

---

# Online Competitive Influence Maximization

---

**Jinhang Zuo**                      **Xutong Liu**                      **Carlee Joe-Wong**  
Carnegie Mellon University    The Chinese University of Hong Kong    Carnegie Mellon University

**John C.S. Lui**  
The Chinese University of Hong Kong

**Wei Chen**  
Microsoft Research

## Abstract

Online influence maximization has attracted much attention as a way to maximize influence spread through a social network while learning the values of unknown network parameters. Most previous works focus on single-item diffusion. In this paper, we introduce a new Online Competitive Influence Maximization (OCIM) problem, where two competing items (e.g., products, news stories) propagate in the same network and influence probabilities on edges are unknown. We adopt a combinatorial multi-armed bandit (CMAB) framework for OCIM, but unlike the non-competitive setting, the important monotonicity property (influence spread increases when influence probabilities on edges increase) no longer holds due to the competitive nature of propagation, which brings a significant new challenge to the problem. We provide a nontrivial proof showing that the Triggering Probability Modulated (TPM) condition for CMAB still holds in OCIM, which is instrumental for our proposed algorithms OCIM-TS and OCIM-OFU to achieve sublinear Bayesian and frequentist regret, respectively. We also design an OCIM-ETC algorithm that requires less feedback and easier offline computation, at the expense of a worse frequentist regret bound. Experimental evaluations demonstrate the effectiveness of our algorithms.

## 1 Introduction

Influence maximization, motivated by viral marketing applications, has been extensively studied since Kempe et al. (2003) formally defined it as a stochastic optimization problem: given a social network  $G$  and a budget  $k$ , how should a set of  $k$  seed nodes in  $G$  be chosen such that the expected number of final activated nodes under a given diffusion model is maximized? They proposed the well-known Independent Cascade (IC) and Linear Threshold (LT) diffusion models, and gave a greedy algorithm that outputs a  $(1 - 1/e - \epsilon)$ -approximate solution for any  $\epsilon > 0$ . However, they only considered a single item (e.g., product, idea) propagating in the network. In reality, different items could propagate concurrently in the same network, interfering with each other and leading to competition during propagation. Several competitive diffusion models (Carnes et al., 2007; Bharathi et al., 2007; Budak et al., 2011; He et al., 2012; Ivanov et al., 2017) have been proposed for this setting. We use a Competitive Independent Cascade (CIC) model (Chen et al., 2013), which extends the classical IC model to multi-item influence diffusion. We consider the competitive influence maximization problem between two items from the “follower’s perspective”: given the seed nodes of the competitor’s item, the follower’s item chooses a set of nodes so as to maximize the expected number of nodes activated by the follower’s item, referred to as the influence spread of the item.

We refer to the above problem as “offline” competitive influence maximization, since the influence probabilities on edges, i.e., the probabilities of an item’s propagation along edges, are known in advance. It can be solved by a greedy algorithm due to submodularity (Chen et al., 2013). However, in many real-world applications, the influence probabilities on edges are unknown. We study the competitive influence maximization in this setting, and call it the Online Competitive Influence Maximization (OCIM) problem. In

Table 1: Summary of the proposed algorithms.

Algorithm	No Prior?	Offline computation	Feedback	Regret
<i>OCIM-TS</i>	×	Standard	Full propagation	Bayes. $O(\sqrt{T \ln T})$
<i>OCIM-OFU</i>	✓	Hard	Full propagation	Freq. $O(\sqrt{T \ln T})$
<i>OCIM-ETC</i>	✓	Standard	Direct out-edges	Freq. $O(T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}})$

OCIM, the influence probabilities on edges need to be learned through repeated influence maximization trials: in each round, given the seed nodes of the competitor, we (i) choose  $k$  seed nodes; (ii) observe the resulting diffusion that follows the CIC model to update our knowledge of the edge probabilities; and (iii) obtain a reward, which is the total number of nodes activated by our item. Our goal is to choose the seed nodes in each round based on previous observations so as to maximize the cumulative reward.

Most previous studies on the online non-competitive influence maximization problem use a combinatorial multi-armed bandit (CMAB) framework (Chen et al., 2016; Wen et al., 2017), an extension of the classical multi-armed bandit problem that captures the tradeoff between exploration and exploitation in sequential decision making. In CMAB, a player chooses a combinatorial action to play in each round, observes a set of arms triggered by this action and receives a reward. The player aims to maximize her cumulative reward over multiple rounds, navigating a tradeoff between exploring unknown actions/arms and exploiting the best known action. CMAB algorithms must also deal with an exponential number of possible combinatorial actions, which makes exploring all actions infeasible.

**Our Contributions.** To the best of our knowledge, we are the first to study the online competitive influence maximization problem. We introduce a general contextual combinatorial multi-armed bandit framework with probabilistically triggered arms (C<sup>2</sup>MAB-T) for OCIM. Within this framework, OCIM presents a new challenge: the key monotonicity property (influence spread increases when influence probabilities on edges increase) no longer holds due to the competitive nature of propagation, and thus upper confidence bound (UCB) based algorithms (Chen et al., 2016; Wen et al., 2017) cannot be directly applied to OCIM. Such non-monotonicity also complicates the analysis of the important Triggering Probability Modulated (TPM) condition for CMAB (Wang and Chen, 2017), and we provide a non-trivial new proof to show it still holds for OCIM. We are the first to identify the OCIM problem as a natural CMAB problem without monotonicity and tackle it from three directions, providing three solutions with different tradeoffs, as shown

in Table 1: OCIM-TS uses standard offline oracles to achieve good Bayesian regret, but requires prior knowledge of edge probabilities; OCIM-OFU has a stronger frequentist regret bound without prior knowledge, but requires harder offline computation; and OCIM-ETC uses standard offline oracles and fewer observations, but leads to a worse frequentist regret bound. None is a perfect solution for OCIM, but we believe their tradeoffs shed light on the challenges involved in solving OCIM and even general CMAB problems without monotonicity. Our regret analysis of OCIM-TS delicately combines the key property of Thompson Sampling (TS) with the TPM condition to tackle non-monotonicity and allows any benchmark (exact, approximate, or even heuristic) oracle; our analysis of OCIM-OFU and OCIM-ETC extends the analysis for CMAB to a new contextual setting (C<sup>2</sup>MAB-T) where the contexts are defined as the feasible sets of super arms and are not bonded with base arms. We also discuss the extension of our framework to settings with more complex competitor actions. Experiments on two real-world datasets demonstrate the effectiveness of our proposed algorithms. Due to the space constraint, we discuss important insights of our proofs and move the complete proofs as well as the results for the general C<sup>2</sup>MAB-T problem to the Appendix.

**Related Work.** Kempe et al. (2003) formally defined the influence maximization problem in their seminal work. Since then, the problem has been extensively studied (Li et al., 2018). Borgs et al. (2014) presented a breakthrough approximation algorithm that runs in near-linear time, which was improved by a series of algorithms (Tang et al., 2015; Nguyen et al., 2016; Tang et al., 2018). A number of studies (Carnes et al., 2007; Bharathi et al., 2007; Budak et al., 2011; He et al., 2012; Lin and Lui, 2015; Ivanov et al., 2017) addressed competitive influence maximization problems where multiple competing sources propagate in the same network. Carnes et al. (2007) proposed the distance-based and wave propagation models, and considered the influence maximization problem from the follower’s perspective. Bharathi et al. (2007) considered the CIC model and gave an algorithm for computing the best response to an opponent’s strategy.

When the influence probabilities of edges are un-

known, the non-competitive online influence maximization problem has been extensively studied (Chen et al., 2016; Wang and Chen, 2017; Wen et al., 2017; Wu et al., 2019; Vaswani et al., 2017; Perrault et al., 2020). Chen et al. (2016) studied the problem under the IC model and proposed a general CMAB framework. We introduce a new contextual extension of CMAB, called C<sup>2</sup>MAB-T, different from the contextual CMAB studied by Chen et al. (2018) and Qin et al. (2014): they consider the context features of all base arms and assume the action space of super arms is a subset of all base arms, while we consider the feasible set of super arms as the context, which is more flexible than a subset of all base arms. Wang and Chen (2017) introduced a triggering probability modulated (TPM) bounded smoothness condition to remove an undesired factor in the regret bound of Chen et al. (2016). Perrault et al. (2020) introduced a budgeted online influence maximization framework, where marketers optimize their seed sets under a budget rather than a cardinality constraint. Our OCIM-TS algorithm is similar to the Combinatorial Thompson Sampling (CTS) algorithm of Wang and Chen (2018). However, CTS requires an exact oracle and has frequentist regret bound, while OCIM-TS allows any benchmark oracle and has Bayesian regret bound. Hüyük and Tekin (2020) studied the Bayesian regret of CTS for CMAB, but they also require an exact oracle and a monotonicity assumption that does not hold for OCIM. Our Bayesian regret analysis is also different from that of Russo and Van Roy (2016): they only study a simple special CMAB problem, while we provide the regret bound for general C<sup>2</sup>MAB-T instances, including the OCIM problem.

## 2 OCIM Formulation

In this section we present the formulation of OCIM. We first introduce the traditional competitive influence maximization problem, and then discuss its online extension where edge probabilities are unknown.

### 2.1 Competitive Independent Cascade Model

We consider a Competitive Independent Cascade (CIC) model, which is an extension of the classical IC model to multi-item influence diffusion. A network is modeled as a directed graph  $G = (V, E)$  with  $n = |V|$  nodes and  $m = |E|$  edges. Every edge  $(u, v) \in E$  is associated with a probability  $p(u, v)$ . There are two items,  $A$  and  $B$ , trying to propagate in  $G$  from their own seed sets  $S_A$  and  $S_B$ . The influence propagation runs as follows: nodes in  $S_A$  (resp.  $S_B$ ) are activated by  $A$  (resp.  $B$ ) at step 0; at each step  $s \geq 1$ , a node  $u$  activated by  $A$  (resp.  $B$ ) in step  $s - 1$  tries to acti-

vate each of its inactive out-neighbors  $v$  to be  $A$  (resp.  $B$ ) with an independent probability  $p(u, v)$  that is the same for  $A$  and  $B$  (i.e., we consider a homogeneous CIC model). The homogeneity assumption is reasonable since typically  $A$  and  $B$  are two items of the same category (thus competing), so they are likely to have similar propagation characteristics.

If two in-neighbors of  $v$  activated by  $A$  and  $B$  respectively both successfully activate  $v$  at step  $s$ , then a tie-breaking rule is applied at  $v$  to determine the final adoption. In this paper, we consider two types of tie-breaking rules: dominance (Budak et al., 2011) and proportional (Chen et al., 2011) tie-breaking rules. Dominance tie-breaking with  $A > B$  (resp.  $B > A$ ) means  $v$  will always adopt  $A$  (resp.  $B$ ) in a competition. Proportional tie-breaking means that if there are  $n_A$  in-neighbors activated by  $A$  and  $n_B$  in-neighbors activated by  $B$  trying to activate  $v$  at the same step, the probability that  $v$  adopts  $A$  (resp.  $B$ ) is  $\frac{n_A}{n_A+n_B}$  (resp.  $\frac{n_B}{n_A+n_B}$ ). The same tie-breaking rule also applies to the case when a node  $u$  is selected both as an  $A$ -seed and a  $B$ -seed. The process stops when no nodes activated at a step  $s$  have inactive out-neighbors.

We consider the follower’s perspective in the optimization task: let  $A$  be the follower and  $B$  be the competitor. Then given  $S_B$ , our goal is to choose at most  $k$  seed nodes in  $G$  as  $S_A$  to maximize the influence spread of  $A$ , denoted as  $\sigma_A(S_A, S_B)$ , which is the expected number of nodes activated by  $A$  after the propagation ends. According to Budak et al. (2011)’s result, the above optimization task under the homogeneous CIC model with the dominance tie-breaking rule has the monotone and submodular properties, and thus can be approximately solved by a greedy algorithm.

### 2.2 OCIM Model

In the online competitive influence maximization (OCIM) problem, the edge probabilities  $p(u, v)$ ’s are unknown and need to be learned: in each round  $t$ , given  $S_B^{(t)}$ , we can choose up to  $k$  seed nodes as  $S_A^{(t)}$ , observe the whole propagation of  $A$  and  $B$  that follows the CIC model, and obtain the reward, which is the number of nodes finally activated by  $A$  in this round. The propagation feedback observed is then used to update the estimates on edge probabilities  $p(u, v)$ ’s, so that we can achieve better influence maximization results in subsequent rounds. Our goal is to accumulate as much reward as possible through this repeated process over multiple rounds.

We introduce a new contextual combinatorial multi-armed bandit framework with probabilistically triggered arms (C<sup>2</sup>MAB-T) for the OCIM problem, which is a contextual extension of CMAB-T from Wang and

Chen (2017). In OCIM, the set of edges  $E$  is the set of (base) arms  $[m] = \{1, \dots, m\}$ , and their outcomes follow  $m$  independent Bernoulli distributions with expectation  $\mu_e = p(u, v)$  for all  $e = (u, v) \in E$ . We denote the independent samples of arms in round  $t$  as  $X^{(t)} = (X_1^{(t)}, \dots, X_m^{(t)}) \in \{0, 1\}^m$ , where  $X_i^{(t)} = 1$  means the  $i$ -th edge is on (or live) and  $X_i^{(t)} = 0$  means the  $i$ -th edge is off (or blocked) in round  $t$ , and thus  $X^{(t)}$  corresponds to the *live-edge graph* (Kempe et al., 2003) in round  $t$ . We consider the seed set of the competitor,  $S_B^{(t)}$ , as the *context* in round  $t$  since it is determined by the competitor and can affect our choice of  $S_A^{(t)}$ . We define  $\mathcal{S}^{(t)} = \left\{ S \mid S = (S_A^{(t)}, S_B^{(t)}), |S_A^{(t)}| \leq k \right\}$  as the action space in round  $t$  and  $S^{(t)} \in \mathcal{S}^{(t)}$  as the real action. We define the triggered arm set  $\tau_t$  as the set of edges reached by the propagation from either  $S_A^{(t)}$  or  $S_B^{(t)}$ . Thus,  $\tau_t$  is the set of edges  $(u, v)$  where  $u$  can be reached from  $S^{(t)}$  by passing through only edges  $e \in E$  with  $X_e^{(t)} = 1$ . The outcomes of  $X_i^{(t)}$  for all  $i \in \tau_t$  are observed as the feedback. We denote the obtained reward in round  $t$  as  $R(S^{(t)}, X^{(t)})$ , which is the number of nodes finally activated by  $A$ . The expected reward  $r_{S^{(t)}}(\boldsymbol{\mu}) = \mathbb{E}[R(S^{(t)}, X^{(t)})]$  is a function of the action  $S^{(t)}$  and the vector  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$ . Note that our framework can also handle dynamic tie-breaking rules over different rounds, by treating the tie-breaking rule as a part of the context. For ease of explanation, we assume a fixed tie-breaking rule in this paper.

The performance of a learning algorithm  $\mathcal{A}$  is measured by its expected regret, which is the difference in expected cumulative reward between always playing the best action and playing actions selected by algorithm  $\mathcal{A}$ . Let  $\text{opt}^{(t)}(\boldsymbol{\mu}) = \sup_{S \in \mathcal{S}^{(t)}} r_{S^{(t)}}(\boldsymbol{\mu})$  denote the expected reward of the optimal action in round  $t$ . Since the offline influence maximization under the CIC model is NP-hard (Budak et al., 2011), we assume that there exists an offline  $(\alpha, \beta)$ -approximation oracle  $\mathcal{O}$ , which takes  $S_B^{(t)}$  and  $\boldsymbol{\mu}$  as inputs and outputs an action  $S^{\mathcal{O},(t)}$  such that  $\Pr\{r_{S^{\mathcal{O},(t)}}(\boldsymbol{\mu}) \geq \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu})\} \geq \beta$ , where  $\alpha$  is the approximation ratio and  $\beta$  is the success probability. Instead of comparing with the exact optimal reward, we use the following  $(\alpha, \beta)$ -approximation *frequentist regret* for  $T$  rounds:

$$Reg_{\alpha, \beta}^{\mathcal{A}}(T; \boldsymbol{\mu}) = \sum_{t=1}^T \alpha \cdot \beta \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - \sum_{t=1}^T r_{S^{\mathcal{A},(t)}}(\boldsymbol{\mu}), \quad (1)$$

where  $S^{\mathcal{A},(t)} := (S_A^{\mathcal{A},(t)}, S_B^{(t)})$  is the action chosen by algorithm  $\mathcal{A}$  in round  $t$ . Here  $S_B^{(t)}$  is the context and  $S_A^{\mathcal{A},(t)}$  is the seed set of item  $A$  chosen by algorithm  $\mathcal{A}$ .

Another way to measure the performance of the algorithm  $\mathcal{A}$  is using *Bayesian regret* (Russo and Van Roy, 2014). Denote the prior distribution of  $\boldsymbol{\mu}$  as  $\mathcal{Q}$  (we

will discuss how to derive  $\mathcal{Q}$  for OCIM in Section 4). When the prior  $\mathcal{Q}$  is given, the corresponding Bayesian regret is defined as:

$$BayesReg_{\alpha, \beta}^{\mathcal{A}}(T) = \mathbb{E}_{\boldsymbol{\mu} \sim \mathcal{Q}} Reg_{\alpha, \beta}^{\mathcal{A}}(T; \boldsymbol{\mu}). \quad (2)$$

We will design algorithms to solve the OCIM problem and bound their achieved Bayesian and frequentist regrets in Section 4 and Section 5, respectively. We also discuss the general C<sup>2</sup>MAB-T problem and its solutions in the Appendix.

### 3 Properties of OCIM

In this section, we first show that the key monotonicity property for CMAB does not hold in OCIM. We then prove that the important Triggering Probability Modulated (TPM) condition still holds, which is essential for the analysis of all proposed algorithms.

#### 3.1 Non-monotonicity

The monotonicity condition given by Wang and Chen (2017) could be stated as follows in the context of OCIM: for any action  $S = (S_A, S_B)$ , for any two expectation vectors  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\mu}' = (\mu'_1, \dots, \mu'_m)$ , we have  $r_S(\boldsymbol{\mu}) \leq r_S(\boldsymbol{\mu}')$  if  $\mu_i \leq \mu'_i$  for all  $i \in [m]$ . Figure 1 shows a simple example of OCIM that does not satisfy the monotonicity condition. The left and right nodes are the seed nodes of  $A$  and  $B$ ; the numbers below edges are influence probabilities. It is easy to calculate that  $r_S(\boldsymbol{\mu}) = \mu_1(1 - \mu_2) + 2$ , for both dominance and proportional tie-breaking rules. Thus, if we increase  $\mu_2$ ,  $r_S(\boldsymbol{\mu})$  will decrease, which is contrary to monotonicity. In general, for every edge  $(u, v)$ , depending on the positions of the  $A$ - and  $B$ -seeds, increasing the influence probability of  $(u, v)$  may benefit the propagation of  $A$  or may benefit the propagation of  $B$  and thus impair the propagation of  $A$ . Thus, the influence spread of  $A$  has intricate connections with the influence probabilities on the edges.

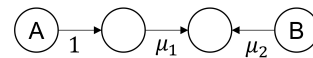


Figure 1: Example of non-monotonicity in OCIM

The lack of monotonicity poses a significant challenge to the OCIM problem. We cannot directly use UCB-type algorithms (Chen et al., 2016), as they will not provide optimistic solutions to bound the regret.

### 3.2 Triggering Probability Modulated (TPM) Bounded Smoothness

The lack of monotonicity further complicates the analysis of the Triggering Probability Modulated (TPM) condition (Wang and Chen, 2017), which is crucial in establishing regret bounds for CMAB algorithms. We use  $p_i^S(\boldsymbol{\mu})$  to denote the probability that the action  $S$  triggers arm  $i$  when the expectation vector is  $\boldsymbol{\mu}$ . The TPM condition in OCIM is given below.

