

Hybrid Combinatorial Multi-armed Bandits with Probabilistically Triggered Arms

Kongchang Zhou

Southern University of Science and Technology

12112825@mail.sustech.edu.cn

Tingyu Zhang

Southern University of Science and Technology

j74638239@gmail.com

Wei Chen

Microsoft Research, Beijing, China

weic@microsoft.com

Fang Kong*

Southern University of Science and Technology

kongf@sustech.edu.cn

Reviewed on OpenReview: <https://openreview.net/forum?id=ZRzB7EVLvi>

Abstract

The problem of combinatorial multi-armed bandits with probabilistically triggered arms (CMAB-T) has been extensively studied. Prior work primarily focuses on either the online setting where an agent learns about the unknown environment through iterative interactions, or the offline setting where a policy is learned solely from logged data. However, each of these paradigms has inherent limitations: online algorithms suffer from high interaction costs and slow adaptation, while offline methods are constrained by dataset quality and lack of exploration capabilities. To address these complementary weaknesses, we propose hybrid CMAB-T, a new framework that integrates offline data with online interaction in a principled manner. Our proposed hybrid CUCB algorithm leverages offline data to guide exploration and accelerate convergence, while strategically incorporating online interactions to mitigate the insufficient coverage or distributional bias of the offline dataset. We provide theoretical guarantees on the algorithm’s regret, demonstrating that hybrid CUCB significantly outperforms purely online approaches when high-quality offline data is available, and effectively corrects the bias inherent in offline-only methods when the data is limited or misaligned. Empirical results further demonstrate the consistent advantage of our algorithm.

1 Introduction

Combinatorial multi-armed bandits with probabilistically triggered arms (CMAB-T) provide a powerful framework for modeling a broad class of real-world sequential decision-making problems, including influence maximization, learning to rank, and large language model cache (Chen et al., 2013; 2016; Wang & Chen, 2017; Wen et al., 2017; Kong et al., 2023; Liu et al., 2023; 2025; Pope et al., 2022; Zhu et al., 2023; Gim et al., 2023; Qu et al., 2024). In these settings, a decision-maker repeatedly selects a combinatorial action, typically a subset of base arms, and receives partial feedback governed by a probabilistic triggering process.

Most existing work on CMAB-T has focused on the *online setting*, where an agent learns through trial and error by interacting with the environment over multiple rounds (Chen et al., 2013; 2016; Wang & Chen, 2017; Wen et al., 2017; Kong et al., 2023; Liu et al., 2023; 2024). While this approach enables

*Corresponding author.

adaptive learning and active exploration, it often incurs high feedback collection costs and suffers from slow convergence—particularly in large-scale or high-stakes domains.

A study (Liu et al., 2025) has begun to explore the *offline setting* for CMAB-T, where the goal is to learn decision policies from pre-collected data logs, thereby avoiding the expense of online interaction. However, offline learning is highly sensitive to the quality and coverage of the logged data. For example, rare but important action combinations may be missing, and distributional shifts between the offline data set and online environment can lead to suboptimal performance. Moreover, the lack of active exploration limits the learner’s ability to gather information about underexplored or high-uncertainty actions.

The limitations of purely online or offline learning motivate the study of hybrid learning methods, which use offline data to warm-start online learning (Sentenac et al., 2025; Zheng et al., 2023; Lee et al., 2021; Song et al., 2023; Li et al., 2023; Bu et al., 2022; Shivaswamy & Joachims, 2012; Chen et al., 2022; Oetomo et al., 2023; Agnihotri et al., 2024; Cheung & Lyu, 2024; Qu et al., 2025). These approaches balance the cost-free nature of offline data with the adaptability of online exploration, often leading to improved sample efficiency in practice. While hybrid methods have been studied in the classical MAB problems, their extension to the general CMAB-T setting remains largely unexplored.

The technical challenge arises when incorporating the offline data into the regret analysis of the online CMAB-T. In particular, we must determine *when* to rely on the pure online observation and *when* the offline data (may be biased) is sufficiently reliable to be used. In the MAB setting, the regret admits a clean decomposition: it can be expressed as the sum over arms of the number of times each suboptimal arm is pulled, multiplied by its corresponding sub-optimality gap. This makes it straightforward to quantify how offline data reduces regret by decreasing the selection count of suboptimal arms (Cheung & Lyu, 2024). But in our considered CMAB-T setting, such gap-based reasoning is no longer directly applicable, where the per-round regret cannot be attributed to individual arms through simple suboptimality gaps. The regret depends on the triggered arms and the combinatorial reward structure, making it much more difficult to define a universal threshold for determining when to use offline data.

To overcome these challenges, our work focuses on the following fundamental questions:

- (1) *How to derive an algorithm that effectively leverages offline data in the online CMAB-T setting?*
- (2) *Can we provide the corresponding theoretical guarantees that offline data leads to measurable improvement compared with purely online algorithms?*

We answer these questions through the following contributions:

Problem Formulation. We formally define the *hybrid CMAB-T* (H-CMAB-T) setting by extending the classical CMAB-T framework to incorporate offline data. In particular, we define the offline dataset as a collection of observations over base arms, and introduce a notion of *bias* based on the discrepancy between the offline and online mean rewards of each arm. This formulation provides a principled basis for assessing when offline data can be beneficial to online learning.

Algorithm Design. We propose a new algorithm *hybrid CUCB* leveraging the biased offline data to improve the classic CUCB algorithm. This algorithm balances offline and online feedback through a dual-UCB mechanism. Specifically, we construct two confidence bounds for each base arm: one purely based on the feedback collected online, and another that hybridizes observations from both the offline data set and online interactions with an explicit bias correction. By selecting the minimum of the two UCB estimates, the algorithm adaptively leverages the offline data based on its quality.

Theoretical Analysis. To overcome challenge from the core difference between MAB and CMAB-T, we draw on the intuition that while the bias may appear at the level of individual arms, the regret in CMAB-T arises from actions that involve multiple arms and triggering mechanisms. Motivated by this, we explore a connection between per-arm bias and action-level regret by considering a hypothetical allocation of the regret to the arms that could be triggered in each round. This perspective allows us to bridge the arm-level discrepancy introduced by offline data and the combinatorial nature of regret in CMAB-T. Leveraging this connection, we construct a threshold condition that determines whether the offline estimates are reliable enough to be used. Finally, We provide both gap-dependent and gap-independent regret bounds. To complement the upper

bound analysis, we further establish a regret lower bound for hybrid CMAB-T, which matches the structure of the data-dependent saving term up to constant factors. This lower bound demonstrates that the regret improvements enabled by offline data are near-optimal, and reveals a fundamental information-theoretic limitation on the extent to which offline data can reduce regret in CMAB-T. Furthermore, our results show that the algorithm achieves improved regret over standard online methods (Wang & Chen, 2017), with a provable *saving term* that depends on the informativeness and reliability of the offline data. Our result recovers the standard online regret when offline data is absent or adversarial, and it matches or improves upon the results of Cheung & Lyu (2024) when the problem reduces to classical MAB.

Empirical Evaluation We complement our theoretical analysis with empirical evaluations. The results consistently demonstrate that hybrid CUCB outperforms both purely online and purely offline baselines, highlighting its adaptability and robustness across varying data conditions.

2 Related Work

Online Bandits. MAB problems have been extensively studied as a foundational model for sequential decision-making under uncertainty (Auer et al., 2002; Bubeck & Cesa-Bianchi, 2012; Lattimore & Szepesvári, 2020). The combinatorial multi-armed bandit (CMAB) framework (Chen et al., 2013) generalizes classical MAB by allowing the learner to select subsets of arms (super arms) in each round, leading to richer modeling power and broader applicability. In particular, the CMAB with probabilistically triggered arms (CMAB-T) framework introduced by Chen et al. (2016); Wang & Chen (2017) captures the settings such as influence maximization, online learning to rank where the reward depends not only on the chosen super arm but also on a random triggering process. This framework has also been extended to incorporate contextual information (Liu et al., 2023). A line of work has established algorithms with theoretical regret guarantees under structural assumptions such as monotonicity and bounded smoothness (Chen et al., 2016; Wang & Chen, 2017; Wen et al., 2017; Liu et al., 2022; 2023; 2024). All these approaches operate entirely in the online setting.

Offline Bandits. Offline learning in bandit and reinforcement learning has gained increasing attention due to the high cost of online exploration and the availability of logged historical data. It has been explored in many bandits settings like the classical MAB (Rashidinejad et al., 2021), contextual MAB (Rashidinejad et al., 2021; Jin et al., 2021; Li et al., 2022) and neural contextual bandits (Nguyen-Tang et al., 2021; 2022). For combinatorial bandits, Liu et al. (2025) recently propose CLCB, the first general framework for offline learning in CMAB problems, which characterizes dataset quality through coverage conditions, and provide near-optimal theoretical guarantees.

Hybrid Bandits. To mitigate the limitations of purely online or offline learning, hybrid methods aim to combine their respective advantages by using offline data to initialize or guide online exploration. Hybrid learning has been studied in various domains, including bandit problems (Shivaswamy & Joachims, 2012; Oetomo et al., 2023; Agnihotri et al., 2024) and reinforcement learning (Song et al., 2023; Qu et al., 2025). Most of these hybrid methods assume that offline data is unbiased and directly compatible with the online environment (Shivaswamy & Joachims, 2012; Song et al., 2023; Oetomo et al., 2023; Agnihotri et al., 2024). In particular, Sentenac et al. (2025) study offline-to-online learning in stochastic MABs and propose an algorithm that adaptively balances pessimistic and optimistic strategies across different online horizons. Qu et al. (2025) assume a strongly biased offline dataset with a lower bound on the discrepancy between offline and online means. Cheung & Lyu (2024) do not require such assumptions and propose an algorithm that adaptively incorporates offline data based on its reliability. To the best of our knowledge, the hybrid learning problem in CMAB-T remains open.

3 Problem Setup

We first introduce the *hybrid combinatorial multi-armed bandits with probabilistically triggered arms* (H-CMAB-T) problem. The H-CMAB-T problem explored in this paper is built upon the standard CMAB-T framework (Wang & Chen, 2017). We begin by reviewing the classical CMAB-T setting, and then introduce how offline data is incorporated in our extension.

The online environment consists of m base arms, represented as random variables X_1, X_2, \dots, X_m , jointly distributed according to an unknown distribution $D^{\text{on}} \in \mathcal{D}$, where D^{on} is supported on $[0, 1]^m$ and \mathcal{D} is the distribution family. For each base arm $i \in [m]$, let $\mu_i^{\text{on}} = \mathbb{E}_{X \sim D^{\text{on}}}[X_i]$ denote its expected value, and define the vector $\mu^{\text{on}} = (\mu_1^{\text{on}}, \dots, \mu_m^{\text{on}}) \in [0, 1]^m$ as the mean vector of all arms. Note that μ^{on} is determined by the underlying distribution D^{on} . The learning process unfolds over discrete rounds $t = 1, 2, \dots, T$. In each round:

1. The learner selects a combinatorial action $S_t \in \mathcal{S}$ based on the previous rounds observation and feedback, where \mathcal{S} is a predefined action space, possibly subject to structural constraints. The combinatorial action S_t is also called “super arm” and in many cases it is a subset of base arms.

2. The environment draws an independent sample $X^{(t)} = (X_1^{(t)}, \dots, X_m^{(t)}) \sim D^{\text{on}}$.

3. Playing action S_t in the environment induces a random subset $\tau_t \subseteq [m]$ of arms to be triggered. The triggering process is stochastic: even given the environment outcome $X^{(t)}$ and the chosen action S_t , the triggered set $\tau_t \subseteq [m]$ may still exhibit randomness. We model this using a *probability triggering function* $D^{\text{trig}}(S, X)$, which defines a distribution over subsets of $[m]$ conditioned on action S and environment realization X . Formally, we assume that for each round t , the triggered set τ_t is independently drawn from $D^{\text{trig}}(S_t, X^{(t)})$, i.e., $\tau_t \sim D^{\text{trig}}(S_t, X^{(t)})$. Moreover, to enable algorithms to estimate μ_i^{on} from observed samples during online learning, we make the following identifiability assumption: the outcome of each arm i does not depend on whether it is triggered. That is, $\mathbb{E}_{X \sim D^{\text{on}}, \tau \sim D^{\text{trig}}(S, X)}[X_i \mid i \in \tau] = \mathbb{E}_{X \sim D^{\text{on}}}[X_i] = \mu_i^{\text{on}}, \forall i \in [m]$.

