

Exploration for Free: How Does Reward Heterogeneity Improve Regret in Cooperative Multi-agent Bandits?

Xuchuang Wang¹ Lin Yang² Yu-Zhen Janice Chen³ Xutong Liu¹ Mohammad Hajiesmaili³ Don Towsley³
John C.S. Lui¹

¹Department of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong

²School of Intelligence Science and Technology, Nanjing University, Jiangsu, China

³College of Information and Computer Sciences, University of Massachusetts Amherst, Massachusetts, USA

Abstract

This paper studies a cooperative multi-agent bandit scenario in which the rewards observed by agents are heterogeneous—one agent’s meat can be another agent’s poison. Specifically, the total reward observed by each agent is the sum of two values: an arm-specific reward, capturing the intrinsic value of the arm, and a privately-known agent-specific reward, which captures the personal preference/limitations of the agent. This heterogeneity in total reward leads to different local optimal arms for agents but creates an opportunity for *free exploration* in a cooperative setting—an agent can freely explore its local optimal arm with no regret and share this free observation with some other agents who would suffer regrets if they pull this arm since the arm is not optimal for them. We first characterize a regret lower bound that captures free exploration, i.e., arms that can be freely explored have no contribution to the regret lower bound. Then, we present a cooperative bandit algorithm that takes advantage of free exploration and achieves a near-optimal regret upper bound which tightly matches the regret lower bound up to a constant factor. Lastly, we run numerical simulations to compare our algorithm with various baselines without free exploration.

1 INTRODUCTION

Multi-armed bandit (MAB) [Lai et al., 1985, Bubeck et al., 2012] is a classic sequential decision making problem. In the stochastic MAB, an agent faces a set $\mathcal{K} := \{1, 2, \dots, K\}$ ($K \in \mathbb{N}^+$) of arms, where each arm k is associated with a reward random variable with unknown mean $\mu(k)$. The agent sequentially pulls arms from \mathcal{K} in $T \in \mathbb{N}^+$ decision rounds and observes the pulled arm rewards. The goal of

the agent is to maximize its total reward over all decision rounds, which is equivalent to minimizing the total *regret*, i.e., the cumulative reward difference between the aggregate reward of the optimal arm k_* with the highest mean and the agent’s sequential choices. To achieve this goal, the agent needs to balance between exploration and exploitation, i.e., either optimistically choose the arm with high uncertainty in reward (exploration), or myopically pull the one with high empirical mean reward (exploitation).

Multi-agent MAB (MA2B) is an extension of the basic MAB, where a group of $M \in \mathbb{N}^+$ agents (denoted as $\mathcal{M} := \{1, 2, \dots, M\}$) pulls arms from the same arm set \mathcal{K} . This model has been studied in various settings, e.g., federated bandits [Shi and Shen, 2021, Shi et al., 2021a, Zhu et al., 2021, Huang et al., 2021], cooperative pure exploration [Hillel et al., 2013, Tao et al., 2019, Karpov et al., 2020], multi-agent MAB with collision [Boursier and Perchet, 2019, Mehrabian et al., 2020, Shi et al., 2021b], and cooperative multi-agent MAB [Landgren et al., 2016, Martínez-Rubio et al., 2019, Wang et al., 2020a,b].

The majority of prior works on MA2B, with a few exceptions (see Appendix A), study a homogeneous reward setting, where the reward distribution of an arm is the same for all agents. The homogeneous reward setting, however, fails to capture agent-specific preferences/limitations. In many real-world applications, the agents represent different clusters of users with specific preferences, or users in different geographical locations with different costs/limits to access the arm set. In such settings, the reward of each arm might be different for different agents. We refer to Section 2.3 for a detailed explanation of various application scenarios.

This paper introduces a multi-agent multi-armed bandits problem with heterogeneous reward (MA2B-HR). In MA2B-HR, the reward observed by an agent consists of two components representing arm- and agent-specific terms. Specifically, when agent $i \in \mathcal{M}$ pulls arm $k \in \mathcal{K}$, the observed reward is $X_t^{(i)}(k) = X_{t,\text{arm}}(k) + X_{t,\text{agent}}^{(i)}(k)$, where $X_{t,\text{arm}}(k)$ is the arm-specific reward with bounded mean

48 $\mu(k) \in (0, b)$ (where b is a positive constant) and $X_{t,\text{agent}}^{(i)}(k)$
 49 is the agent-specific reward with mean $\nu^{(i)}(k)$. We denote
 50 $\omega^{(i)}(k) := \mu(k) + \nu^{(i)}(k)$ as the reward mean of this pull.
 51 In MA2B-HR, we assume both $X_{t,\text{arm}}(k)$ and $X_{t,\text{agent}}^{(i)}(k)$ are
 52 stochastic and independent. The arm-specific reward mean
 53 $\mu(k)$ is not known to agents, and each agent i only privately
 54 knows its own agent-specific mean values $\nu^{(i)}(k), \forall k \in \mathcal{K}$.
 55 Further, in the MA2B-HR setting, the agents can broadcast
 56 the observed values of the arm-specific term in the total
 57 reward (by subtracting the agent-specific reward mean from
 58 the observed reward, i.e., $X_t^{(i)}(k) - \nu^{(i)}(k)$) at no cost. We
 59 note that one may consider other settings for MA2B-HR,
 60 e.g., known vs. unknown and homogeneous vs. heteroge-
 61 neous assumptions for the agent-specific reward. We refer to
 62 Appendix B.1 for a detailed discussion and the connection
 63 of each setting to the prior literature.

64 In MA2B-HR, the reward heterogeneity of agents creates a
 65 counterintuitive opportunity for *free exploration* of a subset
 66 of arms. With heterogeneous rewards among agents, there
 67 might be no global optimal arm(s). In other words, agents
 68 may have different *local* optimal arms, i.e., the arms with
 69 the largest reward mean are different among agents, so the
 70 characterization of the regret of agents becomes more com-
 71 plicated. However, the existence of multiple local optimal
 72 arms poses a surprising opportunity to develop a cooperative
 73 learning algorithm to explore local optimal arms for free
 74 (without cost), share the free observations with others, and
 75 significantly improve the total regret among all agents.

76 While the idea of free exploration is intuitive, designing a
 77 cooperative bandit algorithm that effectively implements
 78 this idea is nontrivial. The main challenge is that the local
 79 optimal arms are unknown in advance to the bandit agents.
 80 Hence, an algorithm should be designed to economically
 81 identify the local optimal arms and assign them to agents
 82 that can freely explore them and prevent other agents from
 83 pulling these arms (with cost).

84 We note that MA2B-HR could be considered as an ex-
 85 tended version of two recent models in the bandits' liter-
 86 ature: action-constrained multi-agent multi-armed bandits
 87 (AC-MA2B) Yang et al. [2022] and grouped K -armed bandits
 88 Baek and Farias [2021]. The idea of free exploration
 89 is applicable to both Yang et al. [2022], Baek and Farias
 90 [2021], however, they did not explicitly utilize free explora-
 91 tion in algorithm design, so they fail to achieve optimal
 92 performance that takes into account the free exploration. A
 93 detailed discussion on both models and their connection to
 94 MA2B-HR, and the significance of our results with respect
 95 to both models are given in Section 1.2.

96 It is worth noting that the high-level idea of free exploration
 97 has been leveraged in some other bandit settings in the liter-
 98 ature [Chen et al., 2018, Shi et al., 2021c]. However, these
 99 works considered the problem of incentivizing exploration;
 100 specifically, they considered a principal, aiming to learn the

global bandit model, offering bonuses to agents to do ex- 101
 102 plorations on the principal's behalf. In these settings, Chen
 103 et al. [2018], Shi et al. [2021c] studied free exploration in
 104 the sense that the principal pays no cost rather than free explora-
 105 tion in cooperation among agents. Hence, these works are
 106 in clear contrast to the idea of free exploration in MA2B-HR
 107 introduced in this paper. A comprehensive comparison to
 108 related works are presented in Appendix A.

