

# Efficient Best Arm Identification in Stochastic Bandits: Beyond $\beta$ -optimality

Arpan Mukherjee      Ali Tajer \*

## Abstract

This paper investigates a hitherto unaddressed aspect of best arm identification (BAI) in stochastic multi-armed bandits in the fixed-confidence setting. Two key metrics for assessing bandit algorithms are computational efficiency and performance optimality (e.g., in sample complexity). In stochastic BAI literature, there have been advances in designing algorithms to achieve optimal performance, but they are generally computationally expensive to implement (e.g., optimization-based methods). There also exist approaches with high computational efficiency, but they have provable gaps to the optimal performance (e.g., the  $\beta$ -optimal approaches in top-two methods). This paper introduces a framework and an algorithm for BAI that achieves optimal performance with a computationally efficient set of decision rules. The central process that facilitates this is a routine for sequentially estimating the optimal allocations up to sufficient fidelity. Specifically, these estimates are accurate enough for identifying the best arm (hence, achieving optimality) but not overly accurate to an unnecessary extent that creates excessive computational complexity (hence, maintaining efficiency). Furthermore, the existing relevant literature focuses on the family of exponential distributions. This paper considers a more general setting of any arbitrary family of distributions parameterized by their mean values (under mild regularity conditions). The optimality is established analytically, and numerical evaluations are provided to assess the analytical guarantees and compare the performance with those of the existing ones.

## 1 Introduction

We consider the problem of best arm identification (BAI) in stochastic multi-armed bandits in the fixed-confidence setting. The bandit instances are assumed to be generated by the single-parameter exponential family (SPEF). In BAI, the objective is to identify the *best* arm (i.e., the arm with the largest mean value) within a pre-specified confidence level with the fewest samples. Performance optimality and computational efficiency are the two central metrics for assessing and comparing different bandit algorithms. This paper focuses on an open aspect of BAI in the fixed-confidence parametric setting, which pertains to achieving optimal performance with a computationally efficient algorithm. In reviewing the existing literature, we will specify these two aspects of the existing algorithms. These will furnish the context to highlight that the current literature on stochastic BAI lacks an algorithm that simultaneously achieves optimal performance and maintains low computational complexity. Subsequently, based on this discussion, we will describe our contributions.

**Fixed-confidence versus Fixed-budget.** BAI was first studied as a pure exploration bandit problem in [1]. Subsequently, it has been investigated in two broad settings: the *fixed-confidence* setting and the *fixed-budget* setting. The

---

\*The authors are with the Department of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, Troy, NY 12180.

goal in the fixed-confidence setting is to identify the best arm within a specified guarantee on the decision confidence while using as few samples as possible to arrive at a decision. Representative studies in the fixed-confidence setting include [2–7]. On the other hand, in the fixed-budget setting, the sampling budget is pre-specified. The goal is to minimize the probability of error in the terminal decision. Some representative studies in this setting include [1, 8, 9]. Our focus is on the fixed-confidence setting, the literature on which is discussed next.

**Bayesian versus Non-Bayesian.** BAI in the fixed-confidence setting can be categorized into Bayesian and non-Bayesian models. Bayesian settings assume a prior distribution on the space of parameters and make arm selection decisions based on the posterior distribution computed from the prior and the observed rewards. In contrast, the non-Bayesian settings do not use posterior sampling for arm selection. Top-two sampling was first introduced in [10] for the Bayesian setting. The principle of top-two sampling involves dynamically, over time, identifying a *leader* and a *challenger* as the top arm candidates. Subsequently, the sampling strategy randomizes between these two arms. The top-two Thompson sampling (TTTS) algorithm, proposed and analyzed in [10, 11], involves sampling the posterior for defining the leader and the challenger. Despite the simplicity of TTTS, it faces the computational challenge of repeatedly sampling from the posterior in defining a challenger. To mitigate this, [11] proposed a computationally efficient alternative called the top-two transportation cost (T3C) algorithm. The empirical performance of T3C was further improved using a penalized transportation cost, promoting exploration, in [12]. However, T3C and its improvement only achieve  $\beta$ -optimality. The nature of a  $\beta$ -optimality guarantee is as follows: if a  $\beta$  fraction of the sampling resources are reserved for the top arm, then these algorithms can determine how to optimally allocate the remaining  $(1 - \beta)$  fraction among the rest of the arms. Hence, these algorithms are said to be only  $\beta$ -optimal, where  $\beta \in (0, 1)$ .

In the non-Bayesian setting, [4] has proposed the track-and-stop (TaS) algorithm for BAI, with optimal performance in the asymptote of diminishing probability of error. More investigations on TaS-based algorithms include [13] and [14]. The TaS sampling strategies, in general, hinge on tracking the optimal allocation of sampling resources over time. Maintaining such allocation in a bandit setting with  $K$  arms necessitates solving  $K$  equations at each time using the bisection method. Hence, these approaches are generally computationally expensive. It was shown in [14] that the tracking procedure could only be performed intermittently at exponentially spaced intervals. This reduces the computational complexity for TaS. Nevertheless, the approach of [14] applies to only linear bandits with Gaussian noise. For stochastic bandits, an asymptotically optimal and computationally efficient alternative was proposed in [15], which is based on the lazy sub-gradient ascent algorithm for arm selection. The results of this study on BAI are limited to only Gaussian settings.

To address the computational challenge, the gamification approach was proposed in [16, 17]. In this approach, BAI is viewed as an unknown two-player game comprising a  $w$  player and a  $\lambda$  player. While the  $w$  player samples a distribution from the probability simplex consisting of possible allocations, the  $\lambda$  player chooses a corresponding bandit instance from the class of instances having a different best arm, with the objective of converging to a saddle point. More recently, a Frank-Wolfe-based algorithm was proposed in [18], which solves BAI using a single iteration of the Frank-Wolfe algorithm. The implementation of this method involves a two-player zero-sum game in each iteration, which requires solving a linear program. To further reduce the computational complexity, the top-two approach has also been used to design algorithms for the non-Bayesian setting. Representative top-two non-Bayesian algorithms include top-two sequential probability ratio test (TT-SPRT) [7, 12, 19], top-two expected improvement (TTEI) [20], and the empirical best leader with an improved transportation cost (EB-TCI) [12]. While efficient, these algorithms achieve optimal performance only when an instance-dependent parameter  $\beta$  is known a priori. Specifically, at their

core, these algorithms identify an optimal allocation of the sampling resources among the arms. These algorithms enjoy  $\beta$ -optimality guarantees. We will discuss and demonstrate empirically that, in some settings, the choice of  $\beta$  can critically affect performance (sample complexity).

**Parametric versus Non-parametric.** Algorithms for BAI can also be categorized based on whether the bandit instance follows *parametric* or *non-parametric* families of distributions. In the case of the non-parametric family, the upper confidence bound (UCB) based approaches have been investigated (e.g., [2, 5]) and shown to be optimal up to constant factors for the family of sub-Gaussian bandits. More recently, the study in [13] has considered the class of distributions satisfying a specific functional boundedness property. This study proposes a tracking-based sampling strategy along with a likelihood ratio-based stopping rule, which was shown to be asymptotically optimal for this specified class of distributions. For parametric bandits, investigations have focused on the single parameter exponential family (e.g., [4, 7, 10, 12, 21]). Some investigations have considered the case of Gaussian bandits with known variances and unknown means, e.g., [11, 20]. In both of these settings, the top-two sampling strategy is  $\beta$ -optimal, while the TaS algorithm is asymptotically optimal, despite being computationally expensive.

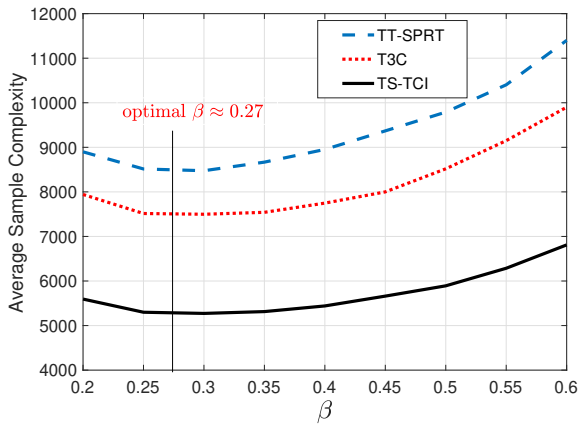


Figure 1: Sample complexity versus  $\beta$ .

**Contributions.** This paper is focused on fixed-confidence, non-Bayesian, and parametric settings. Hence, the directly relevant literature to this scope includes [4, 7, 11, 12, 20]. As discussed, these studies either achieve optimality at the expense of high computational complexity (e.g., [4, 18]) or maintain computational efficiency but exhibit optimality gaps (e.g., [7, 10–12, 20]). We propose an algorithm, referred to as the *transportation cost balancing (TCB)* algorithm, that achieves the optimal sample complexity with computationally efficient sampling and decision rules. Specifically, the TCB algorithm avoids solving an optimization problem to form the arm selection decisions in the next round. Furthermore, compared to the existing class of efficient top-two algorithms [7, 10–12, 20], TCB exhibits optimality in contrast to  $\beta$ -optimality. Leaping from  $\beta$ -optimality to optimality, in some bandit instances, leads to significant improvement in sample complexity. To showcase the gap between a  $\beta$ -optimal solution and an optimal solution, in Figure 1 we demonstrate how the sample complexity of the  $\beta$ -optimal solutions vary with respect to  $\beta$ , demonstrating (i) the sensitivity of the sample complexity to  $\beta$ , and (ii) a substantial gap (e.g., an order of magnitude) between the optimal and the  $\beta$ -optimal guarantees. In the face of not knowing the optimal choice of  $\beta$ ,  $\beta = 0.5$  has

been prescribed as a reasonable choice for the top-two algorithms [10]. However, as expected, this can be noticeably different from the optimal choice in some instances. For instance, Figure 1 empirically shows the sub-optimality of choosing  $\beta = 0.5$  in three existing approaches that ensure  $\beta$ -optimality. In all these cases, the optimal value of  $\beta \approx 0.27$ .

As the second contribution, we also generalize the probability models to any arbitrary class of parametric models (that satisfy certain regularity conditions). These models subsume the exponential family, which is the only parametric class for which algorithms and performance guarantees are available in the literature. For this generalization, we propose a novel concentration inequality for the generalized log-likelihood ratio (GLLR)-based test statistic for BAI [4, 21] that holds for any general parametric bandit instance that satisfies a uniform continuity assumption on the divergence measure of the model, and some mild regularity conditions on the arm distributions.

**Methodology: Transport Cost Balancing (TCB).** For fixed-confidence BAI, for any bandit instance  $\nu$ , the universal lower bound on the average sample complexity is inversely proportional to a problem complexity measure  $\Gamma(\nu)$  (specified later in (21)). Achieving this lower bound is predicated on sequentially determining an optimal allocation of the sampling resources among the arms. While existing optimal algorithms compute these optimal allocations, it is computationally expensive. To mitigate this computational challenge, the top-two sampling rules assign a sampling proportion  $\beta \in (0, 1)$  to the best arm and then determine the allocation of the remaining  $(1 - \beta)$  fraction over the rest. The analysis of the top-two methods shows that this facilitates convergence to the  $\beta$ -optimal allocation. In these methods, the sample complexity is inversely proportional to the transportation cost  $\Gamma_\beta(\nu)$ , where  $\Gamma(\nu) \geq \Gamma_\beta(\nu)$ . Equality holds at the optimal value  $\beta^*$ , which depends on the bandit instance and is unknown a priori.

The central process in our algorithm is a routine for estimating the optimal allocations, including the optimal value of  $\beta$ , up to a sufficient fidelity that enables confidently identifying the best arm. Our algorithm guides its decisions by balancing transportation costs over time. These decisions lead to efficiently estimating the optimal sampling proportion up to sufficient fidelity. The fundamental advantage of our arm selection rules is that they can track the sampling proportions *without* having to compute the optimal sampling proportions at each round. Instead, the sampling proportion is estimated by sampling from the set of under-sampled arms in each round.

## 2 Stochastic BAI Model and Assumption

**Stochastic Model.** Denote the class of probability measures defined on any sample space  $\Omega \subseteq \mathbb{R}$  by  $\mathcal{Q}(\Omega)$ . Let  $\mathcal{P}(\Omega) \subset \mathcal{Q}(\Omega)$  denote the class of probability measures that are parameterized by their mean values, i.e.,

$$\mathcal{P}(\Omega) \triangleq \{\mathbb{P} \in \mathcal{Q}(\Omega) : m(\mathbb{P}) \in \Theta\}, \quad (1)$$

where  $m(\mathbb{P}) \triangleq \mathbb{E}_{\mathbb{P}}[X]$ , and  $\mathbb{E}_{\mathbb{P}}$  denotes the expectation under measure  $\mathbb{P}$ . Furthermore, define  $\mathcal{M} \triangleq \mathcal{P}^{\otimes K}(\Omega)$  as the Cartesian product of  $K$  sets of measures in  $\mathcal{P}(\Omega)$ . We denote the likelihood function associated with measure  $\mathbb{P}$  by  $\pi_{\mathbb{P}}$ . We make the following assumptions on this stochastic model.

1.  $\Theta \subseteq \mathbb{R}$  is a compact parameter space, which is *known* to the learner.
2. The likelihood functions  $\pi_{\mathbb{P}}$  are continuous and twice-differentiable in  $m(\mathbb{P})$  for every  $\mathbb{P} \in \mathcal{P}(\Omega)$ . Furthermore, the log-likelihood function  $\log \pi_{\mathbb{P}}(\cdot | m(\mathbb{P}))$  is concave in  $m(\mathbb{P})$  for any  $m(\mathbb{P}) \in \Theta$ .

3. All distributions in  $\mathcal{P}(\Omega)$  have the same support  $\Omega$ .
4. All the distributions in  $\mathcal{P}(\Omega)$  have *finite* third moments, i.e.,  $\mathbb{E}_{\mathbb{P}}[|X - m(\mathbb{P})|^3] < +\infty$  for every  $\mathbb{P} \in \mathcal{P}(\Omega)$ .
5. For any  $\theta, \theta' \in \Theta$ , we denote the Kullback-Leibler (KL) divergence between  $\mathbb{P}_{\theta} \in \mathcal{P}(\Omega)$  and  $\mathbb{P}_{\theta'} \in \mathcal{P}(\Omega)$  by  $D_{\text{KL}}(\mathbb{P}_{\theta} \parallel \mathbb{P}_{\theta'})$ . Keeping one argument fixed, the KL divergence is assumed to be uniformly continuous in the second argument.
6. For any likelihood function  $\pi_{\mathbb{P}}$ , let us define the Fisher Information (FI) measure as

$$\mathcal{I}_{\mathbb{P}}(\theta) \triangleq -\mathbb{E}_{\mathbb{P}} \left[ \frac{\partial^2}{\partial \theta^2} \log \pi_{\mathbb{P}}(\cdot | \theta) \right]. \quad (2)$$

We assume that  $\mathcal{I}_{\mathbb{P}}(\theta) < +\infty$  for all  $\theta \in \Theta$  and  $\mathbb{P} \in \mathcal{P}(\Omega)$ .

7. There exists  $\sigma^2 > 0$  such that for any  $\mathbb{P} \in \mathcal{P}(\Omega)$ ,

$$\frac{\partial^2 \log \pi_{\mathbb{P}}(x | \theta)}{\partial \theta^2} \leq -\sigma^2, \quad \forall x \in \Omega, \forall \theta \in \Theta. \quad (3)$$

8. For any  $\theta, \theta' \in \Theta$ , we assume that

$$\mathbb{E}_{\mathbb{P}} \left[ \left| \log \frac{\pi_{\mathbb{P}}(X | \theta)}{\pi_{\mathbb{P}}(X | \theta')} \right|^3 \right] < +\infty. \quad (4)$$

Assumption 1 is needed for designing the estimator for the unknown arm means from the observed rewards and is a common assumption in the BAI literature [16, 18]. Assumptions 2 and 3 ensure the existence of the maximum likelihood estimates (MLEs) for the parameters of interest (i.e., the mean values) and that the KL divergence measures between any pair of distributions in  $\mathcal{M}$  are finite. Assumption 4 ensures the almost sure convergence of the sample mean to the ground truth if each arm is sampled sufficiently often. Assumption 5 depicts the uniform continuity of the KL divergence in each argument. We emphasize that this is a mild assumption, and BAI becomes significantly hard without this assumption. Specifically, there exists an impossibility result [13], which states that BAI is impossible for classes of measures that are “KL right dense”, i.e., if the mean values between two measures in the class can be arbitrarily large, even though the KL divergence between the measures is bounded by an arbitrarily small value. Assumption 6 implies the finiteness of the FI measure for all distributions in  $\mathcal{P}(\Omega)$ . The condition in (3) specifies a bound on the second-order derivative of the log-likelihood function over the parameter space  $\Theta$ . For instance, for the exponential family of distributions, (3) is equivalent to the variance being bounded away from zero. Finally, the condition in (4) implies that the log-likelihood ratio corresponding to parameters  $\theta$  and  $\theta'$  has a finite third moment.

**Bandit Model.** Consider a  $K$ -armed stochastic bandit. The rewards of arm  $i \in [K] \triangleq \{1, \dots, K\}$  are generated from  $\mathbb{P}_i \in \mathcal{P}(\Omega)$ . We define  $\mu(i) \triangleq m(\mathbb{P}_i)$ , and denote the likelihood function associated with  $\mathbb{P}_i$  parameterized by the mean value  $\mu(i)$  by  $\pi_i(\cdot | \mu(i))$ . Accordingly, we define the bandit instance  $\nu \triangleq [\mathbb{P}_1, \dots, \mathbb{P}_K]$ . Furthermore, define  $\mathcal{M} \triangleq \mathcal{P}^{\otimes K}(\Omega)$  as the Cartesian product of  $K$  sets of measures in  $\mathcal{P}(\Omega)$ . Finally, let  $(\mathcal{M}, D_{\text{TV}})$  denote the metric space of the distributions in the set  $\mathcal{M}$  endowed with the total variation distance metric  $D_{\text{TV}}$ .

**Sequential Decisions.** At each round  $t \in \mathbb{N}$ , the learner chooses an action  $A_t \in [K]$ , and receives a reward  $X_t \sim \mathbb{P}_{A_t}$ . We denote the sequence of actions, the corresponding rewards, and the filtration generated by the sequence of actions and rewards by the ordered sets

$$\begin{aligned} \mathcal{A}_t &\triangleq \{A_s : s \in [t]\}, \\ \mathcal{X}_t &\triangleq \{X_s : s \in [t]\}, \\ \text{and } \mathcal{F}_t &\triangleq \{A_1, X_1, \dots, A_t, X_t\}. \end{aligned} \quad (5)$$

Denote the set of rewards obtained by selecting an arm  $i \in [K]$  up till time  $t$  by

$$\mathcal{X}_t^i \triangleq \{X_s : s \in [t], A_s = i\}. \quad (6)$$

The objective of the learner is to identify the best arm, which is assumed to be *unique*, and is defined as the arm with the largest mean, i.e.,

$$a^* \triangleq \arg \max_{i \in [K]} \mu(i). \quad (7)$$

For the algorithm design, we use information projection measures defined as follows. For any measure  $\mathbb{P} \in \mathcal{P}(\Omega)$  and  $x \in \mathbb{R}$ , we define

$$d_U(\mathbb{P}, x) \triangleq \inf_{\mathbb{Q} \in \mathcal{P}(\Omega): m(\mathbb{Q}) \leq x} D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}), \quad (8)$$

$$\text{and } d_L(\mathbb{P}, x) \triangleq \inf_{\mathbb{Q} \in \mathcal{P}(\Omega): m(\mathbb{Q}) \geq x} D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}). \quad (9)$$

Specifically, the information measure  $d_U(\mathbb{P}, x)$  for any distribution  $\mathbb{P} \in \mathcal{P}(\Omega)$  and  $x \in \mathbb{R}$  is the minimum KL divergence between  $\mathbb{P}$  and any distribution with mean *at most*  $x$ . Similarly,  $d_L(\mathbb{P}, x)$  measures the KL divergence between  $\mathbb{P}$  and any distribution with mean *at least*  $x$ . In the fixed confidence setting, the goal is to identify the best arm with a pre-specified level of confidence while minimizing the number of samples in making the decision. Let  $\tau$  denote an  $\mathcal{F}$ -adapted stopping time, i.e.,  $\{\tau = t\} \in \mathcal{F}_t$  for every  $t \in \mathbb{N}$ . Corresponding to the stochastic stopping time  $\tau$ , let  $\hat{A}_\tau$  denote the terminal decision of the learner. The  $\delta$ -PAC objective of the learner is formalized next.