4. A non-negative reward $R(S_t, X^{(t)}, \tau_t) \in \mathbb{R}_{\geq 0}$ is revealed to the learner, which is a deterministic function of the chosen action S_t , the sampled instance $X^{(t)}$, and the triggered set τ_t . The expected reward of an action $S \in \mathcal{S}$ is given by $r_S(\mu) := \mathbb{E}[R(S, X, \tau)]$, where the expectation is taken over $X \sim D$ and $\tau \sim D^{\text{trig}}(S, X)$. We emphasize that $r_S(\mu)$ is a function of S and the mean vector μ .

The goal of the learner is to maximize the total expected reward over T rounds, i.e., to design a learning algorithm that selects S_1, \dots, S_T to maximize $\sum_{t=1}^T \mathbb{E}[R(S_t, X^{(t)}, \tau_t)]$.

While the classical CMAB-T framework captures the core structure of combinatorial bandit problems with triggering, it assumes that all learning happens online from scratch. In many practical scenarios, however, a significant amount of data is already available prior to online interaction—collected from historical logs or prior deployments. For example, in *online influence maximization* problem, the organizations often have access to past propagation traces—records of how information spread—which can serve as valuable offline data to accelerate online learning in new deployment scenarios.

Motivated by this, we consider an extension of CMAB-T that incorporates such *offline data*, and investigate how it can be used to improve learning performance. More specifically, the key difference between H-CMAB-T and CMAB-T problem is that before online learning, the player is given an offline data collection \mathcal{B} . It is worth noting that there may be discrepancies between offline data and the online environment. For example, in the OIM problem, due to the characteristics of the product or shifts in user preferences, the diffusion dynamics within social networks can differ. To characterize such phenomenon and avoid misleading of offline data, we consider that the arms in the offline data set and the online setting may have different means. Specifically, the outcomes of m base arms in the offline data set can be represented as random variables Y_1, Y_2, \dots, Y_m , jointly distributed according to an unknown distribution D^{off} and the mean vector of the offline data is $\mu^{\text{off}} = (\mu_1^{\text{off}}, \dots, \mu_m^{\text{off}})$. It is natural that $|\mu_i^{\text{on}} - \mu_i^{\text{off}}| \geq 0$, and equality holds if and only if the offline data is unbiased. Without loss of generality, we denote N_i as the number of the independent observations of arm i . Then the offline data set can be represented as $\mathcal{B} := \{N_i, \{Y_{i,s}\}_{s=1}^{N_i}\}_{i=1}^m$.

Bias control. Besides, to quantify this discrepancy, we adopt the bias control vector $V = (V_1, \dots, V_m)$ as a hyper-parameter which upper bounds the difference between the offline and online means for each arm:

$$|\mu_i^{\text{off}} - \mu_i^{\text{on}}| \leq V_i, \quad \forall i \in [m].$$

Since both means lie in $[0, 1]$, we assume $V_i \in [0, 1]$ for all i . Smaller values of V_i indicate higher alignment between offline and online environments. In settings with prior knowledge—e.g., similar user populations or stable network dynamics—we may set V_i to be small. In fully agnostic cases where no such knowledge is available, we conservatively set $V_i = 1$.

Remark 1. As rigorously shown in Section 3 of [Cheung & Lyu \(2024\)](#), in the presence of biased offline data, no hybrid algorithm in MAB can be guaranteed to outperform a purely online baseline unless some prior knowledge about the bias is available. This theorem highlights that incorporating some form of prior understanding of the bias is not just helpful but fundamentally necessary. To understand this challenge, one can consider the unknown V setting and try to design a hybrid algorithm that learns V during the online interaction. This raises a challenging trade-off: if V is small, estimating it accurately may require excessive online samples, outweighing the benefit of offline data; if V is large, offline estimates are often too biased to be useful, making a pure online strategy preferable. In practice, prior knowledge about the reward-generating process, the data collection mechanism, or the source of distribution shift may provide useful guidance for choosing a tolerable bias level. Thus, selecting V_i can be viewed as a conservative engineering choice that encodes such prior knowledge. Developing methods that adaptively detect bias or learn V_i during online interaction is an important but technically demanding direction, which we leave for future work.

Consequently, based on the above problem formulation, we define an H-CMAB-T instance as a tuple $([m], \mathcal{S}, \mathcal{D}, D^{\text{trig}}, R, \mathcal{B})$. To make the learning problem well-defined and practically solvable, it remains to specify how actions are selected given current estimates of the arm statistics. In many CMAB-T instances, the action space is exponentially large and the underlying optimization problem of selecting the optimal super arm is NP-hard ([Chen et al., 2013; 2016](#)). To decouple the statistical estimation from the combinatorial optimization, prior works commonly assume the access to an *offline oracle* that returns an approximate solution. This allows the learning algorithm to focus on estimating arm statistics while relying on the oracle to select actions.

Offline (α, β) -approximation oracle \mathcal{O} . We assume access to an offline (α, β) -approximation oracle, denoted by \mathcal{O} . This oracle takes as input the mean vector $\mu = (\mu_1, \dots, \mu_m)$ and returns an action $S^\mathcal{O} \in \mathcal{S}$ such that $\mathbb{P}[r_{S^\mathcal{O}}(\mu) \geq \alpha \cdot \text{opt}_\mu] \geq \beta$, where $\alpha \in (0, 1]$ is the approximation ratio, and $\beta \in (0, 1]$ is the success probability. Here, opt_μ denotes the optimal expected reward under mean vector μ , defined as $\text{opt}_\mu := \sup_{S \in \mathcal{S}} r_S(\mu)$. And we assume that opt_μ is bounded for all μ .

Further, the objective of the learner is to minimize the (α, β) -approximation regret defined as below ([Chen et al., 2013; 2016; Wang & Chen, 2017; Wen et al., 2017](#)).

Definition 1 ((α, β) -approximation regret.). *The (α, β) -approximation regret of a learning algorithm \mathcal{A} over T rounds under an H-CMAB-T instance $([m], \mathcal{S}, \mathcal{D}, D^{\text{trig}}, R, \mathcal{B})$ is*

$$\text{Reg}_{\mu^{\text{opt}}, \alpha, \beta}^{\mathcal{A}}(T) := \alpha \cdot \beta \cdot T \cdot \text{opt}_{\mu^{\text{opt}}} - \mathbb{E} \left[\sum_{t=1}^T R(S_t^{\mathcal{A}}, X^{(t)}, \tau_t) \right] = \alpha \cdot \beta \cdot T \cdot \text{opt}_{\mu^{\text{opt}}} - \mathbb{E} \left[\sum_{t=1}^T r_{S_t^{\mathcal{A}}}(\mu^{\text{opt}}) \right],$$

where $S_t^{\mathcal{A}}$ is the action selected by algorithm \mathcal{A} at round t , and the expectation is taken over the randomness of the environment outcomes $\{X^{(t)}\}_{t=1}^T$, the triggered sets $\{\tau_t\}_{t=1}^T$, and the internal randomness of the algorithm.

This notion of regret captures how far the cumulative reward falls short of what could be obtained by always playing a near-optimal action provided by the oracle.

We now introduce several conditions that are used to establish regret guarantees. These conditions are widely adopted in the CMAB literature ([Chen et al., 2016; Wang & Chen, 2017; Wen et al., 2017; Liu et al., 2023; 2025](#)). To facilitate the presentation, we denote $p_i^{D, S}$ as the probability that arm i is triggered when action S is selected in environment D .

Condition 1 (Monotonicity). *We say that a CMAB-T problem instance satisfies monotonicity, if for any action $S \in \mathcal{S}$, for any two distributions $D, D' \in \mathcal{D}$ with expectation vectors $\mu = (\mu_1, \dots, \mu_m)$ and $\mu' = (\mu'_1, \dots, \mu'_m)$, we have $r_S(\mu) \leq r_S(\mu')$ if $\mu_i \leq \mu'_i$ for all $i \in [m]$.*

Condition 2 (1-Norm TPM Bounded Smoothness). *We say that a CMAB-T problem instance satisfies 1-norm TPM bounded smoothness, if there exists $B \in \mathbb{R}^+$ (referred as the bounded smoothness constant) such that, for any two distributions $D, D' \in \mathcal{D}$ with expectation vectors μ and μ' , and any action S , we have $|r_S(\mu) - r_S(\mu')| \leq B \sum_{i \in [m]} p_i^{D, S} |\mu_i - \mu'_i|$.*

The two reward function conditions encode natural intuitions in the CMAB-T setting: Condition 1 reflects monotonicity—if all arm means are higher in one set than another, any action should yield a higher expected

reward; Condition 2 captures the role of triggering probabilities—arms that are triggered more often contribute more to the reward and thus require more accurate mean estimates, while less frequently triggered arms can tolerate greater uncertainty.

4 The Hybrid CUCB Algorithm

Algorithm 1 Hybrid CUCB with Computation Oracle

Require: Valid bias bound V , number of arms m , offline data $\mathcal{B} := \{N_i, \{Y_{i,s}\}_{s=1}^{N_i}\}_{i=1}^m$, horizon T , Oracle

- 1: **for** each arm $i \in [m]$ **do**
- 2: $\hat{\mu}_i^{\text{off}} \leftarrow \frac{1}{N_i} \sum_{s=1}^{N_i} Y_{i,s}$, $T_i \leftarrow 0$, $\hat{\mu}_i^{\text{on}} \leftarrow 0$
- 3: **end for**
- 4: **for** $t = 1, 2, \dots, T$ **do**
- 5: **for** each arm $i \in [m]$ **do**
- 6: $\text{rad}_t(i) \leftarrow \sqrt{\frac{2 \log(4mt^3)}{T_i}}$ $\triangleright = \infty$ if $T_i = 0$
- 7: $\text{rad}_t^{\text{S}}(i) \leftarrow \sqrt{\frac{2 \log(4mt^3)}{N_i + T_i}} + \frac{N_i}{N_i + T_i} V_i$ $\triangleright = \infty$ if $N_i + T_i = 0$
- 8: $\text{UCB}_t(i) \leftarrow \hat{\mu}_i^{\text{on}} + \text{rad}_t(i)$
- 9: $\text{UCB}_t^{\text{S}}(i) \leftarrow \frac{N_i \hat{\mu}_i^{\text{off}} + T_i \hat{\mu}_i^{\text{on}}}{N_i + T_i} + \text{rad}_t^{\text{S}}(i)$
- 10: $\bar{\mu}_i \leftarrow \min \left\{ \text{UCB}_t(i), \text{UCB}_t^{\text{S}}(i), 1 \right\}$
- 11: **end for**
- 12: $S \leftarrow \text{Oracle}(\bar{\mu}_1, \dots, \bar{\mu}_m)$
- 13: Play action S , triggering a set $\tau \subseteq [m]$ of base arms
- 14: **for** each $i \in \tau$ with feedback $X_i^{(t)}$ **do**
- 15: $T_i \leftarrow T_i + 1$ $\hat{\mu}_i^{\text{on}} \leftarrow \hat{\mu}_i^{\text{on}} + (X_i^{(t)} - \hat{\mu}_i^{\text{on}})/T_i$
- 16: **end for**
- 17: **end for**

In this section, we provide an algorithm, hybrid CUCB (Algorithm 1), aiming to leverage *useful* offline data to accelerate the online learning efficiency. The hybrid CUCB algorithm runs as follows. In each round, the algorithm computes two UCB vectors: $\text{UCB}_t = (\text{UCB}_t(1), \dots, \text{UCB}_t(m))$, $\text{UCB}_t^{\text{S}} = (\text{UCB}_t^{\text{S}}(1), \dots, \text{UCB}_t^{\text{S}}(m))$, and then feeds the coordinate-wise minimum two of them into the (α, β) -approximation oracle to select an action.

