
Dynamic Learning in Large Matching Markets

Anand Kalvit¹ and Assaf Zeevi²
Columbia University, New York
{¹akalvit22,²assaf}@gsb.columbia.edu

Abstract

We study a sequential matching problem faced by *large* centralized platforms where “jobs” must be matched to “workers” subject to uncertainty about worker skill proficiencies. Jobs arrive at discrete times with “job-types” observable upon arrival. To capture the “choice overload” phenomenon, we posit an unlimited supply of workers where each worker is characterized by a vector of attributes (aka “worker-types”) drawn from an underlying population-level distribution. The distribution as well as mean payoffs for possible worker-job type-pairs are unobservables and the platform’s goal is to sequentially match incoming jobs to workers in a way that maximizes its cumulative payoffs over the planning horizon. We establish lower bounds on the *regret* of any matching algorithm in this setting and propose a novel rate-optimal learning algorithm that adapts to aforementioned primitives *online*. Our learning guarantees highlight a distinctive characteristic of the problem: achievable performance only has a *second-order* dependence on worker-type distributions; we believe this finding may be of interest more broadly.

1 Introduction

Background and motivation. The problem of sequentially matching “jobs” to “workers” under uncertainty forms the bedrock of many modern operational settings, especially in the online gig economy, see, e.g., applications such as Amazon Mechanical Turk, TaskRabbit, Jobble, and the likes. A simpler instance of the problem dates back to [13] where it is referred to as the sequential stochastic assignment problem (SSAP). A fundamental issue in such settings is that the platform typically is oblivious (at least initially) to the skill proficiencies of individual workers for specific task categories. This complexity is further compounded by the large number of workers usually present on such platforms, tantamount to prohibitively large experimentation costs associated with acquisition of granular information at the level of an individual worker. This issue is commonly mitigated by exploiting structure in the problem (if any), or by positing distributional assumptions on the population of available workers, e.g., workers may be drawn from some distribution \mathcal{D} satisfying certain context-specific desiderata. Such distributional assumptions are vital to designing efficient algorithms for these systems, and as such, traditional literature has largely relied on the availability of ex ante knowledge of \mathcal{D} or certain key aspects thereof (refer to the literature review in §1.2).

Key research question. An important characteristic of the gig economy is that the population of workers may undergo distributional shifts over the course of the platform’s planning horizon. These effects may, many a time, fail to register in a timely manner; as a result, there may be delays in tailoring appropriately the matching algorithm (calibrated typically using available distribution-level information) to the changed environment. This has the potential to cause revenue losses as well as catalyze endogenous worker attrition. Such exigencies necessitate designing algorithms that are *agnostic* to \mathcal{D} and whose performance is *robust* to plausible realizations thereof.

The model at a glance. We consider a finite set of possible job-types (denoted by \mathcal{J}), an assumption we deem appropriate for settings such as those discussed above. In addition, we model workers as exhibiting discrete skill-levels (aka worker-types), indexed by $\{1, \dots, K_j\}$, w.r.t. each job-type $j \in \mathcal{J}$,

and make the simplifying assumption that $(K_j : j \in \mathcal{J})$ is known a priori. It is not unreasonable to make this assumption since it is common, in practice, for platforms to deploy pilot experiments prior to the actual matching phase in order to gather sufficient information on key primitives such as the size and stability of low-dimensional sub-population clusters, if any exist; one can therefore safely assume in settings where such structure exists that $(K_j : j \in \mathcal{J})$ is well-estimated a priori. Furthermore, assumptions pertaining to the finiteness of the set of possible job-types and segmentation of the population of workers based on discretization of skill levels are de rigueur also in the dynamic matching and the broader operations research literature, see, e.g., [7, 2, 17], etc.

While the demand is constituted by sequential job-arrivals (possibly in batches of stochastic size and composition), we posit availability of an unlimited number of workers on the supply side. The latter feature encapsulates the *choice overload* phenomenon characteristic of many large market settings where workers are available in a large number relative to the platform’s planning horizon. To our best knowledge, extant literature on matching under uncertainty is largely limited to “finite” markets (see the literature review in §1.2), and therefore fails to accommodate this important practical consideration. In our setting, the population of workers, albeit large, is governed by a fixed distribution that controls the proportion of each worker-type. Specifically, the K_j distinct worker-types w.r.t. job-type j are distributed according to $\alpha_j := (\alpha_{i,j} : i = 1, \dots, K_j)$, where $\sum_{i=1}^{K_j} \alpha_{i,j} = 1$. We note that this is *one* possible model of a matching market that is closer in spirit to SSAP [13] as well as other related formulations thereof; it differs from other models in the matching literature (see §1.2) in that it tries to capture a salient aspect of large markets, viz., choice overload, as opposed to traditional aspects such as *competition* and *congestion* best elucidated via conventional “finite” market models.

The platform’s goal is to maximize its expected cumulative payoffs over a sequence of n rounds of matching, subject to worker-types w.r.t. job-types and their distributions $\{\alpha_j : j \in \mathcal{J}\}$, as well as mean payoffs for possible worker-job type-pairs being latent attributes. As is the norm in settings with incomplete information and imperfect learning, we reformulate this objective as minimizing the *expected cumulative regret* relative to an oracle that is privy to aforementioned primitives.

On the complexity of the problem. Even with a unique job-type, say $\mathcal{J} = \{j_0\}$, and only one job arriving per period, the ensuing *allocation* problem is challenging to analyze on account of the distribution α_{j_0} and any statistical properties of the rewards being unknown. In the simplest possible formulation, $K_{j_0} = 2$, and the statistical complexity of the corresponding regret minimization problem is governed by three principal primitives: (i) the sub-optimality gap $\underline{\Delta}_{j_0} > 0$ between the mean rewards of the optimal and inferior worker sub-populations; (ii) the probability α_{1,j_0} of sampling an optimal worker from the population; and (iii) the planning horizon n . One may aptly recognize this as an infinitely many-armed bandit problem (where arms are synonymous to workers) with an *arm-reservoir* distribution $(\alpha_{1,j_0}, 1 - \alpha_{1,j_0})$ and a mean reward gap of $\underline{\Delta}_{j_0}$. However, this model differs from the classical literature on infinite-armed bandits in that its arm-reservoir distribution is not endowed with any regularity properties (see the literature review in §1.2), instead we only posit a finite support with cardinality known to the decision maker (in this case, a cardinality of two), absent however, knowledge of the associated probability masses (in this case, $\alpha_{j_0,1}$ and $1 - \alpha_{j_0,1}$). In our setting, absence of information on α_{1,j_0} significantly exacerbates the difficulty of analysis as calibrating *exploration* becomes challenging (on account of a “large” number of arms). In particular, how many arms must one query from the arm-reservoir in order to have at least one optimal arm in the queried set with high probability, is difficult to answer if (a lower bound on) the proportion α_{1,j_0} of optimal arms is unknown. Consequently, any finite consideration set may only contain inferior arms and as a result, any algorithm limited to such a selection will suffer a *linear regret*. One may contrast this setting with its classical two-armed counterpart with gap $\underline{\Delta}_{j_0}$, where finiteness of the set of arms (binary action space) makes it possible to design efficient rate-optimal policies. In our setting, on the other hand, it remains a priori unclear if there even exists a policy capable of achieving *sub-linear regret*. The general matching problem naturally is only harder.

1.1 Contributions

In this work, we resolve several foundational questions pertaining to complexity and achievable performance in the matching problem described earlier, and provide a comprehensive understanding of various other aspects thereof. Among other things, we propose an algorithm that achieves a finite-time instance-dependent expected regret of $\mathcal{O}(\log n)$ after n rounds (the big-Oh subsumes problem-dependent scaling factors encapsulating its fundamental complexity), and prove that this

performance cannot be improved w.r.t. n . While the order of regret and complexity of the problem suggests a great degree of similarity to the classical stochastic finite-armed bandit problem, properties of the performance bounds and salient aspects of algorithm design are quite distinct from the latter, as are the key primitives that determine complexity along with the analysis tools needed to study them. In what follows, we will for expositional reasons assume $\mathcal{J} = \{j_0\}$ and a batch size of 1 (jobs arrive one at a time) whenever $|\mathcal{J}| = 1$. Our theoretical contributions can then be projected along the following axes:

Complexity of regret when $|\mathcal{J}| = 1$. We establish information-theoretic lower bounds on regret that are order-wise tight (in the horizon n) in the instance-dependent setting (Theorem 1). In addition, we establish a *uniform* lower bound on achievable performance (tight in n) that captures explicitly the scaling behavior w.r.t. the fraction α_{1,j_0} of optimal arms (Theorem 2); this is shown via a novel non-information-theoretic proof based entirely on convex analysis.

Algorithm design and achievable performance. We propose a policy (Algorithm 2) that is rate-optimal (in n) for the instance-dependent setting. Our policy relies only on knowledge of K_{j_0} , is agnostic to information pertaining to the reward distributions as well as the distribution $(\alpha_{i,j_0} : i = 1, \dots, K_{j_0})$ of worker-types. Furthermore, the upper bound depends on the distribution of worker-types only in sub-logarithmic terms (see below).

Performance bounds for general \mathcal{J} . Aforementioned results for $|\mathcal{J}| = 1$ and unit batch size are then translated to the general (matching) version of the problem described earlier, where \mathcal{J} can be any arbitrary finite set and jobs may arrive in batches of stochastic size and constitution. In the matching problem, we establish that regret is bounded above by $\sum_{j \in \mathcal{J}} (C_1(\boldsymbol{\mu}_j) \log n + C_2(\boldsymbol{\mu}_j, \boldsymbol{\alpha}_j) \log \log n)$ under our policy tailored to this setting, where the constants $C_1(\cdot), C_2(\cdot, \cdot)$ only depend on their arguments, $\boldsymbol{\mu}_j := (\mu_{i,j} : i = 1, \dots, K_j)^1$ and $\boldsymbol{\alpha}_j := (\alpha_{i,j} : i = 1, \dots, K_j)$ (Theorem 4). When $K_j = 2 \forall j \in \mathcal{J}$, we improve this guarantee to $\sum_{j \in \mathcal{J}} (C_1(\boldsymbol{\mu}_j) \log n + C_2(\boldsymbol{\mu}_j, \boldsymbol{\alpha}_j))$ (Theorem 6). It is noteworthy that the upper bounds depend on $\{\boldsymbol{\alpha}_j : j \in \mathcal{J}\}$ only in sub-logarithmic terms. We believe this finding may be of interest more broadly.

1.2 Literature review

Our work is positioned relative to two major streams of literature; dynamic matching and multi-armed bandits. Below, we briefly survey each of these areas and remark on the distinctions and novelties in our model vis-à-vis the extant body of work.

Dynamic matching under uncertainty. There is a recent line of work on simultaneous learning and matching in bipartite graphs under uncertainty. For a survey of works in this area, see [22, 23, 24, 25, 16, 9, 3], etc. Aforementioned references, by and large, consider an archetypal stable matching problem under uncertainty in preferences where a heterogeneous collection of jobs (represented by nodes on one side of a bipartite graph) must be matched to workers (the other side of the graph) with unknown or noisy preferences over jobs. The matching proceeds iteratively in rounds in a way that meets certain stability criteria at all times as well as ensures that the true preferences are “learnt” at a regret-optimal rate. Cited works, however, differ from our paper fundamentally in that their learning problems are posited over a *finite* set of workers, which allows for sufficient exploration of each; this would be infeasible in our setting owing to a “large” population thereof.

In contrast, [17] considers a stationary setting where a stream of heterogeneous jobs must dynamically be matched to a policy-dependent steady state population of workers in a way that respects capacity constraints on the supply and demand processes. This paper shares basic similarities with our work in studying a “large” population model of workers. Their key technical innovation, however, lies in the way polytope capacity constraints are handled via an intelligent use of shadow prices to create essentially an unconstrained learning problem that may be solved rate-optimally using conventional heuristics. It is noteworthy that the algorithm proposed in aforementioned reference requires *ex ante* knowledge of $\{\boldsymbol{\mu}_j, \boldsymbol{\alpha}_j : j \in \mathcal{J}\}$ in addition to other primitives. Our model, on the other hand, has a richer learning component that is challenging to address as it is, absence of capacity constraints notwithstanding. Our primary contribution here lies in establishing fundamental achievability results for this setting and in the design of novel rate-optimal algorithms that adapt to key primitives online.

¹ $\mu_{i,j}$ denotes the mean reward per match between a worker of type i (w.r.t. job-type j) and a type j job.

Multi-armed bandits. Our problem shares structural similarities with *infinitely many-armed bandits*, modulo heterogeneity and multiplicity of pulls. Infinite-armed bandits involve a fixed *reservoir distribution* over an *uncountable* set of arm-types (possible mean rewards) that may be queried arbitrarily often over the horizon of play. These problems trace their roots to [4] which studied the Bernoulli reward setting under a Uniform (on $[0, 1]$) prior on the mean. Subsequent works have considered more general reservoir distributions, albeit endowed with certain *regularity* properties, see, e.g., [26, 6, 8, 10], etc. In terms of the statistical complexity of regret minimization, such regularity assumptions are tantamount to the minimal achievable regret being polynomial in the horizon (see above references). In contrast, our model is fundamentally simpler owing to a finite set of arm-types despite also being endowed with infinitely many arms. However, unlike cited works, the decision maker in our setting is completely oblivious to the reservoir distribution (or any property thereof) which substantially exacerbates the difficulty of analysis as well as the hardness of the problem itself.

1.3 Organization of the paper

§2 provides a formal description of the problem and §3 contains results pertaining to lower bounds on achievable performance for natural policy classes. §4 deals with design and analysis considerations for rate-optimal policies and also contains our main propositions together with supporting theoretical guarantees. All other discussion (including auxiliary results and proofs) is deferred to the appendix.

2 Problem formulation

Job-arrival process. The platform faces an arrival stream of jobs (i.i.d. in time) given by $\{(\Lambda_{j,t} : j \in \mathcal{J}) : t = 1, 2, \dots\}$, where \mathcal{J} is finite and $\Lambda_{j,t}$ is the number of type j jobs arriving at time t . Types and multiplicities of jobs are perfectly observable upon arrival. We assume that there exists some finite constant $M > 0$ satisfying $\mathbb{P}(\max_{j \in \mathcal{J}} \sup_{t \geq 1} \Lambda_{j,t} \leq M) = 1$. We remark that our algorithms do not require knowledge of M ; the assumption only serves to simplify analysis.²

Supply of workers. We assume that workers are distributed on the unit interval $[0, 1]$ according to some probability distribution \mathcal{D} that is absolutely continuous w.r.t. the Lebesgue measure on $[0, 1]$. Associated with each job-type $j \in \mathcal{J}$, there exists a permutation $\sigma_j := \{\sigma_j(i) : i = 1, \dots, K_j\}$ of $\{1, \dots, K_j\}$, and a sequence of thresholds $0 =: \lambda_{0,j} < \lambda_{1,j} < \dots < \lambda_{K_j-1,j} < \lambda_{K_j,j} := 1$ partitioning the unit interval into K_j disjoint sub-intervals. We posit a payoff model whereby a worker $x \in (\lambda_{i-1,j}, \lambda_{i,j})$ (for some $i \in \{1, \dots, K_j\}$) generates a stochastic reward with mean $\mu_{\sigma_j(i),j}$ upon match with a type j job; it is assumed that the K_j mean rewards adhere to the strict order $\mu_{1,j} > \dots > \mu_{K_j,j}$. We define $\alpha_{i,j} := \mathbb{P}(X \in (\lambda_{\iota(i,j)-1,j}, \lambda_{\iota(i,j),j}))$, where $X \sim \mathcal{D}$ and $\iota(i, j) \in \{1, \dots, K_j\}$ is the unique element satisfying $\sigma_j(\iota(i, j)) = i$, as the probability that a worker sampled at random from \mathcal{D} (equivalently, from the *population*), is i^{th} best for job-type j (generates mean reward $\mu_{i,j}$); such a worker is said to have type i w.r.t. job-type j . Thus, a type 1 worker w.r.t. job-type j is *optimal* for jobs of type j . Note that the model allows for staggered optimality of worker-types associated with different job-types, as Figure 1 below illustrates.

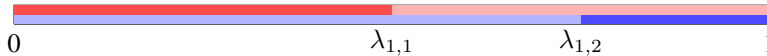


Figure 1: **Possible distribution of worker-types for $\mathcal{J} = \{1, 2\}$ and $K_1 = K_2 = 2$.** The darker shades represent optimal (type 1) workers while the lighter shades represent inferior (type 2) ones w.r.t. each job-type. In this example, no worker can simultaneously be optimal for both job-types.

