionally introduce a right-hand shift to  
 095 the hybrid posterior distribution.

096 Our analysis develops a new regret decomposition tailored  
 097 to median-based TS. In particular, we analyze the behav-  
 098 ior of posterior samples for both suboptimal and optimal  
 099 arms relative to the online sample mean, which serves as an  
 100 anchor, and characterize their joint contributions to regret  
 101 accordingly. This enables us to bound the leading regret  
 102 terms in a way that captures the more favorable behavior  
 103 between online TS and hybrid TS. Specifically, we obtain a  
 104  
 105  
 106  
 107  
 108  
 109

regret upper bound

$$O\left(\sum_{i \in [K] \setminus \{1\}} \Delta_i \left( (\log T / \Delta_i^2 - N_i \max\{0, 1 - 3\omega_i / \Delta_i\})_+ + (\log T / \Delta_i^2 - N_1)_+ \right)\right).$$

where  $T$  is online horizon,  $K$  is the arm set size, arm 1 is the optimal arm,  $\Delta_i$  is arm  $i$ ’s sub-optimality gap compared with 1,  $\omega_i$  is an effective discrepancy related to  $i$ ’s offline and online expected reward,  $N_i$  is the offline sample size of  $i$ ,  $(\cdot)_+$  represents  $\max\{0, \cdot\}$ .

The above regret guarantee formalizes the robustness–efficiency trade-off under distribution shift: it is no worse than the pure online TS (Agrawal & Goyal, 2017), and strictly improves when the distribution shift  $\{\omega_i\}_{i \in [K]}$  is mild. Compared to existing UCB-based algorithms (Cheung & Lyu, 2024), the regret reduction attributable to suboptimal arms estimations is of similar order. Crucially, our analysis further reveals an additional source of improvement that is unique to TS: offline data also reduces regret by accelerating concentration on the optimal arm, as reflected in the dependence on  $N_1$ . Such an advantage is particularly relevant in practice, where offline data are often collected by prior learning policies or expert behavior and therefore tend to contain a large number of samples from the optimal arm. In such cases, Anchor-TS can effectively exploit this abundance of optimal-arm data, whereas UCB-based methods are unable to benefit in the same way. Extensive empirical results further validate our theory. Across all considered settings, Anchor-TS consistently and substantially outperforms UCB-based algorithms. The performance gap becomes even more pronounced when more offline data are on the optimal arm. Importantly, even in the unbiased and pure online settings, the sample-mean–based mechanism helps reduce the excessive exploration caused by the high variance of a single posterior sample in vanilla TS. As a result, it achieves better empirical performance than the corresponding TS baselines in these settings.

## 2. Related Work

**Thompson sampling (TS).** TS is a canonical Bayesian framework for stochastic MAB. It has been widely adopted since the seminal work of (Thompson, 1933), but the establishment of convergent regrets lagged for decades (Agrawal & Goyal, 2012; Kaufmann et al., 2012; Agrawal & Goyal, 2017; Jin et al., 2021; 2023). In many problem settings, although UCB-based algorithms have been applied, researchers continue to devote significant effort to TS-type algorithms (Agrawal & Goyal, 2013; Komiyama et al., 2015; Wang & Chen, 2018; Verstraeten et al., 2020; Zhang et al., 2021), driven by their superior empirical performance

and ease of implementation (Granmo, 2010; Scott, 2010; Chapelle & Li, 2011). More recently, TS-type algorithms are also proposed for offline-to-online learning problems with offline data naturally encoded in the prior (Oetomo et al., 2023; Agnihotri et al., 2024). These approaches typically do not account for distribution shift between offline and online environments.

**Offline-to-online Learning.** Offline-to-online learning aims to accelerate online decision-making by leveraging pre-collected logged data. This paradigm has been extensively explored empirically in both bandit and reinforcement learning (RL) settings (Lee et al., 2022; Nair et al., 2020; Ball et al., 2023; Yu & Zhang, 2023; Nakamoto et al., 2023; Xia et al., 2024), and has also motivated a growing body of theoretical work establishing performance guarantees.

Early theoretical studies primarily focused on the unbiased setting, where offline and online reward distributions coincide. In the bandit setting, Shivaswamy & Joachims (2012) show that historical data can yield constant regret, a result later extended to contextual bandits through informative priors (Oetomo et al., 2023) and meta-algorithmic approaches (Banerjee et al., 2022). More recently, Sentenac et al. (2025) analyze the trade-off between optimism in online learning and pessimism in offline learning in the offline-to-online context. For RL, theoretical frameworks have expanded from tabular Markov decision processes (Xie et al., 2021) to linear models (Wagenmaker & Pacchiano, 2023; Tan & Xu, 2024; Huang et al., 2025) and further to general function approximation (Song et al., 2023), with a common focus on relaxing offline data coverage assumptions.

In practice, offline data often originate from heterogeneous or outdated sources, making it essential to achieve both efficiency gains from informative offline data and robustness to distribution shift. Fundamental lower bounds establish that without a priori knowledge of the distribution shift, no policy can uniformly outperform purely online algorithms (Zhang et al., 2019; Cheung & Lyu, 2024; Zhang et al., 2025; Qu et al., 2025). Consequently, existing works typically incorporate some shift information into the decision process. Such information is encoded through lower bounds (Qu et al., 2025) or upper bounds on the shift (Cheung & Lyu, 2024; Ahn et al., 2025). Our work follows the last line of works. A representative strategy in this line is to construct a hybrid UCB index by combining offline data with a known upper bound on the bias, and to select actions according to the minimum of this hybrid index and a purely online UCB (Cheung & Lyu, 2024). This strategy is applied to a variety of settings, including best arm identification (BAI) (Yang et al., 2025), heterogeneous feedback structures (He et al., 2025), auxiliary rewards (Yin & Fang, 2025), combinatorial MAB (Zhou et al., 2025), and linear bandits (Zhang et al., 2025). To the best of our knowledge, TS-type algorithms

for offline-to-online learning under distribution shift is still open.

**TS with misspecified priors.** Our work is also related to studies on misspecified priors in TS. Liu & Li (2016); Simchowitz et al. (2021) show that when TS is initialized with strongly biased priors, where the prior assigns very little probability to the true model, it can suffer linear regret and perform much worse than methods with uninformative priors. These results demonstrate the negative impact of using an incorrect prior. Compared with this line of work, our setting explicitly models the prior as being induced by offline data. Our setting requires not only robustness to biased priors, but also the ability to achieve regret improvement when offline data are informative.

### 3. Preliminaries

We begin with the classical stochastic multi-armed bandit problem. An agent interacts with an environment over a time horizon  $T$  by repeatedly selecting arms from a finite set  $[K] = \{1, \dots, K\}$ . Each arm  $i \in [K]$  is associated with an unknown reward distribution  $P_i^{(\text{on})}$  supported on  $[0, 1]^1$  and unknown mean  $\mu_i^{(\text{on})}$ . At each round  $t = 1, 2, \dots, T$ , the agent selects an arm  $A(t) \in [K]$  and observes a reward  $R_{A(t)}(t)$  drawn independently from the corresponding online distribution  $P_{A(t)}^{(\text{on})}$ . Without loss of generality, let arm  $1 \in \arg \max_{j \in [K]} \mu_j^{(\text{on})}$  denote the unique<sup>2</sup> optimal arm. For each suboptimal arm  $i \neq 1$ , define the sub-optimality gap  $\Delta_i = \mu_1^{(\text{on})} - \mu_i^{(\text{on})}$ .

In this work, we study a stochastic bandit setting augmented with *offline data*. In addition to the online interaction described above, each arm  $i \in [K]$  is associated with an offline dataset  $S_i = \{X_{i,1}, \dots, X_{i,N_i}\}$  of size  $N_i$ . The samples in  $S_i$  are assumed to be i.i.d. draws from an offline distribution  $P_i^{(\text{off})}$  supported on  $[0, 1]$ , with unknown mean  $\mu_i^{(\text{off})} = \mathbb{E}_{X_i \sim P_i^{(\text{off})}} [X_i]$ .

Importantly, for any arm  $i$ , the offline and online reward distributions are not assumed to coincide. To explicitly model the potential distribution shift between them, we assume that for each arm  $i$ , there exists a known upper bound  $V_i \geq 0$  such that  $|\mu_i^{(\text{off})} - \mu_i^{(\text{on})}| \leq V_i$ . This assumption provides a priori control over the magnitude of the distribution shift and is information-theoretically necessary. Existing works establish fundamental lower bounds that no policy can uni-

<sup>1</sup>Our analysis extends to sub-Gaussian reward distributions by replacing Hoeffding-type concentration bounds with their sub-Gaussian counterparts. The regret order remains unchanged.

<sup>2</sup>The uniqueness assumption is only for the convenience of the analysis. The setting with multiple optimal arms can only decrease the regret as shown in (Agrawal & Goyal, 2012).

formly outperform purely online algorithms without any prior knowledge (Zhang et al., 2019; Cheung & Lyu, 2024; Zhang et al., 2025; Qu et al., 2025) and such upper bound is widely adopted in previous works (Cheung & Lyu, 2024; Zhang et al., 2025; He et al., 2025; Ahn et al., 2025; Yin & Fang, 2025).

The performance of a policy  $\pi$  is measured by its cumulative regret, which quantifies the expected loss relative to always pulling the optimal arm:

$$\begin{aligned} \text{Reg}_\pi(T) &= T\mu_1^{(\text{on})} - \mathbb{E} \left[ \sum_{t=1}^T R_{A(t)}(t) \right] \\ &= \sum_{i \neq 1} \Delta_i \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{A(t) = i\} \right]. \end{aligned} \quad (1)$$

where the expectation is taken over all sources of randomness, including the internal randomness of  $\pi$ , the offline data draw from  $P^{(\text{off})}$ , and the online rewards generated under  $P^{(\text{on})}$ .

**Useful Notations.** For each arm  $i \in [K]$ , let  $T_i(t)$  denote the number of times arm  $i$  is pulled during rounds  $1, \dots, t-1$ , and define  $n_i(t) = T_i(t) + N_i$  as the total number of samples available for arm  $i$ , combining online and offline data. Let  $\hat{\mu}_i^{(\text{off})}$  denote the sample mean of the offline dataset for arm  $i$ , and let  $\hat{\mu}_i^{(\text{on})}(t)$  be the sample mean computed from the online rewards observed up to round  $t-1$ . We further define  $\hat{\mu}_i^{(\text{hyb})}(t)$  as the hybrid sample mean that aggregates both offline samples and online observations for arm  $i$ .

## 4. Algorithm

In the offline-to-online setting, it is natural to maintain two types of estimators: a purely online estimator constructed solely from online observations, and a hybrid estimator that aggregates offline data with online observations. Due to the distribution shift, the hybrid estimators may be biased and can mislead the learning process if used indiscriminately. The key algorithmic challenge is therefore how to balance the use of the hybrid estimator against the purely online estimator. In this section, we present our sample-mean anchored TS (Anchor-TS, Algorithm 1), an efficient and robust mechanism for combining online and hybrid estimators under the TS framework.

For each arm  $i \in [K]$ , the algorithm maintains an online posterior that is updated exclusively based on rewards observed during online interaction, which is the same as the traditional TS algorithm for the purely online setting (Agrawal & Goyal, 2017). At time  $t$ , this posterior is characterized by the online sample mean  $\hat{\mu}_i^{(\text{on})}(t)$  and an associated variance  $\sigma_{i,\text{on}}^2(t)$  (Line 1). An online posterior sample  $\theta_i^{(\text{on})}(t)$  is

**Algorithm 1** Sample-Mean Anchored Thompson sampling (Anchor-TS)

---

- 1: **Input:** Arm set  $[K]$ ; offline sample mean  $\hat{\mu}_i^{(\text{off})}$  and sample size  $N_i$ , bias bound  $V_i, \forall i \in [K]$
- 2: **Initialization:**  $T_i(1) \leftarrow 0, \hat{\mu}_i^{(\text{on})}(1) \leftarrow 0, \hat{\sigma}_{i,\text{on}}^2(1) \leftarrow 1, \hat{\mu}_i^{(\text{hyb})}(1) \leftarrow \hat{\mu}_i^{(\text{off})}, \hat{\sigma}_{i,\text{hyb}}^2(1) \leftarrow 1/(N_i + 1), Z_{i,1} \leftarrow V_i, \forall i \in [K]$
- 3: **for**  $t = 1, \dots$  **do**
- 4:   **for** each arm  $i \in [K]$  **do**
- 5:     Sample  $\theta_i^{(\text{on})}(t)$  from  $\mathcal{N}(\hat{\mu}_i^{(\text{on})}(t), \hat{\sigma}_{i,\text{on}}^2(t))$
- 6:     Sample  $\theta_i^{(\text{hyb})}(t)$  from  $\mathcal{N}(\hat{\mu}_i^{(\text{hyb})}(t) + Z_{i,t}, \hat{\sigma}_{i,\text{hyb}}^2(t))$
- 7:      $\hat{\theta}_i(t) \leftarrow \text{median}\{\hat{\mu}_i^{(\text{on})}(t), \theta_i^{(\text{on})}(t), \theta_i^{(\text{hyb})}(t)\}$
- 8:   **end for**
- 9:   Select arm  $A(t) \leftarrow \arg \max_{i \in [K]} \hat{\theta}_i(t)$  and observe reward  $R_{A(t)}(t)$
- 10:   **// Update online posterior of**  $A(t)$
- 11:    $T_{A(t)}(t+1) \leftarrow T_{A(t)}(t) + 1$
- 12:    $\hat{\mu}_{A(t)}^{(\text{on})}(t+1) \leftarrow \frac{T_{A(t)}(t) \cdot \hat{\mu}_{A(t)}^{(\text{on})}(t) + R_{A(t)}(t)}{T_{A(t)}(t+1) + 1},$   
 $\hat{\sigma}_{A(t),\text{on}}^2(t+1) \leftarrow \frac{1}{T_{A(t)}(t+1) + 1}$
- 13:   **// Update hybrid posterior of**  $A(t)$
- 14:    $\hat{\mu}_{A(t)}^{(\text{hyb})}(t+1) \leftarrow \frac{T_{A(t)}(t+1) \cdot \hat{\mu}_{A(t)}^{(\text{on})}(t+1) + N_{A(t)} \cdot \hat{\mu}_{A(t)}^{(\text{off})}}{T_{A(t)}(t+1) + N_{A(t)} + 1}$
- 15:    $\hat{\sigma}_{A(t),\text{hyb}}^2(t+1) \leftarrow \frac{1}{N_{A(t)} + T_{A(t)}(t+1) + 1},$   
 $Z_{A(t),t+1} \leftarrow \frac{N_{A(t)} V_{A(t)}}{T_{A(t)}(t+1) + N_{A(t)}}$
- 16:   **for**  $i \neq A(t)$  **do**
- 17:     Update  $\hat{\mu}_i^{(\text{on})}(t+1) \leftarrow \hat{\mu}_i^{(\text{on})}(t), \hat{\mu}_i^{(\text{hyb})}(t+1) \leftarrow \hat{\mu}_i^{(\text{hyb})}(t), \hat{\sigma}_{i,\text{on}}^2(t+1) \leftarrow \hat{\sigma}_{i,\text{on}}^2(t), \hat{\sigma}_{i,\text{hyb}}^2(t+1) \leftarrow \hat{\sigma}_{i,\text{hyb}}^2(t), Z_{i,t+1} \leftarrow Z_{i,t}$
- 18:   **end for**
- 19: **end for**

---

drawn from the corresponding Gaussian distribution (Line 1). Because it relies solely on online data, this posterior remains unbiased with respect to the true arm mean, although it may exhibit high variance in the early stages of learning.

