

# Heterogeneous Multi-agent Multi-armed Bandits on Stochastic Block Models

MENGFAN XU, Department of Mechanical and Industrial Engineering, University of Massachusetts Amherst, USA

LIREN SHAN, Toyota Technological Institute at Chicago, USA

FATEMEH GHAFFARI, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, USA

XUCHUANG WANG, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, USA

XUTONG LIU, School of Computer Science, Carnegie Mellon University, USA

MOHAMMAD HAJIESMAILI, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, USA

We study a novel heterogeneous multi-agent multi-armed bandit problem with a cluster structure induced by stochastic block models, influencing not only graph topology, but also reward heterogeneity. Specifically, agents are distributed on random graphs based on stochastic block models - a generalized Erdos-Renyi model with heterogeneous edge probabilities: agents are grouped into clusters (known or unknown); edge probabilities for agents within the same cluster differ from those across clusters. In addition, the cluster structure in stochastic block model also determines our heterogeneous rewards. Rewards distributions of the same arm vary across agents in different clusters but remain consistent within a cluster, unifying homogeneous and heterogeneous settings and varying degree of heterogeneity, and rewards are independent samples from these distributions. The objective is to minimize system-wide regret across all agents. To address this, we propose a novel algorithm applicable to both known and unknown cluster settings. The algorithm combines an averaging-based consensus approach with a newly introduced information aggregation and weighting technique, resulting in a UCB-type strategy. It accounts for graph randomness, leverages both intra-cluster (homogeneous) and inter-cluster (heterogeneous) information from rewards and graphs, and incorporates cluster detection for unknown cluster settings. We derive optimal instance-dependent regret upper bounds of order  $\log T$  under sub-Gaussian rewards. Importantly, our regret bounds capture the degree of heterogeneity in the system (an additional layer of complexity), exhibit smaller constants, scale better for large systems, and impose significantly relaxed assumptions on edge probabilities. In contrast, prior works have not accounted for this refined problem complexity, rely on more stringent assumptions, and exhibit limited scalability.

## 1 Introduction

Multi-armed Bandit (MAB) [6, 7] is an online learning framework in which, during a sequential game, an agent, or decision maker, selects one arm from multiple arms, pulls the arm, and receives the reward observation of the pulled arm from an unknown environment at each time step. The

---

Authors' Contact Information: Mengfan Xu, Department of Mechanical and Industrial Engineering, University of Massachusetts Amherst, Amherst, MA, USA, mengfanxu@umass.edu; Liren Shan, lirensan@ttic.edu, Toyota Technological Institute at Chicago, Chicago, IL, USA; Fatemeh Ghaffari, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, Amherst, MA, USA, fghaffari@umass.edu; Xuchuang Wang, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, Amherst, MA, USA, xuchuangw@gmail.com; Xutong Liu, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA, xutongl@andrew.cmu.edu; Mohammad Hajiesmaili, Manning College of Information & Computer Sciences, University of Massachusetts Amherst, Amherst, MA, USA, hajiesmaili@cs.umass.edu.

objective is to maximize the cumulative received reward by identifying the best arm, a task also known as regret minimization when compared to the ideal case of knowing in advance which arm is the best. Recently, with the rapid development of real-world networks, multi-agent systems have become a major focus, motivating the study of Multi-agent Multi-armed Bandit (MA-MAB) [10, 27, 32, 43, 46, 53, 57, 59, 67]. In this context, multiple agents exist within a system, each with its own arm set, playing a bandit game while communicating with others to exchange bandit information. It is well known that MAB is classified into stochastic (where reward observations come from a time-invariant distribution with reward mean values) [6] and adversarial (where reward observations are arbitrary) [7]. Here, we focus exclusively on the stochastic setting, consistent with the majority of existing work on MA-MAB. For simplicity, we refer to stochastic MA-MAB as MA-MAB for the remainder of this paper.

Depending on the application domain and environment properties, previous work studies several variants of MA-MAB settings. Among all variants, a widely studied one is a cooperative setting, where agents share the same arm set and aim to maximize the overall system’s objective. Depending on how the rewards are generated for different agents, the cooperative MA-MAB can be further categorized into homogeneous and heterogeneous. In a homogeneous setting, the reward mean value of the same arm across different agents is identical. This implies that the locally optimal arm (with respect to an agent’s own reward distribution) is also the globally optimal arm (with respect to the average reward mean values of the same arm across all agents). This scenario has been extensively studied [4, 13, 30–32, 38, 42, 48, 51, 53, 65]. However, it is common that in real applications, the agents’ rewards are heterogeneous. For example, retail companies in different regions may have varying product return rates due to population heterogeneity. To address this, a line of research has focused on the heterogeneous setting [58, 64, 67], where the globally optimal arm can differ from the locally optimal arm. Nonetheless, these prior works assume a fully heterogeneous setting, treating all agents as distinct. It is important to note that, in practice, the setting is not necessarily fully heterogeneous; instead, different degrees of heterogeneity can exist within a heterogeneous setting. This concept, however, has not been well-defined or thoroughly studied, presenting a clear research gap.

Another central aspect of MA-MAB is how agents communicate. In a decentralized setting, agents are distributed on a graph (as vertices connected by edges) and can only communicate if an edge exists between them. In a sequential regime, while time-invariant graphs have been well studied, the advancement of several applications, e.g., wireless IoT networks, motivates the study of time-varying graphs, which introduces additional challenges. Random graphs [23] have been a promising approach to model time-varying graphs in MA-MAB [20, 58] and other areas [17, 40], where the graph is randomly drawn from a distribution by sampling each edge based on a probability, akin to the reward generation process. However, current research in MA-MAB [58] is largely limited to Erdos-Renyi models [23], where the edge probability is homogeneous across all agents. Additionally, the graph generation process is assumed to be independent of the reward generation process, which may not always reflect practical scenarios. These limitations in the studied random graph models for MA-MAB highlight another important research gap.

Notably, a broad family of random graphs is formulated as Stochastic Block Models (SBM) [1, 2, 18], where agents are grouped into clusters, and the edge probability for agents within the same cluster differs from that for agents in different clusters. The existence of cluster structures, which generalize the Erdos-Renyi models as commonly used in MA-MAB, surprisingly but naturally provides a framework for defining the degree of heterogeneity through graph topology. Consequently, considering SBM in the context of MA-MAB holds significant potential for addressing

the aforementioned research gaps. However, SBM has not yet been explored in the cooperative learning context, particularly in MA-MAB, which motivates our work herein.

In this paper, we study the following research problem: *Can we address the heterogeneous multi-agent multi-armed bandit problem on Stochastic Block Models to bridge the following two gaps: 1) varying degrees of heterogeneity and 2) more general random graphs linked to reward dynamics?*

### 1.1 Main Contributions

We provide an affirmative answer to the above question through our main contributions, summarized as follows. First, we formulate a general heterogeneous multi-agent multi-armed bandit problem, where agents are grouped into clusters inspired by the stochastic block model. This cluster structure determines both graph topology and reward similarity. Specifically, the cluster structure introduces heterogeneity in edge probabilities and reward distributions across clusters, while maintaining homogeneity in edge probabilities and reward distributions within clusters, thereby linking reward dynamics with graph dynamics. This framework thus extends random graphs (with homogeneous edge probabilities) from Erdos-Renyi models to general stochastic block models (with heterogeneous edge probabilities). It also incorporates both homogeneity and heterogeneity in reward distributions into a unified setting, characterizing the degree of heterogeneity, which reflects an additional layer of complexity in heterogeneous settings. Existing work on either homogeneous or heterogeneous rewards can be viewed as special cases of this framework, demonstrating its consistency and generalization capability. A more detailed comparison with existing models is provided in Section 2.

Secondly, we propose a learning algorithm tailored to this new formulation that effectively leverages homogeneity to reduce the sample complexity of reward estimations and heterogeneity to efficiently learn the global objective. Specifically, we propose an algorithm, namely UCB-SBM, which consists of a burn-in period to collect local information and a learning period to leverage historical data to improve the arm-pulling strategy. During the burn-in period, the agents randomly pull arms to obtain reward estimators for each arm at an individual level. Then, during the learning period, the agents use a UCB-type strategy to pull the arm with the highest UCB index based on exchanged information and a new weighting technique. They communicate newly designed information to other agents, perform information updates based on newly proposed rules, and novelly run cluster detection using rewards as side information when the cluster structure is unknown.

Our algorithmically technical novelty is as follows. Compared to the most relevant work [58], our UCB-type strategy employs a newly constructed global estimator that integrates both homogeneity (inter-cluster estimators) and heterogeneity (intra-cluster estimators). Additionally, our information transmission involves sending cluster-level estimators instead of individual-level ones, thereby 1) achieving noise reduction in the estimators, and 2) relaxing the assumption on the edge probability, as it requires only one representative in the cluster to exchange cluster-level information, rather than requiring all agent pairs to communicate. The information update process is significantly different in that we construct estimators at the cluster level using a newly proposed weighted sum/average approach, and compute the global information based on the new weight technique over the cluster-level estimators, resulting in three layers of estimators: local, cluster, and global. In contrast, [58] considers only local and global layers. Lastly, the incorporation of a cluster detection method enables us to infer general unknown cluster structures and thus to leverage the cluster structure to design algorithms, which is completely omitted in [58].

Thirdly, we establish precise instance-dependent regret upper bounds for the UCB-SBM algorithm, which are of order  $O(\log T)$ . Additionally, if we examine the coefficient of the regret bounds more

closely, our regret bound accurately captures the relationship between the regret upper bounds and the newly defined degree of heterogeneity,  $C/M$ , which is the ratio between the number of clusters  $C$  and the number of agents  $M$ . Specifically, the regret bound depends linearly on  $C/M$ , reflecting the problem’s complexity. Moreover, this implies that the total regret depends on  $C \leq M$  instead of  $M$ , scaling significantly better with the number of agents  $M$  in large-scale systems. In contrast, the existing algorithm for heterogeneous rewards results in a regret upper bound of order  $M^2$  [58], as it neglects possible cluster structures. This bound may become unmanageable when the number of agents is comparable to the time horizon  $T$ .

Fourthly, our results do not rely on strong assumptions about edge probabilities, making them more broadly applicable compared to prior work. More specifically, the lower bound on the edge probability in our case can be at most  $\frac{e}{e-1} \cdot \frac{C^2}{M^2} \cdot \frac{(C-l-1)!}{(C-2)!} < 1$  for  $1 \leq l \leq C-1$ , while the lower bound on the edge probability in [58] approaches 1 as  $T \rightarrow \infty$ . Furthermore, this lower bound in [58] increases more rapidly with  $M$ , as it depends on  $M$  whereas ours only depends on  $C \leq M$ , and a larger lower bound implies more stringent assumptions on the problem setting as  $M$  grows. Overcoming the limited scalability of problem setting and obtaining a lower bound on the edge probability that is strictly bounded away from 1 is a highly non-trivial yet impactful contribution. Fifth, our results apply to scenarios with both known and unknown cluster assignments, aligning with existing work on Stochastic Block Models. A comprehensive summary of the theoretical results is presented in Table 1. Lastly, through experiments, we demonstrate that our algorithm achieves much lower actual regret (beyond regret bounds), with an improvement of at least 68.69%, highlighting its superior practical effectiveness.

Additionally, we make an independent contribution as follows. The regret bound under the new framework should reflect the degree of heterogeneity, which is essentially highlighted as an open problem in [58]. In that work, the authors numerically observe a dependency of regret on the level of heterogeneity in the problem setting, noting that regret increases monotonically with the level of heterogeneity. This suggests the potential to achieve smaller regret when the degree of heterogeneity is low. However, they do not formally define or analyze this dependency theoretically, leaving a research gap. Moreover, the results in [58] heavily rely on an assumption about the lower bound on the edge probability in the Erdos-Renyi model, which can become quite restrictive when  $T$  is sufficiently large, thereby limiting its practical applicability. How to relax these stringent assumptions remains unexplored, necessitating the development of new methods and analyses—another research gap that our work seeks to address.

*Paper Organization.* The paper is presented as follows. We provide a comparison of our work with existing studies in Section 2. In Section 3, we introduce the notations and formulate the research question. In Section 4, we provide the motivation for the formulation by highlighting some important real-world applications of the problem setting. Section 5 begins by characterizing the framework through a simple case involving a single cluster, where the formulation reduces to homogeneous rewards on random graphs. Subsequently, in Section 6, we extend the framework to the main case involving multiple known clusters, presenting the proposed algorithm and its analysis (with improved regret bounds) under milder assumptions compared to existing work. In Section 7, we illustrate how the algorithm can be adapted to scenarios with multiple unknown clusters. Section 8 demonstrates the numerical performance of the proposed algorithm. Last, we conclude the paper and suggest future research directions in Section 9.

Table 1. Summary of the main results for UCB-SBM.

Cluster	Thm.	Asm. 1*	Asm. 2*	Worst-case <sup>§</sup>	Coef. <sup>‡</sup>
known; $C = 1$	1	$p \in (0, 1]$	$q = p$	N/A	$O(\frac{1}{M})^{\frac{1}{2}}$
known/unknown <sup>¶</sup>	3	$p = 1$	$q > 1 - (\frac{1}{2} - \frac{1}{2}\sqrt{1 - (\frac{\delta}{8CT})^2})^{C^2/M^2}$	1	$O(\frac{C}{M})$
known/unknown <sup>¶</sup>	5	$p = 1$	$q > 1 - (1 - \min\{(\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{8CT})^2}), 1 - \delta(C-1)/8CT\})^{C^2/M^2}$	1	$O(\frac{C}{M})$
known/unknown <sup>¶</sup>	8	$p > \max\{\frac{(C_M - I - 1)!}{(C_M - 2)!} (1 - \frac{\delta(C_M - 1)}{8C_M T}), \frac{(C_M - I - 1)!}{(C_M - 2)!} (\frac{3}{4})^{\frac{1}{2}}\}$	$q > \frac{e}{e-1} \frac{C^2}{M^2} \max\{\frac{(C-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-1)!}{(C-2)!} (\frac{3}{4})^{\frac{1}{2}}\}$	$\frac{e}{e-1} \frac{C^2}{M^2} \frac{(C-1)!}{(C-2)!}$	$O(\frac{C}{M})$

This table assumes reward distributions is sub-Gaussian and the regret bound is of order  $\log T$ . \* Asm. 1 refers to the assumption on the edge probability  $p$  for agents within the same cluster and Asm. 2 refers to the assumption on the edge probability  $q$  for agents belonging to different clusters. † We observe that going from 1 to  $C$  clusters, the regret grows linearly with  $C$ , which represents the dependency between the regret bound and the degree of heterogeneity. ‡ The worst case scenario refers to the case when  $T \rightarrow \infty$ , i.e. the upper bound on the Asm. 2. The value in the existing work [58] is also 1. § We impose additional assumptions on  $p - q$  in the unknown cluster case.

## 2 Related Work

Our proposed model differs significantly from existing work on multi-agent multi-armed bandits. Specifically, we outline these differences in comparison to the existing lines of research on: 1) homogeneous cooperative multi-agent multi-armed bandits, 2) heterogeneous cooperative multi-agent multi-armed bandits, 3) multi-agent multi-armed bandits with clusters of agents, and 4) multi-agent multi-armed bandits with time-varying graphs.

*Homogeneous Cooperative Multi-agent Multi-armed Bandit.* There has been extensive work on cooperative multi-agent bandits, with most studies assuming homogeneous rewards, where the reward distribution for the same arm is identical across all agents [4, 13, 30–32, 38, 42, 48, 51, 53, 65]. In contrast, our model incorporates heterogeneous reward distributions for agents in different clusters, while maintaining homogeneous reward distributions for agents within the same cluster. Notably, when there is only one cluster, our model reduces to the case of the homogeneous reward.

*Heterogeneous Cooperative Multi-agent Multi-armed Bandit.* Although some studies have explored heterogeneous rewards [58, 63, 64], they treat rewards as entirely heterogeneous, without considering the possibility of a framework that bridges heterogeneity and homogeneity, along with the associated problem complexity. In our work, we define and systematically characterize the degree of heterogeneity using the cluster structure. Our model also aligns with existing heterogeneous cases when each agent belongs to a different cluster, and thus, there are  $M$  clusters. In summary, our work bridges the gap between homogeneous and heterogeneous rewards by integrating both paradigms and fully characterizing every possible degree of heterogeneity.

*Clusters by Stochastic Block Models.* One key to our framework is considering cluster structure motivated by the Stochastic Block Model (SBM), which was previously a separate line of work. SBM, introduced by [26], is known as a foundational framework for modeling community (referred to as clusters herein) structures in networks. It has been extensively studied for cluster detection, with detailed analyses providing exact phase transitions and efficient algorithms for different recovery settings [1]. However, it has not yet been coupled with MA-MAB to model and leverage the agent structure to additionally decide on the reward distribution, and thus bridge the gaps. Besides modeling, we also consider scenarios where the cluster assignment is unknown, inspired by the Contextual Stochastic Block Model (CSBM) proposed by [18], which generalizes SBM by incorporating side information—namely node covariates—that depend on cluster assignments. Building on that, recent work provide algorithms to leverage both graph structure and contextual attributes to enhance cluster detection and recovery [3, 12, 18, 19]. Specifically, [3] rigorously study the case where node covariates are generated from a Gaussian Mixture Model (GMM) and propose an algorithm for two-cluster networks. More generally, [12] develop an iterative clustering algorithm and derive the exact recovery threshold for multiple balanced clusters. Notably, none of them consider reward information as side information unique to MA-MAB.

*Multi-agent Multi-armed Bandit with Clusters of Agents.* Another related line of research incorporates cluster structures into multi-armed bandits, commonly referred to as the online clustering of bandits (CLUB) [8, 11, 24, 25, 29, 33–37, 39, 41, 44, 45, 52, 56, 60]. These studies assume that agents can be grouped into clusters, with each group sharing similar reward distributions for each arm, a concept that aligns with our setting. However, there are three significant differences between CLUB and our work. First, while CLUB primarily focuses on contextual bandit scenarios and provides instance-independent regret bounds, our work addresses the canonical multi-agent MAB setting and establishes finer-grained, instance-dependent regret bounds. Second, most CLUB approaches assume a central server within a star-shaped communication graph [11, 41, 60]. To our knowledge, only Korda et al. [29] consider peer-to-peer networks, where agents can communicate with any other agent using a gossip protocol. In contrast, our work involves a more realistic and challenging scenario: communication is constrained by a random communication graph modeled by a stochastic block structure. In this setting, only agents connected by an edge can exchange information, significantly increasing the problem’s complexity. Finally, CLUB aims to identify the optimal arm for each individual agent, whereas our work focuses on finding a *globally* optimal arm across all agents. Consequently, our framework requires each agent not only to learn its own reward distribution but also to infer the reward distributions of other agents. This added complexity is particularly demanding under the constraints of a random communication graph.

*Multi-agent Multi-armed Bandit with Graphs.* Recently, the study of multi-agent bandit problems, where agents are distributed on a graph that constrains their communication, has gained significant attention. Most existing works focus on time-invariant graphs, where the graph remains constant over time [28, 55, 65–67]. However, there is an emerging need to address time-varying graphs, which capture more general scenarios where the graph changes over time, motivated by wireless ad-hoc networks in IoT [47]. It is worth noting that existing work on time-varying graphs either considers Erdos-Renyi graphs with homogeneous edge probabilities [58] or focuses on connected graphs [64], without exploring heterogeneous edge probabilities or the relationship between graph dynamics and reward dynamics. Notably, we are the first to bridge this gap by introducing stochastic block models, which are more general than Erdos-Renyi graphs, and by relating graph topology to reward heterogeneity through a cluster structure. Furthermore, existing work on Erdos-Renyi graphs [58] imposes strong assumptions on edge probabilities, which may be highly impractical. We address this limitation by leveraging cluster information and significantly relaxing these assumptions.

### 3 Problem Formulation

In this section, we introduce the notations and formally present the problem formulation. We start by introducing the notations. Consistent with the basic MAB setting, we consider  $K$  arms, labeled as  $1, 2, \dots, K$ . Let us denote each time step as  $1 \leq t \leq T$ , where  $T$  is the length of the time horizon. Let us denote  $M$  as the number of agents in this multi-agent setting. These agents are distributed on a time-dependent graph  $G_t$  represented by vertex set  $V = \{1, 2, \dots, M\}$  and edge set  $E_t$ . We use  $X_{i,j}^t$  to denote whether an agent pair  $(i, j) \in E_t$ . We use  $\mathcal{N}_m(t)$  to denote the neighbor set of agent  $m$  at time  $t$ , where agent  $j$  is called to be in the neighbor set  $\mathcal{N}_m(t)$  if and only if there is an edge between them, i.e.,  $(m, j) \in E_t$ . The graph  $G_t$  is independent and identically distributed samples from the stochastic block models that extend the Erdos-Renyi model in [58] as described below.

**DEFINITION 1 (STOCHASTIC BLOCK MODELS).** *We consider a stochastic block model, where the set of agents (vertices) with a cluster structure, each agent  $1 \leq i \leq M$ , belongs to a cluster  $c_i \in 1, 2, \dots, C$ . Additionally, there exists an **unknown** probability matrix  $\{p(m, n)\}_{\substack{1 \leq n \leq C \\ 1 \leq m \leq C}}$  associated with the clusters, where  $p(m, n)$  represents the probability of having an edge between an agent pair  $(i, j)$ , where agent*

$i \in m$  and agent  $j \in n$ . Notably,  $p(m, n) \neq p(m, m)$  for  $m \neq n$ , meaning the probability of having an edge between two agents within the same cluster differs from the probability of having an edge between two agents across different clusters, implying heterogeneous random graphs. Then we sample  $X_{i,j}^t$  based on  $\{p(m, n)\}_{1 \leq m \leq C}^{1 \leq n \leq C}$ , for  $\forall i, j \in V$ , and  $E_t = \{(i, j) : X_{i,j}^t = 1, \forall i, j \in V\}$ .

It is worth noting that when  $C = 1$ ,  $X_{i,j}^t$  are sampled according to a Bernoulli distribution with a uniform edge probability, which precisely aligns with the definition of Erdos-Renyi models.