High-level description of the matching problem. Each arriving job may be matched one-to-one to a worker from the available supply. Each match takes at most one period for execution, a matched worker thus frees up before the next lot of jobs arrives. Matched jobs leave the system upon completion and the platform receives a stochastic reward for each completed job; a job that remains unmatched drops out instantaneously. The platform has information neither on individual

²Though our regret bounds will scale linearly with M as we shall later see, proofs do not as such require batch-sizes to be almost surely bounded and can be refined to guarantee $\mathcal{O}(\log n)$ bounds also for $\Lambda_{j,t}$'s supported on \mathbb{Z}_+ under appropriate tail behavior. We do not pursue this line of analysis here purely for expositional reasons.

worker-types w.r.t. job-types nor on their supply distribution, however, it has perfect knowledge of $(K_j : j \in \mathcal{J})$. Subject to this premise, the platform must match incoming jobs to workers in a way that maximizes its expected cumulative payoffs over a sequence of n rounds of matches.

Adaptive control. For any job that arrives at time t , the platform can match it to: (i) a worker that has matched before, (ii) a *new* worker (one without any history of matches) sampled from the population, or (iii) no worker (job is dropped). A policy $\pi := (\pi_1(\cdot, \cdot), \pi_2(\cdot, \cdot), \dots)$ is an *adaptive* rule that prescribes the allocation $\pi_t(\cdot, \cdot)$ at time t . Specifically, $\pi_t(j, k)$ denotes the label of the worker (not its *type*) that gets matched to the k^{th} job of type j arriving at time t (provided there are at least k job-arrivals of type j at t and the k^{th} job is not dropped). Upon match, a $[0, 1]$ -valued stochastic reward with mean $\mu_{\kappa_j(\pi_t(j,k)),j}$ is realized, where $\kappa_j(\pi_t(j,k)) \in \{1, \dots, K_j\}$ encodes the type of worker $\pi_t(j, k)$ w.r.t. job-type j . The realized rewards are independent across matches and in time.

Platform’s objective. The goal of maximizing the expected cumulative payoffs over n rounds is converted to minimizing the expected *regret* relative to a clairvoyant policy that prescribes an “optimal” match for each arriving job. We are thus interested in the following optimization problem

$$\inf_{\pi \in \Pi} \mathbb{E} R_n^\pi := \inf_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=1}^n \sum_{j \in \mathcal{J}: \Lambda_{j,t} \geq 1} \sum_{k=1}^{\Lambda_{j,t}} (\mu_{1,j} - \mu_{\kappa_j(\pi_t(j,k)),j}) \right]. \quad (1)$$

Here, Π is the class of *non-anticipating* policies, i.e., $\pi_{t+1}(\cdot, \cdot)$ is *adapted* to \mathcal{F}_t for all $t \in \{0, 1, \dots\}$, where $\mathcal{F}_t := \sigma\{(\Lambda_s, \pi_s, \mathbf{r}_s) : s = 1, \dots, t\}$ denotes the natural filtration at time t . Here, $\Lambda_s := (\Lambda_{j,s} : j \in \mathcal{J})$, π_s is the set of matches implemented at time s and \mathbf{r}_s is the set of collected rewards. The expectation in (1) is w.r.t. the randomness in job-arrivals, worker supply, policy, and rewards.

Going forward, we will adopt standard terminology from the multi-armed bandit literature and refer to workers as “arms” and jobs as “pulls” interchangeably.

3 Lower bounds for natural policy classes

For a succinct illustration of the statistical complexity of our problem setting, it is conducive to pivot to the paradigmatic case where the arrival stream comprises only of a single job-type (say j) arriving one at a time, and the population of workers is partitioned into $K_j = 2$ clusters with $\alpha_{1,j} \leq 1/2$ (recall that this is the proportion of the optimal worker-type). In this case, one anticipates the problem to be at least as hard as the classical two-armed bandit with a mean reward gap of $\underline{\Delta}_j := \mu_{1,j} - \mu_{2,j} > 0$; this is on account of the infinitely many alternatives available to the decision maker in our setting as opposed to just two. Indeed, we establish this in Theorem 1 via an information-theoretic approach that delicately handles the combinatorial complexity arising due to probabilities over countably many arms (proof is provided in Appendix C). In what follows, an instance of the problem refers to a collection of reward distributions with means $(\mu_{1,j}, \mu_{2,j})$ (and gap $\underline{\Delta}_j$). We will overload the notation for expected regret slightly to emphasize its instance-dependence.

Theorem 1 (Information-theoretic lower bounds on achievable performance) *Suppose that the arrival process comprises only of a single job-type, say $j \in \mathcal{J}$, arriving one at a time, i.e., $\Lambda_{j,t} = 1$ for $t = 1, 2, \dots$. Also suppose that $K_j = 2$ and $\alpha_{1,j} \leq 1/2 - \epsilon$, where $\epsilon \in (0, 1/2)$ is arbitrary. Let Π_{adm} denote the class of admissible³ policies. Then, the following is true under any $\pi \in \Pi_{\text{adm}}$:*

- (i) *For any $\underline{\Delta}_j > 0$, there exists a problem instance ν such that $\mathbb{E} R_n^\pi(\nu) \geq C \log n / \underline{\Delta}_j$ for n large enough⁴, where C is some absolute constant.*
- (ii) *For any $n \in \mathbb{N}$, there exists a problem instance ν such that $\mathbb{E} R_n^\pi(\nu) \geq \epsilon C' \sqrt{n}$, where C' is some absolute constant.*

Distinction from classical MAB. Although the above result bears resemblance to classical information-theoretic lower bounds for finite-armed bandits, it is noteworthy that our setting has a fundamentally greater problem complexity that requires a more nuanced analysis vis-à-vis the finite-armed setting. Traditional proofs, as a result, cannot be translated to our setting in a straightforward manner. To see this, note that when $\alpha_{1,j}$ is high, a new worker is very likely to be optimal;

³This is a rich policy class containing all “reasonable” algorithms; refer to Definition (1) in Appendix C.

⁴The dependence on ϵ is subsumed under “ n large enough.”

in the limit as $\alpha_{1,j} \rightarrow 1$, the problem becomes degenerate as all policies incur zero expected regret. Thus, the problem cannot be harder than the two-armed bandit with gap $\underline{\Delta}_j$ uniformly over all values of $\alpha_{1,j}$. While we conjecture $\alpha_{1,j} < 1$ to be a sufficient condition for the existence of $\Omega(\log n / \underline{\Delta}_j)$ instance-dependent and $\Omega(\sqrt{n})$ instance-independent (minimax) lower bounds (modulo $\alpha_{1,j}$ -dependent scaling factors), there are technical challenges due to probabilistic associations over countably many arms. The restriction to $\alpha_{1,j} < 1/2$ and *admissible* policies is then necessary for tractability of the proof and it remains unclear if this can be generalized further.

Capturing dependence on α_j . Although the instance-dependent lower bound in Theorem 1 is tight in n as we shall later see, it fails to provide actionable insights vis-à-vis α_j . A natural question in our setting is whether and in what manner does the presence of countably many arms (as opposed to finitely many) affect achievable regret. In particular, how does the difficulty associated with finding from the available supply an optimal worker for a given job type (and the dependence on the distribution α_j) come into play. Below, we propose a lower bound that explicitly captures this dependence, albeit with respect to a somewhat restricted policy class.

Theorem 2 (α_j -dependent lower bound) *Suppose that the arrival process comprises only of a single job-type, say $j \in \mathcal{J}$, arriving one at a time, i.e., $\Lambda_{j,t} = 1$ for $t = 1, 2, \dots$. Also suppose that $\alpha_{1,j} \leq 1/2$. Denote by Π_m the class of “memoryless” policies under which the decision to match an incoming job to a new worker at any time $t \in \{1, 2, \dots\}$ is independent of \mathcal{F}_{t-1} . Then, for all problem instances ν with a minimal sub-optimality gap of at least $\underline{\Delta}_j > 0$, one has*

$$\liminf_{n \rightarrow \infty} \inf_{\pi \in \Pi_m} \frac{\mathbb{E}R_n^\pi(\nu)}{\log n} \geq \frac{\underline{\Delta}_j}{4\alpha_{1,j}}.$$

Discussion. The proof is located in Appendix D. Several comments are in order. **(i)** The class Π_m , in particular, includes policies that front-load exploration, i.e., sample upfront a pre-specified number of workers from the population and then deploy a regret minimizing algorithm of choice on the set of workers so obtained. This includes several natural approaches to the problem as we shall later see. **(ii)** The foremost noticeable aspect of Theorem 2 that differs from Theorem 1 is that it provides a uniform lower bound *over all instances* that are at least $\underline{\Delta}_j$ -separated in the mean reward, as opposed to merely establishing their existence. **(iii)** The presence of $\underline{\Delta}_j$ in the numerator in Theorem 2 unlike traditional information-theoretic bounds where $\underline{\Delta}_j$ resides in the denominator suggests that while this bound may be vacuous if $\underline{\Delta}_j$ is “small,” it certainly becomes most relevant when $\underline{\Delta}_j$ is “well-separated.” In that sense, Theorem 1 and 2 provide a tool to separate regimes of $\underline{\Delta}_j$ where one bound captures the dominant effects vis-à-vis the other. **(iv)** A novelty of Theorem 2 lies in its proof, which differs from classical lower bound proofs in that it is based entirely on ideas from convex analysis as opposed to the information-theoretic and change-of-measure techniques hitherto used in the literature.

Remarks. **(i)** It is not impossible to avoid $1/\alpha_{1,j}$ -scaling in the instance-dependent logarithmic regret. We will later show via an upper bound for our algorithm CAB-K that the $\alpha_{1,j}$ -dependence can, in fact, be relegated to sub-logarithmic terms (CAB-K samples new workers from the population *adaptively* based on the sample-history of onboarded workers and therefore does not belong to Π_m). Importantly, this will establish a somewhat surprising fact that the instance-dependent logarithmic bound in Theorem 1 is optimal w.r.t. to its dependence on $\alpha_{1,j}$ (the scaling w.r.t. $\underline{\Delta}_j$, however, may not be best possible as forthcoming upper bounds will suggest). **(ii)** Theorem 2 holds also for any worker supply where the optimal mean reward w.r.t. job-type j is at least $\underline{\Delta}_j$ -separated from the rest, the nature of the set of possible worker-types (countable or uncountable) notwithstanding.

4 Designing adaptive policies for matching

The approach we adopt in this paper directly addresses the fact that there is an unlimited supply of available workers at all times. A natural design then is to tailor sub-routines specific to job-types in \mathcal{J} and instantiate them at the first arrival of each type. Specifically, if jobs of type j arrive at $\{t_1, t_2, \dots\}$, then the platform should call the sub-routine specific to job-type j only at aforementioned times, independent of other job-arrivals. This leads to the meta-algorithm MATCH (see Algorithm 1) for the matching problem. In what follows, ALG refers to an arm-allocation rule w.r.t. a fixed job-type that prescribes *one* arm upon each invocation. ALG can be thought of as a horizon-free sampling strategy for a countably many-armed bandit problem with one pull per period. When multiple jobs (say L) of the same type (say j) arrive at the same time, we instantiate (if necessary) new *parallel threads* of

ALG specific to job-type j so there exist a total of L type j threads when the next lot of jobs arrives. In the following, L_j denotes the running count of the number of parallel threads of ALG for type j jobs.

Algorithm 1 MATCH

```

1: Input: (i)  $\mathcal{J}$ , (ii)  $(K_j : j \in \mathcal{J})$ , and (iii) ALG.
2: Initialization: Set  $L_j = 0$  for each  $j \in \mathcal{J}$ .
3: for  $t \in \{1, 2, \dots\}$  do
4:   for  $j \in \mathcal{J}$  do
5:     if  $\Lambda_{j,t} \geq 1$  then
6:       if  $\Lambda_{j,t} \leq L_j$  then
7:         Match the  $\Lambda_{j,t}$  type  $j$  jobs to the first  $\Lambda_{j,t}$  threads of ALG for type  $j$  jobs.
8:       else
9:         Match the first  $L_j$  type  $j$  jobs to the  $L_j$  available threads of ALG for type  $j$  jobs.
10:      Instantiate  $\Lambda_{j,t} - L_j$  new threads of ALG for the remaining  $\Lambda_{j,t} - L_j$  jobs.
11:      Update  $L_j \leftarrow \Lambda_{j,t}$ .

```

Discussion of MATCH. An immediate observation from Algorithm 1 is that ALG must be *anytime*, i.e., it should not depend on the horizon of play since the number of job arrivals (over the platform’s planning horizon) of each type is not known a priori. Keeping this objective in mind, we shift our focus to designing an arm-allocation rule ALG w.r.t. a fixed job-type, say type j , that: (i) prescribes one pull per period, (ii) depends only on K_j , (iii) is adaptive to the mean reward vector $\boldsymbol{\mu}_j$ and the supply distribution $\boldsymbol{\alpha}_j$, and (iv) is horizon-free. Once such an ALG is designed, its composition with MATCH will transfer learning guarantees to the original matching problem.

4.1 Shifting focus to adaptive sequential sampling strategies tailored to a specific job-type

Going forward, we will assume that jobs belong to a common fixed type and arrive one at a time. With slight abuse of notation, the supply of available workers is characterized by K worker-types with distinct means $\boldsymbol{\mu} := (\mu_i : i = 1, \dots, K)$ adhering to $\mu_1 > \dots > \mu_K$. The maximal and minimal sub-optimality gaps are given by $\bar{\Delta} := \mu_1 - \mu_K$ and $\underline{\Delta} := \mu_1 - \mu_2$ respectively, and the minimal reward gap is $\delta := \min_{1 \leq i < i' \leq K} (\mu_i - \mu_{i'})$. The distribution of worker-types is denoted by $\boldsymbol{\alpha} := (\alpha_i : i = 1, \dots, K)$, where α_i is the probability of sampling a type i arm (characterized by mean μ_i) from the population. The learning horizon is n . The decision maker only knows K and is oblivious to $(\boldsymbol{\mu}, \boldsymbol{\alpha}, n)$.

The specific setting described above was first studied in [18] for $K = 2$ for which a UCB-styled algorithm with $\mathcal{O}(\log n)$ regret was proposed. The analysis of said algorithm leveraged certain concentration properties of the UCB1 policy [1] that were recently discovered in the context of a two-armed bandit with *equal* arm-means (see Theorem 4(i) in [18]). Currently, there is no known algorithm for the general setting with $K > 2$ arm-types. We discuss in the appendix why properties of UCB1 critical to $\mathcal{O}(\log n)$ regret in [18] do not hold for $K > 2$; consequently, a natural adaptation of their algorithm to $K > 2$ will likely fail to generalize $\mathcal{O}(\log n)$ bounds.⁵ To our best knowledge, the general countable-armed bandit (CAB) setting with $K > 2$ types remained open in prior literature.

We close this gap in the literature by proposing a policy CAB-K (see Algorithm 2) that achieves $\mathcal{O}(\log n)$ regret in the general K -typed setting. CAB-K operates based on the Explore-then-Commit principle with adaptive stopping and re-initialization times. It is noteworthy that CAB-K is not horizon-free; in particular, knowledge of the horizon is critical for appropriate calibration of its stopping and re-initialization thresholds. However, this is not a constraining characteristic and we will use a doubling trick dubbed HF (Algorithm 3) to make it horizon-free with an anytime $\mathcal{O}(\log n)$ guarantee.

Discussion of CAB-K. At any time, the algorithm computes two thresholds of $\mathcal{O}(\sqrt{m \log m})$ and $\mathcal{O}(\sqrt{m \log n})$ for the $\binom{K}{2}$ pairwise-difference-of-reward processes, m being the per-arm sample count. If the envelope of said process is dominated by the former threshold, the concerned pair likely contains arms of the same type (equal means). The explanation stems from the Law of the Iterated Logarithm (see [14], Theorem 8.5.2): the envelope of a zero-drift length- m random walk grows as $\mathcal{O}(\sqrt{m \log \log m})$. In the aforementioned scenario, the algorithm discards the *entire* consideration

⁵We will, however, propose a version in the appendix that is asymptotically optimal (achieves $o(n)$ regret).

Algorithm 2 CAB-K

```

1: Input: Horizon of play  $n$ .
2: Set budget  $T = n$ .
3: Initialize new epoch: Query  $K$  new arms; call it consideration set  $\mathcal{A} = \{1, 2, \dots, K\}$ .
4: Play each arm in  $\mathcal{A}$  once; observe rewards  $\{X_{a,1} : a \in \mathcal{A}\}$ .
5: Update budget:  $T \leftarrow T - K$ .
6: Per-arm sample count  $m \leftarrow 1$ .
7: Generate  $\binom{K}{2}$  independent standard Gaussian random variables  $\{\mathcal{Z}_{a,b} : a, b \in \mathcal{A}, a < b\}$ .
8: while  $T \geq K$  do
9:   if  $\exists a, b \in \mathcal{A}, a < b$  s.t.  $|\mathcal{Z}_{a,b} + \sum_{k=1}^m (X_{a,k} - X_{b,k})| < 4\sqrt{m \log m}$  then
10:     Permanently discard  $\mathcal{A}$  and repeat from step (3).
11:   else
12:     if  $|\sum_{k=1}^m (X_{a,k} - X_{b,k})| \geq 4\sqrt{m \log n} \forall a, b \in \mathcal{A}, a < b$  then
13:       Permanently commit to arm  $a^* \in \arg \max_{a \in \mathcal{A}} \{\sum_{k=1}^m X_{a,k}\}$ .
14:     else
15:       Play each arm in  $\mathcal{A}$  once; observe rewards  $(X_{1,m+1}, \dots, X_{K,m+1})$ .
16:        $m \leftarrow m + 1$ .
17:        $T \leftarrow T - K$ .

