

---

# UCB with Clustering and Improved Exploration

---

Anonymous Author(s)

## Abstract

1 In this paper, we present a novel algorithm for the stochastic multi-armed bandit  
2 (MAB) problem. Our proposed Clustered UCB method, referred to as ClusUCB  
3 partitions the arms into clusters and then follows the UCB-Improved strategy with  
4 aggressive exploration factors to eliminate sub-optimal arms, as well as entire  
5 clusters. Through a theoretical analysis, we establish that ClusUCB achieves  
6 a better gap-dependent regret upper bound than UCB1 (Auer et al., 2002) and  
7 UCB-Improved (Auer and Ortner, 2010) and in the worst case matches the gap-  
8 dependent bound of MOSS (Audibert and Bubeck, 2009) and OCUCB (Lattimore,  
9 2015) algorithms. ClusUCB also achieves a gap-independent regret bound of  
10  $O(\sqrt{KT})$  which is better than UCB1 and UCB-Improved, is also comparable  
11 to MOSS and OCUCB and is order optimal. Further, numerical experiments on  
12 test-cases with small gaps between optimal and sub-optimal mean rewards show  
13 that ClusUCB results in lower cumulative regret than several popular UCB variants  
14 as well as MOSS, OCUCB, Thompson sampling and Bayes-UCB.

## 15 1 Introduction

16 In this paper, we consider the stochastic multi-armed bandit problem, a classical problem in sequential  
17 decision making. In this setting, a learning algorithm is provided with a set of decisions (or arms)  
18 with reward distributions unknown to the algorithm. The learning proceeds in an iterative fashion,  
19 where in each round, the algorithm chooses an arm and receives a stochastic reward that is drawn from  
20 a stationary distribution specific to the arm selected. Given the goal of maximizing the cumulative  
21 reward, the learning algorithm faces the exploration-exploitation dilemma, i.e., in each round should  
22 the algorithm select the arm which has the highest observed mean reward so far (*exploitation*), or  
23 should the algorithm choose a new arm to gain more knowledge of the true mean reward of the arms  
24 and thereby avert a sub-optimal greedy decision (*exploration*).

25 Let  $r_i, i = 1, \dots, K$  denote the mean reward of the  $i$ th arm out of the  $K$  arms and  $r^* = \max_i r_i$  the  
26 optimal mean reward. The objective in the stochastic bandit problem is to minimize the cumulative  
27 regret, which is defined as follows:

$$R_T = r^*T - \sum_{i \in A} r_i N_i(T),$$

28 where  $T$  is the number of timesteps,  $N_i(T) = \sum_{m=1}^T I(I_m = i)$  is the number of times the algorithm  
29 has chosen arm  $i$  up to timestep  $T$ . The expected regret of an algorithm after  $T$  timesteps can be  
30 written as

$$\mathbb{E}[R_T] = \sum_{i=1}^K \mathbb{E}[N_i(T)] \Delta_i,$$

31 where  $\Delta_i = r^* - r_i$  denotes the gap between the means of the optimal arm and the  $i$ -th arm.

32 An early work involving a bandit setup is Thompson (1933), where the author deals with the problem  
33 of choosing between two treatments to administer on patients who come in sequentially. Following

the seminal work of Robbins (1952), bandit algorithms have been extensively studied in a variety of applications. From a theoretical standpoint, an asymptotic lower bound for the regret was established in Lai and Robbins (1985). In particular, it was shown that for any consistent allocation strategy, we have  $\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\log T} \geq \sum_{\{i: r_i < r^*\}} \frac{(r^* - r_i)}{D(p_i || p^*)}$ , where  $D(p_i || p^*)$  is the Kullback-Leibler divergence between the reward densities  $p_i$  and  $p^*$ , corresponding to arms with mean  $r_i$  and  $r^*$ , respectively.

There have been several algorithms with strong regret guarantees. For further reference we point the reader to Bubeck et al. (2012). The foremost among them is UCB1 (Auer et al., 2002), which has a regret upper bound of  $O(\frac{K \log T}{\Delta})$ , where  $\Delta = \min_{i: \Delta_i > 0} \Delta_i$ . This result is asymptotically order-optimal for the class of distributions considered. However, the worst case gap independent regret bound of UCB1 can be as bad as  $O(\sqrt{TK \log T})$ . In Audibert and Bubeck (2009), the authors propose the MOSS algorithm and establish that the worst case regret of MOSS is  $O(\sqrt{TK})$  which improves upon UCB1 by a factor of order  $\sqrt{\log T}$ . However, the gap-dependent regret of MOSS is  $O(\frac{K^2 \log(T \Delta^2 / K)}{\Delta})$  and in certain regimes, this can be worse than even UCB1 (see (Audibert and Bubeck, 2009; Lattimore, 2015)). The UCB-Improved algorithm, proposed in Auer and Ortner (2010), is a round-based algorithm<sup>1</sup> variant of UCB1 that has a gap-dependent regret bound of  $O(\frac{K \log T \Delta^2}{\Delta})$ , which is better than that of UCB1. On the other hand, the worst case regret of UCB-Improved is  $O(\sqrt{TK \log K})$ . Recently in Lattimore (2015), the algorithm OCUCB achieves order-optimal gap-dependent regret bound of  $O(\sum_{i=2}^K \frac{\log(T/H_i)}{\Delta_i})$  where  $H_i = \sum_{j=1}^K \min\{\frac{1}{\Delta_i^2}, \frac{1}{\Delta_j^2}\}$  and gap-independent regret bound of  $O(\sqrt{KT})$ . This is the best known bound for the 1-sub-Gaussian distributions in the bandit literature. Moreover, certain powerful algorithms have also been proposed which we will not discuss in detail here for the sake of brevity. These algorithms, like KL-UCB (Garivier and Cappé, 2011), Bayes-UCB (Kaufmann et al., 2012) and Thompson Sampling (Thompson, 1933; Agrawal and Goyal, 2011) are known to perform well empirically and have strong gap-dependent regret guarantees. However, we show that all the aforementioned algorithms fail to take advantage of certain reward structures that our algorithm, by virtue of its implementation, is able to leverage. This discussion is deferred to the contribution section.

The idea of clustering in the bandit framework is not entirely new. In particular, the idea of clustering has been extensively studied in the contextual bandit setup, an extension of the MAB where side information or features are attached to each arm (see Auer (2002); Langford and Zhang (2008); Li et al. (2010); Beygelzimer et al. (2011); Slivkins (2014)). The clustering in this case is typically done over the feature space Bui et al. (2012); Cesa-Bianchi et al. (2013); Gentile et al. (2014), however, in our work we cluster or group the arms.

## 1.1 Our Contribution

We propose a variant of UCB algorithm, called Clustered UCB, henceforth referred to as ClusUCB, that incorporates clustering and an improved exploration scheme. ClusUCB starts with partitioning of arms into small clusters, each having same number of arms. The clustering is done at the start with a prespecified number of clusters. At the end of every round ClusUCB conducts both (individual) arm elimination as well as cluster elimination. This is the first algorithm in bandit literature which uses two simultaneous arm elimination conditions and shows both theoretically and empirically that such an approach is indeed helpful.

The clustering of arms provides two benefits. First, it creates a context where a UCB-Improved like algorithm can be run in parallel on smaller sets of arms with limited exploration, which could lead to fewer pulls of sub-optimal arms with the help of more aggressive elimination of sub-optimal arms. Second, the cluster elimination leads to whole sets of sub-optimal arms being simultaneously eliminated when they are found to yield poor results. These two simultaneous criteria for arm elimination can be seen as borrowing the strengths of UCB-Improved as well as other popular round based approaches.

We will also show that in certain environments ClusUCB is able to take advantage of the underlying structure of the reward distribution of arms that other algorithms fail to take advantage of. We will briefly discuss two of these examples here.

<sup>1</sup>An algorithm is *round-based* if it pulls all the arms equal number of times in each round and then proceeds to eliminate one or more arms that it identifies to be sub-optimal.

84 *1. Bernoulli Distribution with small gaps:* In this environment there are 20 arms with means  $r_{1:12} =$   
85  $0.01$ ,  $r_{13:19} = 0.07$  and  $r_{20}^* = 0.1$ . Here, ClusUCB because of random partitioning of arms into  
86 clusters, will create clusters where there are atleast one arm with means 0.07 and a significant number  
87 of arms with 0.01 means. These clusters behave like independent UCB-Improved algorithms with  
88 improved exploration factors and the arms with means 0.01 are quickly eliminated. Note that since  
89 gaps are very small and the gaps of arms with means 0.07 are very close to the optimal arm, comparing  
90 all arms to the single best performing arm at every timestep will result in fewer arm eliminations.  
91 Hence utilizing the clusters as in ClusUCB results in faster elimination of arms. This is shown in  
92 Experiment 1.

93 *2. Gaussian Distribution with different variances:* In this environment there are 100 arms with means  
94  $r_{1:66} = 0.1$ ,  $\sigma_{1:66}^2 = 0.7$ ,  $r_{67:99} = 0.8$ ,  $\sigma_{67:99}^2 = 0.1$  and  $r_{100}^* = 0.9$ ,  $\sigma_{100}^2 = 0.7$ . Here, the variance  
95 of the optimal arm and arms with mean farthest from the optimal arm are the highest. Whereas, the  
96 arms having mean closest to the optimal arm have lowest variances. In these type of cases, due to  
97 clustering ClusUCB is able to eliminate the arms with means 0.7 quickly because clusters containing  
98 atleast one arm with 0.8 mean behaves as independent UCB-Improved algorithms with improved  
99 exploration factors. This is shown in Experiment 2. Again, note that due to high variance of the  
100 optimal arm, comparing only with the best performing arm at every timestep results in fewer arm  
101 eliminations.

102 Theoretically, while ClusUCB does not achieve the gap-dependent regret bound of OCUCB, the  
103 theoretical analysis establishes that the gap-dependent regret of ClusUCB is always better than that  
104 of UCB-Improved and same as that of MOSS (see Table 1. Moreover, the gap-independent bound of  
105 ClusUCB is of the same order as of MOSS and OCUCB, i.e.,  $O(\sqrt{KT})$ .

Table 1: Regret upper bound of different algorithms

Algorithm	Gap-Dependent	Gap-Independent
ClusUCB	$O\left(\frac{K \log(T\Delta^2/K)}{\Delta}\right)$	$O(\sqrt{KT})$
UCB1	$O\left(\frac{K \log T}{\Delta}\right)$	$O(\sqrt{KT \log T})$
UCB-Imp	$O\left(\frac{K \log(T\Delta^2)}{\Delta}\right)$	$O(\sqrt{KT \log K})$
MOSS	$O\left(\frac{K \log(T\Delta^2/K)}{\Delta}\right)$	$O(\sqrt{KT})$
OCUCB	$O\left(\frac{K \log(T/H)}{\Delta}\right)$	$O(\sqrt{KT})$

106 On two synthetic setups (as discussed before) with small gaps, we observe empirically that ClusUCB  
107 outperforms UCB-ImprovedAuer and Ortner (2010), MOSSAudibert and Bubeck (2009) and  
108 OCUCBLattimore (2015) as well as other popular stochastic bandit algorithms such as UCB-  
109 VAudibert et al. (2009), Median EliminationEven-Dar et al. (2006), Thompson SamplingAgrawal  
110 and Goyal (2011), Bayes-UCBKaufmann et al. (2012) and KL-UCBGarivier and Cappé (2011).