**Definition 1** ( $\delta$ -PAC). *A BAI algorithm is  $\delta$ -PAC if the algorithm has a stopping time  $\tau$  adapted to  $\{\mathcal{F}_t : t \in \mathbb{N}\}$ , and at the stopping time with the terminal decision  $\hat{A}_\tau \in [K]$  it ensures*

$$\mathbb{P}_\nu \{\tau < +\infty, \hat{A}_\tau = a^*\} > 1 - \delta, \quad (10)$$

where  $\mathbb{P}_\nu$  denotes the probability measure induced by the interaction of the BAI algorithm with the bandit instance  $\nu$ .

### 3 Transportation Cost Balancing Algorithm

In this section, we specify (i) a stopping rule that decides when to stop collecting samples and form a confident decision about the best arm, (ii) an arm selection rule that guides the order of sampling arms over time, and (iii) an estimation routine that aims to learn the unknown model parameters of interest (e.g., the mean values). We use different estimators for the arm selection and stopping rules. Specifically, we use the maximum likelihood estimate (MLE) for the stopping rule, whereas we use the sample mean to estimate the arm selection strategy. The central statistic that

guides all three decisions in the TCB algorithm is the GLLR. To formalize the GLLRs, we define  $K(K-1)$  hypotheses  $\{\mathcal{H}_{i,j} : i \in [K]\}$  such that for all  $i \neq j$

$$\mathcal{H}_{i,j} : \mu(i) \geq \mu(j) . \quad (11)$$

Next, we formalize the GLLR test statistic for performing these hypothesis tests, which has been adopted in a wide range of investigations on parametric BAI [4, 7, 11, 16, 18, 20, 21]. At any time  $t \in \mathbb{N}$  and for any arm  $i \in [K]$ , and based on the samples available from this arm, we denote the MLE of  $\mu(i)$  projected on  $\Theta$  by  $\mu_t(i)$ , i.e.,

$$\mu_t(i) \triangleq \arg \max_{\mu \in \Theta} \sum_{s \in [t]} \log \pi_i(X_s | \mu) \cdot \mathbb{1}_{\{A_s=i\}} , \quad (12)$$

where  $\mathbb{1}$  denotes the indicator function. Furthermore, for any two parameters  $\theta, \theta' \in \Theta$ , and for any arm  $i \in [K]$ , let us define

$$d_i(\theta || \theta') \triangleq D_{\text{KL}}(\pi_i(\cdot | \theta) || \pi_i(\cdot | \theta')) . \quad (13)$$

Accordingly, for any pair of arms  $(i, j) \in [K] \times [K]$ , let us define

$$\Lambda_t(i, j) \triangleq \min_{\rho \in \mathbb{R}^K : \rho(i) \leq \rho(j)} \{T_t(i)d_i(\mu_t(i) || \rho(i)) + T_t(j)d_j(\mu_t(j) || \rho(j))\} \mathbb{1}_{\{\mu_t(i) \geq \mu_t(j)\}} , \quad (14)$$

where we have defined

$$T_t(i) \triangleq \sum_{s=1}^t \mathbb{1}_{\{A_s = i\}} , \quad (15)$$

as the counter for the number of times arm  $i \in [K]$  is chosen up to time  $t$ .

**GLLR-based Stopping Rule.** We specify a GLLR-thresholding stopping criterion, which compares the GLLR statistic against a time-varying threshold. When the GLLR statistic exceeds the threshold, the algorithm stops collecting more samples and forms a decision about the best arm. Let us denote the maximum likelihood (ML) decision at time  $t$  about the top arm by  $a_t^{\text{top}}$ , i.e.,

$$a_t^{\text{top}} \in \arg \max_{i \in [K]} \bar{\mu}_t(i) , \quad (16)$$

where  $\bar{\mu}_t(i)$  denotes the sample mean of arm  $i \in [K]$ , projected on to the parameter space  $\Theta$ . Our stopping criterion is based on sufficiently distinguishing between the ML decision  $a_t^{\text{top}}$  and the most likely contender compared to the best arm, which we refer to as the *challenger*. This challenger arm at time  $t$  is the arm closest to the top arm in a GLLR sense, and it is specified by

$$a_t^{\text{ch}} \in \arg \min_{i \in [K] \setminus \{i : \mu_t(i) < \mu_t(a_t^{\text{top}})\}} \Lambda_t(a_t^{\text{top}}, i) . \quad (17)$$

In other words,  $a_t^{\text{ch}}$  denotes the arm that is the top contender to the best arm  $a_t^{\text{top}}$ , where the comparison is made using the likelihood ratio between the arms. The stopping rule compares the GLLR between  $a_t^{\text{top}}$  and  $a_t^{\text{ch}}$  and stops collecting samples when the GLLR exceeds a threshold. The threshold depends on the level of confidence  $\delta$  required on the final decision, and it is denoted by  $\beta_t(\delta)$ . The stopping rule is stated next.

$$\tau \triangleq \inf \{t \in \mathbb{N} : \Lambda_t(a_t^{\text{top}}, a_t^{\text{ch}}) > \beta_t(\delta)\} . \quad (18)$$

The threshold  $\beta_t(\delta)$  is specified in Theorem 1 to ensure the  $\delta$ -PAC guarantee on the decision. Next, we delineate the arm selection rules. For this purpose, we first formalize the *problem complexity* measure, which quantifies the hardness of identification in a BAI instance. Specifically, the problem complexity captures the *minimum* distance between the given bandit instance and any other bandit instance with a different best arm. Clearly, the smaller the problem complexity, the larger the number of samples required for identification. We also state an equivalent representation of the problem complexity, which motivates our arm selection strategies.

**Problem complexity.** Consider a bandit instance  $\nu = [\mathbb{P}_1, \dots, \mathbb{P}_K]$  with the top arm  $a^*$ . Given  $a^*$ , we define an alternative set of bandit instances  $\bar{\nu} = [\bar{\mathbb{P}}_1, \dots, \bar{\mathbb{P}}_K]$  in which the top arm is not  $a^*$ . Specifically,

$$\text{alt}(a^*) \triangleq \{\bar{\nu} \in \mathcal{M} : m(\bar{\mathbb{P}}_{a^*}) \leq \max_{i \neq a^*} m(\bar{\mathbb{P}}_i)\}. \quad (19)$$

Subsequently, given  $a^*$  and  $\text{alt}(a^*)$ , under the weight vector  $\mathbf{w} = [w_1, \dots, w_K] \in \Delta^K$ , where  $\Delta^K$  represents the  $K$ -dimensional probability simplex, we define the problem complexity associated with  $\nu$  as the smallest weighted KL divergence from  $\nu$  to the set  $\text{alt}(a^*)$ . Specifically,

$$\Gamma(\nu, \mathbf{w}) \triangleq \inf_{\bar{\nu} \in \text{alt}(a^*)} \sum_{i \in [K]} w_i D_{\text{KL}}(\mathbb{P}_i \| \bar{\mathbb{P}}_i). \quad (20)$$

Given any weight vector  $\mathbf{w} \in \Delta^K$ ,  $\Gamma(\nu, \mathbf{w})$  captures the hardness of distinguishing  $\nu$  from the closest alternate bandit instance, where the divergence between the arms is weighted by  $\mathbf{w}$ . Finally, we define the problem complexity associated with the bandit instance  $\nu$  as

$$\Gamma(\nu) \triangleq \sup_{\mathbf{w} \in \Delta^K} \Gamma(\nu, \mathbf{w}). \quad (21)$$

It can be readily verified that  $\Gamma(\nu)$  captures the *maximum* hardness in distinguishing  $\nu$  from the closest bandit instance. Accordingly, we define the maximizer weight vector as

$$\mathbf{w}(\nu) \triangleq \arg \sup_{\mathbf{w} \in \Delta^K} \Gamma(\nu, \mathbf{w}). \quad (22)$$

The weight vector  $\mathbf{w}(\nu)$  characterizes the *optimal* allocation in which to sample arms, such that the sample complexity for BAI for the bandit instance  $\nu$  is minimized. Next, we provide an equivalent representation of the problem complexity measure, which facilitates analyzing its key properties.

**Lemma 1.** *The problem complexity  $\Gamma(\nu, \mathbf{w})$  defined in (20) can be equivalently expressed as*

$$\Gamma(\nu, \mathbf{w}) = \min_{i \neq a^*} \inf_{x \in [\mu(i), \mu(a^*)]} \left\{ w_{a^*} d_{\text{U}}(\mathbb{P}_{a^*}, x) + w_i d_{\text{L}}(\mathbb{P}_i, x) \right\}. \quad (23)$$

*Proof.* The proof follows a similar line of arguments as [4, Lemma 3]. For completeness, we provide the proof in Appendix B.1. ■

The expression for  $\Gamma(\nu, \mathbf{w})$  in (23) involves an inner minimization, which is a weighted combination of divergence measures from  $\nu$  to an alternate bandit instance. Note that the inner minimization only depends on the divergence measures for the best arm  $a^*$  and any other arm  $i \neq a^*$ . Furthermore, the outer minimization acts over all the other arms  $i \neq a^*$ , establishing that  $\Gamma(\nu, \mathbf{w})$  is the weighted divergence measure between the bandit instance  $\nu$  and the *closest* alternate bandit instance.



**Transportation cost balancing (TCB).** Designing the arm selection rule consists of two key components. The first component ensures that none of the arms remain under-explored. Specifically, the objective in this phase is to ensure that we have a reasonable estimate of each arm’s mean value, such that our estimates converge to the true mean values if the arm selection rule is allowed to collect samples without stopping. While estimating the mean values is not the goal in BAI, our arm selection rule performs estimation as an intermediate step to form a confident decision about the best arm. The second component of the arm selection rule is to track the optimal proportion  $\mathbf{w}(\boldsymbol{\nu})$  of arm selections defined in (22), which ensures that we minimize the average sample complexity. For this, TaS [4] proposes to compute the optimal sampling proportions at the current mean estimates and track the estimated sampling proportions. However, this is computationally expensive and requires solving  $K$  equations using the bisection method in each round. To circumvent this, we propose a simple sampling mechanism focusing on sampling from the set of under-sampled arms at each instant, in which case, the sampling strategy can converge to the optimal sampling proportions  $\mathbf{w}(\boldsymbol{\nu})$  asymptotically. Next, we describe the sampling rule that combines these two components.

**Under-explored Arms.** At any time  $t$ , the sampling rule defines a set of *under-explored* arms as

$$\mathcal{U}_t \triangleq \left\{ i \in [K] : T_t(i) \leq \left\lceil \sqrt{t/K} \right\rceil \right\}. \quad (24)$$

If the set of under-explored arms is non-empty, indicating that some of the arms are under-explored, the arm selection strategy selects the arm that is sampled the least. Otherwise, when there are no under-explored arms to sample, the goal is to devise a sampling strategy that *estimates* the optimal  $\beta$ , i.e., the optimal allocation of the best arm. For this purpose, our sampling strategy aims to ensure the almost sure convergence of the sampling proportions to the optimal allocation  $\mathbf{w}(\boldsymbol{\nu})$ . We show in Lemma 11 (Appendix C) that this objective is achieved by any sampling rule that eventually *always* samples from the set of under-sampled arms, i.e., the set of arms which are sampled less number of times compared to the optimal allocation. Consequently, we devise a sampling strategy that (eventually) always samples from the set of under-sampled arms. At time  $t$ , if  $\mathcal{U}_t = \emptyset$ , this arm selection rule leverages the MLEs of the arm means to compute an empirical estimate of  $\Gamma(\boldsymbol{\nu}, \mathbf{w})$  defined in Lemma 1. Based on this, the goal is to select the next arm in a way that maximizes the estimate. We show in Appendix C that this is equivalent to sampling from the set of under-sampled arms.

**Transport cost-based Estimation of Allocation.** We provide a few measures that are used to delineate the arm selection rule. Let us denote the distribution of arm  $i \in [K]$  parameterized by the sample mean  $\bar{\mu}_t(i)$  projected on  $\Theta$  by  $\mathbb{P}_{t,i}$ . For any arm  $i \in [K]$ , define the interval  $I_{t,i} \triangleq [\bar{\mu}_t(i), \bar{\mu}_t(a_t^{\text{top}})]$ , which specifies the interval for minimization in the empirical problem complexity defined in (25) next. Based on these, we define the *minimum transportation cost* [11], as the minimum weighted combination of divergence measures  $d_U$  and  $d_L$  defined in (8) and (9) of the top arm  $a_t^{\text{top}}$  and any other arm  $i \neq a_t^{\text{top}}$  as follows.

$$\Gamma_t(\mathbf{w}) \triangleq \min_{i \in [K] \setminus \{a_t^{\text{top}}\}} \min_{x \in I_{t,i}} \left\{ w_{a_t^{\text{top}}} d_U(\mathbb{P}_{t,a_t^{\text{top}}}, x) + w_i d_L(\mathbb{P}_{t,i}, x) \right\}. \quad (25)$$

$\Gamma_t(\mathbf{w})$  in (25) specifies the minimum cost of transporting the currently estimated bandit instance  $\boldsymbol{\nu}_t \triangleq [\mathbb{P}_{t,1}, \dots, \mathbb{P}_{t,K}]$  to an alternate bandit instance for which the best arm is not  $a_t^{\text{top}}$ . Even though we optimize for the information projection measures  $d_U$  and  $d_L$  in  $\mathcal{P}(\Omega)$ , it is worth noting that we may have additional knowledge about the class of bandit instances over which we want to optimize. For example, we may restrict  $\mathcal{P}(\Omega)$  to be the set of distributions in

the single parameter exponential family with the cumulant generating function  $b : \Theta \mapsto \mathbb{R}$ , in which case, we have an explicit closed-form expression for  $\Gamma_t(\mathbf{w})$  [7, 11, 12]. However, our analysis holds for a general class of bandit instances satisfying Assumptions 1-8, and a closed-form expression for  $\Gamma_t(\mathbf{w})$  can be derived based on the bandit instance in consideration. The TCB arm selection rule is based on a *look-ahead* distribution, which navigates the arm selection routine to sample arms in a way that increases the empirical problem complexity. To this end, let us define the look-ahead distribution over the arms  $\mathbf{w}'_t \in \Delta^K$  such that

$$w'_t(i) \triangleq \begin{cases} \frac{T_t(i)}{t} + r_t, & \text{if } i = a_t^{\text{top}} \\ \frac{T_t(i)}{t} - \frac{r_t}{K-1}, & \text{otherwise} \end{cases}, \quad (26)$$

where  $\{r_t : t \in \mathbb{N}\}$  is a sequence of positive real numbers satisfying  $\limsup_{t \uparrow \infty} r_t = 0$ , such that it preserves the property that  $\mathbf{w}'_t \in \Delta^K$ .

**Remark 1.** Note that any choice of the sequence  $\{r_t : t \in \mathbb{N}\}$  satisfying  $\limsup_{t \uparrow \infty} r_t = 0$  is sufficient for the performance guarantees presented in Section 4. However, it is possible that different choices of  $\{r_t : t \in \mathbb{N}\}$  promote convergence in empirical allocations to the optimal one at different rates. However, in this investigation, we are not interested in the rate of convergence in the allocation estimates. Rather, we require the estimates to be sufficiently close to the optimal values at stopping, which ensures  $\delta$ -PAC BAI. Furthermore, since we investigate the asymptotic sample complexity in Theorems 4, it suffices for the sequence  $\{r_t : t \in \mathbb{N}\}$  to converge to 0 asymptotically in  $t$ .

Based on the equivalent form of the problem complexity in Lemma 1, the goal of the sampling strategy is to maximize the lowest information measure  $\Gamma_t(\mathbf{w})$ . This ensures that we sample the under-sampled arms in each round and move closer toward the optimal sampling proportion. To this end, let us define the arm with the lowest information measure as

$$a_t^{\min} \in \arg \min_{i \in [K] \setminus \{i: \mu_t(i) < \mu_t(a_t^{\text{top}})\}} \min_{x \in I_{t,i}} \left\{ \frac{T_t(a_t^{\text{top}})}{t} d_{\text{U}}(\mathbb{P}_{t, a_t^{\text{top}}}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\}. \quad (27)$$

In (27),  $a_t^{\min}$  denotes the arm that minimizes the estimate of the problem complexity at time  $t$ . We note that the arms  $a_t^{\min}$  and  $a_t^{\text{ch}}$  are generally distinct. Even though they occasionally might refer to the same arm, it is not the case. Specifically, given a class of measures, the transportation cost evaluated at the current sampling proportions and the scaled GLLR  $\frac{1}{t} \Lambda_t(a_t^{\text{top}}, a_t^{\text{ch}})$  defined in (14) may have a similar form. However, the transportation cost is evaluated with measures parameterized by the current sample mean, while the GLLR is evaluated with measures parameterized by the constrained MLE. If these estimates are significantly different, we may have different candidates as the arms  $a_t^{\min}$  and  $a_t^{\text{ch}}$ . On the other hand, for special classes, such as the exponential family, the sample mean and the MLE are the same, in which case we may have  $a_t^{\min} = a_t^{\text{ch}}$ . Note that to increase the lowest information measure, we should select either the current best arm  $a_t^{\text{top}}$  or the arm with the lowest information measure  $a_t^{\min}$ . Based on the above definitions, the arm selection for TCB is carried out as follows.

$$A_{t+1} \triangleq \begin{cases} \arg \min_{i \in \mathcal{U}_t} T_t(i), & \text{if } \mathcal{U}_t \neq \emptyset \\ a_t^{\text{top}}, & \text{if } \Gamma_t(\mathbf{w}'_t) > \Gamma_t(\frac{1}{t} \mathbf{T}_t) \text{ and } \mathcal{U}_t = \emptyset \\ a_t^{\min}, & \text{if } \Gamma_t(\mathbf{w}'_t) < \Gamma_t(\frac{1}{t} \mathbf{T}_t) \text{ and } \mathcal{U}_t = \emptyset \end{cases}, \quad (28)$$

where  $\mathbf{T}_t \triangleq [T_t(1), \dots, T_t(K)]$ . The complete procedure is presented in Algorithm 1.