1.1 CONTRIBUTIONS

109 In this paper, we first present the MA2B-HR model and
 110 highlight its real-world applications. Then, we propose
 111 FreeExp, a cooperative algorithm designed to enable free
 112 exploration in the learning process. Finally, we characterize
 113 a regret lower bound that explicitly captures the impact of
 114 free exploration on MA2B-HR, and show that the regret of
 115 FreeExp matches the regret lower bound up to a constant
 116 factor. The contributions of this work are:

Modeling and practical relevance of MA2B-HR: We 117
 118 present the MA2B-HR model in Section 2 and justify its prac-
 119 tical relevance by highlighting several application scenarios
 120 in online advertising, wireless networks, and cloud and edge
 121 resource allocation. We also introduce a new definition for
 122 the suboptimality gap in MA2B-HR as a key parameter to
 123 explicitly characterize the impact of free exploration in the
 124 regret analysis.

Algorithm design: In Section 3, we present FreeExp, a 125
 126 cooperative learning algorithm that tackles MA2B-HR and
 127 implements the idea of free exploration. The high level idea
 128 of FreeExp is that agents judiciously reduce the selection
 129 of arms that are likely to be local optimal for other agents.
 130 Instead, by cooperation, those agents can still get the obser-
 131 vations on those arms from others without regret cost. In
 132 doing so, free exploration of some arms becomes possible
 133 and the cooperative bandit algorithm achieves significant
 134 improvement in regret. A key technique in FreeExp is to
 135 perform periodic pulls of the empirical local optimal arms
 136 (i.e., the arm with the highest empirical mean) while balanc-
 137 ing between exploration and exploitation, which guarantees
 138 that the empirical optimal arm is indeed the ground truth
 139 local optimal arm in most time slots.

Regret analysis: In contrast to the common regret analysis 140
 141 in multi-agent bandits where only the pulled arm matters
 142 regardless of the agent who pull the arm, in MA2B-HR, we
 143 have to address a unique technical challenge since the regret
 144 cost of pulling an arm depends not only on which arm is
 145 pulled, but also on which agent pulls it. In Section 4, we
 146 tackle this challenge and derive a regret lower bound for
 147 MA2B-HR that echos the importance of recognizing free
 148 explorations: arms that can be freely explored only cause
 149 constant regret, instead of the usual logarithmic regret in
 150 MA2B. We derive the regret upper bound of the FreeExp

Table 1: A simple example with three agents and three arms ($b > \mu(1) > \mu(2) > \mu(3) > 0$). The entries of the table show the total reward of each arm for each agent, e.g., $\omega^{(1)}(1) = \mu(1)$ or $\omega^{(3)}(2) = \mu(2) - b < 0$. Arms 1, 2, and 3 are the local optimal arms of agents 1, 2, and 3, respectively. On the right-hand side, denoting $\Delta(i, j) = \mu(i) - \mu(j)$, the regret of our work is compared with a classic non-cooperative algorithm [Auer, 2002] and the works of Yang et al. [2022] and Baek and Farias [2021] as two special cases of MA2B-HR.

	Arm 1	Arm 2	Arm 3		
				UCB [Auer, 2002]	$O\left(\left(\frac{1}{\Delta(1,2)} + \frac{1}{\Delta(1,3)} + \frac{1}{\Delta(2,3)}\right) \log T\right)$
Agent 1	$\mu(1)$	$\mu(2)$	$\mu(3)$	CO-UCB [Yang et al., 2022]	$O\left(\left(\frac{1}{\Delta(1,2)} + \frac{1}{\Delta(2,3)}\right) \log T\right)$
Agent 2	< 0	$\mu(2)$	$\mu(3)$	KL-UCB [Baek and Farias, 2021]	$O(\log \log T)$
Agent 3	< 0	< 0	$\mu(3)$	FreeExp (our work)	$O(1)$

algorithm which matches the regret lower bound up to a constant factor. Deriving this result requires new analysis techniques (see Theorem 4.3’s proof sketch in Section 4 for detail). The tightness of both regret upper and lower bounds reflects the intrinsic property of MA2B-HR where free exploration plays a key role, and that FreeExp is near-optimal. A surprising observation is that in the special cases where every arm is local optimal for at least one agent (reasonable when $M \geq K$), FreeExp achieves an $O(1)$ regret.

Numerical results: In Section 5, we report numerical experiments of comparing our algorithm to several baselines.

1.2 TECHNICAL COMPARISON TO THE PRIOR WORK

In this section, we highlight our contribution in leveraging free exploration by applying our algorithm to the action-constrained MA2B problem (AC-MA2B) which was recently studied by Yang et al. [2022]. In AC-MA2B, each agent $i \in \mathcal{M}$ only pulls from a subset of arms $\mathcal{K}^{(i)} \subset \mathcal{K}$ and its goal is to find the local optimal arm in $\mathcal{K}^{(i)}$. AC-MA2B can be regarded as a special case of MA2B-HR when agent i ’s specific reward $\nu^{(i)}(k)$ for arm k is 0 if $k \in \mathcal{K}^{(i)}$, and $-b$ if $k \notin \mathcal{K}^{(i)}$, where $b > 0$ and $\mu(k) \in (0, b)$ for all arm k (see Remark 2.1 for a formal definition). Since agent i knows its agent-specific reward means, she would never pull arms with $\nu^{(i)}(k) = -b$ and thus is equivalent to only having access to arms in the constrained arm set $\mathcal{K}^{(i)}$. We provide a simple example in Table 1 to illustrate the benefit of free exploration which substantially improves regret as compared to the classic non-cooperative algorithms and the cooperative approach in Yang et al. [2022] as a special case.

Next, we present the theoretical improvement. Recall that the non-cooperative optimal total regret of classic MAB [Lai et al., 1985] for all agents in \mathcal{M} is

$$O\left(\sum_{i \in \mathcal{M}} \sum_{k \in \mathcal{K}^{(i)} \setminus \{k_*^{(i)}\}} \frac{\Delta^{(i)}(k) \log T}{\text{kl}(\mu(k), \mu(k) + \Delta^{(i)}(k))}\right),$$

where the suboptimality gap $\Delta^{(i)}(k) := \mu(k_*^{(i)}) - \mu(k)$

is the difference of reward means between agent i ’s optimal arm $k_*^{(i)}$ and arm k , and $\text{kl}(a, b)$ is the KL-divergence between two Gaussian distributions with means a and b and the same variance (defined later). To improve total regret through cooperation, Yang et al. [2022] proposed cooperative extensions to classic learning algorithms, e.g., UCB [Auer, 2002], which improved the total regret to

$$O\left(\sum_{k \in \cup_i (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})} \frac{\bar{\Delta}(k) \log T}{\text{kl}(\mu(k), \mu(k) + \bar{\Delta}(k))}\right), \quad (1)$$

where $\bar{\Delta}(k)$ denotes the smallest reward mean gap of arm k compared to the local optimal arms (excluding arm k) among agents having access to arm k .

The regret of applying FreeExp to AC-MA2B is

$$O\left(\sum_{k \in \cup_i \mathcal{K}^{(i)} \setminus \cup_i \{k_*^{(i)}\}} \frac{\bar{\Delta}(k) \log T}{\text{kl}(\mu(k), \mu(k) + \bar{\Delta}(k))}\right). \quad (2)$$

The improvement of our result lies in the summation range. Specifically, the summation range $\cup_i \mathcal{K}^{(i)} \setminus \cup_i \{k_*^{(i)}\}$ in (2) is a subset of (1)’s $\cup_i (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})$. The summation range in (2) excludes the regret impact of arms in $\cup_i \{k_*^{(i)}\}$, i.e., arms that are optimal to at least one agent; these arms are freely explored. In contrast, the regret of Yang et al. [2022] in (1) is over $\cup_i (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})$, which counts some arms that are optimal for some agents (and can be freely explored). We note that this improvement can be substantial. Especially, when all arms in \mathcal{K} are locally optimal for some agents, the regret upper in (2) is $O(1)$, e.g., the simple example in Table 1. This implies that capturing the benefit of free exploration requires the development of a completely new cooperative algorithm as explained in Section 3.