111 The rest of the paper is organized as follows: In Section 2 we introduce ClusUCB. In Section 4, we  
112 present the associated regret bounds. In Section 5, we present the numerical experiments and provide  
113 concluding remarks in Section 6. Further proofs of lemmas, corollaries, theorems and propositions  
114 presented in Section 4 are provided in the appendices.

## 115 2 Algorithm: Clustered UCB

116 **Notation.** We denote the set of arms by  $A$ , with the individual arms labeled  $i, i = 1, \dots, K$ . We  
117 denote an arbitrary round of ClusUCB by  $m$ . We denote an arbitrary cluster by  $s_k$ , the subset of arms  
118 within the cluster  $s_k$  by  $A_{s_k}$  and the set of clusters by  $S$  with  $|S| = p \leq K$ . Here  $p$  is a pre-specified  
119 limit for the number of clusters. For simplicity, we assume that the optimal arm is unique and denote  
120 it by  $*$ , with  $s^*$  denoting the corresponding cluster. The best arm in a cluster  $s_k$  is denoted by  $a_{max_{s_k}}$ .  
121 We denote the sample mean of the rewards seen so far for arm  $i$  by  $\hat{r}_i$  and for the true best arm within

---

**Algorithm 1** ClusUCB

---

**Input:** Number of clusters  $p$ , time horizon  $T$ , exploration parameters  $\rho_a, \rho_s$  and  $\psi$ .

**Initialization:** Set  $B_0 := A$ ,  $S_0 = S$  and  $\epsilon_0 := 1$ .

Create a partition  $S_0$  of the arms at random into  $p$  clusters of size up to  $\ell = \left\lceil \frac{K}{p} \right\rceil$  each.

**for**  $m = 0, 1, \dots, \left\lfloor \frac{1}{2} \log_2 \frac{T}{e} \right\rfloor$  **do**

    Pull each arm in  $B_m$  so that the total number of times it has been pulled is  $n_m = \left\lceil \frac{\log(\psi T \epsilon_m^2)}{2\epsilon_m} \right\rceil$ .

**Arm Elimination**

        For each cluster  $s_k \in S_m$ , delete arm  $i \in s_k$  from  $B_m$  if

$$\hat{r}_i + \sqrt{\frac{\rho_a \log(\psi T \epsilon_m)}{2n_m}} < \max_{j \in s_k} \left\{ \hat{r}_j - \sqrt{\frac{\rho_a \log(\psi T \epsilon_m)}{2n_m}} \right\}$$

**Cluster Elimination**

        Delete cluster  $s_k \in S_m$  and remove all arms  $i \in s_k$  from  $B_m$  if

$$\max_{i \in s_k} \left\{ \hat{r}_i + \sqrt{\frac{\rho_s \log(\psi T \epsilon_m)}{2n_m}} \right\} < \max_{j \in B_m} \left\{ \hat{r}_j - \sqrt{\frac{\rho_s \log(\psi T \epsilon_m)}{2n_m}} \right\}.$$

    Set  $\epsilon_{m+1} := \frac{\epsilon_m}{2}$

    Set  $B_{m+1} := B_m$

    Stop if  $|B_m| = 1$  and pull  $i \in B_m$  till  $T$  is reached.

**end for**

---

122 a cluster  $s_k$  by  $\hat{r}_{a_{\max_{s_k}}}$ .  $z_i$  is the number of times an arm  $i$  has been pulled. We assume that the  
123 rewards of all arms are bounded in  $[0, 1]$ .

124 **The algorithm (ClusUCB):** As mentioned in a recent work Liu and Tsuruoka (2016), UCB-Improved  
125 has two shortcomings:

126 (i) A significant number of pulls are spent in early exploration, since each round  $m$  of UCB-Improved  
127 involves pulling every arm an identical  $n_m = \left\lceil \frac{2 \log(T \epsilon_m^2)}{\epsilon_m^2} \right\rceil$  number of times. The quantity  $\epsilon_m$  is  
128 initialized to 1 and halved after every round.

129 (ii) In UCB-Improved, arms are eliminated conservatively, i.e., only after  $\epsilon_m < \frac{\Delta_i}{2}$ , the sub-optimal  
130 arm  $i$  is discarded with high probability. This is disadvantageous when  $K$  is large and the gaps are  
131 identical ( $r_1 = r_2 = \dots = r_{K-1} < r^*$ ) and small.

132 To reduce early exploration, the number  $n_m$  of times each arm is pulled per round in ClusUCB is  
133 lower than that of UCB-Improved and also that of Median-Elimination, which used  $n_m = \frac{4}{\epsilon^2} \log\left(\frac{3}{\delta}\right)$ ,  
134 where  $\epsilon, \delta$  are confidence parameters. To handle the second problem mentioned above, ClusUCB  
135 partitions the larger problem into several small sub-problems using clustering and then performs local  
136 exploration aggressively to eliminate sub-optimal arms within each clusters with high probability.

137 As described in the pseudocode in Algorithm 1, ClusUCB begins with a initial clustering of arms  
138 that is performed by random uniform allocation. The set of clusters  $S$  thus obtained satisfies  $|S| = p$ ,  
139 with individual clusters having a size that is bounded above by  $\ell = \left\lceil \frac{K}{p} \right\rceil$ . Each round of ClusUCB  
140 involves both individual arm as well as cluster elimination conditions. These elimination conditions  
141 are inspired by UCB-Improved. Notice that, unlike UCB-Improved, there is no longer a single point  
142 of reference based on which we are eliminating arms. Instead now we have as many reference points  
143 to eliminate arms as number of clusters formed.

144 The exploration regulatory factor  $\psi$  governing the arm and cluster elimination conditions in ClusUCB  
145 is more aggressive than that in UCB-Improved. With appropriate choice of  $\psi$  and  $\rho_a$  and  $\rho_s$  we can  
146 achieve aggressive elimination even when the gaps  $\Delta_i$  are small and  $K$  is large.

147 In Liu and Tsuruoka (2016), the authors recommend incorporating a factor of  $d_i$  inside the log-term  
 148 of the UCB values, i.e.,  $\max\{\hat{r}_i + \sqrt{\frac{d_i \log T \epsilon_m^2}{2n_m}}\}$ . The authors there examine the following choices  
 149 for  $d_i$ :  $\frac{T}{t_i}$ ,  $\frac{\sqrt{T}}{t_i}$  and  $\frac{\log T}{t_i}$ , where  $t_i$  is the number of times an arm  $i$  has been sampled. Unlike Liu  
 150 and Tsuruoka (2016), we employ cluster as well as arm elimination and establish from a theoretical  
 151 analysis that the choice  $\psi = \frac{T}{K^2}$  helps in achieving a better gap-dependent regret upper bound for  
 152 ClusUCB as compared to UCB-Improved and MOSS (see Corollary 1 in the next section).

### 153 3 Algorithm: Efficient Clustered UCB

---

#### Algorithm 2 EClusUCB

---

**Input:** Number of clusters  $p$ , time horizon  $T$ , exploration parameters  $\rho_a, \rho_s$  and  $\psi$ .

**Initialization:** Set  $m := 0$ ,  $B_0 := A$ ,  $S_0 = S$ ,  $\epsilon_0 := 1$ ,  $M = \lfloor \frac{1}{2} \log_2 \frac{T}{e} \rfloor$ ,  $n_0 = \left\lceil \frac{\log(\psi T \epsilon_0^2)}{2\epsilon_0} \right\rceil$  and  
 $N_0 = K n_0$ .

Create a partition  $S_0$  of the arms at random into  $p$  clusters of size up to  $\ell = \left\lceil \frac{K}{p} \right\rceil$  each.

Pull each arm once

**for**  $t = K + 1, \dots, T$  **do**

Pull arm  $i \in \arg \max_{j \in B_m} \left\{ \hat{r}_j + \sqrt{\frac{\log(\psi T \epsilon_m^2)}{z_j}} \right\}$ , where  $z_j$  is the number of times arm  $j$  has  
 been pulled

**Arm Elimination**

For each cluster  $s_k \in S_m$ , delete arm  $i \in s_k$  from  $B_m$  if

$$\hat{r}_i + \sqrt{\frac{\rho_a \log(\psi T \epsilon_m)}{2z_i}} < \max_{j \in s_k} \left\{ \hat{r}_j - \sqrt{\frac{\rho_a \log(\psi T \epsilon_m)}{2z_j}} \right\}$$

**Cluster Elimination**

Delete cluster  $s_k \in S_m$  and remove all arms  $i \in s_k$  from  $B_m$  if

$$\max_{i \in s_k} \left\{ \hat{r}_i + \sqrt{\frac{\rho_s \log(\psi T \epsilon_m)}{2z_i}} \right\} < \max_{j \in B_m} \left\{ \hat{r}_j - \sqrt{\frac{\rho_s \log(\psi T \epsilon_m)}{2z_j}} \right\}.$$

**if**  $t \geq N_m$  and  $m \leq M$  **then**

**Reset Parameters**

$$\epsilon_{m+1} := \frac{\epsilon_m}{2}$$

$$B_{m+1} := B_m$$

$$n_{m+1} := \left\lceil \frac{\log(\psi T \epsilon_{m+1}^2)}{2\epsilon_{m+1}} \right\rceil$$

$$N_{m+1} := t + |B_{m+1}| n_{m+1}$$

$$m := m + 1$$

Stop if  $|B_m| = 1$  and pull  $i \in B_m$  till  $T$  is reached.

**end if**

**end for**

---

154 **The algorithm (EClusUCB):** One of the principal problems suffered by ClusUCB is that in every  
 155 round it pulls all the arms equal number of time. EClusUCB remedies this by implementing optimistic  
 156 greedy sampling, as done for CCB algorithm (see Liu and Tsuruoka (2016)). As described in the  
 157 pseudocode in Algorithm 2, EClusUCB is almost similar to ClusUCB. It starts with an initial uniform  
 158 clustering of arms and the total number of rounds,  $m = 0, 1, 2, \dots, M$  is also same. Each round of  
 159 EClusUCB consists a total of  $|B_m| n_m$  timesteps and parameters are updated at the end of each round.  
 160 The exploration parameters are also same for both the algorithms. The first major difference with  
 161 ClusUCB is that because of optimistic greedy sampling, EClusUCB only pulls the arm that has the

highest confidence interval at every timestep. Also EClusUCB conducts both individual arm as well as cluster elimination conditions at every timestep.

## 4 Main results

We now state the main result that upper bounds the expected regret of ClusUCB.

**Theorem 1 (Gap dependent regret bound)** For  $T \geq K^{2.4}$ ,  $\rho_a = \frac{1}{2}$ ,  $\rho_s = \frac{1}{2}$  and  $\psi = \frac{T}{K^2}$  the regret  $R_T$  of ClusUCB satisfies

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \Delta_i + 12K + \frac{32 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i} \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + 12K + \frac{64 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i} \right\} \\ & + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} 16K + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} 16K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ \Delta_i > b}} 32K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ 0 < \Delta_i \leq b}} 32K + \max_{i: \Delta_i \leq b} \Delta_i T, \end{aligned}$$

where  $b \geq \sqrt{\frac{e}{T}}$ , and  $A_{s^*}$  is the subset of arms in cluster  $s^*$  containing optimal arm  $a^*$ .

**Proof 1** The proof of this theorem is given in Appendix C.

**Remark:** The most significant term in the bound above is  $\sum_{i \in A: \Delta_i \geq b} \frac{64 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i}$  and hence, the regret upper bound for ClusUCB is of the order  $O\left(\frac{K \log\left(\frac{T\Delta^2}{K}\right)}{\Delta}\right)$ . Since Corollary 1 holds for all  $\Delta \geq \sqrt{\frac{e}{T}}$ , it can be clearly seen that for all  $\sqrt{\frac{e}{T}} \leq \Delta \leq 1$  and  $K \geq 2$ , the gap-dependent bound is better than that of UCB1, UCB-Improved. In the worst case scenario when all the gaps are uniform ClusUCB bound matches that of MOSS and OCUCB (see Table 1).