---

**Algorithm 1** Transportation cost balancing (TCB)
 

---

- 1: **Initialize:**  $t = 0, \mathcal{U}_t = [K], \mu_t(i) = 0 \forall i \in [K], T_t(i) = 0 \forall i \in [K], \Lambda(a_t^{\text{top}}, a_t^{\text{ch}}) = 0, \beta_t(\delta) = 0$
  - 2: **while**  $\Lambda(a_t^{\text{top}}, a_t^{\text{ch}}) \leq \beta_t(\delta)$  **do**
  - 3:    $t \leftarrow t + 1$
  - 4:   Select an arm  $a_t$  specified by (28) and obtain reward  $X_t$
  - 5:   Update  $\mu_t(a_t)$  and  $T_t(a_t)$  using (12)
  - 6:    $a_t^{\text{top}} \leftarrow \arg \max_{i \in [K]} \bar{\mu}_t(i)$
  - 7:   Compute  $a_t^{\text{min}}$  using (27)
  - 8:   For every  $i \in [K]$ , compute  $w'_t(i)$  using (26)
  - 9:   Compute  $\Gamma_t(\mathbf{w}'_t)$  and  $\Gamma_t(\frac{1}{t}\mathbf{T})$  using (25)
  - 10:   Compute  $\Lambda_t(a_t^{\text{top}}, i)$  for every  $i \in [K] \setminus \{a_t^{\text{top}}\}$
  - 11:    $a_t^{\text{ch}} \leftarrow \arg \min_{i \in [K] \setminus \{a_t^{\text{top}}\}} \Lambda_t(a_t^{\text{top}}, i)$
  - 12:   Update  $\beta_t(\delta)$  using (18)
  - 13: **end while**
  - 14: **Output:** Top arm  $a_t^{\text{top}}$
- 

**Improved Transportation Cost Balancing (ITCB).** A recent study by [12] shows that an additional exploration penalty based on the number of times that each arm is chosen improves the empirical performance of top-two algorithms, albeit achieving the same asymptotic optimality guarantee. Specifically, [12] proposes an additive penalty  $\log(T_t(i))$  to the GLLRs for each arm  $i \in [K] \setminus \{a_t^{\text{top}}\}$  to promote further exploration of under-explored arms. Motivated by this observation, we also devise a modified sampling rule that achieves the same optimality guarantee as TCB, with improved empirical performance. To formalize the modified sampling rule, we begin by defining the lowest *penalized* information measure as:

$$\Phi_t(\mathbf{w}) \triangleq \min_{i \in [K] \setminus \{a_t^{\text{top}}\}} \left\{ \min_{x \in I_{t,i}} \left\{ w_{a_t^{\text{top}}} d_{\text{U}}(\mathbb{P}_{t,a_t^{\text{top}}}, x) + w_i d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} + \frac{\log(tw_i)}{t} \right\}, \quad (29)$$

for any  $\mathbf{w} \in \Delta^K$ . Furthermore, let us define the arm having the lowest penalized information measure as

$$b_t^{\text{min}} \triangleq \arg \min_{i \in [K] \setminus \{a_t^{\text{top}}\}} \left\{ \min_{x \in I_{t,i}} \left\{ \frac{T_t(a_t^{\text{top}})}{t} d_{\text{U}}(\mathbb{P}_{t,a_t^{\text{top}}}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} + \frac{\log(T_t(i))}{t} \right\}, \quad (30)$$

Based on this definition, the modified sampling rule in ITCB is specified next.

$$A_{t+1} \triangleq \begin{cases} \arg \min_{i \in \mathcal{U}_t} T_t(i), & \text{if } \mathcal{U}_t \neq \emptyset \\ b_t^{\text{min}}, & \text{if } \Phi_t(\mathbf{w}'_t) > \Phi_t(\frac{1}{t}\mathbf{T}_t) \text{ and } \mathcal{U}_t = \emptyset \\ a_t^{\text{top}}, & \text{if } \Phi_t(\mathbf{w}'_t) < \Phi_t(\frac{1}{t}\mathbf{T}_t) \text{ and } \mathcal{U}_t = \emptyset \end{cases}. \quad (31)$$

Note that main difference in the selection rules (28) and (31) lies in the cost function  $\Phi_t$  which has an additional penalty compared to  $\Gamma_t$ , that promotes additional exploration. The ITCB sampling rule follows the same principle as TCB, with the difference of a penalized cost function.

## 4 Performance Guarantees

This section provides the main results on the performance of TCB and ITCB algorithms. There are two key results that we are interested in proving. First, we want to show that the TCB and ITCB algorithms satisfy the  $\delta$ -PAC guarantee on the decision confidence. Next, we show that the average sample complexities of these algorithms match the known information-theoretic lower bound asymptotically. To this end, we will prove a few properties of TCB and ITCB, which will collectively establish asymptotic optimality in terms of the average sample complexity. We start by stating the guarantee on the probability of error and characterizing  $\beta_t(\delta)$  such that the algorithms are  $\delta$ -PAC. To formalize this result, we define a few quantities that characterize the stopping threshold  $\beta_t(\delta)$ . We denote the FI measure corresponding to arm  $i \in [K]$  evaluated at the current MLE  $\mu_t(i)$  by  $\mathcal{I}_i(\mu_t(i))$ . For any  $i \in [K]$ , let us define

$$\bar{V}_t(i) \triangleq - \sum_{s \in [t]: A_s = i} \left( \frac{\partial^2}{\partial \theta^2} \log \pi_i(X_s | \theta) \right)_{\theta = \mu_t(i)}, \quad (32)$$

which represents the second-order derivative of the log-likelihood function of arm  $i \in [K]$ , evaluated at the ML estimate  $\mu_t(i)$ . Accordingly, for any  $\varepsilon \in \mathbb{R}_+$  and  $i \in [K]$ , we define

$$W_t(\varepsilon, i) \triangleq \int_{\Omega^{\otimes T_t(i)}} \log \left( 1 - 2Q \left( \varepsilon \sqrt{\bar{V}_t(i)} \right) \right) \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i, \quad (33)$$

where  $Q(x)$  denotes the  $Q$  function. Furthermore, let us define

$$W_t(\varepsilon) \triangleq \max_{i \in [K]} \log \mathcal{I}_i(\mu_t(i)) - 2 \min_{i \in [K]} W_t(\varepsilon, i). \quad (34)$$

Note that as  $t \rightarrow \infty$ ,  $\mathcal{I}_i(\mu_t(i))$  converges to the FI measure under the true parameter  $\mu(i)$ , and the second term in (34) converges to 0, given that each arm is sampled sufficiently often. To see why this is true, recall that according to Assumption 7, we have  $\bar{V}_t(i) \geq T_t(i)\sigma^2$  for any  $i \in [K]$ , which can become infinitely large if the arm  $i \in [K]$  is sampled sufficiently often. Hence,  $W_t(\varepsilon)$  converges to the maximum FI measure for the distributions in the bandit instance  $\nu$ . Furthermore, for any  $\varepsilon \in \mathbb{R}_+$  let us define

$$\varepsilon_t \triangleq \max_{i \in [K]} \left\{ \max \{d_i(\mu_t(i) \| \mu_t(i) - \varepsilon), d_i(\mu_t(i) \| \mu_t(i) + \varepsilon)\} \right\}. \quad (35)$$

Note that the quantity  $\varepsilon_t$  can be made arbitrarily small by choosing a sufficiently small  $\varepsilon$  as a consequence of the uniform continuity of the KL divergence measures in Assumption 5. Subsequently, we state the choice of  $\beta_t(\delta)$  that yields a  $\delta$ -PAC guarantee for any BAI algorithm, irrespective of the sampling rule.

**Theorem 1** ( $\delta$ -PAC). *The stopping rule in (18) with the choice of the threshold*

$$\beta_t(\delta) \triangleq W_t(\varepsilon) + t\varepsilon_t + 2 \log \frac{|\Theta|}{\sqrt{2\pi}} + \log \frac{t(K-1)}{2\delta}, \quad (36)$$

along with any arm selection strategy and the decision rule  $\hat{A}_\tau \triangleq a_\tau^{\text{top}}$  is  $\delta$ -PAC for any  $\varepsilon \in \mathbb{R}_+$ , where  $|\Theta|$  denotes the volume of the space of parameters  $\Theta$ .

*Proof.* See Appendix A. ■

The stopping threshold specified in Theorem 1 holds for a general class of bandit instances that satisfy Assumptions 1-8. We observe that the threshold depends on the volume of the parameter space  $\Theta$ , which is uncommon in BAI. However, despite the additional penalty of the order of  $O(\log |\Theta|)$ , our stopping threshold applies a very general class of parameterized bandits for which stopping thresholds do not exist for the test statistic under consideration. Even though we may use the non-parametric GLLR statistic proposed in [12, 13], it is computationally expensive as it requires solving a convex optimization problem in each iteration to compute the test statistic. Naturally, the statistic in [13], designed for non-parametric bandits, does not use the knowledge of the parametric form of the likelihood functions, which it needs to estimate from the rewards. Furthermore, when we specialize to more structured classes of bandit instances, specifically the single parameter exponential family, we can use the tighter thresholds from [21]. The key difficulty in the proof of Theorem 1 for the general case is that the log-likelihood function is generally non-linear in the reward. This is in contrast to the exponential family, which has a linear log-likelihood function in terms of the reward, which facilitates an intelligent mixture martingale construction for designing the stopping threshold [21].

Next, we state the results related to the asymptotic optimality of TCB and ITCB in terms of the average sample complexity. We begin by establishing a few properties of the problem complexity that are useful for characterizing the sample complexity. Specifically, the form of the problem complexity  $\Gamma(\nu)$  in Lemma 1 is instrumental in characterizing the key properties of  $\Gamma(\nu)$  and establishing an upper bound on the average sample complexity of TCB and ITCB. The following lemma characterizes these properties, which include the continuity of the problem complexity in terms of the bandit instance  $\nu \in \mathcal{M}$  in the metric space  $(\mathcal{M}, D_{TV})$  and the characterization of an optimal allocation of samples that maximizes the problem complexity.

**Lemma 2** (Properties of  $\Gamma(\nu)$ ). *The problem complexity  $\Gamma(\nu)$  has the following properties:*

1. Functions  $d_U(\cdot, \cdot)$  and  $d_L(\cdot, \cdot)$  are strictly convex in their second arguments.
2. Problem complexity  $\Gamma : \mathcal{M} \mapsto \mathbb{R}$  and the optimal allocation  $\mathbf{w} : \mathcal{M} \mapsto \Delta^K$  are continuous functions on the metric space  $(\mathcal{M}, D_{TV})$ . Furthermore, an optimal sampling proportion is given by the unique allocation that satisfies:

$$\Gamma_i(\nu, \mathbf{w}) = \Gamma_j(\nu, \mathbf{w}), \quad \forall i, j \neq a^*, \quad (37)$$

where, for any  $i \in [K] \setminus \{a^*\}$ , we have defined

$$\Gamma_i(\nu, \mathbf{w}) \triangleq \inf_{\mathbb{P} \in \mathcal{M}: m(\mathbb{P}_i) \geq m(\mathbb{P}_{a^*})} \left\{ w_{a^*} D_{KL}(\mathbb{P}_{a^*} \| \bar{\mathbb{P}}_{a^*}) + w_i D_{KL}(\mathbb{P}_i \| \bar{\mathbb{P}}_i) \right\}. \quad (38)$$

*Proof.* See Appendix B.2. ■

The continuity of  $\Gamma$  has been previously established by [13] in the metric space  $(\mathcal{M}, W_1)$ , where  $W_1$  denotes the Wasserstein distance metric. However, the continuity was established under the assumption that for any  $\mathbb{P} \in \mathcal{Q}(\Omega)$ , the boundedness assumption  $\mathbb{E}_{\mathbb{P}}[f(|X|)] < B$  holds for any convex and differentiable function  $f$ . The key difference in the proof of Lemma 2 with that of [13, Lemma 4] is that we do not impose the boundedness assumption. Rather, we leverage the properties of the Jensen-Shannon divergence (details in Appendix B.2) to establish the continuity property of the problem complexity.

Leveraging the continuity property of the problem complexity proved in Lemma 2, next, we establish the convergence in the sampling proportions due to the TCB and ITCB sampling rules to that of the optimal proportions. This is a key

theoretical contribution of the paper, which helps establish the asymptotic optimality of the algorithms, compared to the  $\beta$ -optimality achieved by top-two algorithms. Specifically, the top-two algorithms *enforce* the almost sure convergence of allocation for the best arm to  $\beta$ , which accordingly ensures the convergence to the  $\beta$ -optimal allocation. However, such an analysis does not readily extend to that of TCB and ITCB, since we do not enforce the almost sure convergence in allocation of the best arm to  $\beta$ . This makes the analysis of the TCB and ITCB sampling rules significantly more challenging. For details, we refer to Appendix C.

**Theorem 2** (Convergence in sampling proportions). *If the TCB arm selection rule in (28) is allowed to continue sampling without stopping, for any  $\epsilon > 0$ , there exists a stochastic time instant  $N_{\mathbf{w}}^\epsilon$  such that*

$$\left| \frac{T_t(i)}{t} - w_i(\boldsymbol{\nu}) \right| < \epsilon, \quad \forall i \in [K], \forall t \geq N_{\mathbf{w}}^\epsilon. \quad (39)$$

Furthermore,  $\mathbb{E}_{\boldsymbol{\nu}}[N_{\mathbf{w}}^\epsilon] < +\infty$ .

*Proof.* See Appendix C. ■

**Theorem 3** (Convergence in sampling proportions). *If the ITCB arm selection rule in (31) is allowed to continue sampling without stopping, for any  $\epsilon > 0$ , there exists a stochastic time instant  $N_{\mathbf{w}}^\epsilon$  such that*

$$\left| \frac{T_t(i)}{t} - w_i(\boldsymbol{\nu}) \right| < \epsilon, \quad \forall i \in [K], \forall t \geq N_{\mathbf{w}}^\epsilon. \quad (40)$$

Furthermore,  $\mathbb{E}_{\boldsymbol{\nu}}[N_{\mathbf{w}}^\epsilon] < +\infty$ .

*Proof.* See Appendix C. ■

Finally, leveraging Theorem 2 (or 3), we characterize an upper bound on the average sample complexity of the TCB and ITCB algorithms.

**Theorem 4** (Achievable sample complexity). *The TCB and ITCB algorithms, comprising the arm selection rules in (28) and (31) and the stopping rule in (18), satisfy the following upper-bound on the average sample complexity.*

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\boldsymbol{\nu}}[\tau]}{\log(1/\delta)} \leq \frac{1 + \alpha}{\Gamma(\boldsymbol{\nu})}, \quad (41)$$

for any  $\alpha > 0$ .

*Proof.* See Appendix D. ■

The upper bound in Theorem 4 matches the information-theoretic lower bound on the average sample complexity provided by [4] up to any  $\alpha > 0$ . This establishes the asymptotic optimality of the TCB and ITCB algorithms. Note that in the special case that the bandit instance belongs to the single parameter exponential family, we can use the tighter stopping threshold from [21], in which case we can tighten the sample complexity bound in Theorem 4. Specifically, we may achieve a sample complexity bound with  $\alpha = 0$  in (41). The proof follows similar arguments as [7, Theorem 5].

## 5 Numerical Experiments

In this section, we provide numerical evaluations of the different aspects of processes and decisions involved in the TCB and ITCB algorithms. First, we evaluate the convergence of the TCB and ITCB sampling rules in Gaussian and Bernoulli bandit instances. These bandit instances are summarized in Table 1. We have selected two slippage bandit instances (i.e., instances with identical arm means for the sub-optimal arms) and two instances having distinct arm means. All experiments are averaged over  $10^4$  independent trials. We show convergence in the following three senses.

Distribution	$\nu$	$\mathbf{w}(\nu)$
Gaussian	$\nu_1 = [100, 2, 1]$	$[0.4143, 0.2979, 0.2878]$
Gaussian	$\nu_2 = [1.2, 1, 1, 1, 1]$	$[1/3, 1/6, 1/6, 1/6, 1/6]$
Bernoulli	$\nu_3 = [0.8, 0.45, 0.45, 0.45]$	$[0.3854, 0.2049, 0.2049, 0.2049]$
Bernoulli	$\nu_4 = [0.79, 0.29, 0.289]$	$[0.426, 0.2880, 0.2860]$

Table 1: Bandit instances and their associated sampling weights.

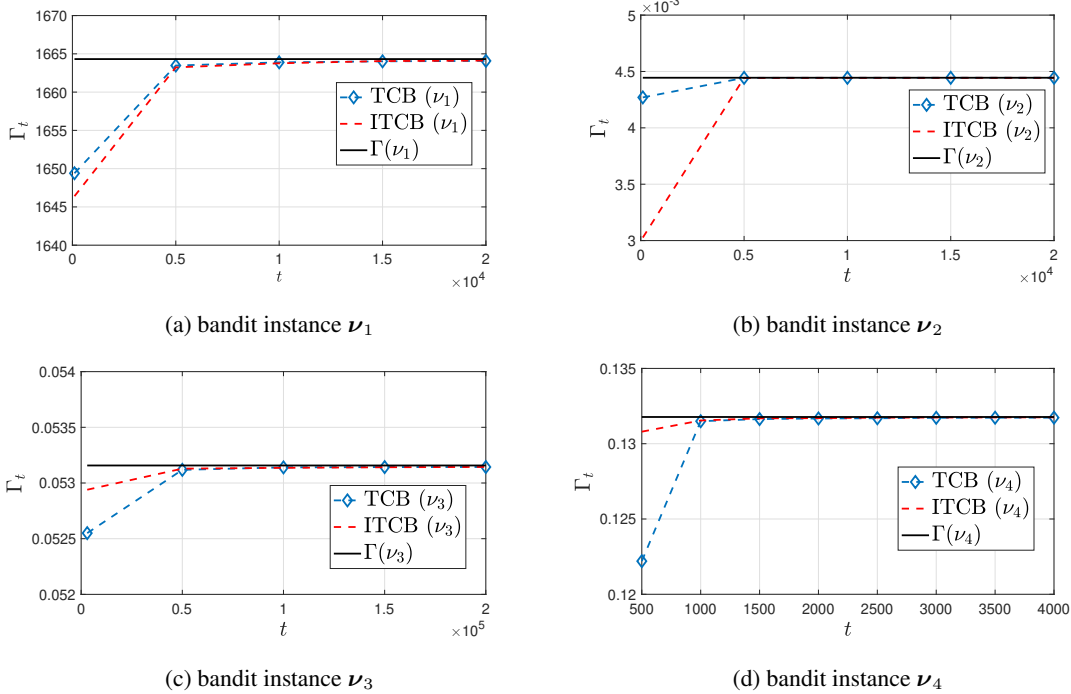


Figure 2: Variations of the estimates of the transportation cost  $\Gamma_t$  over time.

**Transportation cost estimates.** The key idea of the TCB and ITCB arm selection rules is to estimate the problem complexity  $\Gamma(\nu)$ . This is achieved by maximizing the transportation cost  $\Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right)$ , which acts as an estimate of the problem complexity. The analysis of the average sample complexity of the TCB and ITCB algorithms provided in Theorem 4 critically hinges on the convergence in the transportation cost to the problem complexity (see Appendix D for details). To empirically evaluate this, in figures 2a-2d we illustrate the term  $\Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right)$  for the problem instances

specified in Table 1. These empirical results show the convergence of the transportation costs to their optimal values (problem complexity). We note that the range of the transportation cost depends on the mean values of the arms, and consequently, the problem complexity of instance  $\nu_1$  (with  $\mu(1) = 100$ ) is significantly larger than the other bandit instances whose mean values are closer to 1.