```

set and ushers in a new epoch. This is done to avoid the possibility of incurring linear regret should an optimal arm be missing from the consideration set (e.g., when all arms are type 2). In the other scenario that all pairwise-difference-of-reward processes dominate $\mathcal{O}(\sqrt{m \log n})$, the consideration set is likely to contain arms of distinct types (no two have equal means) and the algorithm simply commits to the empirically best arm. Lastly, if neither threshold is crossed (signifying insufficient learning), the sample count for each arm is advanced by one, and the entire process repeats.

Reason for introducing the Gaussian corruption. Centered Gaussian noise is added to all pairwise-difference-of-reward processes in step (9) of CAB-K to avoid the possibility of incurring linear regret should the support of the reward distributions be a “very small” subset of $[0, 1]$. To illustrate this point, suppose that $K = 2$ and the rewards associated with the types are deterministic with a gap of $\underline{\Delta} < 2\sqrt{2 \log 2}$. Then, as soon as the algorithm queries a consideration set containing one arm each of the two types and the per-arm count reaches 2, the difference-of-reward statistic will satisfy $|\sum_{j=1}^2 (X_{1,j} - X_{2,j})| = 2\underline{\Delta} < 4\sqrt{2 \log 2}$ and the consideration set will be discarded. On the other hand, if both arms are of the same type (simultaneously optimal or inferior), the algorithm will still re-initialize as soon as the per-arm count reaches 2.⁶ This will force the algorithm to keep querying new arms from the reservoir at rate that is linear in time, which is tantamount to incurring linear regret in the horizon. The addition of centered Gaussian noise hedges against this risk by guaranteeing that the difference-of-reward process essentially has an infinite support at all times even when the reward distributions might be degenerate. This rids the regret of its fragility w.r.t. the support of reward distributions. The next proposition crystallizes this discussion.

Proposition 1 (Persistence of heterogeneous consideration sets) *Let $\{X_{a,k} : k = 1, 2, \dots\}$ be a collection of independent samples from an arm of type $a \in \{1, 2, \dots, K\} =: \mathcal{A}$, and $\{\mathcal{Z}_{a,b} : a, b \in \mathcal{A}, a < b\}$ be a collection of $\binom{K}{2}$ independently generated standard Gaussians. Then,*

$$\mathbb{P} \left(\bigcap_{m \geq 1} \bigcap_{a, b \in \mathcal{A}, a < b} \left\{ \left| \mathcal{Z}_{a,b} + \sum_{k=1}^m (X_{a,k} - X_{b,k}) \right| \geq 4\sqrt{m \log m} \right\} \right) > \frac{\bar{\Phi}(f(T_0))}{2} =: \beta_{\delta, K} > 0, \quad (2)$$

where $\bar{\Phi}(\cdot)$ is the tail of the standard Gaussian CDF, and $T_0 := \max(\lceil (64/\delta^2) \log^2(\frac{64}{\delta^2}) \rceil, \mathfrak{C}_K)$ with $\mathfrak{C}_K := \inf \left\{ p \in \mathbb{N} : \sum_{m=p}^{\infty} \frac{1}{m^8} \leq \frac{1}{2K^2} \right\}$. Lastly, $f(x) := x + 4\sqrt{x \log x}$ for all $x \geq 1$.

The proof is provided in Appendix E; this meta-result is key to the upper bound stated in Theorem 3.

⁶ $|\sum_{j=1}^m (X_{1,j} - X_{2,j})| = 0$ identically in this case for any $m \in \mathbb{N}$ while $4\sqrt{m \log m} > 0$ only for $m \geq 2$.

Interpretation of $\beta_{\delta,K}$. First of all, note that $\beta_{\delta,K}$ admits a closed-form characterization in terms of standard functions and satisfies $\beta_{\delta,K} > 0$ for $\delta > 0$ with $\lim_{\delta \rightarrow 0} \beta_{\delta,K} = 0$. Secondly, $\beta_{\delta,K}$ depends exclusively on δ and K , and represents a lower bound on the probability that CAB-K will never discard a consideration set containing arms of distinct types.

Theorem 3 (Upper bound on the regret of CAB-K) *For any horizon of play $n \geq 1$, the expected regret of the policy π given by CAB-K is bounded as*

$$\mathbb{E}R_n^\pi \leq \frac{CK^3\bar{\Delta}}{\beta_{\delta,K}} \left(\frac{\log n}{\delta^2} + \frac{1}{K! \prod_{i=1}^K \alpha_i} \right),$$

where $\beta_{\delta,K}$ is as defined in (2) and C is some absolute constant.

Discussion. The dependence on the *minimal reward gap* δ is not an artifact of our analysis but, in fact, reflective of the operating principle of the algorithm. CAB-K keeps querying new consideration sets of size K until it determines with high enough confidence that no two arms have the same type (equal means); this is the genesis of δ in the upper bound. Importantly, equipped just with knowledge of K , it remains unclear if there exists an alternative strategy that does not rely on assessing pairwise differences among the queried arms. Another prominent feature of the upper bound is its $\mathcal{O}(1)$ -dependence on α . This essentially means that the difficulty associated with sampling an optimal arm from an infinite reservoir containing finitely many arm-types is of a *second order*, which starkly contrasts the findings in [12]. Cited paper assumes ex ante knowledge of a lower bound on α_1 (proportion of optimal arms) and posits no structure on sub-optimal arm-types. Under this premise, a *first-order* difficulty of sampling optimal arms from the reservoir is established (to wit, the logarithmic term depends on α_1). However, whether this would hold also for the subset of problems where the reservoir only contains a finite universe of arm-types (with known cardinality) was left open. Theorem 3 essentially settles this problem. The proof of Theorem 3 is provided in Appendix F.

Remarks. (i) Possible improvements. CAB-K, in its present form, indulges in wasteful exploration by discarding entire consideration sets upon re-initialization. It is possible to be parsimonious in this regard and we leave the pursuit of such algorithms to future work. **(ii) More on $\beta_{\delta,K}$.** Notice that when $K = 2$ and $\alpha_1 \leq 1/2$, the upper bound in Theorem 3 assumes the form $C\beta_{\delta,2}^{-1}(\log n/\delta + \delta/\alpha_1)$, where C is some absolute constant. Thus, $\beta_{\delta,2}^{-1}$ captures the relative increase in problem complexity (vis-à-vis the paradigmatic two-armed bandit with gap δ) attributable to an unlimited supply of arms of the two types. To what extent may this factor be shaved off remains an interesting open problem. **(iii) Comparison with lower bounds.** One should also contrast Theorem 3 with the lower bound in Theorem 2; by allowing for policies that query the arm-reservoir *adaptively*, we could achieve a regret performance robust to α (second-order dependence). This also leads to the somewhat remarkable conclusion that the lower bound in Theorem 1 is optimal w.r.t. its dependence on α . **(iv) Anytime guarantees.** A horizon-free version of CAB-K may be obtained by passing it to the HF operator given in Algorithm 3; a logarithmic bound for the resulting composition HF(CAB-K) is stated in Theorem 8.

4.2 Transferring learning guarantees to the matching problem

Theorem 4 (Achievable performance under MATCH \circ HF(CAB-K)) *Denote by π the composition of MATCH with HF(CAB-K). Then, its expected regret after any number $n \geq 1$ of rounds satisfies*

$$\mathbb{E}R_n^\pi \leq CM \sum_{j \in \mathcal{J}} \left[\frac{K_j^3 \bar{\Delta}_j}{\beta_{\delta_j, K_j}} \left(\frac{\log n}{\delta_j^2} + \frac{\log \log(n+2)}{K_j! \prod_{i=1}^{K_j} \alpha_{i,j}} \right) \right],$$

where β_{δ_j, K_j} is as defined in (2) with $\delta \leftarrow \delta_j := \min_{1 \leq i < i' \leq K_j} (\mu_{i,j} - \mu_{i',j})$ and $K \leftarrow K_j$, $\bar{\Delta}_j := \mu_{1,j} - \mu_{K_j,j}$, and C is some absolute constant.

Discussion. The foremost noticeable aspect of Theorem 4 is that achievable regret depends on $\{\alpha_j : j \in \mathcal{J}\}$ (collection of worker-type distributions w.r.t. job-types), surprisingly, only through sub-logarithmic terms. Moreover, when $K_j = 2 \forall j \in \mathcal{J}$, we improve this to an $\mathcal{O}(1)$ -dependence (see Theorem 6 in the appendix). We conjecture that the $\mathcal{O}(\log \log n)$ factor in the upper bound can, in fact, be shaved off also for $K_j > 2$; pursuits in this direction are left to future work. Among other things, characterizing the *minimax* complexity of this setting remains a challenging open problem in light of the multiplicative factors that appear in Theorem 4 (fundamentally different from finite-armed problems). Numerical experiments showing $\mathcal{O}(\log n)$ achievable regret are provided in the appendix.

Acknowledgments and Disclosure of Funding

The authors thank the anonymous reviewers for their constructive feedback on the initial version of this paper. The authors declare an absence of any competing interests, financial or otherwise.

References

- [1] AUER, P., CESA-BIANCHI, N., AND FISCHER, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47, 2-3 (2002), 235–256.
- [2] BANERJEE, S., GOLLAPUDI, S., KOLLIAS, K., AND MUNAGALA, K. Segmenting two-sided markets. In *Proceedings of the 26th International Conference on World Wide Web* (2017), pp. 63–72.
- [3] BASU, S., SANKARARAMAN, K. A., AND SANKARARAMAN, A. Beyond $\log^2(t)$ regret for decentralized bandits in matching markets. In *International Conference on Machine Learning* (2021), PMLR, pp. 705–715.
- [4] BERRY, D. A., CHEN, R. W., ZAME, A., HEATH, D. C., SHEPP, L. A., ET AL. Bandit problems with infinitely many arms. *The Annals of Statistics* 25, 5 (1997), 2103–2116.
- [5] BESSON, L., AND KAUFMANN, E. What doubling tricks can and can’t do for multi-armed bandits. *arXiv preprint arXiv:1803.06971* (2018).
- [6] BONALD, T., AND PROUTIERE, A. Two-target algorithms for infinite-armed bandits with bernoulli rewards. In *Advances in Neural Information Processing Systems* (2013), pp. 2184–2192.
- [7] BUI, L., JOHARI, R., AND MANNOR, S. Clustered bandits. *arXiv preprint arXiv:1206.4169* (2012).
- [8] CARPENTIER, A., AND VALKO, M. Simple regret for infinitely many armed bandits. In *International Conference on Machine Learning* (2015), pp. 1133–1141.
- [9] CEN, S. H., AND SHAH, D. Regret, stability, and fairness in matching markets with bandit learners. *arXiv preprint arXiv:2102.06246* (2021).
- [10] CHAN, H. P., AND HU, S. Infinite arms bandit: Optimality via confidence bounds. *arXiv preprint arXiv:1805.11793* (2018).
- [11] CHAUDHURI, A. R., AND KALYANAKRISHNAN, S. Quantile-regret minimisation in infinitely many-armed bandits. In *UAI* (2018), pp. 425–434.
- [12] DE HEIDE, R., CHESHIRE, J., MÉNARD, P., AND CARPENTIER, A. Bandits with many optimal arms. In *Advances in Neural Information Processing Systems* (2021), vol. 34, pp. 22457–22469.
- [13] DERMAN, C., LIEBERMAN, G. J., AND ROSS, S. M. A sequential stochastic assignment problem. *Management Science* 18, 7 (1972), 349–355.
- [14] DURRETT, R. *Probability: theory and examples*, vol. 49. Cambridge university press, 2019.
- [15] HOEFFDING, W. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association* 58, 301 (1963), 13–30.
- [16] JAGADEESAN, M., WEI, A., WANG, Y., JORDAN, M., AND STEINHARDT, J. Learning equilibria in matching markets from bandit feedback. *Advances in Neural Information Processing Systems* 34 (2021).
- [17] JOHARI, R., KAMBLE, V., AND KANORIA, Y. Matching while learning. *Operations Research* 69, 2 (2021), 655–681.
- [18] KALVIT, A., AND ZEEVI, A. From finite to countable-armed bandits. In *Advances in Neural Information Processing Systems* (2020), vol. 33, pp. 8259–8269.
- [19] KALVIT, A., AND ZEEVI, A. A closer look at the worst-case behavior of multi-armed bandit algorithms. In *Advances in Neural Information Processing Systems* (2021), vol. 34, pp. 8807–8819.
- [20] LAI, T. L., AND ROBBINS, H. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6, 1 (1985), 4–22.
- [21] LATTIMORE, T., AND SZEPESVÁRI, C. *Bandit algorithms*. Cambridge University Press, 2020.

- [22] LIU, L. T., MANIA, H., AND JORDAN, M. Competing bandits in matching markets. In *International Conference on Artificial Intelligence and Statistics* (2020), PMLR, pp. 1618–1628.
- [23] LIU, L. T., RUAN, F., MANIA, H., AND JORDAN, M. I. Bandit learning in decentralized matching markets. *Journal of Machine Learning Research* 22, 211 (2021), 1–34.
- [24] SANKARARAMAN, A., BASU, S., AND SANKARARAMAN, K. A. Dominate or delete: Decentralized competing bandits with uniform valuation. *arXiv preprint arXiv:2006.15166* (2020).
- [25] SANKARARAMAN, A., BASU, S., AND SANKARARAMAN, K. A. Dominate or delete: Decentralized competing bandits in serial dictatorship. In *International Conference on Artificial Intelligence and Statistics* (2021), PMLR, pp. 1252–1260.
- [26] WANG, Y., AUDIBERT, J.-Y., AND MUNOS, R. Algorithms for infinitely many-armed bandits. In *Advances in Neural Information Processing Systems* (2009), pp. 1729–1736.
- [27] ZHU, Y., AND NOWAK, R. On regret with multiple best arms. *Advances in Neural Information Processing Systems* 33 (2020), 9050–9060.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes] Full proofs are provided in the appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Supplementary material: General organization

1. Appendix A contains additional discussion pertaining to the propositions in the main text.
2. Appendix B provides numerical experiments.
3. Appendix C provides the proof of Theorem 1.
4. Appendix D provides the proof of Theorem 2.
5. Appendix E provides the proof of Proposition 1.
6. Appendix F provides the proof of Theorem 3.
7. Appendix G discusses a doubling sub-routine to make algorithms horizon-free.
8. Appendix H provides the proof of Theorem 4.
9. Appendix I discusses a first-order optimal policy for countable-armed bandits with K types.
10. Appendix J discusses auxiliary results pertaining to the aforementioned policy.
11. Appendix K provides more explore-then-commit policies for the countable-armed bandit.

A Additional discussion

We next provide a regret upper bound for a countable-armed bandit policy CAB-K(UCB) that improves upon the UCB-styled algorithm of [18] by ridding its performance guarantee of a certain fragile assumption pertaining to the reward support. The following result, however, is specific to the case of $K = 2$ arm-types. The policy CAB-K(UCB) is stated as Algorithm 4 in Appendix I.

Theorem 5 (Upper bound on the regret of CAB-K(UCB) when $K = 2$) *The expected regret of the policy π given by CAB-2(UCB) after any number $n \geq 1$ of pulls is bounded as*

$$\mathbb{E}R_n^\pi \leq \min \left[\underline{\Delta}n, \frac{C}{\beta_{\underline{\Delta},2}} \left(\frac{\log n}{\underline{\Delta}} + \frac{\underline{\Delta}}{\alpha_1} \right) \right],$$

where $\beta_{\underline{\Delta},2}$ is as defined in (2) with $\delta \leftarrow \underline{\Delta}$ and $K \leftarrow 2$, and C is some absolute constant.

The proof is provided in §I.2.

Limitation of CAB-2(UCB). The performance stated in Theorem 5 together with its anytime property might appear to give an edge to CAB-2(UCB) over CAB-K. However, the former is theoretically disadvantaged in that its logarithmic upper bound is currently not amenable to generalization to $K > 2$. The issue traces its roots to the use of UCB1 as a subroutine; concentration behavior thereof leveraged towards the analysis of CAB-K(UCB) when $K = 2$ fails to hold when $K > 2$ rendering proofs intractable.⁷ This is illustrated via a simple example with $K = 3$ arm-types discussed below.