We now show the the gap-independent regret bound of ClusUCB in Corollary 1.

**Corollary 1 (Gap-independent bound)** Considering the same gap of  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}}$  for all  $i : i \neq *$  and with  $\psi = \frac{T}{K^2}$ ,  $p = \left\lceil \frac{K}{\log K} \right\rceil$ ,  $\rho_a = \frac{1}{2}$  and  $\rho_s = \frac{1}{2}$  and for  $T \geq K^{2.4}$ , we have the following gap-independent bound for the regret of ClusUCB:

$$\mathbb{E}[R_T] \leq 96\sqrt{KT} + 12K^2 + 44K \log K + \frac{64K^3}{K + \log K}$$

**Proof 2** The proof of this corollary is given in Appendix D

**Remarks:** From the above result, we observe that the order of the regret upper bound of ClusUCB is  $O(\sqrt{KT})$ , and this matches the order of MOSS and OUCUCB and is order optimal. This bound is also better than UCB1 and UCB-Improved.

Next, we state the regret upper bound for the special case of ClusUCB when  $p = 1$ , i.e there is a single cluster and there are no cluster elimination condition but only arm elimination condition. We name this algorithm ClusUCB-AE.

**Proposition 1** The regret  $R_T$  for ClusUCB-AE satisfies

$$\mathbb{E}[R_T] \leq \sum_{i \in A: \Delta_i > b} \left\{ 12K + \left( \Delta_i + \frac{32 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i} \right) + 16K \right\} + \sum_{i \in A: 0 < \Delta_i \leq b} 16K + \max_{i \in A: \Delta_i \leq b} \Delta_i T,$$

for all  $b \geq \sqrt{\frac{e}{T}}$ .

**Proof 3** The proof of this proposition is given in Appendix E

## Analysis of elimination error (Why Clustering?)

Let  $\tilde{R}_T$  denote the contribution to the expected regret in the case when the optimal arm  $*$  gets eliminated during one of the rounds of ClusUCB. This can happen if a sub-optimal arm eliminates  $*$  or if a sub-optimal cluster eliminates the cluster  $s^*$  that contains  $*$  – these correspond to cases b2 and b3 in the proof of Theorem 2 (see Section C). As stated before We shall denote variant of ClusUCB that includes arm elimination condition only as ClusUCB-AE while ClusUCB corresponds to Algorithm 1, which uses both arm and cluster elimination conditions. The regret upper bound for ClusUCB-AE is given in Proposition 1.

For ClusUCB-AE, the quantity  $\tilde{R}_T$  can be extracted from the proofs (in particular, case b2 in Appendix E) and simplified to obtain  $\tilde{R}_T = 32K^2$ . Finally, for ClusUCB, the relevant terms from Theorem 2 that corresponds to  $\tilde{R}_T$  can be simplified with  $\rho_a = \frac{1}{2}$ ,  $\rho_s = \frac{1}{2}$ ,  $p = \lceil \frac{K}{\log K} \rceil$  and  $\psi = \frac{T}{K^2}$  (as in Corollary 1 to obtain  $\tilde{R}_T = 32K \log K + \frac{64K^3}{K + \log K}$ . Hence, in comparison to ClusUCB-AE which has an elimination regret bound of  $O(K^2)$ , the elimination error regret bound of ClusUCB is lower and of the order  $O\left(\frac{K^3}{K + \log K}\right)$ . Thus, we observe that clustering in conjunction with improved exploration via  $\rho_a, \rho_s, p$  and  $\psi$  helps in reducing the factor associated with  $K^2$  for the gap-independent error regret bound for ClusUCB. Also in section 5, in experiment 4 we show that ClusUCB outperforms ClusUCB-AE.

Next, we present the main regret upper bound for EClusUCB.

**Theorem 2 (Gap dependent regret bound)** For  $T \geq K^{2.4}$ ,  $\rho_a = \frac{1}{2}$ ,  $\rho_s = \frac{1}{2}$  and  $\psi = \frac{T}{K^2}$  the regret  $R_T$  of EClusUCB satisfies

$$\begin{aligned} \mathbb{E}[R_T] \leq & \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \Delta_i + 4108K^{2.8} + \frac{32 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i} \right\} + \sum_{\substack{i \in A_i, \\ \Delta_i > b}} \left\{ 2\Delta_i + 4108K^{2.8} + \frac{64 \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i} \right\} \\ & + \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} 62K + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} 62K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ \Delta_i > b}} 124K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ 0 < \Delta_i \leq b}} 124K + \max_{i: \Delta_i \leq b} \Delta_i T, \end{aligned}$$

where  $b \geq \sqrt{\frac{e}{T}}$ , and  $A_{s^*}$  is the subset of arms in cluster  $s^*$  containing optimal arm  $a^*$ .

**Proof 4** The proof of this theorem is given in Appendix I.

**Remark:** Thus we see that the regret upper bound of EClusUCB is higher than that of ClusUCB but still it has the same order as  $O\left(\frac{K \log\left(\frac{T\Delta_i^2}{K}\right)}{\Delta_i}\right)$ . Proceeding similarly like Corollary 1 we can show that the gap-independent bound of EClusUCB satisfies an order of  $O\left(\sqrt{KT}\right)$  as that of ClusUCB and is order optimal.

## 5 Simulation experiments

We conduct an empirical performance using cumulative regret as the metric. We implement the following algorithms: KL-UCB Garivier and Cappé (2011), MOSS Audibert and Bubeck (2009), UCB1 Auer et al. (2002), UCB-Improved Auer and Ortner (2010), Median Elimination Even-Dar et al. (2006), Thompson Sampling (TS) Agrawal and Goyal (2011), OCUCB Lattimore (2015), Bayes-UCB (BU) Kaufmann et al. (2012) and UCB-VAudibert et al. (2009)<sup>2</sup>. The parameters of EClusUCB algorithm for all the experiments are set as follows:  $\psi = \frac{T}{K^2}$ ,  $\rho_s = 0.5$ ,  $\rho_a = 0.5$  and  $p = \lceil \frac{K}{\log K} \rceil$  (as in Corollary 1).

**First experiment (Bernoulli with small gaps) :** This is conducted over a testbed of 20 arms in an environment involving Bernoulli reward distributions with expected rewards of the arms  $r_{i \neq *}$  = 0.07 and  $r^*$  = 0.1. These type of cases are frequently encountered in web-advertising domain. The

<sup>2</sup>The implementation for KL-UCB, Bayes-UCB and DMED were taken from Cappe et al. (2012)

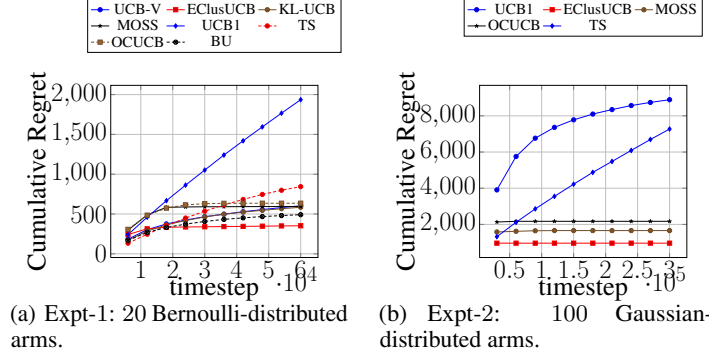


Figure 1: Cumulative regret for various bandit algorithms on two stochastic  $K$ -armed bandit environments.

horizon  $T$  is set to 60000. The regret is averaged over 100 independent runs and is shown in Figure 1(a). EClusUCB, MOSS, UCB1, UCB-V, KL-UCB, TS, BU and DMED are run in this experimental setup and we observe that EClusUCB performs better than all the aforementioned algorithms except TS. Because of the small gaps and short horizon  $T$ , we do not implement UCB-Improved and Median Elimination on this test-case.

**Second experiment (Gaussian with different variances):** This is conducted over a testbed of 100 arms involving Gaussian reward distributions with expected rewards of the arms  $r_{1:33} = 0.7$ ,  $r_{34:99} = 0.8$  and  $r_{100} = 0.9$  with variance set at  $\sigma_{1:33}^2 = 0.7$ ,  $\sigma_{34:99}^2 = 0.1$  and  $\sigma_{100}^2 = 0.7$ . The horizon  $T$  is set for a large duration of  $3 \times 10^5$  and the regret is averaged over 100 independent runs and is shown in Figure 1(b). From the results in Figure 1(b), we observe that EClusUCB outperforms MOSS, UCB1, UCB-Improved and Median-Elimination ( $\epsilon = 0.1, \delta = 0.1$ ). Also the performance of UCB-Improved is poor in comparison to other algorithms, which is probably because of pulls wasted in initial exploration whereas EClusUCB with the choice of  $\psi$ ,  $\rho_a$  and  $\rho_s$  performs much better. Note that the performance of TS is poor and this is in line with the observation in Lattimore (2015) that the worst case regret of TS in Gaussian distributions is  $\Omega(\sqrt{KT \log T})$ .

## 6 Conclusions and future work

From a theoretical viewpoint, we conclude that the gap-dependent regret bound of ClusUCB is lower than UCB1 and UCB-Improved and its gap-independent regret bound is of the same order as MOSS and OCUCB and is also order optimal. From the numerical experiments in specific environments, we observed that EClusUCB outperforms several popular bandit algorithms, including OCUCB, TS and BU which fail to leverage the structure of the rewards. Also ClusUCB is remarkably stable for a large horizon and large number of arms and performs well across different types of distributions. While we exhibited better regret bounds for ClusUCB, it would be interesting future research to improve the theoretical analysis of ClusUCB to achieve the gap-dependent regret bound of OCUCB. This is also one of the first papers to apply clustering in stochastic MAB and another future direction is to use this in contextual or in distributed bandits.



## References

- Agrawal, S. and Goyal, N. (2011). Analysis of thompson sampling for the multi-armed bandit problem. *arXiv preprint arXiv:1111.1797*.
- Audibert, J.-Y. and Bubeck, S. (2009). Minimax policies for adversarial and stochastic bandits. In *COLT*, pages 217–226.
- Audibert, J.-Y., Munos, R., and Szepesvári, C. (2009). Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902.
- Auer, P. (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256.
- Auer, P. and Ortner, R. (2010). Ucb revisited: Improved regret bounds for the stochastic multi-armed bandit problem. *Periodica Mathematica Hungarica*, 61(1-2):55–65.
- Beygelzimer, A., Langford, J., Li, L., Reyzin, L., and Schapire, R. E. (2011). Contextual bandit algorithms with supervised learning guarantees. In *AISTATS*, pages 19–26.
- Bubeck, S., Cesa-Bianchi, N., and Lugosi, G. (2012). Bandits with heavy tail. *arXiv preprint arXiv:1209.1727*.
- Bubeck, S., Munos, R., and Stoltz, G. (2011). Pure exploration in finitely-armed and continuous-armed bandits. *Theoretical Computer Science*, 412(19):1832–1852.
- Bui, L., Johari, R., and Mannor, S. (2012). Clustered bandits. *arXiv preprint arXiv:1206.4169*.
- Cappe, O., Garivier, A., and Kaufmann, E. (2012). pymabandits. <http://mloss.org/software/view/415/>.
- Cesa-Bianchi, N., Gentile, C., and Zappella, G. (2013). A gang of bandits. In *Advances in Neural Information Processing Systems*, pages 737–745.
- Even-Dar, E., Mannor, S., and Mansour, Y. (2006). Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *The Journal of Machine Learning Research*, 7:1079–1105.
- Garivier, A. and Cappé, O. (2011). The kl-ucb algorithm for bounded stochastic bandits and beyond. *arXiv preprint arXiv:1102.2490*.
- Gentile, C., Li, S., and Zappella, G. (2014). Online clustering of bandits. In *ICML*, pages 757–765.
- Kaufmann, E., Cappé, O., and Garivier, A. (2012). On bayesian upper confidence bounds for bandit problems. In *AISTATS*, pages 592–600.
- Lai, T. L. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22.
- Langford, J. and Zhang, T. (2008). The epoch-greedy algorithm for multi-armed bandits with side information. In *Advances in neural information processing systems*, pages 817–824.
- Lattimore, T. (2015). Optimally confident ucb: Improved regret for finite-armed bandits. *arXiv preprint arXiv:1507.07880*.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670. ACM.
- Liu, Y.-C. and Tsuruoka, Y. (2016). Modification of improved upper confidence bounds for regulating exploration in monte-carlo tree search. *Theoretical Computer Science*.