The grouped K -armed bandits model proposed by Baek and Farias [2021] is almost equivalent to AC-MA2B Yang et al. [2022] except for minor differences in how their actions are constrained—the grouped bandits’ action constraint depends on the arrived group while AC-MA2B’s is associates to the agents. Therefore, the grouped bandits model can also be regarded as a special case of our MA2B-HR model. Baek

216 and Farias [2021] proved that the KL-UCB algorithm Cappé
 217 et al. [2013] can address their grouped bandits model with
 218 the regret performance as follows,

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \leq \sum_{k \in \cup_i \mathcal{K}^{(i)} \setminus \cup_i \{k_*^{(i)}\}} \frac{\bar{\Delta}(k)}{\text{kl}(\mu(k), \mu(k) + \bar{\Delta}(k))}.$$

219 We emphasize that the above bound of Baek and Farias
 220 [2021] was in an asymptotic form (i.e., for $T \rightarrow \infty$), while
 221 FreeExp’s regret bound is in a non-asymptotic form (i.e.,
 222 for any time T , see Eq.(10) of Theorem 4.3), which differs
 223 a lot in handling the regret of free arms (see Remark 4.7 for
 224 detail). Here, we pick the toy example in Table 1 to illus-
 225 trate the difference; this can be generalized to any case that
 226 all arms are free arms. In this example, FreeExp attains
 227 the $O(1)$ regret, while KL-UCB’s regret was $o(\log T)$ (or,
 228 $O(\log \log T)$ specifically) [Baek and Farias, 2021]. In Sec-
 229 tion 5, we conduct numerical comparisons to corroborate
 230 the advantage of FreeExp over KL-UCB. Also, we empha-
 231 size that our regret upper bound is proved for the MA2B-HR
 232 model which is more general than Baek and Farias [2021]’s
 233 grouped bandits model.

2 MODEL AND NOTATIONS

234 We first present the multi-agent multi-armed bandits with
 235 heterogeneous rewards problem (MA2B-HR) in Section 2.1
 236 and its performance metric in Section 2.2. In Section 2.4, we
 237 introduce notations related to free exploration to facilitate
 238 our algorithm design and analysis.

2.1 MA2B-HR: THE MULTI-AGENT MULTI-ARMED BANDITS WITH HETEROGENEOUS REWARDS

239 In MA2B-HR, there are $K \in \mathbb{N}^+$ arms and $M \in \mathbb{N}^+$ agents.
 240 Each arm $k \in \mathcal{K} := \{1, 2, \dots, K\}$ is associated with
 241 a Gaussian reward random variable with unknown mean
 242 $\mu(k) \in (0, b)$ and variance σ_1^2 , where b is positive and
 243 known.¹ This is the *arm-specific reward* representing the
 244 intrinsic value of the arm and it is independent of the pref-
 245 erence of the agents. In addition, each agent has its own
 246 private *agent-specific reward* for each arm to capture its
 247 private preference for different arms. The agent-specific re-
 248 ward of agent i for arm k is modelled by a Gaussian random
 249 variable with mean $\nu^{(i)}(k)$ and variance σ_2^2 . The variances
 250 σ_1^2 and σ_2^2 are common for all arms and agents. The agent-
 251 and arm-specific rewards are independent, and both are also
 252 independent across arms \mathcal{K} and time $t = 1, 2, \dots$

253 By pulling an arm k at time t , agent i observes a Gaussian
 254 reward $X_t^{(i)}(k)$ with mean $\omega^{(i)}(k) := \mu(k) + \nu^{(i)}(k)$ and
 255 variance $\sigma_1^2 + \sigma_2^2$. In this paper, we assume that the value

¹If b is unknown, we can set it as an arbitrarily large constant.

of $\nu^{(i)}(k)$ is only known to agent i , but unknown to other
 256 agents, for all agent $i \in \mathcal{M}$. Similar to the basic setting of
 257 stochastic bandits, the arm-specific reward means $\mu(k)$ are
 258 unknown to all agents. We also assume, for each agent i ,
 259 that all mean rewards $\omega^{(i)}(k)$ ($\forall k \in \mathcal{K}$) are different; hence
 260 each agent has a unique optimal arm.
 261

Remark 2.1 (Agent’s local arm set). Observe that $\mu(k) \in$
 262 $(0, b)$. Consequently, if there exist two arms k_1, k_2 such that
 263 $\nu^{(i)}(k_1) \geq \nu^{(i)}(k_2) + b$ for agent $i \in \mathcal{M}$, then
 264

$$\begin{aligned} \omega^{(i)}(k_1) - \omega^{(i)}(k_2) &= (\mu(k_1) + \nu^{(i)}(k_1)) - (\mu(k_2) + \nu^{(i)}(k_2)) \\ &> \mu(k_1) - \mu(k_2) + b > 0, \end{aligned}$$

265 that is, for agent i , the reward mean of arm k_1 is higher than
 266 that of arm k_2 . Therefore, there is no need for agent i to pull
 267 arm k_2 . More generally, we define agent i ’s *local arm set* as
 268 follows, Therefore, agent i ’s local arm set is

$$\mathcal{K}^{(i)} := \left\{ k \in \mathcal{K} : \nu^{(i)}(k) + b > \max_{\ell \in \mathcal{K}} \nu^{(i)}(\ell) \right\},$$

269 and agent i only needs to explore arms in its local arm set.

Another relevant model for reward heterogeneity is contex-
 270 tual bandits [Li et al., 2010]. We discuss it in Appendix B.2.
 271 The MA2B-HR model finds applications in diverse domains,
 272 e.g., online advertising, online shortest path routing, online
 273 cloud and edge resources allocation, and personalized clinical
 274 trial, cf., the detail application scenarios in Appendix 2.3.
 275

2.2 PERFORMANCE METRICS

276 Since rewards are heterogeneous across agents, agents may
 277 have different optimal arms. The goal of each agent is to
 278 find its *local* optimal arm, the one with the largest total
 279 reward, which is the sum of arm- and agent-specific rewards.
 280 Let $k_*^{(i)}$ be the local optimal arm of agent i , i.e., $k_*^{(i)} :=$
 281 $\arg \max_{k \in \mathcal{K}^{(i)}} \omega^{(i)}(k)$. For an algorithm \mathcal{A} , let $J_t^{(i)}(\mathcal{A})$ be
 282 the arm pulled by agent i at time t . The expected regret
 283 of agent i under algorithm \mathcal{A} is the difference between
 284 the aggregate reward of pulling its local optimal arm and
 285 the aggregate reward of pulling arms in an online manner
 286 according to a bandit algorithm, i.e.,

$$\mathbb{E}[\mathbf{R}_T^{(i)}(\mathcal{A})] := T\omega^{(i)}(k_*^{(i)}) - \mathbb{E} \left[\sum_{t=1}^T \omega^{(i)}(J_t^{(i)}(\mathcal{A})) \right],$$

287 where the expectation is taken over the randomness of action
 288 sequence $\{J_1^{(i)}(\mathcal{A}), J_2^{(i)}(\mathcal{A}), \dots\}$.

289 In the MA2B-HR model, agents can cooperate and share in-
 290 formation to accelerate bandit learning. In particular, we as-
 291 sume that each agent can broadcast the arm-specific reward
 292 term (the observed rewards minus the agent-specific reward
 293 mean, $X_t^{(i)}(k) - \nu^{(i)}(k)$) at no cost to all other agents, and

other agents immediately receives the broadcast observations. Note that this basic system model can be extended to include the communication costs, or an underlying topology to govern communication between agents, or agent privacy, etc. We leave these extensions to future works and focus on presenting the key idea of free exploration in this paper. The learning environment is a cooperative one, hence, we consider *aggregate regret* as the performance metric, which is simply the aggregate regret over M agents, i.e.,

$$\mathbb{E}[\mathbf{R}_T(\mathcal{A})] := \sum_{i=1}^M \left(T\omega^{(i)}(k_*^{(i)}) - \mathbb{E} \left[\sum_{t=1}^T \omega^{(i)}(k_t^{(i)}) \right] \right). \quad (3)$$

2.3 APPLICATION SCENARIOS

The heterogeneous and known agent-specific reward means for MA2B-HR is a practically relevant setting and can find applications in diverse domains. The applications mentioned in Yang et al. [2022] and Baek and Farias [2021] can also be handled by MA2B-HR since their models are special cases of MA2B-HR. In the following, we present four motivating application scenarios that MA2B-HR could model. We note that we focus on motivating the arm- and agent-specific rewards. Detailed modeling of each application may require additional effort, which is beyond the scope of this paper.