The vector UCB_t follows the standard CUCB construction (Wang & Chen, 2017) (Line 6 and 8), representing the conventional UCB established with the pure online feedback, where T_i denotes the number of times that arm i has been triggered.

As to H-CMAB-T problem, to effectively leverage offline data while remaining robust to distributional mismatch, we design a hybrid confidence bound UCB_t^{S} that adaptively incorporates offline observations. Intuitively, when the offline mean of an arm is close to its online counterpart, the offline data should be more trusted. Conversely, if the discrepancy between the two is large, the algorithm should rely primarily on online feedback.

Based on this intuition, we construct $\text{UCB}_t^{\text{S}}(i)$ using a weighted empirical mean and a bias-adjusted confidence radius (Line 7 and 9). The empirical mean aggregates offline and online samples proportionally to their counts, while the confidence radius consists of two components: a standard deviation term based on the total offline and online sample size $N_i + T_i$, and a bias penalty scaled by the discrepancy bound V_i . The weight $N_i/(N_i + T_i)$ ensures that the penalty becomes more prominent as more offline data is used.

Finally, by taking the minimum between the two UCB estimates, the algorithm can exploit useful offline data. Intuitively, if N_i is large and V_i is small such that $\text{UCB}_t^{\text{S}}(i) < \text{UCB}_t(i)$, then the offline data is useful for online exploration and the algorithm utilizes the hybrid $\text{UCB}_t^{\text{S}}(i)$. Otherwise, if μ_i^{off} and μ_i^{on} are far apart, then $\text{UCB}_t^{\text{S}}(i)$ becomes large. The algorithm would default to $\text{UCB}_t(i)$, effectively ignoring offline data.

In both cases, the selection rule ensures that the decision is made conservatively, based on the estimated trustworthiness of the offline data. We next provide the regret upper bound for Algorithm 1 in Section 5.

5 Theoretical Analysis

In this section, we provide the theoretical results for hybrid CUCB. We first provide the gap-dependent regret upper bound and the corresponding discussions. The gap-independent regret analysis comes later and the lower bound is discussed in the end. The complete proof is provided in the appendix.

5.1 Gap-Dependent Bound

We first define the reward gaps used in the regret analysis.

Definition 2 (Gap (Wang & Chen, 2017)). *Fix a distribution D and its expectation vector $\boldsymbol{\mu}$. For each action S , we define the gap $\Delta_S = \max(0, \alpha \cdot \text{opt}_{\boldsymbol{\mu}} - r_S(\boldsymbol{\mu}))$. For each arm i , we define*

$$\Delta_{\min}^i = \inf_{S \in \mathcal{S}: p_i^{D,S} > 0, \Delta_S > 0} \Delta_S, \quad \Delta_{\max}^i = \sup_{S \in \mathcal{S}: p_i^{D,S} > 0, \Delta_S > 0} \Delta_S.$$

As a convention, if there is no action S such that $p_i^{D,S} > 0$ and $\Delta_S > 0$, we define $\Delta_{\min}^i = +\infty$, $\Delta_{\max}^i = 0$. Further define $\Delta_{\min} = \min_{i \in [m]} \Delta_{\min}^i$, $\Delta_{\max} = \max_{i \in [m]} \Delta_{\max}^i$.

Let $\tilde{\mathcal{S}} = \{i \in [m] \mid p_i^{\mu,S} > 0\}$ be the set of arms that could be triggered by S . Let $K = \max_{S \in \mathcal{S}} |\tilde{\mathcal{S}}|$. To formally capture the influence of the discrepancy between offline and online environment, we introduce a measure $\omega_i := V_i + \mu_i^{\text{off}} - \mu_i^{\text{on}}$, $i \in [m]$. By the definition of V , we have that $\omega_i \in [0, 2V_i]$. Intuitively, the quantity ω_i allows us to express how much the offline data for arm i deviates from the true online behavior, and plays a key role in determining the extent to which the offline data influences the online learning.

Theorem 1 (Gap-Dependent Regret Bound). *For an H-CMAB-T problem $([m], \mathcal{S}, \mathcal{D}, D^{\text{trig}}, R, \mathcal{B})$ that satisfies monotonicity (Condition 1) and TPM bounded smoothness (Condition 2), the hybrid CUCB algorithm with an input bias control vector V and an (α, β) -approximation oracle achieves an (α, β) -approximate gap-dependent regret bounded by:*

$$\text{Reg}_{\mu^{\text{on}}, \alpha, \beta}(T) \leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{\Delta_{\min}^i} - 8B\sqrt{2N'_i \log(4mT^3)}, 0 \right\} + 4Bm + \frac{\pi^2}{6} \Delta_{\max}, \quad (1)$$

where

$$N'_i = N_i \cdot \max \left\{ 1 - \frac{2BK\omega_i}{\Delta_{\min}^i}, 0 \right\}^2.$$

Following Theorem 1, we now provide a detailed interpretation of the regret bound and its implications for how offline data is used by our algorithm.

A key quantity in the bound is N'_i , which represents the amount of *effectively utilized* offline data for arm i . The multiplicative factor can be interpreted as the *utilization rate* of the offline data. For a fixed online learning setting, the term $2BK/\Delta_{\min}^i$ is constant, so the utilization rate increases as the discrepancy ω_i decreases. When the offline data is unbiased (i.e., $V_i = \omega_i = 0$), we have full utilization: $N'_i = N_i$. In contrast, when $\omega_i \geq \Delta_{\min}^i/(2BK)$, the utilization rate drops to zero, and the offline data is effectively ignored. This reflects our design intuition: offline data that closely matches the online environment should be trusted more and used more aggressively. The result of Theorem 1 recovers the result of CMAB-T (Wang & Chen, 2017) as a special case when $N'_i = 0$ for all i . The setting may correspond to the case where the offline data do not exist (i.e. $N_i = 0$ for all $i \in [m]$) or the case that the offline data is fully misaligned with the online environment.

In general, our regret bound takes the form of the traditional regret in a purely online setting plus a benefit term of order $O(-\sqrt{N'_i})$. One might wonder why the adjustment is of order $O(-\sqrt{N'_i})$ instead of $O(-N'_i)$ in Cheung & Lyu (2024), to make the result easier to compare and cover the hybrid MAB results (Cheung & Lyu, 2024), we have the following result from Theorem 1:

Corollary 1. *For the same problem setting and condition with Theorem 1, we can also bound the regret bound by:*

$$\text{Reg}_{\mu^{\text{on}},\alpha,\beta}(T) \leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{\Delta_{\min}^i} - \frac{\sqrt{2}}{K} N'_i \Delta_{\min}^i, 0 \right\} + 4Bm + \frac{\pi^2}{6} \Delta_{\max}, \quad (2)$$

where

$$N'_i = N_i \cdot \max \left\{ 1 - \frac{2BK\omega_i}{\Delta_{\min}^i}, 0 \right\}^2.$$

Corollary 1 is closely aligned with existing results for hybrid MAB, up to constant factors. In particular, when the problem setting reduces to $B = K = 1$, our bound essentially matches the corresponding hybrid MAB results.

Moreover, we observe that $0 \leq \frac{\log T}{\Delta_{\min}^i} - \sqrt{N'_i \log T} \leq \frac{\log T}{\Delta_{\min}^i} - N'_i \Delta_{\min}^i \leq 2 \left(\frac{\log T}{\Delta_{\min}^i} - \sqrt{N'_i \log T} \right)$, which implies that our bound provides a slightly tighter characterization (within a constant factor). Intuitively, the tighter term $\sqrt{N'_i \log T}$ captures the benefit of offline data from an uncertainty reduction perspective, whereas the classical hybrid MAB analysis interprets the benefit primarily through a reduction in the number of required pulls. Despite this difference, the two bounds are of the same order and differ only by a factor of at most 2, indicating that they are asymptotically equivalent.

It is also worth noting that the analytical approaches are fundamentally different: existing hybrid MAB techniques do not directly extend to the CMAB-T setting, while our analysis naturally applies to the hybrid MAB case. The fact that both bounds coincide up to constant factors further suggests that our result is tight.

5.2 Gap-Independent Bound

We then analyze the gap-independent regret upper bound. We obtain two candidate bounds, denoted as ψ and γ , each derived from a different proof technique. The final regret bound takes the minimum among them.

Theorem 2 (Gap-Independent Regret Bound). *For an H-CMAB-T problem $([m], \mathcal{S}, \mathcal{D}, D^{\text{trig}}, R, \mathcal{B})$ that satisfies monotonicity (Condition 1) and TPM bounded smoothness (Condition 2), the hybrid CUCB algorithm with an input bias control vector V and an (α, β) -approximation oracle achieves an (α, β) -approximate gap-independent regret bounded by:*

$$\text{Reg}_{\mu^{\text{on}},\alpha,\beta}(T) \leq \min\{\psi, \gamma\} + 4Bm + \frac{\pi^2}{6} \Delta_{\max}, \quad (3)$$

where ψ and γ are two candidate bounds derived via distinct proof techniques:

$$\psi = 8\sqrt{2}B\sqrt{\log(4mT^3)} \left(\sum_{i \in [m]} \max \left\{ \sqrt{\frac{KT}{m}} - \sqrt{N''_i}, 0 \right\} + \sqrt{mKT} \right), \quad (4)$$

$$\gamma = 16BKT \sqrt{\frac{2 \log(4mT^3)}{\tau_*}} + BKT \omega_{\max}. \quad (5)$$

Here

$$N''_i = N_i \cdot \max \left\{ 1 - \frac{\omega_i}{4\sqrt{2}} \sqrt{\frac{KT}{m \log(4mT^3)}}, 0 \right\}^2, \quad \omega_{\max} = \max_i \omega_i, \quad (6)$$

and τ_* is defined via

$$\begin{aligned} & \max_{\tau, n} \quad \tau \\ \text{s.t.} \quad & \tau \leq N_i + n(i) \text{ where } \tau \in \mathbb{N}, n(i) \in \mathbb{N}, \forall i, \\ & \sum_{i \in [m]} n(i) \leq KT. \end{aligned}$$

These two upper bounds capture different aspects of how offline data can reduce exploration cost in the H-CMAB-T setting. We will interpret each bound, compare their relative strengths, and highlight how they recover or generalize existing results in the literature as follows.

Formally, the first bound ψ involves the quantity N_i'' , defined analogously to N_i' in the gap-dependent bound, and it is interpreted as the amount of *effectively used* offline data. Similarly, the quantity N_i'' embodies the guiding principle behind our algorithmic design in Section 5.1: the more aligned the offline data is with the online environment, the more confidently and extensively it can be incorporated into the learning process. The setting where $N_i'' = 0$ for all i recovers the pure online CMAB-T problem in (Wang & Chen, 2017), and the resulting bound matches their gap-independent result in order. In this sense, ψ generalizes their analysis by quantifying the potential reduction in regret due to informative offline data via an $O(-\sqrt{N_i''})$ saving term. Moreover, it is worth noting that the use of the $\max\{\cdot, 0\}$ operator implies that ψ ranges between a best-case value (when N_i'' is so large that the $\max\{\cdot, 0\} = 0, \forall i$) and a worst-case value (when $N_i'' = 0, \forall i$) matching the pure online regret bound. Specifically, ψ lies between $8B\sqrt{mKT\log(4mT^3)}$ and $16B\sqrt{mKT\log(4mT^3)}$, depending on the informativeness of the offline data. Therefore, although ψ reflects meaningful offline benefits and can cut down half of the regret at the best case, it does not improve the regret order corresponding to the specific problem parameters.

