DIMINISHING EXPLORATION: A MINIMALIST AP PROACH TO PIECEWISE STATIONARY MULTI-ARMED BANDITS

Anonymous authors

Paper under double-blind review

ABSTRACT

The piecewise-stationary bandit problem is an important variant of the multi-armed bandit problem that further considers abrupt changes in the reward distributions. The main theme of the problem is the trade-off between exploration for detecting environment changes and exploitation of traditional bandit algorithms. While this problem has been extensively investigated, existing works either assume knowledge about the number of change points M or require extremely high computational complexity. In this work, we revisit the piecewise-stationary bandit problem from a minimalist perspective. We propose a novel and generic exploration mechanism, called diminishing exploration, which eliminates the need for knowledge about M and can be used in conjunction with an existing change detection-based algorithm to achieve near-optimal regret scaling. Simulation results show that despite being oblivious of M, equipping existing algorithms with proposed diminishing exploration generally achieves better empirical regret than the traditional uniform exploration.

025 026 027

028 029

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

The multi-armed bandit (MAB) problem, a classic formulation of online decision making, involves a decision-maker facing a set of arms with unknown reward distributions, and the challenge is 031 to determine a learning strategy that maximizes the cumulative reward. MAB encapsulates the fundamental trade-off between exploiting the current known best arm for immediate reward and 033 exploring other arms for potentially discovering better ones. Given the prevalence of such exploration-034 exploitation dilemma in practice, MAB has served as the abstraction of various real-world sequential decision making problems, such as recommender systems (Lu et al., 2010), communication networks (Gupta et al., 2019; Hashemi et al., 2018), and clinical trials (Aziz et al., 2021). In the MAB literature, 036 there are two popular frameworks, namely stochastic bandits and adversarial bandits (Lattimore & 037 Szepesvári, 2020): (i) In a standard stochastic bandit model (Lai et al., 1985; Auer et al., 2002a), the rewards of each arm are drawn independently from its underlying reward distribution, which is assumed stationary (i.e., remains fixed throughout the learning process). (ii) In an adversarial model, 040 the reward distribution of each arm is determined by an adversary and could change abruptly after 041 each time step. To extend and unify the above two frameworks, the *piecewise-stationary bandit* (Yu 042 & Mannor, 2009) incorporates change points into the bandit model, where the reward distribution 043 of each arm could vary abruptly and arbitrarily at each (unknown) change point and remain fixed 044 between two successive change points. Accordingly, the piecewise-stationary bandit framework serves as a more realistic setting for a broad class of applications which are neither fully stationary nor fully adversarial, such as recommender systems (Xu et al., 2020) and dynamic pricing systems 046 (Yu & Mannor, 2009). 047

Piecewise-stationary bandit has been extensively studied from various perspectives, including the passive methods, e.g., forgetting via discounting (Kocsis & Szepesvári, 2006) or a sliding window (Garivier & Moulines, 2011), and the active methods that leverage change-point detectors, e.g., (Liu et al., 2018; Cao et al., 2019). A more comprehensive survey is deferred to Section 1.1. Despite the rich literature, the existing approaches suffer from the following issues: (i) *Required knowledge of the number of change points*: To adapt exploration to the change frequency, most of the existing works require some tuning based on the knowledge about the number of change points (denoted by

059	ALG	TYPE	S.K. FREE	COMPLEXITY	R.B. $\tilde{\mathcal{O}}\left(\cdot\right)$	Reference
060	D-UCB	Р	×	$\mathcal{O}\left(KT\right)$	\sqrt{MT}	(Kocsis & Szepesvári, 2006)
062	SW-UCB	Р	×	$\mathcal{O}\left(KT\right)$	\sqrt{MT}	(GARIVIER & MOULINES, 2011)
063	D-TS	Р	×	$\mathcal{O}\left(KT\right)$	\sqrt{MT}	(QI ET AL., 2023)
064	AdSwitch	Е	\checkmark	$O(KT^4)$	\sqrt{MT}	(AUER ET AL., 2019)
065	META	Е	\checkmark	$O(KT^2)$	\sqrt{ST}	(Suk & Kpotufe, 2023)
0667	ArmSwitch	Е	\checkmark	$\mathcal{O}\left(K^2T^2\right)$	\sqrt{ST}	(Abbasi-Yadkori et al., 2023)
068	MASTER	М	\checkmark	$\mathcal{O}\left(KT\right)$	$\min\left\{\sqrt{MT},\Delta^{1/3}T^{2/3}+\sqrt{T}\right\}$	(WEI & LUO, 2021)
070	CUSUM-UCB	A	×	$O(KT^2)$	\sqrt{MT}	(LIU ET AL., 2018)
071	M-UCB	A	×	$\mathcal{O}\left(KT\right)$	\sqrt{MT}	(CAO ET AL., 2019)
072	GLR-KLUCB	A	\checkmark	$O(KT^2)$	\sqrt{MT}	(BESSON ET AL., 2022)
073	M-UCB (DE)	A	√	$\mathcal{O}\left(KT\right)$	\sqrt{MT}	OURS
074	GLR-UCB (DE)	A	\checkmark	$O(KT^2)$	\sqrt{MT}	OURS
075	M-UCB (DE) ⁺	A	√	$\mathcal{O}(KT)$	\sqrt{ST}	OURS
076	GLR-UCB (DE) ⁺	A	\checkmark	$O(KT^2)$	\sqrt{ST}	OURS

054 Table 1: A summary of the regret bounds of various algorithms (R.B.: Regret Bound, S.K.: Segment Knowledge, (DE): Diminishing exploration version, $(DE)^+$: Diminishing exploration extension version, P: Passive method, A: Active method, E: Elimination approach, M: Multiple instances). The notation S is the total number of times the optimal arm switches to another. 057

078

079

080 M throughout this paper), which could be rather difficult to obtain or estimate in practice. (ii) High computational or algorithmic complexity: AdSwitch (Auer et al., 2019) and its enhanced versions 081 (Suk & Kpotufe, 2022; Abbasi-Yadkori et al., 2023) have been proposed to achieve nearly optimal regret without knowing the number of changes by adopting an elimination approach. However, these 083 approaches are computationally costly as they either relies on a large number of calls of detection 084 mechanism (Auer et al., 2019). On the other hand, MASTER (Wei & Luo, 2021) serves as a generic 085 black-box approach utilizes multiple hierarchical instances to tackle the issue of unknown number of changes and achieve optimal regret guarantees. This additional algorithmic complexity can lead to high overhead and incur high regret, especially for a small number of segments. Moreover, 088 Gerogiannis et al. (2024) emphasizes that MASTER only becomes effective under very large time 089 horizons-a luxury that is often unattainable in practical scenarios. These motivate the need for a minimalist design for piecewise-stationary MAB.

091 In this paper, we answer the above question through the lens 092 of *diminishing exploration*. Our key insight is that one could achieve a proper trade-off between the detection delay and the 094 regret incurred by exploration if the amount of exploration is configured to decrease with the *elapsed time since the latest* 096 *detection*, even without the knowledge of M.

Specifically, the main contribution of this work is to address the 098 piecewise stationary multi-armed bandit problem from a mini-099 malist perspective, reflected in several key aspects:

100 Conceptual Simplicity of the Algorithm: The proposed algo-101 rithm is conceptually very simple, which only involves equipping 102

an active method with a novel diminishing exploration mech-



algorithm. 105

Minimal Knowledge Requirement: The algorithm operates with minimal knowledge of the environ-106 ment, requiring no prior information about the number of change points M, and still can achieve a 107

nearly-optimal regret bound of $O(\sqrt{MKT})$ due to its ability to automatically adapt to the environ-



Figure 1: Regret and computation times.

ment. This balance of minimal assumption and strong performance highlights the algorithm's ability to provide excellent regret results numerically.

Low Complexity: The algorithm maintains one of the lowest possible complexities, designed without adding extra computational burden. The exploration mechanism is efficiently scheduled, requiring only the determination of when to initiate each exploration phase.

To support our statements, we provide evidence in Figure 1, where the mean regret and computation time across different algorithms are compared. It is shown that the proposed diminishing exploration together with M-UCB can achieve the best performance in terms of both regret and complexity. Compared to MASTER, our algorithm achieves significantly lower regret and substantially reduced computational complexity within a non-asymptotic time horizon—a regime in which MASTER typically underperforms, as analyzed in Gerogiannis et al. (2024). This is an initial observation, and we will discuss it in more detail in Section 6.

120 This paper is summarized as follows: (i) We revisit piecewise-stationary bandit problems without 121 the knowledge of the number of changes through the lens of diminishing exploration, which is 122 parameter-free, computationally efficient, and compatible with various change detection methods and 123 bandit algorithms. (ii) We provide a general form that allows any change detector to be combined with 124 diminishing exploration and formally show that M-UCB and GLR-UCB equipped with the proposed 125 diminishing exploration scheme under a properly chosen scheduler enjoy $\mathcal{O}(\sqrt{MKT})$ regret bound. 126 Therefore, the proposed algorithm is nearly optimal in terms of dynamic regret without knowing M. (iii) Through extensive simulations, we corroborate the regret performance of the proposed algorithm 127 and show that it outperforms the existing benchmark methods in empirical regret. 128

129 A summary of the algorithms, which will be reviewed in what follows, and their performance, along 130 with the one proposed in this work, is provided in Table 1. In this table, we highlight the performance 131 of the proposed M-UCB (DE) as evidence of our claim of being an minimalist approach. As shown in the table, it is a nearly optimal algorithm with low computational complexity. In addition, it requires 132 the least knowledge about the environment, offering versatility and ease of extension to M-UCB 133 $(DE)^+$. Last but not least, the performance of GLR-UCB (DE) and its extension GLR-UCB $(DE)^+$ 134 are also provided as an example to demonstrate that the proposed DE can be used in conjunction with 135 active methods, other than M-UCB. 136

137 138 1.1 Related Work

139 Piecewise-Stationary Bandits With Knowledge of Number of Changes. The existing bandit 140 algorithms for the piecewise-stationary setting could be largely divided into two categories: (i) 141 Passive methods: The forgetting mechanism is one widely adopted technique to tackle piecewise 142 stationarity without explicitly detecting the change points. For example, Discounted UCB (Kocsis & 143 Szepesvári, 2006) and Sliding-Window UCB (Garivier & Moulines, 2011) are two important variants 144 of UCB-type algorithms with respective forgetting mechanism, and they both have been shown to 145 achieve $\mathcal{O}(K\sqrt{MT}\log T)$ dynamic regret. Moreover, (Raj & Kalyani, 2017) propose Discounted Thompson Sampling (DTS), which adapts the discounting technique to the Bayesian setup and 146 enjoys good empirical performance despite the lack of any theoretical guarantee. Subsequently, (Qi 147 et al., 2023) provides valuable insights into DTS with theoretical guarantees, demonstrating that the 148 DTS method can achieve $\mathcal{O}(K\sqrt{MT}\log^2 T)$ regret. By adapting the methods originally designed 149 for adversarial bandits, RExp (Besbes et al., 2014) offers another passive strategy by augmenting 150 the classic Exp3 algorithm (Auer et al., 2002b) with restarts and achieve $\mathcal{O}((K \log KV_T)^{1/3}T^{2/3})$ 151 regret, where V_T denotes the total variation budget up to T. (ii) Active methods: MAB algorithms 152 augmented with a change-point detector have been explored quite extensively. One example is the 153 Windowed-Mean Shift Detection (WMD) algorithm (Yu & Mannor, 2009), which offers a generic 154 framework of combining change detectors and standard MAB methods. That said, WMD is designed 155 specifically for the setting with additional side information about the rewards of the unplayed arms 156 and hence is not applicable to the standard piecewise-stationary setting. Allesiardo & Féraud (2015) 157 propose Exp3.R, which augments Exp3 with a change detector that resets Exp3 and thereby achieves 158 $\mathcal{O}(NK\sqrt{T\log T})$ with N denoting the number of changes of the best arm (N = M in the worst 159 case). More recently, Liu et al. (2018) propose change-detection based UCB (CD-UCB), which combines UCB method with off-the-shelf change detectors, such as the classic cumulative sum 160 (CUSUM) procedure and Page-Hinkley test. Cao et al. (2019) propose M-UCB, which augments 161 UCB with a simple change detector that is activated when the sample count of an arm reaches a window size w, comparing the summations of the two halves of the samples within the window. Both CD-UCB and M-UCB could achieve the $O(\sqrt{KMT \log T})$ regret bound. On the other hand, similar ideas have also been studied from a Bayesian perspective (Mellor & Shapiro, 2013; Alami et al., 2017). However, all the methods described above rely on the assumption that M is known. Moreover, piecewise-stationary bandit has also been studied in the constrained setting (Mukherjee, 2022) and the contextual bandit setting (Chen et al., 2019; Zhao et al., 2020).

168 Piecewise-Stationary Bandits Without Knowledge of Number of Changes. Among the existing 169 works, AdSwitch (Auer et al., 2019) and GLR-klUCB (Besson et al., 2022) are the most relevant 170 to ours as they also obviate the need for the knowledge of M. Specifically, AdSwitch achieves 171 $\mathcal{O}(\sqrt{KMT\log T})$ regret via an elimination approach, but at the expense of a high computational 172 complexity incurred by a more sophisticated detection scheme. On the other hand, GLR-klUCB employs the Generalized Likelihood Ratio (GLR) test on the samples to detect changes, offering a 173 more efficient detection framework but could achieve an order-optimal regret bound only in easy 174 problem instances. In contrast, our diminishing exploration is meant to acheive optimal regret without 175 the knowledge of M nor strong assumptions on the segment length. 176

177 Recent studies on enhancing the practical effectiveness of change detection-based algorithms have 178 considered addressing significant changes without having a complete restart for every detected 179 change. When the reward distributions evolve while the optimal arm remains stable, a full restart is deemed too conservative. For instance, (Manegueu et al., 2021) proposed a change point algorithm based on empirical gaps between arms. (Suk & Kpotufe, 2022) expanded on this by quantifying 181 significant shifts at each step, avoiding reliance on non-stationarity knowledge. They employ a 182 sophisticated method to regularly re-explore suboptimal arms, ensuring optimal guarantees for both 183 piecewise-stationary and variation budget assumptions. (Abbasi-Yadkori et al., 2023) also achieved comparable guarantees in scenarios with abrupt changes, albeit with slightly diminished results. (Wei 185 & Luo, 2021) takes a multi-scale approach and maintains multiple competing instances of the base 186 algorithm, which achieves a good theoretical guarantee.

187 188

189

211 212 213

2 PROBLEM FORMULATION

Piecewise-Stationary Bandit Environment. A piecewise-stationary bandit can be described using the tuple $(\mathcal{K}, \mathcal{T}, \{f_{k,t}\}_{k \in \mathcal{K}, t \in \mathcal{T}})$, where \mathcal{K} represents a set of K arms, \mathcal{T} represents a set of T time steps, and $f_{k,t}$ represents the reward distribution of arm $k \in \mathcal{K}$ in time $t \in \mathcal{T}$. Denote by $X_{k,t}$ the reward provided by the environment at the *t*-th time step if the learner selects the *k*-th arm. This reward is drawn from $f_{k,t}$ independently of the rewards obtained in other time steps $t' \in \mathcal{T}$. At time step $t \in \mathcal{T}$, the learner selects A_t , one of the K arms as the action at this time step, and sees a reward $X_{A_t,t}$.