Online Advertising in Social Networks: Online advertising is a classic example of the MAB problem [Tang et al., 2014, Mahadik et al., 2020]. Consider a scenario where there are multiple bandit agents that select ads to be placed on a social platform. Each agent is responsible for a cluster of users with similar interests. The cluster may be constructed based on different criteria, e.g., location, age, etc. Indeed, the popularity of products can differ across different locations or age groups. But the ads (arms) could be selected from a shared pool of available ads. In this scenario, the agent is aware of the personal preferences of users in its cluster, i.e., the agent-specific reward is known. However, the agents need to learn the potential value of ads as well; hence, arm-specific rewards are unknown. Since the learning agents all belong to the same social platform advertising engine, they can cooperate to share arm-specific observations and improve learning performance.

Online Shortest Path Routing in Wireless Networks: Another example is the problem of finding shortest paths in a multi-hop wireless network. Consider a scenario in which multiple learning agents try to learn the shortest paths for different communication sessions. In this scenario, bandit algorithms can be implemented to learn the shortest routing paths [He et al., 2013, Zou et al., 2014, Talebi et al., 2017]. The cost (or latency) of a certain path (arm) depends on the physical condition of the path itself, representing an arm-specific cost unknown to the learning agents. Further, the session of each agent might have its local physical conditions, e.g., distance and the hardware spec of the mobile

device, which is known only to the agent and impacts the overall cost of each path. In this scenario, the former is an arm-specific cost, which is homogeneous and unknown among all agents, while the latter varies across agents and whose mean is privately known to each agent only.

Online Cloud and Edge Resource Allocation: In prior literature, the MAB framework has been used for workload allocation into a pool of cloud/edge servers [Talebi and Proutiere, 2018, Johari et al., 2017, Lattimore et al., 2014, Dagan and Koby, 2018]. In this scenario, the cloud provider may categorize the compute jobs into multiple types, e.g., ML training workload, video processing, financial analytics, etc., and create a learning agent for finding the best server type for them. In this scenario, the arm-specific reward captures the hardware spec of the servers, and the agent-specific reward captures the job-specific hardware requirement of the workload, e.g., video processing is memory-intensive, while finance workload is compute-intensive. In edge scenarios where the workload could be run in multiple locations, the agent-specific reward could be represented as the cost of moving the workload to different locations as well, which is known and heterogeneous for different agents.

Personalized Medicine and Clinical Trial: A classic MAB application is clinical trial [Lai [1987], Villar et al. [2015], Aziz et al. [2021]]. Consider a scenario where patients have different covariates, e.g., age, gender, genomic features, and medical history, and, therefore, should be categorized to several heterogeneous groups, and the doctor should create personalized agents (drug application policies) for every group. In this scenario, the effectiveness of a treatment for a certain patient group depends not only on the treatment itself but also on the patient group’s covariates. For example, the effectiveness of a treatment that disturbs patients’ blood glucose concentrations may be discounted on diabetics. In this scenario, the arm-specific reward captures treatments’ or medicines’ basic effectiveness on a diseases, and the agent-specific reward (or cost) captures the discounted or additional effectiveness due to the patient group features. The latter is known to (or can be well evaluated by) an expert.

2.4 NOTATIONS RELATED TO FREE EXPLORATION

To ease the presentation of FreeExp and its analysis, we introduce some key notations relevant to free exploration. In MA2B-HR, arms that are local optimal for at least one agent can be freely explored. Then, in a cooperative environment, other agents who take these arms as their suboptimal choices can enjoy the freely explored observations of these arms.

Definition 2.2 (Set of free arms). We define the set of free

389 arms \mathcal{K}^{fr} as

$$\mathcal{K}^{\text{fr}} := \{k \in \mathcal{K} : \mathcal{M}_*(k) \neq \emptyset\}, \quad (4)$$

390 where $\mathcal{M}_*(k) := \{i \in \mathcal{M} : k \in \mathcal{K}^{(i)}, k = k_*^{(i)}\}$ is a subset
 391 of agents with arm k as their local optimal arm. Any arm
 392 $k \in \mathcal{K}^{\text{fr}}$ can be freely explored without incurring regret by
 393 any agent in $\mathcal{M}_*(k)$. In the rest of this paper, we refer to
 394 the arms in \mathcal{K}^{fr} as free arms.

395 Recall that in the classic MAB, the difficulty of distinguish-
 396 ing a suboptimal arm k from the optimal arm depends on
 397 $\Delta(k)$ —the reward mean gap between arm k and the optimal
 398 arm k^* . In MA2B-HR, the notion of optimality gap needs to
 399 be redefined since agents may have different local optimal
 400 arms. In the following, we formally define the suboptimality
 401 gap of each arm k as the smallest gap between arm k and
 402 any local optimal arms. A formal definition is given below.

403 **Definition 2.3** (Suboptimality gap). The suboptimality gap
 404 of arm k is defined as

$$\bar{\Delta}(k) := \min_{i \in \mathcal{M}} \Delta^{(i)}(k), \quad (5)$$

405 where $\Delta^{(i)}(k) := \omega^{(i)}(k_*^{(i)}) - \omega^{(i)}(k)$ is the gap between
 406 the mean rewards of arm k and $k_*^{(i)}$ —the local optimal arm
 407 of agent i .

408 All free arms have zero suboptimality gaps, i.e., $\bar{\Delta}(k) =$
 409 $0, \forall k \in \mathcal{K}^{\text{fr}}$. Denote $\bar{i}(k) \in \arg \min_{i \in \mathcal{M}(k)} \Delta^{(i)}(k)$ to
 410 be an agent with the smallest reward gap of arm k (one
 411 can break ties arbitrarily). Then, $\bar{\Delta}(k)$ can be rewritten as
 412 $\bar{\Delta}(k) = \omega^{\bar{i}(k)}(k_*^{\bar{i}(k)}) - \omega^{\bar{i}(k)}(k)$, where for simplicity,
 413 we denote $\omega^{\bar{i}(k)}(k)$ as $\bar{\omega}(k)$, i.e.,

$$\bar{\omega}(k) := \omega^{\bar{i}(k)}(k) = \mu(k) + \nu^{\bar{i}(k)}(k). \quad (6)$$

3 THE FREEEXP ALGORITHM

414 In this section, we present the FreeExp algorithm, which
 415 solves a multi-agent bandit problem in the MA2B-HR model.
 416 Each agent runs its own FreeExp algorithm and cooper-
 417 ates with each other. In Section 4, we demonstrate that with
 418 FreeExp, the reward heterogeneity not only does no harm,
 419 but in fact benefits the cooperative learning by the unique
 420 opportunity of free exploration.