Besides the graph setting, we consider the reward setting characterized by clusters based on stochastic block models. Let  $\mu_k^i$  denote the reward mean value of arm  $1 \leq k \leq K$  for agent  $1 \leq i \leq M$ . Notably, the reward mean values for the same arm are identical for agents within the same cluster, i.e.,  $\mu_k^i = \mu_k^m \doteq \mu_k^{c_m}$ , for agent  $m$  and  $j$  that meet  $c_m = c_i$  while differing for agents in different clusters. This framework effectively bridges the gap between homogeneous and heterogeneous MA-MAB settings. Moreover, we propose a new definition of the degree of heterogeneity as follows.

**DEFINITION 2 (DEGREE OF HETEROGENEITY).** We define  $h_{M,C} = C/M = 1/c_M$  where  $c_M$  represents the average number of agents in one cluster, which is scale-invariant and bounded by 1, i.e.,  $0 < h_{M,C} \leq 1$ .

We argue the rationality as follows. This metric quantifies the variety in the reward/edge distributions across clusters relative to the total number of agents. When  $C = M$ , it is fully heterogeneous, aligning with [58], and when  $C = 1$ , it is fully homogeneous, consistent with [53].

The reward of arm  $k$  at agent  $i$  at time step  $t$ , denoted as  $\{r_k^i(t)\}_{k,i,t}$ , follows a  $\sigma^2$ -sub-Gaussian distribution with a time-invariant mean value  $\{\mu_k^i\}_{k,i}$ .

**REMARK.** While we assume sub-Gaussian reward distributions for illustrative purposes, we highlight that the formulation and results established herein can be extended to sub-exponential cases through straightforward analysis, as this does not require changes to the communication or information update mechanisms. We omit the details and focus on the sub-Gaussian case in this work.

Let  $a_i^t$  represent the arm selected by agent  $i$  at time  $t$ , and let  $n_{i,k}(t)$  denote the number of pulls of arm  $k$  at agent  $i$  up to time  $t$ . We consider a cooperative setting where the goal of all agents is to select the globally optimal arm, defined as  $k^* = \arg \max_k \sum_{i=1}^M \mu_k^i$ . The objective of the system is to maximize the pulls of the globally optimal arm, thereby minimizing regret, which is defined with respect to the globally optimal arm as follows. Formally, the regret and total regret are given by

$$R_T = \frac{1}{M} \sum_{k=1}^K \sum_{i=1}^M \Delta_k n_{i,k}(T), \quad R_T^M = M \cdot R_T = \sum_{k=1}^K \sum_{i=1}^M \Delta_k n_{i,k}(T).$$

respectively, where  $\Delta_k = (\sum_{i=1}^M \mu_{k^*}^i - \sum_{i=1}^M \mu_k^i)/M$ .

## 4 Real-world Applications

In this section, we motivate our problem formulation, which bridges existing gaps by considering agents on stochastic block models, through a range of real-world applications. Here stochastic block models capture agents with specific probabilities of being connected, where agents can observe edges but not the underlying edge probabilities—a scenario that often reflects real-world conditions.

*Collaborative Content Placement in Content Delivery Networks (CDNs).* Online content delivery in content delivery networks (CDNs) is a critical component of modern network applications, including video streaming, web browsing, and software distribution [14, 16, 61]. Unlike traditional architectures that rely on a single central server, CDNs distribute and cache content across multiple edge servers, allowing end users to retrieve data from the nearest server. This distributed architecture significantly reduces latency and enhances the reliability of content delivery. A key challenge in CDNs lies in dynamically placing content across thousands of edge servers to optimize user service-level objectives (e.g., latency, packet loss) and quality of experience (QoE). In this context, each edge server can be modeled as an agent, with its arms representing candidate content placement policies. The reward for each arm corresponds to the number of successfully delivered and precached contents, which ultimately reduce users' loading time. Given the heterogeneity in user preferences and network conditions, edge servers may form clusters where only agents within the same cluster share similar rewards, modeled by a stochastic block model. Furthermore, the large number of edge servers and candidate policies necessitates collaboration among servers to learn optimal policies. However, due to communication bandwidth constraints, servers can only communicate randomly, governed by a random graph. Our framework effectively models this problem, enabling the identification of the *global* optimal content placement policy that maximizes the reward across all edge servers.

*Collaborations in Social Networks.* Examples include scientific collaboration networks of biologists and physicists, where an edge represents a collaboration, defined as co-authorship of one or more scientific articles during the study period, and a cluster refers to working in the same main research area [15]. The collaboration network is highly dynamic, as collaborations change over time and are modeled by time-varying graphs. These scientists may select the most important research topic from a few options, referred to as arms, with the reward of an arm being the impact of the research topic (scientists working in the same area will have the same impact). In the context of collaboration networks of movie actors [22], an edge represents appearing in the same movie, which again varies in movies released at different times and is therefore time-varying, and a cluster refers to club membership. These actors may choose the best club activity among several options, referred to as arms, with the reward of an arm being the engagement in the activity (actors in the same club will have the same reward distribution). Similarly, in a network of directors of Fortune 1000 companies [9], an edge between two directors indicates that they served on the same board that can change over time as the board committee itself may change, and a cluster again refers to club membership. Here, the arms are the start-up candidates for investment, and the reward of an arm is the return on investment (ROI) (directors in the same club will have the same reward distribution).

*Protein-to-Protein in AI-enabled Biology.* In the context of protein-to-protein biology research, AI-enabled proteins embedded in a patient act as agents/nodes, and the physical connections or interactions between proteins are represented as edges within a cell, namely a protein-to-protein interactions network [5]. It is worth noting that these interactions change over time, resulting in time-varying graphs. AI-enabled Proteins functioning similarly in the protein-protein interaction network belong to the same cluster. The task of the proteins is to transport different nutrients (arms) in the body, and the rewards of the arms are the corresponding health conditions of a patient, e.g., blood pressure or blood sugar levels, resulting from different nutrients (the reward distribution for proteins within the same cluster is identical).

*Recommendation in E-commerce with Filtering.* In e-commerce, filtering has been an effective approach where customers utilize others' information to make decisions [21, 49]. In collaborative filtering, customers are represented as nodes/agents, an edge between two customers indicates



similar behaviors, and the most similar customers (those with the highest degrees) form clusters. The arms represent product candidates, and the reward of an arm corresponds to the experience with the product (hence, the reward distribution for the same arm is identical for agents within the same cluster). In item-to-item collaborative filtering, product producers are represented as nodes/agents, an edge signifies similar properties, and the most similar products (again, those with the highest degrees) form clusters. The arms are the warehouse options for the products, and the reward corresponds to the quality of the product after being stored in the warehouse (the reward distributions for the same warehouse are identical for products within the same cluster).

## 5 Warm-up: Single Cluster (Homogeneous Clients)

This section studies the single-cluster multi-agent MAB, focusing on in-cluster learning and serving as a didactic warm-up for the multi-cluster scenarios discussed in later sections. Here, all clients belong to the same cluster and share a homogeneous reward environment. Although there is existing work on homogeneous multi-agent multi-armed bandits [50, 54], our model introduces a time-varying random graph  $G_t$ , which has not yet been studied. In this homogeneous setting, all clients have the same reward distribution for each arm. Consequently, observations from different clients can be combined to improve the estimation of an arm's reward distribution, leading to a more efficient exploration-exploitation trade-off compared to the single-agent case. We first present a simple UCB algorithm in Section 5.1, followed by its regret upper bound in Section 5.2.

### 5.1 Algorithm

Since clients are homogeneous, we propose a simple cooperative Upper Confidence Bounds (UCB) algorithm. Over the whole learning process, every client  $m$  maintains the UCB index for each arm  $k$  as follows,  $u_{k,t}^{(m)} := \hat{\mu}_{k,t}^{(m)} + \sqrt{\log t / \tilde{N}_{k,t}^{(m)}}$ , where  $\tilde{N}_{k,t}^{(m)} = N_{k,t}^{(m)} + \sum_{m' \in \mathcal{M} \setminus \{m\}} N_{k, \tau_t^{(m \leftrightarrow m')}}^{(m')}$  is the total number of observations of arm  $k$  that client  $m$  collects, including its own  $N_{k,t}^{(m)}$  local observations and the  $N_{k, \tau_t^{(m \leftrightarrow m')}}^{(m')}$  observations collected from its neighbors  $m' \in \mathcal{M} \setminus \{m\}$  at the latest communication time slot  $\tau_t^{(m \leftrightarrow m')}$  between these two clients  $m$  and  $m'$  on or before time slot  $t$ . The empirical mean  $\hat{\mu}_{k,t}^{(m)}$  is also the average of all  $\tilde{N}_{k,t}^{(m)}$  observations of arm  $k$  that client  $m$  collects. The clients pull the arm with the highest UCB index, i.e.  $a_m^t = \operatorname{argmax}_k u_{k,t}^m$ , and receive the reward.

### 5.2 Regret Analysis

**THEOREM 1.** *Executing the above algorithm leads to  $\mathbb{E}[R_T] \leq O\left(\sum_{k \neq k^*} \frac{\log T}{M \Delta_k} + \frac{K}{p^{M^2}}\right)$ . (1)*

**PROOF SKETCH.** The full proof is in Appendix E; the main intuition is as follows. Fix a suboptimal arm  $k$ . After the total number of observations for this arm  $k$  exceeds the sample complexity threshold, in expectation, it takes  $\frac{1}{p^{M^2}}$  time slots for all agents to get the information of this arm  $k$ . After that, no more regret will be incurred on this arm  $k$ . As a result, the proof first makes an assumption to reduce the problem to a standard cooperative UCB for homogeneous agents residing on a complete graph with communication delays. Then, we show that this assumption can be fulfilled in the single cluster scenario.  $\square$

The leading  $O(\log T)$  term of the total regret  $R_T^M = M \cdot R_T$  by multiplying  $M$  and (1) is independent of the number of agents  $M$ , highlighting the advantage of multi-agent cooperation. Meanwhile, it is worth noting that in heterogeneous setting in [58], the upper bound of the total regret  $R_T^M$  is of order  $O(M^2 \log T)$  (though it is not tight as illustrated in Section 6), rather than  $O(\log T)$  which

emphasizes the regret reduction by our analysis in homogeneous settings. The second term of (1), however, has a dependence of  $p^{-M}$  where  $p = p(m, m)$ . It suggests the importance of the edge generation probability  $p$ , which will be thoroughly addressed in Section 6 and Section 7.

## 6 Heterogeneous - Multiple Known Clusters

In this section, we consider the general model where the agents are distributed on stochastic block models with multiple clusters, capturing the dependency between reward and graph dynamics. Here, we assume that the cluster information  $\{c_i\}_{i=1}^M$  is known to the agents, and we generalize the results to a more practical setting where the cluster information is unknown in Section 7. We would like to highlight that the probability matrix of the model is always unknown to the agents. The algorithm is presented in Section 6.1, followed by the corresponding regret analyses in Section 6.2.

### 6.1 Algorithm

The newly proposed algorithm is presented as follows and consists of two stages. In the first stage, all agents pull arms one by one without communication to accumulate local information, referred to as the burn-in period. In the second stage, agents use intelligent strategies (based on Upper Confidence Bounds) to pull arms and communicate with one another following the graph structure to collect global information, referred to as the learning period. The corresponding algorithms are provided as Algorithm 2 (see Appendix A) and Algorithm 1, collectively referred to as UCB-SBM (Upper Confidence Bounds for Stochastic Block Models).

We note that there is no difference between the algorithm in the burn-in period herein and that in [58], and thus we show the pseudo code in Appendix A, except that we do not need  $\tau_1$ . The reason is that there is no intelligence during this stage. Specifically at  $t$ , each agent  $m$  pull each arm  $a_m^t = t \bmod K$  one by one and update the average reward as local reward estimators  $\tilde{\mu}_i^m(t)$ . It also updates the edge frequency  $P_t(m, j) = ((t-1)Pt-1(m, j) + X_{m, j}^t)/t$  for each  $j \in V$ , and communicates  $\tilde{\mu}_i^m(t)$  to agent  $j \in \mathcal{N}_m(t)$ . Then at the end of the burn-in period, it outputs the initial global estimator  $\tilde{\mu}_i^m(L+1)$ , which is the weighted average of  $\tilde{\mu}_i^m(L)$  where weights are  $P_t(m, j)$ .

Subsequently, we proceed to the learning stage using either Rule 1 or Rule 2, which define how agents update and aggregate information. Rule 1 is consistent with [58], as it does not consider the cluster structure, whereas Rule 2 is newly proposed and leverages the cluster information. The pseudo-code is presented in Algorithm 1, which includes several stages in the order outlined below.

*Arm selection.* During this stage, the agents use a UCB-based criterion to decide which arm to pull. More specifically, if there is no arm  $i$  such that  $N_{m,i}(t) \leq \tilde{N}_{m,i}(t) - K$ , where  $N$  and  $\tilde{N}$  represent the in-cluster and across-cluster number of pulls, respectively, then agent  $m$  pulls  $a_m^t = \arg \max_i \tilde{\mu}_{m,i}(t) + F(m, i, t)$ , where  $\tilde{\mu}_{m,i}(t)$  is the network-wide estimator for arm  $i$  of agent  $m$  and  $F(m, i, t) = \sqrt{\frac{C_1 \ln t}{N_{m,i}(t)}}$  ( $C_1$  is specified later) quantifies the uncertainty in  $\tilde{\mu}_{m,i}(t)$ . Otherwise, the agents randomly pull an arm by specifying  $a_m^t = t \bmod K$ .

*Transmission.* The agents communicate with their neighbors and integrate information from other agents. Specifically, each agent sends its own information and receives information from agents in its time-dependent neighborhood. The information includes sample counts and reward estimators, covering local, cluster, and global levels, denoted as  $r_i^j(t), N_{j,i}(t), \tilde{N}_{j,i}(t), \tilde{\mu}_i^j(t), \tilde{\mu}_i^j(t)$  defined below.

*Information update.* With such information, agent  $m$  updates estimators as in **Rule 1** or **Rule 2**.

**Rule 1:**  $t_{m,j} = \max_{s \geq \tau_1} \{(m, j) \in E_s\}$  and 0 if such an  $s$  does not exist (2)

$$N_{m,i}(t+1) = n_{m,i}(t+1) = n_{m,i}(t) + \mathbb{1}_{a_m^t=i}, \tilde{N}_{m,i}(t+1) = \max\{N_{m,i}(t+1), \tilde{N}_{j,i}(t), j \in \mathcal{N}_m(t)\}$$

**Algorithm 1:** UCB-SBM: Learning period

---

```

1 Initialization: For each agent  $m$  and arm  $i \in \{1, 2, \dots, K\}$ , we have  $\tilde{\mu}_i^m(L+1)$ ,
    $\tilde{N}_{m,i}(L+1), N_{m,i}(L+1) = n_{m,i}(L)$ ;  $\tau = 1$ , all other values at  $L+1$  are initialized as 0;
2 for  $t = L+1, L+2, \dots, T$  do
3   for each agent  $m$  do // UCB
4     if there is no arm  $i$  such that  $N_{m,i}(t) \leq \tilde{N}_{m,i}(t) - K$  then
5        $a_m^t = \arg \max_i \tilde{\mu}_{m,i}(t) + F(m, i, t)$ 
6     else
7       Randomly sample an arm  $a_m^t$ .
8     end
9     Pull arm  $a_m^t$  and receive reward  $r_{m,a_m^t}(t)$ ;
10  end
11  The environment samples  $G_t = (V, E_t)$  based on SBM; // Env
12  Each agent  $m$  sends  $r_i^m(t), N_{j,i}(t), \tilde{\mu}_i^m(t), \tilde{N}_i^m(t)$  to each agent in  $\mathcal{N}_m(t)$ ;
13  Each agent  $m$  receives  $r_i^j(t), N_{j,i}(t), \tilde{\mu}_i^j(t), \tilde{N}_i^j(t)$  from all agents  $j \in \mathcal{N}_m(t)$  and stores them
   as  $\hat{\mu}_{i,j}^m(t), \hat{N}_{i,j}^m(t), \hat{\mu}_{i,j}^m(t), \hat{N}_{i,j}^m(t)$ ; // Transmission
14  for each agent  $m$  do
15    for  $i = 1, \dots, K$  do
16      Update  $P_t$  for  $1 \leq j \leq M$  by  $P_t(m, j) = (t-1)P_{t-1}(m, j) + X_{m,j}^t/t$ ;
17      Update  $P'_t$  for  $1 \leq j \leq M$  by  $P'_t(m, j) = 1$  if  $P_t(m, j) > 0$  and 0 o.w.;
18      if  $t \bmod \tau = 0$  then
19        Update  $n_{m,i}(t), N_{m,i}(t), \tilde{N}_{m,i}(t)$  and  $\tilde{\mu}_i^m(t)$  based on Rule 1 or Rule 2
20      else
21        Update  $n_{m,i}(t)$  as  $n_{m,i}(t) = n_{m,i}(t) + 1_{a_m^t=i}$ 
22      end
23    end
24  end
25 end

```

---

$$\bar{\mu}_i^m(t+1) = (\bar{\mu}_i^m(t) \cdot n_{m,i}(t) + r_{m,i}(t) \cdot \mathbb{1}_{a_m^t=i}) / n_{m,i}(t+1), P'_t(m, j) = (M-1)/M^2 \text{ if } P_t(m, j) > 0 \text{ and } 0 \text{ otherwise}$$

$$\tilde{\mu}_i^m(t+1) = \sum_{j=1}^M P'_t(m, j) \hat{\mu}_{i,j}^m(t_{m,j}) + d_{m,t} \sum_j \hat{\mu}_{i,j}^m(t_{m,j}) \text{ with } d_{m,t} = (1 - \sum_{j=1}^M P'_t(m, j)) / M$$

**Rule 2:**  $t_{m,j} = \max_{s \geq \tau_1} \{(m, j) \in E_s\}$  and 0 if such an  $s$  does not exist

Local and Cluster sample counts:  $n_{m,i}(t+1) = n_{m,i}(t) + \mathbb{1}_{a_m^t=i}$ ,  $N_{m,i}(t) = N_{c_m,i}(t) = \sum_j n_{j,i}(t_{m,j})$

Global sample counts:  $\tilde{N}_{m,i}(t+1) = \max\{N_i^m(t), \tilde{N}_i^j(t), (m, j) \in E_t\}$

Local estimator:  $\bar{\mu}_i^m(t+1) = (\bar{\mu}_i^m(t) \cdot n_{m,i}(t) + r_{m,i}(t) \cdot \mathbb{1}_{a_m^t=i}) / n_{m,i}(t+1)$

Cluster estimator (local):  $\hat{\mu}_i^{c_m}(t+1) = \hat{\mu}_i^m(t+1) = (\sum_{j \in c_m} \tilde{\mu}_i^j(t_{m,j})) / |c_m|$

Cluster estimator (network-wise):  $\tilde{\mu}_i^{c_m}(t+1) = \tilde{\mu}_i^m(t+1) = (\sum_{j \in c_m} \tilde{\mu}_i^j(t_{m,j})) / |c_m|$

$P_t(c_m, c_j) = (\sum_{s \leq t, p \in c_m, q \in c_j} \mathbb{1}_{(p,q) \in E_s}) / t$ ,  $P'_t(m, j) = (M-1)/M^2$  if  $P_t(c_m, c_j) > 0$  and 0 otherwise (3)

Global estimator:  $\tilde{\mu}_i^m(t+1) =$

$$\sum_{j=1}^M P'_t(m, j) \hat{\mu}_i^j(t_{m,j}) + d_{m,t} \sum_{j \in c_m} \hat{\mu}_i^j(t_{m,j}) + d_{m,t} \sum_{j \notin c_m} \hat{\mu}_i^j(t_{m,j}) \text{ with } d_{m,t} = (1 - \sum_{j=1}^M P'_t(m, j)) / M$$

The difference between Rule 1 and Rule 2 is that Rule 1 is the same as in [58] and leverage no cluster information, whereas Rule 2 is newly proposed herein. Rule 2 considers the stochastic block model structure (agents within a cluster aggregate  $n$  and  $\bar{\mu}$  to obtain  $N$  and  $\hat{\mu}$ ), communicates at the cluster level ( $\hat{\mu}$  instead of  $\bar{\mu}$ ), and utilizes the cluster information to improve the estimators  $\tilde{\mu}$  (with 3 sources: local, cluster, and global information). As shown in Section 6.2, they result in different regret bounds, with Rule 2 achieving a smaller regret and requiring less stringent assumptions.

## 6.2 Regret Analyses

Next, we prove the effectiveness of the proposed algorithm through analyzing the theoretical regret induced by the algorithm. For illustration purposes, let us assume a balanced model where the number of agents in each cluster is the same for all clusters, i.e.,  $|c_i| = \frac{M}{C} \doteq c_M$ . We highlight that the case of imbalanced clusters (with respect to cluster size) can be addressed in our analyses by using the smallest number of agents in a single cluster,  $\min_{1 \leq i \leq M} |c_i|$ , as the universal cluster size.

As a starting point, we consider the regret of Algorithm 2 with Rule 1, which aligns with existing work on heterogeneous rewards without characterizing the cluster structure. A straightforward result is presented below, which is a by-product of Theorem 2 in [58], but with potentially different edge probabilities for different agent pairs, and it reads as follows.