Unlike the stochastic bandit environment, the piecewise-stationary bandit environment has several unknown change points, at which the reward distribution will change. Let us define M as the total number of segments in this piecewise-stationary bandit environment. Mathematically, Mcan be represented as $M := 1 + \sum_{t=1}^{T-1} \mathbf{1}_{\{f_{k,t} \neq f_{k,t+1} \text{ for any } k \in \mathcal{K}\}}$, where the indicator function $\mathbf{1}_{\{f_{k,t} \neq f_{k,t+1} \text{ for any } k \in \mathcal{K}\}}$ represents the occurrence of a change. We denote the time of the *i*-th change point as ν_i , for all $i \in \{1, \dots, M-1\}$ and let $\nu_0 = 0$ and $\nu_M = T$. Moreover, we define $s_i := \nu_i - \nu_{i-1}$ as the segment length of each segment *i*. Within the *i*-th segment, for each $k \in \mathcal{K}$, the reward distribution $f_{k,t}$ are the same for all $t \in [\nu_{i-1} + 1, \nu_i]$; therefore from this point onward, we slightly abuse the notation to simply use $f_k^{(i)}$ and $\mu_k^{(i)}$ to denote the reward distribution and the corresponding expected value, respectively, for arm k in the *i*-th segment.

Regret Minimization. Similar to Liu et al. (2018); Cao et al. (2019), we adopt the *dynamic regret* as the performance metric:

$$\mathcal{R}(T) := \sum_{t=1}^{T} \max_{k \in \mathcal{K}} \mathbb{E}[X_{k,t}] - \mathbb{E}\left[\sum_{t=1}^{T} X_{A_t,t}\right].$$
(1)

214 In the context of piecewise-stationary bandit, our objective, like in other bandit problems, is to 215 minimize regret. Bandit algorithms, such as UCB, are known for solving the tension between exploration and exploitation in stationary bandit problems. The main challenge in piecewise-stationary 216 bandit lies again in the trade-off between exploration and exploitation. To illustrate, in each segment, 217 after sufficiently exploring the environment, a traditional bandit algorithm will (perhaps softly) 218 commit to the current known best arm. This commitment is difficult to break in the presence of 219 unnoticed change. Additional exploration can be invested to identify changes for resetting the 220 algorithm, which inevitably introduces additional regret. This gives rise to a new exploration and exploitation trade-off that investigates how much additional exploration should be conducted. Specifically, the more additional exploration, the quicker the changes are detected, leading to a 222 quicker reset of the algorithm. The goal of this work is to solve this tension by proposing a novel 223 exploration mechanism that can strike an optimal balance between the regret incurred by additional 224 exploration and that due from detection delay. 225

226

3 THE PROPOSED FRAMEWORK: DIMINISHING EXPLORATION

227 In this section, we provide the proposed algorithm in detail. We note that in most of the active 228 algorithms for piecewise-stationary bandit problems, such as Yu & Mannor (2009); Liu et al. (2018); 229 Cao et al. (2019), there is a (periodic or stochastic) uniform exploration scheme which spends a 230 constant fraction of time on exploration for detecting potential changes, together with a traditional 231 algorithm that is capable of attaining near-optimal tradeoff for the traditional bandit problem. We 232 referred to this type of algorithms as a change detection (CD)-based bandit algorithm. Our proposed 233 algorithm is a novel and generic exploration technique, called diminishing exploration, which can be 234 used in conjunction with a CD-based bandit algorithm.

235

236 3.1 DIMINISHING EXPLORATION

237 The motivation of the proposed diminish-238 ing exploration lies in the following two 239 observations about the uniform exploration: 1) The uniform exploration scheme spends 240 a constant fraction of time on exploration, 241 which results in a regret proportional to the 242 configured exploration rate. 2) To deter-243 mine an exploration rate that achieves the 244



Figure 2: Diminishing exploration.

optimal regret scaling, the information about the total number of segments (or change points) is required. To address the above issues, we propose the diminishing exploration scheme as follows. Let us define τ_i as the *i*-th time when the algorithm alarms a change. In the proposed method, a uniform exploration round starts at $u_{i-1}^{(j)}$ for $j \in \{1, 2, ...\}$ with $u_{i-1}^{(1)} = \tau_{i-1} + 1$. i.e., the learner chooses to pull the arm 1, 2, ..., K at time $u_{i-1}^{(j)}, u_{i-1}^{(j)} + 1, ..., u_{i-1}^{(j)} + K - 1$, respectively. The process restarts whenever a new change τ_i is detected.

251 We aim to balance the regret resulting from exploration and that associated with the performance of change detection by dynamically adjusting the exploration rate within a segment. Let $u_{i-1}^{(j)}$ be the start time of the *j*-th uniform exploration session between two consecutive alarms τ_i and τ_{i-1} . In our 252 253 approach, these sessions are designed in such a way that $u_{i-1}^{(j+1)} - u_{i-1}^{(j)}$ is greater than $u_{i-1}^{(j)} - u_{i-1}^{(j-1)}$. This means that the inter-session time within the same time segment increases with j, which in turn 254 255 256 results in reduction in the exploration rate. Specifically, for the *i*-th segment, we choose $u_i^{(1)} =$ 257 $\left[\left(\alpha - K/4\alpha\right)^2\right] \text{ and } u_i^{(j)} = \left[u_i^{(j-1)} + \frac{K}{\alpha}\sqrt{u_i^{(j-1)}} + \frac{K^2}{4\alpha^2}\right], \quad \forall \ 1 \le i \le M, \ j \ge 2, \text{ without the } j \ge 2, \text{ without the$ 258 259 knowledge of M and the parameter α will be chosen later. Clearly, we have $u_i^{(j-1)} + K < u_i^{(j)}$ for every $j \ge 2$; thus, these exploration phases will not overlap. Moreover, it is obvious that the duration 260 261 between two exploration phases $u_i^{(j)} - u_i^{(j-1)} = \mathcal{O}(\sqrt{u_i^{(j-1)}})$ increases with time as Figure 2; hence, the exploration rate decreases. Thus, we term this mechanism *diminishing exploration*. 262 263

264 265 266

3.2 INTEGRATING OFF-THE-SHELF CHANGE DETECTORS WITH DIMINISHING EXPLORATION

The proposed algorithm is given in Algorithm 1, which adopts the proposed diminishing exploration (lines 3-8) and executes traditional UCB (lines 9-12) otherwise. Moreover, the algorithm enters the change detection subroutine (lines 17-19) whenever accumulating sufficient observations for an arm (line 18).

Before concluding this section, we reemphasize that although we selected as an example to employ the change detection algorithm in Algorithm 3 (Cao et al., 2019) and 4 (Besson et al., 2022) in Algorithm 1 together with the proposed diminishing exploration, the diminishing exploration technique can, in fact, be used in conjunction with any CD-based algorithm.

274 275

276

4 REGRET ANALYSIS

277 In this section, to show the effectiveness 278 of the proposed diminishing exploration, 279 we define two sets of events to capture the behavior of the change point detec-281 tion algorithm: The false alarm events 282 are defined as $F_i := \{\tau_i < \nu_i\}, \forall 1 \leq$ $i \leq M - 1$, and $F_0 := \{\tau_0 = 0\};$ 283 284 the event that the detection delay of the *i*-th change is smaller than h_i is de-285 fined as $D_i := \{\tau_i \leq \nu_i + h_i\}, \forall 1 \leq$ 286 $i \leq M-2$, where the choice of h_i de-287 pends on the underlying CD algorithm. 10: 288 We also define $D_0 := \{\tau_0 = 0\}$, and $D_{M-1} := \{\tau_{M-1} \le T\}$. In our regret 289 11: 12: 290 analysis, we also define the following 13: 291 two quantities $\Delta_k^{(i)} := \max_{\tilde{k} \in \mathcal{K}} \left\{ \mu_{\tilde{k}}^{(i)} \right\} -$ 14: 292 $\begin{aligned} \mu_k^{(i)}, &\forall \ 1 \le i \le M, \ k \in \mathcal{K}, \text{ and } \delta_k^{(i)} := \\ \left| \mu_k^{(i+1)} - \mu_k^{(i)} \right|, \quad \forall \ 1 \le i \le M - 1, \ k \in \end{aligned}$ 15: 293 16: 17: 295 18: \mathcal{K} . Furthermore, let $\delta^{(i)} := \max_{k \in \mathcal{K}} \delta_{\iota}^{(i)}$. 296

Theorem 4.1. The Algorithm 1 can be com-

bined with a CD algorithm, which achieves

Algorithm 1 CD-UCB with diminishing exploration **Require:** Positive integer T, K and parameter α 1: Initialize $\tau \leftarrow 0, u \leftarrow \left[\left(\alpha - K/4\alpha \right)^2 \right]$ and $n_k \leftarrow$ $0 \forall k \in \mathcal{K}$: 2: for $t = 1, 2, \ldots, T$ do if $u \leq t - \tau < u + K$ then 3: $A_t \leftarrow (t - \tau) - u + 1$ 4: 5: else if $t - \tau = u + K$ then 6: $u \leftarrow \left[u + \frac{K}{\alpha} \sqrt{u} + \frac{K^2}{4\alpha^2} \right]$ 7: end if 8: for k = 1, ..., K do 9: $\text{UCB}_k \leftarrow \frac{1}{n_k} \sum_{n=1}^{n_k} Z_{k,n} + \sqrt{\frac{2\log(t-\tau)}{n_k}}$ end for $A_t \leftarrow \arg \max_{k \in \mathcal{K}} \mathrm{UCB}_k$ end if Play arm A_t and receive the reward $X_{A_t,t}$. $n_{A_t} \leftarrow n_{A_t} + 1; Z_{A_t, n_{A_t}} \leftarrow X_{A_t, t}$ if CD = True then $\tau \leftarrow t, u \leftarrow 1 \text{ and } n_k \leftarrow 0 \ \forall k \in \mathcal{K}$ end if 19: end for

 (\dot{d})

(2)

the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{i=1}^{M} \tilde{C}_{i}}_{(a)} + \underbrace{2\alpha\sqrt{MT}}_{(b)} + \underbrace{\sum_{i=1}^{M-1} \mathbb{E}\left[\tau_{i} - \nu_{i} \middle| D_{i}\overline{F}_{i}D_{i-1}\overline{F}_{i-1}\right]}_{(c)} + T \sum_{i=1}^{M} \mathbb{P}\left(F_{i} \middle| \overline{F}_{i-1}D_{i-1}\right) + T \sum_{i=1}^{M-1} \mathbb{P}\left(\overline{D}_{i} \middle| \overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right),$$

306 307 308

309

310 311

312

313

314

315

316

317

305

297

298

where $\tilde{C}_i = 8 \sum_{\Delta_k^{(i)} > 0} \frac{\log T}{\Delta_k^{(i)}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \Delta_k^{(i)}.$

To elaborate, let us look into each term in equation 2. As shown in Lemma C.3, term (a) bounds the regret of the UCB algorithm in each stationary segment, given that the CD algorithm successfully detected the previous change. Term (b) bounds the regret incurred from the diminishing exploration, as shown in Lemma C.1. The other two terms bound the regrets incurred in the phase of change detection, whose quantities would depend on the underlying CD algorithm. Specifically, term (c) corresponds to the regret associated with the detection delay while term (d) addresses the regret from unsuccessful detection and false alarm. For a more detailed proof, see Appendix C.

318 319 320

4.1 **Integration with change detectors**

In this section, we will integrate change detectors from M-UCB and GLR-UCB into the framework
 of diminishing exploration. Through theoretical analysis, we will demonstrate that diminishing
 exploration can be extended to other change detectors and achieve a nearly optimal regret bound. All
 the proofs are deferred to Appendix C.

Integration with the change detector of M-UCB. In M-UCB, change detectors are triggered when the sample count of an arm reaches a window size w. The change detector divides the samples in the window into two halves and compares the difference between the two halves' summations. If the result exceeds a threshold b, an alarm is raised. We define $h_0 := 0$ and choose $h_i = \begin{bmatrix} w (K/2\alpha + 1) \sqrt{s_i + 1} + w^2/4 (K/2\alpha + 1)^2 \end{bmatrix}$ and make the following assumption:

Assumption 4.2. The algorithm knows a lower bound $\delta > 0$ such that $\delta \le \min_i \max_{k \in \mathcal{K}} \delta_k^{(i)}$.

Note that Assumption 4.2 is Assumption 1(b) of Cao et al. (2019), which is required for the M-UCB detector to determine good w and b in regret analysis. It is worth noting that almost all schemes that actively detect changes share similar assumptions; however, different algorithms may impose distinct sets of assumptions. This assumption is mild since δ may be statistically derived from historical information. Furthermore, even if the lower bound does not hold, and we occasionally encounter changes with expected reward gaps smaller than the assumed δ , such changes may be perceived as too minor to result in significant regret. In Section 6, this fact will be verified through simulation.

With this assumption, we analyze the regret of Algorithm 3 with w and b given by

$$w = \left(4/\delta^2\right) \cdot \left[\sqrt{\log\left(2KT^2\right)} + \sqrt{\log\left(2T\right)}\right]^2,\tag{3}$$

361

362

338 339

$$b = \left[w \log \left(2KT^2/2 \right) \right]^{1/2}.$$
 (4)

Assumption 4.3. $s_i = \Omega\left(\left(\log KT + \sqrt{K \log KT}\right) \sqrt{s_{i-1}}\right)$.

In particular, if $s_i = \Theta\left(\left(\log KT + \sqrt{K\log KT}\right)^{2(1+\epsilon)}\right)$ for every *i*, Assumption 4.3 holds. This 344 345 assumption essentially posits that the changes are not overly dense, a condition that holds in many 346 practical scenarios. Simple math would then show that given D_{i-1} is true, with this assumption and 347 the proposed diminishing exploration, each arm will have at least w/2 observations before and after 348 a change point. Again, we note that similar assumptions are imposed in other algorithms that actively 349 detect changes, while different algorithms may impose different assumptions. This assumption is 350 necessary with our proof technique, which requires every change to be successfully detected with 351 high probability. In our simulations in Section 6, we will demonstrate that when this assumption is violated, all the considered active methods will experience similar performance degradation due to 352 the overly dense changes and our algorithm is not particularly vulnerable. In fact, in our simulation 353 results in Section 6, we show that our algorithms significantly outperform existing active methods, 354 even when this assumption is violated. 355

Corollary 4.4 (Regret bound of M-UCB). Algorithm 1 integrated with Algorithm 3 with the parameters in (3) and (4) achieves the expected regret upper bound of $\mathcal{O}(\sqrt{KMT \log T})$.

Integration with the change detector of GLR-UCB. In GLR-UCB, the Generalized Likelihood Ratio (GLR) test is employed on the samples to detect changes. i.e., an alarm is raised whenever the GLR statistic exceeds a threshold β given by

$$\beta = 2\mathcal{J}\left(\frac{\log\left(3T^2\right)}{2}\right) + 6\log\left(1 + \log T\right),\tag{5}$$

where the function \mathcal{J} is defined in Appendix C.1.2. Following (Besson et al., 2022), we define $h_0 := 0$

and choose
$$h_i = (\alpha, \epsilon) := \left[2\left(\frac{4}{\left(\delta^{(i)}\right)^2}\beta + 2\right)\left(\frac{K}{2\alpha} + 1\right)\sqrt{s_i + 1} + \left(\frac{4}{\left(\delta^{(i)}\right)^2}\beta + 2\right)^2\left(\frac{K}{2\alpha} + 1\right)^2 \right]$$

and make the following:

and make the following:

Assumption 4.5. $\nu_i - \nu_{i-1} \ge 2 \max\{h_i, h_{i-1}\}$ for all $i \in \{1, \dots, M\}$.