- 295 Robbins, H. (1952). Some aspects of the sequential design of experiments. In *Herbert Robbins*  
296 *Selected Papers*, pages 169–177. Springer.
- 297 Slivkins, A. (2014). Contextual bandits with similarity information. *Journal of Machine Learning*  
298 *Research*, 15(1):2533–2568.
- 299 Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of  
300 the evidence of two samples. *Biometrika*, pages 285–294.

## 301 Appendix

302 The Appendix is organized as follows. First we prove some technical lemmas in Appendix A and  
 303 Appendix B. Next we prove the main theorem in Appendix C. In Appendix D we prove Corollary 1.  
 304 In Appendix E we prove Proposition 1.

### 305 A Proof of Lemma 1

306 **Lemma 1** If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho_a = \frac{1}{2}$  and  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$ , then,

$$\frac{\rho_a m \log(2)}{\log(\psi T) - 2m \log(2)} \leq \frac{3}{2}$$

307 **Proof 5** The proof is based on contradiction. Suppose

$$\frac{\rho_a m \log(2)}{\log(\psi T) - 2m \log(2)} > \frac{3}{2}.$$

308 Then, with  $\psi = \frac{T}{K^2}$  and  $\rho_a = \frac{1}{2}$ , we obtain

$$\begin{aligned} \frac{\rho_a m \log(2)}{\log\left(\frac{T^2}{K^2}\right) - 2m \log(2)} &> \frac{3}{2} \\ \Rightarrow 2\rho_a m \log(2) &> 6 \log\left(\frac{T}{K}\right) - 6m \log(2) \end{aligned}$$

309 This can be further reduced to,

$$\begin{aligned} 6 \log(K) &> 6 \log(T) - 7m \log(2) \\ &\stackrel{(a)}{\geq} 6 \log(T) - \frac{7}{2} \log_2 \left( \frac{T}{e} \right) \log(2) \\ &= 2.5 \log(T) + 3.5 \log_2(e) \log(2) \\ &\stackrel{(b)}{=} 2.5 \log(T) + 3.5 \end{aligned}$$

310 where (a) is obtained using  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$ , while (b) follows from the identity  $\log_2(e) \log(2) = 1$ .

311 Finally, for  $T \geq K^{2.4}$  we obtain,  $6 \log(K) > 6 \log(K) + 3.5$ , which is a contradiction. Hence, for

312  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho = \frac{1}{2}$  and  $m \leq \frac{1}{2} \log_2 \left( \frac{T}{e} \right)$  we have,

$$\frac{\rho m \log(2)}{\log(\psi T) - 2m \log(2)} \leq \frac{3}{2}$$

### 313 B Proof of Lemma 2

314 **Lemma 2** If  $T \geq K^{2.4}$ ,  $\psi = \frac{T}{K^2}$ ,  $\rho_a = \frac{1}{2}$ ,  $m_i = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_i}{4}\}$  and  $c_{m_i} =$

315  $\sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}}}$ , then,  $c_{m_i} < \frac{\Delta_i}{4}$ .

316 **Proof 6** In the  $m_i$ -th round  $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}}}$ . Substituting the value of  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$

317 in  $c_{m_i}$  we get,

$$c_{m_i} \leq \sqrt{\frac{\rho_a \epsilon_{m_i} \log(\psi T \epsilon_{m_i})}{\log(\psi T \epsilon_{m_i}^2)}} \leq \sqrt{\frac{\rho_a \epsilon_{m_i} \log\left(\frac{\psi T \epsilon_{m_i}^2}{\epsilon_{m_i}}\right)}{\log(\psi T \epsilon_{m_i}^2)}}$$

$$\begin{aligned}
&= \sqrt{\frac{\rho_a \epsilon_{m_i} \log(\psi T \epsilon_{m_i}^2) - \rho_a \epsilon_{m_i} \log(\epsilon_{m_i})}{\log(\psi T \epsilon_{m_i}^2)}} \leq \sqrt{\rho_a \epsilon_{m_i} - \frac{\rho_a \epsilon_{m_i} \log(\frac{1}{2^{m_i}})}{\log(\psi T \frac{1}{2^{2m_i}})}} \\
&\leq \sqrt{\rho_a \epsilon_{m_i} + \frac{\rho_a \epsilon_{m_i} \log(2^{m_i})}{\log(\psi T) - \log(2^{2m_i})}} \leq \sqrt{\rho_a \epsilon_{m_i} + \frac{\rho_a \epsilon_{m_i} m_i \log(2)}{\log(\psi T) - 2m_i \log(2)}} \\
&\stackrel{(a)}{\leq} \sqrt{\rho_a \epsilon_{m_i} + \frac{3}{2} \epsilon_{m_i}} < \sqrt{2 \epsilon_{m_i}} < \frac{\Delta_i}{4}
\end{aligned}$$

318 In the above simplification, (a) is obtained using Lemma 1.

## 319 C Proof of Theorem 1

320 **Proof 7** Let  $A' = \{i \in A, \Delta_i > b\}$ ,  $A'' = \{i \in A, \Delta_i > 0\}$ ,  $A'_{s_k} = \{i \in A_{s_k}, \Delta_i > b\}$  and  
321  $A''_{s_k} = \{i \in A_{s_k}, \Delta_i > 0\}$ .  $C_g$  is the cluster set containing max payoff arm from each cluster  
322 in  $g$ -th round. The arm having the true highest payoff in a cluster  $s_k$  is denote by  $a_{\max_{s_k}}$ . Let  
323 for each sub-optimal arm  $i \in A$ ,  $m_i = \min \{m | \sqrt{2\epsilon_m} < \frac{\Delta_i}{4}\}$  and let for each cluster  $s_k \in S$ ,  
324  $g_{s_k} = \min \{g | \sqrt{2\epsilon_g} < \frac{\Delta_{a_{\max_{s_k}}}}{4}\}$ . Let  $\tilde{A} = \{i \in A' | i \in s_k, \forall s_k \in S\}$ . The analysis proceeds by  
325 considering the contribution to the regret in each of the following cases:

326 **Case a:** Some sub-optimal arm  $i$  is not eliminated in round  $\max(m_i, g_{s_k})$  or before, with the optimal  
327 arm  $* \in C_{\max(m_i, g_{s_k})}$ . We consider an arbitrary sub-optimal arm  $i$  and analyze the contribution to  
328 the regret when  $i$  is not eliminated in the following exhaustive sub-cases:

329 **Case a1:** In round  $\max(m_i, g_{s_k})$ ,  $i \in s^*$ .

330 Similar to case (a) of Auer and Ortner (2010), observe that when the following two conditions hold,  
331 arm  $i$  gets eliminated:

$$\hat{r}_i \leq r_i + c_{m_i} \text{ and } \hat{r}^* \geq r^* - c_{m_i}, \quad (1)$$

332 where  $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}}}$ . The arm  $i$  gets eliminated because

$$\begin{aligned}
\hat{r}_i + c_{m_i} &\leq r_i + 2c_{m_i} < r_i + \Delta_i - 2c_{m_i} \\
&\leq r^* - 2c_{m_i} \leq \hat{r}^* - c_{m_i}.
\end{aligned}$$

333 In the above, we have used the fact that  $c_{m_i} = \sqrt{\epsilon_{m_i+1}} < \frac{\Delta_i}{4}$ , from Lemma 2. From the foregoing, we  
334 have to bound the events complementary to that in (1) for an arm  $i$  to not get eliminated. Considering  
335 Chernoff-Hoeffding bound this is done as follows:

$$\begin{aligned}
\mathbb{P}(\hat{r}_i \geq r_i + c_{m_i}) &\leq \exp(-2c_{m_i}^2 n_{m_i}) \\
&\leq \exp(-2 * \frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}} * n_{m_i}) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}}
\end{aligned}$$

336 Along similar lines, we have  $\mathbb{P}(\hat{r}^* \leq r^* - c_{m_i}) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}}$ . Thus, the probability that a sub-  
337 optimal arm  $i$  is not eliminated in any round on or before  $m_i$  is bounded above by  $\left(\frac{2}{(\psi T \epsilon_{m_i})^{\rho_a}}\right)$ .

338 Summing up over all arms in  $A'_{s^*}$  in conjunction with a simple bound of  $T \Delta_i$  for each arm we obtain,

$$\sum_{i \in A'_{s^*}} \left( \frac{2T \Delta_i}{(\psi T \epsilon_{m_i})^{\rho_a}} \right) \leq \sum_{i \in A'_{s^*}} \left( \frac{2T \Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_a}} \right) \stackrel{(a)}{\leq} \sum_{i \in A'_{s^*}} \left( \frac{2T \Delta_i}{(\frac{T^2 \Delta_i^2}{K^2 32})^{\frac{1}{2}}} \right) \leq 8\sqrt{2} \sum_{i \in A'_{s^*}} K$$

339 Here, in (a) we substituted the value  $\rho_a$  and  $\psi$ .

340 **Case a2:** In round  $\max(m_i, g_{s_k})$ ,  $i \in s_k$  for some  $s_k \neq s^*$ .

341 Following a parallel argument like in Case a1, we have to bound the following two events of arm  
342  $a_{\max_{s_k}}$  not getting eliminated on or before  $g_{s_k}$ -th round,

$$\hat{r}_{a_{\max_{s_k}}} \geq r_{a_{\max_{s_k}}} + c_{g_{s_k}} \text{ and } \hat{r}^* \leq r^* - c_{g_{s_k}}$$

343 We can prove using Chernoff-Hoeffding bounds and considering independence of events mentioned  
 344 above, that for  $c_{g_{s_k}} = \sqrt{\frac{\rho_s \log(\psi T \epsilon_{g_{s_k}})}{2n_{g_{s_k}}}}$  and  $n_{g_{s_k}} = \frac{\log(\psi T \epsilon_{g_{s_k}}^2)}{2\epsilon_{g_{s_k}}}$  the probability of the above two  
 345 events is bounded by  $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}})^{\rho_s}}\right)$ .