**Condition 1.** (*1-Norm TPM bounded smoothness*). We say that an OCIM problem instance satisfies 1-norm TPM bounded smoothness, if there exists  $C \in \mathbb{R}^+$  (referred to as the bounded smoothness coefficient) such that, for any two expectation vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ , and any action  $S = (S_A, S_B)$ , we have  $|r_S(\boldsymbol{\mu}) - r_S(\boldsymbol{\mu}')| \leq C \sum_{i \in [m]} p_i^S(\boldsymbol{\mu}) |\mu_i - \mu'_i|$ .

Fortunately, with a more intricate analysis, we are able to show the following TPM condition.

**Theorem 3.1.** *Under both dominance and proportional tie-breaking rules, OCIM instances satisfy the 1-norm TPM bounded smoothness condition with coefficient  $C = \tilde{C}$ , where  $\tilde{C}$  is the maximum number of nodes that any one node can reach in graph  $G$ .*

The proof of the above theorem is one of the key technical contributions of the paper. In the non-competitive setting, an edge coupling method could give a relatively simple proof for the TPM condition.<sup>1</sup> The idea of edge coupling is that for every edge  $e \in E$ , we sample a real number  $X_e \in [0, 1]$  uniformly at random, and determine  $e$  to be live under  $\boldsymbol{\mu}$  if  $X_e \leq \mu_e$  and blocked if  $X_e > \mu_e$ , and similarly for  $\boldsymbol{\mu}'$ . This couples the live-edge graphs  $L$  and  $L'$  under  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$  respectively. In the non-competitive setting, due to the monotonicity property, we only need to consider the TPM condition when  $\boldsymbol{\mu} \geq \boldsymbol{\mu}'$  (coordinate-wise), and this implies that  $L'$  is a subgraph of  $L$ , which significantly simplifies the analysis. However, in the competitive setting, monotonicity does not hold, and we have to show the TPM condition for every pair of  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ . Thus,  $L$  and  $L'$  no longer have the subgraph relationship. In this case, we have to show that for every coupling  $L$  and  $L'$ , for every  $v \in V$  that is activated by  $A$  in  $L$  but not activated by  $A$  in  $L'$ , it is because either (a) some edge  $e = (u, w)$  is live in  $L$  but blocked in  $L'$  while  $u$  is  $A$ -activated (or equivalently  $e$  is  $A$ -triggered); or (b) some edge  $e$  is live in  $L'$  but blocked in  $L$  while  $e$  is  $B$ -triggered. The case (b) is due to the possibility of  $B$  blocking  $A$ 's propagation, a unique scenario in OCIM. The above claim

<sup>1</sup>The original proof in (Wang and Chen, 2017) occupies several pages, but Li et al. (2020) (in their Appendix E) provide a much shorter proof based on edge coupling.

---

#### Algorithm 1 OCIM-TS with offline oracle $\mathcal{O}$

---

- 1: **Input:**  $m, \mathcal{O}$ , Prior  $\mathcal{Q} = \prod_{i \in [m]} \text{Beta}(a_i, b_i)$ .
  - 2: **for**  $t = 1, 2, 3, \dots$  **do**
  - 3:   For each arm  $i \in [m]$ , draw a sample  $\mu_i^{(t)}$  from  $\text{Beta}(a_i, b_i)$ ; let  $\boldsymbol{\mu}^{(t)} = (\mu_1^{(t)}, \dots, \mu_m^{(t)})$ .
  - 4:   Obtain context  $S_B^{(t)}$ .
  - 5:    $S^{(t)} \leftarrow \mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}^{(t)})$ .
  - 6:   Play action  $S^{(t)}$ , which triggers a set  $\tau \subseteq [m]$  of base arms with feedback  $X_i^{(t)}$ 's,  $i \in \tau$ .
  - 7:   **for all**  $i \in \tau$  **do**
  - 8:      $a_i \leftarrow a_i + X_i^{(t)}$ ;  $b_i \leftarrow b_i + 1 - X_i^{(t)}$ .
  - 9:   **end for**
  - 10: **end for**
- 

needs nontrivial inductive proofs for dominance and proportional tie-breaking rules, and then its correctness ensures the TPM condition.

## 4 Bayesian Regret Approach

In our OCIM model, since the samples of base arms follow Bernoulli distributions with mean vector  $\boldsymbol{\mu}$ , we can assume the prior distributions of  $\boldsymbol{\mu}$ ,  $\mathcal{Q}$ , are Beta distributions, where  $\mu_i \sim \text{Beta}(a_i, b_i)$  for all arm  $i$ . Given the prior distributions of all arms, we propose an Online Competitive Influence Maximization-Thompson Sampling (OCIM-TS) algorithm, which is described in Algorithm 1. We initialize the prior distribution of each arm  $i$  to  $\text{Beta}(a_i, b_i)$ . Then we take the context  $S_B^{(t)}$  and the sampled  $\boldsymbol{\mu}^{(t)}$  from prior distributions as inputs to the oracle  $\mathcal{O}$ , and get an output action  $S^{(t)}$ . After taking this action, we get feedback  $X_i^{(t)}$ 's from all triggered arms  $i \in \tau$ , then use them to update the prior distributions of all triggered base arms in  $\tau$ . Let  $\tilde{S} = \{i \in [m] \mid p_i^S(\boldsymbol{\mu}) > 0\}$  be the set of arms that can be triggered by  $S$ . We define  $K = \max_{S \in \mathcal{S}^{(t)}} |\tilde{S}|$  as the largest number of arms that could be triggered by a feasible action. We provide the Bayesian regret bound of OCIM-TS.

**Theorem 4.1.** *The OCIM-TS algorithm has the following Bayesian regret bound with  $\tilde{C}$  as defined in Theorem 3.1:*

$$\text{BayesReg}_{\alpha, \beta}(T) \leq O(\tilde{C} \sqrt{mKT \ln T}). \quad (3)$$

This regret bound essentially matches the distribution-independent frequentist regret bound of OCIM-OFU in the next section. The proof of the above theorem is inspired by the posterior sampling regret decomposition of Russo and Van Roy (2014). However, we combine the key property of posterior sampling with the TPM condition in Theorem 3.1 to tackle non-monotonicity. OCIM-TS can also be applied to gen-

eral C<sup>2</sup>MAB-T problems and allows any benchmark offline oracles (e.g., approximate or heuristic oracles). We provide the Bayesian regret bound of OCIM-TS on general C<sup>2</sup>MAB-T problems in the Appendix.

## 5 Frequentist Regret Approach

Although OCIM-TS can solve the OCIM problem with a standard offline oracle (e.g., TCIM in Lin and Lui (2015)), it requires the prior distribution of the network parameter  $\mu$ , which might not be available in practice. In this section, we first propose the OCIM-OFU algorithm. It achieves good frequentist regret without the prior knowledge, but requires a new oracle to solve a harder offline problem. We then design the OCIM-ETC algorithm, which requires less feedback and easier offline computation, but yields a worse frequentist regret bound.

### 5.1 OCIM-OFU Algorithm

As discussed in Section 3.1, due to the lack of monotonicity, we cannot directly use UCB-type algorithms. However, it is still possible to design bandit algorithms following the principle of Optimism in the Face of Uncertainty (OFU). We first introduce a new offline problem that jointly optimizes for both the seed set  $S^*$  and the optimal influence probability vector  $\mu^*$ , where each dimension of  $\mu^*$ ,  $\mu_i^*$ , is searched within a confidence interval  $c_i$ , for all  $i \in E$ .

$$\begin{aligned} & \underset{S, \mu}{\text{maximize}} && r_S(\mu) \\ & \text{subject to} && |S_A| \leq k, S = (S_A, S_B) \\ & && \mu_i \in c_i, i = 1, \dots, m. \end{aligned} \quad (4)$$

We then define a new offline  $(\alpha, \beta)$ -approximation oracle  $\tilde{\mathcal{O}}$  to solve this problem. Oracle  $\tilde{\mathcal{O}}$  takes  $S_B$  and  $c_i$ 's as inputs and outputs  $\mu^{\tilde{\mathcal{O}}}$  and action  $S^{\tilde{\mathcal{O}}} = (S_A^{\tilde{\mathcal{O}}}, S_B)$ , such that  $\Pr\{r_{S^{\tilde{\mathcal{O}}}(\mu^{\tilde{\mathcal{O}}})} \geq \alpha \cdot r_{S^*(\mu^*)}\} \geq \beta$ , where  $(S^*, \mu^*)$  is the optimal solution for Eq.(4).

With the offline oracle  $\tilde{\mathcal{O}}$ , we propose an algorithm following the principle of Optimism in the Face of Uncertainty (OFU), named OCIM-OFU. The algorithm maintains the empirical mean  $\hat{\mu}_i$  and confidence radius  $\rho_i$  for each edge probability. It uses the lower and upper confidence bounds to determine the range of  $\mu_i$ :  $c_i = [(\hat{\mu}_i - \rho_i)^{0+}, (\hat{\mu}_i + \rho_i)^{1-}]$ , where we use  $(x)^{0+}$  and  $(x)^{1-}$  to denote  $\max\{x, 0\}$  and  $\min\{x, 1\}$  for any real number  $x$ . It feeds  $S_B^{\tilde{\mathcal{O}}}$  and all current  $c_i$ 's into the offline oracle  $\tilde{\mathcal{O}}$  to obtain the action  $S^{(t)} = (S_A^{(t)}, S_B^{\tilde{\mathcal{O}}})$  to play at round  $t$ . The confidence radius  $\rho_i$  is large if arm  $i$  is not triggered often, which leads to a wider search space  $c_i$  to find the optimistic estimate of  $\mu_i$ . We provide its frequentist regret bound.

---

### Algorithm 2 OCIM-OFU with offline oracle $\tilde{\mathcal{O}}$

---

- 1: **Input:**  $m$ , Oracle  $\tilde{\mathcal{O}}$ .
  - 2: For each arm  $i \in [m]$ ,  $T_i \leftarrow 0$ . {maintain the total number of times arm  $i$  is played so far.}
  - 3: For each arm  $i \in [m]$ ,  $\hat{\mu}_i \leftarrow 1$ . {maintain the empirical mean of  $X_{i \cdot}$ }
  - 4: **for**  $t = 1, 2, 3, \dots$  **do**
  - 5:   For each arm  $i \in [m]$ ,  $\rho_i \leftarrow \sqrt{\frac{3 \ln t}{2T_i}}$ . {the confidence radius,  $\rho_i = +\infty$  if  $T_i = 0$ .}
  - 6:   For each arm  $i \in [m]$ ,  $c_i \leftarrow [(\hat{\mu}_i - \rho_i)^{0+}, (\hat{\mu}_i + \rho_i)^{1-}]$ . {the estimated range of  $\mu_i$ .}
  - 7:   Obtain context  $S_B^{(t)}$ .
  - 8:    $S^{(t)} \leftarrow \tilde{\mathcal{O}}(S_B^{(t)}, c_1, c_2, \dots, c_m)$ .
  - 9:   Play action  $S^{(t)}$ , which triggers a set  $\tau \subseteq [m]$  of base arms with feedback  $X_i^{(t)}$ 's,  $i \in \tau$ .
  - 10:   For every  $i \in \tau$  update  $T_i$  and  $\hat{\mu}_i$ :  $T_i = T_i + 1$ ,  $\hat{\mu}_i = \hat{\mu}_i + (X_i^{(t)} - \hat{\mu}_i)/T_i$ .
  - 11: **end for**
- 

**Theorem 5.1.** *The OCIM-OFU algorithm has the following distribution-independent bound (see the Appendix for the distribution-dependent bound) with  $\tilde{\mathcal{C}}$  defined in Theorem 3.1,*

$$\text{Reg}_{\alpha, \beta}(T; \mu) \leq O(\tilde{\mathcal{C}} \sqrt{mKT \ln T})$$

The above regret bound has the typical form of  $\sqrt{T \ln T}$ , indicating that it is tight on the important time horizon  $T$ . In fact, it has the same order as in Wang and Chen (2017)'s for the CMAB problem under monotonicity, despite the fact that the OCIM problem does not enjoy monotonicity, and matches the lower bound of CMAB with general reward functions in (Merlis and Mannor, 2020). This result is due to our non-trivial TPM condition analysis (Theorem 3.1) that shows the same condition as in Wang and Chen (2017)'s setting with monotonicity.

**Computational Efficiency.** We now discuss the computational complexity of implementing the OCIM-OFU algorithm. We show the complexity of the new offline optimization problem in Eq. (4).

**Theorem 5.2.** *The offline problem in Eq.(4) is #P-hard.*

As mentioned before, the original offline problem, i.e., maximizing  $r_S(\mu)$  over  $S$  when fixing  $\mu$ , can be solved by several algorithms (Lin and Lui, 2015) based on submodularity of  $r_S(\mu)$  over  $S$ . A straightforward attempt on the new offline problem in Eq.(4) is to show the submodularity of  $g(S) = \max_{\mu} r_S(\mu)$  over  $S$ , and then to use a greedy algorithm on  $g$  to select  $S$ . Unfortunately, we find that  $g(S)$  is not submodular (see

the Appendix for a counterexample). Implementing the oracle  $\tilde{\mathcal{O}}$  is then a challenge. However, it is possible to design efficient approximate oracles for bipartite graphs, which model the competitive probabilistic maximum coverage problem with applications in online advertising (Chen et al., 2016). The main idea is that we can pre-determine that either the lower or the upper bound of  $c_i$  is optimal and should be chosen as  $\mu_i^*$  depending on the tie-breaking rule, then use existing efficient influence maximization algorithms to get approximate solutions. The competitive propagation in the general graph is much more complicated, but we have a key observation that the optimal solution for the optimization problem in Eq.(4) must occur at the boundaries of the intervals  $c_i$ . Based on that, we discuss solutions for some specific graphs such as trees. See the Appendix for more details.

## 5.2 OCIM-ETC Algorithm

In this section, we propose an OCIM Explore-Then-Commit (OCIM-ETC) algorithm. It has two advantages: first, it does not need the new offline oracle discussed in Sec. 5.1; and second, it requires fewer observations than our other algorithms: instead of the observations of all triggered edges, i.e.,  $\tau$ , it only needs the observations of all direct out-edges of seed nodes.

Like other ETC algorithms (Garivier et al., 2016), OCIM-ETC divides the  $T$  rounds into two phases: an exploration phase and an exploitation phase. In the exploration phase, it chooses each node as the seed node of  $A$  for  $N$  times. The exploration phase thus takes  $\lceil nN/k \rceil$  rounds. In the exploitation phase, it takes  $S_B^{(t)}$  and the empirical means  $\hat{\mu}_i$  as inputs to the oracle  $\mathcal{O}$  mentioned in Sec. 2, then plays the output action  $S^{\mathcal{O},(t)}$ . We give its frequentist regret bound.

**Theorem 5.3.** *The OCIM-ETC algorithm has the following distribution-independent regret bound (see the Appendix for the distribution-dependent bound) with  $\tilde{C}$  defined in Theorem 3.1, when  $N = (\tilde{C}mk)^{\frac{2}{3}}n^{-\frac{4}{3}}T^{\frac{2}{3}}(\ln T)^{\frac{1}{3}}$ ,*

$$\text{Reg}_{\alpha,\beta}(T; \boldsymbol{\mu}) \leq O((\tilde{C}mn)^{\frac{2}{3}}k^{-\frac{1}{3}}T^{\frac{2}{3}}(\ln T)^{\frac{1}{3}}). \quad (5)$$

Although this regret bound is worse than that of the OCIM-OFU algorithm in Theorem 5.1, OCIM-ETC requires easier offline computation and less feedback since it only needs to observe the results of direct out-edges of seed nodes, which shows the tradeoff between regret bound and feedback/computation in OCIM.

## 6 Extension to Probabilistic Seed Distribution for the Competitor

Lin and Lui (2015) extend the offline CIM problem to

a probabilistic setting where the competitor’s seed distribution is known (i.e., the probability of each node being selected as a seed by the competitor). In this section, we extend our algorithms to handle two new settings where the competitor has a probabilistic seed distribution. Note that we need to slightly modify the TPM condition for these settings. We denote the expected reward of follower  $A$  as  $r(S_A, D_B, \boldsymbol{\mu})$ , where  $S_A$  is the seed set of  $A$ ,  $D_B$  is the seed distribution of  $B$ . We use  $p_i(S_A, D_B, \boldsymbol{\mu})$  to denote the probability that either  $S_A$  or  $S_B$  will trigger arm  $i$  when the seed set of  $A$  is  $S_A$ , the seed set of  $B$ ,  $S_B$ , is sampled from  $D_B$ , and the expectation vector is  $\boldsymbol{\mu}$ . The modified TPM condition is given below.

**Condition 2.** *(Modified TPM bounded smoothness). We say that an OCIM problem instance satisfies modified TPM bounded smoothness, if there exists  $C \in \mathbb{R}^+$  such that, for any two expectation vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ , and any seed set  $S_A$  and seed distribution  $D_B$ , we have  $|r(S_A, D_B, \boldsymbol{\mu}) - r(S_A, D_B, \boldsymbol{\mu}')| \leq C \sum_{i \in [m]} p_i(S_A, D_B, \boldsymbol{\mu}) |\mu_i - \mu'_i|$ .*

With the similar analysis of Theorem 3.1, we can show the following TPM condition when the competitor has probabilistic seed distribution.

**Theorem 6.1.** *Under both dominance and proportional tie-breaking rules, OCIM instances satisfy the modified TPM bounded smoothness condition with coefficient  $C = 2\tilde{C}$ , where  $\tilde{C}$  is the maximum number of nodes that any one node can reach in graph  $G$ .*

**Known dynamic seed distribution.** In round  $t$ , the competitor’s seed set  $S_B^{(t)}$  follows a distribution  $D_B^{(t)}$ , i.e.,  $S_B^{(t)} \sim D_B^{(t)}$ . However, the follower only knows  $D_B^{(t)}$  but not  $S_B^{(t)}$  before choosing  $S_A^{(t)}$ . Since our proposed framework has a nice separation between online learning and offline computation, in this setting, only the offline computation part will be affected. Specifically, we can replace the oracle  $\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}^{(t)})$  in OCIM-TS and OCIM-ETC with a new oracle  $\mathcal{O}_{\text{new}}(D_B^{(t)}, \boldsymbol{\mu}^{(t)})$ . For OCIM-OFU, similar to oracle  $\tilde{\mathcal{O}}$ , we need a new oracle  $\tilde{\mathcal{O}}_{\text{new}}$  that takes  $D_B^{(t)}$  and the confidence intervals  $\{c_i\}$  as inputs and outputs  $S_A^{(t)}$ . We can use the TCIM algorithm of (Lin and Lui, 2015) to design  $\mathcal{O}_{\text{new}}$  and  $\tilde{\mathcal{O}}_{\text{new}}$ . Our proposed algorithms will have the same regret bounds as in Theorems 4.1 and 5.1.

**Unknown fixed seed distribution.** In this setting, the seed distribution of the competitor,  $D_B$ , is unknown to the follower but fixed for all rounds. To solve this problem, we introduce a virtual  $B$  seed node  $u_B$ , which connects to each existing node  $u$  with an unknown edge probability  $p(u_B, u)$  equal to the probability of  $u$  being selected as a  $B$  seed. This reduces the

case of probabilistic seed selection to the standard CIC model with a known seed node  $u_B$ . The unknown edge probabilities  $p(u_B, u)$ 's can be learned together with the edge probabilities in the original graph. Therefore, we do not need to know the competitor's seed selection in advance and can learn it over time through the online learning process. Our algorithms will have the same regret guarantees as in Theorems 4.1 and 5.1.