369 Under this assumption, we prove the following:

Corollary 4.6 (Regret bound of GLR-UCB). Algorithm 1 integrated with Algorithm 4 with β in (5) achieves the expected regret upper bound as $O(\sqrt{KMT \log T})$.

Discussion. In some literature, such as Liu et al. (2018) and Cao et al. (2019), the approach involves finding an exploration rate for uniform exploration that balances regret induced by exploration and detection delay, assuming knowledge of M. In GLR-UCB Besson et al. (2022), the exploration rate increases with the number of change detection alarms generated by the algorithm. Compared to other exploration mechanisms, the distinctive feature of *diminishing exploration* is its use of a variable exploration rate within the same segment. The greatest advantage of this approach lies in the fact that it does not require knowledge of M. Moreover, its complexity remains low, and it can be readily applied to other active methods.

378 5 EXTENSION TO DETECTION OF OPTIMAL ARM CHANGES 379

Let us define S as the number of *super-segments*, each of which is the time period be-380 tween two consecutive changes of the optimal arm. Mathematically, S can be repre-381 sented as $S := 1 + \sum_{t=1}^{T-1} \mathbf{1}_{\{\arg\max_{k \in \mathcal{K}} \mu_{k,t} \neq \arg\max_{k \in \mathcal{K}} \mu_{k,t+1}\}}$, where the indicator function 382 $1_{\left\{\arg\max_{k\in\mathcal{K}}\mu_{k,t}\neq\arg\max_{k\in\mathcal{K}}\mu_{k,t+1}\right\}} \text{ represents the occurrence of the optimal arm changing to an-$ 384 other one. We denote the time of the r-th occurrence of the optimal arm changing to another one as 385 ν_r^* , for all $r \in \{1, \dots, S-1\}$ and let $\nu_0^* = 0$ and $\nu_S^* = T$. Moreover, we define $s_r^* := \nu_r^* - \nu_{r-1}^*$ as 386 the super segment length of each segment r, which is the duration for which the optimal arm remains the same. The last, we define τ_r^* as the r-th time when the algorithm alarms an optimal arm changing 387 to another one. We have provided Figure 5 to visually clarify the differences from Section 4. 388

389 Similar to Section 4, we also define two sets of events 390 to capture the behavior of the change point detection 391 algorithm in the version where we only focus on the change of the optimal arm: The false alarm events are 392 defined as $F_r^* := \{\tau_r^* < \nu_r^*\}, \forall 1 \leq r \leq S-1,$ and $F_0^* := \{\tau_0^* = 0\}$; the event that the detection de-393 394 lay of the r-th change is smaller than h_r^* is defined as 395 $D_r^* := \{\tau_r^* \le \nu_r^* + h_r^*\}, \ \forall \ 1 \le r \le S - 2$, where the choice of h_r^* depends on the CD algorithm. We also de-397 fine $D_0^* := \{\tau_0^* = 0\}$, and $D_{S-1}^* := \{\tau_{S-1}^* \le T\}$. In our regret analysis, we also define the following quantities 398 399 $\Delta_{k,t} := \max_{\tilde{k} \in \mathcal{K}} \left\{ \mu_{\tilde{k},t} \right\} - \mu_{k,t}, \ \forall \ 1 \le t \le T, \ k \in \mathcal{K},$ and assume $\Delta_{\min} := \min_{k \in \mathcal{K}} \min_{1 \le t \le T} \Delta_{k,t}$ is known. 400 401 Remark 5.1. In this section, the definition of a false alarm

Algorithm 2 Skipping Mechanism

Require: Two positive integer $n_{\text{skip},k}$ and n_{skip,k^*} , $n_{\text{skip},k}$ observations $\begin{array}{l} \underset{X_{1},\ldots,X_{n_{\mathrm{skip},k}}}{X_{1},\ldots,X_{n_{\mathrm{skip},k}}} & \text{and } n_{\mathrm{skip},k^{*}} & \text{observations} \\ y_{1},\ldots,Y_{n_{\mathrm{skip},k^{*}}} & X_{n_{\mathrm{skip},k^{*}}} & \\ \underset{\ell=1}{\overset{n_{\mathrm{skip},k}}{\sum}} & \sum_{\ell=1}^{n_{\mathrm{skip},k}} X_{\ell}/n_{\mathrm{skip},k} & < \\ & \sum_{\ell=1}^{n_{\mathrm{skip},k}} Y_{\ell}/n_{\mathrm{skip},k^{*}} + \eta \text{ then} \end{array}$ 1: **if** Return True 2: 3: else 4: Return False 5: end if

differs from that in Section 4. Here, a false alarm occurs when we incorrectly claim a change in 403 the optimal arm. Consequently, even if an arm's distribution changes but the optimal arm remains 404 unchanged, any such alarm would still be considered a false alarm. 405

406 407

402

5.1 DIMINISHING EXPLORATION WITH A SKIPPING MECHANISM

408 Here, we introduce a skipping mechanism (see Algorithm 2 for the high-level concept and Algorithm 5 409 in Appendix B for details) to ignore unnecessary alarms. The algorithm takes two sets of observations 410 of size $n_{\text{skip},k}$ and n_{skip,k^*} , respectively, as inputs and checks whether the sample average of the 411 second set is larger than that of the first set by a margin of η . In Appendix B, we present the complete 412 algorithm, where the two sets are samples of our algorithm before and after an alarm of change, respectively. If Algorithm 2 returns true, then this alarm is skipped; otherwise, it declares that the 413 optimal arm has changed and resets the algorithm. Note that having a negative η would reduce 414 miss detection but may also increase false reset. On the other hand, having a positive η would 415 encourage skipping, reducing false reset but leading to higher miss detection. Besides, the optimal 416 η also depends heavily on the underlying CD algorithm, as some CD algorithms cause higher false 417 alarm rates than others. In our numerical (in Appendix F) and analytic results (Theorem C.13 in 418 Appendix C.2), we set $n_{\text{skip},k} = \mathcal{O}(\log T)$ and $n_{\text{skip},k^*} = \mathcal{O}(\log T)$ and demonstrate the effectiveness 419 of the proposed skipping mechanism with $\eta = 0$. When applying the skipping mechanism, minor 420 adjustments may be necessary depending on the specific change detector employed. For instance, 421 for some change detectors, whenever a change in an arm's distribution is identified, regardless of 422 whether it is skipped, the data within the change detector's buffer is cleared. The detailed algorithm is provided in Appendix B. 423

424 425

5.2 INTEGRATION WITH CHANGE DETECTORS

426 In this section, we will integrate change detectors from M-UCB and GLR-klUCB into the extension 427 of diminishing exploration similar to Section 4.

Integration with change detectors of M-UCB. We choose the parameter w as 428

429 430

431

$$w = \left(8/\min\left\{\delta, \Delta_{\min}\right\}^2\right) \cdot \left[\sqrt{\log\left(2KT^2\right)} + \sqrt{\log\left(2T\right)}\right]^2.$$
(6)

The selection of the remaining parameters is the same as in Section 4.

432 **Corollary 5.2** (Regret bound of M-UCB). *Combining Algorithm 1 and 3 with the parameters in Equation 4, and Equation 6 achieves the expected regret upper bound as* $O(\sqrt{KST \log T})$.

Integration with change detectors of GLR-UCB. The selection of parameters is the same as in Section 4. **Corollary 5.3** (Regret bound of GLR-UCB) Combining Algorithm 1 and 4 with β function in

Corollary 5.3 (Regret bound of GLR-UCB). Combining Algorithm 1 and 4 with β function in Equation 5 achieves the expected regret upper bound as $O(\sqrt{KST \log T})$.

The proofs for Corollary 5.2 and 5.3 are in Appendix C.2.

6 SIMULATION RESULTS

437

438

439 440

441

451

452

457

462

463

464

442 In this section, we assess the effectiveness of the proposed diminishing exploration scheme across 443 various dimensions, encompassing regret scaling in M, K, and T, regrets in synthetic environments, 444 and regrets in a real-world scenario. In addition to evaluating M-UCB (Cao et al., 2019) with 445 our diminishing exploration, we also examine a variant of CUSUM-UCB (Liu et al., 2018) that 446 incorporates diminishing exploration with CUSUM-UCB, further highlighting the efficacy of the 447 proposed exploration method. We will compare our approach with M-UCB, CUSUM-UCB, GLR-UCB, Discounted-UCB, Discounted Thompson Sampling, and MASTER. Unless stated otherwise, 448 we report the average regrets over 100 simulation trials, with a zoomed-in view of the figures provided 449 in G. Detailed configuration is provided in Appendix D. 450



Figure 3: Regret in the synthetic environments and under the Yahoo data set.

465 Regret in Each Time Step. In this simulation, we consider a multi-armed bandit problem with 466 T = 20000 time steps and M = 5. Recall that $\mu_k^{(i)}$ represents the expected value for arm k in the *i*-th 467 segment. Here, we set $\mu_{k}^{(i)} = 0.2, 0.5, 0.8$ for i with $(i + k) \mod 3 = 2, 0, 1$, respectively. Figure 3a 468 shows that for both CUSUM-UCB and M-UCB, employing diminishing exploration can effectively 469 reduce the additional regret caused by constant exploration. Moreover, M-UCB with the proposed 470 diminishing exploration achieves the lowest regret. In the figure, the change points are clearly evident 471 by observing the breakpoints in each line. The reason for the overall steeper slope of CUSUM-UCB 472 (both with and without diminishing exploration) is due to the heightened sensitivity of the CUSUM 473 detector itself, resulting in more frequent false alarms.

Regret Scaling in M. Based on the settings outlined above, with the only variation being in the parameter M, Figure 3b illustrates the dynamic regrets across various values of M. In this experiment, adjustments to the exploration parameter settings are required based on the size of M when using a constant exploration rate. However, this is not the case for the proposed diminishing exploration. The result confirms the earlier-discussed rationale that the proposed diminishing exploration can automatically adapt to the environment, resulting in the best regret performance among the algorithms.

Regret Scaling in T. In line with the setting described above, with the only variation being the parameter T, we present the dynamic regrets across different values of T. Observations similar to those made above can again be found in Figure 3c, where the proposed diminishing exploration can effectively reduce the regret.

Regret and Execution Time. Here, we compare the computation time and regret across different algorithms for various choices of M and T with other parameters same as above. Specifically, we

496

501

502

504

505

506

507

508

486 conduct experiments under three scenarioso: one where the environment changes rapidly (M = 50487 and T = 20000), one where the environment changes slowly (M = 5 and T = 20000), and one 488 where the considered time horizon is double (M = 5 and T = 40000). As shown in Figure 4a to 489 4c, despite oblivious of M, our algorithm almost always achieves the lowest computation time and 490 regret in all the scenarios. Moreover, comparing to Master+UCB, another algorithm not requiring the knowledge of M, our algorithm is always significantly better as shown in these figures. Figure 4d 491 plots the ratio of average execution time of Master+UCB to that of M-UCB with our DE for various 492 T. It is shown that the growth rate is faster than $0.5 \log T$, and it appears to become even linear in T 493 as T increases. 494



Figure 4: Regret and computation times.

time of Master+UCB to that of M-UCB (DE).

509 Regret in an Environment Built from a Real-World Dataset. We further utilize the benchmark 510 dataset publicly published by Yahoo! for evaluation. To enhance arm distinguishability in our 511 simulation, we scale up the data by a factor of 10. The number of segments is set to M = 9 and the 512 number of arms is set to K = 6. Figure 3d shows the evolution of dynamic regret. Again, we see 513 that the diminishing exploration scheme could help M-UCB, GLR-UCB and CUSUM-UCB achieve similar or better regret even without knowing M. 514

515 Scenarios when Assumptions are Violated. In Assumption 4.2, we assume the knowledge of 516 δ to select an appropriate w. We emphasize that our settings in many of the above simulations 517 actually violate this assumption. Take Figure 3a for example, w = 200 is chosen, which corresponds 518 to $\delta \approx 0.6$ when back calculating, which is much larger than the actual $\delta = 0.3$ of this scenario. 519 Assumptions 4.3 and 4.5 provide guarantees that the segment length is sufficiently long. However, 520 in the last data point of our Figure 3b (M = 100), these assumptions are clearly violated. In this case, it becomes challenging for the active methods to promptly detect every change point. For 521 algorithms like M-UCB and CUSUM-UCB, which require knowledge of M, their awareness of quick 522 changes causes them to invest more effort into change detection, leading to a very high exploration 523 rate. However, this does not always guarantee successful detection, resulting in very high regret. 524 Diminishing exploration, on the other hand, continues to decrease the exploration rate regardless of 525 M when no changes are detected. This allows more resources to be invested for UCB, which might 526 gradually adapt to the environment's changes, leading to a regret that is lower than that of uniform 527 exploration. 528

Regarding the simulations of AdSwitch, ArmSwitch, and the Meta algorithm, due to the extremely 529 long running time, we have not been able to finish the simulation for T beyond 20000. We alternatively 530 perform comparison with smaller T, whose results are presented in Appendix F. 531

532 7 CONCLUDING REMARKS

533 In this paper, we revisited the piecewise-stationary bandit problem. A novel diminishing exploration 534 mechanism, called diminishing exploration, was proposed that does not require knowledge about the number of stationary segments. When used in conjunction with the M-UCB and GLR-UCB, the 536 proposed diminishing exploration mechanism was rigorously shown to achieve a near optimal regret. 537 Extensive simulations were also provided to show the effectiveness of the proposed mechanism. Regarding the limitations, since the proposed diminishing exploration is employed together with a 538 CD algorithm, the integrated algorithm usually inherits from the CD algorithm a constraint on the length of each segment in order to guarantee near-optimal regret performance.