**THEOREM 2.** *Let us assume that  $\min_{m,n} p(m, n) \geq (\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{MT})^{\frac{2}{M-1}}})$ . For every  $0 < \epsilon < 1$  and  $0 < \delta < \frac{1}{2} + \frac{1}{4} \sqrt{1 - (\frac{\epsilon}{MT})^{\frac{2}{M-1}}}$ , the regret of Algorithm 1 with Rule 1 is upper bounded by with probability  $1 - 7\epsilon$ ,  $E[R_T | A_{\epsilon, \delta}] \leq L + \sum_{i \neq i^*} \Delta_i (\max\{[\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\}) + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K$  where  $A_{\epsilon, \delta} = A_2 \cap A_3$  with  $A_2 = \{\exists t_0, \forall t \geq L, \forall j, \forall m, t + 1 - \min_j t_{m,j} \leq t_0 \leq c_0 \min_i n_{l,i}(t + 1)\}$  and  $A_3 = \{\forall t \geq L, G_t \text{ is connected}\}$ , the length of the burn-in period is explicitly  $L = \max\{\ln T / 2\epsilon / 2\delta^2, 4K \log_2 T / c_0\}$ ,  $c_0 = c_0(K, \min_{i \neq i^*} \Delta_i, M, \epsilon, \delta)$ , and the instance-dependent constant  $C_1 = 8\sigma^2 C = \max\{4^{(M+2)}(1 - \frac{1-c_0}{2(M+2)})^2 / 3M(1-c_0), (M+2)(1 + 4Md_{m,t}^2)\}$ .*

**PROOF SKETCH.** The complete proof is provided in Appendix E; the main proof logic is as follows. The proof of Theorem 1 parallels that of Theorem 2 in [58] for the Erdos-Renyi graph, except that the edge probability is agent-dependent in this case. Interestingly, we find that as long as the minimal edge probability satisfies the condition on the edge probability in the Erdos-Renyi model as specified in [58], the entire proof remains valid. The key observation is that the original analysis relies solely on the lower bound of the edge probability in the Erdos-Renyi model.  $\square$

While the above regret bound depends logarithmically on the time horizon  $T$ , there are two limitations: 1) the assumption on  $\min_{m,n} p_{m,n}$  is stringent and may not always hold in practice, and 2) the total regret ( $M \cdot R_T$ ) depends linearly on  $M$ , which may not scale well in large-scale systems. This is because the analysis does not leverage the homogeneity within clusters, leading to high sample complexity for the reward estimators. To address these limitations, we next present an approach that exploits the cluster structure using Rule 2 in Algorithm 2, which takes advantage of the homogeneity within clusters induced by stochastic block models.

Intuitively, agents using Rule 2 first aggregate the rewards and sample counts of agents within the same cluster, and then communicate these at a cluster level, meaning they only share cluster-wide

information. The aggregation reduces sample complexity because the variance of the averaged estimator is smaller than that of a single agent's estimator. Consequently, this can potentially lead to smaller regret in terms of  $M$ . Additionally, cluster-level communication reduces the need for pairwise (every agent pair) communication. As long as there exists an agent in one cluster and another agent in a different cluster with an edge between them, the two clusters can communicate, rather than requiring every agent in one cluster to be connected to every agent in the other cluster. This approach reduces the connectivity requirements of the graph and relaxes the assumption on  $\min_{m,n} p_{m,n}$ , as a larger  $\min_{m,n} p_{m,n}$  always implies better connectivity (in the high probability sense).

Formally, we consider communication at a cluster level by defining the subgraph generated by the clusters as  $G_t^C$  as follows.

**DEFINITION 3.** A sub-graph  $G_t^C$  is represented by the vertex set  $\{1, 2, \dots, C\}$  of clusters and the edge set  $E_t^C$ , where the pair of clusters  $x$  and  $y$ , namely  $(x, y)$ , belongs to  $E_t^C$  if and only if there exists an agent  $i \in x$  and an agent  $j \in y$  such that  $(i, j) \in E_t$ .

It holds true that the sub-graph  $G_t^C$  is much denser compared to the original graph  $G_t$  as it has a higher probability of having an edge (cluster pair) and thus a lower requirement on the graph topology of  $G_t$ . First, we consider the case where the graph induced by the clusters,  $G_t^C$ , is a connected graph and the edge probability within one cluster is 1, and derive a better regret bound with relaxed assumptions. For illustration purposes, it is natural to assume that the edge probability within the same cluster is 1, and we relax this assumption later (Theorem 8) in this section. The formal statement is summarized as Theorem 3, which reads as follows.

**THEOREM 3.** Let us assume that  $p(m, m) = 1$  for any  $1 \leq m \leq C$ . Let us further assume that  $\min_{m,n} p(m, n) \geq 1 - (\frac{1}{2} - \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon E[R_T | A'_{\epsilon,\delta}] \leq L + \sum_{i \neq i^*} \Delta_i (\max\{\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\}) + 2\pi^2/3P(A_{\epsilon,\delta}) + K^2 + (2M - 1)K = O(C \log T/M)$  where  $A'_{\epsilon,\delta} = A_2 \cap A'_3$ ,  $A_2 = \{\exists t_0, \forall t \geq L, \forall j, \forall m, t + 1 - \min_j t_{m,j} \leq t_0 \leq c_0 \min_i n_{i,i}(t + 1)\}$  and  $A'_3 = \{\forall t \geq L, G_t^C \text{ is connected}\}$ , the length of the burn-in period is explicitly  $L = C/M \max\{\ln \frac{T}{2\epsilon}/2\delta^2, 4K \log_2 T/c_0\}$ ,  $c_0 = c_0(K, \min_{i \neq i^*} \Delta_i, M, \epsilon, \delta)$ , and  $C_1 = \max\{4(M+2)(1 - \frac{1-c_0}{2(M+2)})^2/3M(1-c_0), (M+2)(1 + 4Md_{m,t}^2)/M\}$ .

**PROOF SKETCH.** The full proof is deferred to Appendix E; we present the main logic herein. We note that by Lemma 4, the edge probability  $p(c_m, c_n)$  of the sub-graph  $G_t^C$  is  $1 - (1 - p(m, n))^{M^2/C^2}$ . In other words, as long as  $p(c_m, c_n)$  meets the condition of Theorem 2, i.e.,  $\min_{m,n} p(c_m, c_n) \geq (\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})$ , we achieve the same regret bound, where everything is with respect to the sub-graph instead of the original graph. Subsequently, we derive that the condition is equivalent to  $\min_{m,n} p(m, n) \geq 1 - (\frac{1}{2} - \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$ . Hence, the regret bound in Theorem 2 holds.  $\square$

**LEMMA 4.** For any pair of vertices  $c_m, c_n$  in the sub-graph  $G_t^C$ , the probability that  $c_m$  and  $c_n$  is connected in  $G_t^C$  is  $p(c_m, c_n) = 1 - (1 - p(m, n))^{M^2/C^2}$ .

*Discussion on the total regret.* We novelly derive a regret bound on  $R_T$  that depends on the degree of heterogeneity  $h_{M,C} = \frac{C}{M}$ , reflecting the problem complexity related to the stochastic block model and reward heterogeneity we consider, and highlighting the comprehensiveness of our

regret bound. It is worth noting that our result also resolves an open problem identified in [58], where the authors numerically claimed a dependency of the regret on heterogeneity without formally defining or analyzing it. The advantage of having this explicit dependency compared to the aforementioned established results is as follows. When  $C = 1$ , it is consistent with Theorem 1 in Section 5, further supporting our claim. We note that both the total regret bound ( $M^2 \cdot R_T$ ) in [58] and our result in Theorem 2 depend on  $M^2$ , while the one in Theorem 3 depends linearly on  $C \leq M$ . This demonstrates an improvement over the existing result in [58] when  $C = M$  (no homogeneity) and further implies a significant reduction in the regret bound when  $C < M$ . This improvement is particularly significant in large-scale systems where  $C \ll M$  (a common scenario in real-world applications, e.g., people in different regions where all individuals are agents and their clusters are defined by the regions). It highlights the practical importance of our proposed algorithm and provides practitioners with more effective tools.

*Discussion on the assumption.* Notably, the lower bound on  $\min_{m,n} p_{m,n}$  is  $1 - (\frac{1}{2} - \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$ . When  $C = M$ , i.e. in a fully heterogeneous setting, this lower bound is the same as the lower bound in Theorem 2, i.e. [57], implying consistency. When  $C < M$ , this term is smaller by noting that  $\frac{C^2}{M^2} < 1$ , and thus by the monotone property of the function  $1 - (\frac{1}{2} - \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$  we have  $1 - (\frac{1}{2} - \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2} < 1 - (\frac{1}{2} - \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}) = (\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}) < (\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{MT})^{\frac{2}{M-1}}})$  which suggests that our assumption is less stringent compared to Theorem 2 herein and the original statement of Theorem 2 in [58].

Moreover, we next show that the lower bound on  $p(m, n)$  can be further reduced by modifying the proof, purely from an analytical perspective. This modification also applies to the setting in [58], thereby providing an improvement to the result therein as well, as part of our contribution. The formal statement reads as follows.

**THEOREM 5.** *Let us assume that  $p(m, m) = 1$  for any  $1 \leq m \leq C$ . Let us further assume that  $\min_{m,n} p(m, n) \geq 1 - (1 - \min\{(\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}), 1 - \delta(C-1)/8CT\})^{C^2/M^2}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon E[R_T | A'_{\epsilon, \delta}] \leq L + \sum_{i \neq i^*} (\max(\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)) + 2\pi^2/3P(A_{\epsilon, \delta}) + K^2 + (2M - 1)K) = O(C \log T/M)$  where  $A'_{\epsilon, \delta}$ ,  $L$ ,  $c_0$ ,  $C_1$  are specified in Theorem 3.*

**PROOF SKETCH.** The complete proof is in Appendix E; we introduce the key logic here. The proof mostly follows from the proof of Theorem 3, with the only exception being that the concentration inequality used to prove the following proposition, which in part guarantees the regret bounds by ensuring communication effectiveness given random graphs, considers both Chernoff's Bound (leading to the lower bound  $(\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})$ ) and Chebyshev's inequality (resulting in  $1 - \frac{\delta(C-1)}{8CT}$ ). In contrast, Theorem 3 only considers Chernoff's Bound. The proposition reads as follows, characterizing the probability of event  $A'_3$ . Assume the edge probability  $p(m, n)$  where  $c_m \neq c_n$ ) meets the condition  $1 \geq p(m, n) \geq \min\{(\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}), 1 - \frac{\delta(C-1)}{8CT}\}$ , where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , event  $A'_3$  holds.  $\square$

The assumption on  $\min_{m,n} p(m, n)$  originates from the requirement to establish that  $G_t^C$  is connected with high probability. We observe that, in practice, it may not always be feasible or necessary to have a connected graph  $G_t^C$  at every time step. This assumption has been notoriously hard to

overcome in MA-MAB. The situation becomes even more challenging in our context and in existing work related to random graphs, as the assumption essentially implies that the lower bound on the edge probability in Theorem 3 and existing work [58] can approach 1 when  $T$  is large enough. This heavily constrains the applicability of the established results.

However, this is no longer a concern in what follows, addressed by our new techniques, which represent a significant advancement in this line of work on multi-agent systems with random graphs, and thus highlight our contributions. Surprisingly, we find that as long as the subgraph  $G_t^C$  is  $l$ -periodically connected, based on the following definition, the above regret bound holds, further relaxing the assumption on  $\min_{m,n} p(m, n)$ , which is shown to be strictly bounded away from 1. We introduce the definitions related to  $l$ -periodically connected graphs and present the formal statement below.

**DEFINITION 4 (COMPOSITION OF GRAPHS).** *Let us assume graphs  $G^1, G^2, \dots, G^l$  have the same vertex set  $V$  and possibly different edge set  $E_1, E_2, \dots, E_l$ . Then the composition of these  $l$  graphs  $G^1, G^2, \dots, G^l$ , known as  $G = G_1 \otimes G_2 \otimes G_3 \otimes \dots \otimes G_l$ , is uniquely defined by vertex set  $V$  and edge set  $E$  where  $(i, j) \in E$  if and only if there exist vertex  $v_2, v_3, v_{l-1}$  such that  $(i, v_2) \in E_1, (v_1, v_2) \in E_2, (v_2, v_3) \in E_3, \dots, (v_{l-1}, j) \in E_l$ .*

**DEFINITION 5 ( $l$ -PERIODICALLY CONNECTED [64]).** *A sequence of time-dependent sub-graphs  $G_t^C$  is said to be  $l$ -periodically connected in the sense that the composition of any  $l \geq 1$  consecutive sub-graphs  $G_{t_1}^C, G_{t_1+1}^C, G_{t_1+2}^C, \dots, G_{t_1+l-1}^C$  is a connected graph, formally expressed as  $G = G_{t_1}^C \otimes G_{t_1+1}^C \otimes G_{t_1+2}^C \otimes \dots \otimes G_{t_1+l-1}^C$  is connected.*

It is worth noting that, in our context,  $l$  is at most  $C - 1$  since there are  $C$  vertices. Also, by definition,  $l$  is a positive integer. Next, we demonstrate that any  $l$  within this range specifies a lower bound on  $\min_{m,n} p(m, n)$ , and when this lower bound holds for  $\min_{m,n} p(m, n)$ , the same regret bound as in Theorem 5 also applies, but under much less stringent assumptions.

First, based on the  $l$ -periodical connectivity, we have the following lemma hold which characterizes the relationship between  $p(m, n)$  and  $p(c_m, c_n)$  and is much stronger compared to Lemma 4.

**LEMMA 6.** *For any pair of vertices  $c_m, c_n$  in the sub-graph  $G_t^C$ , the probability that  $c_m$  and  $c_n$  is connected in  $G_t^C$  is  $p(c_m, c_n) \geq (1 - \frac{1}{e}) \min\{1, \frac{M^2}{C^2} \cdot p(m, n)\}$ .*

The formal statement on regret is as follows. Here we assume  $l \geq 2$ , as  $l = 1$  implies connectivity.

**THEOREM 7.** *Let us assume that  $p(m, m) = 1$  for any  $1 \leq m \leq C$ . Given any  $C \geq 4, 2 \leq l \leq C - 1$ , let us further assume that  $\min_{m,n} p(m, n) \geq \frac{e}{e-1} \cdot \frac{C^2}{M^2} \cdot \max\{\frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-l-1)!}{(C-2)!} (\frac{3}{4})^l\}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon E[R_T | A'_{\epsilon, \delta}] \leq L + \sum_{i \neq i^*} \Delta_i (\max\{\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\} + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K) = O(\frac{C \log T}{M})$  where  $A'_{\epsilon, \delta} = A_2 \cap A'_3, A_2 = \{\exists t_0, \forall t \geq L, \forall j, \forall m, t + 1 - \min_j t_{m,j} \leq t_0 \leq c_0 \min_l n_{l,i}(t + 1)\}$  and  $A'_3 = \{\forall t \geq L, G_t^C \text{ is } l\text{-periodically connected}\}$ ,  $L, c_0, C_1$  are specified in Theorem 3.*

**PROOF SKETCH.** We refer to Appendix E for the detailed proof; the intuitive logic is as follows. The relaxation is obtained by proving the following: for any  $m, i, t > L$ , if  $n_{m,i}(t) \geq 2(K^2 + KM + M)$  and the sub-graph  $G_t$  is  $l$ -periodically connected, then we have  $\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t)$ , where the minimum is taken over all clusters, not just the neighbors. This is equivalent to proving that the

agents stay on the same page and achieve consensus, which is guaranteed by  $l$ -periodically connected sub-graphs, not limited to connected sub-graphs. Based on the definition of the composition of graphs, the probability of having an edge is much larger than the edge probability  $p(m, n)$ . Thus, the degree of the composition graph is more likely to exceed  $\frac{C-1}{2}$ , which is a sufficient condition for connectivity, i.e., meeting the condition of being  $l$ -periodically connected.  $\square$

*Choice of  $l$ .* A natural question is how to specify such  $2 \leq l \leq C - 1$ . Since the result in Theorem 7 holds for any  $l$ , an optimal choice of  $l$  is the one that minimizes the lower bound on  $\min_{m,n} p(m, n)$ , which reads as  $l^* = \arg \min_l \frac{e}{e-1} \frac{C^2}{M^2} \max\{\frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-l-1)!}{(C-2)!} (\frac{3}{4})^{\frac{l}{2}}\}$ . Nevertheless, we can always use any  $l$ , which improves the previous results, as stated in the following.

*Comparison with Theorem 3 & 5.* The lower bound on the edge probability in Theorem 7 is given by  $\min_{m,n} p(m, n) \geq \frac{e}{e-1} \frac{C^2}{M^2} \max\{\frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-l-1)!}{(C-2)!} (\frac{3}{4})^{\frac{l}{2}}\}$ . In contrast, the corresponding lower bounds in Theorems 3 and 5 are  $1 - (\frac{1}{2} - \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$  and  $1 - (1 - \min\{(\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}), 1 - \delta(C-1)/8CT\})^{C^2/M^2}$ , respectively. When  $l > 1$ , the lower bound in Theorem 7 can be significantly smaller than the lower bound in Theorem 5. Specifically: - The additional term  $\frac{(C-l-1)!}{(C-2)!}$  is always less than 1 (and substantially so due to factorial decay). - The factor  $C^2/M^2$  is also smaller when  $C < M$ . These differences contribute significantly to relaxing the assumption on edge probabilities. Notably, in Theorems 3 and 5, the lower bounds converge to 1 as  $T \rightarrow \infty$ , implying that the graph becomes fully connected. However, this is no longer a concern in Theorem 7. For the second term, consider the case where  $C = M$  is large. In this scenario:  $\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}} \approx 1 - \frac{1}{2}\frac{\delta}{CT} > 1 - \frac{\delta(C-1)}{8CT}$ . Thus, the lower bound in Theorem 5 becomes approximately:  $1 - \frac{\delta(C-1)}{8CT}^{C^2/M^2} = 1 - \frac{\delta(C-1)}{8CT}$ . This demonstrates the improvement of Theorem 7 over Theorem 5 in this case. However, when  $C$  is small or  $C < M$ , the improvement mainly comes from the additional term in Theorem 7, as the difference between the second terms in Theorems 5 and 7 becomes negligible in comparison.

It is worth noting that the above three results assume that the in-cluster probability  $p(m, m) = 1$ , which may be violated in some cases. To this end, we further relax this assumption by considering more general scenarios. We derive similar results, but only require a lower bound on  $p(m, m)$  by additionally considering the  $l$ -periodically connected sub-graphs induced by the agents within the same cluster (motivated by considering the sub-graphs across clusters), rather than assuming  $p(m, m) = 1$ . Methodologically, in Algorithm 1, we set  $\tau = l$  instead of 1, which implies that the frequency of updating the within cluster information aligns with  $l$ -periodical connectivity. We next present the corresponding theoretical result below.

**THEOREM 8.** *Let us assume that  $\min_m p(m, m) \geq \max\{\frac{(|c_M|-l-1)!}{(|c_M|-2)!} (1 - \frac{\delta(|c_M|-1)}{8c_M T}), \frac{(|c_M|-l-1)!}{(|c_M|-2)!} (\frac{3}{4})^{\frac{l}{2}}\}$  for any  $1 \leq m \leq C$ , and that  $\min_{m \neq n} p(m, n) \geq \frac{e}{e-1} \cdot \frac{C^2}{M^2} \cdot \max\{\frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-l-1)!}{(C-2)!} (\frac{3}{4})^{\frac{l}{2}}\}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon$   $E[R_T | A'_{\epsilon, \delta}] \leq L + \sum_{i \neq i^*} \Delta_i (\max\{\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\} + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K) + l = O(\frac{C \log T}{M})$  where  $A'_{\epsilon, \delta}$  is defined in Theorem 7,  $L, c_0, C_1$  are specified in Theorem 3.*

**PROOF SKETCH.** The complete proof is referred to in Appendix E; here we illustrate the main proof logic. Instead of considering a complete sub-graph within one cluster, we consider  $l$ -periodically connected sub-graphs, which implies that the delay to receive all other agents' information within the



cluster is at most  $l$ . The assumption of  $\min_m p(m, m) \geq \max \left\{ \frac{(|c_M|-l-1)!}{(|c_M|-2)!} \left(1 - \frac{\delta(|c_M|-1)}{8CT}\right), \frac{(|c_M|-l-1)!}{(|c_M|-2)!} \left(\frac{3}{4}\right)^{\frac{1}{2}} \right\}$  guarantees that, with high probability, the sub-graph induced by the agents in one cluster is  $l$ -periodically connected at all times. The proof then follows from the statement for sub-graphs induced by the clusters. The regret incurred in between is at most  $l$ . Otherwise, once an agent collects the information from all agents in the same cluster, the algorithm proceeds as before, where the cluster can be treated as a complete graph (since all information is available). Subsequently, the previous regret bound also holds, which provides the regret bound as stated in Theorem 8.  $\square$

When interpreting the above result in more detail, we surprisingly find that it provides useful insights into the proposed framework, which incorporates both homogeneity and heterogeneity. We next illustrate these insights. If  $|c_M| < C$ , i.e.,  $M \leq C^2$ , the lower bound for the edge probability  $p(m, m)$  within one cluster is larger than that for the edge probability  $p(m, n)$  across clusters where  $m \neq n$ . This implies that for small-scale systems with a large number of clusters, sufficient edges within a cluster (homogeneity) are required to achieve the order reduction in the regret bound. Conversely, when  $|c_M| \gg C$ , i.e.,  $M \gg C^2$ , the reverse holds: the lower bound for the edge probability  $p(m, m)$  within one cluster becomes smaller than that for the edge probability  $p(m, n)$  across clusters where  $m \neq n$ . Hence, for large-scale systems, it is more important to ensure sufficient information is gathered across clusters (heterogeneity) to guarantee the regret bound.

## 7 Heterogeneous - Multiple Unknown Clusters

In practice, the cluster structure is often unknown due to the complexity of real-world data. Simply assuming no cluster structure, as in [58], is an overly simplistic approach that ignores the potential presence of clusters. This neglect not only oversimplifies the problem but also diverts leveraging these structures to design more refined and effective algorithms. To overcome this limitation, and unlike the case where the cluster assignment  $\{i \rightarrow c_i\}$  is known beforehand as in Section 6, we demonstrate that our methodological framework can be extended to scenarios where such information is unavailable. This extension opens the door to tackling more practical and complex real-world problems and providing robust tools. We show that, with only minor modifications, our algorithm seamlessly integrates with a cluster detection algorithm (Section 7.1). We establish the corresponding regret bound under these extended conditions in Section 7.2.

### 7.1 Algorithm

The method involves two steps: first, incorporating a cluster detection method (Algorithm 3) to estimate the cluster structure with reward information given by the burn-in period (Algorithm 2); and second, running Algorithm 1 using this estimation in a plug-and-play fashion.

**7.1.1 Cluster Detection.** To estimate the cluster structure from the graph and the reward information, we use the algorithm for cluster detection in [12], which extends the stochastic block model to contextual stochastic block models by incorporating node covariates as side information.