**Estimates of  $\beta$ .** An optimal sampling rule design for  $\delta$ -PAC BAI is a functional estimation problem in which the learner forms estimates of the problem complexity for arm selection decisions in the form of transportation costs. These transportation costs, among other parameters, are also functions of the sampling proportion assigned to the best arm  $a^*$ . Hence, as established in Theorem 2 and Theorem 3, the TCB and ITCB sampling strategies implicitly *estimate*  $\beta$ , i.e., the optimal sampling proportion for the best arm. To empirically establish the variations of the estimates of the allocation over time, in figures 3a-3d, we illustrate the allocation of sampling resources to the best arm and compare it to the optimal choice of  $\beta = w_{a^*}$  in the four instances specified in Table 1. Similarly to the transportation cost estimates, it is observed that the estimates  $\beta_t$  converge to the optimal value.

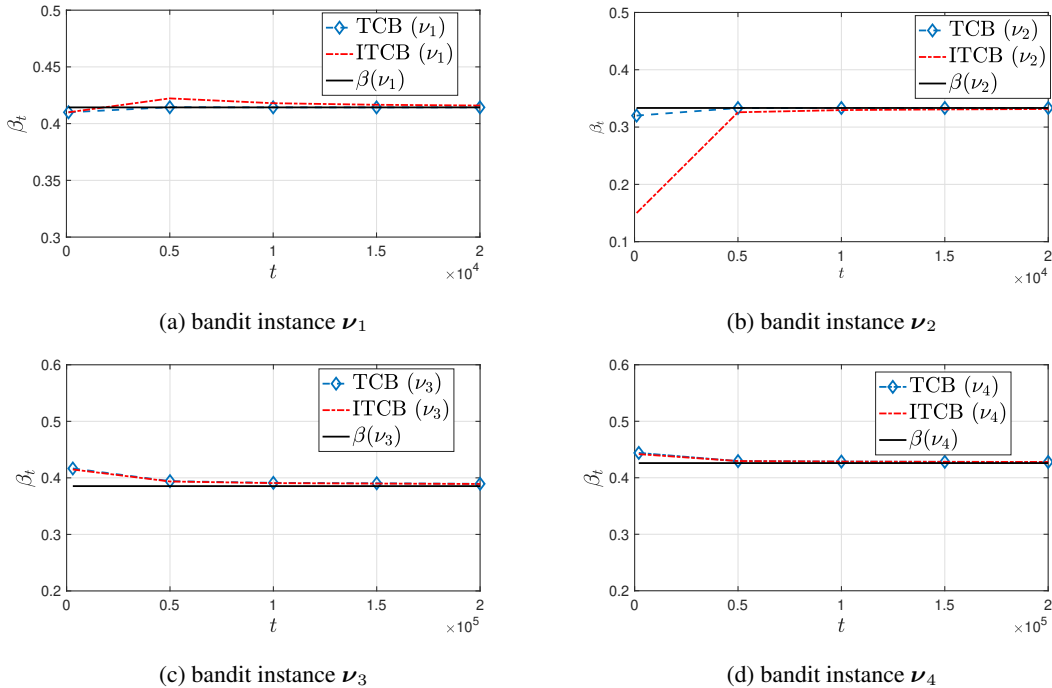


Figure 3: Variations of the estimates  $\beta_t$  over time.

**Sampling over-sampled arms.** The analysis of the TCB and ITCB algorithms show that their sampling rules eventually sample from the set of *under-sampled* arms, i.e., the set of arms that have been sampled fewer times compared to the optimal allocation (for details, see Appendix C). We demonstrate that this is a key property of the TCB and ITCB sampling rules that enables the convergence of the transportation cost to the problem complexity. To show this, we devise an algorithm that preserves the explicit exploration phase of the TCB and ITCB algorithms for convergence in the arm means. However, if the set of under-explored arms is empty, this algorithm always selects the current best arm, i.e.,  $a_t^{\text{top}}$ . It can be readily verified that due to the convergence in mean for all the arms,  $a_t^{\text{top}}$  converges to the best



arm  $a^*$ . Furthermore, the best arm  $a^*$  eventually gets over-sampled since we put the entire allocation on  $a^*$ . Figure 4 shows the deviation of the transportation cost from the problem complexity over time. We observe that the transportation cost diverges from the problem complexity, establishing that sampling from the set of over-sampled arms cannot be optimal.

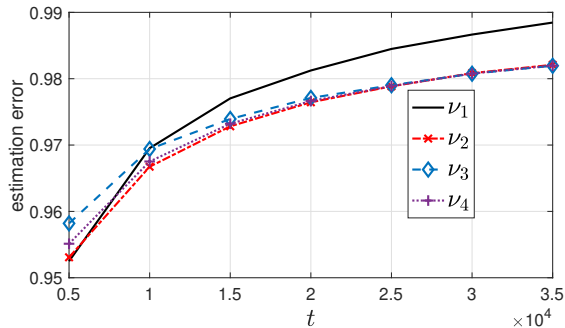


Figure 4: Divergence of transportation cost without controlling under-sampling.

**Sensitivity to  $\beta$ .** Next, we assess the dependency of the state-of-the-art BAI algorithms on  $\beta$ , and compare the empirical performance of the TCB and ITCB algorithms to these algorithms. These state-of-the-art algorithms for BAI include: T3C [11], TT-SPRT [7, 19], TS-TCI [12], EB-TCI [12], and FW [18]. TT-SPRT and EB-TCI use the empirical best arm as the leader. TT-SPRT uses the arm with the closest GLLR statistic to the empirical best arm as the challenger, and EB-TCI includes an additional exploration penalty to the GLLR statistic to determine the challenger. Distinct from TT-SPRT and EB-TCI, T3C and TS-TCI use Thompson sampling from the posterior distribution to identify the leader. T3C selects the challenger as the arm with the closest GLLR statistic to the leader, while TS-TCI uses a penalized GLLR statistic to promote exploration. FW selects the next arm based on a Frank-Wolfe-based update step to solve (22) based on the current mean estimates. Furthermore, we also compare the T3C algorithm with the  $\beta$ -tuning routine from [10]. Note that we have not compared TTTS with  $\beta$ -tuning, owing to its large computation time required in identifying the challenger. We perform our experiments based on two common reward distributions, Bernoulli and Gaussian bandits. All the experiments are averaged over 2000 independent Monte Carlo trials, and we have set  $\delta = 10^{-8}$ . We have the following two sets of experiments.

1. *Bernoulli*. In this experiment, we use a Bernoulli bandit with mean values  $[0.8, 0.5, 0.3, 0.29, 0.06]$  and set  $\delta \triangleq 10^{-8}$ . Figure 5 shows the performance of the TCB and ITCB arm selection rules compared to the state-of-the-art. We observe that the TCB and ITCB arm selection strategies in (28) and (31) are agnostic to the parameter  $\beta$ . Furthermore, ITCB outperforms the top-two sampling strategies for various values of the tuning parameter  $\beta$ , and its performance is comparable to the optimization-based sampling rule FW. Furthermore, the performance of ITCB matches that of TS-TCI and EB-TCI at  $\beta \approx 0.4$ .
2. *Gaussian*. For the next experiment, we take a Gaussian bandit instance with mean values given by  $[1.2, 1, 1, 1, 1]$ . Figure 6 compares the TCB and ITCB algorithms against state-of-the-art BAI algorithms for various values of  $\beta$ . We set  $\delta \triangleq 10^{-8}$  for this experiment. We observe that the ITCB algorithm outperforms the top-two

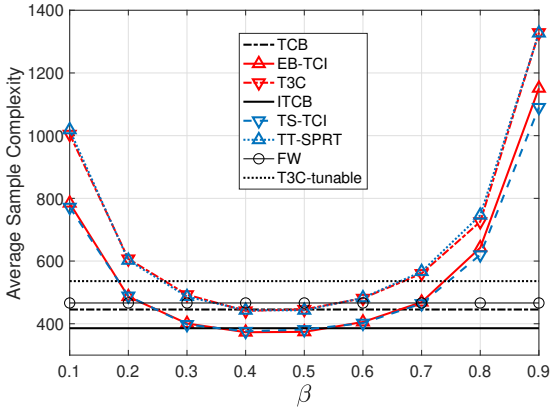


Figure 5: Sensitivity to  $\beta$  in the Bernoulli instance  $[0.8, 0.5, 0.3, 0.29, 0.06]$  with  $\delta = 10^{-8}$ .

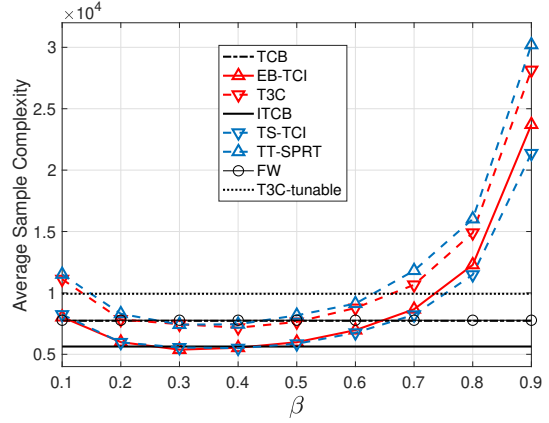


Figure 6: Sensitivity to  $\beta$  in the Gaussian instance  $[1.2, 1, 1, 1, 1]$  with  $\delta = 10^{-8}$ .

algorithms in various regimes of  $\beta$ , and the performances of TS-TCI and EB-TCI match that of ITCB in the range  $0.3 \leq \beta \leq 0.4$ .

## 6 Conclusions

We have investigated the problem of fixed-confidence best-arm identification (BAI) in stochastic multi-arm bandits (MABs). We have designed the transport cost balancing (TCB) algorithm, in which the key decision is finding the optimal allocation of the sampling resources among different arms. This extends the existing strategies that aim to allocate a  $\beta$  fraction of the resources to the best arm and achieve  $\beta$ -optimality, with  $\beta$  remaining a parameter for the sampling routine. In the proposed TCB algorithm, the optimal value of  $\beta$  is also *implicitly* estimated, rendering the algorithm independent of  $\beta$ . As a result, we have established that the proposed TCB algorithm is asymptotically optimal. We have also extended the TCB algorithm by including an additional exploration penalty based on the number of times each arm is chosen. This algorithm, referred to as improved TCB (ITCB), is also shown to achieve asymptotic optimality and improved empirical performance compared to TCB.

## A Proof of Theorem 1

The proof of Theorem 1 is based on upper bounding the GLLR statistic using a mixture martingale construction, which we state in Lemma 4. Subsequently, we use Ville's supermartingale inequality to upper bound the probability of an incorrect terminal decision. The martingale construction closely resembles Laplace's approximation method for approximating the maximum value of any function [22] and is based on an upper bound on the maximum of a sum of twice-differentiable functions, which is stated in Lemma 3.

**Lemma 3.** *For any sequence  $\{g_s : s \in [t]\}$  of twice-differentiable functions  $g : \Theta \mapsto \mathbb{R}$ , where  $\Theta$  is a compact space, let us define the maximizer*

$$\mu_t \triangleq \arg \max_{\rho \in \Theta} \sum_{s \in [t]} g_s(\rho). \quad (42)$$

Furthermore, let us denote  $\eta$  as the uniform distribution over  $\Theta$ . Then, for any  $\varepsilon \in \mathbb{R}_+$ , we have

$$\begin{aligned} \sum_{s \in [t]} g_s(\mu_t) &\leq \log \mathbb{E}_\eta \left[ \exp \left( \sum_{s \in [t]} g_s(\rho) \right) \right] + \frac{1}{2} \log V_t - \log \left( 1 - 2Q(\varepsilon \sqrt{V_t}) \right) \\ &\quad - \min_{\rho \in [\mu_t - \varepsilon, \mu_t + \varepsilon]} \left\{ \sum_{s \in [t]} (g_s(\rho) - g_s(\mu_t)) \right\} + \log \frac{|\Theta|}{\sqrt{2\pi}}, \end{aligned} \quad (43)$$

where we have defined

$$V_t \triangleq - \sum_{s \in [t]} g_s''(\rho) \Big|_{\rho = \mu_t}. \quad (44)$$

*Proof.* Using Taylor's expansion, we obtain

$$\sum_{s \in [t]} g_s(\rho) = \sum_{s \in [t]} g_s(\mu_t) - \frac{1}{2}(\rho - \mu_t)^2 V_t + R(\rho, \mu_t), \quad (45)$$

where we have defined<sup>1</sup>

$$R(\rho, \mu_t) \triangleq \sum_{s \in [t]} \sum_{j=3}^{\infty} \frac{1}{j!} g_s^{(j)}(\rho - \mu_t)^j, \quad (46)$$

and  $g_s^{(j)}(\mu_t)$  denotes the  $j^{\text{th}}$  derivative of the function  $g_s$  evaluated at  $\mu_t$ . Furthermore, we have

$$\begin{aligned} \int_{\Theta} \exp \left( \sum_{s \in [t]} g_s(\rho) \right) d\eta(\rho) &\stackrel{(45)}{=} \frac{1}{|\Theta|} \exp \left( \sum_{s \in [t]} g_s(\mu_t) \right) \int_{\Theta} \exp \left( -\frac{1}{2}(\rho - \mu_t)^2 V_t \right) \cdot \exp(R(\rho, \mu_t)) d\rho \quad (47) \\ &\geq \frac{1}{|\Theta|} \exp \left( \sum_{s \in [t]} g_s(\mu_t) \right) \int_{\mu_t - \varepsilon}^{\mu_t + \varepsilon} \exp \left( -\frac{1}{2}(\rho - \mu_t)^2 V_t \right) \cdot \exp(R(\rho, \mu_t)) d\rho, \end{aligned} \quad (48)$$

<sup>1</sup>Note that we only assume  $g$  to be twice-differentiable. In case the higher order derivatives don't exist, we may simply define  $R(\rho, \mu_t) \triangleq \sum_{s \in [t]} (g_s(\rho) - g_s(\mu_t)) + \frac{1}{2}(\rho - \mu_t)^2 V_t$ .

where the right hand side in (48) is positive due to the fact that  $\mu_t \in \Theta$ . Furthermore, let us define

$$R_t^{\min} \triangleq \min_{\rho \in [\mu_t - \varepsilon, \mu_t + \varepsilon]} R(\rho, \mu_t). \quad (49)$$

Using (49), (48) can be lower bounded as

$$\int_{\Theta} \exp \left( \sum_{s \in [t]} g_s(\rho) \right) d\eta(\rho) \geq \frac{1}{|\Theta|} \exp \left( \sum_{s \in [t]} g_s(\mu_t) + R_t^{\min} \right) \int_{\mu_t - \varepsilon}^{\mu_t + \varepsilon} \exp \left( -\frac{1}{2}(\rho - \mu_t)^2 V_t \right) d\rho. \quad (50)$$

Let us make the following change of variables.

$$y = (\rho - \mu_t)\sqrt{V_t}, \quad y_t = \varepsilon\sqrt{V_t}. \quad (51)$$

We have

$$\int_{\Theta} \exp \left( \sum_{s \in [t]} g_s(\rho) \right) d\eta(\rho) \stackrel{(50)-(51)}{\geq} \frac{1}{|\Theta|} \exp \left( \sum_{s \in [t]} g_s(\mu_t) + R_t^{\min} \right) \cdot \frac{1}{\sqrt{V_t}} \int_{-y_t}^{y_t} \exp \left( -\frac{y^2}{2} \right) dy \quad (52)$$

$$= \frac{1}{|\Theta|} \exp \left( \sum_{s \in [t]} g_s(\mu_t) + R_t^{\min} \right) \cdot \sqrt{\frac{2\pi}{V_t}} \left( 1 - 2Q(\varepsilon\sqrt{V_t}) \right), \quad (53)$$

where  $Q(x)$  denotes the  $Q$  function evaluated at any point  $x \in \mathbb{R}$ . Taking log on both sides of (53) and rearranging, we have

$$\sum_{s \in [t]} g_s(\mu_t) \leq \log \mathbb{E}_{\eta} \left[ \exp \left( \sum_{s \in [t]} g_s(\rho) \right) \right] + \frac{1}{2} \log V_t - \log \left( 1 - 2Q(\varepsilon\sqrt{V_t}) \right) - R_t^{\min} + \log \frac{|\Theta|}{\sqrt{2\pi}}. \quad (54)$$

Furthermore, note that

$$R_t^{\min} \stackrel{(45)}{=} \min_{\rho \in [\mu_t - \varepsilon, \mu_t + \varepsilon]} \left\{ \sum_{s \in [t]} g_s(\rho) + \frac{1}{2}(\rho - \mu_t)^2 V_t \right\} - \sum_{s \in [t]} g_s(\mu_t) \quad (55)$$

$$\geq \min_{\rho \in [\mu_t - \varepsilon, \mu_t + \varepsilon]} \left\{ \sum_{s \in [t]} (g_s(\rho) - g_s(\mu_t)) \right\}. \quad (56)$$

Finally, combining (54) and (56), we obtain (43). ■

**Lemma 4.** *For any arm  $i \in [K]$  and for any  $\varepsilon \in \mathbb{R}_+$ , there exists a non-negative martingale  $M_t(i)$  satisfying  $\mathbb{E}[M_1(i)] = 1$ , such that*

$$\begin{aligned} T_t(i) d_i(\mu_t(i) \| \mu(i)) &\leq \log M_t(i) + \frac{1}{2} \log (T_t(i) \mathcal{I}_i(\mu_t(i))) - W_t(\varepsilon, i) + \log \frac{|\Theta|}{\sqrt{2\pi}} \\ &\quad + T_t(i) \max \{ d_i(\mu_t(i) \| \mu_t(i) - \varepsilon), d_i(\mu_t(i) \| \mu_t(i) + \varepsilon) \}. \end{aligned} \quad (57)$$

*Proof.*

$$T_t(i) d_i(\mu_t(i) \| \mu(i)) = \sum_{s \in [t]: A_s = i} \int_{\Omega} \left( \log \frac{\pi_i(X | \mu_t(i))}{\pi_i(X | \mu(i))} \right) \pi_i(X | \mu_t(i)) dX \quad (58)$$

$$\begin{aligned}
&= \sum_{s \in [t]: A_s = i} \int_{\Omega} (\log \pi_i(X | \mu_t(i))) \pi_i(X | \mu_t(i)) dX \\
&\quad - \sum_{s \in [t]: A_s = i} \int_{\Omega} (\log \pi_i(X | \mu(i))) \pi_i(X | \mu_t(i)) dX .
\end{aligned} \tag{59}$$

Next, we will leverage Lemma 3 using  $g(\rho) = \log \pi_i(X | \rho)$ . Note that

$$\begin{aligned}
&\int_{\Omega^{\otimes T_t(i)}} \left( \log \bar{V}_t(i) \right) \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \\
&\leq \log \int_{\Omega^{\otimes T_t(i)}} \bar{V}_t(i) \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i
\end{aligned} \tag{60}$$

$$= \log \int_{\Omega^{\otimes T_t(i)}} \sum_{s \in [t]: A_s = i} \left( -\frac{\partial^2}{\partial \theta^2} \log \pi_i(X_s | \theta) \right)_{\theta = \mu_t(i)} \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \tag{61}$$

$$= \log \sum_{s \in [t]: A_s = i} \int_{\Omega} \left( -\frac{\partial^2}{\partial \theta^2} \log \pi_i(X | \theta) \right)_{\theta = \mu_t(i)} \pi_i(X | \mu_t(i)) dX \tag{62}$$

$$= \sum_{s \in [t]: A_s = i} \mathcal{I}_i(\mu_t(i)) \tag{63}$$

$$= T_t(i) \mathcal{I}_i(\mu_t(i)) , \tag{64}$$

where (60) is obtained using Jensen's inequality, and (62) is a result of applying the Fubini-Tonelli's theorem. Furthermore, we have

$$\begin{aligned}
&-\int_{\Omega^{\otimes T_t(i)}} \min_{\rho \in [\mu_t(i) - \varepsilon, \mu_t(i) + \varepsilon]} \sum_{s \in [t]: A_s = i} \log \frac{\pi_i(X_s | \rho)}{\pi_i(X_s | \mu_t(i))} \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \\
&\leq -\min_{\rho \in [\mu_t(i) - \varepsilon, \mu_t(i) + \varepsilon]} \int_{\Omega^{\otimes T_t(i)}} \sum_{s \in [t]: A_s = i} \log \frac{\pi_i(X_s | \rho)}{\pi_i(X_s | \mu_t(i))} \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i
\end{aligned} \tag{65}$$

$$= -\min_{s \in [t]: A_s = i} d_i(\rho | \mu_t(i)) \tag{66}$$

$$= \max_{s \in [t]: A_s = i} d_i(\mu_t(i) | \rho) , \tag{67}$$

where (65) follows from Jensen's inequality along with Assumption 2. Next, using Lemma 3 along with (64) and (67), (59) can be upper bounded by

$$\begin{aligned}
T_t(i) d_i(\mu_t(i) | \mu(i)) &\leq \int_{\Omega^{\otimes T_t(i)}} \log \mathbb{E}_{\eta} \left[ \exp \left( \sum_{s \in [t]: A_s = i} \log \pi_i(X_s | \rho) \right) \right] \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \\
&\quad - \int_{\Omega^{\otimes T_t(i)}} \sum_{s \in [t]: A_s = i} (\log \pi_i(X | \mu(i))) \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \\
&\quad + \frac{1}{2} \log (T_t(i) \mathcal{I}_i(\mu_t(i))) - W_t(\varepsilon, i) + \log \frac{|\Theta|}{\sqrt{2\pi}}
\end{aligned}$$

$$\begin{aligned}
& + T_t(i) \cdot \max_{\rho \in [\mu_t(i) - \varepsilon, \mu_t(i) + \varepsilon]} d_i(\mu_t(i) \|\rho) \tag{68} \\
\leq & \log \underbrace{\int_{\Omega^{\otimes T_t(i)}} \mathbb{E}_\eta \left[ \exp \left( \sum_{s \in [t]: A_s = i} \log \frac{\pi_i(X_s | \rho)}{\pi_i(X_s | \mu(i))} \right) \right]}_{\triangleq M_t(i)} \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \\
& + \frac{1}{2} \log (T_t(i) \mathcal{I}_i(\mu_t(i))) - W_t(\varepsilon, i) + \log \frac{|\Theta|}{\sqrt{2\pi}} \\
& + T_t(i) \cdot \max_{\rho \in [\mu_t(i) - \varepsilon, \mu_t(i) + \varepsilon]} d_i(\mu_t(i) \|\rho) , \tag{69}
\end{aligned}$$

where (69) is obtained by using the Jensen's inequality, and  $M_t(i)$  is defined such that  $M_1(i) = 1$ , if  $A_1 \neq i$ .