Technical issues with generalizing logarithmic bounds of CAB-2(UCB) to $K > 2$ types. In the $K = 2$ setting, there are only two possibilities for what a consideration set could be; arms can have means that are either (i) distinct, or (ii) equal. In the former case, an optimal arm is guaranteed to exist in the consideration set and UCB1 will spend the bulk of its sampling effort on it, which is good for regret performance. In the latter scenario, since arms have equal means, UCB1 will split samples approximately equally between the two with high probability (see Theorem 4(i) in [18]); subsequently the consideration set will be discarded within a finite number of samples in expectation (see steps (6) and (7) of CAB-K(UCB)). Contrast this with an alternative setting with $K = 3$ and mean rewards $\mu_1 > \mu_2 > \mu_3$. In this case, CAB-K(UCB) will query consideration sets of size 3. Thus, a query can potentially return one arm with mean μ_2 and two with mean μ_3 . Since an optimal arm (mean μ_1) is missing, the algorithm will incur linear regret on this set; it is therefore imperative to discard it at the earliest. Unfortunately though, UCB1 will invest an overwhelming majority of its sampling effort in the “locally optimal” arm (mean μ_2) and allocate logarithmically fewer samples among the other two. This logarithmic rate of sampling arms with mean μ_3 is proof-inhibiting (vis-à-vis the $K = 2$ case where the rate is linear as previously discussed), making it difficult to theoretically answer if CAB-K(UCB) might still be able to discard the arms within, say, logarithmically many pulls in the horizon. This is an open research question and at the moment, strong $\mathcal{O}(\log n)$

⁷Concentration behavior à la Theorem 4 of [18] is elucidated further in [19].

performance guarantees are available only for $K = 2$; we could only establish asymptotic-optimality ($o(n)$ regret) for $K > 2$ (see Theorem 9). Among other things that remain, identifying the optimal (instance-dependent) scaling factors w.r.t. (μ, α) and the optimal order of minimax regret would be challenging open directions.

More on the inverse scaling w.r.t. $\beta_{\underline{\Delta}, 2}$. The multiplicative factor $\beta_{\underline{\Delta}, 2} \leq 1$ captures the additional complexity of the problem due to the countable nature of arms. Ideally, one would want the sampling strategy to never discard a consideration set of distinct-typed arms (equivalently, the probability in (2) should always be 1). This, however, is not statistically achievable since the rewards are stochastic and types are unobservable. Some “false positives” will be unavoidable (heterogeneous consideration sets incorrectly labeled homogeneous and therefore discarded) causing the aforementioned probability to be bounded away from 1. However, one can show that this probability is bounded away also from 0 (Proposition 1 formalizes this statement), and $\beta_{\underline{\Delta}, 2}$ provides one such lower bound when $K = 2$. As to what the tightest possible lower bound is as a function of $\underline{\Delta}$ remains an open problem at the moment. This would also shed light on the additional complexity of the problem attributable to the countable nature of arms vis-à-vis the classical finite-armed problem. In terms of regret guarantees, positivity alone of $\beta_{\underline{\Delta}, 2}$ suffices to establish instance-dependent rate-optimality (in n) of CAB-K(UCB) when $K = 2$ (see Theorem 1), notwithstanding the structure of its dependence on $\underline{\Delta}$. The worst-case upper bound, however, is polynomially bounded away from the optimal \sqrt{n} minimax rate owing to the presence of $\beta_{\underline{\Delta}, 2}$ in the denominator.

Composing MATCH with CAB-K(UCB) when $K_j = 2 \forall j \in \mathcal{J}$. We are now ready to elevate the performance guarantee of CAB-2(UCB) to the broader setting of the matching problem where each job-type partitions workers into two sub-populations.

Theorem 6 (Achievable performance under MATCH \circ CAB-2(UCB)) *Denote by π the composition of MATCH with CAB-2(UCB). Then, its expected regret after any number $n \geq 1$ of rounds is bounded as*

$$\mathbb{E}R_n^\pi \leq \sum_{j \in \mathcal{J}} \left[\frac{CM}{\beta_{\underline{\Delta}_j, 2}} \left(\frac{\log n}{\underline{\Delta}_j} + \frac{\underline{\Delta}_j}{\alpha_{1,j}} \right) \right],$$

where C is some absolute constant, and $\beta_{\underline{\Delta}_j, 2}$ is as defined in (2) with $\delta \leftarrow \underline{\Delta}_j := \mu_{1,j} - \mu_{2,j}$ and $K_j \leftarrow 2$.

In conclusion. We attempted in this paper to address from first principles the problem of dynamic matching in large markets under type uncertainty on the supply side (worker-types) and size uncertainty on the demand size (batch-sizes). Among things that remain, a predominant open question is how can one incorporate capacity constraints on supply and demand processes into this model, e.g., à la [17]. This would lift our setting to model more realistic scenarios where the population of available workers at any time is “large but finite;” we leave investigations in this direction to future work. Another natural modeling extension is one where workers may interact with the platform *strategically* (this could, e.g., reflect via a time-variant and potentially policy-dependent supply distribution); such *endogeneity* will likely force the platform to align its matching policies also with worker *incentives* in order to plug their potential *attrition*. As to what the fundamental limits of learning and achievable regret might be in such a setting remains an interesting open problem to study. There are also a number of interesting questions more directly on the methodological front and related closely to the technical development in this paper; we leave their pursuit to future work.

B Numerical experiments

We will evaluate the empirical performance of CAB-K (see Algorithm 2) in the countable-armed bandit problem (with one pull per period) characterized by means $\mu_1 > \dots > \mu_K$ and a reservoir distribution $\alpha = (\alpha_i : i = 1, \dots, K)$. Recall that $\underline{\Delta} = \mu_1 - \mu_2$ and $\delta = \min_{1 \leq i < i' \leq K} (\mu_i - \mu_{i'})$. We will conduct experiments for $K = 2$ and $K = 3$.

Experiments. In what follows, the graphs show the performance of different algorithms simulated on synthetic data. The horizon is capped at $n = 10^5$ for $K = 2$ and at $n = 10^4$ for $K = 3$. Each regret plot is averaged over at least 100 independent experiments (sample-paths). The shaded regions indicate standard 95% confidence intervals. For horizon-dependent algorithms, regret is plotted for

discrete values of the horizon n indicated by “*” and interpolated; for anytime algorithms, regret accrued until each $t \in \{1, \dots, n\}$ is plotted.

Baseline policies. We will benchmark the performance of CAB-K against three policies: (i) Sampling-UCB [12], (ii) ETC- $\infty(2)$ [18], and (iii) ETC-RAW (see Algorithm 6 in Appendix K). The first of these, Sampling-UCB, is a UCB-styled policy based on front-loading exploration of *new* arms (Theorem 2 thus applies to this policy). It is, however, noteworthy that Sampling-UCB is predicated on ex ante knowledge of (a lower bound on) the probability α_1 of sampling an optimal arm from the reservoir; we reemphasize that this is not the setting of interest in our paper. Furthermore, its regret scales as $\tilde{O}(\log n / (\alpha_1 \underline{\Delta}))$ (up to poly-logarithmic factors in $1/\underline{\Delta}$), which is inferior in terms of its dependence on α_1 relative to CAB-K and CAB-2(UCB) (see Theorem 3 and 5 respectively). There exist other algorithms as well (see, e.g., [11, 27]) developed for formulations with a prohibitively large number of arms. However, these are either sensitive to certain parametric assumptions on the probability of sampling an optimal arm, or focus on a different notion of regret altogether; both directions remain outside the ambit of our setting.

The second policy ETC- $\infty(2)$ is a non-adaptive explore-then-commit-styled algorithm for reservoirs with $K = 2$ types; this policy requires ex ante knowledge of a lower bound on the difference between the two mean rewards. Although ETC- $\infty(2)$ was originally proposed only for $K = 2$, it is easily generalizable and we present in Algorithm 5 (see Appendix K) a version (ETC- $\infty(K)$) that is adapted to K types.

The last policy ETC-RAW is also based on the explore-then-commit principle and operates using a pre-specified exploration schedule as opposed to an adaptive one à la CAB-K. We do not provide any theoretical performance guarantees for this policy and resort directly to empirical evaluations.

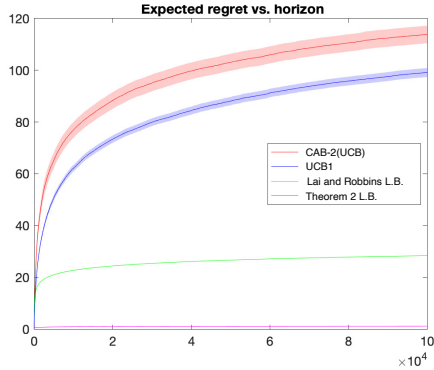


Figure 2: $K = 2$ and $\alpha = (1/2, 1/2)$: Achievable regret in the countable-armed problem vis-à-vis the paradigmatic two-armed bandit.

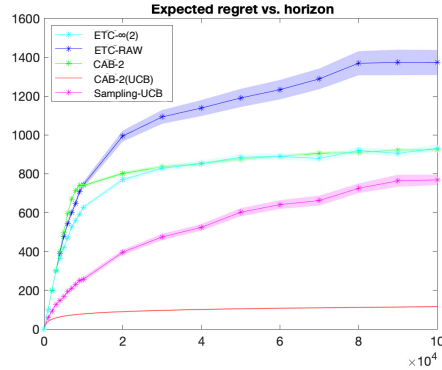


Figure 3: $K = 2$ and $\alpha = (1/2, 1/2)$: An instance with Bernoulli 0.6, 0.4 rewards.

Setup 1 [Figure 2, 3 and 4]. In this setting, we consider $K = 2$ with $\alpha_1 = 0.5$, i.e., two equiprobable arm-types, characterized by Bernoulli(0.6) and Bernoulli(0.4) rewards. Via this setup, we intend to illustrate the difference between the empirical performance achievable in the countable-armed setting vis-à-vis its traditional two-armed counterpart. Refer to Figure 2. The red curve indicates the empirical performance of CAB-2(UCB) in this setting. For reference, the blue one shows the empirical performance of UCB1 [1] in a two-armed bandit with Bernoulli(0.6) and Bernoulli(0.4) rewards; the green curve indicates the best achievable instance-dependent regret [20] in said two-armed configuration. As expected, the regret of CAB-2(UCB) is inflated relative to UCB1. This is owing to the $\beta_{\delta,2} \leq 1$ factor present in the denominator of CAB-2(UCB)’s upper bound; characterizing the sharpest lower bound on the probability in (2) (see Proposition 1) is challenging owing to the limited theoretical tools available to this end and we leave it as an open problem at the moment. Figure 3 shows the empirical performance of the algorithms discussed previously as well as Sampling-UCB initialized with $\alpha_1 = 1/2$ and ETC- $\infty(2)$ initialized with $\underline{\delta} = \delta/2 = 0.1$. Evidently, the (adaptive) explore-then-commit approach in CAB-K outperforms the pre-specified exploration schedule-based approach of ETC-RAW, and performs almost as good as the *gap-aware* approach in ETC- $\infty(2)$. While Sampling-UCB outmatches all explore-then-commit styled approaches, the best performing algorithm

is CAB-K(UCB). Surprisingly, this is despite the fact that the theoretical performance bounds for CAB-K and CAB-K(UCB) are identical (modulo numerical multiplicative constants) when $K = 2$ and $\alpha_1 \leq 0.5$ (see Theorem 3 and 5). A similar hierarchy in performances is also observable in Figure 4, which corresponds to a slightly “easier” instance with $\delta = 0.4$ (as opposed to 0.2) and equiprobable Bernoulli(0.9) and Bernoulli(0.5) rewards.

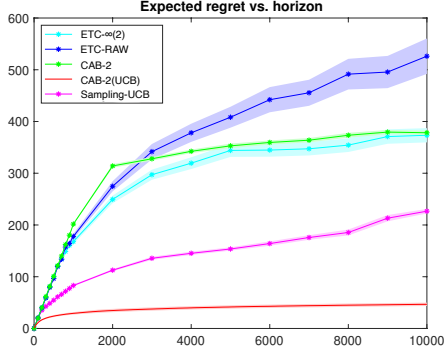


Figure 4: $K = 2$ and $\alpha = (1/2, 1/2)$: An instance with Bernoulli 0.9, 0.5 rewards.

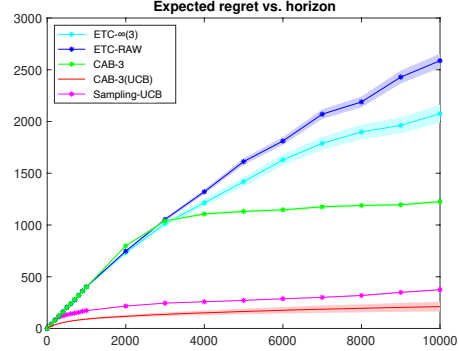


Figure 5: $K = 3$ and $\alpha = (1/3, 1/3, 1/3)$: An instance with Bernoulli 0.9, 0.5, 0.1 rewards.

Setup 2 [Figure 5]. Here, we consider a setting with $K = 3$ arm-types characterized by Bernoulli rewards with means 0.9, 0.5, 0.1, each occurring with probability $1/3$. We compare the performance of ETC-RAW, CAB-K and CAB-K(UCB) with ETC- $\infty(3)$ initialized with $\delta = \delta/2 = 0.2$, and Sampling-UCB initialized with $\alpha_1 = 1/3$. It is noteworthy that despite CAB-K(UCB)’s significantly superior empirical performance relative to aforementioned algorithms, only a weak $o(n)$ bound on its regret is currently available (see Theorem 9) due to reasons discussed earlier in the paper. Investigating best achievable rates under CAB-K(UCB) is an area of active research at the moment.

Remark. For general \mathcal{J} , our approach to the matching problem (see MATCH) involves instantiating horizon-free versions of aforementioned algorithms independently for each job-type upon its first arrival. We have not included in the current version of our paper numerical experiments pertaining to general \mathcal{J} since the complexity of our problem setting, by and large, is encapsulated by $|\mathcal{J}| = 1$ (owing to the design of MATCH), which is already discussed above. We do, however, plan to include comprehensive simulations for the general matching setting as well in future iterations of this paper.

C Proof of Theorem 1

Notation. Since the job-type is fixed (at j), we will with slight abuse of notation drop the subscript j from $(\mu_{1,j}, \mu_{2,j}, \underline{\Delta}_j, \alpha_{1,j}, \kappa_j(\cdot))$. Also, we will denote by π_t the index of the arm played by policy π at time t and by $\kappa(\pi_t) \in \{1, 2\}$ its type. Let $\mathcal{G}_1(x)$ and $\mathcal{G}_2(x)$ be two arbitrary collections of distributions with mean $x \in \mathbb{R}$. The tuple $(\mathcal{G}_1(x), \mathcal{G}_2(x))$ will be referred to as an *instance*.

Since the horizon of play is fixed at n , the decision maker may play at most n distinct arms. Therefore, it suffices to focus only on the sequence of the first n arms that may be played. A *realization* of an instance $\nu = (\mathcal{G}_1(\bar{\mu}_1), \mathcal{G}_2(\bar{\mu}_2))$ is defined as the n -tuple $r \equiv (r_i)_{1 \leq i \leq n}$, where $r_i \in \mathcal{G}_1(\bar{\mu}_1) \cup \mathcal{G}_2(\bar{\mu}_2)$ indicates the reward distribution of arm $i \in \{1, 2, \dots, n\}$. It must be noted that the decision maker need not play every arm in r . Let $i^* := \arg \max_{i \in \{1, 2\}} \bar{\mu}_i$. Then, the distribution over the possible realizations of $\nu = (\mathcal{G}_1(\bar{\mu}_1), \mathcal{G}_2(\bar{\mu}_2))$ in $\{r : r_i \in \mathcal{G}_1(\bar{\mu}_1) \cup \mathcal{G}_2(\bar{\mu}_2), 1 \leq i \leq n\}$ satisfies $\mathbb{P}(r_i \in \mathcal{G}_{i^*}(\bar{\mu}_{i^*})) = \alpha_1$ for all $i \in \{1, \dots, n\}$.