## 7 Experiments

**Datasets and settings.** To validate our theoretical findings, we conduct experiments on two real-world datasets widely used in the influence maximization literature, with detailed statistics summarized in Table 2. First, we use the Yahoo! Search Marketing Advertiser Bidding Data<sup>2</sup> (denoted as Yahoo-Ad), which contains a bipartite graph between 1,000 keywords and 10,475 advertisers. Every entry in the original Yahoo-Ad dataset is a 4-tuple, which represents a "keyword-id" bid by "advertiser-id" at "time-stamp" with "price". We extract advertiser-ids and keyword-ids as nodes, and add an edge if the advertiser bids the keyword at least once. Each edge shows the "who is interested in what" relationship. This dataset will contain 11,475 nodes and 52,567 edges. The motivation of this experiment is to select a set of keywords that is maximally associated to advertisers, which is useful for the publisher to promote keywords to advertisers. We then consider the DM network (Tang et al., 2009) with 679 nodes representing researchers and 3,374 edges representing collaborations between them. We simulate a researcher asking others (i.e.,  $S_A$ ) to spread her ideas while her competitor (i.e.,  $S_B$ ) promotes a competing proposal. We set the parameters of our experiments as the following. For the edge weights, Yahoo-Ad uses the weighted cascade method (Kempe et al., 2003), i.e.  $p(s, t) = 1/deg_-(s)$ , where  $deg_-(s)$  is the in-degree of node  $s$ , and weights for DM are obtained by the learned edge parameters from (Tang et al., 2009). For Bayesian regrets, we set a prior distribution of  $\mu_e \sim Beta(5w_e, 5(1-w_e))$ , where  $w_e$  is the true edge weight as specified above.

We model non-strategic and strategic competitors by selecting the seed set  $S_B$  uniformly at random (denoted as RD) or by running the non-competitive influence maximization algorithm (denoted as IM). In our experiments, we set  $|S_A| = |S_B| = 5$  for Yahoo-Ad and  $|S_A| = |S_B| = 10$  for the DM dataset, and  $B > A$ . Since the optimal solution given the true edge probabilities cannot be derived in polynomial time, for Yahoo-Ad, we use the greedy solution as the optimal baseline, which is a  $(1 - 1/e, 1)$ -approximate solution.

Table 2: Dataset Statistics

Network	$n$	$m$	Average Degree
DM	679	3,374	4.96
Yahoo-Ad	11,475	52,567	4.58

Table 3: Average Running Time (second/round)

Dataset	OCIM-OFU	OCIM-TS	OCIM-ETC	$\epsilon$ -greedy	EMP
Yahoo-Ad	1.221	1.641	0.729	1.244	1.226
DM	1.142	1.195	0.621	1.173	1.125

For the DM dataset, we use the IMM solution as the optimal baseline, which is a  $(1 - 1/e - \epsilon, 1 - n^{-l})$ -approximate solution. For frequentist regrets, we repeat each experiment 50 times and show the average regret with 95% confidence interval. For Bayesian regrets, we draw 5 problem instances according to the prior distributions, conduct 10 experiments in each instance and report the average Bayesian regret over the 50 experiments. Due to the space constraint, results of other settings are provided in the Appendix.

**Algorithms for comparison.** For OCIM-TS, since the true prior distribution is unknown for the frequentist setting, we use the uninformative prior  $Beta(1, 1)$  for each  $\mu_e$ . For OCIM-OFU, we shrink its confidence interval by  $\alpha_\rho$ , i.e.,  $\rho_i \leftarrow \alpha_\rho \sqrt{3 \ln t / 2T_i}$ , to speed up the learning. The role of  $\alpha_\rho$  represents a tradeoff between theoretical guarantees and real-world performance.  $\alpha_\rho \geq 1$  provides theoretical regret bounds for the worst-case (i.e., our algorithms have sublinear regret for any problem instance) and most of the bandit literature gives regret analysis under this condition. However, in practice, we often do not face the worst problem instance. Taking a more aggressive  $\alpha_\rho$  helps speed up the learning empirically (Liu et al., 2021), though the algorithms may incur linear regrets for bad problem instances (which are likely rare in practice), preventing us from achieving worst-case theoretical regret bounds. We compare OCIM-OFU/OCIM-TS to the  $\epsilon$ -Greedy algorithm with parameter  $\epsilon = 0$  (denoted as the EMP algorithm) and  $\epsilon = 0.01$ , which inputs the empirical mean into the offline oracle with  $1 - \epsilon$  probability and otherwise selects  $S_A$  uniformly at random. The results of OCIM-ETC are moved to the Appendix as it requires more rounds to learn than others.

**Running time.** We show the average running times for different algorithms in Table 3. For the Yahoo-Ad dataset, OCIM-ETC is the fastest one as it only needs to call the oracle for one time before the exploitation phase. The running time of OCIM-TS is slower than that of OCIM-OFU because it requires an extra sampling procedure to generate Thompson samples. For

<sup>2</sup><https://webscope.sandbox.yahoo.com>



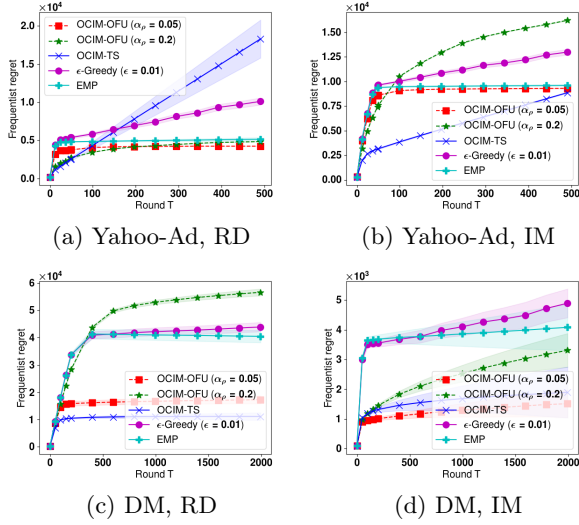


Figure 2: Frequentist regrets of algorithms for bipartite graph Yahoo-Ad and general graph DM.

the DM dataset, all algorithms consume less time since the graph is smaller, but the relative order for different algorithms are preserved.

**Experimental result for frequentist regrets** Figures 2a and 2b show the results for Yahoo-Ad. First, the regret of OCIM-OFU grows sub-linearly with respect to round  $T$  for all  $\alpha_\rho$ , consistent with Theorem 5.1’s regret bound. Second, we can observe that OCIM-OFU is superior to EMP and  $\epsilon$ -Greedy when  $\alpha_\rho = 0.05$ . When  $\alpha_\rho = 0.2$ , OCIM-OFU may have larger regret due to too much exploration. The OCIM-TS algorithm has larger slope in regrets compared to other algorithms. We speculate that such large slope comes from the uninformative prior, which requires more rounds to compensate for the mismatch of the uninformative and the true priors.

The results on the DM dataset are shown in Figs. 2c and 2d. Generally, they are consistent with those on the Yahoo-Ad dataset: OCIM-OFU also grows sub-linearly w.r.t round  $T$ . When  $\alpha_\rho = 0.05$ , OCIM-OFU has smaller regret than all baselines. Moreover, the difference between OCIM-OFU and the baselines for the non-strategic competitor (RD) is more significant than that of the strategic competitor’s (IM), because the non-strategic competitor is less “dominant” and OCIM-OFU can carefully trade off exploration and exploitation to maximize  $A$ ’s influence. OCIM-TS learns faster and achieves better performance in this dataset compared to that in the Yahoo-Ad dataset.

**Experimental result for Bayesian regrets** We show Bayesian regrets of all algorithms in Figure 3. All algorithms except for OCIM-TS have similar curves. OCIM-TS, however, achieves at least two orders of

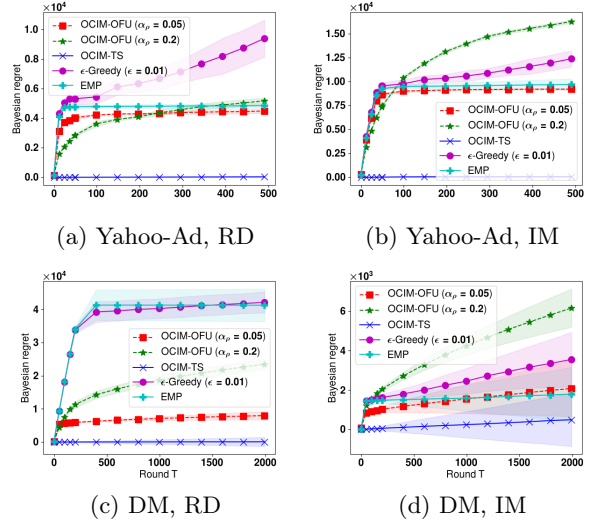


Figure 3: Bayesian regrets of algorithms for bipartite graph Yahoo-Ad and general graph DM.

magnitudes lower regret ( $BayesReg(T) \approx 100$ ) compared with other algorithms. The reason is that OCIM-TS leverages its prior knowledge to quickly converge to the optimal solution, but other algorithms cannot use this knowledge effectively.

## 8 Conclusion and Future Work

In this paper, we formulate the OCIM problem and introduce a general  $C^2$ MAB-T framework for it. We prove that one important condition required by prior CMAB algorithms, the TPM condition, still holds, while the other one, monotonicity, is not satisfied. We propose three algorithms that balance between prior knowledge, offline computation, feedback and regret bound: OCIM-TS relies on prior knowledge and achieves logarithmic Bayesian regret; OCIM-OFU needs to solve a harder offline problem and achieves logarithmic frequentist regret; and OCIM-ETC requires less feedback at the expense of a worse frequentist regret bound. We extend our framework to settings with more complex competitor actions.

This paper initiates the first study on OCIM, and it opens up a number of future directions. One is to design efficient offline approximation algorithms in the competitive setting when edge probabilities take a range of values. Another interesting direction is to study other partial feedback models, e.g. we only observe feedback from edges triggered by  $A$  but not  $B$ . A further direction is to look into distributed online learning, when competitors  $A$  and  $B$  both learn from the propagation and deploy their seeds accordingly.

## Acknowledgements

John C.S. Lui is supported in part by the GRF 14200321.

## References

- Bharathi, S., Kempe, D., and Salek, M. (2007). Competitive influence maximization in social networks. In *International workshop on web and internet economics*, pages 306–311.
- Borgs, C., Brautbar, M., Chayes, J., and Lucier, B. (2014). Maximizing social influence in nearly optimal time. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 946–957.
- Budak, C., Agrawal, D., and El Abbadi, A. (2011). Limiting the spread of misinformation in social networks. In *Proceedings of the 20th international conference on World wide web*, pages 665–674.
- Carnes, T., Nagarajan, C., Wild, S. M., and Van Zuylen, A. (2007). Maximizing influence in a competitive social network: a follower’s perspective. In *Proceedings of the ninth international conference on Electronic commerce*, pages 351–360.
- Chen, L., Xu, J., and Lu, Z. (2018). Contextual combinatorial multi-armed bandits with volatile arms and submodular reward. *Advances in Neural Information Processing Systems*, 31:3247–3256.
- Chen, W., Collins, A., Cummings, R., Ke, T., Liu, Z., Rincon, D., Sun, X., Wang, Y., Wei, W., and Yuan, Y. (2011). Influence maximization in social networks when negative opinions may emerge and propagate. In *Proceedings of the 2011 siam international conference on data mining*, pages 379–390. SIAM.
- Chen, W., Lakshmanan, L. V. S., and Castillo, C. (2013). *Information and Influence Propagation in Social Networks*. Morgan & Claypool Publishers.
- Chen, W., Wang, Y., Yuan, Y., and Wang, Q. (2016). Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *The Journal of Machine Learning Research*, 17(1):1746–1778.
- Garivier, A., Lattimore, T., and Kaufmann, E. (2016). On explore-then-commit strategies. In *Advances in Neural Information Processing Systems*, pages 784–792.
- He, X., Song, G., Chen, W., and Jiang, Q. (2012). Influence blocking maximization in social networks under the competitive linear threshold model. In *Proceedings of the 2012 siam international conference on data mining*, pages 463–474.
- Hüyük, A. and Tekin, C. (2020). Thompson sampling for combinatorial network optimization in unknown environments. *IEEE/ACM Transactions on Networking*, 28(6):2836–2849.
- Ivanov, S., Theocharidis, K., Terrovitis, M., and Karas, P. (2017). Content recommendation for viral social influence. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 565–574.
- Kempe, D., Kleinberg, J., and Tardos, É. (2003). Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 137–146.
- Li, S., Kong, F., Tang, K., Li, Q., and Chen, W. (2020). Online influence maximization under linear threshold model. In *Advances in Neural Information Processing Systems*.
- Li, Y., Fan, J., Wang, Y., and Tan, K. (2018). Influence maximization on social graphs: A survey. *IEEE Trans. Knowl. Data Eng.*, 30(10):1852–1872.
- Lin, Y. and Lui, J. C. (2015). Analyzing competitive influence maximization problems with partial information: An approximation algorithmic framework. *Performance Evaluation*, 91:187–204.
- Liu, X., Zuo, J., Chen, X., Chen, W., and Lui, J. C. (2021). Multi-layered network exploration via random walks: From offline optimization to online learning. In *International Conference on Machine Learning*, pages 7057–7066. PMLR.
- Merlis, N. and Mannor, S. (2020). Tight lower bounds for combinatorial multi-armed bandits. *Proceedings of Thirty Third Conference on Learning Theory*.
- Nguyen, H. T., Thai, M. T., and Dinh, T. N. (2016). Stop-and-stare: Optimal sampling algorithms for viral marketing in billion-scale networks. In *SIGMOD*, pages 695–710.
- Perrault, P., Healey, J., Wen, Z., and Valko, M. (2020). Budgeted online influence maximization. In *International Conference on Machine Learning*, pages 7620–7631. PMLR.
- Qin, L., Chen, S., and Zhu, X. (2014). Contextual combinatorial bandit and its application on diversified online recommendation. In *Proceedings of the 2014 SIAM International Conference on Data Mining*, pages 461–469. SIAM.
- Russo, D. and Van Roy, B. (2014). Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243.

- Russo, D. and Van Roy, B. (2016). An information-theoretic analysis of thompson sampling. *The Journal of Machine Learning Research*, 17(1):2442–2471.
- Tang, J., Sun, J., Wang, C., and Yang, Z. (2009). Social influence analysis in large-scale networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 807–816.
- Tang, J., Tang, X., Xiao, X., and Yuan, J. (2018). Online processing algorithms for influence maximization. In *SIGMOD*, pages 991–1005.
- Tang, Y., Shi, Y., and Xiao, X. (2015). Influence maximization in near-linear time: A martingale approach. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 1539–1554.
- Vaswani, S., Kveton, B., Wen, Z., Ghavamzadeh, M., Lakshmanan, L. V., and Schmidt, M. (2017). Model-independent online learning for influence maximization. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3530–3539.
- Wang, Q. and Chen, W. (2017). Improving regret bounds for combinatorial semi-bandits with probabilistically triggered arms and its applications. In *Advances in Neural Information Processing Systems*, pages 1161–1171.
- Wang, S. and Chen, W. (2018). Thompson sampling for combinatorial semi-bandits. In *International Conference on Machine Learning*, pages 5114–5122. PMLR.
- Wen, Z., Kveton, B., Valko, M., and Vaswani, S. (2017). Online influence maximization under independent cascade model with semi-bandit feedback. In *Advances in neural information processing systems*, pages 3022–3032.
- Wu, Q., Li, Z., Wang, H., Chen, W., and Wang, H. (2019). Factorization bandits for online influence maximization. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 636–646.

## Appendix

### A Proof of Theorem 3.1

*Proof.* Let  $r_S^v(\boldsymbol{\mu})$  be the probability that node  $v$  is activated by  $A$ . From the proof of Lemma 2 in (Wang and Chen, 2017), we know that if for every node  $v$  and every  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$  vectors we have

$$|r_S^v(\boldsymbol{\mu}) - r_S^v(\boldsymbol{\mu}')| \leq \sum_{e \in E} p_e^S(\boldsymbol{\mu}) |\mu_e - \mu'_e|, \quad (6)$$

then Theorem 3.1 is true. Notice that

$$r_S^v(\boldsymbol{\mu}) = \mathbb{E}_{L \sim \boldsymbol{\mu}} [\mathbb{1}\{v \text{ is activated by } A \text{ under } L\}] \quad (7)$$

$$r_S^v(\boldsymbol{\mu}') = \mathbb{E}_{L' \sim \boldsymbol{\mu}'} [\mathbb{1}\{v \text{ is activated by } A \text{ under } L'\}] \quad (8)$$

where  $L$  and  $L'$  are two live-edge graphs sampled under  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ , respectively. As mentioned in Sec. 3.2, we use an edge coupling method to compute the difference between  $r_S^v(\boldsymbol{\mu})$  and  $r_S^v(\boldsymbol{\mu}')$ . Specifically, for each edge  $e$ , suppose we independently draw a uniform random variable  $X_e$  over  $[0, 1]$ , let

$$\begin{aligned} L(e) = L'(e) = 1, & & \text{if } X_e \leq \min(\mu_e, \mu'_e) \\ L(e) = 1, L'(e) = 0, & & \text{if } \mu'_e < X_e < \mu_e \\ L(e) = 0, L'(e) = 1, & & \text{if } \mu_e < X_e < \mu'_e \\ L(e) = L'(e) = 0, & & \text{if } X_e \geq \max(\mu_e, \mu'_e) \end{aligned}$$

where  $L(e)$  represents the live/blocked state of edge  $e$  in live-edge graph  $L$ . Notice that  $L$  and  $L'$  does not have the subgraph relationship. Let  $\mathbf{X} := (X_1, \dots, X_e)$ , the difference can be written as:

$$r_S^v(\boldsymbol{\mu}) - r_S^v(\boldsymbol{\mu}') = \mathbb{E}_{\mathbf{X}} [f(S, L, v) - f(S, L', v)], \quad (9)$$

where  $f(S, L, v) := \mathbb{1}\{v \text{ is activated by } A \text{ under } L\}$ . Since  $f(S, L, v) - f(S, L', v)$  could be 0, 1 or -1, we will discuss these cases separately.

1)  $f(S, L, v) - f(S, L', v) = 0$ .

This will not contribute to the expectation.

2)  $f(S, L, v) - f(S, L', v) = 1$ .

This will occur only if there exists a path such that: under  $L$ ,  $v$  can be activated by  $A$  via this path, while under  $L'$ ,  $v$  cannot be activated by  $A$  via this path. We denote this event as  $\mathcal{E}_1$ . We will show that  $\mathcal{E}_1$  occurs only if at least one of  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  occurs.