421 **High-level idea of FreeExp:** We now explain how
 422 FreeExp implements the idea of free exploration to re-
 423 duce regret. The pivot of FreeExp is the local optimal
 424 (free) arm of each agent, which is unknown in advance. To
 425 address that for an agent i , FreeExp maintains an local
 426 optimal arm estimate $I_t^{(i)}$ of the agent i and an *exploration*
 427 *arm set* $\mathcal{D}_t^{(i)}$ containing arms that might be the ground truth

Algorithm 1 The FreeExp Algorithm (for Agent i)

- 1: **Initialize:** $d_t(k) = 0, \hat{\mu}_t(k) = 0, \hat{\omega}_t^{(i)}(k) := \hat{\mu}_t(k) + \nu^{(i)}(k)$.
 - 2: **for** each time slot t **do**
 - 3: $I_t^{(i)} \leftarrow \arg \max_{k \in \mathcal{K}^{(i)}} \hat{\omega}_t^{(i)}(k)$ {identify the empirical optimal arm}
 - 4: Send $I_t^{(i)}$ to other agents and collect their $I_t^{(j)}$
 - 5: $\mathcal{D}_t^{(i)} \leftarrow \{k \in \mathcal{K}^{(i)} \setminus \{I_t^{(i)}\} : d_t^{(i)}(k) > \hat{\omega}_t^{(i)}(I_t^{(i)})\}$ {choose arms with high KL-UCB}
 - 6: $\mathcal{D}_t^{(i)} \leftarrow \mathcal{D}_t^{(i)} \setminus \{I_t^{(j)} : \forall j \in \mathcal{M}\}$ {take advantage of free exploration}
 - 7: **if** $\mathcal{D}_t^{(i)} = \emptyset$ **then**
 - 8: $J_t^{(i)} \leftarrow I_t^{(i)}$
 - 9: **else**
 - 10: w.p., $\frac{1}{2}, J_t^{(i)} \leftarrow I_t^{(i)}$
 - 11: w.p., $\frac{1}{2}, J_t^{(i)} \leftarrow$ uniformly pick an arm from $\mathcal{D}_t^{(i)}$
 - 12: **end if**
 - 13: Pull arm $J_t^{(i)}$ and receive observations $X_t^{(i)}(J_t^{(i)})$
 - 14: Send observations $X_t^{(i)}(J_t^{(i)}) - \nu^{(i)}(J_t^{(i)})$ to other agents and also collect theirs
 - 15: Update $\hat{\omega}_t^{(i)}(k)$ and $d_t^{(i)}(k)$ for arm k and agent i
 - 16: **end for**
-

local optimal arm and thus need further explorations. To
 utilize free exploration, agent i periodically announces her
 estimated optimal arm $I_t^{(i)}$ to others to discourage other
 agents exploring this arm.

Remark 3.1. We note that some prior works [Combes and Proutiere, 2014, Combes et al., 2015, Wang et al., 2020a], such as the DPE2 algorithm in cooperative MA2B [Wang et al., 2020a], also involved a pivot arm and an exploration arm set in the algorithm design. However, the technical usage of both components in those works is very different from ours. For example, DPE2 estimates the pivot arm to gather all exploration responsibility to a single leader agent, while our usage is relegating/dispersing the free arms to the agents for which they are locally optimal.

Local optimal arm estimate and construction of exploration arm set: Let $n_t(k)$ and $\hat{\mu}_t(k)$ denote the total number of times arm k is pulled up to time t and the empirical mean of these $n_t(k)$ reward observations of arm k among all M agents. Denote $\hat{\omega}_t^{(i)}(k) := \hat{\mu}_t(k) + \nu^{(i)}(k)$ as the empirical reward mean of agent i pulling arm k and it is based on all agents' observations of arm k . FreeExp uses agent i 's *empirical local optimal arm* $I_t^{(i)}$ (the arm with the largest empirical reward mean $\hat{\omega}_t^{(i)}(k)$ of agent i at time t) as an estimate of the pivot. Given this empirical optimal arm as the pivot, the agent either pulls its own empirical optimal arm $I_t^{(i)}$ for free exploration, or explores other arms in $\mathcal{D}_t^{(i)}$ to guarantee the correctness of this estimated pivot. To improve the efficiency of exploring other arms, we con-

456 struct the *exploration arm set* $\mathcal{D}_t^{(i)}$ for each agent i using
 457 the KL-UCB index [Cappé et al., 2013]. The index of arm
 458 k at time slot t is

$$d_t^{(i)}(k) := \sup\{q \geq 0 : n_t(k) \text{kl}(\hat{\omega}_t^{(i)}(k), q) \leq \log t + 4 \log(\log t)\}, \quad (7)$$

459 where $\text{kl}(a, b)$ is the KL-divergence between two Gaussian
 460 distributions with means a and b and same variance $\sigma_1^2 +$
 461 σ_2^2 . The exploration arm set $\mathcal{D}_t^{(i)}$ includes arms whose KL-
 462 UCB indexes $d_t^{(i)}(k)$ are greater than the agent’s highest
 463 empirical mean $\hat{\omega}_t^{(i)}(I_t^{(i)})$ (Line 5) and excludes arms that
 464 are empirically optimal for at least one agent (Line 6)—
 465 discourage agent i exploring others’ local optimal arms.
 466 Note that the agents only share the arm-specific reward to
 467 other, i.e., the agent subtracts the agent-specific reward from
 468 the observed compound reward before sharing (Line 14).

469 **Arm pulling policy:** To guarantee the accuracy of the pivot
 470 estimation (i.e., the empirical optimal arm is correct with
 471 high probability), each agent needs to have enough observa-
 472 tions for her empirically optimal arm. To accomplish this,
 473 `FreeExp` implements an arm pulling policy (Lines 7-11)
 474 as follows: if exploration arm set $\mathcal{D}_t^{(i)}$ is empty, the agent i
 475 pulls the empirical optimal arm $I_t^{(i)}$; if exploration arm set
 476 $\mathcal{D}_t^{(i)}$ is not empty, with probability 1/2, the agent, uniformly
 477 at random picks an arm from $\mathcal{D}_t^{(i)}$ to explore; and with prob-
 478 ability 1/2, pulls her empirical optimal arm—encourage
 479 free explorations of the agent’s empirical optimal arm. This
 480 policy produces sufficient observations of this arm to guaran-
 481 tee fast correction if the current empirical optimal arm is not
 482 the correct one. Let $J_t^{(i)}$ denote the arm selected by agent i
 483 in time slot t under `FreeExp`. We present pseudocode for
 484 `FreeExp` in Algorithm 1.

485 **Remark 3.2** (`NoFreeExp` Algorithm). There is a coun-
 486 terpart algorithm of `FreeExp`, which does not utilize free
 487 exploration, i.e., Algorithm 1 without Line 6. We name it as
 488 `NoFreeExp`. Even without making use of free exploration,
 489 `NoFreeExp` should have a better regret performance than
 490 known baselines, e.g., `CO-UCB`, because `NoFreeExp` is
 491 based on the KL-UCB algorithm, which is theoretically
 492 better than UCB-like algorithms [Cappé et al., 2013].

4 THEORETICAL RESULTS

493 We present our theoretical results and their significance
 494 discussions in this section. The rigorous proofs of these
 495 results are deferred to Appendix C. We first derive a regret
 496 lower bound in Theorem 4.1 which reflects the impact of
 497 free exploration.

498 **Theorem 4.1** (Regret lower bound). *For any consistent*
 499 *policy π (i.e., for any bandit instance ν and any $\alpha > 0$, the*
 500 *policy π always guarantees $\mathbb{E}_{\nu, \pi}[\mathbf{R}_T] = O(T^\alpha)$), the regret*

cost of addressing the MA2B-HR model in T time slots is lower bounded by

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \geq \sum_{k: \bar{\Delta}(k) > 0} \frac{\bar{\Delta}(k)}{\text{kl}(\bar{\omega}(k), \bar{\omega}(k) + \bar{\Delta}(k))}, \quad (8)$$

where $\bar{\Delta}(k)$ defined in (5) is the smallest reward gap of pulling arm k and $\bar{\omega}(k)$ defined in (6) is the reward mean of pulling arm k by the agent who enjoys the smallest gap.

Theorem 4.1’s proof leverages similar techniques of the classic stochastic bandits [Lai et al., 1985]. Since $\bar{\Delta}(k) = 0$ for all free arms $k \in \mathcal{K}^{\text{fr}}$ and vice versa, the regret lower bound can be rewritten as

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \geq \sum_{k \in \mathcal{K} \setminus \mathcal{K}^{\text{fr}}} \frac{\bar{\Delta}(k)}{\text{kl}(\bar{\omega}(k), \bar{\omega}(k) + \bar{\Delta}(k))}. \quad (9)$$

Remark 4.2 (Free arms have no contribution to the asymptotic regret lower bound). Free arms in \mathcal{K}^{fr} contribute at most sub-logarithmic costs to the regret lower bound. In fact, given our finite regret upper bound of `FreeExp` next, free arms only contribute finite regret.