Furthermore, let us define

$$\rho_t \triangleq \arg \max_{\rho \in [\mu_t(i) - \varepsilon, \mu_t(i) + \varepsilon]} \{T_t(i) d_i(\mu_t(i) \|\rho)\} . \tag{70}$$

There exists  $\gamma \in (0, 1)$  such that we have  $\rho_t = \gamma(\mu_t(i) - \varepsilon) + (1 - \gamma)(\mu_t(i) + \varepsilon)$ . Owing to the convexity of KL divergence, we have

$$d_i(\mu_t(i) \|\rho_t) \leq \gamma d_i(\mu_t(i) \|\mu_t(i) - \varepsilon) + (1 - \gamma) d_i(\mu_t(i) \|\mu_t(i) + \varepsilon) \tag{71}$$

$$\leq \max\{d_i(\mu_t(i) \|\mu_t(i) - \varepsilon), d_i(\mu_t(i) \|\mu_t(i) + \varepsilon)\} . \tag{72}$$

Finally, combining (69) and (72), we recover (57). The only thing left to prove is that  $M_t(i)$  is a martingale satisfying  $\mathbb{E}[M_1(i)] = 1$ . Note that if  $A_t \neq i$ , we have  $\mathbb{E}[M_t(i) | \mathcal{F}_{t-1}] = M_{t-1}(i)$ . If  $A_t = i$ , we have

$$\mathbb{E}[M_t(i) | \mathcal{F}_{t-1}] = \mathbb{E} \left[ \int_{\Omega^{\otimes T_t(i)}} \mathbb{E}_\eta \left[ \prod_{s \in [t]: A_s = i} \frac{\pi_i(X_s | \rho)}{\pi_i(X_s | \mu(i))} \right] \prod_{s \in [t]: A_s = i} \pi_i(X_s | \mu_t(i)) d\mathcal{X}_t^i \middle| \mathcal{F}_{t-1} \right] \tag{73}$$

$$= M_{t-1}(i) \cdot \mathbb{E} \left[ \int_{\Omega} \mathbb{E}_\eta \left[ \frac{\pi_i(X | \rho)}{\pi_i(X | \mu(i))} \right] \pi_i(X | \mu_t(i)) dX \right] \tag{74}$$

$$= M_{t-1}(i) \cdot \mathbb{E}_\eta \left[ \int_{\Omega} \mathbb{E} \left[ \frac{\pi_i(X | \rho)}{\pi_i(X | \mu(i))} \right] \pi_i(X | \mu_t(i)) dX \right] \tag{75}$$

$$= M_{t-1}(i) \cdot \mathbb{E}_\eta \left[ \int_{\Omega} \underbrace{\left( \int_{\Omega} \frac{\pi_i(X | \rho)}{\pi_i(X | \mu(i))} \pi_i(X | \mu(i)) dX \right)}_{=1} \pi_i(X | \mu_t(i)) dX \right] \tag{76}$$

$$= M_{t-1}(i) \cdot \mathbb{E}_\eta \left[ \int_{\Omega} \pi_i(X | \mu_t(i)) dX \right] \tag{77}$$

$$= M_{t-1}(i) , \tag{78}$$

which proves that  $M_t(i)$  is a martingale. Furthermore, by definition, if  $A_1 \neq i$ ,  $\mathbb{E}[M_1(i)] = 1$ . Alternatively, if

$A_1 = i$ , we have

$$\mathbb{E}[M_1(i)] = \mathbb{E}_\eta \left[ \int_{\Omega} \underbrace{\left( \int_{\Omega} \frac{\pi_i(X | \rho)}{\pi_i(X | \mu(i))} \pi_i(X | \mu(i)) dX \right)}_{=1} \pi_i(X | \mu_1(i)) dX \right] \quad (79)$$

$$= \mathbb{E}_\eta[1] \quad (80)$$

$$= 1. \quad (81)$$

This concludes the proof of Lemma 4.  $\blacksquare$

Next, we delineate the choice of the threshold  $\beta_t(\delta)$  that facilitates the  $\delta$ -PAC guarantee of the proposed stopping rule in (18). We have

$$\mathbb{P}_\nu(\tau < +\infty, \hat{A}_\tau \neq a^*) = \mathbb{P}_\nu(\exists t \in \mathbb{N}, i \neq a^* : \bar{\mu}_t(i) \geq \max_{j \neq i} \bar{\mu}_t(j), \min_{j \neq i} \Lambda_t(i, j) \geq \beta_t(\delta)) \quad (82)$$

$$\leq \sum_{i \neq a^*} \mathbb{P}_\nu(\exists t \in \mathbb{N} : \bar{\mu}_t(i) \geq \max_{j \neq i} \bar{\mu}_t(j), \min_{j \neq i} \Lambda_t(i, j) \geq \beta_t(\delta)) \quad (83)$$

$$\leq \sum_{i \neq a^*} \mathbb{P}_\nu(\exists t \in \mathbb{N} : \bar{\mu}_t(i) \geq \max_{j \neq i} \bar{\mu}_t(j), \Lambda_t(i, a^*) \geq \beta_t(\delta)) \quad (84)$$

$$\leq \sum_{i \neq a^*} \mathbb{P}_\nu(\exists t \in \mathbb{N} : T_t(i) d_i(\mu_t(i) \| \tilde{\mu}_t(a^*)) + T_t(a^*) d_{a^*}(\mu_t(a^*) \| \tilde{\mu}_t(a^*)) \geq \beta_t(\delta)) \quad (85)$$

$$\leq \sum_{i \neq a^*} \mathbb{P}_\nu(\exists t \in \mathbb{N} : T_t(i) d_i(\mu_t(i) \| \mu(i)) + T_t(a^*) d_{a^*}(\mu_t(a^*) \| \mu(a^*)) \geq \beta_t(\delta)), \quad (86)$$

where (86) uses the definition of  $\Lambda_t(i, a^*)$ . Specifically, using the KKT conditions, we obtain that the minimizer in (14) satisfies the condition that  $\rho(i) = \rho(a^*)$ . Denoting the minimizer by  $\tilde{\mu}_t(a^*)$ , i.e.,

$$\tilde{\mu}_t(a^*) \triangleq \arg \min_{x \in [\mu(i), \mu_t(a^*)]} \{T_t(i) d_i(\mu_t(i) \| x) + T_t(a^*) d_{a^*}(\mu_t(a^*) \| x)\}, \quad (87)$$

we obtain (86). Furthermore  $\mu(i)$  and  $\mu(a^*)$  satisfies the constraint in (14), which yields (86). Next, using Lemma 4, (86) can be upper bounded as

$$\begin{aligned} \mathbb{P}_\nu(\tau < +\infty, \hat{A}_\tau \neq a^*) &\leq \sum_{i \neq a^*} \mathbb{P}_\nu \left( \exists t \in \mathbb{N} : \log \underbrace{M_t(i) M_t(a^*)}_{\triangleq M_t} + \frac{1}{2} \log T_t(i) T_t(a^*) + \frac{1}{2} \log \mathcal{I}_i(\mu_t(i)) \mathcal{I}_{a^*}(\mu_t(a^*)) \right. \\ &\quad + 2 \log \frac{|\Theta|}{\sqrt{2\pi}} + T_t(i) \max\{d_i(\mu_t(i) \| \mu_t(i) - \varepsilon), d_i(\mu_t(i) \| \mu_t(i) + \varepsilon)\} \\ &\quad + T_t(a^*) \max\{d_{a^*}(\mu_t(a^*) \| \mu_t(a^*) - \varepsilon), d_{a^*}(\mu_t(a^*) \| \mu_t(a^*) + \varepsilon)\} \\ &\quad \left. - W_t(\varepsilon, i) - W_t(a^*) \geq \beta_t(\delta) \right) \\ &\leq \sum_{i \neq a^*} \mathbb{P}_\nu \left( \exists t \in \mathbb{N} : \log M_t + \log \frac{t}{2} + \max_{i \in [K]} \mathcal{I}_i(\mu_t(i)) - 2 \cdot \min_{i \in [K]} W_t(\varepsilon, i) \right) \end{aligned} \quad (88)$$

$$+ t \cdot \max_{i \in [K]} \left\{ \max\{d_i(\mu_t(i) \parallel \mu_t(i) - \varepsilon), d_i(\mu_t(i) \parallel \mu_t(i) + \varepsilon)\} + 2 \log \frac{|\Theta|}{\sqrt{2\pi}} \geq \beta_t(\delta) \right\}, \quad (89)$$

where (89) is obtained by using the AM-GM inequality. Next, recalling the definitions of  $W_t(\varepsilon)$  in (34) and  $\varepsilon_t$  in (35), let us set

$$\beta_t(\delta) \triangleq W_t(\varepsilon) + t\varepsilon_t + 2 \log \frac{|\Theta|}{\sqrt{2\pi}} + \log \frac{t(K-1)}{2\delta}. \quad (90)$$

Furthermore, note that  $M_t$  is a martingale. To verify this, WLOG, let us assume that  $A_t = i$ . We have

$$\mathbb{E}[M_t \mid \mathcal{F}_{t-1}] = \mathbb{E}[M_t(i) \cdot M_t(a^*) \mid \mathcal{F}_{t-1}] \quad (91)$$

$$= M_{t-1}(a^*) \cdot \mathbb{E}[M_t(i) \mid \mathcal{F}_{t-1}] \quad (92)$$

$$= M_{t-1}(a^*) \cdot M_{t-1}(i). \quad (93)$$

Furthermore,  $\mathbb{E}[M_1] = E[M_1(a^*) \cdot M_1(i)] = \mathbb{E}[M_1(a^*)] \cdot \mathbb{E}[M_1(i)] = 1$ . Finally, combining (89), (90) and (93), and using Ville's inequality, we obtain

$$\mathbb{P}_\nu \left( \tau < +\infty, \hat{A}_\tau \neq a^* \right) \leq \sum_{i \neq a^*} \mathbb{P}_\nu \left( \exists t \in \mathbb{N} : \log M_t \geq \log \frac{K-1}{\delta} \right) \quad (94)$$

$$\leq \sum_{i \neq a^*} \frac{\delta}{K-1} \mathbb{E}[M_1] \quad (95)$$

$$= \delta. \quad (96)$$

This concludes the proof.

## B Problem Complexity

### B.1 Proof of Lemma 1

Recall that corresponding to a bandit instance  $\nu \in \mathcal{M}$  with the best arm  $a^*$ , the set of alternate bandit instances is defined as

$$\text{alt}(a^*) \triangleq \left\{ \bar{\nu} \in \mathcal{M} : m(\bar{\mathbb{P}}_{a^*}) \leq \max_{i \neq a^*} m(\bar{\mathbb{P}}_i) \right\}. \quad (97)$$

It can be readily verified that the set of alternate bandit instances can be equivalently stated as

$$\text{alt}(a^*) = \bigcup_{i \neq a^*} \left\{ \bar{\nu} \in \mathcal{M} : m(\bar{\mathbb{P}}_i) \geq m(\bar{\mathbb{P}}_{a^*}) \right\}. \quad (98)$$

Using (98), the problem complexity can be simplified as follows:

$$\Gamma(\nu) = \sup_{\mathbf{w} \in \Delta^K} \inf_{\bar{\nu} \in \text{alt}(a^*)} \sum_{i \in [K]} w_i D_{\text{KL}}(\mathbb{P}_i \parallel \bar{\mathbb{P}}_i) \quad (99)$$

$$= \sup_{\mathbf{w} \in \Delta^K} \min_{i \neq a^*} \inf_{\bar{\nu} \in \mathcal{M} : m(\bar{\mathbb{P}}_i) \geq m(\bar{\mathbb{P}}_{a^*})} \left\{ w_{a^*} D_{\text{KL}}(\mathbb{P}_{a^*} \parallel \bar{\mathbb{P}}_{a^*}) + w_i D_{\text{KL}}(\mathbb{P}_i \parallel \bar{\mathbb{P}}_i) \right\}. \quad (100)$$



Let us define

$$\Gamma_i(\boldsymbol{\nu}, \mathbf{w}) \triangleq \inf_{\boldsymbol{\nu} \in \mathcal{M}: m(\mathbb{P}_i) \geq m(\mathbb{P}_{a^*})} \left\{ w_{a^*} D_{\text{KL}}(\mathbb{P}_{a^*} \|\bar{\mathbb{P}}) + w_i D_{\text{KL}}(\mathbb{P}_i \|\bar{\mathbb{P}}) \right\}. \quad (101)$$

It can be readily verified that

$$\Gamma(\boldsymbol{\nu}, \mathbf{w}) \triangleq \min_{i \neq a^*} \Gamma_i(\boldsymbol{\nu}, \mathbf{w}). \quad (102)$$

Note that  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  can be further simplified as:

$$\Gamma_i(\boldsymbol{\nu}, \mathbf{w}) = \inf_{x \in \mathbb{R}} \left\{ w_{a^*} \inf_{\bar{\mathbb{P}} \in \mathcal{P}(\Omega): m(\bar{\mathbb{P}}) \leq x} D_{\text{KL}}(\mathbb{P}_{a^*} \|\bar{\mathbb{P}}) + w_i \inf_{\bar{\mathbb{P}} \in \mathcal{P}(\Omega): m(\bar{\mathbb{P}}) \geq x} D_{\text{KL}}(\mathbb{P}_i \|\bar{\mathbb{P}}) \right\} \quad (103)$$

$$= \inf_{x \in \mathbb{R}} \left\{ w_{a^*} d_{\text{U}}(\mathbb{P}_{a^*}, x) + w_i d_{\text{L}}(\mathbb{P}_i, x) \right\}. \quad (104)$$

Note that  $d_{\text{U}}$  is non-increasing in  $x$  and  $d_{\text{L}}$  is non-decreasing in  $x$ . Thus, when  $x < \mu(i)$ , we have  $d_{\text{U}}(\mathbb{P}_{a^*}, x) > d_{\text{U}}(\mathbb{P}_{a^*}, \mu(i))$ , and  $d_{\text{L}}(\mathbb{P}_i, x) > d_{\text{L}}(\mathbb{P}_i, \mu(i))$ . This implies that the optimizer  $x_i^*$  of (104) should satisfy  $x_i^* \geq \mu(i)$ . Using a similar argument, we can show that  $x_i^* \leq \mu(a^*)$ . Hence, (104) can be rewritten as:

$$\Gamma_i(\boldsymbol{\nu}, \mathbf{w}) = \inf_{x \in [\mu(i), \mu(a^*)]} \left\{ w_{a^*} d_{\text{U}}(\mathbb{P}_{a^*}, x) + w_i d_{\text{L}}(\mathbb{P}_i, x) \right\}. \quad (105)$$

Finally, the problem complexity can be equivalently expressed as:

$$\Gamma(\boldsymbol{\nu}) = \sup_{\mathbf{w} \in \Delta^K} \min_{i \neq a^*} \inf_{x \in [\mu(i), \mu(a^*)]} \left\{ w_{a^*} d_{\text{U}}(\mathbb{P}_{a^*}, x) + w_i d_{\text{L}}(\mathbb{P}_i, x) \right\}. \quad (106)$$

## B.2 Proof of Lemma 2

1. First, we will prove that  $d_{\text{U}}$  and  $d_{\text{L}}$  are strictly convex in  $x$ . For any  $x \in \mathbb{R}$  and  $y \in \mathbb{R}$ , and for any  $\lambda \in [0, 1]$ , let us define

$$z \triangleq \lambda x + (1 - \lambda)y. \quad (107)$$

Furthermore, define

$$\begin{aligned} \eta_x &\triangleq \arg \inf_{\eta \in \mathcal{P}(\Omega): m(\eta) \leq x} D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta), \\ \eta_y &\triangleq \arg \inf_{\eta \in \mathcal{P}(\Omega): m(\eta) \leq y} D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta), \\ \text{and } \eta_z &\triangleq \arg \inf_{\eta \in \mathcal{P}(\Omega): m(\eta) \leq z} D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta). \end{aligned} \quad (108)$$

Furthermore, define  $\kappa_z \triangleq \lambda \eta_x + (1 - \lambda)\eta_y$ . Note that

$$m(\kappa_z) = \lambda m(\eta_x) + (1 - \lambda)m(\eta_y) \quad (109)$$

$$\leq \lambda x + (1 - \lambda)y \quad (110)$$

$$\stackrel{(107)}{=} z. \quad (111)$$

Now,

$$d_{\text{U}}(\mathbb{P}_{a^*}, z) = D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta_z) \quad (112)$$

$$\leq D_{\text{KL}}(\mathbb{P}_{a^*} \|\kappa_z) \quad (113)$$

$$< \lambda D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta_x) + (1 - \lambda) D_{\text{KL}}(\mathbb{P}_{a^*} \|\eta_y) \quad (114)$$

$$= \lambda d_{\text{U}}(\mathbb{P}_{a^*}, x) + (1 - \lambda) d_{\text{L}}(\mathbb{P}_{a^*}, y), \quad (115)$$

where (113) is a result of (111), and (114) is a result of the strict convexity of KL divergence in both arguments. Thus,  $d_{\text{U}}$  is strictly convex in  $x$ . Using a similar argument, we can prove that  $d_{\text{L}}$  is also strictly convex in  $x$ . Thus,  $g_i : \mathcal{M} \times \mathbb{R} \mapsto \mathbb{R}$ , defined as

$$g_i(\boldsymbol{\nu}, x) \triangleq w_{a^*} d_{\text{U}}(\mathbb{P}_{a^*}, x) + w_i d_{\text{L}}(\mathbb{P}_i, x), \quad (116)$$

is strictly convex in its second argument. Thus,  $g_i$  has a unique minimum in  $[\mu(i), \mu(a^*)]$ .