$\mathcal{E}_1^A$ : There exists a path  $u \rightarrow v_1 \rightarrow \dots \rightarrow v_d = v$  such that:

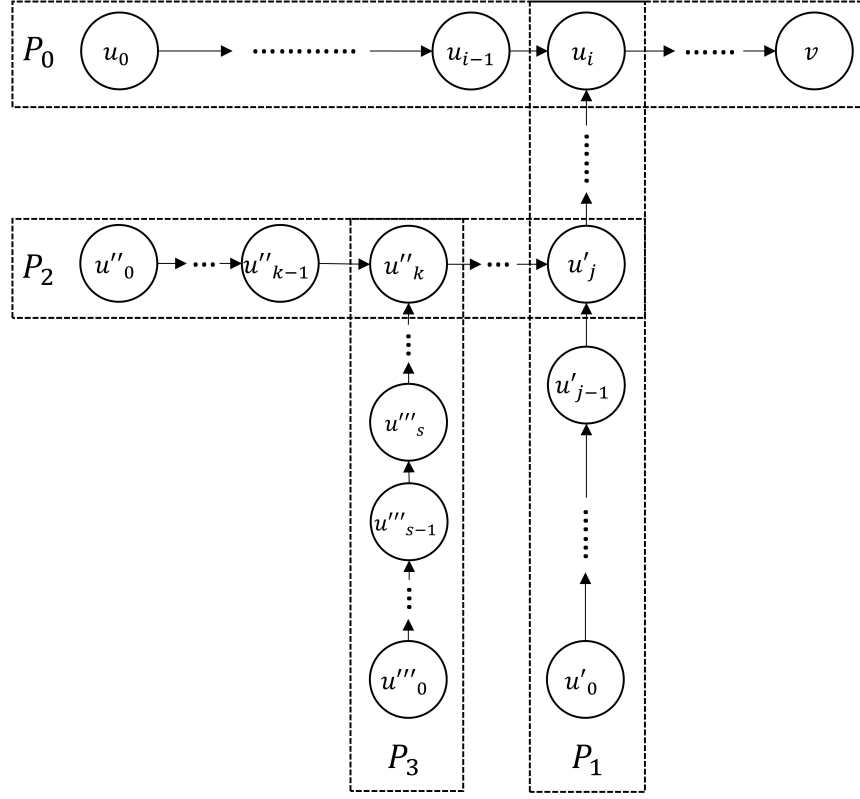
1.  $u$  is activated by  $A$  under both  $L$  and  $L'$
2. edge  $(u, v_1)$  is live under  $L$  but not  $L'$

$\mathcal{E}_1^B$ : There exists a path  $u' \rightarrow v'_1 \rightarrow \dots \rightarrow v'_d = v$  such that:

1.  $u'$  is activated by  $B$  under both  $L$  and  $L'$
2. edge  $(u', v'_1)$  is live under  $L'$  but not  $L$

**Lemma A.1.**  $\mathcal{E}_1$  occurs only if at least one of  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  occurs.

*Proof.* Let us first discuss the relationship between  $\mathcal{E}_1$ ,  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$ . For  $\mathcal{E}_1$ , if  $v$  can be activated by  $A$  under  $L$  but not  $L'$ , it is because either: (a) some edge  $e = (u, w)$  is live in  $L$  but blocked in  $L'$  while  $u$  is  $A$ -activated (or equivalently  $e$  is  $A$ -triggered); or (b) some edge  $e$  is live in  $L'$  but blocked in  $L$  while  $e$  is  $B$ -triggered. The former could be relaxed to  $\mathcal{E}_1^A$ , and the latter could be relaxed to  $\mathcal{E}_1^B$ . Notice that  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  are not mutually exclusive and we are interested in the upper bound of  $\mathbb{P}\{\mathcal{E}_1\}$ .


 Figure 4: Path  $P_0, P_1, P_2$  and  $P_3$ 

Assuming  $\mathcal{E}_1$  is true, consider the shortest path  $P_0 := \{u_0 \rightarrow u_1 \rightarrow \dots \rightarrow u_{l_0} = v\}$  from one seed node of  $A$ ,  $u_0$ , to node  $v$ , such that under  $L$  node  $v$  is activated by  $A$  but under  $L'$  it is not. When  $\mathcal{E}_1$  is true, there must exist a node that is not activated by  $A$  in  $P_0$  under  $L'$ . We denote the first node from  $u_0$  to  $v$  (i.e., closest to  $u_0$ ) in  $P_0$  that is not activated by  $A$  under  $L'$  as  $u_i$ .

Next, let us consider the live/blocked state of edge  $(u_{i-1}, u_i)$ . We already know edge  $(u_{i-1}, u_i)$  is live under  $L$ . If edge  $(u_{i-1}, u_i)$  is blocked under  $L'$ , since  $u_{i-1}$  is activated by  $A$  under both  $L$  and  $L'$ , it directly becomes  $\mathcal{E}_1^A$ . Otherwise, if edge  $(u_{i-1}, u_i)$  is live under  $L'$ , the reason that node  $u_i$  is not activated by  $A$  could only be that it is activated by  $B$ . In this case, there must exist a path  $P_1 := \{u'_0 \rightarrow u'_1 \rightarrow \dots \rightarrow u'_{l_1} = u_i\}$  from one seed node of  $B$ ,  $u'_0$ , to node  $u_i$ , such that  $u_i$  is activated by  $B$  under  $L'$  but not  $L$ . This can only occur when there exists a node that is not activated by  $B$  in  $P_1$  under  $L$ . We denote the first node from  $u'_0$  to  $u'_{l_1}$  (i.e., closest to  $u'_0$ ) in  $P_1$  that is not activated by  $B$  under  $L$  as  $u'_j$ . Notice that when the tie-breaking rule is  $A > B$ , we have  $l_1 < i \leq l_0$  as  $B$  should arrive at  $u_i$  earlier than  $A$ ; when the tie-breaking rule is  $B > A$ , we have  $l_1 \leq i \leq l_0$  as  $B$  should arrive at  $u_i$  no later than  $A$ . We will discuss the case of the proportional tie-breaking rule separately after the discussion of the dominance tie-breaking rules.

Then, let us consider the live/blocked state of edge  $(u'_{j-1}, u'_j)$ . We already know edge  $(u'_{j-1}, u'_j)$  is live under  $L'$ . If edge  $(u'_{j-1}, u'_j)$  is blocked under  $L$ , since  $u'_{j-1}$  is activated by  $B$  under both  $L$  and  $L'$ , it directly becomes  $\mathcal{E}_1^B$ . Otherwise, if edge  $(u'_{j-1}, u'_j)$  is live under  $L$ , the reason that node  $u'_j$  is not activated by  $B$  could only be that it is activated by  $A$ . It also means neither  $\mathcal{E}_1^A$  nor  $\mathcal{E}_1^B$  occurs so far. In this case, there must exist a path  $P_2 := \{u''_0 \rightarrow u''_1 \rightarrow \dots \rightarrow u''_{l_2} = u'_j\}$  from one seed node of  $A$ ,  $u''_0$ , to node  $u'_j$ , such that  $u'_j$  is activated by  $A$  under  $L$  but not  $L'$ . This can only occur when there exists a node that is not activated by  $A$  in  $P_2$  under  $L'$ . We denote the first node from  $u''_0$  to  $u''_{l_2}$  (i.e., closest to  $u''_0$ ) in  $P_2$  that is not activated by  $A$  under  $L'$  as  $u''_k$ . Notice that when  $A > B$ , we have  $l_2 \leq j \leq l_1 < l_0$  as  $A$  should arrive at  $u'_j$  no later than  $B$ ; when  $B > A$ , we have  $l_2 < j \leq l_1 \leq l_0$  as  $A$  should arrive at  $u'_j$  earlier than  $B$ .

Now let us consider the live/blocked state of edge  $(u''_{k-1}, u''_k)$ . We already know edge  $(u''_{k-1}, u''_k)$  is live under  $L$ .

If edge  $(u''_{k-1}, u''_k)$  is blocked under  $L'$ , since  $u''_{k-1}$  is activated by  $A$  under both  $L$  and  $L'$ , it directly becomes  $\mathcal{E}_1^A$ . Otherwise, if edge  $(u''_{k-1}, u''_k)$  is live under  $L'$ , the reason that node  $u''_k$  is not activated by  $A$  could only be that it is activated by  $B$ . In this case, there must exist a path  $P_3 := \{u''_0 \rightarrow u''_1 \rightarrow \dots \rightarrow u''_{l_3} = u''_k\}$  from one seed node of  $B$ ,  $u''_0$ , to node  $u''_k$ , such that  $u''_k$  is activated by  $B$  under  $L'$  but not  $L$ . This can only occur when there exists a node that is not activated by  $B$  in  $P_3$  under  $L$ . We denote the first node from  $u''_0$  to  $u''_{l_3}$  (i.e., closest to  $u''_0$ ) in  $P_3$  that is not activated by  $B$  under  $L$  as  $u''_{l_3}$ . Notice that when  $A > B$ , we have  $l_3 < k \leq l_2 \leq l_1$  as  $B$  should arrive at  $u''_k$  earlier than  $A$ ; when  $B > A$ , we have  $l_3 \leq k \leq l_2 < l_1$  as  $B$  should arrive at  $u''_k$  no later than  $A$ .

Again, let us consider the live/blocked state of edge  $(u'''_{s-1}, u'''_s)$ . We already know edge  $(u'''_{s-1}, u'''_s)$  is live under  $L'$ . If edge  $(u'''_{s-1}, u'''_s)$  is blocked under  $L$ , since  $u'''_{s-1}$  is activated by  $B$  under both  $L$  and  $L'$ , it directly becomes  $\mathcal{E}_1^B$ . Otherwise, if edge  $(u'''_{s-1}, u'''_s)$  is live under  $L$ , similar to the discussion above, we need to consider a new path  $P_4$  with length  $l_4$  and  $l_4 < l_2$ .

For the case of the proportional tie-breaking rule, in addition to the edge coupling, we also need to couple the permutation order (Chen et al., 2011) for each node in  $L$  and  $L'$ . More specific, for each node  $j$ , we randomly permute all of its in-neighbors, then when we need to break a tie on  $j$ , we find its activated neighbor  $i$  that is ordered first in the permutation order, and assign the state of  $i$  as  $j$ 's state. Assuming the same permutation order in  $L$  and  $L'$ , let us consider path  $P_0$  and  $P_1$  again. If  $l_0 = l_1$ , then  $u_i$  must be  $v$ . If  $\mathcal{E}_1^A$  does not occur in  $P_0$ , then the only neighbor of  $v$  in  $P_1$  must be ordered before the only neighbor of  $v$  in  $P_0$  in the permutation order on  $v$ . However, if  $\mathcal{E}_1^B$  does not occur in  $P_1$ , with such permutation order, it is impossible that  $v$  is activated by  $A$  under  $L$  but not  $L'$ . As a result, if neither  $\mathcal{E}_1^A$  nor  $\mathcal{E}_1^B$  occurs in path  $P_0$  and  $P_1$ , we have  $l_2 \leq l_1 < l_0$  in the case of the proportional tie-breaking rule.

To sum up, if neither  $\mathcal{E}_1^A$  nor  $\mathcal{E}_1^B$  occurs in path  $P_0$  and  $P_1$ , we need to check whether they could occur in a new path  $P_2$  shorter than  $P_0$ , and  $P_3$  shorter than  $P_1$ . As a result, we only need to check whether  $\mathcal{E}_1^A$  or  $\mathcal{E}_1^B$  occurs in the path with only one edge. In that case,  $\mathcal{E}_1^A$  or  $\mathcal{E}_1^B$  occurs for sure. Thus, by induction, we conclude that at least one of  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  occurs when considering any path with more than one edge, so  $\mathcal{E}_1$  will occur only if at least one of  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  occurs.  $\square$

Now, let us consider the two events in  $\mathcal{E}_1^A$  for a specific edge  $e = (u, v_1)$ . We find that the first event  $\{u$  is activated by  $A$  under both  $L$  and  $L'\}$ , is independent of the second event  $\{\text{edge } e \text{ is live under } L \text{ but not } L'\}$ , since the live/blocked state of edge  $e$  does not affect the activation of its tail node  $u$ . Also, for edge  $e = (u, v_1)$ , the probability of these two events can be written as

$$\mathbb{P}\{u \text{ is activated by } A \text{ under } L \text{ and } L'\} = \mathbb{P}\{e \text{ is triggered by } A \text{ under } L \text{ and } L'\}, \quad (10)$$

$$\mathbb{P}\{e \text{ is live under } L \text{ but not } L'\} = \begin{cases} \mu_e - \mu'_e & \text{if } \mu_e > \mu'_e \\ 0 & \text{otherwise.} \end{cases} \quad (11)$$

As a result, we have:

$$\mathbb{P}\{\mathcal{E}_1^A\} \leq \sum_{e: \mu_e > \mu'_e} \mathbb{P}\{e \text{ is triggered by } A \text{ under } L \text{ and } L'\}(\mu_e - \mu'_e) \quad (12)$$

Since  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$  are symmetric, we also have:

$$\mathbb{P}\{\mathcal{E}_1^B\} \leq \sum_{e: \mu'_e > \mu_e} \mathbb{P}\{e \text{ is triggered by } B \text{ under } L \text{ and } L'\}(\mu'_e - \mu_e) \quad (13)$$

Combining with Lemma. A.1, we have

$$\mathbb{P}\{\mathcal{E}_1\} \leq \mathbb{P}\{\mathcal{E}_1^A\} + \mathbb{P}\{\mathcal{E}_1^B\} \quad (14)$$

3)  $f(S, \mathbf{w}_1, v) - f(S, \mathbf{w}_2, v) = -1$ .

Similar to the previous case, this will occur only if there exists a path such that: under  $L'$ ,  $v$  can be activated by  $A$  via this path, while under  $L$ ,  $v$  cannot be activated by  $A$  via this path. We denote this event as  $\mathcal{E}_{-1}$ . We show that  $\mathcal{E}_{-1}$  occurs only if at least one of  $\mathcal{E}_{-1}^A$  and  $\mathcal{E}_{-1}^B$  occurs.

$\mathcal{E}_{-1}^A$ : There exists a path  $u \rightarrow v_1 \rightarrow \dots \rightarrow v_d = v$  such that:

1.  $u$  is activated by  $A$  under both  $L$  and  $L'$
2. edge  $(u, v_1)$  is live under  $L'$  but not  $L$

$\mathcal{E}_{-1}^B$ : There exists a path  $u' \rightarrow v'_1 \rightarrow \dots \rightarrow v'_{d'} = v$  such that:

1.  $u'$  is activated by  $B$  under both  $L$  and  $L'$
2. edge  $(u', v'_1)$  is live under  $L$  but not  $L'$

Since they are symmetric with  $\mathcal{E}_1^A$  and  $\mathcal{E}_1^B$ , following the same analysis, we can get

$$\mathbb{P}\{\mathcal{E}_{-1}^A\} \leq \sum_{e: \mu'_e > \mu_e} \mathbb{P}\{e \text{ is triggered by } A \text{ under } L \text{ and } L'\}(\mu'_e - \mu_e) \quad (15)$$

$$\mathbb{P}\{\mathcal{E}_{-1}^B\} \leq \sum_{e: \mu_e > \mu'_e} \mathbb{P}\{e \text{ is triggered by } B \text{ under } L \text{ and } L'\}(\mu_e - \mu'_e) \quad (16)$$

$$\mathbb{P}\{\mathcal{E}_{-1}\} \leq \mathbb{P}\{\mathcal{E}_{-1}^A\} + \mathbb{P}\{\mathcal{E}_{-1}^B\} \quad (17)$$

Combining all cases together, we have:

$$\begin{aligned} |r_S^v(\boldsymbol{\mu}) - r_S^v(\boldsymbol{\mu}')| &= |\mathbb{E}_{\mathbf{X}}[f(S, L, v) - f(S, L', v)]| \\ &\leq |1 \cdot \mathbb{P}\{\mathcal{E}_1\} + (-1) \cdot \mathbb{P}\{\mathcal{E}_{-1}\}| \\ &\leq |1 \cdot (\mathbb{P}\{\mathcal{E}_1^A\} + \mathbb{P}\{\mathcal{E}_1^B\}) + (-1) \cdot (\mathbb{P}\{\mathcal{E}_{-1}^A\} + \mathbb{P}\{\mathcal{E}_{-1}^B\})| \\ &\leq \sum_{e \in E} \mathbb{P}\{e \text{ is triggered by } A \text{ or } B \text{ under } L \text{ and } L'\} |\mu_e - \mu'_e|. \end{aligned} \quad (18)$$

The last inequality above is due to:

$$\begin{aligned} |\mathbb{P}\{\mathcal{E}_1^A\} - \mathbb{P}\{\mathcal{E}_{-1}^B\}| &\leq \sum_{e: \mu_e > \mu'_e} \mathbb{P}\{e \text{ is triggered by } A \text{ or } B \text{ under } L \text{ and } L'\} |\mu_e - \mu'_e| \\ |\mathbb{P}\{\mathcal{E}_1^B\} - \mathbb{P}\{\mathcal{E}_{-1}^A\}| &\leq \sum_{e: \mu'_e > \mu_e} \mathbb{P}\{e \text{ is triggered by } A \text{ or } B \text{ under } L \text{ and } L'\} |\mu_e - \mu'_e| \end{aligned}$$

Notice that Eq.(18) could be relaxed to:

$$\begin{aligned} |r_S^v(\boldsymbol{\mu}) - r_S^v(\boldsymbol{\mu}')| &\leq \sum_{e \in E} \mathbb{P}\{e \text{ is triggered by } A \text{ or } B \text{ under } L\} |\mu_e - \mu'_e| \\ &\leq \sum_{e \in E} p_e^S(\boldsymbol{\mu}) |\mu_e - \mu'_e|. \end{aligned} \quad (19)$$

□

## B Proof of Theorem 4.1

*Proof.* We define  $G^{(t)}$  as the feedback of OCIM in round  $t$ , which includes the outcomes of  $X_i^{(t)}$  for all  $i \in \tau_t$ . We denote by  $\mathcal{F}_{t-1}$  the history  $(S^{(1)}, G^{(1)}, \dots, S^{(t-1)}, G^{(t-1)})$  of observations available to the player when choosing an action  $S^{(t)}$ . For the Bayesian analysis, we assume the mean vector  $\boldsymbol{\mu}$  follows a prior distribution  $\mathcal{Q}$ . In round  $t$ , given  $\mathcal{F}_{t-1}$ , we define the posterior distribution of  $\boldsymbol{\mu}$  as  $\mathcal{Q}^{(t)}$  (i.e.,  $\boldsymbol{\mu}^{(t)} \sim \mathcal{Q}^{(t)}$  where  $\boldsymbol{\mu}^{(t)}$  is given in Alg. 1). As mentioned in Section 4, OCIM-TS allows any benchmark offline oracles, including approximation oracles. We consider a general benchmark oracle  $\mathcal{O}(S_B, \boldsymbol{\mu})$ . As oracle  $\mathcal{O}$  might be a randomized policy (e.g., an  $(\alpha, \beta)$ -approximation oracle with success probability  $\beta$ ), we use a random variable  $\omega \sim \Omega$  to represent all its randomness. In order to discuss the performance of OCIM-TS with oracle  $\mathcal{O}$ , we rewrite the Bayesian regret in Eq.(2) as

$$\text{BayesReg}(T) = \mathbb{E}_{\omega \sim \Omega, \boldsymbol{\mu} \sim \mathcal{Q}} \left[ \sum_{t=1}^T \left( r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu})}(\boldsymbol{\mu}) - r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}) \right) \right]. \quad (20)$$

Notice that  $\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu})$  is the action taken by the player if the true  $\boldsymbol{\mu}$  is known, while  $\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)$  is the real action chosen by OCIM-TS. The original regret definition in Eq.(2) is a special case of Eq.(20) for an  $(\alpha, \beta)$ -approximation oracle, and will focus on this general form in this proof.