Recall that the cumulative pseudo-regret after n plays of a policy π on $\nu = (\mathcal{G}_1(\bar{\mu}_1), \mathcal{G}_2(\bar{\mu}_2))$ is given by $R_n^\pi(\nu) = \sum_{m=1}^n (\bar{\mu}_{i^*} - \bar{\mu}_{\kappa(\pi_t)})$, where $\kappa(\pi_t) \in \{1, 2\}$ indicates the type of the arm played by π at time t . Our goal is to lower bound $\mathbb{E}R_n^\pi(\nu)$, where the expectation is w.r.t. the randomness in π as well as the distribution over the possible realizations of ν . To this end, we define the notion of

expected cumulative regret of π on a realization r of $\nu = (\mathcal{G}_1(\bar{\mu}_1), \mathcal{G}_2(\bar{\mu}_2))$ by

$$S_n^\pi(\nu, r) := \mathbb{E}^\pi \left[\sum_{t=1}^n (\bar{\mu}_{i^*} - \bar{\mu}_{\kappa(\pi_t)}) \right],$$

where the expectation \mathbb{E}^π is w.r.t. the randomness in π . Note that $\mathbb{E}R_n^\pi(\nu) = \mathbb{E}^\nu S_n^\pi(\nu, r)$, where the expectation \mathbb{E}^ν is w.r.t. the distribution over the possible realizations of ν . We define our problem class \mathcal{N}_Δ as the collection of Δ -separated instances given by

$$\mathcal{N}_\Delta := \{(\mathcal{G}_1(\bar{\mu}_1), \mathcal{G}_2(\bar{\mu}_2)) : \bar{\mu}_1 - \bar{\mu}_2 = \Delta, (\bar{\mu}_1, \bar{\mu}_2) \in \mathbb{R}^2\}.$$

Definition 1 (Admissible policy) A policy π is deemed admissible for the problem class \mathcal{N}_Δ if for any instance $\nu \in \mathcal{N}_\Delta$ and any realization r thereof, it satisfies

$$\mathbb{E}^\nu [S_n^\pi(\nu, r) | \mathcal{L}(r) = m] \geq \mathbb{E}^\nu [S_n^\pi(\nu, r) | \mathcal{L}(r) = k] \quad \forall (m, n, k) : 0 \leq m \leq k \leq n, \quad (3)$$

where $\mathcal{L}(r)$ denotes the number of ‘‘optimal’’ arms in realization r , i.e., arms with mean $\bar{\mu}_{i^*}$.

The set of such policies is denoted by $\Pi_{\text{adm}}(\mathcal{N}_\Delta)$. We remark that the condition in (3) is not restrictive since it is only natural that any reasonable policy should incur a larger cumulative regret (in expectation) on realizations with fewer optimal arms.

Fix an arbitrary $\Delta > 0$ and consider an instance $\nu = (\{Q_1\}, \{Q_2\}) \in \mathcal{N}_\Delta$, where Q_1 and Q_2 are unit-variance Gaussian distributions with means μ_1 and μ_2 respectively. Consider an arbitrary realization $r \in \{Q_1, Q_2\}^n$ of ν and let $\mathcal{I} \subseteq \{1, 2, \dots, n\}$ denote the set of inferior arms in r (arms with reward distribution Q_2). Consider another instance $\nu' \in \mathcal{N}_\Delta$ given by $\nu' = (\{\tilde{Q}_1\}, \{Q_1\})$, where \tilde{Q}_1 is another unit variance Gaussian with mean $\mu_1 + \Delta$. Now consider a realization $r' \in \{\tilde{Q}_1, Q_1\}^n$ of ν' that is such that the arms at positions in \mathcal{I} have distribution \tilde{Q}_1 while those at positions in $\{1, 2, \dots, n\} \setminus \mathcal{I}$ have distribution Q_1 . Notice that \mathcal{I} is the set of optimal arms in r' (arms with reward distribution \tilde{Q}_1), implying $\mathcal{L}(r') = |\mathcal{I}|$. Then, the following always holds:

$$S_n^\pi(\nu, r) + S_n^\pi(\nu', r') \geq \left(\frac{\Delta n}{2}\right) \left(\mathbb{P}_{\nu, r}^\pi \left(\sum_{i \in \mathcal{I}} N_i(n) > \frac{n}{2} \right) + \mathbb{P}_{\nu', r'}^\pi \left(\sum_{i \in \mathcal{I}} N_i(n) \leq \frac{n}{2} \right) \right),$$

where $\mathbb{P}_{\nu, r}^\pi(\cdot)$ and $\mathbb{P}_{\nu', r'}^\pi(\cdot)$ denote the probability measures w.r.t. the instance-realization pairs (ν, r) and (ν', r') respectively, and $N_i(n)$ denotes the number of plays up to and including time n of arm $i \in \{1, 2, \dots, n\}$. Using the Bretagnolle-Huber inequality (Theorem 14.2 of [21]), we obtain

$$S_n^\pi(\nu, r) + S_n^\pi(\nu', r') \geq \left(\frac{\Delta n}{4}\right) \exp(-D_{\text{KL}}(\mathbb{P}_{\nu, r}^\pi, \mathbb{P}_{\nu', r'}^\pi)),$$

where $D_{\text{KL}}(\mathbb{P}_{\nu, r}^\pi, \mathbb{P}_{\nu', r'}^\pi)$ denotes the KL-Divergence between $\mathbb{P}_{\nu, r}^\pi$ and $\mathbb{P}_{\nu', r'}^\pi$. Using Divergence decomposition (Lemma 15.1 of [21]), we further obtain

$$S_n^\pi(\nu, r) + S_n^\pi(\nu', r') \geq \left(\frac{\Delta n}{4}\right) \exp\left(-\left(\frac{D_{\text{KL}}(Q_2, \tilde{Q}_1)}{\Delta}\right) S_n^\pi(\nu, r)\right) = \left(\frac{\Delta n}{4}\right) \exp(-2\Delta S_n^\pi(\nu, r)),$$

where the equality follows since \tilde{Q}_1 and Q_2 are unit variance Gaussian distributions with means separated by 2Δ . Next, taking the expectation \mathbb{E}^ν on both the sides above and a direct application of Jensen’s inequality thereafter yields

$$\mathbb{E}R_n^\pi(\nu) + \mathbb{E}^\nu S_n^\pi(\nu', r') \geq \left(\frac{\Delta n}{4}\right) \exp(-2\Delta \mathbb{E}R_n^\pi(\nu)). \quad (4)$$

Consider the $\mathbb{E}^\nu S_n^\pi(\nu', r')$ term in (4). Using a simple change-of-measure argument, we obtain

$$\begin{aligned} \mathbb{E}^\nu S_n^\pi(\nu', r') &= \mathbb{E}^{\nu'} \left[S_n^\pi(\nu', r') \left(\frac{1 - \alpha_1}{\alpha_1}\right)^{2(\mathcal{L}(r') - n/2)} \right] \\ &\leq \mathbb{E}R_n^\pi(\nu') + \mathbb{E}^{\nu'} \left[S_n^\pi(\nu', r') \left(\frac{1 - \alpha_1}{\alpha_1}\right)^{2(\mathcal{L}(r') - n/2)} \mathbb{1}\{\mathcal{L}(r') > n/2\} \right], \end{aligned} \quad (5)$$

where the inequality follows since $\alpha_1 \leq 1/2$. Now consider the second term on the RHS in (5). It follows that

$$\begin{aligned}
& \mathbb{E}^{\nu'} \left[S_n^\pi(\nu', r') \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{2(\mathcal{L}(r') - n/2)} \mathbb{1}_{\{\mathcal{L}(r') > n/2\}} \right] \\
&= \sum_{k > n/2} \mathbb{E}^{\nu'} \left[S_n^\pi(\nu', r') \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{2(\mathcal{L}(r') - n/2)} \mathbb{1}_{\{\mathcal{L}(r') = k\}} \right] \\
&= \sum_{k > n/2} \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{(2k - n)} \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') \mathbb{1}_{\{\mathcal{L}(r') = k\}}] \\
&= \sum_{k > n/2} \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{(2k - n)} \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = k] \mathbb{P}_{\nu'}(\mathcal{L}(r') = k) \\
&= \sum_{k > n/2} \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{(2k - n)} \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = k] \binom{n}{k} \alpha_1^k (1 - \alpha_1)^{(n - k)} \\
&= \alpha_1^n \sum_{k > n/2} \binom{n}{k} \left(\frac{1 - \alpha_1}{\alpha_1} \right)^k \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = k]. \tag{6}
\end{aligned}$$

Recall that $\nu' \in \mathcal{N}_\Delta$ and $\pi \in \Pi_{\text{adm}}(\mathcal{N}_\Delta)$. We have

$$\begin{aligned}
\mathbb{E} R_n^\pi(\nu') &= \mathbb{E}^{\nu'} S_n^\pi(\nu', r') \\
&\geq \sum_{m=1}^k \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = m] \mathbb{P}_{\nu'}(\mathcal{L}(r') = m) \quad (\text{for any } k \leq n) \\
&\geq \mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = k] \mathbb{P}_{\nu'}(\mathcal{L}(r') \leq k). \tag{using (3)} \tag{7}
\end{aligned}$$

Note that under the measure $\mathbb{P}_{\nu'}(\cdot)$, the distribution of $\mathcal{L}(r')$ is Binomial(n, α_1). Therefore, using Markov's inequality, $\mathbb{P}_{\nu'}(\mathcal{L}(r') > k) \leq n\alpha_1/k < 2\alpha_1$ for $k > n/2$ (Note that this bound is non-vacuous since the inequality is strict and $\alpha_1 < 1/2$ by assumption). One then has that $\mathbb{P}_{\nu'}(\mathcal{L}(r') \leq k) > 1 - 2\alpha_1$ for $k > n/2$. Using this observation in (7), we conclude for $k > n/2$ that

$$\mathbb{E}^{\nu'} [S_n^\pi(\nu', r') | \mathcal{L}(r') = k] < \frac{\mathbb{E} R_n^\pi(\nu')}{1 - 2\alpha_1}. \tag{8}$$

Combining (6) and (8), one obtains

$$\mathbb{E}^{\nu'} \left[S_n^\pi(\nu', r') \left(\frac{1 - \alpha_1}{\alpha_1} \right)^{2(\mathcal{L}(r') - n/2)} \mathbb{1}_{\{\mathcal{L}(r') > n/2\}} \right] < \frac{\mathbb{E} R_n^\pi(\nu')}{1 - 2\alpha_1}. \tag{9}$$

Using (5) and (9), one then concludes

$$\mathbb{E}^\nu S_n^\pi(\nu', r') \leq \mathbb{E} R_n^\pi(\nu') \left(1 + \frac{1}{1 - 2\alpha_1} \right) \leq \frac{2\mathbb{E} R_n^\pi(\nu')}{1 - 2\alpha_1}. \tag{10}$$

Finally, from (4) and (10), we have

$$\begin{aligned}
& \mathbb{E} R_n^\pi(\nu) + \left(\frac{2}{1 - 2\alpha_1} \right) \mathbb{E} R_n^\pi(\nu') \geq \frac{\Delta n}{4} \exp(-2\Delta \mathbb{E} R_n^\pi(\nu)) \\
\implies & \mathbb{E} R_n^\pi(\nu) + \mathbb{E} R_n^\pi(\nu') \geq \frac{(1 - 2\alpha_1)\Delta n}{8} \exp(-2\Delta \mathbb{E} R_n^\pi(\nu)) \\
\implies & \tilde{R}_n \geq \frac{(1 - 2\alpha_1)\Delta n}{16} \exp(-2\Delta \tilde{R}_n) \\
\implies & \tilde{R}_n \geq \frac{\epsilon \Delta n}{8} \exp(-2\Delta \tilde{R}_n), \tag{11}
\end{aligned}$$

where $\tilde{R}_n := \max(\mathbb{E}R_n^\pi(\nu), \mathbb{E}R_n^\pi(\nu'))$.

Instance-dependent lower bound

The assertion of the theorem follows from the fact that the inequality (11) is fulfilled only if for any $\varepsilon \in (0, 1)$, \tilde{R}_n satisfies for all n large enough $\tilde{R}_n \geq (1 - \varepsilon) \log n / (2\Delta)$.

Instance-independent (minimax) lower bound

Since $\tilde{R}_n \leq \Delta n$, it follows from (11) that

$$\tilde{R}_n \geq \frac{\varepsilon \Delta n}{16} \exp(-2\Delta^2 n).$$

Setting $\Delta = 1/\sqrt{n}$ now proves the stated assertion. \square

D Proof of Theorem 2

Notation. Again, since the job-type is fixed (at j), we will with slight abuse of notation drop the subscript (j) from $(K_j, \Delta_j, \alpha_{1,j})$. Also, note that the result is stated for the general setting with $K \geq 2$ arm-types where Δ denotes the minimal sub-optimality gap.

Consider now an arbitrary policy $\pi \in \Pi_m$. Denote by A_n^π the number of *distinct* arms played by π until time n . Consider an arbitrary $k \in \{1, 2, \dots, n\}$. Then, conditioned on $A_n^\pi = k$, the expected cumulative regret incurred by π is at least

$$\mathbb{E}[R_n^\pi | A_n^\pi = k] \geq (1 - \alpha_1)\Delta k + (1 - \alpha_1)^k \Delta(n - k) =: f(k). \quad (12)$$

Intuition behind (12). Each of the k arms played during the horizon has at least one pull associated with it. Consider a clairvoyant policy coupled to π that learns the best among the A_n^π arms played by π as soon as each has been pulled exactly once, i.e., after a total of A_n^π pulls. Clearly, the regret incurred by said clairvoyant policy lower bounds $\mathbb{E}R_n^\pi$. Further, since A_n^π is independent of the sample-history of arms, it follows that the A_n^π arms are statistically identical. Thus, conditioned on $A_n^\pi = k$, the expected regret from the first k pulls of the clairvoyant policy is at least $(1 - \alpha_1)\Delta k$. Also, the probability that each of the k arms is inferior-typed is $(1 - \alpha_1)^k$; the clairvoyant policy thus incurs a regret of at least $(1 - \alpha_1)^k \Delta(n - k)$ going forward. This explains the lower bound in (12). Therefore, for any $k \in \{1, 2, \dots, n\}$, we have

$$\mathbb{E}[R_n^\pi | A_n^\pi = k] \geq \min_{k \in \{1, 2, \dots, n\}} f(k) \geq \min_{x \in [0, n]} f(x).$$

We will show that $f(x)$ is strictly convex over $[0, n]$ with $f'(0) < 0$ and $f'(n) > 0$. Then, it would follow that $f(\cdot)$ admits a unique minimizer $x_n^* \in (0, n)$ given by the solution to $f'(x) = 0$. The minimum $f(x_n^*)$ will turn out to be logarithmic in n . Observe that

$$\begin{aligned} f'(x) &= (1 - \alpha_1)\Delta + (1 - \alpha_1)^x \Delta [(n - x) \log(1 - \alpha_1) - 1], \\ f''(x) &= -(1 - \alpha_1)^x \Delta [2 - (n - x) \log(1 - \alpha_1)] \log(1 - \alpha_1). \end{aligned}$$

Since $\Delta > 0$, it follows that $f''(x) > 0$ over $[0, n]$. Further, note that

$$\begin{aligned} f'(0) &= -\alpha_1 \Delta + \Delta n \log(1 - \alpha_1) < 0, \\ f'(n) &= (1 - \alpha_1)\Delta - (1 - \alpha_1)^n \Delta > 0. \end{aligned}$$

Solving $f'(x_n^*) = 0$ for the unique minimizer x_n^* , we obtain

$$\begin{aligned} \left(\frac{1}{1 - \alpha_1}\right)^{x_n^* - 1} - 1 &= (n - x_n^*) \log\left(\frac{1}{1 - \alpha_1}\right) \\ \implies \left(\frac{1}{1 - \alpha_1}\right)^{x_n^*} + x_n^* \log\left(\frac{1}{1 - \alpha_1}\right) &> n \log\left(\frac{1}{1 - \alpha_1}\right) \\ \implies 2 \left(\frac{1}{1 - \alpha_1}\right)^{x_n^*} &> n \log\left(\frac{1}{1 - \alpha_1}\right), \end{aligned}$$

where the last inequality follows using $y > \log y$. Therefore, we have

$$\begin{aligned} \left(\frac{1}{1-\alpha_1}\right)^{x_n^*} &> \frac{n}{2} \log\left(\frac{1}{1-\alpha_1}\right) \\ \implies x_n^* &> \frac{\log n + \log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{\log\left(\frac{1}{1-\alpha_1}\right)}. \end{aligned}$$