2. For establishing the continuity of  $\Gamma : \mathcal{M} \mapsto \mathbb{R}$  and  $\mathbf{w} : \mathcal{M} \mapsto \Delta^K$ , we will leverage the following two lemmas, which provide the sufficient conditions for continuity.

**Lemma 5** (Berge's maximum theorem [23]). *Suppose  $g$  is a continuous function on  $\mathcal{S} \times \Theta$  and  $\mathcal{D} : \Theta \mapsto \mathcal{S}$  is a compact-valued continuous correspondence on  $\Theta$ . Let*

$$g^*(\theta) \triangleq \max_{x \in \mathcal{D}(\theta)} g(x, \theta) \quad \text{and} \quad \mathcal{D}^*(\theta) \triangleq \arg \max_{x \in \mathcal{D}(\theta)} g(x, \theta). \quad (117)$$

*Then,  $g^*$  is a continuous function on  $\Theta$ , and  $\mathcal{D}^*$  is a compact-valued upper semicontinuous correspondence on  $\Theta$ .*

**Lemma 6** ([24]). *Let us denote the generalized Jensen-Shannon (JS) divergence between two measures  $\mathbb{P}_1$  and  $\mathbb{P}_2$  with weight  $\alpha \in (0, 1)$  by*

$$\text{JS}_{\alpha}(\mathbb{P}_1 \|\mathbb{P}_2) \triangleq \alpha D_{\text{KL}}(\mathbb{P}_1 \|\bar{\mathbb{P}}_{\alpha}) + (1 - \alpha) D_{\text{KL}}(\mathbb{P}_2 \|\bar{\mathbb{P}}_{\alpha}), \quad (118)$$

*where we have defined*

$$\bar{\mathbb{P}}_{\alpha} \triangleq \alpha \mathbb{P}_1 + (1 - \alpha) \mathbb{P}_2. \quad (119)$$

*Then,  $\text{JS}_{\alpha}$  is upper-bounded as*

$$\text{JS}_{\alpha}(\mathbb{P}_1 \|\mathbb{P}_2) \leq 1. \quad (120)$$

Now, we show that  $\Gamma(\boldsymbol{\nu})$  and  $\mathbf{w}(\boldsymbol{\nu})$  is continuous in  $\boldsymbol{\nu}$ . For this, let us define the correspondence  $\mathcal{D} : \mathcal{M} \mapsto \Delta^K$  such that for any  $\boldsymbol{\nu} \in \mathcal{M}$ ,  $\mathcal{D}(\boldsymbol{\nu}) \triangleq \Delta^K$ . For any  $\boldsymbol{\nu} \in \mathcal{M}$ ,  $\mathcal{D}(\boldsymbol{\nu})$  is a compact set; hence,  $\mathcal{D}$  is a compact-valued constant correspondence. Finally, we need to show that for each  $i \in [K] \setminus \{a^*\}$ ,  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is continuous in  $\boldsymbol{\nu}$  and  $\mathbf{w}$ , where we have defined  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  in (105). First, note that  $\Gamma_i$  is lower semicontinuous in  $\boldsymbol{\nu}$  due to the lower semicontinuity of KL divergence in both arguments [25]. Next, we will leverage [26, Theorem 5.43], which provides a sufficient condition for the global continuity of convex functions.

**Lemma 7** ([26]). *For a convex function  $f : \mathcal{X} \mapsto \mathbb{R}$  on an open convex subset of a topological vector space, the following statements are equivalent.*

- (a)  *$f$  is bounded above on a neighborhood of some point in  $\mathcal{X}$ .*

(b)  $f$  is upper semicontinuous on  $\mathcal{X}$ .

Note that  $\Gamma_i$  is convex in its first argument since KL divergence is a convex function. Let us denote the interior of the set of distributions  $\mathcal{M}$  by  $\text{int}(\mathcal{M})$ . For any  $\boldsymbol{\eta} \in \text{int}(\mathcal{M})$ , there exists a neighborhood  $\mathcal{N}_r(\boldsymbol{\eta}) \subset \text{int}(\mathcal{M})$ , where we have defined

$$\mathcal{N}_r(\boldsymbol{\eta}) \triangleq \left\{ \boldsymbol{\lambda} \in \text{int}(\mathcal{M}) : D_{\text{TV}}(\boldsymbol{\eta} \parallel \boldsymbol{\lambda}) < r \right\}. \quad (121)$$

Furthermore, for any  $\boldsymbol{\nu} \in \mathcal{N}_r(\boldsymbol{\eta})$  and  $\mathbf{w} \in \Delta^K$ , let us define the distribution

$$\kappa_{a^*,i} \triangleq \frac{w_{a^*} \mathbb{P}_{a^*} + w_i \mathbb{P}_i}{w_{a^*} + w_i}. \quad (122)$$

Expanding  $\Gamma_i$ , we obtain

$$\Gamma_i(\boldsymbol{\nu}, \mathbf{w}) = \inf_{\bar{\boldsymbol{\nu}} \in \mathcal{M}: m(\bar{\mathbb{P}}_{a^*}) \leq m(\mathbb{P}_i)} \left\{ w_{a^*} D_{\text{KL}}(\mathbb{P}_{a^*} \parallel \bar{\mathbb{P}}_{a^*}) + w_i D_{\text{KL}}(\mathbb{P}_i \parallel \bar{\mathbb{P}}_i) \right\} \quad (123)$$

$$\leq w_{a^*} D_{\text{KL}}(\mathbb{P}_{a^*} \parallel \kappa_{a^*,i}) + w_i D_{\text{KL}}(\mathbb{P}_i \parallel \kappa_{a^*,i}) \quad (124)$$

$$= (w_{a^*} + w_i) \left( \frac{w_{a^*}}{w_{a^*} + w_i} D_{\text{KL}}(\mathbb{P}_{a^*} \parallel \kappa_{a^*,i}) + \frac{w_i}{w_{a^*} + w_i} D_{\text{KL}}(\mathbb{P}_i \parallel \kappa_{a^*,i}) \right) \quad (125)$$

$$= (w_{a^*} + w_i) \text{JS}_{\frac{w_{a^*}}{w_{a^*} + w_i}}(\mathbb{P}_{a^*} \parallel \mathbb{P}_i) \quad (126)$$

$$\leq 1, \quad (127)$$

where (127) is a result of Lemma 6. Thus, leveraging Lemma 7, we obtain that  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is upper semicontinuous in  $\boldsymbol{\nu}$ , which proves that  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is continuous in its first argument in  $\mathcal{M}$ . Furthermore,  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is linear in the second argument, and hence, it is continuous. This shows that  $\Gamma_i$  and the correspondence  $\mathcal{D}$  satisfies the conditions in Lemma 5, and we obtain that  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is continuous in  $\boldsymbol{\nu}$  and  $\mathbf{w}$  is upper hemicontinuous in  $\boldsymbol{\nu}$ . Finally, following the same line of arguments as [10, Proposition 7], we can show that for a given  $\boldsymbol{\nu} \in \mathcal{M}$ ,  $\mathbf{w}$  is the unique allocation satisfying

$$\Gamma_i(\boldsymbol{\nu}, \mathbf{w}) = \Gamma_j(\boldsymbol{\nu}, \mathbf{w}), \quad \text{for all } i, j \neq a^*. \quad (128)$$

This shows that  $\mathbf{w}$  is continuous in  $\boldsymbol{\nu}$ .

## C Proof of Theorems 2 and 3

First, we show that the explicit exploration phase ensures that each arm  $i \in [K]$  is sampled sufficiently often, such that the sample mean values converge to the true means. Let us define the time instant  $N_{\boldsymbol{\nu}}^\epsilon$  as

$$N_{\boldsymbol{\nu}}^\epsilon \triangleq \inf \left\{ t \in \mathbb{N} : |\bar{\mu}_s(i) - \mu(i)| < \epsilon, \forall i \in [K], \forall s \geq t \right\}. \quad (129)$$

The stochastic time  $N_{\boldsymbol{\nu}}^\epsilon$  marks the convergence of the sample means to the respective ground truths for every arm  $i \in [K]$ . In the following result, we will show that the TCB and ITCB arm selection strategies ensure that  $N_{\boldsymbol{\nu}}^\epsilon$  has a finite average value. This result is instrumental in showing the convergence in allocation for the TCB and ITCB sampling strategies stated in Theorem 2 and Theorem 3.

**Theorem 5** (Convergence in mean). *Under the TCB and ITCB sampling strategies, we have  $\mathbb{E}_\nu[N_\nu^\epsilon] < +\infty$ .*

*Proof.* We use the notion of  $r$ -quick convergence, which we define below.

**Definition 2** ( $r$ -quick convergence [27]). *Consider the sequence of i.i.d. zero-mean random variables  $\{Z_t : t \in \mathbb{N}\}$ . Let  $\bar{Z}_t \triangleq \frac{1}{t} \sum_{s=1}^t Z_s$  denote the empirical mean. Furthermore, for any  $\epsilon \in \mathbb{R}_+$  define*

$$T_\epsilon \triangleq \sup \{t \in \mathbb{N} : |\bar{Z}_t| > \epsilon\} . \quad (130)$$

*Then,  $\{Z_t : t \in \mathbb{N}\}$  converges  $r$ -quickly for  $r > 0$ , if  $\mathbb{E}[T_\epsilon^r] < +\infty$ .*

We leverage  $r$ -quick convergence for  $r = 2$  to establish the convergence of the sample means of each arm  $i \in [K]$  to the corresponding ground truth values. For this, we first state the necessary and sufficient condition for  $r$ -quick convergence to hold.

**Lemma 8** (Corollary 4, [27]). *The i.i.d. sequence  $\{Z_t : t \in \mathbb{N}\}$  converges  $r$ -quickly if and only if  $\mathbb{E}[|Z_t|^{r+1}] < +\infty$ .*

For any arm  $i \in [K]$ , let us denote the realizations of  $T_t(i)$  by  $\ell_t(i)$ . Furthermore, for any  $s \in [t]$  such that  $A_s = i$ , let us set

$$Z_s(i) \triangleq X_s(i) - \mu(i) , \quad \text{and} \quad \bar{Z}_t(i) \triangleq \frac{1}{T_t(i)} \sum_{s \in [t]: A_s = i} (X_s(i) - \mu(i)) . \quad (131)$$

Let us define

$$T_\epsilon(i) \triangleq \sup \{\ell_t(i) \in \mathbb{N} : |\bar{Z}_t(i)| > \epsilon\} . \quad (132)$$

Using Assumption 4 along with Lemma 8, we have

$$\mathbb{E}_\nu \left[ (T_\epsilon(i))^2 \right] < +\infty . \quad (133)$$

Furthermore, note that owing to the explicit exploration of the TCB and ITCB algorithms, leveraging [4, Lemma 8], for any arm  $i \in [K]$ , we have

$$T_t(i) \geq \sqrt{\frac{t}{K}} - 1 . \quad (134)$$

Accordingly, when arm  $i \in [K]$  has been sampled  $T_\epsilon(i)$  times, the following inequality holds

$$T_\epsilon(i) \geq \sqrt{\frac{N_0}{K}} - 1 , \quad (135)$$

where  $N_0$  denotes the time instant at which the arm  $i \in [K]$  has been sampled  $T_\epsilon(i)$  times. Furthermore, assuming that  $N_0 \geq 4K$ , we have

$$T_\epsilon(i) \geq \sqrt{\frac{N_0}{2K}} , \quad (136)$$

which yields that  $N_0 \leq 2K(T_\epsilon(i))^2$ . Furthermore, from (133) we have  $\mathbb{E}_\nu[(T_\epsilon(i))^2] < +\infty$ . The proof is completed by setting  $N_\nu^\epsilon = N_0$ . ■

Next, let us define the set of *over-sampled* arms as:

$$\mathcal{O}_t^\epsilon \triangleq \left\{ i \in [K] : \frac{T_t(i)}{t} > w_i(\boldsymbol{\nu}) + \epsilon \right\}. \quad (137)$$

Furthermore, we define the set of *under-sampled* arms as:

$$\mathcal{P}_t^\epsilon \triangleq \left\{ i \in [K] : \frac{T_t(i)}{t} < w_i(\boldsymbol{\nu}) + \frac{\epsilon}{2} \right\}. \quad (138)$$

The convergence in allocation for the proposed algorithm is shown in two key steps. First, we prove that if any sampling strategy always samples from the set of under-sampled arms, then the sampling strategy converges to the optimal allocation  $\mathbf{w}(\boldsymbol{\nu})$ . This step is common in the proof for both Theorem 2 and Theorem 3. In the next step, we show that the proposed sampling strategies always sample from the set of under-sampled arms. We show the first step through Lemma 9 and Lemma 10, which we provide next. Essentially, Lemma 9 shows that if the sampling strategy always samples from the set of under-sampled arms, then, after some time, the set of over-sampled arms becomes empty. Lemma 10 then shows that when the set of over-sampled arms is empty, eventually, the allocation for each arm converges to the optimal allocation. The key distinction in the proofs of Theorem 2 and Theorem 3 arises in the next step. In Lemma 11, we show that the TCB arm selection rule stated in (28) always samples from the set of under-sampled arms. In Lemma 12, we show that the ITCB arm selection rule provided in (31) always samples from the set of under-sampled arms. Before stating Lemma 9, let us define the sampling proportion  $\gamma_t \triangleq [\gamma_{t,1}, \dots, \gamma_{t,K}]$  computed at the current MLE  $\boldsymbol{\mu}_t$  as

$$\gamma_t \triangleq \arg \sup_{\mathbf{w} \in \Delta^K} \inf_{\bar{\boldsymbol{\nu}} \in \text{alt}(a^*)} \sum_{i \in [K]} w_i D_{\text{KL}}(\mathbb{P}_{t,i} \| \bar{\mathbb{P}}_i). \quad (139)$$

**Lemma 9.** *There exists a stochastic time  $N^\epsilon \in \mathbb{N}$  such that for all  $t > N^\epsilon$ ,  $\mathcal{O}_t^\epsilon = \emptyset$ , and  $\mathbb{E}[N^\epsilon] < +\infty$ , if the sampling strategy satisfies  $\frac{1}{t}T_t(a_{t+1}) < \gamma_{t,a_{t+1}} + \zeta$  for any  $\zeta \in [0, \frac{\epsilon}{4}]$  and for any  $t > M$ , where  $M$  is a stochastic time satisfying  $\mathbb{E}_{\boldsymbol{\nu}}[M] < +\infty$ .*

*Proof.* Let us define the time instant  $M_1^\epsilon$  such that for all  $t > M_1^\epsilon$  and for all  $i \in [K]$ ,  $|\gamma_{t,i} - w_i(\boldsymbol{\nu})| < \epsilon/8$ . Leveraging the continuity of  $\Gamma$  in Lemma 2 and the convergence of the MLE in Theorem 5, we obtain that  $\mathbb{E}[M_1^\epsilon] < +\infty$ . Furthermore, define  $M_2^\epsilon \triangleq \lceil (8/\epsilon) - 1 \rceil$ , and  $M^\epsilon \triangleq \max\{M, M_1^\epsilon, M_2^\epsilon\}$ . We have the following two cases:

1.  $\mathcal{O}_{M^\epsilon}^\epsilon = \emptyset$ : In this case, we will use induction on  $t$  to show that for all  $t > M^\epsilon$ ,  $\mathcal{O}_t^\epsilon = \emptyset$ . First, by our assumption, for  $t = M^\epsilon$ ,  $\mathcal{O}_{M^\epsilon}^\epsilon = \emptyset$ . Next, assume the inductive hypothesis that for some  $t > M^\epsilon$ ,  $\mathcal{O}_t^\epsilon = \emptyset$ . Then,

$$\frac{T_{t+1}(a_{t+1})}{t+1} = \frac{T_t(a_{t+1}) + 1}{t+1} \quad (140)$$

$$< \frac{T_t(a_{t+1})}{t} + \frac{1}{t+1} \quad (141)$$

$$< \gamma_{t,a_{t+1}} + \frac{1}{t+1} + \zeta \quad (142)$$

$$\leq w_{a_{t+1}}(\boldsymbol{\nu}) + \frac{\epsilon}{8} + \frac{1}{t+1} + \zeta \quad (143)$$

$$\leq w_{a_{t+1}}(\boldsymbol{\nu}) + \frac{\epsilon}{2}, \quad (144)$$

where (142) holds since the sampling strategy satisfies  $\frac{1}{t}T_t(a_{t+1}) < \gamma_{t,a_{t+1}} + \zeta$ , and (144) is obtained using the definition of  $M^\epsilon$ . Hence,  $\mathcal{O}_{t+1}^\epsilon = \emptyset$ , and it concludes the proof.