In parallel, Anchor-TS constructs a hybrid posterior that integrates offline empirical information with online observations. Specifically, the hybrid posterior mean  $\hat{\mu}_i^{(\text{hyb})}(t)$  is obtained by combining the offline sample mean  $\hat{\mu}_i^{(\text{off})}$  and the current online sample mean  $\hat{\mu}_i^{(\text{on})}(t)$ , weighted according to their respective sample sizes (Line 1). The corresponding variance  $\hat{\sigma}_{i,\text{hyb}}^2(t)$  reflects the increased effective sample size enabled by incorporating offline data (Line 1). Notably, to encourage sufficient exploration of the optimal arm 1 and reduce regret, we apply a rightward shift of magnitude  $Z_{i,t}$  to the hybrid posterior distribution of each arm.

A detailed discussion of the intuition behind this design choice will be provided in the end of this section. A hybrid posterior sample  $\theta_i^{(\text{hyb})}(t)$  is then drawn from the resulting Gaussian distribution (Line 1).

At each time step, Anchor-TS computes three indices for every arm: the online sample mean  $\hat{\mu}_i^{(\text{on})}(t)$ , an online posterior sample  $\theta_i^{(\text{on})}(t)$ , and a hybrid posterior sample  $\theta_i^{(\text{hyb})}(t)$ . The arm score is obtained by taking the *median* of these three quantities (Line 1). And Anchor-TS selects the arm with the highest score in each round (Line 1).

### Intuition behind the median aggregation of three indices.

Intuitively, the role of median aggregation is to adaptively select the index that is closest to the true mean rather than to favor large or small posterior realizations like UCB-type algorithms (Cheung & Lyu, 2024). This distinction is intrinsic to TS: posterior samples in TS are neither guaranteed to be optimistic nor pessimistic. Naively selecting the minimum or maximum of multiple indices would either underestimate the optimal arm or overestimating suboptimal arms, leading to uncontrolled regret. The sample mean  $\hat{\mu}^{(\text{on})}$ , estimated from online observations, therefore serves as a natural stabilizing anchor as it concentrates rapidly around the true arm mean. By taking the median of  $\hat{\mu}^{(\text{on})}$ , an online-only and a hybrid posterior sample, Anchor-TS selects the index that is most consistent with this anchor.

When the bias is small, the hybrid posterior, benefiting from a larger effective sample size, tends to concentrate more tightly around the true mean than the online-only posterior. In this regime, the median naturally favors the hybrid posterior sample, allowing Anchor-TS to exploit offline data for faster learning. Conversely, when the offline bias is large, the hybrid posterior may deviate from the true mean, while the online-only posterior remains centered around the unbiased online signal. In this case, the median suppresses the biased hybrid sample and tends to select the online sample instead, yielding robustness against offline-induced distortion. Through this adaptive selection mechanism, Anchor-TS automatically interpolates between efficiency and robustness, leveraging offline data when it is reliable and reverting to online evidence when offline bias is substantial.

**Intuition behind the right-hand shift on the hybrid posterior distribution.** Recall that for TS-type algorithms, the regret can still accumulate when the estimates of suboptimal arms are accurate, as long as the optimal arm is poorly estimated. Specific to the offline-to-online setting, if the offline data underestimate the optimal arm, i.e.,  $\mu_1^{(\text{off})} < \mu_1^{(\text{on})}$ , and the online sample mean of arm 1 happens to be pessimistic, arm 1 may receive too little posterior probability due to the three-index voting mechanism. In this case, it is difficult to correct the estimation error of arm 1 through

additional online observations. As a result, the persistent underestimation of arm 1 amplifies the number of selections of suboptimal arms, leading to increased regret despite their accurate estimation.

Our right-shift operation counteracts this effect by preventing the hybrid posterior distribution from underestimating the online mean of arm 1. This guarantees adequate exploration of the optimal arm and prevents regret from being dominated by prolonged under-exploration of the optimal arm.

Since the identity of the optimal arm is unknown, the right-hand shift is applied uniformly across arms. Importantly, this shift does not compromise the benefit of offline data for suboptimal arms: although overestimation by the hybrid posterior may initially increase their selection frequency, the influence of offline information quickly diminishes as online observations accumulate, and the estimate becomes dominated by online samples through median aggregation.

## 5. Theoretical Results and Discussions

The following theorem summarizes the regret bound of our algorithm:

**Theorem 5.1.** *The cumulative regret of Algorithm 1 can be bounded by*

$$\text{Reg}(T) \leq O\left(\sum_{i \neq 1} \Delta_i \left( \left( \frac{C_1 \log T}{\Delta_i^2} - N_i \left(1 - \frac{3\omega_i}{\Delta_i}\right)_+ \right) + \left( \frac{C_2 \log T}{\Delta_i^2} - N_i \right)_+ + \frac{C_3}{\Delta_i^2} \right)\right).$$

where  $C_1, C_2, C_3$  are constants,  $(\cdot)_+$  represents  $\max\{\cdot, 0\}$ ,  $\omega_i := V_i + \mu_i^{(\text{off})} - \mu_i^{(\text{on})}$  represents the effective discrepancy by adding  $V_i$  in the hybrid posterior distribution.

Due to the space limit, the complete proof of Theorem 5.1 is deferred to Appendix A. In the following, we first provide some discussions and then show the proof sketch.

**Intuition of the regret upper bound.** This upper bound consists of three terms. The first term arises from inaccurate estimation of suboptimal arms, the second from inaccurate estimation of the optimal arm, and the third is a constant term required for convergence of the online sample mean.

This upper bound can be interpreted as the regret in a purely online setting minus the benefit provided by the offline data. When no offline data are available, the regret reduces to the standard purely online regret (Agrawal & Goyal, 2017). When the distribution shift  $\mu_i^{(\text{on})} - \mu_i^{(\text{off})}$  is zero, the offline samples for both the sub-optimal arms and the optimal arm can be viewed as offsetting a portion of the regret. In particular, if the sub-optimal arms have sufficiently many offline

275 samples, their means can be accurately estimated from the  
 276 outset, rendering the corresponding regret term constant.  
 277 Similarly, if the optimal arm has sufficiently many offline  
 278 samples, the regret arising from inaccurate estimation of the  
 279 optimal arm is likewise reduced to a constant.

280 Regarding the bias, for a sub-optimal arm  $i$ , if the offline  
 281 mean is smaller than the online mean and  $V$  is a tight bound,  
 282 i.e.,  $V = \mu_i^{(\text{on})} - \mu_i^{(\text{off})} > 0$ , then  $\omega_i = 0$  and the of-  
 283 fline improvement scales linearly with  $N_i$ . However, if  
 284 the offline mean  $\mu_i^{(\text{off})}$  is larger, the hybrid samples may  
 285 overestimate the true online mean  $\mu_i^{(\text{on})}$ . In this case, the  
 286 contribution of the offline data is attenuated by a discount  
 287 factor  $(1 - 3\omega_i/\Delta_i)$ , which decreases as  $\omega_i$  increases, re-  
 288 sulting in diminishing effective information from the offline  
 289 samples. When  $3\omega_i > \Delta_i$ , the hybrid samples can mislead  
 290 the identification of the optimal arm, and as the result the  
 291 offline benefit vanishes and the regret reduces to that of  
 292 the purely online setting. In contrast, for the optimal arm  
 293 (arm 1), adding the correction term  $V$  to the hybrid samples  
 294 always leads to an overestimation of the true mean  $\mu_1^{(\text{on})}$ .  
 295 Such overestimation consistently favors the identification  
 296 of the optimal arm. Consequently, the contribution of the  
 297 offline data for the optimal arm is not subject to any discount  
 298 factor.  
 299

300  
 301 **Comparison with Cheung & Lyu (2024).** Cheung &  
 302 Lyu (2024) study the same offline-to-online setting as  
 303 ours, but focus on a UCB-based approach. They de-  
 304 rive an upper bound on the sub-optimal  $i$ 's selection time  
 305  $O(C \log T/\Delta_i^2 - N_i \cdot \max\{1 - \omega_i/\Delta_i, 0\}^2)_+$ . The first  
 306 term in our Theorem 5.1, which captures the regret due to  
 307 inaccurate estimation of suboptimal arms, is of the same  
 308 order as this bound. The main difference lies in how the size  
 309 of the offline dataset  $N_i$  is discounted. Specifically, their  
 310 bound uses the discount factor  $(1 - \omega_i/\Delta_i)^2$  whereas ours  
 311 involves the factor  $(1 - 3\omega_i/\Delta_i)$ . This constant coefficient  
 312 discrepancy primarily stems from differences in the analysis  
 313 techniques. In TS, posterior samples may be either opti-  
 314 mistic or pessimistic, which requires partitioning the gap  
 315 between  $\mu_1$  and  $\mu_i$  into three regions, leading to the  $\Delta_i/3$   
 316 factor in our formula. Importantly, our analysis avoids an  
 317 additional squaring of the discount factor.  
 318

319 Compared with this result, our Theorem 5.1 further reveals  
 320 an additional improvement arising from the offline data  $N_1$   
 321 on the optimal arm, as reflected in the second term of our  
 322 bound. This improvement stems from the property that TS  
 323 incurs regret due to inaccurate estimation of both suboptimal  
 324 and optimal arms. When offline data are available for arm 1,  
 325 the estimation of the optimal arm becomes more accurate,  
 326 thereby reducing this source of regret. Such a scenario  
 327 is common in practical applications, since offline data are  
 328 typically collected using expert or near-optimal policies, and  
 329

therefore observations of the optimal arm are often abundant.  
 Our Anchor-TS algorithm is able to exploit this advantage,  
 whereas UCB-based methods do not benefit from offline  
 data on the optimal arm in the same way.

Another aspect worth discussion is the right-shift operation  
 applied to the hybrid samples. Although both the MINUCB  
 algorithm and our approach incorporate a bias-correction  
 term into the hybrid index, the underlying motivations are  
 fundamentally different. In MINUCB, the correction term  
 is introduced to preserve the optimism of the hybrid UCB  
 index, ensuring that it upper-bounds the true reward. In  
 contrast, TS does not intrinsically rely on optimism. The  
 primary purpose of introducing the right-shift in our method  
 is to encourage exploration of arm 1, thereby mitigating  
 the regret caused by underestimation of this arm. While  
 omitting this correction term would in fact improve the first  
 regret term from over-estimation of sub-optimal arms, it  
 would significantly complicate the analysis of  $\hat{\mu}_1^{(\text{on})}$ . In  
 particular, inaccuracies in estimating  $\hat{\mu}_1^{(\text{on})}$  would introduce  
 regret terms that grow exponentially with  $N_1$ , since under-  
 estimation of the hybrid sample  $\theta_1^{(\text{hyb})}$  would make arm 1  
 increasingly unlikely to be selected, leaving little opportu-  
 nity for correction.

**Reduction to the pure online setting and the role of me-  
 dian aggregation.** In the pure online setting, where no  
 offline data are available, the hybrid posterior distribution  
 coincides with the online posterior distribution. In this case,  
 Anchor-TS can be interpreted as sampling two online in-  
 dices,  $\theta_{i,1}^{(\text{on})}$ ,  $\theta_{i,2}^{(\text{on})}$ , and selecting the median among these  
 two samples and the online sample mean  $\hat{\mu}_i^{(\text{on})}$ . We analyze  
 this setting separately and show that this modification pre-  
 serves the classical instance-dependent regret bound, while  
 improving the leading regret order by a factor of  $1/2$ . We  
 further provide empirical results that illustrate this advan-  
 tage. The corresponding theoretical analysis and experimen-  
 tal results are presented in Appendix B. This observation  
 may, to some extent, reflect a similar idea in Jin et al. (2023),  
 who show that TS with less exploration can perform better.  
 But we leverage different ideas to reduce exploration: they  
 divert a fixed probability mass  $\epsilon$  from vanilla TS to select-  
 ing the sample mean, whereas our method incorporates the  
 sample mean in a more adaptive manner through median  
 aggregation.

## 6. Proof Sketch

In this section, we present a proof sketch for Theorem 5.1  
 and highlight the key ideas of our analysis. Our techni-  
 cal contributions operate at two levels: at a high level, we  
 tightly couple the median-based algorithmic design with the  
 regret decomposition, enabling the regret bound to adapt

to the more informative of the online and hybrid posterior distributions; at a more technical level, we develop a refined probabilistic analysis that decouples multiple random indices appearing within probability and expectation operators, allowing their individual regret contributions to be explicitly controlled.

Standard TS analyses consider regret contributions from a single posterior distribution. In contrast, our setting involves two posterior sources, where a naive extension yields additive regret bounds and fails to exploit the more informative one. To address this, we use the online sample mean as an anchor: its estimation errors contribute only a constant-order regret of  $O(1/\Delta_i)$ , allowing us to condition on accurate estimation and isolate the effects of the two posterior distributions.

