

BANDITS WITH ANYTIME KNAPSACKS

Anonymous authors

Paper under double-blind review

ABSTRACT

We consider bandits with anytime knapsacks (BwAK), a novel version of the BwK problem where there is an *anytime* cost constraint instead of a total cost budget. This problem setting introduces additional complexities as it mandates adherence to the constraint throughout the decision-making process. We propose SUAK, an algorithm that utilizes upper confidence bounds to identify the optimal mixture of arms while maintaining a balance between exploration and exploitation. SUAK is an adaptive algorithm that strategically utilizes the available budget in each round in the decision-making process and skips a round when it is possible to violate the anytime cost constraint. In particular, SUAK slightly under-utilizes the available cost budget to reduce the need for skipping rounds. We show that SUAK attains the same problem-dependent regret upper bound of $O(K \log T)$ established in prior work under the simpler BwK framework. Finally, we provide simulations to verify the utility of SUAK in practical settings.

1 INTRODUCTION

Multi-armed bandits (MAB) is one of the fundamental problems in the field of sequential decision-making under uncertainty. In its essence, it is a problem setting where an agent must strategically allocate resources among the arms to maximize cumulative reward over time, navigating the trade-off between gathering information about uncertain arms (exploration) and exploiting known information to optimize immediate rewards (exploitation). This problem finds applications across diverse domains, including reinforcement learning (Intayoad et al., 2020), online advertising (Slivkins, 2013), clinical trials (Villar et al., 2015), and resource allocation (Soare, 2015).

The *Bandits with Knapsacks* (BwK) problem, introduced by Badanidiyuru et al. (2013), is an extension of the classical multi-armed bandit problem, with the additional constraint of limited resource capacity akin to the *knapsack* problem (Tran-Thanh et al., 2012). In this scenario, an agent is confronted with a set of arms, each associated with an *unknown* reward and cost distribution. Unlike the traditional bandit setting, selecting an arm incurs both a reward and a cost here, and the agent’s objective is to maximize the total reward while respecting the total capacity constraint of the knapsack. The BwK problem encapsulates the trade-off between exploration and exploitation while managing resource constraints, presenting a rich framework with applications such as online advertising (Avadhanula et al., 2021; Badanidiyuru et al., 2018), dynamic resource allocation (Kumar & Kleinberg, 2022), and personalized recommendation engines (Yu et al., 2016).

In this paper, we consider a specific variant of this problem, which we name as the bandits with anytime knapsacks (BwAK) problem; where instead of a total cost budget, there is an anytime constraint on the average cost. This problem setting introduces an additional level of complexity as a mixture strategy needs to be employed to be able to pull arms with mean costs higher than the average cost budget without violating the anytime constraint. The main goal of our work is to develop new algorithms for this framework that achieve as much *cumulative* reward as possible.

1.1 APPLICATIONS

The formulation of the anytime constraint considered here has broad applications across various fields. A notable example is inventory management, where a factory produces goods at a constant rate and seeks to maximize revenue by selling to buyers in a marketplace, where bids consisting of price and order size are placed. Our anytime constraint is especially relevant in such scenarios since having a negative inventory is not possible. An important aspect of this constraint is that it introduces

the trade-off between exploiting the available inventory, and skipping a round to accumulate more inventory in an effort to capture bids with higher order sizes. This highlights the added complexity of our problem setting and underscores its broader applicability across a range of settings beyond the standard BwK framework.

Another example is online advertising, where an advertiser sets a daily budget limit to prevent over-spending. In this context, the 'arms' represent different ad campaigns or strategies, each with varying costs that can be selected for the day. The reward can be modeled as the daily revenue generated from clicks or the number of users who subscribe. Further, in portfolio management, our anytime constraint can represent the maximum amount that the customer is willing to invest in a month, and arms can model different investment options. One last example is in satellite systems, where solar panels generate energy and excess energy can be stored in a battery. Here, c can correspond to the energy generated per unit time, and arms can correspond to different tasks that need to be performed, with their rewards reflecting the importance or outcome of the tasks. The costs associated with each arm can then represent the energy consumed to complete the task.

1.2 CONTRIBUTIONS

1. **Formulation:** To our knowledge, this work is the first to consider a multi-armed bandit with knapsacks (BwK) setting with an *anytime* cost constraint.
2. **SUAK Algorithm:** SUAK utilizes the upper confidence bounds to explore the best base that solves the problem, and also strategically under-utilizes the available budget to limit the number of rounds that are skipped when satisfying the anytime cost constraint.
3. **Regret Upper Bound for SUAK:** We provide upper bounds on the expected cumulative regret of SUAK for this problem setting, and establish that it scales as $O(K \log T)$.

Related works is provided in §3.4.

2 PROBLEM STATEMENT

2.1 THE BANDITS WITH ANYTIME KNAPSACKS (BWAK) MODEL

We consider a K -armed stochastic bandit problem with the set of base arms $[K]$, where pulling arm $i \in [K]$ in round t is associated with a random cost, $\rho_i(t)$; drawn from a probability distribution supported in $[0, 1]$ with mean ρ_i , that is independent of the costs of other arms. After pulling arm i in round t , the agent receives a random reward, $r_i(t)$; drawn from a probability distribution supported in $[0, 1]$ with mean μ_i , that is independent of the rewards of other arms. At each round t , the agent has the option of skipping by not pulling any of the K arms. We model this decision by introducing an arm which has a cost and reward of 0, as arm $K + 1$, which is known as the *null arm* in BwK literature. We let $\boldsymbol{\rho} = [\rho_1, \dots, \rho_K, 0]^T$ and $\boldsymbol{\mu} = [\mu_1, \dots, \mu_K, 0]^T$ denote the mean cost vector and the mean reward vector of the arm set $[K + 1]$, respectively. Throughout this paper, we use bold symbols to denote vectors or matrices. Δ_{K+1} is used to denote the $K + 1$ -dimensional probability simplex. We let $i^* := \arg \max_{i \in [K]} \mu_i / \rho_i$ denote the arm with highest mean reward per cost, and let $i^{**} := \arg \max_{i \in [K]} \mu_i$ denote the arm with the highest mean reward. For simplicity, we assume that there is only one arm with highest mean reward per cost and also there is only one arm with highest mean reward.

Let $i(t)$ be the arm pulled by the agent in round t , $r(t)$ represent the reward received in round t , and $c(t)$ represent the cost incurred in round t . Also let $N_i(t)$ denote the total number of times arm i has been pulled up to round t . Further, define $S_c(t) = \sum_{s=1}^t c(s)$ and $\bar{c}(t) = S_c(t)/t$ as the cumulative cost and the average cost incurred until round t . Let $\bar{\rho}_i(t) = \sum_{s=1}^t \rho_i(s) \cdot \mathbb{1}\{i(t) = i\} / N_i(t)$ be the empirical average cost of arm i at round t , and similarly let $\bar{\mu}_i(t) = \sum_{s=1}^t \mu_i(s) \cdot \mathbb{1}\{i(t) = i\} / N_i(t)$ be the empirical average reward. We assume that there is an average cost budget of c per round that cannot be exceeded at any round, which we refer to as the anytime cost constraint. The agent aims to maximize cumulative reward received under this constraint. This can formally be expressed as:

$$\text{maximize } F(t) = \mathbb{E} \left[\frac{1}{t} \sum_{s=1}^t r(s) \right] \quad \text{s.t.} \quad \frac{\sum_{s=1}^u c(s)}{u} \leq c, \forall u \leq t.$$

This setting represents many practical applications as discussed in §1.1.

Linear Relaxation. Following the prior work, we consider the following linear relaxation:

$$\begin{aligned} OPT_{LP}(T) &= \max_{\pi} T \cdot \mu^T \pi \\ \text{s.t. } &\rho^T \pi < c, \\ &\pi_i \geq 0, \forall i \in [K + 1]. \end{aligned} \quad (1)$$

where the vector π represents the policy which defines the fraction of time an arm will be pulled. In any policy, there will be at most two arms that have nonzero π_i values since there are two constraints in the problem. We refer to a set consisting of at most two arms as a base. We denote the set of all possible valid bases (where valid means that the average cost less than or equal to c can be reached through a mixture of arms in the base) as \mathbb{V} . Note that for simplicity, we assume that the arms in a base are ordered so that the higher cost arm appears first. We let $\mathcal{I}^i := \{\mathcal{I} \in \mathbb{V} : i \in \mathcal{I}\}$ denote the set of valid bases that include the arm i . We let π^* denote the optimal solution to (1), and let $r^* := \mu^T \pi^*$ be the optimal reward per round. We also define $\mathcal{I}^* := \{i : \pi_i^* > 0\}$ as the optimal base. The optimal solution of this problem can be divided into three cases. First, if the arm with the highest mean reward has cost less than c , i.e. $\rho_{i^{**}} \leq c$; then the optimal base consists of only this arm; hence $\mathcal{I}^* = \{i^{**}\}$, and $\pi_{i^{**}}^* = 1$. In the second case, if $\rho_{i^{**}} > c$, $\rho_{i^*} > c$, then $\mathcal{I}^* = \{i^*, K + 1\}$, and the optimal solution is $\pi_{i^*}^* = c/\rho_{i^*}$, and $\pi_{K+1}^* = 1 - \pi_{i^*}^*$. In third case, if $\rho_{i^{**}} > c$, $\rho_{i^*} < c$, then optimal base includes two arms which might or might not include i^* or i^{**} .

Let OPT denote the total expected reward of a dynamic policy in T rounds that conforms to a total budget constraint in a total of T rounds as in the standard BwK literature instead of the anytime budget constraint we consider here. It was shown that $OPT_{LP} \geq OPT$ (Badanidiyuru et al., 2013).

Let REF be denote the total expected reward of a dynamic policy in T rounds that conforms to the average cost constraint. This constraint is stricter than the total budget constraint. This can easily be seen as satisfying the anytime constraint in the last round T with $c = B/T$ produces the total budget constraint of B in T rounds. Hence, it holds that $OPT_{LP} \geq OPT \geq REF$. While regret could be defined as the difference expected cumulative reward of SUAK and REF , we choose a stronger regret definition by defining it with respect to OPT_{LP} as $R_T = OPT_{LP} - \mathbb{E}[F(T)] = Tr^* - \mathbb{E}[F(T)]$ so that our results can be compared with prior work on the total budget setting.

3 THE SUAK ALGORITHM

3.1 THE NAIVE APPROACH

Before presenting the SUAK Algorithm, to demonstrate the additional complexities of our problem formulation over the standard BwK setting, and also to serve as a baseline, we present a naive approach which makes it possible to convert any BwK algorithm to our BwAK setting. In this trivial approach, in a given round t , we first simply check if it is possible to violate the anytime constraint, and skip the round if it is the case. Otherwise, we let the BwK algorithm pull an arm. To demonstrate this more concretely, we use the *One Phase Algorithm* in Li et al. (2021), and add skipping behaviour such that a round is t skipped if $S_c(t-1) + 1 > c \cdot t$. The implementation with this skipping rule, which we call as the *One Phase Skip (OPS) Algorithm*, is given in Algorithm 1.

In this algorithm, the initialization phase consists of sampling each arm once while using skips to prevent violation of the constraint. After this phase, we utilize a skipping mechanism in lines 5 - 6, and if the round is not skipped, the algorithm proceeds to solving the linear programming problem in line 8. In this LP, $\mu^U(t)$ is the UCB of arm reward at round t , $\rho^L(t)$ is the LCB of arm cost at round t , and $B_r(t) = cT - S_c(t-1)$ is the total remaining budget in round t . Since UCB values are used, the solution of LP gives the optimistically best policy according to the UCB principle. This policy is normalized to a probability distribution, and the arm is selected using this probability.

To show that this naive approach might suffer a large regret due to large number of skips, we run simulations on the following problem instance with $K + 1 = 4$ arms where $\mu = [0.45, 0.7, 0.8, 0]$; and $\rho = [0.25, 0.75, 0.8, 0]$. Except for the null arm, the arm reward and cost values are independently sampled from a Beta distribution with parameters $\alpha = \mu * 10, \beta = (1 - \mu) * 10$. The average cost budget per round is $c = 0.5$. For SUAK, we take $\omega = 0.143$. We average results from 20 sim-

Algorithm 1 The Naive Approach: One Phase Skip Algorithm

```

1: Input: Average cost target  $c$ , number of rounds  $T$ 
2: Initialize: Sample each arm once while skipping accordingly so that  $\forall t \leq t_{\text{init}}, S_c(t-1) + 1 \leq c \cdot t$ 
3: Set  $t = 1$ 
4: for each round  $t > t_{\text{init}}$  do
5:   if  $S_c(t-1) + 1 > c \cdot t$  then
6:     Skip the round
7:   else
8:     Solve the following LP:
            $\tilde{\pi} = \arg \max_{\pi} \langle \mu^U(t-1), \pi \rangle$ 
           s.t.  $\langle \rho^L(t-1), \pi \rangle \leq B_r(t)$ 
            $\pi \geq \mathbf{0}$ 
9:     Normalize  $\tilde{\pi}$  into a probability and randomly play an arm from this probability
10:  end if
11:  Update  $\rho^L(t), \mu^U(t)$ , and  $B(t)$ 
12:  Update  $t = t + 1$ 
13: end for

```

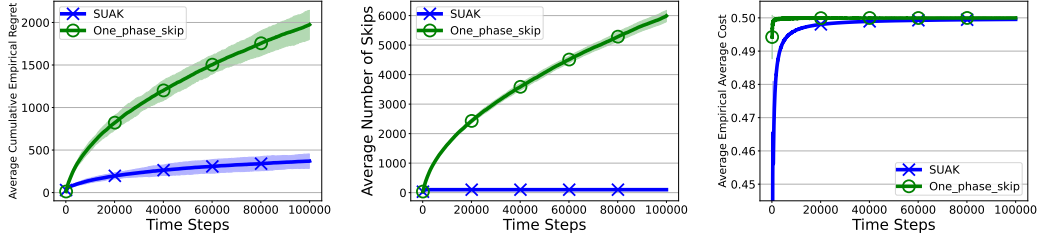


Figure 1: The plots of cumulative empirical regret (Left), number of skips (Middle), and cumulative number of skips (Right); averaged over 20 different trials.

ulation runs, each with 100,000 rounds. We compare the results with results from SUAK which we propose in §3.2. The simulation results are given in Figure 1. The shaded areas in the plots represent error bars with one standard deviation. It can be seen that the number of skips of One Phase Skip algorithm exhibits sublinear growth and it is much higher than the number of skips of SUAK, which results in higher regret compared to SUAK. Hence, this demonstrates that merely adding skips to a BwK algorithm and treating it as a BwAK algorithm is not a viable solution. In view of this, in the next section, we present SUAK, an algorithm that strategically under-utilizes the available budget to reduce the number of skips needed, and achieve smaller regret.