Theorem 4.3 (Regret upper bound for `FreeExp` (Algorithm 1)). *The `FreeExp` algorithm’s regret is upper bounded as follows,*

$$\begin{aligned} \mathbb{E}[\mathbf{R}_T(\mathcal{A})] &\leq 7bM^2K^2(4K + \delta^{-2}) \\ &+ \sum_{k: \bar{\Delta}(k) > 0} \frac{4(\bar{\Delta}(k) - 2\delta)(\log T + 4 \log(\log T))}{\text{kl}(\bar{\omega}(k) + \delta, \bar{\omega}(k) + \bar{\Delta}(k) - \delta)} \end{aligned} \quad (10)$$

where $0 < \delta < \frac{1}{4} \min_{i \in \mathcal{M}, k_1 \neq k_2 \in \mathcal{K}} |\omega^{(i)}(k_1) - \omega^{(i)}(k_2)|$, and that σ_1^2 and σ_2^2 are the variance of arm- and agent-specific Gaussian rewards respectively, and b is an upper bound of arm-specific reward mean $\mu(k)$ for all $k \in \mathcal{K}$.²

If we let $T \rightarrow \infty$ and $\delta \rightarrow 0$ (e.g., $\delta = (\log(\log T))^{-1}$), the above finite-time regret upper bound has the following asymptotical form,

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \leq O \left(\sum_{k: \bar{\Delta}(k) > 0} \frac{\bar{\Delta}(k)}{\text{kl}(\bar{\omega}(k), \bar{\omega}(k) + \bar{\Delta}(k))} \right). \quad (11)$$

Proof sketch and technical challenges. The proof of the regret upper bound in Theorem 4.3 consists of two steps: (i) bound the regret cost of pulling free arms in \mathcal{K}^{fr} , and (ii) other arms outside \mathcal{K}^{fr} . To bound (i), notice that for any free arm k in \mathcal{K}^{fr} , there exists “corresponding” agent(s) that takes arm k as its local optimal and can explore it with no cost. Hence, we only need to count the number of

²One can also obtain a near-optimal regret upper bound if the arm- and agent-specific rewards follow Bernoulli distributions.

532 times that arm k is pulled by agents other than “corresponding”
 533 one(s), which only happens when the “corresponding”
 534 agent’s empirical optimal arm $I_t^{(i)}$ is not its true local opti-
 535 mal arm $k_t^{(i)}$. Such events only occur with finite number
 536 of times even with a very large value of T . The proof of (i)
 537 shares the similar logical flow to that of [Wang et al., 2020b,
 538 Theorem 1]. To proof (ii), however, we need to develop
 539 new techniques for addressing the heterogeneous rewards
 540 in MA2B–HR. Note that in MA2B–HR the suboptimality re-
 541 ward gaps of pulling the same arm depend on the agents
 542 and thus are different. Hence, one cannot bound the cost of
 543 pulling a suboptimal arm k via multiplying the number of
 544 times of pulling the suboptimal arm k by one suboptimality
 545 reward gap as the usual bandits literature did. To address
 546 the challenge, we introduce two new techniques. First, we
 547 respectively count the number of times of the suboptimal
 548 arm pulls by agents (see Lemma C.7 and its proof), and
 549 secondly, we apply an Abel transformation to summing up
 550 the regret costs of all agents on pulling the arm k according
 551 to the order of magnitude of the arm’s reward gaps $\Delta^{(i)}(k)$
 552 for these agents (see Lemma C.8 and its proof).

553 Similar to the regret lower bound’s another expression in (9),
 554 this regret upper bound’s summation range can also be ex-
 555 pressed according to the free arms,

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \leq O \left(\sum_{k \in \mathcal{K} \setminus \mathcal{K}^{\text{fr}}} \frac{\bar{\Delta}(k)}{\text{kl}(\bar{\omega}(k), \bar{\omega}(k) + \bar{\Delta}(k))} \right). \quad (12)$$

556 **Remark 4.4** (Regret optimality of the FreeExp algorithm).
 557 This regret upper bound in (11) matches the regret lower
 558 bound in (8) up to a constant factor, which implies that
 559 both bounds are near-optimal, and therefore the FreeExp
 560 algorithm is near-optimal as well.

561 **Remark 4.5** (Comparison to Yang et al. [2022]’s regret
 562 bounds). Yang et al. [2022] proposed algorithms achieving
 563 regret upper bounds [Yang et al., 2022, Theorems 2 and 4]
 564 for AC–MA2B as follows (adapted to our notations),³

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[\mathbf{R}_T(\mathcal{A})]}{\log T} \leq O \left(\sum_{k \in \cup_{i \in \mathcal{M}} (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})} \frac{\bar{\Delta}(k)}{\text{kl}(\bar{\omega}(k), \bar{\omega}(k) + \bar{\Delta}(k))} \right).$$

565 Note that $\mathcal{K} = \cup_{i \in \mathcal{M}} \mathcal{K}^{(i)}$ and $\mathcal{K}^{\text{fr}} = \cup_{i \in \mathcal{M}} \{k_*^{(i)}\}$. So, we
 566 have $\mathcal{K} \setminus \mathcal{K}^{\text{fr}} \subset \cup_{i \in \mathcal{M}} (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})$. For example, if an

³To express Yang et al. [2022]’s result, we abuse $\bar{\Delta}(k)$ nota-
 tion *once*, where $\bar{\Delta}(k) := \min_{i \in \mathcal{M} \setminus \mathcal{M}_*(k)} \Delta^{(i)}(k)$ —the smallest
 reward mean gap of arm k compared to the local optimal arms
 (excluding arm k) among agents having access to k . The difference
 between this definition and the original one in (5) is that for arm k
 in \mathcal{K}^{fr} this $\bar{\Delta}(k)$ is positive while the original one is zero.

arm $k \in \mathcal{K}^{\text{fr}}$ is also a suboptimal arm for another agent,
 then $k \in \cup_{i \in \mathcal{M}} (\mathcal{K}^{(i)} \setminus \{k_*^{(i)}\})$ but $k \notin \mathcal{K} \setminus \mathcal{K}^{\text{fr}}$. In other
 words, the arm k contributes logarithmic regret costs to
 their upper bound but only contributes finite costs in ours.
 Therefore, their regret upper bound *failed to capture the*
advantage of free exploration and their algorithms did not
 utilize this appealing mechanism.

Remark 4.6 (Special cases with $O(1)$ finite regret in
 MA2B–HR). The regret upper bound in (12) echos the re-
 gret lower bound’s Remark 4.2 that arms in \mathcal{K}^{fr} only cause
 finite $O(1)$ costs in regret. Therefore, if all arms are local
 optimal for some agents, $\mathcal{K} \setminus \mathcal{K}^{\text{fr}} = \emptyset$ (e.g., the example in
 Table 1), then the regret upper bound in (11) becomes $O(1)$,
 i.e., a time horizon independent finite regret.

Remark 4.7 (Comparison to Baek and Farias [2021]). Re-
 call that the set of *free arms* \mathcal{K}^{fr} defined in our Eq.(4) con-
 tains arms that can be freely explored. In our regret upper
 bound, we show that FreeExp’s regret cost due to pulling
 arms in \mathcal{K}^{fr} is $O(1)$, while Baek and Farias [2021]’s regret
 bound was asymptotic with respect to $\log T$, implying that
 KL–UCB’s regret due to pulling arms in \mathcal{K}^{fr} was $o(\log T)$
 (the analysis in Baek and Farias [2021] upper bounds the
 cost for arm set \mathcal{K}^{fr} by $O(\log \log T)$).

Remark 4.8 (Generalization to the homogeneous reward
 setting). If all agents’ local arm sets are the same, then only
 one unique optimal arm can be freely explored (i.e., $|\mathcal{K}^{\text{fr}}| =$
 1) and all other arms would appear in the summation range in
 regret bounds (8) and (11). Then, both the regret upper and
 lower bounds reduce to the ones in classic MABs in Lai et al.
 [1985] (also the same as the optimal bounds of cooperative
 MA2B). This observation highlights the “*generality*” of our
 regret bounds and shows that FreeExp also works for the
 homogeneous reward setting.