540 REFERENCES

547

552

580

581

- Yasin Abbasi-Yadkori, András György, and Nevena Lazić. A new look at dynamic regret for
 non-stationary stochastic bandits. *Journal of Machine Learning Research*, 24(288):1–37, 2023.
- Réda Alami, Odalric Maillard, and Raphael Féraud. Memory bandits: a Bayesian approach for
 the switching bandit problem. In *Proceedings of Conference on Neural Information Processing Systems*, 2017.
- Robin Allesiardo and Raphaël Féraud. Exp3 with drift detection for the switching bandit problem. In *IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pp. 1–7, 2015.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit
 problem. *Machine learning*, 47:235–256, 2002a.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002b.
- Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an unknown number of distribution changes. In *Proceedings of Conference on Learning Theory* (*COLT*), pp. 138–158, 2019.
- Maryam Aziz, Emilie Kaufmann, and Marie-Karelle Riviere. On multi-armed bandit designs for dose-finding clinical trials. *Journal of Machine Learning Research*, 22(1):686–723, 2021.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Stochastic multi-armed-bandit problem with non stationary rewards. *Advances in Neural Information Processing Systems*, 27, 2014.
- Lilian Besson, Emilie Kaufmann, Odalric-Ambrym Maillard, and Julien Seznec. Efficient change point detection for tackling piecewise-stationary bandits. *Journal of Machine Learning Research*, 23(1):3337–3376, 2022.
- Yang Cao, Zheng Wen, Branislav Kveton, and Yao Xie. Nearly optimal adaptive procedure with
 change detection for piecewise-stationary bandit. In *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 418–427, 2019.
- 570 Yifang Chen, Chung-Wei Lee, Haipeng Luo, and Chen-Yu Wei. A new algorithm for non-stationary
 571 contextual bandits: Efficient, optimal and parameter-free. In *Proceedings of Conference on*572 *Learning Theory (COLT)*, pp. 696–726, 2019.
- Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for switching bandit problems. In *Proceedings of International Conference on Algorithmic Learning Theory (ALT)*, pp. 174–188, 2011.
- Argyrios Gerogiannis, Yu-Han Huang, and Venugopal V. Veeravalli. Is prior-free black-box nonstationary reinforcement learning feasible?, 2024. URL https://arxiv.org/abs/2410. 13772.
 - Aditya Gopalan, Braghadeesh Lakshminarayanan, and Venkatesh Saligrama. Bandit quickest changepoint detection. *Advances in Neural Information Processing Systems*, 34:29064–29073, 2021.
- Harsh Gupta, Atilla Eryilmaz, and R Srikant. Link rate selection using constrained thompson sampling. In *Proceedings of IEEE Conference on Computer Communications (INFOCOM)*, pp. 739–747, 2019.
- Morteza Hashemi, Ashutosh Sabharwal, C Emre Koksal, and Ness B Shroff. Efficient beam alignment
 in millimeter wave systems using contextual bandits. In *Proceedings of IEEE Conference on Computer Communications (INFOCOM)*, pp. 2393–2401, 2018.
- Hoda Heidari, Michael J Kearns, and Aaron Roth. Tight policy regret bounds for improving and decaying bandits. In *IJCAI*, pp. 1562–1570, 2016.
- Emilie Kaufmann and Wouter M. Koolen. Mixture martingales revisited with applications to
 sequential tests and confidence intervals. *Journal of Machine Learning Research*, 22(246):1–44,
 2021. URL http://jmlr.org/papers/v22/18-798.html.

594 Levente Kocsis and Csaba Szepesvári. Discounted UCB. In 2nd PASCAL Challenges Workshop, volume 2, pp. 51–134, 2006. 596 Tze Leung Lai, Herbert Robbins, et al. Asymptotically efficient adaptive allocation rules. Advances 597 *in Applied Mathematics*, 6(1):4–22, 1985. 598 Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Cambridge University Press, 2020. 600 Nir Levine, Koby Crammer, and Shie Mannor. Rotting bandits. Advances in Neural Information 601 Processing Systems, 30, 2017. 602 603 Fang Liu, Joohyun Lee, and Ness Shroff. A change-detection based framework for piecewise-604 stationary multi-armed bandit problem. In Proceedings of the AAAI Conference on Artificial 605 Intelligence, volume 32, 2018. 606 Gary Lorden. Procedures for reacting to a change in distribution. The Annals of Mathematical 607 Statistics, pp. 1897–1908, 1971. 608 609 Tyler Lu, Dávid Pál, and Martin Pál. Contextual multi-armed bandits. In Proceedings of International 610 Conference on Artificial Intelligence and Statistics (AISTATS), pp. 485–492, 2010. 611 Anne Gael Manegueu, Alexandra Carpentier, and Yi Yu. Generalized non-stationary bandits. arXiv 612 preprint arXiv:2102.00725, 2021. 613 614 Joseph Mellor and Jonathan Shapiro. Thompson sampling in switching environments with bayesian 615 online change detection. In Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 442-450, 2013. 616 617 Alberto Maria Metelli, Francesco Trovo, Matteo Pirola, and Marcello Restelli. Stochastic rising 618 bandits. In Proceedings of International Conference on Machine Learning (ICML), pp. 15421– 619 15457, 2022. 620 Subhojyoti Mukherjee. Safety aware changepoint detection for piecewise iid bandits. In Proceedings 621 of Uncertainty in Artificial Intelligence (UAI), pp. 1402–1412, 2022. 622 623 Ewan S Page. Continuous inspection schemes. *Biometrika*, 41(1/2):100–115, 1954. 624 Han Qi, Yue Wang, and Li Zhu. Discounted thompson sampling for non-stationary bandit problems. 625 arXiv preprint arXiv:2305.10718, 2023. 626 627 Vishnu Raj and Sheetal Kalyani. Taming non-stationary bandits: A Bayesian approach. arXiv 628 preprint arXiv:1707.09727, 2017. 629 Joe Suk and Samory Kpotufe. Tracking most significant arm switches in bandits. In Conference on 630 Learning Theory, pp. 2160–2182. PMLR, 2022. 631 632 Joe Suk and Samory Kpotufe. Tracking most significant shifts in nonparametric contextual bandits. 633 Advances in Neural Information Processing Systems, 36, 2023. 634 Venugopal V Veeravalli and Taposh Banerjee. Quickest change detection. In Academic Press Library 635 in Signal Processing, volume 3, pp. 209-255. Elsevier, 2014. 636 637 Chen-Yu Wei and Haipeng Luo. Non-stationary reinforcement learning without prior knowledge: 638 an optimal black-box approach. In Mikhail Belkin and Samory Kpotufe (eds.), Proceedings of 639 Thirty Fourth Conference on Learning Theory, volume 134 of Proceedings of Machine Learning Research, pp. 4300-4354. PMLR, 15-19 Aug 2021. URL https://proceedings.mlr. 640 press/v134/wei21b.html. 641 642 Qunzhi Xu, Yajun Mei, and George V Moustakides. Optimum multi-stream sequential change-point 643 detection with sampling control. IEEE Transactions on Information Theory, 67(11):7627–7636, 644 2021. 645 Xiao Xu, Fang Dong, Yanghua Li, Shaojian He, and Xin Li. Contextual-bandit based personalized 646 recommendation with time-varying user interests. In Proceedings of the AAAI Conference on 647 Artificial Intelligence, volume 34, pp. 6518-6525, 2020.

- Jia Yuan Yu and Shie Mannor. Piecewise-stationary bandit problems with side observations. In Proceedings of International Conference on Machine Learning (ICML), pp. 1177–1184, 2009.
- Peng Zhao, Lijun Zhang, Yuan Jiang, and Zhi-Hua Zhou. A simple approach for non-stationary linear bandits. In *Proceedings of International Conference on Artificial Intelligence and Statistics* (AISTATS), pp. 746–755, 2020.
 - Xiang Zhou, Yi Xiong, Ningyuan Chen, and Xuefeng Gao. Regime switching bandits. Advances in Neural Information Processing Systems, 34:4542–4554, 2021.

A CHANGE DETECTION ALGORITHM

660 Change Detector of M-UCB

 The following algorithm is the change detection algorithm for M-UCB (Cao et al., 2019)

Algo	rithm 3 Change Detection of M-UCB: $CD(w, b, Z_1, \dots, Z_w)$
Req 1 1: 1 2:	uire: An even integer w , w observations Z_1, \ldots, Z_w , and a prescribed hreshold $b > 0$ f $\left \sum_{\ell=w/2+1}^{w} Z_{\ell} - \sum_{\ell=1}^{w/2} Z_{\ell} \right > b$ then Return True
3: (else
4:	Return False
5: 0	end if

In this algorithm, one requires w observations as input and check whether the difference between the sample average of the first half and that of the second half exceeds a prescribed threshold b (line 1).

675 Change Detector of GLR-UCB 676 The following definition is the ch

The following definition is the change detection algorithm for GLR-UCB (Besson et al., 2022)

Definition A.1. The Bernoulli GLR change-point detector with threshold function $\beta(n, \epsilon)$ is

$$\tau_{\delta} := \inf \left\{ n \in \mathbb{N}^* : \sup_{s \in [1,n]} \left[s \times \mathrm{kl}\left(\hat{\mu}_{1:s}, \hat{\mu}_{1:n}\right) + (n-s) \times \mathrm{kl}\left(\hat{\mu}_{s+1:n}, \hat{\mu}_{1:n}\right) \right] \ge \beta\left(n, \delta\right) \right\}$$
(7)

Algorithm 4 Change Detection of GLR-UCB: $CD(Y_1, \ldots, Y_{n_k})$

Require: Y_1, \ldots, Y_{n_k} , and a threshold function $\beta(n, \epsilon)$

1: if $\sup_{s \in [1,n]} \left[s \times \operatorname{kl}\left(\sum_{\ell=1}^{s} Y_{\ell}/s, \sum_{\ell=1}^{n} Y_{\ell}/n\right) + (n-s) \times \operatorname{kl}\left(\sum_{\ell=s+1}^{n} Y_{\ell}/(n-s), \sum_{\ell=1}^{n} Y_{\ell}/n\right) \right] \ge \beta(n, \delta)$ then 2: Return True 3: else 4: Return False 5: end if

B THE EXTENDED VERSION ALGORITHM

The skipping mechanism in more detail is as follows:

Algorithm 5 Skipping mechanism (Z, n, n_{skip}, N_I) **Require:** An array Z with elements $Z_{k,t}$ for $k \in \mathcal{K}$ and $t \in \{1, \ldots, n_k\}$; two vectors n and n_{skip} with elements n_k and $n_{\text{skip},k}$ for $k \in \mathcal{K}$; and a positive integer N_I . 1: $k^* \leftarrow \arg \max_{k \in \mathcal{K}} n_k$ 2: if $\exists k \neq k^* \in \mathcal{K}, n_{\text{skip},k} \geq N_I$ and $\sum_{\ell=1}^{n_{\text{skip},k}} Z_{k,n_k-n_{\text{skip},k}+1}/n_{\text{skip},k}$ > $\sum_{\ell=1}^{n_{\text{skip},k^*}} Z_{k,n_{k^*}-n_{\text{skip},k^*}+1}/n_{\text{skip},k^*}+\eta$ then Return False 3: 4: **else** 5: Return True 6: end if As mentioned in Section 5, some adjustments may be required depending on the specific change detector used. Here, we provide the algorithm 6 to illustrate such adjustments. The characteristics of the M-UCB detector allow for a relatively straightforward integration with the skipping mechanism; therefore, our algorithm can be simplified as algorithm 7.

Alg	orithm 6 CD-UCB with diminishing exploration and skipping mechanism
Req	uire: Positive integer T, K , and parameter α, N_I
1:	Initialize $\tau \leftarrow 0, u \leftarrow (\alpha - K/4\alpha)^2 , n_k, n_{\text{CD},A_t} \leftarrow 0 \forall k \in \mathcal{K} \text{ and skip_mode} \leftarrow \text{False:}$
2:	for $t = 1, 2,, T$ do
3:	if skip mode = True then
4:	if $A_{\text{change}} = k^*$ then
5:	$A_t \leftarrow t \mod K$
6:	else
7:	if t is even then
8:	$A_t \leftarrow k^*$
9:	else
10:	$A_t \leftarrow A_{\text{change}}$
11:	end if
12:	end if
13:	$n_{ ext{skip},A_t} \leftarrow n_{ ext{skip},A_t} + 1$
14:	else
15:	if $u \le t - \tau < u + K$ then
16:	$A_t \leftarrow (t - \tau) - u + 1$
17:	else
18:	If $t - \tau = u + K$ then
19:	$u \leftarrow \left u + \frac{\kappa}{\alpha} \sqrt{u} + \frac{\kappa}{4\alpha^2} \right $
20:	end if
21:	for $k = 1, \dots, K$ do
22.	$\text{LICB}_{k} \leftarrow \frac{1}{2} \sum^{n_{k}} Z_{k} + \sqrt{\frac{2 \log(t-\tau)}{2}}$
22.	and for
25:	A (arg max LICP.
24. 25.	$A_t \leftarrow \arg \max_{k \in \mathcal{K}} \operatorname{OCD}_k$
25. 26.	end if
27:	Play arm A_{t} and receive the reward $X_{A_{t}}$
28:	$n_{A_i} \leftarrow n_{A_i} + 1; Z_{A_i} , z_i \leftarrow X_{A_i} t$
29:	$n_{\text{CD}} A_t \leftarrow n_{\text{CD}} A_t + 1; Z_{\text{CD}} A_t n_{\text{CD}} A_t \leftarrow X_{A_t} t$
30:	if $CD = True$ then
31:	if skip_mode = False then
32:	$k^* \leftarrow \arg \max_{k \in \mathcal{K}} n_k$
33:	end if
34:	$A_{\text{change}} \leftarrow A_t$
35:	skip_mode \leftarrow True
36:	$n_{\mathrm{skip},k} \leftarrow 0 \; \forall k \in \mathcal{K}$
37:	$n_{\text{CD},k} \leftarrow 0 \ \forall k \in \mathcal{K}$
38:	end if
39:	If skip_mode = True then f(G) = f(G)
40:	If $Skip(Z, n, n_{skip}, N_I) = False then$
41:	$\tau \leftarrow \iota, u \leftarrow 1 \text{ and } n_k, n_{\text{skip},k} \leftarrow 0 \ \forall k \in \mathcal{K}$
42:	skip_mode \leftarrow raise
45: 11.	the initial formula $\sum n \dots \sum K N_{-}$ or $(A = -h^* \text{ and } \sum n \dots \sum K N_{-})$
44:	If $(A_{\text{change}} - \kappa)$ and $\sum_{k \in \mathcal{K}} n_{\text{skip},k} > K N_I$ of $(A_{\text{change}} \neq \kappa)$ and $\sum_{k \in \mathcal{K}} n_{\text{skip},k} > $ then
15.	$m \dots \leftarrow 0 \forall k \in \mathcal{K}$
45. 46.	$w_{skip,k} \leftarrow 0 \forall n \in \mathbb{N}$ skip mode \leftarrow False
40. 47.	end if
48·	end if
49:	end for
	V

Al	gorithm 7 M-UCB with diminishing exploration and skipping mechanism
Re	quire: Positive integer T, K , and parameter α
1	Initialize $\tau \leftarrow 0, u \leftarrow \left[(\alpha - K/4\alpha)^2 \right]$ and $n_k \leftarrow 0 \ \forall k \in \mathcal{K};$
2	for $t = 1, 2,, T$ do
3	if $u \leq t - \tau < u + K$ then
4	$A_t \leftarrow (t - \tau) - u + 1$
5	else
6	if $t - \tau = u + K$ then
7	$u \leftarrow \left u + \frac{K}{\alpha}\sqrt{u} + \frac{K^2}{4\alpha^2} \right $
8	end if
9	for $k = 1, \dots, K$ do
10	$\text{UCB}_k \leftarrow \frac{1}{n} \sum_{n=1}^{n_k} Z_{k,n} + \sqrt{\frac{2\log(t-\tau)}{n}}$
11	end for
12	$A_t \leftarrow \arg \max_{k \in \mathcal{K}} \mathrm{UCB}_k$
13	end if
14	Play arm A_t and receive the reward $X_{A_t,t}$.
15	$n_{A_t} \leftarrow n_{A_t} + 1; Z_{A_t, n_{A_t}} \leftarrow X_{A_t, t}$
16	if $n_{A_t} \ge w$ then
17	if $CD(w, b, Z_{A_t, n_{A_t}}, w+1)$,, $Z_{A_t, n_{A_t}}$) = True then
18	if $\text{Skip}(Z, n, n, w/2)$ = False then
19	$\tau \leftarrow t, u \leftarrow 1 \text{ and } n_k, n_{\text{skip},k} \leftarrow 0 \ \forall k \in \mathcal{K}$
20	end if
21	end if
22	end if
23	end for

C PROOF DETAIL

This appendix provides the detailed proofs of the results presented in Section 4 and Section 5. Each proof is carefully elaborated to ensure clarity and rigor.