The second bound, γ , is derived via a relaxation of exploration constraints into a covering linear program. The LP solution τ_* appearing in γ satisfies a uniform lower bound $\tau_* \geq KT/m$, which ensures that the first term in γ is always at most the worst case of ψ . It can still be smaller when N_i is large and w_{\max} is small. In some extreme cases where $w_{\max} \leq 1/BKT$ and $N_i \geq (BKT)^2 \log(4mT^3)$, the bound γ tends to be of constant order which is independent of T , highlighting the potential for offline data to fully eliminate exploration cost under perfect alignment. Moreover, γ structurally aligns with recent work on leveraging offline data in the classical MAB setting (Cheung & Lyu, 2024). By setting $K = B = 1$, our H-CMAB-T problem reduces to a hybrid MAB scenario. In this special case, γ recovers (and slightly tightens) Cheung & Lyu (2024): their bound includes a saving term of the form $2TV_{\max}$, whereas ours uses Tw_{\max} with $w_{\max} \leq 2V_{\max}$.

We now compare the two bounds in terms of tightness and interpretability. The bound ψ provides a uniform guarantee and reflects a conservative lower baseline. While it never diverges, it also does not yield a tighter rate even when offline data is abundant. In contrast, γ can become substantially tighter in favorable regimes. When the offline data is highly informative (i.e., large N_i and small ω_i), γ can reduce the regret significantly. For example, in the ideal case of $N_i \geq (BKT)^2 \log(4mT^3)$ and $\omega_{\max} \leq 1/BKT$, the bound tends to be a constant, matching our expectation that regret should vanish when offline information fully resolves arm uncertainty.

Together, these two bounds form a comprehensive characterization of the gap-independent regret in H-CMAB-T. They offer different trade-offs between robustness, interpretability, and tightness, and demonstrate how the size, bias, and coverage of offline data influence the learning performance.

5.3 Lower Bound

Theorem 3 (Regret Lower Bound). *For an H-CMAB-T problem $([m], \mathcal{S}, \mathcal{D}, D^{trig}, R, \mathcal{B})$ that satisfies monotonicity (Condition 1) and TPM bounded smoothness (Condition 2), then any algorithm with an input bias control vector V and an (α, β) -approximation oracle achieves an (α, β) -approximate regret bounded by:*

$$\text{Reg}_{\mu^{\text{on}}, \alpha, \beta}(T) \geq B \sum_{i \in [m]} \left(\frac{cB \log T}{K \Delta_{\min}^i} - \sqrt{cN_i''' \cdot \log T} \right),$$

$$\text{where } N_i''' = N_i \max\left\{1 - \frac{B\omega_i}{K \Delta_{\min}^i}, 0\right\}^2, c = \min\{x \in \{\mu_1, \dots, \mu_m\} : x(1-x)\}$$

The lower bound in Theorem 3 complements our upper bound analysis in several important aspects. First, the bound of N_i matches the gap-dependent saving term in our regret upper bound(Theorem 1) up to constant factors, demonstrating that the improvement achieved by hybrid CUCB is near-optimal in its dependence on the offline data size and bias level.

Second, the lower bound complements existing results for classical hybrid multi-armed bandits. When the CMAB-T problem degenerates to a standard multi-armed bandit with $K = 1$ and deterministic triggering, our analysis yields a gap-dependent lower bound, which is not captured by prior work on hybrid MAB (Cheung & Lyu, 2024), where only gap-independent lower bound is provided. A more detailed comparison is provided in the discussion on Corollary 1.

Finally, Theorem 3 reveals a fundamental limitation of hybrid learning in CMAB-T: even with access to offline observations, the regret cannot be reduced arbitrarily unless the offline data is sufficiently informative and well-aligned with the online environment.

6 Experiments

In this section, we compare our proposed hybrid CUCB with existing CUCB for the pure online setting (Wang & Chen, 2017) and CLCB for the pure offline setting (Liu et al., 2025). To evaluate the performance of CLCB, we first use this algorithm to select an action based on the offline data set and always select this action in the following rounds. For simplicity, we assume that $N_i = N$ and $V_i = V$ for any arm i . Due to the space limit, more details about the reward function and triggering mechanism, as well as the experimental setting and real-world validations, are deferred to appendix.

We evaluate on unbiased offline datasets with varying sizes $N \in \{10, 50, 200\}$. As shown in Figure 1, hybrid CUCB consistently outperforms both online CUCB and offline CLCB. The improvement stems from the warm-start provided by offline data, which reduces early exploration. The advantage becomes more pronounced with larger N , and when N is sufficiently large (e.g., $N = 200$), hybrid CUCB achieves constant regret. Compared to CLCB, the hybrid approach is especially superior when offline data is scarce, since CLCB relies solely on potentially inaccurate offline estimates.

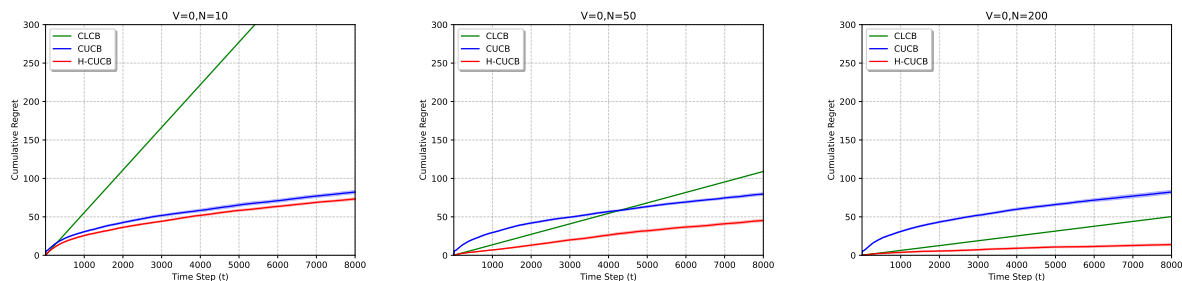


Figure 1: Performance comparison of hybrid CUCB against baselines with unbiased offline data set.

We further evaluate the robustness of the algorithms under distributional bias between the offline and online environments. Specifically, we consider varying levels of bias $V \in \{0.2, 0.3, 0.4\}$, assuming a sufficiently large offline dataset size ($N = 200$) to ensure reliable offline estimates. The results, presented in Figure 2, demonstrate that our hybrid CUCB algorithm consistently outperforms or matches the baseline performance across all tested levels of distributional bias.

7 Conclusion

We introduce H-CMAB-T, a new framework that extends classical CMAB-T by incorporating available offline data into online learning. We propose the hybrid CUCB algorithm, which selectively leverages offline observations via a minimum of two confidence bounds, controlled by a bias-aware mechanism. Theoretically, we established both gap-dependent and gap-independent regret bounds, showing that our method effectively reduces exploration through a data-dependent saving term. Empirical results further corroborate our theoretical findings, demonstrating the effectiveness of the proposed method in benchmark CMAB-T scenarios. The current CMAB-T framework does not naturally handle high-dimensional contexts or side information.

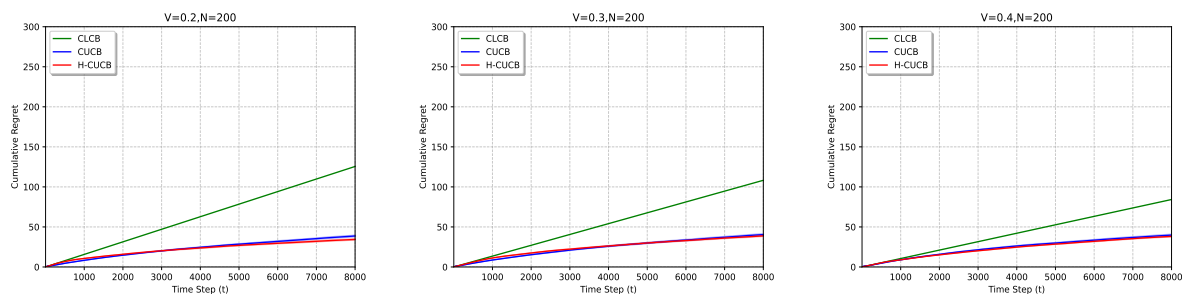


Figure 2: Performance comparison of hybrid CUCB against baselines with the biased offline data set.

Extending hybrid learning to contextual CMAB-T represents a promising direction, with potential for broader applicability in practical scenarios.

Acknowledgements

The corresponding author Fang Kong is supported by National Natural Science Foundation of China (62506150) and Guangdong Basic and Applied Basic Research Foundation (2025A1515011412).

References

- Akhil Agnihotri, Rahul Jain, Deepak Ramachandran, and Zheng Wen. Online bandit learning with offline preference data. *arXiv preprint arXiv:2406.09574*, 2024.
- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- Jinzhi Bu, David Simchi-Levi, and Yunzong Xu. Online pricing with offline data: Phase transition and inverse square law. *Management Science*, 68(12):8568–8588, December 2022. ISSN 0025-1909. doi: 10.48550/arXiv.1910.08693.
- Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. In *Proceedings of the 30th International Conference on Machine Learning (ICML)*, pp. 151–159, 2013.
- Wei Chen, Yajun Wang, Yang Yuan, and Qinshi Wang. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *Journal of Machine Learning Research*, 17(50):1–33, 2016.
- Xinyun Chen, Pengyi Shi, and Shanwen Pu. Data-pooling reinforcement learning for personalized healthcare intervention. *CoRR*, abs/2211.08998, 2022. doi: 10.48550/ARXIV.2211.08998. URL <https://doi.org/10.48550/arXiv.2211.08998>.
- Wang Chi Cheung and Lixing Lyu. Leveraging (biased) information: Multi-armed bandits with offline data. *arXiv preprint arXiv:2405.02594*, 2024. doi: 10.48550/arXiv.2405.02594. Accepted to ICML 2024.
- In Gim, Guojun Chen, Seung Seob Lee, Nikhil Sarda, Anurag Khandelwal, and Lin Zhong. Prompt cache: Modular attention reuse for low-latency inference. *arXiv preprint arXiv:2311.04934*, 2023.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5084–5096. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/jin21e.html>.

- Fang Kong, Jize Xie, Baoxiang Wang, Tao Yao, and Shuai Li. Online influence maximization under decreasing cascade model. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 2197–2204, 2023.
- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. *CoRR*, abs/2107.00591, 2021. URL <https://arxiv.org/abs/2107.00591>.
- Gen Li, Wenhao Zhan, Jason D. Lee, Yuejie Chi, and Yuxin Chen. Reward-agnostic fine-tuning: Provable statistical benefits of hybrid reinforcement learning. *CoRR*, abs/2305.10282, 2023. doi: 10.48550/ARXIV.2305.10282. URL <https://doi.org/10.48550/arXiv.2305.10282>.
- Gene Li, Cong Ma, and Nati Srebro. Pessimism for offline linear contextual bandits using ℓ_p confidence sets. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 20974–20987, 2022.
- Xutong Liu, Jinhang Zuo, Siwei Wang, Carlee Joe-Wong, John Lui, and Wei Chen. Batch-size independent regret bounds for combinatorial semi-bandits with probabilistically triggered arms or independent arms. *Advances in Neural Information Processing Systems*, 35:14904–14916, 2022.
- Xutong Liu, Jinhang Zuo, Siwei Wang, John C.S. Lui, Mohammad Hajiesmaili, Adam Wierman, and Wei Chen. Contextual combinatorial bandits with probabilistically triggered arms. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 22559–22593. PMLR, 23–29 Jul 2023.
- Xutong Liu, Siwei Wang, Jinhang Zuo, Han Zhong, Xuchuang Wang, Zhiyong Wang, Shuai Li, Mohammad Hajiesmaili, John CS Lui, and Wei Chen. Combinatorial multivariate multi-armed bandits with applications to episodic reinforcement learning and beyond. In *International Conference on Machine Learning*, pp. 32139–32172, 2024.
- Xutong Liu, Xiangxiang Dai, Jinhang Zuo, Siwei Wang, Carlee-Joe Wong, John Lui, and Wei Chen. Offline learning for combinatorial multi-armed bandits. *arXiv preprint arXiv:2501.19300*, 2025.
- Thanh Nguyen-Tang, Sunil Gupta, Hung Tran-The, and Svetha Venkatesh. Sample complexity of offline reinforcement learning with deep relu networks. *arXiv preprint arXiv:2103.06671*, 2021.
- Thanh Nguyen-Tang, Sunil Gupta, A Tuan Nguyen, and Svetha Venkatesh. Offline neural contextual bandits: Pessimism, optimization and generalization. In *International Conference on Learning Representations*, 2022.
- Bastian Oetomo, R Malinga Perera, Renata Borovica-Gajic, and Benjamin IP Rubinstein. Cutting to the chase with warm-start contextual bandits. *Knowledge and Information Systems*, 65(9):3533–3565, 2023.
- Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling transformer inference. *arXiv preprint arXiv:2211.05102*, 2022.
- Chengrui Qu, Laixi Shi, Kishan Panaganti, Pengcheng You, and Adam Wierman. Hybrid transfer reinforcement learning: Provable sample efficiency from shifted-dynamics data. In *International Conference on Artificial Intelligence and Statistics*, 2025.
- Guanqiao Qu, Qiyuan Chen, Wei Wei, Zhengyi Lin, Xianhao Chen, and Kaibin Huang. Mobile edge intelligence for large language models: A contemporary survey. *arXiv preprint arXiv:2407.18921*, 2024.
- Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart J. Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *IEEE Transactions on Information Theory*, 68: 8156–8196, 2021.