346 Now, for any round  $g_{s_k}$ , all the elements of  $C_{\max(m_i, g_{s_k})}$  are the respective maximum payoff arms  
 347 of their cluster  $s_k$ ,  $\forall s_k \in S$ , and since clusters are fixed so we can bound the maximum probability  
 348 that a sub-optimal arm  $i \in A'$  and  $i \in s_k$  such that  $a_{\max_{s_k}} \in C_{g_{s_k}}$  is not eliminated on or before  
 349 the  $g_{s_k}$ -th round by the same probability as above. Summing up over all  $p$  clusters and bounding the  
 350 regret for each arm  $i \in A'_{s_k}$  trivially by  $T\Delta_i$ ,

$$\begin{aligned} \sum_{k=1}^p \sum_{i \in A'_{s_k}} \left( \frac{2T\Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_s}} \right) &= \sum_{i \in A'} \left( \frac{2T\Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_s}} \right) \\ &\stackrel{(a)}{\leq} \sum_{i \in A'} \left( \frac{2T\Delta_i}{(\frac{T^2}{K^2} \frac{\Delta_i^2}{32})^{\frac{1}{2}}} \right) = \sum_{i \in A'} (8\sqrt{2}K) \end{aligned}$$

351 Again we obtain (a) by substituting the value of  $\rho_s$  and  $\psi$ .

352 Summing the bounds in Cases a1 – a2 and observing that the bounds in the aforementioned cases  
 353 hold for any round  $C_{\max\{m_i, g_{s_k}\}}$ , we obtain the following contribution to the expected regret from  
 354 case a:

$$\sum_{i \in A'_{s^*}} 8\sqrt{2}K + \sum_{i \in A'} 8\sqrt{2}K \leq \sum_{i \in A'_{s^*}} 12K + \sum_{i \in A'} 12K$$

355 **Case b:** For each arm  $i$ , either  $i$  is eliminated in round  $\max(m_i, g_{s_k})$  or before or there is no optimal  
 356 arm  $*$  in  $C_{\max(m_i, g_{s_k})}$ .

357 **Case b1:**  $*$   $\in C_{\max(m_i, g_{s_k})}$  for each arm  $i \in A'$  and cluster  $s_k \in \tilde{A}$ . The condition in the case  
 358 description above implies the following:

359 (i) each sub-optimal arm  $i \in A'$  is eliminated on or before  $\max(m_i, g_{s_k})$  and hence pulled not more  
 360 than  $n_{m_i}$  number of times.

361 (ii) each sub-optimal cluster  $s_k \in \tilde{A}$  is eliminated on or before  $\max(m_i, g_{s_k})$  and hence pulled not  
 362 more than  $n_{g_{s_k}}$  number of times.

363 Hence, the maximum regret suffered due to pulling of a sub-optimal arm or a sub-optimal cluster is  
 364 no more than the following:

$$\begin{aligned} &\sum_{i \in A'} \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil + \sum_{k=1}^p \sum_{i \in A'_{s_k}} \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{g_{s_k}}^2)}{2\epsilon_{g_{s_k}}} \right\rceil \\ &\stackrel{a}{\leq} \sum_{i \in A'} \Delta_i \left( 1 + \frac{16 \log\left(\psi T \left(\frac{\Delta_i}{2}\right)^4\right)}{\Delta_i^2} \right) + \sum_{i \in A'} \Delta_i \left( 1 + \frac{16 \log\left(\psi T \left(\frac{\Delta_i}{2}\right)^4\right)}{\Delta_i^2} \right) \\ &\stackrel{b}{\leq} \sum_{i \in A'} \left[ 2\Delta_i + \frac{16(\log(\frac{T^2}{K^2} \frac{\Delta_i^4}{1024}) + \log(\frac{T^2}{K^2} \frac{\Delta_i^4}{1024}))}{\Delta_i} \right] \leq \sum_{i \in A'} \left[ 2\Delta_i + \frac{32 \left( \log(\frac{T\Delta_i^2}{K^2}) + \log(\frac{T\Delta_i^2}{K^2}) \right)}{\Delta_i} \right] \end{aligned}$$

365 In the above, the (a) follows since  $\sqrt{2\epsilon_{m_i}} < \frac{\Delta_i}{4}$  and  $\sqrt{2\epsilon_{n_{g_{s_k}}}} < \frac{\Delta_{a_{\max_{s_k}}}}{4}$  and (b) is obtained by  
 366 substituting the values of  $\rho_a, \rho_s$  and  $\psi$ .

367 **Case b2:**  $*$  is eliminated by some sub-optimal arm in  $s^*$   
 368 Optimal arm  $*$  can get eliminated by some sub-optimal arm  $i$  only if arm elimination condition holds,  
 369 i.e.,

$$\hat{r}_i - c_{m_i} > \hat{r}^* + c_{m_i},$$

where, as mentioned before,  $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}}}$ . From analysis in Case a1, notice that, if (1) holds in conjunction with the above, arm  $i$  gets eliminated. Also, recall from Case a1 that the events complementary to (1) have low-probability and can be upper bounded by  $\frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}}$ . Moreover, a sub-optimal arm that eliminates  $*$  has to survive until round  $m_*$ . In other words, all arms  $j \in s^*$  such that  $m_j < m_*$  are eliminated on or before  $m_*$  (this corresponds to case b1). Let, the arms surviving till  $m_*$  round be denoted by  $A'_{s^*}$ . This leaves any arm  $a_b$  such that  $m_b \geq m_*$  to still survive and eliminate arm  $*$  in round  $m_*$ . Let, such arms that survive  $*$  belong to  $A''_{s^*}$ . Also maximal regret per step after eliminating  $*$  is the maximal  $\Delta_j$  among the remaining arms in  $A''_{s^*}$  with  $m_j \geq m_*$ . Let  $m_b = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_b}{4}\}$ . Hence, the maximal regret after eliminating the arm  $*$  is upper bounded by,

$$\begin{aligned}
& \sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{\substack{i \in A''_{s^*}: \\ m_i \geq m_*}} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} \right) \cdot T \max_{\substack{j \in A''_{s^*}: \\ m_j \geq m_*}} \Delta_j \\
& \leq \sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{i \in A''_{s^*}: m_i \geq m_*} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} \right) \cdot T \cdot 4 \sqrt{2\epsilon_{m_*}} \\
& \leq \sum_{m_*=0}^{\max_{j \in A'_{s^*}} m_j} \sum_{i \in A''_{s^*}: m_i \geq m_*} 8\sqrt{2} \left( \frac{T^{1-\rho_a}}{\psi^{\rho_a} \epsilon_{m_*}^{\rho_a - \frac{1}{2}}} \right) \\
& \leq \sum_{i \in A''_{s^*}: m_i \geq m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left( \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} \right) \\
& \leq \sum_{i \in A'_{s^*}} \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_b}} \\
& \leq \sum_{i \in A'_{s^*}} \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} \Delta_i^{2\rho_a - 1}} + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} b^{2\rho_a - 1}} \\
& \leq \sum_{i \in A'_{s^*}} 16K + \sum_{i \in A''_{s^*} \setminus A'_{s^*}} 16K
\end{aligned}$$

**Case b3:**  $s^*$  is eliminated by some sub-optimal cluster. Let  $C'_g = \{a_{\max_{s_k}} \in A' | \forall s_k \in S\}$  and  $C''_g = \{a_{\max_{s_k}} \in A'' | \forall s_k \in S\}$ . A sub-optimal cluster  $s_k$  will eliminate  $s^*$  in round  $g_*$  only if the cluster elimination condition of Algorithm 1 holds, which is the following when  $*$   $\in C_{g_*}$ :

$$\hat{r}_{a_{\max_{s_k}}} - c_{g_*} > \hat{r}^* + c_{g_*}. \quad (2)$$

Notice that when  $*$   $\notin C_{g_*}$ , since  $r_{a_{\max_{s_k}}} > r^*$ , the inequality in (2) has to hold for cluster  $s_k$  to eliminate  $s^*$ . As in case b2, the probability that a given sub-optimal cluster  $s_k$  eliminates  $s^*$  is upper bounded by  $\frac{2}{(\psi T \epsilon_{g_{s^*}})^{\rho_s}}$  and all sub-optimal clusters with  $g_{s_j} < g_*$  are eliminated before round  $g_*$ . This leaves any arm  $a_{\max_{s_b}}$  such that  $g_{s_b} \geq g_*$  to still survive and eliminate arm  $*$  in round  $g_*$ . Let, such arms that survive  $*$  belong to  $C''_g$ . Hence, following the same way as case b2, the maximal regret after eliminating  $*$  is,

$$\sum_{g_*=0}^{\max_{a_{\max_{s_j}} \in C'_g} g_{s_j}} \sum_{\substack{a_{\max_{s_k}} \in C''_g: \\ g_{s_k} \geq g_*}} \left( \frac{2}{(\psi T \epsilon_{g_{s^*}})^{\rho_s}} \right) T \max_{\substack{a_{\max_{s_j}} \in C''_g: \\ g_{s_j} \geq g_*}} \Delta_{a_{\max_{s_j}}}$$

389 Using  $A' \supset C'_g$  and  $A'' \supset C''_g$ , we can bound the regret contribution from this case in a similar  
 390 manner as Case b2 as follows:

$$\begin{aligned} & \sum_{i \in A' \setminus A'_{s^*}} \frac{T^{1-\rho_s} 2^{\rho_s + \frac{7}{2}}}{\psi^{\rho_s} \Delta_i^{2\rho_s - 1}} + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} \frac{T^{1-\rho_s} 2^{\rho_s + \frac{7}{2}}}{\psi^{\rho_s} b^{2\rho_s - 1}} \\ &= \sum_{i \in A' \setminus A'_{s^*}} 16K + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} 16K \end{aligned}$$

391 **Case b4:**  $*$  is not in  $C_{\max(m_i, g_{s_k})}$ , but belongs to  $B_{\max(m_i, g_{s_k})}$ .

392 In this case the optimal arm  $*$  in  $s^*$  is not eliminated, also  $s^*$  is not eliminated. So, for all sub-  
 393 optimal arms  $i$  in  $A'_{s^*}$  which gets eliminated on or before  $\max\{m_i, g_{s_k}\}$  will get pulled no more  
 394 than  $\left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil$  number of times, which leads to the following bound the contribution to the  
 395 expected regret, as in Case b1:

$$\sum_{i \in A'_{s^*}} \left\{ \Delta_i + \frac{32 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right\}$$

396 For arms  $a_i \notin s^*$ , the contribution to the regret cannot be greater than that in Case b3. So the regret  
 397 is bounded by,

$$\sum_{i \in A' \setminus A'_{s^*}} 16K + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} 16K$$

398 The main claim follows by summing the contributions to the expected regret from each of the cases  
 399 above.

## 400 D Proof of Corollary 1

401 **Proof 8** First we recall the definition of Theorem 2 below,

$$\begin{aligned} \mathbb{E}[R_T] &\leq \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} \left\{ \Delta_i + 12K + \frac{32 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right\} + \sum_{\substack{i \in A, \\ \Delta_i > b}} \left\{ 2\Delta_i + 12K + \frac{64 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right\} \\ &+ \sum_{\substack{i \in A_{s^*}, \\ \Delta_i > b}} 16K + \sum_{\substack{i \in A_{s^*}, \\ 0 < \Delta_i \leq b}} 16K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ \Delta_i > b}} 32K + \sum_{\substack{i \in A \setminus A_{s^*}, \\ 0 < \Delta_i \leq b}} 32K + \max_{i: \Delta_i \leq b} \Delta_i T \end{aligned}$$