3.2 THE SUAK ALGORITHM

We propose an algorithm called *Strategic Under-utilization for Anytime Knapsacks* (SUAK) that utilizes upper and lower confidence bounds of the arm rewards and costs using the UCB principle to upper bound the reward that can be obtained from a particular base. SUAK also uses skipping a round to satisfy the anytime cost constraint, and targets an average cost of $c - \log t / (\omega^2 t)$ to limit the number of skips, where ω is defined in Assumption 2. The pseudo-code is provided in Algorithm 2.

The algorithm works as follows. First, SUAK is initialized by sampling each arm once. To prevent SUAK from exceeding the targeted average of $c - \log t / (\omega^2 t)$ during initialization, skipping is employed in rounds $t \leq t_{\text{init}}$ whenever $S_c(t-1) + 1 > c \cdot t - \log(t) / \omega^2$. We define t_{init} as the round where initialization ends (when a sample is obtained from each arm). After this initialization step, the algorithm works as follows. In every round, first the anytime budget constraint is checked. The round is skipped (null arm is pulled) and the algorithm proceeds to the next round if $S_c(t-1) + 1 > c \cdot t$, i.e. if pulling an arm at that round can violate the constraint. Secondly, if there is uncertainty on whether the mean cost of an arm is less than or greater than c , i.e. $\exists l : \rho_l^L(t) \leq c \leq \rho_l^U(t)$ where

$$\rho_l^L(t) := \bar{\rho}_l(t) - 7\sqrt{1.5 \log t / N_l(t)}, \text{ and } \rho_l^U(t) := \bar{\rho}_l(t) + 7\sqrt{1.5 \log t / N_l(t)};$$

then that arm is pulled and SUAK proceeds to the next round. This step is needed for the anytime cost constraint; it ensures whether the mean cost of arm is above or below c is correctly known, which in turn ensures that the base that SUAK selects for that round includes an arm with cost less than c . This step is also needed to establish tighter bounds on the number of times a suboptimal base is selected in the theoretical analysis. To prevent this step from using more than c cost budget per round on average, we define $S_p(t)$ as the sum of all the cost incurred from this step until round $t+1$, and define $N_p(t)$ as the total number of arm pulls due to this step until round $t+1$. The round

Algorithm 2 SUAK: Strategic Under-utilization for Anytime Knapsacks

```

1: Input: Average cost target  $c$ 
2: Initialize: Sample each arm once
   while skipping accordingly so that  $\forall t \leq t_{\text{init}}, S_c(t-1) + 1 \leq c \cdot t - \log t / \omega^2$ 
3: for each round  $t > t_{\text{init}}$  do
4:   if  $S_c(t-1) + 1 > c \cdot t$  then
5:     Skip round  $t$ 
6:   continue
7:   end if
8:   if  $S_p(t-1) + 1 > c \cdot N_p(t-1)$  then
9:     Skip round  $t$ 
10:    continue
11:   end if
12:   if  $\exists l : \varrho_l^L(t) \leq c \leq \varrho_l^U(t)$  then
13:     Pull arm  $i(t) = l$ , observe  $r_{i(t)}(t)$ ,
        $\rho_{i(t)}(t)$ 
14:      $S_p(t) = S_p(t-1) + \rho_{i(t)}(t)$ 
15:      $N_p(t) = N_p(t-1) + 1$ 
16:     continue
17:   end if
18:    $\mathcal{S}_t = \{\pi : \pi \in \Delta_{K+1}, \langle \pi, \boldsymbol{\rho}^L(t-1) \rangle \leq c\}$ 
19:    $\pi(t) = \arg \max_{\pi \in \mathcal{S}_t} \langle \boldsymbol{\mu}^U(t-1), \pi \rangle$ 
20:    $\mathcal{I}_t = \{i : \pi_i(t) > 0\}$ 
21:   if  $|\mathcal{I}_t| = 1$ , i.e.  $\mathcal{I}_t = \{j(t)\}$  then
22:     Pull arm  $j(t)$ 
23:     continue
24:   end if
25:    $j(t), k(t) = j(t), k(t) \in \mathcal{I}_t : \bar{\rho}_{j(t)}(t) > \bar{\rho}_{k(t)}(t)$ 
26:    $b(t) = c \cdot t - S_c(t-1) - \log t / \omega^2$ 
27:   if  $b(t) > \bar{\rho}_{j(t)}(t)$  then
28:      $p(t) = 1 - \omega$ 
29:   else if  $b(t) < \bar{\rho}_{k(t)}(t)$  then
30:      $p(t) = \omega$ 
31:   else
32:      $p_1(t) = \max \left( \frac{b(t) - \bar{\rho}_{k(t)}(t)}{\bar{\rho}_{j(t)}(t) - \bar{\rho}_{k(t)}(t)}, \omega \right)$ 
33:      $p(t) = \min(p_1(t), 1 - \omega)$ 
34:      $i(t) = \begin{cases} j(t) & \text{with probability } p(t), \\ k(t) & \text{otherwise} \end{cases}$ 
35:   end if
36:   Pull arm  $i(t)$ , observe  $r_{i(t)}(t)$ ,  $\rho_{i(t)}(t)$ 
37:   Update  $\boldsymbol{\rho}^L(t)$  and  $\boldsymbol{\mu}^U(t)$ 
38: end for

```

t is skipped if $S_p(t-1) + 1 > c \cdot N_p(t-1)$. The main objective of this skipping mechanism is to decouple the skips needed to satisfy the anytime cost constraint due to regular arm pulls and the skips needed to satisfy the constraint from this step for ease of theoretical analysis; in practice this skipping mechanism can be ignored. Since the expected number of arm pulls from this step is upper bounded by $O(\log T)$, the skips due to this step will also be $O(\log T)$. After this step, the constraint set is constructed as $\mathcal{S}_t = \{\pi : \pi \in \Delta_{K+1}, \langle \pi, \boldsymbol{\rho}^L(t-1) \rangle \leq c\}$; where

$$\begin{aligned} \mu_i^L(t) &:= \text{proj}_{[0,1]}(\bar{\mu}_i(t) - \epsilon_i(t)), & \mu_i^U(t) &:= \text{proj}_{[0,1]}(\bar{\mu}_i(t) + \epsilon_i(t)), \\ \rho_i^L(t) &:= \text{proj}_{[0,1]}(\bar{\rho}_i(t) - \epsilon_i(t)), & \rho_i^U(t) &:= \text{proj}_{[0,1]}(\bar{\rho}_i(t) + \epsilon_i(t)), \end{aligned}$$

are the UCB and LCB values of arm costs and rewards; and $\epsilon_i(t) = \sqrt{3 \log T / N_i(t)}$ is the confidence interval. Hence, the constraint set \mathcal{S}_t includes all policies that have an average cost less than c using the optimistic estimates of arm costs (LCB values of arm costs). The empirically best policy at round t is then found using a linear program (LP) as $\pi(t) = \arg \max_{\pi \in \mathcal{S}_t} \langle \boldsymbol{\mu}^U(t-1), \pi \rangle$. Note that using the UCB of the empirical arm reward along with the LCB of empirical arm costs in \mathcal{S}_t produces an upper confidence bound on the reward of a base. The arms that have nonzero $\pi_i(t)$ values are selected as the empirically optimal base arm set for that round, denoted as \mathcal{I}_t .

If \mathcal{I}_t consists of a single arm, that arm is pulled. Otherwise, $\mathcal{I}_t = \{j(t), k(t)\}$ will consist of two arms; we denote them as $j(t)$, and $k(t)$; where wlog we assume $j(t)$ is the arm with mean cost above c . The available budget at that round with respect to the targeted average cost is found as $b(t) = c \cdot t - S_c(t-1) - \log t / \omega^2$. If the available budget $b(t)$ is greater than $\bar{\rho}_{j(t)}(t)$, arm $j(t)$ is pulled with probability $1 - \omega$, and arm $k(t)$ is pulled otherwise. If $b(t)$ is less than $\bar{\rho}_{k(t)}(t)$, arm $j(t)$ is pulled with probability ω , and arm $k(t)$ is pulled otherwise. If $\bar{\rho}_{k(t)}(t) \leq b(t) \leq \bar{\rho}_{j(t)}(t)$, arm $j(t)$ is pulled with probability $p(t) = \frac{b(t) - \bar{\rho}_{k(t)}(t)}{\bar{\rho}_{j(t)}(t) - \bar{\rho}_{k(t)}(t)}$ clipped at ω from below and $1 - \omega$ from above; and arm $k(t)$ is pulled otherwise. With this design, each arm in a base is pulled with at least ω probability to help explore all arms in a base.

Note that this algorithm is non-stationary as it is adaptive to the available budget at that round. This design is essential as it was shown in Flajolet & Jaillet (2015, Lemma 2) that a non-adaptive design suffers regret of order $\Omega(\sqrt{T})$ even if the optimal solution π^* is known unless all arms consume the same deterministic amount of resources at every round. The main intuition behind this result is that the fluctuation of the available budget around its mean at a round t can be as high as $\Omega(1/\sqrt{t})$.

3.3 ANALYSIS OF SUAK

We now characterize the performance of the SUAK by providing the theoretical upper bound on the expected cumulative regret. We first provide definitions of arm gaps and state a set of mild assumptions that are required for the theoretical analysis. We refer the readers to the Appendix for detailed proofs of the results presented in this section.

Definition 3.1. The reward gap of an arm is defined as $\Delta_i = \mu_{i^*} - \mu_i$.

Definition 3.2. The gap of a base \mathcal{I} is defined as $\Delta_{\mathcal{I}} = r^* - r_{\mathcal{I}}$, where $r_{\mathcal{I}}$ is the reward value of the solution of (1) when only arms in \mathcal{I} are allowed. We also define $\Delta_{\min,i} := \min_{\mathcal{I} \in \mathcal{I}^i \setminus \mathcal{I}^*} \Delta_{\mathcal{I}}$ as the minimum reward gap of bases that include the arm i .

Definition 3.3. The cost gap of an arm is defined as $\delta_i = |\rho_i - c|$.

Assumption 1. We define $\delta_{\min} = \min_{i \in [K]} \delta_i$ as the minimum cost gap, and assume $\delta_{\min} > 0$.

Note that regret depends on the cost gap δ_i since the algorithm needs to be able to correctly identify if the true mean cost of an arm is above or below c . This is in turn needed for the adaptive design since if the empirical average cost is above the targeted cost and if an arm with cost more than c is identified as an arm with cost less than c , pulling that arm might lead to over-consuming the targeted budget. Since regret depends on δ_{\min} , $\delta_{\min} > 0$ is needed so that the regret bound is not unbounded.

Assumption 2. We assume that we are given an $\omega > 0$ such that $\omega \leq \delta_{\min}/(2 + \delta_{\min} - c)$.

Note that for any δ_{\min} or c value, $\delta_{\min}/(2 + \delta_{\min} - c) \geq \delta_{\min}/3$. Assumption 2 is necessary for the adaptive design of the algorithm in meeting the anytime cost constraint, as we use a cost budget under-utilization of $\log t/\omega^2$ at round t in SUAK to be able to achieve theoretical guarantees. Also, in SUAK, we set the minimum fraction of time an arm in a base will be pulled to ω . With this use, ω can be understood as the minimum triggering probability (p^*), in the probabilistic triggering literature discussed in §3.4. Since the fraction of pulls of a particular arm in a given base can be as low as ω , our regret bounds depend on ω as in the worst case a base needs to be selected $1/\omega$ times in expectation to acquire one sample of each arm in the base.

Under these assumptions stated above, we obtain the following upper bound on expected regret.

Theorem 3.1 (Upper Bound on Expected Regret). *Under Assumption 1 and 2; when SUAK is run with a given average cost budget $0 < c \leq 1$, its cumulative expected regret is upper bounded as*

$$R_T \leq \sum_{i=1}^K \frac{96r^*(\frac{\delta_i+1}{\delta_i})^2 \log T}{\omega \Delta_{\min,i}^2} + \frac{202Kr^* \log T}{c\omega^2} + \frac{3\pi^2 r^*}{\delta_{\min}^2} + R_K + r^* t_{in} = O(K \log T) + O(1) \quad (2)$$

where $R_K = 5\pi^2 K^2/3$, and $t_{in} = -W(-\omega^2 c e^{-\omega^2 K})/(\omega^2 c) = O(1)$ is the upper bound on the number of rounds needed for the initialization phase of SUAK; and $W(\cdot)$ is the Lambert function. Also recall that $\Delta_{\min,i}$ is the minimum reward gap among the bases that include the arm i , δ_i is the cost gap of an arm, r^* is the optimal reward, and ω is defined in Assumption 2.

Note that the first term in (2) is related to regret from arm pulls due to selecting a suboptimal base in the round; the second term is related to arm pulls that are used to learn whether the true mean cost of an arm is greater than or less than c , and also the regret resulting from under-utilizing the budget; the third term is due to expected number of times the anytime constraint may be violated; the fourth term is due to suboptimal arm pulls that occur if the confidence bounds do not hold; and the last term is regret from the initialization phase. The proof of Theorem 3.1 is given in §D, and we also provide a brief proof sketch in §3.5.