5 NUMERICAL SIMULATIONS

Baselines: We report results of numerical experiments that
 compare FreeExp to three known cooperative algorithms
 that do not leverage free exploration: (1) CO–UCB and (2)
 CO–KLUCB, extensions of UCB and KLUCB algorithms to
 cooperative multi-agent scenarios proposed by Yang et al.
 [2022] and Baek and Farias [2021] respectively; and (3)
 NoFreeExp, a variant of FreeExp that does not make
 use of free exploration (see Remark 3.2).

Experimental setup: Unless otherwise specified, we con-
 sider a MA2B–HR model with $M = 25$ agents and $K = 50$
 arms. Each arm is associated with a Gaussian distribu-
 tion whose arm-specific mean $\mu(k) \in (0, 1)$ is chosen
 uniformly at random from the click-through-rates of Kag-
 gle’s *Ad-Click* dataset [Avito, 2015] and with variance $1/2$.
 We consider two special cases of agent-specific reward
 means: Case (1) $\nu^{(i)}(k)$ is either 0 or $-1 \forall k \in \mathcal{K}, i \in \mathcal{M}$
 (i.e., AC–MA2B [Yang et al., 2022, Baek and Farias, 2021]

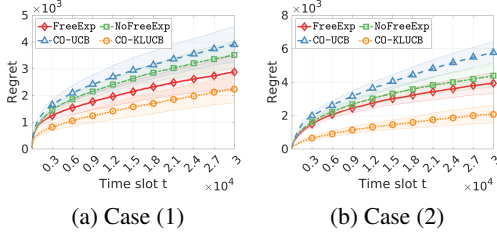


Figure 1: FreeExp vs. baselines

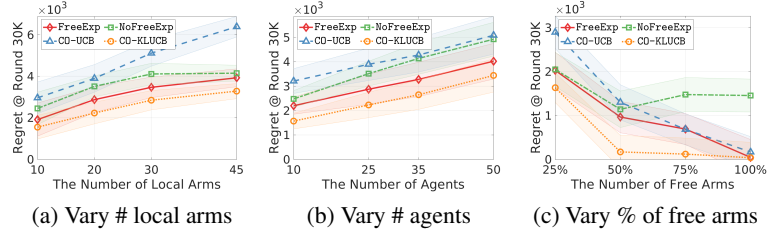


Figure 2: Vary parameters of MA2B-HR

617 where agents have different local arm sets) and Case (2)
 618 $\nu^{(i)}(k) \in (-1/2, 1/2) \forall k \in \mathcal{K}, i \in \mathcal{M}$ (i.e., all agents
 619 have the same local arm sets) as the more general heteroge-
 620 neous reward scenario. The variances of all agent-specific
 621 rewards are set to $1/2$. In the AC-MA2B setting (Case (1)),
 622 for each agent, we randomly select 20 of these 50 arms and
 623 set their agent-specific rewards $\nu^{(i)}(k) = 0$, i.e., as local
 624 arms. The remaining arms' agent-specific rewards is set to
 625 $\nu^{(i)}(k) = -1$. In the heterogeneous reward setting (Case
 626 (2)), all agents have the same 50 arms but different agent-
 627 specific rewards whose means are uniformly and randomly
 628 generated between $(-1/2, 1/2)$ for each arm and agent. All
 629 simulations are averaged over 50 runs and their standard
 630 deviations are plotted as shadow regions.

631 **Experimental results:** In Figures 1a and 1b, we compare
 632 the cumulative regret of all algorithms in Cases (1) and (2).
 633 The notable observations are: (1) Comparison of FreeExp
 634 to NoFreeExp shows that utilizing the free exploration
 635 mechanism can further improve an algorithm's performance.
 636 (2) The KLUCB algorithm outperform our FreeExp algo-
 637 rithm. This is because FreeExp needs to explicitly exclude
 638 arms likely to be local optimal (Line 6) and thus suffers a
 639 high time-independent cost at the beginning, while KLUCB
 640 does not; and the additional cost of FreeExp cannot be
 641 compensated by the advantage of FreeExp in saving cost
 642 on free arms in these two scenarios. Especially, we note that
 643 when the number of free arms are large (e.g., see Figure 2c's
 644 100% free arm case below), the advantage of FreeExp in
 645 saving cost on free arms becomes significant and, therefore,
 646 FreeExp has similar performance to KLUCB.

647 We report the results of varying the number of parameters of
 648 MA2B-HR (Case (1)) in Figure 2. In Figure 2a, we vary the
 649 number of local arms between 10 and 45 and report their cumu-
 650 lative regret at round 30K. All algorithm regrets increase
 651 linearly with respect to the number of local arms. Figure 2b
 652 shows the impact of the number of agents M (from 10
 653 agents to 50) on the regrets. Their regrets also have linear
 654 increasing rate in M , which is due to the fixed per-agent
 655 costs (independent of T). Lastly, we consider an MA2B-HR
 656 consisting of $M = 20$ agents and $K = 20$ arms, and devise
 657 four cases containing $\{5, 10, 15, 20\}$ free arms respectively
 658 (i.e., 25%, 50%, 75%, 100% of all arms are free arms). We
 659 report their regret performance in Figure 2c. The notable

660 observations are: (1) The regret of FreeExp decreases as
 661 the percentage of free arms increases which corroborates
 662 that FreeExp saves the costs due to pulling free arms. (2)
 663 when all (100%) arms are free, FreeExp has similar per-
 664 formance to KLUCB and outperforms other algorithms.

6 CONCLUSION

665 This paper introduced a multi-agent multi-armed bandit
 666 problem with heterogeneous rewards among agents. The
 667 heterogeneous scenario creates a unique opportunity to ex-
 668 plore a subset of arms for free and share the observation
 669 by cooperation, and hence, improve the aggregate regret
 670 significantly. We proposed a cooperative learning algorithm
 671 which would benefit from the free exploration and its regret
 672 is tight up to a constant factor. As a notable special case,
 673 when each arm is a local optimal arm in at least one agent,
 674 the proposed algorithm achieves an $O(1)$ regret.

675 This problem of multi-agent bandits with heterogeneous
 676 reward calls for several interesting follow-up questions, i.e.,
 677 an interesting question is to extend the FreeExp algorithm
 678 with an effective communication protocol. In a distributed
 679 multi-agent setting, cooperation may come with a cost of
 680 communication, and hence the goal is to enhance the coop-
 681 erative algorithms with a communication policies that only
 682 needs sublinear communication times w.r.t. decision rounds
 683 T , while directly extend current algorithm requires $O(T)$
 684 communication times.

Acknowledgements

685 The work of Mohammad Hajiesmaili is supported by
 686 NSF CAREER-2045641, CPS-2136199, CNS-2106299,
 687 and CNS-2102963. The work of Don Towsley is supported
 688 by U.S. Army Research Laboratory under Cooperative
 689 Agreement W911NF-17-2-0196. The work of John C.S. Lui
 690 is supported in part by the RGC GRF 14215722. Lin Yang
 691 is the corresponding author (linyang@nju.edu.cn).
 692