402 Now we know from Bubeck et al. (2011) that the function  $x \in [0, 1] \mapsto x \exp(-Cx^2)$  is decreasing  
 403 on  $\left[\frac{1}{\sqrt{2C}}, 1\right]$  for any  $C > 0$ . So, taking  $C = \left\lfloor \frac{T}{e} \right\rfloor$  and by choosing  $\Delta_i = \Delta = \sqrt{\frac{K \log K}{T}} > \sqrt{\frac{e}{T}}$   
 404 for all  $i : i \neq * \in A$  and substituting  $p = \left\lceil \frac{K}{\log K} \right\rceil$  in the bound of ClusUCB we get,

$$\sum_{i \in A_{s^*} : \Delta_i > b} 12K = 12 \frac{K^2}{p}$$

405 Similarly, for the term,

$$\sum_{i \in A : \Delta_i > b} 12K = 12K^2$$

406 For the term regarding number of pulls,

$$\sum_{i \in A: \Delta_i > b} \frac{64 \log(\frac{T \Delta_i^2}{K})}{\Delta_i} \leq \frac{64K \sqrt{T} \log(T \frac{K \log K}{TK})}{\sqrt{K \log K}} \leq \frac{64\sqrt{KT} \log(\log K)}{\sqrt{\log K}} \\ \stackrel{(a)}{\leq} 64\sqrt{KT}$$

407 Here (a) is obtained by the identity  $\frac{\log \log K}{\sqrt{\log K}} < 1$  for  $K \geq 2$ . Lastly we can bound the error terms  
408 as,

$$\sum_{i \in A_{s^*}: 0 \leq \Delta_i \leq b} 16K = \frac{16K^2}{p} \stackrel{<}{(a)} 16K \log K$$

409 Here we obtain (a) by substituting the value of  $p$ . Similarly for the term,

$$\sum_{i \in A \setminus A_{s^*}: \Delta_i > b} 16K = \frac{16K^2}{p} < 16K \log K$$

410 Also, for all  $b \geq \sqrt{\frac{e}{T}}$ ,

$$\sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} 32K = \left(K - \frac{K}{p}\right) 32K$$

411 Now,  $K - \frac{K}{p} = K \left(\frac{p-1}{p}\right) < K \left(\frac{\frac{K}{\log K} + 1 - 1}{\frac{K}{\log K} + 1}\right) < \frac{K^2}{K + \log K}$ . So, after substituting the  
412 value of  $p = \left\lceil \frac{K}{\log K} \right\rceil$ , we get,

$$\sum_{i \in A \setminus A_{s^*}: 0 < \Delta_i \leq b} 32K = \left(K - \frac{K}{p}\right) 32K < \frac{32K^3}{K + \log K}$$

413 Summing up all the contribution from the individual cases as shown above, the total gap-independent  
414 regret is given by,

$$\mathbb{E}[R_T] \leq 12K \log K + 32\sqrt{KT} + 12K^2 + 64\sqrt{KT} + 32K \log K + \frac{64K^3}{K + \log K}$$

415 So, the total bound for using both arm and cluster elimination cannot be worse than,

$$\mathbb{E}[R_T] \leq 96\sqrt{KT} + 12K^2 + 44K \log K + \frac{64K^3}{K + \log K}$$

## 416 E Proof of Proposition 1

417 **Proof 9** Let  $p = 1$  such that all the arms in  $A$  belongs to a single cluster. Hence, in ClusUCB-  
418 AE there is only arm elimination and no cluster elimination. Let, for each sub-optimal arm  $i$ ,  
419  $m_i = \min \{m | \sqrt{\epsilon_m} < \frac{\Delta_i}{2}\}$ . Also  $\rho_a = \frac{1}{2}$  is a constant in this proof. Let  $A' = \{i \in A : \Delta_i > b\}$   
420 and  $A'' = \{i \in A : \Delta_i > 0\}$ .



421 **Case a: Some sub-optimal arm  $i$  is not eliminated in round  $m_i$  or before and the optimal arm**  
 422  $* \in B_{m_i}$

423 Following the steps of Theorem 2 Case a1, an arbitrary sub-optimal arm  $i \in A'$  can get eliminated  
 424 only when the event,

$$\hat{r}_i \leq r_i + c_{m_i} \text{ and } \hat{r}^* \geq r^* - c_{m_i} \quad (3)$$

425 takes place. So to bound the regret we need to bound the probability of the complementary event of  
 426 these two conditions. Note that  $c_{m_i} = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}}}$ . A sub-optimal arm  $i$  will get eliminated in  
 427 the  $m_i$ -th round because  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  and substituting this in  $c_{m_i}$  and applying Lemma 2  
 428 we get,  $c_{m_i} < \frac{\Delta_i}{4}$ . Hence, for a sub-optimal arm  $i \in A'$ ,

$$\hat{r}_i + c_{m_i} \leq r_i + 2c_{m_i} < r_i + \Delta_i - 2c_{m_i} \leq r^* - 2c_{m_i} \leq \hat{r}^* - c_{m_i}$$

429 Applying Chernoff-Hoeffding bound and considering independence of complementary of the two  
 430 events in 3,

$$\mathbb{P}\{\hat{r}_i \geq r_i + c_{m_i}\} \leq \exp(-2c_{m_i}^2 n_{m_i}) \leq \exp(-2 * \frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_{m_i}} * n_{m_i}) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}}$$

431 Similarly,  $\mathbb{P}\{\hat{r}^* \leq r^* - c_{m_i}\} \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}}$ . Summing the two up, the probability that a sub-optimal  
 432 arm  $i$  is not eliminated on or before  $m_i$ -th round is  $\left(\frac{2}{(\psi T \epsilon_{m_i})^{\rho_a}}\right)$ .

433 Summing up over all arms in  $A'$  and bounding the regret for each arm  $i \in A'$  trivially by  $T\Delta_i$ , we  
 434 obtain

$$\begin{aligned} \sum_{i \in A'} \left( \frac{2T\Delta_i}{(\psi T \epsilon_{m_i})^{\rho_a}} \right) &\leq \sum_{i \in A'} \left( \frac{2T\Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_a}} \right) \leq \sum_{i \in A'} \left( \frac{2^{1+5\rho_a} T^{1-\rho_a} \Delta_i}{\psi^{\rho_a} \Delta_i^{2\rho_a}} \right) \leq \sum_{i \in A'} \left( \frac{2^{1+5\rho_a} T^{1-\rho_a}}{\psi^{\rho_a} \Delta_i^{2\rho_a-1}} \right) \\ &\stackrel{(a)}{\leq} \sum_{i \in A'} \leq 8\sqrt{2}K \end{aligned}$$

435 Here, (a) is obtained by substituting the values of  $\psi$  and  $\rho_a$ .

436 **Case b: Either an arm  $i$  is eliminated in round  $m_i$  or before or else there is no optimal arm**  
 437  $* \in B_{m_i}$

438 **Case b1:  $* \in B_{m_i}$  and each  $i \in A'$  is eliminated on or before  $m_i$**

439 Since we are eliminating a sub-optimal arm  $i$  on or before round  $m_i$ , it is pulled no longer than,

$$\left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil$$

440 So, the total contribution of  $i$  till round  $m_i$  is given by,

$$\begin{aligned} \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil &\leq \Delta_i \left\lceil \frac{\log(\psi T (\frac{\Delta_i}{4\sqrt{2}})^4)}{(\frac{\Delta_i}{4\sqrt{2}})^2} \right\rceil, \text{ since } \sqrt{2\epsilon_{m_i}} < \frac{\Delta_i}{4} \\ &\stackrel{(a)}{\leq} \Delta_i \left( 1 + \frac{32 \log(\frac{T}{K^2} T (\Delta_i)^4)}{\Delta_i^2} \right) \leq \Delta_i \left( 1 + \frac{32 \log(\frac{T \Delta_i^2}{K})}{\Delta_i^2} \right) \end{aligned}$$

441 In the above case, (a) is obtained by substituting the values of  $\psi$  and  $\rho_a$ . Summing over all arms in  
 442  $A'$  the total regret is given by,

$$\sum_{i \in A'} \Delta_i \left( 1 + \frac{32 \log \left( \frac{T \Delta_i^2}{K} \right)}{\Delta_i^2} \right)$$

443 **Case b2: Optimal arm  $*$  is eliminated by a sub-optimal arm**

444 Firstly, if conditions of Case a holds then the optimal arm  $*$  will not be eliminated in round  $m = m_*$   
 445 or it will lead to the contradiction that  $r_i > r^*$ . In any round  $m_*$ , if the optimal arm  $*$  gets eliminated  
 446 then for any round from 1 to  $m_j$  all arms  $j$  such that  $m_j < m_*$  were eliminated according to  
 447 assumption in Case a. Let the arms surviving till  $m_*$  round be denoted by  $A'$ . This leaves any arm  
 448  $a_b$  such that  $m_b \geq m_*$  to still survive and eliminate arm  $*$  in round  $m_*$ . Let such arms that survive  
 449  $*$  belong to  $A'$ . Also maximal regret per step after eliminating  $*$  is the maximal  $\Delta_j$  among the  
 450 remaining arms  $j$  with  $m_j \geq m_*$ . Let  $m_b = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_b}{4}\}$ . Hence, the maximal regret after  
 451 eliminating the arm  $*$  is upper bounded by,

$$\begin{aligned} & \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} \right) \cdot T \max_{j \in A'' : m_j \geq m_*} \Delta_j \\ & \leq \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} \right) \cdot T \cdot 4\sqrt{2}\sqrt{\epsilon_{m_*}} \\ & \leq \sum_{m_*=0}^{\max_{j \in A'} m_j} \sum_{i \in A'' : m_i > m_*} 8\sqrt{2} \left( \frac{T^{1-\rho_a}}{\psi^{\rho_a} \epsilon_{m_*}^{\rho_a - \frac{1}{2}}} \right) \\ & \leq \sum_{i \in A'' : m_i > m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left( \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} \right) \\ & \leq \sum_{i \in A'} \left( \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} \right) + \sum_{i \in A'' \setminus A'} \left( \frac{8\sqrt{2} T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_b}} \right) \\ & \leq \sum_{i \in A'} \left( \frac{4T^{1-\rho_a} * 2^{\rho_a - \frac{1}{2}}}{\psi^{\rho_a} \Delta_i^{8\sqrt{2}\rho_a - 1}} \right) + \sum_{i \in A'' \setminus A'} \left( \frac{8\sqrt{2} T^{1-\rho_a} * 2^{\rho_a - \frac{1}{2}}}{\psi^{\rho_a} b^{2\rho_a - 1}} \right) \\ & \leq \sum_{i \in A'} \left( \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} \Delta_i^{2\rho_a - 1}} \right) + \sum_{i \in A'' \setminus A'} \left( \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} b^{2\rho_a - 1}} \right) \\ & \stackrel{(a)}{\leq} \sum_{i \in A'} 16K + \sum_{i \in A'' \setminus A'} 16K \end{aligned}$$

452 Again (a) is obtained by substituting the values of  $\psi$  and  $\rho_a$ . Summing up **Case a** and **Case b**, the  
 453 total regret is given by,

$$\mathbb{E}[R_T] \leq \sum_{i \in A : \Delta_i > b} \left\{ 12K + \left( \Delta_i + \frac{32 \log \left( \frac{T \Delta_i^2}{K} \right)}{\Delta_i} \right) + 16K \right\} + \sum_{i \in A : 0 < \Delta_i \leq b} 16K + \max_{i \in A : \Delta_i \leq b} \Delta_i T$$

## 454 F Proof of Lemma 3

455 **Lemma 3** If  $m_i = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $c'_i = \sqrt{\frac{\log(\psi T \epsilon_{m_i})}{z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then,

$$\mathbb{P}\{c^{*'} > c'_i\} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}.$$