2.  $|\mathcal{O}_{M^\epsilon}^\epsilon| \geq 1$ : In this case, following the same steps as (142)-(144), we can show that for all  $t > M^\epsilon$ , any  $i \in \mathcal{P}_t^\epsilon$  is not included in  $\mathcal{O}_t^\epsilon$ . Furthermore, for any  $t > M^\epsilon$  and for any  $j \in \mathcal{O}_t^\epsilon$ , let us define  $L_j^\epsilon$  as the time that  $j$  leaves  $\mathcal{O}_t^\epsilon$ , i.e., for all  $t \in \{M^\epsilon, \dots, L_j^\epsilon - 1\}$ ,  $j \in \mathcal{O}_t^\epsilon$ . Next, defining  $L^\epsilon \triangleq \max_{j \in [K]} L_j^\epsilon$ , for all  $t > L^\epsilon$ , we obtain that  $|\mathcal{O}_t^\epsilon| = 0$ . Finally, defining  $N^\epsilon \triangleq \max\{M^\epsilon, L^\epsilon\}$ , we obtain that for all  $t > N^\epsilon$ ,  $\mathcal{O}_t^\epsilon = \emptyset$ . ■

**Lemma 10.** For all  $t > N^{\frac{\epsilon}{K}}$ , the allocation for every arm  $i \in [K]$  satisfies

$$\left| \frac{T_t(i)}{t} - w_i(\boldsymbol{\nu}) \right| \leq \epsilon. \quad (145)$$

*Proof.* We will prove (145) by contradiction. Assume that there exists  $j \in [K]$  such that  $\frac{1}{t}T_t(i) < w_j(\boldsymbol{\nu}) - \epsilon$ . For all  $t > N^{\frac{\epsilon}{K}}$ , leveraging Lemma 9, we have

$$\sum_{i \in [K]} \frac{T_t(i)}{t} = \sum_{i \neq j} \frac{T_t(i)}{t} + \frac{T_t(j)}{t} \quad (146)$$

$$\leq \sum_{i \neq j} \left( w_i(\boldsymbol{\nu}) + \frac{\epsilon}{K} \right) + w_j(\boldsymbol{\nu}) - \epsilon \quad (147)$$

$$= 1 - \frac{\epsilon}{K}, \quad (148)$$

which is a contradiction. Thus, (145) holds for all  $t > N^{\frac{\epsilon}{K}}$ . ■

Next, we show that our proposed sampling strategy always samples from the set of under-sampled arms, which is specified in Lemma 11. Let us define the minimum sub-optimality gap

$$\Delta_{\min} \triangleq \min_{i \in [K] \setminus \{a^*\}} \mu(a^*) - \mu(i). \quad (149)$$

**Lemma 11.** For all  $t > N_\nu^{\Delta_{\min}/4}$ , the TCB sampling rule provided in (28) satisfies

$$\frac{T_t(a_{t+1})}{t} \leq \gamma_{t, a_{t+1}}. \quad (150)$$

*Proof.* Note that by Theorem 5, for all  $t > N_\nu^{\Delta_{\min}/4}$ , the proposed sampling rule satisfies that  $a_t^{\text{top}} = a^*$ . Accordingly, the TCB sampling rule samples between any two arms, either the best arm  $a^*$  or the arm  $a_t^{\min}$ . Furthermore, note that both the arms  $a^*$  and  $a_t^{\min}$  cannot be simultaneously over-sampled. To see why this is true, assume without loss of generality that arm  $a^*$  is over-sampled, i.e.,  $\frac{1}{t}T_t(a^*) > \gamma_{t, a^*}$ .  $\Gamma_t(\mathbf{w})$  is a minimum of linear functions, and hence, it is a concave function [28]. Furthermore,  $\mathbf{w}$  belongs to a compact space  $\Delta^K$ , and  $\Gamma_t(\mathbf{w})$  has a *unique* maxima [10]. We will prove that  $\frac{1}{t}T_t(a_t^{\min}) < \gamma_{t, a_t^{\min}}$  by contradiction. Let us assume that  $\frac{1}{t}T_t(a_t^{\min}) > \gamma_{t, a_t^{\min}}$ . We have

$$\begin{aligned} & \min_{x \in I_{t, a_t^{\min}}} \left\{ \frac{T_t(a^*)}{t} d_{\text{U}}(\mathbb{P}_{t, a^*}, x) + \frac{T_t(a_t^{\min})}{t} d_{\text{L}}(\mathbb{P}_{t, a_t^{\min}}, x) \right\} \\ & \geq \min_{x \in I_{t, a_t^{\min}}} \left\{ \gamma_{t, a^*} d_{\text{U}}(\mathbb{P}_{t, a^*}, x) + \frac{T_t(a_t^{\min})}{t} d_{\text{L}}(\mathbb{P}_{t, a_t^{\min}}, x) \right\} \end{aligned} \quad (151)$$

$$\geq \min_{x \in I_{t, a_t^{\min}}} \left\{ \gamma_{t, a^*} d_{\text{U}}(\mathbb{P}_{t, a^*}, x) + \gamma_{t, a_t^{\min}} d_{\text{L}}(\mathbb{P}_{t, a_t^{\min}}, x) \right\} \quad (152)$$

$$= \Gamma_t(\boldsymbol{\gamma}_t), \quad (153)$$

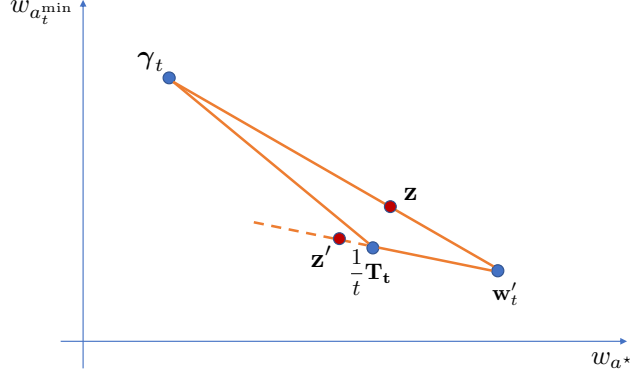


Figure 7: Positions of  $\mathbf{z}$  and  $\mathbf{z}'$  assuming  $\frac{1}{t}T_t(a^*) > \gamma_{t,a^*}$

where (151) and (152) hold due to the fact that  $\Gamma_t$  is an increasing function in each of its arguments, keeping the other argument fixed [10, Lemma 2]. It can be readily verified that (153) is a contradiction, since we obtain that  $\Gamma_t(\frac{1}{t}\mathbf{T}_t) \geq \Gamma_t(\gamma_t)$ , where  $\gamma_t$  is the *unique* maximizer. Thus, we have  $\frac{1}{t}T_t(a_t^{\min}) < \gamma_{t,a_t^{\min}}$ , if  $\frac{1}{t}T_t(a^*) > \gamma_{t,a^*}$ . Next, for any  $t > N_{\nu}^{\Delta_{\min}/4}$ , based on the TCB arm selection strategy, we have the following two cases.

1.  $a_{t+1} = a^*$ : We will prove that  $\frac{1}{t}T_t(a^*) \leq \gamma_{t,a^*}$  by contradiction. We proceed with our assumption that  $\frac{1}{t}T_t(a^*) > \gamma_{t,a^*}$ , and define the points  $\mathbf{z}$  and  $\mathbf{z}'$  such that for any  $\lambda_1, \lambda_2 \in (0, 1)$ ,

$$\mathbf{z} \triangleq \lambda_1 \gamma_t + (1 - \lambda_1) \mathbf{w}'_t, \quad (154)$$

$$\text{and } \frac{1}{t}\mathbf{T}_t = \lambda_2 \mathbf{z}' + (1 - \lambda_2) \mathbf{w}'_t. \quad (155)$$

For a geometric representation of the relative position of these points, we refer to Figure 7. Owing to the concavity of  $\Gamma_t$ , we have

$$\Gamma_t(\mathbf{z}) \geq \lambda_1 \Gamma_t(\gamma_t) + (1 - \lambda_1) \Gamma_t(\mathbf{w}'_t) \quad (156)$$

$$\geq \Gamma_t(\mathbf{w}'_t). \quad (157)$$

Note that as a result of the TCB sampling rule, we have

$$\Gamma_t(\mathbf{w}'_t) > \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right). \quad (158)$$

Accordingly, let us define

$$\epsilon_t \triangleq \Gamma_t(\mathbf{w}'_t) - \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right). \quad (159)$$

Furthermore, for any  $\mathbf{w} \in \Delta^K$ , let  $\nabla\Gamma_t(\mathbf{w})$  denote the sub-gradient of the function  $\Gamma_t$  at  $\mathbf{w}$ . Owing to the concavity of  $\Gamma_t$ , we have

$$\Gamma_t(\mathbf{z}') \geq \Gamma_t(\mathbf{z}) - \langle \nabla\Gamma_t(\mathbf{z}), \mathbf{z}' - \mathbf{z} \rangle \quad (160)$$

$$\geq \Gamma_t(\mathbf{z}) - \|\nabla\Gamma_t(\mathbf{z})\| \|\mathbf{z}' - \mathbf{z}\| \quad (161)$$

$$\stackrel{(157)}{\geq} \Gamma_t(\mathbf{w}'_t) - \|\nabla\Gamma_t(\mathbf{z})\| \|\mathbf{z}' - \mathbf{z}\| \quad (162)$$

$$\stackrel{(159)}{=} \Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) + \epsilon_t - \|\nabla \Gamma_t(\mathbf{z})\| \|\mathbf{z}' - \mathbf{z}\|, \quad (163)$$

where (161) is obtained using the Cauchy–Schwarz inequality. Furthermore, leveraging the concavity of  $\Gamma_t$ , we have

$$\Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) \geq \lambda_2 \Gamma_t(\mathbf{z}') + (1 - \lambda_2) \Gamma_t(\mathbf{w}'_t) \quad (164)$$

$$\stackrel{(163)}{\geq} \lambda_2 \left( \Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) + \epsilon_t - \|\nabla \Gamma_t(\mathbf{z})\| \|\mathbf{z}' - \mathbf{z}\| \right) + (1 - \lambda_2) \Gamma_t(\mathbf{w}'_t), \quad (165)$$

which implies that

$$\Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) \geq \Gamma_t(\mathbf{w}'_t) - \frac{\lambda_2}{1 - \lambda_2} (\|\nabla \Gamma_t(\mathbf{z})\| \|\mathbf{z}' - \mathbf{z}\| - \epsilon_t). \quad (166)$$

Let us set  $\lambda_2 \triangleq O(\epsilon_t^2)$ . Hence, (166) can be rewritten as

$$\Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) \geq \Gamma_t(\mathbf{w}'_t) + O(\epsilon_t^2) \quad (167)$$

$$\stackrel{(159)}{=} \Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) + \epsilon_t + O(\epsilon_t^2), \quad (168)$$

which is a contradiction. This shows that when  $a_{t+1} = a^*$ , we have  $\frac{1}{t} T_t(a_{t+1}) \leq \gamma_{t, a_{t+1}}$ .

2.  $a_{t+1} = a_t^{\min}$ : Let us assume that  $\frac{1}{t} T_t(a_t^{\min}) > \gamma_{t, a_t^{\min}}$ . Furthermore, by our sampling strategy, we have

$$\Gamma_t(\mathbf{w}'_t) < \Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right). \quad (169)$$

Following similar arguments as the case when  $a_{t+1} = a^*$ , leveraging the concavity of  $\Gamma_t$ , we can arrive at a contradiction. Thus, in this case, when  $a_{t+1} = a_t^{\min}$ , we have  $\frac{1}{t} T_t(a_{t+1}) \leq \gamma_{t, a_{t+1}}$ . ■

**Lemma 12.** *There exists a stochastic time  $M_{\text{ITCB}}$  such that  $\mathbb{E}_\nu[M_{\text{ITCB}}] < +\infty$ , and for all  $t > M_{\text{ITCB}}$ , the ITCB sampling rule provided in (31) with the sequence  $r_t = \frac{\epsilon}{t}$  for any  $\epsilon \in \mathbb{R}_+$  satisfies*

$$\frac{T_t(a_{t+1})}{t} \leq \gamma_{t, a_{t+1}}. \quad (170)$$

*Proof.* Let us recall that

$$b_t^{\min} \triangleq \arg \min_{i \in [K] \setminus \{a_t^{\text{top}}\}} \left\{ \min_{x \in I_{t,i}} \left\{ \frac{T_t(a_t^{\text{top}})}{t} d_{\text{U}}(\mathbb{P}_{t, a_t^{\text{top}}}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} + \frac{\log(T_t(i))}{t} \right\}. \quad (171)$$

Note that for all  $t > N_\nu^{\Delta_{\min}/4}$ , we have  $a_t^{\text{top}} = a^*$ . Furthermore, for all  $t > \lceil \frac{1}{K\epsilon^2} \rceil$  and for all  $i \in [K]$ , we almost surely have

$$\frac{\log T_t(i)}{t} \leq \frac{\log(\sqrt{t/K} - 1)}{t} \quad (172)$$

$$\leq \frac{\log(\sqrt{t/K}) + 1}{t} \quad (173)$$



$$\leq \frac{\sqrt{t/K}}{t} \quad (174)$$

$$\leq \epsilon, \quad (175)$$

where (172) is a result of the fact that  $T_t(i) \geq \sqrt{t/K} - 1$  for all  $i \in [K]$  [7, Lemma 4], (174) holds due to the fact that  $\log(1+x) \leq x$ , and (175) is obtained from the fact that  $t > \lceil 1/(K\epsilon^2) \rceil$ . Let us define  $M_3^\epsilon \triangleq \max\{N_{\nu}^{\Delta \min/4}, 1/(K\epsilon^2)\}$ . For a sufficiently small  $\epsilon \in \mathbb{R}_+$ , for any  $t > M_3^\epsilon$ , we almost surely have

$$b_t^{\min} = \arg \min_{i \in [K] \setminus \{a^*\}} \left\{ \min_{x \in I_{t,i}} \left\{ \frac{T_t(a^*)}{t} d_U(\mathbb{P}_{t,a^*}, x) + \frac{T_t(i)}{i} d_L(\mathbb{P}_{t,i}, x) \right\} + \epsilon \right\} \quad (176)$$

$$= a_t^{\min}. \quad (177)$$

We have the following two cases.

1.  $a_{t+1} = a^*$ : Let us assume that  $\frac{1}{t}T_t(a^*) > \gamma_{t,a^*} + \zeta$ . Due to the ITCB sampling strategy in (31), for all  $t > M_3^\epsilon$ , we have

$$\Gamma_t(\mathbf{w}'_t) + \frac{\log(tw'_t(a_t^{\min}))}{t} \geq \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) + \frac{\log(T_t(a_t^{\min}))}{t}, \quad (178)$$

or, equivalently

$$\Gamma_t(\mathbf{w}'_t) - \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) \geq \frac{1}{t} \log\left(1 + \frac{\frac{tr_t}{K-1}}{T_t(a_t^{\min}) - \frac{tr_t}{K-1}}\right). \quad (179)$$

This implies that

$$\Gamma_t(\mathbf{w}'_t) - \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) > 0. \quad (180)$$

Following the same argument as Lemma 11, (180) implies that  $\frac{1}{t}T_t(a^*) \leq \gamma_{t,a^*} + \zeta$  for any  $\zeta \geq 0$ .

2.  $a_{t+1} = a_t^{\min}$ : Let us assume that  $\frac{1}{t}T_t(a_t^{\min}) > \gamma_{t,a_t^{\min}}$ . We will show that this is a contradiction if the condition in the ITCB sampling strategy in (31) holds. Specifically, according to the ITCB sampling strategy, we have

$$\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) - \Gamma_t(\mathbf{w}'_t) \geq \underbrace{-\frac{1}{t} \cdot \log\left(1 + \frac{\frac{tr_t}{K-1}}{T_t(a_t^{\min}) - \frac{tr_t}{K-1}}\right)}_{\triangleq \xi_t}. \quad (181)$$

We may have the following two cases.

- $\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) - \Gamma_t(\mathbf{w}'_t) > 0$ : In this case, following the same line of arguments as in Lemma 11 we obtain that  $\frac{1}{t}T_t(a_t^{\min}) \leq \gamma_{t,a_t^{\min}}$ .
- $\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) \in [\Gamma_t(\mathbf{w}'_t) - \frac{1}{t}\xi_t, \Gamma_t(\mathbf{w}'_t)]$ : Let us define the vector  $\mathbf{e}_i \triangleq [e(1), \dots, e(K)]^\top$ , where, for any  $j \in [K]$  we have defined

$$e(j) \triangleq \begin{cases} -1, & \text{if } j = a^* \\ 1, & \text{if } j \neq a^* \end{cases}. \quad (182)$$

Leveraging the concavity of  $\Gamma_t$ , we have

$$\Gamma_t(\mathbf{w}'_t) \geq \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) - \underbrace{\left\langle \nabla\Gamma_t(\mathbf{w}'_t), \frac{1}{t}\mathbf{T}_t - \mathbf{w}'_t \right\rangle}_{<0}, \quad (183)$$

which implies that

$$\Gamma_t(\mathbf{w}'_t) - \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) \geq \left| \left\langle \nabla\Gamma_t(\mathbf{w}'_t), \frac{1}{t}\mathbf{T}_t - \mathbf{w}'_t \right\rangle \right| \quad (184)$$

$$= r_t \cdot |\langle \nabla\Gamma_t(\mathbf{w}'_t), \mathbf{e} \rangle|. \quad (185)$$

Combining (181) and (185), we obtain

$$\xi_t \geq tr_t |\langle \nabla\Gamma_t(\mathbf{w}'_t), \mathbf{e} \rangle|. \quad (186)$$

Next, let us define the set

$$\mathcal{M}_+ \triangleq \{\mathbf{w} \in \Delta^K : |\langle \nabla\Gamma_t(\mathbf{w}), \mathbf{e} \rangle| > 0\}. \quad (187)$$

It can be readily verified that  $\mathbf{w}'_t \in \mathcal{M}_+$ . Thus, from (186), we obtain that

$$\frac{1}{tr_t}\xi_t \geq \inf_{\mathbf{w} \in \mathcal{M}_+} |\langle \nabla\Gamma_t(\mathbf{w}), \mathbf{e} \rangle|, \quad (188)$$

Next, setting  $r_t = \frac{\epsilon}{t}$ , we obtain

$$\frac{1}{tr_t}\xi_t = \frac{1}{\epsilon} \log \left( 1 - \frac{\epsilon}{K-1} \cdot \frac{1}{T_t(a_t^{\min}) - \frac{\epsilon}{K-1}} \right) \quad (189)$$

$$\leq \frac{1}{\epsilon} \log \underbrace{\left( 1 - \frac{\epsilon}{K-1} \cdot \frac{1}{(\sqrt{t/K} - 1) - \frac{\epsilon}{K-1}} \right)}_{\triangleq g(t)}, \quad (190)$$

where (190) follows from the property of forced exploration that  $T_t(i) \geq \sqrt{t/K} - 1$  for all  $i \in [K]$  [7, Lemma 4]. Next, note that for all  $t > K$ , the function  $g(t)$  is a monotonically decreasing function in  $t$ . Hence, there exists  $M_4^\epsilon \in \mathbb{N}$  such that for all  $t > M_4^\epsilon$ , we have  $g(t) \leq \epsilon^2$ . Hence, for all  $t > M_4^\epsilon$ , we have

$$\frac{1}{tr_t}\xi_t \leq \epsilon. \quad (191)$$

Furthermore, setting

$$\epsilon \triangleq \inf_{\mathbf{w} \in \mathcal{M}_+} |\langle \nabla\Gamma_t(\mathbf{w}), \mathbf{e} \rangle|, \quad (192)$$

it can be readily verified that (186) is a contradiction. This implies that for all  $t > M_4^\epsilon$ , assuming that  $\frac{1}{t}T_t(a_t^{\min}) > \gamma_{t, a_t^{\min}}$ , the ITCB sampling condition (181) does not hold, and hence  $a_{t+1} \neq a_t^{\min}$ . Finally, defining  $M_{\text{ITCB}} \triangleq \max\{M_3^\epsilon, M_4^\epsilon\}$ , it satisfies  $\frac{1}{t}T_t(a_t^{\min}) \leq \gamma_{t, a_t^{\min}}$ . ■

## D Proof of Theorem 4

The upper bound on the average sample complexity is obtained by leveraging the convergence of the empirical problem complexity  $\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right)$  as a result of our sampling strategy, to the true value  $\Gamma(\boldsymbol{\nu})$ , where we have defined  $\mathbf{T}_t \triangleq [T_t(1), \dots, T_t(K)]$ . This is stated in Lemma 13.

**Lemma 13.** *Under TCB and ITCB, for any  $\epsilon \in \mathbb{R}_+$ , there exists  $N_\epsilon$  such that for all  $t \geq N_\epsilon$ , we have*

$$\left| \Gamma(\boldsymbol{\nu}) - \Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) \right| \leq \epsilon, \quad (193)$$

and  $\mathbb{E}_\nu[N_\epsilon] < +\infty$ .

*Proof.* For any  $\epsilon' > 0$ , let us define the time  $N_1^{\epsilon'} \triangleq \max\{N_\nu^{\epsilon'}, N_{\mathbf{w}}^{\epsilon'}\}$ . For all  $t > N_1^{\epsilon'}$ , we have:

1.  $\mu_t(i) \in [\mu(i) - \epsilon', \mu(i) + \epsilon']$  for every arm  $i \in [K]$ .
2.  $\frac{1}{t}T_t(i) \in [w_i(\boldsymbol{\nu}) - \epsilon', w_i(\boldsymbol{\nu}) + \epsilon']$  for every arm  $i \in [K]$ .
3. Let  $\boldsymbol{\nu}_t \triangleq [\mathbb{P}_{t,1}, \dots, \mathbb{P}_{t,K}]$  denote the bandit instance characterized by the mean values  $m(\boldsymbol{\nu}_t) = [\mu_t(1), \dots, \mu_t(K)]$ . As a result of the continuity of  $\Gamma(\boldsymbol{\nu}, \mathbf{w})$  in its first argument established in Lemma 2, for all  $t > N_1^{\epsilon'}$ , there exists  $\epsilon''$  such that we have  $|\Gamma(\boldsymbol{\nu}, \mathbf{w}) - \Gamma(\boldsymbol{\nu}_t, \mathbf{w})| < \epsilon''$ .