- Flore Sentenac, Ilbin Lee, and Csaba Szepesvari. Balancing optimism and pessimism in offline-to-online learning. *arXiv preprint arXiv:2502.08259*, 2025. doi: 10.48550/arXiv.2502.08259. URL <https://arxiv.org/abs/2502.08259>.
- Pannagadatta Shivaswamy and Thorsten Joachims. Multi-armed bandit problems with history. In *Artificial intelligence and statistics*, pp. 1046–1054. PMLR, 2012.
- Yuda Song, Yifei Zhou, Ayush Sekhari, Drew Bagnell, Akshay Krishnamurthy, and Wen Sun. Hybrid rl: Using both offline and online data can make rl efficient. In *International Conference on Learning Representations*, 2023.
- Qinshi Wang and Wei Chen. Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30. Curran Associates, Inc., 2017.
- Zheng Wen, Branislav Kveton, Michal Valko, and Sharan Vaswani. Online influence maximization under independent cascade model with semi-bandit feedback. *Advances in neural information processing systems*, 30, 2017.
- Han Zheng, Xufang Luo, Pengfei Wei, Xuan Song, Dongsheng Li, and Jing Jiang. Adaptive policy learning for offline-to-online reinforcement learning. In Brian Williams, Yiling Chen, and Jennifer Neville (eds.), *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pp. 11372–11380. AAAI Press, 2023. doi: 10.1609/AAAI.V37I9.26345. URL <https://doi.org/10.1609/aaai.v37i9.26345>.
- Banghua Zhu, Ying Sheng, Lianmin Zheng, Clark Barrett, Michael Jordan, and Jiantao Jiao. Towards optimal caching and model selection for large model inference. *Advances in Neural Information Processing Systems*, 36:59062–59094, 2023.

Appendix

A Technical Lemmas

Definition 3 (Event-Filtered Regret (Wang & Chen, 2017)). *For any series of events $\{\mathcal{E}_t\}_{t \geq 1}$ indexed by round number t , we define $Reg_{\mu, \alpha}^A(T, \{\mathcal{E}_t\}_{t \geq 1})$ as the regret filtered by events $\{\mathcal{E}_t\}_{t \geq 1}$, that is, regret is only counted in round t if \mathcal{E}_t happens in round t . Formally,*

$$Reg_{\mu, \alpha}^A(T, \{\mathcal{E}_t\}_{t \geq 1}) = \mathbb{E} \left[\sum_{t=1}^T \mathbb{I}(\mathcal{E}_t) (\alpha \cdot \text{opt}_{\mu} - r_{\mu}(S_t^A)) \right].$$

For convenience, A , α , μ and/or T can be omitted when the context is clear, and we simply use $Reg_{\mu, \alpha}^A(T, \mathcal{E}_t)$ instead of $Reg_{\mu, \alpha}^A(T, \{\mathcal{E}_t\}_{t \geq 1})$.

The regret upper bound relies on considering the following events of accurate estimations by $UCB_t(i)$ and $UCB_t^S(i)$. For every t , define:

$$\mathcal{N}_t = \bigcap_{i \in [m]} (\mathcal{N}_t(i) \cap \mathcal{N}_t^S(i)), \quad \text{where}$$

$$\mathcal{N}_t(i) = \{\mu_i^{\text{on}} \leq UCB_t(i) \leq \mu_i^{\text{on}} + 2 \text{rad}_t(i)\},$$

$$\mathcal{N}_t^S(i) = \left\{ \begin{array}{l} \mu_i^{\text{on}} \leq UCB_t^S(i) \leq \mu_i^{\text{on}} + \text{rad}_t^S(i) \\ + \left[\sqrt{\frac{2 \log(2t/\delta_t)}{N_i + T_{i,t-1}}} + \frac{N_i \cdot (\mu_i^{\text{off}} - \mu_i^{\text{on}})}{N_i + T_{i,t-1}} \right] \end{array} \right\}.$$

Lemma 1. (Cheung & Lyu, 2024) *For the event \mathcal{N}_t defined above, we have $\Pr(\mathcal{N}_t) \geq 1 - 2m\delta_t$. During the part of discussion on regret bound, we set $\delta_t = 1/(2mt^2)$ for $t = 1, 2, \dots, T$.*

In the following, we provide the useful lemmas for the gap-independent regret bound. We firstly define two linear program (IP) and (LP):

$$\begin{aligned} \text{IP} : \quad & \max_{C_T(i), i \in [m]} \sum_{i \in [m]} \sum_{n(i)=1}^{C_T(i)} \sqrt{\frac{1}{N_i + n(i)}} \\ & \text{s.t.} \quad \sum_{i \in [m]} C_T(i) \leq KT, \\ & \quad \quad C_T(i) \in \mathbb{N}^+ \quad \forall i \in [m]. \end{aligned}$$

$$\begin{aligned} \text{LP} : \quad & \max_{\tau, n} \quad \tau \\ & \text{s.t.} \quad \tau \leq N_i + n(i) \quad \forall i \in [m], \\ & \quad \quad \sum_{i \in [m]} n(i) \leq KT, \\ & \quad \quad \tau \geq 0, \quad n(i) \geq 0 \quad \forall i \in [m]. \end{aligned}$$

And we suppose that $(C_T^*(i))_{i \in [m]} \in \mathbb{N}_{\geq 0}^m$ and $(\tau_*, \{n_*(i)\}_{i \in [m]})$ are the solution to (IP) and (LP) correspondingly.

Lemma 2. *For the LP defined above, we have $n_*(i) = \max\{\tau_* - N_i, 0\}$, $\forall i \in [m]$.*

Proof. Since optimal solutions must be feasible, then we have $n_*(i) \geq \max\{\tau_* - N_i, 0\}$, $\forall i \in [m]$. We only need to prove that if \exists such an arm i' that $n_*(i') - \max\{\tau_* - N_{i'}, 0\} = \epsilon > 0$, it will bring into a contradiction statement. In fact, we can construct another solution $(\tau', \{n'(i)\}_{i \in [m]})$ by this immediately:

$$\begin{cases} \tau' = \tau_* + \frac{\epsilon}{m} \\ n'(i) = \begin{cases} n_*(i) + \frac{\epsilon}{m} & \text{if } i \neq i' \\ n_*(i) - \frac{m-1}{m} \cdot \epsilon & \text{if } i = i'. \end{cases} \end{cases}$$

Then we have $\tau' > \tau_*$, which contradicts the optimality of τ_* . \square

Lemma 3. For the LP and IP defined above, we have $C_T^*(i) \leq \max\{\lceil \tau_* \rceil - N_i, 0\}$, $\forall i \in [m]$.

Proof. Suppose that there exists an arm i' such that $C_T^*(i') \geq \max\{\lceil \tau_* \rceil - N_{i'}, 0\} + 1$, then there must exist another arm $i'' \neq i'$ such that $C_T^*(i'') \leq \max\{\lceil \tau_* \rceil - N_{i''}, 0\} - 1$, or we will have:

$$\sum_{i \in [m]} C_T^*(i) > \sum_{i \in [m]} \max\{\lceil \tau_* \rceil - N_i, 0\} \geq \sum_{i \in [m]} \max\{\tau_* - N_i, 0\} \stackrel{(a)}{=} \sum_{i \in [m]} n_*(i) = KT,$$

which contradicts the constraint of (LP). Here, (a) is from lemma 2. As a result, we can construct a feasible solution $\tilde{C}_T(i)_{i \in [m]} \in \mathbb{N}_{\geq 0}^m$ by the existence of two arms i' and i'' that:

$$\tilde{C}_T(i) = \begin{cases} C_T^*(i) - 1 & \text{if } i = i' \\ C_T^*(i) + 1 & \text{if } i = i'' \\ C_T^*(i) & \text{if } i \in [m] \setminus \{i', i''\}. \end{cases}$$

By the property that $C_T^*(i') \geq 1$, $\tilde{C}_T(i') \geq 0$, and $(\tilde{C}_T(i))_{i \in [m]}$ is a feasible solution. But then we have

$$\begin{aligned} & \sum_{i \in [m]} \sum_{n(i)=1}^{\tilde{C}_T(i)} \sqrt{\frac{1}{n(i) + N_i}} - \sum_{i \in [m]} \sum_{n(i)=1}^{C_T^*(i)} \sqrt{\frac{1}{n(i) + N_i}} \\ &= \sqrt{\frac{1}{C_T^*(i'') + N_{i''} + 1}} - \sqrt{\frac{1}{C_T^*(i') + N_{i'}}} \\ &\geq \sqrt{\frac{1}{\lceil \tau_* \rceil}} - \sqrt{\frac{1}{\max\{\lceil \tau_* \rceil, N_{i'}\} + 1}} > 0, \end{aligned}$$

which contradicts the assumed optimality of $(N_T^*(i))_{i \in [m]}$. \square

B Proof of Theorem 1

We define the event

$$F_t = \{r_{S_t}(\bar{\mu}) \leq \alpha \cdot \text{opt}_{\bar{\mu}}\},$$

which captures that the oracle output based on the estimated means $\bar{\mu}$ at round t achieves at least an α -approximation of the optimal reward.

Let the filtration \mathcal{F}_{t-1} represent all the history observed up to and including the decision S_t , formally:

$$\mathcal{F}_{t-1} = (S_1, \tau_1, \{X_{1,i} : i \in \tau_1\}, \dots, S_{t-1}, \tau_{t-1}, \{X_{t-1,i} : i \in \tau_{t-1}\}, S_t).$$

Here, τ_s denotes the triggered set at round s , and $X_{s,i}$ is the observed reward for arm i in round s if triggered. We emphasize that the filtration \mathcal{F}_{t-1} already implicitly incorporates the information from the offline data. In particular, the observations of arm i offline affect the initialization of arm statistics such as rad_t^S and UCB_t^S , which in turn influence the selection of S_t at each round t . Therefore, the subsequent triggered sets τ_t and observed rewards $\{X_{t,i} : i \in \tau_t\}$ are also conditioned on the offline data through the choice of S_t .

The conditional expectation at round t is defined as

$$\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot \mid \mathcal{F}_{t-1}],$$

which aligns with the algorithm's access to the complete history \mathcal{F}_{t-1} when making decisions at round t . Moreover, quantities such as S_t and $\bar{\mu}_{i,t}$ are \mathcal{F}_{t-1} -measurable.

Proof. Since the μ_i^{on} is the actual mean we focus to learn about for every arm i , we set $\mu_i = \mu_i^{\text{on}}$ for every arm i for simplicity. To unify the proofs for the proof of Theorem 1 and the proof of bound ψ in Theorem 2, we introduce a positive parameter M_i for every arm i , which is introduced in Wang & Chen (2017). We also further inherit the definition $M_S := \max_{i \in \tilde{S}} M_i$ for each action S and $M_S = 0$ if $\tilde{S} = \emptyset$ from Wang & Chen (2017).