The key step to derive the Bayesian regret bound of OCIM-TS is to show that the conditional distributions of  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}_t$  given  $\mathcal{F}_{t-1}$  are the same:

$$\mathbb{P}(\boldsymbol{\mu} = \cdot \mid \mathcal{F}_{t-1}) = \mathbb{P}(\boldsymbol{\mu}_t = \cdot \mid \mathcal{F}_{t-1}), \quad (21)$$

which is true since we use Thompson sampling to update the posterior distribution of  $\boldsymbol{\mu}$ . With this finding, we consider the Bayesian regret in Eq.(2):

$$\begin{aligned} & \text{BayesReg}(T) \\ &= \mathbb{E}_{\omega \sim \Omega} \left[ \sum_{t=1}^T \mathbb{E}_{\boldsymbol{\mu} \sim \mathcal{Q}, \boldsymbol{\mu}_t \sim \mathcal{Q}_t} \left[ r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu})}(\boldsymbol{\mu}) - r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}) \right] \right] \end{aligned} \quad (22)$$

$$= \mathbb{E}_{\omega \sim \Omega} \left[ \sum_{t=1}^T \mathbb{E}_{\mathcal{F}_{t-1}} \left[ \mathbb{E}_{\boldsymbol{\mu} \sim \mathcal{Q}, \boldsymbol{\mu}_t \sim \mathcal{Q}_t} \left[ r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu})}(\boldsymbol{\mu}) - r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}) \right] \mid \mathcal{F}_{t-1} \right] \right] \quad (23)$$

$$= \mathbb{E}_{\omega \sim \Omega} \left[ \sum_{t=1}^T \mathbb{E}_{\mathcal{F}_{t-1}} \left[ \mathbb{E}_{\boldsymbol{\mu} \sim \mathcal{Q}, \boldsymbol{\mu}_t \sim \mathcal{Q}_t} \left[ r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}_t) - r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}) \right] \mid \mathcal{F}_{t-1} \right] \right] \quad (24)$$

$$= \mathbb{E} \left[ \sum_{t=1}^T \left[ r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}_t) - r_{\mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)}(\boldsymbol{\mu}) \right] \right], \quad (25)$$

where Eq.(24) comes from applying Eq.(21) to Eq.(23). Let  $S_t = \mathcal{O}(S_B^{(t)}, \boldsymbol{\mu}_t)$  and  $\mathcal{C}_t = \{\boldsymbol{\mu}' : |\mu'_i - \hat{\mu}_{i,t}| \leq \rho_{i,t}, \forall i\}$ , where  $\rho_{i,t} = \sqrt{3 \ln t / 2T_{i,t-1}}$  and  $T_{i,t-1}$  is the total number of times arm  $i$  is played until round  $t$ . We define  $\Delta_{S_t} = r_{S_t}(\boldsymbol{\mu}_t) - r_{S_t}(\boldsymbol{\mu})$  and  $M = \sqrt{576\tilde{C}^2 m K \ln T / T}$ . By Eq.(25), we have

$$\begin{aligned} & \text{BayesReg}(T) \\ &= \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \right] \quad (26) \\ &\leq \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\Delta_{S_t} \geq M, \boldsymbol{\mu}_t \in \mathcal{C}_t, \boldsymbol{\mu} \in \mathcal{C}_t, \mathcal{N}_t^t\} \right]}_{(a)} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\boldsymbol{\mu}_t \notin \mathcal{C}_t\} \right] + \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} \right]}_{(b)} \\ &\quad + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\Delta_{S_t} \leq M\} \right]}_{(c)} + \underbrace{\mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\neg \mathcal{N}_t^t\} \right]}_{(d)} \quad (27) \end{aligned}$$

We can bound these three terms separately. For term (a), when  $\boldsymbol{\mu}_t \in \mathcal{C}_t, \boldsymbol{\mu} \in \mathcal{C}_t$ , we could bound  $|\mu_{i,t} - \mu_i| \leq |\mu_{i,t} - \hat{\mu}_{i,t}| + |\mu_i - \hat{\mu}_{i,t}| \leq 2\rho_{i,t}, \forall i$ . When  $\Delta_{S_t} \geq M$  and  $\mathcal{N}_t^t$  (Definition 7 in (Wang and Chen, 2017)) holds, by the proof of Lemma 5 in (Wang and Chen, 2017), we have  $\Delta_{S_t} \leq \sum_{i \in \tilde{S}_t} \kappa_{j_i, T}(M_i, N_{i, j_i, t-1})$  where  $\tilde{S}_t$  is the set of arms triggered by  $S_t$  and  $\kappa_{j_i, T}(M_i, N_{i, j_i, t-1})$  is defined in (Wang and Chen, 2017). We have



$$\begin{aligned}
 (a) &= \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\Delta_{S_t} \geq M, \boldsymbol{\mu}_t \in \mathcal{C}_t, \boldsymbol{\mu} \in \mathcal{C}_t, \mathcal{N}_t^t\} \right] \\
 &\leq \mathbb{E} \left[ \sum_{t=1}^T \sum_{i \in \mathcal{S}_t} \kappa_{j_i, T}(M_i, N_{i, j_i, t-1}) \right] \\
 &\leq \mathbb{E} \left[ \sum_{i \in [m]} \sum_{j=1}^{+\infty} \sum_{s=0}^{N_{i, j, T-1}} \kappa_{j, T}(M, s) \right] \\
 &\leq 4\tilde{C}m + \sum_{i \in [m]} \frac{576\tilde{C}^2 K \ln T}{M}
 \end{aligned}$$

For term (b), we can observe that  $\mathbb{E}[\mathbb{I}\{\boldsymbol{\mu} \in \mathcal{C}_t\} | \mathcal{F}_{t-1}] = \mathbb{E}[\mathbb{I}\{\boldsymbol{\mu}_t \in \mathcal{C}_t\} | \mathcal{F}_{t-1}]$ , since  $\mathcal{C}_t$  is determined given  $\mathcal{F}_{t-1}$ , and given  $\mathcal{F}_{t-1}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}_t$  follow the same distribution. Since  $\max_{S_t} \Delta_{S_t} \leq n$ , we have

$$\begin{aligned}
 (b) &= \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\boldsymbol{\mu}_t \notin \mathcal{C}_t\} \right] + \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} \right] \\
 &\leq n \left( \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}\{\boldsymbol{\mu}_t \notin \mathcal{C}_t\} \right] + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} \right] \right) \\
 &= n \left( \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}[\mathbb{I}\{\boldsymbol{\mu}_t \notin \mathcal{C}_t\} | \mathcal{F}_{t-1}] \right] \right) + n \left( \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}[\mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} | \mathcal{F}_{t-1}] \right] \right) \\
 &= 2n \left( \mathbb{E} \left[ \sum_{t=1}^T \mathbb{E}[\mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} | \mathcal{F}_{t-1}] \right] \right) \\
 &= 2n \left( \mathbb{E} \left[ \sum_{t=1}^T \mathbb{I}\{\boldsymbol{\mu} \notin \mathcal{C}_t\} \right] \right) \\
 &= 2n \left( \sum_{t=1}^T \mathbb{P}(\boldsymbol{\mu} \notin \mathcal{C}_t) \right) \\
 &\leq \frac{2\pi^2 mn}{3}
 \end{aligned}$$

For term (c), we can bound it by

$$(c) = \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\Delta_{S_t} \leq M\} \right] \leq TM$$

For term (d), similar to Eq.(20) in (Wang and Chen, 2017), we have

$$(d) = \mathbb{E} \left[ \sum_{t=1}^T \Delta_{S_t} \mathbb{I}\{\neg \mathcal{N}_t^t\} \right] \leq \frac{\pi^2}{6} \cdot \sum_{i \in [m]} j_{\max}^i \cdot n$$

Combine them together, we have

$$\text{BayesReg}(T) \leq 4\tilde{C}m + \sum_{i \in [m]} \frac{576\tilde{C}^2 K \ln T}{M} + \frac{2\pi^2 mn}{3} + TM + \frac{\pi^2}{6} \cdot \sum_{i \in [m]} j_{\max}(M) \cdot n$$

where  $j_{\max}(M) = \left\lceil \log_2 \frac{2\tilde{C}K}{M} \right\rceil_0$ . Take  $M = \sqrt{576\tilde{C}^2 m K \ln T / T}$ , we finally get the Bayesian regret bound of TS-OCIM:

$$\text{BayesReg}(T) \leq 12\tilde{C}\sqrt{mKT \ln T} + 2\tilde{C}m + \left( \left\lceil \log_2 \frac{T}{18 \ln T} \right\rceil_0 + 4 \right) \cdot \frac{\pi^2}{6} \cdot n \cdot m.$$

□

## C Proof of Theorem 5.1

*Proof.* We first introduce the following definitions to assist our analysis. Recall that  $\mathcal{S}^{(t)}$  is the action space in round  $t$ . We define the reward gap  $\Delta_S^{(t)} = \max(0, \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - r_S(\boldsymbol{\mu}))$  for all actions  $S \in \mathcal{S}^{(t)}$ . For each base arm  $i$ , we define  $\Delta_{\max}^{i,T} = \max_{t \in [T]} \sup_{S \in \mathcal{S}^{(t)}: p_i^S(\boldsymbol{\mu}) > 0, \Delta_S^{(t)} > 0} \Delta_S^{(t)}$  and  $\Delta_{\min}^{i,T} = \min_{t \in [T]} \inf_{S \in \mathcal{S}^{(t)}: p_i^S(\boldsymbol{\mu}) > 0, \Delta_S^{(t)} > 0} \Delta_S^{(t)}$ . If there is no action  $S$  such that  $p_i^S(\boldsymbol{\mu}) > 0$  and  $\Delta_S^{(t)} > 0$ , we define  $\Delta_{\max}^{i,T} = 0$  and  $\Delta_{\min}^{i,T} = +\infty$ . We define  $\Delta_{\max}^{(T)} = \max_{i \in [m]} \Delta_{\max}^{i,T}$  and  $\Delta_{\min}^{(T)} = \min_{i \in [m]} \Delta_{\min}^{i,T}$ . Let  $\tilde{S} = \{i \in [m] \mid p_i^S(\boldsymbol{\mu}) > 0\}$  be the set of arms that can be triggered by  $S$ . We define  $K = \max_{S \in \mathcal{S}^{(t)}} |\tilde{S}|$  as the largest number of arms could be triggered by a feasible action. We use  $\lceil x \rceil_0$  to denote  $\max\{\lceil x \rceil, 0\}$ . If  $\Delta_{\min}^{(T)} > 0$ , we provide the distribution-dependent bound of the OCIM-OFU algorithm.

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq \sum_{i \in [m]} \frac{576\tilde{C}^2 K \ln T}{\Delta_{\min}^{i,T}} + 4\tilde{C}m + \sum_{i \in [m]} \left( \left\lceil \log_2 \frac{2\tilde{C}K}{\Delta_{\min}^{i,T}} \right\rceil_0 + 2 \right) \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}.$$

To prove the distribution-dependent and the distribution-independent regret bounds, we generally follow the proof of Theorem 1 in Wang and Chen (2017). However, since we extend the original CMAB problem to a new contextual setting where the action space  $\mathcal{S}^{(t)}$  is the context, and monotonicity does not hold in the OCIM setting, we need to modify their analysis to tackle these changes. We introduce a positive real number  $M_i$  for each arm  $i$  and define  $M_{S^{(t)}} = \max_{i \in \tilde{S}^{(t)}} M_i$ . Define

$$\kappa_{j,T}(M, s) = \begin{cases} 4 \cdot 2^{-j} \tilde{C}, & \text{if } s = 0, \\ 2\tilde{C} \sqrt{\frac{72 \cdot 2^{-j} \ln T}{s}}, & \text{if } 1 \leq s \leq \ell_{j,T}(M), \\ 0, & \text{if } s \geq \ell_{j,T}(M) + 1, \end{cases}$$

where

$$\ell_{j,T}(M) = \left\lfloor \frac{288 \cdot 2^{-j} \tilde{C}^2 K^2 \ln T}{M^2} \right\rfloor.$$

Let  $\mathcal{N}_t^s$  be the event that at the beginning of round  $t$ , for every arm  $i \in [m]$ ,  $|\hat{\mu}_{i,t} - \mu_i| \leq 2\rho_{i,t}$ . Let  $\mathcal{H}_t$  be the event that at round  $t$  oracle  $\tilde{\mathcal{O}}$  outputs a solution,  $S^{(t)} = \{S_A^{(t)}, S_B^{(t)}\}$  and  $\boldsymbol{\mu}^{(t)} = (\mu_1^{(t)}, \dots, \mu_m^{(t)})$ , such that  $r_{S^{(t)}}(\boldsymbol{\mu}^{(t)}) < \alpha \cdot r_{S^*}(\boldsymbol{\mu}^*)$ , i.e., oracle  $\tilde{\mathcal{O}}$  fails to output an  $\alpha$ -approximate solution. Let  $\mathcal{N}_t^t$  be the event that the triggering is nice at the beginning of round  $t$  (Definition 7 in (Wang and Chen, 2017)). The following lemma explains how  $\kappa$  contributes to the regret.

**Lemma C.1.** *For any vector  $\{M_i\}_{i \in [m]}$  of positive real numbers and  $1 \leq t \leq T$ , if  $\{\Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}\}, \neg \mathcal{H}_t, \mathcal{N}_t^s$  and  $\mathcal{N}_t^t$  hold, we have*

$$\Delta_{S^{(t)}}^{(t)} \leq \sum_{i \in \tilde{S}^{(t)}} \kappa_{j_i, T}(M_i, N_{i, j_i, t-1}),$$

where  $j_i$  is the index of the TP group with  $S^{(t)} \in \mathcal{S}_{i, j_i}$  (see Definition 5 in (Wang and Chen, 2017)).

*Proof.* By  $\mathcal{N}_t^s$  and  $0 \leq \mu_i \leq 1$  for all  $i \in [m]$ , we have

$$\forall i \in [m], \mu_i \in c_{i,t} = [(\hat{\mu}_{i,t} - \rho_{i,t})^{0+}, (\hat{\mu}_{i,t} + \rho_{i,t})^{1-}]. \quad (28)$$

It means that we have the correct estimated range of  $\mu_i$  for all  $i \in [m]$  at round  $t$ . Combining with  $\neg\mathcal{H}_t$  for the offline oracle  $\tilde{\mathcal{O}}$ , we have

$$r_{S^{(t)}}(\boldsymbol{\mu}^{(t)}) \geq \alpha \cdot r_{S^*}(\boldsymbol{\mu}^*) \geq \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) = r_{S^{(t)}}(\boldsymbol{\mu}) + \Delta_{S^{(t)}}^{(t)}. \quad (29)$$

By the TPM condition in Theorem. 3.1, we have

$$\Delta_{S^{(t)}}^{(t)} \leq r_{S^{(t)}}(\boldsymbol{\mu}^{(t)}) - r_{S^{(t)}}(\boldsymbol{\mu}) \leq \tilde{C} \sum_{i \in [m]} p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i|. \quad (30)$$

We want to bound  $\Delta_{S^{(t)}}^{(t)}$  by bounding  $p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i|$ . We first perform a transformation. Since  $\Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}$ , we have  $\tilde{C} \sum_{i \in [m]} p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i| \geq \Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}$ . Then we have

$$\begin{aligned} \Delta_{S^{(t)}}^{(t)} &\leq \tilde{C} \sum_{i \in [m]} p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i| \\ &\leq -M_{S^{(t)}} + 2\tilde{C} \sum_{i \in [m]} p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i| \\ &\leq 2\tilde{C} \sum_{i \in [m]} \left[ p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i| - \frac{M_i}{2\tilde{C}K} \right]. \end{aligned} \quad (31)$$

In fact, if  $\mathcal{N}_t^s$  holds and  $\mu_i^{(t)} \in c_{i,t}$  for all  $i \in [m]$ ,

$$\forall i \in [m], |\mu_i^{(t)} - \mu_i| \leq 2\rho_{i,t} = 2\sqrt{\frac{3 \ln t}{2T_{i,t-1}}}. \quad (32)$$

So far, all requirements on bounding  $\Delta_{S_i}$  in Lemma 5 from (Wang and Chen, 2017) are also satisfied by  $\Delta_{S^{(t)}}^{(t)}$  of OCIM-OFU algorithm in the OCIM setting without monotonicity. We can then follow the same steps to bound  $p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i|$  in the two cases they considered (combining their Eq.(11)-(13)) and get

$$\begin{aligned} \Delta_{S^{(t)}}^{(t)} &\leq 2\tilde{C} \sum_{i \in [m]} \left[ p_i^{S^{(t)}}(\boldsymbol{\mu}) |\mu_i^{(t)} - \mu_i| - \frac{M_i}{2\tilde{C}K} \right] \\ &\leq \sum_{i \in \tilde{S}^{(t)}} \kappa_{j_i, T}(M_i, N_{i, j_i, t-1}). \end{aligned}$$

□

With Lemma C.1, we can follow the proof of Lemma 6 in (Wang and Chen, 2017) to bound the regret when  $\{\Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}\}, \neg\mathcal{H}_t, \mathcal{N}_t^s$  and  $\mathcal{N}_t^t$  hold.

$$\text{Reg}(\{\Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}\} \wedge \neg\mathcal{H}_t \wedge \mathcal{N}_t^s \wedge \mathcal{N}_t^t) \leq \sum_{i \in [m]} \frac{576\tilde{C}^2 K \ln T}{M_i} + 4\tilde{C}m. \quad (33)$$

Finally, we take  $M_i = \Delta_{\min}^{i, T}$ . If  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}}$ , then  $\Delta_{S^{(t)}}^{(t)} = 0$ , since we have either  $\tilde{S}^{(t)} = \emptyset$  or  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}} \leq M_i$  for some  $i \in \tilde{S}^{(t)}$ . Thus, no regret is accumulated when  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}}$ . Following Eq.(17)-(21) in (Wang and Chen, 2017), we can derive the distribution-dependent regret bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq \sum_{i \in [m]} \frac{576\tilde{C}^2 K \ln T}{\Delta_{\min}^{i, T}} + 4\tilde{C}m + \sum_{i \in [m]} \left( \left\lceil \log_2 \frac{2\tilde{C}K}{\Delta_{\min}^{i, T}} \right\rceil + 2 \right) \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}. \quad (34)$$

To derive the distribution-independent bound, we take  $M_i = M = \sqrt{(576\tilde{C}^2 m K \ln T)/T}$ , follow Eq.(23) in (Wang and Chen, 2017) and get

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq 12\tilde{C}\sqrt{mKT \ln T} + 2\tilde{C}m + \left( \left\lceil \log_2 \frac{T}{18 \ln T} \right\rceil + 2 \right) \cdot \frac{\pi^2}{6} \cdot n \cdot m. \quad (35)$$

□

## D Computational Efficiency of OCIM-OFU

### D.1 Proof of Theorem 5.2

*Proof.* In order to prove Theorem 5.2, we first introduce a new optimization problem denoted as  $P_1$ : given  $S$ , the new problem aims to find the optimal  $\mu_i$  for one edge  $i$  to maximize  $r_S(\boldsymbol{\mu})$ , while fixing the values of all others. The following lemma shows it is #P-hard.

**Lemma D.1.** *Given  $S$  and fixing  $\mu_e$  for all  $e \neq i$ , finding the optimal  $\mu_i \in c_i$  for one edge  $i$  that maximizes  $r_S(\boldsymbol{\mu})$  is #P-hard.*

*Proof.* We prove the hardness of this optimization problem via a reduction from the influence computation problem. We first consider a general graph  $G_0$  with  $n$  nodes and  $m$  edges, where all influence probabilities on edges are set to  $1/2$ . Given  $S_A$ , computing the influence spread of  $A$  in such a graph is #P-hard. Notice that there is no seed set of  $B$  in  $G_0$ . Now let us take one node  $v$  in  $G_0$  and denote its activation probability by  $A$  as  $h_A(G_0, S_A, v)$ . Actually, computing  $h_A(G_0, S_A, v)$  is also #P-hard and we want to show that it can be reduced to our optimization problem in polynomial time.