**DEFINITION 6.** *The Contextual Symmetric Stochastic Block Model (CSSBM) is a special stochastic block model with node covariates. Given the graph size  $M$ , number of clusters  $C$ , edge connection probabilities  $p, q$ ,  $C$  vectors  $\mu^1, \dots, \mu^C \in \mathbb{R}^K$  and variance  $\sigma^2$ , it generates a graph with  $M$  nodes and  $C$  balanced clusters from a stochastic block model with edge probability  $p(m, m) = p$  for any  $m$  and  $p(m, n) = q$  for any  $m \neq n$ . The node covariates are generated from a Gaussian mixture model, where for each node  $i$  in cluster  $c_m$ , its covariates  $v^i$  are sampled from Gaussian distribution  $\mathcal{N}(\mu^m, \sigma^2 I_K)$ .*

Braun et al. [12] proposed an iterative clustering algorithm IR-LSS to recover the cluster structure under CSSBM by incorporating both graph structure and vertex covariates (rewards in our case). We provide the pseudo-code of the IR-LSS algorithm (Algorithm 3) in Appendix A. The IR-LSS algorithm iteratively refines clusters by alternating between the following two steps - 1) Parameter Estimation: Using the current clustering, it estimates model parameters, such as the connectivity probability and covariate means, and 2) Clustering Refinement: It reassigns vertices to clusters by minimizing a least-squares criterion, based on estimated model parameters.

This iterative approach effectively leverages both graph and reward information, addressing cases where clusters are difficult to recover from either source alone, aligning perfectly with MA-MAB.

**7.1.2 Full Algorithm.** The full algorithm first runs Algorithm 2 to accumulate reward information and then executes Algorithm 3 to detect the cluster information, both of which are burn-in steps. Note that for each agent to run the clustering algorithm locally, an additional  $O(M)$  steps are required to propagate the estimated reward information and local edge information across the entire graph. Thus, the new burn-in steps have a total length of  $L_1 = O(L + M)$ , where  $L = O(\frac{C}{M} \cdot K \cdot \log T)$  represents the length of Algorithm 2. Subsequently, the full algorithm proceeds to the learning period by running Algorithm 1, as before, to design UCB-based strategies.

## 7.2 Analyses

Again, let us assume that the cluster structure is balanced for illustration purposes. As before, the imbalanced case can be easily handled by using the smallest cluster size.

We first present the key results regarding the cluster detection algorithm. By using the iterative refinement clustering algorithm by [12], the cluster structure of the CSSBM can be exactly recovered in the following regime efficiently with high probability. First, we define the Signal-to-Noise Ratio (SNR) to be  $\text{SNR} = \frac{1}{8\sigma^2} \min_{m \neq n} \|\mu^m - \mu^n\|_2^2 + \frac{\log M}{C} (\sqrt{p'} - \sqrt{q'})^2$ , where  $p = p' \frac{\log M}{M}$  and  $q = q' \frac{\log M}{M}$ .

**LEMMA 9 (BRAUN ET AL. [12]).** *Consider a CSSBM with  $M$  nodes,  $C$  clusters, and signal-to-noise ratio SNR. Suppose  $\text{SNR} > 2 \log M$  and  $C^3 \leq \text{SNR} \cdot \delta$  for a small constant  $\delta < 1/2$ . Then, with probability at least  $1 - 1/\text{poly}(M)$ , Algorithm 3 exactly recovers the cluster structure.*

Then, we use the graph information and the reward estimation for each agent from the burn-in period with  $L = O(\frac{C}{M} \cdot K \log T)$  rounds as the input for the above clustering algorithm. We assume that the reward for each arm  $i \in [K]$  of agent  $m \in [M]$  is sampled from a Gaussian distribution  $\mathcal{N}(\mu_i^m, \sigma^2)$ . Then, we can recover the cluster structure exactly under the following conditions.

**LEMMA 10.** *Consider the graph is generated from a stochastic block model with  $M$  agents and  $C$  balanced clusters with edge connection probability  $p(m, m) = p$  for all  $m \in [C]$  and  $p(m, n) = q$  for all  $m \neq n$ , where  $p = p' \frac{\log M}{M}$  and  $q = q' \frac{\log M}{M}$ . The reward for each arm  $i \in [K]$  of agent  $m \in [M]$  is sampled from a Gaussian distribution  $\mathcal{N}(\mu_i^m, \sigma^2)$ . Suppose  $\text{SNR} > 2 \log M$  and  $C^3 \leq \text{SNR} \cdot \delta$  for a small constant  $\delta < 1/2$ , where  $\text{SNR} = \frac{C \log T}{8M\sigma^2} \min_{m \neq n} \|\mu^m - \mu^n\|_2^2 + \frac{K \log T \log M}{M} (\sqrt{p'} - \sqrt{q'})^2$ . Then, with probability at least  $1 - 1/\text{poly}(M)$ , Algorithm 3 exactly recovers the cluster structure.*

The proof of Lemma 10 is provided in Appendix D. We would like to emphasize that, given the high probability of detecting the cluster structure, the regret bounds introduced in Section 6 remain valid with certain modifications. Specifically, the probability of the regret bound holds as the product of the original probability and the probability of exact cluster detection. Additionally, the regret

bound includes an extra term of order  $M$ , arising from the burn-in period  $L_1$  during which the cluster detection algorithm is executed.

Formally, we present the following corollary of Theorem 8, removing the assumption of known clusters from the theorem. We show the corollary of Theorem 8 for illustrative purposes, as Theorem 8 represents the most general form of the regret bound under the least stringent assumptions. Likewise, the corollaries of Theorems 2, 3, 5, and 7 also hold, as presented in Appendix C.

We first define event  $A_{\epsilon, \delta, \tau} = A_{\epsilon, \delta} \cap \{c_i = c'_i, \forall 1 \leq i \leq M\}$ , where  $c'_i$  is the cluster label for agent  $i$  recovered by Algorithm 3. It is straightforward to show that  $P(A_{\epsilon, \delta, \tau}) \geq 1 - 7\epsilon - \tau$ , where  $\tau = P(\{c_i \neq c'_i, \exists 1 \leq i \leq M\}) = 1 - 1/\text{poly}(M)$  is the failure probability for the cluster recovery.

**COROLLARY 11 (EXTENSION OF THEOREM 8).** *Let the assumptions in Theorem 8 hold except that  $\{c_i\}$  are unknown. Let the assumptions in Lemma 6 hold. Then the regret of Full Algorithm reads as with probability  $1 - 7\epsilon - 1/\text{poly}(M)$ ,  $E[R_T | A_{\epsilon, \delta, \tau}] \leq L_1 + \sum_{i \neq i^*} (\max\{\frac{C}{M} \cdot \frac{4C_1 \log T}{\Delta_i^2}, 2(K^2 + MK)\}) + \frac{2\pi^2}{3P(A_{\epsilon, \delta, \tau})} + K^2 + (2M - 1)K + l \leq O(\frac{C}{M} \log T)$ .*

**PROOF.** The proof is straightforward, combining the above lemma, which implies that after  $L_1$ , with probability  $1 - \tau$ , the cluster structure is correctly identified. From  $L_1$  to  $T$ , the agents run the same algorithm as before, resulting in the regret bound in Theorem 8, under  $A_{\epsilon, \delta, \tau}$ .  $\square$

This corollary demonstrates that our algorithm and analysis are general, making them applicable to scenarios where the cluster information is unknown and, thus, valuable for many real-world applications, as mentioned in Section 4.

## 8 Numerical Experiments

In this section, we numerically evaluate the performance of the proposed algorithm. Specifically, we compare the regret performance of Algorithm 2 over time with existing benchmark methods studied in the literature. We present the regret curves by computing the algorithms' exact regret as the average over 25 runs, along with the corresponding 95% confidence intervals (CI). We present the results on both synthetic datasets and a real-world dataset in this section. Details about numerical experiments are referred to Appendix B.

The benchmarks include DrFed-UCB [58], GoSInE [13], Gossip\_UCB [67], and Dist\_UCB [64]. Notably, GoSInE and Gossip\_UCB are tailored for time-invariant graphs, whereas Dist\_UCB and DrFed-UCB are recent works designed for time-varying graphs. Among these, we emphasize that DrFed-UCB is the most recent and highly relevant to our framework, as it considers Erdos-Renyi models and pure heterogeneity, while GoSInE focuses solely on homogeneity.

*Benchmark Comparison Results.* The comparison of the regret curves on the synthetic dataset between UCB-SBM (our method) and benchmarks is shown in Fig. 1a, where the shaded area denotes the CI. Here, the  $x$ -axis and  $y$ -axis represent the time  $t$  and the cumulative regret on a log scale  $\log R_t$ , respectively. Among these, UCB-SBM achieves the smallest regret. We observe that UCB-SBM consistently demonstrates significantly smaller regret, showcasing notable improvements and highlighting the advantages of considering homogeneity within clusters and relaxing the assumption on edge probability. More precisely, the improvements in average regret  $R_t$  compared to DrFed-UCB, GoSInE, Gossip\_UCB, and Dist\_UCB are **68.79%**, **79.80%**, **94.14%**, and **94.75%**, respectively. The comparison with DrFed-UCB emphasizes the heavy performance degradation of DrFed-UCB when neglecting the cluster structure. Meanwhile, our algorithm exhibits small

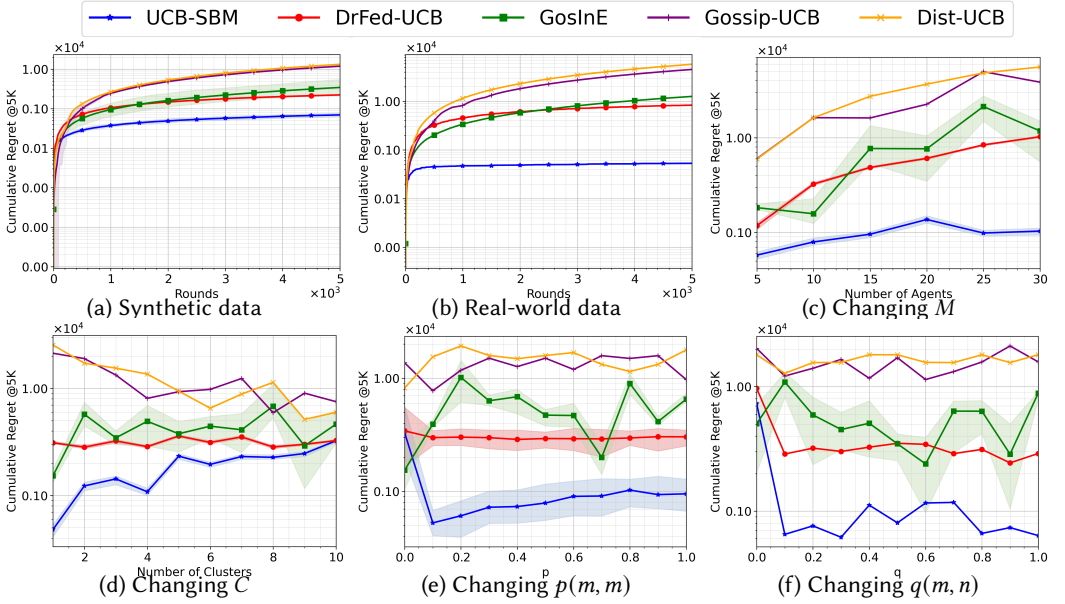


Fig. 1. The regret of different methods across different settings

variances (with GoInE showing the largest), indicating stability even with time-varying graphs. Likewise, we draw similar conclusions from the real-world dataset, as presented in Fig. 1b, wherein our improvement is even more significant.

*Regret Dependency Results.* Moreover, we demonstrate how the actual regret of UCB-SBM depends on several parameters associated with the problem setting, including the number of agents  $M$ , the number of clusters  $C$ , and the parameters  $p = p(m, m)$  and  $q = p(m, n)$  (which also affect  $|p(m, m) - p(m, n)|$ ). While other parameters, such as  $K$  and the difference in mean values  $h$ , are important, they have been studied in [58]. The aforementioned parameters are unique to our problem setting and necessitate examining how the actual regret changes with them beyond the theoretical upper bounds. The results of varying  $M$ ,  $C$ ,  $p$ , and  $q$  are shown in Figs. 1c, 1d, 1e, and 1f, respectively. First, across all possible settings, our algorithm consistently achieves the smallest regret, and its performance does not change dramatically with different parameters, demonstrating both the effectiveness and robustness of the algorithm. In Fig. 1c, we observe that, except for UCB-SBM, all other algorithms exhibit an increasing trend as  $M$  increases, while UCB-SBM remains steady, implying that UCB-SBM scales much better with  $M$  and is thus more practical. In Fig. 1d, UCB-SBM’s regret increases with  $C$ , consistent with the theoretical bound’s dependence on  $C$ . Notably, when  $C = 10 = M$ , i.e., the fully heterogeneous case, UCB-SBM and DrFed-UCB achieve the same regret, validating the consistency of the results. The regret of UCB-SBM increases and then decreases with  $p$  and  $q$ , as shown in Fig. 1e and Fig. 1f. This is possibly because lower bounds on  $p$  and  $q$  are necessary for the theoretical regret bounds to hold, as this pattern holds true for DrFed-UCB. Establishing an explicit dependency of  $R_T$  on  $p$  and  $q$  is left for future exploration.

## 9 Conclusion and Future Work

In this paper, we novelly study the multi-agent multi-armed bandit (MA-MAB) problem, where agents are distributed on random graphs induced by a cluster structure, namely stochastic block models, and their reward mean values also depend on the cluster structure. This introduces both

homogeneity and heterogeneity in edge probabilities and rewards, within and across clusters. The cluster assignment can be either known or unknown. This is the first framework that unifies the existing formulations of both homogeneous MA-MAB (1 cluster) and heterogeneous MA-MAB ( $M$  clusters), smoothly capturing more general cases in between and reflecting different degrees of heterogeneity. Algorithmically, we propose a new method where agents within one cluster aggregate their information to achieve sample complexity reduction, communicate with other clusters to collect heterogeneous information, integrate this information to estimate the globally optimal arm, and pull arms based on newly designed UCB indices. When the cluster assignment is unknown, the agents leverage a cluster detection algorithm to estimate the cluster assignment, and our algorithm operates in a plug-and-play fashion, demonstrating its generalization ability. This approach leads to significantly improved results under less stringent assumptions. Theoretically, we show that the regret bound has a constant reduction of  $\frac{C}{M}$ , uncovering how the regret bound changes with the degree of heterogeneity and improving upon existing work [58], beyond solely  $T$ . Moreover, the assumption on the minimal edge probability of the random graph is significantly relaxed, scaling better with  $T$ . Notably, while the minimal edge probability in existing work can approach 1 as  $T \rightarrow \infty$ , our approach bounds it by much smaller values (e.g.,  $\frac{e}{e-1} \frac{C^2}{M^2}$  and  $\frac{e}{e-1} \frac{C^2}{M^2} \frac{(C-l-1)!}{(C-2)!}$ ). Numerically, we demonstrate the superior performance of the proposed algorithm by comparing it with benchmarks. Consistently, our algorithm shows significant regret improvement, with the relative improvement percentage being at least 68%. We also examine how actual regret changes with parameters unique to the framework, consistent with the theoretical findings.

Moving forward, we identify several promising directions for future work. First, while we assume a balanced cluster structure or use the minimal cluster size to run the algorithm in unbalanced cases, it would be interesting to explore how to fully leverage unbalanced cluster structures instead of relying solely on the minimal cluster size. Additionally, while we assume that the reward distribution is sub-Gaussian (and can be extended to sub-exponential cases), more general heavy-tailed distributions present another direction for future research. Lastly, exploring other types of cluster structures, beyond stochastic block models, and characterizing how regret changes with these structures would be of great interest to both theorists and practitioners.

## References

- [1] E. Abbe. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research*, 18(177):1–86, 2018.
- [2] E. Abbe, A. S. Bandeira, and G. Hall. Exact recovery in the stochastic block model. *IEEE Transactions on information theory*, 62(1):471–487, 2015.
- [3] E. Abbe, J. Fan, and K. Wang. An lp theory of pca and spectral clustering. *The Annals of Statistics*, 50(4):2359–2385, 2022.
- [4] M. Agarwal, V. Aggarwal, and K. Azizzadenesheli. Multi-agent multi-armed bandits with limited communication. *The Journal of Machine Learning Research*, 23(1):9529–9552, 2022.
- [5] E. M. Airoidi, D. M. Blei, S. E. Fienberg, E. P. Xing, and T. Jaakkola. Mixed membership stochastic block models for relational data with application to protein-protein interactions. In *Proceedings of the international biometrics society annual meeting*, volume 15, page 1, 2006.
- [6] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- [7] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- [8] Y. Ban, Y. Qi, T. Wei, L. Liu, and J. He. Meta clustering of neural bandits. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 95–106, 2024.
- [9] S. Battiston and M. Catanzaro. Statistical properties of corporate board and director networks. *The European Physical Journal B*, 38:345–352, 2004.
- [10] I. Bistriz and A. Leshem. Distributed multi-player bandits—a game of thrones approach. *Advances in Neural Information Processing Systems*, 31, 2018.
- [11] E. Blaser, C. Li, and H. Wang. Federated linear contextual bandits with heterogeneous clients. In *International Conference on Artificial Intelligence and Statistics*, pages 631–639. PMLR, 2024.
- [12] G. Braun, H. Tyagi, and C. Biernacki. An iterative clustering algorithm for the contextual stochastic block model with optimality guarantees. In *International Conference on Machine Learning*, pages 2257–2291. PMLR, 2022.
- [13] R. Chawla, A. Sankararaman, A. Ganesh, and S. Shakkottai. The gossiping insert-eliminate algorithm for multi-agent bandits. In *International conference on artificial intelligence and statistics*, pages 3471–3481. PMLR, 2020.
- [14] L. Chen, J. Xu, S. Ren, and P. Zhou. Spatio-temporal edge service placement: A bandit learning approach. *IEEE Transactions on Wireless Communications*, 17(12):8388–8401, 2018.
- [15] M. Cugmas, F. Mali, and A. Žiberna. Scientific collaboration of researchers and organizations: a two-level blockmodeling approach. *Scientometrics*, 125(3):2471–2489, 2020.
- [16] X. Dai, Z. Zhang, P. Yang, Y. Xu, X. Liu, and J. C. Lui. Axiomvision: Accuracy-guaranteed adaptive visual model selection for perspective-aware video analytics. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7229–7238, 2024.
- [17] F. Delarue. Mean field games: A toy model on an Erdős-Renyi graph. *ESAIM: Proceedings and Surveys*, 60:1–26, 2017.
- [18] Y. Deshpande, S. Sen, A. Montanari, and E. Mossel. Contextual stochastic block models. *Advances in Neural Information Processing Systems*, 31, 2018.
- [19] M. Dreveton, F. Fernandes, and D. Figueiredo. Exact recovery and bregman hard clustering of node-attributed stochastic block model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [20] A. Dubey and A. Pentland. Cooperative multi-agent bandits with heavy tails. In *International Conference on Machine Learning*, 2730–2739, 2020.
- [21] Q. Duchemin. Reliable prediction in the markov stochastic block model. *ESAIM: Probability and Statistics*, 27:80–135, 2023.
- [22] A. El Haj. Community detection in multiplex continuous weighted nodes networks using an extension of the stochastic block model. *Computing*, 106(11):3711–3725, 2024.
- [23] P. ERDős and A. R&w. On random graphs i. *Publ. math. debrecen*, 6(290-297):18, 1959.
- [24] C. Gentile, S. Li, and G. Zappella. Online clustering of bandits. In *International conference on machine learning*, pages 757–765. PMLR, 2014.
- [25] C. Gentile, S. Li, P. Kar, A. Karatzoglou, G. Zappella, and E. Etrud. On context-dependent clustering of bandits. In *International Conference on machine learning*, pages 1253–1262. PMLR, 2017.
- [26] P. W. Holland, K. B. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [27] R. Huang, W. Wu, J. Yang, and C. Shen. Federated linear contextual bandits. *Advances in Neural Information Processing Systems*, 34:27057–27068, 2021.
- [28] F. Jiang and H. Cheng. Multi-agent bandit with agent-dependent expected rewards. *Swarm Intelligence*, 1–33, 2023.
- [29] N. Korda, B. Szorenyi, and S. Li. Distributed clustering of linear bandits in peer to peer networks. In *International conference on machine learning*, pages 1301–1309. PMLR, 2016.