C.1 PROOF OF SECTION 4

In this subsection, we present the proofs of Theorem 4.1 in Section 4. In what follows, the first lemma bounds the regret accumulated during the diminishing exploration part of the algorithm. We denote by $R_{\text{DE}}(\tau_{i-1}, \nu_i)$ as the regret caused by the exploration part of the algorithm and by $N_{\text{DE},k}(\tau_{i-1}, \nu_i)$ the number of times that the arm k is selected in the exploration phase from the previous alarm time to the next change point.

Lemma C.1 (Diminishing exploration regret). *If the mean values of the arms remain the same during the time interval* $[\tau_{i-1}, \nu_i)$, *then we have*

$$N_{\text{DE},k}\left(\tau_{i-1},\nu_{i}\right) \leq \frac{2\alpha\sqrt{\nu_{i}-\tau_{i-1}}}{K} + \frac{3}{2},\tag{8}$$

and

$$\mathbb{E}\left[R_{\mathrm{DE}}\left(\tau_{i-1},\nu_{i}\right)\right] \leq 2\alpha\sqrt{\nu_{i}-\tau_{i-1}} + \frac{3}{2}K.$$
(9)

Proof. Recall that $u_i^{(j)}$ is the beginning of the *j*-th uniform exploration session in the *i*-th segment. In Algorithm 1, the initial time of the first exploration session after each τ_i is given by:

$$u_i^{(1)} = \left\lceil \left(\alpha - \frac{K}{4\alpha} \right)^2 \right\rceil,\tag{10}$$

$$u_i^{(j)} = \left[u_i^{(j-1)} + \frac{K}{\alpha} \sqrt{u_i^{(j-1)}} + \frac{K^2}{4\alpha^2} \right] \ge u_i^{(j-1)} + \frac{K}{\alpha} \sqrt{u_i^{(j-1)}} + \frac{K^2}{4\alpha^2}.$$
 (11)

Based on equation 10 and equation 11, one could easily check that the sequence $u_i^{(n)}$ satisfies that for every natural number n,

$$u_i^{(n)} \ge \left(\frac{(2n-3)K}{4\alpha} + \alpha\right)^2.$$
(12)

Let $u_i^{(m)}$ be the last exploration start time in time interval $[\tau_{i-1}, \nu_i)$. Then, we have

$$\mathbb{E}\left[R_{\rm DE}\left(\nu_i - \tau_{i-1}\right)\right] \le mK.\tag{13}$$

Additionally, we have:

$$\nu_{i} - \tau_{i-1} \ge u_{i}^{(m)} \ge \left(\frac{(2m-3)K}{4\alpha} + \alpha\right)^{2} \ge \left(\frac{2\mathbb{E}\left[R_{DE}\left(\nu_{i} - \tau_{i-1}\right)\right] - 3K}{4\alpha} + \alpha\right)^{2}.$$
 (14)

Finally, based on the equation 13 and equation 14, we can conclude that:

$$\mathbb{E}\left[R_{\text{DE}}\left(\nu_{i}-\tau_{i-1}\right)\right] \leq 2\alpha\sqrt{\nu_{i}-\tau_{i-1}} - 2\alpha^{2} + \frac{3}{2}K \leq 2\alpha\sqrt{\nu_{i}-\tau_{i-1}} + \frac{3}{2}K, \quad (15)$$

and

$$N_{\text{DE},k}\left(\nu_{i} - \tau_{i-1}\right) \le \frac{2\alpha\sqrt{\nu_{i} - \tau_{i-1}}}{K} + \frac{3}{2}.$$
(16)

In the following lemma, we aim to explore how long it takes for a given arm to reach a certain number of samples through diminishing exploration.

Lemma C.2 (Samples-time steps transform). When each arm has accumulated n samples, the required time is as follows: If the counting of the n samples begins immediately after a reset, the required time is given by

$$T_{reset} \le \left(\alpha + \frac{(2n-3)K}{4\alpha} + n\right)^2 + K.$$
(17)

However, if there is a delay of t_d time steps after the reset before the counting begins, the required time is given by

$$T_{t_d} \le 2n\left(\frac{K}{2\alpha} + 1\right)\sqrt{t_d + 1} + n^2\left(\frac{K}{2\alpha} + 1\right)^2.$$
(18)

Proof. We can derive the following from Equation 11 in the proof of Lemma C.1:

$$u^{(j)} \le u^{(j-1)} + \frac{k}{\alpha}\sqrt{u^{(j-1)}} + \frac{k^2}{4\alpha^2} + 1 \le \left(\sqrt{u^{(j-1)}} + \frac{k}{2\alpha}\right)^2 + 1 \le \left(\sqrt{u^{(j-1)}} + \frac{k}{2\alpha} + 1\right)^2.$$
(19)

First, let us consider the case where the counting of n samples begins immediately after the reset. From Equation 10, we can derive the following:

$$u^{(1)} \le \left(\alpha - \frac{K}{4\alpha} + 1\right)^2.$$
(20)

Using Equation 19, we can further derive:

$$u^{(2)} \le \left(\alpha - \frac{K}{4\alpha} + 1 + \frac{K}{2\alpha} + 1\right)^2.$$
 (21)

Finally, we obtain:

$$u^{(n)} \le \left(\alpha - \frac{K}{4\alpha} + 1 + (n-1)\left(\frac{K}{2\alpha} + 1\right)\right)^2.$$
 (22)

Here, $u^{(n)}$ represents the starting time of the *n*-th exploration block. The total time required to guarantee that all K arms have been sampled n times is therefore:

$$T_{\text{reset}} = u^{(n)} + K \le \left(\alpha + \frac{(2n-3)K}{4\alpha} + n\right)^2 + K.$$
 (23)

Next, we consider how long it takes for each arm to be sampled n times after t_d time steps following the reset. We first assume that, prior to t_d , each arm has already been sampled x times. For simplicity, we assume an ideal case where the exploration block starts exactly at time $t_d + 1$. Therefore, we have:

$$u^{(x+1)} = t_d + 1. (24)$$

Since, in reality, the exploration block start time may not exactly coincide with $t_d + 1$, we account for the possibility that it could begin at a later time by considering the next exploration block's start time. This allows us to bound the non-ideal case.

Following the same approach as in the first part of the proof, we eventually obtain:

$$u^{(x+n+1)} \le \left(\sqrt{t_d+1} + n\left(\frac{K}{2\alpha} + 1\right)\right)^2.$$
(25)

Finally, we derive that the total time required for each arm to be sampled n times after t_d time steps is:

$$T_{t_d} = u^{(x+n+1)} - u^{(x+1)} \le 2n\left(\frac{K}{2\alpha} + 1\right)\sqrt{t_d + 1} + n^2\left(\frac{K}{2\alpha} + 1\right)^2.$$
 (26)

Define $R(r,s) := \sum_{t=r}^{s} \max_{k \in \mathcal{K}} \mathbb{E}[X_{k,t}] - X_{A_t,t}$ be the regret accumulated during r and s. In the next lemma, we provide an upper bound on the regret accumulated from the (i - 1)-th alarm time to the end of (i-1)-th segment, given that the previous change was successfully detected.

Lemma C.3 (Regret bound with stationary bandit). Consider a stationary bandit interval with $\nu_{i-1} < \tau_{i-1} < \nu_i$. Condition on the successful detection events F_{i-1} and D_{i-1} , the expected regret accumulated during (τ_{i-1}, ν_i) can be bounded by

$$\mathbb{E}\left[R\left(\tau_{i-1},\nu_{i}\right)\big|\overline{F}_{i-1}D_{i-1}\right] \leq \tilde{C} + 2\alpha\sqrt{s_{i}} + T \cdot \mathbb{P}\left(F_{i}\big|\overline{F}_{i-1}D_{i-1}\right),\tag{27}$$

where \tilde{C} is as described in Theorem 4.1.

Proof. For every *i*, we have

$$\mathbb{E}\left[R\left(\tau_{i-1},\nu_{i}\right)\big|\overline{F}_{i-1}D_{i-1}\right] = \mathbb{E}\left[R\left(\tau_{i-1},\nu_{i}\right)\big|F_{i}\overline{F}_{i-1}D_{i-1}\right]\mathbb{P}\left(F_{i}\big|\overline{F}_{i-1}D_{i-1}\right)$$
(28a)

$$+ \mathbb{E}\left[R\left(\tau_{i-1},\nu_{i}\right)\left|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right]\mathbb{P}\left(\overline{F}_{i}\left|\overline{F}_{i-1}D_{i-1}\right)\right.$$
(28b)

$$\leq T \cdot \mathbb{P}\left(F_{i} | \overline{F}_{i-1} D_{i-1}\right) + \mathbb{E}\left[R\left(\tau_{i-1}, \nu_{i}\right) | \overline{F}_{i} \overline{F}_{i-1} D_{i-1}\right].$$
 (28c)

Now, define $N_k(t_1, t_2) := \sum_{t=t_1}^{t_2} \mathbf{1}_{\{A_t=k\}}$ to be the number of times that arm k is selected by Algorithm 1 from t_1 to t_2 . Note that

$$\mathbb{E}\left[R\left(\tau_{i-1},\nu_{i}\right)\big|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right] = \sum_{\Delta_{k}^{(i)}>0} \Delta_{k}^{(i)} \cdot \mathbb{E}\left[N_{k}\left(\tau_{i-1},\nu_{i}\right)\big|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right].$$
(29)

To bound the second term of equation 28c, we further bound $N_k(\tau_{i-1}, \nu_i)$ as follows,

$$N_{k}(\tau_{i-1},\nu_{i}) = \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k,\tau_{i}>\nu_{i},N_{k}(\tau_{i-1},\nu_{i})< l\}} + \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k,\tau_{i}>\nu_{i},N_{k}(\tau_{i-1},\nu_{i})\geq l\}}$$
(30a)

$$\leq l + N_{DE,k} \left(\nu_{i} - \tau_{i-1} \right) + \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\left\{ k = \arg \max_{k \in \mathcal{K}} \text{UCB}_{k \in \mathcal{K}}, \tau_{i} > \nu_{i}, N_{k}(\tau_{i-1}, \nu_{i}) \geq l \right\}}$$
(30b)

$$\leq l + \frac{2\alpha\sqrt{\nu_{i} - \tau_{i-1}}}{K} + \frac{3}{2} + \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\left\{k = \arg\max_{k \in \mathcal{K}} \text{UCB}_{k \in \mathcal{K}}, \tau_{i} > \nu_{i}, N_{k}(\tau_{i-1}, \nu_{i}) \geq l\right\}}$$
(30c)

$$\leq l + \frac{2\alpha\sqrt{s_i}}{K} + \frac{3}{2} + \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\left\{k = \arg\max_{k \in \mathcal{K}} \text{UCB}_{k \in \mathcal{K}}, \tau_i > \nu_i, N_k(\tau_{i-1}, \nu_i) \geq l\right\}}, \quad (30d)$$

Equation 30c follows from Lemma C.1. Setting $l = \left[8 \log T / \left(\Delta_k^{(i)} \right)^2 \right]$, and following the same steps as in the proof of Theorem 1 of Auer et al. (2002a), we arrive at

$$\mathbb{E}\left[N_k\left(\tau_{i-1},\nu_i\right)\Big|\overline{F}_i\overline{F}_{i-1}D_{i-1}\right] \le \frac{2\alpha\sqrt{s_i}}{K} + \frac{8\log T}{\left(\Delta_k^{(i)}\right)^2} + \frac{5}{2} + \frac{\pi^2}{3} + K.$$
(31)

Putting everything together completes the proof.

1000 Theorem 4.1 can then be proved by recursively applying Lemma C.3.

Theorem 4.1 The Algorithm 1 can be combined with a CD algorithm, which achieves the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \sum_{i=1}^{M} \tilde{C}_{i} + 2\alpha\sqrt{MT} + \sum_{i=1}^{M-1} \mathbb{E}\left[\tau_{i} - \nu_{i} \left| D_{i}\overline{F}_{i}D_{i-1}\overline{F}_{i-1} \right] + T\sum_{i=1}^{M} \mathbb{P}\left(F_{i} \left|\overline{F}_{i-1}D_{i-1}\right) + T\sum_{i=1}^{M-1} \mathbb{P}\left(\overline{D}_{i} \left|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right)\right), \quad (32)$$

 where $\tilde{C}_i = 8 \sum_{\Delta_k^{(i)} > 0} \frac{\log T}{\Delta_k^{(i)}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \Delta_k^{(i)}.$

Proof. Recall that $R(r,s) = \sum_{t=r}^{s} \max_{k \in \mathcal{K}} \mathbb{E}[X_{k,t}] - X_{A_t,t}$, then $\mathcal{R}(T) = \mathbb{E}[R(1,T)]$. We have

$$\mathcal{R}(T) = \mathbb{E}[R(1,T)] \tag{33a}$$

$$= \mathbb{E}\left[R\left(1,T\right)\big|\overline{F}_{0}D_{0}\right] \tag{33b}$$

1020
1021
$$\leq \mathbb{E}\left[R\left(1,\nu_{1}\right)\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right] + \mathbb{E}\left[R\left(\nu_{1},T\right)\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right] + T \cdot \mathbb{P}\left(F_{1}\left|\overline{F}_{0}D_{0}\right.\right)\right]$$
(33c)

1022
$$\leq \tilde{C}_1 + 2\alpha\sqrt{(\nu_1 - \nu_0)} + \mathbb{E}\left[R\left(\nu_1, T\right) \middle| \overline{F}_1 \overline{F}_0 D_0\right] + T \cdot \mathbb{P}\left(F_1 \middle| \overline{F}_0 D_0\right), \tag{33d}$$
1023

where equation 33b holds because $\tau_0 = 0$, equation 33c is due to the law of total expectation and some trivial bounds, and equation 33d follows from Lemmas C.3. The third term in equation 33d is then further bounded as follows:

$$1028 \quad \mathbb{E}\left[R\left(\nu_{1},T\right)\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right] \leq \mathbb{E}\left[R\left(\nu_{1},T\right)\left|D_{1}\overline{F}_{1}\overline{F}_{0}D_{0}\right] + T \cdot \left(1 - \mathbb{P}\left(D_{1}\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right)\right)\right.$$
(34a)

$$\leq \mathbb{E} \left[R\left(\nu_{1}, T\right) \middle| D_{1} \overline{F}_{1} \overline{F}_{0} D_{0} \right] + T \cdot \mathbb{P} \left(\overline{D}_{1} \middle| \overline{F}_{1} \overline{F}_{0} D_{0} \right)$$

$$(34b)$$

$$= \mathbb{E}\left[R\left(\tau_{1},T\right)\left|D_{1}\overline{F}_{1}\overline{F}_{0}D_{0}\right] + \mathbb{E}\left[R\left(\nu_{1},\tau_{1}\right)\left|D_{1}\overline{F}_{1}\overline{F}_{0}D_{0}\right] + T \cdot \mathbb{P}\left(\overline{D}_{1}\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right)\right]$$
(34c)