Note that this problem-dependent upper bound order-wise matches the $O(K \log T)$ problem-dependent bound of prior work for the regular BwK setting. However, SUAK is not optimal for a problem-independent bound since regret can be large when the value of ω is small, i.e. if $\omega \leq 1/\log(T)$ assuming the time horizon T is known. For this case, prior work such as Bernasconi et al. (2024b) can be used to achieve $O(\sqrt{KT})$ problem-independent regret in our setting.

3.4 RELATED WORKS

In this section, we provide some of the works that are related to our problem setting. We provide additional related works in Appendix §B.

Bandits with Knapsacks: The BwK problem has been studied before and algorithms that achieve optimal problem-independent regret bounds on the order of $O(\sqrt{K} \cdot OPT)$ have already been developed (Badanidiyuru et al., 2013; Agrawal & Devanur, 2014). However, deriving a problem-dependent lower bound and developing algorithms that achieve this bound are still open questions. One prior work in this regard is by Sankararaman & Slivkins (2021), in which the BwK problem is studied under only one constraint; a one-dimensional cost, and no constraint on time. In this simple setting, there is a single *unique* optimal arm. They propose an algorithm that achieves a regret bound of $O(KG_{LAG}^{-1} \log T)$, where G_{LAG} is defined as the Lagrangian gap of an arm. Compared to this work, our setting is more complex as we have two constraints.

Another notable work in this field is by Flajolet & Jaillet (2015), in which the BwK problem is considered under three different cases of 1, 2; and d constraints. For 2 constraints, which represents a constraint on the total number of rounds and a constraint on total budget where costs of arms are one-dimensional, a regret bound of $O(\lambda^2 K^2 \log T / (\delta_{\min}^3 \Delta) + K^2 \sigma \log T / \delta_{\min}^3)$ is achieved with additional problem-dependent constants where K is the number of arms; δ_{\min} is the minimum distance between the mean cost of arms and the average budget b ; σ is the minimum $1/\mu$ value; $\lambda = 1 + 2\kappa$; and κ is a constant assumed to be known *a priori* such that $|\mu_i - \mu_j| \leq \kappa |\rho_i - \rho_j|$ for any i, j . For the d constraint setting, a regret bound of $O(2^{K+d} \log T)$ is achieved. The 2-constraint setting is similar to our work, as we also have one-dimensional costs, and our anytime cost constraint can be viewed as a total budget constraint that needs to be satisfied in all rounds. It can be seen that the $O(K^2 \log T)$ regret bound in this work is not optimal for its dependence on K^2 . In our work, we reduce this dependence on K^2 to K while considering the more complex BwAK problem. However, our work has an increased dependence on the gap with $1/\Delta^2$ compared to the $1/\Delta$ dependence here.

Another notable prior work is by Li et al. (2021), where a d -dimensional cost vector is considered, with one of the dimensions of the cost vector being time. The optimal solution in this d -dimensional setting can be a base consisting of at most d different arms. They propose a two-phase algorithm where the first phase of the algorithm pulls each arm the same number of times until the suboptimal arms are eliminated. In the second phase, the base with highest upper confidence bound is chosen. This two-phase approach greatly simplifies the theoretical analysis as the number of pulls of each individual arms is the same in the first phase. With this approach, they achieve a regret bound of $O(Kd \log T / (b^3 \Delta^2) + d^4 / (b^2 \min\{\chi^2, \Delta^2\} \min\{1, \sigma^2\}))$, where Δ in their setting is defined as the gap between the reward of the optimal solution per round and the maximum reward that can be obtained per round when one arm (except the null arm) is removed; b is the average cost budget per round; χ is the minimum nonzero value in the optimal policy; σ is a problem dependent constant related to the linear dependency between arms across different constraints. In our work while we have similar dependence on Δ and K , we have $1/\omega\delta^2$ additional dependence on cost gaps of arms. However, our setting (BwAK) is more complicated and we conjecture that these additional terms ω and δ^2 are needed to satisfy the anytime constraint. The BwK setting has also been studied under different problem settings, such as in the adversarial setting (Immorlica et al., 2022), in contextual bandits (Agrawal & Devanur, 2016), under nonstationary distributions (Liu et al., 2022), and in combinatorial bandits (Sankararaman & Slivkins, 2018).

Bandits with Replenishable Knapsacks: In this setting, cost of an arm is allowed to be negative, which allows the knapsack to be replenished. One notable prior work is by Slivkins et al. (2024), where the contextual bandits with linear constraints (CBwLC), a more generalized version of the contextual bandits with knapsacks (CBwK) problem, which allows packing and covering constraints, as well as positive and negative resource consumption, is considered. Their algorithm also works when the initial budget is $B = \Omega(T)$, or $B = o(T)$, compared to the prior work which mostly restricts the initial budget to $B = \Omega(T)$. This is similar to our setting since our problem setting can be reduced to their problem setting by implementing the budget increase as subtracting c from the costs of all arms (this also makes the skip arm in our setting have negative cost c and function as the resource replenishing arm). However, their algorithm is suboptimal in our problem setting with a zero initial budget, as they remark in the discussion of (Slivkins et al., 2024, Theorem 3.6), their proposed algorithm LagrangeCBwLC achieves optimal $O(\sqrt{KT})$ regret when the initial budget $B > \Omega(T)$; and when $B = o(T)$ its regret is suboptimal. This is as expected since Lagrange-based algorithms generally require knowing the ratio T/B , which goes to infinity when $B = o(T)$. In our work, we are only interested in the case where the initial budget is zero, and we consider gap-dependent results instead of the gap-independent results considered here, and we propose an algorithm that achieves an order-optimal $O(K \log T)$ gap-dependent regret bound.

In Bernasconi et al. (2024a), a more general BwK formulation with long-term constraints is considered where the costs can be negative as well as positive. The long-term constraint is defined such that the total consumption of each resource at round T should be less than zero up to small sublinear violations. This is again similar to this setting if we do not allow any violation of the constraint as our problem setting can be reduced to this setting again by subtracting c from the costs of all arms. They show that regret is upper bounded by $O(\sqrt{KT} \log(KT))$ when the EXP3-SIX algorithm is used with the Primal-Dual algorithm based framework that they propose for their problem setting. They also remark that initial $o(T)$ rounds can be skipped to cover the potential violations and implement the long-term constraint as a hard constraint like in our setting. However, they provide an upper bound of $O(\sqrt{KT})$ constraint violations in (Bernasconi et al., 2024a, Corollary 8.2), which suggests that the initial $O(\sqrt{KT})$ rounds would need to be skipped to achieve hard constraints, which would lead to $O(\sqrt{KT})$ gap-independent regret in our problem setting. In our work, we show that we achieve $O(K \log T)$ gap-dependent regret for the same problem setting. In Bernasconi et al. (2024c), they consider the same problem setting as in their prior work Bernasconi et al. (2024a). Instead of a Prior-Dual algorithm based approach, they use a UCB-based approach to optimistically estimate the constraints through a weighted empirical mean of past samples. This approach lets them provide $O(\sqrt{T})$ regret in stochastic settings without assuming Slater’s condition. The upper bound on constraint violations is still $O(\sqrt{KT})$, which would again lead to a $O(\sqrt{KT})$ gap-independent regret in our problem setting.

In Bernasconi et al. (2024b), there exists an arm with a negative expected cost that allows to replenish the budget. This is very similar to our setting as our case can be considered a special case of this setting that starts with zero budget. However, their work cannot be used in our setting as they assume $B = \Omega(T)$ such that $B = T\rho$, and they use the parameter ρ in the Lagrangian function of the Primal-Dual algorithm template that they provide. Further, they only consider instance-independent bounds of $O(\sqrt{KT})$, and do not consider the $O(K \log T)$ instance-dependent bounds that we consider here.

BwK with non-monotonic resource utilization: It is a generalization of the BwK problem where in each round, a vector of resource drifts that can be positive, negative, or zero is observed along with the reward; and the budget of each resource is incremented by this drift amount (Kumar & Kleinberg, 2022). In Kumar & Kleinberg (2022), a three phase algorithm that combines the ideas in Flajolet & Jaillet (2015) and Li et al. (2021) is provided. The algorithm uses the phase one of Li et al. (2021) to identify the optimal arms, then in phase two arms are pulled to shrink the confidence intervals further, and in the third phase, the optimal arms are exploited by sampling from a perturbed distribution to ensure that the budget of each resource stays close to a decreasing sequence of thresholds. While the idea of decreasing sequence of thresholds can be seen as similar to under-budgeting in our algorithm, their problem setting assumes time horizon T to be known, and threshold decays to zero over time as uncertainty decreases; however, in our setting we do not assume knowing T , and we incur regret from under-budgeting as we always under-budget. Their algorithm achieves $O(Km^2 \log T / (\Delta^2 \cdot \min\{\delta_{\text{drift}}^2, \sigma_{\min}^2\}))$, where K is the number of arms, m is the dimension of the cost vector, $\delta_{\text{drift}} > 0$ is the smallest magnitude of the drifts, and σ_{\min} is the smallest singular value of the constraint matrix.

Comparison of our work with the prior work is summarized in Table 1. Due to different gap definitions, and different problem-dependent parameters used, we would like to note that these results are not directly comparable. Also note that a problem-dependent lower bound does not exist for the BwK problem or our BwAK problem. We remark that deriving a lower bound for BwK or BwAK would be an important future work; yet it would be challenging due to the variety of problem-dependent parameters that can be used to define the problem instance.

3.5 PROOF SKETCH

We now present a brief outline of the regret analysis of SUAK, which is provided in §D. In the proof, the regret is first decomposed as follows.

$$R_T \leq R_a(T) + R_b(T) + R_c(T) + R_d(T) + \sum_{t=t_{\text{init}}+1}^T (\mathbb{P}(\mathcal{G}^c(t)) + \mathbb{P}(\mathcal{F}^c(t))) + r^* t_{\text{init}}$$

where $R_a(T)$ is the regret from skips that are used to satisfy the anytime constraint in line 5 of Algorithm 2, and $R_b(T)$ is the regret from skips that are needed while reducing the confidence

Table 1: Comparison of our work with prior work on bandits with knapsacks

Work	Model	Regret Bound
(Li et al., 2021)	Total budget	$O\left(\frac{Kd \log T}{b^3 \Delta^2} + \frac{d^4}{b^2 \min\{\chi^2, \Delta^2\} \min\{1, \sigma^2\}}\right)$
Flajolet & Jaillet (2015)	Total budget	$O\left(\frac{\lambda^2 K^2 \log T}{\delta_{\min}^3 \Delta} + \frac{K^2 \sigma \log T}{\delta_{\min}^3}\right)$
Kumar & Kleinberg (2022)	Total budget and drift	$O\left(\frac{K m^2 \log T}{\Delta^2 \cdot \min\{\delta_{\text{drift}}^2, \sigma_{\min}^2\}}\right)$
Our work	Average budget	$O\left(\frac{K \log T}{\omega \delta^2 \Delta^2} + \frac{K \log T}{w^2}\right) + O(1)$

intervals of the arm costs due to the condition in line 9 of Algorithm 2. $R_c(T)$ is due to pulls needed while reducing the confidence intervals of the arm costs (pulls from line 13 of Algorithm 2), and $R_d(T)$ is due to pulls of arms after selecting a base (pulls from line 36 of Algorithm 2); which includes the selection of suboptimal bases and regret from under-utilization of the cost budget. The terms $\sum_{t=t_{\text{init}}+1}^T (\mathbb{P}(\mathcal{G}^c(t)) + \mathbb{P}(\mathcal{F}^c(t)))$ are due to the probability of confidence bounds not holding, and can be upper bounded as $5\pi^2 K^2/3$. The $r^* t_{\text{init}}$ term is regret due to the initialization phase, we upper bound t_{init} by considering K cost can be incurred in the initialization step in the worst case, and also noticing that the regret per round is upper bounded by r^* .

To upper bound $R_a(T)$, we define t_e as the time the algorithm exceeds the targeted cost of $c - \log t/t$, and we define $t_f + 1$ as the time instant where the algorithm skips. Due to the design of the algorithm, the arm with the lower cost will be pulled with probability $1 - \omega$ between rounds $t_e \leq t \leq t_f$, and the total incurred cost between rounds $t_e \leq t \leq t_f$ needs to exceed the c by at least $\log(t_e)$. We upper bound the probability of this event using standard concentration bounds, and apply a union bound over all possible t_e and t_f values to establish that $R_a(T) \leq 3\pi^2 r^* / (\delta_{\min}^2)$. We upper bound $R_b(T)$ and $R_c(T)$ as follows. Using standard techniques in bandit literature, we show that an arm i will be sampled at most $96 \log T / \delta_i^2$ times to reduce the uncertainty in its cost estimate, and the expected regret per round will be $r^* - \mu_i$. Note that μ_i can be greater than r^* for some arms, but this is balanced by skips. For arms with cost larger than c , we derive $104 \log T / (c \delta_i)$ skips are needed.

We upper bound $R_d(T)$ as follows. For selections of a suboptimal base, using the fact that arms are sufficiently sampled by line 13 of Algorithm 2, we show that a suboptimal base $\mathcal{I} = (i, j)$ can be selected at most $\sum_{i=1}^K 48(\frac{\delta_i+1}{\delta_i})^2 \log T / (\Delta_{i,j}^2)$ times if it is assumed that selection of the base yields a sample of both arms in it. Due to partial observability, it will take $1/\omega$ rounds in expectation to obtain one sample for both arms. Taking this into account, we show that at most $\sum_{i=1}^K 48(\frac{\delta_i+1}{\delta_i})^2 \log T / (\omega \Delta_{\min,i}^2)$ pulls of arm i will occur to satisfy the upper bound on the number of pulls of all bases that include arm i . Using the technique in (Kveton et al., 2015), we derive the worst case regret from this upper bound on the samples of arms. We also upper bound regret from cost under-utilization as $2r^* \log T / (c\omega^2)$ using r^*/c , the optimal reward per cost.