- [30] P. Landgren, V. Srivastava, and N. E. Leonard. On distributed cooperative decision-making in multiarmed bandits. In *2016 European Control Conference*. 243–248. IEEE, 2016.
- [31] P. Landgren, V. Srivastava, and N. E. Leonard. Distributed cooperative decision-making in multiarmed bandits: Frequentist and Bayesian algorithms. In *2016 IEEE 55th Conference on Decision and Control*. 167–172. IEEE, 2016.
- [32] P. Landgren, V. Srivastava, and N. E. Leonard. Distributed cooperative decision making in multi-agent multi-armed bandits. *Automatica*, 125:109445, 2021.
- [33] Q. Li, C. Zhao, T. Yu, J. Wu, and S. Li. Clustering of conversational bandits with posterior sampling for user preference learning and elicitation. *User Modeling and User-Adapted Interaction*, 33(5):1065–1112, 2023.
- [34] S. Li and S. Zhang. Online clustering of contextual cascading bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [35] S. Li, C. Gentile, A. Karatzoglou, and G. Zappella. Online context-dependent clustering in recommendations based on exploration-exploitation algorithms. *ArXiv, abs/1608.03544*, 2016.
- [36] S. Li, A. Karatzoglou, and C. Gentile. Collaborative filtering bandits. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 539–548, 2016.
- [37] S. Li, W. Chen, and K.-S. Leung. Improved algorithm on online clustering of bandits. *arXiv preprint arXiv:1902.09162*, 2019.
- [38] T. Li and L. Song. Privacy-preserving communication-efficient federated multi-armed bandits. *IEEE Journal on Selected Areas in Communications*, 40(3):773–787, 2022.
- [39] Z. Li, M. Liu, X. Dai, and J. Lui. Demystifying online clustering of bandits: Enhanced exploration under stochastic and smoothed adversarial contexts. *arXiv preprint arXiv:2501.00891*, 2025.
- [40] F. W. Lima, A. O. Sousa, and M. Sumuor. Majority-vote on directed Erdős–Rényi random graphs. *Physica A: Statistical Mechanics and its Applications*, 387(14):3503–3510, 2008.
- [41] X. Liu, H. Zhao, T. Yu, S. Li, and J. C. Lui. Federated online clustering of bandits. In *Uncertainty in Artificial Intelligence*, pages 1221–1231. PMLR, 2022.
- [42] D. Martínez-Rubio, V. Kanade, and P. Rebeschini. Decentralized cooperative stochastic bandits. *Advances in Neural Information Processing Systems*, 32, 2019.
- [43] A. Mitra, H. Hassani, and G. Pappas. Exploiting heterogeneity in robust federated best-arm identification. *arXiv preprint arXiv:2109.05700*, 2021.
- [44] T. T. Nguyen and H. W. Lauw. Dynamic clustering of contextual multi-armed bandits. In *Proceedings of the 23rd ACM international conference on conference on information and knowledge management*, pages 1959–1962, 2014.
- [45] S. Pal, A. Suggala, K. Shanmugam, and P. Jain. Blocked collaborative bandits: online collaborative filtering with per-item budget constraints. *Advances in Neural Information Processing Systems*, 36, 2024.
- [46] C. Réda, S. Vakili, and E. Kaufmann. Near-optimal collaborative learning in bandits. In *2022-36th Conference on Neural Information Processing System*, 2022.
- [47] R. Roman, J. Zhou, and J. Lopez. On the features and challenges of security and privacy in distributed internet of things. *Computer networks*, 57(10):2266–2279, 2013.
- [48] A. Sankararaman, A. Ganesh, and S. Shakkottai. Social learning in multi agent multi armed bandits. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 3(3):1–35, 2019.
- [49] N. Stanley, T. Bonacci, R. Kwitt, M. Niethammer, and P. J. Mucha. Stochastic block models with multiple continuous attributes. *Applied Network Science*, 4:1–22, 2019.
- [50] P.-A. Wang, A. Proutiere, K. Ariu, Y. Jedra, and A. Russo. Optimal algorithms for multiplayer multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 4120–4129. PMLR, 2020.
- [51] P.-A. Wang, A. Proutiere, K. Ariu, Y. Jedra, and A. Russo. Optimal algorithms for multiplayer multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 4120–4129. PMLR, 2020.
- [52] Q. Wang, C. Zeng, W. Zhou, T. Li, S. S. Iyengar, L. Shwartz, and G. Y. Grabarnik. Online interactive collaborative filtering using multi-armed bandit with dependent arms. *IEEE Transactions on Knowledge and Data Engineering*, 31(8): 1569–1580, 2019. doi: 10.1109/TKDE.2018.2866041.
- [53] X. Wang, L. Yang, Y.-Z. J. Chen, X. Liu, M. Hajiesmaili, D. Towsley, and J. C. Lui. Achieving near-optimal individual regret & low communications in multi-agent bandits. In *The Eleventh International Conference on Learning Representations*, 2022.
- [54] X. Wang, L. Yang, Y.-Z. J. Chen, X. Liu, M. Hajiesmaili, D. Towsley, and J. C. Lui. Achieve near-optimal individual regret & low communications in multi-agent bandits. In *International Conference on Learning Representations*, 2023.
- [55] Z. Wang, C. Zhang, M. K. Singh, L. Riek, and K. Chaudhuri. Multitask bandit learning through heterogeneous feedback aggregation. In *International Conference on Artificial Intelligence and Statistics*, 1531–1539, 2021.
- [56] J. Wu, C. Zhao, T. Yu, J. Li, and S. Li. Clustering of conversational bandits for user preference learning and elicitation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 2129–2139,

2021.

- [57] M. Xu and D. Klabjan. Regret lower bounds in multi-agent multi-armed bandit. *arXiv preprint arXiv:2308.08046*, 2023.
- [58] M. Xu and D. Klabjan. Decentralized randomly distributed multi-agent multi-armed bandit with heterogeneous rewards. *Advances on Neural Information Processing Systems*, 2023.
- [59] Z. Yan, Q. Xiao, T. Chen, and A. Tajer. Federated multi-armed bandit via uncoordinated exploration. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. 5248–5252. IEEE, 2022.
- [60] H. Yang, X. Liu, Z. Wang, H. Xie, J. C. Lui, D. Lian, and E. Chen. Federated contextual cascading bandits with asynchronous communication and heterogeneous users. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 20596–20603, 2024.
- [61] P. Yang, N. Zhang, S. Zhang, L. Yu, J. Zhang, and X. Shen. Content popularity prediction towards location-aware mobile edge caching. *IEEE Transactions on Multimedia*, 21(4):915–929, 2018.
- [62] W. W. Zachary. An information flow model for conflict and fission in small groups. *Journal of anthropological research*, 33(4):452–473, 1977.
- [63] Y. Zhao, J. Zhao, L. Jiang, R. Tan, D. Niyato, Z. Li, L. Lyu, and Y. Liu. Privacy-preserving blockchain-based federated learning for iot devices. *IEEE Internet of Things Journal*, 8(3):1817–1829, 2020.
- [64] J. Zhu and J. Liu. Distributed multi-armed bandits. *IEEE Transactions on Automatic Control*, 2023.
- [65] J. Zhu, R. Sandhu, and J. Liu. A distributed algorithm for sequential decision making in multi-armed bandit with homogeneous rewards. In *59th IEEE Conference on Decision and Control*. 3078–3083. IEEE, 2020.
- [66] J. Zhu, E. Mulle, C. S. Smith, and J. Liu. Decentralized multi-armed bandit can outperform classic upper confidence bound. *arXiv preprint arXiv:2111.10933*, 2021.
- [67] Z. Zhu, J. Zhu, J. Liu, and Y. Liu. Federated bandit: A gossiping approach. In *Abstract Proceedings of the 2021 ACM SIGMETRICS/International Conference on Measurement and Modeling of Computer Systems*, 3–4, 2021.



## A Pseudo Code of Algorithms

### A.1 Burn-in Peirod

The full algorithm of the burn-in period described in Section 6 is shown as follows.

---

#### Algorithm 2: UCB-SBM: Burn-in period [58]

---

- 1 Initialization: The length of the burn-in period is  $L$ ; the estimates are initialized as  $\bar{\mu}_i^m(0) = 0$ ,  $n_{m,i}(0) = 0$ ,  $\hat{\mu}_{i,j}^m(0) = 0$ , and  $P_0(m, j) = 0$  for any arm  $i$  and agents  $m, j$ ;
  - 2 **for**  $1 < t \leq L$  **do**
  - 3     **for each agent**  $m$  **do**
  - 4         Sample arm  $a_t^m = (t \bmod K)$ ;
  - 5         Receive rewards  $r_{a_t^m}^m(t)$  and update  $n_{m,i}(t) = n_{m,i}(t-1) + \mathbb{1}_{a_t^m=i}$ ;
  - 6         Update the local estimates for any arm  $i$ :  $\bar{\mu}_i^m(t) = \frac{n_{m,i}(t-1)\bar{\mu}_i^m(t-1) + r_{a_t^m}^m(t) \cdot \mathbb{1}_{a_t^m=i}}{n_{m,i}(t-1) + \mathbb{1}_{a_t^m=i}}$ ;
  - 7         Update the maintained matrix  $P_t(m, j) = ((t-1)P_{t-1}(m, j) + X_{m,j}^t)/t$  for each  $j \in V$ ;
  - 8         Send  $\{\bar{\mu}_i^m(t)\}_{i=1}^K$  to all agents in  $\mathcal{N}_m(t)$ ;
  - 9         Receive  $\{\bar{\mu}_i^j(t)\}_{i=1}^K$  from all agents  $j \in \mathcal{N}_m(t)$  and store them as  $\hat{\mu}_{i,j}^m(t)$ .
  - 10     **end**
  - 11 **end**
  - 12 **for each agent**  $m$  **and arm**  $i$  **do**
  - 13     For agent  $1 \leq j \leq M$ , let  $t_{m,j} = \max_{s \geq \tau_1} \{(m, j) \in E_s\}$  or 0 if such  $s$  does not exist
  - 14      $\tilde{\mu}_i^m(L+1) = \sum_{j=1}^M P'_{m,j}(L) \hat{\mu}_{i,j}^m(t_{m,j})$  where  $P'_{m,j}(L) = \frac{1}{M}$  if  $P_L(m, j) > 0$  and 0 o.w.;
  - 15 **end**
- 

### A.2 Clustering Algorithm

The full algorithm of the cluster detection described in Section 7 is presented below.

---

#### Algorithm 3: Iterative Refinement Clustering (IR-LSS) [12]

---

- 1 **Input:** adjacency matrix  $A \in \{0, 1\}^{M \times M}$ , node covariates  $V \in \mathbb{R}^{M \times K}$ , variance  $\sigma^2$ , initial cluster assignment  $Z^{(0)} \in \{0, 1\}^{M \times C}$  and iterations  $T > 1$ ;
  - 2 **Output:** a cluster assignment  $Z \in \{0, 1\}^{M \times C}$
  - 3 **for**  $t = 0, 2, \dots, T-1$  **do**
  - 4     Estimate the model parameters from the current cluster assignment  $Z^{(t)}$ :  $s_n^{(t)} = |c_n^{(t)}|$ ,  $W^{(t)} = Z^{(t)}(D^{(t)})^{-1}$  where  $D^{(t)} = \text{diag}(s_1^{(t)}, \dots, s_C^{(t)})$ ,  $\Pi^{(t)} = (W^{(t)})^\top A W^{(t)}$ , and  $\mu_n^{(t)} = W_n^{(t)} V$ ;
  - 5     Refine clustering by assigning each node  $i$  to the cluster
 
$$z_i^{(t+1)} = \arg \min_{n \in [C]} \|(A_i W^{(t)} + \Pi_n^{(t)}) \sqrt{\Sigma_n^{(t)}}\|_2^2 + \|\mu_n^{(t)} - V_i\|_2^2 / \sigma^2,$$
 where  $\Sigma_n^{(t)} = \frac{M}{C(p^{(t)} - q^{(t)})} \log\left(\frac{p^{(t)}(1-q^{(t)})}{q^{(t)}(1-p^{(t)})}\right) I_C$  with  $p^{(t)} = \sum_{n \in [C]} \Pi_{nn}^{(t)} / C$  and  $q^{(t)} = \sum_{n \neq m} \Pi_{mn}^{(t)} / (C(C-1))$ ;
  - 6     Form the cluster assignment matrix  $Z^{(t+1)}$ ;
  - 7 **end**
-

## B Experiment Details

We present the experimental details in this section. We introduce the datasets, including both the synthetic data in Section B.1 and the real-world dataset, in Section B.2.

### B.1 Synthetic Datasets

We now describe the datasets used in the experiments and provide additional details that can be leveraged to reproduce the experiments discussed in Section 8.

**Synthetic Dataset.** We first examine the performance of the algorithms on a synthetic dataset. The data generation process is as follows: for the results shown in Fig. 1a, we select the number of agents as  $M = 10$ , the number of clusters as  $C = 10$ , the inter-cluster probability as  $p = 0.5$ , and the intra-cluster probability as  $q = 0.5$ . The length of the game is  $5 \cdot 10^5$ .

For Figs. 1c, 1d, 1e, and 1f, we vary  $M$ ,  $C$ ,  $p$ , and  $q$ , respectively, while keeping the other parameters fixed.

### B.2 Real-world Datasets

**Real-world Dataset.** Besides the synthetic dataset, we evaluate our algorithm and the benchmark algorithms on a real-world dataset, as reported in Fig. 1b, using the well-known Zachary’s Karate Club dataset, a widely used benchmark for graph clustering algorithms. This dataset represents the social interactions among 34 members of a university karate club, as observed and documented by Wayne W. Zachary in the 1970s [62]. The dataset consists of 34 nodes (agents), each representing a club member, and 78 unweighted, undirected edges that denote friendships between members.

## C Additional Theoretical Results

In this section, we present the corollaries of Theorem 2, 3, 5, 7 when the cluster structure is unknown. The proof of them follow the proof of Corollary 11 as in the main body, and as a result, we omit the proof steps here. We use  $L + L_1 = O(L + M)$  to denote the length of the burn-in period for the unknown cluster structure setting which contains  $L$  steps for Algorithm 2 and  $L_1 = O(M)$  steps to propagate information to the entire graph for cluster recovery. We use  $\tau = 1/\text{poly}(M)$  to denote the failure probability of cluster recovery.

**COROLLARY 12 (EXTENSION OF THEOREM 2).** *Let us assume that  $\min_{m,n} p(m, n) \geq (\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{MT})^{\frac{2}{M-1}}})$ .*

*For every  $0 < \epsilon < 1$  and  $0 < \delta < \frac{1}{2} + \frac{1}{4}\sqrt{1 - (\frac{\epsilon}{MT})^{\frac{2}{M-1}}}$ , the regret of Algorithm 2 with Rule 1 is upper bounded by with probability  $1 - 7\epsilon - \tau$ ,*

$$E[R_T | A'_{\epsilon, \delta, \tau}] \leq L + L_1 + \sum_{i \neq i^*} \Delta_i (\max \{ \lceil \frac{4C_1 \log T}{\Delta_i^2} \rceil, 2(K^2 + MK) \}) + \frac{2\pi^2}{3P(A'_{\epsilon, \delta, \tau})} + K^2 + (2M - 1)K$$

**COROLLARY 13 (EXTENSION OF THEOREM 3).** *Let us assume that  $p(m, m) = 1$  for any  $1 \leq m \leq C$ .*

*Let us further assume that  $\min_{m,n} p(m, n) \geq \min_{m,n} p(m, n) \geq 1 - (\frac{1}{2} - \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}})^{C^2/M^2}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon - \tau$*

$$E[R_T | A'_{\epsilon, \delta, \tau}] \leq L + L_1 + \sum_{i \neq i^*} \Delta_i (\max \{ \frac{C}{M} \cdot \lceil \frac{4C_1 \log T}{\Delta_i^2} \rceil, 2(K^2 + MK) \}) + \frac{2\pi^2}{3P(A'_{\epsilon, \delta, \tau})} + K^2 + (2M - 1)K + l$$

COROLLARY 14 (EXTENSION OF THEOREM 5). *Let us assume that  $p(m, m) = 1$  for any  $1 \leq m \leq C$ . Let us further assume that  $\min_{m,n} p(m, n) \geq 1 - (1 - \min\{\frac{1}{2} + \frac{1}{2}\sqrt{1 - (\frac{\delta}{CT})^{\frac{2}{C-1}}}, 1 - \delta^{(C-1)}/8CT\})^{C^2/M^2}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon - \tau$*

$$\begin{aligned} & E[R_T | A'_{\epsilon, \delta, \tau}] \\ & \leq L + L_1 + \sum_{i \neq i^*} \Delta_i (\max\{\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\}) + \frac{2\pi^2}{3P(A'_{\epsilon, \delta, \tau})} + K^2 + (2M - 1)K + l \end{aligned}$$

COROLLARY 15 (EXTENSION OF THEOREM 7). *Let us assume that  $\min_m p(m, m) = 1$  for any  $1 \leq m \leq C$ , and that  $\min_{m \neq n} p(m, n) \geq \frac{e}{e-1} \frac{C^2}{M^2} \max\{\frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8CT}), \frac{(C-l-1)!}{(C-2)!} (\frac{3}{4})^{\frac{1}{l}}\}$ . The regret bound of Algorithm 2 with Rule 2 reads as with probability  $1 - 7\epsilon - \tau$*

$$\begin{aligned} & E[R_T | A'_{\epsilon, \delta, \tau}] \\ & \leq L + L_1 + \sum_{i \neq i^*} \Delta_i (\max\{\frac{C}{M} \cdot [\frac{4C_1 \log T}{\Delta_i^2}], 2(K^2 + MK)\}) + \frac{2\pi^2}{3P(A'_{\epsilon, \delta, \tau})} + K^2 + (2M - 1)K + l \end{aligned}$$

## D Proof of Lemmas on Stochastic Block Models

LEMMA (LEMMA 4). *For any pair of vertices  $c_m, c_n$  in the sub-graph  $G_t^C$ , the probability that  $c_m$  and  $c_n$  is connected in  $G_t^C$  is  $p(c_m, c_n) = 1 - (1 - p(m, n))^{M^2/C^2}$ .*

PROOF OF LEMMA 4. Since all clusters have the same size, the clusters  $c_m$  and  $c_n$  have  $\frac{M}{C}$  agents each. Thus, there are  $\frac{M^2}{C^2}$  pairs of vertices between clusters  $c_m$  and  $c_n$ . Since each pair of such vertices is connected with probability  $p(m, n)$  independently, the probability that there is at least one edge between clusters  $c_m$  and  $c_n$  is  $1 - (1 - p(m, n))^{M^2/C^2}$ .  $\square$

LEMMA (LEMMA 6). *For any pair of vertices  $c_m, c_n$  in the sub-graph  $G_t^C$ , the probability that  $c_m$  and  $c_n$  is connected in  $G_t^C$  is  $p(c_m, c_n) \geq (1 - \frac{1}{e}) \min\{1, \frac{M^2}{C^2} \cdot p(m, n)\}$ .*

PROOF OF LEMMA 6. Since all clusters have the same size, the clusters  $c_m$  and  $c_n$  have  $\frac{M}{C}$  agents each. Thus, there are  $\frac{M^2}{C^2}$  pairs of vertices between clusters  $c_m$  and  $c_n$ . Since each pair of such vertices is connected with probability  $p(m, n)$  independently, the probability that there is at least one edge between clusters  $c_m$  and  $c_n$  is  $1 - (1 - p(m, n))^{M^2/C^2}$ . We consider two cases: (1)  $p(m, n) \geq \frac{C^2}{M^2}$ ; (2)  $p(m, n) < \frac{C^2}{M^2}$ .

For  $p(m, n) \geq \frac{C^2}{M^2}$ , we have

$$1 - (1 - p(m, n))^{M^2/C^2} \geq 1 - (1 - \frac{C^2}{M^2})^{M^2/C^2} \geq 1 - \frac{1}{e}.$$

For  $p(m, n) < \frac{C^2}{M^2}$ , we consider  $\frac{1 - (1 - p(m, n))^{M^2/C^2}}{p(m, n) M^2/C^2}$ , which is decreasing with  $p(m, n) \in (0, 1]$ . Then, we have

$$\frac{1 - (1 - p(m, n))^{M^2/C^2}}{p(m, n) M^2/C^2} \geq 1 - (1 - \frac{C^2}{M^2})^{M^2/C^2} \geq 1 - \frac{1}{e}.$$

$\square$

LEMMA (BRAUN ET AL. [12]; LEMMA 9). *Consider a CSSBM with  $M$  nodes,  $C$  clusters, and signal-to-noise ratio SNR. Suppose  $\text{SNR} > 2 \log M$  and  $C^3 \leq \text{SNR} \cdot \delta$  for a small constant  $\delta < 1/2$ . Then, with probability at least  $1 - 1/\text{poly}(M)$ , Algorithm 3 exactly recovers the community structure.*

Then, we use the graph information and the reward estimation for each agent from the burn-in period with  $L = O(\log T)$  rounds as the input for the above clustering algorithm. We assume that the reward for each arm  $i \in [K]$  of agent  $m \in [M]$  is sampled from a Gaussian distribution  $\mathcal{N}(\mu_i^m, \sigma^2)$ . Then, we can recover the cluster structure exactly under the following condition.

LEMMA (LEMMA 10). *Consider the graph is generated from a stochastic block model with  $M$  agents and  $C$  balanced clusters with edge connection probability  $p(m, m) = p$  for all  $m \in [C]$  and  $p(m, n) = q$  for all  $m \neq n$ , where  $p = p' \frac{\log M}{M}$  and  $q = q' \frac{\log M}{M}$ . The reward for each arm  $i \in [K]$  of agent  $m \in [M]$  is sampled from a Gaussian distribution  $\mathcal{N}(\mu_i^m, \sigma^2)$ . Suppose  $\text{SNR} > 2 \log M$  and  $C^3 \leq \text{SNR} \cdot \delta$  for a small constant  $\delta < 1/2$ , where*

$$\text{SNR} = \frac{\log T}{8K\sigma^2} \min_{m \neq n} \|\mu^m - \mu^n\|_2^2 + \frac{\log M \log T}{C} (\sqrt{p'} - \sqrt{q'})^2.$$

Then, with probability at least  $1 - 1/\text{poly}(M)$ , Algorithm 3 exactly recovers the community structure.

PROOF OF LEMMA 10. First, we consider the information from the burn-in period and show that this information can be formed as an instance from the CSSBM. The length of the burn-in period is  $L = O(\log T)$ . Consider the graph  $G_L$  on all  $M$  agents as follows. For any pair of agents  $i, j \in [M]$ , there is an edge connecting  $i, j$  in  $G_L$  if and only if agents  $i$  and  $j$  are connected for at least once in the burn-in period. Thus, for any two agents  $i, j$  in the same cluster, they are connected in  $G_L$  with probability  $p_L = 1 - (1 - p)^L$ . For any two agents  $m, n$  in different clusters, they are connected in  $G_L$  with probability  $q_L = 1 - (1 - q)^L$ . We consider the regime where  $p = p' \log M/M$  and  $q = q' \log M/M$  for constants  $p'$  and  $q'$ . Since edge connection probabilities  $p$  and  $q$  are small, we have  $p_L \approx Lp$  and  $q_L \approx Lq$ . This graph  $G_L$  can be seen as generated from the stochastic block model with  $M$  agents and  $C$  balanced clusters and edge probability  $p(m, m) = p_L$  for  $m \in [C]$  and  $p(m, n) = q_L$  for  $m \neq n$ . For each agent  $m \in [M]$ , we use the reward local estimates  $\bar{\mu}^m(L)$  from the burn-in period as the node covariates. In the burn-in period, each arm is pulled  $L/K$  times for each agent. Since we assume the reward is sampled from a Gaussian distribution with variance  $\sigma^2$ , the reward local estimates  $\bar{\mu}^m(L)$  is a random variable from the Gaussian distribution  $\mathcal{N}(\mu^m, \sigma^2 K/L)$ .

Now, we show that this instance from the CSSBM can be exactly recovered under the given conditions. Note that the Signal-to-noise ratio of the above CSSBM is

$$\text{SNR} = \frac{\log T}{8K\sigma^2} \min_{m \neq n} \|\mu^m - \mu^n\|_2^2 + \frac{\log M \log T}{C} (\sqrt{p'} - \sqrt{q'})^2.$$

By Lemma 9, when  $\text{SNR} > 2 \log M$  and  $C^3 \leq \text{SNR} \cdot \delta$  for a small constant  $\delta < 1/2$ , Algorithm 3 exactly recover the cluster structure with probability at least  $1 - 1/\text{poly}(M)$ .