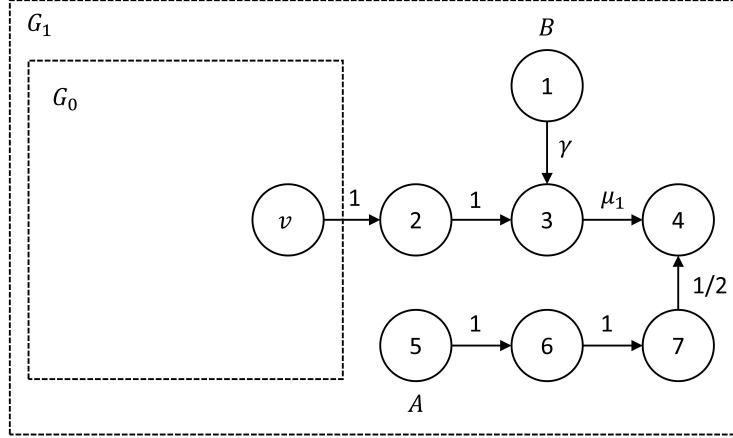


Figure 5: Construction of  $G_1$  based on  $G_0$

We first construct a new graph  $G_1$  based on  $G_0$ . For  $G_1$ , we keep  $G_0$  and  $S_A$  unchanged, then add several nodes and edges as shown in Fig. 5. We add node 1 to the seed set of  $B$  and node 5 to the seed set of  $A$ , so the joint action  $S = \{S_A \cup \{5\}, S_B = \{1\}\}$ . In this new graph  $G_1$ , we consider the optimization problem of finding the optimal  $\mu_1$  (influence probability on edge (3, 4)) within its range  $c_1$  that maximizes  $r_S(\boldsymbol{\mu})$ . Notice that the influence probability  $\gamma$  on edge (1, 3) is a constant and  $\mu_1$  would only affect the activation probability of node 4. We denote the activation probability by  $A$  of node 4 as  $h_A(G_1, S, 4)$ . In order to maximize  $r_S(\boldsymbol{\mu})$ , we only need to maximize  $h_A(G_1, S, 4)$ . It can be written as:

$$h_A(G_1, S, 4) = \frac{1}{2} \left[ (1 - \gamma) \cdot h_A(G_1, S, v) - \gamma \right] \cdot \mu_1 + \frac{1}{2}. \quad (36)$$

It is easy to see  $h_A(G_1, S, 4)$  has a linear relationship with  $\mu_1$ , so the optimal  $\mu_1$  could only be either the lower or upper bound of its range  $c_1$ . Assuming we can solve the optimization problem of finding the optimal  $\mu_1$ , then we can determine the sign of  $\mu_1$ 's coefficient in Eq.(36): if the optimal  $\mu_1$  is the upper bound value in  $c_1$ , we have  $(1 - \gamma) \cdot h_A(G_1, S, v) - \gamma \geq 0$ ; otherwise,  $(1 - \gamma) \cdot h_A(G_1, S, v) - \gamma < 0$ . It means we can answer the question that whether  $h_A(G_1, S, v)$  is larger (or smaller) than  $\frac{\gamma}{1-\gamma}$ . Notice that  $h_A(G_0, S_A, v) = h_A(G_1, S, v)$ , so we can manually change the value of  $\gamma$  to check whether  $h_A(G_0, S_A, v)$  is larger (or smaller) than  $x = \frac{\gamma}{1-\gamma}$  for any  $x \in [0, 1]$ . Recall that all edge probabilities in  $G_0$  are set to  $1/2$ , so the highest precision of  $h_A(G_0, S_A, v)$  should be  $2^{-m}$ . Hence, we can use a binary search algorithm to find the exact value of  $h_A(G_0, S_A, v)$  in at most

$m$  times. It means computing the activation probability of  $v$  in  $G_0$  can be reduced to the optimization problem of finding the optimal  $\mu_1$  in  $G_1$ , which completes the proof.  $\square$

We then show that  $P_1$  is a special case of Eq.(4). The main idea is to relax the constraints  $|S_A| \leq k$ ,  $S = \{S_A, S_B\}$  in Eq.(4) and show that it can find the optimal  $\mu$  for any given  $S$ . Consider a graph  $G$  with  $n$  nodes and a given seed set  $S = \{S_A, S_B\}$ . We construct a new graph  $G'$  by manually add additional  $n + 1$  nodes pointing from each seed node in  $S_A$ . If we can solve the optimization problem Eq.(4) in the new graph  $G'$ , since  $S_A$  must be the optimal seed set of  $A$  and the added nodes will not affect the prorogation in  $G$ , we will also find the optimal  $\mu_i$ 's in the original graph  $G$  for the given  $S$ . Then, it is easy to see  $P_1$  is a special case of Eq.(4) since  $P_1$  only find the optimal  $\mu_i$  for one edge  $i$ . With Lemma D.1, we know Eq.(4) is also #P-hard.  $\square$

## D.2 Non-submodularity of $g(S)$

In Section 5.1, we introduce  $g(S) = \max_{\mu} r_S(\mu)$ , which is an upper bound function of  $r_S(\mu)$  for each  $S$ . If  $g(S)$  is submodular over  $S$ , we can use a greedy algorithm on  $g(S)$  to find an approximate solution. However, the following example in Fig. 6 shows that  $g(S)$  is not submodular.

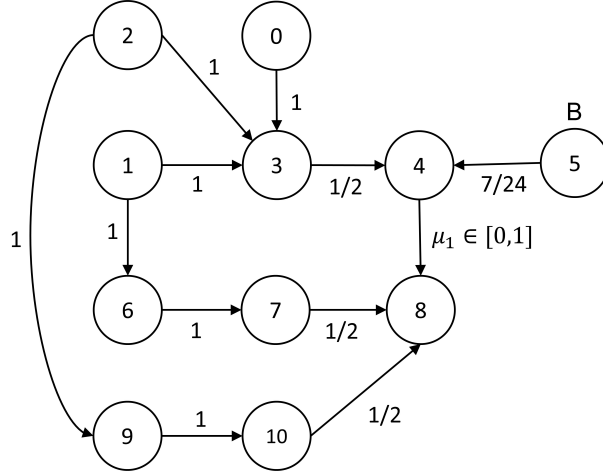


Figure 6: Example showing that  $g(S)$  is not submodular

In Fig. 6, the numbers attached to edges are influence probabilities. Only the influence probability of edge (4, 8) is a variable and we denote it as  $\mu_1$ . We assume  $\mu_1 \in [0, 1]$  and  $S_B = \{5\}$ . Let us consider some choices of  $S_A$ . When  $S_A$  is chosen as  $\{0\}$ ,  $\{0, 1\}$  or  $\{0, 2\}$ , the optimal  $\mu_1$  that maximizes  $r_S(\mu)$  is 1; when  $S_A$  is chosen as  $\{0, 1, 2\}$ , the optimal  $\mu_1$  that maximizes  $r_S(\mu)$  is 0. Based on this observation, we can calculate  $g(S)$  (assuming  $S_B = \{5\}$ ):

$$\begin{aligned} g(S_A = \{0\}) &= 2 + \frac{17}{24}, \\ g(S_A = \{0, 1\}) &= 5 + \frac{17}{24} \times \frac{4}{5}, \\ g(S_A = \{0, 2\}) &= 5 + \frac{17}{24} \times \frac{4}{5}, \\ g(S_A = \{0, 1, 2\}) &= 8 + \frac{17}{24} \times \frac{1}{2} + \frac{3}{4}. \end{aligned}$$

Thus we have

$$g(S_A = \{0, 1\}) + g(S_A = \{0, 2\}) < g(S_A = \{0\}) + g(S_A = \{0, 1, 2\}), \quad (37)$$

which is contrary to submodularity.

### D.3 Bipartite Graph

We consider a weighted bipartite graph  $G = (L, R, E)$  where each edge  $(u, v)$  is associated with a probability  $p(u, v)$ . Given the competitor's seed set  $S_B \subseteq L$ , we need to choose  $k$  nodes from  $L$  as  $S_A$  that maximizes the expected number of nodes activated by  $A$  in  $R$ , where a node  $v \in R$  can be activated by a node  $u \in L$  with an independent probability of  $p(u, v)$ . As mentioned before, if  $A$  and  $B$  are attempting to activate a node in  $L$  at the same time, the result will depend on the tie-breaking rule. If all edge probabilities are fixed, i.e.,  $\boldsymbol{\mu}$  is fixed,  $r_S(\boldsymbol{\mu})$  is still submodular over  $S_A$ , so we can use a greedy algorithm as a  $(1 - 1/e, 1)$ -approximation oracle  $\mathcal{O}_{\text{greedy}}$ . Based on it, let us discuss the new offline optimization problem in Eq.(4) under our two tie-breaking rules: (1)  $A > B$ : since  $B$  will never influence nodes in  $R$  earlier than  $A$  in bipartite graphs, and  $A$  will always win the competition, from  $A$ 's perspective, we can ignore  $S_B$  to choose  $S_A$ . In this case, all edge probabilities should take the maximum values: for all  $i \in E$ ,  $\mu_i$  equals to the upper bound of  $c_i$ , and we then use the oracle  $\mathcal{O}_{\text{greedy}}$  to find  $S_A$ . (2)  $B > A$ : since  $A$  will never influence nodes in  $R$  earlier than  $B$  in bipartite graphs, and  $B$  will always win the competition, all out-edges of  $S_B$ , denoted as  $E_{S_B}$ , should take the minimum probabilities to maximize the influence spread of  $A$ . All the other edges in  $E \setminus E_{S_B}$  should take the maximum probabilities. Formally, for all  $i \in E_{S_B}$ ,  $\mu_i$  equals to the lower bound of  $c_i$ ; for all  $i \in E \setminus E_{S_B}$ ,  $\mu_i$  equals to the upper bound of  $c_i$ . We then use the oracle  $\mathcal{O}_{\text{greedy}}$  to find  $S_A$ . To sum up, in bipartite graphs,  $r_S(\boldsymbol{\mu})$  is optimized by pre-determining  $\boldsymbol{\mu}$  based on the tie-breaking rule, and then using the greedy algorithm to get a  $(1 - 1/e, 1)$ -approximation solution. Since the time complexity of influence computation in the bipartite graph is  $O(m)$ , the time complexity of the offline algorithm is equal to that of the greedy algorithm,  $O(kmn)$ .

### D.4 General Graph

**GraphWe** The competitive propagation in the general graph is much more complicated, so it is hard to pre-determine all edge probabilities as in the bipartite graph case. However, we have a key observation:

**Lemma D.2.** *When fixing the seed set  $S = \{S_A, S_B\}$ , reward  $r_S(\boldsymbol{\mu})$  has a linear relationship with each  $\mu_i$  (when other  $\mu_j$ 's with  $j \neq i$  are fixed). This implies that the optimal solution for the optimization problem in Eq.(4) must occur at the boundaries of the intervals  $c_i$ 's.*

*Proof.* We can expand  $r_S(\boldsymbol{\mu})$  based on the live-edge graph model (Chen et al., 2013a):

$$r_S(\boldsymbol{\mu}) = \sum_L |\Gamma_A(L, S)| \cdot \Pr(L) = \sum_L |\Gamma_A(L, S)| \prod_{e \in E(L)} \mu_e \prod_{e \notin E(L)} (1 - \mu_e), \quad (38)$$

where  $L$  is one possible live-edge graph (each edge  $e \in E$  is in  $L$  with probability  $\mu_e$  and not in  $L$  with probability  $1 - \mu_e$ , and this is independent from other edges),  $\Gamma_A(L, S)$  is the set of nodes activated by  $A$  from seed sets  $S = \{S_A, S_B\}$  under live-edge graph  $L$  and  $E(L)$  is the set of edges that appear in live-edge graph  $L$ . Eq.(38) shows that  $r_S(\boldsymbol{\mu})$  is linear with each  $\mu_i$ , so the optimal  $\mu_i$  must take either the minimum or the maximum value in its range  $c_i$ .  $\square$

Lemma D.2 implies that for any edge  $e$  not reachable from  $B$  seeds, it is safe to always take its upper bound value since it can only help the propagation of  $A$ . This further suggests that if we only have a small number (e.g.  $\log m$ ) of edges reachable from  $B$ , then we can afford enumerating all the boundary value combinations of these edges. For each such boundary setting  $\boldsymbol{\mu}$ , we can use the IMM algorithm (Tang et al., 2014) to design a  $(1 - 1/e - \epsilon, 1 - n^{-l})$ -approximation oracle  $\mathcal{O}_{\text{IMM}}$  with time complexity  $T_{\text{IMM}} = O((k + l)(m + n) \log n / \epsilon^2)$ . We discuss such graphs that satisfy the above condition in directed trees. Specifically, we consider the in-arborescence, where all edges point towards the root. For any node  $u$  in the in-arborescence, there only exists one path from  $u$  to the root; if  $u$  is selected as the seed node of  $B$ , it could only propagate via this path. Hence, if the depth of the in-arborescence is in the order of  $O(\log m)$ , the number of edges reachable from  $S_B$  would be  $O(|S_B| \cdot \log m)$ . In this case, we can use the IMM algorithm for  $O(m^{|S_B|})$  combinations to obtain an approximate solution with time complexity  $O(m^{|S_B|} \cdot T_{\text{IMM}})$ . Examples of such in-arborescences with depth  $O(\log m)$  could be the complete or full binary trees.

For general graphs, designing efficient approximation algorithms for the offline problem in Eq. (4) remains a challenging open problem, due to the joint optimization over  $S$  and  $\boldsymbol{\mu}$  and the complicated function form of  $r_S(\boldsymbol{\mu})$ . Nevertheless, heuristic algorithms are still possible. In the experiment section, we employ the following

heuristic with the  $B > A$  tie-breaking rule: for all outgoing edges from  $B$  seeds, we set their influence probabilities to their lower bound values, while for the rest, we set them to their upper bound values. This setting guarantees that the first-level edges from the seeds are always set correctly, no matter how we select  $A$  seeds. They do not guarantee the correctness of second or higher level edge settings in the cascade, but the impact of those edges to influence spread decays significantly, so the above choice is reasonable as a heuristic.

## E Proof of Theorem 5.3

---

### Algorithm 3 OCIM-ETC with offline oracle $\mathcal{O}$

---

*Proof.* 1: **Input:**  $m, N, T$ , Oracle  $\mathcal{O}$ .

2: For each arm  $i$ ,  $T_i \leftarrow 0$ . {maintain the total number of times arm  $i$  is played so far.}

3: For each arm  $i$ ,  $\hat{\mu}_i \leftarrow 0$ . {maintain the empirical mean of  $X_i$ .}

4: **Exploration phase:**

5: **for**  $t = 1, 2, 3, \dots, \lceil nN/k \rceil$  **do**

6:   Take  $k$  nodes that have not been chosen for  $N$  times as  $S_A$ .

7:   Observe the feedback  $X_i^{(t)}$  for each direct out-edge of  $S_A$ ,  $i \in \tau_{\text{direct}}$ .

8:   For each arm  $i \in \tau_{\text{direct}}$  update  $T_i$  and  $\hat{\mu}_i$ :  $T_i = T_i + 1$ ,  $\hat{\mu}_i = \hat{\mu}_i + (X_i^{(t)} - \hat{\mu}_i)/T_i$ .

9: **end for**

10: **Exploitation phase:**

11: **for**  $t = \lceil nN/k \rceil + 1, \dots, T$  **do**

12:   Obtain context  $S_B^{(t)}$ .

13:    $S^{(t)} \leftarrow \mathcal{O}(S_B^{(t)}, \hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_m)$ .

14:   Play action  $S^{(t)}$ .

15: **end for**

---

The OCIM-ETC algorithm is described in Alg. 3. We utilize the following well-known tail bound in our proof.

**Lemma E.1.** (*Hoeffding's Inequality*) Let  $X_1, \dots, X_n$  be independent and identically distributed random variables with common support  $[0, 1]$  and mean  $\mu$ . Let  $Y = X_1 + \dots + X_n$ . Then for all  $\delta \geq 0$ ,

$$\mathbb{P}\{|Y - n\mu| \geq \delta\} \leq 2e^{-2\delta^2/n}.$$

Let  $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \dots, \hat{\mu}_m)$  be the empirical mean of  $\boldsymbol{\mu}$ . Recall that oracle  $\mathcal{O}$  takes  $S_B^{(t)}$  and  $\hat{\boldsymbol{\mu}}$  as inputs and outputs a solution  $S^{(t)}$ . Let us define event  $\mathcal{F} = \{r_{S^{(t)}}(\hat{\boldsymbol{\mu}}) < \alpha \cdot \text{opt}^{(t)}(\hat{\boldsymbol{\mu}})\}$ , which represents that oracle  $\mathcal{O}$  fails to output an  $\alpha$ -approximate solution, and we know  $\mathbb{P}(\mathcal{F}) < 1 - \beta$ .

With the same definitions in Appendix C, we can decompose the regret as:

$$\begin{aligned} \text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) &\leq \lceil nN/k \rceil \cdot \Delta_{\max}^{(T)} + \sum_{t=T-\lceil nN/k \rceil+1}^T \left[ \alpha\beta \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - \mathbb{E}[r_{S^{(t)}}(\hat{\boldsymbol{\mu}})] \right] \\ &\leq \lceil nN/k \rceil \cdot \Delta_{\max}^{(T)} + \sum_{t=T-\lceil nN/k \rceil+1}^T \left[ \alpha\beta \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - \beta \cdot \mathbb{E}[r_{S^{(t)}}(\hat{\boldsymbol{\mu}}) \mid \neg\mathcal{F}] \right] \\ &\leq \lceil nN/k \rceil \cdot \Delta_{\max}^{(T)} + \sum_{t=T-\lceil nN/k \rceil+1}^T \left[ \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - \mathbb{E}[r_{S^{(t)}}(\hat{\boldsymbol{\mu}}) \mid \neg\mathcal{F}] \right]. \end{aligned} \quad (39)$$

Next, let us rewrite the TPM condition in Theorem 3.1. For any  $S$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ , we have

$$\begin{aligned} |r_S(\boldsymbol{\mu}) - r_S(\boldsymbol{\mu}')| &\leq C \sum_{i \in [m]} p_i^S(\boldsymbol{\mu}) |\mu_i - \mu'_i| \\ &\leq C \sum_{i \in [m]} |\mu_i - \mu'_i| \\ &\leq Cm \cdot \max_{i \in [m]} |\mu_i - \mu'_i|, \end{aligned} \quad (40)$$

where  $C$  is the maximum number of nodes that any one node can reach in graph  $G$ . Let  $S_{\mu}^{*,t}$  denote the optimal action for  $\mu$  in round  $t$ . Under  $\neg\mathcal{F}$ , we have

$$\begin{aligned}
 r_{S^{(t)}}(\hat{\mu}) &\geq \alpha \cdot r_{S_{\mu}^{*,t}}(\hat{\mu}) \\
 &\geq \alpha \cdot r_{S_{\mu}^{*,t}}(\mu) - \alpha \cdot Cm \cdot \max_{i \in [m]} |\mu_i - \hat{\mu}_i| \\
 &\geq r_{S^{(t)}}(\mu) + \Delta_{S^{(t)}}^{(t)} - \alpha \cdot Cm \cdot \max_{i \in [m]} |\mu_i - \hat{\mu}_i|,
 \end{aligned} \tag{41}$$

where the third inequality is due to Eq.(40). Combining Eq.(40) and Eq.(41) together, we have

$$\begin{aligned}
 \Delta_{S^{(t)}}^{(t)} &\leq r_{S^{(t)}}(\hat{\mu}) - r_{S^{(t)}}(\mu) + \alpha \cdot Cm \cdot \max_{i \in [m]} |\mu_i - \hat{\mu}_i| \\
 &\leq (1 + \alpha) \cdot Cm \cdot \max_{i \in [m]} |\mu_i - \hat{\mu}_i|.
 \end{aligned} \tag{42}$$