1032
$$\leq \mathbb{E}\left[R\left(\tau_{1},T\right)\left|\overline{F}_{1}D_{1}\right]+\mathbb{E}\left[\tau_{1}-\nu_{1}\left|\overline{F}_{1}D_{1}\overline{F}_{0}D_{0}\right]+T\cdot\mathbb{P}\left(\overline{D}_{1}\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right)\right.$$
(34d)

$$\leq \mathbb{E}\left[R\left(\tau_{1},T\right)\left|\overline{F}_{1}D_{1}\right] + \mathbb{E}\left[\tau_{1}-\nu_{1}\left|\overline{F}_{1}D_{1}\right] + T \cdot \mathbb{P}\left(\overline{D}_{1}\left|\overline{F}_{1}\overline{F}_{0}D_{0}\right)\right],\tag{34e}$$

where equation 34a applies the law of total expectation and some trivial bounds. From here, we can set up the following recursion:

$$\mathbb{E}\left[R\left(1,T\right)\right] = \mathbb{E}\left[R\left(1,T\right)\big|\overline{F}_{0}D_{0}\right]$$
(35a)

$$\leq \mathbb{E}\left[R\left(\tau_{1},T\right)\left|\overline{F}_{1}D_{1}\right]+\tilde{C}_{1}+2\alpha\sqrt{s_{1}-1}+\mathbb{E}\left[\tau_{1}-\nu_{1}\left|\overline{F}_{1}D_{1}\right]\right]$$
(35b)

$$+ T \cdot \mathbb{P}\left(F_1 | \overline{F}_0 D_0\right) + T \cdot \mathbb{P}\left(\overline{D}_1 | \overline{F}_1 \overline{F}_0 D_0\right)$$
(35c)

$$\leq \mathbb{E}\left[R\left(\tau_{2},T\right)\left|\overline{F}_{2}D_{2}\right]+\sum_{i=1}^{2}\tilde{C}_{i}+2\alpha\sum_{i=1}^{2}\sqrt{s_{i}-1}$$
(35d)

$$+\sum_{i=1}^{2} \mathbb{E}\left[\tau_{i}-\nu_{i}\big|\overline{F}_{i-1}D_{i-1}\right]+T\sum_{i=1}^{2} \mathbb{P}\left(F_{i}\big|\overline{F}_{i-1}D_{i-1}\right)+T\sum_{i=1}^{2} \mathbb{P}\left(\overline{D}_{i}\big|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right) \quad (35e)$$

$$\vdots$$

$$\leq \sum_{i=1}^{M} \tilde{C}_i + 2\alpha \sum_{i=1}^{M} \sqrt{s_i} + \sum_{i=1}^{M-1} \mathbb{E}\left[\tau_i - \nu_i \middle| \overline{F}_{i-1} D_{i-1}\right]$$
(35f)

$$+T\sum_{i=1}^{M} \mathbb{P}\left(F_{i} | \overline{F}_{i-1} D_{i-1}\right) + T\sum_{i=1}^{M-1} \mathbb{P}\left(\overline{D}_{i} | \overline{F}_{i} \overline{F}_{i-1} D_{i-1}\right)$$
(35g)

$$\leq \sum_{i=1}^{M} \tilde{C}_i + 2\alpha\sqrt{MT} + \sum_{i=1}^{M-1} \mathbb{E}\left[\tau_i - \nu_i \big| \overline{F}_{i-1} D_{i-1}\right]$$
(35h)

$$+T\sum_{i=1}^{M} \mathbb{P}\left(F_{i} | \overline{F}_{i-1} D_{i-1}\right) + T\sum_{i=1}^{M-1} \mathbb{P}\left(\overline{D}_{i} | \overline{F}_{i} \overline{F}_{i-1} D_{i-1}\right)$$
(35i)

where equation 35h follows from the Cauchy–Schwarz inequality

$$\left(\sum_{i=1}^{M} \sqrt{s_i}\right)^2 \le \left(\sum_{i=1}^{M} s_i\right) \left(\sum_{i=1}^{M} 1\right) = M \sum_{i=1}^{M} s_i = MT,$$
(36a)

1070 C.1.1 PROOF OF INTEGRATION WITH CHANGE DETECTORS OF M-UCB

1071 With the general regret bound in Theorem 4.1, a regret bound of the M-UCB with the proposed 1072 diminishing exploration can be obtained by bounding $\mathbb{P}(F_i|\overline{F}_{i-1}D_{i-1})$, $\mathbb{P}(D_i|\overline{F}_i\overline{F}_{i-1}D_{i-1})$, and 1073 $\mathbb{E}[\tau_i - \nu_i|\overline{F}_iD_i\overline{F}_{i-1}D_{i-1}]$.

First, in Lemma C.4, we show that the probability of false alarm is very small; thereby, its contribution to the regret is negligible.

Lemma C.4 (Probability of false alarm). Under Algorithm 1 with parameter in equation 3, and equation 4, we have

$$\mathbb{P}\left(F_i | \overline{F}_{i-1} D_{i-1}\right) \le w K \left(1 - \left(1 - \exp\left(-2b^2/w\right)\right)^{\lfloor T/w \rfloor}\right) \le \frac{1}{T}.$$
(37)

Proof. Suppose that at time t, we have gathered w samples of arm $k \in \mathcal{K}$, namely $Y_{k,1}, Y_{k,2}, \ldots, Y_{k,w}$, for change detection in line 17 of Algorithm 1, and we define

$$S_{k,t} = \sum_{\ell=w/2+1}^{w} Y_{k,\ell} - \sum_{\ell=1}^{w/2} Y_{k,\ell}.$$
(38)

Note that $S_{k,t} = 0$ when there is insufficient (less than w) samples to trigger the change detection algorithm. By definition, we have

$$\tau_{k,i} = \inf\{t \ge \tau_{i-1} + w : |S_{k,t}| > b\}.$$
(39)

Given that the events D_{i-1} and F_{i-1} hold, we define $\tau_{k,i}$ as the first detection time of the k-th arm after ν_i . Clearly, $\tau_i = \min_{k \in \mathcal{K}} {\{\tau_{k,i}\}}$ as Algorithm 1 would reset every time a change is detected. Using the union bound, we have

$$\mathbb{P}\left(F_{i} \middle| \overline{F}_{i-1} D_{i-1}\right) = \mathbb{P}\left(\max_{k \in \mathcal{K}} \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k\}} \ge w, F_{i} \middle| \overline{F}_{i-1} D_{i-1}\right)$$
(40a)

$$+ \mathbb{P}\left(\max_{k\in\mathcal{K}}\sum_{t=\tau_{i-1}+1}^{\nu_{i}}\mathbf{1}_{\{A_{t}=k\}} < w, F_{i} \middle| \overline{F}_{i-1}D_{i-1}\right)$$
(40b)

1101
1102
1103
$$= \mathbb{P}\left(F_{i} \middle| \overline{F}_{i-1} D_{i-1}, \max_{k \in \mathcal{K}} \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k\}} \ge w\right)$$
(40c)

1104
1105
1106
1107

$$\cdot \mathbb{P}\left(\max_{k \in \mathcal{K}} \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k\}} \ge w \middle| \overline{F}_{i-1} D_{i-1}\right)$$
(40d)

$$\leq \mathbb{P}\left(F_i \middle| \overline{F}_{i-1} D_{i-1}, \max_{k \in \mathcal{K}} \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k\}} \geq w\right)$$
(40e)

1111
1112
1113
$$\leq \sum_{k=1}^{K} \mathbb{P}\left(\tau_{k,i} \leq \nu_{i} \middle| \overline{F}_{i-1}D_{i-1}, \max_{k' \in \mathcal{K}} \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k'\}} \geq w\right)$$
(40f)

1114
1115
1116
1117
$$\leq \sum_{k=1}^{K} \mathbb{P}\left(\tau_{k,i} \leq \nu_i \left| \overline{F}_{i-1} D_{i-1}, \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k\}} \geq w\right), \quad (40g)$$

where the term in equation 40b is clearly equal to 0 as there will be no false alarm if we do not even have sufficiently many observations to trigger the alarm as suggested by Algorithm 3. Equation 40c and equation 40d hold by the definition of conditional probability, equation 40e is due to the fact that the term in equation 40d is at most one, and equation 40f follows from the union bound. In equation 40g, if $k \neq k'$, we cannot guarantee that $\sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k'\}} \ge w$. Hence, some k might cause the probability in the equation 40f to be zeros.

1124 For any $0 \le j \le w - 1$, define the stopping time

$$\tau_{k,i}^{(j)} := \inf\{t = \tau_{i-1} + j + nw, n \in \mathbb{Z}^+ : |S_{k,t}| > b\}.$$
(41)

1127 1128 Clearly, $\tau_{k,i} = \min\{\tau_{k,i}^{(0)}, \dots, \tau_{k,i}^{(w-1)}\}$. Let us define, for any $0 \le j \le w - 1$,

1129 1130

1131

1125 1126

1084 1085

1089 1090

1094 1095 1096

1099 1100

$$\xi_{k,i}^{(j)} = \frac{\left(\tau_{k,i}^{(j)} - j - \tau_{i-1}\right)}{w}.$$
(42)

1132 1133 Note that condition on the events D_{i-1} and \bar{F}_{i-1} , $\xi_{k,i}^{(j)}$ is a geometric random variable with parameter $p := \mathbb{P}(|S_{k,t}| > b)$, because when fixing j, there is no overlap between the samples in the current window and the next.

$$\begin{array}{ccc} \mathbf{1136} \\ \mathbf{1137} \\ \mathbf{1138} \\ \mathbf{1139} \end{array} & \mathbb{P}\left(\tau_{k,i}^{(j)} = \tau_{i-1} + nw + j \middle| \overline{F}_{i-1} D_{i-1}, \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k\}} \ge w \right)$$

$$= \mathbb{P}\left(\xi_{k,i} = n \middle| \overline{F}_{i-1} D_{i-1}, \sum_{t=\tau_{i-1}+1}^{\nu_i} \mathbf{1}_{\{A_t=k\}} \ge w\right) = p(1-p)^{n-1}.$$
 (43)

Here, the inclusion of subsequent events as conditions should not impact the results, as when entering the change detection algorithm, those events have already occurred. Moreover, by union bound, we have that for any $k \in \mathcal{K}$,

$$\mathbb{P}\left(\tau_{k,i} \leq \nu_{i} \left| \overline{F}_{i-1} D_{i-1}, \sum_{t=\tau_{i-1}+1}^{\nu_{i}} \mathbf{1}_{\{A_{t}=k\}} \geq w \right) \leq w \left(1 - (1-p)^{\lfloor (\nu_{i}-\tau_{i-1})/w \rfloor}\right) \quad (44a)$$

$$\leq w \left(1 - (1-p)^{\lfloor T/w \rfloor}\right). \quad (44b)$$

We further use the McDiarmid's inequality and the union bound to show that

$$p = \mathbb{P}\left(|S_{k,t}| > b\right) = \mathbb{P}\left(S_{k,t} > b\right) + \mathbb{P}\left(S_{k,t} < -b\right)$$
(45a)

$$\leq 2 \cdot \exp\left(-\frac{2b^2}{w}\right). \tag{45b}$$

(44b)

(49)

Using the result in equation 44b and equation 45b into equation 40g,

$$\mathbb{P}\left(F_{i}\big|\overline{F}_{i-1}D_{i-1}\right) \leq \sum_{k=1}^{K} w\left(1 - \left(1 - 2\exp\left(-\frac{2b^{2}}{w}\right)\right)^{\lfloor T/w \rfloor}\right)$$
(46a)

$$= wK\left(1 - \left(1 - 2\exp\left(-\frac{2b^2}{w}\right)\right)^{\lfloor T/w \rfloor}\right).$$
(46b)

Moreover, applying $(1-x)^a > 1 - ax$ for any a > 1 and 0 < x < 1 and plugging the choice of $b = \sqrt{w \log (2KT^2)/2}$ as in equation 4 shows the second inequality.

Lemma C.2 ensures that, with high probability, the detection delay is confined within a tolerable interval. That is, each arm is sampled w/2 times, and using equation 18 from lemma C.2, we select h_i as

$$h_i = \left\lceil w \left(\frac{K}{2\alpha} + 1\right) \sqrt{s_i + 1} + \frac{w^2}{4} \left(\frac{K}{2\alpha} + 1\right)^2 \right\rceil.$$
(47)

Lemma C.5 (Probability of successful detection). Consider a piecewise-stationary bandit envi-ronment. For any $\mu^{(i)}, \mu^{(i+1)} \in [0,1]^K$ with parameters chosen in equation 3 and equation 4 and

$$h_i = \left\lceil w \left(\frac{K}{2\alpha} + 1 \right) \sqrt{s_i + 1} + \frac{w^2}{4} \left(\frac{K}{2\alpha} + 1 \right)^2 \right\rceil,\tag{48}$$

for some $k \in \mathcal{K}, i \geq 1$ and c > 0, under the Algorithm 1, we have

1187
$$\mathbb{P}\left(D_i | \overline{F}_i \overline{F}_{i-1} D_{i-1}\right) \ge 1 - \frac{1}{T}.$$

$$\mathbb{P}\left(D_{i}|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right) = \mathbb{P}\left(\tau_{i} \leq \nu_{i} + h_{i}|\overline{F}_{i}\overline{F}_{i-1}D_{i-1}\right)$$
(50a)

$$\geq \max_{t \in \{\nu_i+1,\dots,\nu_i+h_i\}} \mathbb{P}\left(S_{\tilde{k},t} > b \middle| \overline{F}_i \overline{F}_{i-1} D_{i-1}\right)$$
(50b)

$$\geq \max_{j \in \{0,...,w/2\}} \left(1 - 2 \exp\left(-\frac{(j\left|\delta_{\tilde{k}}^{(i)}\right| - b)^2}{w}\right) \right)$$
(50c)

$$=1 - 2 \exp\left(-\frac{(w|\delta_{\tilde{k}}^{(i)}|/2 - b)^2}{w}\right)$$
(50d)

$$\geq 1 - 2\exp\left(-\frac{wc^2}{4}\right).\tag{50e}$$

where $S_{\tilde{k},t}$ is defined in equation 38, equation 50c follows from the McDiarmid's inequality, and equation 50d is due to the fact that the maximum value is attained when j = w/2. Last, equation 50e is true for any choice of w, b and c such that $\delta_{\tilde{k}}^{(i)} \ge 2b/w + c$ holds. We thus set w and b as in equation 3 and equation 4, respectively, and choose $c = 2\sqrt{\log(2T)/w}$, which leads to $\mathbb{P}(D_i | \overline{F}_i \overline{F}_{i-1} D_{i-1}) \ge 1 - 1/T$.

Lemma C.6 further bounds the expected detection delay in the situation where the change detection algorithm successfully detects the change within the desired interval.