□

## E Proof of Theorems

### E.1 Proof of Theorem 1

PROOF. **Main intuition:** Fix a suboptimal arm  $k$ . After the total number of observations for this arm  $k$  exceeds the sample complexity threshold, in expectation, it takes  $\frac{1}{p^{M^2}}$  time slots for all agents to get the information of this arm  $k$ . After that, no more regret will be incurred on this arm  $k$ .

Below, we present the proof of the algorithm. The proof first makes an assumption to reduce the problem to a standard cooperative UCB for homogeneous agent residing on a complete graph with communication delays. Then, we show that this assumption can be fulfilled in the single cluster scenario.

**Step 1:** Consider a homogeneous multi-agent multi-arm bandit model, where in each time slot, with probability  $q \in (0, 1)$ , a global synchronization of observations would happen among all agents. Applying Wang et al. [54, Lemmas 1 and 2] for cooperative UCB yields the upper bound of pulling times for each suboptimal arm as follows,  $\tilde{N}_{k,t}^{(m)} \leq \frac{8 \log T}{\Delta_k^2}$ , which is a standard property for UCB algorithm. With this stochastic global synchronization, we know that in expectation, the global total pulling times of each suboptimal arm  $k$  is at most  $\frac{8 \log T}{\Delta_k^2} + \frac{M}{q}$ . Therefore, the regret for multi-agent multi-armed bandits with the synchronization probability  $q$  is upper bounded as follows,

$$\mathbb{E}[R_T] \leq O\left(\sum_{k \neq k^*} \frac{\log T}{\Delta_k} + \frac{KM}{q}\right).$$

**Step 2:** Under the single cluster scenario, the communication graph  $G_t$  is generated by the stochastic block model (SBM) with the edge probability  $p$ . With a probability  $p^{M^2}$ , all agents in the cluster are connected in the communication graph  $G_t$ , which fulfills the stochastic synchronization with probability  $q = p^{M^2}$ , which concludes the proof.  $\square$

## E.2 Proof of Theorem 2

PROOF. By specifying  $c = \min_{m,n} p(m, n)$  in the Erdos-Renyi model in [58] and by having  $c > \min_{m,n} p(m, n) \geq (\frac{1}{2} + \frac{1}{2} \sqrt{1 - (\frac{\delta}{MT})^{\frac{2}{M-1}}})$  we have that some key results in [58] hold, which are listed as follows.

First we characterize the graph topology related to the random graph (a reduced form of Erdos-Renyi (E-R) model), which utilizes the above assumption on  $p(m, n)$ .

### Graph connectivity

PROPOSITION 1. Assume  $c$  in setting 1 meets the condition

$$1 \geq c \geq \frac{1}{2} + \frac{1}{2} \sqrt{1 - \left(\frac{\epsilon}{MT}\right)^{\frac{2}{M-1}}},$$

where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , for any  $t > 0$ , the graph  $G_t$  following the E-R model is connected.

Next, we present the result regarding the transmission gap, which guarantees that agents can effectively collect one another's reward information within certain time frame with high probability.

### Explicit transmission gap

PROPOSITION 2. We have that with probability  $1 - \epsilon$ , for any  $t > L$  and any  $m$ , there exists

$$t_0 \geq \frac{\ln\left(\frac{\epsilon}{M^2 T}\right)}{\min_{P(i,j)} \ln(1 - P(i, j))}$$

such that

$$t + 1 - \min_j t_{m,j} \leq t_0, t_0 \leq c_0 \min_l n_{l,i}(t + 1)$$

where  $c_0 = c_0(K, \min_{i \neq i^*} \Delta_i, M, \epsilon, \delta)$ .

By the construction of the estimators based on Rule 1, we derive that the global estimator  $\tilde{\mu}_i^m(t)$  is an unbiased estimator of the underlying true global reward mean value.

### Unbiasedness of the estimator

PROPOSITION 3. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any arm  $i$  and any agent  $m$ , at every time step  $t$ , we have

$$E[\tilde{\mu}_i^m(t) | A_{\epsilon, \delta}] = \mu_i.$$

Again based on Rule 2, we prove that the variance of the global estimator  $\tilde{\mu}_i^m(t)$  decays with  $n_{m,i}(t)$  through the characterization of the moment generating function of  $\tilde{\mu}_i^m(t)$ .

### Variance term

PROPOSITION 4. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . In setting  $s_1, s_2, s_3$  where rewards follow sub-gaussian distributions, for any  $m, i, \lambda$  and  $t > L$  where  $L$  is the length of the burn-in period, the global estimator  $\tilde{\mu}_i^m(t)$  is sub-Gaussian distributed. Moreover, the conditional moment generating function satisfies that with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,

$$\begin{aligned} & E[\exp\{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j n_{j,i}(t)}\right\} \end{aligned}$$

where  $\sigma^2 = \max_{j,i} (\tilde{\sigma}_i^j)^2$  and  $C = \max\left\{\frac{4(M+2)(1-\frac{1-c_0}{2(M+2)})^2}{3M(1-c_0)}, (M+2)(1+4Md_{m,t}^2)/M\right\}$ .

The unbiasedness and decaying variance of the global estimator  $\tilde{\mu}_i^m(t)$  allows us to show how much difference is there between  $\tilde{\mu}_i^m(t)$  and the unknown groundtruth  $\mu_i$  through the following concentration inequality.

### Concentration inequality

PROPOSITION 5. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any  $m, i$  and  $t > L$  where  $L$  is the length of the burn-in period,  $\tilde{\mu}_{m,i}(t)$  satisfies that if  $n_{m,i}(t) \geq 2(K^2 + KM + M)$ , then with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,

$$\begin{aligned} P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{n_{m,i}(t)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}, \\ P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{n_{m,i}(t)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta}) t^2}. \end{aligned}$$

Essentially, the above proposition implies that with high probability, we can identify the globally optimal arm by comparing all arms' global estimators  $\tilde{\mu}_i^m(t)$ . Subsequently, we next show that the number of pulling these globally sub-optimal arms can be upper bounded by the  $\log T$  based on the concentration inequality.

### Number of pulls of sub-optimal arms

Upper bounds on  $E[n_{m,k}(T)|A_{\epsilon,\delta}]$

**PROPOSITION 6.** *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . An arm  $k$  is said to be sub-optimal if  $k \neq i^*$  where  $i^*$  is the unique optimal arm in terms of the global reward, i.e.  $i^* = \arg \max \frac{1}{M} \sum_{j=1}^M \mu_i^j$ . Then when the game ends, for every agent  $m$ ,  $0 < \epsilon < 1$  and  $T > L$ , the expected numbers of pulling sub-optimal arm  $k$  after the burn-in period satisfies with  $P(A_{\epsilon,\delta}) = 1 - 7\epsilon$*

$$\begin{aligned} & E[n_{m,k}(T)|A_{\epsilon,\delta}] \\ & \leq \max \left\{ \left\lceil \frac{4C_1 \log T}{\Delta_i^2} \right\rceil, 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K \\ & \leq O(\log T). \end{aligned}$$

The proof of Proposition 1 - 6 is presented in Appendix in [58] with Erdos-Renyi Models. As a result, we do not repeat the proof details of them herein and refer the proof steps therein.

With these key results, we proceed to bound the regret.

### Regret decomposition

The optimal arm is denoted as  $i^*$  satisfying

$$i^* = \arg \max_i \sum_{m=1}^M \mu_i^m.$$

For the proposed regret, we have that for any constant  $L$ ,

$$\begin{aligned}
R_T &= \frac{1}{M} (\max_i \sum_{t=1}^T \sum_{m=1}^M \mu_i^m - \sum_{t=1}^T \sum_{m=1}^M \mu_{a_t^m}^m) \\
&= \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \\
&\leq \sum_{t=1}^L \left| \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right| + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&\leq L + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&= L + \sum_{t=L+1}^T (\mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m) \\
&= L + ((T-L) \cdot \mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^K n_{m,i}(T) \mu_i^m)
\end{aligned}$$

where the first inequality is by taking the absolute value and the second inequality results from the assumption that  $0 < \mu_i^j < 1$  for any arm  $i$  and agent  $j$ .

Note that  $\sum_{i=1}^K \sum_{m=1}^M n_{m,i}(T) = M(T-L)$  where by definition  $n_{m,i}(T)$  is the number of pulls of arm  $i$  at agent  $m$  from time step  $L+1$  to time step  $T$ , which yields that

$$\begin{aligned}
R_T &\leq L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_{i^*}^m - \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_i^m \\
&= L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&\leq L + \frac{1}{M} \sum_{i=1}^K \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&= L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m).
\end{aligned}$$

where the second inequality uses the fact that  $\sum_{m: \mu_{i^*}^m - \mu_i^m \leq 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \leq 0$  holds for any arm  $i$  and the last equality is true since  $n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) = 0$  for  $i = i^*$  and any  $m$ .

By using Proposition 5 and the fact that

$$R_T \leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m),$$



we derive that

$$\begin{aligned} E[R_T | A_{\epsilon, \delta}] &\leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} E[n_{m,i}(T)] (\mu_{i^*}^m - \mu_i^m) \\ &\leq L + \sum_{i \neq i^*} \Delta_i (\max \{ \lceil \frac{4C_1 \log T}{\Delta_i^2} \rceil, 2(K^2 + MK) \}) + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K \end{aligned}$$

which completes the proof.  $\square$

### E.3 Proof of Theorem 3

PROOF. It is worth noting that our framework implies that only the agents in different clusters have different reward distributions, and the agents in the same cluster stay on the same page by the assumption that  $p(m, m) = 1$ . This indicates that as long as there is an edge between one agent in one cluster and another agent in another agents, the two cluster can exchange the heterogeneous reward distributions. Consequently, it is sufficient to pay attention to the sub-graph with respect to the clusters.

To this end, we establish the following proposition regarding the connectivity of the sub-graph. It is worth noting that the edge probability of this sub-graph, is now  $c = 1 - (1 - p(m, n))^{M^2/C^2}$  by Lemma 4, and the total number of vertex is  $C$  instead of  $M$ .

#### Graph connectivity

PROPOSITION 7. *Assume  $c$  meets the condition*

$$1 \geq c \geq \frac{1}{2} + \frac{1}{2} \sqrt{1 - \left(\frac{\epsilon}{CT}\right)^{\frac{2}{M-1}}},$$

where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , for any  $t > 0$ , the sub-graph  $G_t^C$  following the E-R model is connected.

PROOF OF PROPOSITION 7. The proof of Proposition 7 follows from that of Proposition 1 in Theorem 2, based on [58].  $\square$

Then based on the newly proposed estimator construction as in Rule 2, we have the following results. First, we characterize the consensus regarding arm pulls among the clusters, instead of the agents, since now the communication is on a cluster level.

#### Information delay

LEMMA 16. *For any  $m, i, t > L$ , if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$  and subgraph  $G_t$  induced by the clusters is connected, then we have*

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

where the min is taken over all clusters, not just the neighbors.

PROOF OF LEMMA 7. The proof of this lemma follows from Lemma 3 in [64], with the exception that now the shared arm information is  $N_{m,i}(t)$  instead of  $n_{m,i}(t)$ .

□

Then, we show that the transmission gap with respect to the sub-graph still holds, that clusters can effectively collect one another's reward information within certain time frame with high probability. It is worth noting that when  $p(m, m) = 1$ , the clusters in one cluster obtain such information as well.

### Explicit transmission gap

PROPOSITION 8. *We have that with probability  $1 - \epsilon$ , for any  $t > L$  and any  $m$ , there exists*

$$t_0 \geq \frac{\ln\left(\frac{\epsilon}{M^2 T}\right)}{\min_{P(i,j)} \ln(1 - P(i, j))}$$

such that

$$t + 1 - \min_j t_{m,j} \leq t_0, t_0 \leq c_0 \min_l N_{l,i}(t + 1)$$

where  $c_0 = c_0(K, \min_{i \neq i^*} \Delta_i, M, \epsilon, \delta)$ .

PROOF OF PROPOSITION 8. The proof of this proposition follows from [58], except that the shared information is  $N_{l,i}(t + 1)$  instead of  $n_{l,i}(t + 1)$  since we consider the sub-graph with respect to the clusters (and there is no delay within the cluster).

□

It is straightforward to verify that by the construction of the global estimators based on Rule 2, we have that the global estimator  $\tilde{\mu}_i^m(t)$  is an unbiased estimator of  $\mu_i$ .

### Unbiasedness of the estimator

PROPOSITION 9. *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any arm  $i$  and any agent  $m$ , at every time step  $t$ , we have*

$$E[\tilde{\mu}_i^m(t) | A_{\epsilon, \delta}] = \mu_i.$$

PROOF OF PROPOSITION 9. The proof of this proposition follows from [58].

□

To proceed, we examine the variance of the global estimator  $\tilde{\mu}_i^m(t)$  constructed by Rule through the moment generating function. It is worth noting that the variance is smaller compared to the one based on Rule 1, since we consider cluster-wise information in the decision making and communicate on a cluster-level. More specifically, the upper bound on the moment generating function changes from  $\exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j n_{j,i}(t)}\right\}$  to  $\exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t)}\right\}$ , where  $N_{j,i} = \sum_{m \in c_j} n_{j,i}$  and thus achieves sample complexity reduction.

### Variance term

PROPOSITION 10. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . In setting  $s_1, s_2, s_3$  where rewards follow sub-gaussian distributions, for any  $m, i, \lambda$  and  $t > L$  where  $L$  is the length of the burn-in period, the global estimator  $\tilde{\mu}_i^m(t)$  is sub-Gaussian distributed. Moreover, the conditional moment generating function satisfies that with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,

$$\begin{aligned} & E[\exp \{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \exp \left\{ \frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t)} \right\} \end{aligned}$$

where  $\sigma^2 = \max_{j,i} (\tilde{\sigma}_i^j)^2$  and  $C = \max \left\{ \frac{4(M+2)(1-\frac{1-c_0}{2(M+2)})^2}{3M(1-c_0)}, (M+2)(1+4Md_{m,t}^2)/M \right\}$ .

PROOF OF PROPOSITION 10. The proof is done by induction as in [58].

Based on our construction of  $A'_{\epsilon, \delta}$  and the choice of  $\delta$ , we have that for  $t \geq L$ ,  $|P_t - cE| < \delta < c$  on event  $A_{\epsilon, \delta}$ . It implies that for any  $t \geq L$ ,  $m$  and  $j$ ,  $P_t(m, j) > 0$ , and if  $t = L$

$$P'_t(m, j) = \frac{1}{M} \quad (4)$$

and if  $t > L$

$$P'_t(m, j) = \frac{M-1}{M^2}. \quad (5)$$

Consider the time step  $t \leq L+1$ . The quantity satisfies that (following from [58])

$$\begin{aligned} & E[\exp \{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & = E[\exp \{\lambda(\tilde{\mu}_i^m(L+1) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & = E[\exp \{\lambda(\sum_{j=1}^M P'_{m,j}(L) \hat{\mu}_{i,j}^m(t_{m,j}) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & = E[\exp \{\lambda(\sum_{j=1}^M \frac{1}{M} \hat{\mu}_{i,j}^m(t_{m,j}) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & = E[\exp \{\lambda(\sum_{j=1}^M \frac{1}{M} (\hat{\mu}_{i,j}^m(t_{m,j}) - \mu_i^j)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \Pi_{j=1}^M (E[(\exp \{\lambda(\frac{1}{M} (\hat{\mu}_{i,j}^m(t_{m,j}) - \mu_i^j)\} 1_{A_{\epsilon, \delta}})^M | \sigma(\{n_{m,i}(t)\}_{t,i,m})])^{\frac{1}{M}} \end{aligned} \quad (6)$$

where the third equality holds by (4), the fourth equality uses the definition  $\mu_i = \frac{1}{M} \sum_{i=1}^M \mu_i^j$ , and the last inequality results from the generalized hoeffding inequality and the fact that  $\hat{\mu}_{i,j}^m(t_{m,j}) = \tilde{\mu}_i^j(t_{m,j})$ .

Note that for any agent  $j$ , we have

$$\begin{aligned}
& E[(\exp\{(\lambda \frac{1}{M}(\hat{\mu}_i^j(t_{m,j}) - \mu_i^j)\}1_{A_{\epsilon,\delta}})^M | \sigma(\{n_{m,i}(t)\}_{t,i,m})\})] \\
&= E[\exp\{(\lambda(\hat{\mu}_i^j(t_{m,j}) - \mu_i^j)\}1_{A_{\epsilon,\delta}}) | \sigma(\{n_{m,i}(t)\}_{t,i,m})\}] \\
&= E[\exp\{(\lambda \frac{\sum_s (r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})}\}1_{A_{\epsilon,\delta}}) | \sigma(\{n_{m,i}(t)\}_{t,i,m})\}] \\
&= E[\exp\{\sum_s (\lambda \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})}\}1_{A_{\epsilon,\delta}}) | \sigma(\{n_{m,i}(t)\}_{t,i,m})\}]. \tag{7}
\end{aligned}$$

We observe that based on the reward generation mechanism, given  $s$ ,  $r_i^j(s)$  does not depend on anything else, which implies that

$$\begin{aligned}
& E[(\exp\{(\lambda \frac{1}{M}(\hat{\mu}_i^j(t_{m,j}) - \mu_i^j)\}1_{A_{\epsilon,\delta}})^M | \sigma(\{n_{m,i}(t)\}_{t,i,m})\})] \\
&= \Pi_s E[\exp\{\lambda \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})}\}1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})\}] \\
&= \Pi_s E[\exp\{\lambda \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})}\} | \sigma(\{N_{m,i}(t)\}_{t,i,m})\}] \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})\}] \\
&= \Pi_s E_r[\exp\{\lambda \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})}\}] \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{N_{m,i}(t)\}_{t,i,m})\}] \\
&\leq \Pi_s \exp\{\frac{(\frac{\lambda}{N_{j,i}(t_{m,j})})^2 \sigma^2}{2}\} \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{N_{m,i}(t)\}_{t,i,m})\}] \\
&\leq (\exp\{\frac{(\frac{\lambda}{N_{j,i}(t_{m,j})})^2 \sigma^2}{2}\})^{N_{j,i}(t_{m,j})} \\
&= \exp\{\frac{\lambda^2 \sigma^2}{2 N_{j,i}(t_{m,j})}\} \\
&\leq \exp\{\frac{\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t_{m,j})}\} \tag{8}
\end{aligned}$$

where the first inequality holds by the definition of sub-Gaussian random variables  $r_i^j(s) - \mu_i^j$  with an mean value 0, the second inequality results from  $1_{A_{\epsilon,\delta}} \leq 1$ , and the last inequality uses  $N_{j,i}(t_{m,j}) \geq \min_j N_{j,i}(t_{m,j})$  for any  $j$ .

Therefore, we obtain that

$$\begin{aligned}
(6) &\leq \Pi_{j=1}^M (\exp\{\frac{\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t_{m,j})}\})^{\frac{1}{M}} \\
&= ((\exp\{\frac{\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t_{m,j})}\})^{\frac{1}{M}})^M \\
&= \exp\{\frac{\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t_{m,j})}\}
\end{aligned}$$

which concludes the basis step.