4 SIMULATIONS

We now evaluate the performance of the proposed SUAK Algorithm through simulations. For comparison, we have included the *Primal Dual* and *One Phase* Algorithms in Li et al. (2021); and the *UCB Simplex* Algorithm in Flajolet & Jaillet (2015). We would like to note that while the authors in Li et al. (2021) believe that the *One Phase* Algorithm would be optimal, they leave providing theoretical regret bounds for that algorithm as an open question claiming it would be challenging to do so. Instead, they provide theoretical guarantees for the *Primal Dual* Algorithm, which similar yet expected to have much worse empirical performance compared to the *One Phase* Algorithm. We implement the skip versions of these algorithms as described in §3.1, and we refer them by appending ‘_skip’ after their names.

We perform simulations on the following setting with $K + 1 = 11$ arms where the mean reward and cost vectors are $\mu = [0.2, 0.25, 0.45, 0.4, 0.7, 0.75, 0.8, 0.9, 0.8, 0.7, 0]$; and $\rho = [0.2, 0.25, 0.3, 0.4, 0.6, 0.65, 0.7, 0.75, 0.8, 0.9, 0]$. Except for the null arm, the arm reward and cost values are independently sampled from a Beta distribution with parameters $\alpha = \mu * 10, \beta =$

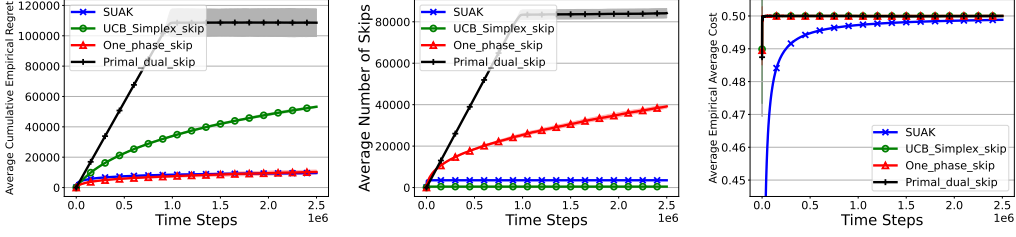


Figure 2: The plots of cumulative empirical regret (Left), number of skips (Middle), and cumulative number of skips (Right); averaged over 20 different trials.

$(1 - \mu) * 10$, where μ represents the mean of the distribution. The average cost budget per round is $c = 0.5$. We perform the simulations for 2.5 million rounds, and average over 20 different trials. The simulation results for this setting are given in Figure 2. The shaded areas in the plots represent error bars with one standard deviation. We use $\omega = 0.0625$ for SUAK.

It can be seen from the top left plot in Figure 2 that the Primal Dual Skip Algorithm performs the worst. This is as expected since the Primal Dual Algorithm is designed for theoretical performance, and pulls every arm the same number of times until finding the optimal solution. The UCB Simplex Skip also performs poorly in simulations. This is since the algorithm assumes knowledge of a constant κ *a priori* such that $|\mu_i - \mu_j| \leq \kappa|\rho_i - \rho_j|$ for any i, j ; and the confidence intervals for the arm rewards are multiplied by a factor of $\lambda = 1 + 2\kappa$. In the simulation setting, $\lambda = 9$; which increases the number of samples needed for exploration.

It can be seen that SUAK performs better compared to the Primal Dual or UCB Simplex algorithms, and also exceeds the performance of One Phase Skip (OPS) after around round 1.7×10^6 . This is since the regret of SUAK concentrates primarily on the initial rounds, which is due to two factors. Skips needed for the under-utilization of the budget by $\log t / \omega^2$, and also pulls from the line 13 of Algorithm 2 mostly concentrate on the initial rounds. After these initial rounds, SUAK can catch up to and eventually surpass the performance of OPS due to higher regret OPS experiences from its high number of skips, verifying the practical utility of SUAK.

In terms of the number of skips, it can be seen from the bottom middle plot in Figure 2 that the number of skips is sublinear for all algorithms, and SUAK has the least number of skips. This demonstrates the effectiveness of SUAK in reducing the number of skips by under-utilizing the available budget. Note that pulls of the null arm originating from the condition in line 5 or line 9 of Algorithm 2 are counted as skips, yet pulls of the null arm when it is in the selected base is not counted as a skip. The plot on the right of Figure 2 shows the incurred average cost. As expected, SUAK starts from a smaller average cost value and approaches the per round cost budget of 0.5 over time due to under-utilization of the budget, and other algorithms are very close to the constraint, and except the UCB Simplex Skip, need to utilize skips to avoid exceeding the constraint.

5 CONCLUDING REMARKS

In this paper, we introduce a previously unexplored setting for the BwK problem, which we call the bandits with anytime knapsacks (BwAK) problem; where we employ a stricter anytime cost constraint instead of a total cost budget. We provide SUAK, a novel algorithm that under-utilizes the available cost budget and uses skipping to limit the probability of violating the anytime cost constraint, and also uses upper confidence bounds to balance exploration and exploitation. SUAK achieves a regret upper bound of $O(K \log T)$ compared to the optimal solution of the linear relaxation version of the problem which does not necessarily obey the anytime cost constraint. This bound is better than the regret upper bound of prior work for the BwK setting on problem-dependent terms in a wide range of problem instances. We provide simulation results to demonstrate the empirical performance of SUAK. Our work opens multiple directions for future research. One interesting future direction is to extend our bandit results to the case where the cost of an arm is a d -dimensional vector. This is a challenging problem as the anytime cost constraint needs to be satisfied in every dimension, which can introduce additional skips, and hence additional regret. Another interesting open direction is the case where the distribution of rewards and costs of arms are stochastic. In this case, satisfying the anytime cost constraint would again be challenging but we conjecture it may be accomplished using a more conservative cost budget under-utilization.

REFERENCES

- Shipra Agrawal and Nikhil Devanur. Linear contextual bandits with knapsacks. *Advances in Neural Information Processing Systems*, 29, 2016.
- Shipra Agrawal and Nikhil R Devanur. Bandits with concave rewards and convex knapsacks. In *Proceedings of the fifteenth ACM conference on Economics and computation*, pp. 989–1006, 2014.
- Vashist Avadhanula, Riccardo Colini Baldeschi, Stefano Leonardi, Karthik Abinav Sankararaman, and Okke Schrijvers. Stochastic bandits for multi-platform budget optimization in online advertising. In *Proceedings of the Web Conference 2021*, pp. 2805–2817, 2021.
- Ashwinkumar Badanidiyuru, Robert D. Kleinberg, and Aleksandrs Slivkins. Bandits with knapsacks. *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pp. 207–216, 2013.
- Ashwinkumar Badanidiyuru, Robert Kleinberg, and Aleksandrs Slivkins. Bandits with knapsacks. *Journal of the ACM (JACM)*, 65(3):1–55, 2018.
- Martino Bernasconi, Matteo Castiglioni, and Andrea Celli. No-regret is not enough! bandits with general constraints through adaptive regret minimization. *arXiv preprint arXiv:2405.06575*, 2024a.
- Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. Bandits with replenishable knapsacks: the best of both worlds. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Martino Bernasconi, Matteo Castiglioni, Andrea Celli, and Federico Fusco. Beyond primal-dual methods in bandits with stochastic and adversarial constraints. *arXiv preprint arXiv:2405.16118*, 2024c.
- Arthur Flajolet and Patrick Jaillet. Logarithmic regret bounds for bandits with knapsacks. *arXiv preprint arXiv:1510.01800*, 2015.
- Yunlong Hou, Vincent YF Tan, and Zixin Zhong. Probably anytime-safe stochastic combinatorial semi-bandits. In *International Conference on Machine Learning*, pp. 13353–13409. PMLR, 2023.
- Alihan Hüyük and Cem Tekin. Thompson sampling for combinatorial network optimization in unknown environments. *IEEE/ACM Transactions on Networking*, 28(6):2836–2849, 2020.
- Nicole Immorlica, Karthik Sankararaman, Robert Schapire, and Aleksandrs Slivkins. Adversarial bandits with knapsacks. *Journal of the ACM*, 69(6):1–47, 2022.
- Wacharawan Intayoad, Chayapol Kamyod, and Punnarumol Temdee. Reinforcement learning based on contextual bandits for personalized online learning recommendation systems. *Wireless Personal Communications*, 115(4):2917–2932, 2020.
- Raunak Kumar and Robert Kleinberg. Non-monotonic resource utilization in the bandits with knapsacks problem. *Advances in Neural Information Processing Systems*, 35:19248–19259, 2022.
- Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight regret bounds for stochastic combinatorial semi-bandits. In *Artificial Intelligence and Statistics*, pp. 535–543. PMLR, 2015.
- Xiaocheng Li, Chunlin Sun, and Yinyu Ye. The symmetry between arms and knapsacks: A primal-dual approach for bandits with knapsacks. In *International Conference on Machine Learning*, pp. 6483–6492. PMLR, 2021.
- Shang Liu, Jiashuo Jiang, and Xiaocheng Li. Non-stationary bandits with knapsacks. *Advances in Neural Information Processing Systems*, 35:16522–16532, 2022.
- Xutong Liu, Jinhang Zuo, Siwei Wang, John CS Lui, Mohammad Hajiesmaili, Adam Wierman, and Wei Chen. Contextual combinatorial bandits with probabilistically triggered arms. In *International Conference on Machine Learning*, pp. 22559–22593. PMLR, 2023.

- Ahmadreza Moradipari, Christos Thrampoulidis, and Mahnoosh Alizadeh. Stage-wise conservative linear bandits. *Advances in neural information processing systems*, 33:11191–11201, 2020.
- Karthik Abinav Sankararaman and Aleksandrs Slivkins. Combinatorial semi-bandits with knapsacks. In *International Conference on Artificial Intelligence and Statistics*, pp. 1760–1770. PMLR, 2018.
- Karthik Abinav Sankararaman and Aleksandrs Slivkins. Bandits with knapsacks beyond the worst case. *Advances in Neural Information Processing Systems*, 34:23191–23204, 2021.
- Aleksandrs Slivkins. Dynamic ad allocation: Bandits with budgets. *arXiv preprint arXiv:1306.0155*, 2013.
- Aleksandrs Slivkins, Xingyu Zhou, Karthik Abinav Sankararaman, and Dylan J. Foster. Contextual bandits with packing and covering constraints: A modular lagrangian approach via regression. 2024. URL <https://arxiv.org/abs/2211.07484>.
- Marta Soare. *Sequential resource allocation in linear stochastic bandits*. PhD thesis, Université Lille 1-Sciences et Technologies, 2015.
- Long Tran-Thanh, Archie Chapman, Alex Rogers, and Nicholas Jennings. Knapsack based optimal policies for budget-limited multi-armed bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pp. 1134–1140, 2012.
- Sofía S Villar, Jack Bowden, and James Wason. Multi-armed bandit models for the optimal design of clinical trials: benefits and challenges. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 30(2):199, 2015.
- Qinshi Wang and Wei Chen. Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. *Advances in Neural Information Processing Systems*, 30, 2017.
- Yifan Wu, Roshan Shariff, Tor Lattimore, and Csaba Szepesvári. Conservative bandits. In *International Conference on Machine Learning*, pp. 1254–1262. PMLR, 2016.
- Baosheng Yu, Meng Fang, and Dacheng Tao. Linear submodular bandits with a knapsack constraint. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.

A TABLE OF NOTATIONS

Below, we provide the notations of some of the terms commonly used throughout the paper.

Table 2: Notations

K	number of arms
$[K + 1]$	set of arms including the <i>null arm</i>
ρ	mean cost vector of the arm set $[K + 1]$
μ	mean reward vector of the arm set $[K + 1]$
\mathbb{V}	the set of all possible valid bases
$c_i(t)$	cost of pulling arm i in round t
ρ_i	mean cost of arm i
$\bar{\rho}_i(t)$	empirical average cost of arm i in round t
$\epsilon_i(t)$	$\epsilon_i(t) = \sqrt{\frac{3 \log t}{N_i(t)}}$, the confidence interval of arm i in round t
$\rho_i^L(t)$	$\rho_i^L(t) = \bar{\rho}_i(t) - \epsilon_i(t)$, the lower confidence bound (LCB) of the cost of arm i in round t
$\mu_i^U(t)$	$\rho_i^L(t) = \bar{\mu}_i(t) + \epsilon_i(t)$, the upper confidence bound (UCB) of the reward of arm i in round t
$r_i(t)$	reward obtained from pulling arm i in round t
$r_{\mathcal{I}}$	reward of base \mathcal{I}
$r_{\mathcal{I}}^U(t)$	upper confidence bound of the reward of base \mathcal{I}
μ_i	mean reward of arm i
π^*	optimal solution of the problem under the linear relaxation
i^*	$i^* := \arg \max_{i \in [K]} \frac{\mu_i}{\rho_i}$
i^{**}	$i^{**} := \arg \max_{i \in [K]} \mu_i$
$S_c(t)$	$S_c(t) := \sum_{s=1}^t c(s)$ cumulative cost amassed until round t
$\bar{c}(t)$	$\bar{c}(t) = S_c(t)/t$ average cost obtained at round t
$N_i(t)$	Total number of times arm i is pulled until round $t + 1$
$N_{(i,j)}(t)$	The number of times base (i, j) is selected until round $t + 1$
$\Delta_{\min, i}$	$\Delta_{\min, i} := \min_{\mathcal{I} \in \mathbb{V} \setminus \mathcal{I}^*} \text{s.t. } i \in \mathcal{I} (\Delta_{\mathcal{I}})$

B ADDITIONAL RELATED WORKS

Conservative and safe bandits: *Conservative bandits* is a framework where an action is considered safe at round t if it keeps the cumulative reward up to round t above a given fraction of a baseline policy. In Wu et al. (2016), a baseline safe arm is given, and they provide an algorithm that utilizes the UCB principle where the baseline arm is pulled if the arm chosen by the UCB principle is not safe. In Moradipari et al. (2020), an anytime constraint is used instead of a constraint on cumulative reward that demands the expected reward of the pulled arm to be greater than a given threshold with high probability. They provide an algorithm that starts with a safe baseline and utilizes confidence regions to explore the other safe arms. *Safe bandits* is a similar framework to conservative bandits where in each round, the agent is required to select an arm with a given property no less than a predetermined (safe) threshold with high probability. For example, in Hou et al. (2023), the agent needs to choose at most K items from a set of L items; and with probability at least $1 - \delta$, the sum of variances of the selected items should not exceed a given threshold, which they call as the anytime-safe constraint. They propose a two step approach where first a set of arms is selected according to the UCB principle, and if this selection exceeds the threshold, the selected arms are split and pulled over multiple rounds. Our work is also similar in the sense that we also implement constraint checking as an additional separate step in SUAK.