We first show that if $\{S_t \geq M_{S_t}\}$, \mathcal{N}_t and $\neg F_t$, and given filtration \mathcal{F}_{t-1} . Recall that $\Delta_S = \max(0, \alpha \cdot \text{opt}_\mu - r_S(\boldsymbol{\mu}))$, we have:

$$\Delta_{S_t} = \mathbb{E}_t[\Delta_{S_t}] \leq \mathbb{E}_t \left[2B \sum_{i \in \tilde{S}_t} \left[p_i^{D, S_t} (\bar{\mu}_{i,t} - \mu_i) - \frac{M_i}{2BK} \right] \right] \quad (7)$$

$$\leq \mathbb{E}_t \left[2B \sum_{i \in \tilde{S}_t} p_i^{D, S_t} \left[(\bar{\mu}_{i,t} - \mu_i) - \frac{M_i}{2BK} \right] \right] \quad (8)$$

$$= \mathbb{E}_t \left[2B \sum_{i \in \tilde{S}_t} \mathbb{I}\{i \in \tau_t\} \left[(\bar{\mu}_{i,t} - \mu_i) - \frac{M_i}{2BK} \right] \right] \quad (9)$$

$$= \mathbb{E}_t \left[2B \sum_{i \in \tau_t} \left[(\bar{\mu}_{i,t} - \mu_i) - \frac{M_i}{2BK} \right] \right], \quad (10)$$

where (7) comes from exactly the equation (11) of Appendix B.3 in Wang & Chen (2017), (8) comes from the fact that $p_i^{D, S_t} \leq 1$ for every arm i , (9) follows from the fact that since the algorithm choose S_t using the information of offline data, then S_t and $\bar{\mu}_{i,t}$ are \mathcal{F}_{t-1} measurable and the only randomness is the triggering set τ_t at round t , which satisfies the conditions of TPE trick in Liu et al. (2023), so we can also use TPE trick (Liu et al., 2023) to replace $p_i^{D, S_t} = \mathbb{E}_t[\mathbb{I}\{i \in \tau_t\}]$, and (10) is the change of notion τ_t .

Then we use κ to describe the concentration for $(\bar{\mu}_{i,t} - \mu_i) - M_i/(2BK)$, the intuition is that we use different UCBs depending on how informative the offline data is.

Case 1 When $\omega_i < M_i/(2BK)$, let

$$\kappa_{T, N_i, \omega_i}(M_i, s) = \begin{cases} 4B, & s = N_i = 0 \\ 4B \sqrt{\frac{2 \log(4mT^3)}{N_i + s}}, & 0 \leq s \leq \ell_{T, N_i, \omega_i}(M_i) \\ 0, & s > \ell_{T, N_i, \omega_i}(M_i) \text{ or } \ell_{T, N_i, \omega_i}(M_i) \leq 0, \end{cases}$$

where

$$\ell_{T, N_i, \omega_i}(M_i) = \frac{64B^2 K^2 \log(4mT^3)}{M_i^2} - N_i \cdot \max\left\{1 - \frac{2BK\omega_i}{M_i}, 0\right\}^2.$$

We first prove when $\omega_i < M_i/(2BK)$, we have $(\bar{\mu}_{i,t} - \mu_i) - M_i/(2BK) \leq \kappa_{T,N_i,\omega_i}(M_i, t)$:

Since $\omega_i < \frac{M_i}{2BK}$, $\max\{1 - \frac{2BK\omega_i}{M_i}, 0\}^2 = \left(1 - \frac{2BK\omega_i}{M_i}\right)^2 > 0$,

1. if $N_i \cdot \left(1 - \frac{2BK\omega_i}{M_i}\right)^2 \geq \frac{32B^2K^2 \log(4mT^3)}{M_i^2}$, then we have:

$$2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}} \leq 2\sqrt{\frac{2\log(4mT^3)}{N_i}} \leq \left(1 - \frac{2BK\omega_i}{M_i}\right) \cdot \frac{M_i}{2BK}. \quad (11)$$

$$\begin{aligned} \Rightarrow \bar{\mu}_{i,t} &\stackrel{(a)}{\leq} \text{UCB}_t^S(i) \stackrel{(b)}{\leq} \mu_i + 2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}} + \frac{N_i}{N_i + T_{i,t-1}} \cdot \omega_i \\ &\stackrel{(c)}{\leq} \mu_i + \left(1 - \frac{2BK\omega_i}{M_i}\right) \cdot \frac{M_i}{2BK} + \omega_i \\ &= \mu_i + \frac{M_i}{2BK}, \end{aligned}$$

$$\Rightarrow \bar{\mu}_{i,t} - \mu_i + \frac{M_i}{2BK} \leq 0 \leq \kappa_{T,N_i,\omega_i}(M_i, t)$$

where (a) comes from the definition of $\bar{\mu}_{i,t}$ in Algorithm 1, (b) follows from the lemma 1 and the definition of ω_i , and (c) follows from (11) and $N_i/(N_i + T_{i,t-1}) \leq 1$.

2. when $N_i \cdot \left(1 - \frac{2BK\omega_i}{M_i}\right)^2 < \frac{32B^2K^2 \log(4mT^3)}{M_i^2}$, then we have:

$$\ell_{T,N_i,\omega_i}(M_i) = \frac{64B^2K^2 \log(4mT^3)}{M_i^2} - N_i \max\left(1 - \frac{2BK\omega_i}{M_i}, 0\right)^2 > \frac{32B^2K^2 \log(4mT^3)}{M_i^2}.$$

(1) when $T_{i,t-1} > \ell_{T,N_i,\omega_i}(M_i) > \frac{32B^2K^2 \log(4mT^3)}{M_i^2}$, we have :

$$\begin{aligned} \bar{\mu}_{i,t} &\stackrel{(a)}{\leq} \text{UCB}_t(i) \stackrel{(b)}{\leq} \mu_i + 2\sqrt{\frac{2\log(4mt^3)}{T_{i,t-1}}} \\ &\stackrel{(c)}{<} \mu_i + \frac{M_i}{2BK}, \end{aligned} \quad (12)$$

$$\Rightarrow \bar{\mu}_{i,t} - \mu_i + \frac{M_i}{2BK} \leq 0 \leq \kappa_{T,N_i,\omega_i}(M_i, t)$$

where (a) comes from the definition of $\bar{\mu}_{i,t}$ in Algorithm 1, (b) follows from lemma 1, and (c) follows from the condition that $T_{i,t-1} > \frac{32B^2K^2 \log(4mT^3)}{M_i^2}$.

(2) when $0 \leq T_{i,t-1} \leq \ell_{T,N_i,\omega_i}(M_i)$, then we have:

$$\begin{aligned} \bar{\mu}_{i,t} - \mu_i - \frac{M_i}{2BK} &\stackrel{(a)}{\leq} \text{UCB}_t^S(i) - \mu_i - \frac{M_i}{2BK} \\ &\stackrel{(b)}{\leq} 2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}} + \frac{N_i}{N_i + T_{i,t-1}} \omega_i - \frac{M_i}{2BK} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{\leq} 2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}} + \frac{N_i}{N_i + T_{i,t-1}} \cdot \frac{M_i}{2BK} - \frac{M_i}{2BK} \\
&= 2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}} + \left(\frac{N_i}{N_i + T_{i,t-1}} - 1\right) \cdot \frac{M_i}{2BK} \\
&\stackrel{(d)}{\leq} 2\sqrt{\frac{2\log(4mT^3)}{N_i + T_{i,t-1}}},
\end{aligned}$$

where (a) comes from the definition of $\bar{\mu}_{i,t}$ in Algorithm 1, (b) follows from lemma 1, (c) follows from the condition that $\omega_i < M_i/(2BK)$ and (d) follows from $N_i/(N_i + T_{i,t-1}) - 1 \leq 0$.

Case 2 When $\omega_i > M_i/(2BK)$, we firstly define:

$$\kappa_T(M_i, s) = \begin{cases} 4B, & s = 0 \\ 4B\sqrt{\frac{2\log(4mT^3)}{s}}, & 1 \leq s \leq \ell_T(M_i) \\ 0, & s > \ell_T(M_i), \end{cases}$$

where

$$\ell_T(M_i) = \frac{32B^2K^2\log(4mT^3)}{M_i^2}.$$

Since $\omega_i \geq \frac{M_i}{2BK}$, $\max\{1 - \frac{2BK\omega_i}{M_i}, 0\}^2 = 0$. Follow the similar analysis in **case 1**, we have:

$$\begin{aligned}
\bar{\mu}_{i,t} - \mu_i - \frac{M_i}{2BK} &\leq \text{UCB}_t(i) - \mu_i - \frac{M_i}{2BK} \\
&\leq 2\sqrt{\frac{2\log(4mT^3)}{T_{i,t-1}}} - \frac{M_i}{2BK} \\
&\stackrel{(a)}{\leq} \kappa_T(M_i, s),
\end{aligned}$$

where (a) follows from: if $T_{i,t-1} > \ell_T(M_i) = \frac{32B^2K^2\log(4mT^3)}{M_i^2}$, then from (12) we have $\bar{\mu}_{i,t} - \mu_i - M_i/(2BK) \leq 0$, and if $T_{i,t-1} \leq \ell_T(M_i)$ we have $2\sqrt{\frac{2\log(4mT^3)}{T_{i,t-1}}} - \frac{M_i}{2BK} \leq 2\sqrt{\frac{2\log(4mT^3)}{T_{i,t-1}}}$.

Notice that if we set $N_i = 0$, then the definition of κ in **case 1** can cover the definition of κ in **case 2**. As a result, we use $\kappa_{T, N_i, \omega_i}(M_i, s)$ for the following statement for simplicity. From the **case 1** and **case 2**, then we can derive the regret into two parts:

$$\begin{aligned}
\text{Reg}(\{S_t \geq M_{S_t}\}, \mathcal{N}_t, \neg F_t) &= \mathbb{E} \left[\sum_{t=1}^T \Delta_{S_t} \right] \\
&\stackrel{(a)}{\leq} \mathbb{E} \left[\sum_{t=1}^T \mathbb{E}_t \left[\sum_{i \in \tau_t} \kappa_{T, N_i, \omega_i}(M_i, T_{i,t-1}) \right] \right] \\
&\stackrel{(b)}{=} \mathbb{E} \left[\sum_{t=1}^T \sum_{i \in \tau_t} \kappa_{T, N_i, \omega_i}(M_i, T_{i,t-1}) \right] \\
&\stackrel{(c)}{=} \mathbb{E} \left[\sum_{i \in [m]} \sum_{s=0}^{T_{-1,i}} \kappa_{T, N_i, \omega_i}(M_i, s) \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{i \in [m]} \sum_{s=0}^{T_{t-1,i}} \kappa_{T,N_i,\omega_i}(M_i, s) \\
&= \sum_{t=1}^T \left[\sum_{i \in \tilde{S}_t, \omega_i \leq M_i/(2BK)} \kappa_{T,N_i,\omega_i}(M_i, t) + \sum_{i \in \tilde{S}_t, \omega_i > M_i/(2BK)} \kappa_T(M_i, t) \right] \\
&:= \underline{A} + \underline{B},
\end{aligned}$$

where (a) follows from the discussion on **case 1** and **case 2**, (b) follows from the tower rule, (c) follows from that $T_{t-1,i}$ is increased by 1 if and only if the arm i is triggered at round t .