Now we proceed to the induction step. Let us assume that for any  $s < t + 1$  where  $t \geq L$ , the following holds

$$\begin{aligned} & E[\exp \{\lambda(\tilde{\mu}_i^m(s) - \mu_i)\} 1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(s)\}_{s,i,m})] \\ & \leq \exp \left\{ \frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(s)} \right\}. \end{aligned} \quad (9)$$

By the derivation of (25) in [58], we derive that

$$\begin{aligned} & E[\exp \{\lambda(\tilde{\mu}_i^m(t+1) - \mu_i)\} 1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(s)\}_{s,i,m})] \\ & \leq \prod_{j=1}^M \left( \exp \left\{ \frac{\lambda^2 (P'_t(m,j))^2 (M+2)^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t_{m,j})} \right\} \right)^{\frac{1}{M+2}} \cdot \\ & \quad \prod_{j \in N_m(t)} \Pi_s \left( E_r \left[ \exp \left\{ \lambda d_{m,t} (M+2) \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t)} \right\} \right] \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \right)^{\frac{1}{M+2}} \cdot \\ & \quad \prod_{j \notin N_m(t)} \Pi_s \left( E_r \left[ \exp \left\{ \lambda d_{m,t} (M+2) \frac{(r_i^j(s) - \mu_i^j)}{N_{j,i}(t_{m,j})} \right\} \right] \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \right)^{\frac{1}{M+2}} \end{aligned} \quad (10)$$

Meanwhile, we derive the following bound for the last two terms by the sub-Gaussian property of  $(r_i^j(s) - \mu_i^j)$

$$\begin{aligned} & E[\exp \{\lambda(\tilde{\mu}_i^m(t+1) - \mu_i)\} 1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(s)\}_{s,i,m})] \\ & \leq \left( \exp \left\{ \frac{\lambda^2 (P'_t(m,j))^2 (M+2)^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t_{m,j})} \right\} \right)^{\frac{M}{M+2}} \cdot \\ & \quad \prod_{j \in N_m(t)} \Pi_s \left( \exp \frac{\lambda^2 d_{m,t}^2 (M+2)^2 \sigma^2}{2N_{j,i}^2(t)} \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \right)^{\frac{1}{M+2}} \cdot \\ & \quad \prod_{j \notin N_m(t)} \Pi_s \left( \exp \frac{\lambda^2 d_{m,t}^2 (M+2)^2 \sigma^2}{2N_{j,i}^2(t_{m,j})} \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{N_{m,i}(t)\}_{t,i,m})] \right)^{\frac{1}{M+2}} \\ & = \left( \exp \left\{ \frac{\lambda^2 (P'_t(m,j))^2 (M+2)^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t_{m,j})} \right\} \right)^{\frac{M}{M+2}} \cdot \\ & \quad \prod_{j \in N_m(t)} \exp \left\{ \frac{N_{j,i}(t)}{M+2} \frac{\lambda^2 d_{m,t}^2 (M+2)^2 \sigma^2}{2N_{j,i}^2(t)} \right\} \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{N_{m,i}(t)\}_{t,i,m})] \cdot \\ & \quad \prod_{j \notin N_m(t)} \exp \left\{ \frac{N_{j,i}(t_{m,j})}{M+2} \frac{\lambda^2 d_{m,t}^2 (M+2)^2 \sigma^2}{2N_{j,i}^2(t_{m,j})} \right\} \cdot E[1_{A_{\epsilon,\delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \end{aligned}$$

Subsequently, we obtain

$$\begin{aligned}
& E[\exp \{ \lambda (\tilde{\mu}_i^m(t+1) - \mu_i) \} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(s)\}_{s,i,m}) ] \\
& \leq (\exp \{ \frac{\lambda^2 (P'_t(m, j))^2 (M+2)^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t_{m,j})} \})^{\frac{M}{M+2}} \cdot \\
& \quad (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)\sigma^2}{2 \min_j N_{j,i}(t)} \})^{|N_m(t)|} \cdot E[1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ] \cdot \\
& \quad (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)\sigma^2}{2 \min_j N_{j,i}(t_{m,j})} \})^{|M-N_m(t)|} \cdot E[1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ] \\
& = E[(\exp \{ \frac{\lambda^2 (P'_t(m, j))^2 M(M+2)}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t_{m,j})} \}) \cdot (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)|N_m(t)|}{2 \min_j N_{j,i}(t)} \}) \\
& \quad \cdot (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)\sigma^2 |M-N_m(t)|}{2 \min_j N_{j,i}(t_{m,j})} \}) 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ] \\
& \leq E[(\exp \{ \frac{\lambda^2 (P'_t(m, j))^2 M(M+2)}{2(1-c_0)} \frac{C\sigma^2}{\min_j N_{j,i}(t+1)} \}) \cdot (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)|N_m(t)|\sigma^2}{2 \frac{L/K}{L/K+1} \min_j N_{j,i}(t+1)} \}) \\
& \quad \cdot (\exp \{ \frac{\lambda^2 d_{m,t}^2 (M+2)|M-N_m(t)|\sigma^2}{2(1-c_0) \min_j N_{j,i}(t+1)} \}) 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ]
\end{aligned}$$

After organizing the terms in the above objective, we have

$$\begin{aligned}
& E[\exp \{ \lambda (\tilde{\mu}_i^m(t+1) - \mu_i) \} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(s)\}_{s,i,m}) ] \\
& = E[(\exp \{ \frac{\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t+1)} \cdot (\frac{C(P'_t(m, j))^2 M(M+2)}{2(1-c_0)} + \\
& \quad \frac{d_{m,t}^2 (M+2)|N_m(t)|}{\frac{L/K}{L/K+1}} + \frac{d_{m,t}^2 (M+2)|M-N_m(t)|}{(1-c_0)}) \} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ] \\
& \leq E[\exp \{ \frac{C\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t+1)} \} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m}) ] \\
& \leq \exp \{ \frac{C\lambda^2 \sigma^2}{2 \min_j N_{j,i}(t+1)} \}
\end{aligned}$$

where the first inequality is true because of the specification of the parameters  $P'_t(m, j)$ ,  $d_{m,t}$ ,  $L$ ,  $c_0$  and  $C$  and the second inequality holds true by the observation that  $1_{A_{\epsilon, \delta}} \leq 1$  and  $\min_j N_{j,i}(t+1) \in \sigma(\{n_{m,i}(t)\}_{t,i,m})$ .

The completion of this induction step subsequently completes the proof of Proposition 10.  $\square$

As in the proof of Theorem 2, we next show how much difference between  $\tilde{\mu}_i^m(t)$  and  $\mu_i$  by establishing the following concentration inequality.

### Concentration inequality

PROPOSITION 11. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any  $m, i$  and  $t > L$  where  $L$  is the length of the burn-in period,  $\tilde{\mu}_{m,i}(t)$  satisfies that if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$ , then with  $P(A_{\epsilon,\delta}) = 1 - 7\epsilon$ ,

$$P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon,\delta}) \leq \frac{1}{P(A_{\epsilon,\delta})} \frac{1}{t^2},$$

$$P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon,\delta}) \leq \frac{1}{P(A_{\epsilon,\delta}) t^2}.$$

PROOF OF PROPOSITION 11. The proof of this proposition follows from [58].

□

Essentially, the above proposition implies that with high probability, we can identify the globally optimal arm by comparing all arms' global estimators  $\tilde{\mu}_i^m(t)$  with smaller sample complexity (by having  $N_{m,i}(t)$  instead of  $n_{m,i}(t)$ ). Subsequently, we next show that the number of pulling these globally sub-optimal arms can be upper bounded by the  $\log T$  based on the concentration inequality.

### Number of pulls of sub-optimal arms

Upper bounds on  $E[n_{m,k}(T) | A_{\epsilon,\delta}]$

PROPOSITION 12. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . An arm  $k$  is said to be sub-optimal if  $k \neq i^*$  where  $i^*$  is the unique optimal arm in terms of the global reward, i.e.  $i^* = \arg \max \frac{1}{M} \sum_{j=1}^M \mu_i^j$ . Then when the game ends, for every agent  $m$ ,  $0 < \epsilon < 1$  and  $T > L$ , the expected numbers of pulling sub-optimal arm  $k$  after the burn-in period satisfies with  $P(A_{\epsilon,\delta}) = 1 - 7\epsilon$

$$E[n_{m,k}(T) | A_{\epsilon,\delta}]$$

$$\leq \max \left\{ \frac{C}{M} \cdot \left[ \frac{4C_1 \log T}{\Delta_i^2} \right], 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K$$

$$\leq O(\log T).$$

PROOF OF PROPOSITION 12. It is easy to observe that by repeating the proof steps of Proposition 6 with  $N_{m,i}(t)$  instead of  $n_{m,i}(t)$ , we obtain that

$$E[N_{m,k}(T) | A_{\epsilon,\delta}]$$

$$\leq \max \left\{ \left[ \frac{4C_1 \log T}{\Delta_i^2} \right], 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K$$

$$\leq O(\log T).$$

It is worth noting that the decision rule of each agent relies on the cluster-wise information, i.e.  $\tilde{\mu}_{m,i}(t)$  and  $N_{m,i}(t)$ , which is the same for all agents within one cluster. That being said, the agents within one cluster are pulling the same arm, formally written as

$$N_{m,k}(T) = |c_M| n_{m,i}(t) = \frac{M}{C} n_{m,i}(t)$$

since we assume a balanced cluster structure.

Subsequently, we derive that

$$\begin{aligned}
& E[n_{m,k}(T)|A_{\epsilon,\delta}] \\
& \leq E\left[\frac{C}{M} \cdot N_{m,k}(T)|A_{\epsilon,\delta}\right] \leq \max\left\{\frac{C}{M} \cdot \left[\frac{4C_1 \log T}{\Delta_i^2}\right], 2(K^2 + MK + M)\right\} + \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K \\
& \leq O(\log T).
\end{aligned}$$

which concludes the proof of Proposition 12. □

Now we are ready to proceed to the proof of the main theorem by characterizing the regret. Likewise, we again perform regret decomposition.

### Regret decomposition

For the proposed regret, we have that for any constant  $L$ ,

$$\begin{aligned}
R_T &= \frac{1}{M} (\max_i \sum_{t=1}^T \sum_{m=1}^M \mu_i^m - \sum_{t=1}^T \sum_{m=1}^M \mu_{a_t^m}^m) \\
&= \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \\
&\leq \sum_{t=1}^L \left| \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right| + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&\leq L + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&= L + \sum_{t=L+1}^T \left( \mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&= L + ((T - L) \cdot \mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^K n_{m,i}(T) \mu_i^m)
\end{aligned}$$

where the first inequality is by taking the absolute value and the second inequality results from the assumption that  $0 < \mu_i^j < 1$  for any arm  $i$  and agent  $j$ .



Note that  $\sum_{i=1}^K \sum_{m=1}^M n_{m,i}(T) = M(T-L)$  where by definition  $n_{m,i}(T)$  is the number of pulls of arm  $i$  at agent  $m$  from time step  $L+1$  to time step  $T$ , which yields that

$$\begin{aligned}
R_T &\leq L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_{i^*}^m - \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_i^m \\
&= L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&\leq L + \frac{1}{M} \sum_{i=1}^K \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&= L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m).
\end{aligned}$$

where the second inequality uses the fact that  $\sum_{m: \mu_{i^*}^m - \mu_i^m \leq 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \leq 0$  holds for any arm  $i$  and the last equality is true since  $n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) = 0$  for  $i = i^*$  and any  $m$ .

By using Proposition 12 and the above inequality, we have that

$$R_T \leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m),$$

Consequently, we obtain that

$$\begin{aligned}
E[R_T | A_{\epsilon, \delta}] &\leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} E[n_{m,i}(T)] (\mu_{i^*}^m - \mu_i^m) \\
&\leq L + \sum_{i \neq i^*} \Delta_i (\max \left\{ \frac{C}{M} \cdot \left[ \frac{4C_1 \log T}{\Delta_i^2} \right], 2(K^2 + MK) \right\} + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M-1)K)
\end{aligned}$$

which concludes the proof.  $\square$

#### E.4 Proof of Theorem 5

PROOF. We first demonstrate that as long as the following lower bound on  $p(m, n)$  holds, then we also have the claims about the connectivity of the sub-graphs and the claim about the delay induced by the random sub-graph.

To this end, we establish the following proposition regarding the connectivity of the sub-graph. It is worth noting that the edge probability of this sub-graph, is now  $c = \frac{e}{e-1} \frac{M^2}{C^2} \cdot \min_{m \neq n} P(m, n)$  based on Lemma 6, and the total number of vertex is  $C$  instead of  $M$ .

We start with the graph connectivity.

##### Graph connectivity

PROPOSITION 13. *Assume  $c$  meets the condition*

$$1 \geq c \geq \left(1 - \frac{\delta(C-1)}{8CT}\right),$$

where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , for any  $t > 0$ , the sub-graph  $G_t^C$  following the E-R model is connected.

PROOF OF PROPOSITION 13. We would like to highlight that this result can help solving the existing problems on random graphs, not limited to the bandit setting studied herein. It is worth noting that for any two clusters  $j_1 \neq j_2 \neq m$ ,  $1_{E_{m,j_1}}$  and  $1_{(m,j_2) \in E_t^C}$  are identical random variables. Meanwhile, we note that the variance of  $1_{(m,j_1) \in E_t^C}$  is no more than 1. Subsequently, based on the Chebyshev's inequality, we obtain

$$\begin{aligned}
& P(d_m \leq \frac{C-1}{2}) \\
&= P(d_m - A_{M-2}^{l-1} p^l (M-1) \leq \frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1)) \\
&\leq P((d_m - A_{M-2}^{l-1} p^l (M-1))^2 \geq (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2) \\
&\leq \frac{\text{Var}(d_m)}{(C-1)^2 (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2} \\
&\leq \frac{(C-1)^2 \cdot A_{M-2}^{l-1} p^l (1 - A_{M-2}^{l-1} p^l)}{(C-1)^2 (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2} \\
&= \frac{1}{(\frac{1}{2} - A_{M-2}^{l-1} p^l)^2} \cdot (1 - A_{M-2}^{l-1} p^l) \\
&\leq 8 \cdot (1 - A_{M-2}^{l-1} p^l) \\
&\leq \frac{\delta}{T}
\end{aligned}$$

when we specify that  $p \geq (1 - \frac{\delta(C-1)}{8T})$ .

Hence, we have

$$P(d_m \leq \frac{C-1}{2}) \leq \frac{\delta}{T}. \quad (11)$$

In other words, with probability at least  $1 - \frac{\delta}{T}$ , we have that  $d_m > \frac{C-1}{2}$ .

It is well known that if  $\delta(G_t) \geq \frac{C-1}{2}$ , then we have that the sub-graph  $G_t^C$  is connected where  $\delta(G_t) = \min_m d_m$ .

As a result, consider the probability and we obtain that

$$\begin{aligned}
& P(\text{graph } G_{t-l+1} \cdots G_t \text{ is connected}) \\
& \geq P(\min_j d_j \geq \frac{C-1}{2}) \\
& = P(\bigcap_j \{d_j \geq \frac{C-1}{2}\}) \\
& = 1 - P(\bigcup_j \{d_j < \frac{C-1}{2}\}) \\
& \geq 1 - \sum_j P(d_j < \frac{C-1}{2}) \\
& = 1 - CP(d_j < \frac{C-1}{2}) \\
& \geq 1 - C\frac{\delta}{T} = 1 - \frac{C\delta}{T}
\end{aligned}$$

where the second inequality holds by the Bonferroni's inequality and the third inequality uses (11).

Consequently, we obtain

$$\begin{aligned}
& P(\text{graph } G_t \text{ is connected}) \\
& = P(\bigcap_t \{G_t \text{ is connected}\}) \\
& \geq 1 - \sum_t P(G_t \text{ is not connected}) \\
& = 1 - \sum_t (1 - P(G_t \text{ is connected})) \\
& \geq 1 - \sum_t (1 - (1 - \frac{C\delta}{T})) = 1 - C\delta
\end{aligned}$$

This indicates that with probability at least  $1 - \epsilon$ , for any time step  $t$ , the corresponding composition graph is connected, which concludes the proof of Proposition 13.  $\square$

Henceforth, we derive the following claim about the graph connectivity, by combining Proposition 13 with Proposition 7, which reads as follows.

**PROPOSITION 14.** *Assume  $c$  meets the condition*

$$1 \geq c \geq \min \left\{ \frac{1}{2} + \frac{1}{2} \sqrt{1 - \left(\frac{\epsilon}{CT}\right)^{\frac{2}{M-1}}}, \left(1 - \frac{\delta(C-1)}{8CT}\right) \right\},$$

*i.e.*

$$1 \geq \min_{m,n} p(m, n) \geq \frac{C^2}{M^2} \min \left\{ \frac{1}{2} + \frac{1}{2} \sqrt{1 - \left(\frac{\epsilon}{CT}\right)^{\frac{2}{M-1}}}, \left(1 - \frac{\delta(C-1)}{8CT}\right) \right\},$$

*where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , for any  $t > 0$ , the sub-graph  $G_t^c$  following the E-R model is connected.*

**PROOF OF PROPOSITION 14.** This is a direct result of merging Proposition 7 and Proposition 13.  $\square$

With the characterization of the graph topology, we next demonstrate the consensus regarding arm pulls among the clusters, instead of the agents, since now the communication is on a cluster level. This is determined by how much information delay we have across the clusters (again, within one cluster, there is no such information delay).

### Information delay

LEMMA 17. *For any  $m, i, t > L$ , if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$  and subgraph  $G_t$  induced by the clusters is connected, then we have*

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

where the min is taken over all clusters, not just the neighbors.

PROOF OF LEMMA 7. The proof of this lemma again follows from Lemma 3 in [64], with the exception that now the shared arm information is  $N_{m,i}(t)$  instead of  $n_{m,i}(t)$ . □

To proceed, we observe that since there are no modifications to the algorithm, the following results on the explicit transmission gap, unbiasedness of the estimator, variance of the estimator, concentration inequality, number of pulls of sub-optimal arms, hold. As a result, we omit the proof steps and refer to the proof of Theorem 3 for details.

### Explicit transmission gap

PROPOSITION 15. *We have that with probability  $1 - \epsilon$ , for any  $t > L$  and any  $m$ , there exists*

$$t_0 \geq \frac{\ln\left(\frac{\epsilon}{M^2 T}\right)}{\min_{P(i,j)} \ln(1 - P(i, j))}$$

such that

$$t + 1 - \min_j t_{m,j} \leq t_0, t_0 \leq c_0 \min_l N_{l,i}(t + 1)$$

where  $c_0 = c_0(K, \min_{i \neq i^*} \Delta_i, M, \epsilon, \delta)$ .

### Unbiasedness of the estimator

PROPOSITION 16. *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any arm  $i$  and any agent  $m$ , at every time step  $t$ , we have*

$$E[\tilde{\mu}_i^m(t) | A_{\epsilon, \delta}] = \mu_i.$$

### Variance term

PROPOSITION 17. *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . In setting  $s_1, s_2, s_3$  where rewards follow sub-gaussian distributions, for any  $m, i, \lambda$  and  $t > L$  where  $L$  is the length of the burn-in period, the global estimator  $\tilde{\mu}_i^m(t)$  is sub-Gaussian distributed. Moreover, the conditional moment generating function satisfies that with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,*

$$\begin{aligned} & E[\exp\{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t)}\right\} \end{aligned}$$

where  $\sigma^2 = \max_{j,i} (\tilde{\sigma}_i^j)^2$  and  $C = \max\{\frac{4(M+2)(1-\frac{1-c_0}{2(M+2)})^2}{3M(1-c_0)}, (M+2)(1+4Md_{m,t}^2)/M\}$ .

### Concentration inequality

PROPOSITION 18. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any  $m, i$  and  $t > L$  where  $L$  is the length of the burn-in period,  $\tilde{\mu}_{m,i}(t)$  satisfies that if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$ , then with  $P(A_{\epsilon,\delta}) = 1 - 7\epsilon$ ,

$$P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon,\delta}) \leq \frac{1}{P(A_{\epsilon,\delta})} \frac{1}{t^2},$$

$$P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon,\delta}) \leq \frac{1}{P(A_{\epsilon,\delta}) t^2}.$$

### Number of pulls of sub-optimal arms

Upper bounds on  $E[n_{m,k}(T) | A_{\epsilon,\delta}]$

PROPOSITION 19. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . An arm  $k$  is said to be sub-optimal if  $k \neq i^*$  where  $i^*$  is the unique optimal arm in terms of the global reward, i.e.  $i^* = \arg \max \frac{1}{M} \sum_{j=1}^M \mu_i^j$ . Then when the game ends, for every agent  $m$ ,  $0 < \epsilon < 1$  and  $T > L$ , the expected numbers of pulling sub-optimal arm  $k$  after the burn-in period satisfies with  $P(A_{\epsilon,\delta}) = 1 - 7\epsilon$

$$E[n_{m,k}(T) | A_{\epsilon,\delta}]$$

$$\leq \max\left\{\frac{C}{M} \cdot \left[\frac{4C_1 \log T}{\Delta_i^2}, 2(K^2 + MK + M)\right], \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K\right\}$$

$$\leq O(\log T).$$

Lastly, by the aforementioned regret decomposition, we obtain

$$E[R_T | A_{\epsilon,\delta}] \leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} E[n_{m,i}(T)] (\mu_{i^*}^m - \mu_i^m)$$

$$\leq L + \sum_{i \neq i^*} \Delta_i \left( \max\left\{\frac{C}{M} \cdot \left[\frac{4C_1 \log T}{\Delta_i^2}, 2(K^2 + MK)\right], \frac{2\pi^2}{3P(A_{\epsilon,\delta})} + K^2 + (2M - 1)K\right\} \right)$$

which concludes the proof.  $\square$

### E.5 Proof of Theorem 7

PROOF. We prove the result following a similar analytical order.

We start by examining the graph connectivity, with respect to the sub-graph induced by the clusters. Instead of requiring that the sub-graph is connected for every time step, it holds true that as long as the sub-graph is  $l$ -periodically connected, we guarantee the same consensus among the clusters (represented by the information delay). And the assumption on  $\min_{m,n} p(m, n)$  in Theorem

5 guarantees that with high probability, the sub-graph is  $l$ -periodically connected at every time step, which is formally presented as follows.

### $l$ -periodically connectivity

PROPOSITION 20. *Let us assume that  $\min_{m,n} p(m, n) \geq \frac{C^2}{M^2} \max \left\{ \frac{(C-l-1)!}{(C-2)!} \left(1 - \frac{\delta(C-1)}{8CT}\right), \frac{(C-l-1)!}{(C-2)!} \left(\frac{3}{4}\right)^{\frac{1}{l}} \right\}$ . Then with probability at least  $1 - \delta$ , for any  $t$ , the sequence starting  $G_t^C$  is a  $l$ -periodically connected graph.*

PROOF OF PROPOSITION 20. We would like to highlight that this result can help solving the existing problems on random graphs, not limited to the bandit setting studied herein. The probability of having  $l$  periodically connected graph, i.e. the composition of  $G_1, G_2, \dots, G_l$  is a connected graph, which means that for any two clusters  $m, n$ , the probability of having a path among them during the  $l$  steps. Formally, let us define  $E_{m,n} = \exists a_1, a_2, \dots, a_{l-1}, s.t. 1_{(m,a_1) \in G_1, (a_1,a_2) \in G_2, \dots, (a_{l-1},n) \in G_l} = 1$

Let us denote  $p = c$  as  $\frac{M^2}{C^2} p_{m,n}$ , which again represents the edge probability of the sub-graph.

$$\begin{aligned} P(E_{m,n}) &= P(\exists a_1, a_2, \dots, a_{l-1}, s.t. 1_{(m,a_1) \in G_1, (a_1,a_2) \in G_2, \dots, (a_{l-1},n) \in G_l} = 1) \\ &= P(\exists a_1, a_2, \dots, a_{l-1}, s.t. 1_{(m,a_1) \in G_1} = 1, 1_{(a_1,a_2) \in G_2} = 1, \dots, 1_{(a_{l-1},n) \in G_l} = 1) \\ &= A_{M-2}^{l-1} p^l \end{aligned}$$

Let us define the degree of cluster  $m$  as  $d_m$ . We then derive that

$$\begin{aligned} E[d_m] &= E\left[\sum_{n=1}^C 1_{n \neq m} \cdot 1_{E_{m,n}}\right] \\ &= (C-1) \cdot A_{M-2}^{l-1} p^l \end{aligned}$$

Likewise, we derive that

$$\begin{aligned} Var(d_m) &\leq (C-1) \sum_{n=1}^C Var(1_{n \neq m} \cdot 1_{E_{m,n}}) \\ &\leq (C-1)^2 \cdot A_{M-2}^{l-1} p^l (1 - A_{M-2}^{l-1} p^l) \end{aligned}$$

Let us assume that

$$A_{M-2}^{l-1} p^l \geq \frac{3}{4}$$

which implies that  $p \geq \left(\frac{4}{3A_{M-2}^{l-1}}\right)^{\frac{1}{l}}$ .