456 **Proof 10** From the definition of  $c'_i$  we know that  $c'_i \propto \frac{1}{z_i}$  as  $\psi$  and  $T$  are constants. Therefore in the  
 457  $m_i$ -th round,

$$\mathbb{P}\{c^{*'} > c'_i\} \leq \mathbb{P}\{z^* < z_i\} \leq \sum_{m=0}^{m_i} \sum_{z_i=1}^{n_{m_i}} \sum_{z^*=1}^{n_{m_i}} \left( \mathbb{P}\{\hat{r}^* < r^* + c^{*'}\} + \mathbb{P}\{\hat{r}_i > r_i + c'_i\} \right)$$

458 Now, applying Chernoff-Hoeffding bound we can show that,

$$\begin{aligned} \mathbb{P}\{\hat{r}^* < r^* + c^{*'}\} &\leq \exp(-2(c^{*'})^2 n^*) \leq \frac{1}{(\psi T \epsilon_{m_i})^2} \\ \mathbb{P}\{\hat{r}_i > r_i + c'_i\} &\leq \exp(-2(c'_i)^2 n_i) \leq \frac{1}{(\psi T \epsilon_{m_i})^2} \end{aligned}$$

459 Hence, summing everything up,

$$\mathbb{P}\{c^{*'} > c'_i\} \leq \sum_{m=0}^{m_i} \sum_{z_i=1}^{n_{m_i}} \sum_{z^*=1}^{n_{m_i}} \frac{2}{(\psi T \epsilon_{m_i})^2} \leq \frac{2T}{(\psi T \epsilon_{m_i})^2} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}.$$

## 460 G Proof of Lemma 4

461 **Lemma 4** If  $m_i = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_i}{4}\}$ ,  $c'_i = \sqrt{\frac{\log(\psi T \epsilon_{m_i})}{z_i}}$  and  $n_{m_i} = \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}}$  then,

$$\mathbb{P}\{z_i > n_{m_i}\} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}.$$

462 **Proof 11** Following a similar argument as in Lemma 3, we can show that in the  $m_i$ -th round,

$$\begin{aligned} \mathbb{P}\{z_i > n_{m_i}\} &\leq \sum_{m=0}^{m_i} \sum_{z_i=1}^{n_{m_i}} \sum_{z^*=1}^{n_{m_i}} \left( \mathbb{P}\{\hat{r}^* < r^* + c^{*'}\} + \mathbb{P}\{\hat{r}_i > r_i + c'_i\} \right) \\ &\leq \sum_{m=0}^{m_i} \sum_{z_i=1}^{n_{m_i}} \sum_{z^*=1}^{n_{m_i}} \frac{2}{(T \psi \epsilon_{m_i})^2} \leq \frac{2T}{(T \psi \epsilon_{m_i})^2} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}. \end{aligned}$$

## 463 H Proof of Lemma 5

464 **Lemma 5** If  $T \geq K^{2.4}$  and  $\Delta_i > \sqrt{\frac{e}{T}}$  then,  $\frac{K^4}{T^2 \Delta_i^3} \leq K^{2.8}$ .

465 **Proof 12** Substituting the value of  $\Delta_i = \sqrt{\frac{e}{T}}$  in the above expression we get,

$$\frac{K^4 T^{\frac{3}{2}}}{T^2 (e)^{\frac{3}{2}}} \leq \frac{K^4}{\sqrt{T}} \stackrel{(a)}{\leq} K^{2.8}$$

466 Here, (a) happens because  $T \geq K^{2.4}$ .

## 467 I Proof of Theorem 2

468 **Proof 13** We follow the same steps as in Theorem 2. We again recall the definition of some of  
 469 our notations. Let  $A' = \{i \in A, \Delta_i > b\}$ ,  $A'' = \{i \in A, \Delta_i > 0\}$ ,  $A'_{s_k} = \{i \in A_{s_k}, \Delta_i > b\}$   
 470 and  $A''_{s_k} = \{i \in A_{s_k}, \Delta_i > 0\}$ .  $C_g$  is the cluster set containing max payoff arm from each cluster  
 471 in  $g$ -th round. The arm having the true highest payoff in a cluster  $s_k$  is denote by  $a_{\max_{s_k}}$ . Let  
 472 for each sub-optimal arm  $i \in A$ ,  $m_i = \min\{m | \sqrt{2\epsilon_m} < \frac{\Delta_i}{4}\}$  and let for each cluster  $s_k \in S$ ,  
 473  $g_{s_k} = \min\{g | \sqrt{2\epsilon_g} < \frac{\Delta_{a_{\max_{s_k}}}}{4}\}$ . Let  $\check{A} = \{i \in A' | i \in s_k, \forall s_k \in S\}$ . Moreover we define  $z_i$  as the

total number of times an arm has been pulled. The analysis proceeds by considering the contribution to the regret in each of the following cases:

**Case a:** Some sub-optimal arm  $i$  is not eliminated in round  $\max(m_i, g_{s_k})$  or before, with the optimal arm  $*$   $\in C_{\max(m_i, g_{s_k})}$ . We consider an arbitrary sub-optimal arm  $i$  and analyze the contribution to the regret when  $i$  is not eliminated in the following exhaustive sub-cases:

**Case a1:** In round  $\max(m_i, g_{s_k})$ ,  $i \in s^*$ .

Note that when the following four conditions hold, arm  $i$  gets eliminated:

$$\hat{r}_i \leq r_i + c_i \text{ and } \hat{r}^* \geq r^* - c^* \text{ and } c^{*'} \leq c_i' \text{ and } z_i \geq n_{m_i}, \quad (4)$$

where  $c_i = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2z_i}}$  and  $c_i' = \sqrt{\frac{\log(\psi T \epsilon_{m_i})}{z_i}}$ . The arm  $i$  gets eliminated because

$$\begin{aligned} \hat{r}_i + c_i &\leq r_i + 2c_i < r_i + \Delta_i - 2c_i \\ &\leq r^* - 2c_i \leq r^* - 2c^* \leq \hat{r}^* - c^*. \end{aligned}$$

In the above, we have used the fact that since  $z_i \geq n_{m_i}$ , so  $c_i = \sqrt{\epsilon_{m_i+1}} < \frac{\Delta_i}{4}$ , from Lemma 2. From the foregoing, we have to bound the events complementary to that in (4) for an arm  $i$  to not get eliminated. Considering Chernoff-Hoeffding bound this is done as follows:

$$\begin{aligned} \mathbb{P}(\hat{r}_i \geq r_i + c_i) &\leq \exp(-2c_i^2 z_i) \\ &\leq \exp(-2 * \frac{\rho_a \log(\psi T \epsilon_{m_i})}{2z_i} * z_i) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}} \end{aligned}$$

Similarly, we have  $\mathbb{P}(\hat{r}^* \leq r^* - c^*) \leq \frac{1}{(\psi T \epsilon_{m_i})^{\rho_a}}$ . Again, applying Lemma 3 we can show that

in the  $m_i$ -th round  $\mathbb{P}\{c^{*'} > c_i'\} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}$ . Similarly, applying Lemma 4 we can show that

$\mathbb{P}\{z_i > n_{m_i}\} \leq \frac{2}{T(\psi \epsilon_{m_i})^2}$ . Thus, the probability that a sub-optimal arm  $i$  is not eliminated in

any round on or before  $m_i$  is bounded above by  $\left(\frac{2}{(\psi T \epsilon_{m_i})^{\rho_a}} + \frac{4}{T(\psi \epsilon_{m_i})^2}\right)$ . Summing up over all

arms in  $A'_{s^*}$  in conjunction with a simple bound of  $T\Delta_i$  for each arm we obtain,

$$\begin{aligned} \sum_{i \in A'_{s^*}} \left( \frac{2T\Delta_i}{(\psi T \epsilon_{m_i})^{\rho_a}} + \frac{4}{T(\psi \epsilon_{m_i})^2} \right) &\leq \sum_{i \in A'_{s^*}} \left( \frac{2T\Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_a}} + \frac{4T\Delta_i}{T(\psi \epsilon_{m_i})^2} \right) \\ &\stackrel{(a)}{\leq} \sum_{i \in A'_{s^*}} \left( \frac{2T\Delta_i}{(\frac{T^2 \Delta_i^2}{K^2 32})^{\frac{1}{2}}} + \frac{4}{(\frac{T \Delta_i^3}{K^2 32})^2} \right) \leq \sum_{i \in A'_{s^*}} \left( 8\sqrt{2}K + \frac{4096K^4}{T^2 \Delta_i^3} \right) \\ &\stackrel{(b)}{\leq} \sum_{i \in A'_{s^*}} (8\sqrt{2}K + 4096K^{2.8}) \end{aligned}$$

Here, in (a) we substituted the value  $\rho_a$  and  $\psi$  and (b) is obtained by applying Lemma 5.

**Case a2:** In round  $\max(m_i, g_{s_k})$ ,  $i \in s_k$  for some  $s_k \neq s^*$ .

Following a parallel argument like in Case a1, we have to bound the following two events of arm  $a_{\max_{s_k}}$  not getting eliminated on or before  $g_{s_k}$ -th round,

$$\hat{r}_{a_{\max_{s_k}}} \geq r_{a_{\max_{s_k}}} + c_{a_{\max_{s_k}}} \text{ and } \hat{r}^* \leq r^* - c^*$$

We can prove using Chernoff-Hoeffding bounds and considering independence of events mentioned

above, that for  $c_{a_{\max_{s_k}}} = \sqrt{\frac{\rho_s \log(\psi T \epsilon_{g_{s_k}})}{2z_{a_{\max_{s_k}}}}}$  and  $z_{a_{\max_{s_k}}} > n_{g_{s_k}} = \frac{\log(\psi T \epsilon_{g_{s_k}}^2)}{2\epsilon_{g_{s_k}}}$  the probability of the

above two events is bounded by  $\left(\frac{2}{(\psi T \epsilon_{g_{s_k}})^{\rho_s}} + \frac{4}{T(\psi \epsilon_{g_{s_k}})^2}\right)$ .

Now, for any round  $g_{s_k}$ , all the elements of  $C_{\max(m_i, g_{s_k})}$  are the respective maximum payoff arms of their cluster  $s_k$ ,  $\forall s_k \in S$ , and since clusters are fixed so we can bound the maximum probability that a sub-optimal arm  $i \in A'$  and  $i \in s_k$  such that  $a_{\max_{s_k}} \in C_{g_{s_k}}$  is not eliminated on or before the  $g_{s_k}$ -th round by the same probability as above. Summing up over all  $p$  clusters and bounding the regret for each arm  $i \in A'_{s_k}$  trivially by  $T\Delta_i$ ,

$$\sum_{k=1}^p \sum_{i \in A'_{s_k}} \left( \frac{2T\Delta_i}{(\psi T \frac{\Delta_i^2}{32})^{\rho_s}} + \frac{4T\Delta_i}{T(\psi \frac{\Delta_i^2}{32})^2} \right) \leq \sum_{i \in A'} (8\sqrt{2}K + 4096K^{2.8})$$

Summing the bounds in Cases a1 – a2 and observing that the bounds in the aforementioned cases hold for any round  $C_{\max\{m_i, g_{s_k}\}}$ , we obtain the following contribution to the expected regret from case a:

$$\sum_{i \in A'_{s^*}} 4108K^{2.8} + \sum_{i \in A'} 4108K^{2.8}$$

**Case b:** For each arm  $i$ , either  $i$  is eliminated in round  $\max(m_i, g_{s_k})$  or before or there is no optimal arm  $*$  in  $C_{\max(m_i, g_{s_k})}$ .