Let us define  $\delta_0 := \frac{\Delta_{\min}^{(T)}}{2Cm}$ . If  $\max_{i \in [m]} |\mu_i - \hat{\mu}_i| < \delta_0$ , then we know  $S^{(t)}$  is at least an  $\alpha$ -approximate solution, such that  $\Delta_{S^{(t)}}^{(t)} = 0$ . Then the regret in Eq.(39) can be written as

$$\begin{aligned}
 \text{Reg}_{\alpha,\beta}(T; \mu) &\leq \lceil nN/k \rceil \cdot \Delta_{\max}^{(T)} + \left( T - \lceil nN/k \rceil \right) \cdot 2m \exp(-2N\delta_0^2) \cdot \Delta_{\max}^{(T)} \\
 &\leq \left( \lceil nN/k \rceil + T \cdot 2m \exp(-2N\delta_0^2) \right) \cdot \Delta_{\max}^{(T)}.
 \end{aligned} \tag{43}$$

The first inequality is obtained by applying the Hoeffding's Inequality (Lemma E.1) and union bound to the event  $\max_{i \in [m]} |\mu_i - \hat{\mu}_i| \geq \delta_0$ . Now we need to choose an optimal  $N$  that minimizes Eq.(43). By taking  $N = \max \left\{ 1, \frac{1}{2\delta_0^2} \ln \frac{4kmT\delta_0^2}{C} \right\} = \max \left\{ 1, \frac{2C^2m^2}{(\Delta_{\min}^{(T)})^2} \ln \left( \frac{kT(\Delta_{\min}^{(T)})^2}{C^2m} \right) \right\}$ , when  $\Delta_{\min}^{(T)} > 0$ , we can get the distribution-dependent bound

$$\text{Reg}_{\alpha,\beta}(T; \mu) \leq \frac{2C^2m^2n\Delta_{\max}^{(T)}}{k(\Delta_{\min}^{(T)})^2} \left( \max \left\{ \ln \left( \frac{kT(\Delta_{\min}^{(T)})^2}{C^2mn} \right), 0 \right\} + 1 \right) + \frac{n}{k} \Delta_{\max}^{(T)}, \tag{44}$$

Next, let us prove the distribution-independent bound. Let  $\mathcal{N}$  denote the event that  $|\hat{\mu}_i - \mu_i| \leq \sqrt{\frac{2 \ln T}{N}}$  for all  $i \in [m]$ . By the Hoeffding's Inequality and union bound, we have

$$\mathbb{P}\{\neg\mathcal{N}\} \leq m \cdot \frac{2}{T^4} \leq \frac{2}{T^3}. \tag{45}$$

When  $\mathcal{N}$  holds, with Eq.(42), we have

$$\Delta_{S^{(t)}}^{(t)} \leq 2Cm \cdot \sqrt{\frac{2 \ln T}{N}}, \tag{46}$$

and the regret in Eq.(39) can be written as

$$\begin{aligned}
 \text{Reg}_{\alpha,\beta}(T; \mu) &\leq \lceil nN/k \rceil \cdot n + \sum_{t=T-\lceil nN/k \rceil+1}^T \Delta_{S^{(t)}}^{(t)} \\
 &\leq \lceil nN/k \rceil \cdot n + O \left( T \cdot Cm \cdot \sqrt{\frac{\ln T}{N}} \right).
 \end{aligned} \tag{47}$$

We can choose  $N$  so as to (approximately) minimize the regret. For  $N = (Cmk)^{\frac{2}{3}} n^{-\frac{4}{3}} T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}}$ , we obtain:

$$\text{Reg}_{\alpha,\beta}(T; \mu) \leq O((Cmn)^{\frac{2}{3}} k^{-\frac{1}{3}} T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}}). \tag{48}$$



To complete the proof, we need to consider both  $\mathcal{N}$  and  $\neg\mathcal{N}$ . As shown in Eq.(45), the probability that  $\neg\mathcal{N}$  occurs is very small, and we have:

$$\begin{aligned} \text{Reg}_{\alpha,\beta}(T; \boldsymbol{\mu}) &= \mathbb{E}[\text{Reg}_{\alpha,\beta}(T; \boldsymbol{\mu}) \mid \mathcal{N}] \cdot \mathbb{P}\{\mathcal{N}\} + \mathbb{E}[\text{Reg}_{\alpha,\beta}(T; \boldsymbol{\mu}) \mid \neg\mathcal{N}] \cdot \mathbb{P}\{\neg\mathcal{N}\} \\ &\leq \mathbb{E}[\text{Reg}_{\alpha,\beta}(T; \boldsymbol{\mu}) \mid \mathcal{N}] + T \cdot n \cdot O(T^{-3}) \\ &\leq O((Cmn)^{\frac{2}{3}} k^{-\frac{1}{3}} T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}}). \end{aligned} \tag{49}$$

□

## F Proof of Theorem 6.1

*Proof.* As mentioned in Section 6, we need to introduce a virtual  $B$  seed node  $u_B$ , which connects to each existing node  $u$  with an unknown edge probability  $p(u_B, u)$  equal to the probability of  $u$  being selected as a  $B$  seed. By adding these virtual nodes and edges, we get a new graph  $G'$  with  $2n$  nodes and  $m + n$  edges. Since  $S_B$  is fixed under  $G'$ , we can follow the same steps in the proof of Theorem 3.1 to show the TPM condition holds under  $G'$ . Note that the maximum number of nodes that any one node can reach in  $G'$  is twice as that in the original graph  $G$ , so the new bounded smoothness coefficient  $C = 2\bar{C}$ . □

## G Additional Experiments

### G.1 Experiments for $A > B$ Tie-breaking Rule

When we consider  $A > B$  in bipartite graphs, we can trivially ignore  $S_B$  to choose  $S_A$  since the influence spread ends in one diffusion round, and OCIM becomes the online influence maximization problem without competition. We show such results in Figure 7. Note that the distribution of  $B$  no longer affects the performance of  $A$  when  $A > B$  and we only use one figure for the IM and RD distribution. For general graphs, we use the same DM dataset and parameter settings described in Sec. 7, and the only difference is that  $A$  now dominates  $B$ . We show the results in Figure 8. Overall, the results and the analysis for  $A > B$  are consistent with  $B > A$ .

### G.2 Experiments for OCIM-ETC

We show the frequentist/Bayesian regret results for the OCIM-ETC algorithm in Figure 9, Figure 10 and Figure 11. In Figure 9, we set exploration phase to be 250 rounds and the experiments show that we suffer linear regrets in both the exploration and the exploitation phase, meaning that the unknown parameters are under-explored. Thus we reset exploration to be 1500 and Figure 11 shows that OCIM-ETC now has constant regret in the exploitation phase. For DM dataset, since the node number and the edge number are less than Yahoo-Ad, we can see constant regrets after 1000 rounds of exploration in Figure 10. Compared with OCIM-OFU/OCIM-TS, OCIM-ETC requires more rounds to learn the unknown influence probabilities and has larger regrets than OCIM-OFU/OCIM-TS, but with sufficient exploration (which is much less than the theoretical requirements  $N = (\tilde{C}m)^{\frac{2}{3}}(nk)^{-\frac{1}{3}}T^{\frac{2}{3}}(\ln T)^{\frac{1}{3}}$  in Theorem 5.3) OCIM-ETC can yield constant regrets during the exploitation phase in our experiments.

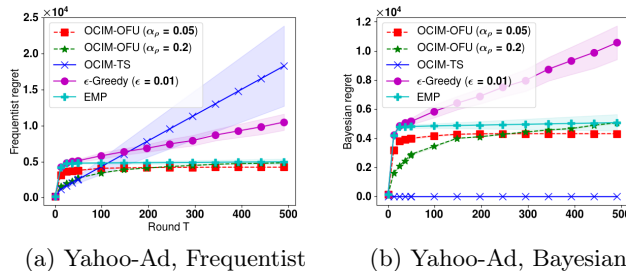


Figure 7: Frequentist/Bayesian regrets of different algorithms for the Yahoo-Ad graph when  $A > B$ .

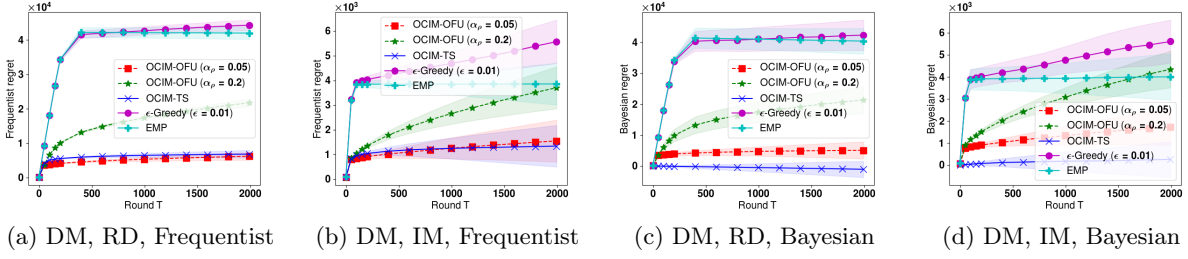
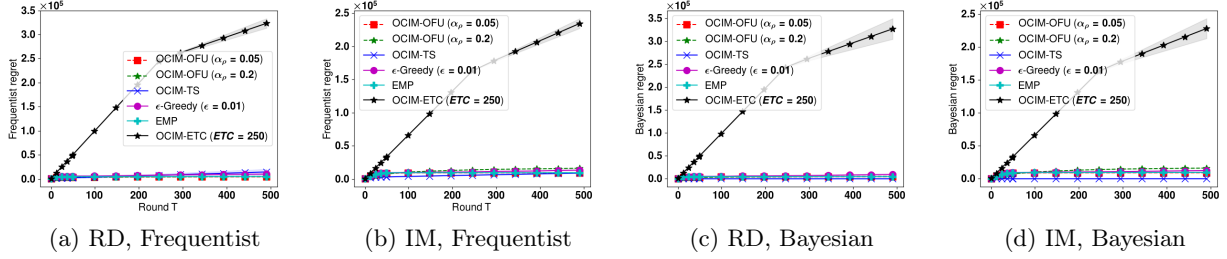

 Figure 8: Frequentist/Bayesian regrets of different algorithms for the general graph DM when  $A > B$ .


Figure 9: Frequentist/Bayesian regrets of OCIM-ETC for the Yahoo-Ad graph.

### G.3 Experiments for Probabilistic Seed Distribution

For the settings where the competitor has unknown fixed seed distribution, we first run the non-competitive influence maximization algorithm for  $S_B$  and get the best 5 seeds on Yahoo-Ad and the best 10 seeds on DM, respectively. We then consider the seed distribution of  $S_B$  as choosing each node from the best seeds with probability 0.5, i.e., the probability that choosing all best 5 seeds on Yahoo-Ad is  $0.5^5$  and the probability that choosing all best 10 seeds on DM is  $0.5^{10}$ . This seed distribution of  $S_B$  is unknown to our algorithms. In our experiments, we set  $|S_A| = 5$  for Yahoo-Ad and  $|S_A| = 10$  for DM, and assume  $B > A$ . Figure 12 shows that OCIM-OFU is still superior to EMP and  $\epsilon$ -Greedy for this setting with more complex competitor actions. We omit the results of OCIM-TS here as it requires the prior knowledge of the competitor's seed distribution. However, as long as the given prior does not differ much from the true prior, OCIM-TS will also achieve good regret results.

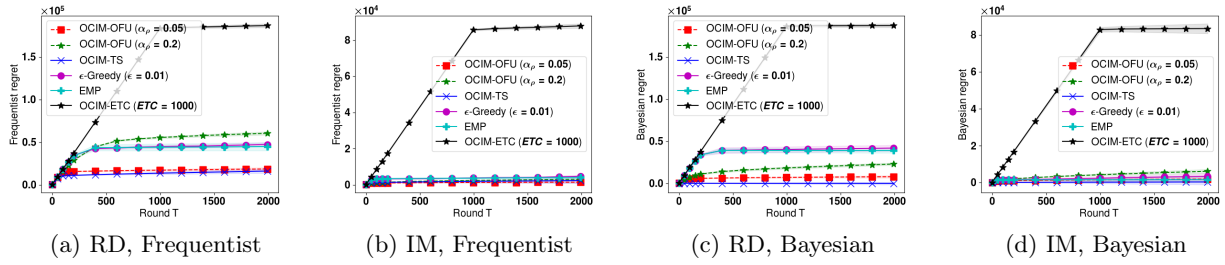


Figure 10: Frequentist/ Bayesian regrets of OCIM-ETC for the DM graph.

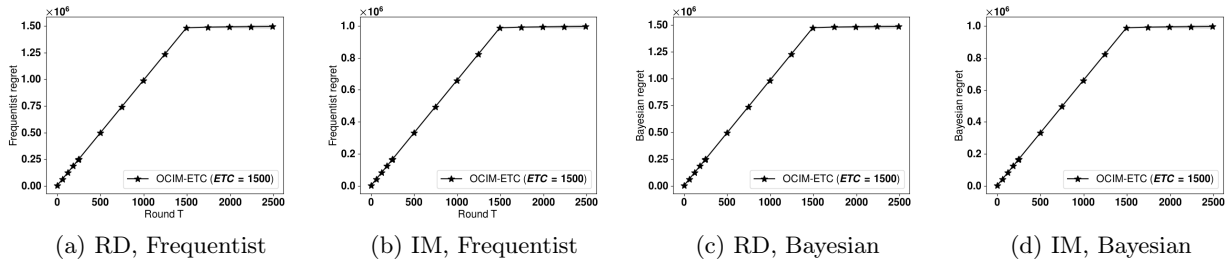


Figure 11: Frequentist/ Bayesian regrets of OCIM-ETC for the Yahoo-Ad graph with 1500 rounds of exploration.

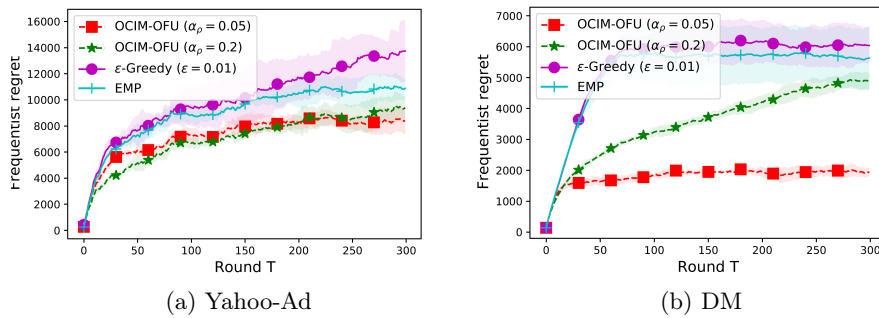


Figure 12: Frequentist regrets for Yahoo-Ad and DM with unknown fixed competitor's seed distribution.

## H Contextual combinatorial multi-armed bandit framework C<sup>2</sup>MAB-T

### H.1 General Framework

We propose a general framework of contextual combinatorial multi-armed bandit with probabilistically triggered arms (C<sup>2</sup>MAB-T), which is a contextual extension of CMAB-T in (Wang and Chen, 2017).

C<sup>2</sup>MAB-T is a learning game between a learning player and an environment. The environment consists of  $m$  random variables  $X_1, \dots, X_m$  called *base arms* following a joint distribution  $D$  over  $[0, 1]^m$ . Distribution  $D$  is chosen by the environment from a class of distributions  $\mathcal{D}$  before the game starts. The player knows  $\mathcal{D}$  but not the actual distribution  $D$  in advance. Different from that in CMAB-T, the environment in C<sup>2</sup>MAB-T also provides *contexts* for the learning agent, which will be discussed in detail later.

The learning process runs in discrete rounds. In round  $t$ , the environment first provides a context,  $\mathcal{S}^{(t)} \subseteq \mathcal{S}$ , to the player, where  $\mathcal{S}$  is the full action space and  $\mathcal{S}^{(t)}$  is a subset of it, representing the current action space in round  $t$ . The player then chooses an action  $S^{(t)} \in \mathcal{S}^{(t)}$  based on the feedback history from previous rounds. The environment also draws an independent sample  $X^{(t)} = (X_1^{(t)}, \dots, X_m^{(t)})$  from the joint distribution  $D$ . When action  $S^{(t)}$  is played on the environment outcome  $X^{(t)}$ , a random subset of arms  $\tau_t \in [m]$  are triggered, and the outcomes of  $X_i^{(t)}$  for all  $i \in \tau_t$  are observed as the feedback to the player.  $\tau_t$  may have additional randomness beyond the randomness of  $X^{(t)}$ . Let  $D_{\text{trig}}(S, X)$  denote a distribution of the triggered subset of  $[m]$  for a given action  $S$  and an environment outcome  $X$ . We assume  $\tau_t$  is drawn independently from  $D_{\text{trig}}(S^{(t)}, X^{(t)})$ . The player obtains a reward  $R(S^{(t)}, X^{(t)}, \tau_t)$  fully determined by  $S^{(t)}, X^{(t)}$  and  $\tau_t$ . A learning algorithm aims at selecting actions  $S^{(t)}$ 's over time based on the past feedback to accumulate as much reward as possible.

For each arm  $i$ , let  $\mu_i = \mathbb{E}_{X \sim D}[X_i]$ . Let  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  denote the expectation vector of arms. We assume that the expected reward  $\mathbb{E}[R(S, X, \tau)]$ , where the expectation is taken over  $X \sim D$  and  $\tau \sim D_{\text{trig}}(S, X)$ , is a function of action  $S$  and the expectation vector  $\boldsymbol{\mu}$  of the arms. Thus, we denote  $r_S(\boldsymbol{\mu}) := \mathbb{E}[R(S, X, \tau)]$ . We assume the outcomes of arms do not depend on whether they are triggered, i.e.,  $\mathbb{E}_{X \sim D, \tau \sim D_{\text{trig}}(S, X)}[X_i \mid i \in \tau] = \mathbb{E}_{X \sim D}[X_i]$ .

The performance of a learning algorithm  $\mathcal{A}$  is measured by its expected regret, which is the difference in expected cumulative reward between always playing the best action and playing actions selected by algorithm  $\mathcal{A}$ . Let  $\text{opt}^{(t)}(\boldsymbol{\mu}) = \sup_{S^{(t)} \in \mathcal{S}^{(t)}} r_{S^{(t)}}(\boldsymbol{\mu})$  denote the expected reward of the optimal action in round  $t$ . We assume that there exists an offline oracle  $\mathcal{O}$ , which takes context  $\mathcal{S}^{(t)}$  and  $\boldsymbol{\mu}$  as inputs and outputs an action  $S^{\mathcal{O},(t)}$  such that  $\Pr\{r_{S^{\mathcal{O},(t)}}(\boldsymbol{\mu}) \geq \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu})\} \geq \beta$ , where  $\alpha$  is the approximation ratio and  $\beta$  is the success probability. Instead of comparing with the exact optimal reward, we take the  $\alpha\beta$  fraction of it and use the following  $(\alpha, \beta)$ -approximation *frequentist regret* for  $T$  rounds:

$$\text{Reg}_{\alpha, \beta}^{\mathcal{A}}(T; \boldsymbol{\mu}) = \sum_{t=1}^T \alpha \cdot \beta \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - \sum_{t=1}^T r_{S^{\mathcal{A},(t)}}(\boldsymbol{\mu}), \quad (50)$$

where  $S^{\mathcal{A},(t)}$  is the action chosen by algorithm  $\mathcal{A}$  in round  $t$ .