Relation to Probabilistic Triggering in Combinatorial Bandits: The feedback obtained from the selected mixture of arms in BwAK problem resembles the feedback model in combinatorial bandits with probabilistically triggered arms. Probabilistic triggering is a special feedback model where when an action is played, a random subset of arms is triggered according to a triggering probability distribution, and the rewards of triggered arms are observed (Wang & Chen, 2017). Since the rewards of arms in a chosen super arm are only observed when that arm is triggered $p^* > 0$ is defined as

the minimum probability that an arm is triggered by any action; it is shown in Wang & Chen (2017, Theorem 3) that the regret lower bound scales with the factor $\frac{1}{p^*}$ for the general combinatorial bandits with probabilistically triggered arms, unless some additional assumptions are made. One such additional assumption is the *triggering probability modulated bounded smoothness* assumption that is used in Wang & Chen (2017). The main rationale behind this assumption is that an arm with a low triggering probability does not have a significant weight on the overall reward of an action, and perturbing its expected mean by a small amount would only cause a marginal change in the expected reward of an action. Leveraging this assumption, they prove regret bounds that are independent of p^* ; but are dependent on B , the bounded smoothness constant. This assumption is also used in many other subsequent work, such as in Hüyük & Tekin (2020). We also have partial observability of arms in our work since the frequency with which an arm is pulled within a given base is contingent upon the costs associated with each arm, such that the average cost of the base converges to c . However, we cannot use the *triggering probability modulated bounded smoothness* assumption in our work, as triggering probabilities of arms are dependent on empirical costs of arms, and subject to change every round. As a result of this, while the actual triggering probabilities are unknown, we have added the additional condition of pulling an arm in a base with at least ω probability in SUAK, and as such, our regret bounds depend on ω as the triggering probability. While time-varying triggering probabilities are considered in Liu et al. (2023) to derive instance-independent bounds that do not depend on the triggering probability, these results are not directly applicable for the instance-dependent bounds that we consider here.

C PRELIMINARIES AND AUXILIARY RESULTS

The following well known properties are used throughout the proofs:

Fact C.1 (Hoeffding’s Inequality). *Let Z_1, Z_2, \dots, Z_n be independent random variables bounded between $a_i \leq Z_i \leq b_i$, then for any $\delta > 0$, we have*

$$\mathbb{P}\left(\frac{\sum_{i=1}^n Z_i}{n} - \mathbb{E}[Z] \geq \delta\right) \leq e^{-\frac{2n^2\delta^2}{\sum_{i=1}^n (b_i - a_i)^2}}.$$

Fact C.2 (Conditional Probabilities). *The probability of an event A can be upper bounded by conditioning on an event B as follows*

$$\mathbb{P}(A) = \mathbb{P}(A, B) + \mathbb{P}(A, B^c) = \mathbb{P}(A|B)\mathbb{P}(B) + \mathbb{P}(A|B^c)\mathbb{P}(B^c) \leq \mathbb{P}(A|B) + \mathbb{P}(B^c).$$

Upper bounds of similar form are used throughout the proof.

C.1 OPTIMAL BASE OF THE OPTIMIZATION PROBLEM

The constrained optimization version of our problem that ignores the anytime cost budget constraint is expressed as follows:

$$\begin{aligned} OPT &= \max_{\pi} \mu^T \pi \\ \text{s.t. } &\rho^T \pi < c, \\ &\sum_{i=1}^{K+1} \pi_i = 1 \\ &\pi_i \geq 0, \forall i \in [K+1]. \end{aligned}$$

Letting $i^{**} := \arg \max_{i \in [K]} \mu_i$, and $i^* := \arg \max_{i \in [K]} \frac{\mu_i}{\rho_i}$, the solution of the problem can be found under three different cases as follows.

Case 1: If $\rho_{i^{**}} \leq c$, then $\mathcal{I}^* = \{i^{**}\}$. Since the cost of the arm with the highest mean reward is less than the cost constraint c , a mixture strategy is not needed and the optimal base includes only this arm.

Case 2: If $\rho_{i^{**}} > c$, $\rho_{i^*} > c$, then the optimal solution is mixing the arm with the highest mean reward per cost with the null arm. Hence, $\mathcal{I}^* = \{i^*, K+1\}$, and the optimal solution is $\pi_{i^*}^* = \frac{c}{\rho_{i^*}}$, and $\pi_{K+1}^* = 1 - \frac{c}{\rho_{i^*}}$. The optimal reward per round is $r^* = \frac{c\mu_{i^*}}{\rho_{i^*}}$.

Case 3: If $\rho_{i^{**}} > c$, $\rho_{i^*} < c$, then the optimal base will be of the form $\mathcal{I}^* = (i, j)$ where $\rho_i > c > \rho_j$; and can be found as:

$$\mathcal{I}^* = \arg \max_{i, j \in [K+1], i \neq j} r(i, j)$$

where

$$\begin{aligned} r(i, j) &= \max_{\pi_i, \pi_j} \mu_i \pi_i + \mu_j \pi_j \\ \text{s.t. } &\rho_i \pi_i + \rho_j \pi_j < c, \\ &\pi_i + \pi_j = 1 \\ &\pi_i \geq 0, \pi_j \geq 0. \end{aligned}$$

is the mean reward of base (i, j) . Note that the optimal base might or might not include i^* or i^{**} .

C.2 CONCENTRATION INEQUALITIES FOR THE CONFIDENCE INTERVALS

In this section, we derive concentration inequalities for the confidence intervals that we use throughout the paper.

Corollary C.3. *For an arm i that is sampled u times up to round t , the following results hold:*

$$\begin{aligned} \mathbb{P}(\mu_i^L(t, u) \geq \mu_i) &\leq t^{-6} \\ \mathbb{P}(\mu_i^U(t, u) \leq \mu_i) &\leq t^{-6} \\ \mathbb{P}(\rho_i^L(t, u) \geq \rho_i) &\leq t^{-6} \\ \mathbb{P}(\rho_i^U(t, u) \leq \rho_i) &\leq t^{-6} \\ \mathbb{P}(\bar{\rho}_i(t, u) \geq \rho_i + \sqrt{1.5 \log t / u}) &\leq t^{-3} \\ \mathbb{P}(\bar{\rho}_i(t, u) \leq \rho_i - \sqrt{1.5 \log t / u}) &\leq t^{-3}. \end{aligned}$$

where $\mu_i^L(t, u)$ is the lower confidence bound of arm i at round t when arm i is sampled u times up to round t ; and other variables are defined similarly.

Proof.

$$\mathbb{P}(\mu_i^L(t, u) \geq \mu_i) = \mathbb{P}\left(\frac{\sum_{s=1}^u r_i(t_{i,s})}{u} - \sqrt{\frac{3 \log t}{u}} \geq \mu_i\right)$$

where $t_{i,s}$ denotes the round in which s^{th} sample of arm i is obtained. Since the samples of arm i are independent across time, this expression can also be written independent of the time instant the sample from arm i was obtained as:

$$\mathbb{P}(\mu_i^L(t, u) \geq \mu_i) = \mathbb{P}\left(\frac{\sum_{s=1}^u \mu_{i,s}}{u} - \mu_i \geq \sqrt{\frac{3 \log t}{u}}\right)$$

where $\mu_{i,s}$ is the s^{th} sample of arm i . The result follows using Fact C.1:

$$\mathbb{P}(\mu_i^L(t, u) \geq \mu_i) = \mathbb{P}\left(\frac{\sum_{s=1}^u \mu_{i,s}}{u} - \mu_i \geq \sqrt{\frac{3 \log t}{u}}\right) \leq e^{-\frac{2u^2 \left(\sqrt{\frac{3 \log t}{u}}\right)^2}{u}} = t^{-6}$$

The other results are proved similarly. \square

Corollary C.4. For a base $\mathcal{I} = (i, j)$ if the arm i is sampled r times; and arm j is sampled s times up to round t , the following results hold:

$$\begin{aligned}\mathbb{P}(r_{\mathcal{I}}^L(t, r, s) > r_{\mathcal{I}}) &\leq 4t^{-6} \\ \mathbb{P}(r_{\mathcal{I}}^U(t, r, s) < r_{\mathcal{I}}) &\leq 4t^{-6}.\end{aligned}$$

where we define $r_{(i,j)}(t, r, s)$ as the empirical mean of the base $\mathcal{I} = (i, j)$ at round t when arm i has been sampled r times, and arm j has been sampled s times up to round t ; $r_{(i,j)}^U(t, r, s)$ as the upper confidence bound of the base $\mathcal{I} = (i, j)$ at round t when arm i has been sampled r times, and arm j has been sampled s times up to round t ; and $r_{(i,j)}^L(t, r, s)$ as the lower confidence bound of the base $\mathcal{I} = (i, j)$ at round t when arm i has been sampled r times, and arm j has been sampled s times up to round t .

Proof. Recall that $r_{(i,j)}$ is the optimal LP solution when only arms in the base $\mathcal{I} = (i, j)$; i.e. i and j ; are allowed. Hence, $r_{(i,j)}$ can be found as:

$$\begin{aligned}r_{(i,j)} &= \max_{\pi_i, \pi_j} \mu_i \pi_i + \mu_j \pi_j \\ \text{s.t. } &\rho_i \pi_i + \rho_j \pi_j < c, \\ &\pi_i + \pi_j = 1 \\ &\pi_i \geq 0, \pi_j \geq 0.\end{aligned}$$

Similar to this, $r_{\mathcal{I}}^U(t, r, s)$ is the solution to the equation below.

$$\begin{aligned}r_{(i,j)}^U(t, r, s) &= \max_{\beta_i, \beta_j} \mu_i^U(t, r) \cdot \beta_i + \mu_j^U(t, s) \cdot \beta_j \\ \text{s.t. } &\rho_i^L(t, r) \cdot \beta_i + \rho_j^L(t, s) \cdot \beta_j < c, \\ &\beta_i + \beta_j = 1 \\ &\beta_i \geq 0, \beta_j \geq 0.\end{aligned}$$

Using the fact that $\mathbb{P}(\mu_i^U(t, r) \geq \mu_i \wedge \mu_j^U(t, s) \geq \mu_j) \geq 1 - 2t^{-6}$ from Corollary C.3, the following holds with at least $1 - 2t^{-6}$ probability:

$$\begin{aligned}r_{(i,j)}^U(t, r, s) &= \mu_i^U(t, r) \cdot \beta_i + \mu_j^U(t, s) \cdot \beta_j \\ &\geq \mu_i \beta_i + \mu_j \beta_j\end{aligned}$$

Using the fact that $\mathbb{P}(\rho_i^L(t, r) \leq \rho_i \wedge \rho_j^L(t, s) \leq \rho_j) \geq 1 - 2t^{-6}$ from Corollary C.3, it can be seen that with probability at least $1 - 2t^{-6}$, $\{\pi_i, \pi_j : \rho_i \pi_i + \rho_j \pi_j < c\} \subset \{\beta_i, \beta_j : \rho_i \beta_i + \rho_j \beta_j < c\}$; i.e. the constraint on π_i and π_j is more restrictive than the constraint on β_i and β_j . Hence, it holds that $\mathbb{P}(\mu_i \beta_i + \mu_j \beta_j \geq \mu_i \pi_i + \mu_j \pi_j) \geq 1 - 2t^{-6}$. Combining this with the fact that $\mathbb{P}(r_{(i,j)}^U(t, r, s) \geq \mu_i \beta_i + \mu_j \beta_j) \geq 1 - 2t^{-6}$, the result follows. Note that the same result also follows if the base \mathcal{I} contains only one arm. This is since a base with only one arm can be viewed as a base with that arm and the null arm where the null arm is never pulled. \square

D PROOF OF THEOREM 3.1

Proof. To prove Theorem 3.1, we start by defining the events where the confidence bounds hold, which are also known as the good events in the bandit literature. First, define

$$\mathcal{G}_{\mathcal{I}}(t) := \left\{ \min_{r \leq t, s \leq t} r_{\mathcal{I}}^U(t, r, s) \geq r^* \wedge \max_{u \leq t, v \leq t} r_{\mathcal{I}}^L(t, u, v) \leq r_{\mathcal{I}} \right\}$$

as the good event for base \mathcal{I} at round t . Hence, $\mathcal{G}_{\mathcal{I}}(t)$ denotes the event where the confidence intervals of both base \mathcal{I} and the optimal base \mathcal{I}^* hold. The events $\mathcal{G}(t)$ and \mathcal{G}_T are defined as

$$\begin{aligned}\mathcal{G}(t) &:= \cap_{\mathcal{I} \in \mathbb{V}} \mathcal{G}_{\mathcal{I}}(t), \\ \mathcal{G}_T &:= \cap_{t=1}^T \mathcal{G}(t).\end{aligned}$$

Further, define

$$\begin{aligned}\mathcal{F}_i(t) &:= \left\{ \min_{r \leq t} \left(\bar{\rho}_i(t, r) + \sqrt{1.5 \log t / r} \right) \geq \rho_i \wedge \max_{r \leq t} \left(\bar{\rho}_i(t, r) - \sqrt{1.5 \log t / r} \right) \leq \rho_i \right\}, \\ \mathcal{F}(t) &:= \cap_{i=1}^K \mathcal{F}_i(t), \text{ and} \\ \mathcal{F}_T &:= \cap_{t=1}^T \mathcal{F}(t)\end{aligned}$$

as the good events for determining the confidence bounds of arm costs for line 13 of Algorithm 2. Using these events, the regret of SUAK can be decomposed as follows.

$$\begin{aligned}R_T &= OPT_{LP} - \mathbb{E}[F(T)] \\ &= T \cdot \boldsymbol{\mu}^T \boldsymbol{\pi}^* - \mathbb{E} \left[\sum_{t=1}^T r(t) \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T r^* - r(t) \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T r^* - r(t) \middle| \mathcal{G}_T, \mathcal{F}_T \right] + \sum_{t=1}^T (\mathbb{P}(\mathcal{G}^c(t)) + \mathbb{P}(\mathcal{F}^c(t)))\end{aligned}$$

We define the following four events based on the behaviour of SUAK:

- $\mathcal{E}_1(t)$: The round is skipped to satisfy the anytime constraint in line 5 of Algorithm 2
- $\mathcal{E}_2(t)$: Round is skipped so that the average cost incurred during pulls that are needed to reduce the confidence interval of the arm cost stay below c in line 9 of Algorithm 2
- $\mathcal{E}_3(t)$: An arm is pulled to reduce the confidence interval of the arm cost in line 13 of Algorithm 2
- $\mathcal{E}_4(t)$: A base is selected and an arm from this base is pulled in line 36 of Algorithm 2.