It is worth noting that for any two clusters  $j_1 \neq j_2 \neq m$ ,  $1_{E_{m,j_1}}$  and  $1_{E_{m,j_2}}$  are dependent but identical random variables. Meanwhile, we note that the variance of  $1_{E_{m,j_1}}$  is no more than 1. Subsequently,

based on the Chebyshev's inequality, we obtain

$$\begin{aligned}
& P(d_m \leq \frac{C-1}{2}) \\
&= P(d_m - A_{M-2}^{l-1} p^l (M-1) \leq \frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1)) \\
&\leq P((d_m - A_{M-2}^{l-1} p^l (M-1))^2 \geq (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2) \\
&\leq \frac{\text{Var}(d_m)}{(C-1)^2 (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2} \\
&\leq \frac{(C-1)^2 \cdot A_{M-2}^{l-1} p^l (1 - A_{M-2}^{l-1} p^l)}{(C-1)^2 (\frac{C-1}{2} - A_{M-2}^{l-1} p^l (M-1))^2} \\
&= \frac{1}{(\frac{1}{2} - A_{M-2}^{l-1} p^l)^2} \cdot (1 - A_{M-2}^{l-1} p^l) \\
&\leq 8 \cdot (1 - A_{M-2}^{l-1} p^l) \\
&\leq \frac{\delta}{T}
\end{aligned}$$

when we specify that  $p \geq \frac{(C-l-1)!}{(C-2)!} (1 - \frac{\delta(C-1)}{8T})$ .

In other words, with probability at least  $1 - \frac{\delta \cdot l}{T}$ , we have that  $d_m > \frac{C-1}{2}$ .

It is well known that if  $\delta(G_t) \geq \frac{C-1}{2}$ , then we have that the composition graph  $G_{t-l+1} \cdot \dots \cdot G_t$  is connected where  $\delta(G_t) = \min_m d_m$ .

As a result, consider the probability and we obtain that

$$\begin{aligned}
& P(\text{graph } G_{t-l+1} \cdot \dots \cdot G_t \text{ is connected}) \\
&\geq P(\min_j d_j \geq \frac{C-1}{2}) \\
&= P(\bigcap_j \{d_j \geq \frac{C-1}{2}\}) \\
&= 1 - P(\bigcup_j \{d_j < \frac{C-1}{2}\}) \\
&\geq 1 - \sum_j P(d_j < \frac{C-1}{2}) \\
&= 1 - CP(d_j < \frac{C-1}{2}) \\
&\geq 1 - C \frac{\epsilon}{T} = 1 - \frac{C\epsilon}{T}
\end{aligned}$$

where the second inequality holds by the Bonferroni's inequality and the third inequality uses (11).

Consequently, we obtain

$$\begin{aligned}
& P(\text{graph } G_t \text{ is connected}) \\
&= P(\cap_t \{G_t \text{ is connected}\}) \\
&\geq 1 - \sum_t P(G_t \text{ is not connected}) \\
&= 1 - \sum_t (1 - P(G_t \text{ is connected})) \\
&\geq 1 - \sum_t (1 - (1 - \frac{C\epsilon}{T})) = 1 - C\epsilon
\end{aligned}$$

This indicates that with probability at least  $1 - \epsilon$ , for any time step  $t$ , the corresponding composition graph of  $G_t^C$  is connected. By definition, we conclude that the sequence  $G_t^C$  is  $l$ -periodically connected, which completes the proof.  $\square$

### Information delay

LEMMA 18. For any  $m, i, t > L$ , if  $n_{m,i}(t) \geq 2(K^2 + KM + M)$  and subgraph  $G_t$  is  $l$ -periodically connected in the sense that the composition of  $l$  consecutive graphs is a connected graph, then we have

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

where the min is taken over all clusters, not just the neighbors.

PROOF OF LEMMA 9. We consider Lemma 10 in [64], and thus establish that

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

We refer the full proof to the proof of Lemma 10 in [64].  $\square$

Then with the information delay results and the same update rule (Rule 2) as in Theorem 3, we can again establish the unbiasedness of the global estimator  $\tilde{\mu}_{m,i}(t)$ , the variance of  $\tilde{\mu}_{m,i}(t)$ , the concentration inequality with respect to  $\tilde{\mu}_{m,i}(t)$ , the upper bound on the number of pulls of sub-optimal arms, which are presented as follows (the proof of them is referred to the proof of Theorem 3).

### Unbiasedness of the estimator

PROPOSITION 21. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any arm  $i$  and any agent  $m$ , at every time step  $t$ , we have

$$E[\tilde{\mu}_i^m(t) | A_{\epsilon, \delta}] = \mu_i.$$

### Variance term



PROPOSITION 22. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . In setting  $s_1, s_2, s_3$  where rewards follow sub-gaussian distributions, for any  $m, i, \lambda$  and  $t > L$  where  $L$  is the length of the burn-in period, the global estimator  $\tilde{\mu}_i^m(t)$  is sub-Gaussian distributed. Moreover, the conditional moment generating function satisfies that with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,

$$\begin{aligned} & E[\exp\{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t)}\right\} \end{aligned}$$

where  $\sigma^2 = \max_{j,i} (\tilde{\sigma}_i^j)^2$  and  $C = \max\left\{\frac{4(M+2)(1-\frac{1-c_0}{2(M+2)})^2}{3M(1-c_0)}, (M+2)(1+4Md_{m,t}^2)/M\right\}$ .

### Concentration inequality

PROPOSITION 23. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any  $m, i$  and  $t > L$  where  $L$  is the length of the burn-in period,  $\tilde{\mu}_{m,i}(t)$  satisfies that if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$ , then with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,

$$\begin{aligned} P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}, \\ P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta}) t^2}. \end{aligned}$$

### Number of pulls of sub-optimal arms

Upper bounds on  $E[n_{m,k}(T) | A_{\epsilon, \delta}]$

PROPOSITION 24. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . An arm  $k$  is said to be sub-optimal if  $k \neq i^*$  where  $i^*$  is the unique optimal arm in terms of the global reward, i.e.  $i^* = \arg \max \frac{1}{M} \sum_{j=1}^M \mu_i^j$ . Then when the game ends, for every agent  $m$ ,  $0 < \epsilon < 1$  and  $T > L$ , the expected numbers of pulling sub-optimal arm  $k$  after the burn-in period satisfies with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$

$$\begin{aligned} & E[n_{m,k}(T) | A_{\epsilon, \delta}] \\ & \leq \max\left\{\frac{C}{M} \cdot \left[\frac{4C_1 \log T}{\Delta_i^2}, 2(K^2 + MK + M)\right] + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K\right. \\ & \left. \leq O(\log T).\right. \end{aligned}$$

As a concluding step, we again use the aforementioned regret decomposition that leads to the following

$$\begin{aligned} E[R_T | A_{\epsilon, \delta}] & \leq L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} E[n_{m,i}(T)] (\mu_{i^*}^m - \mu_i^m) \\ & \leq L + \sum_{i \neq i^*} \Delta_i \left(\max\left\{\frac{C}{M} \cdot \left[\frac{4C_1 \log T}{\Delta_i^2}, 2(K^2 + MK)\right] + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K\right\}\right). \end{aligned}$$

This completes the proof of Theorem 7. □

## E.6 Proof of Theorem 8

PROOF. In a like manner, we approach the regret analysis by decomposing the proof into the following phases, in the order indicated below.

We start with characterizing the sub-graph connectivity (referring to the sub-graphs including the one with respect to the agents within one cluster and the other with respect to the clusters).

### Graph connectivity

PROPOSITION 25. *Assume the edge probabilities meet the condition*

$$1 \geq \min_{m \neq n} p(m, n) \geq \frac{C^2}{M^2} \max \left\{ \frac{(C-l-1)!}{(C-2)!} \left(1 - \frac{\delta(C-1)}{8CT}\right), \frac{(C-l-1)!}{(C-2)!} \left(\frac{3}{4}\right)^{\frac{1}{l}} \right\},$$

and

$$1 \geq \min_{m, m} p(m, m) \geq \max \left\{ \frac{(c_M-l-1)!}{(c_M-2)!} \left(1 - \frac{\delta(c_M-1)}{8c_M T}\right), \frac{(c_M-l-1)!}{(c_M-2)!} \left(\frac{3}{4}\right)^{\frac{1}{l}} \right\},$$

where  $0 < \epsilon < 1$ . Then, with probability  $1 - \epsilon$ , for any  $t > 0$ , the sub-graph  $G_t^C$  induced by the clusters following the E-R model is  $l$ -periodically connected and the sub-graph  $G_t^{c_M}$  induced by the agents within one cluster (since we consider a balanced cluster) is also  $l$ -periodically connected.

PROOF OF PROPOSITION 25. The first part of the statement follows from Proposition 16 in the proof of Theorem 3.

For the second part, we repeat the proof of Proposition by treating the sub-graph as the sub-graph induced by the agents in one cluster, which has  $c_M$  vertex and the corresponding edge set determined by  $\{p(m, m)\}_m$ , instead of the sub-graph induced by the clusters. As a result, we omit the proof steps here and refer to Proposition 16 for a detailed version of the proof. □

Then we consider the information delay due to the randomness in both the cluster-wise sub-graph  $G_t^C$  and the within-cluster sub-graph  $G_t^{c_M}$ .

### Information delay

The following lemma characterize the first case.

LEMMA 19. *For any  $m, i, t > L$ , if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$  and subgraph  $G_t$  is  $l$ -periodically connected in the sense that the composition of  $l$  consecutive graphs is a connected graph, then we have*

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

where the min is taken over all clusters, not just the neighbors.

PROOF OF LEMMA 10. We consider Lemma 10 in [64], and thus establish that

$$\hat{N}_{m,i}(t) \leq 2 \min_j \hat{N}_{j,i}(t).$$

□

Additionally, we have the following proposition

LEMMA 20. For any  $m, i, t > L$ , if  $n_{m,i}(t) \geq 2(K^2 + KM + M)$  and subgraph  $G_t^{c_M}$  is  $l$ -periodically connected, then we have that for any  $t$  and  $m \in c_m$ ,

$$\begin{aligned} N_{m,i}(t) - c_M \cdot l &\leq \min_{m \in c_m} N_{m,i}(t-l) \leq N_{m,i}(t) \\ N_{m,i}(t-l) &\geq \hat{N}_{m,i}(t) - K(K+2M) - c_M \cdot l \end{aligned}$$

where the min is taken over all agents in one cluster.

PROOF OF LEMMA 11. We observe that all agents in one cluster will collect each other's information, after at most  $l$  steps, and only after that, they update the cluster information  $\hat{\mu}_m^i(t)$  that aggregates all agents' information and thus is the same for all agents in one cluster. And in between, the agents use the most recent cluster estimator, which is the same for agents in one cluster.

This implies that  $n_{m,i}(t) = n_{j,i}(t)$  for  $m, j \in c_m$ , as well as  $\min_{m \in c_m} N_{m,i}(t-l) = N_{m,i}(t-l) = N_{j,i}(t-l) \geq N_{j,i}(t) - c_M \cdot l$

Meanwhile, based on Lemma 1 in [64], we have that

$$N_{m,i}(t) \geq \hat{N}_{m,i}(t) - K(K+2M)$$

and thus

$$N_{m,i}(t-l) \geq \hat{N}_{m,i}(t) - K(K+2M) - c_M \cdot l.$$

This completes the proof.

□

## Unbiasedness

PROPOSITION 26. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any arm  $i$  and any agent  $m$ , at every time step  $t$ , we have

$$E[\tilde{\mu}_i^m(t) | A_{\epsilon, \delta}] = \mu_i.$$

PROOF OF PROPOSITION 26. It is worth noting that the information delay of length of  $l$  does not change the expected value of the global estimator, since the delayed estimator is also unbiased, the same as before. Hence, the proof of Proposition 16 for Theorem 3 holds herein, and thus we refer to the proof there.

□

## Variance term

**PROPOSITION 27.** *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . In setting  $s_1, s_2, s_3$  where rewards follow sub-gaussian distributions, for any  $m, i, \lambda$  and  $t > L$  where  $L$  is the length of the burn-in period, the global estimator  $\tilde{\mu}_i^m(t)$  is sub-Gaussian distributed. Moreover, the conditional moment generating function satisfies that with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,*

$$\begin{aligned} & E[\exp\{\lambda(\tilde{\mu}_i^m(t) - \mu_i)\} 1_{A_{\epsilon, \delta}} | \sigma(\{n_{m,i}(t)\}_{t,i,m})] \\ & \leq \exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t-l)}\right\} \end{aligned}$$

where  $\sigma^2 = \max_{j,i} (\tilde{\sigma}_i^j)^2$  and  $C = \max\left\{\frac{4(M+2)(1-\frac{1-c_0}{2(M+2)})^2}{3M(1-c_0)}, (M+2)(1+4Md_{m,t}^2)/M\right\}$ .

**PROOF OF PROPOSITION 27.** Based on Lemma 11, we observe that there is  $l$  delay in the within-cluster information, which implies that the quantity  $N_{m,i}(t)$  is at least  $N_{m,i}(t-l)$ . Also, we note that the term  $\exp\left\{\frac{\lambda^2}{2} \frac{C\sigma^2}{\min_j N_{j,i}(t)}\right\}$  is monotone decreasing in  $N_{j,i}(t)$ , which means that we can use this term  $N_{j,i}(t-l)$  in characterizing the moment generating function of  $\tilde{\mu}_i^m(t)$ , as well as the variance.

With  $N_{m,i}(t-l)$ , we repeat the proof of Proposition 10, following the proof in [58], which concludes the result. □

## Concentration inequality

**PROPOSITION 28.** *Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . For any  $m, i$  and  $t > L$  where  $L$  is the length of the burn-in period,  $\tilde{\mu}_{m,i}(t)$  satisfies that if  $N_{m,i}(t) \geq 2(K^2 + KM + M)$ , then with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$ ,*

$$\begin{aligned} P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t-l)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}, \\ P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t-l)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}. \end{aligned}$$

**PROOF OF PROPOSITION 28.** It is worth noting that in Algorithm 2, we specify the update frequency as  $\tau = l$ , which means between  $t-l$  and  $t$ , the agents in the same cluster do not update the information in order to make sure that they stay on the same page. Also, using Lemma we obtain that with high probability  $1 - \delta$ , after at most  $l$  steps, any two agents in one cluster communicate, and any two cluster communicate. At the end of  $t$ , they already collect all the information of agents within the cluster and they update the information. Based on the concentration inequality we obtained for the cluster-wise information as in the case where the within-cluster graph is a complete graph, we obtain

$$\begin{aligned} P(\tilde{\mu}_{m,i}(t) - \mu_i \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t-l)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}, \\ P(\mu_i - \tilde{\mu}_{m,i}(t) \geq \sqrt{\frac{C_1 \log t}{N_{m,i}(t-l)}} | A_{\epsilon, \delta}) & \leq \frac{1}{P(A_{\epsilon, \delta})} \frac{1}{t^2}. \end{aligned}$$

□

### Number of pulls of sub-optimal arms

PROPOSITION 29. Assume the parameter  $\delta$  satisfies that  $0 < \delta < c = f(\epsilon, M, T)$ . An arm  $k$  is said to be sub-optimal if  $k \neq i^*$  where  $i^*$  is the unique optimal arm in terms of the global reward, i.e.  $i^* = \arg \max \frac{1}{M} \sum_{j=1}^M \mu_i^j$ . Then when the game ends, for every agent  $m$ ,  $0 < \epsilon < 1$  and  $T > L$ , the expected numbers of pulling sub-optimal arm  $k$  after the burn-in period satisfies with  $P(A_{\epsilon, \delta}) = 1 - 7\epsilon$

$$\begin{aligned} & E[n_{m,k}(T)|A_{\epsilon, \delta}] \\ & \leq \max \left\{ \frac{C}{M} \cdot \left[ \frac{4C_1 \log T}{\Delta_i^2} \right], 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3P(A_{\epsilon, \delta})} + K^2 + (2M - 1)K + M \cdot l \\ & \leq O(\log T). \end{aligned}$$

PROOF OF PROPOSITION 29. It is worth noting that by time  $t$ , each agent  $m \in c_m$  has received the information  $\mathcal{F}_{t-l}$  of all agents in the same cluster. In other words,  $N_{m,i}(t) = \sum_{m \in c_m} n_{m,i}(t-l)$ , which implies that  $N_{m,i}(t) = N_{j,i}(t)$  for any agent  $m, j \in c_m$  for any  $t > L$ .

Likewise, we obtain that  $\tilde{\mu}_{m,i}(t) = \tilde{\mu}_{j,i}(t)$  for any  $m, j \in c_m$ , since we update  $\tilde{\mu}_{m,i}(t) = \sum_{j \in c_m} \frac{\hat{\mu}_{j,i}(t-l)}{|C_m|} = \hat{\mu}_{j,i}(t)$ . And the agents across the cluster exchange  $\hat{\mu}_{j,i}(t)$  if they do not belong to the same cluster, and update the estimation towards  $\tilde{\mu}_{j,i}(t)$ .

In light of the UCB decision rule,  $n_{m,i}(t)$  and  $n_{j,i}(t)$  only depend on  $N_{m,i}(t), \tilde{\mu}_{m,i}(t)$ . Therefore, we derive that  $n_{m,i}(t) = n_{j,i}(t)$  for any  $t$  for any  $m, j \in c_m$ , on event  $A_{\epsilon, \delta}$ .

Also, it is worth mentioning that  $n_{m,i}(t) \leq \frac{N_{m,i}(t)}{|C_m|}$ , concluded from the above statement and the fact that the cluster structure is balanced.

By considering 4 different cases regarding the possible values of  $N_{m,i}(t)$  as in [58] and noticing that  $N_{m,i}(t) \leq N_{m,i}(t-l) + c_M \cdot l$  and the fact that  $n_{m,i}(t) \leq \frac{N_{m,i}(t)}{|C_m|}$ , we obtain that

$$\begin{aligned} & E[n_{m,i}(T)|A_{\epsilon, \delta}] \\ & \leq \frac{E[N_{m,i}(T)|A_{\epsilon, \delta}]}{|C_m|} + l \\ & \leq \max \left\{ \left[ \frac{4C_1 \log T}{|C_m| \Delta_i^2} \right], 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3} + K^2 + (2M - 1)K + l. \end{aligned} \quad (12)$$

□

Next, we proceed to the regret decomposition and derive the upper bound on the regret.

### Regret decomposition

For the proposed regret, we have that for any constant  $L$ ,

$$\begin{aligned}
R_T &= \frac{1}{M} (\max_i \sum_{t=1}^T \sum_{m=1}^M \mu_i^m - \sum_{t=1}^T \sum_{m=1}^M \mu_{a_t^m}^m) \\
&= \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \sum_{t=1}^T \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \\
&\leq \sum_{t=1}^L \left| \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right| + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&\leq L + \sum_{t=L+1}^T \left( \frac{1}{M} \sum_{m=1}^M \mu_{i^*}^m - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&= L + \sum_{t=L+1}^T \left( \mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \mu_{a_t^m}^m \right) \\
&= L + ((T-L) \cdot \mu_{i^*} - \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^K n_{m,i}(T) \mu_i^m)
\end{aligned}$$

where the first inequality is by taking the absolute value and the second inequality results from the assumption that  $0 < \mu_i^j < 1$  for any arm  $i$  and agent  $j$ .

Note that  $\sum_{i=1}^K \sum_{m=1}^M n_{m,i}(T) = M(T-L)$  where by definition  $n_{m,i}(T)$  is the number of pulls of arm  $i$  at agent  $m$  from time step  $L+1$  to time step  $T$ , which yields that

$$\begin{aligned}
R_T &\leq L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_{i^*}^m - \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) \mu_i^m \\
&= L + \sum_{i=1}^K \frac{1}{M} \sum_{m=1}^M n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&\leq L + \frac{1}{M} \sum_{i=1}^K \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \\
&= L + \frac{1}{M} \sum_{i \neq i^*} \sum_{m: \mu_{i^*}^m - \mu_i^m > 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m).
\end{aligned}$$

where the second inequality uses the fact that  $\sum_{m: \mu_{i^*}^m - \mu_i^m \leq 0} n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) \leq 0$  holds for any arm  $i$  and the last equality is true since  $n_{m,i}(T) (\mu_{i^*}^m - \mu_i^m) = 0$  for  $i = i^*$  and any  $m$ .

Meanwhile, by the choices of  $\delta$  such that  $\delta < c = f(\epsilon, M, T)$ , we apply Proposition 6 which leads to for any agent  $m$  and arm  $i \neq i^*$ ,

As a result, the upper bound on  $R_T$  can be derived as by taking the conditional expectation over  $R_T$  on  $A_{\epsilon, \delta}$

$$\begin{aligned}
& E[R_T | A_{\epsilon, \delta}] \\
& \leq L + \sum_{i \neq i^*} \sum_{j=1}^C \sum_{m \in c_j} E[n_{m,i}(T) | A_{\epsilon, \delta}] (\mu_{i^*}^m - \mu_i^m) \\
& \leq \sum_{i \neq i^*} \sum_{j=1}^C \sum_{m \in c_j} \max \left\{ \left\lceil \frac{4C_1 \log T}{|C_m| \Delta_i^2} \right\rceil, 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3} + K^2 + (2M - 1)K + l \\
& \leq \sum_{i \neq i^*} \sum_{j=1}^C \max \left\{ \left\lceil \frac{4C_1 \log T}{\Delta_i^2} \right\rceil, 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3} + K^2 + (2M - 1)K + l \\
& \leq \sum_{i \neq i^*} C \max \left\{ \left\lceil \frac{4C_1 \log T}{\Delta_i^2} \right\rceil, 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3} + K^2 + (2M - 1)K + l \tag{13}
\end{aligned}$$

where the second inequality holds by plugging in (12).

Hence, the regret can be upper bounded by

$$\begin{aligned}
& E[R_T | A_{\epsilon, \delta}] \\
& \leq L + \sum_{i \neq i^*} C(\Delta_i + 1) \left( \max \left\{ \left\lceil \frac{C}{M} \cdot \frac{4C_1 \log T}{\Delta_i^2} \right\rceil, 2(K^2 + MK + M) \right\} + \frac{2\pi^2}{3} + K^2 + (2M - 1)K + l \right) \\
& = O(\max\{L, \log T\})
\end{aligned}$$

This completes the proof of Theorem 8.

□