For convenience, define  $x_i := \mu_i^{(\text{on})} + \Delta_i/3$  and  $y_i := \mu_i^{(\text{on})} + 2\Delta_i/3$  as two thresholds,  $E_i^{\mu^{(\text{on})}}(t) := \{\hat{\mu}_i^{(\text{on})}(t) \leq x_i\}$  as the good event that arm  $i$ 's online sample mean is accurate. Then,

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbf{1}\{A(t) = i\} \right] = \sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu^{(\text{on})}}(t)) + \sum_{t=1}^T \Pr(A(t) = i, E_i^{\mu^{(\text{on})}}(t)).$$

The first term can be upper bounded by  $O(1/\Delta_i^2)$  without relying on  $T$ . The second term representing regret from selecting  $i$  when its online sample mean is accurate can be decomposed as

$$\sum_{t=1}^T \Pr(A(t) = i, E_i^{\mu^{(\text{on})}}(t), \hat{\theta}_i(t) > y_i) + \sum_{t=1}^T \Pr(A(t) = i, E_i^{\mu^{(\text{on})}}(t), \hat{\theta}_i(t) \leq y_i). \quad (2)$$

The first term in (2) corresponds to the case where arm  $i$ 's median index is over-estimated. Conditional on the good event  $\{\hat{\mu}_i^{(\text{on})}(t) \leq x_i\}$ , the event  $\{\hat{\theta}_i(t) > y_i\}$  implies that  $\{\theta_i^{(\text{on})}(t) > y_i\} \cap \{\theta_i^{(\text{hyb})}(t) > y_i\}$ . The corresponding regret event is therefore the intersection of two unfavorable events, whose probability is strictly smaller than that induced by either posterior alone,

$$\sum_{t=1}^T \min \left\{ \Pr(A(t) = i, E_i^{\mu^{(\text{on})}}(t), \theta_i^{(\text{on})}(t) > y_i), \Pr(A(t) = i, E_i^{\mu^{(\text{on})}}(t), \theta_i^{(\text{hyb})}(t) > y_i) \right\}.$$

This yields the first term in Theorem 5.1 by taking the minimum of the corresponding regret contributions from the online and hybrid posterior distributions.

The second term in (2) can be transformed to the case where arm 1's median index is under-estimated. Dealing with this term is commonly recognized as a key challenge in TS analyses (Agrawal & Goyal, 2017; Jin et al., 2021; 2023), which is typically upper bounded by  $\mathbb{E}[1/p - 1]$  where  $p$  is the conditional probability (given the history) that  $\hat{\theta}_1 > y_i$ . In vanilla TS,  $p$  is the tail probability of a single posterior sample  $\hat{\theta}_1 := \theta_1^{(\text{on})}$ , so one can directly convert  $\mathbb{E}[1/p]$  into the expectation of a geometric hitting time related to the behavior of the index. In our setting, however, the arm index is the median of three random quantities with distinct concentration behaviors, which prevents a direct reduction to a one-dimensional geometric hitting-time argument as in vanilla TS.

To handle this term, we further condition on the behavior of  $\hat{\mu}_1^{(\text{on})}$ . Under the good event  $\{\hat{\mu}_1^{(\text{on})}(t) > y_i\}$ , it suffices that either the online or the hybrid posterior sample of arm 1 exceeds  $y_i$  for the optimal arm to be selected. This yields a union of favorable events, implying that

$$p \geq \max\{p^{(\text{on})}, p^{(\text{hyb})}\},$$

$$\text{thus } \frac{1}{p} - 1 \leq \min \left\{ \frac{1}{p^{(\text{on})}} - 1, \frac{1}{p^{(\text{hyb})}} - 1 \right\}.$$

where  $p^{(\text{on})}$  and  $p^{(\text{hyb})}$  represent the conditional probability that  $\{\theta_1^{(\text{on})}(t) > y_i\}$  and  $\{\theta_1^{(\text{hyb})}(t) > y_i\}$ , respectively. This leads to the second term in Theorem 5.1, where the right-hand shift of the hybrid posterior distribution ensures that  $p^{(\text{hyb})} > p^{(\text{on})}$  when offline data arm 1 are available.

Another key difficulty arises when  $\hat{\mu}_1^{(\text{on})}(t) < y_i$ . In this regime, the median index can no longer exploit the more favorable of the two posterior samples, and the success probability  $p$  reduces to an intersection event involving both samples. This joint event can significantly shrink  $p$  and amplify regret. Moreover, since  $\hat{\mu}_1^{(\text{on})}(t)$  and  $\hat{\theta}_1(t)$  are statistically coupled, their contributions cannot be bounded separately and combined multiplicatively.

To address this challenge, we first apply Hölder's inequality to decompose the regret term into two components that can be handled separately. The first component involves the success probability  $1/p$  corresponding to a joint hitting event of two posterior samples. Here, the right-hand shift of the hybrid posterior distribution plays a crucial role by preventing this term from becoming excessively large, ensuring that it remains bounded by a constant. The second component captures the probability of inaccurate estimation of the online sample mean, which decreases exponentially as the selection of arm 1. As a result, the combined contribution of these two components can be controlled by an  $O(1/\Delta_i^2)$  upper bound.

## 7. Experiments

In this section, we compare our Anchor-TS with three classes of baselines: online-only methods including standard TS (Agrawal & Goyal, 2017) and standard UCB (Auer et al., 2002); naive offline-to-online methods including hybrid TS and hybrid UCB (Shivaswamy & Joachims, 2012) which use offline data for initialization while treating it as unbiased; and bias-aware offline-to-online method MINUCB (Cheung & Lyu, 2024). For each experiment, we report cumulative regret over  $10k$  rounds, averaged over 50 runs with error bars showing the standard error. We consider a basic setting with  $K = 10$  arms, where the optimal arm has online mean reward 0.8 and all suboptimal arms have mean 0.5, yielding a gap of 0.3. Rewards are drawn from a Gaussian distribution with unit variance. The total number of offline samples is  $2k$ . To investigate the effect of offline data coverage, we consider three coverage patterns: uniform coverage across arms, coverage concentrated on the optimal arm 1 (with 80% of samples), and coverage concentrated on a suboptimal arm 2 (with 80% of samples). We further explore algorithm performance by varying parameters of this basic setting.

**Unbiased offline data.** We first consider a simple setting with zero bias under different sub-optimality gaps  $\Delta \in \{0.3, 0.1\}$ . The results are reported in Figure 1. There are two common observations: methods that leverage offline data consistently outperform purely online methods, and TS-type algorithms typically outperform their UCB-type counterparts.

In the easiest setting with uniformly distributed and sufficiently abundant offline data (Figure 1 left(a)), both hybrid UCB and hybrid TS achieve zero regret. In contrast, Anchor-TS incurs a small constant regret, as it deliberately requires a constant number of online samples to form a reliable anchor. In other settings, Anchor-TS consistently achieves the lowest regret and outperforms hybrid TS, despite both methods leveraging the same offline information. This gap comes from different exploration control: hybrid TS relies on a single high-variance posterior sample and may tend to over-explore, whereas Anchor-TS reduces exploration through sample-mean-based anchoring. A similar phenomenon was also observed in Jin et al. (2023), where incorporating the sample mean into the decision process leads to improved performance.

Compared with UCB-based algorithms, Anchor-TS consistently achieves substantial gains, particularly when offline data are concentrated on the optimal arm. This behavior aligns with our theoretical analysis, which predicts that the benefit of Anchor-TS scales with the amount of offline data available for the optimal arm, while UCB-based methods do not benefit in the same way. This advantage becomes

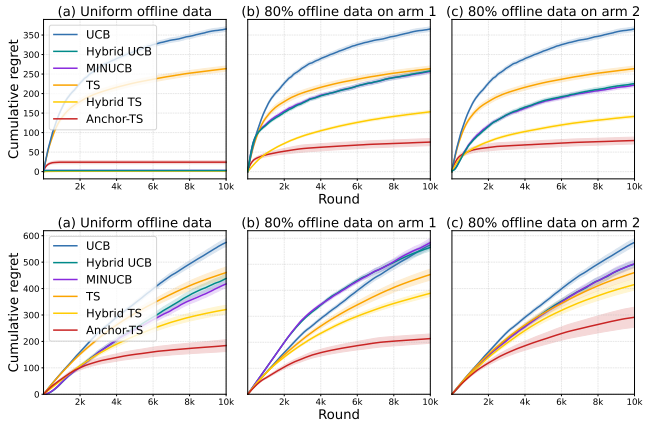


Figure 1. Cumulative regret in the unbiased setting under varying offline coverage regimes and suboptimality gaps (left:  $\Delta = 0.3$ ; right:  $\Delta = 0.1$ ).

even more pronounced in harder settings with a smaller gap where hybrid UCB and MINUCB even underperform pure online UCB (Figure 1 right (b)). This is because the UCB of the optimal arm is excessively reduced when incorporating offline data, forcing other arms to be pulled more often before their sub-optimality can be identified.

**Biased offline data.** We then evaluate the performance of algorithms under distribution shift between offline and online environments. We include both pure online methods and bias-aware approaches as baselines. To ensure that the offline data is sufficiently misleading, we consider a base setting in which  $\mu_1^{(\text{off})} = 0.5$  and  $\mu_i^{(\text{off})} = 0.6$  for  $i \neq 1$ , so that the optimal arm is underestimated in the offline data while the suboptimal arms are overestimated. Figure 2 summarizes the algorithms performances under biased offline data by varying key problem parameters, including the total offline sample size, the real bias level, the hyper-parameter  $V$ , and the arm set size. Unless otherwise specified, the hyperparameter  $V$  is set to the true bias magnitude for each arm. Across all settings, Anchor-TS consistently achieves the lowest regret, demonstrating robustness to offline bias and markedly outperforming UCB-based methods in its ability to effectively utilize offline information.

The top-left panel shows the algorithms performance under different offline sample sizes  $\sum_i N_i$ . Increasing offline data generally improves bias-aware methods, especially under uniform coverage.

The top-right panel studies the effect of the bias level  $\delta$ . For ease of implementation, we fix the offline mean of sub-optimal arms and vary  $\mu_1^{(\text{off})}$  by setting different choices of  $\delta := \mu_1^{(\text{on})} - \mu_1^{(\text{off})}$ . Larger  $\delta$  corresponds to stronger underestimation of the optimal arm in the offline data. Nevertheless, Anchor-TS maintains a stable low-regret profile

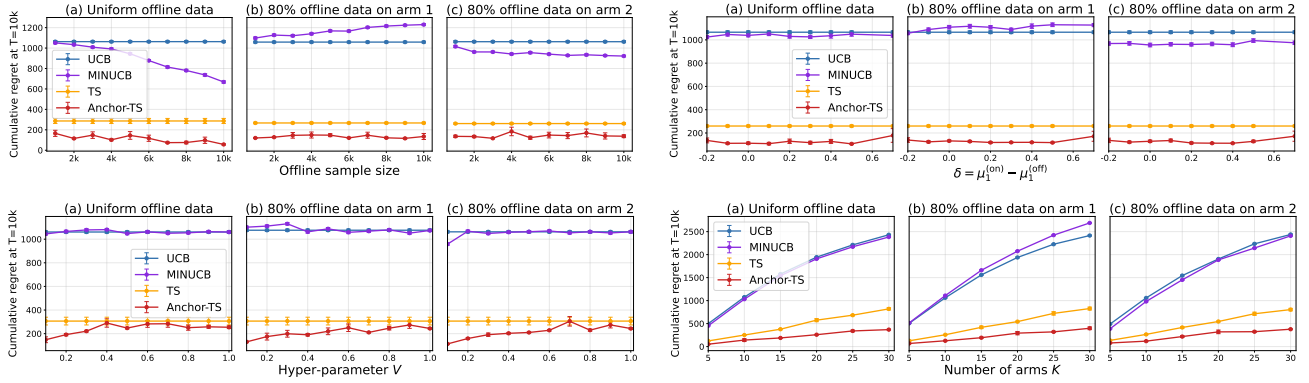


Figure 2. Cumulative regret under varying offline coverage regimes and problem parameters (top-left: total offline sample size; top-right: bias level  $\delta$ ; bottom-left: hyperparameter  $V$ ; bottom-right: number of arms  $K$ ).

and consistently outperforms baselines.

The bottom-left panel examines the impact of hyperparameter  $V$ . For MINUCB and Anchor-TS, we set  $V$  to the maximum of the true bias magnitude and the tested value. When  $V$  is small, tighter prior control allows Anchor-TS to exploit offline data more aggressively. As  $V$  increases, its performance degrades but never falls below pure online TS. Even when offline data are uninformative with large  $V$ , Anchor-TS retains a slight advantage due to sample-mean based exploration.

The bottom-right panel evaluates scalability with respect to the number of arms  $K$ . While regret increases with  $K$  for all methods, UCB-based algorithms grow much faster, whereas Anchor-TS scales more favorably by suppressing unnecessary exploration.

## 8. Conclusion

In this work, we study the TS-type algorithm for offline-to-online learning. Such algorithms face a fundamental challenge in assessing the reliability of the hybrid posterior distribution through comparison with the online posterior, stemming from the lack of inherent optimism in posterior samples. We address this challenge by introducing a novel median-based aggregation mechanism that combines the online sample mean, an online posterior sample, and a hybrid posterior sample, enabling adaptive exploitation of informative offline data. Our theoretical results demonstrate that our algorithm can effectively cope with bias in offline data, while achieving strictly better regret improvement when the offline data are sufficiently informative. Compared with UCB-type algorithms, our TS-based approach can further leverage abundant offline data on the optimal arm to obtain larger performance gains. Moreover, even in the pure online setting, the proposed median-index mechanism yields provable improvement in the leading constant of the regret bound compared with pure online TS. Extensive empirical evalua-

tions further validate our theoretical findings, demonstrating that our algorithm consistently outperforms state-of-the-art baselines across various settings of data coverage and bias magnitudes.

Several promising directions merit further investigation. First, it would be valuable to establish gap-independent regret bounds that explicitly quantify the benefits of offline data. Second, extending the proposed anchor-based framework to contextual TS is an important direction for addressing high-dimensional, real-world decision-making problems.

## Impact Statement

This paper develops theory and algorithms for offline-to-online bandit learning under distribution shift, with experiments only on synthetic instances. The work is methodological and we do not release datasets or models. We do not foresee direct dual-use risks from the theoretical results themselves.