We then compute $\underline{A} + \underline{B}$:

For part \underline{A} :

$$\underline{A} = \sum_{\omega_i \leq M_i/2BK} \sum_{s=0}^{\ell_{T,N_i,\omega_i}(M_i)} \kappa_{T,N_i,\omega_i}(M_i, s).$$

For simplicity we set $N'_i = N_i \cdot \max\{1 - \frac{2BK\omega_i}{M_i}, 0\}^2$. Since when $N'_i \geq \ell_{T,N_i,\omega_i}(M_i)$ $\kappa_{T,N_i,\omega_i}(M_i, s) = 0$, we consider the case that $N'_i < \ell_{T,N_i,\omega_i}(M_i)$:

$$\begin{aligned}
\sum_{s=0}^{\ell_{T,N_i,\omega_i}(M_i)} \kappa_{T,N_i,\omega_i}(M_i, s) &= 4B\sqrt{2\log(4mT^3)} \sum_{s=1}^{\ell_{T,N_i,\omega_i}(M_i)} \frac{1}{\sqrt{N_i + s}} ds + 4B \\
&\stackrel{(a)}{\leq} 4B\sqrt{2\log(4mT^3)} \int_0^{\ell_{T,N_i,\omega_i}(M_i)} \frac{1}{\sqrt{N_i + s}} ds + 4B \\
&\stackrel{(b)}{\leq} 4B\sqrt{2\log(4mT^3)} \int_0^{\ell_{T,N_i,\omega_i}(M_i)} \frac{1}{\sqrt{N'_i + s}} ds + 4B \\
&= \frac{64\sqrt{2}B^2K\log(4mT^3)}{M_i} - 8B\sqrt{2N'_i\log(4mT^3)} + 4B,
\end{aligned}$$

where (a) is by the sum & integral inequality $\int_{L-1}^U f(x)dx \geq \sum_{i=L}^U f(i) \geq \int_L^{U+1} f(x)dx$ for non-increasing function f , (b) follows from $N'_i \leq N_i$ and the monotonicity of integrals.

For part \underline{B} , similar as part \underline{A} , we have:

$$\underline{B} = \sum_{\omega_i > M_i/2BK} \sum_{s=0}^{\ell_T(M_i)} \kappa_T(M_i, s) \leq \sum_{\omega_i > M_i/2BK} \left(\frac{64\sqrt{2}B^2K\log(4mT^3)}{M_i} + 4B \right).$$

Sum up $\underline{A} + \underline{B}$, plus the case that N'_i may be $\geq \ell_{T,N_i,\omega_i}(M_i) := \frac{64B^2K^2\log 4mT^3}{M_i^2}$, we have

$$\text{Reg}(\{S_t \geq M_{S_t}\}, \mathcal{N}_t, -F_t) \leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K\log(4mT^3)}{M_i} - 8B\sqrt{2N'_i\log(4mT^3)}, 0 \right\} + 4Bm. \quad (13)$$

For the gap-dependent bound, take $M_i = \Delta_{\min}^i$, then $\text{Reg}(S_t < M_{S_t}) = 0$. And following Wang & Chen (2017) to handle small probability events $-\mathcal{N}_t$ and F_t we have

$$\text{Reg}(T) \leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K\log(4mT^3)}{\Delta_{\min}^i} - 8B\sqrt{2N'_i\log(4mT^3)}, 0 \right\} + 4Bm + \frac{\pi^2}{6} \Delta_{\max}, \quad (14)$$

where

$$N'_i = N_i \cdot \max \left\{ 1 - \frac{2BK\omega_i}{\Delta_{\min}^i}, 0 \right\}^2.$$

□

C Proof of Corollary 1

Recall that the regret bound in Theorem 1 contains the term

$$\max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{\Delta_{\min}^i} - 8B\sqrt{2N'_i \log(4mT^3)}, 0 \right\}.$$

For the non-trivial case where the above maximum is attained by the first argument, we must have

$$\frac{64\sqrt{2}B^2K \log(4mT^3)}{\Delta_{\min}^i} \geq 8B\sqrt{2N'_i \log(4mT^3)}.$$

Rearranging the inequality gives

$$\sqrt{N'_i} \leq \frac{8\sqrt{2}BK}{\Delta_{\min}^i}.$$

Substituting this upper bound on $\sqrt{N'_i}$ into the square-root term, we obtain

$$8B\sqrt{2N'_i \log(4mT^3)} \geq \frac{\sqrt{2}}{K} N'_i \Delta_{\min}^i,$$

which yields the linear form in N'_i .

Plugging this back into the original regret bound leads to

$$\max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{\Delta_{\min}^i} - \frac{\sqrt{2}}{K} N'_i \Delta_{\min}^i, 0 \right\}.$$

Summing over all arms completes the proof. □

D Proof of Theorem 2

To prove Theorem 2, we present two candidate regret bounds, each derived via a distinct analysis technique. We denote these bounds as ψ and γ , and show that the regret is upper bounded by the minimum of the two.

D.1 Proof of Bound ψ

Proof. We further discuss (13) and (14). For the gap-independent bound, take $M_i = M = \sqrt{64\sqrt{2}mB^2K \log(4mT^3)}/T$, then $\text{Reg}(S_t < M_{S_t}) \leq TM$. (Naturally the N'_i would change correspondingly.) Then we have

$$\begin{aligned} \text{Reg}(T) &\leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{M_i} - 8B\sqrt{2N''_i \log(4mT^3)}, 0 \right\} + \text{Reg}(S_t < M_{S_t}) \\ &\leq \sum_{i \in [m]} \max \left\{ \frac{64\sqrt{2}B^2K \log(4mT^3)}{M_i} - 8B\sqrt{2N''_i \log(4mT^3)}, 0 \right\} + TM \\ &\leq 8\sqrt{2}B\sqrt{\log(4mT^3)} \left(\sum_{i \in [m]} \max \left\{ \sqrt{\frac{KT}{m}} - \sqrt{N''_i}, 0 \right\} + \sqrt{mKT} \right), \end{aligned}$$

where

$$N_i'' = N_i \cdot \max \left\{ 1 - \frac{\omega_i}{4\sqrt{2}} \sqrt{\frac{KT}{m \log(4mT^3)}}, 0 \right\}^2.$$

□

D.2 Proof of Bound γ

Intuition. The key idea behind the γ bound lies in adopting a different perspective for establishing early stopping conditions. In the gap-dependent analysis, the number of times each arm i needs to be triggered is directly related to its gap. However, when such gap information is unavailable, we must seek alternative ways to characterize how offline data effectively reduces the required online exploration for each arm.

To this end, we observe that the regret incurred by an arm depends on both the amount of offline data N_i and the number of times it is triggered online T_i . Motivated by this and [Cheung & Lyu \(2024\)](#), a formulation based on a linear program is obtained, which captures how much exploration can be saved through leveraging informative offline data, even without explicit gap knowledge.

Proof. Under the events \mathcal{N}_t and $\neg F_t$, and given filtration \mathcal{F}_{t-1} , follow the similar analysis from (7) to (10) we have:

$$\begin{aligned} \Delta_{S_t} &= \mathbb{E}_t[\Delta_{S_t}] \leq \mathbb{E}_t \left[B \sum_{i \in \tilde{S}_t} p_i^{D, S_t} (\bar{\mu}_{i,t} - \mu_i) \right] \\ &= \mathbb{E}_t \left[B \sum_{i \in \tilde{S}_t} \mathbb{I}\{i \in \tau_t\} [(\bar{\mu}_{i,t} - \mu_i)] \right] \\ &= \mathbb{E}_t \left[B \sum_{i \in \tau_t} (\bar{\mu}_{i,t} - \mu_i) \right]. \end{aligned} \tag{15}$$

Focusing on the regret analysis on UCB^S then we have:

$$\begin{aligned} \text{Reg}(\{S_t \geq M_{S_t}\}, \mathcal{N}_t, \neg F_t) &= \mathbb{E} \left[\sum_{t=1}^T \Delta_{S_t} \right] \\ &\stackrel{(a)}{\leq} \mathbb{E} \left[\sum_{t=1}^T \mathbb{E}_t \sum_{i \in \tau_t} B(\bar{\mu}_{i,t} - \mu_i) \right] \\ &\stackrel{(b)}{\leq} \mathbb{E} \left[B \sum_{t=1}^T \sum_{i \in \tau_t} (\text{UCB}_t^S(i) - \mu_i) \right] \\ &\leq B \sum_{t=1}^T \sum_{i \in \tau_t} (\text{UCB}_t^S(i) - \mu_i), \end{aligned}$$

where (a) follows from (15), and (b) follows from the tower rule and $\bar{\mu}_{i,t} \leq \text{UCB}_t^S(i)$. Follow from the lemma 1, we have:

$$B \sum_{t=1}^T \sum_{i \in \tau_t} (\text{UCB}_t^S(i) - \mu_i)$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} B \sum_{t=1}^T \sum_{i \in \tau_t} \left(2\sqrt{\frac{2 \log(4mt^3)}{N_i + T_{i,t-1}}} + \frac{N_i}{N_i + T_{i,t-1}} \omega_i \right) \\
&\stackrel{(b)}{\leq} 2B \sum_{t=1}^T \sum_{i \in \tau_t} \sqrt{\frac{2 \log(4mt^3)}{N_i + T_{i,t-1}}} + BKT\omega_{\max} \\
&\stackrel{(c)}{=} 2B \sum_{i \in [m]} \sum_{n(i)=0}^{C_T(i)} \sqrt{\frac{2 \log(4mt^3)}{N_i + n(i)}} + BKT\omega_{\max} \\
&\stackrel{(d)}{\leq} 2B \sum_{i \in [m]} \sum_{n(i)=1}^{C_T(i)} \sqrt{\frac{2 \log(4mt^3)}{N_i + n(i)}} + BKT\omega_{\max} + 2Bm, \tag{16}
\end{aligned}$$

Where (a) is from lemma 1, (b) follows from $N_i + T_{i,t-1} \geq N_i$ and the definition of K , $C_T(i)$ in (c) is an undetermined coefficient discuss next, and (d) considers the case that $n(i) + N_i$ may equal to zero, which in that case we treat the log-term as one as algorithm 1 designed in line 7.

And we use linear program to consider the $C_T(i)$:

$$\begin{aligned}
\sum_{i \in [m]} \sum_{n(i)=1}^{C_T(i)} \frac{1}{\sqrt{N_i + n(i)}} &\stackrel{(a)}{\leq} \sum_{i \in [m]} \sum_{n(i)=1}^{C_T^*(i)} \frac{1}{\sqrt{N_i + n(i)}} \\
&\stackrel{(b)}{\leq} \sum_{i \in [m]} \sum_{n(i)=1}^{\max\{\lceil \tau_* \rceil - N_i, 0\}} \frac{1}{\sqrt{N_i + n(i)}} \\
&\leq \sum_{i \in [m]} \frac{\max\{\lceil \tau_* \rceil - N_i, 0\}}{\lceil \tau_* \rceil} \sum_{t=1}^{\lceil \tau_* \rceil} \frac{1}{\sqrt{t}} \\
&\leq \sum_{i \in [m]} \max\{\lceil \tau_* \rceil - N_i, 0\} \cdot \frac{4}{\sqrt{\tau_*}} \\
&\leq \sum_{i \in [m]} (\max\{\tau_* - N_i, 0\} + 1) \cdot \frac{4}{\sqrt{\tau_*}} \\
&\stackrel{(c)}{\leq} \sum_{i \in [m]} (n_*(i) + 1) \cdot \frac{4}{\sqrt{\tau_*}} \\
&\leq \frac{8KT}{\sqrt{\tau_*}}, \tag{17}
\end{aligned}$$

where (a) comes from the definition of (LP), (b) from lemma 3, (c) follows from the feasibility of $(\tau_*, \{n_*(i)\}_{i \in [m]})$ to (LP).

Combine (16) and (17), and following Wang & Chen (2017) to handle small probability events $\neg \mathcal{N}_t$ and F_t then we the final regret bound:

$$\text{Reg}(T) \leq 16BT \sqrt{\frac{2 \log(4mT^3)}{\tau_*}} + BKT\omega_{\max} + 4Bm + \frac{\pi^2}{6} \Delta_{\max}.$$

□

E Proof of Theorem 3

Step 1. Problem Setup.

We consider a CMAB-T instance with m base arms. An action (super-arm) is a subset $S \subseteq [m]$ with fixed size $|S| = K$.

Triggering distribution. When action S is played at round t , exactly one arm in S is triggered:

$$p_{S,i} = \Pr(i \in \tau_t \mid S_t = S) = \begin{cases} 1/K, & i \in S, \\ 0, & i \notin S. \end{cases}$$

Let T_i denote the number of rounds in which arm i is triggered.