Thus, for any $k \in \{1, 2, \dots, n\}$,

$$\begin{aligned} \mathbb{E}[R_n^\pi | A_n^\pi = k] &\geq f(x_n^*) > (1-\alpha_1)\Delta x_n^* > (1-\alpha_1) \left(\frac{\log n + \log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{\log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta \\ \implies \mathbb{E}R_n^\pi &\geq (1-\alpha_1) \left(\frac{\log n + \log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{\log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta \\ \implies \inf_{\pi \in \Pi_n} \frac{\mathbb{E}R_n^\pi}{\log n} &\geq (1-\alpha_1) \left(\frac{1}{\log\left(\frac{1}{1-\alpha_1}\right)} + \frac{\log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{(\log n) \log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta \\ \implies \inf_{\pi \in \Pi_n} \frac{\mathbb{E}R_n^\pi}{\log n} &\stackrel{(\dagger)}{\geq} (1-\alpha_1) \left(\frac{1-\alpha_1}{\alpha_1} + \frac{\log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{(\log n) \log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta \\ \implies \inf_{\pi \in \Pi_n} \frac{\mathbb{E}R_n^\pi}{\log n} &\geq \frac{(1-\alpha_1)^2 \Delta}{\alpha_1} + (1-\alpha_1) \left(\frac{\log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{(\log n) \log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta \\ \implies \inf_{\pi \in \Pi_n} \frac{\mathbb{E}R_n^\pi}{\log n} &\stackrel{(\ddagger)}{\geq} \frac{\Delta}{4\alpha_1} + (1-\alpha_1) \left(\frac{\log \log\left(\frac{1}{1-\alpha_1}\right) - \log 2}{(\log n) \log\left(\frac{1}{1-\alpha_1}\right)} \right) \Delta, \end{aligned}$$

where (\dagger) follows using $\log y \leq y - 1$, and (\ddagger) since $\alpha_1 \leq 1/2$. Taking the appropriate limit now proves the assertion. \square

E Proof of Proposition 1

Consider the following stopping time:

$$\tau := \inf \left\{ m \in \mathbb{N} : \exists a, b \in \mathcal{A}, a < b \text{ s.t. } \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) < 4\sqrt{m \log m} \right\}.$$

Since $\mathbb{P}\left(\bigcap_{m \geq 1} \bigcap_{a, b \in \mathcal{A}, a < b} \left\{ \left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| \geq 4\sqrt{m \log m} \right\}\right) \geq \mathbb{P}(\tau = \infty)$, it suffices to show that $\mathbb{P}(\tau = \infty)$ is bounded away from 0. To this end, define the following entities:

$$\begin{aligned} \mathfrak{C}_K &:= \inf \left\{ p \in \mathbb{N} : \sum_{m=p}^{\infty} \frac{1}{m^8} \leq \frac{1}{2K^2} \right\}, \\ T_0 &:= \max \left(\left\lceil \left(\frac{64}{\delta^2} \right) \log^2 \left(\frac{64}{\delta^2} \right) \right\rceil, \Lambda_K \right), \\ f(x) &:= x + 4\sqrt{x \log x} \quad \text{for } x \geq 1. \end{aligned}$$

Lemma 1 For any $a, b \in \mathcal{A}$ s.t. $a < b$, it is the case that

$$\{\mathcal{Z}_{a,b} > f(T_0)\} \subseteq \bigcap_{m=1}^{T_0} \left\{ \mathcal{Z} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \geq 4\sqrt{m \log m} \right\}.$$

Proof of Lemma 1. Note that

$$\begin{aligned}
Z_{a,b} &> f(T_0) \\
&= T_0 + 4\sqrt{T_0 \log T_0} \\
&\geq m + 4\sqrt{m \log m} \quad \forall 1 \leq m \leq T_0 \\
&\stackrel{(a)}{\geq} \sum_{j=1}^m (X_{b,j} - X_{a,j}) + 4\sqrt{m \log m} \quad \forall 1 \leq m \leq T_0 \\
\implies Z_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) &\geq 4\sqrt{m \log m} \quad \forall 1 \leq m \leq T_0,
\end{aligned}$$

where (a) follows since the rewards are bounded in $[0, 1]$, i.e., $|X_{a,j} - X_{b,j}| \leq 1$. \square

Lemma 2 For $m \geq T_0$, it is the case that

$$\delta \geq 8\sqrt{\frac{\log m}{m}}.$$

Proof of Lemma 2. First of all, note that $T_0 \geq (64/\delta^2) \log^2(64/\delta^2) \geq 64$ (since $\delta \leq 1$). For $s = (64/\delta^2) \log^2(64/\delta^2)$, one has

$$\delta^2 = \frac{64 \log^2(\frac{64}{\delta^2})}{s} \stackrel{(b)}{\geq} \frac{64 [\log(\frac{64}{\delta^2}) + 2 \log \log(\frac{64}{\delta^2})]}{s} = \frac{64 \log s}{s},$$

where (b) follows since the function $g(x) := x^2 - x - 2 \log x$ is monotone increasing for $x \geq \log 64$ (think of $\log(64/\delta^2)$ as x), and therefore attains its minimum at $x = \log 64$; one can verify that this minimum is strictly positive. Furthermore, since $\log s/s$ is monotone decreasing for $s \geq 64$, it follows that for any $m \geq T_0$,

$$\delta^2 \geq \frac{64 \log m}{m}.$$

\square

Now coming back to the proof of Proposition 1, consider an arbitrary $l \in \mathbb{N}$ such that $l > T_0$. Then,

$$\begin{aligned}
\mathbb{P}(\tau \leq l) &= \mathbb{P}(\tau \leq l, Z > f(T_0)) + \mathbb{P}(\tau \leq l, Z \leq f(T_0)) \\
&\leq \mathbb{P}(\tau \leq l, Z > f(T_0)) + \Phi(f(T_0)).
\end{aligned}$$

Note that

$$\begin{aligned}
\mathbb{P}(\tau \leq l, \mathcal{Z} > f(T_0)) &= \mathbb{P}\left(\bigcup_{m=1}^l \bigcup_{a,b \in \mathcal{A}, a < b} \left\{ \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) < 4\sqrt{m \log m}, \mathcal{Z}_{a,b} > f(T_0) \right\}\right) \\
&\stackrel{(\dagger)}{=} \mathbb{P}\left(\bigcup_{m=T_0}^l \bigcup_{a,b \in \mathcal{A}, a < b} \left\{ \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) < 4\sqrt{m \log m}, \mathcal{Z}_{a,b} > f(T_0) \right\}\right) \\
&\leq \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}\left(\mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) < 4\sqrt{m \log m}, \mathcal{Z}_{a,b} > f(T_0)\right) \\
&= \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}\left(\mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j} - \delta) < -m \left(\delta - 4\sqrt{\frac{\log m}{m}}\right), \mathcal{Z}_{a,b} > f(T_0)\right) \\
&\leq \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}\left(\sum_{j=1}^m (X_{a,j} - X_{b,j} - \delta) < -m \left(\delta - 4\sqrt{\frac{\log m}{m}}\right), \mathcal{Z}_{a,b} > f(T_0)\right) \\
&\leq \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}\left(\sum_{j=1}^m (X_{a,j} - X_{b,j} - \delta) < -4\sqrt{m \log m}, \mathcal{Z}_{a,b} > f(T_0)\right) \\
&\stackrel{(\ddagger)}{=} \bar{\Phi}(f(T_0)) \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}\left(\sum_{j=1}^m (X_{a,j} - X_{b,j} - \delta) < -4\sqrt{m \log m}\right) \\
&\leq \stackrel{(\star)}{\bar{\Phi}(f(T_0))} \sum_{m=T_0}^l \sum_{a,b \in \mathcal{A}, a < b} \frac{1}{m^8} \\
&\leq \bar{\Phi}(f(T_0)) K^2 \sum_{m=T_0}^{\infty} \frac{1}{m^8} \\
&\stackrel{(\ast)}{\leq} \frac{\bar{\Phi}(f(T_0))}{2},
\end{aligned}$$

where (\dagger) follows from Lemma 1, (\ddagger) from Lemma 2, (\star) follows using the Chernoff-Hoeffding bound [15] and finally, (\ast) follows from the definition of T_0 . Therefore, we have

$$\begin{aligned}
\mathbb{P}(\tau \leq l) &\leq \frac{\bar{\Phi}(f(T_0))}{2} + \Phi(f(T_0)) = 1 - \frac{\bar{\Phi}(f(T_0))}{2} \\
\implies \mathbb{P}(\tau > l) &\geq \frac{\bar{\Phi}(f(T_0))}{2}.
\end{aligned}$$

Taking the limit $l \rightarrow \infty$ and appealing to the continuity of probability, we obtain

$$\begin{aligned}
\mathbb{P}(\tau = \infty) &\geq \frac{\bar{\Phi}(f(T_0))}{2} \\
\implies \mathbb{P}\left(\bigcap_{m \geq 1} \bigcap_{a,b \in \mathcal{A}, a < b} \left\{ \left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| \geq 4\sqrt{m \log m} \right\}\right) &\geq \frac{\bar{\Phi}(f(T_0))}{2}.
\end{aligned}$$

□

F Proof of Theorem 3

We will initially assume $\delta > 8\sqrt{\log n/n}$ for technical convenience. In the final step leading up to the asserted bound, we will relax this assumption by offsetting regret appropriately.

Let $\mathcal{A} := \{1, 2, \dots, K\}$. Define the following stopping times:

$$\begin{aligned}\tau_1 &:= \inf \left\{ m \in \mathbb{N} : \exists a, b \in \mathcal{A}, a < b \text{ s.t. } \left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < 4\sqrt{m \log m} \right\}, \\ \tau_2 &:= \inf \left\{ m \in \mathbb{N} : \left| \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| \geq 4\sqrt{m \log n} \forall a, b \in \mathcal{A}, a < b \right\}.\end{aligned}$$

Let R_t denote the cumulative pseudo-regret of CAB-K (calibrated for a horizon of play n) after $t \leq n$ pulls. Let D denote the event that the first batch of K arms queried from the reservoir is “all-distinct,” i.e., no two arms in this batch belong to the same type; let D^c be the complement of this event. Let CI denote the event that the algorithm commits to an inferior-typed arm. Let \tilde{R}_t, \bar{R}_t be independently drawn from the same distribution as R_t . Let $x^+ := \max(x, 0)$ for $x \in \mathbb{R}$. Then, R_n evolves according to the following stochastic recursion:

$$\begin{aligned}R_n &\leq \mathbb{1}\{\mathsf{D}\} \left[\mathbb{1}\{\tau_1 < \tau_2\} \left(\bar{\Delta} \min(K\tau_1, n) + \tilde{R}_{(n-K\tau_1)^+} \right) + \mathbb{1}\{\tau_1 \geq \tau_2\} \left(\bar{\Delta} \min(K\tau_2, n) + \mathbb{1}\{\mathsf{CI}\} \bar{\Delta} (n - K\tau_2)^+ \right) \right] \\ &\quad + \mathbb{1}\{\mathsf{D}^c\} \left[\mathbb{1}\{\tau_1 < \tau_2\} \left(\bar{\Delta} \min(K\tau_1, n) + \bar{R}_{(n-K\tau_1)^+} \right) + \mathbb{1}\{\tau_1 \geq \tau_2\} \bar{\Delta} n \right] \\ &\leq \mathbb{1}\{\mathsf{D}\} \left[\mathbb{1}\{\tau_1 < \tau_2\} \left(\bar{\Delta} \min(K\tau_2, n) + \tilde{R}_n \right) + \mathbb{1}\{\tau_1 \geq \tau_2\} \left(\bar{\Delta} \min(K\tau_2, n) + \mathbb{1}\{\mathsf{CI}\} \bar{\Delta} n \right) \right] \\ &\quad + \mathbb{1}\{\mathsf{D}^c\} \left[\mathbb{1}\{\tau_1 < \tau_2\} \left(\bar{\Delta} \min(K\tau_1, n) + \bar{R}_n \right) + \mathbb{1}\{\tau_1 \geq \tau_2\} \bar{\Delta} n \right] \\ &\leq \mathbb{1}\{\mathsf{D}\} \left[2\bar{\Delta} \min(K\tau_2, n) + \mathbb{1}\{\tau_1 < \tau_2\} \tilde{R}_n + \mathbb{1}\{\mathsf{CI}\} \bar{\Delta} n \right] + \mathbb{1}\{\mathsf{D}^c\} \left[\bar{\Delta} K\tau_1 + \bar{R}_n + \mathbb{1}\{\tau_1 \geq \tau_2\} \bar{\Delta} n \right].\end{aligned}$$

Taking expectations on both sides, one recovers using the independence of \tilde{R}_n, \bar{R}_n that

$$\begin{aligned}\mathbb{E}R_n &\leq \frac{\bar{\Delta}}{\mathbb{P}(\tau_1 \geq \tau_2 | \mathsf{D})} \left[2\mathbb{E}[\min(K\tau_2, n) | \mathsf{D}] + \left(\frac{\mathbb{P}(\mathsf{D}^c)}{\mathbb{P}(\mathsf{D})} \right) \mathbb{E}[K\tau_1 | \mathsf{D}^c] + \left(\mathbb{P}(\mathsf{CI} | \mathsf{D}) + \left(\frac{\mathbb{P}(\mathsf{D}^c)}{\mathbb{P}(\mathsf{D})} \right) \mathbb{P}(\tau_1 \geq \tau_2 | \mathsf{D}^c) \right) n \right],\end{aligned}$$

where $\mathbb{P}(\mathsf{D}) = K! \prod_{i=1}^K \alpha_i$.

F.1 Lower bounding $\mathbb{P}(\tau_1 \geq \tau_2 | \mathsf{D})$

Note that

$$\begin{aligned}\mathbb{P}(\tau_1 < \tau_2 | \mathsf{D}) &= \mathbb{P}(\tau_1 < \tau_2, \tau_2 = \infty | \mathsf{D}) + \mathbb{P}(\tau_1 < \tau_2, \tau_2 < \infty | \mathsf{D}) \\ &\leq \mathbb{P}(\tau_2 = \infty | \mathsf{D}) + \mathbb{P}(\tau_1 < \infty | \mathsf{D}) \\ &= \mathbb{P}(\tau_1 < \infty | \mathsf{D}) \\ &\leq 1 - \beta_{\delta, K},\end{aligned}$$

where the equality in the third step follows since τ_2 is almost surely finite on the event D , and the final inequality is due to Proposition 1. Thus, $\mathbb{P}(\tau_1 \geq \tau_2 | \mathsf{D}) \geq \beta_{\delta, K}$.

F.1.1 Proof that τ_2 is almost surely finite on D

Let $\mathbb{P}_{\mathsf{D}}(\cdot) := \mathbb{P}(\cdot | \mathsf{D})$ be the conditional measure w.r.t. the event D . Let $\mathcal{A} := \{1, 2, \dots, K\}$. Then, by continuity of probability, we have

$$\begin{aligned}\mathbb{P}_{\mathsf{D}}(\tau_2 = \infty) &= \lim_{l \rightarrow \infty} \mathbb{P}_{\mathsf{D}}(\tau_2 > l) \\ &= \lim_{l \rightarrow \infty} \mathbb{P}_{\mathsf{D}} \left(\bigcap_{m=1}^l \bigcup_{a, b \in \mathcal{A}, a < b} \left\{ \left| \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < 4\sqrt{m \log n} \right\} \right) \\ &\leq \lim_{l \rightarrow \infty} \sum_{a, b \in \mathcal{A}, a < b} \mathbb{P}_{\mathsf{D}} \left(\left| \sum_{j=1}^l (X_{a,j} - X_{b,j}) \right| < 4\sqrt{l \log n} \right).\end{aligned}$$

On D , it must be that $|\mathbb{E}[X_{a,j} - X_{b,j}]| \geq \delta$. Without loss of generality, assume that $\mathbb{E}[X_{a,j} - X_{b,j}] \geq \delta$. Then,

$$\begin{aligned} \mathbb{P}_{\mathsf{D}}(\tau_2 = \infty) &\leq \lim_{l \rightarrow \infty} \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathsf{D}} \left(\sum_{j=1}^l (X_{a,j} - X_{b,j}) < 4\sqrt{l \log n} \right) \\ &= \lim_{l \rightarrow \infty} \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathsf{D}} \left(\sum_{j=1}^l (X_{a,j} - X_{b,j} - \delta) < -l \left(\delta - 4\sqrt{\frac{\log n}{l}} \right) \right) \\ &\leq \lim_{l \rightarrow \infty} \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathsf{D}} \left(\sum_{j=1}^l (X_{a,j} - X_{b,j} - \delta) < -4l\sqrt{\log n} \left(\frac{2}{\sqrt{n}} - \frac{1}{\sqrt{l}} \right) \right), \end{aligned}$$

where the last inequality follows since $\delta > 8\sqrt{\log n/n}$ (by assumption). Now, using the Chernoff-Hoeffding bound [15] together with the fact that $-1 \leq X_{a,j} - X_{b,j} \leq 1$, we obtain for $l > n$ and any $a, b \in \mathcal{A}, a < b$ that

$$\begin{aligned} \mathbb{P}_{\mathsf{D}} \left(\sum_{j=1}^l (X_{a,j} - X_{b,j} - \delta) < -4l\sqrt{\log n} \left(\frac{2}{\sqrt{n}} - \frac{1}{\sqrt{l}} \right) \right) &\leq \exp \left[-8l \left(\frac{2}{\sqrt{n}} - \frac{1}{\sqrt{l}} \right)^2 \log n \right] \\ &= \exp \left[-8 \left(\frac{4l}{n} - 4\sqrt{\frac{l}{n}} + 1 \right)^2 \log n \right]. \end{aligned}$$