## References

- Abramowitz, M., Stegun, I. A., and Romer, R. H. Handbook of mathematical functions with formulas, graphs, and mathematical tables. *American Journal of Physics*, 56(10):958–958, 1988.
- Agnihotri, A., Jain, R., Ramachandran, D., and Wen, Z. Online bandit learning with offline preference data for improved rlhf. *arXiv preprint arXiv:2406.09574*, 2024.
- Agrawal, S. and Goyal, N. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, pp. 39–1. JMLR Workshop and Conference Proceedings, 2012.
- Agrawal, S. and Goyal, N. Thompson sampling for contextual bandits with linear payoffs. In *Proceedings of the 30th International Conference on Machine Learning*, pp. 127–135. PMLR, 2013.
- Agrawal, S. and Goyal, N. Near-optimal regret bounds for thompson sampling. *Journal of the ACM (JACM)*, 64(5): 1–24, 2017.
- Ahn, H.-S., Zhang, M., and Zhang, Z. Online decisions with (biased) offline data. *Available at SSRN 5350921*, 2025.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002.
- Ball, P. J., Smith, L., Kostrikov, I., and Levine, S. Efficient online reinforcement learning with offline data. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 1577–1594. PMLR, 2023.
- Banerjee, S., Sinclair, S. R., Tambe, M., Xu, L., and Yu, C. L. Artificial replay: a meta-algorithm for harnessing historical data in bandits. *arXiv preprint arXiv:2210.00025*, 2022.
- Chapelle, O. and Li, L. An empirical evaluation of thompson sampling. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, pp. 2249–2257, 2011.
- Cheung, W. C. and Lyu, L. Leveraging (biased) information: multi-armed bandits with offline data. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 8286–8309. PMLR, 2024.
- Daniel, J. R., Benjamin, V. R., Abbas, K., Ian, O., and Zheng, W. A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–99, 2018.
- Granmo, O.-C. Solving two-armed bernoulli bandit problems using a bayesian learning automaton. *International Journal of Intelligent Computing and Cybernetics*, 3(2): 207–234, 2010.
- He, Q., Wang, M., Liu, X., Wang, Z., and Kong, F. Learning across the gap: Hybrid multi-armed bandits with heterogeneous offline and online data. In *The 39th Annual Conference on Neural Information Processing Systems*, 2025.
- Huang, R., Li, D., Shi, C., Shen, C., and Yang, J. Augmenting online rl with offline data is all you need: A unified hybrid rl algorithm design and analysis. In *The 41st Conference on Uncertainty in Artificial Intelligence*, 2025.
- Jin, T., Xu, P., Shi, J., Xiao, X., and Gu, Q. Mots: Minimax optimal thompson sampling. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 5074–5083. PMLR, 2021.
- Jin, T., Yang, X., Xiao, X., and Xu, P. Thompson sampling with less exploration is fast and optimal. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 15239–15261. PMLR, 2023.
- Kaufmann, E., Korda, N., and Munos, R. Thompson sampling: An asymptotically optimal finite-time analysis. In *International Conference on Algorithmic Learning Theory*, pp. 199–213. Springer, 2012.
- Komiyama, J., Honda, J., and Nakagawa, H. Optimal regret analysis of thompson sampling in stochastic multi-armed bandit problem with multiple plays. In *Proceedings of the 32nd International Conference on Machine Learning*, pp. 1152–1161. PMLR, 2015.
- Lai, T. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1): 4–22, 1985.
- Lattimore, T. and Szepesvári, C. *Bandit algorithms*. Cambridge University Press, 2020.
- Lee, S., Seo, Y., Lee, K., Abbeel, P., and Shin, J. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, pp. 1702–1712. PMLR, 2022.

- 550 Liu, C.-Y. and Li, L. On the prior sensitivity of thompson  
551 sampling. In *International Conference on Algorithmic*  
552 *Learning Theory*, pp. 321–336. Springer, 2016.
- 553 Nair, A., Gupta, A., Dalal, M., and Levine, S. Awac: Accelerating  
554 online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- 555 Nakamoto, M., Zhai, S., Singh, A., Sobol Mark, M., Ma, Y.,  
556 Finn, C., Kumar, A., and Levine, S. Cal-ql: Calibrated  
557 offline rl pre-training for efficient online fine-tuning. *Proceedings of the 37th International Conference on Neural*  
558 *Information Processing Systems*, pp. 62244–62269, 2023.
- 559 Oetomo, B., Perera, R. M., Borovica-Gajic, R., and Rubinstein,  
560 B. I. Cutting to the chase with warm-start contextual  
561 bandits. *Knowledge and Information Systems*, 65(9):  
562 3533–3565, 2023.
- 563 Qu, C., Shi, L., Panaganti, K., You, P., and Wierman, A.  
564 Hybrid transfer reinforcement learning: Provable sample  
565 efficiency from shifted-dynamics data. In *International*  
566 *Conference on Artificial Intelligence and Statistics*, pp.  
567 1054–1062. PMLR, 2025.
- 568 Scott, S. L. A modern bayesian look at the multi-armed bandit.  
569 *Applied Stochastic Models in Business and Industry*,  
570 26(6):639–658, 2010.
- 571 Sentenac, F., Lee, I., and Szepesvari, C. Balancing optimism  
572 and pessimism in offline-to-online learning. *arXiv preprint arXiv:2502.08259*, 2025.
- 573 Shivaswamy, P. and Joachims, T. Multi-armed bandit problems  
574 with history. In *Artificial Intelligence and Statistics*,  
575 pp. 1046–1054. PMLR, 2012.
- 576 Simchowitz, M., Tosh, C., Krishnamurthy, A., Hsu, D., Lykouris,  
577 T., Dudík, M., and Schapire, R. Bayesian decision-making  
578 under misspecified priors with applications to meta-learning.  
579 In *Proceedings of the 35th International Conference on Neural*  
580 *Information Processing Systems*, pp. 26382–26394, 2021.
- 581 Song, Y., Zhou, Y., Sekhari, A., Bagnell, D., Krishnamurthy,  
582 A., and Sun, W. Hybrid rl: Using both offline and online  
583 data can make rl efficient. In *The 11th International*  
584 *Conference on Learning Representations*, 2023.
- 585 Tan, K. and Xu, Z. A natural extension to online algorithms  
586 for hybrid rl with limited coverage. *Reinforcement*  
587 *Learning Journal*, 1, 2024.
- 588 Thompson, W. R. On the likelihood that one unknown  
589 probability exceeds another in view of the evidence of  
590 two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 591 Verstraeten, T., Bargiacchi, E., Libin, P. J., Helsen, J., Roijers,  
592 D. M., and Nowé, A. Multi-agent thompson sampling for  
593 bandit applications with sparse neighbourhood structures. *Scientific Reports*, 10(1):6728, 2020.
- 594 Wagenmaker, A. and Pacchiano, A. Leveraging offline data in  
595 online reinforcement learning. In *Proceedings of the 40th*  
596 *International Conference on Machine Learning*, pp.  
597 35300–35338. PMLR, 2023.
- 598 Wang, S. and Chen, W. Thompson sampling for combinatorial  
599 semi-bandits. In *Proceedings of the 35th International*  
600 *Conference on Machine Learning*, pp. 5114–5122. PMLR, 2018.
- 601 Xia, Y., Xie, Z., Yu, T., Zhao, C., and Li, S. Toward joint  
602 utilization of absolute and relative bandit feedback for  
603 conversational recommendation. *User Modeling and User-*  
604 *Adapted Interaction*, 34(5):1707–1744, 2024.
- 605 Xie, T., Jiang, N., Wang, H., Xiong, C., and Bai, Y. Policy  
606 finetuning: Bridging sample-efficient offline and online  
607 reinforcement learning. In *35th Conference on Neural*  
608 *Information Processing Systems*, pp. 27395–27407, 2021.
- 609 Yang, L., Tan, V. Y., and Cheung, W. C. Best arm identification  
610 with possibly biased offline data. *arXiv preprint arXiv:2505.23165*, 2025.
- 611 Yin, Z. and Fang, Z. Multi-armed bandits with biased and  
612 heteroscedastic auxiliary rewards. In *Proceedings of the 34th*  
613 *ACM International Conference on Information and Knowledge*  
614 *Management*, pp. 3899–3908, 2025.
- 615 Yu, Z. and Zhang, X. Actor-critic alignment for offline-to-online  
616 reinforcement learning. In *Proceedings of the 40th*  
617 *International Conference on Machine Learning*, pp.  
618 40452–40474. PMLR, 2023.
- 619 Zhang, C., Agarwal, A., Iii, H. D., Langford, J., and Negahban,  
620 S. Warm-starting contextual bandits: Robustly combining  
621 supervised and bandit feedback. In *Proceedings of the 36th*  
622 *International Conference on Machine Learning*, pp. 7335–7344. PMLR, 2019.
- 623 Zhang, W., Zhou, D., Li, L., and Gu, Q. Neural thompson  
624 sampling. In *The 9th International Conference on Learning*  
625 *Representations*, 2021.
- 626 Zhang, Y., Zhu, R., and Xie, Q. Contextual online pricing  
627 with (biased) offline data. In *The 39th Annual Conference*  
628 *on Neural Information Processing Systems*, 2025.
- 629 Zhou, K., Zhang, T., Chen, W., and Kong, F. Hybrid combinatorial  
630 multi-armed bandits with probabilistically triggered arms. *arXiv preprint arXiv:2512.21925*, 2025.

## A. Proof of Theorem 5.1

We first introduce some useful definitions that will be used in the proof.

For each sub-optimal arm  $i \neq 1$ , we define two thresholds:  $x_i := \mu_i^{(\text{on})} + \Delta_i/3$  and  $y_i := \mu_i^{(\text{on})} + 2\Delta_i/3 = \mu_1^{(\text{on})} - \Delta_i/3$ . It holds that  $x_i < y_i$  and  $y_i - x_i = \Delta_i/3$ .

Define the confidence radius  $L_t := 2 \log(2t/\delta_t)$  with uncertainty level  $\delta_t = 1/(t^2 \Delta_i^2)$ .

Let  $\mathcal{F}_t := \cup_{i=1}^K S_i \cup \{(A(\tau), R_{A(\tau)}(\tau)) : 1 \leq \tau \leq t\}$  to represent the history comprising the offline dataset and the sequence of online observations collected up to round  $t$ .

For the optimal arm 1 and arm  $i$ , define  $p_{i,t} := \Pr(\hat{\theta}_1(t) > y_i \mid \mathcal{F}_{t-1})$  to represent the conditional probability that the median-index is larger than  $y_i$ . Similarly, define  $p_{i,t}^{(\text{on})} := \Pr(\theta_1^{(\text{on})}(t) > y_i \mid \mathcal{F}_{t-1})$  and  $p_{i,t}^{(\text{hyb})} := \Pr(\theta_1^{(\text{hyb})}(t) > y_i \mid \mathcal{F}_{t-1})$ .

Further define the following good events representing the estimators are not far from their centers.

**Definition A.1** (Good events  $E_i^{\mu(\text{on})}(t), E_i^{\mu(\text{hyb})}(t), E_i^{\theta(\text{hyb})}(t), \mathcal{G}_i(t), E_1^{\mu(\text{on})}(t)$ ). For  $i \neq 1$ , define the following good event

$$\begin{aligned} E_i^{\mu(\text{on})}(t) &:= \left\{ \hat{\mu}_i^{(\text{on})}(t) \leq x_i \right\}, \\ E_i^{\mu(\text{hyb})}(t) &:= \left\{ \left| \hat{\mu}_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})} \right| \leq \sqrt{\frac{L_t}{N_i(t) + T_i} + \frac{T_i V_i}{N_i(t) + T_i}} \right\}, \\ E_i^{\theta(\text{hyb})}(t) &:= \left\{ \left| \theta_i^{(\text{hyb})}(t) - \hat{\mu}_i^{(\text{hyb})}(t) \right| \leq \sqrt{\frac{L_t}{N_i(t) + T_i} + \frac{T_i V_i}{N_i(t) + T_i}} \right\}. \end{aligned}$$

and the global good event

$$\mathcal{G}_i(t) := E_i^{\mu(\text{on})}(t) \cap E_i^{\mu(\text{hyb})}(t) \cap E_i^{\theta(\text{hyb})}(t).$$

Similarly, for optimal arm 1, define

$$E_1^{\mu(\text{on})}(t) := \left\{ \hat{\mu}_1^{(\text{on})}(t) > y_i \right\}.$$

In the following, we provide the detailed proof of Theorem 5.1.

*Proof of Theorem 5.1.* We first analyze the regret by bounding the number of selections of each sub-optimal arm  $i \neq 1$ .

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{A(t) = i\} \right] &= \sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu(\text{on})}(t)) + \sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu(\text{hyb})}(t)) \\ &\quad + \sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\theta(\text{hyb})}(t)) + \sum_{t=1}^T \Pr(A(t) = i, \mathcal{G}_i(t)) \end{aligned} \quad (3)$$

$$\leq O\left(\frac{1}{\Delta_i^2}\right) + \sum_{t=1}^T \Pr(A(t) = i, \mathcal{G}_i(t)) \quad (4)$$

$$\leq O\left(\frac{1}{\Delta_i^2}\right) + \sum_{t=1}^T \Pr(\underbrace{A(t) = i, \mathcal{G}_i(t), \hat{\theta}_i(t) > y_i}_{\mathcal{E}_{1,t}})$$

$$+ \sum_{t=1}^T \Pr(\underbrace{A(t) = i, \mathcal{G}_i(t), \hat{\theta}_i(t) \leq y_i}_{\mathcal{E}_{2,t}})$$

$$\leq O\left(\frac{1}{\Delta_i^2}\right) + \sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) + \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1}\{E_1^{\mu(\text{on})}(\tau_k + 1)\} \right]$$

$$+ \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \quad (5)$$

$$\leq \frac{C_3}{\Delta_i^2} + \left( \frac{C_1 \log T}{\Delta_i^2} - N_i \max \left\{ \left( 1 - \frac{3\omega_i}{\Delta_i} \right), 0 \right\} \right)_+ + c \left( \frac{C_2 \log T}{\Delta_i^2} - N_1 \right)_+. \quad (6)$$

Here (4) follows from Lemma D.4, Lemma A.2, and Lemma A.3 for bounding the bad events. (5) is due to Lemma A.5. (6) comes from the upper bound for  $\mathcal{E}_1$  in Lemma A.4 and the upper bound for two terms of  $\mathcal{E}_2$  in Lemma A.6 and Lemma A.9.  $C_1, C_2, C_3$  are constant terms independent of the problem parameters,  $c = \max \left\{ e^{11}, \exp \left\{ \frac{28+16\sqrt{3}}{N_1+1} \right\} \right\} + 5$  is a constant smaller than the coefficient in the pure online TS (Agrawal & Goyal, 2017) and decreases as  $N_1$  increases.  $O$  hides constant terms.

The final regret can then be obtained by

$$\begin{aligned} \text{Reg}(T) &= \sum_{i \neq 1} \Delta_i \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} \{ A(t) = i \} \right] \\ &\leq \sum_{i \neq 1} \Delta_i \left( \frac{C_3}{\Delta_i^2} + \left( \frac{C_1 \log T}{\Delta_i^2} - N_i \max \left\{ \left( 1 - \frac{3\omega_i}{\Delta_i} \right), 0 \right\} \right)_+ + c \left( \frac{C_2 \log T}{\Delta_i^2} - N_1 \right)_+ \right). \end{aligned}$$

□

We next introduce the lemmas used in the above proof and provide the proof of lemmas in Appendix C. The first two lemmas establish constant upper bounds on the probability of bad events.