Another way to measure the performance of the algorithm  $\mathcal{A}$  is using *Bayesian regret*. Denote the prior distribution of  $\boldsymbol{\mu}$  as  $\mathcal{Q}$ . When the prior  $\mathcal{Q}$  is given, the corresponding Bayesian regret is defined as:

$$\text{BayesReg}_{\alpha, \beta}^{\mathcal{A}}(T) = \mathbb{E}_{\boldsymbol{\mu} \sim \mathcal{Q}} \text{Reg}_{\alpha, \beta}^{\mathcal{A}}(T; \boldsymbol{\mu}). \quad (51)$$

Note that the contextual combinatorial bandit problem is also studied in (Chen et al., 2018; Qin et al., 2014). They consider the context features of all bases arms, which can affect their expected outcomes in each round, and assume the action space of super arms is a subset of  $[m]$ . However, we do not bond the context with base arms and consider the feasible set of super arms,  $\mathcal{S}^{(t)}$ , as the context, which is more flexible than a subset of  $[m]$ . Besides, we are the first to consider probabilistically triggered arms in the contextual combinatorial bandit problem.

### H.2 Monotonicity and Triggering Probability Modulated Condition

In order to guarantee the theoretical regret bounds, we consider two conditions given in (Wang and Chen, 2017). The first one is monotonicity, which is stated below.

**Condition 3.** (*Monotonicity*). We say that a C<sup>2</sup>MAB-T problem instance satisfies monotonicity, if for any action  $S$ , for any two expectation vectors  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  and  $\boldsymbol{\mu}' = (\mu'_1, \dots, \mu'_m)$ , we have  $r_S(\boldsymbol{\mu}) \leq r_S(\boldsymbol{\mu}')$  if  $\mu_i \leq \mu'_i$  for all  $i \in [m]$ .

---

**Algorithm 4** Contextual CUCB with offline oracle  $\mathcal{O}$ ,  $C^2$ -UCB
 

---

- 1: **Input:**  $m$ , Oracle  $\mathcal{O}$ .
  - 2: For each arm  $i \in [m]$ ,  $T_i \leftarrow 0$ . {maintain the total number of times arm  $i$  is played so far.}
  - 3: For each arm  $i \in [m]$ ,  $\hat{\mu}_i \leftarrow 1$ . {maintain the empirical mean of  $X_i$ .}
  - 4: **for**  $t = 1, 2, 3, \dots$  **do**
  - 5: For each arm  $i \in [m]$ ,  $\rho_i \leftarrow \sqrt{\frac{3 \ln t}{2T_i}}$ . {the confidence radius,  $\rho_i = +\infty$  if  $T_i = 0$ .}
  - 6: For each arm  $i \in [m]$ ,  $\bar{\mu}_i = \min\{\hat{\mu}_i + \rho_i, 1\}$ . {the upper confidence bound.}
  - 7: Obtain context  $\mathcal{S}^{(t)}$ .
  - 8:  $S^{(t)} \leftarrow \mathcal{O}(\mathcal{S}^{(t)}, \bar{\mu}_1, \bar{\mu}_2, \dots, \bar{\mu}_m)$ .
  - 9: Play action  $S^{(t)}$ , which triggers a set  $\tau \subseteq [m]$  of base arms with feedback  $X_i^{(t)}$ 's,  $i \in \tau$ .
  - 10: For every  $i \in \tau$  update  $T_i$  and  $\hat{\mu}_i$ :  $T_i = T_i + 1$ ,  $\hat{\mu}_i = \hat{\mu}_i + (X_i^{(t)} - \hat{\mu}_i)/T_i$ .
  - 11: **end for**
- 

The second condition is Triggering Probability Modulated (TPM) Bounded Smoothness. We use  $p_i^S(\boldsymbol{\mu})$  to denote the probability that the action  $S$  triggers arm  $i$  when the expectation vector is  $\boldsymbol{\mu}$ . The TPM condition in  $C^2$ MAB-T is given below.

**Condition 4.** (*1-Norm TPM bounded smoothness*). We say that a  $C^2$ MAB-T problem instance satisfies 1-norm TPM bounded smoothness, if there exists  $C \in \mathbb{R}^+$  (referred as the bounded smoothness coefficient) such that, for any two expectation vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ , and any action  $S$ , we have  $|r_S(\boldsymbol{\mu}) - r_S(\boldsymbol{\mu}')| \leq C \sum_{i \in [m]} p_i^S(\boldsymbol{\mu}) |\mu_i - \mu'_i|$ .

### H.3 Regret Bounds with Monotonicity

For the general  $C^2$ MAB-T problem that satisfies both monotonicity (Condition 3) and TPM bounded smoothness (Condition 4), we introduce a contextual version of the CUCB algorithm (Wang and Chen, 2017), which is described in Algorithm 4. Recall that  $\mathcal{S}^{(t)}$  is the action space in round  $t$ . We define the reward gap  $\Delta_S^{(t)} = \max(0, \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) - r_S(\boldsymbol{\mu}))$  for all actions  $S \in \mathcal{S}^{(t)}$ . For each arm  $i$ , we define  $\Delta_{\min}^{i,T} = \min_{t \in [T]} \inf_{S \in \mathcal{S}^{(t)}: p_i^S(\boldsymbol{\mu}) > 0, \Delta_S^{(t)} > 0} \Delta_S^{(t)}$  and  $\Delta_{\max}^{i,T} = \max_{t \in [T]} \sup_{S \in \mathcal{S}^{(t)}: p_i^S(\boldsymbol{\mu}) > 0, \Delta_S^{(t)} > 0} \Delta_S^{(t)}$ . If there is no action  $S$  such that  $p_i^S(\boldsymbol{\mu}) > 0$  and  $\Delta_S^{(t)} > 0$ , we define  $\Delta_{\min}^{i,T} = +\infty$  and  $\Delta_{\max}^{i,T} = 0$ . We define  $\Delta_{\min}^{(T)} = \min_{i \in [m]} \Delta_{\min}^{i,T}$  and  $\Delta_{\max}^{(T)} = \max_{i \in [m]} \Delta_{\max}^{i,T}$ . Let  $\tilde{\mathcal{S}} = \{i \in [m] \mid p_i^S(\boldsymbol{\mu}) > 0\}$  be the set of arms that can be triggered by  $S$ . We define  $K = \max_{S \in \mathcal{S}^{(t)}} |\tilde{\mathcal{S}}|$  as the largest number of arms could be triggered by a feasible action. We use  $[x]_0$  to denote  $\max\{[x], 0\}$ . Contextual CUCB ( $C^2$ -UCB) has the following regret bounds.

**Theorem H.1.** For the Contextual CUCB algorithm  $C^2$ -UCB (Algorithm 4) on an  $C^2$ MAB-T problem satisfying 1-norm TPM bounded smoothness (Condition 4) with bounded smoothness constant  $C$ , (1) if  $\Delta_{\min}^{(T)} > 0$ , we have a distribution-dependent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq \sum_{i \in [m]} \frac{576C^2 K \ln T}{\Delta_{\min}^{i,T}} + 4Cm + \sum_{i \in [m]} \left( \left\lceil \log_2 \frac{2CK}{\Delta_{\min}^{i,T}} \right\rceil + 2 \right) \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}, \quad (52)$$

and (2) we have a distribution-independent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq 12C\sqrt{mKT \ln T} + 2Cm + \left( \left\lceil \log_2 \frac{T}{18 \ln T} \right\rceil_0 + 2 \right) \cdot m \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}.$$

*Proof.* We first show that Lemma 5 in (Wang and Chen, 2017) still holds for Contextual CUCB algorithm in the  $C^2$ MAB-T problem. Let  $\mathcal{N}_t^s$  be the event that at the beginning of round  $t$ , for every arm  $i \in [m]$ ,  $|\hat{\mu}_{i,t} - \mu_i| \leq \rho_{i,t}$ . Let  $\mathcal{H}_t$  be the event that at round  $t$  oracle  $\mathcal{O}$  fails to output an  $\alpha$ -approximate solution. In Lemma 5 from (Wang and Chen, 2017), it assumes that  $\mathcal{N}_t^s$  and  $\neg \mathcal{H}_t$  hold, then we have

$$r_{S^{(t)}}(\bar{\boldsymbol{\mu}}_t) \geq \alpha \cdot \text{opt}^{(t)}(\bar{\boldsymbol{\mu}}_t) \geq \alpha \cdot \text{opt}^{(t)}(\boldsymbol{\mu}) = r_{S^{(t)}}(\boldsymbol{\mu}) + \Delta_{S^{(t)}}^{(t)}. \quad (53)$$

By the TPM condition, we have

$$\Delta_{S^{(t)}}^{(t)} \leq r_{S^{(t)}}(\bar{\boldsymbol{\mu}}_t) - r_{S^{(t)}}(\boldsymbol{\mu}) \leq C \sum_{i \in [m]} p_i^{S^{(t)}}(\boldsymbol{\mu}) |\bar{\mu}_{i,t} - \mu_i|, \quad (54)$$

---

**Algorithm 5** C<sup>2</sup>-TS with offline oracle  $\mathcal{O}$

---

- 1: **Input:**  $m$ , Prior  $\mathcal{Q}$ , Oracle  $\mathcal{O}$ .
  - 2: **Initialize** Posterior  $\mathcal{Q}_1 = \mathcal{Q}$
  - 3: **for**  $t = 1, 2, 3, \dots$  **do**
  - 4:   Draw a sample  $\boldsymbol{\mu}^{(t)}$  from  $\mathcal{Q}_t$ .
  - 5:   Obtain context  $\mathcal{S}^{(t)}$
  - 6:    $S^{(t)} \leftarrow \mathcal{O}(\mathcal{S}^{(t)}, \boldsymbol{\mu}^{(t)})$ .
  - 7:   Play action  $S^{(t)}$ , which triggers a set  $\tau \subseteq [m]$  of base arms with feedback  $X_i^{(t)}$ 's,  $i \in \tau$ .
  - 8:   Update posterior  $\mathcal{Q}_{t+1}$  using  $X_i^{(t)}$  for all  $i \in \tau$ .
  - 9: **end for**
- 

---

**Algorithm 6** C<sup>2</sup>-OFU with offline oracle  $\tilde{\mathcal{O}}$

---

- 1: **Input:**  $m$ , Oracle  $\tilde{\mathcal{O}}$ .
  - 2: For each arm  $i \in [m]$ ,  $T_i \leftarrow 0$ . {maintain the total number of times arm  $i$  is played so far.}
  - 3: For each arm  $i \in [m]$ ,  $\hat{\mu}_i \leftarrow 1$ . {maintain the empirical mean of  $X_i$ .}
  - 4: **for**  $t = 1, 2, 3, \dots$  **do**
  - 5:   For each arm  $i \in [m]$ ,  $\rho_i \leftarrow \sqrt{\frac{3 \ln t}{2T_i}}$ . {the confidence radius,  $\rho_i = +\infty$  if  $T_i = 0$ .}
  - 6:   For each arm  $i \in [m]$ ,  $c_i \leftarrow [(\hat{\mu}_i - \rho_i)^{0+}, (\hat{\mu}_i + \rho_i)^{1-}]$ . {the estimated range of  $\mu_i$ .}
  - 7:   Obtain context  $\mathcal{S}^{(t)}$ .
  - 8:    $S^{(t)} \leftarrow \tilde{\mathcal{O}}(\mathcal{S}^{(t)}, c_1, c_2, \dots, c_m)$ .
  - 9:   Play action  $S^{(t)}$ , which triggers a set  $\tau \subseteq [m]$  of base arms with feedback  $X_i^{(t)}$ 's,  $i \in \tau$ .
  - 10:   For every  $i \in \tau$  update  $T_i$  and  $\hat{\mu}_i$ :  $T_i = T_i + 1$ ,  $\hat{\mu}_i = \hat{\mu}_i + (X_i^{(t)} - \hat{\mu}_i)/T_i$ .
  - 11: **end for**
- 

which is in the same form of Eq.(10) in (Wang and Chen, 2017). Hence, we can follow the remaining proof of its Lemma 5. With Lemma 5, we can follow the proof of Lemma 6 in (Wang and Chen, 2017) to bound the regret when  $\Delta_{S^{(t)}}^{(t)} \geq M_{S^{(t)}}$ , where  $M_{S^{(t)}} = \max_{i \in \tilde{S}^{(t)}} M_i$  and  $M_i$  is a positive real number for each arm  $i$ . Finally, we take  $M_i = \Delta_{\min}^{i,T}$ . If  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}}$ , then  $\Delta_{S^{(t)}}^{(t)} = 0$ , since we have either  $\tilde{S}^{(t)} = \emptyset$  or  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}} \leq M_i$  for some  $i \in \tilde{S}^{(t)}$ . Thus, no regret is accumulated when  $\Delta_{S^{(t)}}^{(t)} < M_{S^{(t)}}$ . Following Eq.(17)-(22) in (Wang and Chen, 2017), we can derive the distribution-dependent and distribution-independent regret bounds shown in the theorem.  $\square$

#### H.4 Regret Bounds without Monotonicity

As discussed in Section 3, OCIM is an example of C<sup>2</sup>MAB-T that satisfies the TPM condition but not monotonicity. For the general C<sup>2</sup>MAB-T problem without monotonicity, we proposed two algorithms, C<sup>2</sup>-TS, C<sup>2</sup>-OFU, that can still achieve logarithmic Bayesian and frequentist regrets respectively. We also present C<sup>2</sup>-ETC that has a tradeoff between feedback requirement and regret bound.

C<sup>2</sup>-TS is described in Algorithm 5. Different from OCIM-TS, we input a general prior  $\mathcal{Q}$  (which depends on  $\mathcal{D}$  and might not be Beta distributions anymore) and update the posterior distribution  $\mathcal{Q}_t$  accordingly. With the same definitions in H.3 and  $\delta_{\max}^{(T)} = \max_{\boldsymbol{\mu}} \Delta_{\max}^{(T)}$ , it has the following Bayesian regret bound.

**Theorem H.2.** *For the C<sup>2</sup>-TS (Algorithm 5) on an C<sup>2</sup>MAB-T problem satisfying 1-norm TPM bounded smoothness (Condition 4) with bounded smoothness constant  $C$ , we have the Bayesian regret bound*

$$\text{BayesReg}_{\alpha,\beta}(T) \leq 12C\sqrt{mKT \ln T} + 2Cm + \left(\lceil \log_2 \frac{T}{18 \ln T} \rceil_0 + 4\right) \cdot m \cdot \frac{\pi^2}{6} \cdot \delta_{\max}^{(T)}, \quad (55)$$

C<sup>2</sup>-OFU is described in Algorithm 6. Similar to OCIM-OFU, it requires an offline oracle  $\tilde{\mathcal{O}}$  that takes the context  $\mathcal{S}^{(t)}$  and  $c_i$ 's (ranges of  $\mu_i$ 's) as inputs and outputs an approximate solution  $S^{(t)}$ . With such an oracle, C<sup>2</sup>-OFU has the following frequentist regret bounds.

**Theorem H.3.** *For the C<sup>2</sup>-OFU (Algorithm 6) on an C<sup>2</sup>MAB-T problem satisfying 1-norm TPM bounded*

---

**Algorithm 7** C<sup>2</sup>-ETC with offline oracle  $\mathcal{O}$ 


---

- 1: **Input:**  $m, k, N, T$ , Oracle  $\mathcal{O}$ .
  - 2: For each arm  $i$ ,  $T_i \leftarrow 0$ . {maintain the total number of times arm  $i$  is played so far.}
  - 3: For each arm  $i$ ,  $\hat{\mu}_i \leftarrow 0$ . {maintain the empirical mean of  $X_i$ .}
  - 4: **Exploration phase:**
  - 5: **for**  $t = 1, 2, 3, \dots, \lceil mN/k \rceil$  **do**
  - 6:   Obtain context  $\mathcal{S}^{(t)}$ .
  - 7:   Play action  $S^{(t)} \in \mathcal{S}^{(t)}$ , which contains  $k$  base arms that have not been chosen for  $N$  times.
  - 8:   Observe the feedback  $X_i^{(t)}$  for each base arm in  $S^{(t)}$ ,  $i \in \tau_{\text{direct}}$ .
  - 9:   For each arm  $i \in \tau_{\text{direct}}$  update  $T_i$  and  $\hat{\mu}_i$ :  $T_i = T_i + 1, \hat{\mu}_i = \hat{\mu}_i + (X_i^{(t)} - \hat{\mu}_i)/T_i$ .
  - 10: **end for**
  - 11: **Exploitation phase:**
  - 12: **for**  $t = \lceil mN/k \rceil + 1, \dots, T$  **do**
  - 13:   Obtain context  $\mathcal{S}^{(t)}$ .
  - 14:    $S^{(t)} \leftarrow \mathcal{O}(\mathcal{S}^{(t)}, \hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_m)$ .
  - 15:   Play action  $S^{(t)}$ .
  - 16: **end for**
- 

smoothness (Condition 4) with bounded smoothness constant  $C$ , (1) if  $\Delta_{\min}^{(T)} > 0$ , we have a distribution-dependent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq \sum_{i \in [m]} \frac{576C^2 K \ln T}{\Delta_{\min}^{i, T}} + 4Cm + \sum_{i \in [m]} \left( \left\lceil \log_2 \frac{2CK}{\Delta_{\min}^{i, T}} \right\rceil_0 + 2 \right) \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}, \quad (56)$$

and (2) we have a distribution-independent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq 12C\sqrt{mKT \ln T} + 2Cm + \left( \left\lceil \log_2 \frac{T}{18 \ln T} \right\rceil_0 + 2 \right) \cdot m \cdot \frac{\pi^2}{6} \cdot \Delta_{\max}^{(T)}.$$

Besides C<sup>2</sup>-TS and C<sup>2</sup>-OFU, we also provide a general explore-then-commit algorithm C<sup>2</sup>-ETC, as described in Algorithm 7. In the general setting,  $\tau_{\text{direct}}$  is defined as the set of base arms that is deterministically triggered by the action in question. C<sup>2</sup>-ETC is simple and only requires feedback from directly triggered arms, but it has a worse regret bound and requires the following condition besides Condition 4.

**Condition 5.** For some  $k \geq 1$ , given any context  $\mathcal{S}^{(t)}$  and any set  $S' \subseteq [m]$  with  $|S'| = k$ , there exists  $S \in \mathcal{S}^{(t)}$  such that  $p_i^S(\boldsymbol{\mu}) = 1$  for every  $i \in |S'|$ .

With such a condition, C<sup>2</sup>-ETC has the following frequentist regret bounds.

**Theorem H.4.** For the C<sup>2</sup>-ETC (Algorithm 7) on an C<sup>2</sup>MAB-T problem satisfying Condition 5 and 1-norm TPM bounded smoothness (Condition 4) with bounded smoothness constant  $C$ , (1) if  $\Delta_{\min}^{(T)} > 0$ , when  $N = \max \left\{ 1, \frac{2C^2 m^2}{(\Delta_{\min}^{(T)})^2} \ln \left( \frac{kT(\Delta_{\min}^{(T)})^2}{C^3 m} \right) \right\}$ , we have a distribution-dependent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq \frac{m}{k} \Delta_{\max}^{(T)} + \frac{2C^2 m^3 \Delta_{\max}^{(T)}}{k(\Delta_{\min}^{(T)})^2} \left( \max \left\{ \ln \left( \frac{kT(\Delta_{\min}^{(T)})^2}{C^2 m^2} \right), 0 \right\} + 1 \right) \quad (57)$$

and (2) when  $N = (Ck)^{\frac{2}{3}} m^{-\frac{2}{3}} T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}}$ , we have a distribution-independent bound

$$\text{Reg}_{\alpha, \beta}(T; \boldsymbol{\mu}) \leq O(C^{\frac{2}{3}} m^{\frac{4}{3}} k^{-\frac{1}{3}} T^{\frac{2}{3}} (\ln T)^{\frac{1}{3}}). \quad (58)$$

The proofs of Theorem H.2, H.3 and H.4 generally follow the same steps in Appendix B, C and E.