Summing over $a, b \in \mathcal{A}, a < b$ and taking the limit $l \rightarrow \infty$ proves the stated assertion. \square

F.2 Upper bounding $\mathbb{E}[\min(K\tau_2, n) | \mathsf{D}]$

Let $\mathbb{P}_{\mathsf{D}}(\cdot) := \mathbb{P}(\cdot | \mathsf{D})$ be the conditional measure w.r.t. the event D . Let $\mathcal{A} := \{1, 2, \dots, K\}$. Then,

$$\begin{aligned} \mathbb{E}[\min(K\tau_2, n) | \mathsf{D}] &= K \mathbb{E} \left[\min \left(\tau_2, \frac{n}{K} \right) | \mathsf{D} \right] \\ &\leq K \mathbb{E}[\min(\tau_2, n) | \mathsf{D}] \\ &\leq K + K \sum_{k=2}^n \mathbb{P}_{\mathsf{D}}(\tau_2 \geq k) \\ &\leq K + K \sum_{k=1}^n \mathbb{P}_{\mathsf{D}}(\tau_2 \geq k+1) \\ &\leq K + K \sum_{k=1}^n \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathsf{D}} \left(\left| \sum_{j=1}^k (X_{a,j} - X_{b,j}) \right| < 4\sqrt{k \log n} \right). \end{aligned}$$

On \mathcal{D} , it must be that $|\mathbb{E}[X_{a,j} - X_{b,j}]| \geq \delta$. Without loss of generality, assume that $\mathbb{E}[X_{a,j} - X_{b,j}] \geq \delta$. Then,

$$\begin{aligned} \mathbb{E}[\min(K\tau_2, n) | \mathcal{D}] &\leq K + K \sum_{k=1}^n \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathcal{D}} \left(\sum_{j=1}^k (X_{a,j} - X_{b,j}) < 4\sqrt{k \log n} \right) \\ &= K + K \sum_{k=1}^n \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathcal{D}} \left(\sum_{j=1}^k (X_{a,j} - X_{b,j} - \delta) < -k \left(\delta - 4\sqrt{\frac{\log n}{k}} \right) \right) \\ &\leq K + \frac{32K^3 \log n}{\delta^2} + K \sum_{k=\lceil \frac{64 \log n}{\delta^2} \rceil}^n \sum_{a,b \in \mathcal{A}, a < b} \mathbb{P}_{\mathcal{D}} \left(\sum_{j=1}^k (X_{a,j} - X_{b,j} - \delta) < \frac{-k\delta}{2} \right), \end{aligned}$$

where the last step follows since $\delta > 8\sqrt{\log n/n}$ (by assumption) implies $n > 64 \log n/\delta^2$, and $k \geq 64 \log n/\delta^2$ implies $\delta - 4\sqrt{\log n/k} \geq \delta/2$. Finally, using the Chernoff-Hoeffding inequality [15] together with the fact that $|X_{a,j} - X_{b,j}| \leq 1$, one obtains

$$\mathbb{E}[\min(K\tau_2, n) | \mathcal{D}] \leq K + \frac{32K^3 \log n}{\delta^2} + \frac{K^3}{2} \sum_{k=\lceil \frac{64 \log n}{\delta^2} \rceil}^n \exp\left(\frac{-\delta^2 k}{8}\right) \leq \frac{64K^3 \log n}{\delta^2}.$$

F.3 Upper bounding $\mathbb{E}[K\tau_1 | \mathcal{D}^c]$

The event \mathcal{D}^c will be implicitly assumed and we will drop the conditional argument for notational simplicity. Let $\mathcal{A} := \{1, 2, \dots, K\}$. Without loss of generality, suppose that arm 1 and 2 belong to the same type. Then,

$$\begin{aligned} \mathbb{E}[K\tau_1 | \mathcal{D}^c] &= K + K \sum_{k \geq 2} \mathbb{P}(\tau_1 \geq k) \\ &= K + K \sum_{k \geq 1} \mathbb{P}(\tau_1 \geq k + 1) \\ &= K + K \sum_{k \geq 1} \mathbb{P} \left(\bigcap_{m=1}^k \bigcap_{a,b \in \mathcal{A}, a < b} \left\{ \left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m} \right\} \right) \\ &\leq K + K \sum_{k \geq 1} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^k (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{k \log k} \right). \end{aligned}$$

Since $\mathcal{Z}_{a,b}$ is a standard Gaussian, and the increments $X_{1,j} - X_{2,j}$ are zero-mean sub-Gaussian with variance proxy 1, it follows from the Chernoff-Hoeffding concentration bound [15] that

$$\mathbb{E}[K\tau_1 | \mathcal{D}^c] \leq K + 2K \sum_{k \geq 1} \frac{1}{k^4} = \left(1 + \frac{\pi^4}{45}\right) K < 4K.$$

F.4 Upper bounding $\mathbb{P}(\text{CI} | \mathcal{D})$

Let $\mathbb{P}_{\mathcal{D}}(\cdot) := \mathbb{P}(\cdot | \mathcal{D})$ be the conditional measure w.r.t. the event \mathcal{D} . Let $\mathcal{A} := \{1, 2, \dots, K\}$ and without loss of generality, suppose that arm 1 is optimal (mean μ_1). Then,

$$\begin{aligned}
\mathbb{P}(\text{CI}|\mathbb{D}) &\leq \mathbb{P}_{\mathbb{D}} \left(\bigcup_{b=2}^K \left\{ \sum_{j=1}^{\tau_2} (X_{1,j} - X_{b,j}) \leq -4\sqrt{\tau_2 \log n} \right\} \right) \\
&\leq \sum_{b=2}^K \sum_{k=1}^n \mathbb{P}_{\mathbb{D}} \left(\sum_{j=1}^k (X_{1,j} - X_{b,j}) \leq -4\sqrt{k \log n} \right) + \sum_{b=2}^K \mathbb{P}_{\mathbb{D}}(\tau_2 > n) \\
&\leq \sum_{b=2}^K \sum_{k=1}^n \frac{1}{n^8} + \frac{K^3}{n^8} \\
&\leq \frac{K}{n^7} + \frac{K^3}{n^8},
\end{aligned}$$

where the second-to-last step follows using the Chernoff-Hoeffding inequality [15].

F.5 Upper bounding $\mathbb{P}(\tau_1 \geq \tau_2 | \mathbb{D}^c)$

Let $\mathbb{P}_{\mathbb{D}^c}(\cdot) := \mathbb{P}(\cdot | \mathbb{D}^c)$ be the conditional measure w.r.t. the event \mathbb{D}^c . Let $\mathcal{A} := \{1, 2, \dots, K\}$. On \mathbb{D}^c , there exist 2 arms in \mathcal{A} that belong to the same type; without loss of generality suppose that these arms are indexed by 1, 2. Then,

$$\begin{aligned}
\mathbb{P}(\tau_1 \geq \tau_2 | \mathbb{D}^c) &\leq \mathbb{P}(\tau_1 > n | \mathbb{D}^c) + \mathbb{P}(\tau_2 \leq n | \mathbb{D}^c) \\
&\leq \frac{2}{n^4} + \mathbb{P}_{\mathbb{D}^c}(\tau_2 \leq n | \mathbb{D}^c) \\
&\leq \frac{2}{n^4} + \mathbb{P}_{\mathbb{D}^c} \left(\bigcup_{m=1}^n \bigcap_{a,b \in \mathcal{A}, a < b} \left\{ \left| \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| \geq 4\sqrt{m \log n} \right\} \right) \\
&\leq \frac{2}{n^4} + \sum_{m=1}^n \mathbb{P}_{\mathbb{D}^c} \left(\left| \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log n} \right) \\
&\leq \frac{2}{n^4} + \frac{2}{n^7}, \tag{13}
\end{aligned}$$

where the last step follows using the Chernoff-Hoeffding bound [15].

F.6 Putting everything together

Combining everything, one finally obtains that when $\delta > 8\sqrt{\log n/n}$,

$$\mathbb{E}R_n \leq \frac{CK^3\bar{\Delta}}{\beta_{\delta,K}} \left(\frac{\log n}{\delta^2} + \frac{1}{\mathbb{P}(\mathbb{D})} \right),$$

where $\beta_{\delta,K}$ is as defined in (2), $\mathbb{P}(\mathbb{D}) = K! \prod_{i=1}^K \alpha_i$, and C is some absolute constant. When $\delta \leq 8\sqrt{\log n/n}$, regret is at most $\bar{\Delta}n \leq 64\bar{\Delta}/\delta^2 \log n$. Thus, the aforementioned bound, in fact, holds generally for some large enough absolute constant C . \square

G Exponential doubling for “anytime” algorithms with logarithmic regret

Below, we propose an exponential doubling trick that can be used to make algorithms horizon-free while preserving their instance-dependent regret upper bounds in the leading order. The reader is referred to [5] for a comprehensive survey of standard doubling tricks used in the literature.

Algorithm 3 HF(\mathfrak{A}) (EXPONENTIAL DOUBLING)

- 1: **Input:** Algorithm \mathfrak{A} .
 - 2: **Generate doubling sequence:** $T_{-1} = 0; T_i = 2^{2^i}$ for $i = 0, 1, \dots$
 - 3: **for** $i \in \{0, 1, \dots\}$ **do**
 - 4: Specify as input to \mathfrak{A} a horizon of play of $T_i - T_{i-1}$.
 - 5: Restart \mathfrak{A} at time $t = T_{i-1} + 1$ and run until $t = T_i$.
-

Theorem 7 (Preservation of logarithmic regret under HF(\mathfrak{A})) Consider an algorithm \mathfrak{A} (possibly horizon-dependent) with an expected regret of $\mathbb{E}R_n^{\mathfrak{A}} \leq C \log n + D$ at the end of a horizon of $n \geq 1$ plays, where C and D are non-negative constants. Then, the expected regret of the horizon-free policy π given by HF(\mathfrak{A}) after any number $n \geq 1$ of plays is bounded as

$$\mathbb{E}R_n^\pi \leq 80 [C \log n + D \log \log(n+2)].$$

Discussion. The proof is provided in Appendix G.1. It is possible to improve the multiplicative constants in Theorem 7 via a finer calibration of the doubling sequence $(T_i : i = 0, 1, \dots)$. However, our goal simply is to design horizon-free policies for the matching problem that are rate-optimal in n modulo constant multiplicative factors and to that end, Theorem 7 serves its purpose. We now state an anytime upper bound on the regret of the horizon-free version of CAB-K.

Theorem 8 (Upper bound on the regret of HF(CAB-K)) The expected regret of the policy π given by HF(CAB-K) after any number $n \geq 1$ of plays is bounded as

$$\mathbb{E}R_n^\pi \leq \frac{CK^3\bar{\Delta}}{\beta_{\delta,K}} \left(\frac{\log n}{\delta^2} + \frac{\log \log(n+2)}{K! \prod_{i=1}^K \alpha_i} \right),$$

where $\beta_{\delta,K}$ is as defined in (2) and C is some absolute constant.

G.1 Proof of Theorem 7

In what follows, \log corresponds to the natural logarithm unless a base is explicitly specified otherwise. Define $\Lambda_n := \min \{i \in \mathbb{N} : T_i \geq n\} = \lceil \log_2 \log_2 n \rceil$. Since the algorithm \mathfrak{A} is restarted at times $\{T_i + 1 : i = -1, 0, 1, \dots\}$, where $T_{-1} := 0$, it follows that the expected cumulative regret of the horizon-free policy π given by HF(\mathfrak{A}) after n plays is bounded as

$$\begin{aligned} \mathbb{E}R_n^\pi &\leq \sum_{i=0}^{\Lambda_n} \mathbb{E}R_{T_i - T_{i-1}}^{\mathfrak{A}} \leq \sum_{i=0}^{\Lambda_n} (C \log T_i + D) = \sum_{i=0}^{\Lambda_n} \left(\frac{C}{\log_2 e} \log_2 T_i + D \right) \\ &\leq \sum_{i=0}^{\Lambda_n} (C2^i + D) \\ &= C2^{\Lambda_n+1} + D(\Lambda_n + 1) - C \\ &= 4C \log_2 n + D \log_2 \log_2 n + 2D - C \\ &\leq 8C \log n + 2D \log \log n + 3D - C \\ &\leq 8C \log n + 80D \log \log(n+2). \end{aligned}$$

□

G.2 Proof of Theorem 8

This proof is a direct application of Theorem 7. □

H Proof of Theorem 4

Note that there are at most $M |\mathcal{J}|$ active threads of HF(CAB-K) at any time. Since HF(CAB-K) is horizon-free, it naturally follows that the regret incurred under π is dominated by that in the scenario where all $M |\mathcal{J}|$ threads are active at each $t \in \{1, \dots, n\}$. The assertion is now immediate. □

I A first-order optimal algorithm for countable-armed bandits with $K \geq 2$

Algorithm 4 CAB-K(UCB): Nested UCB1 for K types

- 1: **Initialize new epoch** (resets clock $t \leftarrow 0$): Query K new arms; call it set $\mathcal{A} = \{1, 2, \dots, K\}$.
 - 2: Play each arm in \mathcal{A} once; observe rewards $\{X_{a,1} : a \in \mathcal{A}\}$.
 - 3: Minimum per-arm sample count $m \leftarrow 1$.
 - 4: Generate $\binom{K}{2}$ independent standard Gaussian random variables $\{\mathcal{Z}_{a,b} : a, b \in \mathcal{A}, a < b\}$.
 - 5: **for** $t \in \{K+1, K+2, \dots\}$ **do**
 - 6: **if** $\exists a, b \in \mathcal{A}, a < b$ s.t. $\left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < 4\sqrt{m \log m}$ **then**
 - 7: Permanently discard \mathcal{A} and repeat from step (1).
 - 8: **else**
 - 9: Play arm $a_t \in \arg \max_{a \in \mathcal{A}} \left(\frac{\sum_{j=1}^{N_a(t-1)} X_{a,j}}{N_a(t-1)} + \sqrt{\frac{2 \log(t-1)}{N_a(t-1)}} \right)$.
 - 10: Observe reward $X_{a_t, N_{a_t}(t)}$.
 - 11: **if** $m < \min_{a \in \mathcal{A}} N_a(t)$ **then**
 - 12: $m \leftarrow m + 1$.
-

Theorem 9 (Upper bound on the regret of CAB-K(UCB)) *The expected regret of the policy π given by CAB-K(UCB) after any number $n \geq 1$ of plays is bounded as*

$$\mathbb{E}R_n^\pi \leq \frac{CK}{\beta_{\delta,K}} \left(\frac{\log n}{\underline{\Delta}} + \bar{\Delta} \right) + o\left(\frac{\bar{\Delta}n}{\beta_{\delta,K} \prod_{i=1}^K \alpha_i} \right),$$

where C is some absolute constant, $\beta_{\delta,K}$ is as defined in (2), and the little-Oh is asymptotic in n and only hides multiplicative factors in K .

The proof is provided in §I.1.

Remark 1 *The upper bound is logarithmic in n for $K = 2$; refer to Theorem 5. Whether this would hold also for general K remains an open problem.*

I.1 Proof of Theorem 9

Let $\mathcal{A} := \{1, 2, \dots, K\}$ be the collection of K arms queried during the first epoch. Consider an arbitrary $l \in \mathbb{N}$ s.t. $l \geq K$ and define the following:

$$M_l := \min_{a \in \mathcal{A}} N_a(l),$$

$$\tau := \inf \left\{ l \geq K : \exists a, b \in \mathcal{A}, a < b \text{ s.t. } \left| \mathcal{Z}_{a,b} + \sum_{j=1}^{M_l} (X_{a,j} - X_{b,j}) \right| < 4\sqrt{M_l \log M_l} \right\}, \quad (14)$$

where $N_a(l)$ denotes the sample count from arm a under UCB1 until time l . Note that τ marks the termination of epoch 1.