**Lemma A.2.** For sub-optimal arm  $i \neq 1$ , the second term in (3) can be bounded by

$$\sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu(\text{hyb})}(t)) \leq \sum_{t=1}^T \delta_t = O \left( \frac{1}{\Delta_i^2} \right).$$

**Lemma A.3.** For sub-optimal arm  $i \neq 1$ , the third term in (3) can be bounded by

$$\sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\theta(\text{hyb})}(t)) \leq \sum_{t=1}^T \frac{\delta_t}{t} = O \left( \frac{1}{\Delta_i^2} \right).$$

The following lemma establishes an upper bound on  $\sum_{t=1}^T \Pr(\mathcal{E}_{1,t})$  in the regret, corresponding to the regret incurred when the global good event  $\mathcal{G}_i(t)$  holds and the index  $\hat{\theta}_i(t)$  of arm  $i$  is inaccurate.

**Lemma A.4.** For sub-optimal arm  $i \neq 1$ ,

$$\sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) \leq \left( 36 \frac{\log(T\Delta_i^2 + e^6)}{\Delta_i^2} - N_i \max \left\{ \left( 1 - \frac{3\omega_i}{\Delta_i} \right), 0 \right\} \right)_+ + \frac{1}{\Delta_i^2}.$$

The following lemma establishes an upper bound on  $\sum_{t=1}^T \Pr(\mathcal{E}_{2,t})$  in the regret. This term can be transformed into the case where the index  $\hat{\theta}_1(t)$  of arm 1 is inaccurate, which may lead to the selection of arm  $i$ .

**Lemma A.5.** For all  $t, i \neq 1$  and all instantiations  $F_{t-1}$  of  $\mathcal{F}_{t-1}$ ,

$$\Pr(A(t) = i, \mathcal{G}_i(t), \hat{\theta}_i(t) \leq y_i \mid F_{t-1}) \leq \frac{1 - p_{i,t}}{p_{i,t}} \Pr(A(t) = 1, \mathcal{G}_i(t), \hat{\theta}_i(t) \leq y_i \mid F_{t-1}).$$

Further,

$$\sum_{t=1}^T \Pr(\mathcal{E}_{2,t}) \leq \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] + \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right].$$

**Lemma A.6.** For any  $t \geq 0$ ,

$$\sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \leq \sum_{k=0}^{T-1} \min \left\{ \mathbb{E} \left[ \frac{1}{p_{i, \tau_k + 1}^{(\text{on})}} - 1 \right], \mathbb{E} \left[ \frac{1}{p_{i, \tau_k + 1}^{(\text{hyb})}} - 1 \right] \right\}.$$

We then prove a similar lemma corresponding to Lemma D.6 for the hybrid estimator.

**Lemma A.7.** Let  $\tau_k$  be the time of the  $k$ -th play of arm 1. Then

$$\mathbb{E} \left[ \frac{1}{p_{i, \tau_k + 1}^{(\text{hyb})}} - 1 \right] \leq \begin{cases} \max \left\{ e^{11}, \exp \left( \frac{28 + 16\sqrt{3}}{N_1 + 1} \right) \right\} + 5, & \forall k, \\ \frac{5}{T \Delta_i^2}, & k > L_i^{(\text{hyb})}(T), \end{cases}$$

where  $L_i^{(\text{hyb})}(T) = \left( \frac{288 \ln(T \Delta_i^2 + e^6)}{\Delta_i^2} - N_1 \right)_+$ .

**Lemma A.8.**

$$\sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \leq C_1 \left( \frac{\log(T \Delta_i^2 + e^{32})}{\Delta_i^2} - N_1 \right)_+ + C_2 \frac{1}{\Delta_i^2},$$

for suitable constants  $C_1, C_2 > 0$  independent of  $T$ .

*Proof of Lemma A.8.* Combining Lemma A.6, Lemma D.6 and Lemma A.7 we get the bound.  $\square$

**Lemma A.9.**

$$\sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1 - p_{i, \tau_k + 1}}{p_{i, \tau_k + 1}} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \leq \frac{C}{\Delta_i^2},$$

for some constant  $C$ .

## B. Theoretical Analysis and Experiments in the Pure Online Setting

In this section, we provide a tighter analysis to the pure online setting, where no offline data are available. In this case, the hybrid posterior coincides with the online posterior, and Anchor-TS reduces to selecting arms according to the median index

$$\text{median} \left\{ \theta_{i,1}^{(\text{on})}(t), \theta_{i,2}^{(\text{on})}(t), \hat{\mu}_i^{(\text{on})} \right\},$$

where  $\theta_{i,1}^{(\text{on})}$  and  $\theta_{i,2}^{(\text{on})}$  are i.i.d. samples from the standard online posterior of arm  $i$  at time  $t$ . We show that this median rule preserves the classical instance-dependent regret guarantee while improving the leading logarithmic term by a factor  $1/2$  compared with vanilla TS.

### B.1. $1/2$ improvement on the leading $T$ -term

We reuse the notation in Appendix A. Recall that the  $\log T$  contribution in Theorem 5.1 arises from two parts: (i) controlling over-optimistic indices of suboptimal arms (captured by the term  $\sum_{t=1}^T \Pr(\mathcal{E}_{1,t})$  as in Lemma A.4), and (ii) controlling under-exploration of the optimal arm (via the ratio term in Lemma A.5). We bound these two parts separately.

**Suboptimal-arm term.** Recall that for a suboptimal arm  $i \neq 1$ ,

$$\sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) = \sum_{t=1}^T \Pr \left( A(t) = i, \mathcal{G}_i(t), \hat{\theta}_i(t) > y_i \right). \quad (7)$$

Conditioned on the global good event  $\mathcal{G}_i(t)$ , we have  $\hat{\mu}_i^{(\text{on})} < x_i < y_i$ . Therefore, the event  $\{\hat{\theta}_i(t) > y_i\}$  can only happen when both posterior samples exceed  $y_i$ , i.e.,

$$\left\{ \text{median} \left\{ \theta_{i,1}^{(\text{on})}(t), \theta_{i,2}^{(\text{on})}(t), \hat{\mu}_i^{(\text{on})}(t) \right\} > y_i \right\} = \left\{ \theta_{i,1}^{(\text{on})}(t) > y_i \cap \theta_{i,2}^{(\text{on})}(t) > y_i \right\}.$$

Since  $\theta_{i,1}^{(\text{on})}(t), \theta_{i,2}^{(\text{on})}(t)$  are i.i.d. conditional on  $F_{t-1}$ ,

$$\Pr(\hat{\theta}_i(t) > y_i \mid \mathcal{F}_{t-1}) = \Pr(\theta_i^{(\text{on})}(t) > y_i \mid \mathcal{F}_{t-1})^2. \quad (8)$$

Using (8) squares the Gaussian tail probability (Lemma D.2) in the vanilla TS analysis (Lemma D.4) and yields a halved exploration threshold as

$$\begin{aligned} \sum_{t=1}^T \Pr\left(T_i(t) > \frac{1}{2}L_i(T), G_i(t), \hat{\theta}_i(t) > y_i\right) &\leq \sum_{t=\frac{1}{2}L_i(T)}^T \left(\exp\left\{-\frac{T_i(t)(y_i - x_i)^2}{2}\right\}\right)^2 \\ &\leq \sum_{t=1}^T \left(\exp\left\{-\frac{\frac{1}{2}L_i(t)(y_i - x_i)^2}{2}\right\}\right)^2 \\ &\leq \sum_{t=1}^T \frac{1}{T\Delta_i^2} \\ &\leq \frac{1}{\Delta_i^2}, \end{aligned} \quad (9)$$

where the threshold  $L_i(T)$  is the same as in Lemma D.4.

Thus take (9), we get

$$\begin{aligned} \sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) &\leq \sum_{t=1}^T \Pr\left(A(t) = i, T_i(t) \leq \frac{1}{2}L_i(T), G_i(t), \hat{\theta}_i(t) > y_i\right) \\ &\quad + \sum_{t=1}^T \Pr\left(A(t) = i, T_i(t) > \frac{1}{2}L_i(T), G_i(t), \hat{\theta}_i(t) > y_i\right) \\ &\leq \frac{1}{2}L_i(T) + \sum_{t=1}^T \Pr\left(T_i(t) > \frac{1}{2}L_i(T), G_i(t), \hat{\theta}_i(t) > y_i\right) \\ &\leq \frac{1}{2}L_i(T) + \frac{1}{\Delta_i^2}, \end{aligned}$$

which matches the classical bound but with a factor 1/2 on the leading logarithmic term.

**Optimal-arm term.** Next we bound the term arising from Lemma A.5, namely

$$\sum_{t=1}^T \Pr(\mathcal{E}_{2,t}) \leq \sum_{k=0}^{T-1} \mathbb{E}\left[\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\left\{E_1^{\mu^{(\text{on})}}(\tau_k + 1)\right\}\right] + \sum_{k=0}^{T-1} \mathbb{E}\left[\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\left\{\neg E_1^{\mu^{(\text{on})}}(\tau_k + 1)\right\}\right]. \quad (10)$$

We mainly focus on the first term since the second term upper bound does not depend on  $T$ . Under  $E_1^{\mu^{(\text{on})}}(\tau_k + 1)$ , we have  $\hat{\mu}_1^{(\text{on})}(\tau_k + 1) > y_i$ , hence

$$p_{i,\tau_k+1} = \Pr(\hat{\theta}_1(\tau_k + 1) > y_i \mid \mathcal{F}_{\tau_k}) = \Pr(\theta_{1,1}^{(\text{on})}(\tau_k + 1) > y_i \cup \theta_{1,2}^{(\text{on})}(\tau_k + 1) > y_i \mid \mathcal{F}_{\tau_k}).$$

Let  $p_{i,\tau_k+1}^{(\text{on})} := \Pr(\theta_1^{(\text{on})}(\tau_k + 1) > y_i \mid \mathcal{F}_{\tau_k})$ . Then

$$\begin{aligned} p_{i,\tau_k+1} &= 1 - \left(1 - p_{i,\tau_k+1}^{(\text{on})}\right)^2 \\ &= 2p_{i,\tau_k+1}^{(\text{on})} - \left(p_{i,\tau_k+1}^{(\text{on})}\right)^2. \end{aligned}$$

Therefore,

$$\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} = \frac{(1 - p_{i,\tau_k+1}^{(\text{on})})^2}{p_{i,\tau_k+1}^{(\text{on})}(2 - p_{i,\tau_k+1}^{(\text{on})})} = \frac{1 - p_{i,\tau_k+1}^{(\text{on})}}{p_{i,\tau_k+1}^{(\text{on})}} \cdot \frac{1 - p_{i,\tau_k+1}^{(\text{on})}}{2 - p_{i,\tau_k+1}^{(\text{on})}} \leq \frac{1}{2} \cdot \frac{1 - p_{i,\tau_k+1}^{(\text{on})}}{p_{i,\tau_k+1}^{(\text{on})}}, \quad (11)$$

where we used  $(1 - p)/(2 - p) \leq 1/2$  for  $p \in [0, 1]$ . Plugging (11) into the first term of (10) directly yields a factor-1/2 reduction on the leading logarithmic contribution in the optimal-arm analysis. The second term in (10) is handled identically to Appendix A and remains a lower-order constant term.

## B.2. Experiments in the pure online setting

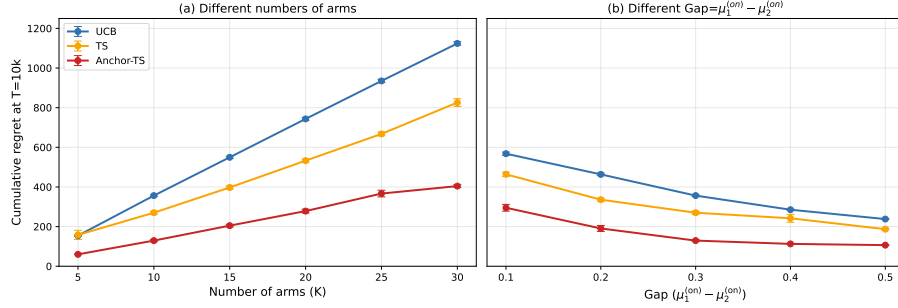


Figure 3. Cumulative regret of our Anchor-TS and baselines in the pure online setting with different arm set size  $K$  and sub-optimality gap  $\Delta$ .

We examine the pure online setting in Figure 3, where no offline data are available. In this case, MIN-UCB, hybrid UCB, and hybrid TS reduce to their respective pure online variants. The experimental setting follows Section 7. We test different choices of arm set size  $K$  and sub-optimality gap  $\Delta = \mu_1^{(\text{on})} - \mu_2^{(\text{on})}$ . Anchor-TS makes decisions based on the median of two posterior samples and the empirical mean from online data. As shown in the figure, Anchor-TS consistently achieves the lowest cumulative regret across all tested settings. This demonstrates that even without offline data, the median-based decision rule employed by Anchor-TS effectively suppresses over-exploration caused by the randomness of posterior sampling, thereby exhibiting greater learning efficiency compared with vanilla TS and UCB-based methods.