Using these events, regret can be decomposed as

$$R_T \leq R_a(T) + R_b(T) + R_c(T) + R_d(T) + \sum_{t=t_{\text{init}}+1}^T (\mathbb{P}(\mathcal{G}^c(t)) + \mathbb{P}(\mathcal{F}^c(t))) + t_{\text{init}}$$

where

$$\begin{aligned}R_a(T) &:= \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_1(t)\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ R_b(T) &:= \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_2(t)\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ R_c(T) &:= \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_3(t)\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ R_d(T) &:= \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_4(t), \mathcal{E}_3^c(t)\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right]\end{aligned}$$

Note that these four events $\mathcal{E}_1(t), \dots, \mathcal{E}_4(t)$ are mutually exclusive, any pair of these events cannot happen at the same time. Also note that in the definition of $R_d(T)$, we explicitly define the event as $\mathbb{1}\{\mathcal{E}_4(t), \mathcal{E}_3^c(t)\}$ to highlight that $\mathcal{E}_4(t)$ can happen only under $\mathcal{E}_3^c(t)$, i.e. when the confidence

intervals of arm costs have been reduced enough that whether the cost of an arm is greater or less than c is correctly known. This property is essential in satisfying the anytime constraint as it enables to pull an arm whose true mean cost is less than c if the targeted cost budget is exceeded.

Each term in the regret can be upper bounded as below.

Corollary D.1. $t_{\text{init}} \leq -\frac{1}{\omega^2 c} W\left(-\omega^2 c e^{-\omega^2 K}\right) = O(1)$, where $W(\cdot)$ is the Lambert function.

Proof of this result is provided in §D.2.

Lemma D.2. *The regret from skips needed to prevent violating the anytime constraint can be upper bounded as*

$$R_a(T) \leq \frac{3\pi^2 r^*}{\delta_{\min}^2}.$$

Proof of this result is provided in §D.6.

Lemma D.3. *The regret from pulls needed to reduce the confidence intervals of arm costs can be upper bounded as*

$$R_b(T) + R_c(T) \leq \sum_{i=1}^K \frac{96(r^* - \mu_i) \log T}{\delta_i^2} + \sum_{i: \rho_i > c} \frac{104r^* \log T}{c\delta_i} \leq \frac{200Kr^* \log T}{c\delta_{\min}^2}$$

Proof of this result is provided in §D.5.

Lemma D.4. *The regret from arm pulls due to line 36 of Algorithm 2, i.e. regret from pulls of an arm in a selected base, can be upper bounded as*

$$R_d(T) \leq \sum_{i=1}^K \frac{96r^* \log T}{\omega \Delta_{\min, i}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 + \frac{2r^* \log T}{c\omega^2}$$

Proof of this result is provided in §D.4.

Corollary D.5. *It holds that*

$$\sum_{t=t_{\text{init}}+1}^T \mathbb{P}(\mathcal{G}^c(t)) \leq \frac{4\pi^2 K^2}{3}$$

Proof of this result is provided in §D.1

Corollary D.6. *It holds that*

$$\sum_{t=t_{\text{init}}+1}^T \mathbb{P}(\mathcal{F}^c(t)) \leq \frac{\pi^2 K}{3}$$

Proof of this result is provided in §D.3.

Note that to present the main result in a simpler way, we use define R_K as

$$\sum_{t=t_{\text{init}}+1}^T \mathbb{P}(\mathcal{G}^c(t)) + \sum_{t=t_{\text{init}}+1}^T \mathbb{P}(\mathcal{F}^c(t)) \leq \frac{\pi^2 K}{3} + \frac{4\pi^2 K^2}{3} \leq R_K := \frac{5\pi^2 K^2}{3}$$

Combining all these results, it can be seen that

$$R_T \leq \sum_{i=1}^K \frac{96(\frac{\delta_i+1}{\delta_i})^2 \log T}{\omega \Delta_{\min, i}} + \frac{200Kr^* \log T}{c\delta_{\min}^2} + \frac{2r^* \log T}{c\omega^2} + \frac{3\pi^2 r^*}{\delta_{\min}^2} + R_K + r^* t_{\text{in}}$$

$$\begin{aligned}
&\leq \sum_{i=1}^K \frac{96(\frac{\delta_i+1}{\delta_i})^2 \log T}{\omega \Delta_{\min,i}} + \frac{202K r^* \log T}{c \delta_{\min}^2} + R_K + r^* t_{\text{in}} \\
&= O(K \log T) + O(1)
\end{aligned}$$

□

D.1 PROOF OF COROLLARY D.5

Proof.

$$\begin{aligned}
\sum_{t=1}^T \mathbb{1} \{ \mathcal{G}_{\mathcal{I}}^c(t) \} &= \sum_{t=1}^T \mathbb{1} \left\{ \min_{r \leq t, s \leq t} r_{\mathcal{I}^*}^U(t, r, s) < r^* \vee \max_{u \leq t, v \leq t} r_{\mathcal{I}}^L(t, u, v) > r_{\mathcal{I}} \right\} \\
&\leq \sum_{t=1}^T \sum_{r=1}^t \sum_{s=1}^t \sum_{u=1}^t \sum_{v=1}^t \mathbb{1} \{ r_{\mathcal{I}^*}^U(t, r, s) < r^* \vee r_{\mathcal{I}}^L(t, u, v) > r_{\mathcal{I}} \}
\end{aligned}$$

By the monotonicity of expectation, it holds that

$$\begin{aligned}
\sum_{t=1}^T \mathcal{G}_{\mathcal{I}}^c(t) &\leq \sum_{t=1}^T \sum_{r=1}^t \sum_{s=1}^t \sum_{u=1}^t \sum_{v=1}^t \mathbb{P} (r_{\mathcal{I}^*}^U(t, r, s) < r^* \vee r_{\mathcal{I}}^L(t, u, v) > r_{\mathcal{I}}) \\
&\leq \sum_{t=1}^T \sum_{r=1}^t \sum_{s=1}^t \sum_{u=1}^t \sum_{v=1}^t 8t^{-6} \\
&\leq \sum_{t=1}^T 8t^{-2} \leq \frac{4\pi^2}{3}
\end{aligned} \tag{3}$$

where we used Corollary C.4 in (3). The result follows using $\mathcal{G}(t) = \cap_{\mathcal{I} \in \mathcal{V}} \mathcal{G}_{\mathcal{I}}(t)$.

$$\begin{aligned}
\sum_{t=1}^T \mathcal{G}(t) &= \sum_{t=1}^T \bigcup_{\mathcal{I} \in \mathcal{V}} \mathcal{G}_{\mathcal{I}}^c(t) \\
&\leq \sum_{\mathcal{I} \in \mathcal{V}} \sum_{t=1}^T \mathcal{G}_{\mathcal{I}}^c(t) = \frac{4\pi^2 K^2}{3}
\end{aligned}$$

□

D.2 PROOF OF COROLLARY D.1

Proof. At round t_{init} , each arm will have been sampled once, and the incurred average cost will be less than or equal to the targeted average cost $c - \frac{\log t}{\omega^2 t}$. Hence,

$$\frac{\sum_{i=1}^K \rho_i(t_{\text{init},i})}{t_{\text{init}}} \leq c - \frac{\log t_{\text{init}}}{t_{\text{init}} \omega^2}$$

needs to hold.

$$\frac{\sum_{i=1}^K \rho_i(t_{\text{init},i})}{t_{\text{init}}} \leq \frac{K}{t_{\text{init}}} \leq c - \frac{\log t_{\text{init}}}{t_{\text{init}} \omega^2}$$

From this expression, it can be seen that $t_{\text{init}} \leq t_{\text{in}}$, where t_{in} is the solution of

$$t_{\text{in}} = \frac{K}{c} + \frac{\log t_{\text{in}}}{\omega^2 c},$$

which can be written explicitly as

$$t_{\text{in}} = -\frac{1}{\omega^2 c} W \left(-\omega^2 c e^{-\omega^2 K} \right)$$

where W is the Lambert function.

□

D.3 PROOF OF COROLLARY D.6

Proof.

$$\begin{aligned} \sum_{t=1}^T \mathbb{1} \{ \mathcal{F}_i^c(t) \} &= \sum_{t=1}^T \mathbb{1} \left\{ \min_{r \leq t} \left(\bar{\rho}_i(t, r) + \sqrt{1.5 \log t/r} \right) < \rho_i \vee \max_{r \leq t} \left(\bar{\rho}_i(t, r) - \sqrt{1.5 \log t/r} \right) > \rho_i \right\} \\ &\leq \sum_{t=1}^T \sum_{r=1}^t \mathbb{1} \left\{ \left(\bar{\rho}_i(t, r) + \sqrt{1.5 \log t/r} \right) < \rho_i \vee \left(\bar{\rho}_i(t, r) - \sqrt{1.5 \log t/r} \right) > \rho_i \right\} \end{aligned}$$

By the monotonicity of expectation, it holds that

$$\begin{aligned} \sum_{t=1}^T \mathbb{1} \{ \mathcal{F}_i^c(t) \} &\leq \sum_{t=1}^T \sum_{r=1}^t \mathbb{P} \left(\left(\bar{\rho}_i(t, r) + \sqrt{1.5 \log t/r} \right) < \rho_i \vee \left(\bar{\rho}_i(t, r) - \sqrt{1.5 \log t/r} \right) > \rho_i \right) \\ &\leq \sum_{t=1}^T \sum_{r=1}^t 2t^{-3} \\ &\leq \sum_{t=1}^T 2t^{-2} \leq \frac{\pi^2}{3} \end{aligned} \tag{4}$$

where we used Corollary C.4 in (4). The result follows using $\mathcal{F}(t) = \cap_{i=1}^K \mathcal{F}_i(t)$.

$$\begin{aligned} \sum_{t=1}^T \mathcal{F}(t) &= \sum_{t=1}^T \bigcup_{i=1}^K \mathcal{F}_i^c(t) \\ &\leq \sum_{i=1}^K \sum_{t=1}^T \mathcal{F}_i^c(t) = \frac{\pi^2 K}{3} \end{aligned}$$

□

D.4 PROOF OF LEMMA D.4

Regret will be incurred in $R_d(T)$ under two different ways; one is selecting a suboptimal base; and the other is the under-utilization of the cost budget. Under-utilization causes regret since even if the selected base is optimal, less reward can be obtained when the targeted average budget is less than c . First, we start by ignoring under-utilization (assume we target an anytime cost budget of ct), the regret from under-utilization will be added separately.