Let R_n denote the cumulative pseudo-regret of CAB-K(UCB) after n pulls (superscript π suppressed for notational convenience). Let S_n denote the cumulative pseudo-regret of UCB1 after n pulls in a K -armed bandit with means $\mu_1 > \mu_2 > \dots > \mu_K$. Let D denote the event that the K arms queried in epoch 1 have distinct types (no two belong to the same type). Similarly, let OPT denote the event that the K arms have optimal types (type 1). Let $\tilde{R}_n, \hat{R}_n, \bar{R}_n$ be independently drawn from the same distribution as R_n . Then, the evolution of R_n satisfies

$$\begin{aligned} R_n &\leq \mathbb{1}\{\mathsf{D}\} \left[S_{\min(\tau, n)} + \tilde{R}_{(n-\tau)^+} \right] + \mathbb{1}\{\mathsf{D}^c \setminus \mathsf{OPT}\} \left[\bar{\Delta} \min(\tau, n) + \hat{R}_{(n-\tau)^+} \right] + \mathbb{1}\{\mathsf{OPT}\} \bar{R}_{(n-\tau)^+} \\ &\stackrel{(t)}{\leq} \mathbb{1}\{\mathsf{D}\} \left[S_n + \mathbb{1}\{\tau < n\} \tilde{R}_n \right] + \mathbb{1}\{\mathsf{D}^c \setminus \mathsf{OPT}\} \bar{\Delta} \min(\tau, n) + \mathbb{1}\{\mathsf{D}^c \setminus \mathsf{OPT}\} \hat{R}_n + \mathbb{1}\{\mathsf{OPT}\} \bar{R}_n \\ &\leq \mathbb{1}\{\mathsf{D}\} \left[S_n + \mathbb{1}\{\tau < \infty\} \tilde{R}_n \right] + \mathbb{1}\{\mathsf{D}^c \setminus \mathsf{OPT}\} \bar{\Delta} \min(\tau, n) + \mathbb{1}\{\mathsf{D}^c \setminus \mathsf{OPT}\} \hat{R}_n + \mathbb{1}\{\mathsf{OPT}\} \bar{R}_n, \end{aligned}$$

where (\dagger) follows since CAB-K(UCB) is agnostic to n , and hence the pseudo-regret R_n is weakly increasing in n . Taking expectations on both sides, one recovers using the independence of $\tilde{R}_n, \hat{R}_n, \bar{R}_n$ that

$$\begin{aligned} \mathbb{E}R_n &\leq \frac{1}{\mathbb{P}(\tau = \infty | \mathbf{D})} \left[\mathbb{E}S_n + \left(\frac{\bar{\Delta} \mathbb{P}(\mathbf{D}^c \setminus \text{OPT}) \mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}]}{\mathbb{P}(\mathbf{D})} \right) \right] \\ &\leq \frac{1}{\beta_{\delta, K}} \left[\mathbb{E}S_n + \left(\frac{\bar{\Delta} \mathbb{P}(\mathbf{D}^c \setminus \text{OPT}) \mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}]}{\mathbb{P}(\mathbf{D})} \right) \right], \end{aligned} \quad (15)$$

where $\mathbb{P}(\mathbf{D}) = K! \prod_{i=1}^K \alpha_i$, and the last inequality follows using Lemma 3 with $\beta_{\delta, K}$ as defined in (2). We know that $\mathbb{E}S_n \leq CK (\log n / \underline{\Delta} + \bar{\Delta})$ for some absolute constant C [1]. The rest of the proof is geared towards showing that $\mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}] = o(n)$.

I.1.1 Proof of $\mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}] = o(n)$

Let $\mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}}(\cdot) := \mathbb{P}(\cdot | \mathbf{D}^c \setminus \text{OPT})$ be the conditional measure w.r.t. the event $\mathbf{D}^c \setminus \text{OPT}$. On $\mathbf{D}^c \setminus \text{OPT}$, there exist two arms in the consideration set \mathcal{A} that belong to the same type. Without loss of generality, suppose that these are indexed by 1, 2. Then,

$$\begin{aligned} \mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}] &\leq K + \sum_{k=K+1}^n \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}}(\tau \geq k) \\ &\leq K + \sum_{k=K}^n \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}}(\tau \geq k+1) \\ &= K + \sum_{k=K}^n \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}} \left(\bigcap_{l=1}^k \bigcap_{a, b \in \mathcal{A}, a < b} \left\{ \left| \mathcal{Z}_{a,b} + \sum_{j=1}^{M_l} (X_{a,j} - X_{b,j}) \right| \geq 4\sqrt{M_l \log M_l} \right\} \right) \\ &\leq K + \sum_{k=K}^n \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_k} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_k \log M_k} \right) \\ &= K + \sum_{k=K}^n \sum_{m=1}^k \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_k} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_k \log M_k}, M_k = m \right) \\ &\stackrel{(\dagger)}{=} K + \sum_{k=K}^n \sum_{m=f(k)}^k \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_k} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_k \log M_k}, M_k = m \right) \\ &\leq K + \sum_{k=K}^n \sum_{m=f(k)}^k \mathbb{P}_{\mathbf{D}^c \setminus \text{OPT}} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m} \right) \\ &\stackrel{(\ddagger)}{\leq} K + 2 \sum_{k=K}^n \sum_{m=f(k)}^k \frac{1}{m^4} \\ &= K + 2 \sum_{k=K}^n \left(\frac{1}{(f(k))^4} + \sum_{m=f(k)+1}^k \frac{1}{m^4} \right) \\ &\leq K + 2 \sum_{k=K}^n \left(\frac{1}{(f(k))^4} + \frac{1}{3(f(k))^3} \right), \end{aligned}$$

where (\dagger) follows from Lemma 4, and (\ddagger) using the Chernoff-Hoeffding bound [15]. Since $f(k)$ is monotone non-decreasing and coercive in k , it follows that $\mathbb{E}[\min(\tau, n) | \mathbf{D}^c \setminus \text{OPT}] = o(n)$, where the little-Oh only hides dependence on K . \square

I.2 Proof of Theorem 5

We know from (15) that

$$\mathbb{E}R_n \leq \frac{1}{\beta_{\Delta,2}} \left[\mathbb{E}S_n + \left(\frac{\Delta \mathbb{P}(\mathcal{D}^c \setminus \text{OPT}) \mathbb{E}[\min(\tau, n) | \mathcal{D}^c \setminus \text{OPT}]}{\mathbb{P}(\mathcal{D})} \right) \right],$$

where τ is as defined in (14), $\mathbb{P}(\mathcal{D}) = 2\alpha_1(1 - \alpha_1)$, $\mathbb{P}(\mathcal{D}^c \setminus \text{OPT}) = (1 - \alpha_1)^2$, and $\mathbb{E}S_n \leq C(\log n/\underline{\Delta} + \underline{\Delta})$ for some absolute constant C [1]. The rest of the proof is geared towards showing that $\mathbb{E}[\min(\tau, n) | \mathcal{D}^c \setminus \text{OPT}] \leq \mathbb{E}[\tau | \mathcal{D}^c \setminus \text{OPT}] \leq C'$ (for some absolute constant C'). Going forward, we will assume that the probabilities are implicitly conditional to avoid overloading notation. Then, note that

$$\begin{aligned} \mathbb{E}[\tau | \mathcal{D}^c \setminus \text{OPT}] &= 1 + \sum_{k \geq 2} \mathbb{P}(\tau \geq k) \\ &= 1 + \sum_{k \geq 2} \mathbb{P} \left(\bigcap_{l=1}^{k-1} \left\{ \left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_l} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_l \log M_l} \right\} \right) \\ &\leq 1 + \sum_{k \geq 1} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_k} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_k \log M_k} \right) \\ &= 1 + \sum_{k \geq 1} \sum_{m=1}^k \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{M_k} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{M_k \log M_k}, N_1(k) = m \right) \\ &= 1 + \sum_{k \geq 1} \sum_{1 \leq m \leq k/2} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m}, N_1(k) = m \right) \\ &\quad + \sum_{k \geq 1} \sum_{k/2 < m \leq k} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^{(k-m)} (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{(k-m) \log(k-m)}, N_1(k) = m \right) \\ &= 1 + \sum_{k \geq 1} \sum_{1 \leq m \leq k/2} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m}, N_1(k) = m \right) \\ &\quad + \sum_{k \geq 1} \sum_{1 \leq m \leq k/2} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m}, N_2(k) = m \right) \\ &\leq 1 + 2 \sum_{k \geq 1} \sum_{\theta k \leq m \leq k/2} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m} \right) \\ &\quad + \sum_{k \geq 1} [\mathbb{P}(N_1(k) \leq \theta k) + \mathbb{P}(N_2(k) \leq \theta k)], \end{aligned}$$

where $\theta = 1/2 - \sqrt{15}/8$. Using Theorem 4(i) of [18] with $\epsilon = \sqrt{15}/8$, one obtains

$$\begin{aligned} \mathbb{E}[\tau | \mathcal{D}^c \setminus \text{OPT}] &\leq 1 + 2 \sum_{k \geq 1} \sum_{\theta k \leq m \leq k/2} \mathbb{P} \left(\left| \mathcal{Z}_{1,2} + \sum_{j=1}^m (X_{1,j} - X_{2,j}) \right| \geq 4\sqrt{m \log m} \right) + 16 \sum_{k \geq 1} \frac{1}{k^2} \\ &\leq 1 + 4 \sum_{k \geq 1} \sum_{\theta k \leq m \leq k/2} \frac{1}{m^4} + 16 \sum_{k \geq 1} \frac{1}{k^2} \\ &\leq 1 + \frac{4}{\theta^4} \sum_{k \geq 1} \frac{1}{k^3} + 16 \sum_{k \geq 1} \frac{1}{k^2} \\ &=: C'. \end{aligned}$$

□

J Auxiliary results used in the analysis of CAB-K(UCB)

Lemma 3 (Persistence of heterogeneous consideration sets) Consider a K -armed bandit with rewards bounded in $[0, 1]$ and means $\mu_1 > \mu_2 > \dots > \mu_K$. Let $\{X_{a,j} : j = 1, 2, \dots\}$ denote the rewards collected from arm $a \in \{1, \dots, K\} =: \mathcal{A}$ by UCB1 [1]. Let $\{\mathcal{Z}_{a,b} : a, b \in \mathcal{A}, a < b\}$ be a collection of $\binom{K}{2}$ independent standard Gaussian random variables. Let $N_a(n)$ be the sample count of arm a under UCB1 until time n . Define

$$M_l := \min_{a \in \mathcal{A}} N_a(l),$$

$$\tau := \inf \left\{ l \geq K : \exists a, b \in \mathcal{A}, a < b \text{ s.t. } \left| \mathcal{Z}_{a,b} + \sum_{j=1}^{M_l} (X_{a,j} - X_{b,j}) \right| < 4\sqrt{M_l \log M_l} \right\}.$$

Then, $\mathbb{P}(\tau = \infty) > \beta_{\delta, K}$, where $\beta_{\delta, K}$ is as defined in (2).

Lemma 4 (Path-wise lower bound on the arm-sampling rate of UCB1) Consider a K -armed bandit with rewards bounded in $[0, 1]$. Let $N_a(n)$ be the sample count of arm $a \in \{1, \dots, K\} =: \mathcal{A}$ under UCB1 [1] until time n . Then, for all $n \geq K$,

$$M_n := \min_{a \in \mathcal{A}} N_a(n) \geq f(n),$$

where $(f(n) : n = K, K + 1, \dots)$ is some deterministic monotone non-decreasing integer-valued sequence satisfying $f(K) = 1$ and $f(n) \rightarrow \infty$ as $n \rightarrow \infty$.

J.1 Proof of Lemma 3

Suppose that there exists a sample-path on which some non-empty subset of arms $\mathfrak{A} \subset \mathcal{A}$ receives a bounded number of pulls asymptotically in the horizon of play. Also suppose that \mathfrak{A} is the maximal such subset, i.e., each arm in $\mathcal{A} \setminus \mathfrak{A}$ is played infinitely often asymptotically on said sample-path. This implies that the UCB score of any arm in $\mathcal{A} \setminus \mathfrak{A}$ is of the order $o(\sqrt{\log t})$ at time t (the empirical mean term remains bounded in $[0, 1]$ and therefore can be ignored). At the same time, the boundedness hypothesis implies that the UCB score of any arm in \mathfrak{A} will grow as $\Omega(\sqrt{\log t})$. Thus, for t large enough, UCB scores of arms in \mathfrak{A} will start to dominate those in $\mathcal{A} \setminus \mathfrak{A}$ and the algorithm will end up playing an arm from \mathfrak{A} at some point, thus increasing the cumulative sample-count of arms in \mathfrak{A} by 1. As t grows further, one can replicate the preceding argument an arbitrary number of times to conclude that \mathfrak{A} receives an unbounded number of pulls on the sample-path under consideration, thereby contradicting the boundedness hypothesis. Therefore, it must be the case that each arm in \mathcal{A} is played infinitely often on every sample-path. Consequently, $M_n = \min_{a \in \mathcal{A}} N_a(n) \rightarrow \infty$ as $n \rightarrow \infty$ on every sample-path.

Now since $(M_n : n = K, K + 1, \dots)$ is an integer-valued process (starting from $M_K = 1$) with unit increments (wherever they exist), it follows that on every sample-path, τ , in fact, weakly dominates the stopping time τ' given by

$$\tau' := \inf \left\{ m \in \mathbb{N} : \exists a, b \in \mathcal{A}, a < b \text{ s.t. } \left| \mathcal{Z}_{a,b} + \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < 4\sqrt{m \log m} \right\}. \quad (16)$$

Therefore, $\mathbb{P}(\tau = \infty) \geq \mathbb{P}(\tau' = \infty) > \beta_{\delta, K}$; the last inequality follows from Proposition 1. \square

J.2 Proof of Lemma 4

Suppose that $\mathcal{S}_n = \{(N_a(n) : a \in \mathcal{A})\}$ denotes the set of possible sample-count realizations under UCB1 when the horizon of play is n . Define $f(n) := \min_{(N_a(n) : a \in \mathcal{A}) \in \mathcal{S}_n} \min_{a \in \mathcal{A}} N_a(n)$. Since \mathcal{S}_n is finite, aforementioned minimum is attained at some $(N_a^*(n) : a \in \mathcal{A}) \in \mathcal{S}_n$. Note that $(N_a^*(n) : a \in \mathcal{A})$ is not a random vector as it corresponds to a specific set of sample-paths (possibly non-unique) on which $\min_{a \in \mathcal{A}} N_a(n)$ is minimized. Therefore, $f(n) = \min_{a \in \mathcal{A}} N_a^*(n)$ is deterministic. We have already established in the proof of Lemma 3 that for each $a \in \mathcal{A}$, $N_a(n) \rightarrow \infty$ as $n \rightarrow \infty$ on every sample-path. In particular, this also implies $N_a^*(n) \rightarrow \infty$ as $n \rightarrow \infty$. Thus, we have established the existence of a sequence $f(n)$ satisfying the assertions of the lemma. \square

K More Explore-then-Commit policies for countable-armed bandits

Algorithm 5 ETC- ∞ (K)

- 1: **Input:** (i) Horizon of play $n \geq K$, (ii) A lower bound $\underline{\delta} \in (0, \delta]$ on the minimal reward gap δ .
 - 2: Set budget $T = n$.
 - 3: **Initialize new epoch:** Query K new arms; call it consideration set $\mathcal{A} = \{1, 2, \dots, K\}$.
 - 4: Set exploration duration $L = \lceil 2\underline{\delta}^{-2} \log n \rceil$.
 - 5: $m \leftarrow \min(L, \lfloor T/K \rfloor)$.
 - 6: Play each arm in \mathcal{A} m times; observe rewards $\{(X_{1,j}, \dots, X_{K,j}) : j = 1, \dots, m\}$.
 - 7: Update budget: $T \leftarrow T - Km$.
 - 8: **if** $\exists a, b \in \mathcal{A}, a < b$ s.t. $\left| \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < \underline{\delta}m$ **then**
 - 9: Permanently discard \mathcal{A} , and repeat from step (3).
 - 10: **else**
 - 11: Permanently commit to arm $a^* \in \arg \max_{a \in \mathcal{A}} \left\{ \sum_{j=1}^m X_{a,j} \right\}$.
-

Algorithm 6 ETC-RAW

- 1: **Input:** Horizon of play $n \geq K$.
 - 2: Set budget $T = n$; set epoch counter $k = 1$.
 - 3: **Initialize new epoch:** Query K new arms; call it consideration set $\mathcal{A} = \{1, 2, \dots, K\}$.
 - 4: Set exploration duration $L = \lceil e^{2\sqrt{k}} \log n \rceil$.
 - 5: $m \leftarrow \min(L, \lfloor T/K \rfloor)$.
 - 6: Play each arm in \mathcal{A} m times; observe rewards $\{(X_{1,j}, \dots, X_{K,j}) : j = 1, \dots, m\}$.
 - 7: Update budget: $T \leftarrow T - Km$.
 - 8: **if** $\exists a, b \in \mathcal{A}, a < b$ s.t. $\left| \sum_{j=1}^m (X_{a,j} - X_{b,j}) \right| < 2me^{-\sqrt{k}}$ **then**
 - 9: Permanently discard \mathcal{A} .
 - 10: $k \leftarrow k + 1$.
 - 11: Repeat from step (3).
 - 12: **else**
 - 13: Permanently commit to arm $a^* \in \arg \max_{a \in \mathcal{A}} \left\{ \sum_{j=1}^m X_{a,j} \right\}$.
-