693 References

- 694 Peter Auer. Using confidence bounds for exploitation-
695 exploration trade-offs. *Journal of Machine Learning*
696 *Research*, 3(Nov):397–422, 2002.
- 697 Avito. *Avito Context Ad Clicks*, 2015. <https://www.kaggle.com/c/avito-context-ad-clicks>.
698
- 699 Maryam Aziz, Emilie Kaufmann, and Marie-Karelle Riviere.
700 On multi-armed bandit designs for dose-finding clinical
701 trials. *Journal of Machine Learning Research*, 22(1-38):
702 4, 2021.
- 703 Jackie Baek and Vivek Farias. Fair exploration via axiomatic
704 bargaining. *Advances in Neural Information Processing*
705 *Systems*, 34:22034–22045, 2021.
- 706 Etienne Boursier and Vianney Perchet. Sic-mmab: synchro-
707 nisation involves communication in multiplayer multi-
708 armed bandits. *Advances in Neural Information Process-*
709 *ing Systems*, 32, 2019.
- 710 Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret anal-
711 ysis of stochastic and nonstochastic multi-armed bandit
712 problems. *Foundations and Trends® in Machine Learn-*
713 *ing*, 5(1):1–122, 2012.
- 714 Olivier Cappé, Aurélien Garivier, Odalric-Ambrym Mail-
715 lard, Rémi Munos, and Gilles Stoltz. Kullback-leibler
716 upper confidence bounds for optimal sequential alloca-
717 tion. *The Annals of Statistics*, pages 1516–1541, 2013.
- 718 Bangrui Chen, Peter Frazier, and David Kempe. Incentiviz-
719 ing exploration by heterogeneous users. In *Conference*
720 *On Learning Theory*, pages 798–818. PMLR, 2018.
- 721 Richard Combes and Alexandre Proutiere. Unimodal band-
722 its: Regret lower bounds and optimal algorithms. In
723 *International Conference on Machine Learning*, pages
724 521–529. PMLR, 2014.
- 725 Richard Combes, Stefan Magureanu, Alexandre Proutiere,
726 and Cyrille Laroche. Learning to rank: Regret lower
727 bounds and efficient algorithms. In *Proceedings of the*
728 *2015 ACM SIGMETRICS International Conference on*
729 *Measurement and Modeling of Computer Systems*, pages
730 231–244, 2015.
- 731 Yuval Dagan and Crammer Koby. A better resource alloca-
732 tion algorithm with semi-bandit feedback. In *Algorithmic*
733 *Learning Theory*, pages 268–320. PMLR, 2018.
- 734 Ting He, Dennis Goeckel, Ramya Raghavendra, and Don
735 Towsley. Endhost-based shortest path routing in dynamic
736 networks. In *Proc. IEEE INFOCOM*, pages 2202–2210,
737 2013.
- Eshcar Hillel, Zohar S Karnin, Tomer Koren, Ronny Lempel,
and Oren Somekh. Distributed exploration in multi-armed
bandits. *Advances in Neural Information Processing*
Systems, 26, 2013.
- Ruiquan Huang, Weiqiang Wu, Jing Yang, and Cong Shen.
Federated linear contextual bandits. *Advances in Neural*
Information Processing Systems, 34, 2021.
- Ramesh Johari, Vijay Kamble, and Yash Kanoria. Matching
while learning. In *Proceedings of the 2017 ACM Con-*
ference on Economics and Computation, pages 119–119,
2017.
- Nikolai Karpov, Qin Zhang, and Yuan Zhou. Collaborative
top distribution identifications with limited interaction.
In *2020 IEEE 61st Annual Symposium on Foundations of*
Computer Science (FOCS), pages 160–171. IEEE, 2020.
- Tze Leung Lai. Adaptive treatment allocation and the multi-
armed bandit problem. *The Annals of Statistics*, pages
1091–1114, 1987.
- Tze Leung Lai, Herbert Robbins, et al. Asymptotically
efficient adaptive allocation rules. *Advances in applied*
mathematics, 6(1):4–22, 1985.
- Peter Landgren, Vaibhav Srivastava, and Naomi Ehrich
Leonard. Distributed cooperative decision-making in
multiarmed bandits: Frequentist and bayesian algorithms.
In *2016 IEEE 55th Conference on Decision and Control*
(CDC), pages 167–172. IEEE, 2016.
- Tor Lattimore, Koby Crammer, and Csaba Szepesvári. Op-
timal resource allocation with semi-bandit feedback. In
Proceedings of the Thirtieth Conference on Uncertainty
in Artificial Intelligence, pages 477–486, 2014.
- Lihong Li, Wei Chu, John Langford, and Robert E Schapire.
A contextual-bandit approach to personalized news arti-
cle recommendation. In *Proceedings of the 19th inter-*
national conference on World wide web, pages 661–670,
2010.
- Kanak Mahadik, Qingyun Wu, Shuai Li, and Amit Sabne.
Fast distributed bandits for online recommendation sys-
tems. In *Proceedings of the 34th ACM international*
conference on supercomputing, pages 1–13, 2020.
- David Martínez-Rubio, Varun Kanade, and Patrick Rebes-
chini. Decentralized cooperative stochastic bandits. *Ad-*
vances in Neural Information Processing Systems, 32,
2019.
- Abbas Mehrabian, Etienne Boursier, Emilie Kaufmann, and
Vianney Perchet. A practical algorithm for multiplayer
bandits when arm means vary among players. In *Interna-*
tional Conference on Artificial Intelligence and Statistics,
pages 1211–1221. PMLR, 2020.

- 786 Chengshuai Shi and Cong Shen. Federated multi-armed
787 bandits. In *Proceedings of the 35th AAAI Conference on*
788 *Artificial Intelligence (AAAI)*, 2021.
- 789 Chengshuai Shi, Cong Shen, and Jing Yang. Federated
790 multi-armed bandits with personalization. In *International*
791 *Conference on Artificial Intelligence and Statistics*,
792 pages 2917–2925. PMLR, 2021a.
- 793 Chengshuai Shi, Wei Xiong, Cong Shen, and Jing Yang. Het-
794 erogeneous multi-player multi-armed bandits: Closing the
795 gap and generalization. *Advances in Neural Information*
796 *Processing Systems*, 34, 2021b.
- 797 Chengshuai Shi, Haifeng Xu, Wei Xiong, and Cong Shen.
798 (almost) free incentivized exploration from decentralized
799 learning agents. *Advances in Neural Information Process-*
800 *ing Systems*, 34:560–571, 2021c.
- 801 Mohammad Sadegh Talebi and Alexandre Proutiere. Learn-
802 ing proportionally fair allocations with low regret. *Pro-*
803 *ceedings of the ACM on Measurement and Analysis of*
804 *Computing Systems*, 2(2):1–31, 2018.
- 805 Mohammad Sadegh Talebi, Zhenhua Zou, Richard Combes,
806 Alexandre Proutiere, and Mikael Johansson. Stochastic
807 online shortest path routing: The value of feedback. *IEEE*
808 *Transactions on Automatic Control*, 63(4):915–930, 2017.
- 809 Liang Tang, Yexi Jiang, Lei Li, and Tao Li. Ensemble
810 contextual bandits for personalized recommendation. In
811 *Proceedings of the 8th ACM Conference on Recommender*
812 *Systems*, pages 73–80, 2014.
- 813 Chao Tao, Qin Zhang, and Yuan Zhou. Collaborative learn-
814 ing with limited interaction: Tight bounds for distributed
815 exploration in multi-armed bandits. In *2019 IEEE 60th*
816 *Annual Symposium on Foundations of Computer Science*
817 *(FOCS)*, pages 126–146. IEEE, 2019.
- 818 Sofia S Villar, Jack Bowden, and James Wason. Multi-armed
819 bandit models for the optimal design of clinical trials:
820 benefits and challenges. *Statistical science: a review*
821 *journal of the Institute of Mathematical Statistics*, 30(2):
822 199, 2015.
- 823 Po-An Wang, Alexandre Proutiere, Kaito Ariu, Yassir Jedra,
824 and Alessio Russo. Optimal algorithms for multiplayer
825 multi-armed bandits. In *International Conference on*
826 *Artificial Intelligence and Statistics*, pages 4120–4129.
827 PMLR, 2020a.
- 828 Yuanhao Wang, Jiachen Hu, Xiaoyu Chen, and Liwei Wang.
829 Distributed bandit learning: Near-optimal regret with ef-
830 ficient communication. In *8th International Conference*
831 *on Learning Representations, ICLR 2020, Addis Ababa,*
832 *Ethiopia, April 26-30, 2020*, 2020b.
- Lin Yang, Yu-Zhen Janice Chen, Mohammad Hajiesmaili, 833
John C.S. Lui, and Don Towsley. Distributed bandits with 834
heterogeneous agents. In *Proceedings of The IEEE* 835
International Conference on Computer Communications 836
2022, 2022. 837
- Zhaowei Zhu, Jingxuan Zhu, Ji Liu, and Yang Liu. Fed- 838
erated bandit: A gossiping approach. In *Abstract Pro-* 839
ceedings of the 2021 ACM SIGMETRICS/International 840
Conference on Measurement and Modeling of Computer 841
Systems, pages 3–4, 2021. 842
- Zhenhua Zou, Alexandre Proutiere, and Mikael Johansson. 843
Online shortest path routing: The value of information. 844
In *2014 American Control Conference*, pages 2142–2147. 845
IEEE, 2014. 846