For SUAK to select a suboptimal base $\mathcal{I} = (i, j)$ in round t , the following needs to hold.

$$r_{\mathcal{I}}^U(t) \geq r_{\mathcal{I}^*}^U(t)$$

Under the good event $\mathcal{G}(t)$, $r^* \leq r_{\mathcal{I}^*}^U(t)$; $r_{\mathcal{I}} \geq r_{\mathcal{I}}^L(t)$; and hence $r_{\mathcal{I}}^U(t) \leq r_{\mathcal{I}}^{U,U}(t)$ holds where

$$\begin{aligned} r_{\mathcal{I}}^{U,U}(t) &= \max_{\beta_i(t), \beta_j(t)} (\mu_i + 2\epsilon_i(t))\beta_i(t) + (\mu_j + 2\epsilon_j(t))\beta_j(t) \\ \text{s.t. } &(\rho_i - 2\epsilon_i(t))\beta_i(t) + (\rho_j - 2\epsilon_j(t))\beta_j(t) < c, \\ &\beta_i(t) + \beta_j(t) = 1 \\ &\beta_i(t) \geq 0, \beta_j(t) \geq 0. \end{aligned}$$

Hence, the condition for SUAK to select a suboptimal base $\mathcal{I} = (i, j)$ in round t can be written as:

$$r^* \leq r_{\mathcal{I}^*}^U(t) \leq r_{\mathcal{I}}^U(t) \leq r_{\mathcal{I}}^{U,U}(t)$$

This means that the suboptimal base $\mathcal{I} = (i, j)$ would not be selected in round t if $r_{\mathcal{I}}^{U,U}(t) \leq r^*$. To analyze this, we note that when $\rho_i - 4\epsilon_i(t) > c > \rho_j$, the value of $r_{(i,j)}^{U,U}(t)$ can be written as:

$$r_{(i,j)}^{U,U}(t) = (\mu_i + 2\epsilon_i(t))\beta_i(t) + (\mu_j + 2\epsilon_j(t))\beta_j(t)$$

where

$$\beta_i(t) = \frac{c - \rho_j + 2\epsilon_j(t)}{\rho_i - \rho_j + 2(\epsilon_j(t) - \epsilon_i(t))}, \quad \beta_j(t) = \frac{\rho_i - c - 2\epsilon_i(t)}{\rho_i - \rho_j + 2(\epsilon_j(t) - \epsilon_i(t))}$$

Note that given $\mathcal{F}_i(t)$, $\rho_i - 4\epsilon_i(t) > c$ holds whenever SUAK selects a base due to the design of SUAK. If this was not the case, SUAK would pull arm i due to condition in line 13 and would not select a base at that round. Also note that we only consider the case where the arms i, j in the base satisfy $\rho_i > c > \rho_j$. The proof for the case where \mathcal{I} contains a single arm i with $\rho_i < c$ is much simpler and can be done similarly. Furthermore, under the event $\mathcal{F}_i(t) \wedge \mathcal{F}_j(t)$, it is not possible for SUAK to select a base $\mathcal{I} = (i, j)$, where $\rho_i > c, \rho_j > c$.

$r_{(i,j)}^{U,U}(t)$ can be written in terms of $r_{(i,j)}$ as follows:

$$\begin{aligned} r_{(i,j)}^{U,U}(t) &= (\mu_i + 2\epsilon_i(t))\beta_i(t) + (\mu_j + 2\epsilon_j(t))\beta_j(t) \\ &= \mu_i\beta_i(t) + 2\epsilon_i(t)\beta_i(t) + \mu_j\beta_j(t) + 2\epsilon_j(t)\beta_j(t) \\ &= \mu_i\beta_i(t) + \mu_j\beta_j(t) + 2\epsilon_i(t)\beta_i(t) + 2\epsilon_j(t)\beta_j(t) \\ &= \mu_i\pi_i + \mu_j\pi_j + \mu_i(\beta_i(t) - \pi_i) + \mu_j(\beta_j(t) - \pi_j) + 2\epsilon_i(t)\beta_i(t) + 2\epsilon_j(t)\beta_j(t) \\ &= \mu_i\pi_i + \mu_j\pi_j + \mu_i(\beta_i(t) - \pi_i) + \mu_j(1 - \beta_i(t) - 1 + \pi_i) + 2\epsilon_i(t)\beta_i(t) + 2\epsilon_j(t)\beta_j(t) \\ &= r_{(i,j)} + (\mu_i - \mu_j) \cdot (\beta_i(t) - \pi_i) + 2\epsilon_i(t)\beta_i(t) + 2\epsilon_j(t)\beta_j(t) \end{aligned}$$

The expression $\beta_i(t) - \pi_i$ can be written as

$$\begin{aligned} \beta_i(t) - \pi_i &= \frac{c - \rho_j + 2\epsilon_j(t)}{\rho_i - \rho_j + 2(\epsilon_j(t) - \epsilon_i(t))} - \frac{c - \rho_j}{\rho_i - \rho_j} \\ &= \frac{2\epsilon_j(t) \cdot (\rho_i - c) + 2\epsilon_i(t) \cdot (c - \rho_j)}{(\rho_i - \rho_j + 2\epsilon_j(t) - 2\epsilon_i(t)) \cdot (\rho_i - \rho_j)} \end{aligned}$$

Hence,

$$r_{(i,j)}^{U,U}(t) = r_{(i,j)} + (\mu_i - \mu_j) \cdot \frac{2\epsilon_j(t) \cdot (\rho_i - c) + 2\epsilon_i(t) \cdot (c - \rho_j)}{(\rho_i - \rho_j + 2\epsilon_j(t) - 2\epsilon_i(t)) \cdot (\rho_i - \rho_j)} + 2\epsilon_i(t) + 2\epsilon_j(t).$$

Hence, the base $\mathcal{I} = (i, j)$ can be selected in round t if

$$\begin{aligned} r^* - r_{(i,j)} &= \Delta_{(i,j)} \leq (\mu_i - \mu_j) \cdot \frac{2\epsilon_j(t) \cdot (\rho_i - c) + 2\epsilon_i(t) \cdot (c - \rho_j)}{(\rho_i - \rho_j + 2\epsilon_j(t) - 2\epsilon_i(t)) \cdot (\rho_i - \rho_j)} + 2\epsilon_i(t) + 2\epsilon_j(t) \\ &= \frac{(2\epsilon_j(t) \cdot \delta_i + 2\epsilon_i(t) \cdot \delta_j) \cdot (\mu_i - \mu_j)}{(\delta_i + \delta_j + 2\epsilon_j(t) - 2\epsilon_i(t)) \cdot (\delta_i + \delta_j)} + 2\epsilon_i(t) + 2\epsilon_j(t) \end{aligned}$$

where we used the cost gaps $\delta_i = |\rho_i - c|$. Under the event $\mathcal{F}_i(t)$, it is known that $\epsilon_i(t) < \delta_i/4$ holds due to the design of SUAK. Hence, the condition can be written as

$$\Delta_{(i,j)} \leq 2 \frac{(2\epsilon_j(t) \cdot \delta_i + 2\epsilon_i(t) \cdot \delta_j) \cdot (\mu_i - \mu_j)}{(\delta_i + \delta_j) \cdot (\delta_i + \delta_j)} + 2\epsilon_i(t) + 2\epsilon_j(t)$$

Define $N_{(i,j)}(t)$ as the number of times the base $\mathcal{I} = (i, j)$ is selected by SUAK up to round t . Since SUAK selects an arm in a base with at least ω probability, an individual arm in a base $\mathcal{I} = (i, j)$

will be pulled at least $\omega \cdot N_{(i,j)}(t)$ times in expectation (greater than is used as an arm can also be pulled if through other bases that include that arm), hence

$$\begin{aligned}\mathbb{E}[N_i(t)] &\geq \omega \cdot N_{(i,j)}(t) \\ \mathbb{E}[N_j(t)] &\geq \omega \cdot N_{(i,j)}(t)\end{aligned}$$

Defining

$$\epsilon_{(i,j)}(t) := \sqrt{\frac{3 \log t}{N_{(i,j)}(t)}},$$

it can be seen that $\sqrt{\omega} \cdot \epsilon_i(t) \leq \epsilon_{(i,j)}(t)$, $\sqrt{\omega} \cdot \epsilon_j(t) \leq \epsilon_{(i,j)}(t)$. Using this

$$\begin{aligned}\sqrt{\omega} \cdot \Delta_{(i,j)} &\leq 2 \frac{2\epsilon_{(i,j)}(t) \cdot (\delta_i + \delta_j) \cdot (\mu_i - \mu_j)}{(\delta_j + \delta_i) \cdot (\delta_i + \delta_j)} + 4\epsilon_{(i,j)}(t) \\ &\leq \frac{4\epsilon_{(i,j)}(t) \cdot (\mu_i - \mu_j)}{\delta_j + \delta_i} + 4\epsilon_{(i,j)}(t) \\ &= 4\epsilon_{(i,j)}(t) \cdot \left(1 + \frac{\mu_i - \mu_j}{\delta_i + \delta_j}\right) \\ &= 4\sqrt{\frac{3 \log t}{N_{(i,j)}(t)}} \left(1 + \frac{\mu_i - \mu_j}{\delta_i + \delta_j}\right)\end{aligned}$$

Using this, under the method described above, it can be seen that a suboptimal base $\mathcal{I} = (i, j)$ can be pulled at most

$$N_{(i,j)}(T) \leq \frac{48 \log T}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{\mu_i - \mu_j}{\delta_i + \delta_j}\right)^2$$

times in T rounds under the good event $\cap_{t=1}^T \mathcal{G}(t)$. Using the inequality above, it can also be seen that the following two inequalities hold.

$$\begin{aligned}N_{(i,j)}(T) &\leq \frac{48 \log T}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 \\ N_{(i,j)}(T) &\leq \frac{48 \log T}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{1}{\delta_j}\right)^2\end{aligned}$$

Note that information on an individual arm i can be obtained from any base $\mathcal{I} : i \in \mathcal{I}$, not just the base $\mathcal{I} = (i, j)$. As selecting and pulling arms from one base might reduce the need to pull the other base, simply summing the upper bounds of $N_{(i,j)}(T)$ values of all bases to find the total number of arm pulls needed would lead to an over-count. To prevent this kind of over-count, we consider an upper bound on the number of pulls of individual arms; and for this regard, we define the following event.

$$\mathcal{B}_t = \{\mathcal{I}_t \in \mathbb{V} \setminus \mathcal{I}^*\} \cap \left\{ \exists i \in \mathcal{I}_t : N_i(t) \leq \frac{48 \log t}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 \right\}$$

It can be seen that the base $\mathcal{I}_t \in \mathbb{V} \setminus \mathcal{I}^*$ cannot be chosen under the event \mathcal{B}_t^c . Further, define

$$\mathcal{B}_{i,t} := \mathcal{B}_t \cap \left\{ i \in \mathcal{I}_t, N_i(t) \leq \frac{48 \log t}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 \right\}$$

as the event that the arm i is not observed *sufficiently often* under event \mathcal{B}_t . Then, it can be seen that

$$\mathbb{1}\{\mathcal{B}_t, \Delta_{\mathcal{I}_t} > 0\} \leq \sum_{i=1}^K \mathbb{1}\{\mathcal{B}_{i,t}, \Delta_{\mathcal{I}_t} > 0\}.$$

Using this, regret can be bounded as:

$$\begin{aligned} R_d(T) &= \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_4(t), \mathcal{E}_3^c(t)\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ &\leq \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \sum_{\mathcal{I} \in \mathbb{V} \setminus \mathcal{I}^*} \mathbb{1}\{\mathcal{I}_t = \mathcal{I}\} \cdot r^* \middle| \mathcal{G}(T) \right] \\ &\quad + \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \mathbb{1}\{\mathcal{E}_4(t), \mathcal{E}_3^c(t)\} \cdot \mathbb{1}\{\mathcal{I}_t = \mathcal{I}^*\} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ &:= R_{d,1}(T) + R_{d,2}(T) \end{aligned}$$

where $R_{d,1}(T)$ is the gap from selecting a suboptimal base, and $R_{d,2}(T)$ is the regret due to under-budgeting. Starting with $R_{d,1}(T)$,

$$R_{d,1}(T) = \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \sum_{\mathcal{I} \in \mathbb{V} \setminus \mathcal{I}^*} \mathbb{1}\{\mathcal{I}_t = \mathcal{I}\} \cdot r^* \middle| \mathcal{G}(T) \right] \leq \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \sum_{i=1}^K \mathbb{1}\{\mathcal{B}_{i,t}\} \cdot r^* \middle| \mathcal{G}(T) \right]$$

where the gap is taken as r^* . Let $\mathcal{S}_i := \{\Delta_{\mathcal{I}} : \mathcal{I} \in \mathbb{V} \setminus \{\mathcal{I}^*\}, i \in \mathcal{I}\}$ be the set of gaps of suboptimal bases that include arm i . Also let $\sigma_{i,1} \geq \dots \geq \sigma_{i,|\mathcal{S}_i|}$ be the gaps of the bases in \mathcal{S}_i ordered from the one with largest gap to the smallest one. Note that $|\mathcal{S}_i|$ is the number of valid bases that contain the arm i ; and since two arms that have mean costs larger than c do not form a valid base, $|\mathcal{S}_i| \leq K$ will hold.

$$\begin{aligned} R_{d,1}(T) &\leq \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}_i|} \mathbb{1}\{\mathcal{B}_{i,t}, \Delta_{\mathcal{I}_t} = \sigma_{i,j}\} \cdot r^* \middle| \mathcal{E}(T) \right] \\ &\leq \mathbb{E} \left[\sum_{t=t_{\text{init}}+1}^T \left(\sum_{i=1}^K \sum_{j=1}^{|\mathcal{S}_i|} \mathbb{1}\left\{ i \in \mathcal{I}_t, N_i(t) \leq \frac{48 \log t}{\omega \Delta_{(i,j)}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 \right\} \cdot r^* \right) \middle| \mathcal{E}(T) \right] \end{aligned}$$

To proceed, as in (Kveton et al., 2015), we consider the worst case, i.e. the way with which the samples of arm i are obtained with the highest regret possible. The key idea is that this worst case occurs when first the base with highest gap is repeatedly selected to obtain samples of arm i until this base can no longer be selected, and then selecting the base with the highest gap among the remaining bases, and then repeatedly selecting that base, and so on. Since all bases have the same regret per sample, it can be seen that regret from samples for arm i will be bounded by

$$R_{d,1,i}(T) \leq \frac{96r^* \log T}{\omega \Delta_{\min,i}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2$$

□

We now upper bound $R_{d,2}(T)$, the regret from under-utilizing the cost budget. Note that we only consider the regret for under-budgeting from the selections of the suboptimal base since the upper

bound on the number of selections of a suboptimal base that we considered in $R_{d,1}(T)$ is not affected by under-budgeting.

Using similar arguments as in D.6, it can be shown that the probability that the empirical cost average of the pulls from the optimal base at round T being less than $c - 2 \log T / \omega^2 T$ is a constant. Because of this, the under-utilization of the cost budget is upper bounded by $2 \log T / \omega^2 T$. Hence,

$$R_{d,2}(T) \leq \frac{2r^* \log T}{c\omega^2}.$$

This is since the playing the optimal action could have obtained $\frac{2r^* \log T}{c\omega^2}$ reward with the unspent budget of $2 \log T / \omega^2 T$ as the optimal reward per unit cost is r^*/c .