## C. Proof of Lemmas

*Proof of Lemma A.2.* We expand the difference between the empirical hybrid mean and the true mean of the online distribution in  $E_i^{\mu^{(\text{hyb})}}(t)$  as follows:

$$\begin{aligned} \hat{\mu}_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})} &= \left( \hat{\mu}_i^{(\text{hyb})}(t) - \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i} \right) + \left( \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i} - \mu_i^{(\text{on})} \right) \\ &= \left( \hat{\mu}_i^{(\text{hyb})}(t) - \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i} \right) + \frac{N_i}{T_i(t) + N_i} (\mu_i^{(\text{off})} - \mu_i^{(\text{on})}). \end{aligned}$$

Using  $|\mu_i^{(\text{off})} - \mu_i^{(\text{on})}| \leq V_i$ , we obtain

$$|\hat{\mu}_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})}| \leq \left| \frac{T_i(t)\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{T_i(t) + N_i} - \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i} \right| + \frac{N_i V_i}{T_i(t) + N_i}. \quad (12)$$

Combining the definition of  $\neg E_i^{\mu^{(\text{hyb})}}(t)$

$$|\hat{\mu}_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})}| > \sqrt{\frac{L_t}{T_i(t) + N_i}} + \frac{N_i V_i}{T_i(t) + N_i},$$

with (12), we can derive the following relation:

$$\neg E_i^{\mu(\text{hyb})}(t) \subseteq \left\{ \left| \frac{T_i(t)\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{T_i(t) + N_i} - \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i} \right| > \sqrt{\frac{L_t}{T_i(t) + N_i}} \right\}. \quad (13)$$

Noting that in (13),

$$\mathbb{E} \left[ \frac{T_i(t)\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{T_i(t) + N_i} \right] = \frac{T_i(t)\mu_i^{(\text{on})}(t) + N_i\mu_i^{(\text{off})}}{T_i(t) + N_i}.$$

Given  $T_i(t) = s$ , by Chernoff-Hoeffding bounds (Lemma D.1), we have

$$\Pr \left( \left| \frac{s\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{s + N_i} - \mathbb{E} \left[ \frac{s\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{s + N_i} \right] \right| > \sqrt{\frac{L_t}{s + N_i}} \right) \leq 2e^{-(s+N_i)\frac{L_t}{s+N_i}} \leq \frac{\delta_t}{t}.$$

Thus

$$\begin{aligned} \Pr \left( \neg E_i^{\mu(\text{hyb})}(t) \right) &= \Pr \left( \exists s \in \{0, \dots, t-1\} : T_i(t) = s \text{ and } \neg E_i^{\mu(\text{hyb})}(t) \right) \\ &\leq \sum_{s=0}^{t-1} \Pr \left( \left| \frac{s\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{s + N_i} - \mathbb{E} \left[ \frac{s\hat{\mu}_i^{(\text{on})}(t) + N_i\hat{\mu}_i^{(\text{off})}}{s + N_i} \right] \right| > \sqrt{\frac{L_t}{s + N_i}} \right) \\ &\leq \sum_{s=0}^{t-1} \frac{\delta_t}{t} \leq \delta_t. \end{aligned}$$

Then we have

$$\sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu(\text{hyb})}(t)) \leq \sum_{t=1}^T \Pr \left( \neg E_i^{\mu(\text{hyb})}(t) \right) \leq \sum_{t=1}^T \delta_t = O\left(\frac{1}{\Delta_i^2}\right).$$

□

*Proof of Lemma A.3.* Recall that

$$\theta_i^{(\text{hyb})}(t) \sim \mathcal{N} \left( \hat{\mu}_i^{(\text{hyb})}(t) + \frac{N_i V_i}{T_i(t) + N_i}, \frac{1}{T_i(t) + N_i + 1} \right),$$

thus  $\mathbb{E}[\theta_i^{(\text{hyb})}(t)] = \hat{\mu}_i^{(\text{hyb})}(t) + \frac{N_i V_i}{T_i(t) + N_i}$ . By the triangle inequality,

$$|\theta_i^{(\text{hyb})}(t) - \hat{\mu}_i^{(\text{hyb})}(t)| \leq |\theta_i^{(\text{hyb})}(t) - \mathbb{E}[\theta_i^{(\text{hyb})}(t)]| + |\mathbb{E}[\theta_i^{(\text{hyb})}(t)] - \hat{\mu}_i^{(\text{hyb})}(t)|$$

Thus

$$\begin{aligned} \neg E_i^{\theta(\text{hyb})}(t) &= \left\{ |\theta_i^{(\text{hyb})}(t) - \hat{\mu}_i^{(\text{hyb})}(t)| > \sqrt{\frac{L_t}{T_i(t) + N_i} + \frac{N_i V_i}{T_i(t) + N_i}} \right\} \\ &\subseteq \left\{ |\theta_i^{(\text{hyb})}(t) - \mathbb{E}[\theta_i^{(\text{hyb})}(t)]| \geq \sqrt{\frac{L_t}{T_i(t) + N_i}} \right\}. \end{aligned} \quad (14)$$

Conditional on the filtration  $\mathcal{F}_t$ ,  $\hat{\mu}_i^{(\text{hyb})}(t)$  and  $T_i(t)$  are deterministic. Thus the right-hand of (14) can be bounded by Gaussian tail inequality (Lemma D.2) as

$$\Pr(\neg E_i^{\theta(\text{hyb})}(t) | \mathcal{F}_t) \leq 2 \exp \left( -\frac{1}{2} \frac{L_t}{T_i(t) + N_i} \cdot (T_i(t) + N_i + 1) \right) \leq \frac{\delta_t}{t}.$$

Then we have

$$\begin{aligned}
 \sum_{t=1}^T \Pr(A(t) = i, \neg E_i^{\mu(\text{hyb})}(t)) &= \sum_{t=1}^T \mathbb{E} \left[ \Pr(A(t) = i, \neg E_i^{\mu(\text{hyb})}(t) | \mathcal{F}_t) \right] \\
 &\leq \sum_{t=1}^T \mathbb{E} \left[ \Pr(\neg E_i^{\mu(\text{hyb})}(t) | \mathcal{F}_t) \right] \\
 &\leq \sum_{t=1}^T \frac{\delta_t}{t} \\
 &= O\left(\frac{1}{\Delta_i^2}\right).
 \end{aligned}$$

□

*Proof of Lemma A.4.*  $\hat{\theta}_i(t) > y_i$  indicates that the median among three indices is above  $y_i$ , which implies that both the online and hybrid indexes are above  $y_i$  based on the event  $E_i^{\mu(\text{on})}(t)$ , i.e.,

$$\{\hat{\theta}_i(t) > y_i\} \subseteq \{\theta_i^{(\text{on})}(t) > y_i\} \cap \{\theta_i^{(\text{hyb})}(t) > y_i\} \text{ when } \mathcal{G}_i(t) \text{ holds.}$$

Therefore

$$\begin{aligned}
 \Pr(\mathcal{E}_{1,t}) &= \Pr(A(t) = i, \mathcal{G}_i(t), \hat{\theta}_i(t) > y_i) \\
 &\leq \min \left\{ \Pr(A(t) = i, \mathcal{G}_i(t), \theta_i^{(\text{on})}(t) > y_i), \Pr(A(t) = i, \mathcal{G}_i(t), \theta_i^{(\text{hyb})}(t) > y_i) \right\}.
 \end{aligned}$$

Then

$$\begin{aligned}
 \sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) &\leq \sum_{t=1}^T \min \left\{ \Pr(A(t) = i, \mathcal{G}_i(t), \theta_i^{(\text{on})}(t) > y_i), \Pr(A(t) = i, \mathcal{G}_i(t), \theta_i^{(\text{hyb})}(t) > y_i) \right\} \\
 &\leq \min \left\{ \underbrace{\sum_{t=1}^T \Pr(A(t) = i, E_i^{\mu(\text{on})}(t), \theta_i^{(\text{on})}(t) > y_i)}_{S_1}, \underbrace{\sum_{t=1}^T \Pr(A(t) = i, \mathcal{G}_i(t), \theta_i^{(\text{hyb})}(t) > y_i)}_{S_2} \right\}. \quad (15)
 \end{aligned}$$

Note that  $S_1$  is exactly the same as the corresponding term in vanilla TS (Agrawal & Goyal, 2017). By Lemma D.5,

$$S_1 \leq \frac{2 \log(T \Delta_i^2)}{(y_i - x_i)^2} + \frac{1}{\Delta_i^2} = \frac{18 \log(T \Delta_i^2)}{\Delta_i^2} + \frac{1}{\Delta_i^2}. \quad (16)$$

Next we bound term  $S_2$ . Based on the good event  $\mathcal{G}_i(t)$  and triangle inequality,

$$\begin{aligned}
 \theta_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})} &= \theta_i^{(\text{hyb})}(t) - \mu_i^{(\text{hyb})} + \mu_i^{(\text{hyb})} - \mu_i^{(\text{on})} \\
 &\leq 2 \sqrt{\frac{L_t}{T_i(t) + N_i}} + \frac{N_i V_i}{T_i(t) + N_i} + \frac{N_i}{T_i(t) + N_i} (\mu_i^{(\text{off})} - \mu_i^{(\text{on})}) \\
 &= 2 \sqrt{\frac{L_t}{T_i(t) + N_i}} + \frac{N_i \omega_i}{T_i(t) + N_i}. \quad (17)
 \end{aligned}$$

Conditional on the event  $\{\theta_i^{(\text{hyb})}(t) > y_i\}$ ,

$$\theta_i^{(\text{hyb})}(t) - \mu_i^{(\text{on})} > y_i - \mu_i^{(\text{on})} = 2\Delta_i/3. \quad (18)$$

Combining (17) and (18),

$$\begin{aligned}
 S_2 &\leq \sum_{t=1}^T \Pr \left( A(t) = i, 2\sqrt{\frac{L_t}{T_i(t) + N_i}} + \frac{N_i \omega_i}{T_i(t) + N_i} \geq \frac{2\Delta_i}{3} \right) \\
 &\leq \sum_{t=1}^T \Pr \left( A(t) = i, 2\sqrt{\frac{L_t}{T_i(t) + N_i}} \geq \frac{\Delta_i}{3} \text{ or } \frac{N_i \omega_i}{T_i(t) + N_i} \geq \frac{\Delta_i}{3} \right) \\
 &= \sum_{t=1}^T \Pr \left( A(t) = i, T_i(t) + N_i \leq \frac{36L_t}{\Delta_i^2} \text{ or } T_i(t) + N_i \leq \frac{3N_i \omega_i}{\Delta_i} \right) \\
 &\leq \sum_{t=1}^T \Pr \left( A(t) = i, T_i(t) + N_i \leq \max \left\{ \frac{36L_t}{\Delta_i^2}, \frac{3N_i \omega_i}{\Delta_i} \right\} \right) \\
 &\leq \left( \max \left\{ 36 \frac{L_T}{\Delta_i^2}, 3 \frac{N_i \omega_i}{\Delta_i} \right\} - N_i \right)_+ \leq \left( 36 \frac{\log(T\Delta_i^2)}{\Delta_i^2} - N_i \left( 1 - 3 \frac{\omega_i}{\Delta_i} \right) \right)_+ \tag{19}
 \end{aligned}$$

Above all, based on (15), (16) and (19),

$$\begin{aligned}
 \sum_{t=1}^T \Pr(\mathcal{E}_{1,t}) &\leq \min \{S_1, S_2\} \\
 &\leq \left( 36 \frac{\log(T\Delta_i^2)}{\Delta_i^2} - N_i \max \left\{ \left( 1 - 3 \frac{\omega_i}{\Delta_i} \right), 0 \right\} \right)_+ + \frac{1}{\Delta_i^2}. \tag{20}
 \end{aligned}$$

□

*Proof of Lemma A.5.* Our goal is to demonstrate

$$\Pr(A(t) = i, \hat{\theta}_i(t) \leq y_i \mid \mathcal{G}_i(t), F_{t-1}) \leq \frac{1 - p_{i,t}}{p_{i,t}} \Pr(A(t) = 1, \hat{\theta}_i(t) \leq y_i \mid \mathcal{G}_i(t), F_{t-1}). \tag{21}$$

For sub-optimal arm  $i$  to be chosen given this constraint, every other arm  $j$  must satisfy  $\hat{\theta}_j(t) \leq \hat{\theta}_i(t) \leq y_i$ . Besides, noting that given instantiation  $F_{t-1}$ , the posterior distribution between optimal arm 1 and other sub-optimal arms are independent, and  $\Pr(\hat{\theta}_1(t) \leq y_i \mid \mathcal{G}_i(t), F_{t-1}) = \Pr(\hat{\theta}_1(t) \leq y_i \mid F_{t-1})$ . Therefore,

$$\begin{aligned}
 \text{LHS of (21)} &\leq \Pr(\hat{\theta}_j(t) \leq y_i, \forall j \mid \mathcal{G}_i(t), F_{t-1}) \\
 &= \Pr(\hat{\theta}_1(t) \leq y_i \mid F_{t-1}) \cdot \Pr(\hat{\theta}_j(t) \leq y_i, \forall j \neq 1 \mid \mathcal{G}_i(t), F_{t-1}) \\
 &= (1 - p_{i,t}) \cdot \Pr(\hat{\theta}_j(t) \leq y_i, \forall j \neq 1 \mid \mathcal{G}_i(t), F_{t-1}).
 \end{aligned}$$

Next, consider the probability of selecting the optimal arm 1 under the same conditions. Arm 1 is chosen if its median index  $\hat{\theta}_1(t)$  exceeds all others. Thus we have

$$\begin{aligned}
 \Pr(A(t) = 1 \mid \mathcal{G}_i(t), F_{t-1}) &\geq \Pr(\hat{\theta}_1(t) > y_i \geq \hat{\theta}_j(t), \forall j \mid \mathcal{G}_i(t), F_{t-1}) \\
 &= \Pr(\hat{\theta}_1(t) > y_i \mid F_{t-1}) \\
 &\quad \cdot \Pr(\hat{\theta}_j(t) \leq y_i, \forall j \neq 1 \mid \mathcal{G}_i(t), F_{t-1}) \\
 &= p_{i,t} \cdot \Pr(\hat{\theta}_j(t) \leq y_i, \forall j \neq 1 \mid \mathcal{G}_i(t), F_{t-1}).
 \end{aligned}$$

Combining the above two inequalities, we get the first result in the lemma.

For the term  $\mathcal{E}_{2,t}$ , taking expectations and summing over  $t$  yields

$$\begin{aligned} \sum_{t=1}^T \Pr(\mathcal{E}_{2,t}) &= \sum_{t=1}^T \mathbb{E} \left[ \Pr(\mathcal{E}_{2,t} \mid \mathcal{F}_{t-1}) \right] \\ &\leq \sum_{t=1}^T \mathbb{E} \left[ \frac{1-p_{i,t}}{p_{i,t}} \Pr(A(t) = 1, \mathcal{G}_i(t) \mid \mathcal{F}_{t-1}) \right] \end{aligned} \quad (22)$$

$$= \sum_{t=1}^T \mathbb{E} \left[ \mathbb{E} \left[ \frac{1-p_{i,t}}{p_{i,t}} \mathbb{1}\{A(t) = 1, \mathcal{G}_i(t)\} \mid \mathcal{F}_{t-1} \right] \right] \quad (23)$$

$$= \sum_{t=1}^T \mathbb{E} \left[ \frac{1-p_{i,t}}{p_{i,t}} \mathbb{1}\{A(t) = 1, \mathcal{G}_i(t)\} \right]. \quad (24)$$

(22) is derived from the first result in the lemma, (23) above uses that  $p_{i,t}$  is fixed given  $\mathcal{F}_{t-1}$ .