Reward model. When arm i is triggered, its outcome is $X_i^{(t)} \sim \text{Bern}(\mu_i)$. The reward is defined as

$$R(S, X, \tau) = B \cdot \sum_{i \in \tau} X_i = B \cdot X_{j^*},$$

where j^* is the unique triggered arm.

Thus the expected reward of action S under mean vector μ is

$$r_S(\mu) = \mathbb{E}[R(S, X, \tau)] = B \sum_{j \in S} p_{S,j} \mu_j = \frac{B}{K} \sum_{j \in S} \mu_j.$$

Verification of TPM smoothness. For any two mean vectors μ, μ' ,

$$\begin{aligned} |r_S(\mu) - r_S(\mu')| &= \left| \frac{B}{K} \sum_{j \in S} (\mu_j - \mu'_j) \right| \\ &\leq \frac{B}{K} \sum_{j \in S} |\mu_j - \mu'_j| = B \sum_j p_{S,j} |\mu_j - \mu'_j|. \end{aligned}$$

Therefore this instance satisfies the TPM condition with constant B .

Gaps. Assume the optimal base arm is arm 1, and its mean dominates:

$$\mu_1 = \mu_2 = \dots = \mu_K \geq \mu_{K+1} \geq \dots \geq \mu_m.$$

The optimal action is $S = \{1, \dots, K\}$ with reward

$$r_S(\mu) = \frac{B}{K} \cdot K \mu_1 = B \mu_1.$$

For any suboptimal arm $i > K$, let $A_i \subseteq [K]$ be an arbitrary subset such that $|A_i| = K - 1$, and define

$$S_i := A_i \cup \{i\}.$$

That is, S_i consists of $K - 1$ optimal base arms and one suboptimal arm i . Since

$$\mu_1 = \mu_2 = \dots = \mu_K,$$

the expected reward of S_i does not depend on the particular choice of A_i . And as a result, its expected reward is

$$r_{S_i}(\mu) = \frac{B}{K} ((K - 1)\mu_1 + \mu_i).$$

Thus the action gap is

$$\Delta_{S_i} = r_S(\mu) - r_{S_i}(\mu) = \frac{B}{K} (\mu_1 - \mu_i) = \frac{B}{K} g_i,$$

where we define the base-arm gap $g_i := \mu_1 - \mu_i$.

Regret decomposition. Let N_i^{act} be the number of times the algorithm selects action S_i . Then the regret contributed by arm i is

$$\text{Reg}_i(T) = \Delta_{S_i} \cdot \mathbb{E}_\nu[N_i^{\text{act}}] = \frac{B}{K} g_i \mathbb{E}_\nu[N_i^{\text{act}}].$$

Because arm i is triggered with probability $1/K$ under S_i , we have

$$\mathbb{E}_\nu[T_i] = \frac{1}{K} \mathbb{E}_\nu[N_i^{\text{act}}] \Leftrightarrow \mathbb{E}_\nu[N_i^{\text{act}}] = K \cdot \mathbb{E}_\nu[T_i].$$

Thus,

$$\boxed{\text{Reg}_i(T) = B g_i \cdot \mathbb{E}_\nu[T_i].}$$

This establishes the link between regret and the required online triggering counts.

Step 2. Construct Two Instances.

For arm i , we consider two environments ν and $\nu^{(i)}$:

$$\text{online: } \mu_i^{\text{on}} \mapsto \mu_i^{\text{on}} + g_i, \quad \text{and for all } j \neq i, \mu_j^{\text{on}} \text{ unchanged.}$$

$$\text{offline: } \begin{cases} g_i \leq 2V_i: & \text{the offline means can fully align,} \\ g_i > 2V_i: & \text{match as closely as possible, but cannot fully align.} \end{cases}$$

Since,

$$\mu_i^{\text{off},(i)} \geq \mu_i^{\text{on}} + g_i - V_i, \quad \mu_i^{\text{off}} \leq \mu_i^{\text{on}} + V_i.$$

Hence the difference cannot be smaller than

$$\mu_i^{\text{off},(i)} - \mu_i^{\text{off}} = g_i - 2V_i, \quad \text{i.e., } \mu_i^{\text{off}} = \mu_i^{\text{on}} + V_i, \quad \mu_i^{\text{off},(i)} = \mu_i^{\text{on}} + g_i - V_i.$$

Trigger- i outcomes are Bernoulli(μ_i) in both online and offline data (with different means under the two environments).

KL terms.

$$\text{KL}(P_\nu^{\text{on}} \parallel P_{\nu'}^{\text{on}}) = \mathbb{E}_\nu[T_i] \cdot \text{KL}(\text{Bern}(\mu_i^{\text{on}}) \parallel \text{Bern}(\mu_i^{\text{on}} + g_i)).$$

$$\text{KL}(P_\nu^{\text{off}} \parallel P_{\nu'}^{\text{off}}) = N_i \cdot \text{KL}(\text{Bern}(\mu_i^{\text{off}}) \parallel \text{Bern}(\mu_i^{\text{off},(i)})).$$

By standard bandit lower bounds or the Bretagnolle–Huber inequality, we have

$$\mathbb{E}_\nu[T_i] \cdot \text{KL}(\text{Bern}(\mu_i^{\text{on}}) \parallel \text{Bern}(\mu_i^{\text{on}} + g_i)) + N_i \cdot \text{KL}(\text{Bern}(\mu_i^{\text{off}}) \parallel \text{Bern}(\mu_i^{\text{off},(i)})) \geq \log \frac{1}{2\delta} \geq \log T,$$

Here, we set $\delta = P(\text{error}) + Q(\text{error}) \leq \frac{1}{2T}$, following standard arguments in bandit lower bounds.

Since

$$\text{KL}(\text{Bern}(\mu_i^{\text{on}}) \parallel \text{Bern}(\mu_i^{\text{on}} + g_i)) \leq \frac{g_i^2}{\mu_i^{\text{on}}(1 - \mu_i^{\text{on}})} \quad (\text{from the property of the Bernoulli distribution})$$

we obtain

$$\mathbb{E}_\nu[T_i] \cdot \frac{g_i^2}{\mu_i^{\text{on}}(1 - \mu_i^{\text{on}})} + N_i \cdot \frac{(g_i - 2V_i)^2}{\mu_i^{\text{on}}(1 - \mu_i^{\text{on}})} \geq \log T.$$

Therefore,

$$\mathbb{E}_\nu[T_i] \geq \frac{\mu_i^{\text{on}}(1 - \mu_i^{\text{on}}) \log T}{g_i^2} - N_i \cdot \left(\frac{g_i - 2V_i}{g_i}\right)^2 := \frac{c \log T}{g_i^2} - N_i''',$$

where

$$N_i''' = N_i \max\left\{1 - \frac{2V_i}{g_i}, 0\right\}^2, c \text{ is a universal constant.}$$

Step3. Summary

Finally, we have:

$$\begin{aligned} \text{Reg}_i(T) &\geq B \cdot g_i \cdot \mathbb{E}_\nu[T_i] \\ &\geq B \left(\frac{c \log T}{g_i} - N_i''' \cdot g_i \right) \\ &\geq B \left(\frac{c \log T}{g_i} - \sqrt{c N_i''' \cdot \log T} \right), \quad \left(N_i''' < \frac{c \log T}{g_i^2} \Rightarrow g_i < \sqrt{\frac{c \log T}{N_i'''}} \right). \end{aligned}$$

Notice that $\Delta_{\min}^i = \inf_{S \in \mathcal{S}: p_i^{D,S} > 0, \Delta_S > 0} \Delta_S = \Delta_{S_i} = \frac{B g_i}{K}$, and $\omega_i = 2V_i$, then we have:

$$\begin{aligned} \text{Reg}(T) &= \sum_{i \in [m]} \text{Reg}_i(T) \geq B \sum_{i \in [m]} \left(\frac{c \log T}{g_i} - \sqrt{c N_i''' \cdot \log T} \right) \\ &= B \sum_{i \in [m]} \left(\frac{c B \log T}{K \Delta_{\min}^i} - \sqrt{c N_i''' \cdot \log T} \right), \end{aligned}$$

$$\text{where } N_i''' = N_i \max\left\{1 - \frac{B \omega_i}{K \Delta_{\min}^i}, 0\right\}^2, c = \min\{x \in \{\mu_1, \dots, \mu_m\} : x(1-x)\}$$

F Experimental Details and Real-World Validation

We compare our proposed hybrid CUCB with existing CUCB for the pure online setting (Wang & Chen, 2017) and CLCB for the pure offline setting (Liu et al., 2025). We mainly focus on the task of online learning to rank for the considered CMAB-T problem, where the agent selects k from m base arms. The outcome distribution of each base arm is Bernoulli. We set $m = 10$ and $k = 5$. All results are averaged over 20 runs, and the error bar is defined as the standard deviation divided by $\sqrt{20}$. The triggering process and reward function are introduced as below following existing literature (Chen et al., 2016; Liu et al., 2025):

- Triggering process: The super arm S_t is a permutation over k arms. The environment would check the Bernoulli outcome from the first to the last one. If the first arm has outcome 1, then the triggering stops. Otherwise, the environment would check the second arm. Similar process continues until one arm has the outcome 1. All arms ranked before this arm are observed with outcome 0 and this arm is observed with outcome 1. The following arms have no observations.

- Reward function: The reward function is defined as:

$$r(S_t, \mu) = 1 - \prod_{i \in S_t} (1 - \mu_i).$$

We evaluate the performances of algorithms in both unbiased ($V = 0$) and biased ($V \neq 0$) environments. For the unbiased case, we generate μ_i^{on} uniformly in the interval $(0, 0.5)$ and set $\mu_i^{\text{off}} = \mu_i^{\text{on}}$. For the biased setting, we test different values of discrepancy $V \in \{0.2, 0.3, 0.4\}$. To ensure both μ_i^{on} and μ_i^{off} fall into interval $[0, 1]$ when evaluating different values of V_i , we generate μ_i^{on} uniformly in the interval $(0.4, 0.5)$ and uniformly choose $V_i = V$ or $V_i = -V$. We set $\mu_i^{\text{off}} = \mu_i^{\text{on}} + V_i$.

Finally, we validate the performance of our hybrid CUCB algorithm on a real-world dataset. Specifically, we use the MovieLens dataset, where we randomly select 10 movies as the arms and split the data into two disjoint parts to represent online and offline feedback. The bias level V is computed as the mean difference between the two parts. As shown in Figure 3, our algorithm consistently outperforms or matches the baselines across different offline dataset sizes. Notably, hybrid CUCB achieves significantly lower regret compared to CLCB, while maintaining performance comparable to CUCB even under distributional shift. These results highlight the robustness of hybrid CUCB in practical settings, demonstrating that the algorithm can effectively leverage real-world offline data despite inherent biases and variability.

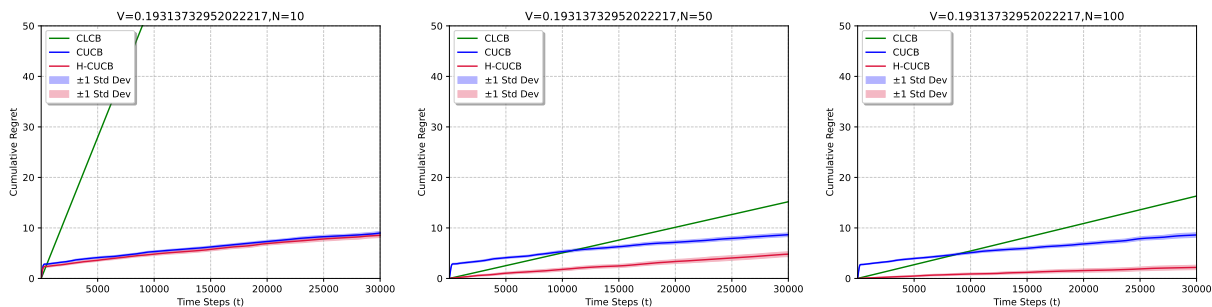


Figure 3: Performance comparison of hybrid CUCB against baselines in a real-world dataset.