**Case b1:**  $*$   $\in C_{\max(m_i, g_{s_k})}$  for each arm  $i \in A'$  and cluster  $s_k \in \tilde{A}$ . The condition in the case description above implies the following:

(i) each sub-optimal arm  $i \in A'$  is eliminated on or before  $\max(m_i, g_{s_k})$  and hence pulled not more than  $z_i \leq n_{m_i}$  number of times.

(ii) each sub-optimal cluster  $s_k \in \tilde{A}$  is eliminated on or before  $\max(m_i, g_{s_k})$  and hence pulled not more than  $z_{a_{\max_{s_k}}} \leq n_{g_{s_k}}$  number of times.

Hence, proceeding similarly to case b1 in Theorem 2, we can show that the maximum regret suffered due to pulling of a sub-optimal arm or a sub-optimal cluster is no more than the following:

$$\begin{aligned} & \sum_{i \in A'} \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil + \sum_{k=1}^p \sum_{i \in A'_{s_k}} \Delta_i \left\lceil \frac{\log(\psi T \epsilon_{g_{s_k}}^2)}{2\epsilon_{g_{s_k}}} \right\rceil \\ & \stackrel{a}{\leq} \sum_{i \in A'} \Delta_i \left( 1 + \frac{16 \log\left(\psi T \left(\frac{\Delta_i}{2}\right)^4\right)}{\Delta_i^2} \right) + \sum_{i \in A'} \Delta_i \left( 1 + \frac{16 \log\left(\psi T \left(\frac{\Delta_i}{2}\right)^4\right)}{\Delta_i^2} \right) \\ & \stackrel{b}{\leq} \sum_{i \in A'} \left[ 2\Delta_i + \frac{16(\log(\frac{T^2}{K^2} \frac{\Delta_i^4}{1024}) + \log(\frac{T^2}{K^2} \frac{\Delta_i^4}{1024}))}{\Delta_i} \right] \leq \sum_{i \in A'} \left[ 2\Delta_i + \frac{32 \left( \log(\frac{T\Delta_i^2}{K}) + \log(\frac{T\Delta_i^2}{K}) \right)}{\Delta_i} \right] \end{aligned}$$

In the above, the (a) follows since  $\sqrt{2\epsilon_{m_i}} < \frac{\Delta_i}{4}$  and  $\sqrt{2\epsilon_{n_{g_{s_k}}}} < \frac{\Delta_{a_{\max_{s_k}}}}{4}$  and (b) is obtained by substituting the values of  $\rho_a, \rho_s$  and  $\psi$ .

**Case b2:**  $*$  is eliminated by some sub-optimal arm in  $s^*$

Optimal arm  $*$  can get eliminated by some sub-optimal arm  $i$  only if arm elimination condition holds, i.e.,

$$\hat{r}_i - c_i > \hat{r}^* + c^*,$$

where, as mentioned before,  $c_i = \sqrt{\frac{\rho_a \log(\psi T \epsilon_{m_i})}{2n_i}}$ . From analysis in Case a1, notice that, if (4) holds in conjunction with the above, arm  $i$  gets eliminated. Also, recall from Case a1 that the events complementary to (4) have low-probability and can be upper bounded by  $\left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} + \frac{4}{T(\psi \epsilon_{m_*})^2} \right)$ .

Moreover, a sub-optimal arm that eliminates  $*$  has to survive until round  $m_*$ . In other words, all arms  $j \in s^*$  such that  $m_j < m_*$  are eliminated on or before  $m_*$  (this corresponds to case b1). Let, the arms surviving till  $m_*$  round be denoted by  $A'_{s^*}$ . This leaves any arm  $a_b$  such that  $m_b \geq m_*$  to still survive and eliminate arm  $*$  in round  $m_*$ . Let, such arms that survive  $*$  belong to  $A''_{s^*}$ . Also

527 maximal regret per step after eliminating  $*$  is the maximal  $\Delta_j$  among the remaining arms in  $A_{s^*}''$  with  
 528  $m_j \geq m_*$ . Let  $m_b = \min\{m \mid \sqrt{2\epsilon_m} < \frac{\Delta_b}{4}\}$ . Hence, the maximal regret after eliminating the arm  $*$   
 529 is upper bounded by,

$$\begin{aligned}
 & \max_{j \in A_{s^*}'} \sum_{m_*=0}^{m_j} \sum_{\substack{i \in A_{s^*}'' : \\ m_i \geq m_*}} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} + \frac{4}{T(\psi \epsilon_{m_*})^2} \right) \cdot T \max_{\substack{j \in A_{s^*}'' : \\ m_j \geq m_*}} \Delta_j \\
 & \leq \sum_{m_*=0}^{\max_{j \in A_{s^*}'} m_j} \sum_{i \in A_{s^*}'' : m_i \geq m_*} \left( \frac{2}{(\psi T \epsilon_{m_*})^{\rho_a}} + \frac{4}{T(\psi \epsilon_{m_*})^2} \right) \cdot T \cdot 4\sqrt{2\epsilon_{m_*}} \\
 & \leq \sum_{m_*=0}^{\max_{j \in A_{s^*}'} m_j} \sum_{i \in A_{s^*}'' : m_i \geq m_*} 8\sqrt{2} \left( \frac{T^{1-\rho_a}}{\psi^{\rho_a} \epsilon_{m_*}^{\rho_a - \frac{1}{2}}} + \frac{2}{\psi^2 \epsilon_{m_*}^{2 - \frac{1}{2}}} \right) \\
 & \leq \sum_{i \in A_{s^*}'' : m_i \geq m_*} \sum_{m_*=0}^{\min\{m_i, m_b\}} \left( \frac{8\sqrt{2}T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} + \frac{16\sqrt{2}}{\psi^2 2^{-\frac{3}{2}m_*}} \right) \\
 & \leq \sum_{i \in A_{s^*}'} \frac{8\sqrt{2}T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_*}} + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} \frac{8\sqrt{2}T^{1-\rho_a}}{\psi^{\rho_a} 2^{-(\rho_a - \frac{1}{2})m_b}} + \sum_{i \in A_{s^*}'} \frac{16\sqrt{2}}{\psi^2 2^{-\frac{3}{2}m_*}} + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} \frac{16\sqrt{2}}{\psi^2 2^{-\frac{3}{2}m_b}} \\
 & \leq \sum_{i \in A_{s^*}'} \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} \Delta_i^{2\rho_a - 1}} + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} \frac{T^{1-\rho_a} 2^{\rho_a + \frac{7}{2}}}{\psi^{\rho_a} b^{2\rho_a - 1}} + \sum_{i \in A_{s^*}'} \frac{46}{\psi^2 \Delta_i^{\frac{3}{2}}} + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} \frac{46}{\psi^2 b^{\frac{3}{2}}} \\
 & \stackrel{(a)}{\leq} \sum_{i \in A_{s^*}'} 16K + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} 16K + \sum_{i \in A_{s^*}'} \frac{46K^2 T^{\frac{3}{2}}}{T^2 e^{\frac{3}{2}}} + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} \frac{46K^2 T^{\frac{3}{2}}}{T^2 e^{\frac{3}{2}}} \\
 & \stackrel{(b)}{\leq} \sum_{i \in A_{s^*}'} 16K + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} 16K + \sum_{i \in A_{s^*}'} 46K + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} 46K \leq \sum_{i \in A_{s^*}'} 62K + \sum_{i \in A_{s^*}'' \setminus A_{s^*}'} 62K
 \end{aligned}$$

530 Here, in (a) we substitute  $\sqrt{\frac{e}{T}} < b < \Delta_i$  and (b) happens because  $T \geq K^{2.4}$ .

531 **Case b3:**  $s^*$  is eliminated by some sub-optimal cluster. Let  $C_g' = \{a_{\max_{s_k}} \in A' \mid \forall s_k \in S\}$  and  
 532  $C_g'' = \{a_{\max_{s_k}} \in A'' \mid \forall s_k \in S\}$ . A sub-optimal cluster  $s_k$  will eliminate  $s^*$  in round  $g_*$  only if the  
 533 cluster elimination condition of Algorithm 1 holds, which is the following when  $*$   $\in C_{g_*}$ :

$$534 \quad \hat{r}_{a_{\max_{s_k}}} - c_{a_{\max_{s_k}}} > \hat{r}^* + c^*. \quad (5)$$

534 Notice that when  $*$   $\notin C_{g_*}$ , since  $r_{a_{\max_{s_k}}} > r^*$ , the inequality in (5) has to hold for cluster  $s_k$  to  
 535 eliminate  $s^*$ . As in case b2, the probability that a given sub-optimal cluster  $s_k$  eliminates  $s^*$  is upper  
 536 bounded by  $\left( \frac{2}{(\psi T \epsilon_{g_s^*})^{\rho_s}} + \frac{4}{T(\psi \epsilon_{g_s^*})^2} \right)$  and all sub-optimal clusters with  $g_{s_j} < g_*$  are eliminated  
 537 before round  $g_*$ . This leaves any arm  $a_{\max_{s_b}}$  such that  $g_{s_b} \geq g_*$  to still survive and eliminate arm  $*$   
 538 in round  $g_*$ . Let, such arms that survive  $*$  belong to  $C_g''$ . Hence, following the same way as case b2,  
 539 the maximal regret after eliminating  $*$  is,

$$\begin{aligned}
 & \max_{\substack{j \in A_{s^*}' \\ a_{\max_{s_j}} \in C_g'}} \sum_{g_*=0}^{g_{s_j}} \sum_{\substack{a_{\max_{s_k}} \in C_g'' : \\ g_{s_k} \geq g_*}} \left( \frac{2}{(\psi T \epsilon_{g_s^*})^{\rho_s}} + \frac{4}{T(\psi \epsilon_{g_s^*})^2} \right) T \max_{\substack{a_{\max_{s_j}} \in C_g'' : \\ g_{s_j} \geq g_*}} \Delta_{a_{\max_{s_j}}}
 \end{aligned}$$

540 Using  $A' \supset C_g'$  and  $A'' \supset C_g''$ , we can bound the regret contribution from this case in a similar  
 541 manner as Case b2 as follows:

$$\sum_{i \in A' \setminus A_{s^*}'} 62K + \sum_{i \in A'' \setminus A' \cup A_{s^*}'} 62K$$

542 **Case b4:**  $*$  is not in  $C_{\max(m_i, g_{s_k})}$ , but belongs to  $B_{\max(m_i, g_{s_k})}$ .

543 In this case the optimal arm  $*$   $\in s^*$  is not eliminated, also  $s^*$  is not eliminated. So, for all sub-  
 544 optimal arms  $i$  in  $A'_{s^*}$  which gets eliminated on or before  $\max\{m_i, g_{s_k}\}$  will get pulled no more  
 545 than  $\left\lceil \frac{\log(\psi T \epsilon_{m_i}^2)}{2\epsilon_{m_i}} \right\rceil$  number of times, which leads to the following bound the contribution to the  
 546 expected regret, as in Case b1:

$$\sum_{i \in A'_{s^*}} \left\{ \Delta_i + \frac{32 \log\left(\frac{T \Delta_i^2}{K}\right)}{\Delta_i} \right\}$$

547 For arms  $a_i \notin s^*$ , the contribution to the regret cannot be greater than that in Case b3. So the regret  
 548 is bounded by,

$$\sum_{i \in A' \setminus A'_{s^*}} 62K + \sum_{i \in A'' \setminus A' \cup A'_{s^*}} 62K$$

549 The main claim follows by summing the contributions to the expected regret from each of the cases  
 550 above.