Combining $R_{d,1}(T)$ and $R_{d,2}(T)$, the regret $R_d(T)$ can be upper bounded as

$$\begin{aligned} R_d(T) &= R_{d,1}(T) + R_{d,2}(T) \\ &\leq \sum_{i=1}^K \frac{96r^* \log T}{\omega \Delta_{\min,i}^2} \cdot \left(1 + \frac{1}{\delta_i}\right)^2 + \frac{2r^* \log T}{c\omega^2} \end{aligned}$$

D.5 PROOF OF LEMMA D.3

In round t , arm i might be pulled due to line 13 of Algorithm 2 if $\varrho_t^L(t) \leq c$ and $c \leq \varrho_t^U(t)$. Without loss of generality, we consider the case where $\rho_i > c$. The case where $\rho_i < c$ can be derived similarly. Under the good event $\mathcal{F}(t)$; $\rho_i \leq \bar{\rho}_i(t) + \sqrt{1.5 \log t / N_i(t)}$, and $\rho_i \geq \bar{\rho}_i(t) - \sqrt{1.5 \log t / N_i(t)}$ holds. Hence, $c \leq \varrho_t^U(t)$ will always hold under $\mathcal{F}(t)$ since

$$\varrho_t^U(t) = \bar{\rho}_i(t) + 7\sqrt{\frac{1.5 \log t}{N_i(t)}} \geq \bar{\rho}_i(t) + \sqrt{\frac{1.5 \log t}{N_i(t)}} \geq \rho_i \geq c.$$

Therefore, arm i will be pulled in round t if $\varrho_t^L(t) \leq c$. This condition can be written as

$$\bar{\rho}_i(t) - 7\sqrt{\frac{1.5 \log t}{N_i(t)}} \leq c$$

Using $\rho_i - \sqrt{1.5 \log t / N_i(t)} \leq \bar{\rho}_i(t)$, the following will hold under $\mathcal{F}(t)$.

$$\rho_i - 8\sqrt{\frac{1.5 \log t}{N_i(t)}} \leq \bar{\rho}_i(t) - 7\sqrt{\frac{1.5 \log t}{N_i(t)}} \leq c$$

Hence, arm i will be pulled in round t if

$$\rho_i - c = \delta_i \leq 8\sqrt{\frac{1.5 \log t}{N_i(t)}}$$

It can be concluded from here that the arm i can be pulled at most

$$N_i(T) \leq \frac{96 \log T}{\delta_i^2} \tag{5}$$

times in T rounds under the good event $\cap_{t=1}^T \mathcal{F}(t)$. Also note that when

$$\bar{\rho}_i(t) - 7\sqrt{\frac{1.5 \log t}{N_i(t)}} \geq c$$

holds, using $\rho_i + \sqrt{\frac{1.5 \log t}{N_i(t)}} \geq \bar{\rho}_i(t)$, it can be seen that

$$\rho_i - 6\sqrt{\frac{1.5 \log t}{N_i(t)}} \geq \bar{\rho}_i(t) - 7\sqrt{\frac{1.5 \log t}{N_i(t)}} \geq c$$

holds. Hence, the following holds at any round t under the event $\cap_{t=1}^T \mathcal{F}(t)$.

$$\frac{\delta_i}{6} \geq \sqrt{\frac{1.5 \log t}{N_i(t)}}$$

Thus, it can be seen that $\frac{\delta_i}{4} \geq \epsilon_i(t)$ at any round t . This outcome is used in the proof of Lemma D.4 in §D.4.

Using (5), the regret $R_c(T)$ can be upper bounded as

$$\begin{aligned} R_c(T) &= \mathbb{E} \left[\sum_{t=1}^T \mathbb{1} \{ \mathcal{E}_3(t) \} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ &\leq \sum_{i=1}^K \frac{96 \log T}{\delta_i^2} \cdot (r^* - \mu_i) \end{aligned}$$

Note that for some arms $\rho_i > c$, this regret term might be negative if $r^* \leq \mu_i$ holds. However, these arms will also cause skips, and the overall regret of these pulls will be reflected in $R_b(T) + R_c(T)$.

Now, we upper bound $R_b(T)$. It can be seen that on expectation only pulls from arms $\rho_i > c$ can lead to skips due to exceeding the average budget c . First, the total cost incurred from pulling arm an $i : \rho_i > c$ can be expressed as $\sum_{s=1}^T c_i(s) \cdot \mathbb{1} \{ i(t) = i \}$. Skipping is used to reduce the average cost incurred from these pulls to c . The number of skips needed to reduce the empirical average cost of an arm i to c , which we denote as $N_i^s(T)$, can be found through the following relation

$$c = \frac{\sum_{s=1}^T c_i(s) \cdot \mathbb{1} \{ i(t) = i \}}{N_i(T) + N_i^s(T)}$$

To proceed, note that $N_i(T) \leq \frac{96 \log T}{\delta_i^2}$ by (5). If $N_i(T) = \frac{96 \log T}{\delta_i^2}$, using Hoeffding's Inequality (Fact C.1); we can upper bound

$$\mathbb{P} \left(\sum_{s=1}^T c_i(s) \cdot \mathbb{1} \{ i(t) = i \} \leq N_i(T) \cdot (\rho_i + \delta_i/12) \right) \leq e^{-2 \frac{96 \log T}{\delta_i^2} \left(\frac{\delta_i}{12} \right)^2} \leq \frac{1}{T}$$

Hence, for the case where $N_i(T) = \frac{96 \log T}{\delta_i^2}$,

$$\begin{aligned} c \cdot (N_i(T) + N_i^s(T)) &= \sum_{s=1}^T c_i(s) \cdot \mathbb{1} \{ i(t) = i \} \leq \left(\rho_i + \frac{\delta_i}{12} \right) \cdot N_i(T) \\ c \cdot N_i^s(T) &\leq \left(\rho_i + \frac{\delta_i}{12} - c \right) \cdot N_i(T) = \frac{13\delta_i}{12} \cdot N_i(T) \end{aligned}$$

From this equation, it can be concluded that

$$N_i^s(T) \leq \frac{104 \log T}{c \delta_i} + 1$$

Note that we while used $N_i(T) = \frac{96 \log T}{\delta_i^2}$ to derive this result, the upper bound still holds for the case $N_i(T) \leq \frac{96 \log T}{\delta_i^2}$ as well since if $N_i(T)$ is decreased, it will also cause $N_i^s(T)$ to decrease. Also note that the term 1 is for the case where the Hoeffding's Inequality does not hold (with probability upper bounded by $1/T$).

To derive the upper bound on regret, we multiply the upper bound on the number of skips with r^* , hence

$$R_b(T) \leq r^* \cdot N_i^s(T) + O(1) \leq \frac{104r^* \log T}{c\delta_i} + 1$$

And the total regret $R_b(T) + R_c(T)$ can be upper bounded as

$$R_b(T) + R_c(T) \leq \sum_{i:\rho_i > c} \frac{104r^* \log T}{c\delta_i} + \sum_{i=1}^K \frac{96 \log T}{\delta_i^2} \cdot (r^* - \mu_i) + 1$$

□

D.6 PROOF OF LEMMA D.2

Define $\mathcal{Z}_t := \{\exists i \in \mathcal{I}_t : \rho_i < c\}$ as the event that the selected base is correctly identified, i.e. there exists an arm that has a cost less than c in the base. Assume that t_f is the round where $\bar{c}(t_f) + 1/(t_f + 1) > c$, which means that round $t_f + 1$ will be skipped due to the condition in the algorithm. Also define t_e as the latest round $t_e < t_f$ where $\bar{c}(t) \leq c - \frac{\log t}{\omega^2 t}$. Given $\mathcal{Z}_{t_e, t_f} := \cap_{s=t_e}^{t_f} \mathcal{Z}_s$, the algorithm will pull the arm with cost less than c in the base for all rounds $t_e < t < t_f$, and it can be seen that $t_f - t_e \geq \frac{\log t_e}{\omega^2(1-c)}$. This is since accumulated cost at round t_f is ct_f , and in round t_e , it is $ct_e - \log t_e/\omega^2$. Also using the fact that cost observed in a round is upper bounded by 1, it holds that $c(t_f - t_e) + \frac{\log t_e}{\omega^2} \leq t_f - t_e$. Note that we ignore the pulls that occur from determining if the cost of an arm is greater than or less than c (line 10 in Algorithm 1) since that part has its own skipping rule to limit the average cost attained from these pulls to c .

We define Z_i as the cost observed in round i where $t_e < i \leq t_f$ given the event \mathcal{Z}_i . This means that $\mathbb{E}[Z] \leq c - \delta_{\min}$ since under these circumstances, the algorithm will pull the arm with empirical average cost lower than c in the base with probability $1 - w$; and given the event $\mathcal{E}_3^c(t)$, the true mean cost of this arm will be less than $c - \delta_{\min}$. Further, with probability w , the arm with the higher cost will be pulled whose mean cost is bounded by 1. Defining \mathcal{S}_{t_e, t_f} , where $t_f \geq t_e + \frac{\log t_e}{\omega^2(1-c)}$ as the event that round t_f is skipped when the target budget started to be exceeded at round t_e , its probability can be upper bounded as:

$$\begin{aligned} \mathbb{P}(\mathcal{S}_{t_e, t_f}) &\leq \mathbb{P}\left(\sum_{s=t_e}^{t_f} Z_s \geq (t_f - t_e) \cdot c + \log(t_e)\right) \\ &= \mathbb{P}\left(\sum_{s=t_e}^{t_f} Z_s - (t_f - t_e) \cdot \mathbb{E}[Z] \geq (t_f - t_e) \cdot (c - \mathbb{E}[Z]) + \log(t_e)\right) \\ &\leq \mathbb{P}\left(\sum_{s=t_e}^{t_f} Z_s - (t_f - t_e) \cdot \mathbb{E}[Z] \geq (t_f - t_e) \cdot (c - \mathbb{E}[Z])\right) \end{aligned}$$

Since only arms whose cost is less than or equal to $c - \delta_{\min}$ are pulled with probability $1 - \omega$, and with probability w the cost is bounded by 1; the expected mean of the random variable Z_i will be less than or equal to $(c - \delta_{\min}) \cdot (1 - \omega) + \omega$, i.e. $\mathbb{E}[Z] \leq c - \delta_{\min} + \omega(1 + \delta_{\min} - c)$. Hence,

$$\mathbb{P}(\mathcal{S}_{t_f}) \leq \mathbb{P}\left(\sum_{s=t_e}^{t_f} Z_s - (t_f - t_e) \cdot \mathbb{E}[Z] \geq (t_f - t_e) \cdot (\delta_{\min} - \omega \cdot (1 + \delta_{\min} - c))\right)$$

$$\leq e^{-2(\delta_{\min} - \omega \cdot (1 + \delta_{\min} - c))^2 (t_f - t_e)}$$

Then, the total expected number of skips can be upper bounded as:

$$S(T) \leq \mathbb{E} \left[\sum_{t_e=1}^T \sum_{t_f=t_e + \left\lceil \frac{\log t_e}{\omega^2(1-c)} \right\rceil}^T \mathbb{1} \{ \mathcal{S}_{t_e, t_f} \} \right] \quad (6)$$

$$\begin{aligned} &= \sum_{t_e=1}^T \sum_{t_f=t_e + \left\lceil \frac{\log t_e}{\omega^2(1-c)} \right\rceil}^T \mathbb{P}(\mathcal{S}_{t_e, t_f}) \\ &\leq \sum_{t_e=1}^T \sum_{r=\left\lceil \frac{\log t_e}{\omega^2(1-c)} \right\rceil}^{\infty} e^{-2(\delta_{\min} - \omega \cdot (1 + \delta_{\min} - c))^2 r} \end{aligned} \quad (7)$$

$$\begin{aligned} &\leq \sum_{t_e=1}^T \frac{e^{-2(\delta_{\min} - \omega \cdot (1 + \delta_{\min} - c))^2 \cdot \frac{\log t_e}{\omega^2(1-c)}}}{1 - e^{-2(\delta_{\min} / (2 + \delta_{\min}))^2}} \\ &\leq \sum_{t_e=1}^T \frac{2(2 + \delta_{\min})^2}{\delta_{\min}^2} e^{-2(\delta_{\min} - \omega \cdot (1 + \delta_{\min} - c))^2 \cdot \frac{\log t_e}{\omega^2(1-c)}} \end{aligned} \quad (8)$$

The inequality in (6) is due to the fact that the upper limit of the summation over t_e is upper bounded by T . Fact C.1 is used in (7), and the fact that $e^{-2x} \leq 1 - x/2$, $\forall 0 \leq x \leq 1$ is used in (8). Using $\omega \leq \frac{\delta_{\min}}{2 + \delta_{\min} - c}$, it can be seen that

$$\begin{aligned} S(T) &\leq \sum_{t_e=1}^T \frac{18}{\delta_{\min}^2} e^{-2 \cdot \log t_e} \\ &\leq \frac{18}{\delta_{\min}^2} \sum_{t_e=1}^T \frac{1}{t_e^2} \leq \frac{3\pi^2}{\delta_{\min}^2} \end{aligned}$$

Using this, the result follows.

$$\begin{aligned} R_a(T) &= \mathbb{E} \left[\sum_{t=1}^T \mathbb{1} \{ \mathcal{E}_1(t) \} \cdot (r^* - r(t)) \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ &\leq r^* \cdot \mathbb{E} \left[\sum_{t=1}^T \mathbb{1} \{ \mathcal{E}_1(t) \} \middle| \mathcal{G}_T, \mathcal{F}_T \right] \\ &\leq r^* \cdot S(T) \leq \frac{3\pi^2 r^*}{\delta_{\min}^2} \end{aligned}$$

□