Thus, for all  $t > N_1^{\epsilon'}$ , we have

$$\Gamma_t\left(\frac{1}{t}\mathbf{T}\right) = \min_{i \neq a^*} \min_{x \in I_{t,i}} \left\{ \frac{T_t(a^*)}{t} d_{\text{U}}(\mathbb{P}_{t,a^*}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} \quad (194)$$

$$\leq \min_{i \neq a^*} \min_{x \in I_{t,i}} \left\{ (w_{a^*}(\boldsymbol{\nu}) + \epsilon') d_{\text{U}}(\mathbb{P}_{t,a^*}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} \quad (195)$$

$$\leq \min_{i \neq a^*} \min_{x \in I_{t,i}} \left\{ (w_{a^*}(\boldsymbol{\nu}) + \epsilon' d_{\text{U}}(\mathbb{P}_{t,a^*}, x)) + (w_i(\boldsymbol{\nu}) + \epsilon') d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} \quad (196)$$

$$= \Gamma(\boldsymbol{\nu}_t, \mathbf{w}(\boldsymbol{\nu})) + O(\epsilon') \quad (197)$$

$$\leq \Gamma(\boldsymbol{\nu}) + \underbrace{\epsilon''}_{\triangleq \epsilon} + O(\epsilon'), \quad (198)$$

where (194) follows from the fact that  $a_t^{\text{top}} = a^*$  for all  $t > N_1^{\epsilon'}$ , (195) and (196) follow from the fact that  $\Gamma_i(\boldsymbol{\nu}, \mathbf{w})$  is an increasing function in each coordinate  $w_i$ , keeping the other coordinates fixed [10, Lemma 2], and (198) follows from the fact that  $t > N_1^{\epsilon'}$ . Following similar steps as (194)-(198), we can show that

$$\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right) \geq \Gamma(\boldsymbol{\nu}) - \epsilon. \quad (199)$$

The proof is completed by setting  $N_\epsilon \triangleq N_1^{\epsilon'}$ . ■

Next, we investigate the relationship between the empirical problem complexity  $\Gamma_t\left(\frac{1}{t}\mathbf{T}_t\right)$  and the normalized test statistic  $\frac{1}{t}\Lambda_t(a_t^{\text{top}}, a_t^{\text{ch}})$ . For this, we leverage the convergence of the test statistic to the log-likelihood ratio. For any arm  $i \in [K]$  and parameters  $\theta, \theta' \in \Theta$ , let us define

$$\text{nLLR}_t(i, \theta, \theta') \triangleq \frac{1}{T_t(i)} \sum_{s \in [t]: A_s = i} \log \frac{\pi_i(X_s | \theta)}{\pi_i(X_s | \theta')}. \quad (200)$$

Furthermore, for any  $\epsilon \in \mathbb{R}_+$  let us define the time instant

$$N_{\text{KL}}^\epsilon(i, \theta, \theta') \triangleq \sup \{t \in \mathbb{N} : |\text{nLLR}_t(i, \theta, \theta') - d_i(\theta \|\theta')| > \epsilon\}. \quad (201)$$

Together with Assumption 8, and following the same line of arguments as Theorem 5, we have

$$\mathbb{E}_\nu [N_{\text{KL}}^\epsilon(i, \theta, \theta')] < +\infty. \quad (202)$$

Let us define

$$M_{\text{KL}}^\epsilon \triangleq \max_{i \in [K]} N_{\text{KL}}^\epsilon(i, \mu_t(i), \tilde{\mu}_t(i)), \quad (203)$$

$$\bar{M}_{\text{KL}}^\epsilon \triangleq \max_{i \in [K]} N_{\text{KL}}^\epsilon(i, \bar{\mu}_t(i), \tilde{\mu}_t(i)), \quad (204)$$

$$\bar{N}_{\text{KL}}^\epsilon \triangleq \max\{M_{\text{KL}}^\epsilon, \bar{M}_{\text{KL}}^\epsilon, N_\nu^{\Delta_{\min}/4}\}. \quad (205)$$

For all  $t > \bar{N}_{\text{KL}}^\epsilon$ , we have

$$\frac{1}{t} \Lambda_t(a_t^{\text{top}}, a_t^{\text{ch}}) = \min_{i \neq a^*} \frac{1}{t} \Lambda_t(a^*, i) \quad (206)$$

$$= \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} d_{a^*}(\mu_t(a^*) \|\tilde{\mu}_t(i)) + \frac{T_t(i)}{t} d_i(\mu_t(i), \tilde{\mu}_t(i)) \right\} \quad (207)$$

$$\geq \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} (\text{nLLR}_t(a^*, \mu_t(a^*), \tilde{\mu}_t(i)) - \epsilon) + \frac{T_t(i)}{t} (\text{nLLR}_t(i, \mu_t(i), \tilde{\mu}_t(i)) - \epsilon) \right\} \quad (208)$$

$$\geq \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} \text{nLLR}_t(a^*, \mu_t(a^*), \tilde{\mu}_t(i)) + \frac{T_t(i)}{t} \text{nLLR}_t(i, \mu_t(i), \tilde{\mu}_t(i)) \right\} - \epsilon \quad (209)$$

$$\geq \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} \text{nLLR}_t(a^*, \bar{\mu}_t(a^*), \tilde{\mu}_t(i)) + \frac{T_t(i)}{t} \text{nLLR}_t(i, \bar{\mu}_t(i), \tilde{\mu}_t(i)) \right\} - \epsilon \quad (210)$$

$$\geq \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} (d_{a^*}(\bar{\mu}_t(a^*) \|\tilde{\mu}_t(i)) - \epsilon) + \frac{T_t(i)}{t} (d_i(\bar{\mu}_t(i) \|\tilde{\mu}_t(i)) - \epsilon) \right\} - \epsilon \quad (211)$$

$$\geq \min_{i \neq a^*} \left\{ \frac{T_t(a^*)}{t} d_{a^*}(\bar{\mu}_t(a^*) \|\tilde{\mu}_t(i)) + \frac{T_t(i)}{t} d_i(\bar{\mu}_t(i) \|\tilde{\mu}_t(i)) \right\} - 2\epsilon \quad (212)$$

$$\geq \min_{i \neq a^*} \min_{x \in \mathbb{R}} \left\{ \frac{T_t(a^*)}{t} d_{a^*}(\bar{\mu}_t(a^*) \|x) + \frac{T_t(i)}{t} d_i(\bar{\mu}_t(i) \|x) \right\} - 2\epsilon \quad (213)$$

$$\geq \min_{i \neq a^*} \min_{x \in I_{t,i}} \left\{ \frac{T_t(a^*)}{t} d_{\text{U}}(\mathbb{P}_{t,a^*}, x) + \frac{T_t(i)}{t} d_{\text{L}}(\mathbb{P}_{t,i}, x) \right\} - 2\epsilon \quad (214)$$

$$= \Gamma_t \left( \frac{1}{t} \mathbf{T}_t \right) - 2\epsilon. \quad (215)$$

Next, we proceed with the proof of Theorem 4. Expanding the time instance just before stopping, we have

$$\tau - 1 = (\tau - 1) \mathbb{1}_{\{\tau - 1 \leq N_2^\epsilon\}} + (\tau - 1) \mathbb{1}_{\{\tau - 1 > N_2^\epsilon\}} \quad (216)$$

$$\leq N_2^\epsilon + (\tau - 1) \mathbb{1}_{\{\tau - 1 > N_2^\epsilon\}}, \quad (217)$$

where we have defined  $N_2^\epsilon \triangleq \max\{N_\epsilon, \bar{N}_{\text{KL}}^\epsilon, N_2\}$ , and  $N_2$  will be specified later. Leveraging Lemma 13, along with the fact that at  $\tau - 1$ ,  $\Lambda_{\tau-1}(a_{\tau-1}^{\text{top}}, a_{\tau-1}^{\text{ch}}) \leq \beta_{\tau-1}(\delta)$ , if  $\tau - 1 > N_2^\epsilon$  we have

$$\Gamma(\nu) - 3\epsilon \leq \Gamma_{\tau-1} \left( \frac{1}{\tau-1} \mathbf{T}_{\tau-1} \right) - 2\epsilon \stackrel{(215)}{\leq} \frac{\Lambda_{\tau-1}(a_{\tau-1}^{\text{top}}, a_{\tau-1}^{\text{ch}})}{\tau-1} \leq \frac{\beta_{\tau-1}(\delta)}{\tau-1}. \quad (218)$$

Furthermore, leveraging the stopping threshold  $\beta_t(\delta)$  defined in Theorem 1, if  $\tau - 1 > N_2^\epsilon$  we have

$$\begin{aligned}
(\tau - 1)(\Gamma(\boldsymbol{\nu}) - 3\epsilon) &\leq \max_{i \in [K]} \log \mathcal{I}_i(\mu_{\tau-1}(i)) - 2 \cdot \min_{i \in [K]} W_{\tau-1}(i) \\
&\quad + (\tau - 1) \cdot \max_{i \in [K]} \{ \max\{d_i(\mu_{\tau-1}(i) \parallel \mu_{\tau-1}(i) - \epsilon), d_i(\mu_{\tau-1}(i) \parallel \mu_{\tau-1}(i) + \epsilon)\} \} \\
&\quad + 2 \log \frac{|\Theta|}{\sqrt{2\pi}} + \log \frac{\tau - 1}{2} + \log \frac{K - 1}{\delta}. \tag{219}
\end{aligned}$$

Let us define

$$\mathcal{I}_{\max} \triangleq \max_{i \in [K]} \max_{\theta \in \Theta} \mathcal{I}_i(\theta). \tag{220}$$

Furthermore, for any  $t \in \mathbb{N}$  and for all  $i \in [K]$  we have

$$W_t(\epsilon, i) = \int_{\Omega^{\otimes T_t(i)}} \log \left( 1 - 2Q \left( \epsilon \sqrt{\bar{V}_t(i)} \right) \right) \prod_{s \in [t]: A_s = i} \pi_i(X_s \mid \mu_t(i)) d\mathcal{X}_t^i \tag{221}$$

$$\log \left( 1 - 2Q \left( \epsilon \sqrt{T_t(i)} \mathcal{I}_i(\mu_t(i)) \right) \right) \geq \log \left( 1 - 2Q \left( \epsilon \sigma \sqrt{T_t(i)} \right) \right) \tag{222}$$

$$\geq \log \left( 1 - 2Q \left( \epsilon \sigma \sqrt{\sqrt{t/K} - 1} \right) \right), \tag{223}$$

where (222) is obtained using Assumption 7, and (223) is follows from the explicit exploration property of the TCB and ITCB algorithms. Hence, defining

$$N_2 \triangleq K \left( \left( \frac{1}{\epsilon \sigma} Q^{-1} \left( \frac{1}{4} \right) \right)^2 + 1 \right)^2, \tag{224}$$

for all  $t > N_2$  and for all  $i \in [K]$  we have

$$-\log \left( 1 - 2Q \left( \epsilon \sigma \sqrt{T_t(i)} \right) \right) \leq \log 2. \tag{225}$$

Next, note that as a result of Assumption 5, we have

$$\max_{i \in [K]} \{ \max\{d_i(\mu_{\tau-1}(i) \parallel \mu_{\tau-1}(i) - \epsilon), d_i(\mu_{\tau-1}(i) \parallel \mu_{\tau-1}(i) + \epsilon)\} \} \leq \zeta(\epsilon), \tag{226}$$

where  $\zeta(\epsilon) \in \mathbb{R}_+$  can be made arbitrarily small by appropriately choosing  $\epsilon$ . Hence, using (220), (225) and (226), for  $\tau_1 > N_2^\epsilon$  we can rearrange (219) as follows.

$$(\tau - 1)(\Gamma(\boldsymbol{\nu}) - 3\epsilon - \zeta(\epsilon)) \leq \log \mathcal{I}_{\max} + 2 \log \sqrt{\frac{2}{\pi}} |\Theta| + \log \frac{\tau - 1}{2} + \log \frac{K - 1}{\delta}. \tag{227}$$

To upper bound (227), we leverage [4, Lemma 18] which gives

$$\tau \stackrel{(217)}{\leq} N_2^\epsilon + \frac{1}{(\Gamma(\boldsymbol{\nu}) - \epsilon - \zeta(\epsilon))} \cdot \left( \log \frac{(K - 1) \mathcal{I}_{\max} |\Theta|^2 e}{(\Gamma(\boldsymbol{\nu}) - \epsilon - \zeta(\epsilon)) \delta} + \log \log \frac{(K - 1) \mathcal{I}_{\max} |\Theta|^2}{(\Gamma(\boldsymbol{\nu}) - \epsilon - \zeta(\epsilon))} \right) + 1. \tag{228}$$

Next, taking expectation on both sides, dividing by  $\log(1/\delta)$  and taking the limit of  $\delta \rightarrow 0$ , we have

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}_{\boldsymbol{\nu}}[\tau]}{\log(1/\delta)} \leq \frac{1}{(\Gamma(\boldsymbol{\nu}) - 3\epsilon - \zeta(\epsilon))}. \tag{229}$$

Taking infimum with respect to  $\epsilon$  in (229), we have

$$\lim_{\delta \rightarrow 0} \frac{\mathbb{E}_{\nu}[\tau]}{\log(1/\delta)} \leq \frac{1}{\left(\Gamma(\nu) - \zeta(\epsilon)\right)} \quad (230)$$

$$= \frac{1 + \alpha}{\Gamma(\nu)}, \quad (231)$$

where we have defined  $\alpha \triangleq \zeta(\epsilon)/(\Gamma(\nu) - \zeta(\epsilon))$ , and  $\alpha$  can be made arbitrarily small by choosing a sufficiently small  $\epsilon$ . This completes the proof.

## References

- [1] S. Bubeck, R. Munos, and G. Stoltz, “Pure exploration in multi-armed bandits problems,” in *Proc. International Conference on Algorithmic Learning Theory*, Porto, Portugal, October 2009.
- [2] V. Gabillon, M. Ghavamzadeh, and A. Lazaric, “Best arm identification: A unified approach to fixed budget and fixed confidence,” in *Proc. Advances in Neural Information Processing Systems*, Lake Tahoe, NV, December 2012.
- [3] S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone, “PAC subset selection in stochastic multi-armed bandits,” in *Proc. International Conference on Machine Learning*, Madison, WI, June 2012.
- [4] A. Garivier and E. Kaufmann, “Optimal best arm identification with fixed confidence,” in *Proc. Conference on Learning Theory*, New York, NY, June 2016.
- [5] L. Xu, J. Honda, and M. Sugiyama, “A fully adaptive algorithm for pure exploration in linear bandits,” in *Proc. International Conference on Artificial Intelligence and Statistics*, Lanzarote, Canary Islands, April 2018.
- [6] K. Jamieson, M. Malloy, R. Nowak, and S. Bubeck, “‘lil’ UCB : An optimal exploration algorithm for multi-armed bandits,” in *Proc. Conference on Learning Theory*, Barcelona, Spain, June 2014.
- [7] A. Mukherjee and A. Tajer, “SPRT-based efficient best arm identification in stochastic bandits,” *IEEE Journal on Selected Areas in Information Theory (accepted for publication)*, June 2023.
- [8] M. Hoffman, B. Shahriari, and N. Freitas, “On correlation and budget constraints in model-based bandit optimization with application to automatic machine learning,” in *Proc. International Conference on Artificial Intelligence and Statistics*, Reykjavik, Iceland, April 2014.
- [9] J. Katz-Samuels, L. Jain, Z. Karnin, and K. G. Jamieson, “An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits,” in *Proc. Advances in Neural Information Processing Systems*, Virtual, December 2020.
- [10] D. Russo, “Simple bayesian algorithms for best-arm identification,” *Operations Research*, vol. 68, no. 6, pp. 1625–1647, April 2020.
- [11] X. Shang, R. de Heide, P. Menard, E. Kaufmann, and M. Valko, “Fixed-confidence guarantees for Bayesian best-arm identification,” in *Proc. International Conference on Artificial Intelligence and Statistics*, Sicily, Italy, August 2020.

- [12] M. Jourdan, R. Degenne, D. Baudry, R. de Heide, and E. Kaufmann, “Top two algorithms revisited,” in *Proc. Advances in Neural Information Processing Systems*, New Orleans, LA, December 2022.
- [13] S. Agrawal, S. Juneja, and P. Glynn, “Optimal  $\delta$ -correct best-arm selection for heavy-tailed distributions,” in *Proc. International Conference on Algorithmic Learning Theory*, San Diego, CA, February 2020.
- [14] Y. Jedra and A. Proutiere, “Optimal best-arm identification in linear bandits,” in *Proc. Advances in Neural Information Processing Systems*, Virtual, December 2020.
- [15] P. Ménard, “Gradient ascent for active exploration in bandit problems,” *arXiv 1905.08165*, 2019.
- [16] R. Degenne, W. M. Koolen, and P. Ménard, “Non-asymptotic pure exploration by solving games,” in *Proc. Advances in Neural Information Processing Systems*, Vancouver, Canada, December 2019.
- [17] R. Degenne, P. Menard, X. Shang, and M. Valko, “Gamification of pure exploration for linear bandits,” in *Proceedings of the 37th International Conference on Machine Learning*, Vienna, Austria, July 2020.
- [18] P.-A. Wang, R.-C. Tzeng, and A. Proutiere, “Fast pure exploration via Frank-Wolfe,” in *Proc. Advances in Neural Information Processing Systems*, virtual, December 2021.
- [19] A. Mukherjee and A. Tajer, “SPRT-based best arm identification in stochastic bandits,” in *Proc. International Symposium on Information Theory*, Espoo, Finland, June 2022.
- [20] C. Qin, D. Klabjan, and D. Russo, “Improving the expected improvement algorithm,” in *Proc. Advances in Neural Information Processing Systems*, Long Beach, CA, December 2017.
- [21] E. Kaufmann and W. M. Koolen, “Mixture martingales revisited with applications to sequential tests and confidence intervals,” *Journal of Machine Learning Research*, vol. 22, no. 246, pp. 1–44, 2021.
- [22] T. Lattimore and C. Szepesvári, *Bandit Algorithms*, Cambridge University Press, Cambridge, UK, 2020.
- [23] R. K. Sundaram, *A First Course in Optimization Theory*, Cambridge University Press, Cambridge, UK, June 1996.
- [24] J. Lin, “Divergence measures based on the Shannon entropy,” *IEEE Transactions on Information Theory*, vol. 37, no. 1, pp. 145–151, January 1991.
- [25] E. Posner, “Random coding strategies for minimum entropy,” *IEEE Transactions on Information Theory*, vol. 21, no. 4, pp. 388–391, July 1975.
- [26] C. D. Aliprantis and K. C. Border, *Infinite Dimensional Analysis: A Hitchhiker’s Guide*, Springer-Verlag, Berlin, Germany, 2006.
- [27] Y. S. Chow and H. Teicher, *Probability Theory Independence, Interchangeability, Martingales*, Springer, 1978.
- [28] S. P. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, Cambridge, UK, 2004.