Let  $\tau_k$  be the time of the  $k$ -th pull of arm 1 ( $k \geq 1$ ) and set  $\tau_0 := 0$ . Between two consecutive pulls of arm 1, neither its posteriors nor their empirical mean estimations change, so both  $p_{i,t}$  and  $E_1^{\mu(\text{on})}(t)$  is invariant within each block  $\{\tau_k + 1, \dots, \tau_{k+1}\}$ . Summing up the sum in (24) block wise we obtain

$$\begin{aligned} \text{RHS of (24)} &\leq \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1-p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \sum_{t=\tau_k+1}^{\tau_{k+1}} \mathbb{1}\{A(t) = 1, \mathcal{G}_i(t)\} \right] \\ &\leq \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1-p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \right] \\ &\leq \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1-p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\{E_1^{\mu(\text{on})}(\tau_k + 1)\} \right] + \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{1-p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\{\neg E_1^{\mu(\text{on})}(\tau_k + 1)\} \right], \end{aligned} \quad (25)$$

(25) is derived because each block contains at most one round with  $A(t) = 1$  and  $E_1^{\mu(\text{on})}(t)$  coincides with  $E_1^{\mu(\text{on})}(\tau_k + 1)$  throughout the block.  $\square$

*Proof of Lemma A.6.* For any time  $t$ , define  $\text{MIX}_t := \max\{p_{i,t}^{(\text{on})}, p_{i,t}^{(\text{hyb})}\}$ . We first prove that when  $E_1^{\mu(\text{on})}(t)$  holds,

$$p_{i,t} = \Pr(\hat{\theta}_1(t) > y_i \mid \mathcal{F}_{t-1}) \geq \text{MIX}_t. \quad (26)$$

Recall that event  $E_1^{\mu(\text{on})}(t) = \{\hat{\mu}_1^{(\text{on})}(t) > y_i\}$ . This implies with  $E_1^{\mu(\text{on})}(t)$  being true,

$$\left( \{\theta_1^{(\text{hyb})}(t) > y_i\} \cup \{\theta_1^{(\text{on})}(t) > y_i\} \right) \subseteq \{\hat{\theta}_1(t) > y_i\}.$$

Taking conditional probabilities given  $\mathcal{F}_{t-1}$  and using that  $\theta_1^{(\text{on})}(t)$  and  $\theta_1^{(\text{hyb})}(t)$  are independent of  $\hat{\mu}_1^{(\text{on})}(t)$  conditional on  $\mathcal{F}_{t-1}$ , we obtain

$$p_{i,t} \geq \Pr \left( \{\theta_1^{(\text{hyb})}(t) > y_i\} \cup \{\theta_1^{(\text{on})}(t) > y_i\} \mid \mathcal{F}_{t-1} \right) \geq \max\{p_{i,t}^{(\text{on})}, p_{i,t}^{(\text{hyb})}\} = \text{MIX}_t \text{ on } E_1^{\mu(\text{on})}(t).$$

Further, using  $\frac{1-p}{p} = \frac{1}{p} - 1$  and the monotonicity of  $x \mapsto 1/x$  on  $(0, 1]$ , (26) implies that on  $E_1^{\mu(\text{on})}(t)$

$$\frac{1-p_{i,t}}{p_{i,t}} \leq \frac{1-\text{MIX}_t}{\text{MIX}_t} = \frac{1}{\text{MIX}_t} - 1,$$

and multiplying by  $\mathbb{1}\{E_1^{\mu(\text{on})}(t)\}$  gets the point-wise bound

$$\frac{1-p_{i,t}}{p_{i,t}} \mathbb{1}\{E_1^{\mu(\text{on})}(t)\} \leq \left( \frac{1}{\text{MIX}_t} - 1 \right) \mathbb{1}\{E_1^{\mu(\text{on})}(t)\} \leq \frac{1}{\text{MIX}_t} - 1.$$

The last inequality comes from  $\mathbb{1}\{E_1^{\mu(\text{on})}(t)\} \leq 1$ .

Since  $\text{MIX}_t = \max\{p_{i,t}^{(\text{on})}, p_{i,t}^{(\text{hyb})}\}$ ,

$$\frac{1}{\text{MIX}_t} - 1 \leq \min\left\{\frac{1}{p_{i,t}^{(\text{on})}} - 1, \frac{1}{p_{i,t}^{(\text{hyb})}} - 1\right\}. \quad (27)$$

Applying (26)–(27) at time  $\tau_k + 1$ , we obtain that for each  $k$ ,

$$\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\{E_1^{\mu(\text{on})}(\tau_k + 1)\} \leq \min\left\{\frac{1}{p_{i,\tau_k+1}^{(\text{on})}} - 1, \frac{1}{p_{i,\tau_k+1}^{(\text{hyb})}} - 1\right\}, \quad (28)$$

and hence

$$\sum_{k=0}^{T-1} \mathbb{E}\left[\frac{1 - p_{i,\tau_k+1}}{p_{i,\tau_k+1}} \mathbb{1}\{E_1^{\mu(\text{on})}(\tau_k + 1)\}\right] \leq \sum_{k=0}^{T-1} \min\left\{\mathbb{E}\left[\frac{1}{p_{i,\tau_k+1}^{(\text{on})}} - 1\right], \mathbb{E}\left[\frac{1}{p_{i,\tau_k+1}^{(\text{hyb})}} - 1\right]\right\}. \quad (29)$$

□

*Proof of Lemma A.7.* Recall that  $p_{i,t}^{(\text{hyb})}$  denotes the probability that  $\theta_1^{(\text{hyb})}(t)$  exceeds  $y_i$  given  $F_{t-1}$ , and for the algorithm with Gaussian priors we have

$$\theta_1^{(\text{hyb})}(t) \sim \mathcal{N}\left(\hat{\mu}_1^{(\text{hyb})}(t) + \frac{N_1 V_1}{T_1(t) + N_1}, \frac{1}{T_1(t) + N_1 + 1}\right).$$

Given  $F_{\tau_k}$ , let  $\Theta_k$  denote a Gaussian random variable sampled from the above gaussian distribution with  $T_1(t) = k$ . Let  $G_k$  be the geometric random variable denoting the number of consecutive independent trials until a sample of  $\Theta_k$  becomes greater than  $y_i$ . Then

$$p_{i,\tau_k+1} = \Pr(\Theta_k > y_i \mid \mathcal{F}_{\tau_k})$$

and hence

$$\mathbb{E}\left[\frac{1}{p_{i,\tau_k+1}}\right] = \mathbb{E}[\mathbb{E}[G_k \mid \mathcal{F}_{\tau_k}]] = \mathbb{E}[G_k].$$

We first bound  $\mathbb{E}[G_k]$  by a constant for all  $k$ .

Fix any integer  $r \geq 1$ . Let  $z = \sqrt{\ln r}$  and let the random variable  $\text{MAX}_r$  denote the maximum of  $r$  independent samples of  $\Theta_k$ . For brevity, write  $\hat{\mu}_1^{(\text{hyb})} = \hat{\mu}_1^{(\text{hyb})}(\tau_k + 1)$ . Then, for any integer  $r \geq 1$ ,

$$\begin{aligned} \Pr(G_k \leq r) &\geq \Pr(\text{MAX}_r > y_i) \\ &\geq \Pr\left(\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right) \\ &= \mathbb{E}\left[\mathbb{E}\left[\mathbb{1}\left\{\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right\} \mid \mathcal{F}_{\tau_k}\right]\right] \\ &= \mathbb{E}\left[\mathbb{1}\left\{\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right\}\right] \\ &= \Pr\left(\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \mid \mathcal{F}_{\tau_k}\right). \end{aligned} \quad (30)$$

Given  $F_{\tau_k}$ ,  $\Theta_k$  is Gaussian with distribution  $\mathcal{N}\left(\hat{\mu}_1^{(\text{hyb})}(t) + \frac{N_1 V_1}{k + N_1}, 1/(k + N_1 + 1)\right)$ , using a Gaussian tail (Lemma D.2) we have

$$\begin{aligned} \Pr\left(\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \mid \mathcal{F}_{\tau_k}\right) &\geq 1 - \left(1 - \frac{1}{\sqrt{2\pi}} \frac{z}{z^2 + 1} e^{-z^2/2}\right)^r \\ &= 1 - \left(1 - \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\ln r}}{\ln r + 1} \frac{1}{\sqrt{r}}\right)^r. \end{aligned}$$

For sufficiently large  $r$  (in particular, for all  $r \geq e^{11}$ ), we have

$$\left(1 - \frac{1}{\sqrt{2\pi}} \frac{\sqrt{\ln r}}{\ln r + 1} \frac{1}{\sqrt{r}}\right)^r \leq \exp\left(-\frac{\sqrt{r}}{4\pi r \ln r}\right) \leq \frac{1}{r^2},$$

and hence for all  $r \geq e^{11}$ ,

$$\Pr\left(\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{j + N_1 + 1}} + \frac{N_1 V_1}{j + N_1} \mid \mathcal{F}_{\tau_j}\right) \geq 1 - \frac{1}{r^2}. \quad (31)$$

Substituting (31) into (30) yields, for  $r \geq e^{11}$ ,

$$\Pr(G_k \leq r) \geq \left(1 - \frac{1}{r^2}\right) \Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right).$$

Next we obtain a lower bound on the second term using Chernoff–Hoeffding bound (Lemma D.1).

$$\begin{aligned} \Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right) &\geq \Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq \mu_1^{(\text{on})}\right) \\ &\geq \Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} \geq \frac{k \cdot \mu_1^{(\text{on})} + N_1 \mu_1^{(\text{off})}}{k + N_1}\right), \end{aligned}$$

where the last inequality is from a simple decomposition for the hybrid mean like (12).

Applying the Chernoff bound at time  $t = \tau_k + 1$  (so that  $k_1(t) = k$ ), we get, for any  $x > 0$ ,

$$\Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{1}{k + N_1} + \frac{x}{\sqrt{k + N_1 + 1}} \geq \mu_1^{(\text{hyb})}\right) \geq 1 - e^{-2x^2}.$$

Here, the term  $\frac{1}{k + N_1 + 1}$  was added to  $\hat{\mu}_1^{(\text{hyb})}$  to adjust for the fact that  $\hat{\mu}_1^{(\text{hyb})}$  is not simply average of the past  $k + N_1$  samples, instead, it is the sum of past  $k + N_1$  samples divided by  $k + N_1 + 1$ . Now, we use  $x := z - \frac{1}{\sqrt{k + N_1 + 1}} \geq z - \frac{1}{\sqrt{N_1 + 1}}$  for all  $k \geq 0$ , we obtain

$$\begin{aligned} \Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq \mu_1\right) &\geq 1 - \exp\left(-2\left(z - \frac{1}{\sqrt{N_1 + 1}}\right)^2\right) \\ &= 1 - \frac{1}{r^2} \exp\left(\frac{4}{\sqrt{N_1 + 1}} \sqrt{\log r} - \frac{2}{N_1 + 1}\right). \end{aligned}$$

Since  $y_i \leq \mu_1$ , this further implies

$$\Pr\left(\hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} \geq y_i\right) \geq 1 - \frac{1}{r^2} \exp\left(\frac{4}{\sqrt{N_1 + 1}} \sqrt{\log r} - \frac{2}{N_1 + 1}\right). \quad (32)$$

Noting that for  $r \geq \exp\left(\frac{28 + 16\sqrt{3}}{N_1 + 1}\right)$ , we have

$$\frac{1}{r^2} \exp\left(\frac{4}{\sqrt{N_1 + 1}} \sqrt{\log r} - \frac{2}{N_1 + 1}\right) \leq \frac{1}{r^{1.5}} \quad (33)$$

and combining this with (32), (32) and (33) yields, for all  $r \geq \max\left\{e^{11}, \exp\left(\frac{28 + 16\sqrt{3}}{N_1 + 1}\right)\right\}$ ,

$$\Pr(G_k \leq r) \geq 1 - \frac{1}{r^2} - \frac{1}{r^{1.5}}.$$

We can now bound  $\mathbb{E}[G_k]$ :

$$\begin{aligned}
 \mathbb{E}[G_k] &= \sum_{r=0}^{\infty} \Pr(G_k \geq r) \\
 &= 1 + \sum_{r=1}^{\infty} \Pr(G_k \geq r) \\
 &\leq 1 + e^{11} + \sum_{r \geq 1} \left( \frac{1}{r^2} + \frac{1}{r^{1.5}} \right) \\
 &\leq 1 + \max \left\{ e^{11}, \exp \left( \frac{28 + 16\sqrt{3}}{N_1 + 1} \right) \right\} + 2 + 2.7 \\
 &\leq \max \left\{ e^{11}, \exp \left( \frac{28 + 16\sqrt{3}}{N_1 + 1} \right) \right\} + 6.
 \end{aligned} \tag{34}$$

Hence

$$\mathbb{E} \left[ \frac{1}{p_{i, \tau_k + 1}} - 1 \right] = \mathbb{E}[G_k] - 1 \leq \max \left\{ e^{11}, \exp \left( \frac{28 + 16\sqrt{3}}{N_1 + 1} \right) \right\} + 5 \quad \text{for all } k.$$

Next we derive a tighter bound for large  $k$ .

$$k > L_i^{(\text{hyb})}(T) = \left( \frac{288 \ln(T\Delta_i^2 + e^6)}{\Delta_i^2} - N_1 \right)_+.$$

Again fix  $r \geq 1$ , let  $z = \sqrt{\ln r}$ , and define  $\text{MAX}_r$  as before. Then

$$\begin{aligned}
 \Pr(G_k \leq r) &\geq \Pr(\text{MAX}_r > y_i) \\
 &\geq \Pr \left( \text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} - \frac{\Delta_i}{6} \geq y_i \right) \\
 &= \mathbb{E} \left[ \mathbf{1} \left\{ \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} + \frac{\Delta_i}{6} \geq \mu_1^{(\text{on})} \right\} \right]
 \end{aligned} \tag{35}$$

$$\cdot \Pr \left( \text{MAX}_r > \hat{\mu}_1^{(\text{hyb})} + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} - \frac{\Delta_i}{6} \mid \mathcal{F}_{\tau_k} \right), \tag{36}$$

where we used that  $y_i = \mu_1 - \Delta_i/3$ .