Lemma C.6 (Expected detection delay). Consider a piecewise-stationary bandit environment. For any $\mu^{(i)}, \mu^{(i+1)} \in [0,1]^K$ with parameters chosen in equation 3 and equation 4 and

$$h_i = \left[w \left(\frac{K}{2\alpha} + 1 \right) \sqrt{s_i + 1} + \frac{w^2}{4} \left(\frac{K}{2\alpha} + 1 \right)^2 \right],\tag{51}$$

for some $K \in \mathcal{K}, i \geq 1$ and c > 0, under the Algorithm 1, we have

$$\mathbb{E}\left[\tau_{i}-\nu_{i}\middle|\overline{F}_{i}D_{i}\overline{F}_{i-1}D_{i-1}\right] \leq h_{i}.$$
(52)

Proof. For any $1 \le i \le M$, we have

$$\mathbb{E}\left[\tau_{i}-\nu_{i}\big|\overline{F}_{i}D_{i}G_{i}\overline{F}_{i-1}D_{i-1}\right]=\sum_{j=1}^{h_{i}}\mathbb{P}\left(\tau_{i}\geq\nu_{i}+j\big|\overline{F}_{i}D_{i}G_{i}\overline{F}_{i-1}D_{i-1}\right)\leq h_{i}.$$
(53a)

 Proof.

Plugging the bounds in Lemmas C.4, C.5 and C.6 into Theorem 4.1 shows the following regret bound in Corollary 4.4.

1229 Corollary 4.4 Combining Algorithm 1 and 3 with the parameters in Equation 3, and Equation 41230 achieves the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{i=1}^{M} \tilde{C}_{i}}_{(a)} + \underbrace{2\alpha\sqrt{MT}}_{(b)} + \underbrace{w\left(\frac{K}{2\alpha}+1\right)\sqrt{M\left(T+M\right)}}_{(c)} + \underbrace{\frac{w^{2}M}{4}\left(\frac{K}{2\alpha}+1\right)^{2}}_{(c)} + \underbrace{\frac{2M}{4}\left(\frac{K}{2\alpha}+1\right)^{2}}_{(c)} + \underbrace{\frac{2M}{4}\left(\frac{K}{2\alpha}+1\right)^{2}}_{($$

1240 where $\tilde{C}_i = 8 \sum_{\Delta_k^{(i)} > 0} \frac{\log T}{\Delta_k^{(i)}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \Delta_k^{(i)}$. By setting $\alpha = c\sqrt{K \log (KT)}$ for some constant c, the expected regret is upper-bounded by $\mathcal{O}(\sqrt{KMT \log T})$.

C.1.2 PROOF OF INTEGRATION WITH CHANGE DETECTORS OF GLR-UCB

First, we introduce the function \mathcal{J} , originally introduced by Kaufmann & Koolen (2021),

$$\mathcal{J}(x) := 2\tilde{g}\left(\frac{g^{-1}\left(1+x\right) + \ln\left(\pi^2/3\right)}{2}\right),\tag{55}$$

where $g^{-1}(y)$ is the inverse function of $g(y) := y - \ln(y)$ defined for $y \ge 1$, and for any $x \ge 0$, $\tilde{q}(x) := e^{1/g^{-1}(x)}q^{-1}(x)$ if $x \ge q^{-1}(1/\ln(3/2))$ and $\tilde{q}(x) = (3/2)(x - \ln(\ln(3/2)))$ otherwise. We select the threshold function

$$\beta(n,\epsilon) := 2\mathcal{J}\left(\frac{\log\left(3n\sqrt{n}/\epsilon\right)}{2}\right) + 6\log\left(1+\log n\right),\tag{56}$$

and define h_i in successful detection events D_i to be $h_i := h_i(\alpha, \epsilon)$ with $h_0(\alpha, \epsilon) := 0$ and for i > 0,

$$h_{i}(\alpha,\epsilon) := \left[2\left(\frac{4}{\left(\delta^{(i)}\right)^{2}}\beta\left(T,\epsilon\right) + 2\right)\left(\frac{K}{2\alpha} + 1\right)\sqrt{s_{i} + 1} + \left(\frac{4}{\left(\delta^{(i)}\right)^{2}}\beta\left(T,\epsilon\right) + 2\right)^{2}\left(\frac{K}{2\alpha} + 1\right)^{2} \right],$$
(57)

which guarantees that with the proposed diminishing exploration, we will observe

$$\left\lceil \frac{4}{\left(\delta^{(i)}\right)^2} \beta\left(\frac{3}{2}s_i, \epsilon\right) + 1 \right\rceil,\tag{58}$$

post-change samples for each $k \in \mathcal{K}$ after ν_i . We analyze the GLR-UCB with diminishing exploration under the following assumption:

Assumption C.7. $\nu_i - \nu_{i-1} \ge 2 \max\{h_i, h_{i-1}\}$ for all $i \in \{1, ..., M\}$.

Following the proof of Lemma 8 in Besson et al. (2022), we can show the following lemma:

Lemma C.8. Under assumption C.7 and Equation 57, it holds that

$$\sum_{i=1}^{M} \mathbb{P}\left(F_{i} \middle| \overline{F}_{i-1} D_{i-1}\right) + \sum_{i=1}^{M-1} \mathbb{P}\left(\overline{D}_{i} \middle| \overline{F}_{i} \overline{F}_{i-1} D_{i-1}\right) \le \epsilon \left(K+1\right) M.$$
(59)

Plugging equation 57 and Lemma C.8 into Theorem 4.1 shows the following Corollary 4.6.

Corollary 4.6 Combining Algorithm 1 and 4 with β function in equation 5 achieves the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{i=1}^{M} \tilde{C}_{i}}_{(a)} + \underbrace{2\alpha\sqrt{MT}}_{(b)} + \underbrace{\left(\frac{4}{\left(\delta^{(i)}\right)^{2}}\beta\left(T,\epsilon\right) + 2\right)^{2}\left(\frac{K}{2\alpha} + 1\right)^{2}M}_{(c)}$$

$$+\underbrace{2\left(\frac{4}{\left(\delta^{(i)}\right)^{2}}\beta\left(T,\epsilon\right)+2\right)\left(\frac{K}{2\alpha}+1\right)\sqrt{M\left(T+M\right)}}_{(c)}+\underbrace{\epsilon\left(K+1\right)M}_{(d)},\quad(60)$$

where $\tilde{C}_i = 8 \sum_{\Delta_k^{(i)} > 0} \frac{\log T}{\Delta_k^{(i)}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \Delta_k^{(i)}$. By setting $\alpha = c\sqrt{K\log(KT)}$ for some constant c and $\epsilon = 1/\sqrt{T}$, the expected regret is upper-bounded by $\mathcal{O}(\sqrt{KMT \log T})$.

1296 C.2 PROOF OF SECTION 5

This subsection covers the proofs of the extended results discussed in Section 5. These extensions include advanced integrations and new theoretical insights.

Corollary C.9. We can extend Lemma C.1 to the case where we only care about the optimal arm changing to another one. Then, the number of times that arm k is selected in the exploration phase and the diminishing exploration regret during the time interval $[\tau_{r-1}^*, \nu_r^*)$ would be

$$N_{\text{DE},k}\left(\tau_{r-1}^{*},\nu_{r}^{*}\right) \leq \frac{2\alpha\sqrt{\nu_{r}^{*}-\tau_{r-1}^{*}}}{K} + \frac{3}{2},\tag{61}$$

1309 and

$$\mathbb{E}\left[R_{\rm DE}\left(\tau_{r-1}^{*},\nu_{r}^{*}\right)\right] \le 2\alpha\sqrt{\nu_{r}^{*}-\tau_{r-1}^{*}} + \frac{3}{2}K.$$
(62)

Corollary C.10. We can extend Lemma C.3 to the case where we only care about the optimal arm changing to another one. Then, the number of times that arm k is selected in the exploration phase and the diminishing exploration regret during $[\tau_{r-1}^*, \nu_r^*)$ can be bounding by

$$\mathbb{E}\left[R\left(\tau_{r-1}^{*},\nu_{r}^{*}\right)\Big|\overline{F}_{r-1}^{*}D_{r-1}^{*}\right] \leq \tilde{C} + 2\alpha\sqrt{s_{r}^{*}} + T \cdot \mathbb{P}\left(F_{r}^{*}\Big|\overline{F}_{r-1}^{*}D_{r-1}^{*}\right),\tag{63}$$

Lemma C.11. The false alarm of the super segment can be bounded as follow:

$$\mathbb{P}\left(F_{r}^{*} \middle| \overline{F}_{r-1}^{*} D_{r-1}^{*}\right) \leq \sum_{k=1}^{K} \mathbb{P}\left(\overline{Ignore}, \text{ optimal arm no change } \middle| arm k \text{ alarm, } \overline{F}_{r-1}^{*} D_{r-1}^{*}\right)$$
(64)

Proof. Conditioning on $\overline{F}_{r-1}^* D_{r-1}^*$ holds, we have the event

$$F_r^* = \bigcup_{k=1}^K \left\{ \overline{Ignore}, \text{ optimal arm no change, arm } k \text{ alarm} \right\}.$$
 (65)





Using the union bound,

1352
1353
$$\mathbb{P}\left(F_r^* \middle| \overline{F}_{r-1}^* D_{r-1}^*\right) = \mathbb{P}\left(\bigcup_{k=1}^K \left\{\overline{Ignore}, \text{ optimal arm no change, arm } k \text{ alarm}\right\} \middle| \overline{F}_{r-1}^* D_{r-1}^*\right)$$
1355

$$\leq \sum_{k=1}^{K} \mathbb{P}\left(\overline{Ignore}, \text{ optimal arm no change, arm } k \text{ alarm} \Big| \overline{F}_{r-1}^* D_{r-1}^* \right)$$
(66b)

$$=\sum_{k=1}^{K} \mathbb{P}\left(\overline{Ignore}, \text{ optimal arm no change} \middle| \operatorname{arm} k \text{ alarm, } \overline{F}_{r-1}^{*} D_{r-1}^{*}\right)$$
(66c)

$$\times \mathbb{P}\left(\operatorname{arm} k \operatorname{alarm} \middle| \overline{F}_{r-1}^* D_{r-1}^* \right) \tag{66d}$$

$$\leq \sum_{k=1}^{K} \mathbb{P}\left(\overline{Ignore}, \text{ optimal arm no change} \middle| \text{arm } k \text{ alarm, } \overline{F}_{r-1}^* D_{r-1}^* \right)$$
(66e)

Lemma C.12. The miss detection probability of the super segment can be bounded as follow:

$$\mathbb{P}\left(\overline{D}_{r}^{*} \middle| \overline{F}_{r}^{*} \overline{F}_{r-1}^{*} D_{r-1}^{*} \right) \leq \sum_{k=1}^{K} \mathbb{P}\left(\overline{D}_{r,k} \middle| \overline{F}_{r}^{*} \overline{F}_{r-1}^{*} D_{r-1}^{*} \right)$$
(67)

Proof. We make some modifications to event D_i in Section 4 and extend $D_{r,k}$ as the arm k alarm 1379 in the tolerate delay (regardless of whether to ignore or not) after the r-th optimal arm change. 1380 Conditioning on $\overline{F}_r^* \overline{F}_{r-1}^* D_{r-1}^*$ holds, we have the event

$$\overline{D}_{r}^{*} = \bigcup_{k=1}^{K} \overline{D}_{r,k}.$$
(68)

Using the union bound, we can get the result as Equation 67.

(66a)

Theorem C.13 (General form of regret bound). *Insert Algorithm 2 into Algorithm 1 and combine with a CD algorithm, which achieves the expected regret upper bound as follows:*

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \sum_{\substack{r=1\\(a)}}^{S} \tilde{C}_{r}^{*} + \underbrace{2\alpha\sqrt{ST}}_{(b)} + \underbrace{\sum_{r=1}^{S-1} \left(\mathbb{E}\left[\tau_{r}^{*} - \nu_{r}^{*} \middle| D_{i}^{*}\overline{F}_{i}^{*}D_{i-1}^{*}\overline{F}_{i-1}\right] + d_{I,r}\right)}_{(c)} + \underbrace{T\sum_{r=1}^{S-1} \mathbb{P}\left(F_{r}^{*} \middle| \overline{F}_{r-1}^{*}D_{r-1}^{*}\right) + \mathbb{P}\left(\overline{D}_{r}^{*} \middle| \overline{F}_{r}^{*}\overline{F}_{r-1}^{*}D_{r-1}^{*}\right)}_{(d)} + \underbrace{K\left(M-1\right)N_{I} + T\sum_{i=1}^{M} \mathbb{P}\left(F_{i} \middle| \overline{F}_{i-1}D_{i-1}\right)}_{(e)},$$
(69)

1400
1401 where
$$\tilde{C}_r^* = 8 \sum_{\min_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t} > 0} \frac{\log T}{\min_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \max_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t},$$

and term (e) represents the cost incurred during the period of deciding whether to skip, which
 includes the additional exploration cost for collecting sufficient observations for the skipping
 mechanism. This cost applies to all the changes detected, including actual changes and false alarms.

Moreover, we can transform Equation 69 into another form using lemma C.11 and C.12 as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{r=1}^{S} \tilde{C}_{r}^{*}}_{(a)} + \underbrace{2\alpha\sqrt{ST}}_{(b)} + \underbrace{\sum_{r=1}^{S-1} \left(\mathbb{E}\left[\tau_{r}^{*} - \nu_{r}^{*} \middle| D_{i}^{*}\overline{F}_{i}^{*}D_{i-1}^{*}\overline{F}_{i-1}^{*}\right] + d_{I,r}\right)}_{(c)}$$

$$\begin{array}{c} \begin{array}{c} \begin{array}{c} 1410\\ 1411\\ 1412\\ 1413\\ 1414 \end{array} \\ \begin{array}{c} +T\sum_{r=1}^{S-1}\sum_{k=1}^{K} \mathbb{P}\left(\overline{Ignore}, optimal \ arm \ no \ change \left| arm \ k \ alarm, \ \overline{F}_{r-1}^{*}D_{r-1}^{*}\right) + T\sum_{r=1}^{S-1}\sum_{k=1}^{K} \mathbb{P}\left(\overline{D_{r,k}} \left| \overline{F}_{r}^{*}\overline{F}_{r-1}^{*}D_{r-1}^{*}\right) \right), \end{array} \\ \begin{array}{c} (d) \end{array}$$

1416 Lemma C.14. Suppose an arm $k \in \mathcal{K}$ changes at time ν and raises an alarm at time τ , but the **1417** optimal arm is the same one, and we choose a variable N_I such that $N_k(\nu, \tau) + N_I \geq \frac{4\xi \log T}{\Delta_{\min}^2}$ and **1418** $\xi = 1$, we have

$$\mathbb{P}\left(\overline{Ignore}, optimal \ arm \ no \ change \left| arm \ k \ alarm \ , \overline{F}_{r-1}^* D_{r-1}^* \right) \le \frac{2}{T^2}$$
(71)

1423 *Proof.* Condition on arm k alarms, \overline{F}_{r-1}^* and D_r^* hold, the event

1424
1425
$$\{\overline{Ignore}, \text{ optimal arm no change}\} \subset \left\{\hat{\mu}_k \ge \mu_k + \sqrt{\frac{\xi \log T}{N_k (\nu, \tau) + N_I}}\right\}$$
(72a)
1426 $\left(\sqrt{-\xi \log T}\right)$

$$\cup \left\{ \hat{\mu}_{k^*} \le \mu_{k^*} - \sqrt{\frac{\xi \log T}{N_{k^*}(\nu, \tau) + N_I}} \right\}$$
(72b)

$$\cup \left\{ \mu_{k^*} - \mu_k < 2\sqrt{\frac{\xi \log T}{N_k \left(\nu, \tau\right) + N_I}}, N_k \left(\nu, \tau\right) + N_I \ge \frac{4\xi \log T}{\Delta_{\min}^2} \right\}$$
(72c)

1434 The third event will vanish because

$$2\sqrt{\frac{\xi\log T}{N_k(\nu,\tau)+N_I}} \le 2\sqrt{\frac{\xi\log T\Delta_{\min}^2}{4\xi\log T}} = \Delta_{\min} \le \mu_{k^*} - \mu_k \tag{73}$$

Equation 73 substitutes the latter term of event 72c into the former term. As a result of the substitution, it is determined that this event cannot occur, hence the probability is zero. Therefore, we only need to consider events 72a and 72b. Using the Chernoff-Hoeffding bound, we can obtain

$$\mathbb{P}\left(\hat{\mu}_k \ge \mu_k + \sqrt{\frac{\xi \log T}{N_k \left(\nu, \tau\right) + N_I}}\right) \le T^{-2\xi}.$$
(74)

$$\mathbb{P}\left(\hat{\mu}_{k^*} \le \mu_{k^*} - \sqrt{\frac{\xi \log T}{N_{k^*}(\nu, \tau) + N_I}}\right) \le T^{-2\xi}$$
(75)

If we choose $\xi = 1$, then

$$\mathbb{P}\left(\hat{\mu}_{k} \ge \mu_{k} + \sqrt{\frac{\xi \log T}{N_{k}\left(\nu, \tau\right) + N_{I}}}\right) + \mathbb{P}\left(\hat{\mu}_{k^{*}} \le \mu_{k^{*}} - \sqrt{\frac{\xi \log T}{N_{k^{*}}\left(\nu, \tau\right) + N_{I}}}\right) \le \frac{2}{T^{2}}$$
(76a)

1455 C.2.1 PROOF OF INTEGRATION WITH CHANGE DETECTORS OF M-UCB

1457 Assumption C.15. The algorithm knows a lower bound $\delta > 0$ such that $\delta \le \min_i \max_{k \in \mathcal{K}} \delta_k^{(i)}$. Assumption C.16. $s_r^* = \Omega\left(\left(\log KT + \sqrt{K \log KT}\right) \sqrt{s_{r-1}^*}\right)$.

1458 In particular, if
$$s_r^* = \Theta\left(\left(\log KT + \sqrt{K\log KT}\right)^{2(1+\epsilon)}\right)$$
 for every *i*, Assumption C.16 holds.
1460

Corollary 5.2 Combining Algorithm 1 and 3 with the parameters in Equation 6, and Equation 4 achieves the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{r=1}^{S} \tilde{C}_{r}^{*}}_{(a)} + \underbrace{2\alpha\sqrt{ST}}_{(b)} + \underbrace{w\left(\frac{K}{2\alpha}+1\right)\sqrt{S\left(T+S\right)}}_{(c)} + \underbrace{\frac{w^{2}S}{4}\left(\frac{K}{2\alpha}+1\right)^{2}}_{(c)} + \underbrace{\frac{2S}{(d)}}_{(c)}, \quad (77)$$

where

 $\tilde{C}_r^* = 8 \sum_{\min_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t} > 0} \frac{\log T}{\min_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t}} + \left(\frac{5}{2} + \frac{\pi^2}{3} + K\right) \sum_{k=1}^K \max_{\nu_{r-1}^* \le t \le \nu_r^*} \Delta_{k,t}.$ By setting $\alpha = c_{\sqrt{K}\log{(KT)}}$ for some constant c, the expected regret is upper-bounded by $\mathcal{O}(\sqrt{KST\log T}).$

Proof. We can substitute lemma C.6, lemma C.4, and lemma C.5 into terms (c) and (d) of Equation 70 respectively, and term (e) could be zero because w (window size) samples are sufficient to determine whether to ignore. We don't need to make an additional decision interval to ensure that the number of samples is sufficient for ignoring.

C.2.2 PROOF OF INTEGRATION WITH CHANGE DETECTORS OF GLR-UCB

Corollary 5.3 Combining Algorithm 1 and 4 with β function in Equation 5 achieves the expected regret upper bound as follows:

$$\mathbb{E}\left[R\left(1,T\right)\right] \leq \underbrace{\sum_{\substack{r=1\\(a)}}^{S} \tilde{C}_{r}^{*} + \underbrace{2\alpha\sqrt{ST}}_{(b)} + \underbrace{\frac{2KS}{T} + \epsilon S(K+1)}_{(d)} + \underbrace{K(M-1)\frac{4\log T}{\Delta_{\min}^{2}} + \epsilon MK}_{(e)} + \underbrace{\left(\frac{4}{\delta^{2}}\beta\left(T,\epsilon\right) + 2\right)^{2} \left(\frac{K}{2\alpha} + 1\right)^{2} S + 4\left(\frac{2}{\delta^{2}}\beta\left(T,\epsilon\right) + 1\right) \left(\frac{K}{2\alpha} + 1\right) \sqrt{S\left(T+S\right)}}_{(c)}, \quad (78)$$

where

where $\tilde{C}_{r}^{*} = 8 \sum_{\min_{\nu_{r-1}^{*} \le t \le \nu_{r}^{*}} \Delta_{k,t} > 0} \frac{\log T}{\min_{\nu_{r-1}^{*} \le t \le \nu_{r}^{*}} \Delta_{k,t}} + \left(\frac{5}{2} + \frac{\pi^{2}}{3} + K\right) \sum_{k=1}^{K} \max_{\nu_{r-1}^{*} \le t \le \nu_{r}^{*}} \Delta_{k,t}.$ By setting $\alpha = c_{\sqrt{K \log (KT)}}$ for some constant c and $\epsilon = 1/\sqrt{T}$, the expected regret is upper-bounded by $\mathcal{O}(\sqrt{KST \log T})$.

Proof. We can substitute Equation 57, and lemma C.8 into terms (c) and (d) of Equation 70 respec-tively.

D ALGORITHMS AND PARAMETERS TUNING

In this appendix, we provide an explanation of our parameter selection. For M-UCB, the window size w is set to 200 unless otherwise specified; however, for the last data point (M = 100) in Figure 3b, we chose w = 50 due to the limitations inherent to change detection. We compute the change detection threshold $b_{\text{M-UCB}} = \sqrt{w/2 \log (2KT^2)}$ following the original formulation in Cao et al. (2019). Additionally, the uniform exploration rate $\gamma_{\text{M-UCB}} = \sqrt{MK \log T/T}$ is determined as initially stated in Besson et al. (2022).Concerning CUSUM-UCB, we adhere to Liu et al. (2018) by

1512 fixing $\epsilon = 0.1$, setting the change detection threshold $b_{\text{CUSUM-UCB}} = \log (T/M - 1)$, and establishing 1513 the uniform exploration rate $\gamma_{\text{CUSUM-UCB}} = \sqrt{MK \log T/T}$ as initially stated in Besson et al. (2022). 1514 Additionally, in CUSUM-UCB, the change point detection involves averaging the first H samples, 1515 where H is set to 100. For GLR-UCB (Besson et al., 2022), we set $\gamma_{m,\text{GLR-UCB}} = \sqrt{mK \log T/T}$, 1516 where m is the number of alarms. We utilize the threshold function $\beta(n, \delta) = \log (n^{3/2}/\delta)$ and 1517 set $\delta = 1/\sqrt{T}$. In our setup, for both the diminishing versions of M-UCB and CUSUM-UCB, we 1518 follow the parameter selection approach described earlier, except for the choice of the exploration 1519 rate. In this context, we opt for $\alpha = 1$. For passive methods, including DUCB (Garivier & Moulines, 1520 2011) and DTS (Qi et al., 2023), we use a discounting factor $\gamma = 0.75$. MASTER, on the other hand, 1521 follows the theoretical settings outlined in Wei & Luo (2021) and is categorized as an active method. 1522

To evaluate the scalability of our methods, we conducted three sets of scaling experiments. For scaling in t (Figure 3a), scaling in M (Figure 3b), and scaling in T (Figure 3c), The parameters of DUCB, DTS and MASTER follow the theoretical settings outlined in Garivier & Moulines (2011), Qi et al. (2023) and Wei & Luo (2021), respectively.

Table 2: Parameter	Selection	for Active a	and Passive	Methods
--------------------	-----------	--------------	-------------	---------

Method	Parameters	References	
M-UCB	- Window size $w = 200$ (default), $w = 50$ for $M = 100$ (Figure 3b). - Threshold: $b_{\text{M-UCB}} = \sqrt{w/2 \log (2KT^2)}$.		
	- Exploration rate: $\gamma_{M-UCB} = \sqrt{MK \log T/T}$.	Cao et al. (2019),	
CUSUM-UCB	- Fixed $\epsilon = 0.1$. - Threshold: $b_{\text{CUSUM-UCB}} = \log (T/M - 1)$. - Exploration rate: $\gamma_{\text{CUSUM-UCB}} = \sqrt{MK \log T/T}$.		
	- Change detection based on averaging first $H = 100$ samples.	Liu et al. (2018),	
GLR-UCB	- Exploration rate: $\gamma_{m,\text{GLR-UCB}} = \sqrt{mK \log T/T}$. Threshold function: $\beta(n, \delta) = \log (m^{3/2}/\delta) \cdot \delta = 1/\sqrt{T}$.	Passon at al. (2022)	
	- The shold function: $\beta(n, 0) = \log(n + 70), 0 = 1/\sqrt{1}$.	Besson et al. (2022)	
MASTER	- All parameters follow theoretical settings.	Wei & Luo (2021)	
DUCB (Passive)	- Discounting factor $\gamma = 0.75$.(Figure 3a and 3d) - Discounting factor follows theoretical setting. (Figure 3b and 3c)	Garivier & Moulines (2011)	
DTS (Passive)	- Discounting factor $\gamma = 0.75$. (Figure 3a and 3d) - Discounting factor follows theoretical setting. (Figure 3b and 3c)	Qi et al. (2023)	

1544 1545

1546

1527

1528

E ADDITIONAL RELATED WORK

Structured Non-Stationary Bandits. Another related line of research works is non-stationary bandit 1547 with structured reward changes, which are typically motivated by the dynamic behavior of real-world 1548 applications. For example, Heidari et al. (2016) proposes the *Rising bandit* problem, where the 1549 rewards are assumed to be a non-decreasing and concave function of the current time index and the 1550 number of pulls. Subsequently, this model is extended to the stochastic setting by (Metelli et al., 1551 2022). Another related setting is the *Rotting bandit* (Levine et al., 2017), where the expected reward 1552 is a non-increasing function of the number of pulls. Moreover, Zhou et al. (2021) studies the regime 1553 switching bandit, where the rewards are jointly controlled by an underlying finite-state Markov chain. 1554 However, the algorithms tailored to the above customized formulations are not directly applicable to 1555 the general piecewise-stationary MAB problem.

1556 Bandit Quickest Change Detection. Since the seminal works (Page, 1954; Lorden, 1971), the 1557 quickest change detection (QCD) problem, which involves identifying the change of distribution at an 1558 unknown time with minimal delay, has been a well studied detection problem of stochastic processes 1559 (Veeravalli & Banerjee, 2014). Bandit QCD, a variant of QCD problem recently proposed by (Gopalan 1560 et al., 2021), adds another layer of complexity to the conventional QCD by considering bandit 1561 feedback. A concurrent work (Xu et al., 2021) also studies a similar setting, namely multi-stream 1562 QCD under sampling control, and proposes a myopic sampling policy that achieves a second-order asymptotically optimal detection delay. Despite the above bandit QCD methods focusing mainly on 1563 achieving low detection delay rather than characterizing regret bounds, these recent progress could 1564 nicely complement the studies of piecewise stationary bandits. 1565

1566 F ADDITIONAL SIMULATIONS

1567 1568

Regret Scaling in K. We considered an environment with T = 20000 and M = 5 for various 1569 K, aiming to showcase dynamic regrets versus different K. In this experiment, expected rewards 1570 are generated randomly. Specifically, we randomly generated 5 instances, averaging each instance 1571 over 50 times. Figure 6a demonstrates that our method is not limited to working only in simple 1572 environments with small K values but is adaptable to a broader range of scenarios.

1573 Comepare to AdSwitch, ArmSwitch and Meta Algorithm. These algorithms are indeed com-1574 putationally quite complex, evidenced by the time complexity of $O(KT^4)$ in Auer et al. (2019), 1575 that of $O(K^2T^2)$ in Abbasi-Yadkori et al. (2023), and the fact that recursive calls of the base algo-1576 rithm are needed by the algorithm in Suk & Kpotufe (2023). Besides, while being able to achieve 1577 near-optimal regret bound asymptotically, these elimination-based algorithms generally would not 1578 perform well when the time horizon T is small, shown as Figure 6b and 6c. Figure 6d compares our extension, which incorporates a skipping mechanism, with ArmSwitch, where the latter focuses on 1579 tracking the most significant arm switches. For our setup, we defined $\mu_1^{(i)} = 0.8, 0.2$ for *i* where $(i+1) \mod 4 = \{2,3\}, \{0,1\}$, and $\mu_2^{(i)} = 0.4, 0.6$ for *i* where $(i+2) \mod 2 = 0, 1$, as well as 1580 1581 1582 $\mu_2^{(i)} = 0.4, 0.6$ for i where $(i+3) \mod 2 = 0, 1$. In our parameter settings, we set $N_I = 50$ and 1583 $\alpha = 1$. The results clearly show that our performance significantly exceeds that of the ArmSwitch. The above discussion precisely constitutes the main reason we did not include these algorithms in 1585 our experimental comparison, as such algorithms with $T = 20000 \sim 100000$ would take too long to 1586 finish while showing results for small T may look unfair for those excellent algorithms.

1587 Figure 7. To explore the applicability of our approach in more general nonstationary settings, we 1588 further conducted experiments in a slowly changing environment, where the two arms followed 1589 sinusoidal oscillations as considered in (Qi et al., 2023). This setup represents a gradual, continuous 1590 shift, distinct from the abrupt changes in our primary focus.

1591 Performance of Active and Passive Methods: Surprisingly, some active methods perform reasonably 1592 well in this setting. However, their effectiveness is highly dependent on the choice of change detector. 1593 For example, M-UCB struggles in this environment, as it is not designed for smooth transitions. 1594 Among the passive methods (e.g., Discounted-klUCB and Discounted-klUCB-TS), setting a discount 1595 factor of 0.99 yields strong results, emphasizing their adaptability to slowly changing environments. 1596

Figure 8. Our empirical results further support the importance of the skipping mechanism. As 1597 illustrated in figure 8, the inclusion of the skipping mechanism significantly improves the performance 1598 of the algorithm in regret. By applying the proposed diminishing exploration and the skipping 1599 mechanism, we can improve the empirical regrets of various active methods, including M-UCB, CUSUM-UCB, and GLR-UCB. 1601

Figure 9. We retain the mean reward structure in figure 3d but introduce controlled variations with an even larger horizon T = 54000. These experiments reveal a clear advantage of our DE framework, 1603 particularly in more dynamic scenarios. 1604

- 1608 1609 1610 1611 1612
- 1613
- 1614
- 1615
- 1616
- 1617 1618
- 1619









1782 A ZOOMED-IN VIEW OF FIGURES G 1783

1784

Below are the zoomed-in versions of Figure 3 and Figure 4.



34