By the definition of  $L_i^{(\text{hyb})}(T)$  we have

$$k + N_1 \geq \frac{288 \ln(T\Delta_i^2 + e^6)}{\Delta_i^2},$$

and hence

$$\frac{\sqrt{2 \ln(T\Delta_i^2 + e^6)}}{\sqrt{k + N_1 + 1}} \leq \frac{\Delta_i}{12}.$$

Therefore, for all  $r \leq (T\Delta_i^2 + e^6)^2$ ,

$$\frac{z}{\sqrt{k + N_1 + 1}} - \frac{\Delta_i}{6} = \frac{\sqrt{\log r}}{\sqrt{k + N_1 + 1}} - \frac{\Delta_i}{6} \leq -\frac{\Delta_i}{12}.$$

Using the upper tail bound for a Gaussian random variable (Lemma D.2) we obtain, for any realization  $F_{\tau_k}$ ,

$$\Pr \left( \Theta_k > \hat{\mu}_1^{(\text{hyb})}(\tau_k + 1) + \frac{N_1 V_1}{k + N_1} - \frac{\Delta_i}{12} \mid F_{\tau_k} \right) \geq 1 - \frac{1}{2} \exp \left( -\frac{(k+1)\Delta_i^2}{288} \right) \geq 1 - \frac{1}{2(T\Delta_i^2 + e^{32})},$$

which implies

$$\Pr\left(\text{MAX}_r > \hat{\mu}_1^{(\text{hyb})}(\tau_k + 1) + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} - \frac{\Delta_i}{6} \mid \mathcal{F}_{\tau_k}\right) \geq 1 - \frac{1}{2^r (T \Delta_i^2 + e^6)^r}.$$

Moreover, for any  $t \geq \tau_k + 1$  we have  $T_1(t) \geq k$ , and by Chernoff–Hoeffding bound (Lemma D.1),

$$\begin{aligned} \Pr\left(\hat{\mu}_1^{(\text{hyb})}(t) + \frac{z}{\sqrt{k + N_1 + 1}} + \frac{N_1 V_1}{k + N_1} - \frac{\Delta_i}{6} \geq y_i\right) &\geq \Pr\left(\hat{\mu}_1^{(\text{hyb})}(t) \geq \mu_1^{(\text{hyb})} - \frac{\Delta_i}{6}\right) \\ &\geq 1 - \exp\left(-\frac{2(T_1(t) + N_1)\Delta_i^2}{36}\right) \\ &\geq 1 - \frac{1}{(T \Delta_i^2 + e^6)^{16}}. \end{aligned}$$

Let  $T' = (T \Delta_i^2 + e^6)^2$ . Therefore,

$$\begin{aligned} \mathbb{E}[G_k] &\leq \sum_{r=1}^{\infty} \Pr(G_k \geq r) \\ &\leq 1 + \sum_{r=1}^{T'} \Pr(G_k \geq r) + \sum_{r=T'+1}^{\infty} \Pr(G_k \geq r) \\ &\leq 1 + \sum_{r=1}^{T'} \left(\frac{1}{(2\sqrt{T'})^r} + \frac{1}{(T')^8}\right) + \sum_{r=T'+1}^{\infty} \left(\frac{1}{r^2} + \frac{1}{r^{1.5}}\right) \\ &\leq 1 + \frac{1}{\sqrt{T'}} + \frac{1}{(T')^7} + \frac{2}{T'} + \frac{3}{\sqrt{T'}} \\ &\leq 1 + \frac{5}{T \Delta_i^2 + e^6}. \end{aligned}$$

The above upper bound shows that  $\mathbb{E}\left[\frac{1}{p_{i, \tau_k+1}^{(\text{hyb})}}\right] - 1 = \mathbb{E}[G_k] - 1 \leq \frac{5}{T \Delta_i^2}$  for  $k > L_i^{(\text{hyb})}(T)$ .  $\square$

*Proof of Lemma A.9.* On  $E_1^{\mu}(\text{on}) (\tau_k + 1)$  we have  $\hat{\mu}_1^{(\text{on})}(\tau_k + 1) \leq y_i$ , hence

$$\left\{ \text{median}\{\theta_1^{(\text{on})}(\tau_k+1), \hat{\mu}_1^{(\text{on})}(\tau_k+1), \theta_1^{(\text{hyb})}(\tau_k+1)\} > y_i \right\} = \left\{ \theta_1^{(\text{on})}(\tau_k+1) > y_i, \theta_1^{(\text{hyb})}(\tau_k+1) > y_i \right\}. \quad (37)$$

We denote the right-hand of (37) as

$$p_k^\wedge := \Pr\left(\theta_1^{(\text{on})}(\tau_k+1) > y_i, \theta_1^{(\text{hyb})}(\tau_k+1) > y_i \mid \mathcal{F}_{\tau_k}\right).$$

Conditioned on  $\mathcal{F}_{\tau_k}$ , draw i.i.d. copies

$$\left(\Theta_{k,\ell}^{(\text{on})}, \Theta_{k,\ell}^{(\text{hyb})}\right)_{\ell \geq 1} \stackrel{\text{i.i.d.}}{\sim} \left(\theta_1^{(\text{on})}(\tau_k+1), \theta_1^{(\text{hyb})}(\tau_k+1)\right) \mid \mathcal{F}_{\tau_k},$$

and define the (conditional) geometric hitting time

$$G_k := \min \left\{ \ell \geq 1 : \Theta_{k,\ell}^{(\text{on})} > y_i, \Theta_{k,\ell}^{(\text{hyb})} > y_i \right\}.$$

Then  $G_k \mid \mathcal{F}_k \sim \text{Geom}(p_k^\wedge)$  (supported on  $\{1, 2, \dots\}$ ), so

$$\frac{1}{p_k^\wedge} = \mathbb{E}[G_k \mid \mathcal{F}_k].$$

Fix any  $\alpha \in (1, \frac{3}{2})$  and let  $\beta := \frac{\alpha}{\alpha-1}$ . By Jensen's inequality (since  $x \mapsto x^\alpha$  is convex for  $\alpha > 1$ ),

$$\left(\frac{1}{p_k^\wedge}\right)^\alpha = (\mathbb{E}[G_k | \mathcal{F}_k])^\alpha \leq \mathbb{E}[G_k^\alpha | \mathcal{F}_k], \quad \Rightarrow \quad \mathbb{E}\left[\left(\frac{1}{p_k^\wedge}\right)^\alpha\right] \leq \mathbb{E}[G_k^\alpha].$$

Next we bound  $\mathbb{E}[G_k^\alpha]$  uniformly in  $k$ . For any integer  $r \geq 1$ ,

$$\{G_k > r\} \subseteq \left\{ \max_{1 \leq \ell \leq r} \Theta_{k,\ell}^{(\text{on})} \leq y_i \right\} \cup \left\{ \max_{1 \leq \ell \leq r} \Theta_{k,\ell}^{(\text{hyb})} \leq y_i \right\}.$$

Recall in Lemma D.6 and Lemma A.7 where we prove a constant upper bound for all  $k \geq 0$ , for  $r \geq r_0 = e^{64}$ , we have

$$\begin{aligned} \Pr\left(\max_{1 \leq \ell \leq r} \Theta_{k,\ell}^{(\text{on})} \leq y_i\right) &\leq \frac{1}{r^2} + \frac{1}{r^{3/2}}, \\ \Pr\left(\max_{1 \leq \ell \leq r} \Theta_{k,\ell}^{(\text{hyb})} \leq y_i\right) &\leq \frac{1}{r^2} + \frac{1}{r^{3/2}}. \end{aligned}$$

Therefore,

$$\Pr(G_k > r) \leq \frac{2}{r^2} + \frac{2}{r^{3/2}} \leq \frac{4}{r^{3/2}}.$$

For integer-valued  $X \geq 1$ , we first note that

$$(r+1)^\alpha - r^\alpha \leq \alpha(r+1)^{\alpha-1} \leq \alpha 2^{\alpha-1} r^{\alpha-1},$$

Then we have

$$\begin{aligned} X^\alpha &= \sum_{r=1}^{X-1} ((r+1)^\alpha - r^\alpha) + 1 \\ &\leq \sum_{r=1}^X ((r+1)^\alpha - r^\alpha) \\ &\leq \alpha 2^{\alpha-1} \sum_{r=1}^X r^{\alpha-1} \\ &\leq \alpha 2^{\alpha-1} \sum_{r=1}^{\infty} r^{\alpha-1} \mathbb{1}\{r \leq X\} \\ &= \alpha 2^{\alpha-1} \sum_{r=1}^{\infty} r^{\alpha-1} \mathbb{1}\{X \geq r\} \end{aligned}$$

Taking exception on both sides, let  $C_\alpha = \alpha 2^{\alpha-1}$ , we have

$$\mathbb{E}[X^\alpha] \leq C_\alpha \sum_{r=1}^{\infty} r^{\alpha-1} \mathbb{P}(X \geq r),$$

hence with  $X = G_k$ ,

$$\mathbb{E}[G_k^\alpha] \leq C_\alpha \left( \sum_{r=1}^{r_0-1} r^{\alpha-1} + \sum_{r=r_0}^{\infty} r^{\alpha-1} \cdot \frac{4}{r^{3/2}} \right) = C_\alpha \left( C_{r_0} + 4 \sum_{r=r_0}^{\infty} r^{\alpha-\frac{5}{2}} \right) \leq M_\alpha < \infty,$$

where the last inequality uses  $\alpha < \frac{3}{2}$  so that  $\alpha - \frac{5}{2} < -1$ , and we use  $C_{r_0}$  to denote  $\sum_{r=1}^{r_0-1} r^{\alpha-1}$  as a constant.

Consequently,

$$\sup_{k \geq 0} \mathbb{E} \left[ \left( \frac{1}{p_k^\wedge} \right)^\alpha \right] \leq M_\alpha < \infty. \quad (38)$$

By Hölder's inequality (Lemma D.3) with exponents  $(\alpha, \beta)$ ,

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{p_{i, \tau_k+1}} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] &= \mathbb{E} \left[ \frac{1}{p_k^\wedge} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \\ &\leq \left( \mathbb{E} \left[ \left( \frac{1}{p_k^\wedge} \right)^\alpha \right] \right)^{1/\alpha} \cdot \Pr(\neg E_1^{\mu(\text{on})}(\tau_k + 1))^{1/\beta} \\ &\leq M_\alpha^{1/\alpha} \Pr(\neg E_1^{\mu(\text{on})}(\tau_k + 1))^{1/\beta}. \end{aligned}$$

Noting that, by approximation, taking  $r_0 = e^{64}$ ,  $(C_{r_0})^{1/\alpha} \approx e^{64}$ .

Finally, since at time  $\tau_k + 1$  the online empirical mean of arm 1 is computed from exactly  $k$  i.i.d. Gaussian samples with mean  $\mu_1$  and variance 1, we have the standard tail bound

$$\Pr(\neg E_1^{\mu(\text{on})}(\tau_k + 1)) = \mathbb{P}(\hat{\mu}_{1,k}^{(\text{on})} \leq y_i) \leq \exp\left(-\frac{k(\mu_1 - y_i)^2}{2}\right) = \exp\left(-\frac{k\Delta_i^2}{18}\right),$$

where we used  $y_i = \mu_1 - \Delta_i/3$ . Therefore,

$$\sum_{k=1}^{T-1} \mathbb{E} \left[ \frac{1}{p_{i, \tau_k+1}} \mathbb{1} \left\{ \neg E_1^{\mu(\text{on})}(\tau_k + 1) \right\} \right] \leq M_\alpha^{1/\alpha} \sum_{k=1}^{\infty} \exp\left(-\frac{k\Delta_i^2}{18\beta}\right) \leq \frac{C}{\Delta_i^2}, \quad (39)$$

for some constant  $C > 0$  depending only on  $\alpha$  (hence on  $\beta$ ), which completes the bound for the  $\mathbb{I}\{\neg E_1^{\mu(\text{on})}(\tau_k + 1)\}$ -part of Lemma A.5.  $\square$

## D. Useful Lemmas

**Lemma D.1** (Chernoff-Hoeffding Bound). *Let  $X_1, \dots, X_n$  be independent random variables in  $[0, 1]$  with  $\mathbb{E}[X_i] = \mu_i$  (not necessarily equal). Let  $X = \frac{1}{n} \sum_{i=1}^n X_i$ ,  $\mu = \mathbb{E}[X] = \frac{1}{n} \sum_{i=1}^n \mu_i$ . Then for any  $0 < \epsilon < 1 - \mu$ ,*

$$\Pr(X \geq \mu + \epsilon) \leq e^{-2n\epsilon^2},$$

and, for any  $0 < \epsilon < \mu$ ,

$$\Pr(X \leq \mu - \epsilon) \leq e^{-2n\epsilon^2}.$$

**Lemma D.2** (Concentration inequality of Gaussian variables (Abramowitz et al., 1988)). *For a Gaussian distributed random variable  $Z$  with mean  $\mu$  and variance  $\sigma^2$ , for any  $x > 0$ ,*

$$\Pr(Z > \mu + x\sigma) \geq \frac{1}{\sqrt{2\pi}} \frac{x}{x^2 + 1} e^{-x^2/2}.$$

**Lemma D.3** (Hölder's inequality). *Let  $p, q \in (1, \infty)$  be conjugate exponents such that  $\frac{1}{p} + \frac{1}{q} = 1$ . Let  $X$  and  $Y$  be real-valued random variables defined on the same probability space. If  $\mathbb{E}[|X|^p] < \infty$  and  $\mathbb{E}[|Y|^q] < \infty$ , then*

$$\mathbb{E}[|XY|] \leq (\mathbb{E}[|X|^p])^{1/p} (\mathbb{E}[|Y|^q])^{1/q}.$$

**Lemma D.4** (Lemma 2.15 in (Agrawal & Goyal, 2017)). *For any sub-optimal arm  $i \neq 1$ ,*

$$\sum_{t=1}^T \Pr(A(t) = i, \hat{\mu}_i^{(\text{on})}(t) \leq x_i) \leq \frac{1}{d(x_i, y_i)} + 1 \leq \frac{9}{2\Delta_i^2} + 1,$$

where  $d(a, b) = a \ln \frac{a}{b} + (1-a) \ln \frac{(1-a)}{(1-b)}$ .

1430 **Lemma D.5** (Lemma 2.16 in (Agrawal & Goyal, 2017)). For any sub-optimal arm  $i \neq 1$ ,

1431

1432

1433

1434

$$\sum_{t=1}^T \Pr(A(t) = i, \theta_i^{(\text{on})}(t) > y_i, \hat{\mu}_i^{(\text{on})}(t) \leq x_i) \leq L_i(T) + \frac{1}{\Delta_i^2},$$

1435

1436

where  $L_i(T) \geq \frac{2 \log(T \Delta_i^2)}{(y_i - x_i)^2}$ .

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

$$\mathbb{E} \left[ \frac{1}{p_{i, \tau_j+1}^{(\text{on})}} - 1 \right] \leq \begin{cases} e^{64} + 5 & \forall j, \\ \frac{5}{T \Delta_i^2}, & j > L_i(T), \end{cases}$$

where  $L_i^{(\text{on})}(T) = \frac{288 \log(T \Delta_i^2 + e^{32})}{\Delta_i^2}$ .